

Hardware Accelerator for the Training of Neural Networks

James Erik Groving Meade
s173476



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.

GitHub Repository Organization

As the scope of this thesis was fairly wide, several different technologies and programming languages were used in the project. As such, this section was added to help guide the reader around the organization of the GitHub repository¹. The organization of the GitHub repository is present on the homepage of the repository. It may also be read in textform by viewing the README.md file.

¹<https://github.com/erikgroving/NeuralNetworkHardwareAccelerator/>

Contents

Abstract	i
Preface	iii
Acknowledgements	v
GitHub Repository Organization	vii
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	23
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	37
4.7.3	Softmax Layer	39
4.8	PS – FPGA Communication	40
4.8.1	AXI Implementation for the PS	40
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	Experiment Data	81
B	Hardware Model Code	85
B.1	Source Files	85
B.1.1	sys_defs.vh	85
B.1.2	neural_net_top.sv	87
B.1.3	fc0_layer.sv	105
B.1.4	fc0_weight_bram_controller.sv	113
B.1.5	fc_kernel.sv	114
B.1.6	fc_scheduler.sv	117
B.1.7	interlayer_activation_buffer.sv	119
B.1.8	fc1_layer.sv	121
B.1.9	previous_layer_gradient_adder.sv	130
B.1.10	fc1_weight_bram_controller.sv	133
B.1.11	fc2_layer.sv	135
B.1.12	fc2_weight_bram_controller.sv	143
B.1.13	softmax.sv	145
B.2	Testbenches	148
B.2.1	neural_net_top_tb.sv	148
B.2.2	softmax_tb.sv	149
B.2.3	fc1_scheduler_tb.sv	151
B.2.4	fc1_layer_tb.sv	153

C Processing System Code	155
C.1 train.c	155
C.2 inference_only.c	160
C.3 parse_mnist.c	163
D Hardware Testing Code	165
D.1 fpga_forward_backward_pass_test.py	165
E Weight Generation	173
E.1 weight_coeff.py	173
F Software Model	175
F.1 Header Files	175
F.1.1 net.h	175
F.1.2 layer.h	177
F.1.3 convolutional.h	178
F.1.4 fullyconnected.h	179
F.1.5 pooling.h	180
F.1.6 neuron.h	181
F.1.7 parse_data.h	182
F.2 Source Files	183
F.2.1 main.cpp	183
F.2.2 net.cpp	187
F.2.3 convolutional.cpp	193
F.2.4 fullyconncted.cpp	197
F.2.5 pooling.cpp	199
F.2.6 neuron.cpp	201
F.2.7 parse_data.cpp	203
F.3 Testing Files	205
F.3.1 gradient_check_test.cpp	205
F.3.2 conv_test.cpp	209
F.3.3 fullyconnected_test.cpp	210
F.3.4 neuron_test.cpp	211
G PyTorch Model	213
G.1 mnist_model.py	213
G.2 parse_data.py	216
Bibliography	217

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

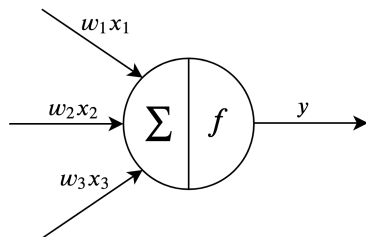


Figure 2.1: A neuron with 3 inputs; bias term omitted for simplicity.

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

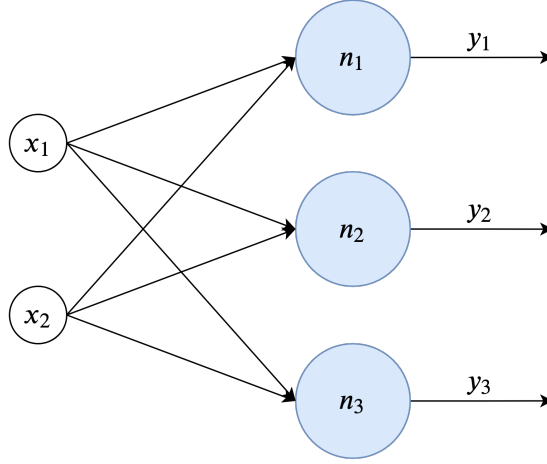


Figure 2.2: A fully-connected layer with 3 neurons, each receiving an input vector of size 2 from the input layer.

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.

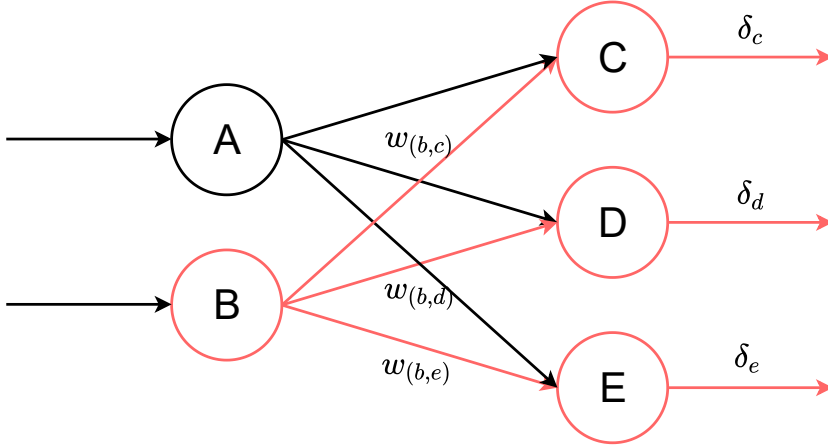


Figure 2.3: Example of backpropagating the loss gradient to the previous layer. Values used in backpropagating the loss to neuron B shown in red.

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.

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

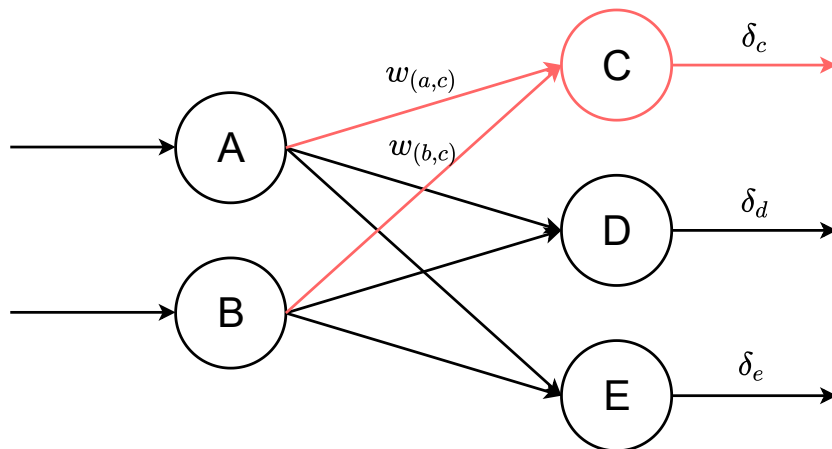


Figure 2.4: Example of computing weight gradients. Relevant values shown in red.

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.

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 Appendix F as well as online on GitHub in the SWModel folder.

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

Name	Type	Description
in	uint32_t	Size of the input to the neural network.
out	uint32_t	Size of the output of the neural network.
bs	uint32_t	Size of the batch size to be used when training the net.
lr	double	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> .
momentum	double	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.

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

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.

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 `dim_i` \times `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 confi-

guration 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, 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 Appendix F.

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

```

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.

(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

```

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

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 3×3 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.

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 ¹. 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

```

¹<https://github.com/google/googletest>

```
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.

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

```
> ./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. The source code can be found in the *FPGA* folder of the GitHub repository or in Appendix B.

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.

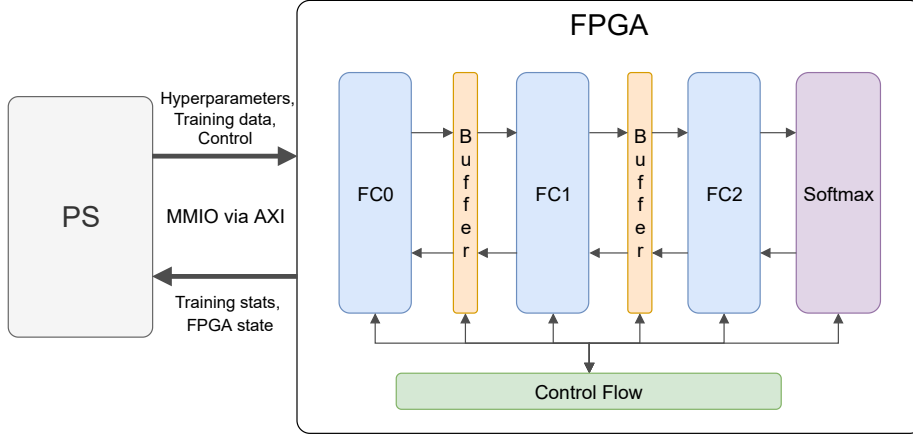


Figure 4.1: Architecture of the hardware accelerator

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.

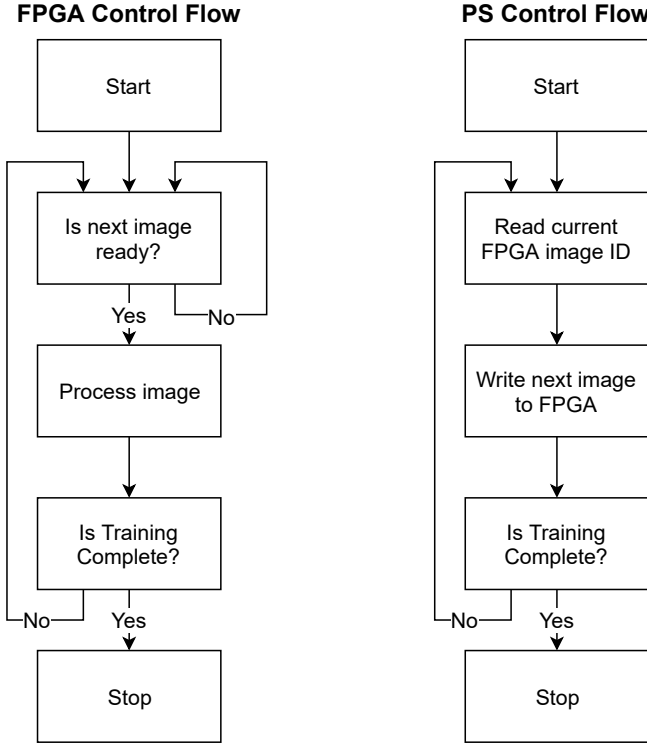


Figure 4.2: The high-level control flow of the training process

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 B, 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.

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

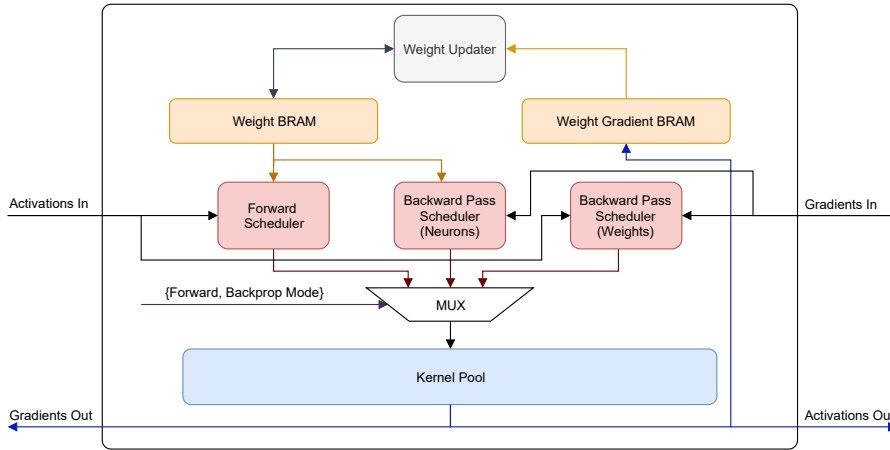


Figure 4.3: Architecture of the fully connected layer

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 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

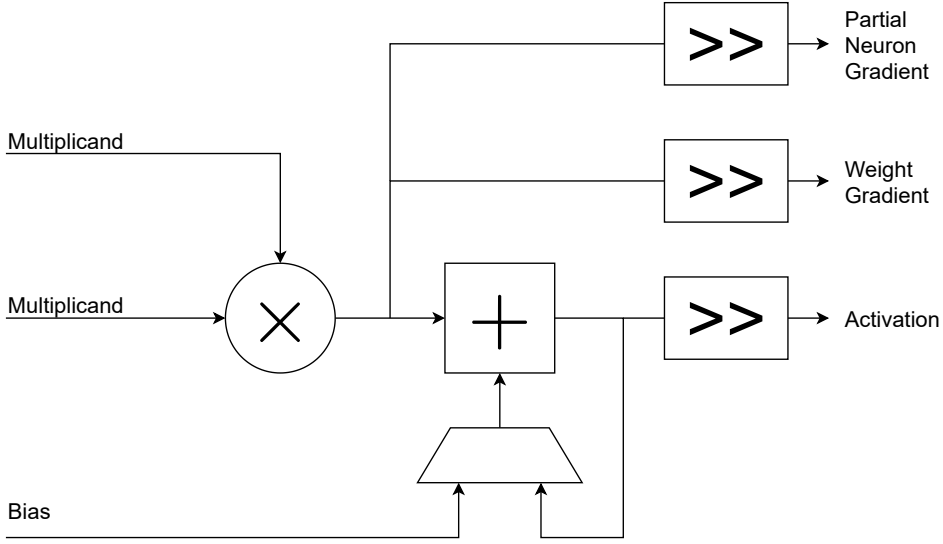


Figure 4.4: Architecture of the kernel, saturation checking not shown.

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.

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.

Layer	# Kernels
FC0	196
FC1	16
FC2	2
Total	214

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

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

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.

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

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.

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

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

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.

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

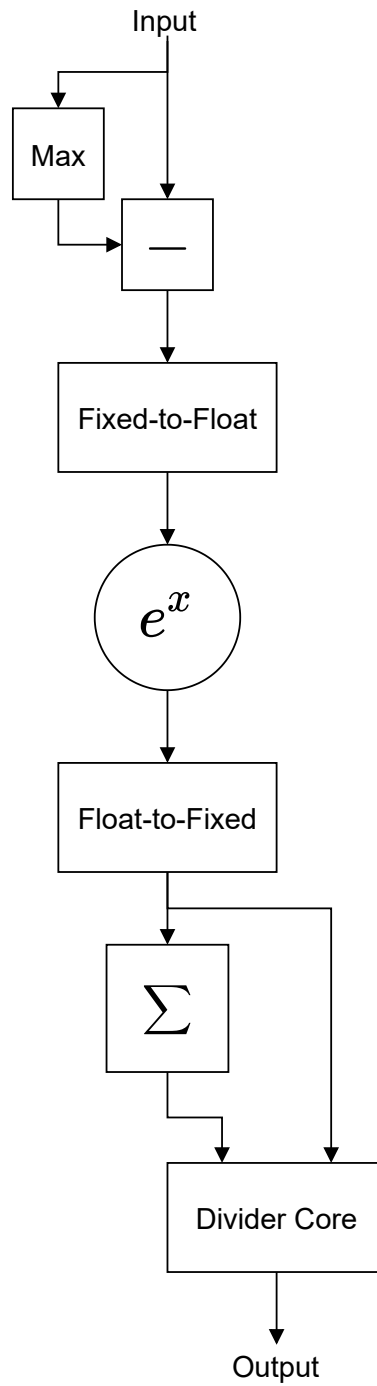


Figure 4.6: Architecture of the softmax layer

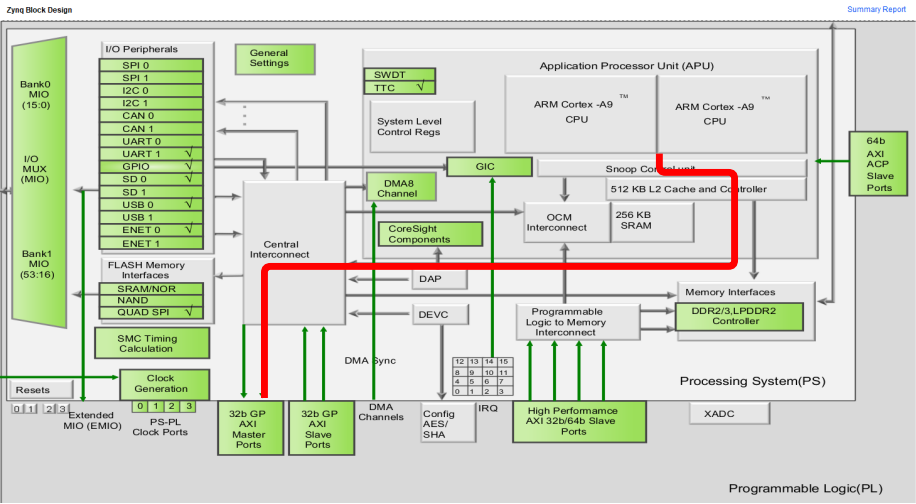


Figure 4.7: Communication from the PS to the FPGA

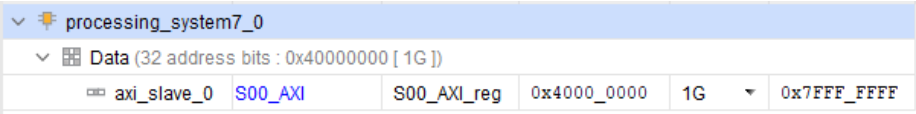


Figure 4.8: Specifying the Address Range for the AXI Bus

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 Appendix C 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 1000 Loading MNIST images...
2 1000 Loading complete!
3
4 1000 EPOCH 1
5 1000 Training Images: 2950/5000
6 Accuracy: 58.999997%
7 1000Test 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 1000 EPOCH 2
14 1000 Training Images: 3623/5000
15 Accuracy: 72.460002%
16 1000Test 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.

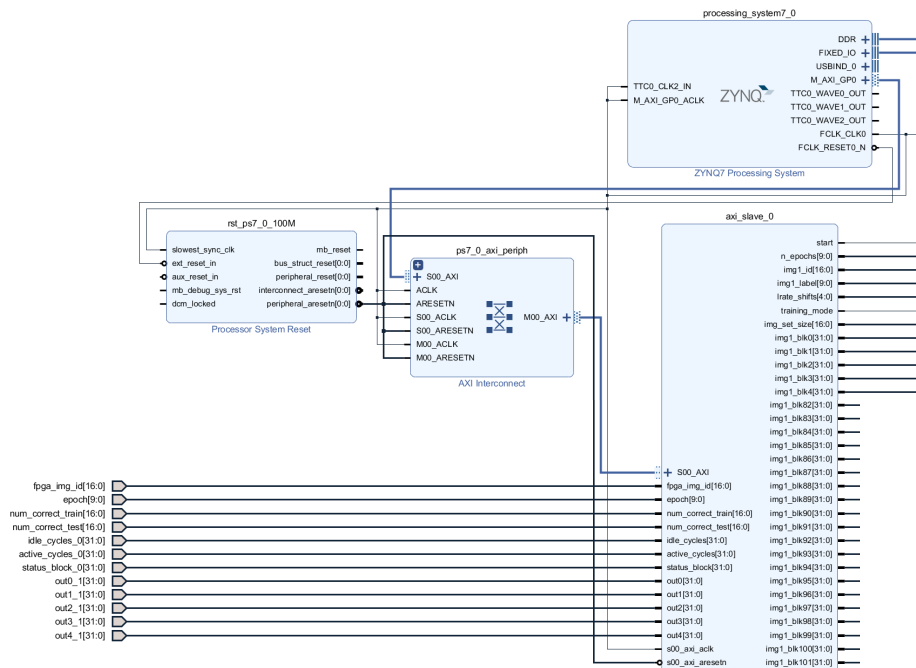


Figure 4.9: Block diagram to establish connection from the FPGA side to the AXI bus from the PS

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 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 can be found on GitHub in the *ip_repo/axi_slave_1.0/hdl* folder.

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.

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

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.

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 B.

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 D 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

```

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

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. All experiment data can be found in Appendix A.

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.

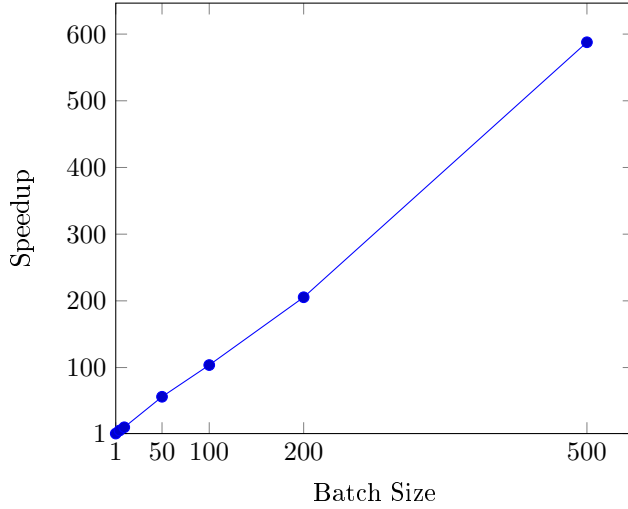


Figure 6.1: Speedup for 1 training epoch when training using different batch sizes for PyGPU. Speedup is calculated using the PyGPU model with a batch size of 1 as the reference.

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.

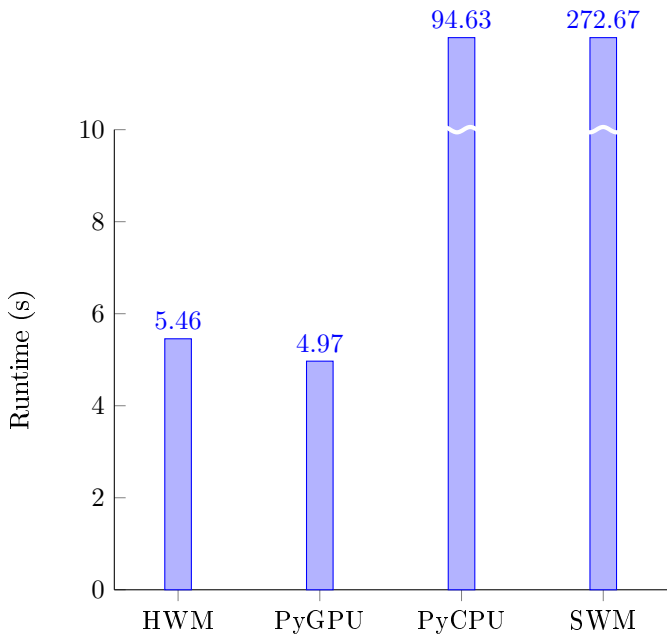


Figure 6.2: Training runtime for various network models

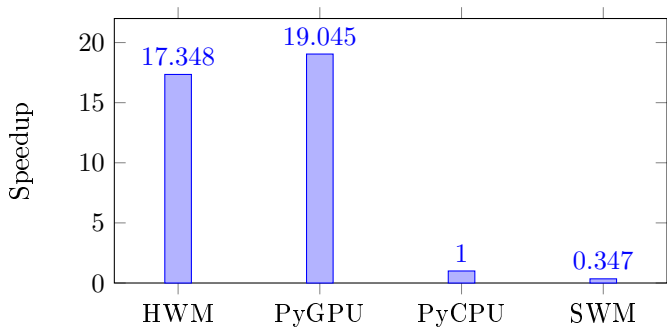


Figure 6.3: Training speedup using the PyCPU as a baseline

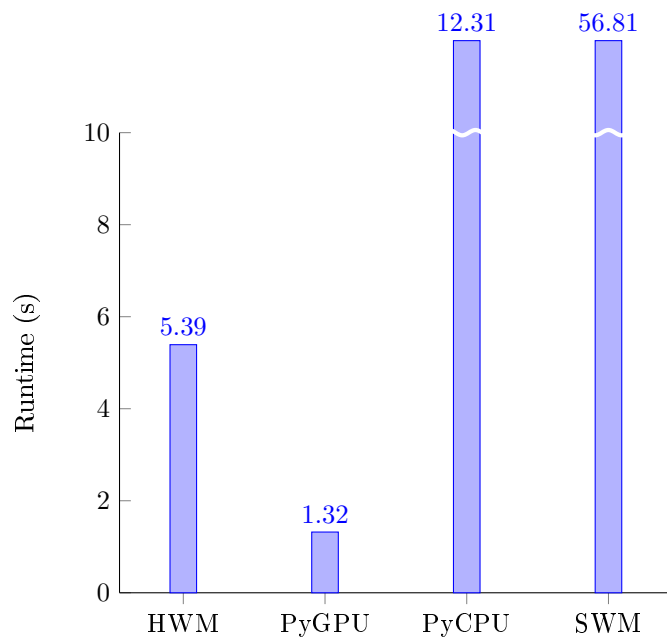


Figure 6.4: Inference runtime for various network models

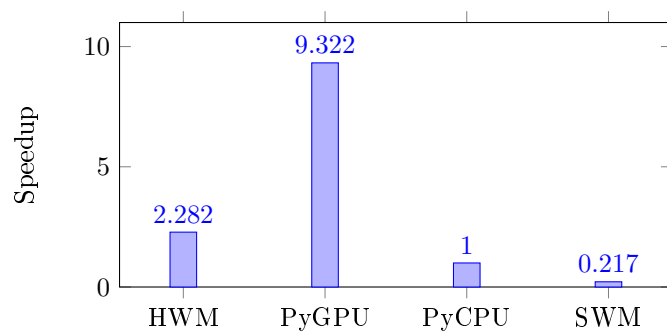


Figure 6.5: Inference speedup using the PyCPU as a baseline

	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

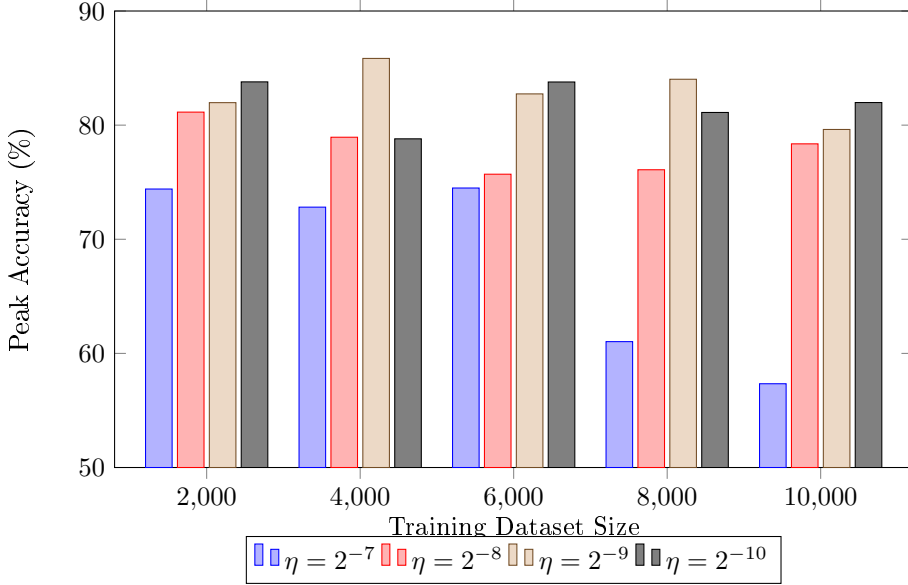


Figure 6.6: Maximum training accuracy reached for various learning rate and training set sizes.

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.

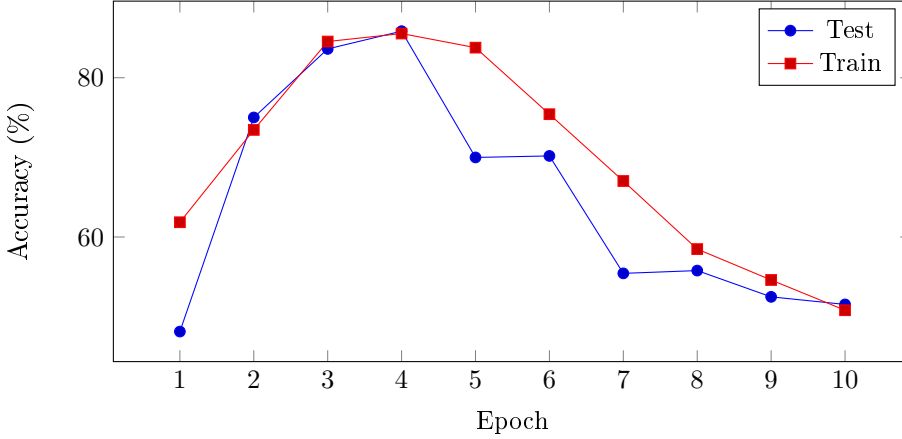


Figure 6.7: Epoch-by-epoch training data for an HWM configuration. Clearly visible degradation of accuracy instead of convergence after epoch 4.

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

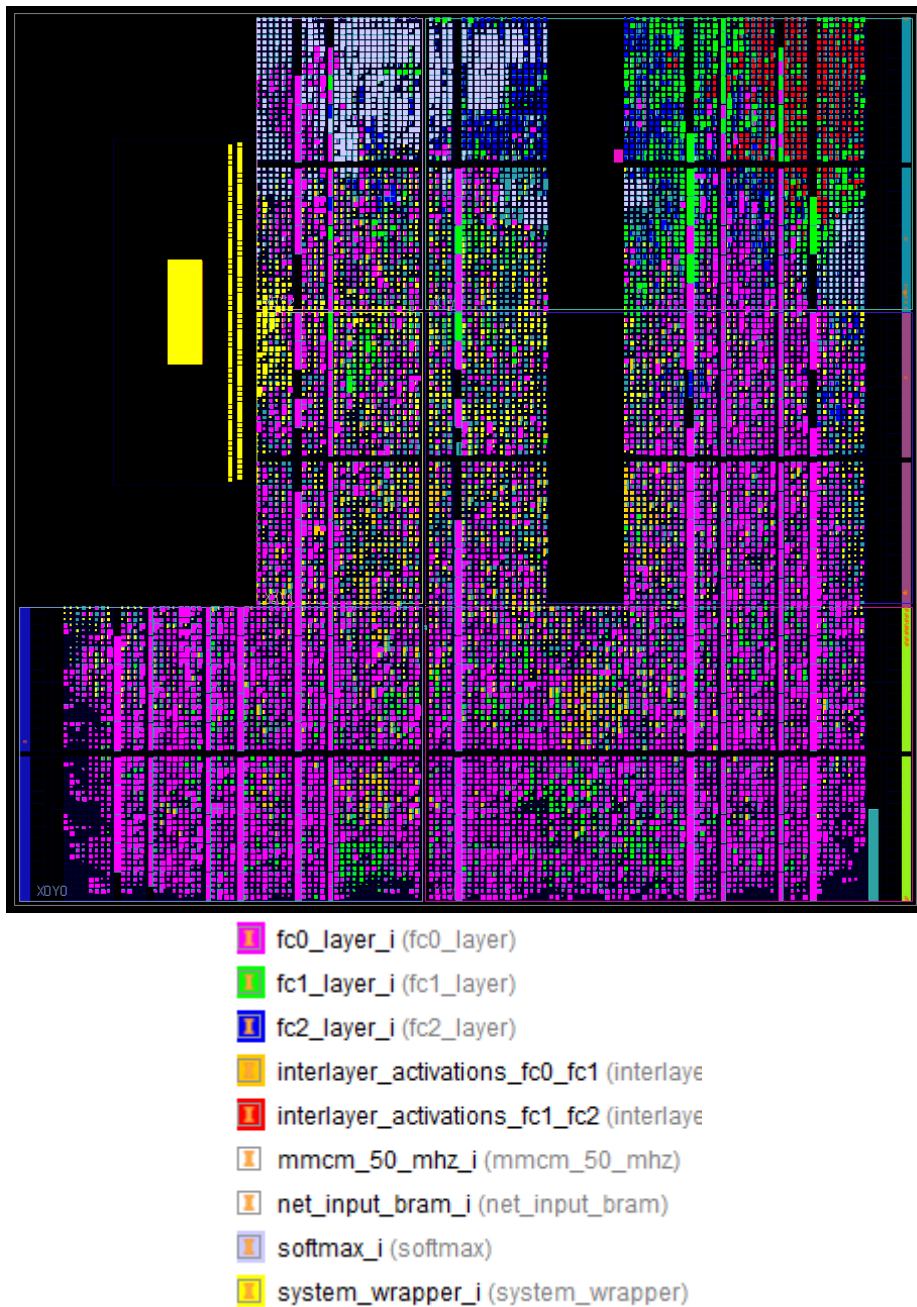


Figure 6.8: The implemented design of the hardware model

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.

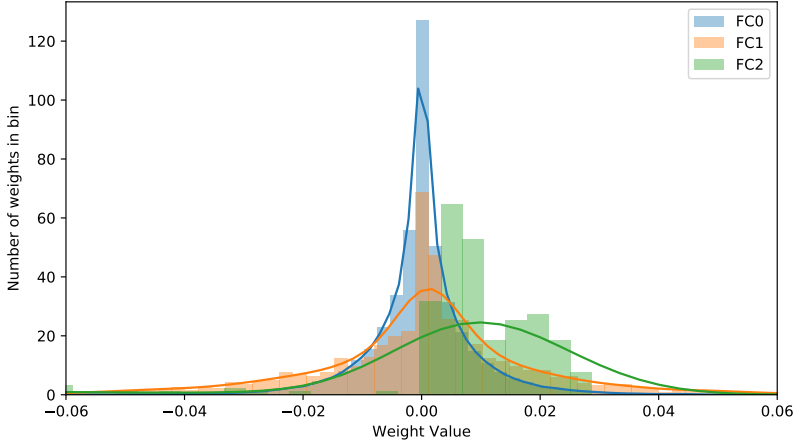


Figure 7.1: Weight gradient distribution for the 3 fully-connected layers after performing a single backward pass on the network using the input used to verify the hardware model from Chapter 5.

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

Experiment Data

Appendix A contains the experiment data used in the Results chapter of the thesis.

Training: Time Per Epoch (60,000 images)

Model	Runtime (s)	Speedup
HWM	5.455	17.348
SWM	272.67	0.347
PyCPU	94.633	1
PyGPU (Batch size: 50)	4.969	19.05

Training: Time Per Epoch. Model: PyGPU. Varying Batch Sizes
(60,000 images)

Batch Size	Runtime (s)	Speedup
1	279.14	1
5	49.892	5.594965125
10	26.747	10.43646016
50	4.9689	56.17822858
100	2.6911	103.7285868
200	1.359	205.4039735
500	0.479	582.7640919

Inference Performance (70,000 images)

Model	Runtime (s)	Speedup
HWM	5.392	2.28208457
SWM	56.808	0.216606816
PyCPU	12.305	1
PyGPU (Batch size: 50)	1.32	9.321969697

Epoch-by-epoch Data for Unstable Training (70,000 images)

Epoch	Test Accuracy (%)	Train Accuracy (%)
1	48.12	61.85
2	75.012	73.45
3	83.62	84.55
4	85.85	85.55
5	69.99	83.78
6	70.18	75.42
7	55.43	67.03
8	55.78	58.48
9	52.48	54.6
10	51.53	50.8

Maximum Accuracy Results for Varying Learning Rates and Training Dataset Sizes

LRate (Shifts)	Dataset Size	Top Test (%)	Train (%)
7	1000	74.779	74.2
7	2000	74.401	77.5
7	3000	74.061	73.866
7	4000	72.814	65.85
7	5000	79.351	68.06
7	6000	74.49	68.78
7	7000	70.641	69.957
7	8000	61.03	69.22
7	9000	60.957	68.077
7	10000	57.34	53.4
8	1000	81.714	87.9
8	2000	81.14	84.25
8	3000	82.696	79.5
8	4000	78.94	77.82
8	5000	78.64	79.68
8	6000	75.698	70.1
8	7000	74.581	71.629
8	8000	76.088	72.47
8	9000	71.04	72.77
8	10000	78.36	73.21
9	1000	67.01	66.3
9	2000	81.97	84.95
9	3000	79.31	79.57
9	4000	85.845	85.549
9	5000	85.789	85.26
9	6000	82.737	84.2
9	7000	85.56	85.65
9	8000	84.02	85.69
9	9000	83.6	85.2
9	10000	79.615	84.34
10	1000	74.4	84.2
10	2000	83.79	86.5
10	3000	83.19	82.87
10	4000	78.8	77.72
10	5000	85.3	83.44
10	6000	83.78	81.78
10	7000	81.3	82.94
10	8000	81.11	84.425
10	9000	79.1	83.2
10	10000	81.98	80.59

APPENDIX B

Hardware Model Code

This appendix contains the SystemVerilog code used to implement the hardware accelerator. Display statements have been omitted as they serve no functional purpose to the implementation nor is there any particular insight to be gained from them. They may of course still be viewed at the GitHub repository. Note that some style has been modified make code fit on the page.

B.1 Source Files

B.1.1 sys_defs.vh

```
1 `ifndef __SYS_DEFS_VH__
2 `define __SYS_DEFS_VH__
3
4 // Precision defines
5 `define PREC 18
6 `define MULT_BITS 36
7 `define ACT_INT_BITS 6
8 `define ACT_FRAC_BITS 12
9 `define GRAD_INT_BITS 1
10 `define GRAD_FRAC_BITS 17
11 `define ONE 18'h1_ffff
12 `define MAX_VAL 18'h1_ffff
```

```

13 `define MIN_VAL                18'h2_0000
14
15 // FC0 defines
16 `define FC0_N_KERNELS          196
17 `define FC0_PORT_WIDTH         98
18 `define FC0_NEURONS             98
19 `define FC0_FAN_IN              10'd784
20 `define FC0_KERNEL_FAN_IN      10'd392
21 `define FC0_MID_PTR_OFFSET     10'd784
22 `define FC0_ADDR                10
23 `define FC0_BIAS_ADDR           1
24
25 // FC1 defines
26 `define FC1_N_KERNELS           16
27 `define FC1_ADDR                 10
28 `define FC1_PORT_WIDTH           8
29 `define FC1_PORT_WIDTH_TIMES2    16
30 `define FC1_PORT_WIDTH_TIMES3    24
31 `define FC1_BRAM                  1
32 `define FC1_NEURONS              64
33 `define FC1_BIAS_ADDR             2
34 `define FC1_FAN_IN               10'd98
35 `define FC1_STEP2                10'd196
36 `define FC1_STEP3                10'd294
37 `define FC1_MID_PTR_OFFSET       10'd392
38 `define FC1_MID_PTR_END          10'd784
39 `define FC1_HALF_NEURONS         32
40
41 // FC2 defines
42 `define FC2_BRAM                  1
43 `define FC2_NEURONS              10
44 `define FC2_FAN_IN               64
45 `define FC2_N_KERNELS            2
46 `define FC2_ADDR                 10
47 `define FC2_BIAS_ADDR             3
48 `define FC2_MID_PTR_OFFSET       320
49 `define FC2_HALF_NEURONS         5
50
51
52 // Backward pass defines
53 `define FC0_LOOPS                 1
54
55 `define FC1_MODE_SWITCH           4
56 `define FC1_LOOPS                 8
57
58 `define FC2_MODE_SWITCH           5
59 `define FC2_LOOPS                 10
60
61 `endif

```

B.1.2 neural_net_top.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3
4  module neural_net_top(
5      inout                [14:0]DDR_addr ,
6      inout                [2:0]DDR_ba ,
7      inout                DDR_cas_n ,
8      inout                DDR_ck_n ,
9      inout                DDR_ck_p ,
10     inout                DDR_cke ,
11     inout                DDR_cs_n ,
12     inout                [3:0]DDR_dm ,
13     inout                [31:0]DDR_dq ,
14     inout                [3:0]DDR_dqs_n ,
15     inout                [3:0]DDR_dqs_p ,
16     inout                DDR_odt ,
17     inout                DDR_ras_n ,
18     inout                DDR_reset_n ,
19     inout                DDR_we_n ,
20     inout                FIXED_IO_dds_vrn ,
21     inout                FIXED_IO_dds_vrp ,
22     inout                [53:0]FIXED_IO_mio ,
23     inout                FIXED_IO_ps_clk ,
24     inout                FIXED_IO_ps_porb ,
25     inout                FIXED_IO_ps_srstb ,
26
27     input                rst ,
28     input                [7: 0] sw_in ,
29     input                clock_in ,
30     output logic        [7: 0] led_o
31 );
32
33     logic                fab_clk;
34     logic                clk;
35     logic                forward;
36     // Logics for the fc0 layer
37     logic                fc0_start;
38     logic [1: 0][`PREC - 1: 0] fc0_activation_i;
39     logic                fc0_valid;
40     logic                fc0_valid_i;
41     logic [`FC0_NEURONS - 1: 0][`PREC - 1: 0] fc0_activation_o ;
42     logic [`FC0_NEURONS - 1: 0][6: 0] fc0_neuron_id_o ;
43     logic                fc0_valid_act_o;
44     logic                fc0_busy;
45     logic [`FC0_NEURONS - 1: 0][`PREC - 1: 0] fc0_gradients;
46     logic                fc0_grad_valid;
47
48
49     // Logics for the fc1 layer
50     logic                fc1_start;
51     logic [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0] fc1_activation_i ;
52     logic                fc1_valid_i;
53     logic [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0] fc1_activation_o;

```

```

54 logic [`FC1_N_KERNELS - 1: 0][5: 0]          fc1_neuron_id_o ;
55 logic                                          fc1_valid_act_o;
56 logic                                          fc1_buff_rdy;
57 logic                                          fc1_busy;
58 logic                                          fc1_grad_valid;
59
60 // Logics for the fc2 layer (the last fc layer)
61 logic                                          fc2_start;
62 logic                                          fc2_buff_rdy;
63 logic [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0] fc2_activation_i;
64 logic                                          fc2_valid_i;
65 logic                                          fc2_busy;
66
67 logic [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0] fc2_activation_o;
68 logic [`FC2_N_KERNELS - 1: 0][3: 0]          fc2_neuron_id_o;
69 logic                                          fc2_valid_o;
70 logic [`FC2_NEURONS - 1: 0][`PREC - 1: 0]    fc2_act_o_buf;
71 logic                                          fc2_buf_valid;
72
73 // Backward pass logics
74 logic [`FC0_N_KERNELS - 1: 0][`PREC - 1: 0] fc0_b_gradient_i;
75 logic [`FC0_N_KERNELS - 1: 0][`PREC - 1: 0] fc0_b_activation_i;
76 logic [9: 0]                                fc0_b_activation_id_i
77 ;
78 logic [9: 0]                                fc0_b_activation_id_o
79 ;
80 logic                                          fc0_b_valid_i;
81 logic                                          fc0_b_start;
82 logic                                          fc0_b_start_r;
83 logic [3: 0]                                fc0_loops;
84 logic [`FC0_NEURONS - 1: 0][`PREC - 1: 0]    fc0_gradients_i;
85 logic                                          fc0_gradients_rdy;
86 logic [6: 0]                                fc0_n_loop_offset;
87 logic                                          fc0_bp_done;
88 logic                                          fc0_update;
89 logic                                          fc0_update_done;
90
91 logic [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0] fc1_b_gradient_i;
92 logic [`PREC - 1: 0]                          fc1_b_activation_i;
93 logic [6: 0]                                fc1_b_activation_id_i
94 ;
95 logic [6: 0]                                fc1_b_activation_id_o
96 ;
97 logic [`FC1_N_KERNELS - 1: 0][5: 0]          fc1_b_neuron_id_i;
98 logic                                          fc1_b_valid_i;
99 logic                                          fc1_b_start;
100 logic                                          fc1_b_start_r;
101 logic [3: 0]                                fc1_loops;
102 logic [`FC1_NEURONS - 1: 0][`PREC - 1: 0]    fc1_gradients;
103 logic                                          fc1_gradients_rdy;
104 logic [5: 0]                                fc1_n_offset;
105 logic [5: 0]                                fc1_n_loop_offset;

```

```

105     logic          fc1_bp_mode;
106     logic          fc1_bp_done;
107     logic          fc1_update;
108     logic          fc1_update_done;
109
110
111
112     logic [ `FC2_N_KERNELS - 1: 0 ][ `PREC - 1: 0 ] fc2_b_gradient_i;
113     logic [ `PREC - 1: 0 ] fc2_b_activation_i;
114     logic [5: 0] fc2_b_activation_id_i
115     ;
116     logic [5: 0] fc2_b_activation_id_o
117     ;
118     logic [ `FC2_N_KERNELS - 1: 0 ][3: 0] fc2_b_neuron_id_i;
119     logic fc2_b_valid_i;
120     logic fc2_b_start;
121     logic fc2_b_start_r;
122     logic [3: 0] fc2_loops;
123     logic [ `FC2_NEURONS - 1: 0 ][ `PREC - 1: 0 ] fc2_gradients;
124     logic [ `FC2_N_KERNELS - 1: 0 ][ `PREC - 1: 0 ] fc2_gradients_i;
125     logic fc2_gradients_rdy;
126     logic [3: 0] fc2_n_offset;
127     logic fc2_bp_mode;
128     logic fc2_bp_done;
129     logic fc2_update;
130     logic fc2_update_done;
131
132     logic [7: 0] img1_unpacked[784];
133     logic new_img;
134     logic [9:0] epoch;
135     logic [16:0] img_id;
136     logic [4: 0] lrate_shifts;
137     logic [4: 0] lrate_shifts_bus;
138     logic [31: 0] active_cycles;
139     logic [31: 0] idle_cycles;
140     logic training_mode;
141     logic training_mode_bus;
142     logic [16:0] img_set_size;
143
144     logic [16:0] img1_id;
145     logic [16: 0] prev_img_id;
146     logic [9:0] img1_label;
147     logic [9:0] n_epochs;
148     logic [16:0] num_correct_test;
149     logic [16:0] num_correct_train;
150     logic start;
151     logic start_bus;
152
153     // Layer States
154     logic [2: 0] fc0_state;
155     logic [2: 0] next_fc0_state;
156     logic [2: 0] fc1_state;
157     logic [2: 0] next_fc1_state;
158     logic [2: 0] fc2_state;
159     logic [2: 0] next_fc2_state;

```

```

158     logic all_idle;
159
160
161     logic [9: 0]                input_addr;
162     logic [`PREC - 1: 0]        net_input_bram_dout_a
163     ;
164     logic [`PREC - 1: 0]        net_input_bram_dout_b
165     ;
166     logic [`PREC - 1: 0]        input_data_a;
167     logic [`PREC - 1: 0]        input_data_b;
168     logic [9: 0]                img_label;
169     logic                        img_rdy;
170     logic                        epoch_fin;
171     logic                        correct;
172     logic [12: 0]               fc0_ptr_a;
173     logic [12: 0]               fc0_ptr_b;
174     logic [9: 0]                fc0_addr_a;
175     logic [9: 0]                fc0_addr_b;
176     logic [`FC2_NEURONS - 1: 0][`PREC - 1: 0] fc2_out;
177     logic [4: 0][`PREC - 1: 0]  max1;
178     logic [2: 0][`PREC - 1: 0]  max2;
179     logic [1: 0][`PREC - 1: 0]  max3;
180     logic [`PREC - 1: 0]        max4;
181     logic [`PREC - 1: 0]        max;
182     logic [4: 0]                max_valid;
183     logic [7: 0]                led_o_r;
184     logic                        sm_valid_o;
185     logic [`FC2_NEURONS - 1: 0][`PREC - 1: 0] sm_grad_o;
186     logic [31: 0]               status_block;
187
188     localparam sf      = 2.0**-12.0;
189     localparam sf2     = 2.0**-17.0;
190
191     // Backward pass states
192     localparam WEIGHT_MODE = 0;
193     localparam NEURON_MODE = 1;
194
195     // Layer states
196     localparam FORWARD    = 1;
197     localparam WAITING    = 2;
198     localparam BACKWARD   = 3;
199     localparam UPDATE     = 4;
200     localparam IDLE       = 5;
201
202     mmcm_50_mhz mmcm_50_mhz_i (
203         .clk_in1(fab_clk),
204         // .clk_in1(clock_in),
205         .clk_out1(clk)
206     );
207
208     logic [7: 0] sw_i; // so simulation uses net_input_bram
209     assign sw_i = sw_in;
210

```



```

211
212     assign start          = sw_i[0] ? 1'b1 : start_bus;
213     assign training_mode = sw_i[0] ? 1'b1 : training_mode_bus;
214     assign forward       = fc0_state == FORWARD || fc1_state ==
        FORWARD || fc2_state == FORWARD;
215     assign all_idle      = (fc0_state == IDLE) && (fc1_state == IDLE)
        && (fc2_state == IDLE);
216     assign img_rdy       = (img1_id == (img_id + 1'b1)) | (img1_id ==
        0 && img_id == img_set_size);
217     assign new_img       = start & all_idle & img_rdy;
218     assign epoch_fin     = sw_i[0] ? 1'b0 : epoch == n_epochs;
219
220     logic reset_i;
221     logic reset;
222     always_ff @(posedge clk) begin
223         reset_i      <= rst;
224         lrate_shifts <= sw_i[0] ? 5'd7 : lrate_shifts_bus;
225     end
226
227     always_ff @(posedge clk) begin
228         if (reset || !start) begin
229             idle_cycles    <= 0;
230             active_cycles  <= 0;
231         end
232         else begin
233             idle_cycles    <= idle_cycles + all_idle;
234             active_cycles  <= all_idle ? active_cycles : active_cycles +
                1'b1;
235         end
236     end
237
238     BUFG BUFG_reset(.I(reset_i), .0(reset));
239
240     always_ff @(posedge clk) begin
241         if (reset) begin
242             input_addr <= 0;
243             fc0_start  <= 0;
244         end
245         else if (fc0_state == FORWARD & !fc0_start && ~epoch_fin) begin
246             fc0_start  <= 1'b1;
247             input_addr <= 0;
248         end
249         else if (fc0_state == FORWARD & fc0_start) begin
250             input_addr <= input_addr + 1'b1;
251         end
252         else begin
253             fc0_start  <= 1'b0;
254             input_addr <= 0;
255         end
256     end
257
258
259
260

```

```

261 | assign fc0_addr_a = (forward) ? input_addr << 1 :
    |         fc0_b_activation_id_i << 1;
262 | assign fc0_addr_b = fc0_addr_a + 1'b1;
263 |
264 |
265 | net_input_bram net_input_bram_i (
266 |     .addra(fc0_addr_a),
267 |     .clka(clk),
268 |     .dina(18'b0),
269 |     .douta(net_input_bram_dout_a),
270 |     .ena(1'b1),
271 |     .wea(1'b0),
272 |
273 |     .addrb(fc0_addr_b),
274 |     .clkb(clk),
275 |     .dinb(18'b0),
276 |     .doutb(net_input_bram_dout_b),
277 |     .enb(1'b1),
278 |     .web(1'b0)
279 | );
280 |
281 |
282 |
283 |
284 | always_ff @(posedge clk) begin
285 |     if (reset) begin
286 |         prev_img_id <= img_set_size;
287 |     end
288 |     else begin
289 |         prev_img_id <= img_id;
290 |     end
291 |
292 |
293 |     if (reset || (img_id == 0 && prev_img_id != 0)) begin
294 |         num_correct_train <= 0;
295 |         num_correct_test  <= 0;
296 |     end
297 |     else if (correct) begin
298 |         num_correct_train <= (~training_mode) ?
299 |             num_correct_train : num_correct_train + 1'b1;
300 |         num_correct_test  <= (training_mode) ?
301 |             num_correct_test  : num_correct_test  + 1'b1;
302 |     end
303 |
304 |     if (reset) begin
305 |         epoch <= 0;
306 |     end
307 |     else if (img_id == 0 && prev_img_id != 0) begin
308 |         epoch <= epoch + 1'b1;
309 |     end
310 | end
311 |
312 |
313 | always_comb begin
314 |     input_data_a <= sw_i[0] ? net_input_bram_dout_a :

```

```

315         {6'b0, img1_unpacked[fc0_addr_a], 4'b0};
316     input_data_b <= sw_i[0] ? net_input_bram_dout_b :
317         {6'b0, img1_unpacked[fc0_addr_b], 4'b0};
318 end
319
320
321 always_ff @(posedge clk) begin
322     if (reset) begin
323         fc0_valid          <= 0;
324         fc0_valid_i        <= 0;
325     end
326     else begin
327         fc0_valid          <= fc0_start;
328         fc0_valid_i        <= fc0_valid;
329         fc0_activation_i   <= {input_data_b, input_data_a};
330     end
331 end
332
333
334 assign fc0_b_activation_i = {{`FC0_NEURONS{input_data_b}}, {`
    FC0_NEURONS{input_data_a}}};
335 assign fc0_gradients_rdy = fc0_grad_valid;
336 // Start when backward is good and gradients are ready. Only do
    backprop once
337 assign fc0_b_start = fc0_state == BACKWARD;
338 bit [7: 0] q, r;
339 always_ff @(posedge clk) begin
340     if (reset) begin
341         fc0_b_start_r     <= 1'b0;
342     end
343     else begin
344         fc0_b_start_r     <= fc0_b_start;
345     end
346
347     // Loop over fan in
348     if (reset) begin
349         fc0_loops         <= 0;
350     end
351     else if (fc0_state != BACKWARD) begin
352         fc0_loops         <= 0;
353     end
354     else if (fc0_b_activation_id_i == (`FC0_KERNEL_FAN_IN - 1))
        begin
355         fc0_loops         <= fc0_loops + 1'b1;
356     end
357
358     if (reset) begin
359         fc0_b_activation_id_i <= 0;
360     end
361     else if (fc0_state != BACKWARD) begin
362         fc0_b_activation_id_i <= 0;
363     end
364     else if (fc0_b_start) begin
365         fc0_b_activation_id_i <= (fc0_b_activation_id_i == (`
            FC0_KERNEL_FAN_IN - 1'b1)) ?

```

```

366             0 : fc0_b_activation_id_i + 1'b1;
367     end
368
369     for (q = 0, r = `FC0_PORT_WIDTH; q < `FC0_PORT_WIDTH; q=q+1, r=
370         r+1) begin
371         fc0_gradients_i[q]    <= fc0_gradients[q];
372     end
373     fc0_b_activation_id_o    <= fc0_b_activation_id_i << 1;
374 end
375 always_comb begin
376     case(fc0_state)
377     FORWARD:
378         next_fc0_state = fci_buff_rdy & training_mode
379             ? WAITING :
380             fci_buff_rdy & ~training_mode
381             ? IDLE : FORWARD;
382     WAITING:
383         next_fc0_state = (fc0_gradients_rdy) ? BACKWARD :
384             WAITING;
385     BACKWARD:
386         next_fc0_state = (fc0_bp_done) ? UPDATE : BACKWARD;
387     UPDATE:
388         next_fc0_state = (fc0_update_done) ? IDLE : UPDATE;
389     IDLE:
390         next_fc0_state = (new_img | sw_i[0]) ? FORWARD : IDLE
391         ;
392     default:
393         next_fc0_state = IDLE;
394     endcase
395 end
396 always_ff @(posedge clk) begin
397     if (reset) begin
398         fc0_state <= IDLE;
399     end
400     else begin
401         fc0_state <= next_fc0_state;
402     end
403 end
404
405 assign fc0_update = fc0_state == UPDATE;
406 // FC0
407 fc0_layer fc0_layer_i (
408     // inputs
409     .clk(clk),
410     .rst(reset),
411     .forward(forward),
412     .update(fc0_update),
413     .activations_i(fc0_activation_i),
414     .valid_i(fc0_valid_i & forward),
415     .lrate_shifts(lrate_shifts),
416
417     // backward pass inputs
418     .b_gradient_i(fc0_gradients_i),
419     .b_activation_i(fc0_b_activation_i),
420     .b_activation_id(fc0_b_activation_id_o),

```

```

418     .b_valid_i(fc0_b_start_r),
419
420     // outputs
421     .activation_o(fc0_activation_o),
422     .neuron_id_o(fc0_neuron_id_o),
423     .valid_act_o(fc0_valid_act_o),
424     .fc0_busy(fc0_busy),
425     .bp_done(fc0_bp_done),
426     .update_done(fc0_update_done)
427 );
428
429 always_ff @(posedge clk) begin
430     if (reset) begin
431         fc1_start    <= 1'b0;
432     end
433     else begin
434         fc1_start    <= fc1_state == FORWARD & fc1_buff_rdy;
435     end
436 end
437 interlayer_activation_buffer
438 #(.N_KERNELS_I(`FC0_NEURONS),
439  .N_KERNELS_O(`FC1_N_KERNELS),
440  .ID_WIDTH(7),
441  .BUFF_SIZE(`FC0_NEURONS),
442  .LOOPS(4))
443 interlayer_activations_fc0_fc1 (
444     // inputs
445     .clk(clk),
446     .rst(reset),
447
448     .start(fc1_start),
449     .activation_i(fc0_activation_o),
450     .neuron_id_i(fc0_neuron_id_o),
451     .valid_act_i(fc0_valid_act_o & forward),
452     .b_ptr(fc1_b_activation_id_i),
453     // outputs
454     .activation_o(fc1_activation_i),
455     .valid_o(fc1_valid_i),
456
457     .b_act_o(fc1_b_activation_i),
458
459     .buff_rdy(fc1_buff_rdy)
460 );
461
462
463
464 assign fc1_gradients_rdy  = fc1_grad_valid;
465 assign fc1_n_offset      = (fc1_loops >= `FC1_MODE_SWITCH) ?
466     fc1_loops - 4 : fc1_loops;
467 // Start when backward is good and gradients are ready. Only do
468     backprop once
469 assign fc1_b_start      = fc1_state == BACKWARD;
470 bit [5: 0] o, p;
471 always_ff @(posedge clk) begin
472     if (reset) begin

```

```

471     fc1_b_start_r    <= 1'b0;
472     fc1_bp_mode     <= 1'b0;
473 end
474 else begin
475     fc1_b_start_r    <= fc1_b_start;
476     fc1_bp_mode     <= fc1_loops >= `FC1_MODE_SWITCH ? WEIGHT_MODE
        : NEURON_MODE;
477 end
478
479 // Loop over fan in
480 if (reset) begin
481     fc1_loops     <= 0;
482 end
483 else if (fc1_state != BACKWARD) begin
484     fc1_loops     <= 0;
485 end
486 else if (fc1_b_activation_id_i == (`FC0_NEURONS - 1)) begin
487     fc1_loops     <= fc1_loops + 1'b1;
488 end
489
490 if (reset) begin
491     fc1_b_activation_id_i <= 0;
492 end
493 else if (fc1_state != BACKWARD) begin
494     fc1_b_activation_id_i <= 0;
495 end
496 else if (fc1_b_start) begin
497     fc1_b_activation_id_i <= (fc1_b_activation_id_i == (`
        FC1_FAN_IN - 1'b1)) ?
498         0 : fc1_b_activation_id_i + 1'b1;
499 end
500
501 for (p = 0, o = `FC1_PORT_WIDTH; p < `FC1_PORT_WIDTH; p=p+1, o=
    o+1) begin
502     fc1_gradients_i[p]    <= fc1_gradients[(fc1_n_offset << 3) +
        p];
503     fc1_gradients_i[o]    <= fc1_gradients[((fc1_n_offset << 3) +
        p) | 6'd32];
504     fc1_b_neuron_id_i[p]  <= (fc1_n_offset << 3) + p;
505     fc1_b_neuron_id_i[o]  <= ((fc1_n_offset << 3) + p) | 6'd32;
506 end
507 fc1_b_activation_id_o    <= fc1_b_activation_id_i;
508 end
509
510 always_comb begin
511     case(fc1_state)
512     FORWARD:
513         next_fc1_state = fc2_buff_rdy & training_mode
514             ? WAITING :
515             fc2_buff_rdy & ~training_mode
516             ? IDLE : FORWARD;
517     WAITING:
518         next_fc1_state = (fc1_gradients_rdy) ? BACKWARD :
519             WAITING;
519     BACKWARD:

```

```

520         next_fc1_state = (fc1_bp_done)      ? UPDATE   : BACKWARD;
521     UPDATE:
522         next_fc1_state = (fc1_update_done)   ? IDLE     : UPDATE;
523     IDLE:
524         next_fc1_state = (new_img | sw_i[0])  ? FORWARD  : IDLE
525         ;
526     default:
527         next_fc1_state = IDLE;
528 endcase
529 end
530 always_ff @(posedge clk) begin
531     if (reset) begin
532         fc1_state <= IDLE;
533     end
534     else begin
535         fc1_state <= next_fc1_state;
536     end
537 end
538 assign fc1_update = fc1_state == UPDATE;
539 // FC1
540 fc1_layer fc1_layer_i (
541     // inputs
542     .clk(clk),
543     .rst(reset),
544     .forward(forward),
545     .update(fc1_update),
546     .activations_i(fc1_activation_i),
547     .valid_i(fc1_valid_i & forward),
548     .lrate_shifts(lrate_shifts),
549
550     // backward pass inputs
551     .b_gradient_i(fc1_gradients_i),
552     .b_activation_i({`FC1_N_KERNELS{fc1_b_activation_i}}),
553     .b_activation_id(fc1_b_activation_id_o),
554     .b_neuron_id_i(fc1_b_neuron_id_i),
555     .b_valid_i(fc1_b_start_r),
556     .bp_mode(fc1_bp_mode),
557
558     // outputs
559     .activation_o(fc1_activation_o),
560     .neuron_id_o(fc1_neuron_id_o),
561     .valid_act_o(fc1_valid_act_o),
562     .fc1_busy(fc1_busy),
563     .bp_done(fc1_bp_done),
564     .update_done(fc1_update_done),
565
566     // backward pass outputs
567     .pl_gradients(fc0_gradients),
568     .pl_grad_valid(fc0_grad_valid)
569 );
570
571
572 always_ff @(posedge clk) begin
573     if (reset) begin

```

```

574         fc2_start    <= 1'b0;
575     end
576     else begin
577         fc2_start    <= fc2_state == FORWARD & fc2_buff_rdy;
578     end
579 end
580
581
582 interlayer_activation_buffer
583 #( .N_KERNELS_I(`FC1_N_KERNELS),
584   .N_KERNELS_O(`FC2_N_KERNELS),
585   .ID_WIDTH(6),
586   .BUFF_SIZE(`FC1_NEURONS),
587   .LOOPS(`FC2_NEURONS))
588 interlayer_activations_fc1_fc2 (
589     // inputs
590     .clk(clk),
591     .rst(reset),
592
593     .start(fc2_start),
594     .activation_i(fc1_activation_o),
595     .neuron_id_i(fc1_neuron_id_o),
596     .valid_act_i(fc1_valid_act_o & forward),
597     .b_ptr(fc2_b_activation_id_i),
598     // outputs
599
600     .activation_o(fc2_activation_i),
601     .valid_o(fc2_valid_i),
602
603     .b_act_o(fc2_b_activation_i),
604
605     .buff_rdy(fc2_buff_rdy)
606 );
607
608 always_comb begin
609     case(fc2_state)
610         FORWARD:
611             next_fc2_state =      fc2_buf_valid & training_mode
612                               ? WAITING   :
613                               fc2_buf_valid & ~training_mode
614                               ? IDLE      : FORWARD;
615         WAITING:
616             next_fc2_state = (fc2_gradients_rdy) ? BACKWARD : WAITING
617                               ;
618         BACKWARD:
619             next_fc2_state = (fc2_bp_done)          ? UPDATE : BACKWARD;
620         UPDATE:
621             next_fc2_state = (fc2_update_done)      ? IDLE   : UPDATE;
622         IDLE:
623             next_fc2_state = (new_img | sw_i[0]) ? FORWARD : IDLE;
624         default:
625             next_fc2_state = IDLE;
626     endcase
627 end
628 always_ff @(posedge clk) begin

```



```

628     if (reset) begin
629         fc2_state    <= IDLE;
630     end
631     else begin
632         fc2_state    <= next_fc2_state;
633     end
634 end
635
636
637
638
639 assign fc2_n_offset = (fc2_loops >= `FC2_MODE_SWITCH) ? fc2_loops
    - 5 : fc2_loops;
640
641 // Start when backward is good and gradients are ready. Only do
    backprop once
642 assign fc2_b_start = fc2_state == BACKWARD;
643 always_ff @(posedge clk) begin
644     if (reset) begin
645         fc2_b_start_r    <= 1'b0;
646         fc2_bp_mode      <= 1'b0;
647     end
648     else begin
649         fc2_b_start_r    <= fc2_b_start;
650         fc2_bp_mode      <= fc2_loops >= `FC2_MODE_SWITCH ? WEIGHT_MODE
            : NEURON_MODE;
651     end
652
653 // Loop over fan in
654 if (reset) begin
655     fc2_loops    <= 0;
656 end
657 else if (fc2_state != BACKWARD) begin
658     fc2_loops    <= 0;
659 end
660 else if (fc2_b_activation_id_i == (`FC1_NEURONS - 1)) begin
661     fc2_loops    <= fc2_loops + 1'b1;
662 end
663
664
665 if (reset) begin
666     fc2_b_activation_id_i <= 0;
667 end
668 else if (fc2_state != BACKWARD) begin
669     fc2_b_activation_id_i <= 0;
670 end
671 else if (fc2_b_start) begin
672     fc2_b_activation_id_i <= fc2_b_activation_id_i + 1'b1;
673 end
674 fc2_gradients_i    <= {fc2_gradients[fc2_n_offset + 5],
    fc2_gradients[fc2_n_offset]};
675 fc2_b_neuron_id_i    <= {fc2_n_offset + 5, fc2_n_offset};
676 fc2_b_activation_id_o <= fc2_b_activation_id_i;
677 end
678

```

```

679     assign fc2_update = fc2_state == UPDATE;
680     // FC2, fed directly from FC1 due to the small size
681     fc2_layer fc2_layer_i (
682         // inputs
683         .clk(clk),
684         .rst(reset),
685         .forward(forward),
686         .update(fc2_update),
687         .activations_i(fc2_activation_i),
688         .valid_i(fc2_valid_i & forward),
689         .lrate_shifts(lrate_shifts),
690
691         // backward pass inputs
692         .b_gradient_i(fc2_gradients_i),
693         .b_activation_i({fc2_b_activation_i, fc2_b_activation_i}),
694         .b_activation_id(fc2_b_activation_id_o),
695         .b_neuron_id_i(fc2_b_neuron_id_i),
696         .b_valid_i(fc2_b_start_r),
697         .bp_mode(fc2_bp_mode),
698
699         // outputs
700         .activation_o(fc2_activation_o),
701         .neuron_id_o(fc2_neuron_id_o),
702         .valid_act_o(fc2_valid_o),
703         .fc2_busy(fc2_busy),
704         .bp_done(fc2_bp_done),
705         .update_done(fc2_update_done),
706
707         // backward pass outputs
708         .p1_gradients(fc1_gradients),
709         .p1_grad_valid(fc1_grad_valid)
710     );
711
712
713
714
715
716
717     bit [`FC2_N_KERNELS - 1: 0] m;
718     logic prev_fc2_buf_valid;
719     always_ff @(posedge clk) begin
720         if (reset) begin
721             prev_fc2_buf_valid <= 0;
722             fc2_act_o_buf <= 0;
723         end
724         else begin
725             prev_fc2_buf_valid <= fc2_buf_valid;
726             for (m = 0; m < `FC2_N_KERNELS; m=m+1) begin
727                 if (fc2_valid_o && forward) begin
728                     fc2_act_o_buf[fc2_neuron_id_o[m]] <= fc2_activation_o[m];
729                 end
730             end
731         end
732     end

```

```

733     if (reset) begin
734         fc2_buf_valid    <= 1'b0;
735     end
736     else if (fc2_valid_o) begin
737         fc2_buf_valid    <= fc2_neuron_id_o[`FC2_N_KERNELS - 1] == `
            FC2_NEURONS - 1;
738     end
739     else if (fc2_state == IDLE) begin
740         fc2_buf_valid    <= 1'b0;
741     end
742 end
743
744 always @(posedge clk) begin
745     if (fc2_buf_valid) begin
746         fc2_out <= fc2_act_o_buf;
747     end
748 end
749
750
751
752
753 // LED Logic
754 bit [3: 0] k;
755 bit [3: 0] j, t;
756 always_ff @(posedge clk) begin
757     if (reset) begin
758         max    <= 0;
759         max_valid    <= 0;
760     end
761     else if ({fc2_buf_valid, prev_fc2_buf_valid} == 2'b10) begin
762         for (k = 0; k < 5; k=k+1) begin
763             max1[k] <= $signed(fc2_act_o_buf[2*k]) > $signed(
                fc2_act_o_buf[2*k+1]) ?
764                 fc2_act_o_buf[2*k] : fc2_act_o_buf[2*k + 1];
765         end
766         max_valid    <= {max_valid[3: 0], 1'b1};
767     end
768     else begin
769         max_valid[0]    <= 1'b0;
770
771         max2[0]    <= $signed(max1[0]) > $signed(max1[1]) ? max1
            [0] : max1[1];
772         max2[1]    <= $signed(max1[2]) > $signed(max1[3]) ? max1
            [2] : max1[3];
773         max2[2]    <= max1[4];
774         max_valid[1]    <= max_valid[0];
775
776         max3[0]    <= $signed(max2[0]) > $signed(max2[1]) ? max2
            [0] : max2[1];
777         max3[1]    <= max2[2];
778         max_valid[2]    <= max_valid[1];
779
780         max4    <= $signed(max3[0]) > $signed(max3[1]) ? max3
            [0] : max3[1];
781         max_valid[3]    <= max_valid[2];

```

```

782
783         max                <= max4;
784         max_valid[4]       <= max_valid[3];
785
786
787     end
788     if (reset) begin
789         led_o_r            <= 0;
790         correct             <= 1'b0;
791     end
792     else if (max_valid[4]) begin
793         correct <= fc2_act_o_buf[img_label] == max;
794         for (t = 0; t < `FC2_NEURONS; t=t+1) begin
795             if (fc2_act_o_buf[t] == max && t != img_label) begin
796                 correct <= 1'b0;
797             end
798         end
799         for (j = 0; j < 8; j=j+1) begin
800             led_o_r[j] <= fc2_act_o_buf[j] == max;
801         end
802     end
803     else begin
804         correct <= 1'b0;
805     end
806     led_o[7:0] <= led_o_r[7: 0];
807 end
808
809
810 softmax softmax_i (
811     .clk(clk),
812     .reset(reset),
813     .start(max_valid[4]),
814     .max(max),
815     .act_in(fc2_act_o_buf),
816
817     .valid_o(sm_valid_o),
818     .grad_o(sm_grad_o)
819 );
820
821 bit [3: 0] u;
822 always_ff @(posedge clk) begin
823     if (reset) begin
824         fc2_gradients_rdy <= 0;
825     end
826     else if (all_idle) begin
827         fc2_gradients_rdy <= 1'b0;
828     end
829     else if (sm_valid_o) begin
830         fc2_gradients_rdy <= 1'b1;
831     end
832
833     if (sm_valid_o) begin
834         for (u = 0; u < `FC2_NEURONS; u=u+1) begin
835             fc2_gradients[u] <= (fc2_act_o_buf[img_label] == `MIN_VAL)
836                 ? 0 :

```

```

836             sm_grad_o[u];
837         end
838         fc2_gradients[img_label] <= (fc2_act_o_buf[img_label] == `
            MAX_VAL) ? 0 :
839             $signed(sm_grad_o[img_label]) - $signed(`ONE)
            ;
840     end
841 end
842
843
844 assign status_block = {5'b0, led_o_r, fc0_state, fc1_state,
            fc2_state, forward, fc0_start,
845             fc1_start, fc2_start, fc0_busy, fc1_busy, fc2_busy,
            new_img,
846             all_idle, img_rdy};
847
848
849 logic [31:0]img1_blk0_0;
850 logic [31:0]img1_blk100_0;
851 logic [31:0]img1_blk101_0;
852 ...
853
854
855 system_wrapper system_wrapper_i
856     (DDR_addr,
857      DDR_ba,
858      DDR_cas_n,
859      DDR_ck_n,
860      DDR_ck_p,
861      DDR_cke,
862      DDR_cs_n,
863      DDR_dm,
864      DDR_dq,
865      DDR_dqs_n,
866      DDR_dqs_p,
867      DDR_odt,
868      DDR_ras_n,
869      DDR_reset_n,
870      DDR_we_n,
871      fab_clk,
872      FIXED_IO_dds_vrn,
873      FIXED_IO_dds_vrp,
874      FIXED_IO_mio,
875      FIXED_IO_ps_clk,
876      FIXED_IO_ps_por_b,
877      FIXED_IO_ps_srstb,
878      active_cycles,
879      epoch,
880      img_id,
881      idle_cycles,
882      img1_blk0_0,
883      img1_blk100_0,
884      ...
885      img1_id,
886      img1_label,

```

```

887     img_set_size,
888     lrate_shifts_bus,
889     n_epochs,
890     num_correct_test,
891     num_correct_train,
892     {fc2_out[1][17:2], fc2_out[0][17:2]},
893     {fc2_out[3][17:2], fc2_out[2][17:2]},
894     {fc2_out[5][17:2], fc2_out[4][17:2]},
895     {fc2_out[7][17:2], fc2_out[6][17:2]},
896     {fc2_out[9][17:2], fc2_out[8][17:2]},
897     start_bus,
898     status_block,
899     training_mode_bus);
900
901     always_ff @(posedge clk) begin
902         if (reset) begin
903             img_id      <= img_set_size;
904             img_label   <= 0;
905         end
906         else if (new_img) begin
907             img_id      <= img1_id;
908             img_label   <= img1_label;
909         end
910         if (new_img) begin
911             img1_unpacked[0] <= img1_blk0_0[7:0];
912             img1_unpacked[1] <= img1_blk0_0[15:8];
913             img1_unpacked[2] <= img1_blk0_0[23:16];
914             img1_unpacked[3] <= img1_blk0_0[31:24];
915             img1_unpacked[4] <= img1_blk1_0[7:0];
916         ...
917         end
918     end
919
920 endmodule

```

B.1.3 fc0_layer.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3
4  module fc0_layer(
5      input  clk,
6      input  rst,
7      input  forward,
8      input  update,
9      input  [1: 0][`PREC - 1: 0] activations_i,
10     input  valid_i,
11     input  [4: 0] lrate_shifts,
12
13     input  [`FC0_NEURONS - 1: 0][`PREC - 1: 0] b_gradient_i,
14     input  [`FC0_N_KERNELS - 1: 0][`PREC - 1: 0] b_activation_i,
15     input  [9: 0] b_activation_id,
16     input  b_valid_i,
17
18     output logic [`FC0_NEURONS - 1: 0][`PREC - 1: 0] activation_o
19     ,
20     output logic [`FC0_NEURONS - 1: 0][6: 0] neuron_id_o,
21     output logic valid_act_o,
22     output logic fc0_busy,
23     output logic bp_done,
24     output logic update_done
25 );
26
27 logic  [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_in_a;
28 logic  [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_in_b;
29 logic  [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_out_a;
30 logic  [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_out_b;
31
32 logic  [`FC0_N_KERNELS - 1: 0][`PREC - 1: 0] weights;
33 logic  [`FC0_ADDR - 1: 0] head_ptr;
34 logic  [`FC0_ADDR - 1: 0] mid_ptr;
35 logic  [`FC0_ADDR - 1: 0] addr_a;
36 logic  [`FC0_ADDR - 1: 0] addr_b;
37 logic  [`FC0_BIAS_ADDR - 1: 0] bias_ptr;
38
39 logic  [1: 0][`PREC - 1: 0] sch_activations;
40 logic  [1: 0][`PREC - 1: 0] sch_valid;
41 logic  [1: 0][`PREC - 1: 0] bram_activations;
42 logic  [1: 0][`PREC - 1: 0] bram_valid;
43 logic  [1: 0][`PREC - 1: 0] kern_activations;
44 logic  [1: 0][`PREC - 1: 0] kern_valid;
45
46 logic  [1: 0][`PREC - 1: 0] bias;
47 logic  [1: 0][`PREC - 1: 0] kern_bias;
48 logic  [255: 0] bias_container;
49 logic  [1: 0][`PREC - 1: 0] sch_has_bias;
50 logic  [1: 0][`PREC - 1: 0] bram_has_bias;
51 logic  [1: 0][`PREC - 1: 0] kern_has_bias;
52 logic  [1: 0][`PREC - 1: 0] neuron_id;
53 logic  [1: 0][`PREC - 1: 0] kern_neuron_id;

```

```

53 logic    [`FCO_N_KERNELS - 1: 0]                last_weight;
54
55 logic    [`FCO_N_KERNELS - 1: 0]                valid;
56 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  kern_activation_o
57 ;
58 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  activation_o_rel;
59 logic    [`FCO_N_KERNELS - 1: 0][6: 0]          kern_neuron_id_o;
60
61 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_gradient;
62 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_gradient_pl;
63 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_kern_grad;
64 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_act;
65 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_act_pl;
66 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_kern_act;
67
68 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  b_kern_grad_o;
69 logic    [`FCO_N_KERNELS - 1: 0]                b_kern_valid_o;
70 logic    [2: 0]                                b_valid;
71 logic    [3: 0][9: 0]                          b_act_id;
72
73 logic                                           b_weight_we;
74
75 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  kern_mult1;
76 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  kern_mult2;
77 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  weight_grad_o;
78 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  weight_grad;
79 logic    [1: 0][9: 0]
80     fc0_weight_grad_addr;
81 logic    [1: 0][9: 0]
82     fc0_weight_grad_addr_offset;
83 logic    [`FCO_NEURONS - 1: 0]                act_o_sign;
84 logic    [`FCO_N_KERNELS - 1: 0][`PREC: 0]
85     update_weights_sat;
86 logic    [`FCO_N_KERNELS - 1: 0][`PREC - 1: 0]  update_weights;
87 logic    [10: 0]                                update_ptr;
88 logic    [9: 0]                                update_addr_a;
89 logic    [9: 0]                                update_addr_b;
90 logic    [9: 0]                                w_addr_a;
91 logic    [9: 0]                                w_addr_b;
92 logic    [9: 0]                                wg_addr_a;
93 logic    [9: 0]                                wg_addr_b;
94 logic                                           w_we;
95 logic                                           wg_we;
96
97 logic                                           sch_valid_i;
98 localparam WEIGHT_MODE = 0;
99 localparam NEURON_MODE = 1;
100 logic bp_mode;
101 assign bp_mode = WEIGHT_MODE;
102 always_ff @(posedge clk) begin
103     if (rst) begin
104         sch_valid    <= 0;
105     end
106     else begin

```



```

104         sch_valid      <= valid_i;
105     end
106     sch_activations <= activations_i;
107 end
108
109 assign sch_valid_i = (forward) ? valid_i : b_valid_i;
110
111 // Scheduler for the fully connected layer
112 fc_scheduler #(.ADDR(`FCO_ADDR), .BIAS_ADDR(`FCO_BIAS_ADDR),
113 .MID_PTR_OFFSET(`FCO_KERNEL_FAN_IN), .FAN_IN(`FCO_FAN_IN))
114 fc0_scheduler_i (
115     //inputs
116     .clk(clk),
117     .rst(rst),
118     .forward(forward),
119     .valid_i(sch_valid_i),
120
121     //outputs
122     .head_ptr(head_ptr),
123     .mid_ptr(mid_ptr),
124     .bias_ptr(bias_ptr),
125     .has_bias(sch_has_bias)
126 );
127
128
129
130
131 always_ff @(posedge clk) begin
132     if (rst) begin
133         bram_valid      <= 0;
134         bram_has_bias    <= 0;
135         fc0_busy         <= 0;
136     end
137     else begin
138         bram_valid      <= sch_valid;
139         bram_has_bias    <= sch_has_bias;
140         fc0_busy         <= valid_i;
141     end
142     bram_activations    <= sch_activations;
143 end
144
145
146
147
148 always_ff @(posedge clk) begin
149     if (rst) begin
150         update_ptr <= 0;
151     end
152     else if (update) begin
153         update_ptr <= update_ptr + 1'b1;
154     end
155     else begin
156         update_ptr <= 0;
157     end
158 end

```

```

159
160
161 assign update_done      = update_ptr == 11'd783;
162 assign update_addr_a    = update_ptr[10: 1] << 1;
163 assign update_addr_b    = update_addr_a + 1'b1;
164 assign w_addr_a         = (update) ? update_addr_a : addr_a;
165 assign w_addr_b         = (update) ? update_addr_b : addr_b;
166 assign wg_addr_a        = (update) ? update_addr_a :
    fc0_weight_grad_addr[0];
167 assign wg_addr_b        = (update) ? update_addr_b :
    fc0_weight_grad_addr[1];
168 assign w_we             = (update) ? update_ptr[0] : 1'b0; //
    write when odd
169 assign wg_we            = (update) ? 1'b0 : b_weight_we;
170 assign addr_a           = (head_ptr << 1);
171 assign addr_b           = (head_ptr << 1) + 1'b1;
172 bit [7: 0] a,c;
173 always_comb begin
174     weight_grad = 0;
175     for (a = 0, c = `FC0_PORT_WIDTH; a < `FC0_PORT_WIDTH; a = a + 1,
        c=c+1) begin
176         case(lrate_shifts)
177             5'd7: begin
178                 weight_grad[a] = {{7{weight_grad_o[a][`PREC - 1]}}, {
                    weight_grad_o[a][`PREC - 1: 7]}};
179                 weight_grad[c] = {{7{weight_grad_o[c][`PREC - 1]}}, {
                    weight_grad_o[c][`PREC - 1: 7]}};
180             end
181
182             5'd9: begin
183                 weight_grad[a] = {{9{weight_grad_o[a][`PREC - 1]}}, {
                    weight_grad_o[a][`PREC - 1: 9]}};
184                 weight_grad[c] = {{9{weight_grad_o[c][`PREC - 1]}}, {
                    weight_grad_o[c][`PREC - 1: 9]}};
185             end
186             5'd11: begin
187                 weight_grad[a] = {{11{weight_grad_o[a][`PREC - 1]}}, {
                    weight_grad_o[a][`PREC - 1: 11]}};
188                 weight_grad[c] = {{11{weight_grad_o[c][`PREC - 1]}}, {
                    weight_grad_o[c][`PREC - 1: 11]}};
189             end
190             5'd10: begin
191                 weight_grad[a] = {{10{weight_grad_o[a][`PREC - 1]}}, {
                    weight_grad_o[a][`PREC - 1: 10]}};
192                 weight_grad[c] = {{10{weight_grad_o[c][`PREC - 1]}}, {
                    weight_grad_o[c][`PREC - 1: 10]}};
193             end
194             default: begin
195                 weight_grad[a] = {{8{weight_grad_o[a][`PREC - 1]}}, {
                    weight_grad_o[a][`PREC - 1: 8]}};
196                 weight_grad[c] = {{8{weight_grad_o[c][`PREC - 1]}}, {
                    weight_grad_o[c][`PREC - 1: 8]}};
197             end
198         endcase

```

```

199         update_weights_sat[a]    = $signed(data_out_a[a]) - $signed(
200             weight_grad[a]);
201         update_weights_sat[c]    = $signed(data_out_b[a]) - $signed(
202             weight_grad[c]);
203     end
204     end
205     bit [7: 0] d;
206     always_comb begin
207         for (d = 0; d < `FC0_N_KERNELS; d=d+1) begin
208             if (update_weights_sat[d][`PREC:`PREC - 1] == 2'b01) begin
209                 update_weights[d] = `MAX_VAL;
210             end
211             else if (update_weights_sat[d][`PREC:`PREC - 1] == 2'b10)
212                 begin
213                     update_weights[d] = `MIN_VAL;
214                 end
215             else begin
216                 update_weights[d] = update_weights_sat[d][`PREC - 1: 0];
217             end
218         end
219     end
220     // BRAM for the weights of the fully connected layer
221     fc0_weight_bram_controller fc0_weight_bram_controller_i (
222         // inputs
223         .clk(clk),
224         .rst(rst),
225
226         .addr_a(w_addr_a),
227         .data_in_a(update_weights[97: 0]),
228         .en_a(1'b1),
229         .we_a(w_we),
230
231         .addr_b(w_addr_b),
232         .data_in_b(update_weights[195: 98]),
233         .en_b(1'b1),
234         .we_b(w_we),
235
236         // outputs
237         .data_out_a(data_out_a),
238         .data_out_b(data_out_b),
239         .neuron_id(neuron_id)
240     );
241
242     assign b_weight_we = &b_kern_valid_o;
243
244     assign fc0_weight_grad_addr_offset[0] = 0;
245     assign fc0_weight_grad_addr_offset[1] =
246         fc0_weight_grad_addr_offset[0] + 1'b1;
247     assign fc0_weight_grad_addr[0] =
248         fc0_weight_grad_addr_offset[0] + b_act_id[3];
249     assign fc0_weight_grad_addr[1] =
250         fc0_weight_grad_addr_offset[1] + b_act_id[3];

```

```

248
249     assign bp_done = fc0_weight_grad_addr[1] == `FC0_FAN_IN - 1'b1;
250
251     fc0_weight_gradients fc0_weight_gradients_i (
252         .addr_a(wg_addr_a),
253         .clka(clk),
254         .dina(b_kern_grad_o[97: 0]),
255         .douta(weight_grad_o[97: 0]),
256         .ena(1'b1),
257         .wea(wg_we),
258
259         .addrb(wg_addr_b),
260         .clkb(clk),
261         .dinb(b_kern_grad_o[195: 98]),
262         .doutb(weight_grad_o[195: 98]),
263         .enb(1'b1),
264         .web(wg_we)
265     );
266
267     assign bias = 0;
268
269
270     always_ff @(posedge clk) begin
271         if (rst) begin
272             kern_valid      <= 0;
273             kern_has_bias   <= 0;
274         end
275         else begin
276             kern_valid      <= bram_valid;
277             kern_has_bias   <= bram_has_bias;
278         end
279         kern_activations <= {{`FC0_NEURONS{bram_activations[1]}}, {`
280             FC0_NEURONS{bram_activations[0]}}};
281         kern_bias        <= 0; //bias;
282         kern_neuron_id    <= {2{neuron_id}};
283         weights           <= {data_out_b, data_out_a};
284     end
285
286     assign kern_mult1    =    (forward) ? weights           : b_kern_grad;
287
288     assign kern_mult2    =    (forward) ? kern_activations : b_kern_act
289     ;
290
291     // Computational kernel for the fully connected layer
292     genvar i;
293     generate
294         for (i = 0; i < `FC0_N_KERNELS; i=i+1) begin
295             fc_kernel #(FAN_IN(`FC0_KERNEL_FAN_IN), .ID_WIDTH(7))
296                 fc_kernel_i (
297                     // input
298                     .clk(clk),
299                     .rst(rst),
300                     .activation_i(kern_mult2[i]),
301                     .weight(kern_mult1[i]),

```

```

300         .bias(18'b0),
301         .neuron_id_i(kern_neuron_id[i]),
302         .has_bias(kern_has_bias),
303         .valid_i(kern_valid),
304         .b_valid_i(b_valid[2]),
305         .bp_mode(bp_mode),
306         // output
307         .b_gradient_o(b_kern_grad_o[i]),
308         .b_valid_o(b_kern_valid_o[i]),
309         .activation_o(kern_activation_o[i]),
310         .neuron_id_o(kern_neuron_id_o[i]),
311         .valid_o(valid[i])
312     );
313 end
314 endgenerate
315
316 bit [7: 0] b;
317 always_ff @(posedge clk) begin
318     if (&valid) begin
319         for (b = 0; b < `FC0_NEURONS; b = b + 1) begin
320             act_o_sign[neuron_id_o[b]] <= activation_o_rel[b][`PREC -
321                 1];
322         end
323     end
324 end
325
326 assign valid_act_o = &valid;
327 assign neuron_id_o = kern_neuron_id_o[`FC0_NEURONS - 1: 0];
328
329 bit [8: 0] m, n;
330 always_comb begin
331     for (m = 0, n = `FC0_NEURONS; m < `FC0_NEURONS; m=m+1, n=n+1)
332         begin
333             activation_o_rel[m] = $signed(kern_activation_o[m]) + $signed
334                 (kern_activation_o[n]);
335             activation_o[m] = activation_o_rel[m][`PREC - 1] ? 0 :
336                 activation_o_rel[m];
337         end
338     end
339 end
340
341 bit [7: 0] q, w;
342 // Backward pass logic
343 always_ff @(posedge clk) begin
344     for (q = 0, w = `FC0_NEURONS; q < `FC0_NEURONS; q = q + 1, w =
345         w+1) begin
346         b_gradient[q] <= act_o_sign[q] ? 0 : b_gradient_i[q];
347         b_gradient[w] <= act_o_sign[q] ? 0 : b_gradient_i[q];
348     end
349     b_gradient_pl <= b_gradient;
350     b_kern_grad <= b_gradient_pl;
351
352     b_act <= b_activation_i;
353     b_act_pl <= b_act;

```

```
350     b_kern_act    <= b_act_pl;  
351  
352  
353     b_act_id      <= {b_act_id[2:0], b_activation_id};  
354     b_valid       <= {b_valid[1: 0], b_valid_i};  
355     end  
356 endmodule
```

B.1.4 fc0_weight_bram_controller.sv

```

1  `timescale 1ns / 1ps
2
3  `include "sys_defs.vh"
4
5  module fc0_weight_bram_controller (
6      input                                clk,
7      input                                rst,
8
9      input  [`FC0_ADDR - 1: 0]            addr_a,
10     input  [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_in_a,
11     input                                en_a,
12     input                                we_a,
13
14     input  [`FC0_ADDR - 1: 0]            addr_b,
15     input  [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_in_b,
16     input                                en_b,
17     input                                we_b,
18
19     output logic [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_out_a,
20     output logic [`FC0_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_out_b,
21     output logic [`FC0_PORT_WIDTH - 1: 0][6: 0] neuron_id
22 );
23
24
25 bit [6: 0] i, j;
26 always_ff @(posedge clk) begin
27     for (i = 0; i < `FC0_PORT_WIDTH; i=i+1) begin
28         neuron_id[i] <= i;
29     end
30 end
31
32 fc0_weights_bram fc0_weights_bram_i (
33     .addra(addr_a),
34     .clka(clk),
35     .dina(data_in_a),
36     .douta(data_out_a),
37     .ena(en_a),
38     .wea(we_a),
39
40     .addrb(addr_b),
41     .clkb(clk),
42     .dinb(data_in_b),
43     .doutb(data_out_b),
44     .enb(en_b),
45     .web(we_b)
46 );
47
48 endmodule

```

B.1.5 fc_kernel.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3
4  module fc_kernel #(
5      parameter FAN_IN = 0,
6      parameter ID_WIDTH = 0
7  )(
8      input                clk,
9      input                rst,
10     input  [`PREC - 1: 0] activation_i,
11     input  [`PREC - 1: 0] weight,
12     input  [`PREC - 1: 0] bias,
13     input  [ID_WIDTH - 1: 0] neuron_id_i,
14     input                has_bias,
15     input                valid_i,
16     input                b_valid_i,
17     input                bp_mode,
18
19     output logic  [`PREC - 1: 0] b_gradient_o,
20     output logic                b_valid_o,
21     output logic  [`PREC - 1: 0] activation_o,
22     output logic  [ID_WIDTH - 1: 0] neuron_id_o,
23     output logic                valid_o
24 );
25
26     logic [31: 0]                dsp_o;
27
28     logic [ID_WIDTH - 1: 0]      neuron_id;
29     logic [ID_WIDTH - 1: 0]      prev_neuron_id_i;
30     logic                        valids;
31     logic [31: 0]                kernel_in;
32
33
34     logic [35: 0]                mult_res;
35     logic [33: 0]                mac_res;
36
37     logic                        last;
38     logic                        prev_valid_i;
39     logic [8: 0]                cnt;
40
41
42     localparam WEIGHT_MODE = 0;
43     localparam NEURON_MODE = 1;
44
45
46     always_ff @(posedge clk) begin
47         if (valid_i) begin
48             cnt      <= (cnt == FAN_IN - 1) ? 0 : cnt + 1'b1;
49             last     <= cnt == FAN_IN - 1;
50         end
51         else begin
52             cnt      <= 0;
53             last     <= cnt == FAN_IN - 1;

```



```

54         end
55     end
56
57     always_ff @(posedge clk) begin
58         prev_neuron_id_i    <= neuron_id_i;
59     end
60
61     always_ff @(posedge clk) begin
62         if (bp_mode == WEIGHT_MODE) begin
63             if ({mult_res[35], &mult_res[34: 30]} == 2'b10) begin
64                 // negative saturation
65                 b_gradient_o    <= `MIN_VAL;
66             end
67             else if ({mult_res[35], |mult_res[34: 30]} == 2'b01)
68                 begin
69                 // positive saturation
70                 b_gradient_o    <= `MAX_VAL;
71             end
72             else begin
73                 b_gradient_o    <= mult_res[29: 12];
74             end
75         end
76         else begin
77             if ({mult_res[35], mult_res[34]} == 2'b10) begin
78                 // negative saturation
79                 b_gradient_o    <= `MIN_VAL;
80             end
81             else if ({mult_res[35], mult_res[34]} == 2'b01) begin
82                 // positive saturation
83                 b_gradient_o    <= `MAX_VAL;
84             end
85             else begin
86                 b_gradient_o    <= mult_res[34: 17];
87             end
88         end
89         b_valid_o            <= b_valid_i;
90     end
91
92
93
94     always_ff @(posedge clk) begin
95         if (last) begin
96             activation_o    <= dsp_o[31: 14];
97             neuron_id_o    <= prev_neuron_id_i;
98             valid_o        <= 1'b1;
99         end
100        else begin
101            valid_o            <= 1'b0;
102        end
103    end
104
105    assign kernel_in    = has_bias ? {14'b0, bias} : dsp_o;
106    assign mult_res      = $signed(weight) * $signed(activation_i);

```

```
107     assign mac_res      = $signed(mult_res[35:3]) + $signed(  
108         kernel_in);  
109     always_ff @(posedge clk) begin  
110         if ({mac_res[33], &mac_res[32: 31]} == 2'b10) begin  
111             // negative saturation  
112             dsp_o    <= 32'h8000_0000;  
113         end  
114         else if ({mac_res[33], |mac_res[32: 31]} == 2'b01) begin  
115             // positive saturation  
116             dsp_o    <= 32'h7FFF_FFFF;  
117         end  
118         else begin  
119             dsp_o    <= mac_res[31: 0];  
120         end  
121     end  
122 endmodule
```

B.1.6 fc_scheduler.sv

```

1  `timescale 1ns / 1ps
2
3  module fc_scheduler #(
4      parameter ADDR          = 0,
5      parameter BIAS_ADDR     = 0,
6      parameter MID_PTR_OFFSET = 0,
7      parameter FAN_IN        = 0
8  )(
9      input                clk,
10     input                rst,
11     input                forward,
12     input                valid_i,
13
14     output logic          [ADDR - 1: 0] head_ptr,
15     output logic          [ADDR - 1: 0] mid_ptr,
16     output logic          [BIAS_ADDR - 1: 0] bias_ptr,
17     output logic          has_bias
18 );
19
20     logic                start;
21     logic [ADDR - 1: 0] h_thresh;
22     logic [ADDR - 1: 0] next_head_ptr;
23     logic [ADDR - 1: 0] next_mid_ptr;
24     logic [BIAS_ADDR - 1: 0] next_bias_ptr;
25     logic                prev_forw;
26     logic                mode_switch;
27
28
29     assign h_thresh          = MID_PTR_OFFSET - 2;
30     assign mode_switch       = prev_forw ^ forward;
31
32     assign next_head_ptr = (mode_switch || !start) ? 0 :
33                          (!valid_i) ? head_ptr : head_ptr + 1'b1;
34     assign next_mid_ptr = (mode_switch || !start) ? MID_PTR_OFFSET
35                          :
36                          (!valid_i) ? mid_ptr : mid_ptr + 1'b1;
37     assign next_bias_ptr = (mode_switch || !start) ? 0 :
38                          (!valid_i) ? bias_ptr : bias_ptr + 1'b1;
39
40
41     always_ff @(posedge clk) begin
42         head_ptr    <= next_head_ptr;
43         mid_ptr     <= next_mid_ptr;
44         prev_forw   <= forward;
45     end
46
47     logic [ADDR - 1: 0] bias_cntr;
48     always_ff @(posedge clk) begin
49         if (rst) begin
50             bias_cntr <= 0;
51         end
52         else if (valid_i && forward) begin

```

```

53         bias_cntr    <= (bias_cntr == FAN_IN - 1) ? 0 :
           bias_cntr + 1'b1;
54     end
55     else begin
56         bias_cntr    <= 0;
57     end
58 end
59
60
61 always_ff @(posedge clk) begin
62     if (rst) begin
63         has_bias      <= 0;
64         bias_ptr      <= 0;
65     end
66     else if (valid_i && bias_cntr == 0 && forward) begin
67         has_bias      <= 1'b1;
68         bias_ptr      <= next_bias_ptr;
69     end
70     else begin
71         has_bias      <= 1'b0;
72         bias_ptr      <= bias_ptr;
73     end
74 end
75
76 always_ff @(posedge clk) begin
77     if (rst) begin
78         start        <= 1'b0;
79     end
80     else if (valid_i && !start) begin
81         start        <= 1'b1;
82     end
83     else if (valid_i && head_ptr == h_thresh) begin
84         start        <= 1'b0;
85     end
86     else if (mode_switch) begin
87         start        <= 1'b0;
88     end
89 end
90
91 endmodule

```

B.1.7 interlayer_activation_buffer.sv

```

1  `timescale 1ns / 1ps
2
3  module interlayer_activation_buffer #(
4      parameter N_KERNELS_I = 0,
5      parameter N_KERNELS_O = 0,
6      parameter ID_WIDTH = 0,
7      parameter BUFF_SIZE = 0,
8      parameter LOOPS = 0
9  )(
10     input                clk,
11     input                rst,
12     input                start,
13     input [N_KERNELS_I - 1: 0][`PREC - 1: 0] activation_i,
14     input [N_KERNELS_I - 1: 0][ID_WIDTH - 1: 0] neuron_id_i,
15     input                valid_act_i,
16     input [ID_WIDTH - 1: 0] b_ptr,
17
18
19     output logic [N_KERNELS_O - 1: 0][`PREC - 1: 0] activation_o,
20     output logic                valid_o,
21     output logic [`PREC - 1: 0] b_act_o,
22     output logic                buff_rdy
23 );
24     logic [ID_WIDTH - 1: 0] buff_ptr;
25     logic [BUFF_SIZE - 1: 0][`PREC - 1: 0] buffer;
26     logic read_o;
27     logic [LOOPS: 0] loop_cnt;
28
29
30     bit [ID_WIDTH: 0] i;
31     always_ff @(posedge clk) begin
32         if (rst) begin
33             buff_rdy <= 0;
34         end
35         else if (valid_act_i) begin
36             if (!read_o && neuron_id_i[N_KERNELS_I - 1] ==
37                 BUFF_SIZE - 1) begin
38                 buff_rdy <= 1'b1;
39             end
40             if (valid_act_i) begin
41                 for (i = 0; i < N_KERNELS_I; i=i+1) begin
42                     buffer[neuron_id_i[i]] <= activation_i[i];
43                 end
44             end
45             if (read_o) begin
46                 buff_rdy <= 1'b0;
47             end
48         end
49
50     always_ff @(posedge clk) begin
51         if (rst) begin
52             read_o <= 1'b0;

```

```

53         buff_ptr      <= 0;
54     end
55     else if (buff_rdy && start && !read_o) begin
56         read_o        <= 1'b1;
57         buff_ptr      <= 0;
58     end
59     else if (read_o) begin
60         read_o        <= ~((buff_ptr == (BUFF_SIZE - 1'b1)) & (
61             loop_cnt == LOOPS - 1));
62         buff_ptr      <= (buff_ptr == (BUFF_SIZE - 1'b1)) ? 0 :
63             buff_ptr + 1'b1;
64     end
65     if (rst) begin
66         loop_cnt      <= 0;
67     end
68     else if (~read_o) begin
69         loop_cnt      <= 0;
70     end
71     else if (buff_ptr == BUFF_SIZE - 1'b1) begin
72         loop_cnt      <= loop_cnt + 1'b1;
73     end
74 end
75 bit [ID_WIDTH - 1: 0] j;
76 always_ff @(posedge clk) begin
77     for (j = 0; j < N_KERNELS_0; j=j+1) begin
78         activation_o[j] <= buffer[buff_ptr];
79     end
80     valid_o <= read_o;
81 end
82
83 always_ff @(posedge clk) begin
84     b_act_o <= buffer[b_ptr];
85 end
86 endmodule

```

B.1.8 fc1_layer.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3
4  module fc1_layer(
5      input                                clk,
6      input                                rst,
7      input                                forward,
8      input                                update,
9      input  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  activations_i
10     ,
11     input                                valid_i,
12     input  [4: 0]                        lrate_shifts,
13
14     input  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  b_gradient_i,
15     input  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  b_activation_i,
16     input  [6: 0]                        b_activation_id,
17     input  [`FC1_N_KERNELS - 1: 0][5: 0]          b_neuron_id_i
18     ,
19     input                                b_valid_i,
20     input                                bp_mode,
21
22     output logic  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  activation_o,
23     output logic  [`FC1_N_KERNELS - 1: 0][5: 0]          neuron_id_o,
24     output logic                                valid_act_o,
25     output logic                                fci_busy,
26     output logic                                bp_done,
27     output logic                                update_done,
28
29     output logic  [`FC0_NEURONS - 1: 0][`PREC - 1: 0]  pl_gradients,
30     output logic                                pl_grad_valid
31 );
32
33
34     logic  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0]  data_in_a;
35     logic  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0]  data_in_b;
36     logic  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0]  data_out_a;
37     logic  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0]  data_out_b;
38
39     logic  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  weights;
40     logic  [`FC1_ADDR - 1: 0]                      head_ptr;
41     logic  [`FC1_ADDR - 1: 0]                      mid_ptr;
42     logic  [`FC1_BIAS_ADDR - 1: 0]                  bias_ptr;
43
44     logic  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  sch_activations;
45     logic                                sch_valid;
46     logic  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  bram_activations;
47     logic                                bram_valid;
48     logic  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]  kern_activations;
49     logic                                kern_valid;

```

```

50
51 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   bias;
52 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   kern_bias;
53 logic   [255: 0]                                   bias_container;
54 logic                                       sch_has_bias;
55 logic                                       bram_has_bias;
56 logic                                       kern_has_bias;
57 logic   [`FC1_N_KERNELS - 1: 0][5: 0]           neuron_id;
58 logic   [`FC1_N_KERNELS - 1: 0][5: 0]           kern_neuron_id;
59 logic   [`FC1_N_KERNELS - 1: 0]                 last_weight;
60
61 logic   [`FC1_N_KERNELS - 1: 0]                 valid;
62
63 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   kern_activation_o
        ;
64
65 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_gradient;
66 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_gradient_pl;
67 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_kern_grad;
68 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_act;
69 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_act_pl;
70 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_kern_act;
71
72 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   b_kern_grad_o;
73 logic   [`FC1_N_KERNELS - 1: 0]                 b_kern_valid_o;
74 logic   [2: 0]                                   b_valid;
75 logic   [3: 0][6: 0]                             b_act_id;
76 logic   [3: 0][`FC1_N_KERNELS - 1: 0][5: 0]     b_neuron_id;
77
78 logic                                       b_kern_valid;
79 logic                                       b_weight_we;
80
81 logic                                       sch_bp_mode;
82 logic                                       bram_bp_mode;
83 logic                                       kern_bp_mode;
84 logic                                       kern_bp_mode_o;
85
86 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   kern_mult1;
87 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   kern_mult2;
88 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   weight_grad;
89 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   weight_grad_o;
90 logic   [1: 0][9: 0]                             fc1_weight_grad_addr;
91 logic   [1: 0][9: 0]                             fc1_weight_grad_addr_offset;
92 logic   [`FC1_NEURONS - 1: 0]                   act_o_sign;
93 logic   [`FC1_N_KERNELS - 1: 0][`PREC: 0]       update_weights_sat;
94 logic   [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0]   update_weights;
95
96 logic   [10: 0]                                   update_ptr;
97 logic   [9: 0]                                   update_addr_a;
98 logic   [9: 0]                                   update_addr_b;
99 logic   [9: 0]                                   w_addr_a;
100 logic   [9: 0]                                   w_addr_b;

```



```

101 logic [9: 0]                wg_addr_a;
102 logic [9: 0]                wg_addr_b;
103 logic                        w_we;
104 logic                        wg_we;
105 logic                        sch_valid_i;
106
107 localparam WEIGHT_MODE = 0;
108 localparam NEURON_MODE = 1;
109
110 always_ff @(posedge clk) begin
111     if (rst) begin
112         sch_valid      <= 0;
113         sch_bp_mode    <= 0;
114     end
115     else begin
116         sch_valid      <= valid_i;
117         sch_bp_mode    <= bp_mode;
118     end
119     sch_activations <= activations_i;
120 end
121
122
123 assign sch_valid_i = (forward) ? valid_i : b_valid_i & bp_mode ==
    NEURON_MODE;
124
125 // Scheduler for the fully connected layer
126 fc_scheduler #(.ADDR(`FC1_ADDR), .BIAS_ADDR(`FC1_BIAS_ADDR),
127 .MID_PTR_OFFSET(`FC1_MID_PTR_OFFSET), .FAN_IN(`FC1_FAN_IN))
128     fc1_scheduler_i (
129 //inputs
130     .clk(clk),
131     .rst(rst),
132     .forward(forward),
133     .valid_i(sch_valid_i),
134
135 //outputs
136     .head_ptr(head_ptr),
137     .mid_ptr(mid_ptr),
138     .bias_ptr(bias_ptr),
139     .has_bias(sch_has_bias)
140 );
141
142
143
144
145 always_ff @(posedge clk) begin
146     if (rst) begin
147         bram_activations <= 0;
148         bram_valid      <= 0;
149         bram_has_bias    <= 0;
150         fc1_busy         <= 0;
151         bram_bp_mode     <= 0;
152     end
153     else begin
154         bram_activations <= sch_activations;

```

```

155     bram_valid          <= sch_valid;
156     bram_has_bias       <= sch_has_bias;
157     fc1_busy            <= valid_i;
158     bram_bp_mode        <= sch_bp_mode;
159 end
160 end
161
162
163
164
165 always_ff @(posedge clk) begin
166     if (rst) begin
167         update_ptr <= 0;
168     end
169     else if (update) begin
170         update_ptr <= update_ptr + 1'b1;
171     end
172     else begin
173         update_ptr <= 0;
174     end
175 end
176
177
178 assign update_done      = update_ptr == 11'd783;
179 assign update_addr_a    = update_ptr[10: 1] << 1;
180 assign update_addr_b    = update_addr_a + 1'b1;
181 assign w_addr_a         = (update) ? update_addr_a : head_ptr;
182 assign w_addr_b         = (update) ? update_addr_b : mid_ptr;
183 assign wg_addr_a        = (update) ? update_addr_a :
184     fc1_weight_grad_addr[0];
185 assign wg_addr_b        = (update) ? update_addr_b :
186     fc1_weight_grad_addr[1];
187 assign w_we             = (update) ? update_ptr[0] : 1'b0; //
188 write when odd
189 assign wg_we           = (update) ? 1'b0 : b_weight_we;
190
191 bit [4: 0] a,c;
192 always_comb begin
193     for (a = 0, c = `FC1_PORT_WIDTH; a < `FC1_PORT_WIDTH; a = a + 1,
194         c=c+1) begin
195         case(lrate_shifts)
196             5'd7: begin
197                 weight_grad[a] = {{7{weight_grad_o[a][`PREC - 1]}}, {
198                     weight_grad_o[a][`PREC - 1: 7]}};
199                 weight_grad[c] = {{7{weight_grad_o[c][`PREC - 1]}}, {
200                     weight_grad_o[c][`PREC - 1: 7]}};
201             end
202             5'd9: begin
203                 weight_grad[a] = {{9{weight_grad_o[a][`PREC - 1]}}, {
204                     weight_grad_o[a][`PREC - 1: 9]}};
205                 weight_grad[c] = {{9{weight_grad_o[c][`PREC - 1]}}, {
206                     weight_grad_o[c][`PREC - 1: 9]}};
207             end
208             5'd11: begin

```

```

202         weight_grad[a] = {{11{weight_grad_o[a][`PREC - 1]}}}, {
203             weight_grad_o[a][`PREC - 1: 11]}};
204         weight_grad[c] = {{11{weight_grad_o[c][`PREC - 1]}}}, {
205             weight_grad_o[c][`PREC - 1: 11]}};
206     end
207     5'd10: begin
208         weight_grad[a] = {{10{weight_grad_o[a][`PREC - 1]}}}, {
209             weight_grad_o[a][`PREC - 1: 10]}};
210         weight_grad[c] = {{10{weight_grad_o[c][`PREC - 1]}}}, {
211             weight_grad_o[c][`PREC - 1: 10]}};
212     end
213     default: begin
214         weight_grad[a] = {{8{weight_grad_o[a][`PREC - 1]}}}, {
215             weight_grad_o[a][`PREC - 1: 8]}};
216         weight_grad[c] = {{8{weight_grad_o[c][`PREC - 1]}}}, {
217             weight_grad_o[c][`PREC - 1: 8]}};
218     end
219     endcase
220     update_weights_sat[a] = $signed(data_out_a[a]) - $signed(
221         weight_grad[a]);
222     update_weights_sat[c] = $signed(data_out_b[a]) - $signed(
223         weight_grad[c]);
224 end
225 end
226 bit [7: 0] d;
227 always_comb begin
228     for (d = 0; d < `FC1_N_KERNELS; d=d+1) begin
229         if (update_weights_sat[d][`PREC:`PREC - 1] == 2'b01) begin
230             update_weights[d] = `MAX_VAL;
231         end
232         else if (update_weights_sat[d][`PREC:`PREC - 1] == 2'b10)
233             begin
234                 update_weights[d] = `MIN_VAL;
235             end
236         else begin
237             update_weights[d] = update_weights_sat[d][`PREC - 1: 0];
238         end
239     end
240 end
241 end
242
243 // BRAM for the weights of the fully connected layer
244 fc1_weight_bram_controller fc1_weight_bram_controller_i (
245     // inputs
246     .clk(clk),
247     .rst(rst),
248
249     .addr_a(w_addr_a),
250     .data_in_a(update_weights[7: 0]),
251     .en_a(1'b1),
252     .we_a(w_we),
253
254     .addr_b(w_addr_b),
255     .data_in_b(update_weights[15: 8]),
256     .en_b(1'b1),

```

```

248     .we_b(w_we),
249
250     // outputs
251     .data_out_a(data_out_a),
252     .data_out_b(data_out_b),
253     .neuron_id(neuron_id)
254 );
255
256
257 /*
258 biases_fc1_blk_mem_gen_1 biases_fc1_blk_mem_gen_1_i (
259     .addra(bias_ptr),
260     .clka(clk),
261     .dina(256'b0),
262     .douta(bias),
263     .ena(1'b1),
264     .wea(1'b0)
265 );*/
266 assign bias = 0;
267 assign b_weight_we = &b_kern_valid_o & kern_bp_mode_o ==
    WEIGHT_MODE;
268
269 assign fc1_weight_grad_addr_offset[0] = ({6'b0, b_neuron_id
    [3][0][5:3]} << 6) +
270     ({6'b0, b_neuron_id[3][0][5:3]} << 5) +
271     ({6'b0, b_neuron_id[3][0][5:3]} << 1);
272 assign fc1_weight_grad_addr_offset[1] =
    fc1_weight_grad_addr_offset[0] + `FC1_MID_PTR_OFFSET;
273
274 assign fc1_weight_grad_addr[0] = fc1_weight_grad_addr_offset[0] +
    b_act_id[3];
275 assign fc1_weight_grad_addr[1] = fc1_weight_grad_addr_offset[1] +
    b_act_id[3];
276
277 fc1_weight_gradients fc1_weight_gradients_i (
278     .addra(wg_addr_a),
279     .clka(clk),
280     .dina(b_kern_grad_o[7: 0]),
281     .douta(weight_grad_o[7: 0]),
282     .ena(1'b1),
283     .wea(wg_we),
284
285     .addrb(wg_addr_b),
286     .clkb(clk),
287     .dinb(b_kern_grad_o[15:8]),
288     .doutb(weight_grad_o[15:8]),
289     .enb(1'b1),
290     .web(wg_we)
291 );
292
293
294 previous_layer_gradient_adder previous_layer_gradient_adder_i (
295     // inputs
296     .clk(clk),
297     .rst(rst),

```

```

298     .forward(forward),
299     .grad_i(b_kern_grad_o),
300     .neuron_id_i(b_act_id[3]),
301     .valid_i(&b_kern_valid_o),
302     .bp_mode_i(kern_bp_mode_o),
303
304     // outputs
305     .pl_gradients(pl_gradients),
306     .pl_grad_valid(pl_grad_valid)
307 );
308
309 always_ff @(posedge clk) begin
310     if (rst) begin
311         kern_valid      <= 0;
312         kern_has_bias   <= 0;
313     end
314     else begin
315         kern_valid      <= bram_valid;
316         kern_has_bias   <= bram_has_bias;
317     end
318     kern_activations    <= bram_activations;
319     kern_bias           <= 0; // bias;
320     kern_neuron_id      <= neuron_id;
321     kern_bp_mode        <= bram_bp_mode;
322     kern_bp_mode_o      <= kern_bp_mode;
323     weights             <= {data_out_b, data_out_a};
324 end
325
326
327 // 3 modes of use in kernel
328 // forward: weight * activations
329 // weight gradient: gradient * activations
330 // neuron gradient: weight * gradient
331 assign kern_mult1 = (forward) ? weights :
332                     (bram_bp_mode == WEIGHT_MODE) ? b_kern_grad : weights
333                     ;
334
335 assign kern_mult2 = (forward) ? kern_activations :
336                     (bram_bp_mode == WEIGHT_MODE) ? b_kern_act :
337                     b_kern_grad;
338
339 // Computational kernel for the fully connected layer
340 genvar i;
341 generate
342     for (i = 0; i < `FC1_N_KERNELS; i=i+1) begin
343         fc_kernel #(FAN_IN(`FC1_FAN_IN), .ID_WIDTH(6)) fc_kernel_i (
344             // input
345             .clk(clk),
346             .rst(rst),
347             .activation_i(kern_mult2[i]),
348             .weight(kern_mult1[i]),
349             .bias(18'b0/*kern_bias[i]*/),
350             .neuron_id_i(kern_neuron_id[i]),
351             .has_bias(kern_has_bias),

```

```

351     .valid_i(kern_valid),
352     .b_valid_i(b_valid[2]),
353     .bp_mode(bp_mode),
354     // output
355     .b_gradient_o(b_kern_grad_o[i]),
356     .b_valid_o(b_kern_valid_o[i]),
357     .activation_o(kern_activation_o[i]),
358     .neuron_id_o(neuron_id_o[i]),
359     .valid_o(valid[i])
360 );
361 end
362 endgenerate
363
364 bit [6: 0] b;
365 always_ff @(posedge clk) begin
366     if (rst) begin
367         act_o_sign <= 0;
368     end
369     else if (&valid) begin
370         for (b = 0; b < `FC1_N_KERNELS; b = b + 1) begin
371             act_o_sign[neuron_id_o[b]] <= kern_activation_o[b][`PREC
                 - 1];
372         end
373     end
374 end
375
376 bit [10: 0] j;
377 always_comb begin
378     for (j = 0; j < `FC1_N_KERNELS; j=j+1) begin
379         activation_o[j] = kern_activation_o[j][`PREC - 1] ? 0 :
                 kern_activation_o[j];
380     end
381 end
382
383 assign bp_done = wg_we && wg_addr_a == `FC1_MID_PTR_OFFSET - 1;
384 assign valid_act_o = &valid;
385
386
387 bit [5: 0] q;
388 // Backward pass logic
389 always_ff @(posedge clk) begin
390     for (q = 0; q < `FC1_N_KERNELS; q = q + 1) begin
391         b_gradient[q] <= act_o_sign[b_neuron_id_i[q]] ? 0 :
                 b_gradient_i[q];
392     end
393     b_gradient_pl <= b_gradient;
394     b_kern_grad <= b_gradient_pl;
395
396     b_act <= b_activation_i;
397     b_act_pl <= b_act;
398     b_kern_act <= b_act_pl;
399
400
401     b_act_id <= {b_act_id[2:0], b_activation_id};
402     b_neuron_id <= {b_neuron_id[2:0], b_neuron_id_i};

```

```
403     b_valid      <= {b_valid[1: 0], b_valid_i};
404     end
405 endmodule
```

B.1.9 previous_layer_gradient_adder.sv

```

1  `timescale 1ns / 1ps
2
3  module previous_layer_gradient_adder (
4      input                                clk,
5      input                                rst,
6      input                                forward,
7      input  [`FC1_N_KERNELS - 1: 0][`PREC - 1: 0] grad_i,
8      input  [6: 0] neuron_id_i,
9      input  valid_i,
10     input  bp_mode_i,
11
12     output logic  [`FC0_NEURONS - 1: 0][`PREC - 1: 0] pl_gradients,
13     output logic
14 );
15
16 localparam WEIGHT_MODE = 0;
17 localparam NEURON_MODE = 1;
18
19 logic [7: 0][`PREC - 1: 0] stage1_grad;
20 logic stage1_valid;
21 logic [6: 0] stage1_neuron_id;
22 logic stage1_bp_mode;
23
24 logic [3: 0][`PREC - 1: 0] stage2_grad;
25 logic stage2_valid;
26 logic [6: 0] stage2_neuron_id;
27 logic stage2_bp_mode;
28
29 logic [1: 0][`PREC - 1: 0] stage3_grad;
30 logic stage3_valid;
31 logic [6: 0] stage3_neuron_id;
32 logic stage3_bp_mode;
33
34 logic [`PREC - 1: 0] stage4_grad;
35 logic stage4_valid;
36 logic prev_stage4_valid;
37 logic [6: 0] stage4_neuron_id;
38 logic stage4_bp_mode;
39
40 logic prev_bp_mode_i;
41
42 bit [4: 0] i, j, m, n, o, p;
43 always_ff @(posedge clk) begin
44     prev_bp_mode_i <= bp_mode_i;
45
46     // Stage 1 of Adder
47     if (rst) begin
48         stage1_grad <= 0;
49         stage1_valid <= 0;
50         stage1_neuron_id <= 0;
51         stage1_bp_mode <= 0;
52     end
53     else begin

```



```

54     for (i = 0, n = 0; i < 16; i = i + 2, n = n + 1) begin
55         stage1_grad[n] <= $signed(grad_i[i]) + $signed(grad_i[i +
           1]);
56     end
57     stage1_valid      <= valid_i & (bp_mode_i == NEURON_MODE);
58     stage1_neuron_id <= neuron_id_i;
59     stage1_bp_mode   <= bp_mode_i;
60 end
61
62 // Stage 2 of Adder
63 if (rst) begin
64     stage2_grad      <= 0;
65     stage2_valid     <= 0;
66     stage2_neuron_id <= 0;
67     stage2_bp_mode   <= 0;
68 end
69 else begin
70     for (j = 0, o = 0; j < 8; j = j + 2, o = o + 1) begin
71         stage2_grad[o] <= $signed(stage1_grad[j]) + $signed(
           stage1_grad[j + 1]);
72     end
73     stage2_valid      <= stage1_valid;
74     stage2_neuron_id <= stage1_neuron_id;
75     stage2_bp_mode   <= stage1_bp_mode;
76 end
77
78
79 // Stage 3 of Adder
80 if (rst) begin
81     stage3_grad      <= 0;
82     stage3_valid     <= 0;
83     stage3_neuron_id <= 0;
84     stage3_bp_mode   <= 0;
85 end
86 else begin
87     for (m = 0, p = 0; m < 4; m = m + 2, p = p + 1) begin
88         stage3_grad[p] <= $signed(stage2_grad[m]) + $signed(
           stage2_grad[m + 1]);
89     end
90     stage3_valid      <= stage2_valid;
91     stage3_neuron_id <= stage2_neuron_id;
92     stage3_bp_mode   <= stage2_bp_mode;
93 end
94
95
96 // Stage 4 of Adder
97 if (rst) begin
98     stage4_grad      <= 0;
99     stage4_valid     <= 0;
100    prev_stage4_valid <= 0;
101    stage4_neuron_id <= 0;
102    stage4_bp_mode   <= 0;
103 end
104 else begin

```

```

105     stage4_grad      <= $signed(stage3_grad[0]) + $signed(
        stage3_grad[1]);
106     stage4_valid     <= stage3_valid;
107     prev_stage4_valid <= stage4_valid;
108     stage4_neuron_id <= stage3_neuron_id;
109     stage4_bp_mode    <= stage3_bp_mode;
110 end
111
112
113 // Stage 5
114 if (rst || forward) begin
115     pl_gradients <= 0;
116 end
117 else if (stage4_valid & stage4_bp_mode == NEURON_MODE) begin
118     pl_gradients[stage4_neuron_id] <= $signed(pl_gradients[
        stage4_neuron_id]) + $signed(stage4_grad);
119 end
120
121 if (rst) begin
122     pl_grad_valid <= 0;
123 end
124 else if ({prev_stage4_valid, stage4_valid} == 2'b10) begin
125     pl_grad_valid <= 1'b1;
126 end
127 else if (forward) begin
128     pl_grad_valid <= 1'b0;
129 end
130 end
131 endmodule

```

B.1.10 fc1_weight_bram_controller.sv

```

1  `timescale 1ns / 1ps
2
3  `include "sys_defs.vh"
4
5  module fc1_weight_bram_controller (
6      input                                clk,
7      input                                rst,
8
9      input  [`FC1_ADDR - 1: 0]           addr_a,
10     input  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_in_a,
11     input                                en_a,
12     input                                we_a,
13
14     input  [`FC1_ADDR - 1: 0]           addr_b,
15     input  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_in_b,
16     input                                en_b,
17     input                                we_b,
18
19     output logic  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_out_a,
20     output logic  [`FC1_PORT_WIDTH - 1: 0][`PREC - 1: 0] data_out_b,
21     output logic  [`FC1_N_KERNELS - 1: 0][5: 0] neuron_id
22
23 );
24
25 bit  [`FC1_PORT_WIDTH - 1: 0] i, j;
26 always_ff @(posedge clk) begin
27     for (i = 0, j = 8; i < `FC1_PORT_WIDTH; i=i+1, j=j+1) begin
28         if (addr_a < `FC1_FAN_IN) begin
29             neuron_id[i]    <= i;
30             neuron_id[j]    <= i + `FC1_HALF_NEURONS;
31         end
32         else if (addr_a < `FC1_STEP2) begin
33             neuron_id[i]    <= i + `FC1_PORT_WIDTH;
34             neuron_id[j]    <= i + `FC1_PORT_WIDTH + `
35                 FC1_HALF_NEURONS;
36         end
37         else if (addr_a < `FC1_STEP3) begin
38             neuron_id[i]    <= i + `FC1_PORT_WIDTH_TIMES2;
39             neuron_id[j]    <= i + `FC1_PORT_WIDTH_TIMES2 + `
40                 FC1_HALF_NEURONS;
41         end
42         else begin
43             neuron_id[i]    <= i + `FC1_PORT_WIDTH_TIMES3;
44             neuron_id[j]    <= i + `FC1_HALF_NEURONS + `
45                 FC1_PORT_WIDTH_TIMES3;
46         end
47     end
48 end
49
50 fc1_weights_bram_0 fc1_weights_bram_0_i (
51     .addra(addr_a),
52     .clka(clk),
53     .dina(data_in_a),

```

```
51         .douta(data_out_a),  
52         .ena(en_a),  
53         .wea(we_a),  
54  
55         .addrb(addr_b),  
56         .clkb(clk),  
57         .dina(data_in_b),  
58         .doutb(data_out_b),  
59         .enb(en_b),  
60         .web(we_b)  
61     );  
62  
63 endmodule
```

B.1.11 fc2_layer.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3
4  module fc2_layer(
5      input                                clk,
6      input                                rst,
7      input                                forward,
8      input                                update,
9      input  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  activations_i
10     ,
11     input                                valid_i,
12     input  [4: 0]                        lrate_shifts,
13
14     input  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  b_gradient_i,
15     input  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  b_activation_i,
16     input  [5: 0]                                b_activation_id,
17     input  [`FC2_N_KERNELS - 1: 0][3: 0]          b_neuron_id_i
18     ,
19     input                                b_valid_i,
20     input                                bp_mode,
21
22     output logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  activation_o,
23     output logic  [`FC2_N_KERNELS - 1: 0][3: 0]          neuron_id_o,
24     output logic                                valid_act_o,
25     output logic                                fc2_busy,
26     output logic                                bp_done,
27     output logic                                update_done,
28
29     output logic  [`FC1_NEURONS - 1: 0][`PREC - 1: 0]  pl_gradients,
30     output logic                                pl_grad_valid
31 );
32
33
34     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  data_in;
35     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  data_out;
36
37     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  weights;
38     logic  [`FC2_ADDR - 1: 0]                    head_ptr;
39     logic  [`FC2_ADDR - 1: 0]                    mid_ptr;
40     logic  [`FC2_BIAS_ADDR - 1: 0]                bias_ptr;
41
42     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  sch_activations;
43     logic                                sch_valid;
44     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  bram_activations;
45     logic                                bram_valid;
46     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  kern_activations;
47     logic                                kern_valid;
48
49     logic  [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0]  bias;

```

```

50 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] kern_bias;
51 logic sch_has_bias;
52 logic bram_has_bias;
53 logic kern_has_bias;
54 logic [ `FC2_N_KERNELS - 1: 0][3: 0] neuron_id;
55 logic [ `FC2_N_KERNELS - 1: 0][3: 0] kern_neuron_id;
56 logic [ `FC2_N_KERNELS - 1: 0] last_weight;
57
58
59 logic [ `FC2_N_KERNELS - 1: 0] valid;
60
61 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_gradient;
62 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_gradient_pl;
63 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_kern_grad;
64 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_act;
65 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_act_pl;
66 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_kern_act;
67
68 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] b_kern_grad_o;
69 logic [ `FC2_N_KERNELS - 1: 0] b_kern_valid_o;
70 logic [2: 0] b_valid;
71 logic [3: 0][5: 0] b_act_id;
72 logic [3: 0][`FC2_N_KERNELS - 1: 0][3: 0] b_neuron_id;
73
74 logic b_kern_valid;
75 logic b_weight_we;
76
77 logic sch_bp_mode;
78 logic bram_bp_mode;
79 logic kern_bp_mode;
80 logic kern_bp_mode_o;
81
82 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] kern_mult1;
83 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] kern_mult2;
84 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] weight_grad;
85 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] weight_grad_o;
86 logic [ `FC2_N_KERNELS - 1: 0][9: 0]
    fc2_weight_grad_addr;
87
88 logic [ `FC2_N_KERNELS - 1: 0][`PREC : 0]
    update_weights_sat;
89 logic [ `FC2_N_KERNELS - 1: 0][`PREC - 1: 0] update_weights;
90
91 logic prev_b_kern_valid;
92 ;
93 logic [10: 0] update_ptr;
94 logic [9: 0] update_addr_a;
95 logic [9: 0] update_addr_b;
96 logic [9: 0] w_addr_a;
97 logic [9: 0] w_addr_b;
98 logic [9: 0] wg_addr_a;
99 logic [9: 0] wg_addr_b;
100 logic w_we;
101 logic wg_we;

```

```

102     logic                                                    sch_valid_i;
103
104     localparam WEIGHT_MODE = 0;
105     localparam NEURON_MODE = 1;
106
107     always_ff @(posedge clk) begin
108         if (rst) begin
109             prev_b_kern_valid    <= 0;
110         end
111         else begin
112             prev_b_kern_valid    <= &b_kern_valid_o;
113         end
114     end
115
116     always_ff @(posedge clk) begin
117         if (rst) begin
118             sch_activations <= 0;
119             sch_valid       <= 0;
120             sch_bp_mode     <= 0;
121         end
122         else begin
123             sch_activations <= activations_i;
124             sch_valid       <= valid_i;
125             sch_bp_mode     <= bp_mode;
126         end
127     end
128
129
130     assign sch_valid_i = (forward) ? valid_i : b_valid_i & bp_mode ==
        NEURON_MODE;
131
132     // Scheduler for the fully connected layer
133     fc_scheduler #(.ADDR(`FC2_ADDR), .BIAS_ADDR(`FC2_BIAS_ADDR),
134     .MID_PTR_OFFSET(`FC2_MID_PTR_OFFSET), .FAN_IN(`FC2_FAN_IN))
135     fc2_scheduler_i (
136         //inputs
137         .clk(clk),
138         .rst(rst),
139         .forward(forward),
140         .valid_i(sch_valid_i),
141
142         //outputs
143         .head_ptr(head_ptr),
144         .mid_ptr(mid_ptr),
145         .bias_ptr(bias_ptr),
146         .has_bias(sch_has_bias)
147     );
148
149
150
151     always_ff @(posedge clk) begin
152         if (rst) begin
153             bram_activations <= 0;
154             bram_valid       <= 0;
155             bram_has_bias    <= 0;

```

```

156         fc2_busy          <= 0;
157         bram_bp_mode      <= 0;
158     end
159     else begin
160         bram_activations  <= sch_activations;
161         bram_valid        <= sch_valid;
162         bram_has_bias     <= sch_has_bias;
163         fc2_busy          <= valid_i;
164         bram_bp_mode      <= sch_bp_mode;
165     end
166 end
167
168
169
170 always_ff @(posedge clk) begin
171     if (rst) begin
172         update_ptr <= 0;
173     end
174     else if (update) begin
175         update_ptr <= update_ptr + 1'b1;
176     end
177     else begin
178         update_ptr <= 0;
179     end
180 end
181
182
183 assign update_done      = update_ptr == 11'd639;
184 assign update_addr_a    = update_ptr[10: 1] << 1;
185 assign update_addr_b    = update_addr_a + 1'b1;
186 assign w_addr_a         = (update) ? update_addr_a : head_ptr;
187 assign w_addr_b         = (update) ? update_addr_b : mid_ptr;
188 assign wg_addr_a        = (update) ? update_addr_a :
189     fc2_weight_grad_addr[0];
189 assign wg_addr_b        = (update) ? update_addr_b :
190     fc2_weight_grad_addr[1];
190 assign w_we             = (update) ? update_ptr[0] : 1'b0; //
191     write when odd
191 assign wg_we            = (update) ? 1'b0 : b_weight_we;
192
193 always_comb begin
194     case(lrate_shifts)
195         5'd7: begin
196             weight_grad[0] = {{7{weight_grad_o[0][`PREC - 1]}}, {
197                 weight_grad_o[0][`PREC - 1: 7]}};
197             weight_grad[1] = {{7{weight_grad_o[1][`PREC - 1]}}, {
198                 weight_grad_o[1][`PREC - 1: 7]}};
198         end
199
200         5'd9: begin
201             weight_grad[0] = {{9{weight_grad_o[0][`PREC - 1]}}, {
202                 weight_grad_o[0][`PREC - 1: 9]}};
202             weight_grad[1] = {{9{weight_grad_o[1][`PREC - 1]}}, {
203                 weight_grad_o[1][`PREC - 1: 9]}};
203         end
204     end

```



```

204     5'd11: begin
205         weight_grad[0] = {{11{weight_grad_o[0][`PREC - 1]}}, {
                weight_grad_o[0][`PREC - 1: 11]}};
206         weight_grad[1] = {{11{weight_grad_o[1][`PREC - 1]}}, {
                weight_grad_o[1][`PREC - 1: 11]}};
207     end
208     5'd10: begin
209         weight_grad[0] = {{10{weight_grad_o[0][`PREC - 1]}}, {
                weight_grad_o[0][`PREC - 1: 10]}};
210         weight_grad[1] = {{10{weight_grad_o[1][`PREC - 1]}}, {
                weight_grad_o[1][`PREC - 1: 10]}};
211     end
212     default: begin
213         weight_grad[0] = {{8{weight_grad_o[0][`PREC - 1]}}, {
                weight_grad_o[0][`PREC - 1: 8]}};
214         weight_grad[1] = {{8{weight_grad_o[1][`PREC - 1]}}, {
                weight_grad_o[1][`PREC - 1: 8]}};
215     end
216 endcase
217
218 update_weights_sat[0]    = $signed(data_out[0]) - $signed(
                weight_grad[0]);
219 update_weights_sat[1]    = $signed(data_out[1]) - $signed(
                weight_grad[1]);
220 end
221
222 bit [7: 0] d;
223 always_comb begin
224     for (d = 0; d < `FC2_N_KERNELS; d=d+1) begin
225         if (update_weights_sat[d][`PREC :`PREC - 1] == 2'b01) begin
226             update_weights[d]    = `MAX_VAL;
227         end
228         else if (update_weights_sat[d][`PREC :`PREC - 1] == 2'b10)
                begin
229             update_weights[d]    = `MIN_VAL;
230         end
231         else begin
232             update_weights[d]    = update_weights_sat[d][`PREC - 1: 0];
233         end
234     end
235 end
236
237 // BRAM for the weights of the fully connected layer
238 fc2_weight_bram_controller fc2_weight_bram_controller_i (
239     // inputs
240     .clk(clk),
241     .rst(rst),
242
243     .addr_a(w_addr_a),
244     .data_in_a(update_weights[0]),
245     .en_a(1'b1),
246     .we_a(w_we),
247
248     .addr_b(w_addr_b),
249     .data_in_b(update_weights[1]),

```

```

250     .en_b(1'b1),
251     .we_b(w_we),
252
253     // outputs
254     .data_out(data_out),
255     .neuron_id(neuron_id)
256
257 );
258 /*
259     biases_fc2_blk_mem biases_fc2_blk_mem_i (
260         .addra(bias_ptr),
261         .clka(clk),
262         .dina(32'b0),
263         .douta({bias[1], bias[0]}),
264         .ena(1'b1),
265         .wea(1'b0)
266     );
267 */
268 assign b_weight_we = &b_kern_valid_o & kern_bp_mode_o ==
    WEIGHT_MODE;
269 assign fc2_weight_grad_addr[0] = ({6'b0, b_neuron_id[3][0]} << 6)
    + b_act_id[3];
270 assign fc2_weight_grad_addr[1] = ({6'b0, b_neuron_id[3][1]} << 6)
    + b_act_id[3];
271
272 fc2_weight_gradients fc2_weight_gradients_i (
273     .addra(wg_addr_a),
274     .clka(clk),
275     .dina(b_kern_grad_o[0]),
276     .douta(weight_grad_o[0]),
277     .ena(1'b1),
278     .wea(wg_we),
279
280     .addrb(wg_addr_b),
281     .clkb(clk),
282     .dinb(b_kern_grad_o[1]),
283     .doutb(weight_grad_o[1]),
284     .enb(1'b1),
285     .web(wg_we)
286 );
287
288
289 bit [2: 0] z;
290 always_ff @(posedge clk) begin
291
292
293     // Calculating gradients for the neurons of the previous layer
294     if (rst || forward) begin
295         pl_gradients <= 0;
296     end
297     else if (&b_kern_valid_o & kern_bp_mode_o == NEURON_MODE) begin
298         pl_gradients[b_act_id[3]] <= $signed(pl_gradients[b_act_id
299             [3]]) +
300             $signed(b_kern_grad_o[0]) +
301             $signed(b_kern_grad_o[1]);

```

```

301     end
302
303     if (rst) begin
304         pl_grad_valid    <= 0;
305     end
306     else if (&b_kern_valid_o & {kern_bp_mode_o, kern_bp_mode} == 2'
307         b10) begin
308         pl_grad_valid    <= 1'b1;
309     end
310     else if (forward) begin
311         pl_grad_valid    <= 1'b0;
312     end
313 end
314
315
316 always_ff @(posedge clk) begin
317     if (rst) begin
318         kern_valid        <= 0;
319         kern_has_bias     <= 0;
320     end
321     else begin
322         kern_valid        <= bram_valid;
323         kern_has_bias     <= bram_has_bias;
324     end
325     kern_activations     <= bram_activations;
326     kern_bias            <= 0;//bias;
327     kern_neuron_id       <= neuron_id;
328     kern_bp_mode         <= bram_bp_mode;
329     kern_bp_mode_o       <= kern_bp_mode;
330     weights              <= data_out;
331 end
332
333
334
335 // 3 modes of use in kernel
336 // forward: weight * activations
337 // weight gradient: gradient * activations
338 // neuron gradient: weight * gradient
339 assign kern_mult1 = (forward) ? weights :
340     (bram_bp_mode == WEIGHT_MODE) ? b_kern_grad : weights
341     ;
342
343 assign kern_mult2 = (forward) ? kern_activations :
344     (bram_bp_mode == WEIGHT_MODE) ? b_kern_act :
345     b_kern_grad;
346
347 // Computational kernel for the fully connected layer
348 genvar i;
349 generate
350     for (i = 0; i < `FC2_N_KERNELS; i=i+1) begin
351         fc_kernel #(FAN_IN(`FC2_FAN_IN), .ID_WIDTH(4)) fc_kernel_i (
352             // input
353             .clk(clk),
354             .rst(rst),

```

```

353     .activation_i(kern_mult2[i]),
354     .weight(kern_mult1[i]),
355     .bias(kern_bias[i]),
356     .neuron_id_i(kern_neuron_id[i]),
357     .has_bias(kern_has_bias),
358     .valid_i(kern_valid),
359     .b_valid_i(b_valid[2]),
360     .bp_mode(bp_mode),
361     // output
362     .b_gradient_o(b_kern_grad_o[i]),
363     .b_valid_o(b_kern_valid_o[i]),
364     .activation_o(activation_o[i]),
365     .neuron_id_o(neuron_id_o[i]),
366     .valid_o(valid[i])
367 );
368 end
369 endgenerate
370 assign bp_done = wg_we && wg_addr_a == `FC2_MID_PTR_OFFSET - 1;
371 assign valid_act_o = &valid;
372
373
374 // Backward pass logic
375 always_ff @(posedge clk) begin
376     b_gradient      <= b_gradient_i;
377     b_gradient_pl   <= b_gradient;
378     b_kern_grad     <= b_gradient_pl;
379
380     b_act           <= b_activation_i;
381     b_act_pl        <= b_act;
382     b_kern_act      <= b_act_pl;
383
384
385     b_act_id        <= {b_act_id[2:0], b_activation_id};
386     b_neuron_id     <= {b_neuron_id[2:0], b_neuron_id_i};
387     b_valid         <= {b_valid[1: 0], b_valid_i};
388 end
389 endmodule

```

B.1.12 fc2_weight_bram_controller.sv

```

1  `timescale 1ns / 1ps
2
3  `include "sys_defs.vh"
4
5  module fc2_weight_bram_controller (
6      input                                clk,
7      input                                rst,
8
9      input  [`FC2_ADDR - 1: 0]            addr_a,
10     input  [`PREC - 1: 0]                data_in_a,
11     input  en_a,
12     input  we_a,
13
14     input  [`FC2_ADDR - 1: 0]            addr_b,
15     input  [`PREC - 1: 0]                data_in_b,
16     input  en_b,
17     input  we_b,
18
19     output logic [`FC2_N_KERNELS - 1: 0][`PREC - 1: 0] data_out,
20     output logic [`FC2_N_KERNELS - 1: 0][3: 0] neuron_id
21 );
22
23
24     logic [`FC2_BRAM - 1: 0][`PREC - 1: 0] data_out_a;
25     logic [`FC2_BRAM - 1: 0][`PREC - 1: 0] data_out_b;
26
27     assign data_out = {data_out_b, data_out_a};
28
29     bit [`FC2_BRAM - 1: 0] i, j;
30     always_ff @(posedge clk) begin
31         if (rst) begin
32             neuron_id <= 0;
33         end
34         else begin
35             for (i = 0, j = 1; i < `FC2_BRAM; i=i+1, j=j+1) begin
36                 neuron_id[i] <= addr_a[9: 6]; // abuse the
37                                         fact that fan in is a power of 2
38                 neuron_id[j] <= addr_a[9: 6] + `FC2_HALF_NEURONS
39                 ;
40             end
41         end
42     end
43
44     fc2_weights_bram fc2_weights_bram_i (
45         .addra(addr_a),
46         .clka(clk),
47         .dina(data_in_a),
48         .douta(data_out_a[0]),
49         .ena(en_a),
50         .wea(we_a),
51
52         .addrb(addr_b),
53         .clkb(clk),

```

```
52         .dinb(data_in_b),  
53         .doutb(data_out_b[0]),  
54         .enb(en_b),  
55         .web(we_b)  
56     );  
57  
58     endmodule
```

B.1.13 softmax.sv

```

1  `timescale 1ns / 1ps
2
3  module softmax(
4      input                                clk,
5      input                                reset,
6      input                                start,
7      input  [`PREC - 1: 0]               max,
8      input  [`FC2_NEURONS - 1: 0][`PREC - 1: 0] act_in,
9
10     output logic                          valid_o,
11     output logic  [`FC2_NEURONS - 1: 0][`PREC - 1: 0] grad_o
12 );
13
14
15     logic  [`FC2_NEURONS - 1: 0][23: 0]  act_in_norm;
16     logic  [`FC2_NEURONS - 1: 0][23: 0]  fixed_exp_res;
17     logic  [`FC2_NEURONS - 1: 0][31: 0]  act_in_norm_float;
18     logic  [31: 0]                        float_o;
19     logic  [31: 0]                        float_exp_o;
20     logic                                      float_valid_o;
21     logic                                      float_exp_valid_o;
22     logic  [23: 0]                        fixed_exp_o;
23     logic  [23: 0]                        fixed_exp_sum;
24     logic                                      fixed_exp_valid_o;
25     logic  [3: 0]                        fp_in_ptr;
26     logic  [3: 0]                        fixed_exp_ptr;
27     logic  [3: 0]                        div_ptr;
28     logic                                      in_prog;
29     logic  [47: 0]                       div_o;
30     logic                                      div_valid_i;
31     logic                                      div_valid_o;
32
33     bit [3: 0] i;
34     always_ff @(posedge clk) begin
35         if (start) begin
36             for (i = 0; i < `FC2_NEURONS; i=i+1) begin
37                 act_in_norm[i] <= $signed(act_in[i]) - $signed(max);
38             end
39         end
40     end
41
42     always_ff @(posedge clk) begin
43         if (~in_prog) begin
44             fp_in_ptr <= 4'b0;
45         end
46         if (in_prog && fp_in_ptr != `FC2_NEURONS) begin
47             fp_in_ptr <= fp_in_ptr + 1'b1;
48         end
49     end
50
51     fp_to_float fp_to_float_i (
52         .s_axis_a_tdata(act_in_norm[fp_in_ptr]),
53         .s_axis_a_tvalid(in_prog & fp_in_ptr != `FC2_NEURONS),

```

```

54     .aclk(clk),
55
56     .m_axis_result_tdata(float_o),
57     .m_axis_result_tvalid(float_valid_o)
58 );
59
60
61 float_exp float_exp_i (
62     .s_axis_a_tdata(float_o),
63     .s_axis_a_tvalid(float_valid_o),
64     .aclk(clk),
65
66     .m_axis_result_tdata(float_exp_o),
67     .m_axis_result_tvalid(float_exp_valid_o)
68 );
69
70 float_to_fp float_to_fp_i (
71     .s_axis_a_tdata(float_exp_o),
72     .s_axis_a_tvalid(float_exp_valid_o),
73     .aclk(clk),
74
75     .m_axis_result_tdata(fixed_exp_o),
76     .m_axis_result_tvalid(fixed_exp_valid_o)
77 );
78
79 always_ff @(posedge clk) begin
80     if (~in_prog) begin
81         fixed_exp_sum    <= 0;
82     end
83     if (fixed_exp_valid_o) begin
84         fixed_exp_sum    <= $signed(fixed_exp_sum) + $signed(
85             fixed_exp_o);
86     end
87 end
88
89 always_ff @(posedge clk) begin
90     if (~in_prog) begin
91         fixed_exp_ptr    <= 4'b0;
92         fixed_exp_res[fixed_exp_ptr] <= 0;
93     end
94     else if (fixed_exp_valid_o) begin
95         fixed_exp_ptr    <= fixed_exp_ptr + 1'b1;
96         fixed_exp_res[fixed_exp_ptr] <= fixed_exp_o;
97     end
98 end
99
100 always_ff @(posedge clk) begin
101     if (~in_prog) begin
102         div_ptr <= 0;
103     end
104     else if (div_ptr != `FC2_NEURONS && fixed_exp_ptr == `
105         FC2_NEURONS) begin
106         div_ptr <= div_ptr + 1'b1;

```



```

107     end
108
109     assign div_valid_i = (fixed_exp_ptr == `FC2_NEURONS) & (div_ptr
110         != `FC2_NEURONS);
111
112     fixed_divider fixed_divider_i (
113         .s_axis_divisor_tdata(fixed_exp_sum),
114         .s_axis_divisor_tvalid(div_valid_i),
115         .s_axis_dividend_tdata(fixed_exp_res[div_ptr]),
116         .s_axis_dividend_tvalid(div_valid_i),
117         .aclk(clk),
118         .m_axis_dout_tdata(div_o),
119         .m_axis_dout_tvalid(div_valid_o)
120     );
121
122     );
123
124     logic [3: 0] grad_ptr;
125     always_ff @(posedge clk) begin
126         if (~in_prog) begin
127             grad_ptr <= 0;
128         end
129         else if (div_valid_o) begin
130             grad_ptr <= grad_ptr + 1'b1;
131         end
132     end
133 end
134
135 always_ff @(posedge clk) begin
136     if (div_valid_o) begin
137         grad_o[grad_ptr] <= div_o[`PREC - 1: 0];
138     end
139
140     valid_o <= ((grad_ptr == `FC2_NEURONS) & in_prog);
141
142 end
143
144 always_ff @(posedge clk) begin
145     if (reset) begin
146         in_prog <= 1'b00;
147     end
148     else if (start) begin
149         in_prog <= 1'b1;
150     end
151     else if (valid_o) begin
152         in_prog <= 1'b0;
153     end
154 end
155 endmodule

```

B.2 Testbenches

B.2.1 neural_net_top_tb.sv

```
1  `timescale 1ns / 1ps
2
3  module neural_net_top_tb(
4      );
5
6      logic clock;
7      logic reset;
8
9      neural_net_top neural_net_top_i (
10         .clock_in(clock),
11         .rst(reset),
12         .sw_in(8'h01),
13         .led_o()
14     );
15
16
17     always begin
18         #5
19         clock = ~clock;
20     end
21
22     initial begin
23         clock = 1'b0;
24         reset = 1'b1;
25         @(negedge clock);
26         reset = 1'b1;
27         @(negedge clock);
28         @(negedge clock);
29         @(negedge clock);
30         @(negedge clock);
31         @(negedge clock);
32         reset = 1'b0;
33
34         #100000;
35         $finish;
36     end
37 endmodule
```

B.2.2 softmax_tb.sv

```
1  `timescale 1ns / 1ps
2
3  module softmax_tb(
4
5      );
6
7      logic                                clock;
8      logic                                reset;
9      logic                                start;
10     logic [15: 0]                        max;
11     logic [`FC2_NEURONS - 1: 0][15: 0] act_in;
12
13     logic valid_o;
14     logic [`FC2_NEURONS - 1: 0][15: 0] grad_o;
15
16
17     assign act_in = {
18         16'h1234,
19         16'h0735,
20         16'hdf28,
21         16'hf801,
22         16'hf206,
23         16'h1842,
24         16'h1842,
25         16'h2182,
26         16'h0321,
27         16'h0a18
28     };
29
30     assign max = 16'h2182;
31
32     softmax softmax_i (
33         .clk(clock),
34         .reset(reset),
35         .start(start),
36         .max(max),
37         .act_in(act_in),
38
39         .valid_o(),
40         .grad_o()
41     );
42
43     always begin
44         #5
45         clock = ~clock;
46     end
47
48     initial begin
49         clock = 1'b0;
50         reset = 1'b1;
51         @(negedge clock);
52         reset = 1'b1;
53     end
```

```
54         @(negedge clock);
55         reset = 1'b0;
56         @(negedge clock);
57         @(negedge clock);
58         start = 1'b1;
59         @(negedge clock);
60         start = 1'b0;
61
62
63         #100000;
64         $finish;
65     end
66 endmodule
```

B.2.3 fc1_scheduler_tb.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3
4
5  module fc1_scheduler_tb(
6
7      );
8
9      logic clock;
10     logic reset;
11     logic forward;
12     logic has_bias;
13
14     logic    [`FC1_N_KERNELS - 1: 0]    valid_i;
15
16     logic    [`FC1_ADDR - 1: 0]          head_ptr;
17     logic    [`FC1_ADDR - 1: 0]          mid_ptr;
18     logic    [`FC1_BIAS_ADDR - 1: 0]     bias_ptr;
19
20
21
22
23     fc_scheduler #( .ADDR(`FC1_ADDR), .BIAS_ADDR(`FC1_BIAS_ADDR),
24     .MID_PTR_OFFSET(`FC1_MID_PTR_OFFSET), .FAN_IN(`FC1_FAN_IN))
25     fc1_scheduler_i (
26         .clk(clock),
27         .rst(reset),
28         .forward(forward),
29         .valid_i(&valid_i),
30
31         .head_ptr(head_ptr),
32         .mid_ptr(mid_ptr),
33         .bias_ptr(bias_ptr),
34         .has_bias(has_bias)
35     );
36
37     always begin
38         #5 clock = ~clock;
39     end
40
41
42     initial begin
43         clock = 1'b0;
44         reset = 1'b1;
45         forward = 1'b1;
46         valid_i = 0;
47         @(negedge clock);
48         reset = 1'b1;
49         @(negedge clock);
50         reset = 1'b0;
51         valid_i = `{`FC1_N_KERNELS{1'b1}};
52
53     end

```

```
54 |  
55 |  
56 |  
57 | endmodule
```

B.2.4 fc1_layer_tb.sv

```

1  `timescale 1ns / 1ps
2  `include "sys_defs.vh"
3  module fc1_layer_tb(
4
5      );
6
7      logic clock;
8      logic reset;
9      logic forward;
10
11     logic [`FC1_N_KERNELS - 1: 0][15: 0]    activations_i;
12     logic [`FC1_N_KERNELS - 1: 0]            valid_i;
13     logic [`FC1_N_KERNELS - 1: 0][15: 0]    activation_o;
14     logic [`FC1_N_KERNELS - 1: 0][4: 0]      neuron_id_o;
15     logic [`FC1_N_KERNELS - 1: 0]            valid_act_o;
16
17
18
19
20     fc1_layer fc1_layer_i (
21         // inputs
22         .clk(clock),
23         .rst(reset),
24         .forward(forward),
25         .activations_i(activations_i),
26         .valid_i(valid_i),
27
28         // outputs
29         .activation_o(activation_o),
30         .neuron_id_o(neuron_id_o),
31         .valid_act_o(valid_act_o)
32     );
33
34     always begin
35         #5
36         clock = ~clock;
37     end
38
39     initial begin
40         clock = 1'b0;
41         reset = 1'b1;
42         forward = 1'b1;
43         valid_i = 0;
44         @(negedge clock);
45         reset = 1'b1;
46         @(negedge clock);
47         reset = 1'b0;
48         valid_i = {'FC1_N_KERNELS{1'b1}};
49         activations_i = {'FC1_N_KERNELS{16'h2000}};
50
51         #100000;
52         $finish;
53     end

```

54

`endmodule`

APPENDIX C

Processing System Code

This appendix contains the source code for the programs that run on the PS. It should be noted that *inference_only.c* and *train.c* are quite similar. This code was written in C and cross-compiled for ARM.

C.1 train.c

```
1  #include <stdio.h>
2  #include <fcntl.h>
3  #include <sys/mman.h>
4  #include <stdint.h>
5  #include "parse_mnist.h"
6  #include <unistd.h>
7  #include <math.h>
8  #include <string.h>
9  #include <sys/time.h>
10 #include <time.h>
11
12 #define FORWARD      1
13 #define WAITING      2
14 #define BACKWARD     3
15 #define UPDATE       4
16 #define IDLE         5
17 #define SET_SIZE     70000
```

```

18 #define TRAIN_SIZE 4000
19 #define START_LRATE 9
20
21 typedef struct ddr_data {
22     // written to by fpga      Offset Desc
23     uint32_t fpga_img_id;      // 0    fpga image ptr
24     uint32_t epoch;            // 1
25     uint32_t num_correct_train; // 2
26     uint32_t num_correct_test;  // 3
27     uint32_t idle_cycles;       // 4
28     uint32_t active_cycles;     // 5
29     uint32_t status;            // 6    contains status info
30
31     // written to by arm
32     uint32_t start;             // 7    start looping
33     uint32_t n_epochs;          // 8    upper limit on epochs
34     uint32_t learning_rate;     // 9    # of right shifts
35     uint32_t training_mode;     // 10   train or just forward
36     pass
37     uint32_t img_set_size;       // 11   size of dataset
38     uint32_t img_id;             // 12   arm image ptr
39     uint32_t img_label;         // 13
40     uint32_t img[196];          // 14
41     int16_t out[10];
42 } ddr_data_t;
43
44 void state_enc_to_str(uint32_t state, char* enc);
45 void parse_mnist_data(char* filename, uint32_t** mnist_images);
46 void print_debug_data(volatile ddr_data_t* ddr_ptr);
47
48 int main() {
49     uint32_t magic_number;
50     uint32_t id, test_idx, epoch, corr_tr, corr_test;
51     uint32_t** train_images;
52     uint32_t** test_images;
53     uint32_t* train_labels;
54     uint32_t* test_labels;
55
56     int handle = open("/dev/mem", O_RDWR | O_SYNC);
57     ddr_data_t* ddr_ptr = mmap(NULL, 134217728, PROT_READ |
58         PROT_WRITE, MAP_SHARED, handle, 0x40000000);
59
60     uint32_t* ptr = (uint32_t*)ddr_ptr;
61     magic_number = ptr[400];
62     printf("### Checking Magic Number\n");
63     if (magic_number != 0xFADEDBEE) {
64         printf("### Memory was read incorrectly.\n");
65         return -1;
66     }
67     printf("### Magic number: %08x\n", magic_number);
68     printf("### Magic number successfully read.\n");
69
70     // Load MNIST images into memory

```

```

71     printf("### Loading MNIST images...\n");
72     train_images = parse_mnist_images("data/train-images.idx3-ubyte")
73     ;
74     train_labels = parse_mnist_labels("data/train-labels.idx1-ubyte")
75     ;
76     test_images = parse_mnist_images("data/t10k-images.idx3-ubyte");
77     test_labels = parse_mnist_labels("data/t10k-labels.idx1-ubyte");
78     printf("### Loading complete!\n");
79
80     struct timespec sleep;
81     sleep.tv_sec = 0;
82     sleep.tv_nsec = 1000;
83
84     // Start training!
85     ddr_ptr->start = 0;
86     usleep(100);
87     ddr_ptr->start = 1;
88     ddr_ptr->n_epochs = 5;
89     ddr_ptr->learning_rate = START_LRATE;
90     ddr_ptr->training_mode = 1;
91     ddr_ptr->img_set_size = SET_SIZE - 1;
92     struct timeval start, end;
93     gettimeofday(&start, NULL);
94     do {
95         id = (ddr_ptr->fpga_img_id + 1) % SET_SIZE;
96         epoch = ddr_ptr->epoch;
97         // Print data if epoch just finished
98         if ((id == 0) && epoch != 0) {
99             gettimeofday(&end, NULL);
100             corr_tr = ddr_ptr->num_correct_train;
101             corr_test = ddr_ptr->num_correct_test;
102             printf("\n\n### EPOCH %d\n### Training Images"
103                 ": %d/%d\nAccuracy: %f%%\n"
104                 "###Test Images: %d/%d\n"
105                 "Accuracy: %f%%\n", epoch, corr_tr, TRAIN_SIZE,
106                 (float)(corr_tr/(float)TRAIN_SIZE) * 100., corr_test,
107                 SET_SIZE - TRAIN_SIZE,
108                 ((float)corr_test/(float)(SET_SIZE - TRAIN_SIZE)) * 100.)
109                 ;
110
111             uint32_t active = ddr_ptr->active_cycles;
112             uint32_t idle = ddr_ptr->idle_cycles;
113             printf("Active Cycles: %d\t Idle Cycles: %d\n", active, idle)
114             ;
115             printf("Active Cycle Percentage: %f%%\n", 100.*((float)active
116                 / ((float)idle + (float)active));
117             printf("Elapsed time: %.5f seconds\n", (end.tv_sec - start.
118                 tv_sec) + ((end.tv_usec - start.tv_usec) * 1e-6));
119             gettimeofday(&start, NULL);
120         }
121     }
122
123     ddr_ptr->training_mode = (id < TRAIN_SIZE);

```

```

119     if (id < 60000) {
120         memcpy((void*)ddr_ptr->img, train_images[id], sizeof(uint32_t
121             ) * 196);
122         ddr_ptr->img_label = train_labels[id];
123     }
124     else {
125         test_idx = id - 60000;
126         memcpy((void*)ddr_ptr->img, test_images[test_idx], sizeof(
127             uint32_t) * 196);
128         ddr_ptr->img_label = test_labels[test_idx];
129     }
130     nanosleep(&sleep, NULL);
131     ddr_ptr->img_id = id;
132 } while (epoch < ddr_ptr->n_epochs);
133
134 void state_enc_to_str (uint32_t state, char* enc) {
135     if (state == IDLE) {
136         sprintf(enc, "IDLE");
137     }
138     else if (state == FORWARD) {
139         sprintf(enc, "FORWARD");
140     }
141     else if (state == WAITING) {
142         sprintf(enc, "WAITING");
143     }
144     else if (state == BACKWARD) {
145         sprintf(enc, "BACKWARD");
146     }
147     else if (state == UPDATE) {
148         sprintf(enc, "UPDATE");
149     }
150 }
151
152 void print_debug_data(volatile ddr_data_t* ddr_ptr) {
153     uint32_t start, fpga_img_id, img_id, img_label;
154     uint32_t status;
155     uint32_t led_o_r, fc0_state, fc1_state, fc2_state, forward,
156         fc0_start, fc1_start;
157     uint32_t fc2_start, fc0_busy, fc1_busy, fc2_busy, new_img,
158         all_idle, img_valid;
159     char fc0_state_str[40];
160     char fc1_state_str[40];
161     char fc2_state_str[40];
162     uint32_t corr_tr, corr_test;
163     float output[10];
164
165     printf("\n@@@ CURRENT STATE \n");
166     fpga_img_id = ddr_ptr->fpga_img_id;
167     img_id = ddr_ptr->img_id;
168     img_label = ddr_ptr->img_label;
169     start = ddr_ptr->start;

```

```
170 // parse the status data
171 status      = ddr_ptr->status;
172 img_valid   = status & 0x1;
173 all_idle    = (status >> 1) & 0x1;
174 new_img     = (status >> 2) & 0x1;
175 fc2_busy    = (status >> 3) & 0x1;
176 fc1_busy    = (status >> 4) & 0x1;
177 fc0_busy    = (status >> 5) & 0x1;
178 fc2_start   = (status >> 6) & 0x1;
179 fc1_start   = (status >> 7) & 0x1;
180 fc0_start   = (status >> 8) & 0x1;
181 forward     = (status >> 9) & 0x1;
182 fc2_state   = (status >> 10) & 0x7;
183 fc1_state   = (status >> 13) & 0x7;
184 fc0_state   = (status >> 16) & 0x7;
185 led_o_r     = (status >> 19) & 0xFF;
186 corr_tr     = ddr_ptr->num_correct_train;
187 corr_test   = ddr_ptr->num_correct_test;
188 state_enc_to_str(fc0_state, fc0_state_str);
189 state_enc_to_str(fc1_state, fc1_state_str);
190 state_enc_to_str(fc2_state, fc2_state_str);
191
192 float max_out    = -100;
193 int max_out_id   = 0;
194 for (int i = 0; i < 10; i++) {
195     output[i] = (float)(ddr_ptr->out[i]) / pow(2, 10);
196     if (output[i] > max_out) {
197         max_out = output[i];
198         max_out_id = i;
199     }
200 }
201
202 printf("fpga_img_id: %d\t\timg1_id: %d\n", fpga_img_id, img_id);
203 printf("img1_label: %d\t\tmax_out: %d\t\tled_o: %08x\n",
        img_label, max_out_id, led_o_r);
204 printf("Output:\n");
205 for (int i = 0; i < 10; i++) {
206     printf("%d: %f\n", i, output[i]);
207 }
208
209
210 }
```

C.2 inference_only.c

```

1  #include <stdio.h>
2  #include <fcntl.h>
3  #include <sys/mman.h>
4  #include <stdint.h>
5  #include "parse_mnist.h"
6  #include <unistd.h>
7  #include <math.h>
8  #include <sys/time.h>
9  #include <time.h>
10 #include <string.h>
11
12 #define FORWARD 1
13 #define WAITING 2
14 #define BACKWARD 3
15 #define UPDATE 4
16 #define IDLE 5
17 #define SET_SIZE 70000
18 #define TRAIN_SIZE 70000
19
20 typedef struct ddr_data {
21     // written to by fpga           Offset   Desc
22     uint32_t fpga_img_id;           // 0      fpga image ptr
23     uint32_t epoch;                 // 1
24     uint32_t num_correct_train;     // 2
25     uint32_t num_correct_test;      // 3
26     uint32_t idle_cycles;           // 4
27     uint32_t active_cycles;         // 5
28     uint32_t status;                // 6      contains status info
29
30     // written to by arm
31     uint32_t start;                 // 7      start looping
32     uint32_t n_epochs;              // 8      upper limit on epochs
33     uint32_t learning_rate;         // 9      # of right shifts
34     uint32_t training_mode;         // 10     train or just forward
35     pass
36     uint32_t img_set_size;           // 11     size of dataset
37     uint32_t img_id;                // 12     arm image ptr
38     uint32_t img_label;             // 13
39     uint32_t img[196];              // 14
40     int16_t out[10];
41 } ddr_data_t;
42
43 void state_enc_to_str(uint32_t state, char* enc);
44 void parse_mnist_data(char* filename, uint32_t** mnist_images);
45 void print_debug_data(volatile ddr_data_t* ddr_ptr);
46
47 int main() {
48     uint32_t magic_number;
49     uint32_t id, test_idx, epoch, corr_tr, corr_test;
50     uint32_t** train_images;
51     uint32_t** test_images;

```

```

52  uint32_t* train_labels;
53  uint32_t* test_labels;
54
55  int handle = open("/dev/mem", O_RDWR | O_SYNC);
56  volatile ddr_data_t* ddr_ptr = mmap(NULL, 134217728, PROT_READ |
    PROT_WRITE, MAP_SHARED, handle, 0x40000000);
57
58
59  uint32_t* ptr = (uint32_t*)ddr_ptr;
60  magic_number = ptr[400];
61  printf("### Checking Magic Number\n");
62  if (magic_number != 0xFADEDBEE) {
63      printf("### Memory was read incorrectly.\n");
64      return -1;
65  }
66  printf("### Magic number: %08x\n", magic_number);
67  printf("### Magic number successfully read.\n");
68
69  // Load MNIST images into memory
70  printf("### Loading MNIST images...\n");
71  train_images = parse_mnist_images("data/train-images.idx3-ubyte")
    ;
72  train_labels = parse_mnist_labels("data/train-labels.idx1-ubyte")
    ;
73  test_images = parse_mnist_images("data/t10k-images.idx3-ubyte");
74  test_labels = parse_mnist_labels("data/t10k-labels.idx1-ubyte");
75  printf("### Loading complete!\n");
76
77  struct timespec sleep;
78  sleep.tv_sec = 0;
79  sleep.tv_nsec = 1000;
80
81  // Start training!
82  ddr_ptr->start = 0;
83  usleep(10);
84  ddr_ptr->start = 1;
85  ddr_ptr->n_epochs = 2;
86  ddr_ptr->training_mode = 0;
87  ddr_ptr->img_set_size = SET_SIZE - 1;
88  struct timeval start, end;
89  gettimeofday(&start, NULL);
90  do {
91      id = (ddr_ptr->fpga_img_id + 1) % SET_SIZE;
92      epoch = ddr_ptr->epoch;
93      // Print data if epoch just finished
94      if ((id == 0) && epoch != 0) {
95          gettimeofday(&end, NULL);
96
97          corr_tr = ddr_ptr->num_correct_train;
98          corr_test = ddr_ptr->num_correct_test;
99          printf("\nImages"
100              " : %d/%d\nAccuracy: %f%%\n", corr_test, 70000,
101              ((float)corr_test/70000.) * 100.);
102
103

```

```
104     uint32_t active = ddr_ptr->active_cycles;
105     uint32_t idle = ddr_ptr->idle_cycles;
106     printf("Active Cycles: %d\t Idle Cycles: %d\n", active, idle)
107     ;
108     printf("Active Cycle Percentage: %f%%\n", ((float)active / ((
109         float)idle + (float)active));
110     printf("Elapsed time: %.5f seconds\n\n", (end.tv_sec - start.
111         tv_sec) + ((end.tv_usec - start.tv_usec) * 1e-6));
112     gettimeofday(&start, NULL);
113 }
114
115 if (id < 60000) {
116     memcpy((void*)ddr_ptr->img, train_images[id], sizeof(uint32_t
117         ) * 196);
118     ddr_ptr->img_label = train_labels[id];
119 }
120 else {
121     test_idx = id - 60000;
122     memcpy((void*)ddr_ptr->img, test_images[test_idx], sizeof(
123         uint32_t) * 196);
124     ddr_ptr->img_label = test_labels[test_idx];
125 }
126
127 nanosleep(&sleep, NULL);
128 ddr_ptr->img_id = id;
129
130 } while (epoch < ddr_ptr->n_epochs);
131 }
```


C.3 parse_mnist.c

```

1  #include <stdio.h>
2  #include <stdint.h>
3  #include <stdlib.h>
4
5  uint32_t** parse_mnist_images(char* filename) {
6      FILE* f;
7      uint8_t* u8;
8      uint8_t** images;
9      uint32_t** images32;
10     uint32_t magic_number;
11     uint32_t n_items;
12     uint32_t rows;
13     uint32_t cols;
14
15     f = fopen(filename, "rb");
16
17     fread(&magic_number, 4, 1, f);
18     fread(&n_items, 1, 4, f);
19     fread(&rows, 1, 4, f);
20     fread(&cols, 1, 4, f);
21
22     u8 = (uint8_t*)&magic_number;
23     magic_number = u8[3] + (u8[2] << 8) + (u8[1] << 16) + (u8[0] <<
24         24);
25     u8 = (uint8_t*)&n_items;
26     n_items = u8[3] + (u8[2] << 8) + (u8[0] << 16) + (u8[0] << 24);
27     u8 = (uint8_t*)&rows;
28     rows = u8[3] + (u8[2] << 8) + (u8[0] << 16) + (u8[0] << 24);
29     u8 = (uint8_t*)&cols;
30     cols = u8[3] + (u8[2] << 8) + (u8[0] << 16) + (u8[0] << 24);
31
32     images = malloc(sizeof(uint8_t*) * n_items);
33
34     for (int i = 0; i < n_items; i++) {
35         images[i] = malloc(sizeof(uint8_t) * 784);
36         fread(images[i], 784, 1, f);
37     }
38
39     images32 = malloc(sizeof(uint32_t*) * n_items);
40     for (int i = 0; i < n_items; i++) {
41         images32[i] = malloc(sizeof(uint32_t) * 196);
42         for (int j = 0; j < 196; j++) {
43             images32[i][j] = images[i][4*j] + (images[i][4*j+1] << 8) +
44                 (images[i][4*j+2] << 16) + (images[i][4*j+3] << 24);
45         }
46     }
47
48     fclose(f);
49     return images32;
50 }
51
52 uint32_t* parse_mnist_labels(char* filename) {

```

```
52  uint8_t* ptr;
53  uint8_t* labels;
54  uint32_t* labels32;
55  uint32_t magic_number;
56  uint32_t n_items;
57  uint8_t label;
58  FILE* f;
59  f = fopen(filename, "rb");
60
61  fread(&magic_number, 4, 1, f);
62  ptr = (uint8_t*)&magic_number;
63  magic_number = ptr[3] + (ptr[2] << 8) + (ptr[1] << 16) + (ptr[0] <<
    << 24);
64  fread(&n_items, 1, 4, f);
65  ptr = (uint8_t*)&n_items;
66  n_items = ptr[3] + (ptr[2] << 8) + (ptr[0] << 16) + (ptr[0] <<
    24);
67  labels = malloc(sizeof(uint8_t) * n_items);
68  labels32 = malloc(sizeof(uint32_t) * n_items);
69  for (int i = 0; i < n_items; i++) {
70      fread(&labels[i], 1, 1, f);
71  }
72
73  for (int i = 0; i < n_items; i++) {
74      labels32[i] = labels[i];
75  }
76
77  fclose(f);
78  return labels32;
79 }
```

APPENDIX D

Hardware Testing Code

This Appendix contains the code that was use to verify the simulated output of the hardware model. This code was written in Python.

D.1 fpga_forward_backward_pass_test.py

```
1 import glob
2 import math
3 import random
4 import matplotlib.pyplot as plt
5 import matplotlib as mpl
6 import numpy as np
7 import seaborn as sns
8 from itertools import chain
9
10 integer_bits = 6
11 int_bits_grad = 1
12
13 # Set activations in
14 fc0_fan_in = 28*28
15 activations_i = []
16 f = open('../FPGA/FPGA.srcs/sources_1/ip/rand_input_6.12.coe')
17 next(f)
18 next(f)
19 for line in f:
```

```

20     act = 0
21     bit_val = 2 ** (integer_bits - 1)
22
23     if (line[0] == '1'):
24         act -= bit_val
25     bit_val /= 2.
26
27     for bit in line[1:]:
28         if (bit == '1'):
29             act += bit_val
30         bit_val /= 2.
31     activations_i.append(act)
32
33 # FC0 layer
34 fc0_n_neurons = 98
35 fc0_fan_in = 28*28
36 fc0_bram = 1
37 fc0_neurons_per_bram = fc0_n_neurons / fc0_bram
38 fc0_neurons = [[] for i in range(fc0_n_neurons)]
39 n_offset = 0
40
41
42 # read fc0
43 path = '../FPGA/FPGA.srcs/sources_1/ip/fc0_weights_1.17.coe'
44 for fname in glob.glob(path):
45     print(fname)
46     f = open(fname, 'r')
47     next(f)
48     next(f)
49     for line in f:
50         for i in range(fc0_n_neurons):
51             curr_neuron = i
52             st_bit = i * 18
53             end_bit = (i + 1) * 18
54
55             bit_str = line[st_bit: end_bit]
56             weight_val = 0
57             bit_val = 2 ** (int_bits_grad - 1)
58
59             if (bit_str[0] == '1'):
60                 weight_val -= bit_val
61             bit_val /= 2.
62
63             for bit in bit_str[1:]:
64                 if (bit == '1'):
65                     weight_val += bit_val
66                 bit_val /= 2.
67             fc0_neurons[fc0_n_neurons - (curr_neuron + 1)].append(
68                 weight_val)
69
70 # FC1 layer
71 fc1_n_neurons = 64
72 fc1_fan_in = 98
73 fc1_bram = 8

```

```
74 fcl_neurons_per_bram = fcl_n_neurons / fcl_bram
75 fcl_neurons = [[] for i in range(fcl_n_neurons)]
76 n_offset = 0
77
78
79 # read fcl
80 path = '../FPGA/FPGA.srscs/sources_1/ip/fcl_weights2_1.17.coe'
81 for fname in glob.glob(path):
82     print(fname)
83     f = open(fname, 'r')
84     next(f)
85     next(f)
86     for line in f:
87         for i in range(8):
88             curr_neuron = n_offset + (7 - i)
89             st_bit = i * 18
90             end_bit = (i + 1) * 18
91
92             bit_str = line[st_bit: end_bit]
93             weight_val = 0
94             bit_val = 2 ** (int_bits_grad - 1)
95
96             if (bit_str[0] == '1'):
97                 weight_val -= bit_val
98             bit_val /= 2.
99
100             for bit in bit_str[1:]:
101                 if (bit == '1'):
102                     weight_val += bit_val
103                 bit_val /= 2.
104             fcl_neurons[curr_neuron].append(weight_val)
105         if len(fcl_neurons[n_offset]) == fcl_fan_in:
106             n_offset += 8
107
108 # FC2 layer
109 fc2_n_neurons = 10
110 fc2_fan_in = fcl_n_neurons
111 fc2_bram = 1
112 fc2_neurons_per_bram = fc2_n_neurons / fc2_bram
113 fc2_neurons = [[] for i in range(fc2_n_neurons)]
114 n_offset = 0
115
116
117 # read fc2
118 path = '../FPGA/FPGA.srscs/sources_1/ip/fc2_weights_1.17.coe'
119 for fname in glob.glob(path):
120     print(fname)
121     f = open(fname, 'r')
122     next(f)
123     next(f)
124     curr_neuron = n_offset
125     for line in f:
126         weight_val = 0
127         bit_val = 2 ** (int_bits_grad - 1)
128
```

```

129     if (line[0] == '1'):
130         weight_val -= bit_val
131     bit_val /= 2.
132
133     for bit in line[1:]:
134         if (bit == '1'):
135             weight_val += bit_val
136             bit_val /= 2.
137     fc2_neurons[curr_neuron].append(weight_val)
138     if len(fc2_neurons[curr_neuron]) == fc2_fan_in:
139         curr_neuron += fc2_bram
140     n_offset += 1
141
142 fc0_output = []
143 for neuron in fc0_neurons:
144     n_out = 0
145     for j in range(len(activations_i[0: fc0_fan_in])):
146         n_out += activations_i[j + 0*fc0_fan_in] * neuron[j]
147     fc0_output.append(max(n_out, 0))
148
149
150 fc1_output = []
151 for neuron in fc1_neurons:
152     n_out = 0
153     for j in range(len(fc0_output)):
154         n_out += fc0_output[j] * neuron[j]
155     fc1_output.append(max(n_out, 0))
156
157 fc2_output = []
158 for neuron in fc2_neurons:
159     n_out = 0
160     for j in range(len(fc1_output)):
161         n_out += fc1_output[j] * neuron[j]
162     fc2_output.append(n_out)
163
164 gradients = []
165 sm_output = []
166 sm_sum = 0.
167 max_v = max(fc2_output)
168 for output in fc2_output:
169     sm_sum += math.exp(output - max_v)
170
171 for output in fc2_output:
172     sm_output.append(math.exp(output - max_v) / sm_sum)
173
174 for output in sm_output:
175     gradients.append(output)
176
177 print(gradients[0])
178 gradients[0] -= 1
179 loss = -math.log(gradients[0] + 1.)
180
181
182
183 weight_grad = [[] for i in range(fc2_n_neurons)]

```

```

184 for i in range(len(fc2_neurons)):
185     for j in range(len(fc2_neurons[i])):
186         weight_grad[i].append(gradients[i] * fc1_output[j])
187
188
189 fc1_grad = [0 for i in range(fc1_n_neurons)]
190 for i in range(len(fc2_neurons)):
191     for j in range(len(fc2_neurons[i])):
192         fc1_grad[j] += gradients[i] * fc2_neurons[i][j]
193
194 fc1_w_grad = [[] for i in range(fc1_n_neurons)]
195 for i in range(len(fc1_neurons)):
196     for j in range(len(fc1_neurons[i])):
197         de_dnet = fc1_output[i] > 0
198         fc1_w_grad[i].append(de_dnet * fc1_grad[i] * fc0_output[j])
199
200 fc0_grad = [0 for i in range(fc0_n_neurons)]
201 for i in range(len(fc1_neurons)):
202     for j in range(len(fc1_neurons[i])):
203         de_dnet = fc1_output[i] > 0
204         fc0_grad[j] += (de_dnet * fc1_grad[i]) * fc1_neurons[i][j]
205
206 fc0_w_grad = [[] for i in range(fc0_n_neurons)]
207 for i in range(len(fc0_neurons)):
208     for j in range(len(fc0_neurons[i])):
209         de_dnet = fc0_output[i] > 0
210         fc0_w_grad[i].append(de_dnet * fc0_grad[i] * activations_i[j])
211
212 # conduct 10 gradient check tests
213 for i in range(10):
214
215     #pick random weight in layer 0
216     n_idx = random.randint(0, len(fc0_neurons) - 1)
217     w_idx = random.randint(0, len(fc0_neurons) - 1)
218     print('Calculated gradient:\t' + str(fc0_w_grad[n_idx][w_idx]))
219     eps = 1e-4
220     fc0_neurons[n_idx][w_idx] += eps
221     fc0_output = []
222     for neuron in fc0_neurons:
223         n_out = 0
224         for j in range(len(activations_i[0: fc0_fan_in])):
225             n_out += activations_i[j + 0*fc0_fan_in] * neuron[j]
226         fc0_output.append(max(n_out, 0))
227
228
229     fc1_output = []
230     for neuron in fc1_neurons:
231         n_out = 0
232         for j in range(len(fc0_output)):
233             n_out += fc0_output[j] * neuron[j]
234         fc1_output.append(max(n_out, 0))
235
236     fc2_output = []
237     for neuron in fc2_neurons:
238         n_out = 0

```

```

239     for j in range(len(fc1_output)):
240         n_out += fc1_output[j] * neuron[j]
241     fc2_output.append(n_out)
242
243     gradients = []
244     sm_output = []
245     sm_sum = 0.
246     max_v = max(fc2_output)
247     for output in fc2_output:
248         sm_sum += math.exp(output - max_v)
249
250     for output in fc2_output:
251         sm_output.append(math.exp(output - max_v) / sm_sum)
252
253     for output in sm_output:
254         gradients.append(output)
255
256
257     gradients[0] -= 1
258     loss2 = -math.log(gradients[0] + 1.)
259
260
261     fc0_neurons[n_idx][w_idx] -= (2 * eps)
262     fc0_output = []
263     for neuron in fc0_neurons:
264         n_out = 0
265         for j in range(len(activations_i[0: fc0_fan_in])):
266             n_out += activations_i[j + 0*fc0_fan_in] * neuron[j]
267         fc0_output.append(max(n_out, 0))
268
269
270     fc1_output = []
271     for neuron in fc1_neurons:
272         n_out = 0
273         for j in range(len(fc0_output)):
274             n_out += fc0_output[j] * neuron[j]
275         fc1_output.append(max(n_out, 0))
276
277     fc2_output = []
278     for neuron in fc2_neurons:
279         n_out = 0
280         for j in range(len(fc1_output)):
281             n_out += fc1_output[j] * neuron[j]
282         fc2_output.append(n_out)
283
284     gradients = []
285     sm_output = []
286     sm_sum = 0.
287     max_v = max(fc2_output)
288     for output in fc2_output:
289         sm_sum += math.exp(output - max_v)
290
291     for output in fc2_output:
292         sm_output.append(math.exp(output - max_v) / sm_sum)
293

```



```

294     for output in sm_output:
295         gradients.append(output)
296
297
298     gradients[0] -= 1
299     loss3 = -math.log(gradients[0] + 1.)
300     print('Numerical gradient:\t' + str((loss2 - loss3) / (2*eps)) +
          '\n')
301
302
303 print('\n--- FC0 OUT ---')
304 print('Neuron\t\tActivation')
305 for i in range(len(fc0_output)):
306     print(str(i) + "\t\t\t" + str(fc0_output[i]))
307
308 print('\n--- FC1 OUT ---')
309 print('Neuron\t\tActivation')
310 for i in range(len(fc1_output)):
311     print(str(i) + "\t\t\t" + str(fc1_output[i]))
312
313 print('\n--- FC2 OUT ---')
314 print('Neuron\t\tActivation')
315 for i in range(len(fc2_output)):
316     print(str(i) + "\t\t\t" + str(fc2_output[i]))
317
318 print('\n--- SOFTMAX OUT ---')
319 print('Neuron\t\tActivation')
320 for i in range(len(sm_output)):
321     print(str(i) + "\t\t\t" + str(sm_output[i]))
322
323 print('\n--- FC2 NEURON GRADIENTS ---')
324 for i in range(len(gradients)):
325     print(str(i) + ": " + str(gradients[i]))
326
327
328 print('\n--- FC2 WEIGHT GRADIENTS ---')
329 for i in range(len(weight_grad)):
330     print("Neuron " + str(i))
331     for j in range(len(weight_grad[i])):
332         print(str(j) + ": " + str(weight_grad[i][j]))
333
334 print('\n--- FC1 NEURON GRADIENTS ---')
335 for i in range(len(fc1_grad)):
336     print(str(i) + ": " + str(fc1_grad[i]))
337
338 print('\n--- FC1 WEIGHT GRADIENTS ---')
339 for i in range(64):
340     print("Neuron " + str(i))
341     for j in range(98):
342         print(str(j) + ": " + str(fc1_w_grad[i][j]))
343
344
345 print('\n--- FC0 NEURON GRADIENTS ---')
346 for i in range(len(fc0_grad)):
347     print(str(i) + ": " + str(fc0_grad[i]))

```

```

348
349 print('\n--- FC0 WEIGHT GRADIENTS ---')
350 for i in range(98):
351     print("Neuron " + str(i))
352     for j in range(784):
353         print(str(j) + ": " + str(fc0_w_grad[i][j]))
354
355
356 #Plot weight gradient distribution
357 no_zero_fc0 = []
358 no_zero_fc1 = []
359 no_zero_fc2 = []
360
361 for i in range(len(fc0_w_grad)):
362     for grad in fc0_w_grad[i]:
363         if grad != 0.0:
364             no_zero_fc0.append(grad)
365
366 for i in range(len(fc1_w_grad)):
367     for grad in fc1_w_grad[i]:
368         if grad != 0.0:
369             no_zero_fc1.append(grad)
370
371 for i in range(len(weight_grad)):
372     for grad in weight_grad[i]:
373         if grad != 0.0:
374             no_zero_fc2.append(grad)
375
376 print('Fc0: ' + str(len(no_zero_fc0) / (fc0_n_neurons * fc0_fan_in)
377 ))
378 print('Fc1: ' + str(len(no_zero_fc1) / (fc1_n_neurons * fc1_fan_in)
379 ))
380 print('Fc2: ' + str(len(no_zero_fc2) / (fc2_n_neurons * fc2_fan_in)
381 ))
382
383 no_zero_fc0 = np.asarray(no_zero_fc0)
384 no_zero_fc1 = np.asarray(no_zero_fc1)
385 no_zero_fc2 = np.asarray(no_zero_fc2)
386
387 sns.distplot(no_zero_fc0, label='FC0', bins = 70, norm_hist=True)
388 sns.distplot(no_zero_fc1, label='FC1', bins = 70, norm_hist=True)
389 sns.distplot(no_zero_fc2, label='FC2', bins = 70, norm_hist=True)
390 plt.xlabel('Weight Value')
391 plt.ylabel('Number of weights in bin')
392 plt.legend()
393 plt.xlim(-0.06, 0.06)
394 plt.show()

```

APPENDIX E

Weight Generation

This appendix contains the code that was used to generate the Xilinx coefficient files to initialize the weight BRAMs in the hardware model using He Initialization. This code was written in Python.

E.1 weight_coeff.py

```
1  from random import seed
2  from random import gauss
3  import math
4
5  n_neurons = 64
6  params_per_neuron = 98
7  r_width = 8
8  fan_in = 98
9
10 def intToBinaryString(x, l):
11     str = ""
12     neg = False
13     if x < 0:
14         x += 2 ** (l - 1)
15         neg = True
16
17
```

```
18     while x:
19         if int(x) & 1:
20             str = "1" + str
21         else:
22             str = "0" + str
23         x = int(x / 2)
24
25     if neg:
26         str = "1" + str
27     while (len(str) != 1):
28         str = "0" + str
29     return str
30
31 params = []
32 binary_params = []
33
34 for i in range(n_neurons):
35     for j in range(params_per_neuron):
36         param = gauss(0, math.sqrt(2/(fan_in - 1)))
37         params.append(param)
38
39 for p in params:
40     b = int(p * (2**17))
41     binary_params.append(intToBinaryString(b, 18))
42 print(params[0])
43 print(binary_params[0])
44
45 contents = "memory_initialization_radix=2;\n"
46         nmemory_initialization_vector=\n"
47 cnt = 0
48 for b in binary_params:
49     contents += str(b)
50     cnt += 1
51     if cnt == r_width:
52         contents += ",\n"
53         cnt = 0
54 contents = contents[:-2] + ";"
55
56 f = open('output.coe', 'w')
57 f.write(contents)
58 f.close()
```

APPENDIX F

Software Model

This appendix contains the code that was used to implement the software model for this project. The software model was written in C++.

F.1 Header Files

F.1.1 net.h

```
1  #ifndef __NET_H
2  #define __NET_H
3
4  #include <vector>
5  #include <stdint.h>
6  #include "layer.h"
7
8  class Net {
9      std::vector<Layer*> layers;
10     std::vector< std::vector< std::vector<double> > > activations;
11     std::vector< std::vector<double> > batch_output;
12     std::vector< std::vector<double> > ol_gradient;
13     uint32_t batch_size;
14     uint32_t input_size;
15     uint32_t output_size;
```

```

16     double learning_rate;
17     double momentum;
18
19 public:
20
21     std::vector< std::vector<double> > operator() (std::vector <
22         std::vector<double> > input);
23     void addLayer(Layer*);
24     std::vector< std::vector<double> > inference(std::vector< std:::
25         vector<double> > input);
26     double computeLossAndGradients(std::vector<int> labeled);
27     void backpropLoss();
28     void update();
29     void clearSavedData();
30
31     std::vector<double> convLogitToProb(std::vector<double> logits)
32         ;
33     std::vector<double> getPredictions();
34
35     const uint32_t getBatchSize() const { return batch_size; }
36     std::vector< std::vector<double> > getOlGradient() const {
37         return ol_gradient; }
38     std::vector< std::vector< std::vector<double> > >
39         getActivations() const { return activations; }
40     std::vector<Layer*> getLayers() { return layers; }
41     void setLearningRate(double lr) { learning_rate = lr; };
42     const double getLearningRate() const& { return learning_rate;
43         };
44
45     Net(uint32_t in, uint32_t out, uint32_t bs, double lr, double
46         momentum);
47     Net(const Net& net);
48     ~Net();
49 };
50
51 #endif

```

F.1.2 layer.h

```
1  #ifndef __LAYER_H
2  #define __LAYER_H
3
4  #include "neuron.h"
5  #include <vector>
6
7  #define FULLY    0
8  #define CONV     1
9  #define POOL     2
10
11 class Layer {
12 public:
13     bool last_layer;
14     virtual void forward(std::vector<double>) = 0;
15     virtual void forward(std::vector<double>, bool) = 0;
16     virtual std::vector< std::vector<double> > backward(std::vector
        < std::vector<double> >,
17         std::vector< std::vector<double> >,
18         std::vector< std::vector<double> >) = 0;
19     virtual void updateWeights(double lr, double momentum) = 0;
20     virtual void clearData() = 0;
21     virtual const std::vector<double>& getOutput() = 0;
22     virtual std::vector<Neuron>& getNeurons() = 0;
23     virtual void setNeurons (const std::vector<Neuron>& n) = 0;
24     virtual int getType() = 0;
25     virtual ~Layer() {};
26 };
27
28 #endif
```

F.1.3 convolutional.h

```

1  #ifndef __CONVOLUTIONAL_H
2  #define __CONVOLUTIONAL_H
3
4  #include <vector>
5  #include <stdint.h>
6  #include "layer.h"
7  #include "neuron.h"
8
9  class ConvLayer : public Layer {
10
11     uint32_t dim;
12     uint32_t filt_size;
13     uint32_t stride;
14     uint32_t padding;
15     uint32_t in_channels;
16     uint32_t out_channels;
17     uint32_t dim_o;
18
19     std::vector<Neuron> neurons;
20     std::vector<double> output;
21
22 public:
23     ConvLayer (uint32_t dim, uint32_t filt_size, uint32_t stride,
24               uint32_t padding, uint32_t in_channels, uint32_t
25               out_channels);
26     ConvLayer(const ConvLayer& A);
27     ~ConvLayer ();
28
29     void forward(std::vector<double>);
30     void forward(std::vector<double>, bool in);
31     std::vector< std::vector<double> > backward (std::vector< std::
32     vector<double> >,
33     std::vector< std::vector<double> >,
34     std::vector< std::vector<double> >);
35     void updateWeights(double lr, double momentum);
36     void clearData();
37     std::vector<Neuron>& getNeurons() { return neurons; };
38     const std::vector<double>& getOutput() { return output; };
39     const uint32_t getDim() const { return dim; }
40     const uint32_t getInChannels() const { return in_channels; }
41     const uint32_t getOutChannels() const { return out_channels; }
42     const uint32_t getFiltSize() const { return filt_size; }
43     std::vector<double> getWindowPixels(const std::vector<double>&
44     input, int row, int col);
45     int getType() { return CONV; }
46
47     void setNeurons (const std::vector<Neuron>& n) {neurons = n;}
48
49 };
50
51 #endif

```


F.1.4 fullyconnected.h

```

1  #ifndef __FULLYCONNECTED_H
2  #define __FULLYCONNECTED_H
3
4  #include <stdint.h>
5  #include <vector>
6
7  #include "neuron.h"
8  #include "layer.h"
9
10 class FullyConnected : public Layer {
11
12     uint32_t input_size;
13     uint32_t output_size;
14     std::vector<Neuron> neurons;
15     std::vector<double> output;
16
17 public:
18     FullyConnected(uint32_t in, uint32_t out);
19     FullyConnected(const FullyConnected& x) {
20         input_size = x.input_size;
21         output_size = x.output_size;
22         neurons = x.neurons;
23     }
24     ~FullyConnected() {}
25
26     void forward(std::vector<double> input);
27     void forward(std::vector<double> input, bool last_layer);
28
29     std::vector< std::vector<double> > backward(
30         std::vector< std::vector<double> > gradients_ps,
31         std::vector< std::vector<double> > in_activations,
32         std::vector< std::vector<double> > out_activations);
33
34     void updateWeights(double lr, double momentum);
35     void clearData();
36
37     void setNeurons (const std::vector<Neuron>& n) {neurons = n;}
38
39     const std::vector<double>& getOutput() { return output; }
40     std::vector<Neuron>& getNeurons() { return neurons; }
41     int getType() { return FULLY; }
42 };
43
44 #endif

```

F.1.5 pooling.h

```

1  #ifndef __POOLING_H
2  #define __POOLING_H
3
4  #include <stdint.h>
5  #include <vector>
6  #include "layer.h"
7  #include "neuron.h"
8
9  class PoolingLayer : public Layer {
10     uint32_t dim_i;
11     uint32_t dim_o;
12     uint32_t channels;
13     std::vector<Neuron> placeholder;
14     std::vector<double> output;
15
16 public:
17     PoolingLayer(uint32_t d_i, uint32_t d_o, uint32_t c) :
18         dim_i(d_i), dim_o(d_o), channels(c),
19         output(std::vector<double>(d_o * d_o * c, 0)) {last_layer =
20             false;}
21     PoolingLayer(const PoolingLayer& p) {
22         dim_i = p.dim_i;
23         dim_o = p.dim_o;
24         channels = p.channels;
25         output = p.output;
26     }
27     ~PoolingLayer() {};
28
29     void forward(std::vector<double>);
30     void forward(std::vector<double>, bool);
31     std::vector< std::vector<double> > backward(std::vector< std::
32         vector<double> >,
33         std::vector< std::vector<double> >,
34         std::vector< std::vector<double> >);
35
36     std::vector<double> getWindowPixels (const std::vector<double>&
37         input,
38         uint32_t ch, uint32_t row, uint32_t col);
39
40     void updateWeights(double lr, double momentum) {};
41     void clearData() {}
42     const std::vector<double>& getOutput() { return output; }
43     std::vector<Neuron>& getNeurons() { return placeholder; }
44     void setNeurons (const std::vector<Neuron>& n) {};
45     int getType() { return POOL; };
46 };
47 #endif

```

F.1.6 neuron.h

```
1  #ifndef __NEURON_H
2  #define __NEURON_H
3
4  #include <stdint.h>
5  #include <vector>
6
7  class Neuron {
8      std::vector<double> weights;
9      std::vector<double> gradient_per_weight;
10     std::vector<double> momentum_per_weight;
11     double offset_gradient;
12     double offset_momentum;
13     double offset;
14     double de_dnet;
15     uint32_t fan_in;
16     double net;
17     double activation;
18
19 public:
20
21     Neuron(uint32_t in);
22     Neuron(const Neuron& n);
23     ~Neuron();
24
25     void initWeights(); // He initialization
26     double computeNet(std::vector<double> input);
27     double computeActivation();
28     void calculateGradient(double grad, std::vector<double> act_in,
29                           double act_out, bool last_layer);
30     void updateWeights(double lr, double momentum);
31     void clearBackwardData();
32
33     const double& getActivation() { return activation; }
34     const double& getSensitivity() { return de_dnet; }
35     void setOffset(double offset) { this->offset = offset; }
36     const double& getOffset() { return offset; }
37     void setWeights(std::vector<double> weights) { this->weights =
        weights; }
38     const std::vector<double>& getWeights() { return weights; }
39     const std::vector<double>& getGradients() { return
        gradient_per_weight; }
40 };
41
42 #endif
```

F.1.7 parse_data.h

```
1  #ifndef __PARSE_DATA_H
2  #define __PARSE_DATA_H
3
4  #include <string>
5  #include <vector>
6  #include <stdint.h>
7
8  std::vector<int> readLabels(std::string filename);
9  std::vector< std::vector<double> > readImages(std::string filename)
10     ;
11 #endif
```

F.2 Source Files

F.2.1 main.cpp

```

1  #include <iostream>
2  #include <random>
3  #include <time.h>
4  #include <chrono>
5
6  #include "convolutional.h"
7  #include "fullyconnected.h"
8  #include "pooling.h"
9  #include "parse_data.h"
10 #include "layer.h"
11 #include "net.h"
12
13 double printAccuracy(Net& net, std::vector< std::vector<double> >&
    in, std::vector<int>& out);
14 void trainNet(Net& net, std::vector< std::vector<double> >& in, std
    ::vector<int>& out,
15     std::vector< std::vector<double> >& in_test, std::vector<int>&
    out_test, int n_epochs,
16     int epochs_per_change, double geometric_rate);
17 int main () {
18
19     std::cout << "Running software model...\n";
20
21     std::vector< std::vector<double> > trainX;
22     std::vector<int> trainY;
23     std::vector< std::vector<double> > testX;
24     std::vector<int> testY;
25     trainX = readImages("data/train-images.idx3-ubyte");
26     trainY = readLabels("data/train-labels.idx1-ubyte");
27     testX = readImages("data/t10k-images.idx3-ubyte");
28     testY = readLabels("data/t10k-labels.idx1-ubyte");
29     int n_epochs = 50;
30
31     int input_size = 28*28;
32     int output_size = 10;
33     int batch_size = 1;
34     double momentum = 0.9;
35     double lr = 0.001;
36     Net net(input_size, output_size, batch_size, lr, momentum);
37
38     trainX = std::vector< std::vector<double> > (trainX.begin(),
    trainX.begin() + 60000);
39     trainY = std::vector<int> (trainY.begin(), trainY.begin() +
    60000);
40
41
42     testX = std::vector< std::vector<double> > (testX.begin(),
    testX.begin() + 10000);

```

```

43     testY = std::vector<int> (testY.begin(), testY.begin() + 10000)
44         ;
45     /*
46     Convolutional configuration
47
48     Layer* conv1 = new ConvLayer(28, 3, 1, 1, 1, 8);
49     Layer* pool1 = new PoolingLayer(28, 14, 8);
50     Layer* conv2 = new ConvLayer(14, 3, 1, 1, 8, 16);
51     Layer* pool2 = new PoolingLayer(14, 7, 16);
52     Layer* fc1 = new FullyConnected(16*7*7, 64);
53     Layer* fc2 = new FullyConnected(64, 10);
54
55     net.addLayer(conv1);
56     net.addLayer(pool1);
57     net.addLayer(conv2);
58     net.addLayer(pool2);
59     net.addLayer(fc1);
60     net.addLayer(fc2);
61 */
62
63     Layer* fc1 = new FullyConnected(28*28, 98);
64     Layer* fc2 = new FullyConnected(98, 64);
65     Layer* fc3 = new FullyConnected(64, 10);
66
67     net.addLayer(fc1);
68     net.addLayer(fc2);
69     net.addLayer(fc3);
70
71     trainNet(net, trainX, trainY, testX, testY, n_epochs, 25, .1);
72
73     printAccuracy(net, testX, testY);
74 }
75
76 void trainNet(Net& net, std::vector< std::vector<double> >& in, std
77 ::vector<int>& out,
78     std::vector< std::vector<double> >& in_test, std::
79     vector<int>& out_test, int n_epochs,
80     int epochs_per_change, double geometric_rate) {
81     std::cout << "Starting Accuracy" << std::endl;
82     printAccuracy(net, in_test, out_test);
83     std::cout <<std::endl;
84     clock_t start, end, diff;
85     start = clock();
86     for (int i = 0; i <= n_epochs; i++) {
87         double train_loss = 0.0;
88         int batch_size = net.getBatchSize();
89         int lb = 0;
90         int ub = batch_size;
91         int size = in.size();
92         while (ub <= size) {
93             /* Get the batch */
94             std::vector< std::vector<double> >::iterator startX =
95                 in.begin() + lb;

```

```

94         std::vector< std::vector<double> >::iterator endX = in.
           begin() + ub;
95         std::vector<int>::iterator startY = out.begin() + lb;
96         std::vector<int>::iterator endY = out.begin() + ub;
97
98         std::vector< std::vector<double> > in_batch(startX,
           endX);
99         std::vector<int> out_batch(startY, endY);
100        /* Train by batch size! */
101        net(in_batch);
102        train_loss += net.computeLossAndGradients(out_batch);
103
104        net.backpropLoss();
105        net.update();
106        net.clearSavedData();
107
108        lb += batch_size;
109        ub += batch_size;
110    }
111    end = clock();
112    diff = end - start;
113    std::cout << "Epoch: " << i << std::endl;
114    std::cout << "\n--- Training Stats ---\n";
115    train_loss = printAccuracy(net, in, out);
116    std::cout << "Loss: " << train_loss / (double)in.size() <<
           std::endl;
117    std::cout << "\n--- Test Stats ---\n";
118    double test_loss = printAccuracy(net, in_test, out_test);
119    std::cout << "Elapsed time: " << (float)diff /
           CLOCKS_PER_SEC << std::endl;
120    std::cout << "Loss: " << test_loss / (double)in_test.size()
           << std::endl << std::endl;
121    if ( (i + 1) % epochs_per_change == 0) {
122        std::cout << "Learning rate changed from " << net.
           getLearningRate();
123        net.setLearningRate(net.getLearningRate() *
           geometric_rate);
124        std::cout << " to " << net.getLearningRate() << std::
           endl << std::endl;
125    }
126    }
127 }
128
129 double printAccuracy(Net& net, std::vector< std::vector<double> >&
    in, std::vector<int>& out) {
130     auto result = net(in);
131     int corr = 0;
132     for (size_t i = 0; i < result.size(); i++) {
133         int max_idx = 0;
134         double max = result[i][0];
135         for (size_t j = 1; j < result[i].size(); j++) {
136             if (result[i][j] > max) {
137                 max_idx = j;
138                 max = result[i][j];
139             }

```

```
140     }
141     if (max_idx == out[i]) {
142         corr++;
143     }
144 }
145 double loss = net.computeLossAndGradients(out);
146 net.clearSavedData();
147 std::cout << "Total correct: " << corr << " / " << result.size
148     () << std::endl;
149 std::cout << "Accuracy: " << (double)corr / result.size() <<
150     std::endl;
151 return loss;
152 }
```


F.2.2 net.cpp

```

1  #include "net.h"
2  #include <string>
3  #include <math.h>
4  #include <algorithm>
5  #include <iostream>
6  #include "convolutional.h"
7  #include "fullyconnected.h"
8
9  void Net::addLayer(Layer* layer) {
10     if (layers.size()) {
11         layers[layers.size() - 1]->last_layer = false;
12     }
13     layers.push_back(layer);
14     layers[layers.size() - 1]->last_layer = true;
15 }
16
17 std::vector< std::vector<double> > Net::operator() (std::vector<
    std::vector<double> > input) {
18     return inference(input);
19 }
20
21 std::vector< std::vector<double> > Net::inference(std::vector< std
    ::vector<double> > input) {
22     batch_output = std::vector< std::vector<double> >();
23     activations = std::vector< std::vector< std::vector<double> > >
        >();
24     for (size_t i = 0; i < layers.size() + 1; i++) {
25         activations.push_back(std::vector< std::vector<double> >())
            ;
26     }
27
28     for (std::vector<double> in : input) {
29         if (in.size() != input_size) {
30             std::cout << "Input size does not match, expected: " +
                std::to_string(input_size) +
31             ", got: " + std::to_string(in.size()) << std::endl;
32             exit(1);
33         }
34
35         for (size_t i = 0; i < layers.size(); i++) {
36             Layer*&l = layers[i];
37             /*double max = *(std::max_element(in.begin(), in.end()))
                );
38             for (double& e : in) {
39                 e /= max;
40             }*/
41             activations[i].push_back(in);
42             if (i == layers.size() - 1) {
43                 l->forward(in, true);
44             }
45             else {
46                 l->forward(in);
47             }

```

```

48         in = l->getOutput();
49     }
50
51
52     if (in.size() != output_size) {
53         std::cout << "Output size does not match, expected: " +
            std::to_string(output_size) + ", got: " + std::
            to_string(in.size());
54         std::cout << std::endl;
55         exit(1);
56     }
57     activations[layers.size()].push_back(in);
58     std::vector<double> output = convLogitToProb(in);
59     //std::vector<double> output = in;
60     batch_output.push_back(output);
61 }
62 return activations[layers.size()];
63 //return batch_output;
64 }
65
66 double Net::computeLossAndGradients(std::vector<int> labeled) {
67     if (labeled.size() != batch_output.size()) {
68         std::cout << "Labeled data size does not match the net's
            output size, expected: " +
69             std::to_string(batch_output.size()) + ", got: "
            + std::to_string(labeled.size()) << std::
            endl;
70         std::cout << std::endl;
71         exit(1);
72     }
73     ol_gradient = std::vector< std::vector<double> > ();
74     double loss = 0;
75     // Compute cross entropy loss for each output
76     // CrossEntropy loss -q(x) * log(p(x))
77     // q(x) is true distribution, so it is 1 for our labeled data
    on the correct sample
78     for (size_t i = 0; i < labeled.size(); i++) {
79         std::vector<double> gradient(output_size, 0);
80         unsigned short label = labeled[i];
81
82         for (size_t j = 0; j < output_size; j++) {
83             //double f_out = (tanh(batch_output[i][j] / 2.))/2. +
            0.5;
84             if (j == label) {
85                 //gradient[j] = batch_output[i][j] * (1 -
                    batch_output[i][j]);
86                 gradient[j] = batch_output[i][j] - 1;
87
88                 //gradient[j] = (f_out - 1) * (f_out*(1-f_out));
89             }
90             else {
91                 //gradient[j] = -batch_output[i][j] * batch_output[
                    i][label];
92                 gradient[j] = batch_output[i][j];
93                 //gradient[j] = (f_out) * (f_out*(1-f_out));

```

```

94         }
95     }
96     loss += -log(batch_output[i][label]);
97     ol_gradient.push_back(gradient);
98 }
99
100
101 // Compute mean square error
102 /*for (size_t i = 0; i < labeled.size(); i++) {
103     unsigned short label = labeled[i];
104     std::vector<double> gradient(output_size, 0);
105     for (size_t j = 0; j < output_size; j++) {
106         double err, grad;
107         if (j == label) {
108             grad = 1 - batch_output[i][j];
109             err = 0.5 * pow(grad, 2);
110         }
111         else {
112             grad = 0 - batch_output[i][j];
113             err = 0.5 * pow(grad, 2);
114         }
115         gradient[j] = grad;
116         loss += err;
117     }
118     ol_gradient.push_back(gradient);
119 }*/
120
121 return loss;
122 }
123
124 // Backpropagate the gradients of the error
125 void Net::backpropLoss() {
126     std::vector< std::vector<double> > gradients = ol_gradient;
127     std::vector< std::vector<double> > sens;
128     // Outer layer gradients is just the loss
129     for (int i = layers.size() - 1; i >= 0; i--) {
130         //std::cout << "Grad len: " << gradients[0].size() << "
131             Layer len: " << layers[i]->getNeurons().size() << "
132             Prev: " << activations[i].size() << " Next: " <<
133             activations[i+1].size() << std::endl;
134         sens = layers[i]->backward(gradients, activations[i],
135             activations[i + 1]);
136         // fully connected gradients
137         // if fully connected check, on layers[i]
138         gradients = std::vector< std::vector<double> >();
139         Layer* l = layers[i];
140         std::vector<Neuron> neurons = l->getNeurons();
141
142         if (l->getType() == POOL) {
143             gradients = sens;
144         }
145         else {
146             for (size_t j = 0; j < sens.size(); j++) {
147                 if (l->getType() == FULLY) {

```

```

144         // The gradient of neuron i in prev layer is
145         // the sum of the weights[i] * de_dnet of all
146         // neurons in layer j
147         std::vector<double> grad(neurons[0].getWeights
148             ().size(), 0);
149         for (size_t k = 0; k < neurons.size(); k++) {
150             std::vector<double> weights = neurons[k].
151                 getWeights();
152             for (size_t l = 0; l < weights.size(); l++)
153                 grad[l] += sens[j][k] * weights[l];
154         }
155         gradients.push_back(grad);
156     }
157     else if (l->getType() == CONV) {
158         ConvLayer* cl = dynamic_cast<ConvLayer* >(l);
159         // If the sensitivities of i + 1 layer were
160         // from a convolution, then the
161         // neurons for layer i only need to do weights[
162         // i] * de_dnet for
163         // the relevant windows that the activation was
164         // in
165         int dim = cl->getDim();
166         int in_chan = cl->getInChannels();
167         int out_chan = cl->getOutChannels();
168
169         int num_neurons = dim * dim * in_chan; //
170         // amount of gradients to give previous layer
171         std::vector<double> grad(num_neurons, 0);
172
173         for (int k = 0; k < num_neurons; k++) {
174             // for each window it goes to... need to
175             // know which weight to use
176             int chan = k / (dim * dim);
177             int row = (k - (chan * dim * dim)) / dim;
178             int col = (k - (chan * dim * dim + row *
179                 dim)) % dim;
180             int filt_size = cl->getFiltSize();
181             int filt_sq = filt_size * filt_size;
182             // Iterate over the neurons in the window
183             // for this gradient
184             int dim_sq = dim * dim;
185             int start_row = row - (filt_size / 2);
186             int end_row = row + (filt_size / 2);
187             int start_col = col - (filt_size / 2);
188             int end_col = col + (filt_size / 2);
189
190             /*std::cout << "Row: " << row << " Col: "
191                 << col << " Start row: " << start_row
192                 <<
193                 " End row: " << end_row << " Start col: "
194                 << start_col << " End col: " << end_col <<
195                 " j: " << j

```

```

184         << " Sens size: " << sens.size() << std::
185             endl;*/
186         for (int o = 0; o < out_chan; o++) {
187             int count = 0;
188             for (int m = start_row; m <= end_row; m
189                 ++){
190                 for (int n = start_col; n <=
191                     end_col; n++) {
192                     if (m < 0 || m >= dim || n < 0
193                         || n >= dim) {
194                         count++;
195                         continue;
196                     }
197                     int o_neur_idx = o * dim_sq + m
198                         * dim + n;
199                     int filt_offset = (filt_sq) - (
200                         count + 1);
201                     int weight_idx = (chan *
202                         filt_sq) + filt_offset;
203
204                     grad[k] += sens[j][o_neur_idx]
205                         * neurons[o_neur_idx].
206                         getWeights()[weight_idx];
207                     count++;
208                 }
209             }
210             gradients.push_back(grad);
211         }
212     }
213 }
214
215 void Net::update() {
216     double effective_lr = learning_rate / batch_size;
217     for (int i = layers.size() - 1; i >= 0; i--) {
218         layers[i]->updateWeights(effective_lr, momentum);
219     }
220 }
221
222 std::vector<double> Net::convLogitToProb(std::vector<double> logits
223 ) {
224     double sum = 0;
225     double max = *std::max_element(logits.begin(), logits.end());
226     for (auto l : logits) {
227         sum += exp(l - max);
228     }
229     std::vector<double> prob;
230
231     for (auto l : logits) {
232         prob.push_back(exp(l - max) / sum);
233     }

```

```

229     return prob;
230 }
231
232 std::vector<double> Net::getPredictions() {
233     std::vector<double> preds;
234     for (size_t i = 0; i < batch_output.size(); i++) {
235         int pred_class = 0;
236         double pred_max = batch_output[i][0];
237         for (size_t j = 1; j < output_size; j++) {
238             if (batch_output[i][j] > pred_max) {
239                 pred_class = j;
240             }
241         }
242         preds.push_back(pred_class);
243     }
244     return preds;
245 }
246
247 void Net::clearSavedData() {
248     activations = std::vector< std::vector< std::vector<double> >
249         >();
250     batch_output = std::vector< std::vector<double> >();
251     ol_gradient = std::vector< std::vector<double> >();
252     /*for (int i = layers.size() - 1; i >= 0; i--) {
253         layers[i]->clearData();
254     }*/
255     for (Layer* l : layers) {
256         l->clearData();
257     }
258 }
259
260 Net::Net(uint32_t in, uint32_t out, uint32_t bs, double lr, double
    moment) {
261     layers = std::vector<Layer*>();
262     input_size = in;
263     output_size = out;
264     batch_size = bs;
265     learning_rate = lr;
266     momentum = moment;
267 }
268
269 Net::Net(const Net& net) {
270     layers = net.layers;
271     input_size = net.input_size;
272     output_size = net.output_size;
273 }
274
275 Net::~Net() {
276     for (size_t i = 0; i < layers.size(); i++) {
277         delete layers[i];
278     }
279 }

```

F.2.3 convolutional.cpp

```

1  #include "convolutional.h"
2
3  #include <iostream>
4
5  void ConvLayer::forward(std::vector<double> input) {
6      uint32_t h_steps = 1 + ((dim + (padding * 2) - filt_size) /
7          stride);
8
9      if (input.size() != h_steps * h_steps * in_channels) {
10         std::cout << "Wrong input size for convolutional layer\n";
11         exit(1);
12     }
13
14     unsigned start = (filt_size / 2) - padding;
15
16     for (unsigned int i = 0; i < out_channels; i++) { // channel of
17         output
18         for (unsigned int k = 0; k < dim_o; k++) { // row
19             for (unsigned int l = 0; l < dim_o; l++) { // column
20                 int row = k * stride + start;
21                 int col = l * stride + start;
22                 std::vector<double> pixels = getWindowPixels(input,
23                     row, col);
24                 int out_idx = i * dim_o * dim_o + k * dim_o + l;
25                 neurons[out_idx].computeNet(pixels);
26                 output[out_idx] = neurons[out_idx].
27                     computeActivation();
28             }
29         }
30     }
31 }
32
33 void ConvLayer::forward(std::vector<double> input, bool in) {
34     forward(input);
35 }
36
37 std::vector<double> ConvLayer::getWindowPixels(const std::vector<
38     double>& input, int row, int col) {
39     std::vector<double> pixels;
40
41     int start_row = row - (filt_size / 2);
42     int end_row = row + (filt_size / 2);
43     int start_col = col - (filt_size / 2);
44     int end_col = col + (filt_size / 2);
45     int dim_sq = dim * dim;
46
47     for (unsigned int ch = 0; ch < in_channels; ch++) {
48         for (int i = start_row; i <= end_row; i++) {
49             for (int j = start_col; j <= end_col; j++) {
50                 int idx = ch * dim_sq + i * dim + j;
51                 if (i < 0 || i >= (int)dim || (j < 0 || j >= (int)
52                     dim)) {

```

```

48         // in padding region, just push 0
49         pixels.push_back(0);
50     }
51     else {
52         pixels.push_back(input[idx]);
53     }
54 }
55 }
56 }
57 return pixels;
58 }
59
60 std::vector< std::vector<double> > ConvLayer::backward (std::vector
    < std::vector<double> > gradients,
61
62                                     std::vector< std:::
        vector<double>
        >
        in_activations,
        std::vector< std:::
        vector<double>
        >
        out_activations
        ) {
63
64     std::vector< std::vector<double> > sensitivity;
65     int dim_sq = dim_o * dim_o;
66
67     unsigned start = (filt_size / 2) - padding;
68     for (size_t i = 0; i < gradients.size(); i++) {
69         std::vector<double> single_sens;
70         for (size_t j = 0; j < out_channels; j++) {
71             for (size_t k = 0; k < dim_o; k++) {
72                 for (size_t l = 0; l < dim_o; l++) {
73                     // Neuron at row k, column l, in output channel
74                     j
75                     int idx = l + k * dim_o + j * dim_sq;
76                     int row = k * stride + start;
77                     int col = l * stride + start;
78                     std::vector<double> in_act = getWindowPixels(
79                         in_activations[i], row, col);
80                     neurons[idx].calculateGradient(gradients[i][idx],
81                         in_act, out_activations[i][idx],
82                         last_layer);
83                     single_sens.push_back(neurons[idx].
84                         getSensitivity());
85                 }
86             }
87         }
88         sensitivity.push_back(single_sens);
89     }
90     return sensitivity;
91 }
92
93 void ConvLayer::updateWeights(double lr, double momentum) {

```



```

90     for (Neuron& n : neurons) {
91         n.updateWeights(lr, momentum);
92     }
93 }
94
95
96
97 void ConvLayer::clearData() {
98     for (Neuron& n : neurons) {
99         n.clearBackwardData();
100     }
101 }
102
103
104 ConvLayer::ConvLayer(uint32_t d, uint32_t fsize, uint32_t str,
105     uint32_t pad, uint32_t in_ch, uint32_t out_ch) {
106     last_layer = false;
107     dim = d;
108     in_channels = in_ch;
109     out_channels = out_ch;
110     filt_size = fsize;
111     stride = str;
112     padding = pad;
113     uint32_t weights_per_fmap = filt_size * filt_size * in_channels
114         ;
115     uint32_t steps = 1 + ((dim + (padding * 2) - filt_size) /
116         stride);
117     uint32_t num_neurons = steps * steps * out_channels;
118     dim_o = steps;
119
120
121     output.resize(num_neurons);
122     neurons.reserve(num_neurons);
123
124     for (uint32_t i = 0; i < num_neurons; i++) {
125         Neuron n(weights_per_fmap);
126         n.initWeights();
127         neurons.push_back(n);
128     }
129 }
130
131 ConvLayer::ConvLayer(const ConvLayer& a) {
132     dim = a.dim;
133     in_channels = a.in_channels;
134     out_channels = a.out_channels;
135     filt_size = a.filt_size;
136     neurons = a.neurons;
137     output = a.output;
138 }
139
140 ConvLayer::~ConvLayer() {
141

```

142

}

F.2.4 fullyconncted.cpp

```

1  #include "fullyconnected.h"
2
3  #include <iostream>
4
5  void FullyConnected::forward(std::vector<double> input) {
6      for (size_t i = 0; i < output_size; i++) {
7          neurons[i].computeNet(input);
8          output[i] = neurons[i].computeActivation();
9      }
10 }
11
12 void FullyConnected::forward(std::vector<double> input, bool
    last_layer) {
13     if (last_layer) {
14         for (size_t i = 0; i < output_size; i++) {
15             output[i] = neurons[i].computeNet(input);
16         }
17     }
18     else {
19         forward(input);
20     }
21 }
22
23 std::vector< std::vector<double> > FullyConnected::backward(
24     std::vector< std::vector<double> > gradients,
25     std::vector< std::vector<double> >
        in_activations,
26     std::vector< std::vector<double> >
        out_activations) {
27     std::vector< std::vector<double> > sensitivity;
28     for (size_t i = 0; i < gradients.size(); i++) {
29         std::vector<double> sngl_sens;
30         for (size_t j = 0; j < neurons.size(); j++) {
31             neurons[j].calculateGradient(gradients[i][j],
                in_activations[i], out_activations[i][j],
                last_layer);
32             sngl_sens.push_back(neurons[j].getSensitivity());
33         }
34         sensitivity.push_back(sngl_sens);
35     }
36     return sensitivity;
37 }
38
39 void FullyConnected::updateWeights(double lr, double momentum) {
40     for (Neuron& n : neurons) {
41         n.updateWeights(lr, momentum);
42     }
43 }
44
45 void FullyConnected::clearData() {
46     for (Neuron& n : neurons) {
47         n.clearBackwardData();
48     }

```

```
49 }  
50  
51 FullyConnected::FullyConnected(uint32_t in, uint32_t out) :  
    input_size(in), output_size(out) {  
52     last_layer = false;  
53     neurons.reserve(out);  
54     output.resize(out);  
55     for (size_t i = 0; i < out; i++) {  
56         Neuron n(in);  
57         n.initWeights();  
58         neurons.push_back(n);  
59     }  
60 }
```

F.2.5 pooling.cpp

```

1  #include "pooling.h"
2  #include <iostream>
3  #include <algorithm>
4
5  void PoolingLayer::forward(std::vector<double> in) {
6      for (size_t c = 0; c < channels; c++) {
7          for (size_t i = 0; i < dim_i - 1; i += 2) {
8              for (size_t j = 0; j < dim_i - 1; j += 2) {
9                  int row_o = i / 2;
10                 int col_o = j / 2;
11                 int out_idx = c * dim_o * dim_o + row_o * dim_o +
                             col_o;
12
13                 std::vector<double> pixels = getWindowPixels(in, c,
14                     i, j);
15                 double max = *std::max_element(pixels.begin(),
16                     pixels.end());
17                 output[out_idx] = max;
18             }
19         }
20     }
21
22 void PoolingLayer::forward(std::vector<double> in, bool first) {
23     forward(in);
24 }
25
26 std::vector<double> PoolingLayer::getWindowPixels(const std::vector<
27     double>& input,
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
uint32_t ch,
uint32_t
row,
uint32_t
col) {
    int offset = dim_i * dim_i * ch + dim_i * row + col;
    std::vector<double> pixels;
    pixels.push_back(input[offset]);
    pixels.push_back(input[offset + 1]);
    pixels.push_back(input[offset + dim_i]);
    pixels.push_back(input[offset + dim_i + 1]);
    return pixels;
}

std::vector< std::vector<double> > PoolingLayer::backward(
    std::vector< std::vector<double> > gradients,
    std::vector< std::vector<double> >
        in_activations,
    std::vector< std::vector<double> >
        out_activations) {
    std::vector< std::vector<double> > sensitivity;
    for (size_t i = 0; i < gradients.size(); i++) {
        std::vector<double> sngl_sens(channels * dim_i * dim_i, 0);

```

```

44     for (size_t c = 0; c < channels; c++) {
45         for (size_t j = 0; j < dim_o; j++) {
46             for (size_t k = 0; k < dim_o; k++) {
47                 int offset_o = c * dim_o * dim_o + j * dim_o +
48                     k;
49                 int row_i = j * 2;
50                 int col_i = k * 2;
51                 int offset_i = c * dim_i * dim_i + row_i *
52                     dim_i + col_i;
53                 std::vector<double> pixels = getWindowPixels(
54                     in_activations[i], c, row_i, col_i);
55
56                 double max = *std::max_element(pixels.begin(),
57                     pixels.end());
58
59                 // write the gradients for each pixel
60                 std::vector<double> window_gradients(4, 0);
61                 for (size_t l = 0; l < pixels.size(); l++) {
62                     if (pixels[l] == max) {
63                         window_gradients[l] = gradients[i][
64                             offset_o];
65                     }
66                 }
67                 sngl_sens[offset_i] = window_gradients[0];
68                 sngl_sens[offset_i + 1] = window_gradients[1];
69                 sngl_sens[offset_i + dim_i] = window_gradients
70                     [2];
71                 sngl_sens[offset_i + dim_i + 1] =
72                     window_gradients[3];
73             }
74         }
75     }
76     sensitivity.push_back(sngl_sens);
77 }
78 return sensitivity;
79 }

```

F.2.6 neuron.cpp

```
1  #include "neuron.h"
2
3  #include <random>
4  #include <chrono>
5  #include <math.h>
6  #include <iostream>
7
8  Neuron::Neuron(uint32_t in) {
9      fan_in = in;
10     gradient_per_weight = std::vector<double> (fan_in, 0);
11     momentum_per_weight = std::vector<double> (fan_in, 0);
12     offset_gradient = 0;
13     offset_momentum = 0;
14 }
15
16 Neuron::Neuron(const Neuron& n) {
17     weights = n.weights;
18     gradient_per_weight = n.gradient_per_weight;
19     momentum_per_weight = n.momentum_per_weight;
20     offset_gradient = n.offset_gradient;
21     offset_momentum = n.offset_momentum;
22     offset = n.offset;
23     fan_in = n.fan_in;
24     net = n.net;
25     activation = n.activation;
26 }
27
28 Neuron::~Neuron() {
29 }
30
31
32 /**
33  * Uses He initialization
34  */
35 void Neuron::initWeights() {
36     weights.reserve(fan_in);
37     std::normal_distribution<double> distribution(0, sqrt(2. / (
38         fan_in)));
39
40     static unsigned seed = std::chrono::system_clock::now().
41         time_since_epoch().count();
42     static std::default_random_engine generator(seed);
43
44     for (size_t i = 0; i < fan_in; i++) {
45         weights.push_back(distribution(generator));
46     }
47     offset = distribution(generator);
48 }
49
50 /**
51  * Compute net for the neuron
52  */
53 double Neuron::computeNet(std::vector<double> input) {
```

```

52     if (input.size() != weights.size()) {
53         std::cerr << "Input size did not match size of weights.
           Input size: " <<
54             input.size() << " Weight size: " << weights
               .size() << std::endl;
55         exit(1);
56     }
57
58     // Dot product and sum offset
59     net = offset;
60     for (size_t i = 0; i < fan_in; i++) {
61         net += input[i] * weights[i];
62     }
63
64     return net;
65 }
66
67 /**
68  * Compute output for neuron
69  */
70 double Neuron::computeActivation() {
71     activation = std::max(net, 0.);    // ReLU
72     return activation;
73 }
74
75 void Neuron::calculateGradient(double grad, std::vector<double>
act_in, double act_out, bool last_layer) {
76     double dact_dnet = (!last_layer) ? (act_out > 0) : 1;
77     de_dnet = grad * dact_dnet;
78
79     for (size_t i = 0; i < fan_in; i++) {
80         gradient_per_weight[i] += (de_dnet * act_in[i]);
81     }
82     offset_gradient += de_dnet;
83 }
84
85 void Neuron::updateWeights(double lr, double momentum) {
86     // update weights
87     for (size_t i = 0; i < fan_in; i++) {
88         momentum_per_weight[i] = momentum * momentum_per_weight[i]
           + (lr * -gradient_per_weight[i]);
89         weights[i] += momentum_per_weight[i];
90     }
91     offset_momentum = offset_momentum * momentum + (lr * -
           offset_gradient);
92     offset += offset_momentum;
93     clearBackwardData();
94 }
95
96 void Neuron::clearBackwardData() {
97     gradient_per_weight = std::vector<double>(fan_in, 0);
98     offset_gradient = 0;
99 }

```


F.2.7 parse_data.cpp

```
1  #include "parse_data.h"
2  #include <fstream>
3  #include <iostream>
4
5  std::vector< std::vector<double> > readImages(std::string s) {
6      std::ifstream f(s, std::ios::binary | std::ios::in);
7
8      std::vector< std::vector<double> > data;
9      uint8_t pixel;
10     uint32_t rows;
11     uint32_t cols;
12
13     uint8_t buf[4];
14
15
16     f.read((char*)&buf, 4); // magic number
17     f.read((char*)&buf, 4); // number of images
18     f.read((char*)&buf, 4); // number of rows
19     rows = buf[3] + (buf[2] << 8) + (buf[1] << 16) + (buf[0] << 24)
20     ;
21     f.read((char*)&buf, 4); // number of cols
22     cols = buf[3] + (buf[2] << 8) + (buf[1] << 16) + (buf[0] << 24)
23     ;
24
25     while(!f.eof()) {
26         // read an image
27         std::vector<double> img;
28         for (unsigned i = 0; i < rows * cols; i++) {
29             f.read((char*)&pixel, 1);
30             double p = ((double)pixel / 255.);
31             img.push_back(p);
32         }
33         if (!f.eof()) {
34             data.push_back(img);
35         }
36     }
37     f.close();
38     return data;
39 }
40
41 std::vector<int> readLabels(std::string s) {
42     std::ifstream f(s, std::ios::binary | std::ios::in);
43     std::vector<int> labels;
44     uint8_t label;
45     int32_t res;
46     f.read((char*)&res, 4); // magic number
47     f.read((char*)&res, 4); // number of items
48     f.read((char*)&label, 1);
49     while(!f.eof()) {
50         labels.push_back((int)label);
51         f.read((char*)&label, 1);
52     }
53 }
```

```
52     f.close();  
53     return labels;  
54 }
```

F.3 Testing Files

F.3.1 gradient_check_test.cpp

```

1  #include <gtest/gtest.h>
2
3  #include <iostream>
4  #include <random>
5  #include <chrono>
6
7  #include "../src/convolutional.h"
8  #include "../src/fullyconnected.h"
9  #include "../src/pooling.h"
10 #include "../src/parse_data.h"
11 #include "../src/layer.h"
12 #include "../src/net.h"
13
14 TEST(GradientTest, FCGradientCheck) {
15
16     int input_size = 100;
17     int output_size = 2;
18     int batch_size = 1;
19     double momentum = 0.9;
20     double lr = 0.001;
21     Net net(input_size, output_size, batch_size, lr, momentum);
22
23
24     Layer* fc1 = new FullyConnected(input_size, 98);
25     Layer* fc2 = new FullyConnected(98, 64);
26     Layer* fc3 = new FullyConnected(64, output_size);
27
28     net.addLayer(fc1);
29     net.addLayer(fc2);
30     net.addLayer(fc3);
31
32     // Generate random input and output labels
33     std::vector< std::vector<double> > in;
34     std::vector<int> out;
35
36     std::uniform_real_distribution<double> distribution(-1.0, 1.0);
37     std::uniform_int_distribution<int> distribution_out(0,
38         output_size - 1);
39     static unsigned seed = std::chrono::system_clock::now().
40         time_since_epoch().count();
41     static std::default_random_engine generator(seed);
42     int test_size = 10;
43     for (int i = 0; i < test_size; i++) {
44         std::vector<double> smpl;
45         for (int j = 0; j < input_size; j++) {
46             smpl.push_back(distribution(generator));
47         }
48         in.push_back(smpl);
49         out.push_back(distribution_out(generator));
50     }
51 }

```

```

48     }
49
50     int grad_tests = 100;
51     int num_layers = 3;
52     srand(time(0));
53     for (int i = 0; i < grad_tests; i++) {
54         double sigma = pow(10, -4);
55         auto out_ = net(in);
56         net.computeLossAndGradients(out);
57
58         net.backpropLoss();
59
60         int l_idx = (rand() % num_layers);
61         FullyConnected* fc = (FullyConnected* )net.getLayers()[
            l_idx];
62         std::vector<Neuron>& neurons = fc->getNeurons();
63         int n_idx = (rand() % neurons.size());
64         auto neuron = neurons[n_idx];
65         auto weights = neuron.getWeights();
66         int w_idx = (rand() % neuron.getWeights().size());
67         double grad = neuron.getGradients()[w_idx];
68         net.clearSavedData();
69
70
71         // + sigma loss
72         weights[w_idx] += sigma;
73         neurons[n_idx].setWeights(weights);
74         fc->setNeurons(neurons);
75         net(in);
76         double loss_plus = net.computeLossAndGradients(out);
77         net.clearSavedData();
78
79         // - sigma loss
80         weights[w_idx] -= (sigma + sigma);
81         neurons[n_idx].setWeights(weights);
82         fc->setNeurons(neurons);
83         net(in);
84         double loss_minus = net.computeLossAndGradients(out);
85         net.clearSavedData();
86
87         weights[w_idx] += sigma;
88         neurons[n_idx].setWeights(weights);
89         fc->setNeurons(neurons);
90
91         double num_grad = (loss_plus - loss_minus) / (2 * sigma);
92
93         double rel = std::max(num_grad > 0 ? num_grad : -num_grad,
            grad > 0 ? grad : -grad);
94         rel = rel == 0. ? 1 : rel;
95         double diff = (num_grad - grad) / rel;
96         diff = (diff > 0) ? diff : -diff;
97         ASSERT_LE(diff, 1e-7);
98     }
99 }
100

```

```

101 TEST(GradientTest, ConvGradientCheck) {
102
103     int input_size = 8*8;
104     int output_size = 2;
105     int batch_size = 1;
106     double momentum = 0.9;
107     double lr = 0.001;
108     Net net(input_size, output_size, batch_size, lr, momentum);
109
110     Layer* conv1 = new ConvLayer(8, 3, 1, 1, 1, 3);
111     Layer* pool1 = new PoolingLayer(8, 4, 3);
112     Layer* conv2 = new ConvLayer(4, 3, 1, 1, 3, 6);
113     Layer* fc1 = new FullyConnected(4*4*6, output_size);
114
115     net.addLayer(conv1);
116     net.addLayer(pool1);
117     net.addLayer(conv2);
118     net.addLayer(fc1);
119
120     // Generate random input and output labels
121     std::vector< std::vector<double> > in;
122     std::vector<int> out;
123
124     std::uniform_real_distribution<double> distribution(-1.0, 1.0);
125     std::uniform_int_distribution<int> distribution_out(0,
126         output_size - 1);
127     static unsigned seed = std::chrono::system_clock::now().
128         time_since_epoch().count();
129     static std::default_random_engine generator(seed);
130     int test_size = 10;
131     for (int i = 0; i < test_size; i++) {
132         std::vector<double> smpl;
133         for (int j = 0; j < input_size; j++) {
134             smpl.push_back(distribution(generator));
135         }
136         in.push_back(smpl);
137         out.push_back(distribution_out(generator));
138     }
139
140     int grad_tests = 100;
141     int num_layers = 3;
142     srand(time(0));
143     for (int i = 0; i < grad_tests; i++) {
144         double sigma = pow(10, -4);
145         auto out_ = net(in);
146         net.computeLossAndGradients(out);
147         net.backpropLoss();
148         int l_idx = (rand() % num_layers);
149         while (l_idx == 1) { // don't select the pooling layer
150             l_idx = (rand() % num_layers);
151         }
152         Layer* l = (FullyConnected*)net.getLayers()[l_idx];
153         std::vector<Neuron>& neurons = l->getNeurons();
154         int n_idx = (rand() % neurons.size());
155         auto neuron = neurons[n_idx];

```

```
154     auto weights = neuron.getWeights();
155     int w_idx = (rand() % neuron.getWeights().size());
156     double grad = neuron.getGradients()[w_idx];
157     net.clearSavedData();
158
159     // + sigma loss
160     weights[w_idx] += sigma;
161     neurons[n_idx].setWeights(weights);
162     l->setNeurons(neurons);
163     net(in);
164     double loss_plus = net.computeLossAndGradients(out);
165     net.clearSavedData();
166
167     // - sigma loss
168     weights[w_idx] -= (sigma + sigma);
169     neurons[n_idx].setWeights(weights);
170     l->setNeurons(neurons);
171     net(in);
172     double loss_minus = net.computeLossAndGradients(out);
173     net.clearSavedData();
174
175     weights[w_idx] += sigma;
176     neurons[n_idx].setWeights(weights);
177     l->setNeurons(neurons);
178
179     double num_grad = (loss_plus - loss_minus) / (2 * sigma);
180
181     double rel = std::max(num_grad > 0 ? num_grad : -num_grad,
182                          grad > 0 ? grad : -grad);
182     rel = rel == 0. ? 1 : rel;
183     double diff = (num_grad - grad) / rel;
184     diff = (diff > 0) ? diff : -diff;
185     ASSERT_LE(diff, 1e-7);
186 }
187 }
```

F.3.2 conv_test.cpp

```
1  #include <gtest/gtest.h>
2
3  #include "../src/convolutional.h"
4
5  TEST(ConvTest, TestForward) {
6      ConvLayer conv1(2, 3, 1, 1, 2, 2);
7
8      std::vector<double> input = {    2, 3,
9                                      1, 4,
10
11                                      3, 1,
12                                      5, 0 };
13
14      std::vector<Neuron> neurons;
15
16      for (double i = 0; i < 8; i++) {
17          Neuron n(18);
18          std::vector<double> weights = { 5, 1, 2,
19                                          3, 3, 2,
20                                          4, 1, 1,
21
22                                          2, 3, 5,
23                                          0, 1, 2,
24                                          4, 2, 1};
25          n.setWeights(weights);
26          n.setOffset(4);
27          neurons.push_back(n);
28      }
29      conv1.setNeurons(neurons);
30      conv1.forward(input);
31
32      std::vector<double> outputs = conv1.getOutput();
33
34      ASSERT_EQ(outputs[0], 36.);
35      ASSERT_EQ(outputs[1], 48);
36      ASSERT_EQ(outputs[2], 42);
37      ASSERT_EQ(outputs[3], 41);
38      ASSERT_EQ(outputs[4], 36.);
39      ASSERT_EQ(outputs[5], 48);
40      ASSERT_EQ(outputs[6], 42);
41      ASSERT_EQ(outputs[7], 41);
42  }
```

F.3.3 fullyconnected_test.cpp

```
1  #include <gtest/gtest.h>
2
3  #include "../src/fullyconnected.h"
4
5  TEST(FCTest, TestForward) {
6      FullyConnected fc(3, 4);
7
8      std::vector<Neuron> neurons;
9      for (double i = 0; i < 4; i++) {
10         Neuron n(3);
11         std::vector<double> weights = {i, i, i};
12         n.setWeights(weights);
13         n.setOffset(i + 1.5);
14         neurons.push_back(n);
15     }
16     std::vector<double> input = {1, 2, 3};
17     fc.setNeurons(neurons);
18     fc.forward(input);
19
20     std::vector<double> outputs = fc.getOutput();
21
22     ASSERT_EQ(outputs[0], 0 + 1.5);
23     ASSERT_EQ(outputs[1], 1 + 2 + 3 + 1 + 1.5);
24     ASSERT_EQ(outputs[2], 2 * 1 + 2 * 2 + 2 * 3 + 2 + 1.5);
25     ASSERT_EQ(outputs[3], 3 * 1 + 3 * 2 + 3 * 3 + 3 + 1.5);
26 }
```


F.3.4 neuron_test.cpp

```
1  #include <gtest/gtest.h>
2
3  #include "../src/neuron.h"
4
5  TEST(NeuronTest, InitWeights) {
6      int fan_in = 11;
7      Neuron n(fan_in);
8      n.initWeights();
9      ASSERT_EQ(n.getWeights().size(), fan_in);
10 }
11
12 TEST(NeuronTest, SetWeightsAndGetOutput) {
13     int fan_in = 5;
14     Neuron n(fan_in);
15
16     std::vector<double> weights;
17     std::vector<double> input;
18     double offset = 13;
19     for (int i = 0; i < fan_in; i++) {
20         weights.push_back(i + 1);
21         input.push_back(2 * i + 1);
22     }
23
24     n.setOffset(offset);
25     n.setWeights(weights);
26
27     double result = 1 * 1 + 2 * 3 + 3 * 5 + 4 * 7 + 5 * 9 + 13;
28
29     n.computeNet(input);
30
31     ASSERT_EQ(n.computeActivation(), result);
32 }
```


APPENDIX G

PyTorch Model

This appendix contains the code for the PyTorch version of the implemented neural network. This code was written in Python.

G.1 mnist_model.py

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torch.optim as optim
5 from torchvision import datasets, transforms
6 from torch.autograd import Variable
7 import numpy as np
8 import parse_data
9 from timeit import default_timer as timer
10 from sklearn.model_selection import train_test_split
11 import csv
12
13 class Net(nn.Module):
14
15     def __init__(self):
16         super(Net, self).__init__()
17         self.fc0 = nn.Linear(28*28, 2*7*7)
18         self.fc1 = nn.Linear(2*7*7, 64)
19         self.fc2 = nn.Linear(64, 10)
```

```

20
21     def forward(self, x):
22         in_size = x.size(0)
23         x = x.view(in_size, -1)
24         x = F.relu(self.fc0(x))
25         x = F.relu(self.fc1(x))
26         x = self.fc2(x)
27         return F.log_softmax(x, dim=0)
28
29 use_cuda = torch.cuda.is_available()
30 #use_cuda = False
31 device = torch.device("cuda:0" if use_cuda else "cpu")
32
33 X_train, Y_train, X_test, Y_test = parse_data.loadData('mnist_train
    .csv', 'mnist_test.csv', device)
34
35 def finalTrainAndTest():
36     start = timer()
37     n_epochs = 100
38     net = Net()
39     net.cuda()
40     lrate = 0.01
41     momen = 0.9
42     criterion = nn.CrossEntropyLoss()
43     optimizer = optim.SGD(net.parameters(), lr=lrate, momentum=0.0)
44     #Test for number of epochs we found with above function
45     for i in range(n_epochs):
46         running_loss = 0.0
47
48         batch_s = 1
49         lb = 0
50         ub = batch_s
51
52         if i == 15:
53             lrate = 1e-3
54         elif i == 30:
55             lrate = 1e-4
56         elif i == 45:
57             lrate = 1e-5
58
59         for g in optimizer.param_groups:
60             g['lr'] = lrate
61
62         while ub <= len(X_train):
63             optimizer.zero_grad()
64             output = net(X_train[lb: ub])
65             loss = criterion(output, Y_train[lb: ub])
66             loss.backward()
67             optimizer.step()
68             lb += batch_s
69             ub += batch_s
70             running_loss += loss.item()
71
72         num_correct = 0
73         val_guess = net(X_test)

```

```
74         loss = criterion(val_guess, Y_test)
75
76     for j in range(len(Y_test)):
77         if torch.argmax(val_guess[j]) == Y_test[j]:
78             num_correct += 1
79
80     acc = (num_correct / len(Y_test))
81
82     test_loss = 0.0 + loss.item()
83
84     print("Epoch: " + str(i) + ": " + str(acc))
85     print("Training loss: " + str((batch_s * running_loss) /
86         len(X_train)))
87     print("Test loss: " + str(test_loss))
88     end = timer()
89     print("Training time: " + str(end - start) + " seconds\n")
90 finalTrainAndTest()
```

G.2 parse_data.py

```
1 import torch.utils.data as data_utils
2 import torch
3 import numpy as np
4 import pandas as pd
5
6 def loadData(train_filename, test_filename, device):
7     train = pd.read_csv(train_filename, skiprows=0).values
8     trainX = train[:, 1:].reshape(train.shape[0], 1, 28, 28).astype('
9         float32')
10     X_train = trainX / 255.0
11
12     y_train = train[:, 0]
13
14     print(X_train.shape)
15     print(y_train.shape)
16
17     test = pd.read_csv(test_filename, skiprows=0).values
18     testX = test[:, 1:].reshape(test.shape[0], 1, 28, 28).astype('
19         float32')
20     X_test = testX / 255.0
21
22     y_test = test[:, 0]
23
24     X_train_tsr = torch.from_numpy(X_train)
25     Y_train_tsr = torch.from_numpy(y_train)
26     X_test_tsr = torch.from_numpy(X_test)
27     Y_test_tsr = torch.from_numpy(y_test)
28
29     if device == torch.device("cuda:0"):
30         X_train_tsr = X_train_tsr.cuda()
31         Y_train_tsr = Y_train_tsr.cuda()
32         X_test_tsr = X_test_tsr.cuda()
33         Y_test_tsr = Y_test_tsr.cuda()
34
35     return X_train_tsr, Y_train_tsr, X_test_tsr, Y_test_tsr
```

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.