

1. Abstract

This project implemented a neural network calculation system on an FPGA using FSM structures. For complex networks this is a difficult computational task in itself with a significant amount of ongoing research. However, interesting usage applications can be solved on smaller network topologies. We achieved our core goal by interfacing this structure with a computer via USB-UART, where it was used to compute the propagation values for an artificial neural network.

2. Symbol Key

T: Maximum number of multiplication operations per layer calculation step

M: Bit width of weights

N: Bit width of inputs

I, W: Maximum number of nodes per layer

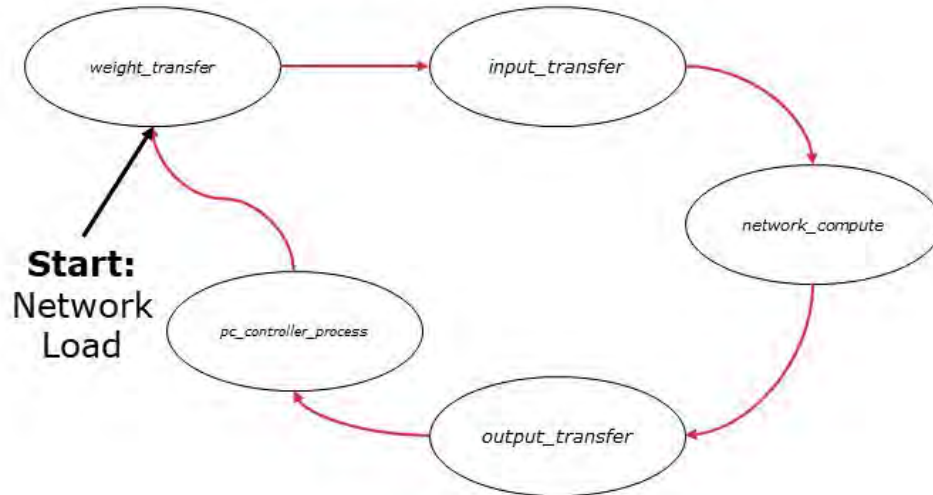
D: Network depth (number of layers)

P: Parallelism. For a given layer, this is the maximum number of nodes whose output can be computed in a single timestep.

our visualization and testing frameworks along with a performance analysis of our final implementation.

3.1. State Machine

Forward Propagation FSM



The diagram above describes the implemented FSM which computes the outputs of a multi-layer network on the FPGA. We begin in the *weight_transfer* communication state. In this state the weights for the entire neural network (in this case a single layer) are sent across the USB-UART interface from the PC to the FPGA Communication Controller module. The FPGA then stores these weights through the Memory Controller module. Along with these weights additional network topology information is sent. This currently includes the layer's parallelism, "P", and maximum node index.

Once all network information is transferred successfully, we enter the *input_transfer* state. In this state the network input layer values are sent to the FPGA via the UART interface in a similar fashion to the weights. These inputs are then stored in the Input Pool module as opposed to the Memory Controller module.

When all values for the input layer have been transferred, we begin the *network_compute* step. This state contains the bulk of our logic and computation. The parallelism "P", discussed earlier, defines the number of nodes that can be computed in parallel for a given layer. In order to perform parallel computation of the forward propagation values for a given set of active nodes, the Layer Controller module reads in the appropriate inputs from the Input Pool and weights from the Memory Controller, and forwards these values to the Layer Computation module. The number of nodes that can be computed in parallel "P" is given by the expression T / I , where T represents the total number of DSP multipliers, and I represents the number of inputs to any given node in the layer. After this computation is complete, the outputs are written to the output pool. When multiple layers are being computed we signal a transfer from the output pool to input pool and computation restarts. When all final layer outputs are in the output pool, the layer controller signals a transition to the *output_transfer* state.

In the *output_transfer* state, the communication module reads the layer's outputs from the output pool, and sends them to the PC as described in the protocol section. Upon reception of all output data the *pc_controller* can begin using the outputs received.

3.2. Communication Protocol (Josh)

3.2.1. Overview

In order to perform inference on the FPGA we must have the ability to transfer data between the Host PC and FPGA. A common protocol is shared between the PC and FPGA in order to encode all the necessary data. Specifically, we have a process running on a computer that feeds network topology information, including weights, and network inputs to the FPGA. This process is referred to as the *pc_controller*. On the FPGA side the *communication_module* receives this information and decodes it for use. The module also handles the transfer of network outputs back to the PC.

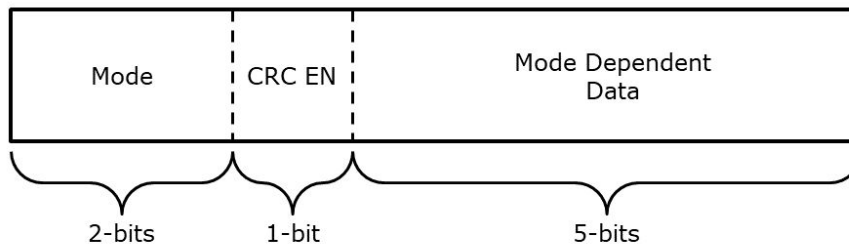
Our serial communication protocol has two actors: a sender and a receiver. For the *weight_transfer* and *input_transfer* states, *pc_controller* is the sender and the *communication_module* is the receiver. The roles are reversed in the *output_transfer* state.

3.2.2. Packet Streams

The protocol itself is defined as follows. All data is sent as a packet stream. A stream is defined to be 1 header packet and 9 data packets where a packet is a single byte. This stream abstraction allows for data to be sent in independent chunks, and any data that would span multiple chunks is sent across multiple packet streams. Note that any leftover space in a packet stream after data insertion is zero-padded. The decision to make the packet stream 10 packets or 80 bits is that it aligns well with our weight bit-width of 24 bits and input/output bit-width of 16 bits. Excluding the header, there are 72 data bits per stream. This means that a single stream can contain 3 weight values with no wasted space or overflow. For inputs there is an overflow of 8 bits per stream ($72 \text{ data bits} \% 16 \text{ bit input/outputs} = 8 \text{ bits}$). This slight overlap does not add much complexity to input encoding as all data can still be handled on a byte-level granularity.

3.2.3. Packet Stream Headers

Packet Stream Header



Header packets contain metadata for the packet stream as a whole. The first 2 bits specify the mode for the following data packets. The supported modes are, *weight*, *input*, *output*, and *debug*. The next bit is whether CRC data for the stream should be expected immediately after the stream ends¹. The final 5 bits may contain additional metadata dependent on the mode of the stream.

Mode 0 is weight transfer mode, which indicates the data packets are weight packets. The mode dependent header data is the layer id of the weights. Note that the communication protocol for weight transfer contains additional requirements so network metadata can be sent as data packets. Specifically, in the first packet stream for each layer (i.e. whenever the layer id portion of the weight header changes), the first 6 data packets must be the parallelism and maximum node id for the layer respectively. Note that 3 data packets per information is 24-bits the same width as a weight. This ensures that no special encode or decode logic must exist to locate this data within a stream, however the *pc_controller* and *communication_module* must know to treat the first two “weights” for each layer specially. In addition the first actual weight for every node is the node’s threshold rather than a weight connecting it to the previous layer.

Mode 1 is input transfer mode and mode 2 is output transfer mode. These types contain the same data, however in input transfer mode data flows from *pc_controller* to *communication_module* and in the opposite direction in output transfer mode. In this mode the dependent data is just zeroed. The input packets flow PC to FPGA and output packets from FPGA to PC.

Mode 3 is debug mode, indicating debug packets which are used to transfer debug data from FPGA to PC. In this mode the dependent data is an error code. Different error codes define different interpretations of the debug data packets.

As described in annotation 1, crc check is currently disabled, however this the protocol still defines its use as follows. If `CRC_EN=1` every packet stream will be followed by a corresponding CRC bit sequence. This sequence will be checked on the receiver by recalculating the CRC of the concatenation [packet_stream, CRC]. If this does not yield 0 an ERROR header packet is sent. Otherwise a SUCCESS packet is sent (All 0 bits). Note that this means communication with `CRC_EN=1` is necessarily two-way as confirm messages must be sent from

¹ Note that the `CRC_EN` bit is always disabled in our current implementation as we did not encounter any errors during data transfer. However, the bit was left as it allows for easy extension in the future if some error correction must be implemented.

the receiver after every packet stream from the sender. Note that an LED on the FPGA is lit if any CRC error is detected or error packet received. If any error packet is received by the sender, it restarts transmission of the current state from the beginning.

3.3. Communication Module (FPGA)

This module will implement the FPGA end of the communication protocol described above. The input will be the TXD port (C4) of the FTDI FT232R USB-UART bridge. The output is the RXD port (D4). This module will operate at 115,200 baud. Note that the FT232R bridge is rated for up to 12MBaud, but we used a lower baud rate to simplify debugging as well as to better demonstrate the performance implications of a high communication overhead. These implications will be discussed further in section 7.

The implementation is similar to Lab 2B and Lab 5C. For this reason the uart interface was abstracted away in the modules *uart_rx* and *uart_tx* for receiving and transmission respectively. *uart_rx* was Sebastian's Lab5C module modified to buffer an entire packet stream before sending it for decode. *uart_tx* was based on lab2b, but it takes in a packet stream and transmits at the necessary baudrate. Note that the 16x supersample clock generator was also abstracted into the *uart_clkdiv* module. The *communication_module* itself handles the decoding and encoding of packet streams for the previously described uart interfaces.

For decoding weight streams, every 3 data packets are combined to form a single 24-bit weight. The current layer_id is extracted from the header and also output from this module. The module must keep some state regarding the last layer received in order to decode the network metadata described in the communication protocol section. If it detects that it has received the first packet stream for a new layer, it will treat the first two weights as the layer's parallelism and maximum node id respectively. These values are stored into output registers and sent to the *parallelism_pool* and *max_node_id_pool*. Note that the fact the first weight per node is a threshold is completely ignored by the communication module to cut down on complexity. It simply assumes every node has $I+1$ weights. Once all weights for a single node are received the communication module asserts a signal that the weight buffer is full. This signal is handled by the memory controller as described in section 3.4. This process continues as long as weight packets are received.

Input decoding is much simpler than weight decoding. Even though a single input may span multiple packet streams, all streams are buffered contiguously in the order received. This places the split input bytes consecutively. The communication module expects I inputs, so once it receives this many it asserts to the input pool that it should read the input buffer.

Note that sending an invalid number of weight or input bits such that it only partially fills the communication module weight or input buffer is disallowed by the communication protocol.

Once a debug packet with all 0 data packets is received, the communication module asserts that transfer in is completed and the layer controller begins working.

This module will also handle the sending of output data from the FPGA to the computer once the *layer_controller* asserts that the last layer's outputs are in the output pool. These outputs will be read from the output pool, translated to packet streams, and each stream is then sent through *uart_tx*.

The main lesson learned through this project is in regards to this module. As I was testing communication it would have been very useful to create a test that combined only this module with the memory controller and input pool, and then used this test actually integrated with the

computer. I wrote a test similar to this closer to the end of the project that would simply send back to the computer the weights or inputs it had received. This tests not only the communication and decoding aspect, but also it's interaction with the elements it needs to store persistent data to. This exposed a couple of bugs that would have been nice to find much earlier in the process. My main takeaway from this experience is to begin testing modules closer to the the environment they will be used in as soon as possible rather than relying on tests in isolation to prove correctness of the system as a whole.

3.4. Memory Controller (Josh)

As we have defined our maximum network width and depth we can derive how much memory is needed for the storage of weights. As described in the symbol key, we have a maximum network width of 24 and depth of 4. Every node has $24 \text{ weights} + 1 \text{ threshold}$ worth of bits to store through the memory controller. Multiplying all of this by a weight width of 24-bits yields:

$$24\text{weightwidth} * (24\text{netwidth} * 4\text{netdepth} * (24\text{netwidth} + 1\text{threshold})) = 57,600 \text{ bits}$$

The Nexys Video has 13 MBits of fast BRAM, so this is well within the limitations of the board. As described in section 3.2, the communication module sends the weights and thresholds a single node at a time ($I \text{ weights} + 1 \text{ threshold}$). This yields a write port width of:

$$24\text{weightwidth} * (24\text{netwidth} + 1) = 600 \text{ bits}$$

We synthesized a BRAM IP core with this port width, and depth defined as the maximum number of nodes in a network. In regards to the memory controller implementation itself it has two distinct operating modes: writing and reading.

For writing the memory controller keeps an internal counter tracking the number of nodes that have been written so far. It then writes to this index in BRAM. The resulting memory layout is such that any reads will be from consecutive indices, as all nodes within a layer will have been stored one after another.

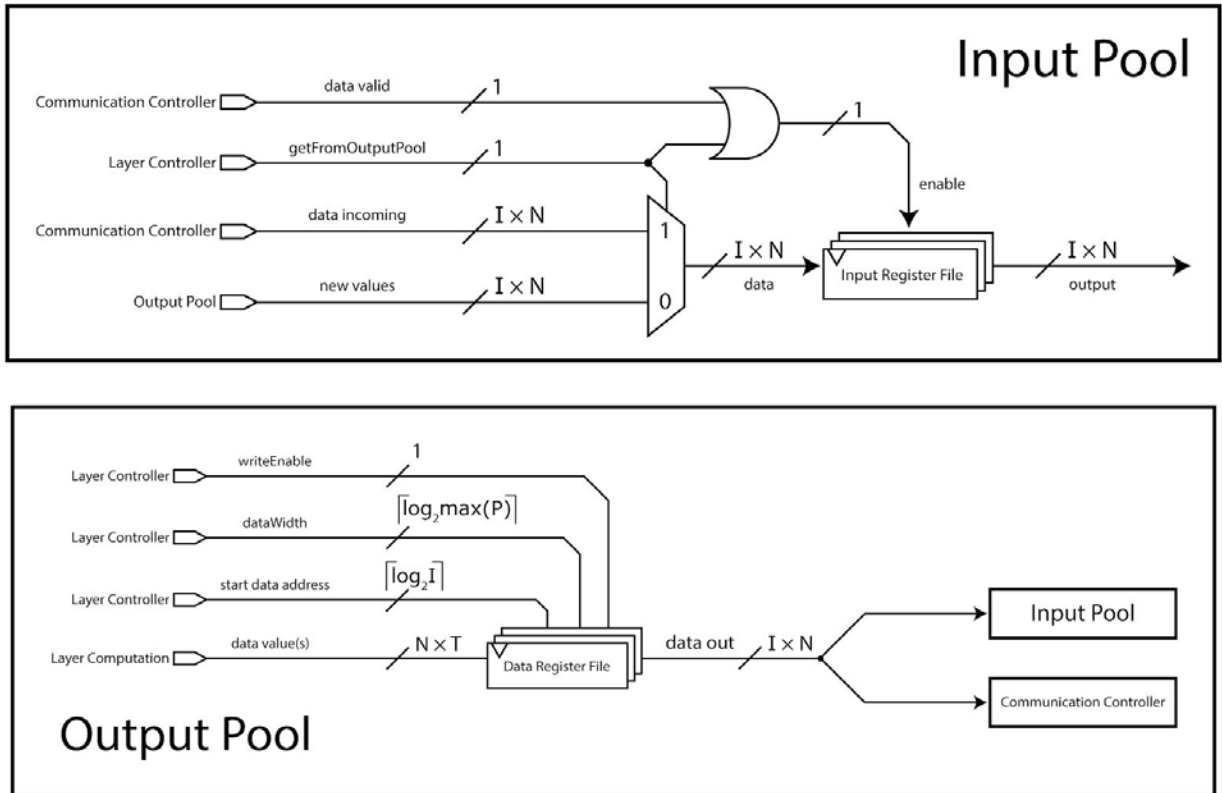
For reading the layer computation module expects to receive weights and thresholds for P (parallelism) nodes in the current layer. For this reason two output shift registers are created *read_weights* and *read_thresholds*. In order to read all necessary values from memory, a single layer read request is broken up into a BRAM read request for P nodes in the layer. The initial read address is calculated from the signals *read_layer_id* and *read_node_id*. Since data is written layer-by-layer we can reconstruct the start index for a given layer and node as:

$$\text{read_layer_id} * \text{MAX_NET_WIDTH} + \text{read_node_id}$$

Once a single node read response is returned from BRAM, the first weight is shifted into the output threshold shift register and the remaining weights shifted into the output weight shift register. P defines the parallelism of the current layer, but if $\text{MAX_NET_WIDTH} \% P \neq 0$ then it is possible that less than P reads must be done for the final nodes in a layer. For this reason the layer controller checks the current node id being read against the maximum node id for the layer. Note that this id verification functionality is not needed in the current implementation given that

we have enough DSP slices to support calculating an entire layer per computation step. Thus, P will never cause a read past max_node_id as $P = LayerWidth$. Once all consecutive BRAM reads are completed the output shift registers are full and a $read_rdy$ signal is asserted.

3.5. Pool Module (Josh)



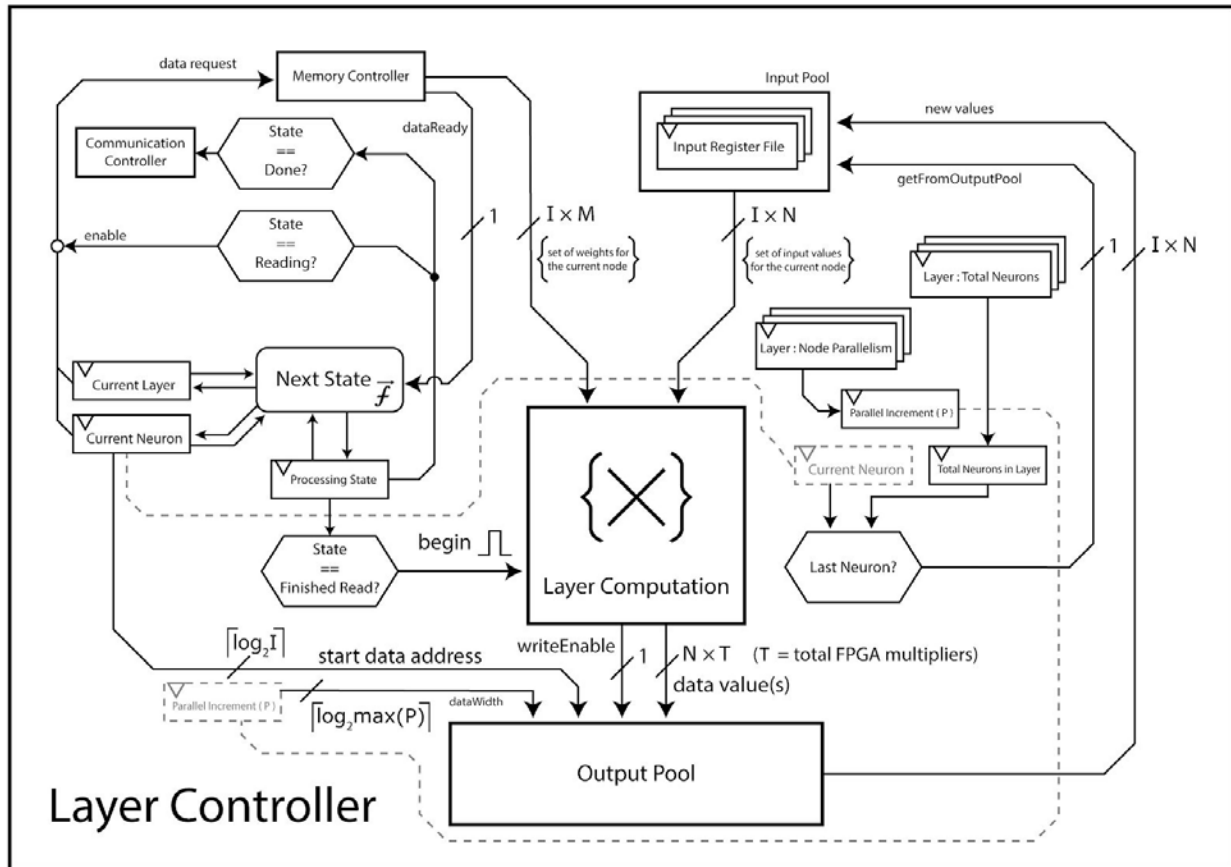
The Input Pool module encapsulates the input data buffer that is used to calculate values in the Layer Computation module. The write inputs to this register file come from the Output Pool and Communication Controller modules, given that the buffer can store initial values received from the host computer or hidden layer values stored in the Output Pool buffer after computation is complete. In the latter case, these output values are passed to the input pool since they are required to calculate the next layer of neuron values. The register file data write is enabled when either the data valid signal from the communication controller or the getFromOutputPool signal from the Layer Controller is asserted.

In addition to the input pool, two other pools are written to by the communication module, but for the purpose of storing network metadata. Specifically, a parallelism pool and max node id pool were created. As described in the communication section, these are values that occur once per received layer. The index to which the data should be written is defined by the current $layer_id$ output of the communication module. The communication module asserts the $write_en$ signal to these pools once per node received. As the $layer_id$ is the same for all nodes

in the same layer, multiple writes of the same data occur for every layer, but this has no negative effects. To read from these pools the *index* input is used to extract a value from the internal buffers and store it in an output register. Reads are initiated by the layer controller in order to read from *memory controller* as this module needs to know the parallelism and maximum node id of the layer being read.

The Output Pool serves as the buffer for Layer Computation outputs. As opposed to the Input Pool, which modifies the entire register file on a write operation, the output pool uses the data start address and data width values from the Layer Controller to determine which registers to write to when writeEnable is asserted. This occurs because the output pool was designed to allow for partial computation of layers where some outputs will be sent after one computation step, and the rest of the outputs will be sent across some number of computation timesteps. Note, however, that although the output pool supports this functionality, the rest of the system currently always calculates a full layer per computation step. This change occurred as with the increased number of DSP slices on the Nexys Video, we had enough multipliers to do so. The data from this register file is forwarded to the Input Pool upon layer computation completion and the Communication Controller once all layers finish.

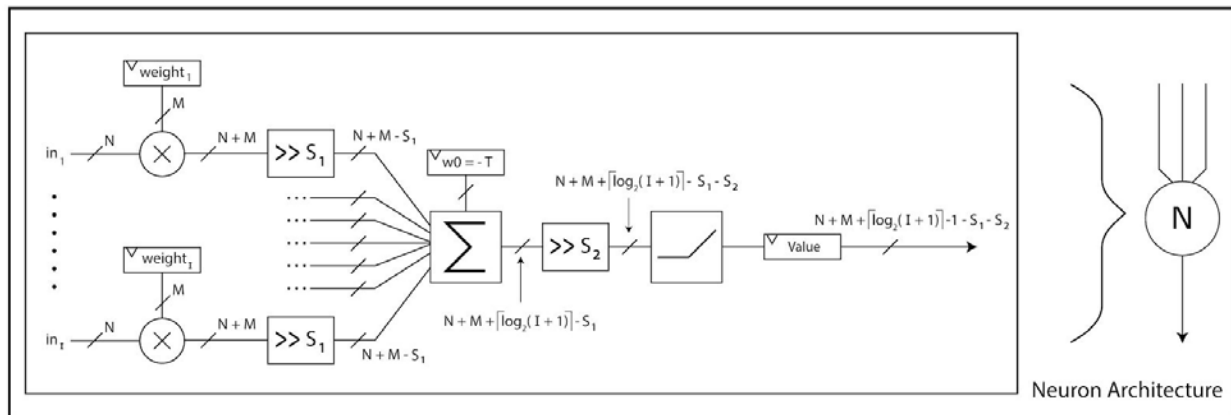
3.6. Layer Controller (Sebastian)

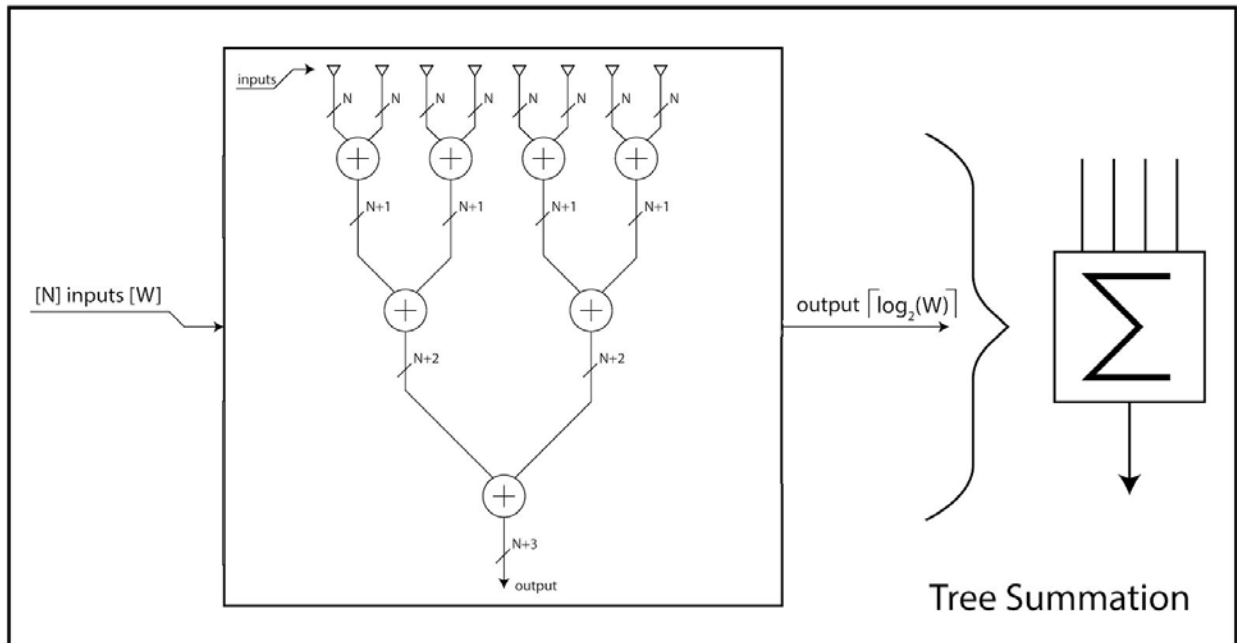
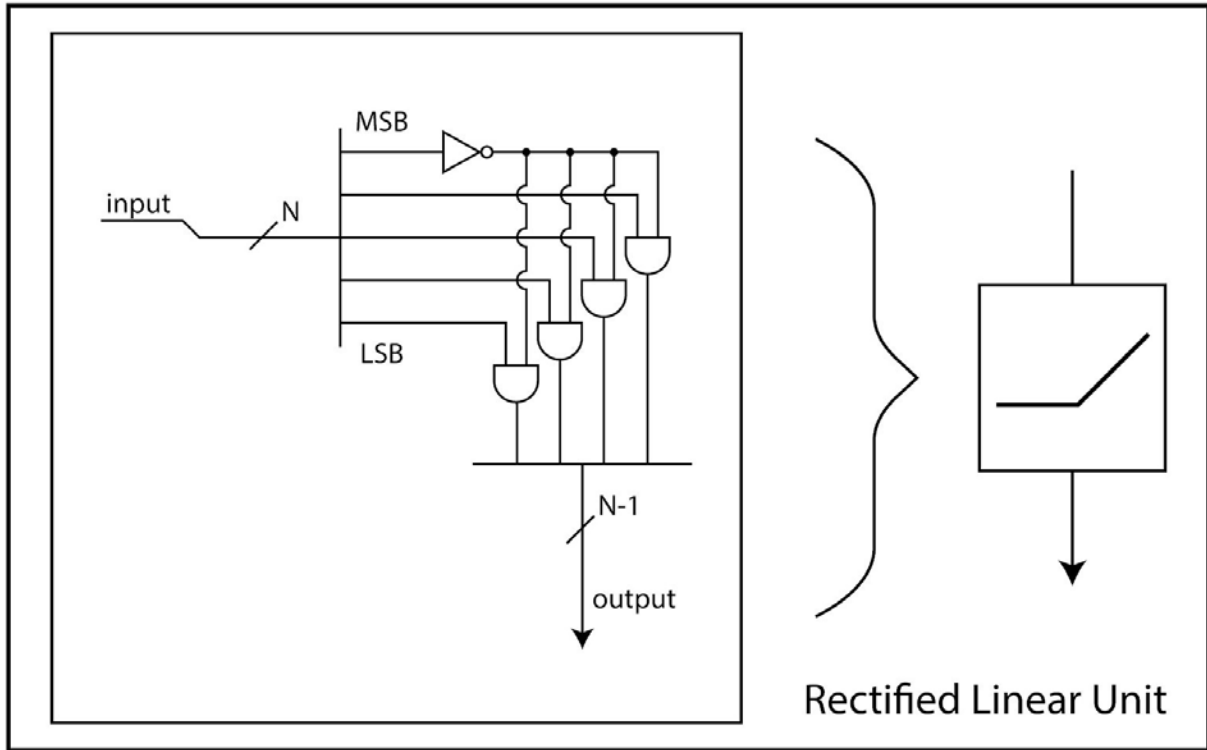


The Layer Controller module acts as the primary control FSM for calculating the neural network layers. To begin, the network depth is loaded from the Communication Controller via the *commCTL_networkDepth* bus and corresponding ready signal. During this data latching the requests for layer architecture data and neuron weight data for the primary layer are initiated via *memCTL_getLayerDataRequest* and *memCTL_getNeuronDataRequest* signals. These values are latched once the requests have valid data on their respective buses and active ready signals. The Layer Controller then waits until the Input Pool has neuron input data ready for calculation, which forms the *S_WAITING_INPUTS* state of the FSM. Once these conditions are met, the Layer Controller sends a pulse to the Layer Computation module (*layerComputation_startPulse*) to initiate the calculations for the layer. The Layer Controller FSM waits for the layer computation to finish (state *S_WAITING_LAYERCOMP*), and if the current layer being considered is the last layer, a signal is asserted for the communication controller to fetch data outputs from the Output Pool (*commCTL_isStateDone*). If the layer being processed is an intermediate ("hidden") layer, the layer counter is incremented, and requests for the next layer architecture and neuron weights are sent to the memory controller. The FSM then returns to the input waiting state as previously delineated.

As inputs, the Layer Controller receives the network depth, parallelism values, maximum neuron index for the given layer, as well as flags from other modules indicating the ready status of weights, inputs, and valid data outputs. The Layer Controller outputs the data ready signal to the Input Pool so that values are cycled from the Output Pool, parameters for the Memory Controller to fetch the appropriate data values, request signals to assert data fetch, as well as flags to the Communication Controller and Layer Computation modules regarding data-flow validity and timing, respectively.

3.7. Layer Computation (Sebastian)





The Layer Computation module serves as the central data processor for the project. The inputs of this module include the incoming weights, thresholds, and inputs for the neurons in the current layer of computation, as well as a start pulse from the Layer Controller. To begin, the set of weights and inputs are multiplied in pairs, effectively calculating the dot product between the two value vectors. These values are then right shifted by a quantity “S1”, in order to reduce the

hardware needed for summation as well as reduce the probability of an overflow result. The results of this shift are stored in the *multiplicationResult* bit array. The values are then added together using an array of one tree adder per neuron *TreeAdder24(0-23)*, and then right shifted by another quantity “S2” to reduce the bit width of the output values to be the same as the input values, given the cyclic implementation of calculating layers in this hardware implementation. Once all values are computed, they are sent to the Output Pool for further data routing along with an enable signal to indicate that valid data is being sent across the *outputPool_finishedVals* bus. Once the values are latched by the Output Pool, the Layer Computation module waits for the next set of weights and inputs to be sent, and initiates the next round of calculation once the *startPulse* signal from Layer Controller is asserted. The module uses a FSM to keep track of arithmetic steps, consisting of states *S_WAITING*, *S_SUM*, and *S_DONE*.

4. Neural Net Implementation / Visualization (Sebastian)

The visualization software was critical to showcasing the functionality of the system. This code interfaced with the Lab 6 code provided by staff in 6.034 to use our FPGA solution. An encoding function was created to translate the Lab API to a format consisting of list of lists for layer inputs, weights, and thresholds ("neuronal biases"), that could then be sent over USB UART communication. The biggest challenge in modifying the lab code was changing the activation function from a sigmoid curve to a rectified linear unit, in order to match the FPGA neural net implementation. The data sets used would quickly cause the well-known "dying ReLU" problem to occur, at which point no further training could be performed. Forcing data inputs and weights to be integers to match the FPGA implementation caused this problem to worsen. Solving this issue took a large portion of time, proving to be non-trivial. The solution used was implementing leaky ReLUs as opposed to rigid ReLUs, and then iteratively reducing the leakage factor to determine hyperparameters such as training rate that would allow for the net to be successfully trained with regular ReLUs. In addition, rather than training on the FPGA using only integers, the training was instead performed on the host computer using floating point arithmetic. Once the training was complete the network weights and thresholds were converted into integers. The FPGA would then perform inference on the training space, and the data received was displayed on the screen in a two-dimensional chart.

5. PC_Controller (Josh)

The implementation of the PC side of the communication protocol was done near the end of the project similar to the visualization code. The code for it can be seen in the *app.py* and *uart.py* files. Since the FPGA side was finished at this point the implementation in python was not very difficult. The underlying UART implementation is based off of the code given in lab5c with slight updates to support user-specified COM port selection through the command line. The communication protocol was implemented on top of this. In summary, weights and inputs were stored as multidimensional python int arrays, sent from the neural network implementation described above. Two modifications were required on these arrays before transfer. First, all non-layer arrays were zero-padded to length “I” in order to provide a constant length shared by the FPGA and PC. Once this was done, the multidimensional arrays would be flattened. At this point the “send_data_stream” method would be called specifying the data type being sent and the corresponding data array. Python would then iterate over this array creating and sending packet streams. Doing so involved converting the data array to bytes where the number of bytes per

entry is stream-type dependent as specified in the communication protocol (i.e. weights are 3 bytes and inputs 2-bytes). Once the entire array is in bytes form the code takes every 72 data bytes, prepends a header byte, and transfers it using pyserial. If there are less than 72 data bytes the stream is zero-padded.

6. Top Level Module Implementation & Debugging (Josh)

Implementation of the top-level module *_top* involved modifying the constraints file with all I/O ports necessary for implementation and debugging as well as connecting the modules described above. As we had a well-defined block diagram, connecting the modules took very little time.

The main hurdles came while debugging all of the modules integrated together at the top level, taking more time than implementing the modules themselves. Though all modules worked in isolation, once connected there were some mismatches in control signals and data layout as some modules were not changed to reflect design updates since the proposal. Along with this, connecting the modules would sometimes expose internal bugs within modules that had not been tested for prior. This full system testing was done partially in simulation using the *full_test.v* testbench. This testbench supported sending packet streams across the *uart_rx* interface and then having the final network outputs flow all the way through a dummy *uart_rx* interface to mimic the PC receiving the outputs. This test would check all aspects of the design, which was very useful. However, when debugging the entire system rather than single modules it was much harder to track down where incorrect outputs were coming from.

app.py contained testing code in addition to the code necessary to support visualization. This eased the process of constructing meaningful test data as *app.py* provides python constructs, while in the simulation testbench, *full_test*, packet streams had to be manually constructed.

while computation takes about 10 clock cycles with a 1ns clock. Increasing the baudrate would help relieve this issue.

Appendix A: Verilog Implementation

(Note: lines prefaced with “...” are continuations of code in the previous line, and do not denote line breaks in the source code file.)

Top Level Module

```

1 `timescale 1ns / 1ps
2 `include "master_params.vh"
3
4 module _top (
5     input CLK_100MHZ,
6     input UART_TX_IN,
7     input CPU_RESETN,
8     output UART_RX_OUT,
9     output [6:0] LED
10    );
11
12
13    // STATE
14    localparam STATE_INIT = 0; // PC is sending weights and
... inputs
15    localparam STATE_CALC = 1; // FPGA is calculating network
... outputs
16    localparam STATE_OUT = 2;
17    reg state;
18    reg network_inputs_rec;
19    reg network_weights_rec;
20
21
22
23    ////////// Communication module wires
24    wire [`PARAM_E_EM-1:0] debug_outputs;
25    wire [`PARAM_I_N-1:0] comm_inputs;
26    wire [`PARAM_I_M1-1:0] comm_layer_weights;
27    wire comm_input_rdy, comm_weight_rdy;
28    wire comm_transfer_in_done, comm_transfer_out_done;
29    wire network_outputs_done, debug_outputs_done;
30    wire [`BTW_MAX_NET_DEPTH-1:0] comm_layer_id;
31    wire [`BTW_MAX_PARALLELISM-1:0] comm_layer_parallelism;
32    wire [`BTW_MAX_NET_WIDTH-1:0] comm_layer_max_node_id;
33    wire [`BTW_MAX_NET_DEPTH-1:0] comm_network_depth;
34    wire [2:0] comm_debug_state;
35
36    ////////// Memory Controller Wires
37    wire mem_rd_layer_en, mem_rd_neuron_en;
38    wire mem_rd_weights_rdy;
39    wire [`PARAM_I_I_M-1:0] mem_rd_weights;

```



```

40     wire [`PARAM_I_M-1:0] mem_rd_thresholds;
41
42     /////// Output Pool Wires
43     wire cntrl_rd_new_layer; // Pulse to update
... cntrl_parallelism and cntrl_max_node_id
44     wire [`PARAM_I_N-1:0] output_pool_data;
45
46     /////// Parallel Pool Wires
47     wire [`BTW_MAX_PARALLELISM-1:0] parallelism_rd;
48     wire parallelism_pool_rdy;
49
50     /////// Max Node Pool Wires
51     wire [`BTW_MAX_NET_WIDTH-1:0] max_node_id_rd;
52     wire max_node_id_pool_rdy;
53
54     /////// Input Pool Wires
55     wire [`PARAM_I_N-1:0] input_pool_rd;
56     wire input_pool_rd_rdy;
57
58
59     ///// Layer Computation wires
60     wire computation_done, computation_start;
61     wire [`PARAM_I_N-1:0] computation_output;
62
63     /////// Layer Controller wires
64     wire cntrl_done_all;
65     wire [`BTW_MAX_NET_DEPTH-1:0] cntrl_cur_layer_id;
66     wire [`BTW_MAX_NET_WIDTH-1:0] cntrl_cur_node_id;
67     wire [`BTW_MAX_PARALLELISM-1:0] cntrl_parallelism;
68     wire [`BTW_MAX_NET_WIDTH-1:0] cntrl_max_node_id;
69
70     // System wires
71     wire sysclk = CLK_100MHZ;
72     wire btn_cpu_resetrn;
73     debounce debounce_reset(.clock(sysclk), .reset(0),
... .noisy(~CPU_RESETRN), .clean(btn_cpu_resetrn));
74     wire reset = btn_cpu_resetrn | comm_transfer_out_done;
75
76     // DEBUG wires
77     assign LED[0] = comm_debug_state[0];
78     assign LED[1] = comm_debug_state[1];

```

```

79     assign LED[2] = comm_transfer_in_done;
80     assign LED[3] = cntrl_done_all;
81     assign LED[4] = comm_transfer_out_done;
82     assign LED[5] = output_pool_data[7:0] != 0;
83     assign LED[6] = input_pool_rd[7:0] == 1;
84
85
86
87     memory_controller memory_controller_instance(.clk(sysclk),
... .reset(reset),
88         .write_weights_rdy(comm_weight_rdy),
... .rd_layer_en(mem_rd_layer_en), .rd_node_en(mem_rd_neuron_en),
89         .write_layer_weights(comm_layer_weights),
... .write_layer_id(comm_layer_id),
90         .read_layer_id(cntrl_cur_layer_id),
... .read_node_id(cntrl_cur_node_id),
... .read_parallelism(cntrl_parallelism),
91         .read_weights(mem_rd_weights),
... .read_thresholds(mem_rd_thresholds),
... .read_weights_rdy(mem_rd_weights_rdy));
92
93     communication_module comm_module_instance(.clk(sysclk),
... .reset(reset), .uart_tx_in(UART_TX_IN),
... .uart_rx_out(UART_RX_OUT),
94         .network_outputs(output_pool_data),
... .debug_outputs(debug_outputs),
95         .debug_outputs_done(debug_outputs_done),
96         .network_inputs(comm_inputs),
... .input_rdy(comm_input_rdy),
97         .layer_weights(comm_layer_weights),
... .weight_rdy(comm_weight_rdy), .layer_id(comm_layer_id),
98         .layer_parallelism(comm_layer_parallelism),
... .layer_max_node_id(comm_layer_max_node_id),
99         .network_depth(comm_network_depth),
... .in_transfer_done(comm_transfer_in_done),
100        .start_output_transfer(cntrl_done_all),
... .out_transfer_done(comm_transfer_out_done),
101        .debug_state(comm_debug_state));
102
103
104     LayerController layer_controller_instance(.clk(sysclk),

```

```

104.. .reset(reset),
105         .commCTL_networkDepth(comm_network_depth),
... .commCTL_networkDepthReady(comm_transfer_in_done),
106         .parallelPool_parallelism(parallelism_rd),
... .parallelPool_parallelismReady(parallelism_pool_rdy),
107
... .neuronIndicesPool_maxNeuronIndex(max_node_id_rd),
... .neuronIndicesPool_maxNeuronIndexReady(max_node_id_pool_rdy),
108
... .memCTL_weightsAreReady(mem_rd_weights_rdy),
... .inputPool_inputsAreReady(1),
... .layerComputation_dataReady(computation_done),
109         .outputPool_dataWidth(cntrl_parallelism),
... .outputPool_maxNeuronIndex(cntrl_max_node_id),
110
... .inputPool_getFromOutputPool(cntrl_rd_new_layer),
111         .memCTL_currentLayer(cntrl_cur_layer_id),
... .memCTL_currentNeuron(cntrl_cur_node_id),
... .memCTL_parallelism(cntrl_parallelism),
112
... .memCTL_getLayerDataRequest(mem_rd_layer_en),
... .memCTL_getNeuronDataRequest(mem_rd_neuron_en),
113         .commCTL_isStateDone(cntrl_done_all),
... .layerComputation_startPulse(computation_start)
114     );
115     LayerComputation layer_computation_instance (.clk(sysclk),
... .reset(reset), .startPulse(computation_start),
116         .incomingWeights(mem_rd_weights),
... //topmost weight is "w0 = - T" Seperated out by memory
... controller and sent in threshold port //this is currently for
... just one node, if we want to do parallelism this needs to be
... P{ }'d (consecutive weights must be sent for each node)
117         .incomingThresholds(mem_rd_thresholds),
118         .incomingInputs(input_pool_rd),
119
... .outputPool_finishedVals(computation_output),
120         .outputPool_writeEnable(computation_done)
121     );
122
123     /////// Input pool wires
124     input_pool input_pool_instance(.clk(sysclk),

```

```

124.. .reset(reset),
125         .comm_write_en(comm_input_rdy),
... .cntrl_pool_transfer(cntrl_rd_new_layer),
126         .comm_data(comm_inputs),
... .output_pool_data(output_pool_data),
127         .read_data(input_pool_rd),
... .read_data_rdy(input_pool_rd_rdy));
128
129     //// Parallelism pool
130     pool#(.DATA_WIDTH(`BTW_MAX_PARALLELISM))
... parallelism_pool_instance(.clk(sysclk), .reset(reset),
131         .write_en(comm_weight_rdy),
... .write_data(comm_layer_parallelism),
132         .read_en(mem_rd_layer_en),
... .index(cntrl_cur_layer_id), .read_data(parallelism_rd),
133         .data_rdy(parallelism_pool_rdy));
134
135     //// Max Node Index pool
136     pool#(.DATA_WIDTH(`BTW_MAX_NET_WIDTH))
... max_id_pool_instance(.clk(sysclk), .reset(reset),
137         .write_en(comm_weight_rdy),
... .write_data(comm_layer_max_node_id),
138         .read_en(mem_rd_layer_en),
... .index(cntrl_cur_layer_id), .read_data(max_node_id_rd),
139         .data_rdy(max_node_id_pool_rdy));
140
141
142     //// Output Pool
143     output_pool output_pool_instance(.clk(sysclk),
... .reset(reset),
144         .write_en(computation_done),
... .parallelism(cntrl_parallelism),
145         .max_node_id(cntrl_max_node_id),
146         .write_data(computation_output),
... .cur_outputs(output_pool_data));
147
148 endmodule
149

```

Communication Module


```

1  `include "master_params.vh"
2
3  module communication_module(
4      input clk, //CLK_100MHZ
5      input reset,
6      input uart_tx_in, //UART_TX_IN,
7      input start_output_transfer,
8      input [`PARAM_I_N-1:0] network_outputs,
9      input [`PARAM_E_EM-1:0] debug_outputs,
10     input debug_outputs_done,
11     output uart_clk, // Output subdivided clock signal for
... testing
12     output uart_rx_out, //UART_RX_OUT,
13     output reg input_rdy,
14     output reg weight_rdy,
15     output reg in_transfer_done,
16     output reg out_transfer_done,
17     output reg [`PARAM_I_N-1:0] network_inputs,
18     output reg [`PARAM_I_M1-1:0] layer_weights, // We output a
... single node's weights +cutoff per write
19     output reg [`BTW_MAX_NET_DEPTH-1:0] layer_id,
20     output reg [`BTW_MAX_PARALLELISM-1:0] layer_parallelism,
21     output reg [`BTW_MAX_NET_WIDTH-1:0] layer_max_node_id,
22     output reg [`BTW_MAX_NET_DEPTH-1:0] network_depth,
23     output [2:0] debug_state
24 );
25
26     localparam STATE_HEADER_WAIT = 0;
27     localparam STATE_CRC_WAIT = 1;
28     localparam STATE_IN_DONE = 2;
29     localparam STATE_OUTPUT_TRANSFER = 3;
30     localparam STATE_DEBUG_TRANSFER = 4;
31
32     //// Note: Modes from packet header. Must match comm.
... protocol
33     localparam MODE_WEIGHT = 2'd0;
34     localparam MODE_INPUT = 2'd1;
35     localparam MODE_OUTPUT = 2'd2;
36     localparam MODE_DEBUG = 2'd3;
37
38     parameter VERIFY_STOP = 1; // Whether UART_RX module

```

```

38... should verify a stop bit before accepting a packet
39
40
41     //// UART RX Setup
42     wire clk_serial_sample;
43     wire verify_stop = VERIFY_STOP;
44     // Contains a valid, complete packet stream
45     wire [`BITS_PER_STREAM-1:0] packet_stream;
46     wire packet_stream_rdy;
47     //TDO synch reset signal with sample slk
48     uart_clkdiv uart_clkdiv_instance(.clk(clk), .reset(reset),
... .clk_out(clk_serial_sample));
49     uart_rx uart_rx_instance(.clk(clk),
... .uart_clk(clk_serial_sample), .reset(reset),
... .signal(uart_tx_in), .verifyStopBit(1),
... .data_out(packet_stream), .ready(packet_stream_rdy));
50     assign uart_clk = clk_serial_sample;
51     ///// UART TX Setup
52
53     wire transmit_done;
54     reg transmit_en;
55     reg start_output_transfer_pulse;
56     reg last_transmit_en;
57     wire tx_reset = reset | (last_transmit_en == 0 &&
... transmit_en == 1); //(start_output_transfer_pulse;
58     reg [`TOTAL_BITS_PER_STREAM-1:0] transmit_stream;
59     uart_tx uart_tx_instance(.clk(clk), .reset(tx_reset),
... .en(transmit_en), .start_data(transmit_stream),
... .xmit_data(uart_rx_out), .done(transmit_done));
60     // Read in outputs stream by stream
61     // Read in debug stream by stream
62
63     //// Internal Setup
64     reg [2:0] state;
65     reg [`CLOG2(`STREAMS_PER_OUTPUTS+1)-1:0]
... output_streams_left;
66     reg [`PARAM_I_N-1:0] cur_network_outputs;
67     reg [`CLOG2(`PARAM_I_M1):0] weight_bits_left;
68     reg [`CLOG2(`PARAM_I_N):0] input_bits_left;
69
70     wire [1:0] header_mode = packet_stream[1:0];

```

```

71     wire          header_crc = packet_stream[2];
72     wire [4:0] header_data = packet_stream[7:3];
73     assign debug_state = {layer_parallelism[0],
... layer_max_node_id[0]};
74
75     reg prev_layer_id;
76     wire layer_first_weight_stream = (header_data !=
... prev_layer_id);
77
78     integer i;
79     function [`TOTAL_BITS_PER_STREAM-1:0] data_to_stream
... (input [`BITS_PER_STREAM-1:0] data);
80         for(i = 0; i < `PACKETS_PER_STREAM; i = i + 1) begin
81             data_to_stream[8*(i+1)+2*i -: 8] = data[8*(i+1)-1
... -: 8];
82             data_to_stream[8*(i+1)+2*i+1] = 1;
83             data_to_stream[8*i+2*i] = 0;
84         end
85     endfunction
86
87     always @(posedge clk) begin
88         last_transmit_en <= transmit_en;
89         if (reset) begin
90             state <= STATE_HEADER_WAIT;
91             weight_bits_left <= `PARAM_I_M1;
92             input_bits_left <= `PARAM_I_N;
93             input_rdy <= 0;
94             weight_rdy <= 0;
95             in_transfer_done <= 0;
96             prev_layer_id = -1;
97             transmit_en <= 0;
98             last_transmit_en <= 0;
99             network_inputs <= 0;
100            layer_weights <= 0;
101            network_depth <= 0;
102            output_streams_left <= 0;
103            out_transfer_done <= 0;
104            start_output_transfer_pulse <= 0;
105            transmit_stream <= -1;
106            layer_parallelism <= 0;
107            layer_max_node_id <= 0;

```



```

108         end else if (start_output_transfer_pulse ||
... output_streams_left > 0) begin
109             if (transmit_en == 0 &&
... start_output_transfer_pulse == 1) begin
110                 start_output_transfer_pulse <= 0;
111                 // Starting new stream transfer
112                 transmit_stream <=
... data_to_stream({cur_network_outputs[`BITS_PER_STREAM-9:0],
... 6'd0, MODE_OUTPUT});
113                 cur_network_outputs <= cur_network_outputs >>
... (`BITS_PER_STREAM-8);
114                 transmit_en <= 1;
115             end else if (transmit_en == 1 && transmit_done ==
... 1) begin
116                 // transmission done
117                 transmit_en <= 0;
118                 output_streams_left <= output_streams_left -
... 1;
119                 if (output_streams_left == 1) begin
120                     // Last packet just finished
121                     out_transfer_done <= 1;
122                 end else begin
123                     // Intermediate packet finished, so reset
... uart_tx for the next one
124                     start_output_transfer_pulse <= 1;
125                 end
126             end
127         end else if (start_output_transfer &&
... out_transfer_done == 0) begin
128             // Start transfer process by reading in outputs
129             //start_output_transfer_pulse <= 1;
130             cur_network_outputs <= network_outputs >>
... (`BITS_PER_STREAM-8);
131             output_streams_left <= `STREAMS_PER_OUTPUTS;
132             transmit_stream <=
... data_to_stream({network_outputs[`BITS_PER_STREAM-9:0], 6'd0,
... MODE_OUTPUT});
133             transmit_en <= 1;
134         end else if (packet_stream_rdy && transmit_en == 0)
... begin
135             case (state)

```



```

161... current layer
162         if (weight_bits_left <
...  `BITS_PER_STREAM - 8) begin
163             layer_weights <=
... {packet_stream[8+`STREAM_WEIGHT_OVERLAP-1:8],
... layer_weights[`PARAM_I_M1-1:`STREAM_WEIGHT_OVERLAP]};
164             //weight_bits_left <=
... weight_bits_left - `STREAM_WEIGHT_OVERLAP;
165         end else begin
166             layer_weights <=
... {packet_stream[`BITS_PER_STREAM-1:8],
... layer_weights[`PARAM_I_M1-1:`BITS_PER_STREAM-8]};
167             //weight_bits_left <=
... weight_bits_left - (`BITS_PER_STREAM - 8);
168         end
169         weight_rdy <= 1;
170         weight_bits_left <= `PARAM_I_M1;
171     end
172     layer_id <= header_data;
173 end
174 MODE_INPUT: begin
175     if (input_bits_left >= `BITS_PER_STREAM -
... 8) begin
176         // We are transferring an entire
... packet and must wait for at least 1 more stream to read the
... whole layer
177         network_inputs <=
... {packet_stream[`BITS_PER_STREAM-1:8],
... network_inputs[`PARAM_I_N-1:`BITS_PER_STREAM-8]};
178         input_bits_left <= input_bits_left -
... (`BITS_PER_STREAM - 8);
179     end else begin
180         // This is the last stream in the
... current layer
181         if (input_bits_left < `BITS_PER_STREAM
... - 8) begin
182             network_inputs <=
... {packet_stream[8+`STREAM_INPUT_OVERLAP-1:8],
... network_inputs[`PARAM_I_N-1:`STREAM_INPUT_OVERLAP]};
183         end else begin
184             network_inputs <=

```

```

184... {packet_stream[`BITS_PER_STREAM-1:8],
... network_inputs[`PARAM_I_N-1:`BITS_PER_STREAM-8]};
185         end
186         input_rdy <= 1;
187         input_bits_left <= `PARAM_I_N;
188     end
189     end
190     MODE_DEBUG: begin
191         network_depth <= layer_id;
192         in_transfer_done <= 1;
193     end
194     default: ; // TODO: This is some kind of
... transmission error as we only have 3 valid input modes
195     endcase
196     end
197     STATE_CRC_WAIT: begin
198         // Begin parallel CRC check by passing crc
... bits to checker and latching the response
199         state <= STATE_HEADER_WAIT;
200     end
201     // TODO: Output states
202
203     default: ;
204     endcase
205     end else begin
206         if (input_rdy) input_rdy <= 0;
207         if (weight_rdy) weight_rdy <= 0;
208     end
209     end
210
211 endmodule
212

```

```

1  `ifndef CONSTANT_FUNCS_H
2  `define CONSTANTS_FUNCS_H
3
4  // Constant clog2 hack as verilog has a $clog2 bug that breaks
... its use in localparameter and initial blocks
5  `define CLOG2(x) \
6      ((x <= 2) ? 1 : \
7       (x <= 4) ? 2 : \
8       (x <= 8) ? 3 : \
9       (x <= 16) ? 4 : \
10      (x <= 32) ? 5 : \
11      (x <= 64) ? 6 : \
12      (x <= 128) ? 7 : \
13      (x <= 256) ? 8 : \
14      (x <= 512) ? 9 : \
15      (x <= 1024) ? 10 : \
16      (x <= 2048) ? 11 : \
17      (x <= 4096) ? 12 : \
18      (x <= 8192) ? 13 : \
19      (x <= 16384) ? 14 : \
20      (x <= 32768) ? 15 : \
21      (x <= 65536) ? 16 : \
22      -1)
23
24  `define CEIL(x,y) (((x)+(y)-1)/(y))
25
26  `define MIN(x,y) ((x) < (y) ? (x) : (y))
27
28  `endif

```



```

1  `include "master_params.vh"
2
3  /*
4  we're doing T mults in a step
5  those T mults can be broken down into a set of
... {inputs,weights}
6  therefore total number of inputs we have = total number of
... weights sent = T
7  how those inputs and weights are grouped is determined by P =
... T / I
8  */
9  module LayerComputation
10     (
11     input clk,
12     input reset,
13     input startPulse,
14     input [`PARAM_I_I_M-1:0] incomingWeights,
15     input [`PARAM_I_M-1:0] incomingThresholds,
16     input [`PARAM_I_N-1:0] incomingInputs,
17
18     output reg [`PARAM_I_N-1:0] outputPool_finishedVals,
19     output reg outputPool_writeEnable
20     );
21
22     wire [`PARAM_I_I_N-1:0] inputDuplicates =
... {`MAX_NET_WIDTH{incomingInputs[`PARAM_I_N-1:0]}};
23
24     // Parameters //
25
26     localparam S_WAITING = 2'd0;
27     localparam S_SUM = 2'd1;
28     localparam S_DONE = 2'd2;
29     // Debugging
30     //assign outputPool_finishedVals = incomingInputs;
31     //assign outputPool_writeEnable = 1;
32
33     // Registers //
34
35     reg [1:0] state = S_WAITING; //"Processing State"
36
37     reg [`PARAM_N_M_S1_I_I-1:0] multiplicationResult;

```

```

38
39     reg sumStartPulse = 0;
40     wire sumOutputReady;
41     wire signed [`PARAM_N_M_S1+3:0] sumOutput0;
42     wire signed [`PARAM_N_M_S1+3:0] sumOutput1;
43     wire signed [`PARAM_N_M_S1+3:0] sumOutput2;
44     wire signed [`PARAM_N_M_S1+3:0] sumOutput3;
45     wire signed [`PARAM_N_M_S1+3:0] sumOutput4;
46     wire signed [`PARAM_N_M_S1+3:0] sumOutput5;
47     wire signed [`PARAM_N_M_S1+3:0] sumOutput6;
48     wire signed [`PARAM_N_M_S1+3:0] sumOutput7;
49     wire signed [`PARAM_N_M_S1+3:0] sumOutput8;
50     wire signed [`PARAM_N_M_S1+3:0] sumOutput9;
51     wire signed [`PARAM_N_M_S1+3:0] sumOutput10;
52     wire signed [`PARAM_N_M_S1+3:0] sumOutput11;
53     wire signed [`PARAM_N_M_S1+3:0] sumOutput12;
54     wire signed [`PARAM_N_M_S1+3:0] sumOutput13;
55     wire signed [`PARAM_N_M_S1+3:0] sumOutput14;
56     wire signed [`PARAM_N_M_S1+3:0] sumOutput15;
57     wire signed [`PARAM_N_M_S1+3:0] sumOutput16;
58     wire signed [`PARAM_N_M_S1+3:0] sumOutput17;
59     wire signed [`PARAM_N_M_S1+3:0] sumOutput18;
60     wire signed [`PARAM_N_M_S1+3:0] sumOutput19;
61     wire signed [`PARAM_N_M_S1+3:0] sumOutput20;
62     wire signed [`PARAM_N_M_S1+3:0] sumOutput21;
63     wire signed [`PARAM_N_M_S1+3:0] sumOutput22;
64     wire signed [`PARAM_N_M_S1+3:0] sumOutput23;
65
66     // Tree Adders //
67
68     /* module TreeAdder24
69     (
70     input clk,
71     input reset,
72     input [`PARAM_N_M_S1_I-1:0] data,
73     input inputReady,
74     output outputReady,
75     output reg [`PARAM_N_M_S1+3:0] dataOut
76     ); */
77
78     TreeAdder24 n0

```

```

78... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*0)+:
... `PARAM_N_M_S1_I], sumStartPulse, sumOutputReady, sumOutput0);
79   TreeAdder24 n1
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*1)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput1);
80   TreeAdder24 n2
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*2)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput2);
81   TreeAdder24 n3
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*3)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput3);
82   TreeAdder24 n4
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*4)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput4);
83   TreeAdder24 n5
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*5)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput5);
84   TreeAdder24 n6
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*6)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput6);
85   TreeAdder24 n7
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*7)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput7);
86   TreeAdder24 n8
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*8)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput8);
87   TreeAdder24 n9
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*9)+:
... `PARAM_N_M_S1_I], sumStartPulse, , sumOutput9);
88   TreeAdder24 n10
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*10)+: `
... PARAM_N_M_S1_I], sumStartPulse, , sumOutput10);
89   TreeAdder24 n11
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*11)+: `
... PARAM_N_M_S1_I], sumStartPulse, , sumOutput11);
90   TreeAdder24 n12
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*12)+: `
... PARAM_N_M_S1_I], sumStartPulse, , sumOutput12);
91   TreeAdder24 n13
... (clk, reset, multiplicationResult[(`PARAM_N_M_S1_I*13)+: `
... PARAM_N_M_S1_I], sumStartPulse, , sumOutput13);

```



```

92     TreeAdder24 n14
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*14)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput14);
93     TreeAdder24 n15
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*15)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput15);
94     TreeAdder24 n16
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*16)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput16);
95     TreeAdder24 n17
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*17)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput17);
96     TreeAdder24 n18
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*18)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput18);
97     TreeAdder24 n19
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*19)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput19);
98     TreeAdder24 n20
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*20)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput20);
99     TreeAdder24 n21
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*21)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput21);
100    TreeAdder24 n22
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*22)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput22);
101    TreeAdder24 n23
...   (clk,reset,multiplicationResult[(`PARAM_N_M_S1_I*23)+:
...   PARAM_N_M_S1_I],sumStartPulse,,sumOutput23);
102
103    function signed [`PARAM_N_M_S1:0]shift1 (input signed
...   [`PARAM_N_M:0] data);
104        shift1 = data[`PARAM_N_M-1:`SHIFT_1];
105    endfunction
106
107    function signed [`PARAM_N_M_LG_S1_S2-1:0] shift2 (input
...   signed [`PARAM_N_M_LG_S1-1:0] data);
108        shift2 = data[`PARAM_N_M_LG_S1-1:`SHIFT_2];
109    endfunction
110

```



```

111     function signed [`INPUT_WIDTH-1:0] relu (input signed
...   [`PARAM_N_M_LG_S1_S2-1:0] data);
112         relu = data[`PARAM_N_M_LG_S1_S2-1] ? 'b0 :
... data[`PARAM_N_M_LG_S1_S2-2:0];
113     endfunction
114
115     wire signed [`PARAM_N_M_S1-1:0] test_mul, test_mul2,
... test_mul5;
116     wire [(2*`PARAM_N_M_S1)-1:0] test_mul3;
117     wire signed [`PARAM_N_M_S1_I-1:0] test_mul4;
118     wire signed [`PARAM_N_M_S1+3:0] test_sum;
119     wire signed [`PARAM_N_M_S1+10:0] test_sum3;
120     wire signed [`WEIGHT_WIDTH-1:0] test_sum1;
121     wire signed [`PARAM_N_M_LG_S1_S2-1:0] test_sum2;
122     wire signed [`PARAM_N_M_S1-1:0] muls [`PARAM_I_I-1:0];
123     // Loop indices //
124
125     integer i_dotProd = 0;
126     genvar i_test;
127     integer test_idx = 574; // index into multiplication for
... loop
128     integer test_idx_2 = 23; // index into corresponding adder
... tree
129
130     // Sequential Logic //
131     /**The first is wrong? Copies an extra msb from
... inputDuplicates*/
132     assign test_mul = {{(`PARAM_N_M_S1-`INPUT_WIDTH){'b0}},
... inputDuplicates[(test_idx*`INPUT_WIDTH)+:`INPUT_WIDTH]};
133     assign test_mul2 =
... $signed(incomingWeights[(test_idx*`WEIGHT_WIDTH)+:`
... WEIGHT_WIDTH]);//{{(`PARAM_N_M_S1-`INPUT_WIDTH){'b0}},
... inputDuplicates[15:0]};
134     assign test_mul3[0+:`PARAM_N_M_S1] =
... shift1($signed(inputDuplicates[(test_idx*`INPUT_WIDTH)+:`
... INPUT_WIDTH]) *
... $signed(incomingWeights[(test_idx*`WEIGHT_WIDTH)+:`
... WEIGHT_WIDTH]));
135     assign test_mul4 =
... multiplicationResult[(`PARAM_N_M_S1_I*test_idx_2)+:`
... PARAM_N_M_S1_I];// input to sum23

```

```

136     assign test_mul5 =
... multiplicationResult[(`PARAM_N_M_S1*test_idx)+:`PARAM_N_M_S1];
137     assign test_sum = sumOutput23;
138     assign test_sum1 =
... incomingThresholds[(`WEIGHT_WIDTH*test_idx_2)+:`WEIGHT_WIDTH];
139     assign test_sum2 = shift2($signed({1'b0, sumOutput23})) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*test_idx_2)+:`
... WEIGHT_WIDTH]);
140     assign test_sum3 = relu(shift2($signed({1'b0,
... sumOutput23})) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*test_idx_2)+:`
... WEIGHT_WIDTH]));
141
142     generate
143         for(i_test = 0; i_test < `PARAM_I_I; i_test = i_test +
... 1) begin
144             assign muls[i_test] =
... multiplicationResult[(i_test*`PARAM_N_M_S1)+:`PARAM_N_M_S1];
145         end
146     endgenerate
147
148
149
150     always @(posedge clk) begin
151         if(reset) begin
152             state <= S_WAITING;
153             outputPool_finishedVals <= -1;
154             outputPool_writeEnable <= 0;
155         end
156         else case(state)
157             S_WAITING: begin
158                 outputPool_writeEnable <= 0;
159                 if(startPulse) begin //start multiplication
160                     //indexing scheme = [(i*WIDTH)-1 :
... (i-1)*WIDTH ]
161                     for(i_dotProd = 0; i_dotProd < `PARAM_I_I;
... i_dotProd = i_dotProd + 1) begin
162
... multiplicationResult[(i_dotProd*`PARAM_N_M_S1)+:`PARAM_N_M_S1]
... <=
... shift1($signed(inputDuplicates[(i_dotProd*`INPUT_WIDTH)+:`

```

```

162... INPUT_WIDTH]) *
... $signed(incomingWeights[(i_dotProd*`WEIGHT_WIDTH)+: `
... WEIGHT_WIDTH]));
163
... //{{(`PARAM_N_M_S1-`INPUT_WIDTH){'b0}},
... inputDuplicates[(i_dotProd*`INPUT_WIDTH)+: `INPUT_WIDTH]};//
... shift1($signed(inputDuplicates[(i_dotProd*`INPUT_WIDTH)+: `
... INPUT_WIDTH]) *
... $signed(incomingWeights[(i_dotProd*`WEIGHT_WIDTH)+: `
... WEIGHT_WIDTH]));
164                                     //(multiplicationResult)i = (input)i *
... (weight)i
165                                     end
166                                     sumStartPulse <= 1;
167                                     state <= S_SUM;
168                                 end
169                             end
170                             S_SUM: begin
171                                 sumStartPulse <= 0;
172                                 if (sumOutputReady) begin
173                                     // testing weight input by sending back
... truncated weights
174
... outputPool_finishedVals[(`INPUT_WIDTH*0)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput0}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*0)+: `WEIGHT_WIDTH]))
... ;
175
... outputPool_finishedVals[(`INPUT_WIDTH*1)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput1}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*1)+: `WEIGHT_WIDTH]))
... ;
176
... outputPool_finishedVals[(`INPUT_WIDTH*2)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput2}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*2)+: `WEIGHT_WIDTH]))
... ;
177
... outputPool_finishedVals[(`INPUT_WIDTH*3)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput3}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*3)+: `WEIGHT_WIDTH]))

```



```

177... ;
178
... outputPool_finishedVals[(`INPUT_WIDTH*4)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput4}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*4)+: `WEIGHT_WIDTH]))
... ;
179
... outputPool_finishedVals[(`INPUT_WIDTH*5)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput5}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*5)+: `WEIGHT_WIDTH]))
... ;
180
... outputPool_finishedVals[(`INPUT_WIDTH*6)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput6}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*6)+: `WEIGHT_WIDTH]))
... ;
181
... outputPool_finishedVals[(`INPUT_WIDTH*7)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput7}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*7)+: `WEIGHT_WIDTH]))
... ;
182
... outputPool_finishedVals[(`INPUT_WIDTH*8)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput8}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*8)+: `WEIGHT_WIDTH]))
... ;
183
... outputPool_finishedVals[(`INPUT_WIDTH*9)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput9}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*9)+: `WEIGHT_WIDTH]))
... ;
184
... outputPool_finishedVals[(`INPUT_WIDTH*10)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput10}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*10)+: `WEIGHT_WIDTH]))
... );
185
... outputPool_finishedVals[(`INPUT_WIDTH*11)+: `INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput11}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*11)+: `WEIGHT_WIDTH]))
... );

```

```
186 ... outputPool_finishedVals[(`INPUT_WIDTH*12)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput12}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*12)+:`WEIGHT_WIDTH])
... );
187 ... outputPool_finishedVals[(`INPUT_WIDTH*13)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput13}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*13)+:`WEIGHT_WIDTH])
... );
188 ... outputPool_finishedVals[(`INPUT_WIDTH*14)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput14}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*14)+:`WEIGHT_WIDTH])
... );
189 ... outputPool_finishedVals[(`INPUT_WIDTH*15)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput15}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*15)+:`WEIGHT_WIDTH])
... );
190 ... outputPool_finishedVals[(`INPUT_WIDTH*16)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput16}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*16)+:`WEIGHT_WIDTH])
... );
191 ... outputPool_finishedVals[(`INPUT_WIDTH*17)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput17}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*17)+:`WEIGHT_WIDTH])
... );
192 ... outputPool_finishedVals[(`INPUT_WIDTH*18)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput18}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*18)+:`WEIGHT_WIDTH])
... );
193 ... outputPool_finishedVals[(`INPUT_WIDTH*19)+:`INPUT_WIDTH] <=
... relu(shift2($signed({1'b0, sumOutput19}))) +
... $signed(incomingThresholds[(`WEIGHT_WIDTH*19)+:`WEIGHT_WIDTH])
... );
194
```

```

194... outputPool_finishedVals[(`INPUT_WIDTH*20)+: `INPUT_WIDTH] <=
...   relu(shift2($signed({1'b0, sumOutput20}))) +
...   $signed(incomingThresholds[(`WEIGHT_WIDTH*20)+: `WEIGHT_WIDTH])
... );
195
...   outputPool_finishedVals[(`INPUT_WIDTH*21)+: `INPUT_WIDTH] <=
...   relu(shift2($signed({1'b0, sumOutput21}))) +
...   $signed(incomingThresholds[(`WEIGHT_WIDTH*21)+: `WEIGHT_WIDTH])
... );
196
...   outputPool_finishedVals[(`INPUT_WIDTH*22)+: `INPUT_WIDTH] <=
...   relu(shift2($signed({1'b0, sumOutput22}))) +
...   $signed(incomingThresholds[(`WEIGHT_WIDTH*22)+: `WEIGHT_WIDTH])
... );
197
...   outputPool_finishedVals[(`INPUT_WIDTH*23)+: `INPUT_WIDTH] <=
...   relu(shift2($signed({1'b0, sumOutput23}))) +
...   $signed(incomingThresholds[(`WEIGHT_WIDTH*23)+: `WEIGHT_WIDTH])
... );
198           outputPool_writeEnable <= 1;
199           state <= S_WAITING;
200           end
201       end
202       default: state <= S_WAITING;
203   endcase
204 end
205 endmodule

```

Layer Controller


```

1 `include "master_params.vh"
2
3 module LayerController
4     (
5         input clk,
6         input reset,
7
8         input [`BTW_MAX_NET_DEPTH-1:0] commCTL_networkDepth,
9         input commCTL_networkDepthReady,
10
11        input [`BTW_MAX_PARALLELISM-1:0] parallelPool_parallelism,
12        input parallelPool_parallelismReady,
13
14        input [`BTW_MAX_NET_WIDTH-1:0]
... neuronIndicesPool_maxNeuronIndex,
15        input neuronIndicesPool_maxNeuronIndexReady,
16
17        input memCTL_weightsAreReady,
18        input inputPool_inputsAreReady,
19
20        input layerComputation_dataReady,
21
22        output [`BTW_MAX_PARALLELISM-1:0] outputPool_dataWidth,
23        output [`BTW_MAX_NET_WIDTH-1:0] outputPool_maxNeuronIndex,
24
25        output inputPool_getFromOutputPool,
26
27        output [`BTW_MAX_NET_DEPTH-1:0] memCTL_currentLayer,
28        output [`BTW_MAX_NET_WIDTH-1:0] memCTL_currentNeuron,
29        output [`BTW_MAX_PARALLELISM-1:0] memCTL_parallelism,
30        output reg memCTL_getLayerDataRequest, //for layer params
31        output reg memCTL_getNeuronDataRequest, //for neuron
... weights, needs to keep in mind the parallelism above
32
33        output commCTL_isStateDone,
34        output reg layerComputation_startPulse // "begin" pulse
35    );
36
37    // Parameters //
38
39    localparam S_WAITING_INPUTS = 2'd0;

```

```

40     localparam S_WAITING_LAYERCOMP = 2'd1;
41     localparam S_DONE = 2'd2;
42
43     //indices to architectureReady:
44     localparam ARCH_PARALLELISM = 0;
45     localparam ARCH_MAXNEURON = 1;
46
47     // Registers //
48
49     reg [1:0] state = S_WAITING_INPUTS; //"Processing State"
50
51     reg networkDepthLoaded = 0;
52     reg [`BTW_MAX_NET_DEPTH-1:0] networkDepth =
... {(`CLOG2(`BTW_MAX_NET_DEPTH)){1'b0}}; //NOTE: actually needs
... to be depth -1, since again, it's 0 indexed
53
54     //microarchitectural ready states, all need to be 1 to
... enable the system
55     reg [1:0] architectureReady = 2'b0;
56
57     reg [`BTW_MAX_NET_DEPTH-1:0] currentLayer =
... {(`CLOG2(`BTW_MAX_NET_DEPTH)){1'b0}}; // set to 0
58     reg [`BTW_MAX_NET_WIDTH-1:0] currentNeuron =
... {(`CLOG2(`BTW_MAX_NET_WIDTH)){1'b0}}; // set to 0
59
60     reg [`BTW_MAX_NET_WIDTH-1:0] maxNeuronIndex; //"Total
... Neurons in Layer"
61     reg [`BTW_MAX_PARALLELISM-1:0] parallelism; //"Parallel
... Increment (P)"
62
63     // Wires //
64
65     // We add parallelism because this should be true if we
... are calculating last neuron
66     // currentNeuron is first neuron in set being calculated
67     wire isLastNeuron = (currentNeuron + parallelism >=
... maxNeuronIndex) & architectureReady[ARCH_MAXNEURON]; //sent to
... inputPool
68     wire isLastLayer = ((currentLayer == networkDepth-1) ||
... networkDepth == 0)& networkDepthLoaded;
69

```



```

70 // Assign statements //
71
72 assign commCTL_isStateDone = state == S_DONE;
73
74 assign inputPool_getFromOutputPool = isLastNeuron &
... layerComputation_dataReady & (~isLastLayer);
75 assign outputPool_dataWidth = parallelism; //todo check if
... P valid here or?
76 assign outputPool_maxNeuronIndex = maxNeuronIndex; //todo
... check if M valid here or?
77
78 assign memCTL_currentLayer = currentLayer;
79 assign memCTL_currentNeuron = currentNeuron;
80 assign memCTL_parallelism = parallelism;
81
82 // nextState Functions //
83
84 //We've sent the requests and have our layer architecture,
... wait for inputs to be ready
85 task doState_WaitingInputs;
86 begin
87     if(memCTL_weightsAreReady && inputPool_inputsAreReady)
... begin
88         layerComputation_startPulse <= 1;
89         state <= S_WAITING_LAYERCOMP;
90     end
91
92     //reset previous requests
93     memCTL_getLayerDataRequest <= 0;
94     memCTL_getNeuronDataRequest <= 0;
95 end
96 endtask
97
98 //We've sent a request to the layer computation module,
... wait for it to finish calculating
99 task doState_WaitingLayerComp;
100 begin
101     layerComputation_startPulse <= 0;
102     if(layerComputation_dataReady) begin //if the
... computation is done
103         if(isLastNeuron) begin

```

```

104         if(isLastLayer) state <= S_DONE;
105         else begin
106             //update layer
107             currentLayer <= currentLayer + 1;
108             currentNeuron <=
... {(`CLOG2(`BTW_MAX_NET_WIDTH)){1'b0}}; // set to 0
109
110             //request next layer architecture
111             memCTL_getLayerDataRequest <= 1; //is set
... to 0 in next state
112
113             //request weights for the new neuron-layer
... combo
114             memCTL_getNeuronDataRequest <= 1; //is set
... to 0 in next state
115             state <= S_WAITING_INPUTS; //wait for
... weights to refresh
116         end
117
118         //reset layer architecture flags
119         architectureReady <= 2'b0;
120     end
121     else begin
122         currentNeuron <= currentNeuron + parallelism;
123
124         //request weights for the new neuron-layer
... combo
125         memCTL_getNeuronDataRequest <= 1; //is set to
... 0 in next state
126         state <= S_WAITING_INPUTS; //wait for weights
... to refresh
127     end
128 end
129 endtask
130
131
132 // Sequential Logic //
133
134 always @(posedge clk) begin
135     if(reset) begin
136         state <= S_WAITING_INPUTS;

```

```

137     currentLayer <= 0;
138     currentNeuron <= 0;
139     architectureReady <= 2'b0;
140     memCTL_getLayerDataRequest <= 0;
141     memCTL_getNeuronDataRequest <= 0;
142     layerComputation_startPulse <= 0;
143     networkDepthLoaded <= 0;
144     end
145     else if(~networkDepthLoaded) begin //global parameter
... needs to be set
146         if(commCTL_networkDepthReady) begin
147             networkDepth <= commCTL_networkDepth;
148             networkDepthLoaded <= 1; //never gets unset
149
150             //send initial architecture (layer data)
... request
151             memCTL_getLayerDataRequest <= 1;
152             memCTL_getNeuronDataRequest <= 1;
153         end
154     end
155     else if(!(&{architectureReady})) begin //if any bit of
... architectureReady is 0 (i.e. if the architecture is not ready)
156         //reset previous request
157         memCTL_getLayerDataRequest <= 0;
158         memCTL_getNeuronDataRequest <= 0;
159
160         //Load neuron indices if RDY
161         if(neuronIndicesPool_maxNeuronIndexReady) begin
162             maxNeuronIndex <=
... neuronIndicesPool_maxNeuronIndex;
163             architectureReady[ARCH_MAXNEURON] <= 1;
164         end
165
166         //Load parallelism values if RDY
167         if(parallelPool_parallelismReady) begin
168             parallelism <= parallelPool_parallelism;
169             architectureReady[ARCH_PARALLELISM] <= 1;
170         end
171     end
172     else case(state)
173         S_WAITING_INPUTS: doState_WaitingInputs;

```

```
174         S_WAITING_LAYERCOMP: doState_WaitingLayerComp;
175         S_DONE: ; //intentionally empty state, waits for
... external reset to begin next computation
176         default: state <= S_WAITING_INPUTS;
177     endcase
178 end
179 endmodule
```



```

1 `ifndef MASTER_PARAMS_H
2 `define MASTER_PARAMS_H
3
4 `include "constant_funcs.vh"
5
6 // Basic parameters
7 `define NODES_PER_STREAM 9 //for communication module
8 // `define WRITE_CNTRL 15 // TODO: Width of address wire from
... comm module to
9 `define MAX_NET_WIDTH 24 // I
10 `define MAX_NET_DEPTH 4 // D
11 `define INPUT_WIDTH 16 // N
12 `define WEIGHT_WIDTH 24 // M
13 `define SHIFT_1 3 //S1
14 `define SHIFT_2 1 //S2
15 `define DEBUG_WIDTH 8 // E
16 `define MAX_DEBUG_MSGS 9 // E_max
17 `define BITS_PER_STREAM 80
18 `define DSP_SLICES 740 // Based on # DSP slices
19 // Same as PARAM_I_I_M1 `define MAX_WEIGHT_BITS_PER_STEP
... `MAX_NET_WIDTH*`MAX_NET_WIDTH*`WEIGHT_WIDTH+`MAX_NET_WIDTH //
... This is the logical max weights. With 24 width, 4 depth, we
... know the output of min. `MIN(`DSP_SLICES * `WEIGHT_WIDTH,
... `MAX_NET_WIDTH*`MAX_NET_WIDTH*`WEIGHT_WIDTH+`MAX_NET_WIDTH)
20 `define MAX_PARALLELISM `MAX_NET_WIDTH //theoretical max is
... $sqrt(`DSP_SLICES) = 27, but 24 now as we only do entire layers
... in a compute steps
21
22 // Bit widths
23 `define BTW_MAX_NET_WIDTH $clog2(`MAX_NET_WIDTH)
24 `define BTW_MAX_NET_DEPTH $clog2(`MAX_NET_DEPTH)
25 `define BTW_MAX_PARALLELISM $clog2(`MAX_PARALLELISM)
26 `define BTW_MAX_WEIGHTS_PER_STEP
... $clog2(`MAX_WEIGHT_BITS_PER_STEP)
27 `define BTW_BITS_PER_STREAM $clog2(`BITS_PER_STREAM)
28
29 // Derived Parameters
30 `define PARAM_I_M (`MAX_NET_WIDTH*`WEIGHT_WIDTH) // I*M
31 `define PARAM_I_N (`MAX_NET_WIDTH*`INPUT_WIDTH) // I*N
32 `define PARAM_I_D_N
... (`MAX_NET_WIDTH*`MAX_NET_DEPTH*`INPUT_WIDTH) // I*N

```

```

33 `define PARAM_E_EM (`DEBUG_WIDTH*`MAX_DEBUG_MSGS) // E*E_max
34 // Subtract 24 to account for header, parallelism, max_node
... data in first data packet
35 `define STREAM_WEIGHT_OVERLAP
... ((`PARAM_I_M1-(`BITS_PER_STREAM-(8+24+24))) %
... (`BITS_PER_STREAM-8) == 0 ? `BITS_PER_STREAM-8 :
... (`PARAM_I_M1-(`BITS_PER_STREAM-(8+24+24))) %
... (`BITS_PER_STREAM-8)) // Overlap into last packet stream for
... weight data
36 `define NODES_PER_STREAM_OVERLAP `STREAM_WEIGHT_OVERLAP % 8
37 `define STREAM_INPUT_OVERLAP (`PARAM_I_N % (`BITS_PER_STREAM-8)
... == 0) ? `BITS_PER_STREAM-8 : (`PARAM_I_N %
... (`BITS_PER_STREAM-8)) // Overlap into last packet stream for
... inputs data
38 `define STREAMS_PER_OUTPUTS (`CEIL(`PARAM_I_N,
... `BITS_PER_STREAM-8))
39 `define PACKETS_PER_STREAM (`BITS_PER_STREAM/8)
40 `define TOTAL_BITS_PER_STREAM (`PACKETS_PER_STREAM*10) // Bits
... per stream including start and stop bits
41 // Derived Parameters
42 `define PARAM_I_I (`MAX_NET_WIDTH * `MAX_NET_WIDTH)
... //interconnect of 2 layers
43 `define PARAM_I_I_N (`MAX_NET_WIDTH*`PARAM_I_N)
44 `define PARAM_I_M1 (`MAX_NET_WIDTH*(`WEIGHT_WIDTH+1)) //
... I*(M+1) := Weight/Threshold bits per node
45 `define PARAM_I_I_M (`MAX_NET_WIDTH*`PARAM_I_M) // I^2*(M) :=
... Weight bits per layer
46 `define PARAM_I_I_M1 (`MAX_NET_WIDTH*`PARAM_I_M1) //
... I^2*(M+1) := Weight/Threshold bits per layer
47 `define PARAM_N_M (`INPUT_WIDTH + `WEIGHT_WIDTH) // N+M
48 `define PARAM_N_M_S1 (`PARAM_N_M - `SHIFT_1) // N+M-S1
49 `define PARAM_N_M_S1_I (`PARAM_N_M_S1 * `MAX_NET_WIDTH) //
... I(N+M-S1) //all the weights for one node
50 `define PARAM_N_M_S1_I_I (`PARAM_N_M_S1_I * `MAX_NET_WIDTH) //
... I(N+M-S1) //all the weights for one layer
51 `define PARAM_N_T (`INPUT_WIDTH*`NUM_MULTIPLIERS)
52
53
54 `define PARAM_N_M_LG_S1 (`PARAM_N_M + `CLOG2(`MAX_NET_WIDTH+1)
... - `SHIFT_1)
55 `define PARAM_N_M_LG_S1_S2 (`PARAM_N_M +

```

```
55... `CLOG2(`MAX_NET_WIDTH+1) - `SHIFT_1 - `SHIFT_2)
56
57 // BRAM write widths
58 `define BRAM_WIDTH `PARAM_I_M1 // For now we read in a single
... node's weights per step
59 //`define BRAM_BTW_ADDR
... $clog2(`MAX_NET_DEPTH*(`MAX_NET_WIDTH*`MAX_NET_WIDTH*`
... WEIGHT_WIDTH +`MAX_NET_WIDTH))
60 `define BRAM_BTW_ADDR $clog2(`MAX_NET_DEPTH*`MAX_NET_WIDTH)
61
62 `endif
63
```



```

1 `timescale 1ns / 1ps
2
3 `include "master_params.vh"
4
5
6 // Current implementation reads and writes a single node's
... weights per timestep
7 // During writing phase, expects all nodes in a layer to have
... weights sent in-order
8 // TODO: Do we have to buffer write requests?
9 module memory_controller(
10     input clk,
11     input reset,
12     input write_weights_rdy,
13     input rd_layer_en,
14     input rd_node_en,
15     input [`PARAM_I_M1-1:0] write_layer_weights, // Note that
... we currently expect a single node's weights + cutoff per
... write.
16     input [`BTW_MAX_NET_DEPTH-1:0] write_layer_id,
17     input [`BTW_MAX_NET_DEPTH-1:0] read_layer_id,
18     input [`BTW_MAX_NET_WIDTH-1:0] read_node_id,
19     input [`BTW_MAX_PARALLELISM-1:0] read_parallelism,
20     output reg [`PARAM_I_I_M-1:0] read_weights,
21     output reg [`PARAM_I_M-1:0] read_thresholds,
22     output reg read_weights_rdy);
23
24
25     localparam READS_PER_REQ = `MAX_NET_WIDTH; // Layer
... controller will send for each computation step. Currently this
... is all nodes in a layer.
26     localparam WRITES_PER_NODE = 1; // Port width set to write
... all of a node's data at a time
27     localparam STATE_WAITING = 0; // We are not waiting for
... BRAM to finish reading. Inputs are monitored
28     localparam STATE_READING = 1; // BRAM is reading into
... output buffer. All inputs ignored except reset
29     localparam STATE_WRITING = 2;
30
31
32     //wire bram_start_addr = read_layer_id * `PARAM_I_M1 +

```



```

32... read_node_id*`WEIGHT_WIDTH + read_node_id;
33   wire [`BRAM_BTW_ADDR-1:0] bram_start_addr = read_layer_id
... * `MAX_NET_WIDTH + read_node_id;
34   reg [`BRAM_BTW_ADDR-1:0] bram_cur_addr;
35   reg [$clog2(READS_PER_REQ):0] bram_reads_left;
36   wire [`BRAM_WIDTH-1:0] cur_rd_weights;
37
38   reg [1:0] state;
39   reg last_wr_layer;
40   reg cur_wr_layer;
41   reg [`BRAM_BTW_ADDR-1:0] bram_wr_offset;
42   reg [$clog2(WRITES_PER_NODE):0] bram_writes_left;
43   reg [`PARAM_I_M1-1:0] cur_wr_weights;
44   wire dummy_bram_busy;
45
46   wire bram_en = bram_writes_left > 0 || bram_reads_left >
... 0;
47
48   blk_width24_depth4 bram_inst (
49     .clka(clk),           // input wire clka
50     .rsta(reset),        // input wire rsta
51     .ena(bram_en),       // input wire ena
52     .wea(bram_writes_left > 0), // input wire
... [0 : 0] write enable
53     .addra(bram_cur_addr), // input wire [14 : 0] addra
54     .dina(cur_wr_weights[`BRAM_WIDTH-1:0]), // input wire
... [374 : 0] dina
55     .douta(cur_rd_weights), // output wire [374 : 0] douta
56     .rsta_busy(dummy_bram_busy) // dummy reset busy. Assume
... we don't care about this for now
57   );
58
59   always @(posedge clk) begin
60     if (reset) begin
61       read_weights_rdy <= 0;
62       state <= STATE_WAITING;
63       cur_wr_layer <= 0;
64       last_wr_layer <= 1;
65       bram_writes_left <= 0;
66       bram_reads_left <= 0;
67       bram_cur_addr <= 0;

```

```

68         end else begin
69             case (state)
70             STATE_WAITING: begin
71                 if (write_weights_rdy) begin
72                     // Write layer weights `PARAM_I_M bits to
... BRAM
73                     // Reset layer offset if we are writing
... nodes in a new layer
74                     if (last_wr_layer != write_layer_id) begin
75                         bram_wr_offset <= 0;
76                         last_wr_layer <= write_layer_id;
77                     end
78                     cur_wr_layer <= write_layer_id;
79                     cur_wr_weights <= write_layer_weights;
80                     bram_writes_left <= WRITES_PER_NODE;
81                     state <= STATE_WRITING;
82                 end else if (rd_node_en) begin
83                     // Addresses BRAM using layer_id, node_id,
... and parallelism
84                     bram_cur_addr <= bram_start_addr;
85                     bram_reads_left <= READS_PER_REQ;
86                     // Prefill output weight and thresholds to
... account for layers that do not read to fill the buffer
87                     // NOTE: that this neuron id check only
... works because we do a full layer per computation step
88                     if (read_node_id == 0) begin
89                         read_weights <= 0;
90                         read_thresholds <= 0;
91                     end
92                     state <= STATE_READING;
93                 end else begin
94                     if (read_weights_rdy) read_weights_rdy <=
... 0;
95                     end
96                 end
97             STATE_READING: begin
98                 // delay first read from bram by a cycle due
... to delay
99                 if (bram_reads_left != READS_PER_REQ) begin
100                     if (bram_reads_left == 0) begin
101                         state <= STATE_WAITING;

```

```

102         read_weights_rdy <= 1;
103     end
104     // least significant weight of every node
... is the cutoff
105         read_weights <=
... {cur_rd_weights[`BRAM_WIDTH-1:`WEIGHT_WIDTH],
... read_weights[`PARAM_I_I_M-1:`BRAM_WIDTH-`WEIGHT_WIDTH]};
106         read_thresholds <=
... {cur_rd_weights[`WEIGHT_WIDTH-1:0],
... read_thresholds[`PARAM_I_M-1:`WEIGHT_WIDTH]};
107         bram_cur_addr <= bram_cur_addr +
... 1; // `BRAM_WIDTH;
108     end
109     bram_reads_left <= bram_reads_left - 1;
110 end
111 STATE_WRITING: begin
112     if (bram_writes_left == 1) state <=
... STATE_WAITING;
113     cur_wr_weights <= cur_wr_weights >>
... `BRAM_WIDTH;
114     bram_writes_left <= bram_writes_left - 1;
115     bram_cur_addr <= bram_cur_addr +
... 1; // `BRAM_WIDTH;
116     end
117     default: ;
118     endcase
119 end
120 end
121
122 endmodule
123

```



```

1 `timescale 1ns / 1ps
2
3 `include "master_params.vh"
4
5 // General pool module to support input pool and
... parallelism/max_node_id pools
6 // This can be split into two specialized pool as some
... functionality only exists for the former/latter
7 module pool #(DATA_WIDTH = `BTW_MAX_PARALLELISM) (
8     input clk,
9     input reset,
10    input write_en, // Initial write enable
11    input [DATA_WIDTH-1:0] write_data, // Initial write data
12    input read_en, // read enable. asserts index is valid
13    input [`BTW_MAX_NET_DEPTH-1:0] index, // write/read index
14    output reg data_rdy,
15    output reg [DATA_WIDTH-1:0] read_data // output data
16 );
17
18    reg [(`MAX_NET_WIDTH*DATA_WIDTH)-1:0] storage;
19    always @(posedge clk) begin
20        if (reset) begin
21            storage <= 0;
22            data_rdy <= 0;
23        end else begin
24            if (write_en) begin
25                storage[index*DATA_WIDTH +: DATA_WIDTH] <=
... write_data;
26            end else if (read_en) begin
27                read_data <= storage[index*DATA_WIDTH +:
... DATA_WIDTH];
28                data_rdy <= 1;
29            end else if (data_rdy) begin
30                data_rdy <= 0;
31            end
32        end
33    end
34 endmodule
35
36 module input_pool(
37     input clk,

```



```

38     input reset,
39     input comm_write_en, // Initial write enable
40     input cntrl_pool_transfer, // Transfer output pool to
... input pool
41     input [`PARAM_I_N-1:0] comm_data, // Initial write data
42     input [`PARAM_I_N-1:0] output_pool_data, // Intermediate
... data to overwrite
43     output reg [`PARAM_I_N-1:0] read_data, // output data
44     output reg read_data_rdy
45 );
46
47     always @(posedge clk) begin
48         if (reset) begin
49             read_data <= 0;
50             read_data_rdy <= 0;
51         end else begin
52             if (comm_write_en | cntrl_pool_transfer) begin
53                 read_data <= (cntrl_pool_transfer) ?
... output_pool_data : comm_data;
54                 read_data_rdy <= 1;
55             end else if (read_data_rdy) read_data_rdy <= 0;
56         end
57     end
58 endmodule
59
60 module output_pool(
61     input clk,
62     input reset,
63     input write_en,
64     input [`BTW_MAX_PARALLELISM-1:0] parallelism,
65     input [`BTW_MAX_NET_WIDTH-1:0] max_node_id,
66     input [`PARAM_I_N-1:0] write_data,
67     output reg [`PARAM_I_N-1:0] cur_outputs,
68     output reg done_layer // Calculated by layer controller.
... Output here should matche for debugging
69 );
70
71     integer i;
72     reg [`BTW_MAX_NET_WIDTH-1:0] cur_write_node;
73     always @(posedge clk) begin
74         if (reset) begin

```

```

75         cur_outputs <= 1;
76         cur_write_node <= 0;
77         done_layer <= 0;
78         end else if (write_en) begin
79             //cur_outputs <= {{(`PARAM_I_N-24-8){1'b1}},
... {3'b0, cur_write_node}, {3'b0, parallelism}, {3'b0,
... max_node_id}, 8'b00};
80             //done_layer <= 1;
81             // partial layer output write to cur_outputs based
... on data and parallelism
82             for (i = 0; i < `MAX_NET_WIDTH; i = i + 1) begin
83                 if (i >= cur_write_node && i < cur_write_node
... + parallelism && i <= max_node_id) begin
84                     // Update new node outputs
85                     cur_outputs[i*`INPUT_WIDTH +:
... `INPUT_WIDTH] <=
... write_data[`INPUT_WIDTH*(i-cur_write_node)+:`INPUT_WIDTH];
86                     end else if (i >= cur_write_node) begin
87                         // Zero future node outputs. This handles
... zeroing outputs that are never written
88                         cur_outputs[i*`INPUT_WIDTH +:
... `INPUT_WIDTH] <= 'b0;
89                     end
90                 end
91                 if (cur_write_node + parallelism <= max_node_id)
... begin
92                     cur_write_node <= cur_write_node +
... parallelism;
93                 end else begin
94                     cur_write_node <= 0;
95                     done_layer <= 1;
96                 end
97             end else if (done_layer) begin
98                 done_layer <= 0;
99             end
100         end
101     endmodule
102

```

Tree Adder

```

1 `include "master_params.vh"
2
3 /*
4
5 SPECIFICATION:
6
7 assert inputReady
8 next cycle data has to be ready to latch (must be held only
... for that cycle)
9 when outputReady is asserted the data will be available on
... dataOut
10 (same clock cycle) and will remain until inputReady is
... asserted again
11
12 reset does not need to be called before this, and there is no
... mechanism
13 to clear garbage data because it is not needed in the use-case
14
15 */
16
17 module TreeAdder24
18     (
19     input clk,
20     input reset,
21     input [`PARAM_N_M_S1_I-1:0] data,
22     input inputReady,
23     output outputReady,
24     output reg [`PARAM_N_M_S1+3:0] dataOut
25     );
26
27     //Tree-layout adder, log delay
28     //sumA layer 0
29     reg signed [`PARAM_N_M_S1:0] sumA00;
30     reg signed [`PARAM_N_M_S1:0] sumA01;
31     reg signed [`PARAM_N_M_S1:0] sumA02;
32     reg signed [`PARAM_N_M_S1:0] sumA03;
33     reg signed [`PARAM_N_M_S1:0] sumA04;
34     reg signed [`PARAM_N_M_S1:0] sumA05;
35     reg signed [`PARAM_N_M_S1:0] sumA06;
36     reg signed [`PARAM_N_M_S1:0] sumA07;
37

```



```

38 //sumA layer 1
39 reg signed [`PARAM_N_M_S1+1:0] sumA10;
40 reg signed [`PARAM_N_M_S1+1:0] sumA11;
41 reg signed [`PARAM_N_M_S1+1:0] sumA12;
42 reg signed [`PARAM_N_M_S1+1:0] sumA13;
43
44 //sumA layer 2
45 reg signed [`PARAM_N_M_S1+2:0] sumA20;
46 reg signed [`PARAM_N_M_S1+2:0] sumA21;
47
48 //sumA layer 3
49 reg signed [`PARAM_N_M_S1+3:0] sumA3;
50
51
52 //sumB latch 0
53 reg signed [`PARAM_N_M_S1:0] sumB00;
54 reg signed [`PARAM_N_M_S1:0] sumB01;
55 reg signed [`PARAM_N_M_S1:0] sumB02;
56 reg signed [`PARAM_N_M_S1:0] sumB03;
57
58 //sumB latch 1
59 reg signed [`PARAM_N_M_S1:0] sumB10;
60 reg signed [`PARAM_N_M_S1:0] sumB11;
61 reg signed [`PARAM_N_M_S1:0] sumB12;
62 reg signed [`PARAM_N_M_S1:0] sumB13;
63
64
65 //sumB layer 1
66 reg signed [`PARAM_N_M_S1+1:0] sumB20;
67 reg signed [`PARAM_N_M_S1+1:0] sumB21;
68
69 //sumB layer 2
70 reg signed [`PARAM_N_M_S1+2:0] sumB3;
71
72 reg active = 0;
73
74 reg [3:0] cycleCount = 3'b0;
75 assign outputReady = cycleCount == 3'd5;
76
77 always @(posedge clk) begin
78     if(reset | inputReady) begin

```



```

79         cycleCount <= 3'b0;
80         active <= inputReady;
81     end
82     else if(active) begin
83         cycleCount <= cycleCount + 1;
84         if(cycleCount == 3'd4) active <= 0; //will be 5
... next round, we are done
85
86         //sumA data latch
87         sumA00 <= data[`PARAM_N_M_S1 *(0) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(1) +: `PARAM_N_M_S1];
88         sumA01 <= data[`PARAM_N_M_S1 *(2) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(3) +: `PARAM_N_M_S1];
89         sumA02 <= data[`PARAM_N_M_S1 *(4) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(5) +: `PARAM_N_M_S1];
90         sumA03 <= data[`PARAM_N_M_S1 *(6) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(7) +: `PARAM_N_M_S1];
91         sumA04 <= data[`PARAM_N_M_S1 *(8) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(9) +: `PARAM_N_M_S1];
92         sumA05 <= data[`PARAM_N_M_S1 *(10) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(11) +: `PARAM_N_M_S1];
93         sumA06 <= data[`PARAM_N_M_S1 *(12) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(13) +: `PARAM_N_M_S1];
94         sumA07 <= data[`PARAM_N_M_S1 *(14) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(15) +: `PARAM_N_M_S1];
95
96         //sumB data latch
97         sumB00 <= data[`PARAM_N_M_S1 *(16) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(17) +: `PARAM_N_M_S1];
98         sumB01 <= data[`PARAM_N_M_S1 *(18) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(19) +: `PARAM_N_M_S1];
99         sumB02 <= data[`PARAM_N_M_S1 *(20) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(21) +: `PARAM_N_M_S1];
100        sumB03 <= data[`PARAM_N_M_S1 *(22) +:
... `PARAM_N_M_S1] + data[`PARAM_N_M_S1 *(23) +: `PARAM_N_M_S1];
101
102        //sumB layer 1 ( latch to keep pipeline uniform )
103        sumB10 <= sumB00;
104        sumB11 <= sumB01;
105        sumB12 <= sumB02;
106        sumB13 <= sumB03;

```

```
107
108 //sumA layer 1
109 sumA10 <= sumA00+sumA01;
110 sumA11 <= sumA02+sumA03;
111 sumA12 <= sumA04+sumA05;
112 sumA13 <= sumA06+sumA07;
113
114 //sumA layer 2
115 sumA20 <= sumA10+sumA11;
116 sumA21 <= sumA12+sumA13;
117
118 //sumA layer 3
119 sumA3 <= sumA20+sumA21;
120
121 //sumB layer 2
122 sumB20 <= sumB10 + sumB11;
123 sumB21 <= sumB12 + sumB13;
124
125 //sumB layer 3
126 sumB3 <= sumB20 + sumB21;
127
128 //final Sum (layer 4)
129 dataOut <= sumA3 + sumB3;
130     end
131   end
132 endmodule
133
```

UART

```

1 `include "master_params.vh"
2
3 // Modifications from base:
4 // Parameterized uart interface to support variable packet
... stream size
5
6 //115200 * 16 oversampling = 1843200 samples / sec
7 //100000000 hz / N = 1843200 hz
8 //N ~~~ 54 (54.253)
9 // TODO: test max baud
10 module uart_clkdiv //to account for non blocking assignment
... timing
11     (input clk,
12     input reset,
13     output clk_out);
14
15     localparam TARGET = 'd54;
16     reg [8:0] count; //9 bits needed to store values up to and
... including TARGET
17
18     assign clk_out = (count == 9'b0);
19     always @(posedge clk) begin
20         if (reset) begin
21             count <= 0;
22         end else begin
23             if(count == TARGET-1) count <= 9'b0;
24             else count <= count + 1;
25         end
26     end
27 endmodule
28
29 // 10 packets per stream. 8 bits per packet
30 `define BITS_PER_STREAM 80
31 //2ms hold constraint:
32 //115200 hz ^-1 = 5.425 E -7 sec / cycle
33 //5.425 E -7 * NUM_CYCLES = .002 sec
34 //VERIFY_SAMPLES = NUM_CYCLES ~~~ 230 ( = 230.4 ) (8 bits
... needed to store)
35
36 module uart_rx
37     (input clk,

```

```

38     input uart_clk,
39     input reset,
40     input signal,
41     input verifyStopBit, //for debugging purposes, negative
... logic
42     output reg [`BITS_PER_STREAM-1:0] data_out,
43     output reg ready = 0); //THIS IS A PULSE
44
45     localparam VERIFY_SAMPLES = 'd230;
46
47     localparam S_VERIFY_HIGH = 0;
48     localparam S_WAITING = 1;
49     localparam S_LISTENING = 2;
50     localparam S_VERIFYSTOP = 3;
51
52     reg [1:0] state;
53     reg [7:0] tempData;
54     reg [7:0] timingCounter;
55     reg [2:0] bitCounter; //counts how many bits sampled
56     reg [3:0] cycleCounter; //counts how many cycles have
... passed
57     reg [`BITS_PER_STREAM-1:0] data;
58
59
60     always @(posedge clk) begin
61         if (reset) begin
62             state <= 2'b0;
63             tempData <= 8'b0;
64             timingCounter <= 12'b0;
65             bitCounter <= 3'b0;
66             cycleCounter <= 3'b0;
67             data <= 'b0;
68             data_out <= -1;
69             ready <= 0;
70         end else if (uart_clk) begin
71             case(state)
72                 S_VERIFY_HIGH: begin
73                     ready <= 0;
74                     if(signal) begin
75                         if(timingCounter == VERIFY_SAMPLES)
... begin //will be 307 on next cycle

```



```

76             state <= S_WAITING; //we have
... verified we are not in the middle of data transmission, go to
... wait state
77             if(cycleCounter >= 'd9) begin
78                 cycleCounter <= 3'b0;
79             end
80             timingCounter <= 12'b0;
81         end
82         else timingCounter <= timingCounter +
... 1;
83     end
84     else timingCounter <= 12'b0;
85 end
86 S_WAITING: begin
87     if(~signal) begin //signal has gone low
88         if(timingCounter == 'd7) begin //will
... be 8 the next cycle, start sampling
89             state <= S_LISTENING;
90             timingCounter <= 'd0;
91         end
92         else timingCounter <= timingCounter +
... 1;
93     end
94     else timingCounter <= 12'b0; //likely
... noise, keep waiting for data transmission start
95     end
96     S_LISTENING: begin
97         if(timingCounter == 12'd15) begin
98             timingCounter <= 'd0;
99             tempData <= {signal,tempData[7:1]};
100            if(bitCounter == 'd7) begin // this
... cycle is the 8th bit, we're done
101                state <= S_VERIFYSTOP;
102                bitCounter <= 3'b0;
103            end
104            else bitCounter <= bitCounter + 1;
105        end
106        else timingCounter <= timingCounter + 1;
107    end
108    S_VERIFYSTOP: begin
109        if(timingCounter == 12'd15) begin

```

```

110             timingCounter <= 12'b0;
111             if(verifyStopBit | signal) begin//stop
... bit found, data valid (or debug enabled)
112                 data <=
... {tempData,data[`BITS_PER_STREAM-1:8]};
113                 tempData <= 0;
114                 if(cycleCounter == 'd9) begin
115                     ready <= 1; //this is the 10th
... byte, send ready pulse
116                     data_out <=
... {tempData,data[`BITS_PER_STREAM-1:8]};
117                     state <= S_VERIFY_HIGH;
118                     end
119                     else state <= S_WAITING;
120                     cycleCounter <= cycleCounter + 1;
121                     end
122                     else state <= S_WAITING;
123                     end
124                     else timingCounter <= timingCounter +1;
125                     end
126                     default: state <= S_VERIFY_HIGH;
127                 endcase
128             end else if (ready == 1) ready <= 0;
129         end
130     endmodule
131
132 module uart_tx(
133     input clk,
134     input reset,
135     input en,
136     input [`TOTAL_BITS_PER_STREAM-1:0] start_data,
137     output reg xmit_data,
138     output reg done);
139
140     localparam SEND_RATE = 16;
141     localparam DELAY = 100; // UART cycles to wait before
... sending data. Possible bug in delay from computer send inputs
... to waiting for outputs
142     localparam N = 868;
143
144     reg signed [`TOTAL_BITS_PER_STREAM-1:0] shift_data;

```

```

145     reg [`CLOG2(`TOTAL_BITS_PER_STREAM)-1:0] cntnr;
146     reg [9:0] downsampler; // xmit_clk is the 16x sample clock
... from receiver. Need to send slower.
147     reg [9:0] delay;
148     reg [9:0] trans_cntr; // 100 Mhz * n = 115200 hz => n =
... 100000000/115200 = 868
149
150     always @(posedge clk) begin
151         if (reset) begin
152             xmit_data <= 1;
153             cntnr <= `TOTAL_BITS_PER_STREAM;
154             shift_data <= start_data;
155             done <= 0;
156             downsampler <= 0;//SEND_RATE;
157             delay <= DELAY;
158             trans_cntr <= N-1;
159         end else if (trans_cntr == 0) begin//(xmit_clk && cntnr
... > 0) begin
160             trans_cntr <= N-1;
161             if (delay == 0) begin
162                 //if (downsampler == 0) begin
163                     //downsampler <= SEND_RATE;
164                     cntnr <= cntnr - 1;
165                     xmit_data <= shift_data[0];
166                     shift_data <= shift_data >>> 1;
167                     if (cntnr == 1) done <= 1;
168                 // end else begin
169                     // downsampler <= downsampler - 1;
170                 // end
171             end else
172                 delay <= delay - 1;
173         end else if (done) begin
174             done <= 0;
175         end else if (en) begin
176             trans_cntr <= trans_cntr - 1;
177         end
178     end
179
180 endmodule
181

```

Appendix B: Python Visualization Code

(Note: Solutions to the 6.034 lab assignment have been intentionally omitted. These are not key files in interfacing the FPGA neural network with the “training.py” script that is given as part of the 6.034 starting lab material, which has been modified to suit this project.)

Main Script

(app.py)


```

1 from comm import pc_uart
2 from ctypes import c_uint8
3 import random
4 import sys
5 import time
6 import labinterface
7
8 TESTPRINT = False
9 def TP(*args):
10     if TESTPRINT:
11         TP(*args)
12 # Matched with communication_module.v for MODE values expected
... in packet headers
13 class stream_types(): #ENUM
14     WEIGHTS = 0
15     INPUTS = 1
16     OUTPUT = 2
17     DEBUG = 3
18
19 MAX_NET_WIDTH = 24
20 MAX_NET_DEPTH = 4
21 NUM_LAYERS= 1
22 CRC_EN = 0
23 WEIGHT_WIDTH = 24
24 INPUT_WIDTH = 16
25 NUM_NODES = [15]
26 parallelism = [15]
27
28 # ensure this is 0 padded in actual implementation
29 test_weights = [NUM_NODES[i] * [(MAX_NET_WIDTH+1)*[1]] for i
... in range(len(NUM_NODES))]
30 max_node = [NUM_NODES[i]-1 for i in range(len(NUM_NODES))]
31 # 24 nodes, each has weight array [cutoff=1, weight_0 = 1,
... ..., weight_23 = 1]
32 test_inputs = 12*[0x20,0x20] #range(15)
33 uart_ifc = None
34
35 def make_header(data_type, header_data):
36     header = (data_type | (CRC_EN << 2) | (header_data <<
... 3)).to_bytes(1, byteorder='little')
37     return header

```

```

38
39 def get_size(data_type):
40     if data_type == stream_types.WEIGHTS:
41         return 3
42     if data_type == stream_types.INPUTS:
43         return 2
44     if data_type == stream_types.DEBUG:
45         return 1
46
47 def ceil_div(a, b):
48     return -(-a // b)
49
50 def create_send_stream(data_type, header_data, data_bytes):
51     global uart_ifc
52     header_bytes = make_header(data_type, header_data)
53     packet_stream = header_bytes + data_bytes
54     uart_ifc.send_stream(packet_stream)
55
56 # def send_data_array(data_type, header_data, data_arr):
57 #     data_i = 0
58 #     # check data being sent is a multiple of packet size
59 #     packets_per_num = get_size(data_type)
60 #     increment = 9 // packets_per_num
61 #     overlap = 9 % packets_per_num
62 #
63 #     while data_i + increment <= len(data_arr):
64 #         data_bytes = [i.to_bytes(get_size(data_type),
... byteorder='little') for i in
... data_arr[data_i:data_i+increment]]
65 #         if overlap > 0:
66 #             overlap_start = data_i + increment
67 #             overlap_end = data_i + increment + overlap
68 #             data_bytes += [i.to_bytes(get_size(data_type),
... byteorder='little') for i in
69 #
... data_arr[overlap_start:overlap_end]]
70 #             data_i += overlap
71 #             TP(str(data_bytes))
72 #             create_send_stream(data_type, header_data,
... data_bytes)
73 #             data_i += increment

```

```

74 #     # handle leftover
75 #     leftover = len(data_arr) - data_i - 1
76 #     if leftover > 0:
77 #         TP(leftover)
78 #         data = data_arr[data_i:data_i+leftover] +
... (9-leftover)*[0]
79 #         TP(str(data))
80 #         create_send_stream(data_type, header_data, data)
81
82 def send_data_array(data_type, header_data, data_arr):
83     # check data being sent is a multiple of packet size
84     data_bytes = b''.join([i.to_bytes(get_size(data_type),
... bytearray='little', signed=True) for i in data_arr])
85     TP("DATA (" +str(len(data_arr)) + ") = " + str(data_arr))
86     TP("DATA_BYTES (" + str(len(data_bytes)) + ") = " +
... str(data_bytes))
87     overlap = len(data_bytes) % 9
88     i = 0
89     while i + 9 < len(data_bytes):
90         stream_data_bytes = data_bytes[i:i+9]
91         TP("STREAM_DATA_BYTES = " + str(stream_data_bytes))
92         create_send_stream(data_type, header_data,
... stream_data_bytes)
93         i += 9
94     # handle leftover
95     if overlap > 0:
96         TP("overlap = " + str(overlap))
97         data = data_bytes[-overlap:] + (9-overlap)*b'\x00'
98         TP(str(data))
99         create_send_stream(data_type, header_data, data)
100 def get_outputs():
101     global uart_ifc
102     outputs = b''
103     while len(outputs) < ceil_div(MAX_NET_WIDTH*INPUT_WIDTH,
... 8):
104         streams = uart_ifc.rec_stream()
105         if streams:
106             for stream in [streams[i:i+10] for i in range(0,
... len(streams), 10)]:
107                 TP(str(stream[1:].hex()))
108                 outputs += stream[1:]

```



```

109         else:
110             return None
111         #grouped_bytes = [[outputs[i+1:i+2] + outputs[i:i+1]] for
... i in range(0,len(outputs),2)]
112         grouped_bytes = [[outputs[i:i+1] + outputs[i+1:i+2]] for i
... in range(0,len(outputs),2)]
113         TP('{{}'.format(', '.join(x[0].hex() for x in
... grouped_bytes)))
114         outputs = [int.from_bytes(b''.join(g), byteorder='little',
... signed=True) for g in grouped_bytes]
115         return outputs[0:24]
116 def main():
117     global uart_ifc
118     if len(sys.argv) < 2:
119         TP("Please pass part of port name to connect to. E.g.
... COM7")
120         return -1
121     uart_ifc = pc_uart()
122     if uart_ifc.setup(str(sys.argv[1])) == False:
123         TP("Exiting....")
124         return -1
125     # WEIGHTS
126     try:
127         TP("Generating weights")
128         #while(True):
129         #    uart_ifc.send_stream(b''.join([int(0).to_bytes(1,
... byteorder='little', signed=True)] + 9*[int(0).to_bytes(1,
... byteorder='big')]))
130         #    time.sleep(0.5)
131         ## generate weight array
132         weight_data = [NUM_LAYERS*[]]
133         for i in range(NUM_LAYERS):
134             weight_data[i] = [parallelism[0], max_node[0]]
135             for j in range(NUM_NODES[i]):
136                 weight_data[i] += test_weights[i][j]
137         TP("Sending weights" + str(test_weights))
138         for layer_id in range(NUM_LAYERS):
139             send_data_array(stream_types.WEIGHTS, layer_id,
... weight_data[layer_id])
140         TP("Sending inputs")
141         send_data_array(stream_types.INPUTS, 0, test_inputs)

```



```

142     TP("Sending debug send_done packet")
143     send_data_array(stream_types.DEBUG, 0, [0])
144     TP("Waiting for outputs...")
145     # Wait for outputs
146     outputs = get_outputs()
147     TP("outputs = " + str(outputs))
148     finally:
149         TP("Closing serial port")
150         uart_ifc.close()
151 def verify_TestNet():
152     #inputs = [12*[1<<5,1<<5]]
153     inputs = [[0,0]]
154     #weights = 2*[7 * [(24)*[1]] + (24-7)* [24*[0]] for i in
... range(1)]
155     weights = [[[-27, 40, 0], [-27, 40, 0], [-27, 40, 0]],
... [[24, 24, 24], [24, 24, 24], [0, 0, 0]], [[40, 40, 0], [0, 0,
... 0], [0, 0, 0]]]
156     #thresh = 2*[24 * [1] + (24-7) * [24*[0]] for i in
... range(1)]
157     thresh = [[17, 17, 17], [-6, -6, 0], [-11, 0, 0]]
158     depth = len(thresh)
159     inputs = [l + (24-len(l))*[0] for l in inputs]
160     weights = [l + (24-len(l))*[0] for n in weights for l in
... n]
161     thresh = [l + (24-len(l))*[0] for l in thresh]
162     newNet = labinterface.TestNet(inputs, weights, thresh,
... depth)
163     TP("Running on FPGA")
164     run(inputs, weights, thresh, len(weights[0]), depth-1)
165     newNet.printFstLayer()
166     TP('-----')
167     newNet.printResults()
168 def run_rand():
169     global uart_ifc
170     if len(sys.argv) < 2:
171         TP("Please pass part of port name to connect to. E.g.
... COM7")
172         return -1
173     uart_ifc = pc_uart()
174     if uart_ifc.setup(str(sys.argv[1])) == False:
175         TP("Exiting....")

```

```

176     return -1
177 # WEIGHTS
178 try:
179     newNet = labinterface.TestNet()
180     TP("Generating weights")
181     #while(True):
182     #    uart_ifc.send_stream(b''.join([int(0).to_bytes(1,
... byteorder='little', signed=True)] + 9*[int(0).to_bytes(1,
... byteorder='big')]))
183     #    time.sleep(0.5)
184     ## generate weight array
185     # reverse weight to have correct ordering relative to
... lsb
186     #weight_data.reverse()
187     for layer_id in range(NUM_LAYERS):
188         combined = [num_nodes, num_nodes-1] # parallelism
... and max node
189         for i,l in enumerate(newNet.weights[layer_id]):
190             combined += [newNet.thresholds[layer_id][i]] +
... l
191             TP(str(combined))
192             TP("Sending weights: " + str(combined))
193             send_data_array(stream_types.WEIGHTS, layer_id,
... combined)
194             TP("Sending inputs")
195             newNet.printFstLayer()
196             send_data_array(stream_types.INPUTS, 0,
... newNet.values[0])
197             TP("Sending debug send_done packet")
198             send_data_array(stream_types.DEBUG, 0, [0])
199             TP("Waiting for outputs...")
200             # Wait for outputs
201             outputs = get_outputs()
202             TP("outputs = " + str(outputs))
203             TP("Actual Results: ")
204             newNet.printResults()
205             TP("Results match? " + newNet.checkResults(outputs))
206
207 finally:
208     TP("Closing serial port")
209     uart_ifc.close()

```

```

210 def run(inputs, weights, thresholds, num_nodes, depth,
... port_name = "COM8"):
211     single_layer = depth == 1
212     max_node = [num_nodes-1 for i in range(depth)]
213     # WEIGHTS
214     if uart_ifc is None:
215         setup()
216     try:
217         TP("Generating weights")
218         # reverse weight to have correct ordering relative to
... lsb
219         #weight_data.reverse()
220         for layer_id in range(depth):
221             combined = [num_nodes, num_nodes-1] # parallelism
... and max node
222             for i,l in enumerate(weights[layer_id]):
223                 combined += [thresholds[layer_id][i]] + l +
... (24-len(l)) * [0]
224                 TP("Sending weights: " + str(combined))
225                 send_data_array(stream_types.WEIGHTS, layer_id,
... combined)
226                 TP("Sending inputs")
227                 inputs = inputs + (24-len(inputs)) * [0]
228                 TP(str(inputs))
229                 send_data_array(stream_types.INPUTS, 0, inputs)
230                 TP("Sending debug send_done packet")
231                 send_data_array(stream_types.DEBUG, 0, [0])
232                 TP("Waiting for outputs...")
233                 # Wait for outputs
234                 outputs = []
235                 if not single_layer:
236                     for i in range(MAX_NET_DEPTH):
237                         outputs.append(get_outputs())
238                 else:
239                     outputs.append(get_outputs())
240                 print("outputs (final first) = " + str(outputs))
241                 return outputs
242
243     except:
244         TP("Closing serial port")
245         uart_ifc.close()

```



```

246 def setup(port_name = "COM8"):
247     global uart_ifc
248     if port_name == "":
249         TP("Please pass part of port name to connect to. E.g.
... COM7")
250         return -1
251     uart_ifc = pc_uart()
252     if uart_ifc.setup(port_name) == False:
253         TP("Exiting....")
254         return -1
255 if __name__ == "__main__":
256     setup()
257     main()
258     #t = labinterface.TestNet.run_working(0)
259     #t[0].printResults()
260     #verify_TestNet()
261     #run_rand()
262

```



```

1 import serial.tools.list_ports
2 import sys
3 import time
4 from time import sleep
5
6 TESTPRINT = False
7 def TP(*args):
8     if TESTPRINT:
9         TP(*args)
10 #Version 2.7 or Above?
11 if sys.version_info[0] >2:
12     version3 = True
13     kwargs = {'newline':''}
14 else:
15     version3 = False
16     kwargs = {}
17 def get_usb_port(port_name):
18     usb_port = list(serial.tools.list_ports.grep("USB-Serial
... Controller"))
19     if len(usb_port) == 1:
20         print("Automatically found USB-Serial Controller:
... {}".format(usb_port[0].description))
21         return usb_port[0].device
22     else:
23         ports = list(serial.tools.list_ports.comports())
24         port_dict = {i:[ports[i],ports[i].vid] for i in
... range(len(ports))}
25         usb_id=None
26         for p in port_dict:
27             print("{}: {} (Vendor ID:
... {})".format(p,port_dict[p][0],port_dict[p][1]))
28             if port_name in port_dict[p][0]:
29                 print("\tAbove Port Matches Argument:
... "+str(port_name))
30                 usb_id = p
31         if usb_id== None:
32             return False
33         else:
34             print("USB-Serial Controller: Device {}".format(p))
35             return port_dict[usb_id][0].device
36 class pc_uart:

```

```

37 def __init__(self):
38     self._serial = None
39 def setup(self, port_name):
40     s = get_usb_port(port_name)
41     if s:
42         self._serial = serial.Serial(port = s,
43                                     baudrate=115200,
44                                     parity=serial.PARITY_NONE,
45                                     stopbits=serial.STOPBITS_ONE,
46                                     bytesize=serial.EIGHTBITS,
47                                     timeout=None)
48         print(self._serial)
49         print("Serial Connected!")
50         if self._serial.isOpen():
51             print(self._serial.name + ' is open...')
52         return True
53     else:
54         print("USB-Serial Controller Not Found")
55         return False
56 def rec_stream(self):
57     data = []
58     data = self._serial.read(10)
59     TP("stream: " + str(data.hex()))
60     if len(data) > 0:
61         return data
62     else:
63         print("ERR: Received invalid packet stream")
64         return None
65     #sleep(3.0)
66 def send_stream(self, stream):
67     assert len(stream) == 10
68     TP("STREAM_BYTES: " + str(stream))
69     time.sleep(0.005) # hold constraint
70     TP("# packets sent: " +
... str(self._serial.write(stream)))
71 def close(self):
72     self._serial.close()
73

```

```

1 import lab6, random, app, sys, time
2
3 def NNGet(net):
4     """NeuralNet object -> lists"""
5     neuronNameToPos = {}
6     posToNeuronName = {}
7     neurons = [[]]
8     currLayerI = 0
9     currNodeI = 0
10    maxNodeI = 0
11
12    #get neuron list
13    currLayer = neurons[0]
14    for neuron in net.topological_sort(True):
15        for parentNode in net.get_incoming_neighbors(neuron):
16            if parentNode in neurons[currLayerI]: #next layer
17                found
18                    currLayerI += 1
19                    maxNodeI = max(maxNodeI, currNodeI)
20                    currNodeI = 0
21                    neurons.append([])
22                    break
23            if neuron != -1:
24                neurons[currLayerI].append(neuron)
25                neuronNameToPos[neuron] = (currLayerI, currNodeI)
26                posToNeuronName[(currLayerI, currNodeI)] = neuron
27                currNodeI += 1
28
29    #for layer in neurons: #zero pad to max size
30    weights = [[[]]*maxNodeI for _ in range(maxNodeI)] for _
31    ... in range(currLayerI)]
32    thresholds = [[[]]*maxNodeI for _ in range(currLayerI)]
33
34    #Reference (copied from below)
35    #values[1][2] -> layer, neuron
36    #weights[1][2][3] -> layer, neuron, prevNeuron
37    #thresholds[1][2] -> layer, neuron
38
39    for wire in net.get_wires():
40        n0 = wire.startNode
41        n1_layerI, n1_nodeI = neuronNameToPos[wire.endNode]

```



```

40     n1_layerI -= 1 #shift array data left by 1, since
... there is nothing before the 0th layer, and this indexes by
... prev layer
41     if n0 == -1: thresholds[n1_layerI][n1_nodeI] =
... -wire.weight
42     else:
43         n0_layerI, n0_nodeI = neuronNameToPos[n0]
44         assert n0_layerI == n1_layerI #no +1 because we
... did a -1 above to n1_layerI
45         weights[n1_layerI][n1_nodeI][n0_nodeI] =
... wire.weight
46
47     return
... (neurons,weights,thresholds,maxNodeI,currLayerI,
... posToNeuronName)
48 def FPGAFwdProp(net,inputs,runOnFPGA):
49     """ return (final output, dict {neuron -> output}
50     this function has to stall until the final outputs are
... ready"""
51
52     if runOnFPGA is False: return
... lab6.forward_prop(net,inputs)
53     else:
54         neurons,weights,thresholds, numNodes, numLayers,
... posToNeuronName = NNGet(net)
55
56         inputLayer = [int(i) for i in inputs.values()]
57         DEBUG = True
58         if DEBUG:
59             print("n",neurons)
60             print("w",weights)
61             print("t",thresholds)
62             print("i",inputLayer)
63             print("params",numNodes,numLayers)
64             print("pos -> neuron name",posToNeuronName)
65
66         #_ = input("CALL FPGA HERE") #todo remove this
67         # in case we only have single layer
68         print(inputLayer)
69         netOutputs = [inputLayer]
70         port_name = "COM8"

```



```

71         if True:
72             cur_inputs = inputLayer
73             for layer in range(numLayers):
74                 outputs = app.run(cur_inputs, weights[layer:],
... thresholds[layer:], numNodes, 1, port_name)
75                 if type(outputs) == int:
76                     break
77                 netOutputs.append(outputs[0])
78                 cur_inputs = outputs[0]
79                 #time.sleep(0.05)
80         else:
81             netOutputs = app.run(inputLayer, weights,
... thresholds, numNodes, numLayers-1,
... port_name)#FPGA_FUNCTIONCALL() #todo, returns a 1d list of raw
... neuron outputs
82             if type(netOutputs) == int:
83                 print("FPGA forward propagation returned error: "
... + str(netOutputs))
84                 sys.exit("Please kill and restart after fixing the
... error")
85             #make 2d list
86             print(str(netOutputs))
87             #netOutputs = [netOutputs[i:i+numNodes] for i in
... range(0, len(netOutputs), numNodes)]
88             print(str(netOutputs))
89             print(str(posToNeuronName))
90             #make into dict
91             values = {}#{neuron:output}
92             for layerI in range(len(netOutputs)):
93                 for neuronI in range(len(netOutputs[layerI])):
94                     if (layerI,neuronI) in posToNeuronName:
95                         name = posToNeuronName[(layerI,neuronI)]
96                         values[name] = netOutputs[layerI][neuronI]
97
98             return (values[net.get_output_neuron()],values)
99 class working_tests():
100     inputs = [
101         [12*[1<<5,1<<5]],
102     ]
103     weights = [
104         2*[5 * [(24)*[1]] for i in range(1)],

```

```

105     ]
106     thresholds = [
107         2*[5 * [1] for i in range(1)],
108     ]
109
110 class TestNet():
111     layerWidth = 24
112     inputBitWidth = 16
113     weightBitWidth = 24
114     shift1 = 3
115     shift2 = 1
116     def __init__(self, inputs = None, weights = None,
117 ... thresholds = None, depth=2):
118         self.layerDepth = depth
119         #create values and initialize first layer:
120 ... values[1][2] -> layer, neuron
121         if inputs is None:
122             self.values =
123 ... [[self.getRandomSigned(self.inputBitWidth) for _ in
124 ... range(self.layerWidth)]]
125         else:
126             self.values = inputs
127         #create weights: weights[1][2][3] -> layer, neuron,
128 ... prevNeuron
129         if weights is None:
130             self.weights =
131 ... [[[self.getRandomSigned(self.weightBitWidth) for _ in
132 ... range(self.layerWidth)] for _ in range(self.layerWidth)] for _
133 ... in range(self.layerDepth)]
134         else:
135             self.weights = weights
136         #create thresholds for each neuron: thresholds[1][2]
137 ... -> layer, neuron
138         if thresholds is None:
139             self.thresholds =
140 ... [[self.getRandomSigned(self.weightBitWidth) for _ in
141 ... range(self.layerWidth)] for _ in range(self.layerDepth)]
142         else:
143             self.thresholds = thresholds
144         print(self.weights)
145         print(self.values)

```

```

135         print(self.thresholds)
136
137         for i in range(1,self.layerDepth):
138             currLayer = []
139             self.values.append(currLayer)
140             for currNeuron in range(self.layerWidth):
141
142 ... currLayer.append(max(0,(sum([(self.weights[i][currNeuron][
143 ... prevN] * self.values[i-1][prevN]) >> self.shift1 for prevN in
144 ... range(self.layerWidth)]) >>
145 ... self.shift2)+self.thresholds[i][currNeuron]))
146         # Returns instances of this class containing either a
147 ... specified or all working test cases
148         def run_working(self, test_num = -1):
149             # specific test case chosen
150             new_inst = []
151             if test_num >= 0:
152                 new_inst = [TestNet(inputs =
153 ... working_tests.inputs[test_num], weights =
154 ... working_tests.weights[test_num],
155 ...                               thresholds =
156 ... working_tests.thresholds[test_num])]
157             else:
158                 for i in range(len(working_tests.inputs)):
159
160 ... new_inst.append(TestNet(working_tests.inputs[i],
161 ... working_tests.weights[i], working_tests.thresholds[i]))
162             return new_inst
163
164         @staticmethod
165         def getRandomSigned(n): return
166 ... random.randint(-(2**(n-1)),(2**(n-1))-1)
167
168         def getFstLayer(self): return self.values[0]
169
170         def checkResults(self,actualResult): return actualResult
171 ... == self.values[-1]
172
173         def printLayer(self,n): print "[" + ",".join([str(i) for i
174 ... in self.values[n]]) + "]"
175         def printFstLayer(self): self.printLayer(0)

```

Training

(training.py, modified version of 6.034 starting code)


```

1 #!/usr/bin/env python3
2 # MIT 6.034 Lab 6: Neural Nets
3 # This file originally written by Joel Gustafson and Kenny
... Friedman
4
5 from __future__ import print_function
6 from sys import argv
7 from random import random, randint, shuffle
8 from time import sleep
9 from matplotlib import pyplot, cm
10 import numpy
11 # import _tkinter
12 from lab6 import *
13 from neural_net_api import *
14 import labinterface
15 from math import ceil
16 AMP = 200 #positive data amplitude
17 STATIC_W = 30 #static weight used
18 ACCURACY = -20
19 MAX_STEPS = 50
20
21 colormap = cm.get_cmap("plasma")
22 def multi_accuracy(desired_outputs, actual_outputs):
23     pairs = []
24     actual_outputs.sort()
25     for d_o in sorted(desired_outputs):
26         a_o = actual_outputs.pop(0)
27         while a_o[0] != d_o[0]:
28             a_o = actual_outputs.pop(0)
29         pairs.append((d_o, a_o))
30     total = sum(accuracy(a[1], b[1]) for a, b in pairs)
31     return float(total) / len(pairs)
32 # Multi-point forward propagation
33 def multi_forward_prop(net, resolution=1, runOnFPGA = False):
34     outputs = []
35     data = []
36     for i in range(resolution * 5):
37         y = float(i) / resolution
38         line = []
39         for j in range(resolution * 5):
40             x = float(j) / resolution

```



```

41         #result = forward_prop(net,{'x': x, 'y': y})[0]
42         result = labinterface.FPGAForwardProp(net,{'x': x,
... 'y': y},runOnFPGA)[0]
43         if i % resolution == 0 and j % resolution == 0:
44             outputs.append((x, y), result)
45             line.append(result)
46         data.append(line)
47         data = numpy.array(data)
48         pyplot.pcolor(data, cmap=colormap)
49         # pyplot.pcolor(data)
50         pyplot.pause(0.0001)
51         sleep(0.05)
52         return (sorted(outputs), data)
53 # Backward propagation
54 def multi_update_weights(net, desired_outputs, width=5,
... height=5):
55     shuffle(desired_outputs)
56     for desired_output in desired_outputs:
57         input_values = {'x': desired_output[0][0], 'y':
... desired_output[0][1]}
58         #neuron_outputs = forward_prop(net, input_values)[1]
59         neuron_outputs =
... labinterface.FPGAForwardProp(net,input_values,False)[1]
60         net = update_weights(net, input_values,
... desired_output[1], neuron_outputs)
61         #TP([w.weight for w in net.get_wires()])
62         return net
63 def multi_back_prop(net, desired_outputs, resolution=1):
64     actual_outputs = multi_forward_prop(net, resolution)[0]
65     c = 0
66     current_accuracy = multi_accuracy(desired_outputs,
... actual_outputs)
67     TP("ii",c, current_accuracy)
68     while current_accuracy < ACCURACY:
69         net = multi_update_weights(net, desired_outputs)
70         actual_outputs = multi_forward_prop(net,
... resolution)[0]
71         c += 1
72         current_accuracy = multi_accuracy(desired_outputs,
... actual_outputs)
73         TP("ii",c, current_accuracy)

```

```

74         if c > MAX_STEPS: break
75     return net
76 # Define neural nets
77 def get_small_nn(w=None):
78     if w is None:
79         f = lambda: STATIC_W # 10 * (0.5 - random())
80         w = [f() for n in range(9)]
81     return NeuralNet(['x', 'y', -1], ['A', 'B', 'C']) \
82         .join('x', 'A', w[0]).join('x', 'B', w[1]) \
83         .join('y', 'A', w[2]).join('y', 'B', w[3]) \
84         .join('A', 'C', w[4]).join('B', 'C', w[5]) \
85         .join(-1, 'A', w[6]).join(-1, 'B', w[7]).join(-1, 'C',
... w[8]) \
86         .join('C', NeuralNet.OUT)
87 def get_medium_nn(w=None):
88     if w is None:
89         f = lambda: STATIC_W #10 * (0.5 - random())
90         w = [f() for n in range(20)]
91     return NeuralNet(['x', 'y', -1], list('ABCDEF')) \
92         .join('x', 'A', w[0]).join('x', 'B', w[1]).join('y',
... 'A', w[2]) \
93         .join('y', 'B', w[3]).join('y', 'C', w[4]).join('x',
... 'C', w[5]) \
94         .join(-1, 'A', w[6]).join(-1, 'B', w[7]).join(-1, 'C',
... w[8]) \
95         .join('A', 'E', w[9]).join('A', 'D', w[10]).join('B',
... 'D', w[11]) \
96         .join(-1, 'D', w[12]).join(-1, 'E', w[13]).join(-1,
... 'F', w[14]) \
97         .join('B', 'E', w[15]).join('C', 'E', w[16]).join('C',
... 'D', w[17]) \
98         .join('D', 'F', w[18]).join('E', 'F', w[19]).join('F',
... NeuralNet.OUT)
99 def get_large_nn():
100     w = lambda: STATIC_W #10*(0.5-random())
101     nn = NeuralNet(['x', 'y', -1],
... list('ABCDEFGHIJKLMN OPQRSTUVWXYZ'))
102     #first layer: A-J (10 neurons)
103     for n1 in 'ABCDEFGHIJ':
104         nn.join('x', n1, w()).join('y', n1, w())
105         #second layer: K-T (10 neurons)

```

```

106         for n2 in 'KLMNOPQRST':
107             nn.join(n1, n2, w())
108     #third layer: U-Y (5 neurons)
109     for n3 in 'UVWXY':
110         for n2 in 'KLMNOPQRST':
111             nn.join(n2, n3, w())
112         #final layer: Z (1 neuron)
113         nn.join(n3, 'Z', w())
114     #define Z as output neuron
115     nn.join('Z', NeuralNet.OUT)
116
117     return nn
118 nets = {'small': get_small_nn, 'medium': get_medium_nn,
119 ... 'large': get_large_nn}
120 # Define data sets
121 # horizontal
122 # - - - - -
123 # - - - - -
124 # - - - - -
125 # + + + + +
126 # + + + + +
127 horizontal =
128 ... sorted([(0,0),AMP),((0,1),AMP),((0,2),AMP),((0,3),0),((0,4),0
129 ... ),((1,0),AMP),
130 ... ((1,1),AMP),((1,2),AMP),((1,3),0),((1,4),0),((2,0),AMP),((2,1)
131 ... ,AMP),
132 ... ((2,2),AMP),((2,3),0),((2,4),0),((3,0),AMP),((3,1),AMP),((3,2)
133 ... ,AMP),
134 ... ((3,3),0),((3,4),0),((4,0),AMP),((4,1),AMP),((4,2),AMP),((4,3)
135 ... ,0),((4,4),0)])
136 # diagonal
137 # + + + + -
138 # + + + - -
139 # + + - - -
140 # + - - - -
141 # - - - - -
142 diagonal =
143 ... sorted([(0,0),0),((0,1),AMP),((0,2),AMP),((0,3),AMP),((0,4),1

```



```

136... ), ((1,0),0),
137
... ((1,1),0), ((1,2),AMP), ((1,3),AMP), ((1,4),AMP), ((2,0),0), ((2,1)
... ,0),
138
... ((2,2),0), ((2,3),AMP), ((2,4),AMP), ((3,0),0), ((3,1),0), ((3,2),0
... ),
139
... ((3,3),0), ((3,4),AMP), ((4,0),0), ((4,1),0), ((4,2),0), ((4,3),0),
... ((4,4),0]])
140 # stripe
141 # - - - - +
142 # - - - + -
143 # - - + - -
144 # - + - - -
145 # + - - - -
146 stripe =
... sorted([(0,0),AMP], [(0,1),0], [(0,2),0], [(0,3),0], [(0,4),0], [(
... 1,0),0],
147
... [(1,1),AMP], [(1,2),0], [(1,3),0], [(1,4),0], [(2,0),0], [(2,1),0],
148
... [(2,2),AMP], [(2,3),0], [(2,4),0], [(3,0),0], [(3,1),0], [(3,2),0],
149
... [(3,3),AMP], [(3,4),0], [(4,0),0], [(4,1),0], [(4,2),0], [(4,3),0],
... [(4,4),AMP]])
150 # checkerboard
151 # - - + +
152 # - - + +
153 #
154 # + + - -
155 # + + - -
156 checkerboard =
... sorted([(0,0),AMP], [(1,0),AMP], [(0,1),AMP], [(1,1),AMP], [(3,3)
... ,AMP], [(4,3),AMP],
157
... [(3,4),AMP], [(4,4),AMP], [(0,3),0], [(1,3),0], [(0,4),0],
158
... [(1,4),0], [(3,0),0], [(3,1),0], [(4,0),0], [(4,1),0]])
159 # letterL
160 # + -

```



```

161 # + -
162 # + -
163 # + - - - -
164 # - + + + +
165 letterL =
... sorted([(0,0),0],[(1,0),AMP],[(2,0),AMP],[(3,0),AMP],[(4,0),
... AMP],[(0,1),AMP],
166
... [(1,1),0],[(2,1),0],[(3,1),0],[(4,1),0],[(0,2),AMP],[(0,3),AMP
...
),
167
... [(0,4),AMP],[(1,2),0],[(1,3),0],[(1,4),0]])
168 # moat
169 # - - - - -
170 # - - - - -
171 # - + - - -
172 # - - - - -
173 # - - - - -
174 moat =
... sorted([(0,0),0],[(0,1),0],[(0,2),0],[(0,3),0],[(0,4),0],[(1,
... 4),0],
175
... [(2,4),0],[(3,4),0],[(4,4),0],[(4,3),0],[(4,2),0],[(4,1),0],
176
... [(4,0),0],[(3,0),0],[(2,0),0],[(1,0),0],[(2,2),AMP]])
177 training_data = {'horizontal': horizontal, 'diagonal':
... diagonal, 'stripe': stripe,
178
... 'checkerboard': checkerboard, 'letterL':
... letterL, 'moat': moat}
179 # Main function for training and heatmap
180 def start_training(data=None, net=None, resolution=None,
... fullscreen=False):
181     if data == None:
182         print('defaulting to diagonal training dataset')
183         train = diagonal
184     elif data in training_data:
185         train = training_data[data]
186     else:
187         print('training dataset not found, defaulting to
... diagonal')
188         train = diagonal

```

```

189     if net == None:
190         print('defaulting to medium net')
191         net_fn = get_medium_nn
192     elif net in nets:
193         net_fn = nets[net]
194     else:
195         print('net not found, defaulting to medium net')
196         net_fn = get_medium_nn
197     if resolution == None:
198         print('defaulting to resolution of 1')
199         resolution = 1
200     else:
201         try:
202             resolution = int(resolution)
203             assert resolution > 0
204         except Exception:
205             print('invalid resolution, defaulting to 1')
206             resolution = 1
207     pyplot.ion
208     pyplot.show()
209     if fullscreen:
210         time.sleep(1)
211     nn = net_fn()
212     #print('\nInitial neural net:\n', nn)
213     #print('\nIter, Accuracy:')
214     try: nn = multi_back_prop(nn, train, resolution)
215     except _tkinter.TclError as e:
216         print('\nException caught: ', e, '\n')
217         Athena_ssh_error_message = ("If you are running this
... on Athena "
218             + "over ssh, try sshing again using the -X flag,
... which allows "
219             + "Athena to display GUI windows on your local
... desktop. "
220             + "If you want to see the original stack trace
... instead of this "
221             + "error, find the line in training.py that raises
... this "
222             + "RuntimeError and replace it with 'raise e'.")
223         raise RuntimeError(Athena_ssh_error_message)
224     pyplot.ioff()

```

```

225     print('\nTrained neural net:\n', nn)
226
227     #trained NN, now we make it into integers
228     nn.integerize()
229     print(nn)
230
231     data = multi_forward_prop(nn, resolution, True)[1]
232     pyplot.clf()
233     pc = pyplot.pcolor(data, cmap=colormap)
234     pyplot.colorbar(pc)
235     pyplot.show()
236 if __name__ == "__main__":
237     train = "diagonal"
238     net = "medium"
239     resolution = 1
240     if '-data' in argv:
241         train = argv[argv.index('-data') + 1]
242     else:
243         print('defaulting to diagonal training dataset')
244     if '-net' in argv:
245         net = argv[argv.index('-net') + 1]
246     else:
247         print('defaulting to medium net')
248     if '-resolution' in argv:
249         resolution = argv[argv.index('-resolution') + 1]
250     else:
251         print('defaulting to resolution of 1')
252     start_training(train, net, resolution)
253

```