# Deep Neural Network Accelerator using Verilog

Name: Sounder Rajendran
Unityid: srajend4
StudentID: 200475996

| | | |
|---|---|---|
| Delay (ns to run provided example): 28400 ns<br>Clock period: 5<br># cycles": 5680 cycles | Logic Area: ($um^2$) 12070.2820 $um^2$ | 1/(delay x area) ($ns^{-1}.um^{-2}$)<br><br>2.917 x $10^{-9}$ $ns^{-1}.um^{-2}$ |
| Delay (TA provided example. TA to complete) | | 1/(delay x area) (TA) |

**Abstract:**
A deep neural network accelerator circuit is designed using Verilog, keeping in mind the inherent parallel nature of silicon circuits and parallel algorithms that can be utilized in design of multi-stage and multi-state hardware. Convolution layer, max-pooling layer and ReLu activation layers are implemented, that take in an image matrix containing pixel values and a kernel matrix of desired values, to apply a filter, extract useful information (edge pixels, feature point extraction, etc.) from the image and decrease the overall size of the image while keeping the useful information intact. The approach taken to design, interfaces used, implementation, and results achieved are discussed. A combination of testbench using Verilog, and python based golden output creation and verification using modelsim timing diagrams were used. A performance per area factor of 2.917e[-9] is achieved, which is calculated using synopsis design_vision synthesis tool.

# Deep Neural Network Accelerator using Verilog

Sounder Rajendran

## 1. Introduction:

An efficient multi-stage artificial neural network is designed using Verilog.

Input images are transformed into input arrays and are convoluted against a kernel acting like a filter in order to extract useful information from the image. The hardware accelerator that is designed, processes these image arrays inside a hidden layer. The hidden layer consists of a convolution layer, ReLu, and a max-pooling layer. In the convolution layer (whose outputs are fitted into a range using a linear activation function known as ReLu; ReLu is perfect for activating image bits as the range of pixels does not fall beneath 0), the output of the convolution is then aggregated, which helps to reduce the overall image size, while still preserving the valuable information obtained from convolution. The aggregation algorithm used is max-pooling, which keeps the most intense bit within a chosen group of bits.

The neural network accelerator was designed in a manner such that two convolution elements are being processed every six clock cycles (after all inputs were available for fetching) and thereby we get one-half of the max pooling output every six clock cycles on average.

Usually, convolution and max pooling layers require multiple rows of input from the input and kernel matrices before they can arrive at a single output element. Not waiting for the complete availability of inputs, finding the temporary values of convolution and max pooling, and iteratively arriving at the complete elements of convolution and max-pooling helps transform the general serial algorithm of the convolution layer and max-pooling layer into a highly parallel implementation without spending too many cycles waiting for a complete set of input values.

The algorithms utilized, parallel and pipelined data fetching and processing, the microarchitecture of design, the design choices made, and techniques used to overcome the constraints that allowed arrival at the microarchitecture are discussed in the forthcoming paragraphs.

## 2. Micro-architecture:

The design of any hardware, apart from its own execution algorithms, is restricted by other hardware it depends on. The neural network accelerator that is to be designed depends on input SRAM and weights SRAM to fetch its input and kernel matrices upon which convolution and max pooling are to be done, and it also depends on output SRAM to store its output. The SRAMs are configured in such a way that they contain two blocks of data on each address line which is arranged in a linear fashion. Therefore, in order to obtain all necessary inputs to compute one element of the convolution matrix, we have to have 3 rows and 3 columns of input elements. Since the 9 elements of the kernel remain unchanged, we fetch them once, store them in a register array and use them whenever required. As for the input elements, we fetch data from one address at a time. The number of data elements available is found by using the size value given as the first element. Therefore, we fetch the size value of an input matrix first and use it to calculate the number

of convolution elements. To calculate a particular convolution element, convolution_matrix_element (p, q), the input elements received (3x3) are paired with corresponding kernel elements (3x3) for one-on-one multiplication-accumulation. The MAC function is done till we arrive at the complete convolution output. Since we get two data at the same time, we multiply-accumulate the two inputs at the same time. For every second input received, we have a stray input value that isn't useful at any stage of convolution for the convolution element being calculated. To make use of this value and not discard it, we calculate the value of convolution_matrix_element (p, q+1). Therefore, we calculate 4 multiplication and 2 accumulation values in parallel with two sets of data lines containing 4 data elements. The parallel MAC hardware is shown below. The temporary conv elements calculated are considered complete when we process all 9 elements required to complete convolution calculation. Then the values are transferred from temp registers to the max-pooling stage.

Fig 1: The parallel algorithm

The two values that are received at the max pooling stage, every cycle is compared with each other and the max is stored (for an even row) in a temporary max_pool register with 31 flipflops. For every odd row, the calculated temp_max_pooled elements are additionally compared with the temp_max pooled values stored in the 31 temp_max_pooled registers, and finally, we arrive at the final max_pooled output. In total, this takes 6 cycles to fetch

every required element from input SRAM to achieve two convolution elements, and thereby one-half of the max pooled value. Therefore, for every complete max pooled output, we spend 12 cycles on average. This figure varies when we access edge values and because of state transitions during different input fetches.

**3. Interface specification**

The neural network accelerator is designed to receive data from input SRAM, which gives one data block containing two data values per address, for every cycle. At the end of computations, the hardware is designed to write to output SRAM one data line at a time, every clock cycle.

Optionally given scratchpad SRAM isn't utilized to save area. Weights SRAM is accessed sequentially at the beginning of input fetch state to get kernel values. The kernel values are then stored in a register array to be accessed in parallel over multiple convolution stages.

The control signals:

All control signals are 1 bit wide.

*output reg dut_busy* is set to true when the hardware accelerator accepts dut_run signal and starts data fetching.

*input wire dut_run* is received as input trigger to start the accelerator.

*input wire reset_b* active low reset. When signal is low, we set all controls to initial state.

*input wire clk* gives the clock signal for sequential logic.

SRAM interface

Input, output and weights SRAM have similar interfaces.

*sram_write_enable* is set to high when there is enough data to be written to output SRAM. This is a status signal and so it is one bit wide.

*sram_write_addresss* is set to write address when there is data to be written to output SRAM. This is 12 bits wide.

*sram_write_data* is set to 16-bit data packet containing two data values that is to be written to address mentioned in sram_write_address. It is 16 bits wide.

*sram_read_address* is used to set 12-bit address when trying to fetch data from SRAM.

*input_sram_read_data* is a 16-bit wide reg used to read a line of data addressed using sram_read_address. Each data line read contains 2 data blocks.
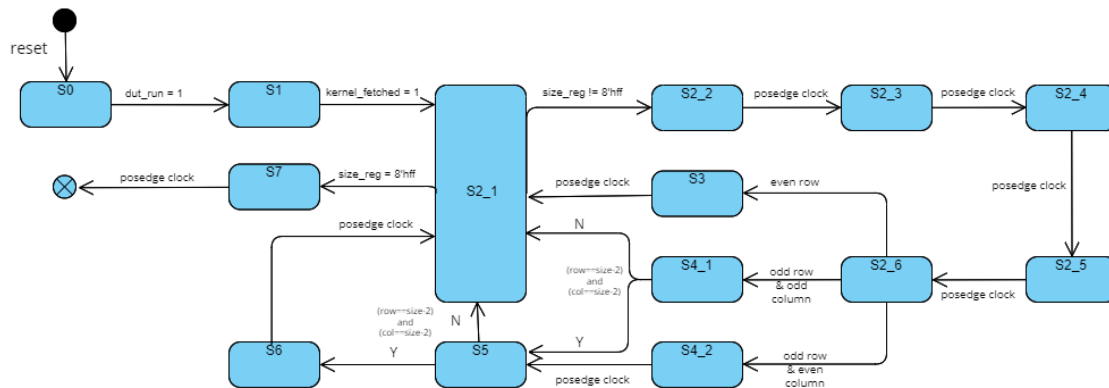
## 4. Technical Implementation:



Fig 2: Finite State Machine Diagram

An FSM with 15 states is designed keeping in mind each state must last for minimum number of clock cycle (to make control easier). The hardware design is split into controller and datapath.
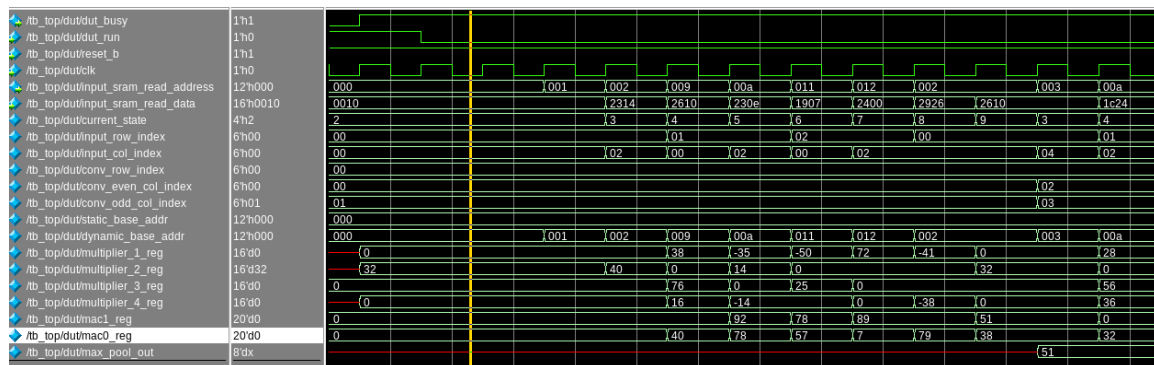
## Timing Diagram:



Fig 3: Timing diagram showing calculation of one element of max-pooling layer

## 5. Verification:

The make verify-564-base is run on Grendel to verify the correctness of the entire circuit.



Fig 4: make verify-564-base command results to verify output values using modelsim

To verify circuit output in parts, the intermediate timing diagram is analyzed and reg+wire outputs are compared against the intermediate values obtained by running.a python script that gives outputs after each stage of the hardware.



Fig 5: Python script generating outputs that should be present after each stage

**6. Results Achieved and Conclusion:**

A deep neural network accelerator with convolution layer and max-pooled layer with ReLu activation function has been developed using parallel utilization of data retrieved from SRAMs. The number of clock cycles to complete calculation of sample input of matrix sizes ranging from taken 8x8, 16x16, 32x32 and 64x64 elements is found to be 5680 cycles with a clock period of 5ns. This gives a delay of 28400 ns for the entire hardware accelerator to run the given input. This result was achieved on a silicon base with a cell area of 12070.3 $um^2$ that gives a performance/area parameter of 2.917 x 10-9 $ns^{-1}.um^{-2}$.

Detailed analysis of timing diagram revealed improvements can be made in reducing number of adders and multipliers that were idle in certain states, which could've been handled different using resource sharing. The resultant algorithm might have been more complex and thus increased the total clock cycles taken to achieve the results, but would've given an improvement in power and area.