

Report on SystemVerilog Implementation of One Layer Model

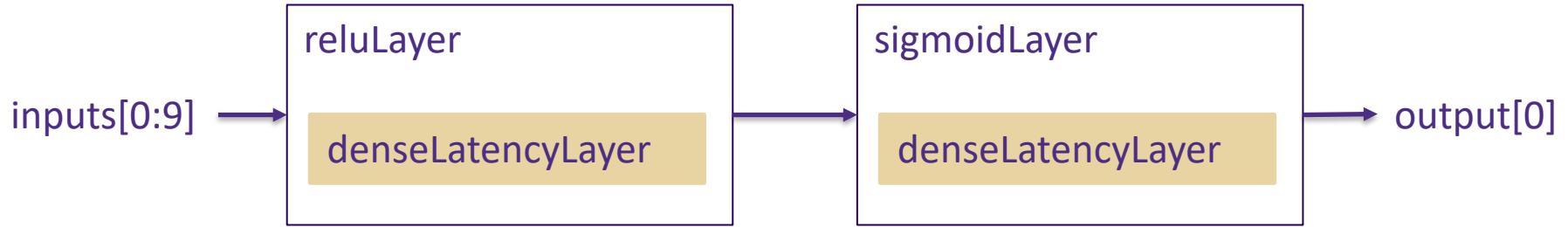


Scott Hauck, Geoffrey Jones, Anatoliy Martynyuk, Matthew Bavier,
Caroline Johnson , Oleh Kondratyuk, Trinh Nguyen, Aidan Short, Jan Silva

UNIVERSITY *of* WASHINGTON



Data Flow of SystemVerilog (SV) Implementation



- > Inputs are obtained from HLS and processed through Python script
- > Every module is pipelined (`denseLatencyLayer` , `reluLayer`, `sigmoidLayer`, `One_Layer_NN`)
- > `denseLatencyLayer` uses an adder tree

denseLatencyLayer

- > **Dense layer**
 - Each neuron in this layer receives input from all neurons from the previous layer
- > **Module computes the following: $\text{dot}(\text{input}, \text{weight}) + \text{bias}$**
 - Via matrix multiplication
 - Use adder tree to add results from matrix multiplication in order to calculate dot product
- > **This module is used by both reluLayer and sigmoidLayer**



reluLayer



- > **output = activation(dot(input, weight) + bias)**
 - Dot product and bias computation is done within denseLatencyLayer
 - reluLayer takes the output of denseLatencyLayer and performs activation function
- > **Input size: array of size 10**
- > **Output size: array of size 32**

sigmoidLayer

- > **output = activation(dot(input, weight) + bias)**
 - Dot product and bias computation is done within denseLatencyLayer
 - sigmoidLayer takes the output of denseLatencyLayer and use a lookup table to perform the activation function
- > **Input size: array of size 32**
- > **Output size: a single number**

Timing Comparison Between HLS and SV Implementation

HLS

- > 3.989 ns period = ~250 MHz

Clock	Target	Estimated	Uncertainty
ap_clk	5.00 ns	3.989 ns	1.35 ns

SV

- > Through pipelining and implementing adder tree instead of simple +, we were able to increase the clock speed
- > Up to 500 MHz

Timing of SV Implementation

Design Timing Summary

Setup

Worst Negative Slack (WNS): 0.071 ns
Total Negative Slack (TNS): 0.000 ns
Number of Failing Endpoints: 0
Total Number of Endpoints: 15214

Hold

Worst Hold Slack (WHS): 0.057 ns
Total Hold Slack (THS): 0.000 ns
Number of Failing Endpoints: 0
Total Number of Endpoints: 15214

Pulse Width

Worst Pulse Width Slack (WPWS): 0.161 ns
Total Pulse Width Negative Slack (TPWS): 0.000 ns
Number of Failing Endpoints: 0
Total Number of Endpoints: 15052

All user specified timing constraints are met.

Name	Waveform	Period (ns)	Frequency (MHz)
clk_p	{0.000 3.200}	6.400	156.250
clk_out1_clk_wiz_0	{0.000 1.000}	2.000	500.000
clkfbout_clk_wiz_0	{0.000 16.000}	32.000	31.250

Latency SV Implementation



- > Latency is roughly 56ns
- > 28 clock cycles at 500MHz
- > 8.75 clock cycles at 156.25MHz input clock

Utilization Comparison Between HLS and SV

	HLS	SV
BRAM	1	0.5
DSP	70	352
FF	13260	14651
LUT	15728	5972

* Each value is the total used in the One_Layer_NN

Utilization of HLS

Name	BRAM_18K	DSP	FF	LUT	URAM
DSP	-	-	-	-	-
Expression	-	-	0	12611	-
FIFO	-	-	-	-	-
Instance	-	70	3150	630	-
Memory	1	-	0	0	-
Multiplexer	-	-	-	2487	-
Register	-	-	10110	-	-
Total	1	70	13260	15728	0
Available SLR	2160	2760	663360	331680	0
Utilization SLR (%)	~0	2	1	4	100
Available	4320	5520	1326720	663360	0
Utilization (%)	~0	1	~0	2	0

Utilization of SV

Name	Slice LUTs (433200)	Slice Registers (866400)	Slice (108300)	LUT as Logic (433200)	LUT as Memory (174200)	Block RAM Tile (1470)	DSPs (3600)	Bonded IOB (850)	IBUFDS (816)	BUFGCTRL (32)	PLLE2_ADV (20)
One_Layer_NN	5972	14651	3662	5939	33	0.5	352	190	1	2	1
> dout (DFF_2D_parameterized1)	0	17	5	0	0	0	0	0	0	0	0
> PLL (clk_wiz_0)	0	0	0	0	0	0	0	0	1	2	1
> relu (reluLayer)	5420	13510	3382	5391	29	0	320	0	0	0	0
> sig_layer (sigmoidLayer)	552	1124	339	548	4	0.5	32	0	0	0	0

Resource	Utilization	Available	Utilization %
LUT	5972	433200	1.38
LUTRAM	33	174200	0.02
FF	14651	866400	1.69
BRAM	0.50	1470	0.03
DSP	352	3600	9.78
IO	190	850	22.35
PLL	1	20	5.00

Extra: Timing of SV Implementation w/out using BRAM

Design Timing Summary

Setup

Worst Negative Slack (WNS): 0.084 ns
Total Negative Slack (TNS): 0.000 ns
Number of Failing Endpoints: 0
Total Number of Endpoints: 16135

Hold

Worst Hold Slack (WHS): 0.054 ns
Total Hold Slack (THS): 0.000 ns
Number of Failing Endpoints: 0
Total Number of Endpoints: 16135

Pulse Width

Worst Pulse Width Slack (WPWS): 0.390 ns
Total Pulse Width Negative Slack (TPWS): 0.000 ns
Number of Failing Endpoints: 0
Total Number of Endpoints: 15982

All user specified timing constraints are met.

Name	Waveform	Period (ns)	Frequency (MHz)
clk_p	{0.000 3.200}	6.400	156.250
clk_out1_clk_wiz_0	{0.000 1.032}	2.065	484.375
clkfbout_clk_wiz_0	{0.000 16.000}	32.000	31.250