CNN Accelerator for Gesture Recognition on DE1-SoC FPGA

CSEE 4840 Embedded System

Presenters: Yangfan Wang (yw4415) Fengze Zhong (fz2393) Xincheng Yu (xy2654)

May 14, 2025



Project Introduction

Target: Use FPGA to accelerate the inference process of a CNN model



Figure 1: 10 Gestures for recognition



The CNN contains 4 convolutional layers, 3 max pooling layers and 2 fully connected layers



Figure 2: CNN Architecture





Figure 3: CNN Architecture (detailed)



Golden Model

Dataset:

Sign Language Digits Dataset, which contains 2,180 color images of hand gestures representing digits from 0 to 9.

Platform:

PyTorch

Data Augmentation Techniques:

Random rotation, Center cropping, and

Adjustments to brightness and contrast

Data Split:

80% for training and 20% for validation.

Accuracy:

93% on the original dataset







System Architecture

The system we build contains the software control and the hardware acceleration.

 The hardware side will be purely for storage and computation including the kernel convolution, max pooling, ReLU activation and fully connected layer.



Figure 5: FPGA and HPS Block Diagram



Software

- With an image pending for recognition, the software will firstly downsample, resize it to 32x32 and do zero padding.
- Then it will follow the structure of the CNN, load weights and biases from files to the hardware according to the current layer.



Figure 6: Example of 32x32 matrix zero padding at the input



Software

- compute_network():
 - Main function that runs the network
- run_conv():
 - Set convolution layer parameters
 - Allocate memory for output
- run_fc():
 - Set fully connected layer parameters
 - Allocate memory for output
- memory_ioctl()
 - CONV_WRITE & FC_WRITE:
 - copy_from_user: configuration
 - copy_from_user: intput_data, weight, bias
 - copy_to_user: output_data



Hardware Top Level







```
conv_output_wr_addr:
curr_output_y * (num_cols - 'd2) + curr_output_x + curr_dst_chan * (num_cols - 'd2) * (num_rows - 'd2)
```

```
curr_pixel_pos + num_cols * num_rows * curr_src_chan
```

```
conv_data_rd_addr:
```

```
9 * curr_src_chan + 9 * num_src_chans * curr_dst_chan + curr_kernel_id
```

```
conv_weight_rd_addr:
```

 tb_iop/cnn_interface_inst/conv_inst/mem_rd tb_iop/cnn_interface_inst/conv_inst/mem_wr_data_o tb_iop/cnn_interface_inst/conv_inst/mem_wr_addr_o tb_iop/cnn_interface_inst/conv_inst/mem_wr 	1'h0 16'hxxxx 19'h00010 1'h1	0000 Obeef										D00 beef						00
/tb_top/cnn_interface_inst/conv_inst/start	1'h0																	
#	16'h0001							-AV-				- D-4					00	01
	16'h0001																00	01
# /tb_top/cnn_interface_inst/conv_inst/curr_output_x	16'h0000						1				100	-	-		100		00	00
	16'h0000						(in										00	00
Hotop/cnn_interface_inst/conv_inst/curr_src_chan	16'd0									111111			IIIII					
🖅 🎝 /tb_top/cnn_interface_inst/conv_inst/curr_dst_chan	16'd0	<u>1 12 13</u>	14 15 X	5 X7 X8	(9) 10) :	11 <u>12</u>)	13 (14)	15 16	<u>17 18 18 18 18 18 18 18 </u>	19 20	21 22	23 24	125 12	5 (27)	28 29	(30 (3	1 0	
🔩 /tb top/cnn interface inst/conv inst/valid	1'h0																	

1 Obeef

Initial Weight Load

Hardware Top Level

Convolution

/tb top/cnn interface inst/conv inst/mem rd addr o

Hardware Top Level

Fully Connected

+	🗇 /tb_top/cnn_interface_inst/conv_wr_data	16'h0000	0000	1 20000	<u>i 10000</u>	1 10000
1	/tb_top/cnn_interface_inst/conv_wr	1'h0			ومعمد المعتمان والمعط	
+	/tb_top/cnn_interface_inst/conv_wr_addr	19'hObeef	Obeef	l lobeet	, jObeef	1 10beef
+	/tb_top/cnn_interface_inst/conv_rd_data	16'h0060	0060			
	/tb_top/cnn_interface_inst/conv_rd	1'h0				
+	/tb_top/cnn_interface_inst/conv_rd_addr	19'h0beef	Obeef	X X X X 0beef X X X X X X X X X X X	I (Objeef) I I I I I I I I I	, j j0beef
	/tb_top/cnn_interface_inst/valid	1'h1				

fc_weight_rd_addr:

```
curr_fc_w_row * num_dst_chans + curr_fc_w_col
```

fc_data_rd_addr: curr_fc_w_row

fc_wr_addr: curr_fc_w_col



Address Computation

Pixel Position:

curr_kernel_pos - num_cols - 'd1
curr_kernel_pos - num_cols
curr_kernel_pos - num_cols + 'd1
curr_kernel_pos - 'd1
curr_kernel_pos + 'd1
curr_kernel_pos + 'd1
curr_kernel_pos + num_cols - 'd1
curr_kernel_pos + num_cols
curr_kernel_pos + num_cols + 'd1;



Address Computation

pool_rd_addr:

```
(curr_pool_size * pool_size) + curr_pool_stride_y) * (num_cols - 'd2) + (curr_pool_x * pool_size + curr_pool_stride_x + curr_dst_chan * (num_cols - 'd2) * (num_rows - 'd2)
```

pool_wr_addr:

```
(curr_pool_y) * (num_cols - 'd2) / pool_size + (curr_pool_x)
```



Hardware Component



Figure 8: MAC Unit



Programmable Registers

0	num_src_chans	8	data_starting_addr[15:0]
1	num dst chans	9	data_starting_addr[31:16]
2	num cols	10	weight_starting_addr[15:0]
3	num rows	11	weight_starting_addr[31:16]
4	fc _	12	bias_starting_addr[15:0]
5	do pool	13	bias_starting_addr[31:16]
6	pool size	14	output_starting_addr[15:0]
7	pool stride	15	output_starting_addr[31:16]
		16	start

Read to any address will receive the LAYER_DONE signal from HW

This is how software tells hardware what to do and receive feedback from hardware



HW/SW Interface

Memory Space

Data
Weight
Bias
Output

Software does the zero padding for each layer



Verification

- Implemented Python Codes
 - Fully Connected Layer
 - $\circ \quad \text{Convolution} \quad$
 - Pooling

Generating Tests Case and Expected Output



• The main constraint for this project is the BRAM needed to store the source data, weight, bias and output. All the data involved will be represented in Q8.8 fixed point using 2 bytes.

72 KB
146 KB
68 KB
305 KB
33 KB
1.4 KB

Table 1. memory requirement of each layer

- BRAM size must be larger than the maximum layer memory requirement
- BRAM is set to 340 KB



+	
; Flow Summary	
+	.+
; Flow Status	; Successful - Wed May 14 14:40:42 2025
; Quartus Prime Version	; 21.1.0 Build 842 10/21/2021 SJ Lite Edition
; Revision Name	; soc_system
; Top-level Entity Name	; soc_system_top
; Family	; Cyclone V
; Device	; 5CSEMA5F31C6
; Timing Models	; Final
; Logic utilization (in ALMs)	; 1,498 / 32,070 (5 %)
; Total registers	; 1301
; Total pins	; 362 / 457 (79 %)
; Total virtual pins	; 0
; Total block memory bits	; 2,785,424 / 4,065,280 (69 %)
; Total DSP Blocks	; 21 / 87 (24 %)
; Total HSSI RX PCSs	; 0
; Total HSSI PMA RX Deserializers	; 0
; Total HSSI TX PCSs	; 0
; Total HSSI PMA TX Serializers	; 0
; Total PLLs	; 0 / 6 (0 %)
; Total DLLs	; 1 / 4 (25 %)
+	+

Fmax ~ 45MHz



Thanks for Listening

Any Questions?



May 14, 2025