

# **CNN Plus Accelerator IP**

IP Version: v2.0.0

# **User Guide**

FPGA-IPUG-02115-1.6

December 2024



#### **Disclaimers**

Lattice makes no warranty, representation, or guarantee regarding the accuracy of information contained in this document or the suitability of its products for any particular purpose. All information herein is provided AS IS, with all faults, and all associated risk is the responsibility entirely of the Buyer. The information provided herein is for informational purposes only and may contain technical inaccuracies or omissions, and may be otherwise rendered inaccurate for many reasons, and Lattice assumes no obligation to update or otherwise correct or revise this information. Products sold by Lattice have been subject to limited testing and it is the Buyer's responsibility to independently determine the suitability of any products and to test and verify the same. LATTICE PRODUCTS AND SERVICES ARE NOT DESIGNED, MANUFACTURED, OR TESTED FOR USE IN LIFE OR SAFETY CRITICAL SYSTEMS, HAZARDOUS ENVIRONMENTS, OR ANY OTHER ENVIRONMENTS REQUIRING FAIL-SAFE PERFORMANCE, INCLUDING ANY APPLICATION IN WHICH THE FAILURE OF THE PRODUCT OR SERVICE COULD LEAD TO DEATH, PERSONAL INJURY, SEVERE PROPERTY DAMAGE OR ENVIRONMENTAL HARM (COLLECTIVELY, "HIGH-RISK USES"). FURTHER, BUYER MUST TAKE PRUDENT STEPS TO PROTECT AGAINST PRODUCT AND SERVICE FAILURES, INCLUDING PROVIDING APPROPRIATE REDUNDANCIES, FAIL-SAFE FEATURES, AND/OR SHUT-DOWN MECHANISMS. LATTICE EXPRESSLY DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY OF FITNESS OF THE PRODUCTS OR SERVICES FOR HIGH-RISK USES. The information provided in this document is proprietary to Lattice Semiconductor, and Lattice reserves the right to make any changes to the information in this document or to any products at any time without notice.

#### **Inclusive Language**

This document was created consistent with Lattice Semiconductor's inclusive language policy. In some cases, the language in underlying tools and other items may not yet have been updated. Please refer to Lattice's inclusive language FAQ 6878 for a cross reference of terms. Note in some cases such as register names and state names it has been necessary to continue to utilize older terminology for compatibility.



# **Contents**

Contents	3
Acronyms in This Document	5
1. Introduction	6
1.1. Quick Facts	6
1.2. Features	6
2. Functional Descriptions	8
2.1. Overview	8
2.2. Interface Descriptions	11
2.2.1. Control and Status Interface	14
2.2.2. Input Data Interface	15
2.2.3. Result and Debug Interface	15
2.2.4. DRAM Interface	16
2.2.5. Command FIFO Interface	16
2.3. Attribute Summary	17
2.4. Clock Domain	21
2.5. Reset Behavior	21
2.6. Register Description	22
2.7. Operation Sequence	22
2.7.1. Command Format	22
2.7.2. Input Data Format	22
2.7.3. Output Data Format	23
2.8. Supported Commands	23
3. IP Core Generation, Simulation, and Validation	24
3.1. Licensing the IP	24
3.2. Generation and Synthesis	24
3.3. Running Functional Simulation	27
3.4. Hardware Evaluation	27
4. Hardware Validation	28
5. Ordering Part Number	29
Appendix A. Resource Utilization	30
References	31
Technical Support Assistance	32
Revision History	33



# **Figures**

Figure 2.1. Functional Block Diagram of CNN Plus Accelerator (Compact CNN Type)	8
Figure 2.2. Order of Layers for CNN Plus Accelerator	
Figure 2.3. Functional Block Diagram of CNN Plus Accelerator (Optimized CNN Type)	9
Figure 2.4. Functional Block Diagram of CNN Plus Accelerator (Extended CNN Type)	
Figure 2.5. Additional Order of Layers for CNN Plus Accelerator (Extended CNN Type)	
Figure 2.6. CNN Plus Accelerator IP Core Interface Diagram	11
Figure 2.7. Control and Status Interface Timing Diagram	14
Figure 2.8. General Purpose Output Sample Application	15
Figure 2.9. Result and Debug Interface Timing Diagram	16
Figure 2.10. Command FIFO Interface Timing Diagram	16
Figure 2.11. Clock Domain Diagram	21
Figure 2.12. Reset Timing Diagram	21
Figure 2.13. Command Format	22
Figure 3.1. Module/IP Block Wizard	24
Figure 3.2. Configure User Interface of CNN Plus Accelerator IP Core	25
Figure 3.3. Check Generating Result	26
Tables	
Table 1.1. Quick Facts	6
Table 2.1. CNN Plus Accelerator IP Core Signal Descriptions	11
Table 2.2. Attributes for EngineConfig Tab	
Table 2.3. Attributes for LUT-Activation-Config Tab (Scale-Module)	
Table 2.4. Attributes Table for LUT-Activation-Config Tab (FC-Module)	
Table 2.5. Attributes Descriptions	18
Table 2.6. Shifting of Data Byte Based on Byte Shift and Byte Mode	
Table 2.7. Supported Devices and Attributes	
Table 3.1. Generated File List	
Table A.1. Performance and Resource Utilization <sup>1</sup>	30

5



# **Acronyms in This Document**

A list of acronyms used in this document.

Acronym	Definition
CNN	Convolutional Neural Network
DRAM	Dynamic Random Access Memory
FC	Fully Connected
FIFO	First In, First Out
FPGA	Field-Programmable Gate Array
LMMI	Lattice Memory Mapped Interface
LSQ	Learned Step Size
LUT	Look-Up Table
MSB	Most Significant Bit
OPN	Ordering Part Number



## 1. Introduction

The Lattice Semiconductor CNN Plus Accelerator IP Core is a calculation engine for Deep Neural Network with fixed point weight. It calculates full layers of Neural Network including convolution layer, pooling layer, batch normalization layer, and full connect layer by executing sequence code with weight value, which is generated by Lattice sensAl™ Neural Network Compiler. The engine is optimized for convolutional neural network, so it can be used for vision-based application such as classification or object detection and tracking. The IP Core does not require an extra processor; it can perform all required calculations by itself.

The design is implemented in Verilog HDL. It can be targeted to CrossLink™-NX, Certus™-NX, and CertusPro™-NX FPGA devices implemented using the Lattice Radiant™ software Place and Route tool integrated with the Synplify Pro® synthesis tool.

## 1.1. Quick Facts

Table 1.1 presents a summary of the CNN Plus Accelerator IP Core.

**Table 1.1. Quick Facts** 

IP Requirements	Supported FPGA Family	CrossLink-NX, Certus-NX, CertusPro-NX		
	Targeted Devices	LIFCL-40, LIFCL-33, LIFCL-33U, LIFCL-17, LFD2NX-40, LFD2NX-17, LFCPNX-50, LFCPNX-100		
Resource Utilization	Supported User Interface	AXI4, LMMI (Lattice Memory Mapped Interface) Native interfaces as described in Interface Descriptions section.		
	Resources	See Appendix A. Resource Utilization.		
Design Tool Support	Lattice Implementation	IP Core v1.0.x – Lattice Radiant Software 2.0 and 2.1 IP Core v1.1.x – Lattice Radiant Software 2.0 or later IP Core v1.2.x – Lattice Radiant Software 2.0 or later IP Core v1.4.x – Lattice Radiant Software 2022.1 or later IP Core v2.0.0 – Lattice Radiant Software 2023.1 or later		
	Synthesis	Lattice Synthesis Engine Synopsys Synplify® Pro for Lattice		
	Simulation	For a list of supported simulators, see the Lattice Radiant software user guide.		

## 1.2. Features

The key features of the CNN Plus Accelerator IP Core include:

- Selectable three implementation types:
  - Compact CNN
  - Optimized CNN
  - Extended CNN
- Selectable AXI4 or FIFO interface
- Support for convolution layer, max pooling layer, global average pooling layer, batch normalization layer, and full connect layer
- Configurable bit width of activation (16/8-bit)
- Configurable number of memory blocks for tradeoff between resource and performance
- Optimized for 3 × 3 2D convolution calculation
- Support for 2 × 2 max pooling
- Support for 2 × 2 max pooling/unpooling with max argument Extended CNN only
- Support for global average pooling by full connect engine
- Configurable input byte mode: signed, unsigned and disable only Compact CNN supports disable mode
- Partial DRAM access



- Configurable maximum burst length (32, 256) for AXI4 Interface
- Supports paired and quadruple convolution engines to improve performance Optimized/Extended CNN only
- Configurable Line Buffer Size (512, 1024, 2048) Optimized/Extended CNN only
- · Supports general purpose output signal for controlling external logic through command code
- Supports higher precision convolution calculation to improve overall result accuracy Optimized/Extended CNN only
- Provides option for additional LUT reduction with configurable reduction factor using DSP slices Optimized CNN only
- Added the Focus layer support Optimized CNN only
- Added the LSQ quantization support for the Scale an FC layer Optimized CNN only
- Supports external LUT file for the Scale and FC layer Optimized CNN only
- Supports 3 modes of FC layer which are Normal, 4BIT4CH, and 4BIT8CH Optimized CNN only
- Supports Stride2 for 3 × 3 convolution Optimized CNN only
- Supports 4-bit quantized activations (with 8 bit weights) mode for convolution in Optimized CNN to save internal memory and facilitate fitting bigger models in embedded mode.



## 2. Functional Descriptions

#### 2.1. Overview

8

The CNN Plus Accelerator IP Core performs a series of calculations per command sequence that is generated by the Lattice Neural Network Compiler tool. Commands must be written at the DRAM address specified by the code\_base\_addr\_i signal, which is accessible through AXI BUS. Input data may be read from the DRAM at a pre-defined address or directly written through the input data write port. After command code and input data are available, the CNN Plus Accelerator IP Core starts calculation at the rising edge of the start signal. During calculation, intermediate data and final result may be transferred to the DRAM or fed out through the result write port. All operations are fully programmable by command code.

The CNN Plus Accelerator IP Core offers three types of implementations: the compact CNN type, which is suitable for small FPGA devices due to its low utilization, the optimized CNN type, which can perform four convolution calculations in parallel, making it suitable for high-speed applications, and the extended CNN type, which offers the same features as optimized CNN plus additional support for max pooling/unpooling with max argument. The block diagram of each type of implementation is shown in the following figures: Figure 2.1 for compact CNN type, Figure 2.3 for optimized CNN type, and Figure 2.4 for extended CNN type. The sizes of the Activation Data Storage and the Conv. Scratch Storage are configurable by using the Memory Type and Scratch Pad Memory Size attributes respectively. Refer to Table 2.5 for the description of these attributes.

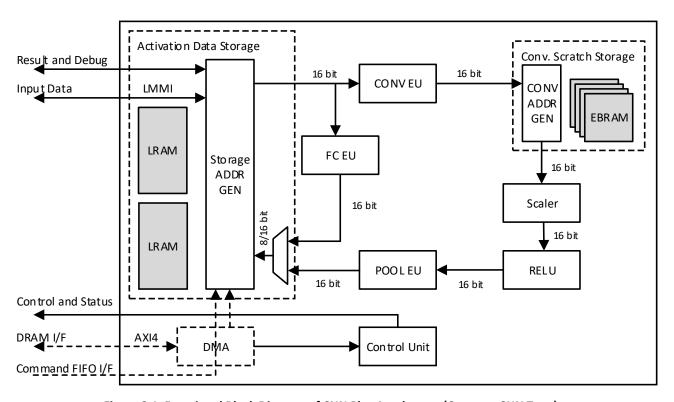


Figure 2.1. Functional Block Diagram of CNN Plus Accelerator (Compact CNN Type)

In Figure 2.1, the input data is stored in Activation Data Storage through LMMI. It can feed data to either the Convolution Engine or the Full Connect Engine. The result of Full Connect Engine goes back to Activation Data Storage. For the other path, the result of Convolution Engine is buffered to Convolution Scratch Storage and is transmitted to Scaler, which performs batch normalization. The Scaler can be bypassed by setting unity value, that is multiply by 1 and add by 0.

That result goes to RELU and then to POOL EU. Both RELU and POOL EU has bypass option. Lastly, the result of POOL EU goes back to Activation Data Storage.

FPGA-IPLIG-02115-1 6

9



The order of layers for CNN configuration is shown in Figure 2.2, which is similar for both Compact CNN and Optimized CNN implementation type. The only difference is that the Optimized CNN as shown in Figure 2.3 can perform four parallel calculations efficiently. Follow this when designing your network model. The layers inside the parenthesis have bypass option.

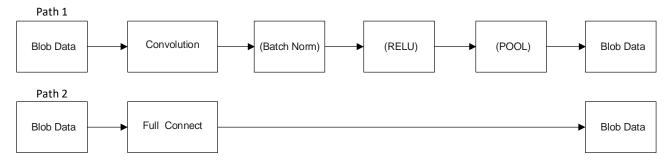


Figure 2.2. Order of Layers for CNN Plus Accelerator

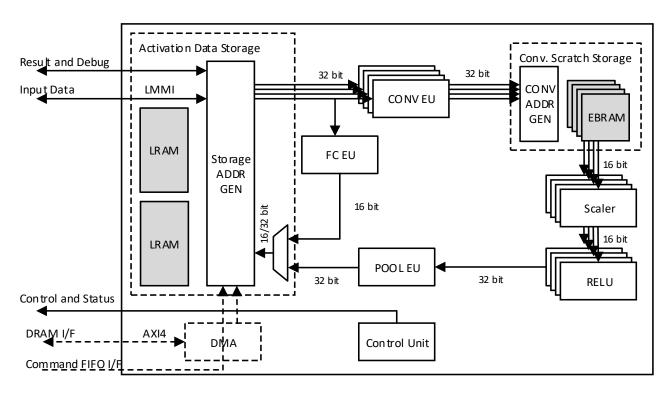


Figure 2.3. Functional Block Diagram of CNN Plus Accelerator (Optimized CNN Type)

In Figure 2.4, two additional computing functions are present in Extended CNN type. First, the max pooling/unpooling that processes the input data in reference to the input max argument. Second, the addition of RELU result and the input data. These new functions also provide additional order of layers as shown in Figure 2.5.



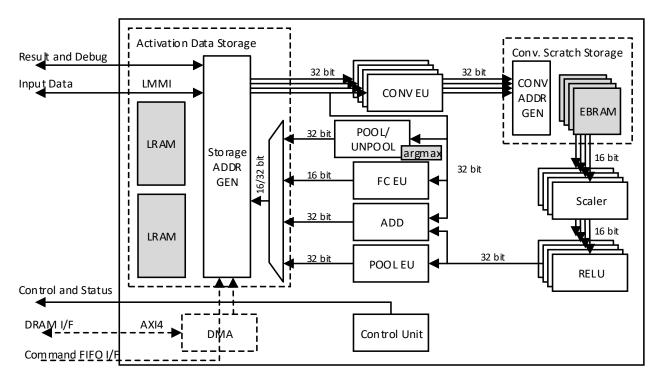


Figure 2.4. Functional Block Diagram of CNN Plus Accelerator (Extended CNN Type)

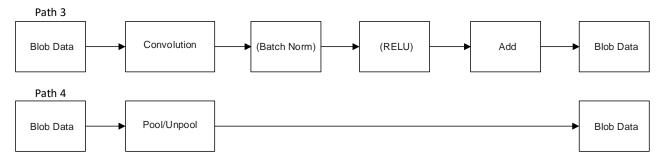


Figure 2.5. Additional Order of Layers for CNN Plus Accelerator (Extended CNN Type)



## 2.2. Interface Descriptions

Figure 2.6. shows the interface diagram for the CNN Plus Accelerator IP Core. The diagram shows all of the available ports for the IP core.

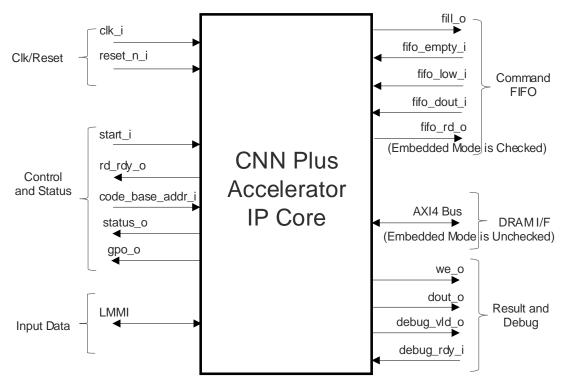


Figure 2.6. CNN Plus Accelerator IP Core Interface Diagram

Table 2.1. CNN Plus Accelerator IP Core Signal Descriptions

Pin Name Direction Function Description				
	Clock/Reset			
clk_i				
reset_n_i	Input	Active low system reset that is synchronous to clk signal and is asynchronous to aclk signal.  [0] – Resets all ports and sets internal registers to their default values.  [1] – Reset is NOT active		
		Control and Status		
code_base_addr_i[31:0]	Input This signal specifies the base/start address that is read by CNN Plus Accelerator IP ( to get the command code. External logic should write the command code to this address. This signal must be set before start of operation.			
start_i	Input	ut Start execution signal. Level sensitive. Must deassert after rd_rdy_o going 0.		
rd_rdy_o  Output  Output  Output  Output  Ready signal  [0] - Engine is busy/running.  [1] - Engine is idle and ready to get input.  External logic should write input data to internal memory only during rd_rdy_		[0] – Engine is busy/running.		
status_o	Output	Debug information  [0] – If bit value is 1, engines are running except full connect engine.  [1] – If bit value is 1, full connect engine is running.  [2] – If bit value is 1, AXI bus transfer is ongoing.  [3] – If bit value is 1, during engine is running.		



Pin Name	Direction	Function Description	
		[4] – If bit value is 1, command FIFO is reading.	
		[5] – If bit value is 1, during DRAM access.	
		[6] – If bit value is 1, DRAM command is issued.	
		[7] – If bit value is 1, engine is waiting for FIFO fill.	
[24.0]	0	General Purpose Output signal. Use for communication from firmware to outside block	
gpo_o[31:0]	Output	such as informing the type of current output, or informing the firmware version.	
		Input Data (LMMI)	
lmmi_request_i	Input	Start transaction	
lmmi_wr_rdn_i	Input	Write = HIGH, Read = LOW	
Immi_offset_i	Input	Address offset for accessing the internal LRAMs in Activation Data Storage.	
lmmi_wdata_i	Input	Write data	
lmmi_ready_o	Output	Target Ready signal. When Immi_ready_o and Immi_request_i are both HIGH, it indicates write transaction is completed and address phase of ready transition is completed.	
lmmi_rdata_valid_o	Output	Indicates read transaction is complete and Immi_rdata_o contains valid data.	
lmmi_rdata_o	Output	Read data	
	-	Result and Debug	
		Write enable of result, indicates result data is valid	
we_o	Output	[0] – Result data is NOT valid	
		[1] – Result data is valid	
dout_o[15:0]	Output	IF we_o is asserted, dout_o[15:0] contains result data.	
uout_0[13.0]	Output	IF debug_vld_o is asserted, dout_o[15:0] contains debug data.	
		Ready signal for one burst read of debug data	
		0 – External logic is not ready to receive 1 burst of debug data.	
debug_rdy_i	Input	1 – External logic is ready to receive 1 burst of debug data.	
		The recommendation is to connect to 1 if debug feature is not used so that the	
		operation is not halted when the IP Core sends debug data.	
		Indicates that dout_o[15:0] is not result data, but debug data.	
debug_vld_o	Output	0 – Debug data is NOT valid.	
		1 – Debug data is valid	
aclk	Input	DRAM I/F (AXI4 Bus)¹  AXI4 clock signal. Fully asynchronous from clk. Recommend to use DRAM system clock.	
	<del></del>		
axi4_awid_o[7:0]	Output	AXI4 write address channel, write address ID signal  Constant output: 0x10	
avi4 awaddr o[21:0]	Output	AXI4 write address channel, write address signal	
axi4_awaddr_o[31:0]	<u> </u>	AXI4 write address channel, write address signal  AXI4 write address channel, region identifier signal	
axi4_awregion_o[3:0]	Output		
	Outrout	Constant output: 0x0 (default)	
axi4_awlen_o[7:0]	Output	AXI4 write address channel, burst length signal	
axi4_awsize_o[2:0]	Output	AXI4 write address channel, burst size signal  Constant output: 3'b011 (8 bytes per beat)	
:4bet -[1.0]	Outro		
axi4_awburst_o[1:0]	Output	AXI4 write address channel, burst type signal	
	0	Constant output: 2'b01 (INCR)	
axi4_awlock_o	Output	AXI4 write address channel, lock type signal	
ovid owes sha - [2:0]	O + +	Constant output: 1'b0 (Normal Access)	
axi4_awcache_o[3:0]	Output	AXI4 write address channel, memory type signal	
and annual -[2 0]	0	Constant output: 0x00 (Device Non-bufferable)	
axi4_awprot_o[2:0]	Output	AXI4 write address channel, protection type signal	
	0	Constant output: 3'b010 (Data, Non-Secure Access, Unpriviledged)	
axi4_awqos_o[3:0]	Output	AXI4 write address channel, quality of service signal	
	0	Constant output: 0x0 (no QoS scheme)	
axi4_awvalid_o	Output	AXI4 write address channel, write address valid signal	



Pin Name	Direction	Function Description		
axi4 awready i	Input	AXI4 write address channel, write address ready signal		
axi4_wid_o[7:0]	Output	AXI4 write data channel, write ID tag signal		
		Constant output: 0x10		
axi4_wdata_o[ 63:0]	Output	AXI4 write data channel, write data signal		
axi4_wstrb_o[ 7:0]	Output	AXI4 write data channel, write strobe signal		
		Constant output: 0xFF		
axi4_wlast_o	Output	AXI4 write data channel, write last signal		
axi4_wvalid_o	Output	AXI4 write data channel, write valid signal		
axi4_wready_i	Input	AXI4 write data channel, write ready signal		
axi4_bid_i[7:0]	Input	AXI4 write response channel, response ID tag signal		
axi4_bresp_i[1:0]	Input	AXI4 write response channel, write response signal		
axi4_bvalid_i	Input	AXI4 write response channel, write response valid signal		
axi4_bready_o	Output	AXI4 write response channel, response ready signal		
		Constant output: 1'b1 (AXI4 write response channel is ignored)		
axi4_arid_o[7:0]	Output	AXI4 read address channel, read address ID signal		
axi4_araddr_o[31:0]	Output	AXI4 read address channel, read address signal		
axi4_arregion_o[3:0]	Output	AXI4 read address channel, region identifier signal		
		Constant output: 0x0 (default)		
axi4_arlen_o[7:0]	Output	AXI4 read address channel, burst length signal		
axi4_arsize_o[2:0]	Output	AXI4 read address channel, burst size signal		
axi4_arburst_o[1:0]	Output	AXI4 read address channel, burst type signal		
axi4_arlock_o	Output	AXI4 read address channel, lock type signal		
		Constant output: 1'b0 (Normal Access)		
axi4_arcache_o[3:0]	Output	AXI4 read address channel, memory type signal		
axi4_arprot_o[2:0]	Output	AXI4 read address channel, protection type signal		
axi4_arqos_o[3:0]	Output	AXI4 read address channel, quality of service signal		
		Constant output: 0x0 (no QoS scheme)		
axi4_arvalid_o	Output	AXI4 read address channel, read address valid signal		
axi4_arready_i	Input	AXI4 read address channel, read address ready signal		
axi4_rid_i[7:0]	Input	AXI4 Read data channel, read ID tag signal		
axi4_rdata_i[ 63:0]	Input	AXI4 Read data channel, read data signal		
axi4_rresp_i[1:0]	Input	AXI4 Read data channel, read response signal		
axi4_rlast_i	Input	AXI4 Read data channel, read last signal		
axi4_rvalid_i	Input	AXI4 Read data channel, read valid signal		
axi4_rready_o	Output	AXI4 Read data channel, read ready signal		
	1	FIFO I/F <sup>2</sup>		
fill_o	Output	Request for filling the external command FIFO		
		0 – Request to flush and fill the external command FIFO External command fifo must be flushed and the command code generated by		
		Lattice sensAl Neural Network Compiler software must be loaded to external		
		command FIFO.		
		1 – Normal FIFO operation		
fifo_empty_i	Input	Indicates whether external command FIFO is empty or not.		
		0 – External FIFO is not empty		
CC I		1 – External FIFO is empty		
fifo_low_i	Input	Indicates whether external command FIFO level is low or not.  0 – FIFO level is not low		
		1 – FIFO level is not low		
		If fifo level is not low, fifo must consecutively provide one full set of full connect		
		weights, that is, by one 32-bit data per two cycles.		



Pin Name	Direction	Function Description	
fifo_dout_i[31:0]	Input	xternal command FIFO read data	
fifo_rd_o	Output	Read request to external command FIFO  0 – No Read request  1 – Read request to external FIFO is active  Value of i_fifo_dout is sampled when i_fifo_empty = 0 and o_fifo_rd = 1.	

#### Notes:

- DRAM I/F (AXI4 Bus) is available when Embedded Mode is unchecked.
- 2. FIFO I/F is available when Embedded Mode is checked.

#### 2.2.1. Control and Status Interface

After reset or when engine is idle, rd rdy o is high. During this state, external logic may perform the following:

- 1. Write input data through Input Data interface.
- 2. Set the start address of command code to code\_base\_addr signal\_i.

After the above steps, external logic must assert start\_i signal. Engine starts execution when it gets start\_i = 1 and rd\_rdy\_o goes 0. During execution, each bit of status\_o indicates activity of sub calculation engine or AXI BUS. After finishing execution, that is by getting the finish command, the CNN Plus Accelerator IP Core asserts rd\_rdy\_o and waits for the next execution. Repeat from asserting start i as shown in Figure 2.7.

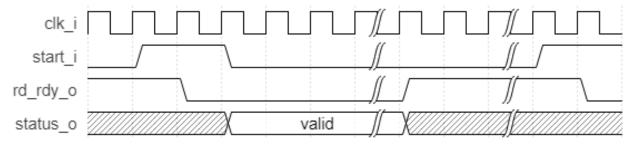


Figure 2.7. Control and Status Interface Timing Diagram

#### 2.2.1.1. General Purpose Output

The general-purpose output signal, gpo\_o is available since version v1.1.0 of the IP Core. This signal is controlled by the Lattice sensAl Neural Network Compiler software. One possible application of this signal is shown in Figure 2.8.In this example, different post processing operation needs to be performed on certain outputs. The gpo\_o signal may be used to select which post process to perform.



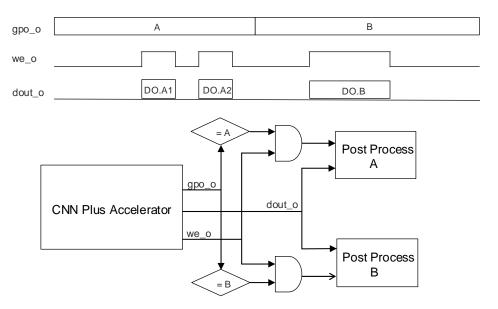


Figure 2.8. General Purpose Output Sample Application

#### 2.2.2. Input Data Interface

Input data can be written to the DRAM by external logic. In this case, loading from the DRAM for input data must be in command codes. This is done by enabling the *Store Input* option and setting the target DRAM address in the Project Implementation window of the Lattice sensAl Neural Network Compiler software. In this case, the external logic needs to store the image data to target the DRAM address.

In addition, if input data is small enough to hold in internal memory, writing to the DRAM and reading back to the CNN Plus Accelerator IP Core may be a waste of cycle time and energy. In that case, the *Store Input* option should be disabled and the external logic should write input data to the internal memory of the CNN Plus Accelerator IP Core during idle state (Immi\_ready\_o is high). The Input Data Interface is LMMI. Refer to Lattice Memory Mapped Interface and Lattice Interrupt Interface (FPGA-UG-02039) for more information on LMMI and its timing diagram. The read latency from Immi\_request\_i to Immi\_ready\_o is three clock cycles. Memory ID and memory address are mapped to the Immi\_offset\_i as described in Table 2.1. They should be matched to command code. There is no required order or rule for writing the input data to the internal memory; any random access including memory ID is acceptable. Overwriting of same address is also accepted. However, you should ensure that the input data format is satisfied before asserting start i signal. Refer to Input Data Format section for details.

#### 2.2.3. Result and Debug Interface

The result, that is, final Blob data of neural network and the debug information are outputted to external logic through this Result and Debug Interface. When we\_o signal is asserted, the dout\_o signal contains the valid final Blob data. The external logic must have enough buffer to sample the result data whenever we\_o asserts. This is usually is a single burst series of 16-bit data based on the command code. On the other hand, when debug\_vld\_o is asserted, the dout\_o signal contains a vaid debug information. The IP Core keeps transmitting available debug information when debug\_rdy\_i is asserted. The IP Core halts until all debug information are transmitted. Thus, when debug information is not used, you should set debug\_rdy\_i to a fixed 1'b1 value so that it can flush-out debug information to proceed operation. The we\_o and debug\_vld\_o signals never assert at the same time. Figure 2.9 shows an example timing diagram of Result and Debug interface with debug\_rdy\_i fixed to 1'b1.



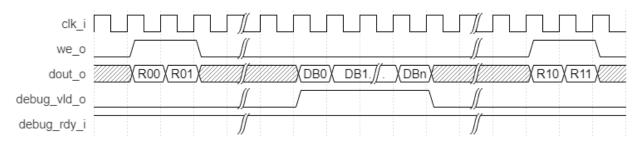


Figure 2.9. Result and Debug Interface Timing Diagram

#### 2.2.4. DRAM Interface

Command code must be written in the DRAM before the execution of CNN Plus Accelerator IP Core. Input data may be written in the DRAM too. During the execution of CNN Plus Accelerator IP Core, it reads command code from the DRAM and performs calculation with internal sub-execution engine per command code. Intermediate data may be transferred from/to the DRAM per command code.

Refer to AXI4 Protocol Specification for the timing diagram of DRAM Interface.

#### 2.2.5. Command FIFO Interface

External command FIFO must provide the command code that is generated by the Lattice sensAl Neural Network Compiler software. When the fill\_o signal is at logic 0, the external logic should load the command code to external command FIFO. After the start\_i control signal is asserted, the fill\_o signal go to logic 1 and normal FIFO operation proceeds as shown in. When fifo\_empty\_i is low, the fifo\_dout\_i signal contains a valid command. The command is latched/sampled when fifo\_empty\_i is low and fifo\_rd\_o is high. The external command FIFO holds current command (value of fifo\_dout\_i) when fifo\_empty\_i and fifo\_rd\_o are both low.

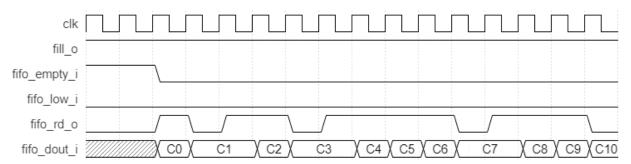


Figure 2.10. Command FIFO Interface Timing Diagram



## 2.3. Attribute Summary

The configurable attributes of the CNN Plus Accelerator IP Core are shown in Table 2.2, Table 2.3, and Table 2.4, and are described in Table 2.5. The attributes can be configured through the IP Catalog's Module/IP Block wizard of the Lattice Radiant software. For the supported devices, the applicable attributes are listed in Table 2.7.

Table 2.2. Attributes for EngineConfig Tab

Attribute	Selectable Values	Default	Dependency on Other Attributes
Embedded Mode <sup>1</sup>	Checked, Unchecked	Unchecked	-
Machine Learning Type	COMPACT_CNN, OPTIMIZED_CNN, EXTENDED_CNN	OPTIMIZED_CNN	_
Convolution Engine <sup>2</sup>	SINGLE_CONV, DUAL_CONV, QUAD_CONV	SINGLE_CONV	Only for <i>Machine Learning Type</i> = OPTIMIZED_CNN or EXTENDED_CNN
Memory Type	SINGLE_LRAM, DUAL_LRAM, QUAD_LRAM	DUAL_LRAM	
LRAM Enable Output Register <sup>2</sup>	Checked, Unchecked	Checked	Only for Machine Learning Type = COMPACT_CNN or OPTIMIZED_CNN
Scratch Pad Memory Size	1K, 4K, 8K, 16K	1K	v1.0.0: Machine Learning Type = COMPACT_CNN
	1K, 2K, 4K, 8K		v1.0.0: Machine Learning Type = OPTIMIZED_CNN v1.1.x: Same value for both Machine Learning Type 8K is not supported for Scratch Pad Memory Mode = OCTA
Scratch Pad Memory Mode <sup>2</sup>	QUAD, OCTA	QUAD	Only for <i>Machine Learning Type</i> = COMPACT_CNN
Line Buffer Size <sup>1</sup>	512, 1024, 2048	512	Machine Learning Type = OPTIMIZED_CNN
Byte Mode	UNSIGNED, SIGNED, DISABLE	UNSIGNED	The DISABLE option is only available when Machine Learning Type = COMPACT_CNN
Byte Shift	b000, b001, b010, b011, b100, b101, b110, b111	b101	v1.0.0: Fixed to b101 Only for <i>Machine Learning Type</i> = COMPACT_CNN
Maximum Burst Length	32, 256	32	Embedded Mode is unchecked
Maximum Argument Size <sup>2</sup>	4096, 8192	4096	Machine Learning Type = EXTENDED_CNN
FC Layer Feature	NORMAL, 4BIT4CH, 4BIT8CH	NORMAL	Only for Machine Learning Type = OPTIMIZED_CNN
Convolution Higher Precision	Checked, Unchecked	Unchecked	Only for Machine Learning Type = OPTIMIZED_CNN or EXTENDED_CNN
LUT Reduction <sup>3</sup>	Checked, Unchecked	Unchecked	Machine Learning Type = OPTIMIZED_CNN
LUT Reduction Factor	1, 2, 3, 4	3	Only for LUT Reduction = Checked

#### Notes:

- 1. The attributes *Embedded Mode* and *Line buffer Size* are available from v1.1.0.
- The attributes Convolution Engine, LRAM Enable Output Register, Scratch Pad Memory Mode and Maximum Argument Size are available from v1.2.0.
- 3. The attribute LUT Reduction can only be used in LIFCL40, LIFCL33, LFD2NX40, LFCPNX50, and LFCPNX100 devices.



## Table 2.3. Attributes for LUT-Activation-Config Tab (Scale-Module)

		•	
Attribute	Selectable Values	Default	Dependency on Other Attributes
Enable activation function in scale module	Checked, Unchecked	Unchecked	Only for Machine Learning Type = OPTIMIZED_CNN
Lookup table address width	12–16	12	_
Lookup table data width	8, 16	8	_
Scale Lookup table init file relative path	File path	empty	
Enable MSB clipping	Checked, Unchecked	Unchecked	_

#### Table 2.4. Attributes Table for LUT-Activation-Config Tab (FC-Module)

Attribute	Selectable Values	Default	Dependency on Other Attributes
Enable activation function in FC module	Checked, Unchecked	Unchecked	Only for <i>Machine Learning Type</i> = OPTIMIZED_CNN
Lookup table address width	12–16	12	_
Lookup table data width	8, 16	8	_
FC Lookup table init file relative path	File path	empty	_
Enable MSB clipping	Checked, Unchecked	Unchecked	_

## **Table 2.5. Attributes Descriptions**

Attribute	Description		
Embedded Mode	Selects between AXI4 I/F (unchecked) and command FIFO I/F (checked). When command FIFO I/F is enabled, the external logic feeds the command codes to this interface.		
Machine Learning Type	This option allows you to choose between the following two implementation options:  COMPACT_CNN – Low footprint, similar to Compact CNN Accelerator for thunderplus but enhanced with external DRAM access capability and supports larger network.  OPTIMIZED_CNN – High performance – four parallel engines that efficiently calculate 3 × 3 convolution, depth-wise convolution and 1 × 1 convolution.  EXTENDED_CNN – Same performance as OPTIMIZED_CNN plus support to max pooling/unpooling with max argument and adding of RELU result and input data.		
Convolution Engine	Specifies the number of convolution engines.  SINGLE_CONV – Use Single Convolution Engine  DUAL_CONV – Use Paired Convolution Engine  QUAD_CONV – Use Quadruple Convolution Engine  This option is only for <i>Machine Learning Type</i> == OPTIMIZED_CNN or EXTENDED_CNN.		
Memory Type	Specifies the memory size of Activation Data Storage.  SINGLE_LRAM – Total size is 64 kB.  DUAL_LRAM – Total size is 128 kB.  QUAD_LRAM – Total size is 256 kB.  Set this attribute based on required size of peak Blob data in the Neural Network Compiler:  • Use SINGLE_LRAM when network blob data ≤ 64 kB  • Use DUAL_LRAM when network blob data > 64 kB and ≤ 128 kB  • Use QUAD_LRAM when network blob data > 128 kB and ≤ 256 kB		
LRAM Enable Output Register	Enables the use of output register for LRAM  Unchecked – Disable the use of output register  Checked – Enable the use of output register  This option is only for <i>Machine Learning Type</i> == COMPACT_CNN or OPTIMIZED_CNN.		
Scratch Pad Memory Size	Configures the memory size of each EBRAM in Convolution Scratch Storage. 1K – Size of each EBRAM is 2 kB (1024 × 2 bytes).		



Attribute	Description
	2K – Size of each EBRAM is 4 kB (2048 × 2 bytes).
	4K – Size of each EBRAM is 8 kB (4096 × 2 bytes).
	8K – Size of each EBRAM is 16 kB (8192 × 2 bytes).
	The Lattice sensAI Neural Network Compiler software schedules data access to Convolution Scratch Storage based on the size of this memory and size of the input image It is recommended to allocate more size to this memory when your design and target device allows it.
Scratch Pad Memory Mode	Specifies the total size of scratch pad
	QUAD – 2x of scratch pad
	OCTA – 4x of scratch pad
	This option is only for <i>Machine Learning Type</i> == COMPACT_CNN.
Line Buffer Size	Selects the size of line buffer.
	512 – Select this when the input image is ≤ QVGA (320×240).
	1024 – Select this when the input image is VGA
	2048 – Select this when the input image size is ≥ VGA
Dotte Manda	
Byte Mode	Specifies the byte mode of input data.  SIGNED – Input data is signed 16-bit/8-bit data.
	UNSIGNED – Input data is signed 16-bit data or unsigned 8-bit data.
	DISABLE – Input data is unsigned 16-bit data only. This option saves LUT because the byte mode support is not implemented. This option is only for <i>Machine Learning Type</i> == COMPACT_CNN.
Byte Shift	Specifies the shift amount for converting byte data and 16-bit data. Conversion occurs
	when storing the result of convolution, pooling, and full connect layers into the Activation Data Storage as Byte Mode (signed or unsigned) and reading out the byte data to the 16-bit dout_o signal. Refer to Table 2.6. for the shifting of data byte based on Byte Shift and Byte Mode. This is fixed to default value (b101) for IP Core v1.0 and sensAl v3.0.
Maximum Burst Length	Specifies the maximum burst length of AXI4 bus.
Maximum Barse Length	This should be set less than or equal to the maximum burst length that is supported by
	the connected subordinate device/memory.
Maximum Argument Size	Specifies the number of output data which can be processed by max pooling/unpooling without updating table. This option is only for <i>Machine Learning Type</i> == EXTENDED_CNN
FC Layer Feature	Normal – 8-bit weight, 8-bit data
,	4BIT4CH – 4-bit weight, 4-channel/8-bit Data-in
	4BIT8CH – 4-bit weight, 8-channel/8-bit Data-in
Convolution Higher Precision	Improves convolution result accuracy by using higher precision internal calculations.
convolution riigher recision	Unchecked – Lower internal precision, consumes less LUT and EBR resources
	Checked – Higher internal precision, consumes more LUT and EBR resources
LUT Reduction	Enables the <i>LUT Reduction</i> feature which requires availability of unused DSP slices.
	Unchecked – No changes to LUTs and DSP slices consumed to perform the calculation.
	Checked – Reduces LUTs consumed to perform the calculation, increases DSP slices to
	perform the calculation.
	This option is only for <i>Machine Learning Type</i> == OPTIMIZED_CNN and is not supported of LIFCL17 and LFD2NX-17 devices.
LUT Reduction Factor	Specifies the degree of LUT reduction. The higher the factor, the higher the consumption
	of DSP slices. Select the factor based on LUT/DSP slices utilization tradeoff and availability of unused DSP slices in the device used.
	1 – Sets LUT reduction factor to 1, the lowest degree of LUT reduction.
	2 – Sets LUT reduction factor to 2.
	3 – Sets LUT reduction factor to 3.
	4 – Sets LUT reduction factor to 4, the highest degree of LUT reduction.
Enable activation function in FC/Scale module	Enables/disables the LUT activation table based on file.
Lookup table address width	LUT length.

© 2024 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



Attribute	Description
Lookup Table init file relative path	LUT file path.
Enable MSB clipping	Enables/disables the clipping of the output bit over data.

Table 2.6. Shifting of Data Byte Based on Byte Shift and Byte Mode

Byte Shift	Byte Mode == SIGNED	Byte Mode == UNSIGNED
b000	{Byte[7:0], 8'b0}	{1'b0, Byte[6:0], 8'b0}
b001	{Byte[7], Byte[7:0], 7'b0}	{1'b0, Byte[7:0], 7'b0}
b010	{2{Byte[7]}, Byte[7:0], 6'b0}	{2'b0, Byte[7:0], 6'b0}
b011	{3{Byte[7]}, Byte[7:0], 5'b0}	{3'b0, Byte[7:0], 5'b0}
b100	{4{Byte[7]}, Byte[7:0], 4'b0}	{4'b0, Byte[7:0], 4'b0}
b101	{5{Byte[7]}, Byte[7:0], 3'b0}	{5'b0, Byte[7:0], 3'b0}
b110	{6{Byte[7]}, Byte[7:0], 2'b0}	{6'b0, Byte[7:0], 2'b0}
b111	{7{Byte[7]}, Byte[7:0], 1'b0}	{7'b0, Byte[7:0], 1'b0}

**Note:** Byte[7:0] is the result of convolution, pooling, and full connect layers. Byte[7:0] = 8'h00 when sign of 16-bit value is negative and Byte[7:0] = 8'hFF when saturated.

**Table 2.7. Supported Devices and Attributes** 

Attribute	LIFCL-17/LFD2NX-17	LIFCL-40/LFD2NX-40	LFCPNX-100
Machine Learning Type == COMPACT_CNN	0	0	0
Machine Learning Type == COMPACT_CNN			
Scratch Pad Memory Size == 1K/2K/4K	0	0	0
Scratch Pad Memory Mode == QUAD			
Machine Learning Type == COMPACT_CNN			
Scratch Pad Memory Size == 8K	X	0	0
Scratch Pad Memory Mode == QUAD			
Machine Learning Type == COMPACT_CNN			
Scratch Pad Memory Size == 1K/2K	0	0	0
Scratch Pad Memory Mode == OCTA			
Machine Learning Type == COMPACT_CNN			
Scratch Pad Memory Size == 4K	X	0	0
Scratch Pad Memory Mode == OCTA			
Machine Learning Type == OPTIMIZED_CNN	0	0	0
Machine Learning Type == OPTIMIZED_CNN	X	0	0
Use Paired Convolution Engine == Checked	^	0	U
Machine Learning Type == OPTIMIZED_CNN	X	Х	0
Use Quadrupled Convolution Engine == Checked	^	^	U
Machine Learning Type == EXTENDED_CNN	0	0	0
Machine Learning Type == EXTENDED_CNN	X	0	0
Use Paired Convolution Engine == Checked	^	0	U
Machine Learning Type == EXTENDED_CNN	X	Х	0
Use Quadrupled Convolution Engine == Checked	^	^	U
Memory Type == SINGLE_LRAM	0	0	0
Memory Type == DUAL_LRAM	0	0	0
Memory Type == QUAD_LRAM	0	Χ	0



### 2.4. Clock Domain

The clk\_i and aclk\_i domains for Compact CNN implementation type are shown in Figure 2.11. ,the aclk\_i domain is limited to the AXI4 interface. The difference in clock is absorbed by Lattice Dual Clocked FIFO IP (FIFO\_DC) and is implemented in the DMA sub block. This is the same for the Optimized CNN implementation type.

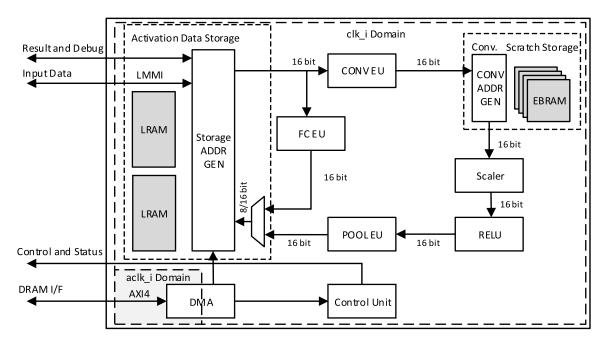


Figure 2.11. Clock Domain Diagram

## 2.5. Reset Behavior

When resetn signal asserts, output ports return to logic 0 in the next cycle. When resetn deasserts, output ready signals assert in the next cycle. A timing diagram of reset during AXI4 access is shown as an example in Figure 2.12. Only a few output signals are shown in this figure. The clk and aclk signals are 50% out-of-phase to show asynchronous relationship. The minimum resetn assert period is one cycle of the slower clock between clk\_i and aclk\_i.

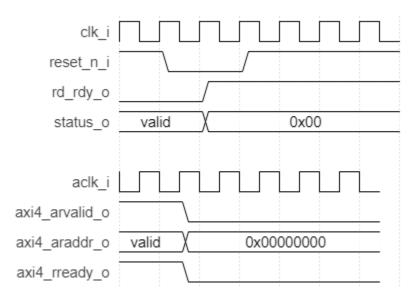


Figure 2.12. Reset Timing Diagram



Some AXI4 output signals are constant outputs; these are not affected by reset. Refer to Table 2.1 for the AXI4 output signals that are constant.

## 2.6. Register Description

CNN Plus Accelerator IP Core has no user-configurable register.

## 2.7. Operation Sequence

Operation sequence must be executed in the following steps. debug\_rdy\_i is tied to 1 in this example.

- Assert Reset.
- Deassert Reset. start i must be deasserted.
- 3. Write command sequence code, which is generated by the Lattice sensAl Neural Network Compiler software into the DRAM starting at the address specified by code base addr i signal.
- 4. Check that rd\_rdy\_o is high, if not, go back to step 1.
- 5. Write input data into the DRAM at the proper address, which is decided by command sequence, or directly write into the internal memory block of the CNN Plus Accelerator IP Core through the input data LMMI interface.
- 6. Assert start\_i and check rd\_rdy\_o. The rd\_rdy\_o signal deasserts to 0 after asserting start\_i.
- 7. Deassert start i.
- Check if we o if code has direct output commands. Collect dout o while we o == 1.
- 9. Repeat from step 5.

#### 2.7.1. Command Format

Command is a sequence of 32-bit data with or without additional parameters or weights as shown in Figure 2.13. . It should be loaded at DRAM address specified by code\_base\_addr\_i signal before execution. The command is generated by the Lattice sensAl Neural Network Compiler software. For more information, refer to Lattice sensAl Neural Network Compiler Software User Guide (FPGA-UG-02052).

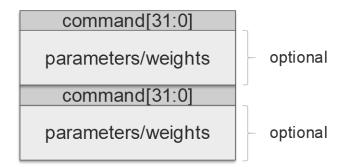


Figure 2.13. Command Format

#### 2.7.2. Input Data Format

22

Input data is a sequence of 8-bit or 16-bit data. Memory index and address are decided by Neural Network. Therefore, the external block should process input raw data and write input data to the Lattice CNN Plus Accelerator IP Core through the input data write interface. The input data interface of this IP core is 16-bit wide. If input data is 8-bit (Byte Mode is SIGNED/UNSIGNED), use only the lower 8-bit of the input data interface and set the upper 8-bit to 8'h00.

For example, face detection neural network may take  $32 \times 32$  of R, G, B planes at memory index 0 with address 0x0000 for Red plane, 0x0400 for Green plane and 0x0800 for Blue plane. Another example is object detection neural network may take  $90 \times 90$  of R, G, B planes, which are assigned to memory index 0, 1, and 2 respectively. Because memory assignment is defined by neural network, external block should handle input raw data and write it to proper position of internal memory of the CNN Plus Accelerator IP Core.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.

FPGA-IPLIG-02115-1 6



Writing input data to DRAM and using Load command to fetch input data are also possible in the case of large input data. The IP core expects data in little-endian order.

The input data must agree with the Byte Mode attribute settings. Refer to Table 2.5 for details.

### 2.7.3. Output Data Format

Output data is a sequence of 16-bit data which is controlled by commands. The amount of data is also decided by Neural Network, that is, by output blobs. External block should interpret output sequence and generate usable information. For example, face detection outputs 2-beat burst (two consecutive) of 16-bit data, the first is confidence of non-face while the second one is confidence of face. Whenever the latter is larger than the former, conclusion is Face. The IP core outputs data in little-endian order.

## 2.8. Supported Commands

Command sequences are generated by Lattice sensAl Neural Network Compiler Software (FPGA-UG-02052).



## 3. IP Core Generation, Simulation, and Validation

This section provides information on how to generate the CNN Plus Accelerator IP Core using the Lattice Radiant software and how to run simulation and synthesis. For more details on the Lattice Radiant software, refer to the Lattice Radiant Software 2.0 User Guide.

## 3.1. Licensing the IP

An IP core-specific device-specific license string is required to enable full, unrestricted use of the Lattice CNN Plus Accelerator IP Core in a complete, top-level design. You may refer to the instructions on how to obtain licenses for Lattice IP cores at Lattice Software Licensing web page.

You may download and generate the CNN Plus Accelerator IP Core and fully evaluate the core through functional simulation and implementation (synthesis, map, place, and route) without an IP license string. The CNN Plus Accelerator IP Core supports Lattice's IP hardware evaluation capability, which makes it possible to create versions of the IP core which operate in hardware for a limited time (approximately four hours) without requiring an IP license string. See Hardware Evaluation section for further details. However, a license string is required to generate bitstream file that does not include the hardware evaluation timeout limitation.

Note: All IP has a license whether in eval mode or full mode. Difference is license string.

## 3.2. Generation and Synthesis

The Lattice Radiant software allows you to customize and generate modules and IPs and integrate them into the device's architecture. The procedure for generating the CNN Plus Accelerator IP Core in Lattice Radiant software is described below.

To generate the CNN Plus Accelerator IP Core:

- 1. Create a new Lattice Radiant software project or open an existing project.
- In the IP Catalog tab, double-click on CNN\_Plus\_Accelerator under IP, DSP category. The Module/IP Block Wizard
  opens as shown in Figure 3.1. Enter values in the Component name and the Create in fields and click Next.

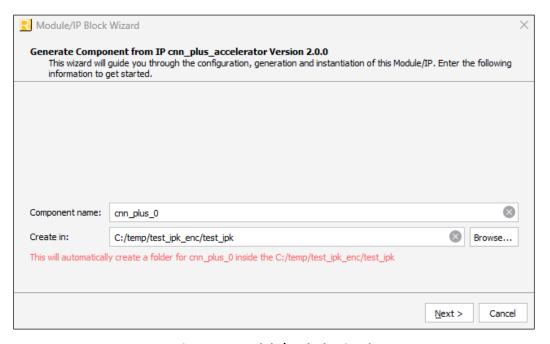


Figure 3.1. Module/IP Block Wizard



3. In the module's dialog box of the **Module/IP Block Wizard** window, customize the selected CNN Plus Accelerator IP Core using drop-down menus and check boxes. As a sample configuration, see Figure 3.2. For configuration options, see the Attribute Summary section.

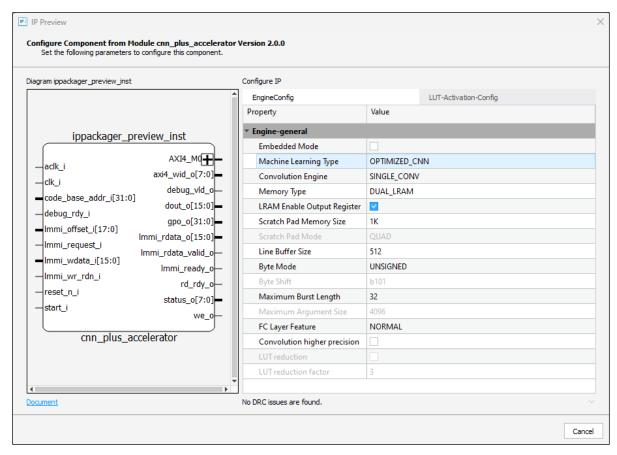


Figure 3.2. Configure User Interface of CNN Plus Accelerator IP Core



4. Click **Generate**. The **Check Generating Result** dialog box opens, showing design block messages and results as shown in Figure 3.3.

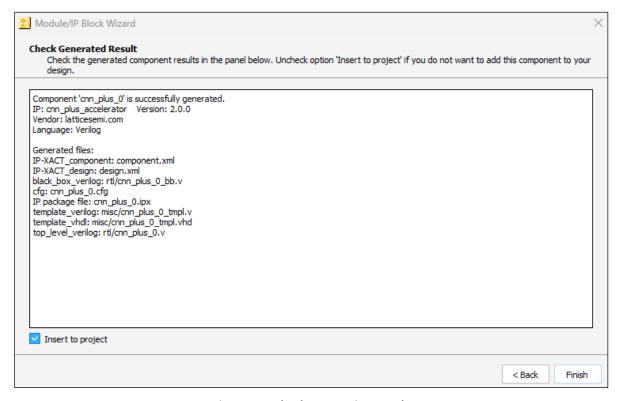


Figure 3.3. Check Generating Result

5. Click the **Finish** button. All the generated files are placed under the directory paths in the **Create in** and the **Component name** fields shown in Figure 3.1.

The generated CNN Plus Accelerator IP Core package includes the closed-box (<Component name>\_bb.v) and instance templates (<Component name>\_tmpl.v/vhd) that can be used to instantiate the core in a top-level design. An example RTL top-level reference source file (<Component name>.v) that can be used as an instantiation template for the IP core is also provided. You may also use this top-level reference as the starting template for the top-level for their complete design. The generated files are listed in Table 3.1.

**Table 3.1. Generated File List** 

Attribute	Description
<component name="">.ipx</component>	Contains the information on the files associated to the generated IP.
<component name="">.cfg</component>	Contains the parameter values used in IP configuration.
component.xml	Contains the ipxact:component information of the IP.
design.xml	Documents the configuration parameters of the IP in IP-XACT 2014 format.
rtl/ <component name="">.v</component>	Provides an example RTL top file that instantiates the IP core.
rtl/ <component name="">_bb.v</component>	Provides the synthesis closed-box.
misc/ <component name="">_tmpl.v misc /<component name="">_tmpl.vhd</component></component>	Provide instance templates for the IP core.



## 3.3. Running Functional Simulation

The CNN Plus Accelerator IP does not have a functional simulation testbench.

### 3.4. Hardware Evaluation

The CNN Plus Accelerator IP Core supports Lattice's IP hardware evaluation capability, which makes it possible to create versions of the IP core that operate in hardware for a limited period of time (approximately four hours) without requiring the purchase of an IP license. It may also be used to evaluate the core in hardware in user-defined designs. Choose Project > Active Strategy > Translate Design Settings. The hardware evaluation capability may be enabled/disabled in the Strategy dialog box. It is enabled by default.



## 4. Hardware Validation

This IP has been validated using the *CertusPro-NX MobileNet Object Classification on VVML Board* reference design on VVML CPNX-100 board, and *CrossLink-NX Scene Segmentation* reference design on VVML LIFCL-40 board. Refer to Object Classification Reference Design and Scene Segmentation Reference Design sections within the Lattice sensAl Stack for these reference designs and details on the configuration used for the validation.



## 5. Ordering Part Number

The Ordering Part Numbers (OPN) for CNN Plus Accelerator IP Core are the following:

- CNNPLUS-ACCEL-CNX-US CNN Plus Accelerator IP for CrossLink-NX Single Seat Annual License
- CNNPLUS-ACCEL-CNX-UT CNN Plus Accelerator IP for CrossLink-NX Multi-site Perpetual License
- CNNPLUS-ACCEL-CTNX-US CNN Plus Accelerator IP for Certus-NX Single Seat Annual License
- CNNPLUS-ACCEL-CTNX-UT CNN Plus Accelerator IP for Certus-NX Multi-site Perpetual License
- CNNPLUS-ACCEL-CPNX-US CNN Plus Accelerator IP for CertusPro-NX Single Seat Annual License
- CNNPLUS-ACCEL-CPNX-UT CNN Plus Accelerator IP for CertusPro-NX Multi-site Perpetual License



## **Appendix A. Resource Utilization**

Table A.1 shows configuration and resource utilization for LFCPNX-100-9BBG484I using Lattice Synthesis Engine of Lattice Radiant software.

#### Table A.1. Performance and Resource Utilization<sup>1</sup>

Configuration	FPGA	clk, aclk (MHz)²	Registers	LUTs	LRAMs	EBRs	DSP MULT (in terms of 18x18 MULT)
ML_TYPE = OPTIMIZED_CNN	LFCPNX-100	96, 96	6636	9168	4	90	25
SCRATCH_SIZE: 8K							
LB_SIZE : 2048							
MEM_TYPE: QUAD_LRAM							
CONV_ENGINE : SINGLE_CONV							
ARGMAX_SIZE: 8192							
LRAM_ORG:1							
LONG_CONV: 1							
FC_LAYER_FEATURE: NORMAL							
(Both Scale and FC LUT activations are							
enabled)							

#### Notes:

- 1. Performance may vary when using a different software version or targeting a different device density or speed grade.
- 2. Fmax is generated when the FPGA design only contains the CNN Plus Accelerator IP Core. These values may be reduced when user logic is added to the FPGA design.



## **References**

- Certus-NX web page
- CertusPro-NX web page
- CrossLink-NX web page
- Lattice Radiant Software web page
- Lattice sensAl Stack web page
- Lattice Solutions IP Cores web page
- Lattice Insights web page for Lattice Semiconductor training courses and learning plans



# **Technical Support Assistance**

Submit a technical support case through www.latticesemi.com/techsupport. For frequently asked questions, refer to the Lattice Answer Database at www.latticesemi.com/Support/AnswerDatabase.



# **Revision History**

#### Revision 1.6, IP v2.0.0, December 2024

Section	Change Summary
All	<ul> <li>Removed <i>Core</i> from the document title.</li> <li>Added the IP version information on the cover page.</li> </ul>
	<ul> <li>Updated all instances of Machine Learning to Machine Learning.</li> </ul>
Acronyms in This Document	Added the following acronyms:
Actoriyins in this bocument	Fully Connected (FC)
	Learned Step Size (LSQ)
	Look-Up Table (LUT)
	Most Significant Bit (MSB)
Introduction	Updated Resources and Lattice Implementation in Table 1.1. Quick Facts.
	Updated the Features section.
Attribute Summary	<ul> <li>Renamed Table 2.2. Attributes Table to Table 2.2. Attributes for EngineConfig Tab.</li> <li>Updated Default for the Machine Learning Type and LRAM Enable Output Register attributes, and added the FC Layer Feature attribute to Table 2.2. Attributes for EngineConfig Tab.</li> <li>Added Table 2.3. Attributes for LUT-Activation-Config Tab (Scale-Module) and Table 2.4. Attributes Table for LUT-Activation-Config Tab (FC-Module), and adjusted the heading number of remaining tables.</li> <li>Added the following attributes to Table 2.5. Attributes Descriptions:         <ul> <li>FC Layer Feature</li> <li>Enable activation function in FC/Scale module</li> <li>Lookup table address width</li> <li>Lookup Table init file relative path</li> <li>Enable MSB clipping</li> </ul> </li> </ul>
Core Generation, Simulation,	Updated all figures in the Generation and Synthesis section.
and Validation	
Ordering Part Number	Updated this section.
Resource Utilization	Updated Table A.1. Performance and Resource Utilization <sup>1</sup> .
References	Added the Lattice Solutions IP Cores web page.
	Made editorial fixes.

#### Revision 1.5, December 2023

Section	Change Summary
All	<ul> <li>Changed the document name from CNN Plus Accelerator IP Core – Lattice Radiant Software to CNN Plus Accelerator IP Core.</li> <li>Standardized the Lattice sensAl naming.</li> <li>Minor adjustments to ensure the document is consistent with Lattice Semiconductor's inclusive language policy.</li> </ul>
Disclaimers	Updated boilerplate.
Inclusive Language	Added boilerplate.
Introduction	<ul> <li>Added LIFCL-33, LIFCL-33U and LFCPNX-50 as targeted devices in Table 1.1.</li> <li>Added new features in Features section.</li> </ul>
Attribute Summary	Updated Table 2.2 and Table 2.3.
Core Generation, Simulation, and Validation	Updated Figure 3.1, Figure 3.2, and Figure 3.3.
Hardware Validation	Added this section.
Ordering Part Number	Updated the header number.

© 2024 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



Section	Change Summary
References	<ul> <li>Removed in latticesemi.com from existing references.</li> <li>Added references to Lattice Radiant Software, Lattice sensAl Stack and Lattice Insights web pages.</li> </ul>
Technical Support Assistance	Added link to the Lattice Answer Database.

### Revision 1.4, March 2022

Section	Change Summary
Introduction	Removed "Dynamic support for various 1D convolution from 1 to 9 taps" from the list of features.

### Revision 1.3, October 2021

Section	Change Summary		
All	Minor adjustments in formatting across the document.		
Acronyms in This Document	Added this section.		
Introduction	<ul> <li>Added CertusPro-NX as supported FPGA Family and LFCPNX-100 as target device in Table 1.1.</li> <li>Added new features in Features section.</li> </ul>		
Functional Description	<ul> <li>Added extended CNN type in Overview.</li> <li>Updated Figure 2.1 and Figure 2.3.</li> <li>Added Figure 2.4 and Figure 2.5 for extended CNN type block diagram and order of layers.</li> <li>Updated Table 2.2 and</li> <li>Table 2.3 for the new attributes.</li> <li>Added Table 2.5 for supported devices and attributes.</li> </ul>		
Core Generation Simulation and Validation	Updated Figure 3.1, Figure 3.2, and Figure 3.3.		
Ordering Part Number	Updated content.		
References	Added reference to the CertusPro-NX web page.		
Resource Utilization	Added new rows in Table A.1 to update the resource utilization result and changed target device to LFCPNX-100-9BBG484I.		

#### Revision 1.2, May 2021

Section	Change Summary	
Functional Description	Updated Input Data Format section.	
Resource Utilization	Updated resource utilization result.	

#### Revision 1.1, December 2020

Section	Change Summary
Introduction	<ul> <li>Added Certus-NX as supported FPGA Family, LFD2NX-40 and LFD2NX-17 as target devices in Quick Facts section.</li> <li>Added new features in Features section.</li> </ul>
Functional Description	<ul> <li>Added gpo_o and Command FIFO interface in Figure 2.1, Figure 2.4, and Table 2.1.</li> <li>Added new sections: General Purpose Output and Command FIFO Interface.</li> <li>Updated Table 2.2 and Table 2.3 for the new attributes.</li> </ul>
Core Generation Simulation and Validation	Updated Figure 3.1, Figure 3.2, and Figure 3.3.
Ordering Part Number	Updated content.
References	Added Certus-NX Web page in latticesemi.com.
Resource Utilization	Updated resource utilization result.

© 2024 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



#### Revision 1.0, May 2020

Section	Change Summary
All	Initial release



www.latticesemi.com