

DPUCAHX8H for Convolutional Neural Networks

PG367 (v1.0) July 22, 2021



Table of Contents

Chapter 1: Introduction.....	3
Features.....	3
IP Facts.....	4
Chapter 2: Overview.....	5
Core Overview.....	5
Hardware Architecture.....	6
Development Tools.....	7
Example System with DPUCAHX8H.....	9
Vitis AI Development Kit.....	11
Chapter 3: Product Specification.....	12
Port Descriptions.....	12
Register Space.....	13
Interrupts.....	18
Chapter 4: DPU Configuration.....	19
Configuration Options.....	20
Chapter 5: Development Flow.....	21
Customizing and Generating the Core in Shell Mode with Vitis Flow.....	21
Chapter 6: Upgrading.....	27
Chapter 7: Additional Resources and Legal Notices.....	28
Xilinx Resources.....	28
Documentation Navigator and Design Hubs.....	28
References.....	28
Revision History.....	29
Please Read: Important Legal Notices.....	29

Introduction

The Xilinx[®] Deep Learning Processor Unit (DPU) is a series of soft IP for convolutional neural networks acceleration. The DPUCAHX8H is a high throughput CNN inference IP for Alveo[™] cards with high bandwidth memory (HBM). It runs with a set of efficiently optimized instructions and it can support most convolutional neural networks, such as VGG, ResNet, GoogLeNet, YOLO, SSD, FPN, etc.

Features

- Supports one AXI slave interface for accessing configuration and status registers.
- Supports one AXI master interface for code fetch.
- Supports two AXI master interface for model parameters loading.
- Supports 1-5 AXI master interface for accessing input/output/intermediate feature map stored in the HBM.
- Supports all AXI master interfaces with 256-bit width.
- DPU functionality includes the following:
 - Configurable number of processing engines (PE).
 - Convolution and deconvolution
 - Max pooling
 - Average pooling
 - ReLU, ReLU6, and Leaky ReLU
 - Concat
 - Elementwise-sum
 - Dilation
 - Reorg
 - Fully connected layer
 - Batch normalization

- Split

IP Facts

LogiCORE™ IP Facts Table	
Core Specifics	
Supported Device Family	Alveo™ U280 and U50/U50LV Data Center accelerator cards
Supported User Interfaces	AXI4-Lite CSR Interface
Resources	Chapter 4: DPU Configuration
Provided with Core	
Design Files	Encrypted RTL
Example Design	Verilog
Test Bench	Not Provided
Constraints File	Xilinx Constraints File
Simulation Model	Not Provided
Supported S/W Driver ¹	Xilinx® Runtime (XRT)
Tested Design Flows ²	
Design Entry	Vitis™ unified software platform
Simulation	N/A
Synthesis	Vivado® Synthesis
Support	
Xilinx Support web page	

Notes:

1. Vitis™ AI development flow.
2. For the supported versions of the tools, see the Linux OS and driver support information are available from the *Vitis Unified Software Platform Documentation: Application Acceleration Development* ([UG1393](#)).

Overview

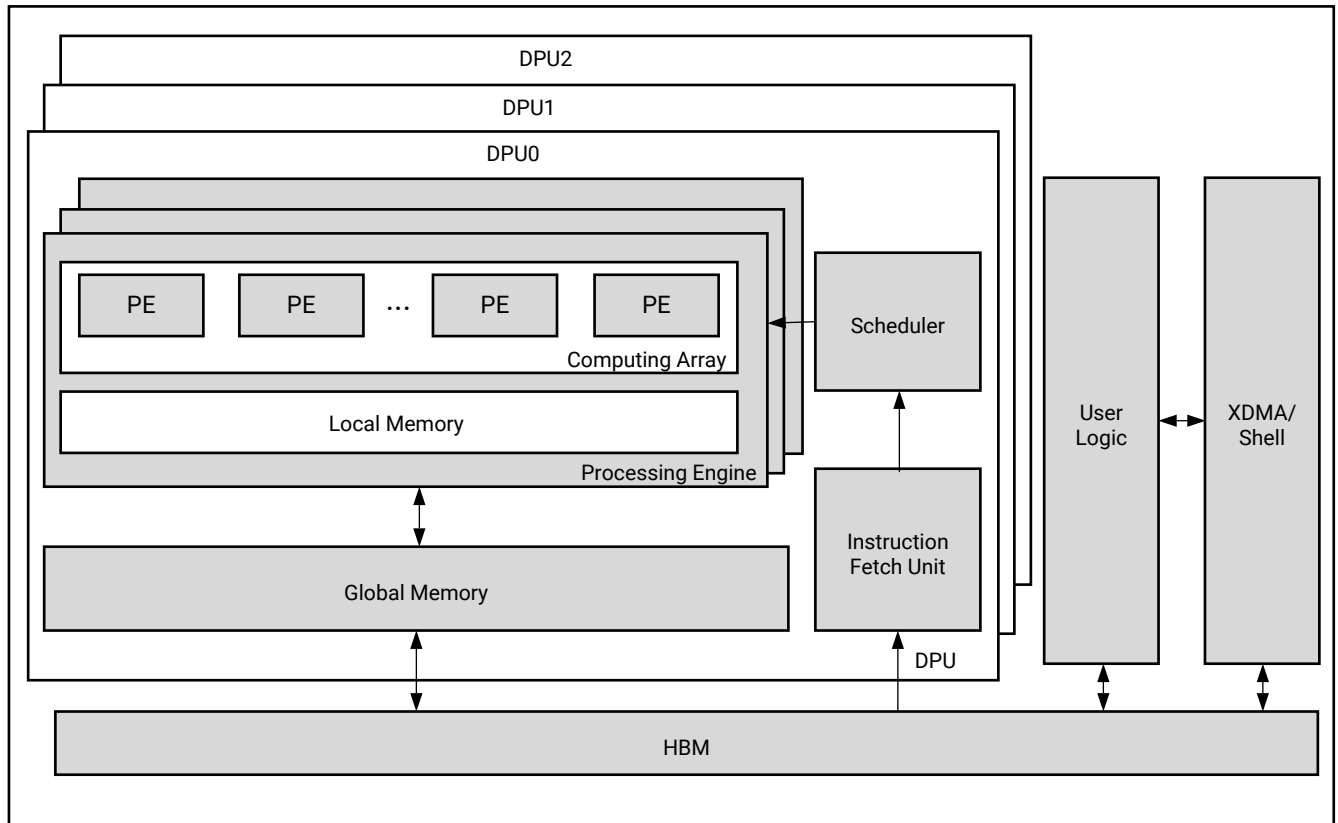
Core Overview

The Xilinx[®] DPUCAHX8H is a programmable DPU core optimized for convolutional neural networks, mainly for high throughput applications. The core includes a high-performance scheduler module, a hybrid computing array module, an instruction fetch module, and a frame buffer module. It uses a specialized instruction set that allows efficient implementation of many convolutional neural networks. Some examples of convolutional neural networks which have been deployed include VGG, ResNet, GoogLeNet, YOLO, SSD, FPN, and many others.

The DPUCAHX8H is implemented in the programmable logic (PL) of the Alveo[™] U280 and U50/U50LV Data Center accelerator cards.

The following figure shows the top-level block diagram of the DPUCAHX8H:

Figure 1: DPU Top-Level Block Diagram



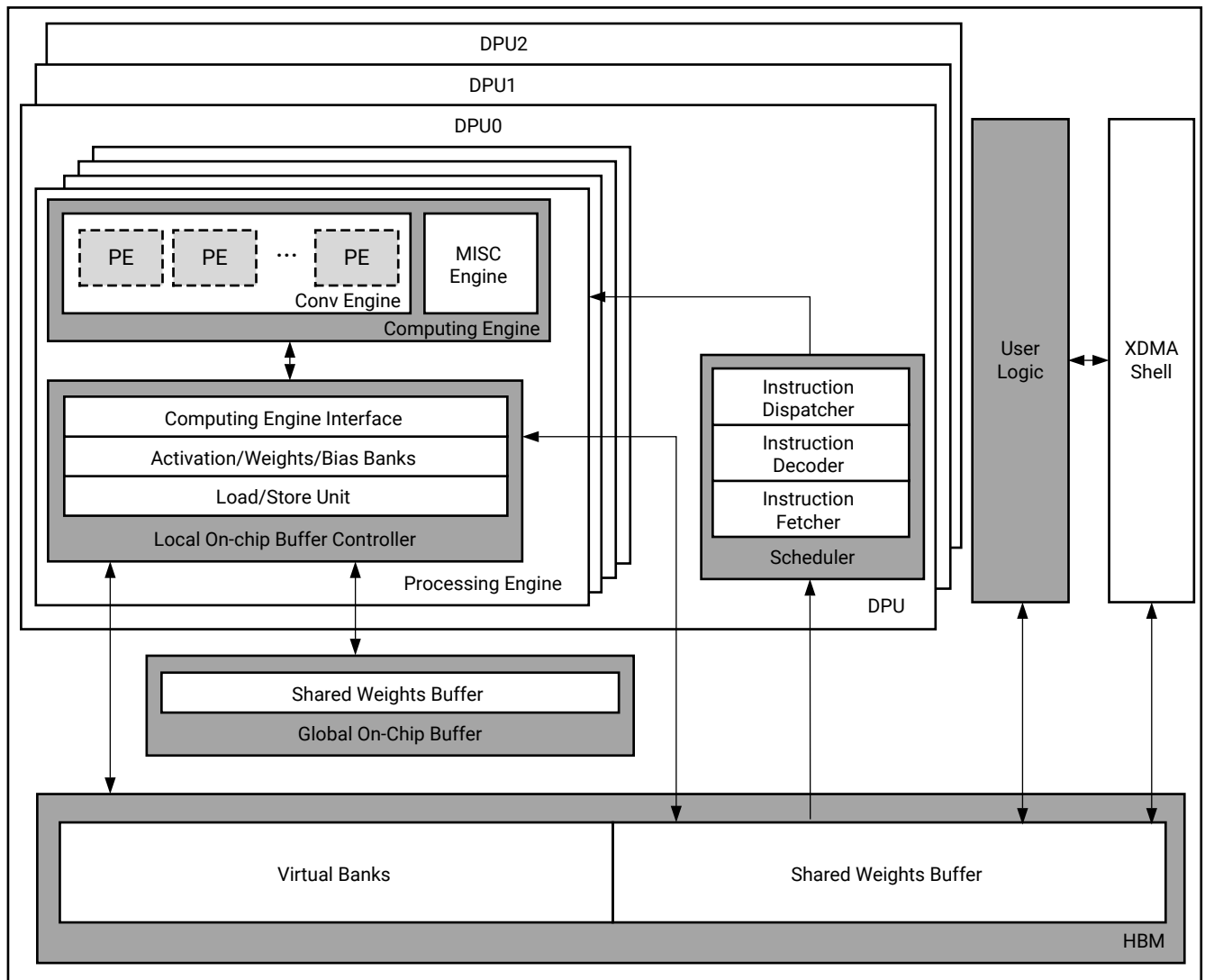
X23531-062221

Hardware Architecture

The detailed hardware architecture of the DPUCAHX8H is shown in the following figure. Each implementation has one to three DPU cores, and each DPU has one to five processing engines. The number of cores and PEs/cores are chosen based on throughput needs versus FPGA resource usage. The HBM memory space is divided into virtual banks and system memory. The virtual banks are used to store temporary data and the system memory is used to store instructions, input images, output results, and user data. After starting up, the DPU fetches model instructions from system memory to control the operation of the computing engine. The model instructions are generated by the Vitis AI compiler (running on the host server) which performs substantial optimizations.

On-chip memory is used to buffer weights, bias, intermediate data, and output data to achieve high throughput and efficiency. The local buffer is private to each PE; the global buffer is shared by the PEs in the same DPU core. A deeply pipelined design is used for the computing engine. PEs which include the conv engine, depthwise conv engine, and misc logic take full advantage of the fine-grained building blocks such as multipliers, adders, and accumulators in Xilinx® devices.

Figure 2: DPU Hardware Architecture



X23530-062221

Development Tools

The Vitis™ Integrated Design Environment (IDE) version 2020.2 is required to integrate the DPU into your projects. Contact your local sales representative if the project requires an older version of the Vitis™ software platform.

Note: For timing closure issues, use Vivado® Design Suite 2020.2 instead of 2021.1.

Device Resources

The DPUCAHX8H can only be deployed on the Alveo U50/U50LV and U280 cards.

For more information on resource utilization, see [Chapter 4: DPU Configuration](#).

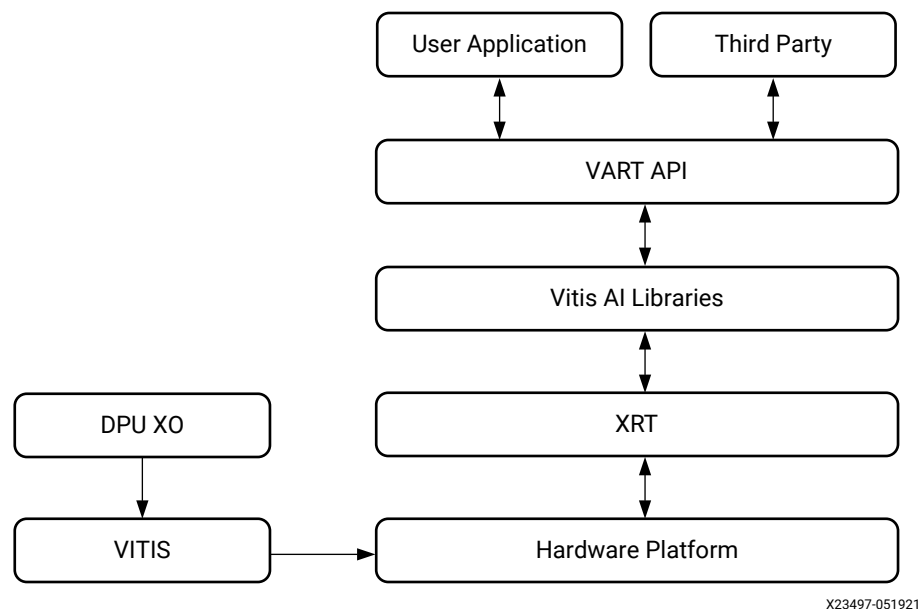
DPU Development Flow

The DPU is available with the software development stack of Vitis AI development kit. Free developer resources can be obtained from the [Xilinx website](#).

The *Vitis AI User Guide* ([UG1414](#)) describes how to use the DPU for deploying machine learning applications with the Vitis AI tools. The development flow for DPU applications is summarized in the following steps and shown in the following figure.

1. Use the Vitis tool to generate the bitstream.
2. Download the bitstream to the target board. For instructions on how to set up a running environment for DPU applications, refer to the *Vitis AI User Guide* ([UG1414](#)).

Figure 3: DPU Development Flow

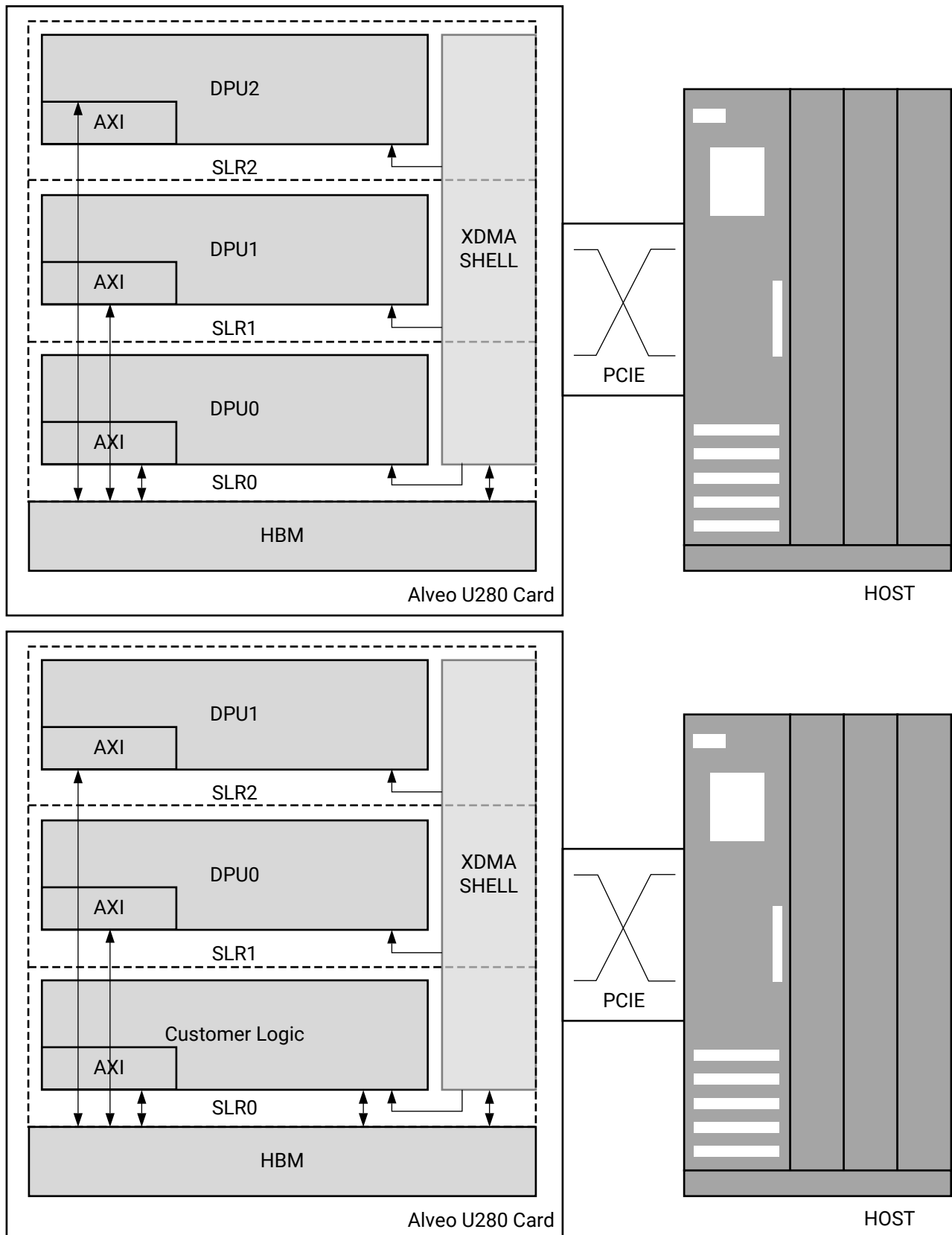


X23497-051921

Example System with DPUCAHX8H

The following figure shows two example system block diagrams with the Alveo U280 Data Center accelerator card which includes an UltraScale+ XCU280 FPGA and a PCIe[®] port. The card is inserted into the PCIe slot of the host server. The first example does not have user logic; all SLRs are used by DPU cores for better performance. In the other example, one SLR is reserved for user logic. You can implement a self-defined pre- or post-processing logic on this die. The DPU cores are integrated into the system through AXI interfaces that connect to the HBM, and the whole system is integrated into the server through a PCIe interconnect. It can be used to perform deep learning inference tasks such as image classification, object detection, and semantic segmentation.

Figure 4: Example System with Integrated DPU



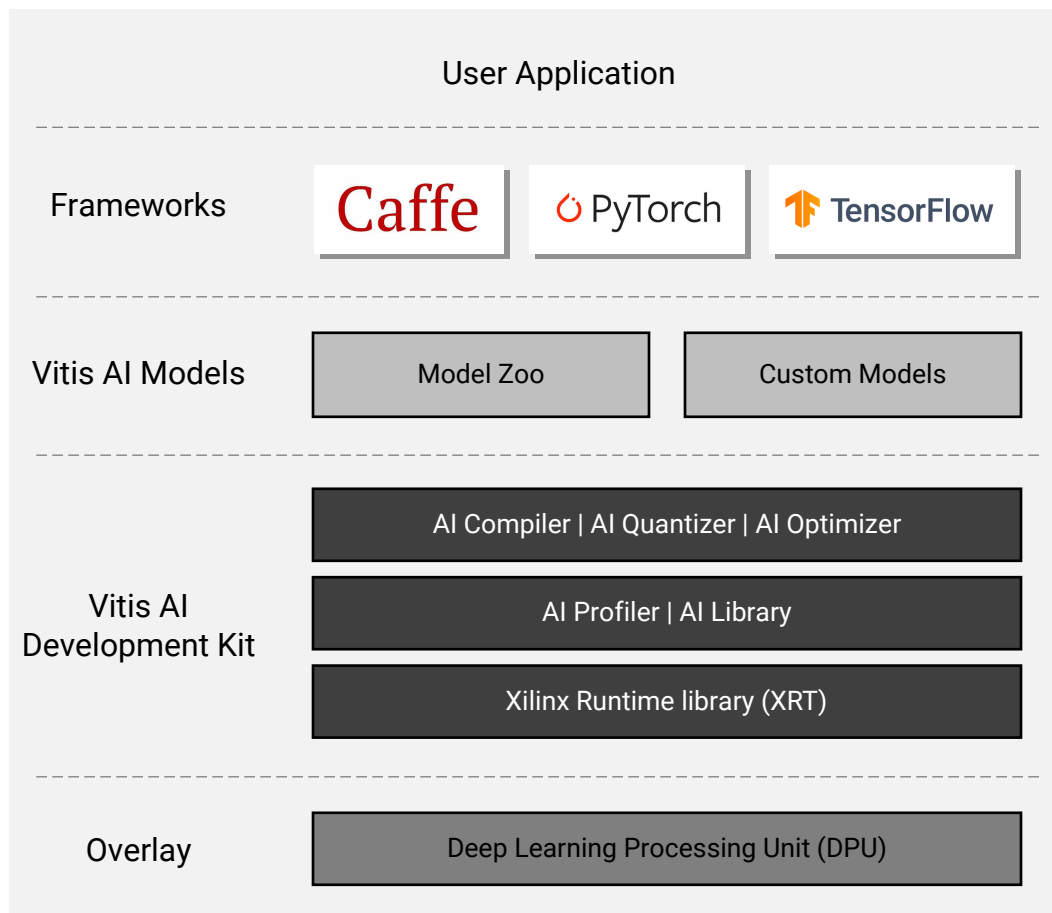
X23532-062721

Vitis AI Development Kit

The Vitis™ AI development environment is used for AI inference on Xilinx® hardware platforms. It consists of optimized IP cores, tools, libraries, models, and example designs.

As shown in the following figure, the Vitis AI development kit consists of AI Compiler, AI Quantizer, AI Optimizer, AI Profiler, AI Library, and Xilinx Runtime Library (XRT).

Figure 5: Vitis AI Stack



X24893-120920

For more information of the Vitis AI development kit, see the *Vitis AI User Guide* ([UG1414](#)).

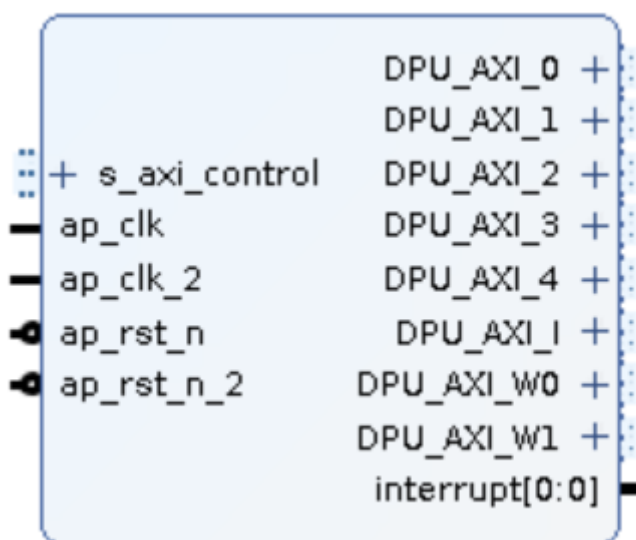
The Vitis AI development kit can be freely downloaded from [here](#).

Product Specification

Port Descriptions

The top-level interfaces of the DPU with five processing engines is shown in the following figure:

Figure 6: DPU IP Port



DPU Signals

The following table lists the 5-engine DPU I/O signals and their function descriptions.

Table 1: DPU Signals

Signal Name	Interface Type	Width	I/O	Description
s_axi_control	Memory mapped AXI slave interface	32	I/O	32-bit memory mapped AXI interface for registers.

Table 1: DPU Signals (cont'd)

Signal Name	Interface Type	Width	I/O	Description
ap_clk	Clock	1	I	Kernel clock. The frequency of the clock should match the clock of the DPU core. The supported frequencies are 300 MHz, 275 MHz, and 250 MHz.
ap_clk_2	Clock	1	I	Reference clock for mmcm in the DPU. It is set to 100 MHz.
ap_rst_n	Reset	1	I	Active-Low reset for DPU general logic.
ap_rst_n_2	Reset	1	I	Unused in the core.
DPU_AXI_0	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU engine0.
DPU_AXI_1	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU engine1.
DPU_AXI_2	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU engine2.
DPU_AXI_3	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU engine3.
DPU_AXI_4	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU engine4.
DPU_AXI_I	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU instructions.
DPU_AXI_W0	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU parameters.
DPU_AXI_W1	Memory mapped AXI master interface	256	I/O	256-bit memory mapped AXI interface for DPU parameters.
interrupt	Interrupt	1	O	Active-High interrupt output from DPU.

Notes:

1. The DPU_AXI_1~ DPU_AXI_4 options are shown depending on the number of DPU engines.

Register Space

The DPU IP implements registers in programmable logic. The following tables show the DPU IP registers. These registers are accessible from the APU through the s_axi_control interface.

DPU Control Registers

The DPU control registers are used to start a DPU core, waiting for task finish and then clear DPU status. The details of control registers are shown in the following table.

Table 2: DPU Control Registers

Register	Address Offset	Width	Type	Description
reg_ap_control	0x000	32	r/w	Bit 0: ap_start (read/write/clear on handshake) Bit 1: ap_done (read/clear on read) Bit 2: ap_idle (read) Others: reserved
Global interrupt enable register (GIER)	0x004	32	r/w	Bit 0: global interrupt enable Others: reserved
IP interrupt enable register (IPIER)	0x008	32	r/w	Bit 0: channel 0 (ap_done) Others: reserved
IP interrupt status register (IPISR)	0x00c	32	r/w	Bit 0: channel 0 (ap_done) (read/toggle on write) Others: reserved
reg_dpu_start	0x010	32	r/w	Bit [0]: enable DPU to start
reg_finish_clr	0x018	32	r/w	Bit [0]: clear reg_finish_sts
reg_finish_sts	0x080	32	r	Bit [0]: indicate DPU has finished. The DPU finish signal is also output as DPU interrupt to trigger XDMA or custom logic. The DPU finish is a level and asynchronous signal.

DPU Configuration Registers

The DPU configuration registers are used to indicate instruction address, common address and mean value settings.

The reg_instr_addr register is used to indicate the instruction address of the DPU core.

The reg_base_addr register is used to indicate the address of input image and parameters for the DPU in external memory. The width of a DPU base address is 33 bits so it can support an address space up to 8 GB. All registers are 32-bit wide, so two registers are required to represent a 33-bit wide base address. The reg_base_addr0_l register represents the lower 32 bits of base_address0 in DPU core0, and reg_base_addr0_h represents the upper 1 bit of base_address0. There are eight groups of DPU base addresses for each DPU core, and therefore there are 40 groups of DPU base addresses for up to five DPU cores.

The details of configuration registers are shown in the following figure.

Table 3: DPU Configuration Registers

Register	Address Offset	Width	Type	Description
reg_instr_addr_l	0x140	32	r/w	The lower 32 bits of instruction address of DPU. 4 KB aligned.
reg_instr_addr_h	0x144	32	r/w	The lower 1-bit in the register represent the upper 1-bit of instruction address of DPU. 4 KB aligned.
reg_engine0_base_addr_0_l	0x100	32	r/w	The lower 32 bits of base address0 of DPU engine0.
reg_engine0_base_addr_0_h	0x104	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address0 of DPU engine0.
reg_engine0_base_addr_1_l	0x108	32	r/w	The lower 32 bits of base address1 of DPU engine0.
reg_engine0_base_addr_1_h	0x10c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address1 of DPU engine0.
reg_engine0_base_addr_2_l	0x110	32	r/w	The lower 32 bits of base address2 of DPU engine0.
reg_engine0_base_addr_2_h	0x114	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address2 of DPU engine0.
reg_engine0_base_addr_3_l	0x118	32	r/w	The lower 32 bits of base address3 of DPU engine0.
reg_engine0_base_addr_3_h	0x11c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address3 of DPU engine0.
reg_engine0_base_addr_4_l	0x120	32	r/w	The lower 32 bits of base address4 of DPU engine0.
reg_engine0_base_addr_4_h	0x124	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address4 of DPU engine0.
reg_engine0_base_addr_5_l	0x128	32	r/w	The lower 32 bits of base address5 of DPU engine0.
reg_engine0_base_addr_5_h	0x12c	32	r/w	The lower 1 bit in the register represent the upper 1 bit of base address5 of DPU engine0.
reg_engine0_base_addr_6_l	0x130	32	r/w	The lower 32 bits of base address6 of DPU engine0.
reg_engine0_base_addr_6_h	0x134	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address6 of DPU engine0.
reg_engine0_base_addr_7_l	0x138	32	r/w	The lower 32 bits of base address7 of DPU engine0.
reg_engine0_base_addr_7_h	0x13c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address7 of DPU engine0.
reg_engine1_base_addr_0_l	0x200	32	r/w	The lower 32 bits of base address0 of DPU engine1.
reg_engine1_base_addr_0_h	0x204	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address0 of DPU engine1.
reg_engine1_base_addr_1_l	0x208	32	r/w	The lower 32 bits of base address1 of DPU engine1.
reg_engine1_base_addr_1_h	0x20c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address1 of DPU engine1.
reg_engine1_base_addr_2_l	0x210	32	r/w	The lower 32 bits of base address2 of DPU engine1.
reg_engine1_base_addr_2_h	0x214	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address2 of DPU engine1.
reg_engine1_base_addr_3_l	0x218	32	r/w	The lower 32 bits of base address3 of DPU engine1.
reg_engine1_base_addr_3_h	0x21c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address3 of DPU engine1.
reg_engine1_base_addr_4_l	0x220	32	r/w	The lower 32 bits of base address4 of DPU engine1.

Table 3: DPU Configuration Registers (cont'd)

Register	Address Offset	Width	Type	Description
reg_engine1_base_addr_4_h	0x224	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address4 of DPU engine1.
reg_engine1_base_addr_5_l	0x228	32	r/w	The lower 32 bits of base address5 of DPU engine1.
reg_engine1_base_addr_5_h	0x22c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address5 of DPU engine1.
reg_engine1_base_addr_6_l	0x230	32	r/w	The lower 32 bits of base address6 of DPU engine1.
reg_engine1_base_addr_6_h	0x234	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address6 of DPU engine1.
reg_engine1_base_addr_7_l	0x238	32	r/w	The lower 32 bits of base address7 of DPU engine1.
reg_engine1_base_addr_7_h	0x23c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address7 of DPU engine1.
reg_engine2_base_addr_0_l	0x300	32	r/w	The lower 32 bits of base address0 of DPU engine2.
reg_engine2_base_addr_0_h	0x304	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address0 of DPU engine2.
reg_engine2_base_addr_1_l	0x308	32	r/w	The lower 32 bits of base address1 of DPU engine2.
reg_engine2_base_addr_1_h	0x30c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address1 of DPU engine2.
reg_engine2_base_addr_2_l	0x310	32	r/w	The lower 32 bits of base address2 of DPU engine2.
reg_engine2_base_addr_2_h	0x314	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address2 of DPU engine2.
reg_engine2_base_addr_3_l	0x318	32	r/w	The lower 32 bits of base address3 of DPU engine2.
reg_engine2_base_addr_3_h	0x31c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address3 of DPU engine2.
reg_engine2_base_addr_4_l	0x320	32	r/w	The lower 32 bits of base address4 of DPU engine2.
reg_engine2_base_addr_4_h	0x324	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address4 of DPU engine2.
reg_engine2_base_addr_5_l	0x328	32	r/w	The lower 32 bits of base address5 of DPU engine2.
reg_engine2_base_addr_5_h	0x32c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address5 of DPU engine2.
reg_engine2_base_addr_6_l	0x330	32	r/w	The lower 32 bits of base address6 of DPU engine2.
reg_engine2_base_addr_6_h	0x334	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address6 of DPU engine2.
reg_engine2_base_addr_7_l	0x338	32	r/w	The lower 32 bits of base address7 of DPU engine2.
reg_engine2_base_addr_7_h	0x33c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address7 of DPU engine2.
reg_engine3_base_addr_0_l	0x400	32	r/w	The lower 32 bits of base address0 of DPU engine3.
reg_engine3_base_addr_0_h	0x404	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address0 of DPU engine3.
reg_engine3_base_addr_1_l	0x408	32	r/w	The lower 32 bits of base address1 of DPU engine3.
reg_engine3_base_addr_1_h	0x40c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address1 of DPU engine3.
reg_engine3_base_addr_2_l	0x410	32	r/w	The lower 32 bits of base address2 of DPU engine3.

Table 3: DPU Configuration Registers (cont'd)

Register	Address Offset	Width	Type	Description
reg_engine3_base_addr_2_h	0x414	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address2 of DPU engine3.
reg_engine3_base_addr_3_l	0x418	32	r/w	The lower 32 bits of base address3 of DPU engine3.
reg_engine3_base_addr_3_h	0x41c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address3 of DPU engine3.
reg_engine3_base_addr_4_l	0x420	32	r/w	The lower 32 bits of base address4 of DPU engine3.
reg_engine3_base_addr_4_h	0x424	32	r/w	The lower 1 bit in the register represent the upper 1 bit of base address4 of DPU engine3.
reg_engine3_base_addr_5_l	0x428	32	r/w	The lower 32 bits of base address5 of DPU engine3.
reg_engine3_base_addr_5_h	0x42c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address5 of DPU engine3.
reg_engine3_base_addr_6_l	0x430	32	r/w	The lower 32 bits of base address6 of DPU engine3.
reg_engine3_base_addr_6_h	0x434	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address6 of DPU engine3.
reg_engine3_base_addr_7_l	0x438	32	r/w	The lower 32 bits of base address7 of DPU engine3.
reg_engine3_base_addr_7_h	0x43c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address7 of DPU engine3.
reg_engine4_base_addr_0_l	0x500	32	r/w	The lower 32 bits of base address0 of DPU engine4.
reg_engine4_base_addr_0_h	0x504	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address0 of DPU engine4.
reg_engine4_base_addr_1_l	0x508	32	r/w	The lower 32 bits of base address1 of DPU engine4.
reg_engine4_base_addr_1_h	0x50c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address1 of DPU engine4.
reg_engine4_base_addr_2_l	0x510	32	r/w	The lower 32 bits of base address2 of DPU engine4.
reg_engine4_base_addr_2_h	0x514	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address2 of DPU engine4.
reg_engine4_base_addr_3_l	0x518	32	r/w	The lower 32 bits of base address3 of DPU engine4.
reg_engine4_base_addr_3_h	0x51c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address3 of DPU engine4.
reg_engine4_base_addr_4_l	0x520	32	r/w	The lower 32 bits of base address4 of DPU engine4.
reg_engine4_base_addr_4_h	0x524	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address4 of DPU engine4.
reg_engine4_base_addr_5_l	0x528	32	r/w	The lower 32 bits of base address5 of DPU engine4.
reg_engine4_base_addr_5_h	0x52c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address5 of DPU engine4.
reg_engine4_base_addr_6_l	0x530	32	r/w	The lower 32 bits of base address6 of DPU engine4.
reg_engine4_base_addr_6_h	0x534	32	r/w	The lower 1 bit in the register represent the upper 1 bit of base address6 of DPU engine4.
reg_engine4_base_addr_7_l	0x538	32	r/w	The lower 32 bits of base address7 of DPU engine4.
reg_engine4_base_addr_7_h	0x53c	32	r/w	The lower 1-bit in the register represent the upper 1-bit of base address7 of DPU engine4.

DPU Debug Registers

The DPU debug registers are used to indicate the processing cycles for task.

The details of debug registers are shown in the following table.

Table 4: DPU Debug Registers

Register	Address Offset	Width	Type	Description
reg_prof_value	0x0a8	32	r	Indicates the cycle counter of DPU processing time. Saturation counting.

Interrupts

As all DPU cores work synchronously, all DPU engines generate an interrupt to signal the completion of a task. A high state on `reg_dpu_start` or `ap_start` signals the start of a DPU task. At the end of the task, the DPU generates an interrupt and bit0 in IPISR and `reg_finish_sts` is set to 1.

To support DPU interrupt, the DPU implements the following registers:

- **Global Interrupt Enable Register (GIER):** Provides the master enable/disable for the interrupt output to the processor or Interrupt Controller. See Global Interrupt Enable Register (GIER) in [Table 2](#) for more details.
- **IP Interrupt Enable Register (IPIER):** Implements the independent interrupt enable bit for each channel. See IP Interrupt Enable (IPIER) and IP Status Registers (IPISR) in [Table 2](#) for more details.
- **IP Interrupt Status Register (IPISR):** Implements the independent interrupt status bit for each channel. The IPISR provides Read and Toggle-On-Write access. The Toggle-On-Write mechanism allows interrupt service routines to clear one or more ISR bits using a single write transaction. The IPISR can also be manually set to generate an interrupt for testing purposes. See IP Interrupt Enable (IPIER) and IP Status Registers (IPISR) in [Table 1](#) for additional details.

DPU Configuration

There is an option to determine the number of DPU engines that will be instantiated in a single DPU IP. The deep neural network features and the associated parameters supported by the DPU are shown in the following table.

Table 5: Deep Neural Network Features and Parameters Supported by DPU

Features	Description (channel_parallel=16)	
conv2d	Kernel Sizes	kernel_w: [1, 16] kernel_h: [1, 16]
	Strides	stride_w: [1, 4] stride_h: [1, 4]
	Pad_left/Pad_right	[0, (kernel_w - 1) * dilation_w + 1]
	Pad_top/Pad_bottom	[0, (kernel_h - 1) * dilation_h + 1]
	In Size	kernel_w * kernel_h * ceil(input_channel / channel_parallel) <= 2048
	Out Size	output_channel <= 256 * channel_parallel
	Activation	ReLU, LeakyReLU, ReLU6
	Dilation	dilation * input_channel <= 256 * channel_parallel
transposed-conv2d	Kernel Sizes	kernel_w: [1, 16] kernel_h: [1, 16]
	Strides	kernel_w: [1, 16] kernel_h: [1, 16]
	Pad_left/Pad_right	[1, kernel_w-1]
	Pad_top/Pad_bottom	[1, kernel_h-1]
	Out Size	output_channel <= 256 * channel_parallel
	Activation	ReLU, LeakyReLU, ReLU6
max-pooling	Strides	kernel_w: [1, 8] kernel_h: [1, 8]
	Kernel Sizes	W: [1, 8] H: [1, 8]
	Pad_left/Pad_right	[1, kernel_w-1]
	Pad_top/Pad_bottom	[1, kernel_h-1]
average-pooling	Kernel Sizes	kernel_w: [1, 8] kernel_h: [1, 8] kernel_w==kernel_h
	Strides	kernel_w: [1, 8] kernel_h: [1, 8]
	Pad_left/Pad_right	[1, kernel_w-1]
	Pad_top/Pad_bottom	[1, kernel_h-1]

Table 5: Deep Neural Network Features and Parameters Supported by DPU (cont'd)

Features	Description (channel_parallel=16)	
elementwise-sum	Input channel	input_channel <= 256 * channel_parallel
	Activation	ReLU
Fully Connected	Input Channel	input_channel <= 16*16*16
Concat	Network-specific limitation related to the size of feature maps, quantization results, and compiler optimizations.	

Configuration Options

The DPUCAHX8H can be configured with some predefined options which include the number of processing engines, frequency, and card type. These options are used to choose how many engines to implement and the frequency of the IP.

Number of Processing Engines

A range of one to five engines can be selected in one DPU IP. Multiple DPU engines can be used to achieve higher performance at the cost of increased resource utilization.

Resource Utilization

The resource utilization of the DPUCAHX8H with different configurations on the Alveo U50LV card is shown in the following table.

Table 6: DPUCAHX8H Utilization

Architecture	LUT	Register	Block RAM	UltraRAM	DSP
DPUCAHX8H 4PE	193675	338966	143.5	256	2080
DPUCAHX8H 5PE	238433	420289	147.5	320	2598

Performance

The following table shows the peak performance of the DPU on different devices.

Table 7: DPU_EU Performance (GOPS) on Different Device

Device	DPU Configuration	Frequency (MHz)	Peak Performance
U50	2 cores (3 ENGINES + 3 ENGINES)	300M	7373
U50LV	2 cores (5 ENGINES + 5 ENGINES)	275M	11264
U280	3 Cores (4 ENGINES + 5 ENGINES + 5 ENGINES)	300M	17203

Development Flow

Customizing and Generating the Core in Shell Mode with Vitis Flow

The following sections describe the development flow on how to use the DPU IP with the Vitis™ IDE.

Generating .xo Files

The .xo file is a format of IP that can be used by the shell in the Vitis flow. DPU IP files are released as .xo files.

1. Download the DPUCAHX8H_xo_gen_flow.tar package to generate DPU .xo files. The package includes DPU-related encrypted RTL and timing constraint files.
2. cd to the path: DPU_v3e_xo_gen_flow.
3. Run the following command to generate the DPU .xo files with required options.

```
vivado -mode tcl -source gen_DPUCAHX8H_ENGINE_xo.tcl -tclargs  
[1-5]ENGINE FREQ card(u50/u50lv/u280)
```

For example, to generate 5ENGINE 275M xo for the U50LV implementation, run the following command:

```
vivado -mode tcl -source gen_DPUCAHX8H_ENGINE_xo.tcl -tclargs 5ENGINE  
275 u50lv
```

The interface of the DPU xo is as following figures (3ENGINE, 4ENGINE, 5ENGINE). Each ENGINE uses one AXI port, so the DPU_AXI_* port number is the same as ENGINE number of the xo.

Figure 7: Ports of DPU of 3ENGINE



Figure 8: Ports of DPU of 4ENGINE

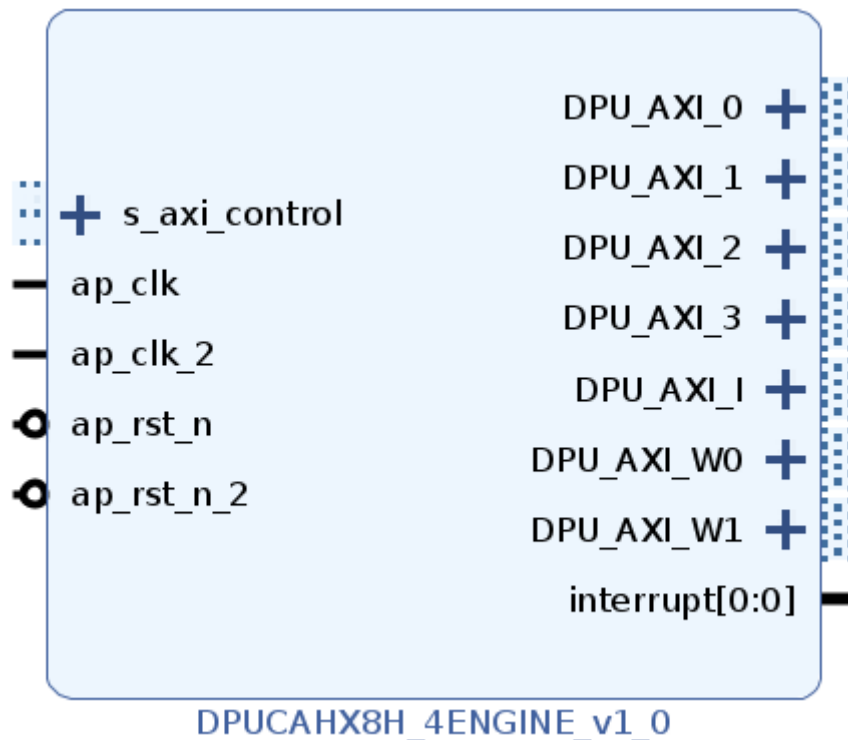
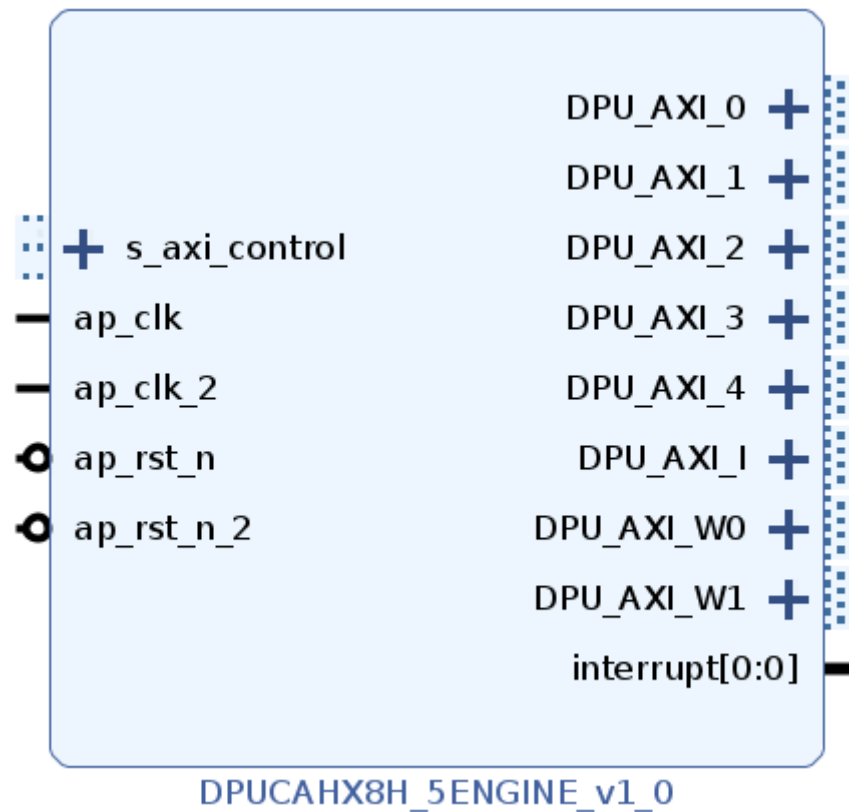


Figure 9: Ports of DPU of 5ENGINE



Adding the DPU xo Files to the Implementation Flow

To add the DPU xo files to the implementation flow, follow these steps:

Note: 5ENGINE at SLR0+ 5 ENGINE at SLR1 is used in this example.

1. Add following constraints into `cons.ini` file:

```
[connectivity]
nk=DPUCAHX8H_5ENGINE:1:dpu_0
nk=DPUCAHX8H_5ENGINE:1:dpu_1
```

2. Run the `v++` command with the `--config "cons.ini"` option.

Configuring the HBM

The Alveo™ U50, U50LV, and U280 cards support HBM. This section describes how to customize the HBM IP to meet DPU requirements and improve performance.

Although each AXI port connected to the HBM controller (HMSS) can access all the DDR memory banks within the HBM, the efficiency to access the DDR memory bank directly connected or within one switch is much higher than to access the DDR memory bank which is far from the AXI port. To get better access performance, each feature map AXI port must be constrained at a DDR memory bank to avoid latency caused by inefficient horizontal access.

For the DPUCAHX8L_A, 0~2 GB of space should not be used by your kernel. For the DPUCAHX8L_B, it is 2~4 GB.

To configure the HBM, follow these steps:

1. Add the sp constraints to the `cons.ini` file. In this example, the connectivity of each AXI port covers full HBM-range access.

```
[connectivity]
nk=DPUCAHX8H_5ENGINE:1:dpu_0
nk=DPUCAHX8H_5ENGINE:1:dpu_1
sp=dpu_0.DPU_AXI_0:HBM[00:31]
sp=dpu_0.DPU_AXI_1:HBM[00:31]
sp=dpu_0.DPU_AXI_4:HBM[00:31]
sp=dpu_1.DPU_AXI_0:HBM[00:31]
sp=dpu_1.DPU_AXI_1:HBM[00:31]
sp=dpu_1.DPU_AXI_4:HBM[00:31]
sp=dpu_0.DPU_AXI_I0:HBM[00:31]
sp=dpu_1.DPU_AXI_I0:HBM[00:31]
sp=dpu_0.DPU_AXI_2:HBM[00:31]
sp=dpu_0.DPU_AXI_3:HBM[00:31]
sp=dpu_1.DPU_AXI_2:HBM[00:31]
sp=dpu_1.DPU_AXI_3:HBM[00:31]
sp=dpu_0.DPU_AXI_W0:HBM[00:31]
sp=dpu_0.DPU_AXI_W1:HBM[00:31]
sp=dpu_1.DPU_AXI_W0:HBM[00:31]
sp=dpu_1.DPU_AXI_W1:HBM[00:31]
```

2. Use the `--config "cons.ini"` command to the `v++` command.
3. Put the following code at `sys_link_post.tcl` for one core:

```
hbm_memory_subsystem::force_host_port 28 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_0] 0 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_1] 1 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_4] 2 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_0] 3 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_1] 4 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_4] 5 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_I0] 6 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_I0] 7 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_2] 16 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_3] 17 1 [get_bd_cells hmss_0]
```



```
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_2] 18 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_3] 19 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_W0] 20 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_0/
DPU_AXI_W1] 21 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_W0] 22 1 [get_bd_cells hmss_0]
hbm_memory_subsystem::force_kernel_port [get_bd_intf_pins /dpu_1/
DPU_AXI_W1] 23 1 [get_bd_cells hmss_0]
```

4. Add the following lines to the `cons.ini` file.

```
[advanced]
param=compiler.userPostSysLinkTcl=*/sys_link_post.tcl
```

Making HBM Connections

DPUCAHX8H can be deployed with one or two cores on the Alveo U50/U50LV card or with one, two, or three cores on the Alveo U280 card.. The two core are named core0, core 1, and core2. You can locate core0 on SLR0, core1 on SLR1, and core2 on SLR2. To get the best performance of HBM, DPU AXI ports should be connected with the following rules.

Table 8: core0 Connection with HBM

DPU AXI Interfaces	HBM AXI Interfaces
DPU_AXI_0	AXI_00
DPU_AXI_1	AXI_01
DPU_AXI_2	AXI_02
DPU_AXI_3	AXI_03 (unused if PE number <4)
DPU_AXI_4	AXI_04 (unused if PE number <5)
DPU_AXI_I	AXI_05
DPU_AXI_W0	AXI_06
DPU_AXI_W1	AXI_07

Table 9: core1 Connection with HBM

DPU AXI Interfaces	HBM AXI Interfaces
DPU_AXI_0	AXI_08
DPU_AXI_1	AXI_09
DPU_AXI_2	AXI_10
DPU_AXI_3	AXI_11 (unused if PE number <4)
DPU_AXI_4	AXI_12 (unused if PE number <5)
DPU_AXI_I	AXI_13
DPU_AXI_W0	AXI_14
DPU_AXI_W1	AXI_15

Table 10: core2 Connection with HBM

DPU AXI Interfaces	HBM AXI Interfaces
DPU_AXI_0	AXI_16
DPU_AXI_1	AXI_17
DPU_AXI_2	AXI_18
DPU_AXI_3	AXI_19 (unused if PE number <4)
DPU_AXI_4	AXI_20 (unused if PE number <5)
DPU_AXI_I	AXI_21
DPU_AXI_W0	AXI_22
DPU_AXI_W1	AXI_23

Generating the Bitstream

Run the following command to generate the bitstream:

```
v++ -t hw --platform $platform_name --temp_dir $direction_name -l --config
"cons.ini" -o dpu.xclbin DPUCAHX8H_*ENGINE.xo
```

If more than one xo is used, add the xo file names in the command as shown in the following example for implementing the 4ENGINE at SLR0+ 5 ENGINE at SLR1 on the Alveo U50LV card:

```
v++ -t hw --platform xilinx_u50lv_gen3x4_xdma_2_202010_1 --save-temps --
temp_dir imp_dir -l --config "deephini" -o dpu.xclbin
DPUCAHX8H_4ENGINE.xo DPUCAHX8H_5ENGINE.xo
```

Upgrading

This appendix is not applicable for the first release of the core.

Additional Resources and Legal Notices

Xilinx Resources

For support resources such as Answers, Documentation, Downloads, and Forums, see [Xilinx Support](#).

Documentation Navigator and Design Hubs

Xilinx[®] Documentation Navigator (DocNav) provides access to Xilinx documents, videos, and support resources, which you can filter and search to find information. To open DocNav:

- From the Vivado[®] IDE, select **Help** → **Documentation and Tutorials**.
- On Windows, select **Start** → **All Programs** → **Xilinx Design Tools** → **DocNav**.
- At the Linux command prompt, enter `docnav`.

Xilinx Design Hubs provide links to documentation organized by design tasks and other topics, which you can use to learn key concepts and address frequently asked questions. To access the Design Hubs:

- In DocNav, click the **Design Hubs View** tab.
- On the Xilinx website, see the [Design Hubs](#) page.

Note: For more information on DocNav, see the [Documentation Navigator](#) page on the Xilinx website.

References

These documents provide supplemental material useful with this guide:

1. Vitis AI User Guide (UG1414)

Revision History

The following table shows the revision history for this document.

Section	Revision Summary
07/22/2021 Version 1.0	
Initial release.	N/A

Please Read: Important Legal Notices

The information disclosed to you hereunder (the "Materials") is provided solely for the selection and use of Xilinx products. To the maximum extent permitted by applicable law: (1) Materials are made available "AS IS" and with all faults, Xilinx hereby DISCLAIMS ALL WARRANTIES AND CONDITIONS, EXPRESS, IMPLIED, OR STATUTORY, INCLUDING BUT NOT LIMITED TO WARRANTIES OF MERCHANTABILITY, NON-INFRINGEMENT, OR FITNESS FOR ANY PARTICULAR PURPOSE; and (2) Xilinx shall not be liable (whether in contract or tort, including negligence, or under any other theory of liability) for any loss or damage of any kind or nature related to, arising under, or in connection with, the Materials (including your use of the Materials), including for any direct, indirect, special, incidental, or consequential loss or damage (including loss of data, profits, goodwill, or any type of loss or damage suffered as a result of any action brought by a third party) even if such damage or loss was reasonably foreseeable or Xilinx had been advised of the possibility of the same. Xilinx assumes no obligation to correct any errors contained in the Materials or to notify you of updates to the Materials or to product specifications. You may not reproduce, modify, distribute, or publicly display the Materials without prior written consent. Certain products are subject to the terms and conditions of Xilinx's limited warranty, please refer to Xilinx's Terms of Sale which can be viewed at <https://www.xilinx.com/legal.htm#tos>; IP cores may be subject to warranty and support terms contained in a license issued to you by Xilinx. Xilinx products are not designed or intended to be fail-safe or for use in any application requiring fail-safe performance; you assume sole risk and liability for use of Xilinx products in such critical applications, please refer to Xilinx's Terms of Sale which can be viewed at <https://www.xilinx.com/legal.htm#tos>.

AUTOMOTIVE APPLICATIONS DISCLAIMER

AUTOMOTIVE PRODUCTS (IDENTIFIED AS "XA" IN THE PART NUMBER) ARE NOT WARRANTED FOR USE IN THE DEPLOYMENT OF AIRBAGS OR FOR USE IN APPLICATIONS THAT AFFECT CONTROL OF A VEHICLE ("SAFETY APPLICATION") UNLESS THERE IS A SAFETY CONCEPT OR REDUNDANCY FEATURE CONSISTENT WITH THE ISO 26262 AUTOMOTIVE SAFETY STANDARD ("SAFETY DESIGN"). CUSTOMER SHALL, PRIOR TO USING OR DISTRIBUTING ANY SYSTEMS THAT INCORPORATE PRODUCTS, THOROUGHLY TEST SUCH SYSTEMS FOR SAFETY PURPOSES. USE OF PRODUCTS IN A SAFETY APPLICATION WITHOUT A SAFETY DESIGN IS FULLY AT THE RISK OF CUSTOMER, SUBJECT ONLY TO APPLICABLE LAWS AND REGULATIONS GOVERNING LIMITATIONS ON PRODUCT LIABILITY.

Copyright

© Copyright 2021 Xilinx, Inc. Xilinx, the Xilinx logo, Alveo, Artix, Kintex, Spartan, Versal, Virtex, Vivado, Zynq, and other designated brands included herein are trademarks of Xilinx in the United States and other countries. PCI, PCIe, and PCI Express are trademarks of PCI-SIG and used under license. All other trademarks are the property of their respective owners.