# NVMe IP Core

October 12, 2020

Product Specification

Г

Rev3.7



## **Design Gateway Co.,Ltd**

E-mail: ip-sales@design-gateway.com URL: www.design-gateway.com

## **Features**

- NVMe host controller for access one NVMe SSD without CPU and external memory
- Include 256-Kbyte RAM to be data buffer
- Simple user interface by dgIF typeS
- Support six commands, i.e. Identify, Shutdown, Write, Read, SMART and Flush
- Supported NVMe device
  - Base Class Code:01h (mass storage), Sub Class Code:08h (Non-volatile), Programming Interface:02h (NVMHCI)
  - MPSMIN (Memory Page Size Minimum): 0 (4Kbyte)
  - MDTS (Maximum Data Transfer Size): At least 5 (128 Kbyte) or 0 (no limitation)
  - LBA unit: 512 bytes or 4096 bytes
- User clock frequency must be more than or equal to PCIe clock (125MHz for Gen2, 250MHz for Gen3)
- Operating with Integrated Block for PCI Express from Xilinx by using 4-lane PCIe Gen2 or Gen3 (128bit bus interface)
- One NVMe IP connects to one NVMe SSD directly
- Available reference design
  - 1-ch demo on AC701, KC705, VC707, VC709, ZC706, Zynq Mini-ITX, KCU105, KCU116, ZCU106, and VCU118 with AB18-PCIeX16, AB17-M2FMC, and AB16-PCIeXOVR adapter board
  - 1-ch with DDR for transferring as sustain rate on KCU105
  - 2-ch/4-ch RAID0 demo on KCU105 and VCU118
- Customized service for following features
  - Additional NVMe commands
  - RAM size or RAM type (URAM) modification

### Table 1: Example Implementation Statistics for 7-Series device (PCIe Gen2/Gen3)

Family	Example Device	Fmax (MHz)	Slice Regs	Slice LUTs	Slices	ЮВ	BRAMTile <sup>1</sup>	PLL	GTX	Design Tools
Virtex-7	XC7VX690TFFG1761-2	300	4549	3204	1591	-	66	-	-	Vivado2017.4
Virtex-7	XC7VX485TFFG1761-2	300	4549	3203	1638	-	66	-	-	Vivado2017.4
Zynq-7000	XC7Z045FFG900-2	300	4549	3201	1532	-	66	-	-	Vivado2017.4
Kintex-7	XC7K325TFFG900-2	300	4549	3204	1574	-	66	-	-	Vivado2017.4
Artix-7	XC7A200TFBG676-2	225	4549	3089	1500	-	66	-	-	Vivado2017.4

Core Facts						
Pr	ovided with Core					
Documentation	Reference Design Manual					
	Demo Instruction Manual					
Design File Formats	Encrypted Netlist					
Instantiation Templates	VHDL					
Reference Designs &	Vivado Project,					
Application Notes	See Reference Design Manual					
Additional Items	Demo on AC701, KC705, VC707,					
	VC709, ZC706, Zynq Mini-ITX,					
	KCU105, KCU116,					
	ZCU106, VCU118					
Support						

Support Provided by Design Gateway Co., Ltd.

Family	Example Device	Fmax (MHz)	CLB Regs	CLB LUTs	CLB	юв	BRAMTile <sup>1</sup>	PLL	GTX	Design Tools
Kintex-Ultrascale	XCKU040FFVA1156-2E	375	4165	2718	814	-	66	-	-	Vivado2017.4
Kintex-Ultrascale+	XCKU5P-FFVB676-2E	400	4165	2737	834	-	66	-	-	Vivado2017.4
Zynq-Ultrascale+	XCZU7EV-FFVC1156-2E	400	4165	2732	880	-	66	-	-	Vivado2017.4
Virtex-Ultrascale+	XCVU9P-FLGA2104-2L	400	4165	2734	815	-	66	-	-	Vivado2017.4

<b>Table 2: Example Implementation Statistics</b>	s for Ultrascale/Ultrascale+ device (	PCle Gen3)
---	---------------------------------------	------------

# Applications



Figure 1: NVMe IP Application

NVMe IP Core integrated with Integrated Block for PCI Express (PCIe hard IP) from Xilinx is ideal to access NVMe SSD without CPU and external memory such as DDR. NVMe IP with PCIe hard IP can connect to one NVMe SSD. When transfer speed must be increased, multiple NVMe IPs and PCIe hard IPs can be applied to connect as RAID0 system as shown in Figure 1. By using two SSDs, transfer speed of RAID0 is double. In video data stream recording or replaying system, RAID0 system is applied to match with data stream bandwidth. For more specific applications, we also provides alternative IP as follows.

**Random Access NVMe IP Core**– To access NVMe SSD with multiple commands, individual address for each. Recommended for the application which requires the access in non-contiguous area, stream data recording without specified length or application that one SSD can be accessed by multiple users.

https://dgway.com/raNVMe-IP\_X\_E.html

NVMe IP Core for PCIe Switch- For access one or multiple NVMe SSDs via PCIe switch.

https://dgway.com/NVMe-IP X E.html

**NVMe IP Core with PCIe Gen 3 Soft IP** – When the selected FPGA does not have enough PCIe hard IP for the application.

NVMe IP Core with PCIe Gen 4 Soft IP - To access with PCIe Gen4 NVMe SSD.

https://dgway.com/NVMeG4-IP\_X\_E.html

### **General Description**



Figure 2: NVMe IP Block Diagram

NVMe IP implements as host controller to access NVMe SSD following NVM express standard. Physical interface of NVMe SSD is PCIe. The lower layer hardware is implemented by using Integrated Block for PCI Express (PCIe hard IP) from Xilinx.

NVMe IP supports six NVMe commands, i.e. Identify, Shutdown, Write, Read, SMART and Flush command by using two user interface groups. First is Control interface for transferring command and the parameters. Another is Data interface for transferring data when the command must have the data transferring. The Control interface and Data interface for Write/Read command use dgIF typeS format. Control interface of dgIF typeS has start address and transfer length with asserting the request signal while Data interface of dgIF typeS is the FIFO interface.

SMART and Flush command require the specific interface, called Custom command interface, which consists of Ctm I/F for control path and Ctm RAM I/F for data path. Furthermore, Identify command has its own data interface, named Iden RAM I/F, as shown in Figure 2.

During initialization process or running some commands, error signal may be asserted by NVMe IP if some abnormal conditions are found. The IP includes the error status to check the more details of error condition. To recover error status, NVMe IP and SSD must be reset.

There is one limitation about clock frequency of user logic. Transmit packet to PCIe hard IP must be sent continuously until end of packet. Therefore, data must be valid every clock between start of packet and end of packet. To support this feature, user logic clock frequency must be more than or equal to PCIe clock frequency (250 MHz for PCIe Gen3, 125MHz for PCIe Gen2) to have the bandwidth of transmit logic higher than or equal to PCIe hard IP bandwidth.

The reference design on FPGA evaluation boards are available for evaluation before purchasing.

# **Functional Description**

Figure 3 shows operation flow of NVMe IP after IP reset is de-asserted.



Figure 3: NVMe IP Operation Flow

The operation of NVMe IP is described as follows.

- 1) IP waits until PCIe is ready by monitoring Linkup status from PCIe IP core.
- IP begins the initialization process by setup PCIe and NVMe registers. After that, the IP enters to the Idle state to wait new command request from user. If some errors are detected during setup process, the IP changes to the Inactive state.
- The 1<sup>st</sup> command from user must be Identify command (UserCmd="000") to update LBASize (disk capacity) and LBAMode (LBA unit=512 bytes or 4 Kbytes).
- 4) The last command before power down the system must be Shutdown command (UserCmd="001"). This command is recommended to guarantee SSD powered down in a good sequence. Without Shutdown command, Write data in SSD cannot be guaranteed. After finishing Shutdown command, NVMe IP and SSD enter to the Inactive state and do not receive any commands until the IP is reset.
- 5) For Write command (UserCmd="010"), the maximum data size for one command is 128 Kbytes. If total length from user is more than 128 Kbytes, repeat step 5a) 5b) for many times.
  - a) The IP waits until Write data, sent by user, is enough for one command (transfer size of one command in NVMe IP is 128 Kbytes, except the last loop which could be less than 128 Kbytes).
  - b) The IP sends Write command to SSD and then waits until the status is returned from SSD. The IP changes to the Idle state when total data is completely transferred. Otherwise, the IP goes back to step 5a) to send the next Write command.
- 6) Similar to Write command, Read command (UserCmd="011") must be sent to the SSD many times when total length from user is more than 128 Kbytes. As a result, step 6a) 6b) are repeated many times.
  - a) The IP waits until free space of data buffer in NVMe IP is enough for one command. If remaining transfer size is equal to zero, the IP skips to step c).
  - b) The IP sends Read command to SSD. After that, goes back to step 6a) to check the remaining transfer size and the free space of data buffer.
  - c) IP waits until all data are completely transferred from data buffer to user logic and then changes to the Idle state. Therefore, data buffer is empty after finishing Read command.
- 7) For SMART command (UserCmd="100"), 512-byte data is returned after finishing the operation.
  - a) IP sends Get Log Page command to read SMART/Health information from SSD.
  - b) 512-byte data is returned from the SSD. The IP forwards the data through Custom command RAM interface (CtmRamAddr=0x000 0x01F).
- 8) For Flush command (UserCmd="110"), there is no data transferring during the operation.
  - a) IP sends Flush command to SSD.
  - b) IP waits until SSD returns status to complete the operation.

To design NVMe host controller, NVMe IP implements two protocols, i.e. NVMe protocol for interfacing with user and PCIe protocol for interfacing with PCIe hard IP. The details of the hardware inside NVMe IP are described as follows.

#### NVMe

NVMe IP supports six commands, i.e. Identify, Write, Read, SMART, Flush and Shutdown command which can split in two command types, i.e. Admin command and NVM command. NVMe IP supports three Admin commands, i.e. Identify, Shutdown and SMART command and supports three NVM commands, i.e. Write, Read and Flush command. After finishing the command, the status returned from the SSD is latched to AdmCompStatus signal when running Admin command or IOCompStatus signal when running NVM command.

The parameters of Write or Read command are set by Control interface of dgIF typeS while the parameters of SMART or Flush command are set by CtmSubm of Ctm interface. Data interface for Write or Read command is transferred by FIFO interface, a part of dgIF typeS, which is finally connected with 256-Kbyte buffer inside the IP. The data interface of other commands has its own interface, i.e. Identify RAM for Identify command and Custom RAM for SMART command.

The details of each submodule are described as follows.

#### NVMe Host Controller

NVMe host controller is the main controller in NVMe IP. The operation is split into two phases. First is the initialization phase which is once run after the system is boot up for setting NVMe register inside the SSD. After finishing the initialization phase, the next phase is the command operation phase. The order of transmitted and received packet is controlled by NVMe host controller.

To operate the command, the parameters of each command are prepared in Command Parameter and then forwarded to AsyncCtrl. After finishing the command, the status packet returned from SSD is monitored by NVMe host controller to check if no error is found. When the command needs to transfer data such as Write, Read, SMART and Identify command, NVMe host controller must also handle with NVMe data controller.

#### • Command Parameter

This module is designed to prepare command packet for SSD and also decode status packet returned from SSD.

Typically, the command consists of 16 Dwords (1 Dword = 32-bit). When running Identify, Shutdown, Write or Read command, all 16 Dwords are created by Command parameter, following the user inputs on dgIF typeS. When running SMART and Flush command, all 16 Dwords are directly loaded via CtmSubmDW0-CtmSubmDW15 of Ctm interface.

#### Data Buffer

256-Kbyte simple dual port RAM is implemented by BlockRAM to be data buffer. The buffer stores data transferred between UserLogic and SSD during Write and Read command operation.

#### • NVMe Data Controller

This module is operated when the command must transfer the data, i.e. Identify, SMART, Write and Read command. There are three data interfaces for transferring with the SSD, i.e. FIFO interface with 256-Kbyte buffer when running Write or Read command, Custom command RAM interface when running SMART command or Identify RAM interface when running Identify command.

When running Write or Read command, the address of the data buffer is controlled by NVMe data controller.

### PCle

The PCIe standard is the outstanding lower layer protocol for very high speed application. The NVMe standard is the protocol which is run over PCIe protocol. In the initialization process, NVMe layer is setup after finishing PCIe layer setup. Two modules are designed to support PCIe protocol, i.e. PCIe controller and AsyncCtrl. More details of each module are described as follows.

#### PCle Controller

During initialization process, PCIe controller sets up PCIe environment of SSD via CFG interface. After that, PCIe packet is created or decoded via 128-bit Tx/Rx AXI4-Stream. The command and data packet from NVMe module are converted to be PCIe packet by PCIe controller. On the other hand, the received PCIe packet is decoded and converted to be NVMe packet for NVMe module by this module.

#### AsyncCtrl

AsyncCtrl includes asynchronous registers and buffers to support signal and data clockcrossing. Most logics in NVMe IP run on user clock domain while PCIe hard IP runs on PCIe clock domain. AXI4-stream interface of PCIe hard IP must transfer data of each packet continuously, so the user bandwidth must be greater than or equal to PCIe bandwidth by running at higher or the same clock frequency of PCIe clock.

# **User Logic**

This module could be designed by using small state machine to send the commands with assigning the parameters for each command. For Write/Read command, the command parameters are address and transfer size. Data interface for Write/Read command can directly connect to FIFO while data output from SMART and Identify command can directly connect to simple dual port RAM. RAM size depends on the data size transferring in each command, but data width of all commands is 128-bit. Data size of Identify command is 8-Kbyte while data size of SMART command is 512-byte.

# **Integrated Block for PCI Express**

Integrated Block for PCI Express is the hard IP provided by Xilinx for some Xilinx FPGAs. The maximum number of SSDs connecting to one FPGA device is limited by the number of PCIe hard IPs in FPGA device. By using NVMe IP, one PCIe hard IP can connect to one NVMe SSD. More details of PCIe hard IP are described in following document.

PG054: 7 Series FPGAs Integrated Block for PCI Express PG023: Virtex-7 FPGA Gen3 Integrated Block for PCI Express PG156: UltraScale Devices Gen3 Integrated Block for PCI Express

PG213: UltraScale+ Devices Integrated Block for PCI Express



Figure 4: Integrated Block for PCI Express (UltraScale+)

# Core I/O Signals

Descriptions of all signal I/O are provided in Table 3 - Table 5.

# Table 3: User logic I/O Signals (Synchronous to Clk signal)

Signal	Dir	Description
	-	Control I/F of dgIF typeS
RstB	In	Synchronous reset signal. Active low. De-assert to '1' when Clk signal is stable.
Clk	In	System clock for running NVMe IP. The frequency must be more than or equal to PCIeClk which
		is output from PCIe hard IP (250 MHz for PCIe Gen3, 125 MHz for PCIe Gen2).
UserCmd[2:0]	In	User Command. Valid when UserReq='1'.
		("000": Identify, "001": Shutdown, "010": Write SSD, "011": Read SSD,
		"100": SMART, "110": Flush, "101"/"111": Reserved)
UserAddr[47:0]	In	Start address to write/read SSD in 512-byte unit. Valid when UserReq='1'.
		In case LBA unit = 4 Kbyte, UserAddr[2:0] must be always set to "000" to align 4-Kbyte unit.
		In case LBA unit = 512 byte, it is recommended to set UserAddr[2:0]="000" to align 4-Kbyte size
		(SSD page size). Write/Read performance of most SSDs is reduced when start address is not
		aligned to page size.
UserLen[47:0]	In	Total transfer size to write/read SSD in 512-byte unit. Valid from 1 to (LBASize-UserAddr).
		In case LBA unit = 4 Kbyte, UserLen[2:0] must be always set to "000" to align 4-Kbyte unit.
		Valid when UserReq='1'.
UserReq	In	Assert to '1' to send the new command request and de-assert to '0' after IP starts the operation
		by asserting UserBusy to '1'. This signal can be asserted when the IP is Idle (UserBusy='0').
		Command parameters (UserCmd, UserAddr, UserLen and CtmSubmDW0-DW15) must be valid
		and stable during UserReq='1'. UserAddr and UserLen are inputs for Write/Read command while
		CtmSubmDW0-DW15 are inputs for SMART/Flush command.
UserBusy	Out	Assert to '1' when IP is busy. New request must not be sent (UserReq to '1') when IP is busy.
LBASize[47:0]	Out	Total capacity of SSD in 512-byte unit. Default value is 0.
		This value is valid after finishing Identify command.
LBAMode	Out	LBA unit size of SSD ('0': 512 bytes, '1': 4 Kbytes). Default value is 0.
		This value is valid after finishing Identify command.
UserError	Out	Error flag. Assert to '1' when UserErrorType is not equal to 0.
		The flag can be cleared by asserting RstB to '0'.
UserErrorType[31:0]	Out	Error status.
		[0] – Error when PCIe class code is not correct.
		[1] – Error from CAP (Controller capabilities) register which may be caused from
		- MPSMIN (Memory Page Size Minimum) is not equal to 0.
		- NVM command set flag (bit 37 of CAP register) is not set to 1.
		- DSTRD (Doorbell Stride) is not 0.
		- MQES (Maximum Queue Entries Supported) is more than or equal to 7.
		More details of each register can be checked from NVMeCAPReg signal.
		[2] – Error when Admin completion entry is not received until timeout.
		[3] - Error when status register in Admin completion entry is not 0 or phase tag/command ID is
		invalid. Please see more details from AdmCompStatus signal.
		[4] – Error when IO completion entry is not received until timeout.
		[5] – Error when status register in IO completion entry is not 0 or phase tag is invalid. Please see
		more details from IOCompStatus signal.
		[6] – Error when Completion TLP packet size is not correct.

Signal	Dir	Description
		Control I/F of dgIF typeS
UserErrorType[31:0]	Out	[7] – Error when PCIe hard IP detects Error correction code (ECC) error from the internal buffer.
		[8] – Error from Unsupported Request (UR) flag in Completion TLP packet.
		[9] – Error from Completer Abort (CA) flag in Completion TLP packet.
		[15:10] – Reserved
		[16] - Error from unsupport LBA unit (LBA unit is not equal to 512 bytes or 4 Kbytes)
		[31:17] – Reserved
		Note: Timeout period of bit[2]/[4] is set from TimeOutSet input.
		Data I/F of dgIF typeS
UserFifoWrCnt[15:0]	In	Write data counter of Receive FIFO. Used to check full status.
		If FIFO size signal is less than 16 bits, please fill '1' to upper bit.
UserFifoWrEn	Out	Asserted to '1' to write data to Receive FIFO during running Read command.
UserFifoWrData[127:0]	Out	Write data bus of Receive FIFO. Valid when UserFifoWrEn='1'.
UserFifoRdCnt[15:0]	In	Read data counter of Transmit FIFO. Used to check data size stored in FIFO.
		If FIFO size signal is less than 16 bits, please fill '0' to upper bit.
UserFifoEmptv	In	The signal is unused for this IP.
UserFifoRdEn	Out	Asserted to '1' to read data from Transmit FIFO during running Write command.
UserFifoRdData[127:0]	In	Read data returned from Transmit FIFO.
		Valid in the next clock after UserFifoRdEn is asserted to '1'.
		NVMe IP Interface
IPVesion[31:0]	Out	IP version number
TestPin[31:0]	Out	Reserved to be IP Test point
TimeOutSet[31:0]	In	Timeout value to wait completion from SSD. Time unit is equal to 1/(Clk frequency).
		When TimeOutSet is set to 0. Timeout function is disabled.
PCIel inkup	In	Asserted to '1' when I TSSM state of PCIe hard IP is in L0 State
AdmCompStatus[15:0]	Out	Status output from Admin Completion Entry
/ amoompotatao[10.0]	out	[0] – Set to '1' when Phase tag or Command ID in Admin Completion Entry is invalid.
		[15:1] – Status field value of Admin Completion Entry
IOCompStatus[15:0]	Out	Status output from IO Completion Entry
	• • •	[0] – Set to '1' when Phase tag in IQ Completion Entry is invalid
		[15:1] – Status field value of IO Completion Entry
NVMeCAPReg[31:0]	Out	The parameter value of the NVMe capability register when UserErrorType[1] is asserted to '1'
	out	[15:0] – MOES (Maximum Queue Entries Supported)
		[19:16] – DSTRD (Doorbell Stride)
		[20] – NVM command set flag
		[24:21] – MPSMIN (Memory Page Size Minimum)
		[31:25] – Undefined
IdenWrEn	Out	Assert to '1' for sending data output from Identify command.
IdenWrDWEn[3:0]	Out	Dword (32-bit) enable of IdenWrData, Valid when IdenWrEn='1'.
		'1': this dword data is valid. '0': this dword data is not available.
		Bit[0], [1], [2] and [3] corresponds to IdenWrData[31:0], [63:32], [95:64] and [127:96]
		respectively.
IdenWrAddr[8:0]	Out	Index of IdenWrData in 128-bit unit. Valid when IdenWrEn='1'.
		0x000-0x0FF is 4Kbyte Identify controller data,
		0x100-0x1FF is 4Kbyte Identify namespace data.
IdenWrData[127:0]	Out	4Kbyte Identify controller data or Identify namespace data. Valid when IdenWrEn='1'.

Signal	Dir	Description			
	Custom interface				
CtmSubmDW0[31:0] -	In	16 Dwords of Submission queue entry for SMART/Flush command.			
CtmSubmDW15[31:0]		DW0: Command Dword0, DW1: Command Dword1,, and DW15: Command Dword15.			
		These inputs must be valid and stable during UserReq='1' and UserCmd="100" (SMART) or			
		"110" (Flush).			
CtmCompDW0[31:0] -	Out	4 Dwords of Completion queue entry, output from SMART/Flush command.			
CtmCompDW3[31:0]		DW0: Completion Dword0, DW1: Completion Dword1,, and DW3: Completion Dword3			
CtmRamWrEn	Out	Asserted to '1' for sending data output from custom command such as SMART command.			
CtmRamWrDWEn[3:0]	Out	Dword (32-bit) enable of CtmRamWrData. Valid when CtmRamWrEn='1'.			
		'1': this dword data is valid, '0': this dword data is not available.			
		Bit[0], [1], [2] and [3] corresponds to CtmRamWrData[31:0], [63:32], [95:64] and [127:96]			
		respectively.			
CtmRamAddr[8:0]	Out	Index of CtmRamWrData when SMART data is received. Valid when CtmRamWrEn='1'.			
		(Optional) Index to request data input through CtmRamRdData for customized custom			
		commands.			
CtmRamWrData[127:0]	Out	512-byte data output from SMART command. Valid when CtmRamWrEn='1'.			
CtmRamRdData[127:0]	In	(Optional) Data input for customized custom commands.			

# Table 4: Physical I/O Signals for PCIe Gen3 Hard IP (Synchronous to PCIeCIk)

Signal	Dir	Description
	Ť	PCIe Gen3 hard IP
PCIeRstB	In	Synchronous reset signal. Active low.
		De-assert to '1' when PCIe hard IP is not in reset state.
PCIeClk	In	Clock output from PCIe hard IP (250 MHz for PCIe Gen3).
		Configuration Management Interface
PCIeCfgDone	In	Read/Write operation complete. Asserted for 1 cycle when operation completes.
PCIeCfgRdEn	Out	Read enable. Asserted to '1' for a read operation.
PCIeCfgWrEn	Out	Write enable. Asserted to '1' for a write operation.
PCIeCfgWrData[31:0]	Out	Write data which is used to configure the Configuration and Management registers.
PCIeCfgByteEn[3:0]	Out	Byte enable for write data, where bit[0] corresponds to PCIeCfgWrData[7:0], and so on.
PCIeCfgAddr[18:0]	Out	Read/Write Address.
		Requester Request Interface
PCIeMtTxData[127:0]	Out	Requester request data bus.
PCIeMtTxKeep[3:0]	Out	Bit [i] indicates that Dword [i] of PCIeMtTxData contains valid data.
PCIeMtTxLast	Out	Asserts this signal in the last cycle of a TLP to indicate the end of the packet.
PCIeMtTxReady[3:0]	In	Asserts to accept data. Data is transferred when both PCIeMtTxValid and PCIeMtTxReady
		are asserted in the same cycle.
PCIeMtTxUser[59:0]	Out	Requester request user data. Valid when PCIeMtTxValid is high.
PCIeMtTxValid	Out	Asserts to drive valid data on PCIeMtTxData bus. NVMe IP keeps the valid signal asserted
		during the transfer of a packet.
	-	Completer Request Interface
PCIeMtRxData[127:0]	In	Receive data from PCIe hard IP.
PCIeMtRxKeep[3:0]	In	Bit [i] indicates that Dword [i] of PCIeMtRxData contains valid data.
PCIeMtRxLast	In	Asserts this signal in the last beat of a packet to indicate the end of the packet.
PCIeMtRxReady	Out	Indicates that NVMe IP is ready to accept data.
PCIeMtRxUser[74:0]	In	Sideband information for the TLP being transferred. Valid when PCIeMtRxValid is high.
PCIeMtRxValid	In	Assert when PCIe hard IP drives valid data on PCIeMtRxData bus. PCIe hard IP keeps the
		valid signal asserted during the transfer of packet.
	1	Completer Completion Interface
PCIeSITxData[127:0]	Out	Completion data from NVMe IP.
PCIeSITxKeep[3:0]	Out	Bit [i] indicates that Dword [i] of PCIeSITxData contains valid data.
PCIeSITxLast	Out	Asserts this signal in the last cycle of a packet to indicate the end of the packet.
PCIeSITxReady[3:0]	In	Indicates that PCIe hard IP is ready to accept data.
PCIeSITxUser[32:0]	Out	Sideband information for the TLP being transferred. Valid when PCIeSITxValid is high.
PCIeSITxValid	Out	Asserts to drive valid data on PCIeSITxData bus. NVMe IP keeps the valid signal asserted
		during the transfer of a packet.
	1	Requester Completion Interface
PCIeSIRxData[127:0]	In	Receive data from PCIe hard IP.
PCIeSIRxKeep[3:0]	In	Bit [i] indicates that Dword [i] of PCIeSIRxData contains valid data.
PCIeSIRxLast	In	Asserts this signal in the last beat of a packet to indicate the end of the packet.
PCIeSIRxReady	Out	Indicates that NVMe IP is ready to accept data.
PCIeSIRxUser[84:0]	In	Sideband information for the TLP being transferred. Valid when PCIeSIRxValid is high.
PCIeSIRxValid	In	Asserts when PCIe hard IP drives valid data on PCIeSIRxData bus. PCIe hard IP keeps the
		valid signal asserted during the transfer of packet.

# Table 5: Physical I/O Signals for PCIe Gen2 Hard IP (Synchronous to PCIeCIk)

Signal	Dir	Description
		PCIe Gen2 hard IP
PCleRstB	In	Synchronous reset signal. Active low.
		De-assert to '1' when PCIe hard IP is not in reset state.
PCIeClk	In	Clock output from PCIe hard IP (125 MHz for PCIe Gen2).
	_	Configuration Management Interface
PCIeCfgDone	In	Read/Write operation complete. Asserted for 1 cycle when operation completes.
PCIeCfgRdEn	Out	Read enable. Asserted to '1' for a read operation.
PCIeCfgWrEn	Out	Write enable. Asserted to '1' for a write operation.
PCIeCfgWrData[31:0]	Out	Write data which is used to configure the Configuration and Management registers.
PCIeCfgByteEn[3:0]	Out	Byte enable for write data, where bit[0] corresponds to PCIeCfgWrData[7:0], and so on.
PCIeCfgAddr[9:0]	Out	Read/Write Address.
		PCIe Transmit Interface
PCleTxData[127:0]	Out	Transmit data to PCIe hard IP.
PCIeTxKeep[15:0]	Out	Bit[i] indicates that Byte[i] of PCIeTxData contains valid data.
PCIeTxLast	Out	Asserts this signal in the last cycle of a TLP to indicate the end of the packet.
PCIeTxReady[3:0]	In	Asserts to accept data. Data is transferred when both PCIeTxValid and PCIeTxReady are
		asserted in the same cycle.
PCIeTxUser[3:0]	Out	Bit[3]: Transmit source discontinue.
		Bit[2]: Transmit streamed.
		Bit[1]: Transmit error forward.
		Bit[0]: Transmit ECRC Generate.
		Always set to 0000b.
PCIeTxValid	Out	Indicates that NVMe IP is presenting valid data on PCIeTxData.
	T	PCIe Receive Interface
PCIeRxData[127:0]	In	Receive data from PCIe hard IP. Valid only PCIeRxValid is also asserted.
PCIeRxKeep[15:0]	In	Bit[i] indicates that Byte[i] of PCIeRxData contains valid data.
PCIeRxLast	In	Receive End-of-Frame. Valid only if PCIeRxValid is also asserted.
PCIeRxReady	Out	Indicates that NVMe IP is ready to accept data.
PCIeRxUser[21:0]	In	Bit[0]: Receive ECRC error.
		Bit[14:13]: Indicates the start of a new packet header in PCIeRxData.
		Bit[21]: Indicates the end of a packet in PCIeRxData.
		Other bits are ignored in NVMe IP.
PCIeRxValid	In	Assert when PCIe hard IP drives valid data on PCIeRxData bus. PCIe hard IP keeps the valid
		signal asserted during the transfer of packet.

# **Timing Diagram**

### Initialization



Figure 5: Timing diagram during initialization process

The step of the initialization process is as follows.

- 1) Wait until Clk is stable and then de-asserts RstB to '1' to start IP initialization.
- 2) PCIe hard IP de-asserts PCIeRstB to '1' after finishing PCIe reset sequence. PCIe hard IP is ready to interface with the application layer.
- Assert PCIeLinkup to '1' after LTSSM state of PCIe hard IP is L0 state. Though LTSSM state is run on PCIeClk, PCIeLinkup must be generated on Clk domain. Asynchronous register must be applied to generate PCIeLinkup. After that, NVMe IP starts initialization process.
- 4) UserBusy is deasserted to '0' after NVMe IP completes initialization process.

After finishing above sequences, NVMe IP is ready to receive the command from user.

### Control interface of dgIF typeS

dgIF typeS signals are split into two groups. First group is control interface for sending command with the parameters and monitoring the status. Second group is data interface for transferring data stream in both directions.



### Figure 6: Control Interface of dgIF typeS timing diagram

- 1) Before sending new command to the IP, UserBusy must be equal to '0' to confirm that IP is the Idle state.
- 2) Command and the parameters such as UserCmd, UserAddr and UserLen must be valid when asserting UserReq to '1' for sending the new command request.
- 3) IP asserts UserBusy to '1' after starting the new command operation.
- 4) After UserBusy is asserted to '1', UserReq is de-asserted to '0' to finish the current request. New parameters for the next command could be prepared on the bus. UserReq for the new command must not be asserted to '1' until the current command operation is finished.
- 5) UserBusy is de-asserted to '0' after the command operation is completed. New command request could be sent by asserting UserReq to '1'.

Note: The number of parameters using in each command is different.

- Write and Read command: Use UserCmd, UserAddr and UserLen.
- SMART and Flush command: Use UserCmd and CtmSubmDW0-DW15.
- Identify and Shudown command: Use only UserCmd.

### Data interface of dgIF typeS

Data interface of dgIF typeS is applied for transferring data stream when operating Write command or Read command. The interface is compatible to general FIFO interface. 16-bit FIFO read data counter (UserFifoRdCnt) shows total data stored in the FIFO before transferring as a burst. The burst size is 512 bytes or 32 cycles of 128-bit data.



Figure 7: Transmit FIFO Interface for Write command

In Write command, data is read from Transmit FIFO until total data are transferred completely. The details to transfer data are described as follows.

- 1) Before starting a new burst transfer, UserFifoRdCnt[15:5] is monitored. The IP waits until at least 512-byte data is available in Transmit FIFO (UserFifoRdCnt[15:5] is not equal to 0).
- 2) The IP asserts UserFifoRdEn to '1' for 32 clock cycles to read 512-byte data from Transmit FIFO.
- 3) UserFifoRdData is valid in the next clock cycle after asserting UserFifoRdEn to '1'. 32 data are continuously transferred.
- 4) UserFifoRdEn is de-asserted to '0' after reading the 32<sup>th</sup> data. Repeat step 1) − 4) to transfer the next 512-byte until total data size is equal to the transfer size in the command.
- 5) After total data is completely transferred, UserBusy is de-asserted to '0'.



Figure 8: Receive FIFO Interface for Read command

In Read command, data is transferred from SSD to Receive FIFO until total data are completely transferred. The details to transfer data are as follows.

- 1) Before starting the new burst transmission, UserFifoWrCnt[15:6] is monitored. The IP waits until the free space of Receive FIFO is enough (UserFifoWrCnt[15:6] is not equal to all 1 or 1023). After received data from the SSD is more than or equal to 512 bytes, the new burst transmission begins.
- 2) The IP asserts UserFifoWrEn to '1' for 32 clock cycles to transfer 512-byte data from the data buffer to user logic.
- 3) After finishing transferring 512-byte data, UserFifoWrEn is de-asserted to '0'. Repeat step 1) 3) to transfer the next 512-byte data until total data size is equal to the transfer size in the command.
- 4) After total data is completely transferred, UserBusy is de-asserted to '0'.

### IdenCtrl/IdenName

It is recommened to send Identify command to the IP as the first command after system boots up. This command updates the necessary information of SSD, i.e. total capacity (LBASize) and LBA unit size (LBAMode). The SSD information is applied to be the limitation of the input parameters when operating Write or Read command, described as follows.

- 1) The sum of the address (UserAddr) and transfer length (UserLen) of Write and Read command must not be more than total capacity (LBASize) of the SSD.
- 2) If LBAMode of the active SSD is equal to '1' (LBA unit size is 4 Kbytes), the three lower bit (bit[2:0]) of UserAddr and UserLen must be always equal to '0' to align 4-Kbyte unit.



Figure 9: Identify command timing diagram

The details when running Identify command are shown as follows.

- 1) Send Identify command to the IP (UserCmd="000" and UserReq='1').
- 2) The IP asserts UserBusy to '1' after running Identify command.
- 3) 4096-byte Identify controller data is returned to user. IdenWrAddr is equal to 0-255 with asserting IdenWrEn. Also, IdenWrData and IdenWrDWEn are valid at the same clock as IdenWrEn='1'.
- 4) 4096-byte Identify namespace data is returned. IdenWrAddr is equal to 256-511. IdenWrAddr[8] can be applied to check data type which is Identify controller data or Identify namespace data.
- 5) UserBusy is de-asserted to '0' after finishing the Identify command.
- 6) LBASize and LBAMode of the SSD are simultaneously updated.



Figure 10: IdenWrDWEn timing diagram

IdenWrDWEn is 4-bit signal to be valid signal of 32-bit data. Some SSDs do not return 4-Kbyte Identify controller data and Identify namespace data continuously, but it returns only one dword (32-bit) at a time. Therefore, one bit of IdenWrDWEn is asserted to '1' in the write cycle to write 32-bit data, as shown in Figure 10. IdenWrDWEn[0], [1], [2] and [3] corresponds to IdenWrData[31:0], [63:32], [95:64] and [127:96] respectively.

### Shutdown

Shutdown command is recommended to send as the last command before the system is powered down. When Shutdown command is issued, SSD flushes the data from the internal cache to flash memory. After the command is issued, NVMe IP and SSD are not available until the system is powered down. If the SSD is powered down without Shutdown command, the total count of unsafe shutdowns, read by SMART command, is increased.



Figure 11: Shutdown command timing diagram

The details when running Shutdown command are shown as follows.

- 1) Before sending the command request, the IP must be in the Idle state (UserBusy='0'). To send Shutdown command, user asserts UserReq to '1' with UserCmd="001".
- 2) UserBusy is asserted to '1' after NVMe IP runs Shutdown command.
- 3) UserReq is de-asserted to '0' to clear the current request when UserBusy is asserted to '1'.
- 4) UserBusy is de-asserted to '0' when the SSD is completely shut down. After that, the IP does not receive any command requested from user.

### SMART

SMART command is the command to check the SSD health. After sending SMART command, 512-byte health information is returned from the SSD. SMART command loads the parameters from CtmSubmDW0-DW15 signals on Custom command interface. User sets 16-dword data as constant value for SMART command. After that, the SMART data is returned via CtmRAM port as shown in Figure 12.



Figure 12: SMART command timing diagram

The details when running SMART command are shown as follows.

- 1) Before sending the command request, the IP must be in the Idle state (UserBusy='0'). All input parameters must be stable when UserReq is asserted to '1' for sending the request. CtmSubmDW0-DW15 is set as constant value by following value for SMART command. CtmSubmDW0  $= 0 \times 0000 0002$ CtmSubmDW1 = 0xFFFF FFFF CtmSubmDW2 – CtmSubmDW5  $= 0 \times 0000 0000$ CtmSubmDW6 = 0x2000 0000CtmSubmDW7 – CtmSubmDW9  $= 0 \times 0000 0000$ CtmSubmDW10  $= 0 \times 007F 0002$ CtmSubmDW11 – CtmSubmDW15  $= 0x0000_{0000}$
- 2) Assert UserBusy to '1' after NVMe IP runs SMART command.
- 3) UserReq is de-asserted to '0' to clear the current request. Next, user logic can change the input parameters for the next command request.
- 4) 512-byte SMART data is returned on CtmRamWrData signal with asserting CtmRamWrEn to '1'. CtmRamAddr is equal to 0-31 to be data index of 512-byte data. When CtmRamAddr=0, byte0-15 of SMART data is valid on CtmRamWrData. CtmRamWrDWEn is dword enable for each 32-bit CtmRamWrData. If CtmRamWrDWEn="1111", all 128-bit CtmRamWrData are valid.
- 5) UserBusy is de-asserted to '0' when finishing SMART command.



Figure 13: CtmRamWrDWEn timing diagram

Similar to Identify command, some SSDs do not return 512-byte data continuously but returns only one dword (32-bit) at a time. Therefore, one bit of CtmRamWrDWEn is asserted to '1' in the write cycle to be the valid signal of 32-bit CtmRamWrData. CtmRamWrDWEn[0], [1], [2] and [3] corresponds to CtmRamWrData[31:0], [63:32], [95:64] and [127:96] respectively.

### Flush

Most SSDs accelerate write performance by storing write data to cache before flushing to the flash memory by the SSD controller. If power is down unexpectedly, the data in the cache may be lost and not stored to the flash memory. Flush command is the command to force the SSD controller to flush data from the cache. After sending Flush command, all data in previous Write command can be guaranteed.



Figure 14: Flush command timing diagram

The details for running Flush command are shown as follows.

 Before sending the command request, the IP must be in the Idle state (UserBusy='0'). All input parameters must be stable when UserReq is asserted to '1' for sending the request. CtmSubmDW0-DW15 is set as constant value by following value for Flush command. CtmSubmDW0 = 0x0000\_0000 CtmSubmDW1 = 0x0000\_0001

CtmSubmDW2 – CtmSubmDW15 = 0x0000\_0000

- 2) UserBusy is asserted to '1' after NVMe IP runs Flush command.
- 3) UserReq is de-asserted to '0' to clear the current request. Next, user logic can change the input parameters for the next command request.
- 4) UserBusy is de-asserted to '0' when finishing Flush command.





Figure 15: Error flag timing diagram

When the error is found during running initialization process or operating some commands, UserError flag is asserted to '1'. UserErrorType is read to check the error type. NVMeCAPReg, AdmCompStatus and IOCompStatus are valid for monitoring error details after UserError is asserted to '1'.

When the error is found during running initialization process, it is recommended to read NVMeCAPReg to check capability of NVMe SSD. When the error is found during operating the command, it is recommended to read AdmCompStatus or IOCompStatus, depending on the command type, which is decoded from the received packet.

UserError flag is cleared by RstB signal only. After the failure is solved, RstB is asserted to '0' to clear the error flag.

### **Verification Methods**

The NVMe IP Core functionality was verified by simulation and also proved on real board design by using AC701/KC705/VC707/VC709/ZC706/Zynq Mini-ITX/KCU105/ZCU106/VCU118 evaluation board.

### **Recommended Design Experience**

Experience design engineers with a knowledge of Vivado Tools should easily integrate this IP into their design.

## **Ordering Information**

This product is available directly from Design Gateway Co., Ltd. Please contact Design Gateway Co., Ltd. for pricing and additional information about this product using the contact information on the front page of this datasheet.

## **Revision History**

Revision	Date	Description
1.0	2-Jun-2016	Initial Release
1.1	15-Jun-2016	Add KCU105 support
1.2	5-Sep-2016	Add ZC706 support
1.3	9-Sep-2016	Add KC705 support
1.4	27-Oct-2016	Support VC709 and Zynq Mini-ITX
1.5	9-Dec-2016	Modify buffer to be BRAM and modify user interface to dgIF typeS
2.0	7-Jun-2017	Change PCIe interface to connect to PCIe hard IP
2.1	30-Nov-2017	Support ZCU106 board
3.0	19-July-2018	Support Shutdown, SMART, and Flush command
3.1	10-Oct-2018	Support Memory Write Request which is not aligned to 128-bit unit
3.2	27-Nov-2018	Add AC701 support
3.3	2-May-2019	Add AB18-PCIeX16 support
3.4	28-Jan-2020	Add AB17-M2FMC adaptor board support for KCU105 and ZCU106
3.5	26-Jun-2020	Add alternative IP description
3.6	28-Aug-2020	Add KCU116 support and update company info
3.7	12-Oct-2020	Correct information