

NVMe IP Core

July 17, 2018

Product Specification

Rev3.0



Design Gateway Co.,Ltd

54 BB Building 14th Fl., Room No.1402 Sukhumvit
21 Rd. (Asoke), Klongtoey-Nua, Wattana,
Bangkok 10110
Phone: 66(0)2-261-2277
Fax: 66(0)2-261-2290
E-mail: ip-sales@design-gateway.com
URL: www.design-gateway.com

Features

- Implement application layer to access NVMe SSD without CPU usage
- Simple user interface by dgIF typeS
- User clock frequency must be more than or equal to PCIe clock (125 MHz for PCIe Gen2, 250 MHz for PCIe Gen3)
- Direct connect to Integrated Block for PCI Express from Xilinx by using 128-bit bus interface
- Include 256 Kbyte RAM to be data buffer
- Support six commands, i.e. Identify, Shutdown, Write, Read, SMART, and Flush
- Additional command support as optional
- 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 byte or 4096 byte
- Reference design
 - KC705, VC707, VC709, ZC706, KCU105, ZCU106 with AB16-PCIeXOVR adapter board
 - VCU118 with AB17-M2FMC adapter board
 - Zynq Mini-ITX

Core Facts	
Provided with Core	
Documentation	Reference Design Manual Demo Instruction Manual
Design File Formats	Encrypted Netlist
Constraints Files	Constraint file example in reference design project
Instantiation Templates	VHDL
Reference Designs & Application Notes	Vivado Project, See Reference Design Manual
Additional Items	Demo on KC705, VC707, VC709, ZC706, Zynq Mini-ITX, KCU105, ZCU106, VCU118
Support	
Support Provided by Design Gateway Co., Ltd.	

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

Family	Example Device	Fmax (MHz)	Slice Regs	Slice LUTs	Slices	IOB	BRAMTile ¹	PLL	GTX	Design Tools
Virtex-7	XC7VX690TFFG1761-2	300	4169	2828	1423	-	59	-	-	Vivado2017.4
Virtex-7	XC7VX485TFFG1761-2	300	4159	3465	1446	-	59	-	-	Vivado2017.4
Zynq-7000	XC7Z045FFG900-2	300	4159	3461	1506	-	59	-	-	Vivado2017.4
Kintex-7	XC7K325TFFG900-2	300	4159	3463	1556	-	59	-	-	Vivado2017.4

July 17, 2018

Table 2: Example Implementation Statistics for Ultrascale/Ultrascale+ device (PCIe Gen3)

Family	Example Device	Fmax (MHz)	CLB Regs	CLB LUTs	CLB	IOB	BRAMTile ¹	PLL	GTX	Design Tools
Kintex-Ultrascale	XCKU040FFVA1156-2E	400	4170	2724	772	-	59	-	-	Vivado2017.4
Zynq-Ultrascale+	XCZU7EV-FFVC1156-2E	400	4170	2670	790	-	59	-	-	Vivado2017.4
Virtex-Ultrascale+	XCVU9P-FLGA2104-2L	400	4170	2675	761	-	59	-	-	Vivado2017.4

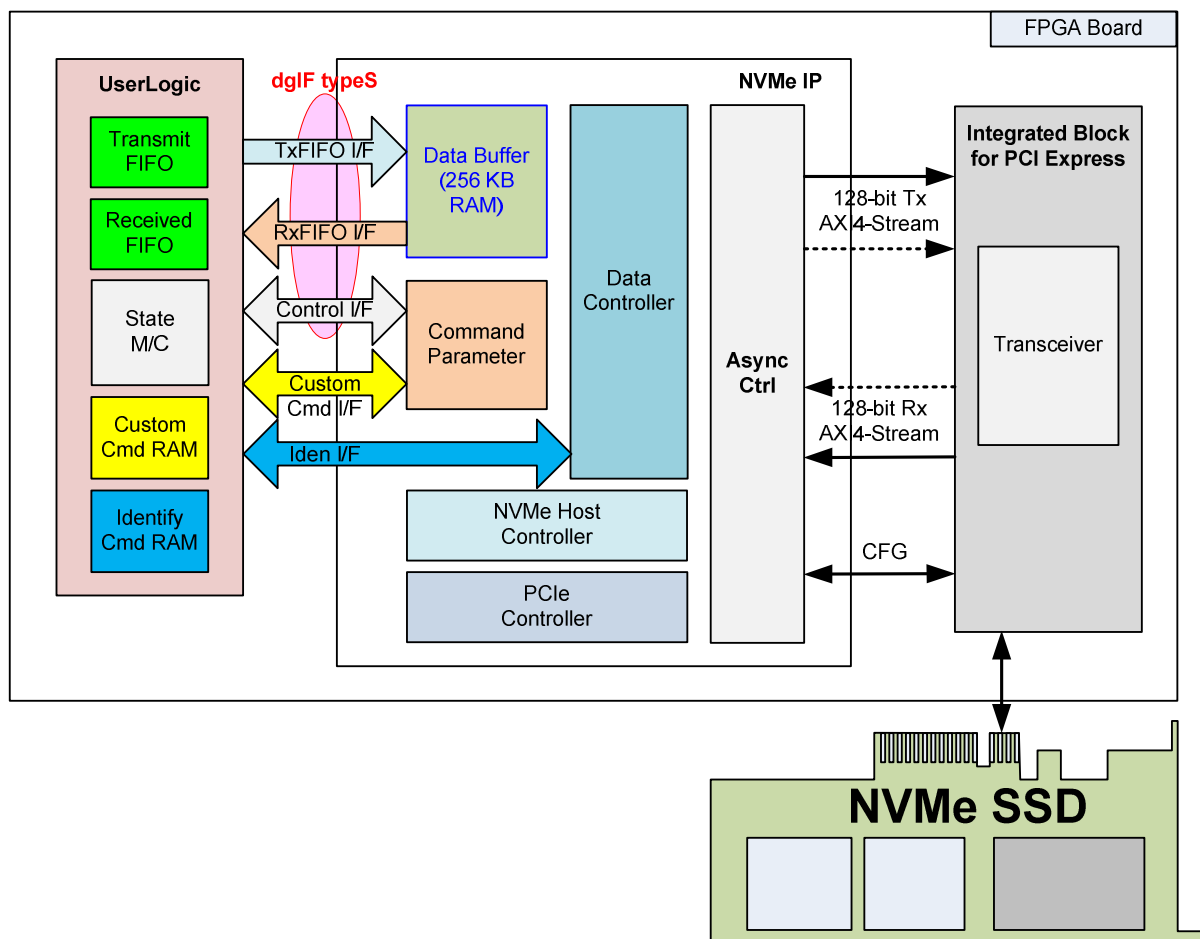


Figure 1: NVMe IP Block Diagram

Applications

NVMe IP Core integrated with Integrated Block for PCI Express from Xilinx is ideal to access NVMe SSD without CPU or DDR. 256 Kbyte buffer implemented by BlockRAM is included in NVMe IP Core to be data buffer between user logic and NVMe SSD. This solution is recommended for the application which requires high capacity storage at ultra high speed performance. Small size system can be designed by using M.2 SSD package.

General Description

NVMe IP implements as host controller to access NVMe SSD following NVM express standard. Physical interface of NVMe SSD is PCIe, so the hardware of lower layer is implemented by using Integrated Block for PCI Express from Xilinx. NVMe IP supports six NVMe commands, i.e. Identify, Shutdown, Write, Read, SMART, and Flush command.

256 Kbyte RAM is included in NVMe IP to be data buffer between User Logic and NVMe SSD when operating Write and Read command. By using big buffer size, the system achieves good performance to write and read SSD.

The user interface of NVMe IP is simply designed by using dgIF typeS. dgIF typeS consists of two interfaces, i.e. command interface and data interface. Input parameters of Write and Read command (received through dgIF typeS) are start address and transfer length. Data stream of Write and Read command are transferred through TxFIFO and RxFIFO interface of dgIF typeS.

To support SMART and Flush command, Custom command interface is additional designed to receive command parameters. Also, SMART data is transferred through Custom command RAM interface.

From Integrated Block for PCI Express limitation, clock frequency of user logic must be more than or equal to PCIe clock frequency (250 MHz for PCIe Gen3, 125 MHz for PCIe Gen2). Error signal will be asserted with the error status if IP detects the abnormal condition during packet transferring.

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

Functional Description

Figure 2 shows operation sequence of NVMe IP after IP reset is released.

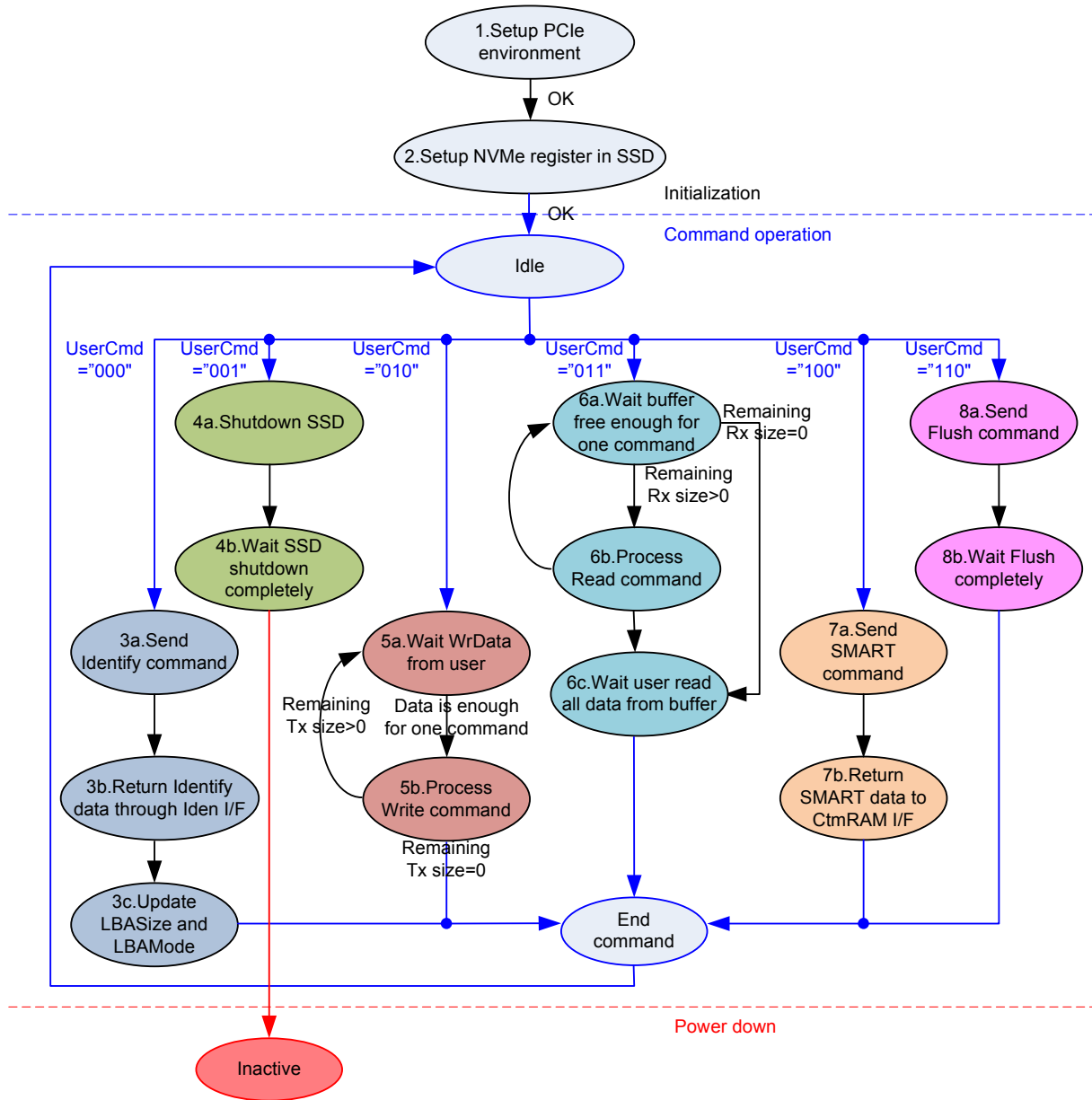


Figure 2: NVMe IP Operation Flow

The operation of NVMe IP is described as follows.

- 1) IP sets up PCIe environment for NVMe operation.
- 2) IP sets up NVMe parameters in SSD. After complete initialization process, IP is in idle state to wait new command request from user.
- 3) The 1st command from user must be Identify command (UserCmd="000") to update LBASize (disk capacity) and LBAMode (LBA unit=512 byte or 4 Kbyte).
- 4) The last command before power down the system must be Shutdown command (UserCmd="001"). This command is used to shutdown SSD in good sequence. Without Shutdown command, Write data in SSD cannot be guaranteed. After running Shutdown command, NVMe IP and SSD are Inactive state.
- 5) In case of Write command (UserCmd="010"),
 - IP waits until user sends write data to NVMe IP is much enough for one command (Transfer size of one command in NVMe IP is 128 Kbyte or less).
 - IP sends Write command to NVMe SSD.
 - IP waits status from SSD to confirm that all data in one command have been received completely. Remaining transfer size is decreased by the transmit size.
 - If remaining transfer size is not zero, IP will continue to check write data size from user in NVMe IP. If data size is much enough, the next command will be sent.
 - If remaining transfer size is zero, IP will go back to Idle state.
- 6) In case of Read command (UserCmd="011"),
 - IP checks free space of data buffer in NVMe IP.
 - If free space is much enough for one command, IP will send Read command to SSD.
 - IP waits until received size is equal to transfer size in one command. After that, the data is ready for user to read. Remaining transfer size is decreased by the received size.
 - If remaining transfer size is not zero, IP will check free space size for sending next command.
 - If remaining transfer size is zero, IP will go back to Idle state.
- 7) In case of SMART command (UserCmd="100"),
 - IP sends Get Log Page command to read SMART/Health information from SSD.
 - 512-byte data is returned from SSD and forwarded to user through Custom Command RAM interface (CtmRamAddr=0x000 – 0x01F).
- 8) In case of Flush command (UserCmd="110"),
 - IP sends Flush command to SSD.
 - IP waits until SSD returns status to complete the operation.

From above sequences, NVMe IP consists of three controllers, i.e. PCIe controller, NVMe host controller, and Data controller. After system power-on, PCIe controller sets up PCIe environment for connecting with SSD. Next, NVMe host controller writes and reads NVMe register of SSD following NVMe specification to complete initialization phase.

NVMe host controller controls the packet sequence of each user command following NVMe standard. The input parameters from user are received to create Command packet. Data controller is responsible to create and decode all packets transferring with SSD.

More details of each module in NVMe IP are described as follows.

PCIe

NVMe protocol uses PCIe standard to be physical interface and lower layer protocol, so the initialization sequence and lower layer communication are designed by PCIe controller.

- **PCIe Controller**

This module is designed to check PCIe device class, to set BAR address, and to enable master mode. The essential parameters to set up PCIe environment are mapped to configuration space area of SSD which can be accessed through 128-bit Tx/Rx AXI4-Stream.

Also, Integrated Block for PCI Express needs configuration space setup through CFG interface. PCIe controller includes the logic to initialize configuration space in Integrated Block for PCI Express.

NVMe

Following NVMe standard, there are two command types, i.e. Admin command and NVM command. Admin command is controlled through Admin submission queues and Admin completion queues. NVM command is controlled through I/O submission queues and I/O completion queues.

In case of Identify/Shutdown/Write/Read command, Submission queue entry has been prepared and status value within Completion queue entry has been monitored by NVMe block to check error status. 256 Kbyte data buffer is used to store data transferring between user logic and SSD when operating Write/Read command. Data output from Identify command is transferred through Identify interface.

For SMART/Flush command, Submission queue entry is prepared by user through Custom command interface. NVMe block still monitors error status within Completion queue entry like other commands. For more information, Completion queue entry of SMART/Flush command is also mapped to be output signal through Custom command interface. Data output from SMART command is transferred to user through Custom command RAM interface.

The details of each module within NVMe block are described as follows.

- **NVMe Host Controller**

The sequence of NVMe IP (shown in Figure 2) is controlled by NVMe host controller. The first step is sending Submission queue entry to SSD to start new command operation. If the command requires data transferring, NVMe host controller needs to prepare and send data pointer for using in the command. Data pointer for each command is defined to different interface, i.e. data buffer for Write/Read command, Identify interface for Identify command, and Custom command RAM interface for SMART command. The command is completed after Completion queue entry is received from SSD. Busy flag of user interface is de-asserted to '0' when the command is completed.

Otherwise, the logic of NVMe host controller needs to create multiple Write/Read command when total transfer length from user is more than 128 Kbyte size (Maximum transfer size of one NVM command is 128 Kbyte).

Busy flag of most command is de-asserted to '0' after SSD returns no error status in Completion queue entry to NVMe IP, except Read command. Busy flag of Read command is de-asserted to '0' after user reads all data from NVMe IP (Data buffer is empty).

- **Command Parameter**

This module is designed to prepare Submission queue entry for all commands and decode status value of Completion queue entry.

Submission queue entry size is 16 dwords (1 dword=32-bit). When running Identify, Shutdown, Write, and Read command, Command parameter module creates the value of all 16 dwords of Submission queue entry following Command interface of dgIF typeS. When running SMART and Flush command, 16 dwords of Submission queue entry are loaded from user input through Custom command interface.

Completion queue entry size is 4 dwords. The value of Completion queue entry is monitored by Command parameter for all commands. For SMART and Flush command, 4 dwords of Completion queue entry are latched to be output of Custom command interface.

- **Data Controller**

This module is operated when the command requires data transferring such as Write and Read command. Data packet request is created by SSD after receiving new command. Data controller decodes the address in SSD request packet to check data type such as Write/Read data, Identify command data, SMART data, and pointer list. Each data type is mapped to different address. When SSD sends read data request, the logic inside Data controller selects returned data following the address in the request.

- **AsyncCtrl**

NVMe IP is designed to run user logic in another clock domain which is not PCIe clock. The limitation is that User clock frequency must be higher than or equal to PCIe clock because transmit packet to Integrated Block for PCIe must be always ready between start of frame and end of frame. Using low frequency, data packet in each frame could not be sent to Integrated Block for PCIe continuously.

- Data Buffer**
 256 Kbyte simple dual port RAM is implemented by BlockRAM. This RAM is used to be data buffer transferring between UserLogic and SSD in Write and Read command.

User Logic

To set user parameters such as command, address, and transfer size, simple logic with small state machine could be designed in User Logic. Data stream for Write/Read command are connected to FIFO. Data output from SMART and Identify command are connected to simple dual port RAM. Data size returning from Identify command is 8 Kbyte while data size from SMART command is 512 byte. Bus width of all data interface are 128-bit.

Integrated Block for PCI Express

The integrated hard block PCIe solution in the device is provided by Xilinx. The maximum number of SSDs connecting to one FPGA device is limited by the numbers of Integrated Block for PCIe. More details of Integrated Block for PCIe 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

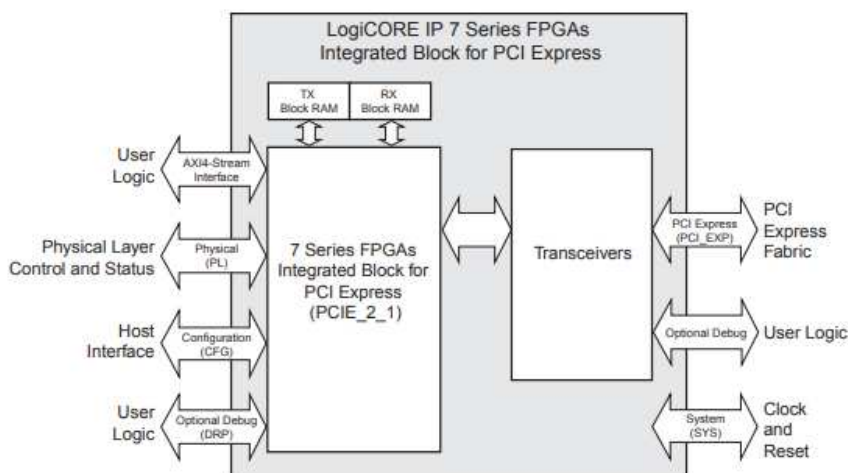


Figure 3 Integrated Block for PCI Express (7 Series FPGA)

Core I/O Signals

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

Table 3: Core I/O Signals

Signal	Dir	Description
Control I/F of dgIF typeS (Synchronous to Clk)		
RstB	In	Synchronous reset signal. Active low. Release to '1' when Clk signal is stable.
Clk	In	System clock for running NVMe IP. The frequency of Clk must be more than or equal to PCIeClk which is output from Integrated Block for PCIe (125 MHz for PCIe Gen2, 250 MHz for PCIe Gen3).
UserCmd[2:0]	In	User Command ("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. 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 in the request 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.
UserReq	In	Assert to '1' to request the new command. Can be asserted when the IP is Idle (UserBusy='0'). Command parameter (UserCmd, UserAddr, UserLen, 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. UserReq will not be allowed to '1' if IP is still busy status.
LBASize[47:0]	Out	Total capacity of SSD in 512-byte unit. Default value is 0. This value is valid after user sends Identify command.
LBAMode	Out	LBA unit size ('0': 512byte, '1': 4 Kbyte). Default value is 0. This value is valid after user sends Identify command.
UserError	Out	Error flag. Assert when UserErrorType is not equal to 0. The flag can be cleared by asserting RstB signal.
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 returned 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 returned 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.

Signal	Dir	Description
Control I/F of dgIF typeS (Synchronous to Clk)		
UserErrorType[31:0]	Out	[6] – Error when Completion TLP packet size is not correct. [7] – Error when Avalon-ST PCIe Hard IP detects 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. [10] – Error when Length[1:0] in Memory Write Request TLP packet is not equal to 0 (not aligned to 128-bit unit). [11] - Error when Address[3:2] in Memory Write Request TLP packet is not equal to 0 (not aligned to 128-bit unit). [15:12] – Reserved [16]- Error from unsupport LBA unit (LBA unit is not equal to 512 byte or 4 Kbyte) [31:17] – Reserved Note: Timeout period of bit[2]/[4] is set from TimeOutSet input.
Data I/F of dgIF typeS (Synchronous to Clk)		
UserFifoWrCnt[15:0]	In	Write data counter of Received FIFO. Used to check full status. If total FIFO size is less than 16-bit, please fill '1' to upper bit.
UserFifoWrEn	Out	Assert to '1' to write data valid of Received FIFO.
UserFifoWrData[127:0]	Out	Write data bus of Received FIFO. Synchronous to UserFifoWrEn.
UserFifoRdCnt[15:0]	In	Read data counter of Transmit FIFO. Used to check total numbers of data stored in FIFO. If FIFO size signal is less than 16-bit, please fill '0' to upper bit.
UserFifoEmpty	In	FIFO empty flag of Transmit FIFO to check available data status. This signal is not used.
UserFifoRdEn	Out	Assert to '1' to read data from Transmit FIFO.
UserFifoRdData[127:0]	In	Read data returned from Transmit FIFO. Valid in the next clock after UserFifoRdEn is asserted.
NVMe IP Interface (Synchronous to Clk)		
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). If TimeOutSet is equal to 0, Timeout will be disabled.
PCleLinkup	In	Asserted to '1' when LTSSM state of Integrated Block for PCIe is in L0 State.
AdmCompStatus[15:0]	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	[0] – Set to '1' when Phase tag in IO Completion Entry is invalid. [15:1] – Status field value of IO Completion Entry
NVMeCAPReg[31:0]	Out	Some parts of NVMe capability register output from SSD. [15:0] – MQES (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' when IdenWrAddr and IdenWrData are valid on the bus.
IdenWrAddr[8:0]	Out	Index of IdenWrData in 128-bit unit. Synchronous to IdenCtrlWrEn. 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. Synchronous to IdenWrEn.

Signal	Dir	Description
NVMe IP Custom interface (Synchronous to Clk)		
CtmSubmDW0[31:0] – CtmSubmDW15[31:0]	In	16 Dwords of Submission queue entry for SMART/Flush command. DW0: Command Dword0, DW1: Command Dword1, ..., and DW15: Command Dword15. These inputs must be valid and stable when UserReq='1' and UserCmd="100" (SMART) or "110" (Flush).
CtmCompDW0[31:0] – CtmCompDW3[31:0]	Out	4 Dwords of Completion queue entry output from SMART/Flush command. DW0: Completion Dword0, DW1: Completion Dword1, ..., and DW3: Completion Dword3
CtmRamWrEn	Out	Assert to '1' when CtmRamAddr and CtmRamWrData are valid on the bus (found in SMART command).
CtmRamAddr[8:0]	Out	Index of CtmRamWrData when command returns SMART data. (Optional) Index to request data input through CtmRamRdData to support other custom commands.
CtmRamWrData[127:0]	Out	512-byte data output from SMART command.
CtmRamRdData[127:0]	In	(Optional) Data input for other custom commands.

Signal	Dir	Description
Integrated Block for PCIe (Synchronous to PCIeClk) for PCIe Gen3 only		
PCleRstB	In	Synchronous reset signal. Active low. Release to '1' when Integrated Block for PCIe is not in reset state.
PCleClk	In	Clock output from Integrated Block for PCIe (250 MHz for PCIe Gen3).
Configuration Management Interface (for PCIe Gen3 only)		
PCleCfgDone	In	Read/Write operation complete. Asserted for 1 cycle when operation completes.
PCleCfgRdEn	Out	Read enable. Asserted for a read operation.
PCleCfgWrEn	Out	Write enable. Asserted for a write operation.
PCleCfgWrData[31:0]	Out	Write data which is used to configure the Configuration and Management registers.
PCleCfgByteEn[3:0]	Out	Byte enable for write data, where bit[0] corresponds to PCleCfgWrData[7:0], and so on.
PCleCfgAddr[18:0]	Out	Read/Write Address.
Requester Request Interface (for PCIe Gen3 only)		
PCleMtTxData[127:0]	Out	Requester request data bus.
PCleMtTxKeep[3:0]	Out	Bit i indicates that Dword i of PCleMtTxData contains valid data.
PCleMtTxLast	Out	Asserts this signal in the last cycle of a TLP to indicate the end of the packet.
PCleMtTxReady[3:0]	In	Asserts to accept data. Data is transferred when both PCleMtTxValid and PCleMtTxReady are asserted in the same cycle.
PCleMtTxUser[59:0]	Out	Requester request user data. Valid when PCleMtTxValid is high.
PCleMtTxValid	Out	Asserts to drive valid data on PCleMtTxData bus. NVMe IP keeps the valid signal asserted during the transfer of a packet.
Completer Request Interface (for PCIe Gen3 only)		
PCleMtRxData[127:0]	In	Receive data from Integrated Block for PCIe.
PCleMtRxKeep[3:0]	In	Bit i indicates that Dword i of PCleMtRxData contains valid data.
PCleMtRxLast	In	Asserts this signal in the last beat of a packet to indicate the end of the packet.
PCleMtRxReady	Out	Indicates that NVMe IP is ready to accept data.
PCleMtRxUser[74:0]	In	Sideband information for the TLP being transferred. Valid when PCleMtRxValid is high.
PCleMtRxValid	In	Assert when Integrated Block for PCIe drives valid data on PCleMtRxData bus. Integrated Block for PCIe keeps the valid signal asserted during the transfer of packet.
Completer Completion Interface (for PCIe Gen3 only)		
PCleSITxData[127:0]	Out	Completion data from NVMe IP.
PCleSITxKeep[3:0]	Out	Bit i indicates that Dword i of PCleSITxData contains valid data.
PCleSITxLast	Out	Asserts this signal in the last cycle of a packet to indicate the end of the packet.
PCleSITxReady[3:0]	In	Indicates that Integrated Block for PCIe is ready to accept data.
PCleSITxUser[32:0]	Out	Sideband information for the TLP being transferred. Valid when PCleSITxValid is high.
PCleSITxValid	Out	Asserts to drive valid data on PCleSITxData bus. NVMe IP keeps the valid signal asserted during the transfer of a packet.
Requester Completion Interface (for PCIe Gen3 only)		
PCleSIRxData[127:0]	In	Receive data from Integrated Block for PCIe.
PCleSIRxKeep[3:0]	In	Bit i indicates that Dword i of PCleSIRxData contains valid data.
PCleSIRxLast	In	Asserts this signal in the last beat of a packet to indicate the end of the packet.
PCleSIRxReady	Out	Indicates that NVMe IP is ready to accept data.
PCleSIRxUser[84:0]	In	Sideband information for the TLP being transferred. Valid when PCleSIRxValid is high.
PCleSIRxValid	In	Asserts when Integrated Block for PCIe drives valid data on PCleSIRxData bus. Integrated Block for PCIe keeps the valid signal asserted during the transfer of packet.

Signal	Dir	Description
Integrated Block for PCIe (Synchronous to PCIeClk) for PCIe Gen2 only		
PCleRstB	In	Synchronous reset signal. Active low. Release to '1' when Integrated Block for PCIe is not in reset state.
PCleClk	In	Clock output from Integrated Block for PCIe (125 MHz for PCIe Gen2).
Configuration Management Interface (for PCIe Gen2 only)		
PCleCfgDone	In	Read/Write operation complete. Asserted for 1 cycle when operation completes.
PCleCfgRdEn	Out	Read enable. Asserted for a read operation.
PCleCfgWrEn	Out	Write enable. Asserted for a write operation.
PCleCfgWrData[31:0]	Out	Write data which is used to configure the Configuration and Management registers.
PCleCfgByteEn[3:0]	Out	Byte enable for write data, where bit[0] corresponds to PCleCfgWrData[7:0], and so on.
PCleCfgAddr[9:0]	Out	Read/Write Address.
PCIe Transmit Interface (for PCIe Gen2 only)		
PCleTxData[127:0]	Out	Transmit data to Integrated Block for PCIe.
PCleTxKeep[15:0]	Out	Transmit data strobe. Determine which data bytes are valid on PCleTxData.
PCleTxLast	Out	Transmit End-of-Frame. Valid only along with assertion of PCleTxValid.
PCleTxReady[3:0]	In	Indicates that Integrated Block for PCIe is ready to accept data. The simultaneous assertion of PCleTxValid and PCleTxReady marks the successful transfer of one data beat on PCleTxData.
PCleTxUser[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.
PCleTxValid	Out	Indicates that NVMe IP is presenting valid data on PCleTxData.
PCIe Receive Interface (for PCIe Gen2 only)		
PCleRxData[127:0]	In	Receive data from Integrated Block for PCIe. Valid only PCleRxValid is also asserted.
PCleRxKeep[15:0]	In	Receive data strobe. Determine which data bytes are valid on PCleRxData.
PCleRxLast	In	Receive End-of-Frame. Valid only if PCleRxValid is also asserted.
PCleRxReady	Out	Indicates that NVMe IP is ready to accept data on PCleRxData. The simultaneous assertion of PCleRxValid and PCleRxReady marks the successful transfer of one data beat on PCleRxData.
PCleRxUser[21:0]	In	Bit[0]: Receive ECRC error. Bit[14:13]: Indicates the start of a new packet header in PCleRxData. Bit[21]: Indicates the end of a packet in PCleRxData. Other bits are ignored in NVMe IP.
PCleRxValid	In	Indicates that the core is presenting valid data on PCleRxData.

Timing Diagram

Initialization

The sequence of the initialization process is as follows.

- 1) RstB is released by user after Clk is stable.
- 2) NVMe IP waits until both PCIeRstB and PCIeLinkup are asserted to '1'. PCIeLinkup is asserted to '1' when PCIe Hard IP completes power sequence (LTSSM State is in L0 state). PCIeRstB is de-asserted to '1' after PCIe Hard IP is not in reset state and ready to interface with application layer. PCIeRstB is generated in PCIeClk domain.
- 3) NVMe IP starts initialization process.
- 4) UserBusy is deasserted to '0' after NVMe IP completes initialization process.

After complete above sequence, NVMe IP is ready to receive the command from user.

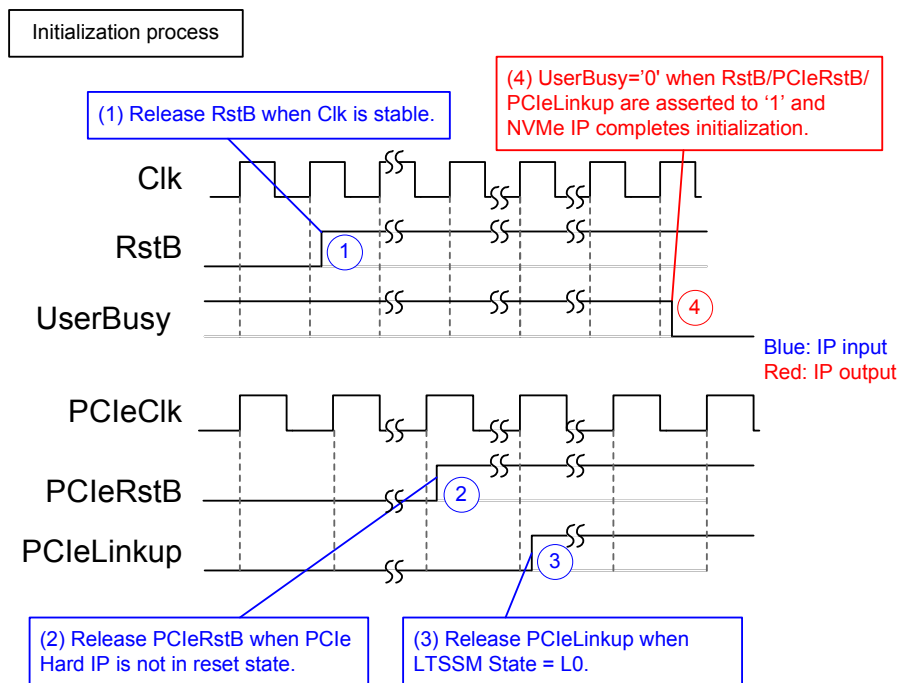


Figure 4: UserBusy after system boot-up

dgIF typeS

dgIF typeS signal is split into two interfaces, i.e. command interface and data interface. Figure 5 shows timing diagram of command interface of dgIF typeS. Before sending new command to the IP, UserBusy must be monitored to confirm that IP is Idle. Command parameters (UserCmd, UserAddr, UserLen, and CtmSubmDW0-DW15) must be valid and stable during asserting UserReq='1'. After receiving new request from user, UserBusy changes status from '0' to '1'. Next, UserReq could be de-asserted to '0' and user logic can prepare the next command parameter on the bus.

Note: 1) UserAddr and UserLen value are not used in Identify, Shutdown, SMART, and Flush command.
2) CtmSubmDW0 – DW15 are not used in Identify, Shutdown, Write, and Read command.

For data interface, Transmit FIFO is read for Write command, while Received FIFO is written for Read command. Timing diagram of data interface is shown in Figure 6 and Figure 7.

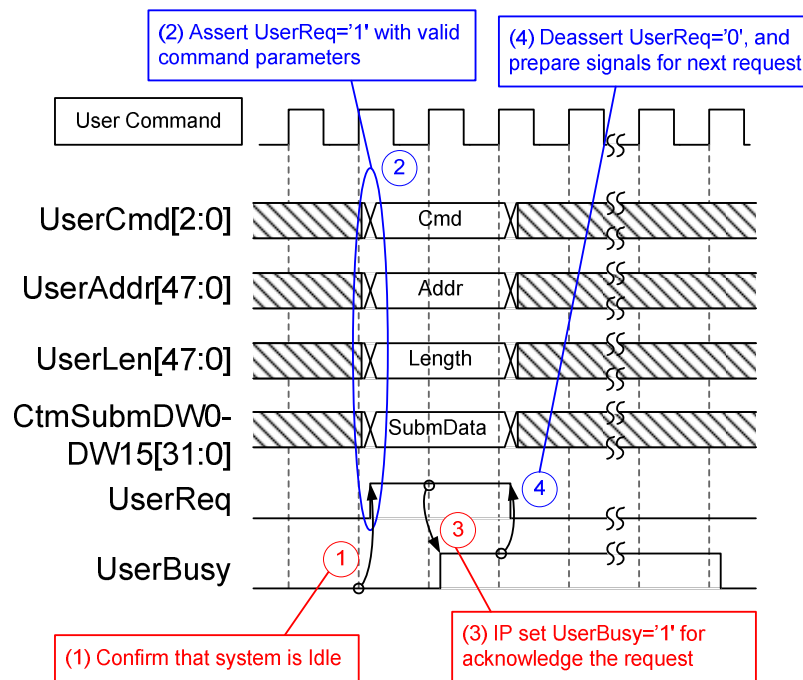


Figure 5: Command Interface of dgIF typeS Timing diagram

For Write command, 128-bit data from Transmit FIFO is stored to data buffer within NVMe IP. DMA Engine in NVMe IP monitors UserFifoRdCnt signal until it indicates that data in Transmit FIFO is equal to or more than 512 bytes. After that, UserFifoRdEn is asserted to '1' for 32 clock cycles to read 512-byte data, as shown in Figure 6. Similar to general FIFO timing diagram, UserFifoRdData is valid in the next clock after UserFifoRdEn is asserted to '1'.

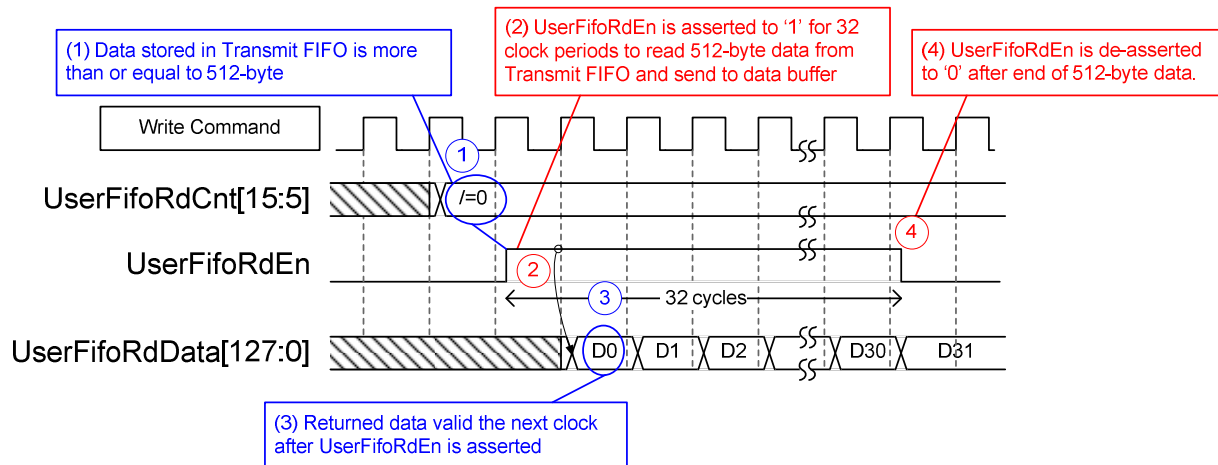


Figure 6: Transmit FIFO Interface for Write command

For Read command, UserFifoWrEn is asserted to '1' with the valid value of UserFifoWrData to send received data from data buffer. Similar to Write command, UserFifoWrCnt is monitored to check free space of Received FIFO is much enough (more than or equal to 1024-byte) before transferring 512-byte data to FIFO.

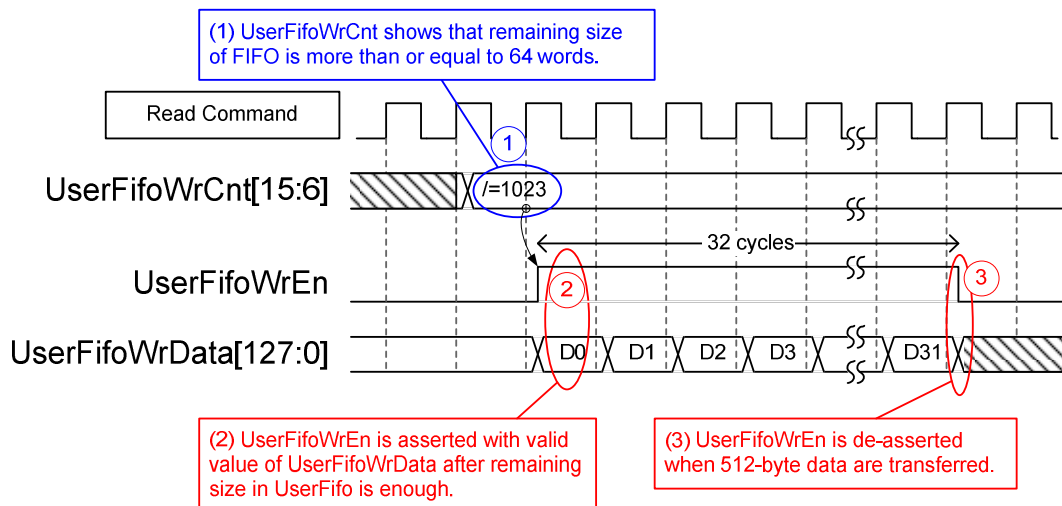


Figure 7: Received FIFO Interface for Read command

IdenCtrl/IdenName

Before sending Write or Read command to IP, user should send Identify command firstly to update LBASize and LBAMode output. LBASize value is used in User Logic to confirm that the sum of address and length in Write/Read command is not out-of-range. If LBAMode='1' (LBA unit = 4 Kbyte), the address and length input for Write/Read command must be aligned to 4Kbyte.

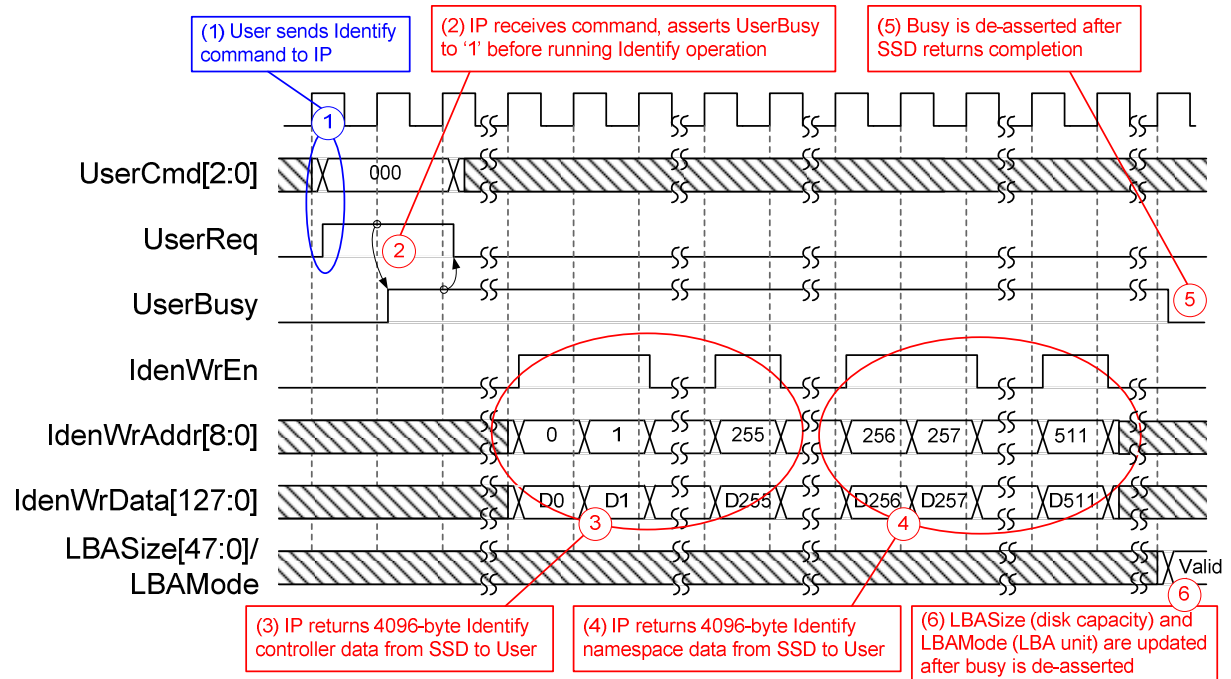


Figure 8: LBASize/LBAMode updated after Identify command

Figure 8 shows timing diagram of Identify command. UserCmd and UserReq are sent out when UserBusy='0'. UserAddr and UserLen input are not used for Identify command. After Identify command is sent to SSD, 4096-byte Identify Controller data and 4096-byte Identify Namespace data are returned. IdenWrAddr is equal to 0 – 255 during sending Identify Controller data and equal to 256 – 511 during sending Identify Namespace data. 8-byte Identify Controller data or Identify Namespace data is valid on IdenWrData bus for each clock cycle, synchronous to IdenWrAddr and IdenWrEn signal. Finally, LBASize and LBAMode are updated when UserBusy is de-asserted to '0' from Identify command.

Shutdown

Before system is power down, it is recommended to send Shutdown command to NVMe IP. After issue this command, SSD shutdowns in good sequence. During shutdown operation, Cache in SSD is flushed to the memory and all write data sending to SSD could be guaranteed. After issue this command, NVMe IP and SSD cannot be used until system is power down.

To send Shutdown command, UserCmd is set to "001". UserAddr, UserLen, and CtmSubmDW0-DW15 are not used in this command.

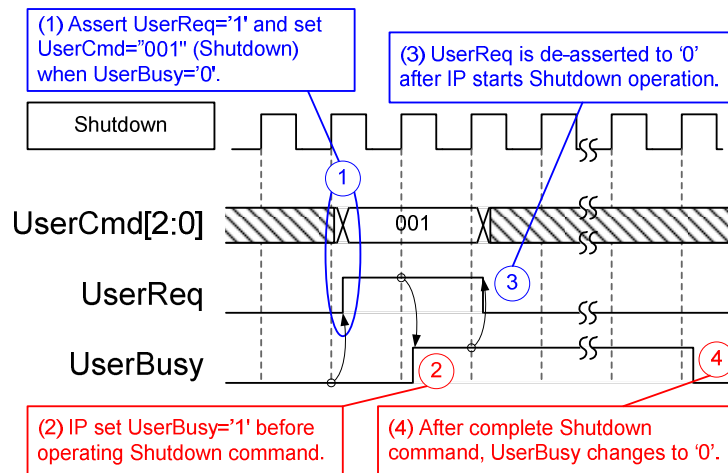


Figure 9: Example to send Shutdown command

SMART

To check SSD health, SMART command is applied to get 512-byte information through Custom command interface. To run SMART command, user must set UserCmd="100", set 16-Dword data to CtmSubmDW0-DW15, and assert UserReq to '1' at the same clock cycle. UserAddr and UserLen input are not used in SMART command. After that, SSD returns 512-byte SMART data through CtmRam (Custom command RAM) interface. 128-bit SMART data is valid on CtmRamWrData bus at the same clock cycle as CtmRamWrEn='1'. CtmRamAddr is used to be data index of 32 x 128-bit SMART data.

The value of CtmSubmDW0-DW15 for SMART command is fixed as follows.

CtmSubmDW0	= 0x0000_0002
CtmSubmDW1	= 0xFFFF_FFFF
CtmSubmDW2 – CtmSubmDW5	= 0x0000_0000
CtmSubmDW6	= 0x2000_0000
CtmSubmDW7 – CtmSubmDW9	= 0x0000_0000
CtmSubmDW10	= 0x007F_0002
CtmSubmDW11 – CtmSubmDW15	= 0x0000_0000

As shown in Figure 10, SMART data is sent through CtmRam interface. CtmRamWrEn may not be asserted to '1' continuously. The last data is returned when CtmRamAddr=0x1F. Finally, UserBusy is de-asserted to '0'.

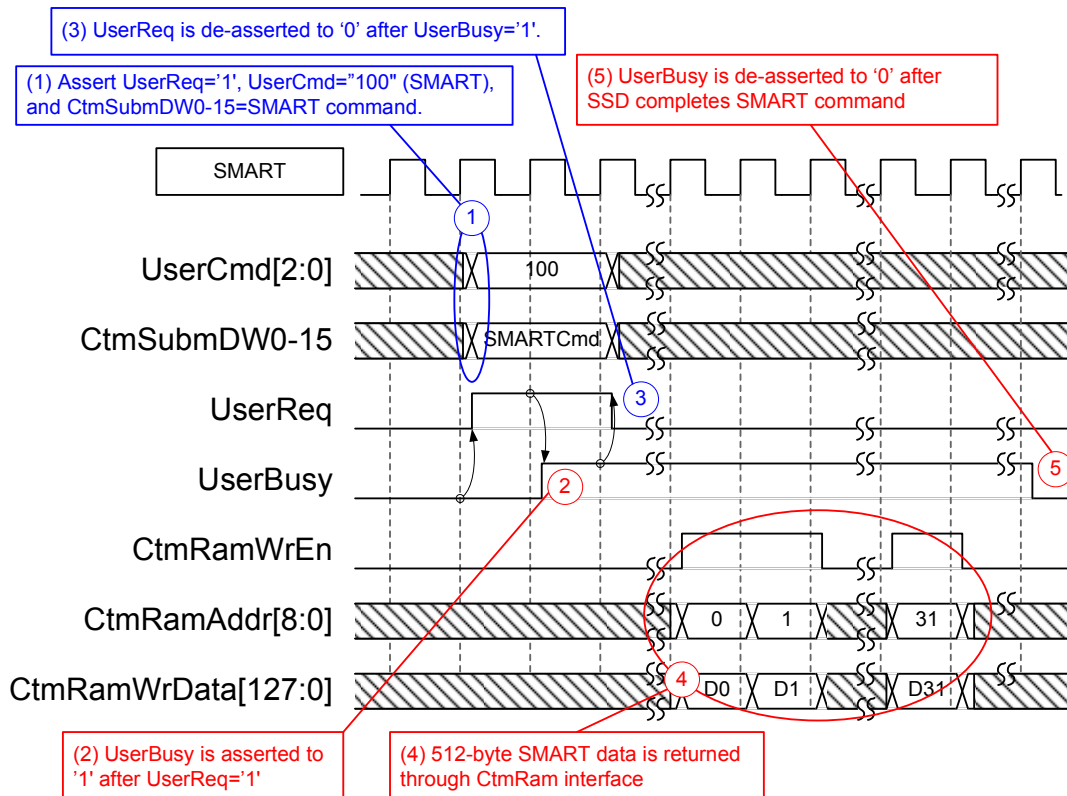


Figure 10: Example to send SMART command

Flush

Generally, when user writes data to SSD, the data is stored to internal cache of SSD firstly. If SSD still not flush data from cache to memory and system is power down immediately, Write data in cache will be lost. To guarantee that data is not lost, there are two solutions.

First solution is sending Shutdown command before power down the system. This command is required before shutdown the system. Second solution is sending Flush command to SSD. Using this command, data in cache is flushed to the memory. Flush command could be used when system is Idle to back up data from cache to the memory. After using Flush command, SSD still be active status.

Similar to SMART command, Flush command must be sent through CtmSubmDW0-DW15 by using following value.

CtmSubmDW0	= 0x0000_0000
CtmSubmDW1	= 0x0000_0001
CtmSubmDW2 – CtmSubmSW15	= 0x0000_0000

UserAddr and UserLen inputs are not used for Flush command.

As shown in Figure 11, after sending Flush command, user logic waits until UserBusy changes from '1' to '0' to confirm that write data in cache has already flushed.

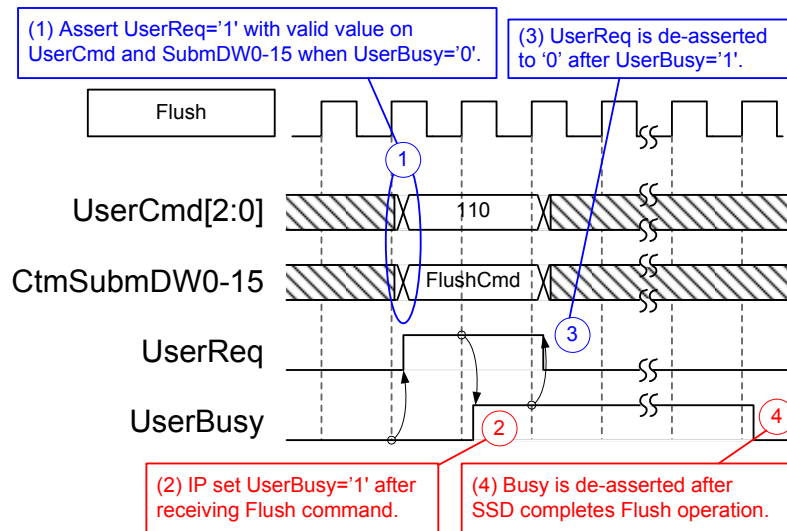


Figure 11: Example to send Flush command

Error

During normal operation, UserError and all bits of UserErrorType signal are always 0. UserError is generated by OR condition of each-bit of UserErrorType. If some bit of UserErrorType is set to '1', UserError will be asserted to '1'. UserError could be cleared by asserting RstB to be '0' only, as shown in Figure 12.

If AdmCompStatus or IOCompStatus value has error condition, UserErrorType bit[3]/[5] will be set. User can see more details of the error by reading AdmCompStatus and IOCompStatus value.

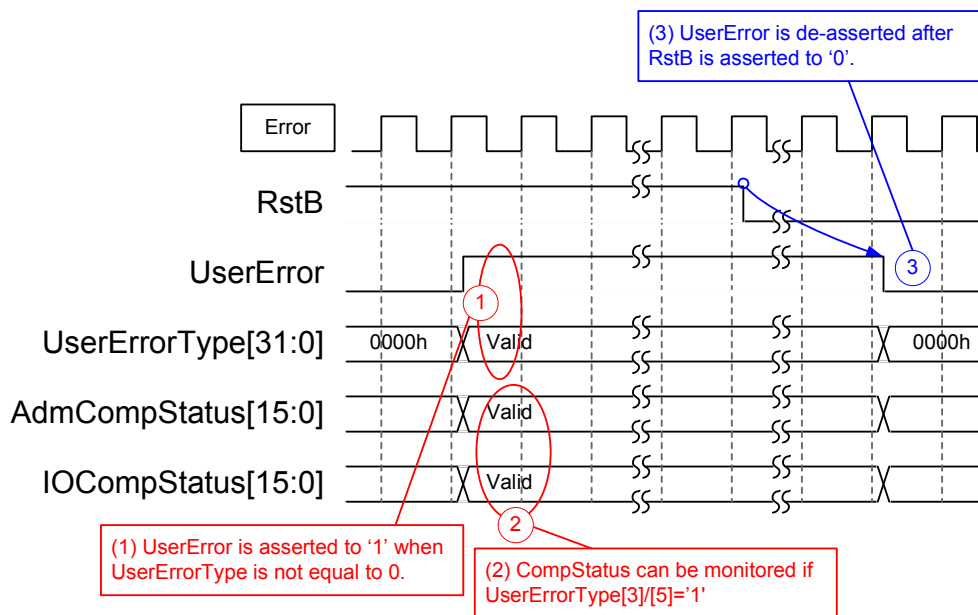


Figure 12: Error flag Timing diagram

Verification Methods

The NVMe IP Core functionality was verified by simulation and also proved on real board design by using 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	Jun-2-2016	Initial Release
1.1	Jun-15-2016	Add KCU105 support
1.2	Sep-5-2016	Add ZC706 support
1.3	Sep-9-2016	Add KC705 support
1.4	Oct-27-2016	Support VC709 and Zynq Mini-ITX
1.5	Dec-9-2016	Modify buffer to be BRAM and modify user interface to dgIF typeS
2.0	Jun-7-2017	Change PCIe interface to connect to Integrated Block for PCIe
2.1	Nov-30-2017	Support ZCU106 board
3.0	July-19-2018	Support Shutdown, SMART, and Flush command