

User's Manual

NetXtreme-E

Broadcom[®] NetXtreme-C and NetXtreme-E

USER'S MANUAL

Revision History

<i>Revision</i>	<i>Date</i>	<i>Change Description</i>
NetXtreme-E-UG100	2/26/18	Initial release for 20.6.

© 2018 by Broadcom. All rights reserved.

Broadcom[®], the pulse logo, Connecting everything[®], Avago Technologies, and the A logo are among the trademarks of Broadcom and/or its affiliates in the United States, certain other countries and/or the EU. The term “Broadcom” refers to Broadcom Limited and/or its subsidiaries. For more information, please visit www.broadcom.com.

Broadcom reserves the right to make changes without further notice to any products or data herein to improve reliability, function, or design. Information furnished by Broadcom is believed to be accurate and reliable. However, Broadcom does not assume any liability arising out of the application or use of this information, nor the application or use of any product or circuit described herein, neither does it convey any license under its patent rights nor the rights of others.

Table of Contents

Regulatory and Safety Approvals	7
Regulatory	7
Safety	7
Electromagnetic Compatibility (EMC).....	8
Electrostatic Discharge (ESD) Compliance.....	8
FCC Statement.....	8
Functional Description	9
Network Link and Activity Indication	14
BCM957402AXXXX/BCM957412AXXXX	14
BCM957404AXXXX/BCM957414AXXXX	15
BCM957406AXXXX/BCM957416AXXXX	16
BCM957414M4140D.....	17
BCM957412M4120D.....	18
BCM957416M4160	19
Features	20
Software and Hardware Features	20
Virtualization Features.....	21
VXLAN.....	22
NVGRE/GRE/IP-in-IP/Geneve	22
Stateless Offloads	22
RSS.....	22
TPA	22
Header-Payload Split.....	22
UDP Fragmentation Offload.....	22
Stateless Transport Tunnel Offload.....	23
Multiqueue Support for OS.....	23
NDIS VMQ	23
VMWare NetQueue.....	23
KVM/Xen Multiqueue	23
SR-IOV Configuration Support Matrix	23
SR-IOV	24
Network Partitioning (NPAR).....	24
RDMA over Converge Ethernet – RoCE	24
Supported Combinations.....	25
NPAR, SR-IOV, and RoCE	25
NPAR, SR-IOV, and DPDK.....	25
Unsupported Combinations.....	25

Installing the Hardware	26
Safety Precautions	26
System Requirements	26
Hardware Requirements	26
Preinstallation Checklist	26
Installing the Adapter	27
Connecting the Network Cables	27
Supported Cables and Modules	27
Copper	28
SFP+	28
SFP28	28
Software Packages and Installation	29
Supported Operating Systems	29
Installing Drivers	29
Windows	29
<i>Dell DUP</i>	29
<i>GUI Install</i>	29
<i>Silent Install</i>	29
<i>INF Install</i>	29
Linux	30
<i>Module Install</i>	30
<i>Linux Ehtool Commands</i>	31
VMware	32
Firmware Update	33
Dell Update Package	33
Windows	33
Linux	33
Windows Driver Advanced Properties and Event Log Messages	34
Driver Advanced Properties	34
Event Log Messages	35
Teaming	37
Windows	37
Linux	37
System-level Configuration	38
UEFI HII Menu	38
Main Configuration Page	38
Firmware Image Properties	38
Device Level Configuration	38
NIC Configuration	38

iSCSI Configuration.....	38
Comprehensive Configuration Management.....	39
Device Hardware Configuration	39
MBA Configuration Menu	39
iSCSI Boot Main Menu.....	39
Auto-Negotiation Configuration	40
<i>Operational Link Speed</i>	43
<i>Firmware Link Speed</i>	43
<i>Auto-negotiation Protocol</i>	43
<i>Windows Driver Settings</i>	43
<i>Linux Driver Settings</i>	43
<i>ESXi Driver Settings</i>	44
FEC Auto-Negotiation	45
Link Training.....	46
Media Auto Detect.....	47
iSCSI Boot	49
Supported Operating Systems for iSCSI Boot	49
Setting up iSCSI Boot.....	49
Configuring the iSCSI Target	49
Configuring iSCSI Boot Parameters.....	50
MBA Boot Protocol Configuration	50
iSCSI Boot Configuration	51
<i>Static iSCSI Boot Configuration</i>	51
<i>Dynamic iSCSI Boot Configuration</i>	52
Enabling CHAP Authentication	53
Configuring the DHCP Server to Support iSCSI Boot.....	54
DHCP iSCSI Boot Configurations for IPv4.....	54
<i>DHCP Option 17, Root Path</i>	54
<i>DHCP Option 43, Vendor-Specific Information</i>	54
<i>Configuring the DHCP Server</i>	55
DHCP iSCSI Boot Configuration for IPv6.....	55
<i>DHCPv6 Option 16, Vendor Class Option</i>	55
<i>DHCPv6 Option 17, Vendor-Specific Information</i>	55
<i>Configuring the DHCP Server</i>	56
VXLAN: Configuration and Use Case Examples	56
SR-IOV: Configuration and Use Case Examples	57
Linux Use Case.....	57
Windows Case	58
VMWare SRIOV Case.....	59

NPAR – Configuration and Use Case Example	62
Features and Requirements	62
Limitations	62
Configuration	62
Notes on Reducing NIC Memory Consumption	65
RoCE – Configuration and Use Case Examples	66
Linux Configuration	66
Requirements	66
BNXT_RE Driver Dependencies	66
Installation	67
Limitations	67
Known Issues	67
Windows	68
Kernel Mode	68
Verifying RDMA	68
User Mode	69
VMware ESX	70
Limitations	70
BNXT RoCE Driver Requirements	70
Installation	70
Configuring Paravirtualized RDMA Network Adapters	71
<i>Configuring a Virtual Center for PVRDMA</i>	71
<i>Tagging vmknics for PVRDMA on ESX Hosts</i>	71
<i>Setting the Firewall Rule for PVRDMA</i>	71
<i>Adding a PVRDMA Device to the VM</i>	71
<i>Configuring the VM on Linux Guest OS</i>	72
DCBX – Data Center Bridging	73
QoS Profile – Default QoS Queue Profile	73
DCBX Mode = Enable (IEEE only)	74
DCBX Willing Bit	74
Frequently Asked Questions	77

Regulatory and Safety Approvals

The following sections detail the Regulatory, Safety, Electromagnetic Compatibility (EMC), and Electrostatic Discharge (ESD) standard compliance for the NetXtreme-E Network Interface Card.

Regulatory

Table 1: Regulatory Approvals

Item	Applicable Standard	Approval/Certificate
CE/European Union	EN 62368-1:2014	CB report and certificate
UL/USA	IEC 62368-1 ed. 2	CB report and certificate
CSA/Canada	CSA 22.2 No. 950	CSA report and certificate.
Taiwan	CNS14336 Class B	–

Safety

Table 2: Safety Approvals

Country	Certification Type/Standard	Compliance
International	CB Scheme ICES 003 – Digital Device UL 1977 (connector safety) UL 796 (PCB wiring safety) UL 94 (flammability of parts)	Yes

Electromagnetic Compatibility (EMC)

Table 3: Electromagnetic Compatibility

Standard / Country	Certification Type	Compliance
CE/EU	EN 55032:2012/AC:2013 Class B EN 55024:2010 EN 61000-3-2:2014 EN 61000-3-3:2013	CE report and CE DoC
FCC/USA	CFR47, Part 15 Class B	FCC/IC DoC and EMC report referencing FCC and IC standards
IC/Canada	ICES-003 Class B	FCC/IC DoC and report referencing FCC and IC standards
ACA/Australia, New Zealand	AS/NZS CISPR 22:2009 +A1:2010 / AS/NZS CISPR 32:2015	ACA certificate RCM Mark
BSMI/Taiwan	CNS13438 Class B	BSMI certificate
BSMI/Taiwan	CNS15663	BSMI certificate
MIC/S. Korea	KN32 Class B KN35	Korea certificate MSIP Mark
VCCI/Japan	V-3/2014/04 (take effect 3/31/2015)	Copy of VCCI on-line certificate

Electrostatic Discharge (ESD) Compliance

Table 4: ESD Compliance Summary

Standard	Certification Type	Compliance
EN55024:2010 (EN 61000-4-2)	Air/Direct discharge	Yes

FCC Statement

This equipment has been tested and found to comply with the limits for a Class B digital device, pursuant to Part 15 of the FCC Rules. These limits are designed to provide reasonable protection against harmful interference in a residential installation. This equipment generates uses and can radiate radio frequency energy and, if not installed and used in accordance with the instructions, may cause harmful interference to radio communications. However, there is no guarantee that interference will not occur in a particular installation. If this equipment does cause harmful interference to radio television reception, which can be determined by turning the equipment off and on, the user is encouraged to try to correct the interference by one or more of the following measures:

- Reorient or relocate the receiving antenna.
- Increase the separation between the equipment and receiver.
- Consult the dealer or an experienced radio/TV technician for help.



Note: Changes or modifications not expressly approved by the manufacture responsible for compliance could void the user's authority to operate the equipment.

Functional Description

Dell supports 10GBase-T, 10G SFP+, and 25G SFP28 Network Interface Cards (NICs). These NICs are described in [Table 5](#).

Table 5: Functional Description

Network Interface Card	Description
BCM957402A4020DLPC/BCM957402A4020DC/BCM957412A4120D/BCM957412M4120D	
Speed	Dual Port 10 Gbps Ethernet
PCI-E	Gen 3 x8 ^a
Interface	SFP+ for 10 Gbps
Device	Broadcom BCM57402/BCM57412 10 Gbps MAC controller with Integrated dual channel 10 Gbps SFI transceiver.
NDIS Name	Broadcom NetXtreme E-Series Dual-port 10Gb SFP+ Ethernet PCIe Adapter
UEFI Name	Broadcom Dual 10Gb SFP+ Ethernet
BCM57404A4041DLPC/BCM57404A4041DC/BCM957414A4141D/BCM957414M4140D	
Speed	Dual Port 25 Gbps or 10 Gbps Ethernet
PCI-E	Gen 3 x8 ^a
Interface	SFP28 for 25 Gbps and SFP+ for 10 Gbps
Device	Broadcom BCM57404/BCM57414 25 Gbps MAC controller with Integrated dual channel 25 Gbps SFI transceiver.
NDIS Name	Broadcom NetXtreme E-Series Dual-port 25 Gb SFP28 Ethernet PCIe Adapter
UEFI Name	Broadcom Dual 25 Gb SFP 28 Ethernet
BCM957406A4060DLPC/BCM957406A4060DC/BCM957416A4160D/BCM957416M4160	
Speed	Dual Port 10GBase-T Ethernet
PCI-E	Gen 3 x8 ^a
Interface	RJ45 for 10 Gbps and 1 Gbps
Device	Broadcom BCM57406/BCM57416 10 Gbps MAC controller with Integrated dual channel 10GBase-T transceiver.
NDIS Name	Broadcom NetXtreme E-Series Dual-port 10GBASE-T Ethernet PCIe Adapter
UEFI Name	Broadcom Dual 10GBASE-T Ethernet

- a. The NIC supports PCI-E Gen 3, Gen 2, and Gen 1 speeds, however, PCI Gen 3 is recommended to achieve nominal throughput when 2 ports of 25G links transmit and receive traffic at the same time.

Figure 1: BCM957402A4020DC, BCM957412A4120D Network Interface Card

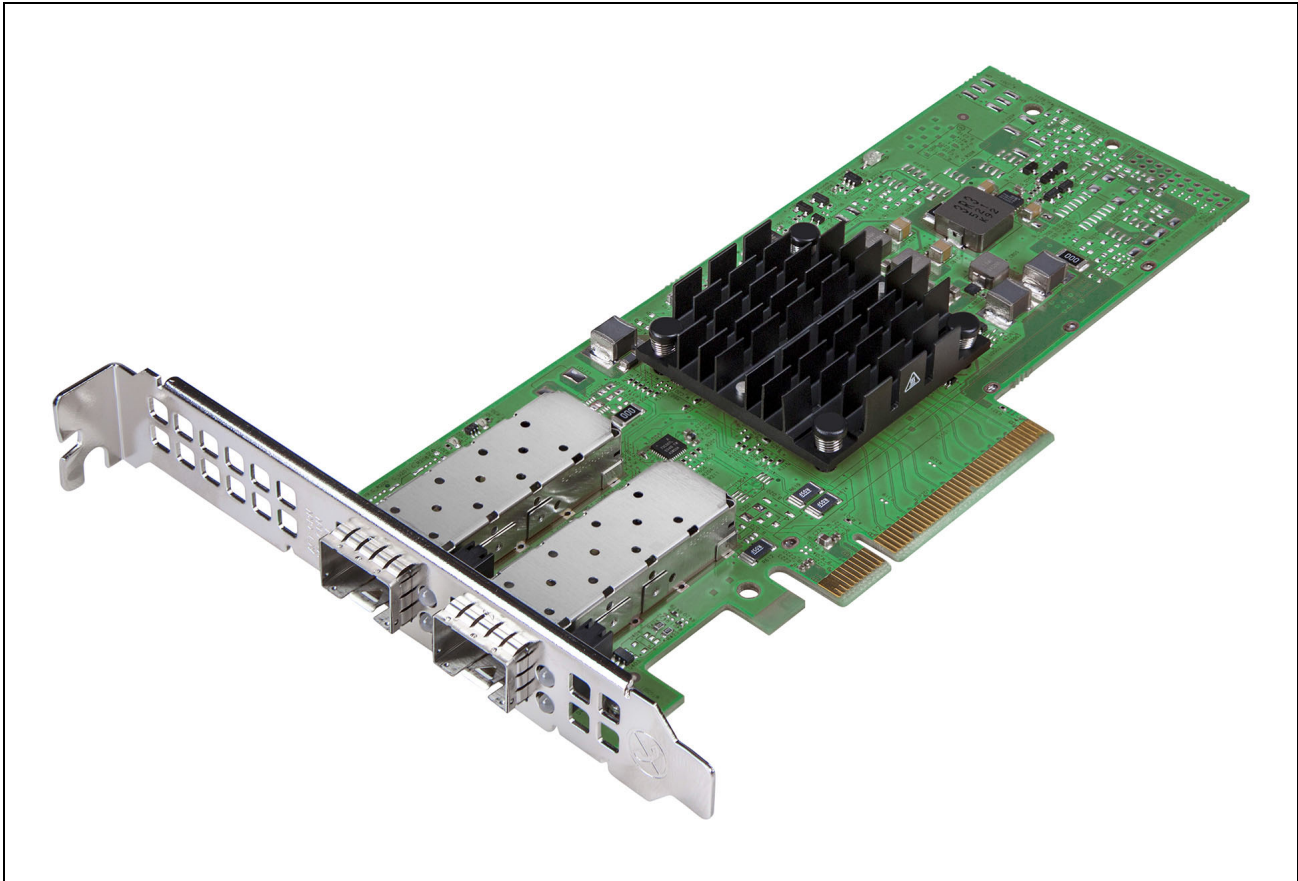


Figure 2: BCM957404A4041DLPC, BCM957414A4141D Network Interface Card

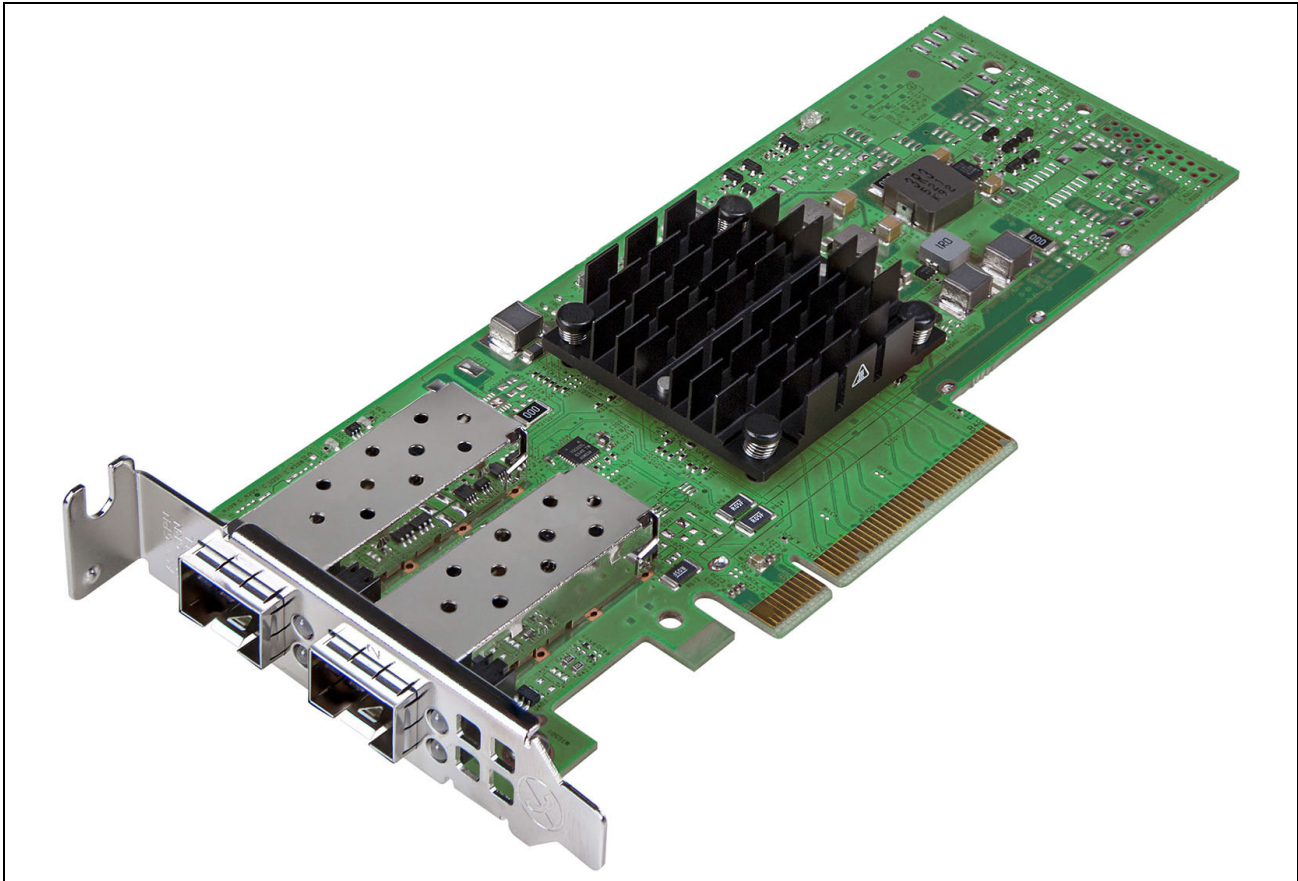


Figure 3: BCM957406A4060DLPC, BCM957416A4160D Network Interface Card

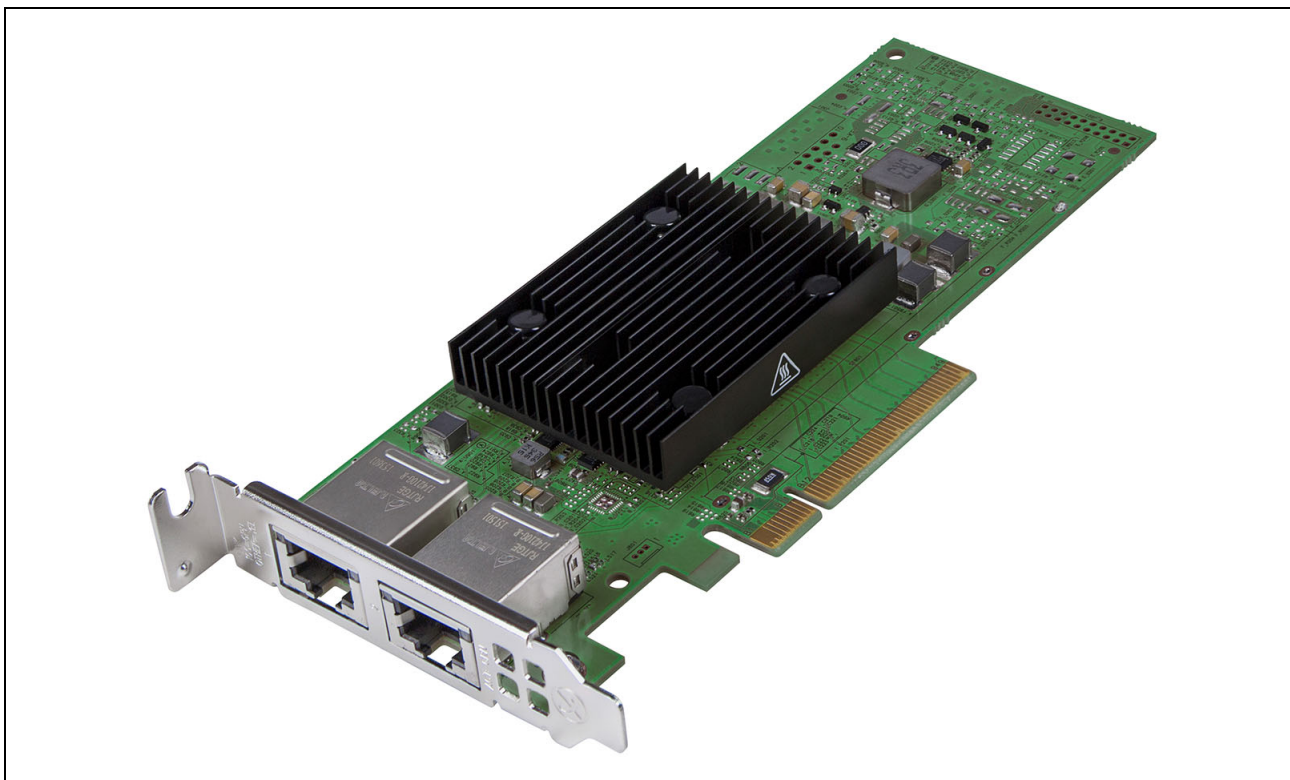


Figure 4: BCM957414M4140D Network Daughter Card (rNDC)

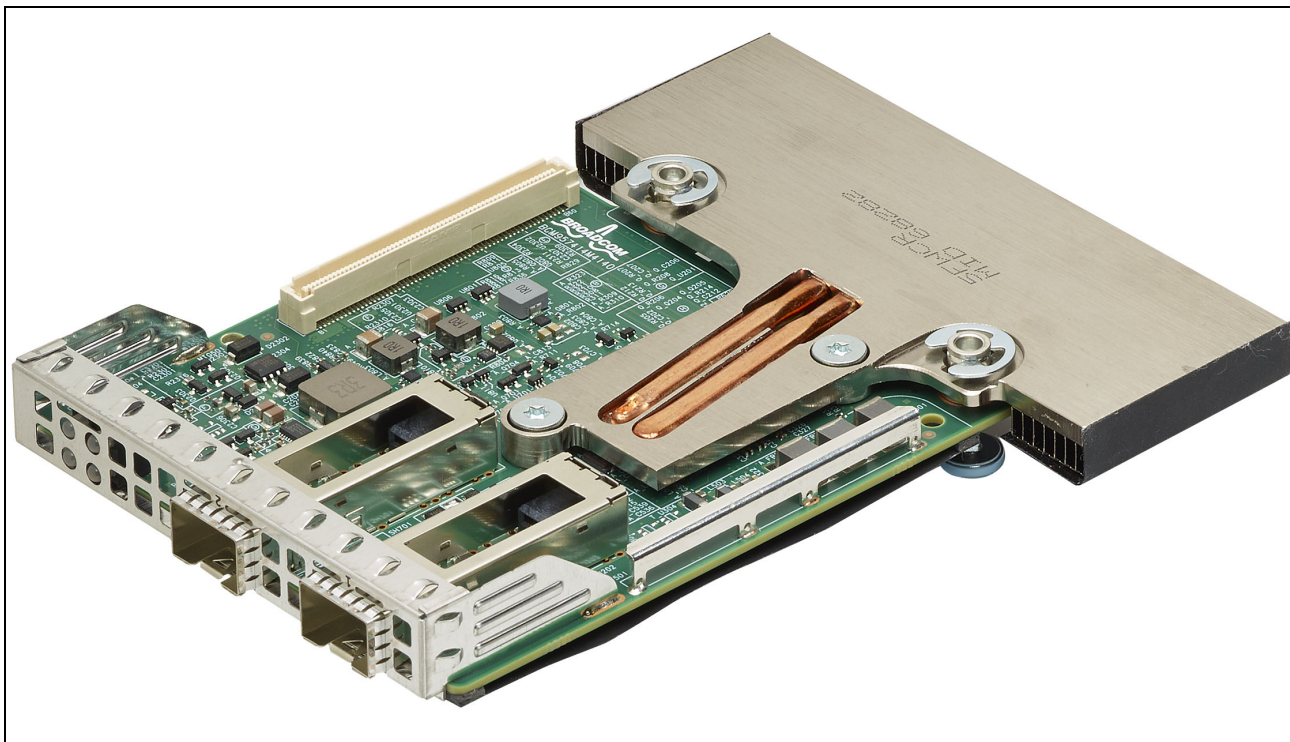


Figure 5: BCM957412M4120D Network Daughter Card (rNDC)

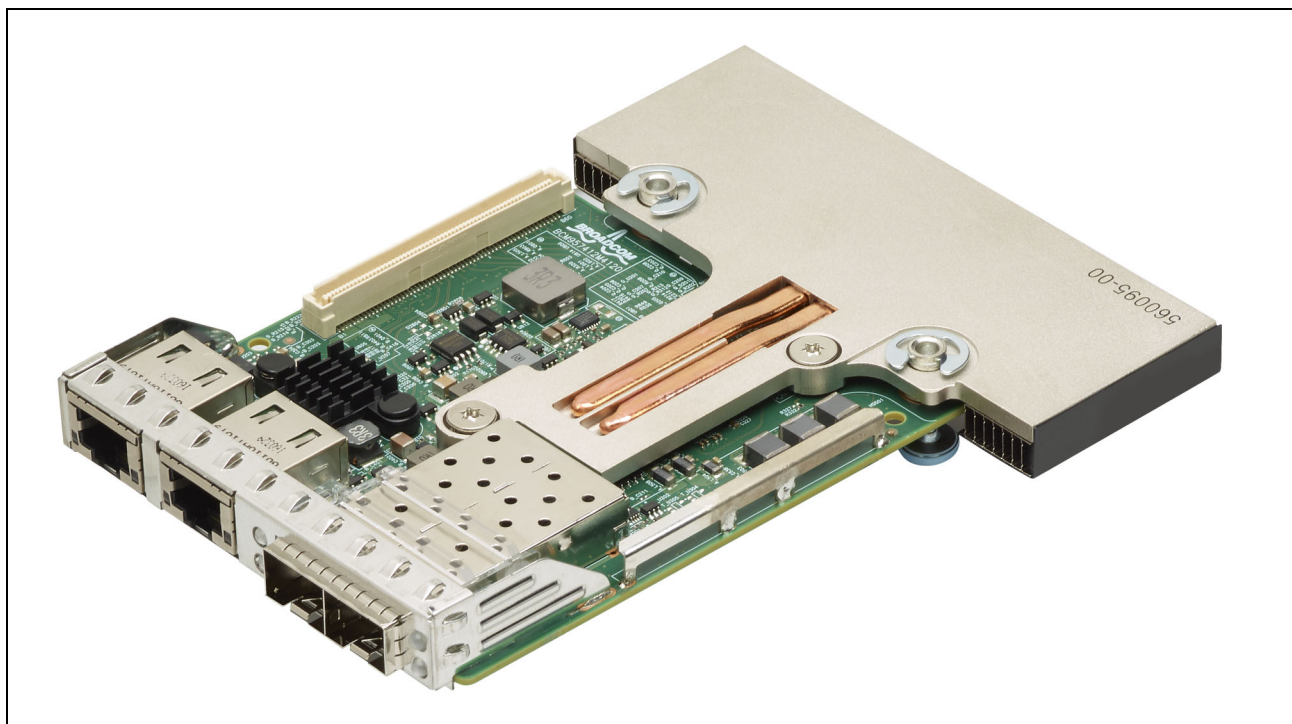
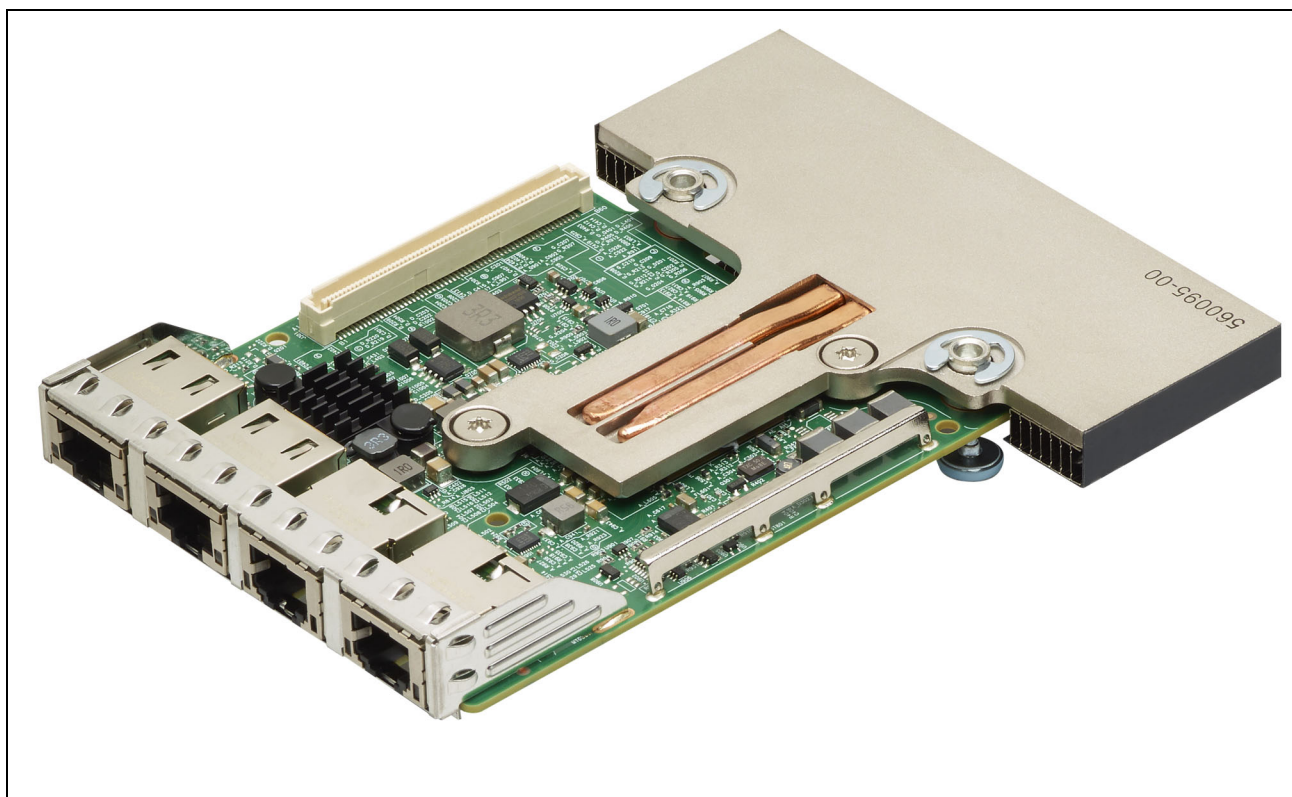


Figure 6: BCM957416M4160 Network Daughter Card (rNDC)



Network Link and Activity Indication

BCM957402AXXX/BCM957412AXXX

The SFP+ port has two LEDs to indicate traffic activities and link speed. The LEDs are visible through the cutout on the bracket as shown in [Figure 7](#). The LED functionality is described in [Table 6](#).

Figure 7: BCM957402AXXX/BCM957412AXXX Activity and Link LED Locations

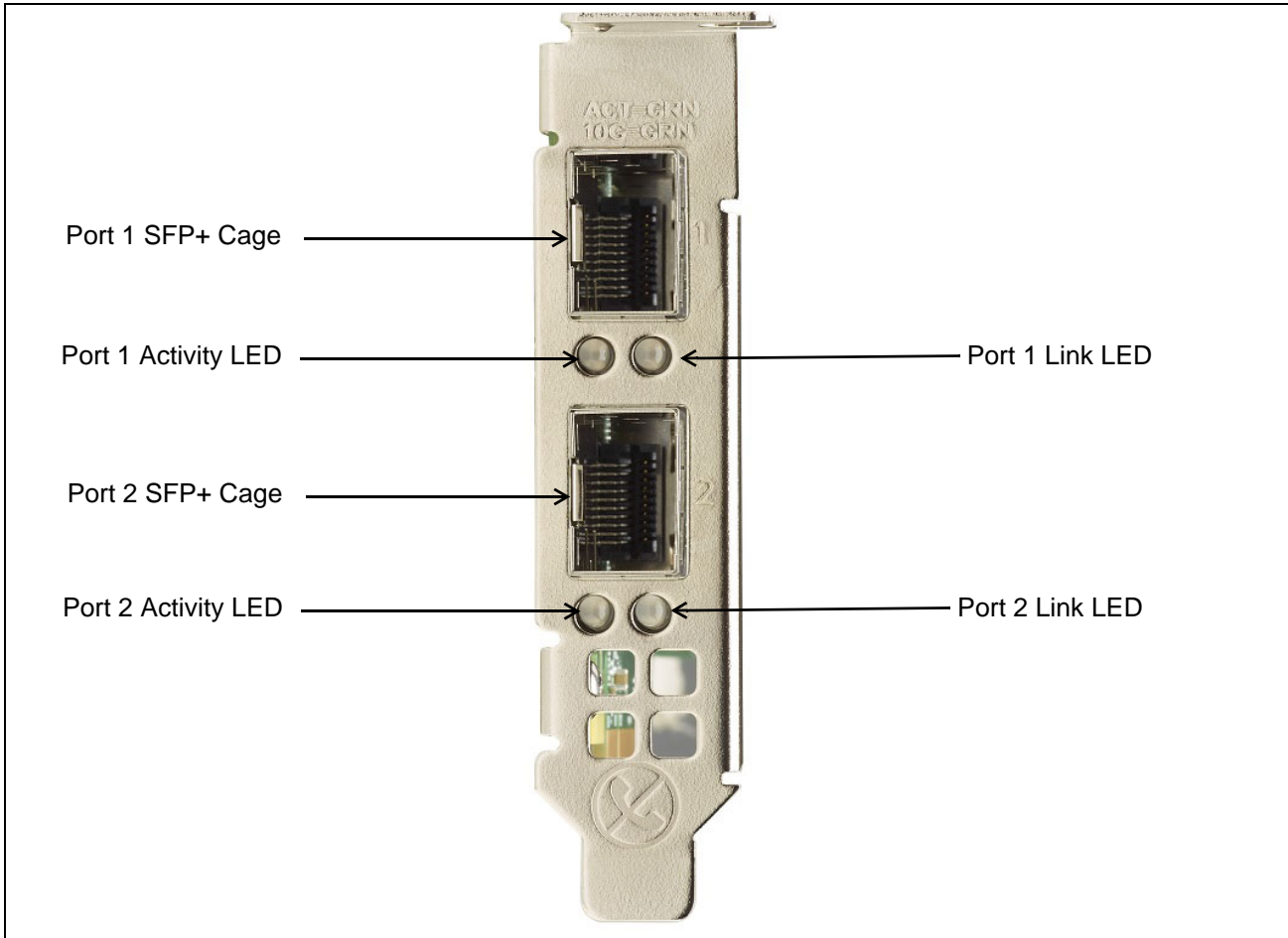


Table 6: BCM957402AXXX/BCM957412AXXX Activity and Link LED Locations

LED Type	Color/Behavior	Note
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 10 Gbps

BCM957404AXXX/BCM957414AXXX

The SFP28 port has two LEDs to indicate traffic activities and link speed. The LEDs are visible through the cutout on the bracket as shown in [Figure 8](#). The LED functionality is described in [Table 7](#).

Figure 8: BCM957404AXXX/BCM957414AXXX Activity and Link LED Locations

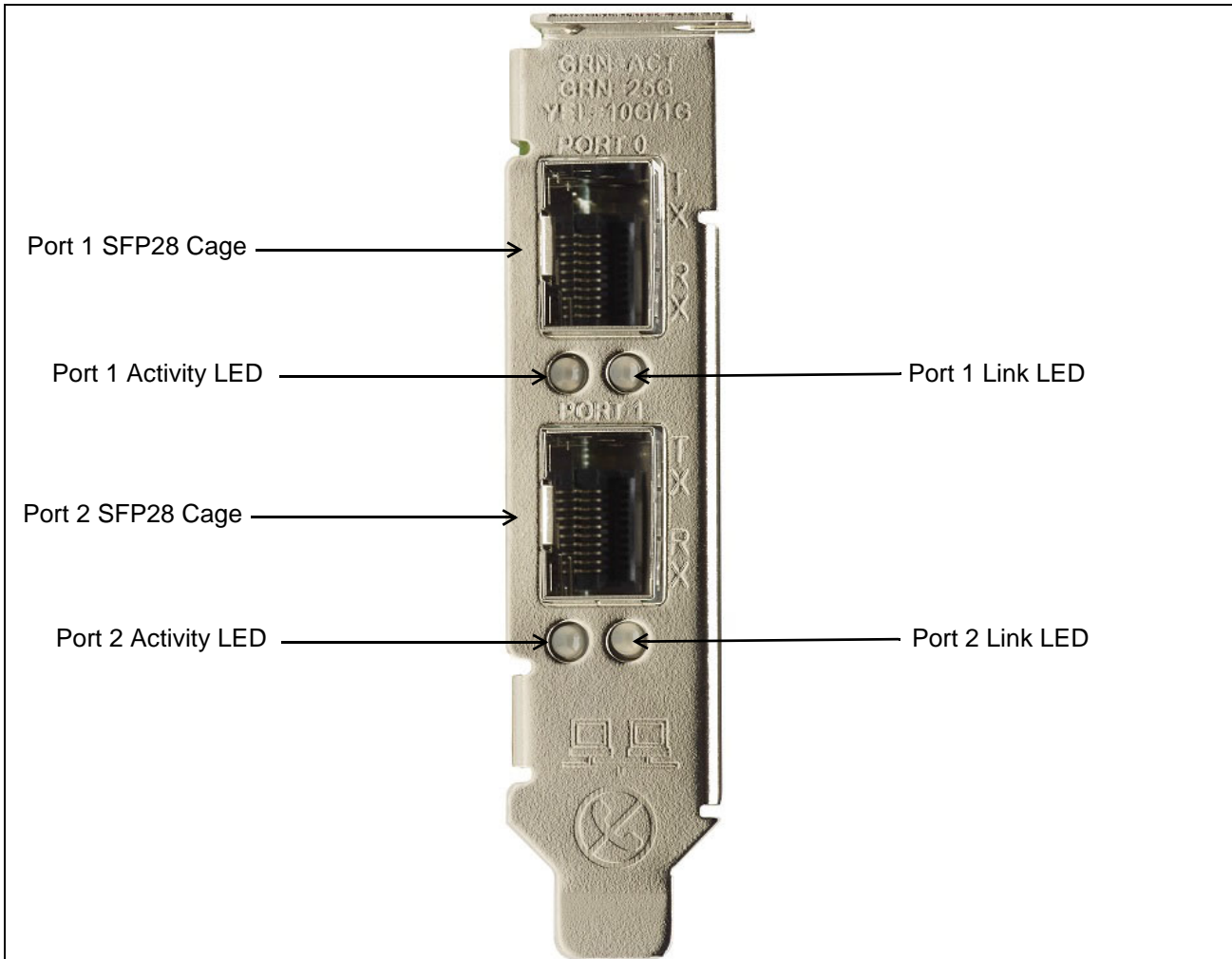


Table 7: BCM957404AXXX/BCM957414AXXX Activity and Link LED Locations

LED Type	Color/Behavior	Note
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 25 Gbps
	Yellow	Linked at 10 Gbps

BCM957406AXXX/BCM957416AXXX

The RJ-45 port has two LEDs to indicate traffic activities and link speed. The LEDs are visible through the cutout on the bracket as shown in [Figure 9](#). The LED functionality is described in [Table 8](#).

Figure 9: BCM957406AXXX/BCM957416AXXX Activity and Link LED Locations

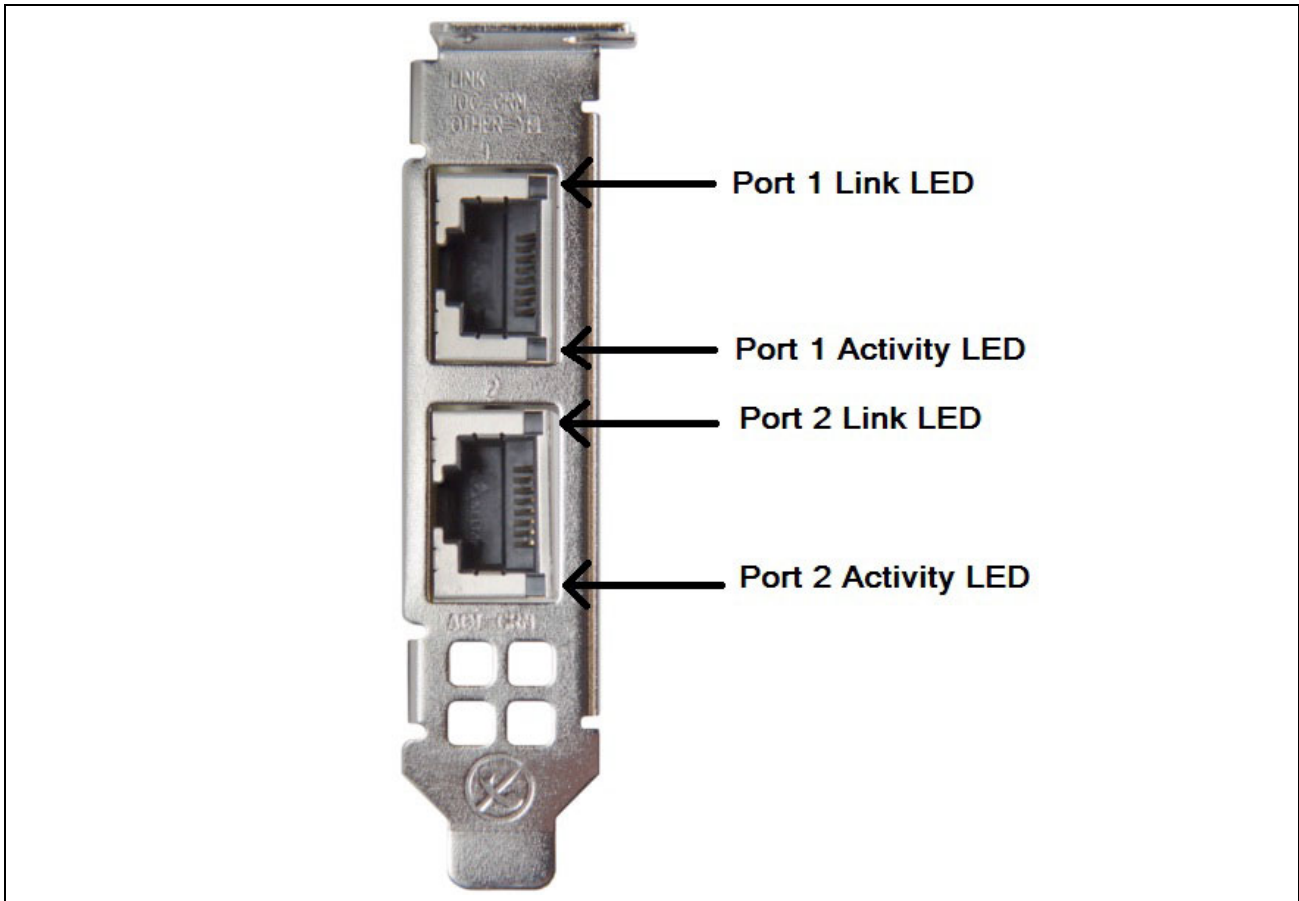


Table 8: BCM957406AXXX/BCM957416AXXX Activity and Link LED Locations

LED Type	Color/Behavior	Notes
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 10 Gbps
	Amber	Linked at 1 Gbps

BCM957414M4140D

The SFP28 port has two LEDs to indicate traffic activities and link speed. The LEDs are visible through the cutout on the bracket as shown in [Figure 10](#).

Figure 10: BCM957414M4140D Network Daughter Card (rNDC) Activity and Link LED Locations

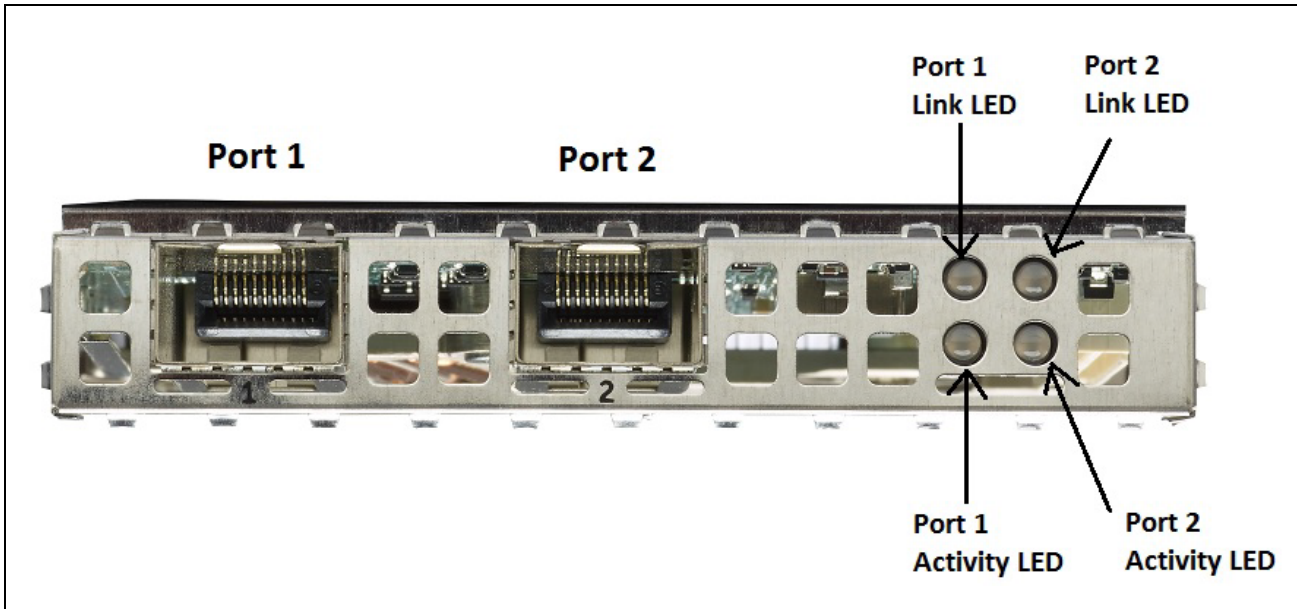


Table 9: BCM957414M4140D Network Daughter Card (rNDC) Activity and Link LED Locations

LED Type	Color/Behavior	Notes
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 25 Gbps
	Yellow	Linked at 10 Gbps

BCM957412M4120D

This rNDC has SFP+ and RJ-45 ports, each with two LEDs to indicate traffic activities and link speed. The LEDs are visible as shown in [Figure 11](#).

Figure 11: BCM957412M4120D Network Daughter Card (rNDC) Activity and Link LED Locations

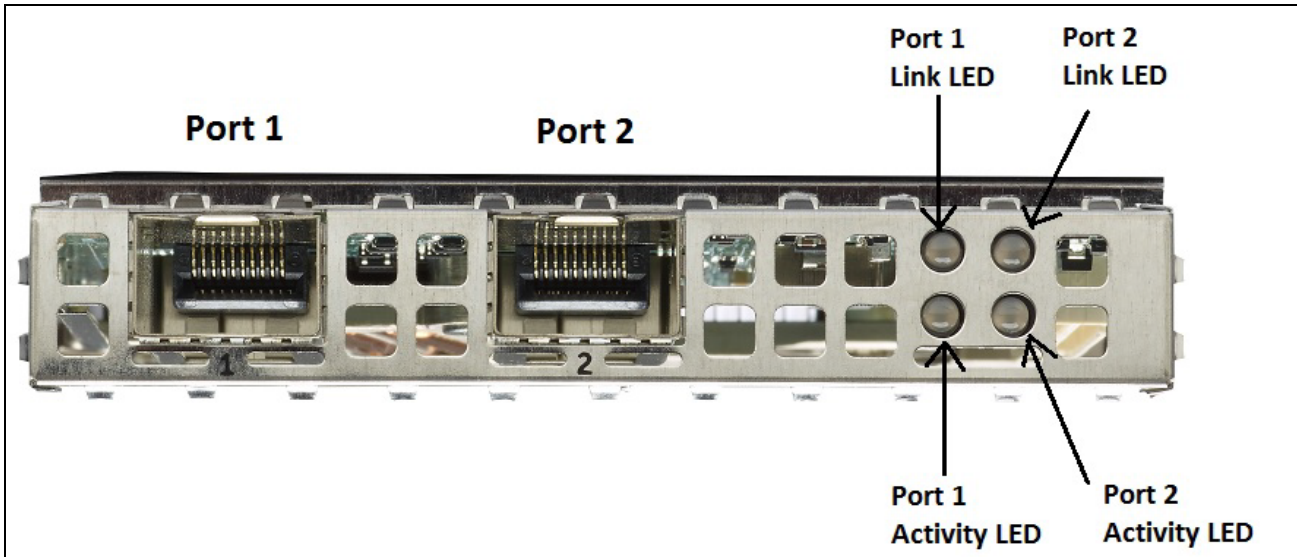


Table 10: BCM957412M4120D Network Daughter Card (rNDC) Activity and Link LED Locations SFP+ Port 1 and 2

LED Type	Color/Behavior	Notes
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 10 Gbps

Table 11: 1000BaseT Port 3 and 4

LED Type	Color/Behavior	Notes
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 1 Gbps
	Amber	Linked at 10/100 Mbps

BCM957416M4160

This rNDC has 10GBaseT and 1000BaseT RJ-45 ports, each with two LEDs to indicate traffic activities and link speed. The LEDs are visible as shown in [Figure 12](#).

Figure 12: BCM957416M4160 Network Daughter Card (rNDC) Activity and Link LED Locations

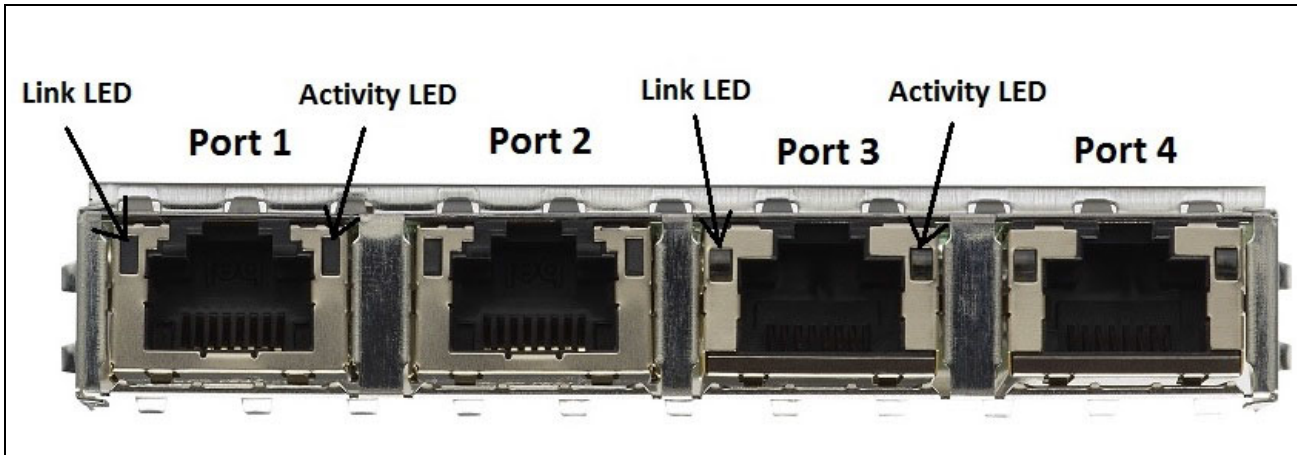


Table 12: BCM957416M4160 Network Daughter Card (rNDC) Activity and Link LED Locations 10GBaseT Port 1 and 2

LED Type	Color/Behavior	Notes
Activity	Off	No Activity
	Green blinking	Traffic Flowing Activity
Link	Off	No Link
	Green	Linked at 10 Gbps
	Amber	Linked at 1 Gbps

Features

Refer to the following sections for device features.

Software and Hardware Features

Table 13 provides a list of host interface features.

Table 13: Host Interface Features

Feature	Details
Host Interface	PCIe v3.0 (Gen 3: 8 GT/s; Gen 2: 5 GT/s; Gen 1: 2.5 GT/s).
Number of PCIe lanes	PCI-E Edge connector: x8.
Vital Product Data (VPD)	Supported.
Alternate Routing ID (ARI)	Supported.
Function Level Reset (FLR)	Supported.
Advanced Error Reporting	Supported.
PCIe ECNs	Support for TLP Processing Hints (TPH), Latency Tolerance Reporting (LTR), and Optimized Buffer Flush/Fill (OBFF).
MSI-X Interrupt vector per queue	1 per RSS queue, 1 per NetQueue, 1 per Virtual Machine Queue (VMQ).
IP Checksum Offload	Support for transmit and receive side.
TCP Checksum Offload	Support for transmit and receive side.
UDP Checksum Offload	Support for transmit and receive side.
NDIS TCP Large Send Offload	Support for LSOV1 and LSOV2.
NDIS Receive Segment Coalescing (RSC)	Support for Windows environments.
TCP Segmentation Offload (TSO)	Support for Linux and VMware environments.
Large Receive Offload (LRO)	Support for Linux and VMware environments.
Generic Receive Offload (GRO)	Support for Linux and VMware environments.
Receive Side Scaling (RSS)	Support for Windows, Linux, and VMware environments. Up to 8 queues/port supported for RSS.
Header-Payload Split	Enables the software TCP/IP stack to receive TCP/IP packets with header and payload data split into separate buffers. Supports Windows, Linux, and VMware environments.

Table 13: Host Interface Features (Cont.)

Feature	Details
Jumbo Frames	Supported.
iSCSI boot	Supported.
NIC Partitioning (NPAR)	Supports up to eight Physical Functions (PFs) per port, or up to 16 PFs per silicon. This option is configurable in NVRAM.
RDMA over Converge Ethernet (RoCE)	The BCM5741X supports RoCE v1/v2 for Windows, Linux, and VMware.
Data Center Bridging (DCB)	The BCM5741X supports DCBX (IEEE and CEE specification), PFC, and AVB.
NCSI (Network Controller Sideband Interface)	Supported.
Wake on LAN (WOL)	Supported on rNDC with 10GBase-T, SFP+, and SFP28 interfaces.
PXE boot	Supported.
UEFI boot	Supported.
Flow Control (Pause)	Supported.
Auto negotiation	Supported.
802.1q VLAN	Supported.
Interrupt Moderation	Supported.
MAC/VLAN filters	Supported.

Virtualization Features

Table 14 lists the virtualization features of the NetXtreme-E.

Table 14: Virtualization Features

Feature	Details
Linux KVM Multiqueue	Supported.
VMware NetQueue	Supported.
NDIS Virtual Machine Queue (VMQ)	Supported.
Virtual eXtensible LAN (VXLAN) – Aware stateless offloads (IP/UDP/TCP checksum offloads)	Supported.
Generic Routing Encapsulation (GRE) – Aware stateless offloads (IP/UDP/TCP checksum offloads)	Supported.
Network Virtualization using Generic Routing Encapsulation (NVGRE) – Aware stateless offloads	Supported.
IP-in-IP aware stateless offloads (IP/UDP/TCP checksum offloads)	Supported
SR-IOV v1.0	128 Virtual Functions (VFs) for Guest Operating Systems (GOS) per device. MSI-X vector per VF is set to 16.
MSI-X vector port	74 per port default value (two port configuration). 16 per VF and is configurable in HII and CCM.

VXLAN

A Virtual eXtensible Local Area Network (VXLAN), defined in IETF RFC 7348, is used to address the need for overlay networks within virtualized data centers accommodating multiple tenants. VXLAN is a Layer 2 overlay or tunneling scheme over a Layer 3 network. Only VMs within the same VXLAN segment can communicate with each other.

NVGRE/GRE/IP-in-IP/Geneve

Network Virtualization using GRE (NVGRE), defined in IETF RFC 7637, is similar to a VXLAN.

Stateless Offloads

RSS

Receive Side Scaling (RSS) uses a Toeplitz algorithm which uses 4 tuple match on the received frames and forwards it to a deterministic CPU for frame processing. This allows streamlined frame processing and balances CPU utilization. An indirection table is used to map the stream to a CPU.

Symmetric RSS allows the mapping of packets of a given TCP or UDP flow to the same receive queue.

TPA

Transparent Packet Aggregation (TPA) is a technique where received frames of the same 4 tuple matched frames are aggregated together and then indicated to the network stack. Each entry in the TPA context is identified by the 4 tuple: Source IP, destination IP, source TCP port, and destination TCP port. TPA improves system performance by reducing interrupts for network traffic and lessening CPU overhead.

Header-Payload Split

Header-payload split is a feature that enables the software TCP/IP stack to receive TCP/IP packets with header and payload data split into separate buffers. The support for this feature is available in both Windows and Linux environments. The following are potential benefits of header-payload split:

- The header-payload split enables compact and efficient caching of packet headers into host CPU caches. This can result in a receive side TCP/IP performance improvement.
- Header-payload splitting enables page flipping and zero copy operations by the host TCP/IP stack. This can further improve the performance of the receive path.

UDP Fragmentation Offload

UDP Fragmentation Offload (UFO) is a feature that enables the software stack to offload fragmentation of UDP/IP datagrams into UDP/IP packets. The support for this feature is only available in the Linux environment. The following is a potential benefit of UFO:

- The UFO enables the NIC to handle fragmentation of a UDP datagram into UDP/IP packets. This can result in the reduction of CPU overhead for transmit side UDP/IP processing.

Stateless Transport Tunnel Offload

Stateless Transport Tunnel Offload (STT) is a tunnel encapsulation that enables overlay networks in virtualized data centers. STT uses IP-based encapsulation with a TCP-like header. There is no TCP connection state associated with the tunnel and that is why STT is stateless. Open Virtual Switch (OVS) uses STT.

An STT frame contains the STT frame header and payload. The payload of the STT frame is an untagged Ethernet frame. The STT frame header and encapsulated payload are treated as the TCP payload and TCP-like header. The IP header (IPv4 or IPv6) and Ethernet header are created for each STT segment that is transmitted.

Multiqueue Support for OS

NDIS VMQ

The NDIS Virtual Machine Queue (VMQ) is a feature that is supported by Microsoft to improve Hyper-V network performance. The VMQ feature supports packet classification based on the destination MAC address to return received packets on different completion queues. This packet classification combined with the ability to DMA packets directly into a virtual machine's memory allows the scaling of virtual machines across multiple processors.

Refer to the ["Windows Driver Advanced Properties and Event Log Messages"](#) on page 34 for information on VMQ.

VMWare NetQueue

The VMware NetQueue is a feature that is similar to Microsoft's NDIS VMQ feature. The NetQueue feature supports packet classification based on the destination MAC address and VLAN to return received packets on different NetQueues. This packet classification combined with the ability to DMA packets directly into a virtual machine's memory allows the scaling of virtual machines across multiple processors.

KVM/Xen Multiqueue

KVM/Multiqueue returns the frames to different queues of the host stack by classifying the incoming frame by processing the received packet's destination MAC address and or 802.1Q VLAN tag. The classification combined with the ability to DMA the frames directly into a virtual machine's memory allows scaling of virtual machines across multiple processors.

SR-IOV Configuration Support Matrix

- Windows VF over Windows hypervisor
- Windows VF and Linux VF over VMware hypervisor
- Linux VF over Linux KVM

SR-IOV

The PCI-SIG defines optional support for Single-Root IO Virtualization (SR-IOV). SR-IOV is designed to allow access of the VM directly to the device using Virtual Functions (VFs). The NIC Physical Function (PF) is divided into multiple virtual functions and each VF is presented as a PF to VMs.

SR-IOV uses IOMMU functionality to translate PCI-E virtual addresses to physical addresses by using a translation table.

The number of Physical Functions (PFs) and Virtual Functions (VFs) are managed through the UEFI HII menu, the CCM, and through NVRAM configurations. SRIOV can be supported in combination with NPAR mode.

Network Partitioning (NPAR)

The Network Partitioning (NPAR) feature allows a single physical network interface port to appear to the system as multiple network device functions. When NPAR mode is enabled, the NetXtreme-E device is enumerated as multiple PCIe physical functions (PF). Each PF or “partition” is assigned a separate PCIe function ID on initial power on. The original PCIe definition allowed for eight PFs per device. For Alternative Routing-ID (ARI) capable systems, Broadcom NetXtreme-E adapters support up to 16 PFs per device. Each partition is assigned its own configuration space, BAR address, and MAC address allowing it to operate independently. Partitions support direct assignment to VMs, VLANs, etc., just as any other physical interface.



Note: In the **System Setup > Device Settings > [Broadcom 5741x Device] > Device Level Configuration** page, the user can enable **NParEP** to allow the NXE adapter to support up to 16 PFs per device. For 2 port devices, this means up to eight PFs for each port.

RDMA over Converge Ethernet – RoCE

Remote Direct Memory Access (RDMA) over Converge Ethernet (RoCE) is a complete hardware offload feature in the BCM5741X that allows RDMA functionality over an Ethernet network. RoCE functionality is available in user mode and kernel mode application. RoCE Physical Functions (PF) and SRIOV Virtual Functions (VF) are available in single function mode and in multi-function mode (NIC Partitioning mode). Broadcom supports RoCE in Windows, Linux, and VMWare.

Please refer to the following links for RDMA support for each operating system:

Windows

[https://technet.microsoft.com/en-us/library/jj134210\(v=ws.11\).aspx](https://technet.microsoft.com/en-us/library/jj134210(v=ws.11).aspx)

Redhat Linux

https://access.redhat.com/documentation/en-US/Red_Hat_Enterprise_Linux/7/html/Networking_Guide/ch-Configure_InfiniBand_and_RDMA_Networks.html

VMware

<https://pubs.vmware.com/vsphere-65/index.jsp?topic=%2Fcom.vmware.vsphere.networking.doc%2FGUID-4A5EBD44-FB1E-4A83-BB47-BBC65181E1C2.html>

Supported Combinations

The following sections describe the supported feature combinations for this device.

NPAR, SR-IOV, and RoCE

[Table 15](#) provides the supported feature combinations of NPAR, SR-IOV, and RoCE.

Table 15: NPAR, SR-IOV, and RoCE

SW Feature	Notes
NPAR	Up to 8 PFs or 16 PFs
SR-IOV	Up to 128 VFs (total per chip)
RoCE on PFs	Up to 4 PFs
RoCE on VFs	Valid for VFs attached to RoCE-enabled PFs
Host OS	Linux, Windows, ESXi (no vRDMA support)
Guest OS	Linux and Windows
DCB	Up to two COS per port with non-shared reserved memory

NPAR, SR-IOV, and DPDK

[Table 16](#) provides the supported feature combinations of NPAR, SR-IOV, and DPDK.

Table 16: NPAR, SR-IOV, and DPDK

SW Feature	Notes
NPAR	Up to 8 PFs or 16 PFs
SR-IOV	Up to 128 VFs (total per chip)
DPDK	Supported only as a VF
Host OS	Linux
Guest OS	DPDK (Linux)

Unsupported Combinations

The combination of NPAR, SR-IOV, RoCE, and DPDK is not supported.

Installing the Hardware

Safety Precautions



Caution! The adapter is being installed in a system that operates with voltages that can be lethal. Before removing the cover of the system, observe the following precautions to protect yourself and to prevent damage to the system components:

- Remove any metallic objects or jewelry from your hands and wrists.
- Make sure to use only insulated or nonconducting tools.
- Verify that the system is powered OFF and unplugged before you touch internal components.
- Install or remove adapters in a static-free environment. The use of a properly grounded wrist strap or other personal antistatic devices and an antistatic mat is strongly recommended.

System Requirements

Before you install the Broadcom NetXtreme-E Ethernet adapter, verify that the system meets the requirements listed for the operating system.

Hardware Requirements

Refer to the following list of hardware requirements:

- Dell 13G systems that meet operating system requirements.
- Dell 13G systems that support NetXtreme-E Ethernet cards.
- One open PCI-E Gen 3 x8 slot or an open PCIe Gen3 rNDC slots for a NIC Adapter with rNDC form factor.
- 4 GB memory or more (32 GB or more is recommended for virtualization applications and nominal network throughput performance).

Preinstallation Checklist

Refer to the following list before installing the NetXtreme-E device.

1. Verify that the server meets the hardware and software requirements listed in [“System Requirements”](#).
2. Verify that the server is using the latest BIOS.
3. If the system is active, shut it down.
4. When the system shutdown is complete, turn off the power and unplug the power cord.
5. Holding the adapter card by the edges, remove it from its shipping package and place it on an antistatic surface.
6. Check the adapter for visible signs of damage, particularly on the card edge connector. Never attempt to install a damaged adapter.

Installing the Adapter

The following instructions apply to installing the Broadcom NetXtreme-E Ethernet adapter (add-in NIC) into most servers. Refer to the manuals that are supplied with the server for details about performing these tasks on this particular server.

1. Review the “[Safety Precautions](#)” on page 26 and “[Preinstallation Checklist](#)” before installing the adapter. Ensure that the system power is OFF and unplugged from the power outlet, and that proper electrical grounding procedures have been followed.
2. Open the system case and select any empty PCI Express Gen 3 x8 slot.
3. Remove the blank cover-plate from the slot.
4. Align the adapter connector edge with the connector slot in the system.
5. Secure the adapter with the adapter clip or screw.
6. Close the system case and disconnect any personal antistatic devices.



Note: For NIC Adapters with rNDC form factor, locate an open RNDC slot or remove replace exiting default rNDC with the NetXtreme rNDC.

Connecting the Network Cables

Dell Ethernet switches are productized with SFP+/SFP28/QSFP28 ports that support up to 100 Gbps. These 100 Gbps ports can be divided into 4 x 25 Gbps SFP28 ports. QSFP ports can be connected to SFP28 ports using 4 x 25G SFP28 breakout cables.

Supported Cables and Modules

Table 17: Supported Cables and Modules

<i>Optical Module</i>	<i>Dell Part Number Adapters</i>	<i>Description</i>
FTLX8571D3BCL-DL	3G84K	BCM957404A4041DLPC, BCM957404A4041DC, BCM957402A4020DLPC, BCM957402A4020DC
FCLF-8521-3	8T47V	BCM957406A4060DLPC, BCM957406A4060DC
FTLF8536P4BCL	P7D7R	BCM957404A4041DLPC, BCM957404A4041DC
Methode DM7051	PGYJT	BCM57402X, BCM57404X, BCM57412X, BCM57414X
FTLX1471D3BCL-FC	RN84N	BCM57402X, BCM57404X, BCM57412X, BCM57414X
FTLX8574D3BNL	N8TDR	BCM57402X, BCM57404X, BCM57412X, BCM57414X
FTLF8536P4BNL-FC	HHHHC	BCM57402X, BCM57404X, BCM57412X, BCM57414X

Table 17: Supported Cables and Modules (Cont.)

Optical Module	Dell Part Number	Adapters	Description
FTLX8574D3BCL-FC or PLRXPLSCS43811	WTRD1	BCM57402X, BCM57404X, BCM57412X, BCM57414X	10 Gbps-SR SFP+ Transceiver

Note:

1. Direct Attach Cables (DAC) that conform to IEEE standards can be connected to the adapter.
2. Dell part HHHHC and N8TDR are required for the BCM957414M4140D.

Copper

The BCM957406AXXXX, BCM957416AXXXX, and BCM957416XXXX adapters have two RJ-45 connectors used for attaching the system to a CAT 6E Ethernet copper-wire segment.

SFP+

The BCM957402AXXXX, BCM957412AXXXX, and BCM957412MXXXX adapters have two SFP+ connectors used for attaching the system to a 10 Gbps Ethernet switch.

SFP28

The BCM957404AXXXX, BCM957414XXXX, and BCM957414AXXXX adapters have two SFP28 connectors used for attaching the system to a 100 Gbps Ethernet switch.

Software Packages and Installation

Refer to the following sections for information on software packages and installation.

Supported Operating Systems

Table 18 provides a list of supported operating systems.

Table 18: Supported Operating System List

OS Flavor	Distribution
Windows	Windows 2012 R2 or above
Linux	Redhat 6.9, Redhat 7.1 or above SLES 11 SP 4, SLES 12 SP 2 or above
VMWare	ESXi 6.0 U3 or above

Installing Drivers

Refer to the following sections for driver installation.

Windows

Dell DUP

Broadcom NetXtreme E series controller drivers can be installed using the driver DUP. The installer is provided in x64 executable format.

GUI Install

When the file is executed, a dialog box appears requesting user input. The installer supports the driver only option.

Silent Install

The executable can be silently executed using the command shown below.

Example:

```
Network_Driver_<version>.EXE /s /driveronly
```

INF Install

The Dell DUP is used to install drivers for Broadcom NetXtreme-E Ethernet controllers. Use the following command to extract the driver INF files from the Dell DUP:

```
Network_Driver_<version>.EXE /s /v"EXTRACTDRIVERS=c:\dell\drivers\network"
```

Once the files are extracted, INF installation is executed through "upgrade driver" functionality using the Device Manager (devmgmt.msc). Open the Device Manager, select the desired NIC, right click and select the upgrade driver to update it.

Linux

The Linux drivers are supplied in RPM, KMP, and source code format. To build the device driver from the source code using Linux, refer to the following example:

1. As a root user, log on to the Linux system.
2. scp or cp the driver tar ball on to the Linux system. A typical example is:

```
cp /var/run/media/usb/bnxt_en-<version>.tar.gz /root/
```

3. Execute the following command:

```
tar -zxvf /root/bnxt_en-<version>.tar.gz
```

4. Execute the following command:

```
cd bnxt_en-<version>
```

5. Execute the following command:

```
make; make install; modprobe -r bnxt_en; modprobe bnxt_en
```

For RDMA functionality, install both the `bnxt_en` and `bnxt_re` driver. Use `netxtreme-bnxt_en-<version>.tar.gz` instead of `bnxt_en-<version>.tar.gz`.

Module Install

RHEL

The driver image can be installed with one of the following options:

- Mount the `bnxt_en-x.x.x-rhelYuZ-x86_64-dd.iso` image using the Dell iDRAC Virtual console.
- Mount the `bnxt_en-x.x.x-rhelYuZ-x86_64-dd.iso` image from a CD/DVD.
- Copy the `bnxt_en-x.x.x-rhelYuZ-x86_64-dd.iso` image to a USB device and mount the device.

Start the OS install, push the tab key, and enter "linux dd". Continue the installation until the driver disk is requested and select the `bnxt_en` driver.

SLES

The driver image can be installed with one of the following options:

- Mount the `bnxt_en-x.x.x-rhelYuZ-x86_64-dd.iso` image using the Dell iDRAC Virtual console.
- Mount the `bnxt_en-x.x.x-rhelYuZ-x86_64-dd.iso` image from a CD/DVD.
- Extract the `bnxt_en-x.x.x-rhelYuZ-x86_64-dd.iso` image, copy the contents to a USB device, and mount the USB device.

Linux Ethtool Commands



Note: In [Table 19](#), ethX should be replaced with the actual interface name.

Table 19: Linux Ethtool Commands

Command	Description
<code>ethtool -s ethX speed 25000 autoneg off</code>	Set the speed. If the link is up on one port, the driver will not allow the other port to be set to an incompatible speed.
<code>ethtool -i ethX</code>	Output includes Package version, NIC BIOS version (boot code).
<code>ethtool -k ethX</code>	Show offload features.
<code>ethtool -K ethX tso off</code>	Turn off TSO.
<code>ethtool -K ethX gro off lro off</code>	Turn off GRO / LRO.
<code>ethtool -g ethX</code>	Show ring sizes.
<code>ethtool -G ethX rx N</code>	Set Ring sizes.
<code>ethtool -S ethX</code>	Get statistics.
<code>ethtool -l ethX</code>	Show number of rings.
<code>ethtool -L ethX rx 0 tx 0 combined M</code>	Set number of rings.
<code>ethtool -C ethX rx-frames N</code>	Set interrupt coalescing. Other parameters supported are: rx-usecs, rx-frames, rx-usecs-irq, rx-frames-irq, tx-usecs, tx-frames, tx-usecs-irq, tx-frames-irq.
<code>ethtool -x ethX</code>	Show RSS flow hash indirection table and RSS key.
<code>ethtool -s ethX autoneg on speed 10000 duplex full</code>	Enable Autoneg (see "Auto-Negotiation Configuration" on page 40 for more details)
<code>ethtool --show-eee ethX</code>	Show EEE state.
<code>ethtool --set-eee ethX eee off</code>	Disable EEE.
<code>ethtool --set-eee ethX eee on tx-lpi off</code>	Enable EEE, but disable LPI.
<code>ethtool -L ethX combined 1 rx 0 tx 0</code>	Disable RSS. Set the combined channels to 1.
<code>ethtool -K ethX ntuple off</code>	Disable Accelerated RFS by disabling ntuple filters.
<code>ethtool -K ethX ntuple on</code>	Enable Accelerated RFS.
<code>Ethtool -t ethX</code>	Performs various diagnostic self-tests.
<code>echo 32768 > /proc/sys/net/core/rps_sock_flow_entries</code>	Enable RFS for Ring X.
<code>echo 2048 > /sys/class/net/ethX/queues/rx-X/rps_flow_cnt</code>	
<code>sysctl -w net.core.busy_read=50</code>	This sets the time to busy read the device's receive ring to 50 usecs. For socket applications waiting for data to arrive, using this method can decrease latency by 2 or 3 usecs typically at the expense of higher CPU utilization.
<code>echo 4 > /sys/bus/pci/devices/0000:82:00.0/sriov_numvfs</code>	Enable SR-IOV with four VFs on bus 82, Device 0 and Function 0.

Table 19: Linux Etool Commands (Cont.)

Command	Description
<code>ip link set ethX vf 0 mac 00:12:34:56:78:9a</code>	Set VF MAC address.
<code>ip link set ethX vf 0 state enable</code>	Set VF link state for VF 0.
<code>ip link set ethX vf 0 vlan 100</code>	Set VF 0 with VLAN ID 100.

VMware

The ESX drivers are provided in VMware standard VIB format.

1. To install the Ethernet and RDMA driver, issue the commands:

```
$ esxcli software vib install --no-sig-check -v <bnxtnet>-<driver version>.vib
```

```
$ esxcli software vib install --no-sig-check -v <bnxtroce>-<driver version>.vib
```

2. A system reboot is required for the new driver to take effect.

Other useful VMware commands shown in [Table 20](#).



Note: In [Table 20](#), `vmnicX` should be replaced with the actual interface name.



Note: `$ kill -HUP $(cat /var/run/vmware/vmkdevmgr.pid)` This command is required after `vmkload_mod bnxtnet` for successful module bring up.

Table 20: VMware Commands

Command	Description
<code>esxcli software vib list grep bnx</code>	List the VIBs installed to see whether the bnxt driver installed successfully.
<code>esxcfg-module -l bnxtnet</code>	Print module info on to screen.
<code>esxcli network get -n vmnicX</code>	Get vmnicX properties.
<code>esxcfg-module -g bnxtnet</code>	Print module parameters.
<code>esxcfg-module -s 'multi_rx_filters=2 disable_tap=0 max_vfs=0,0 RSS=0'</code>	Set the module parameters.
<code>vmkload_mod -u bnxtnet</code>	Unload bnxtnet module.
<code>vmkload_mod bnxtnet</code>	Load bnxtnet module.
<code>esxcli network nic set -n vmnicX -D full -S 25000</code>	Set the speed and duplex of vmnicX.
<code>esxcli network nic down -n vmnicX</code>	Disable vmnicX.
<code>esxcli network nic up -n vmnic6</code>	Enable vmnicX.
<code>bnxtnetcli -s -n vmnic6 -S "25000"</code>	Set the link speed. Bnxtnetcli is needed for older ESX versions to support the 25G speed setting.

Firmware Update

The NIC firmware can be updated using one of the following methods:

- Using the Dell Update Package (DUP) when the system is in the OS booted state. This method only applies to the Windows and Linux operating systems.
- Using the Dell iDRAC – Lifecycle Controller. This method can be used regardless of the operating system. If the system is running VMware, use the Lifecycle Controller to upgrade the firmware.

Refer to product support page at <http://www.dell.com/support>

Dell Update Package

Refer to the following sections to use the Dell Update Package (DUP):

Windows

Broadcom NetXtreme-E series controller firmware can be upgraded using the Dell DUP package. The executable is provided in standard Windows x64 executable format. Double-click on the file to execute it.

DUP packages can be downloaded from <http://support.dell.com>

Linux

The Dell Linux DUP is provided in x86_64 executable format. Use the standard Linux `chmod` to update the execute permission and run the executable. Refer to the following example:

1. Login to Linux.
2. `scp` or `cp` the DUP executable on to file system. A typical example is:

```
cp /var/run/media/usb/Network_Firmware_<version>.BIN /root/
```

3. Execute the following command:

```
chmod 755 Network_Firmware_<version>.BIN
```

4. Execute the following command:

```
./Network_Firmware_<version>.BIN
```

A reboot is needed to activate the new firmware.

Windows Driver Advanced Properties and Event Log Messages

Driver Advanced Properties

The Windows driver advanced properties are shown in [Table 21](#).

Table 21: Windows Driver Advanced Properties

Driver Key	Parameters	Description
Encapsulated Task offload	Enable or Disable	Used for configuring NVGRE encapsulated task offload.
Energy Efficient Ethernet	Enable or Disable	EEE enabled for Copper ports and Disabled for SFP+ or SFP28 ports. This feature is only enabled for the BCM957406A4060 adapter.
Flow control	TX or RX or TX/RX enable	Configure flow control on RX or TX or both sides.
Interrupt Moderation	Enable or Disable	Default Enabled. Allows frames to be batch processed by saving CPU time.
Jumbo packet	1514 or 4088 or 9014	Jumbo packet size.
Large Send offload V2 (IPV4)	Enable or Disable	LSO for IPV4.
Large Send offload V2 (IPV6)	Enable or Disable	LSO for IPV6.
Locally Administered Address	User entered MAC address.	Override default hardware MAC address after OS boot.
Max Number of RSS Queues	2, 4 or 8.	Default is 8. Allows user to configure Receive Side Scaling queues.
Priority and VLAN	Priority and VLAN Disable, Priority enabled, VLAN enabled, Priority and VLAN enabled.	Default Enabled. Used for configuring 802.1Q and 802.1P.
Receive Buffer (0=Auto)	Increments of 500.	Default is Auto.
Receive Side Scaling	Enable or Disable.	Default Enabled
Receive Segment Coalescing (IPV4)	Enable or Disable.	Default Enabled
Receive Segment Coalescing (IPV6)	Enable or Disable.	Default Enabled
RSS load balancing profile	NUMA scaling static, Closest processor, Closest processor static, conservative scaling, NUMA scaling.	Default NUMA scaling static.
Speed and Duplex	1 Gbps or 10 Gbps or 25 Gbps or Auto Negotiation.	10 Gbps Copper ports can Auto negotiate speeds, whereas 25 Gbps ports are set to forced speeds.

Table 21: Windows Driver Advanced Properties (Cont.)

Driver Key	Parameters	Description
SR-IOV	Enable or Disable.	Default Enabled. This parameter works in conjunction with HW configured SR-IOV and BIOS configured SR-IOV setting.
TCP/UDP checksum offload IPV4	TX/RX enabled, TX enabled or RX Enabled or offload disabled.	Default RX and TX enabled.
TCP/UDP checksum offload IPV6	TX/RX enabled, TX enabled or RX Enabled or offload disabled.	Default RX and TX enabled.
Transmit Buffers (0=Auto)	Increment of 50.	Default Auto.
Virtual Machine Queue	Enable or Disable.	Default Enabled.
VLAN ID	User configurable number.	Default 0.

Event Log Messages

Table 22 provides the Event Log messages logged by the Windows NDIS driver to the Event Logs.

Table 22: Windows Event Log Messages

Message ID	Comment
0x0001	Failed Memory allocation.
0x0002	Link Down Detected.
0x0003	Link up detected.
0x0009	Link 1000 Full.
0x000A	Link 2500 Full.
0x000b	Initialization successful.
0x000c	Miniport Reset.
0x000d	Failed Initialization.
0x000E	Link 10Gb successful.
0x000F	Failed Driver Layer Binding.
0x0011	Failed to set Attributes.
0x0012	Failed scatter gather DMA.
0x0013	Failed default Queue initialization.
0x0014	Incompatible firmware version.
0x0015	Single interrupt.

Table 23: Event Log Messages

0x0016	Firmware failed to respond within allocated time.
0x0017	Firmware returned failure status.
0x0018	Firmware is in unknown state.
0x0019	Optics Module is not supported.

Table 23: Event Log Messages (Cont.)

0x001A	Incompatible speed selection between Port 1 and Port 2. Reported link speeds are correct and might not match Speed and Duplex setting.
0x001B	Incompatible speed selection between Port 1 and Port 2. Link configuration became illegal.
0x001C	Network controller configured for 25Gb full-duplex link.
0x0020	RDMA support initialization failed.
0x0021	Device's RDMA firmware is incompatible with this driver.
0x0022	Doorbell BAR size is too small for RDMA.
0x0023	RDMA restart upon device reset failed.
0x0024	RDMA restart upon system power up failed
0x0025	RDMA startup failed. Not enough resources.
0x0026	RDMA not enabled in firmware.
0x0027	Start failed, a MAC address is not set.
0x0028	Transmit stall detected. TX flow control will be disabled from now on.

Teaming

Windows

The Broadcom NetXtreme-E devices installed on Dell platforms can participate in NIC teaming functionality using the Microsoft teaming solution. Refer to Microsoft public documentation described in the following link: <https://www.microsoft.com/en-us/download/details.aspx?id=40319>

Microsoft LBFO is a native teaming driver that can be used in the Windows OS. The teaming driver also provides VLAN tagging capabilities.

Linux

Linux bonding is used for teaming under Linux. The concept is loading the bonding driver and adding team members to the bond which would load-balance the traffic.

Use the following steps to setup Linux bonding:

1. Execute the following command:

```
modprobe bonding mode="balance-alb". This will create a bond interface.
```

2. Add bond clients to the bond interface. An example is shown below:

```
ifenslave bond0 ethX; ifenslave bond0 ethY
```

3. Assign an IPV4 address to bond the interface using `ifconfig bond0 IPV4Address netmask NetMask up`. The IPV4Address and NetMask are an IPV4 address and the associated network mask.



Note: IPV4 address should be replaced with the actual network IPV4 address. NetMask should be replaced by the actual IPV4 network mask.

4. Assign an IPV6 address to bond the interface using `ifconfig bond0 IPV6Address netmask NetMask up`. The IPV6Address and NetMask are an IPV6 address and the associated network mask.



Note: IPV6 address should be replaced with the actual network IPV6 address. NetMask should be replaced by the actual IPV6 network mask.

Refer to the Linux Bonding documentation for advanced configurations.

System-level Configuration

Refer to the following sections for information on system-level NIC configuration.

UEFI HII Menu

Broadcom NetXtreme-E series controllers can be configured for pre-boot, iscsi and advanced configuration such as SR-IOV using HII (Human Interface) menu.

To configure the settings, during system boot select F2 -> System Setup -> Device Settings. Select the desired network adapter for viewing and changing the configuration.

Main Configuration Page

This page displays current network link status, PCI-E Bus:Device:Function, MAC address of the adapter and the Ethernet device.

10GBaseT card will allow the user to enable or disable Energy Efficient Ethernet (EEE).

Firmware Image Properties

Main configuration page -> Firmware Image properties displays, Family version which consists of version numbers of controller BIOS, Multi Boot Agent (MBA), UEFI, iSCSI and Comprehensive Configuration Management (CCM) version numbers.

Device Level Configuration

Main configuration page -> Device level configuration allows the user to enable SR-IOV mode, number of virtual functions per physical function, MSI-X vectors per Virtual function, and the Max number of physical function MSI-X vectors.

NIC Configuration

NIC configuration -> Legacy boot protocol is used to select and configure PXE, iSCSI, or disable legacy boot mode. The boot strap type can be Auto, int18h (interrupt 18h), int19h (interrupt 19h), or BBS.

MBA and iSCSI can also be configured using CCM. Legacy BIOS mode uses CCM for configuration. The hide setup prompt can be used for disabling or enabling the banner display.

VLAN for PXE can be enabled or disabled and the VLAN ID can be configured by the user. Refer to the [“Auto-Negotiation Configuration” on page 40](#) for details on link speed setting options.

iSCSI Configuration

iSCSI boot configuration can be set through the Main configuration page -> iSCSI configuration. Parameters such as IPV4 or IPV6, iSCSI initiator, or the iSCSI target can be set through this page.

Refer to ["iSCSI Boot" on page 49](#) for detailed configuration information.

Comprehensive Configuration Management

Preboot configuration can be configured using the Comprehensive Configuration Management (CCM) menu option. During the system BIOS POST, the Broadcom banner message will be displayed with an option to change the parameters through the Control-S menu. When Control-S is pressed, a device list will be populated with all the Broadcom network adapters found in the system. Select the desired NIC for configuration.

Device Hardware Configuration

Parameters that can be configured using this section are the same as the HII menu "Device level configuration".

MBA Configuration Menu

Parameters that can be configured using this section are the same as HII menu "NIC configuration".

iSCSI Boot Main Menu

Parameters that can be configured using this section are same as HII menu "iSCSI configuration".

Auto-Negotiation Configuration



Note: In NPAR (NIC partitioning) devices where one port is shared by multiple PCI functions, the port speed is preconfigured and cannot be changed by the driver.

The Broadcom NetXtreme-E controller supports the following auto-negotiation features:

- Link speed auto-negotiation
- Pause/Flow Control auto-negotiation
- FEC – Forward Error Correction auto-negotiation



Note: Regarding link speed AN, when using SFP+, SFP28 connectors, use DAC or Multimode Optical transceivers capable of supporting AN. Ensure that the link partner port has been set to the matching auto-negotiation protocol. For example, if the local Broadcom port is set to IEEE 802.3by AN protocol, the link partner must support AN and must be set to IEEE 802.3by AN protocol.



Note: For Dual ports NetXtreme-E network controllers, 10 Gbps and 25 Gbps are not supported combination of link speed.

The supported combination of link speed settings for two ports NetXtreme-E network controller are shown in [Table 24 on page 41](#).

Table 24: Supported Combination of Link Speed Settings

Port1 Link Speed Setting	Port 2 Link Setting									
	Forced 1G	Forced 10G	Forced 25G	AN Enabled {1G}	AN Enabled {10G}	AN Enabled {25G}	AN Enabled {1/10G}	AN Enabled {1/25G}	AN Enabled {10/25G}	AN Enabled {1/10/25G}
Forced 1G	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN
	P2: no AN	P2: no AN	P2: no AN	P2: {1G}	P2: AN {10G}	P2: AN {25G}	P2: AN {1/10G}	P2: AN {1/25G}	P2: AN {10/25G}	P2: AN {1/10/25G}
Forced 10G	P1: no AN	P1: no AN	Not supported	P1: no AN	P1: no AN	Not supported	P1: no AN	P1: no AN	P1: no AN	P1: no AN
	P2: no AN	P2: no AN		P2: {1G}	P2: {10G}		P2: AN {1/10G}	P2: AN {1G}	P2: AN {10G}	P2: AN {1/10G}
Forced 25G	P1: no AN	Not supported	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN	P1: no AN
	P2: no AN		P2: no AN	P2: no AN	P2: no AN	P2: AN {1G}	P2: AN {1/25G}	P2: AN {25G}	P2: AN {1/25G}	
AN Enabled {1G}	P1: {1G}	P1: {1G}	P1: {1G}	P1: AN {1G}	P1: AN {1G}	P1: AN {1G}	P1: AN {1G}	P1: AN {1G}	P1: AN {1G}	P1: AN {1G}
	P2: no AN	P2: no AN	P2: no AN	P2: AN {1G}	P2: AN {10G}	P2: AN {25G}	P2: AN {1/10G}	P2: AN {1/25G}	P2: AN {10/25G}	P2: AN {1/10/25G}
AN Enabled {10G}	P1: AN {10G}	P1: AN {25G}	Not supported	P1: AN {10G}	P1: AN {10G}	Not supported	P1: AN {25G}	P1: AN {10G}	P1: AN {10G}	P1: AN {10G}
	P2: no AN	P2: no AN		P2: AN {1G}	P2: AN {10G}		P2: AN {1G}	P2: AN {10G}	P2: AN {10G}	P2: AN {1/10G}
AN Enabled {25G}	P1: AN {25G}	Not supported	P1: AN {25G}	P1: AN {25G}	Not supported	P1: AN {25G}	P1: AN {1/10G}	P1: AN {25G}	P1: AN {25G}	P1: AN {25G}
	P2: no AN		P2: no AN	P2: AN {1G}		P2: AN {25G}	P2: AN {1/10G}	P2: AN {1/25G}	P2: AN {25G}	P2: AN {1/25G}
AN Enabled {1/10G}	P1: AN {1/10G}	P1: AN {1/10G}	P1: AN {1G}	P1: AN {1/10G}	P1: AN {1/10G}	P1: AN {1/10G}	P1: AN {1/25G}	P1: AN {1G}	P1: AN {1/10G}	P1: AN {1/10G}
	P2: no AN	P2: no AN	P2: no AN	P2: AN {1G}	P2: AN {10G}	P2: AN {25G}	P2: AN {1/10G}	P2: AN {1G}	P2: AN {10G}	P2: AN {1/10G}
AN Enabled {1/25G}	P1: AN {1/25G}	P1: {1G}	P1: AN {1/25G}	P1: AN {1/25G}	P1: AN {1G}	P1: AN {1/25G}	P1: AN {10/25G}	P1: AN {1/25G}	P1: AN {1/25G}	P1: AN {1/25G}
	P2: no AN	P2: no AN	P2: no AN	P2: AN {1G}	P2: AN {10G}	P2: AN {25G}	P2: AN {1/10G}	P2: AN {1/25G}	P2: AN {25G}	P2: AN {1/25G}
AN Enabled {10/25G}	P1: AN {10/25G}	P1: {10G}	P1: AN {25G}	P1: AN {10/25G}	P1: AN {10G}	P1: AN {25G}	P1: AN {1/10/25G}	P1: AN {25G}	P1: AN {10/25G}	P1: AN {10/25G}
	P2: no AN	P2: no AN	P2: no AN	P2: AN {1G}	P2: AN {10G}	P2: AN {25G}	P2: AN {1/10G}	P2: AN {25G}	P2: AN {10/25G}	P2: AN {1/10/25G}
AN Enabled {1/10/25G}	P1: AN {1/10/25G}	P1: {1/10G}	P1: AN {1/25G}	P1: AN {1/10/25G}	P1: AN {1/10G}	P1: AN {1/25G}	P1: AN {1/10/25G}	P1: AN {1/25G}	P1: AN {1/10/25G}	P1: AN {1/10/25G}
	P2: no AN	P2: no AN	P2: no AN	P2: AN {1G}	P2: AN {10G}	P2: AN {25G}	P2: AN {1/10G}	P2: AN {1/25G}	P2: AN {10/25G}	P2: AN {1/10/25G}



Note: Use a supported optical transceiver or direct attached copper cable to achieve 1 Gbps link speed.

- P1 – port 1 setting
- P2 – port 2 setting
- AN – auto-negotiation
- No AN – forced speed

- {link speed} – expected link speed
- AN {link speeds} – advertised supported auto-negotiation link speeds.

The expected link speeds based on the local and link partner settings are shown in [Table 25](#).

Table 25: Expected Link Speeds

Local Speed Settings	Link Partner Speed Settings									
	Forced 1G	Forced 10G	Forced 25G	AN Enabled {1G}	AN Enabled {10G}	AN Enabled {25G}	AN Enabled {1/10G}	AN Enabled {1/25G}	AN Enabled {10/25G}	AN Enabled {1/10/25G}
Forced 1G	1G	No link	No link	No link	No link	No link	No link	No link	No link	No link
Forced 10G	No link	10G	No link	No link	No link	No link	No link	No link	No link	No link
Forced 25G	No link	No link	25G	No link	No link	No link	No link	No link	No link	No link
AN {1G}	No link	No link	No link	1G	No link	No link	1G	1G	No link	1G
AN {10G}	No link	No link	No link	No link	10G	No link	10G	No link	10G	10G
AN {25G}	No link	No link	No link	No link	No link	25G	No link	25G	25G	25G
AN {1/10G}	No link	No link	No link	1G	10G	No link	10G	1G	10G	10G
AN {1/25G}	No link	No link	No link	1G	No link	25G	1G	25G	25G	25G
AN {10/25G}	No link	No link	No link	No link	10G	25G	10G	25G	25G	25G
AN {1/10/25G}	No link	No link	No link	1G	10G	25G	10G	25G	25G	25G



Note: 1 Gbps link speed for SFP+/SFP28 is currently not support in this release.

To enable link speed auto-negotiation, the following options can be enabled in system BIOS HII menu or in CCM:



Note: Media Auto Detect must be enabled when enabling link speed auto-negotiation when connected via an optical transceiver module/cable.

System BIOS->Device Settings->NetXtreme-E NIC->Device Level Configuration

Operational Link Speed

This option configures the link speed used by the preboot (MBA and UEFI) drivers (Linux, ESX), OS driver, and firmware. This setting is overridden by the driver setting in the OS present state. The Windows driver (bnxnd_.sys) uses link speed settings in the driver .inf file.

Firmware Link Speed

This option configures the link speed used by the firmware when the device is in D3.

Auto-negotiation Protocol

This is the supported auto-negotiation protocol used to negotiate the link speed with the link partner. This option must match the AN protocol setting in the link partner port. The Broadcom NetXtreme-E NIC supports the following auto-negotiation protocols: IEEE 802.3by, 25G/50G consortiums and 25G/50G BAM. By default, this option is set to IEEE 802.3by and falls back to 25G/50G consortiums.

Link speed and Flow Control/Pause must be configured in the driver in the host OS.

Windows Driver Settings

To access the Windows driver settings:

Open **Windows Device Manager** -> **Broadcom NetXtreme E Series adapter** -> **Advanced Properties** -> **Advanced tab**

Flow Control = Auto-Negotiation

This enables Flow Control/Pause frame AN.

Speed and Duplex = Auto-Negotiation

This enables link speed AN.

Linux Driver Settings



Note: For 10GBase-T NetXtreme-E network adapters, auto-negotiation must be enabled.



The 25G advertisement is a newer standard first defined in the 4.7 kernel's ethtool interface. To fully support these new advertisement speeds for auto-negotiation, a 4.7 (or newer) kernel and a newer ethtool utility (version 4.8) are required.

ethtool -s eth0 speed 25000 autoneg off

This command turns off auto-negotiation and forces the link speed to 25 Gbps.

ethtool -s eth0 autoneg on advertise 0x0

This command enables auto-negotiation and advertises that the device supports all speeds: 1G, 10G, 25G.

The following are supported advertised speeds.

- 0x020 – 1000baseT Full
- 0x1000 – 10000baseT Full
- 0x80000000 – 25000baseCR Full

ethtool -A eth0 autoneg on|off

Use this command to enable/disable Pause frame auto-negotiation.

ethtool -a eth0

Use this command to display the current flow control auto-negotiation setting.

ESXi Driver Settings



Note: For 10GBase-T NetXtreme-E network adapters, auto-negotiation must be enabled. Using forced speed on 10GBase-T adapter results in esxcli command failure.



Note: VMWare does not support 25G speeds in ESX6.0. In this case, use the second utility (BNXTNETCLI) to set 25G speed. For ESX6.0U2, the 25G speed is supported.

\$ esxcli network nic get -n <iface>

This command shows the current speed, duplex, driver version, firmware version and link status.

\$ esxcli network nic set -S 10000 -D full -n <iface>

This command sets the forced speed to 10 Gbps.

\$ esxcli network nic set -a -n <iface>

This enables link speed auto-negotiation on interface <iface>.

\$ esxcli network nic pauseParams list

Use this command to get Pause Parameters list.

\$ esxcli network nic pauseParams set --auto <1/0> --rx <1/0> --tx <1/0> -n <iface>

Use this command to set Pause Parameters.



Note: Flow control/pause auto-negotiation can be set only when the interface is configured in link speed auto-negotiation mode.

FEC Auto-Negotiation

To enable/disable Link FEC auto-negotiation, the following options can be enabled in the system BIOS HII menu or in CCM:

- **System BIOS->Device Settings->NetXtreme-E NIC->Device Level Configuration**

The FEC Auto-negotiation uses two parameters during negotiation exchange: FEC capable and FEC request.

If the NIC advertises it as FEC Auto-negotiation capable, then the FEC settings are driven by the Switch. It can then link up with a switch which enables FEC or a switch which disables FEC.

Example:

- switch – capable=1, request = 1 then, the link is a FEC link.
- switch – capable= N/A, request= 0 = then the FEC is disabled.

For the NetXtreme-E Ethernet controllers, only Base-R FEC (CL74) is supported. [Table 26](#) shows all supported configurations with a link partner.

Table 26: Supported FEC Configurations for the BCM5730X/BCM5740X

Local SFEC Setting	Link Partner FEC Setting			
	Force Speed No FEC	Force Speed Base-R FEC CL74	AN (None)	AN (None, Base-R)
Force Speed No FEC	Link w/ no FEC	No link	No link	No link
Force Speed Base-R FEC CL74	No link	Base-R FEC CL74	No link	No link
AN (None)	No link	No link	Link w/ no FEC	Base-R FEC CL74
AN (None, Base-R)	No link	No link	Base-R FEC CL74	Base-R FEC CL74



Note: For force speed, it must be the same speed setting on both sides.



Note: AN {None} means AN advertises the Base-R capable bit. Set the F0 bit on IEEE802.3by and the F2 bit on Consortium.



Note: AN {None, Base-R} means the AN advertises the Base-R capable bit and requested bit. Set the F0 and F1 bits on IEEE802.3by and the F2 and F4 bits on Consortium.

For the NetXtreme-E Ethernet controller, the FEC supports Base-R FEC (CL74) and RS-FEC (CL91/CL108). [Table 27](#) shows all supported configurations with a link partner.

Table 27: Supported FEC Configurations for the BCM5741X

Local FEC Setting	Link Partner FEC Setting					
	Force Speed No FEC	Force Speed Base-R FEC CL74	Force Speed RS-FEC CL91/CL108	AN (None)	An (None, Base-R)	AN (None, Base-R, RS)
Force No FEC	Link w/no FEC	No link	No link	No link	No link	No link
Force Speed Base-R FEC CL74	No link	Base-R FEC CL74	No link	No link	No link	No link
Force RS-FEC CL91/CL108	No link	No link	RS-FEC CL91/CL108	No link	No link	No link
AN (None)	No link	No link	No link	Link w/ no FEC	Base-R FEC CL74	RS-FEC CL91/CL108
AN (None, Base-R)	No link	No link	No link	Base-R FEC CL74	Base-R FEC CL74	RS-FEC CL91/CL108
AN (None, Base-R, RS)	No link	No link	No link	RS-FEC CL91/CL108	RS-FEC CL91/CL108	RS-FEC CL91/CL108



Note: For force speed, it must be the same speed setting on both sides.



Note: AN {None} means AN advertises the Base-R capable bit. Set the F0 bit on IEEE802.3by and the F2 bit on Consortium.



Note: AN {None, Base-R} means the AN advertises the Base-R capable bit and requested bit. Set the F0 and F1 bits on IEEE802.3by and the F2 and F4 bits on Consortium.

Link Training

Link training allows both endpoints, the Broadcom adapter and the other side, to adjust power settings and other tuning parameters to maximize the reliability and efficiency of the communication channel between the two devices. The goal is to eliminate the need for channel-specific tuning between different cable lengths and types. Link training is performed for CR/KR speeds and it precedes auto-negotiation. Link training is operational when auto-negotiation is enabled. The link policy automatically disables link training if link training does not result in a link up with the link partner. This link policy ensures compatibility with a link partner that does not support link training.

[Table 28](#) shows the relationship between media type and speed for the BCM5730X, BCM5740X, and BCM5741X Ethernet controllers.

Table 28: Link Training Relationship Between Media Type and Speed

Local and Media Cable Type	Link Partner Link Training Setting		
	Force Speed Link Training Disabled	Force Speed Link Training Enabled	AN (Auto-Link Training)
Force Speed DAC (SFP+/SFP28/QSFP28)	Link without Link Training	Link with Link Training	Link with Link Training
Force Speed optical (Transceiver/AOC)	Link without Link Training	Link without Link Training	N/A
AN DAC (SFP+SFP28/QSFP28)	No Link	No Link	Link with Link Training

Media Auto Detect

Since parallel detection is not supported in SerDes, the firmware has implemented a method to enhance link detection called **Media Auto Detect**. This feature is controlled by CCM/HII as shown in [Figure 13](#).

When **Media Auto Detect** is enabled, the link policy follows the state machine (see [Figure 13](#)) to establish a link with a link partner. This behavior is dependent on the media type. For DAC cables, the method falls back to different force modes and stops when there is a link-up.

NOTE: Parallel detection is not supported when auto-negotiating at 25G or 10GbE. This means that if one side is auto-negotiating and the other side is not, the link will not come up.

Figure 13: State Machine for the Media Auto Detect Feature

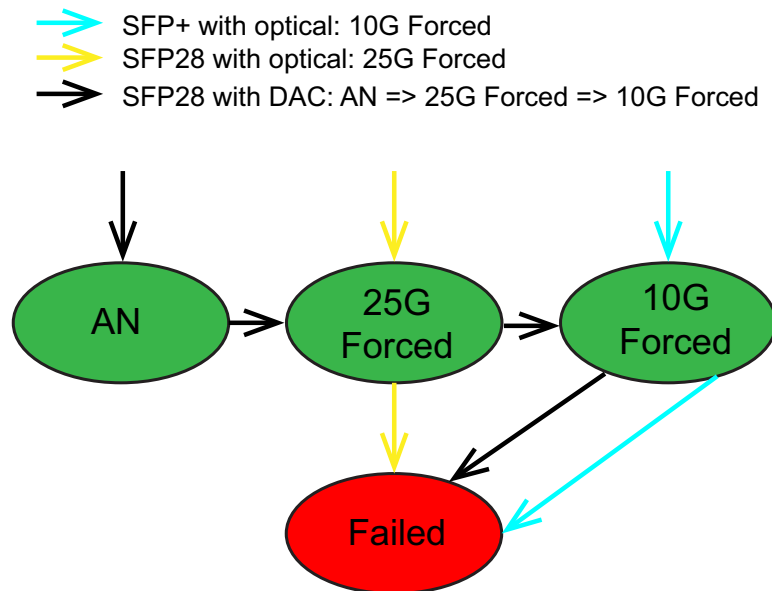


Table 29 and Table 30 show the link results with the **Media Auto Detect** feature enabled.

Table 29: Media Auto Detect for the BCM5730X and BCM5740X

		Link Partner Link Training Setting	
Link Partner Settings		Media Auto Detect	
Speed	FEC	No FEC	Base-R FEC
10G	No FEC	Link	Link with Base-R
	Base-R	No Link	Link with Base-R
25G	No FEC	Link	Link
	Base-R	No Link	Link with Base-R
AN	Disabled	Link	Link
	Auto FEC	Link	Link with Base-R
	Base-R	Link with Base-R	Link with Base-R

Table 30: Media Auto Detect for the BCM5741X

		Link Partner Link Training Setting		
Link Partner Settings		Media Auto Detect		
Speed	FEC	No FEC	Base-R FEC	RS-FEC
10G	No FEC	Link	Link with Base-R	No Link
	Base-R	No Link	Link with Base-R	No Link
	RS	No Link	No Link	Link with RS-FEC
25G	No FEC	Link	No Link	No Link
	Base-R	No Link	Link with Base-R	No Link
	RS	No Link	No Link	Link with RS-FEC
AN	Disabled	Link	Link	Link
	Auto FEC	Link	Link with Base-R	Link with RS-FEC
	Base-R	Link with Base-R	Link with Base-R	Link with RS-FEC
	RS	Link with RS	Link with RS-FEC	Link with RS-FEC

iSCSI Boot

Broadcom NetXtreme-E Ethernet adapters support iSCSI boot to enable the network boot of operating systems to diskless systems. iSCSI boot allows a Windows, Linux, or VMware operating system to boot from an iSCSI target machine located remotely over a standard IP network.

Supported Operating Systems for iSCSI Boot

The Broadcom NetXtreme-E Gigabit Ethernet adapters support iSCSI boot on the following operating systems:

- Windows Server 2012 and later 64-bit
- Linux RHEL 7.1 and later, SLES11 SP4 or later
- VMware 6.0 U2

Setting up iSCSI Boot

Refer to the following sections for information on setting up iSCSI boot.

Configuring the iSCSI Target



Note: Windows 2016 (or older) may have an outdated inbox driver which only supports forced Linkspeed. To avoid installation issues such as the inability to link up to a mismatching link partner, it is highly recommended that the user customize the Windows Installation image/media with the latest available NIC driver which supports Linkspeed Auto-negotiation.

Configuring the iSCSI target varies per the target vendor. For information on configuring the iSCSI target, refer to the documentation provided by the vendor. The general steps include:

1. Create an iSCSI target.
2. Create a virtual disk.
3. Map the virtual disk to the iSCSI target created in [Step 1 on page 49](#).
4. Associate an iSCSI initiator with the iSCSI target.
5. Record the iSCSI target name, TCP port number, iSCSI Logical Unit Number (LUN), initiator Internet Qualified Name (IQN), and CHAP authentication details.
6. After configuring the iSCSI target, obtain the following:
 - Target IQN
 - Target IP address
 - Target TCP port number
 - Target LUN
 - Initiator IQN
 - CHAP ID and secret

Configuring iSCSI Boot Parameters

Configure the Broadcom iSCSI boot software for either static or dynamic configuration. Refer to [Table 31](#) for configuration options available from the General Parameters menu. [Table 31](#) lists parameters for both IPv4 and IPv6. Parameters specific to either IPv4 or IPv6 are noted.

Table 31: Configuration Options

Option	Description
TCP/IP parameters via DHCP	This option is specific to IPv4. Controls whether the iSCSI boot host software acquires the IP address information using DHCP (Enabled) or use a static IP configuration (Disabled).
IP Autoconfiguration	This option is specific to IPv6. Controls whether the iSCSI boot host software will configure a stateless link-local address and/or stateful address if DHCPv6 is present and used (Enabled). Router Solicit packets are sent out up to three times with 4 second intervals in between each retry. Or use a static IP configuration (Disabled).
iSCSI parameters via DHCP	Controls whether the iSCSI boot host software acquires its iSCSI target parameters using DHCP (Enabled) or through a static configuration (Disabled). The static information is entered through the iSCSI Initiator Parameters Configuration screen.
CHAP Authentication	Controls whether the iSCSI boot host software uses CHAP authentication when connecting to the iSCSI target. If CHAP Authentication is enabled, the CHAP ID and CHAP Secret are entered through the iSCSI Initiator Parameters Configuration screen.
DHCP Vendor ID	Controls how the iSCSI boot host software interprets the Vendor Class ID field used during DHCP. If the Vendor Class ID field in the DHCP Offer packet matches the value in the field, the iSCSI boot host software looks into the DHCP Option 43 fields for the required iSCSI boot extensions. If DHCP is disabled, this value does not need to be set.
Link Up Delay Time	Controls how long the iSCSI boot host software waits, in seconds, after an Ethernet link is established before sending any data over the network. The valid values are 0 to 255. As an example, a user may need to set a value for this option if a network protocol, such as Spanning Tree, is enabled on the switch interface to the client system.
Use TCP Timestamp	Controls if the TCP Timestamp option is enabled or disabled.
Target as First HDD	Allows specifying that the iSCSI target drive will appear as the first hard drive in the system.
LUN Busy Retry Count	Controls the number of connection retries the iSCSI Boot initiator will attempt if the iSCSI target LUN is busy.
IP Version	This option specific to IPv6. Toggles between the IPv4 or IPv6 protocol. All IP settings will be lost when switching from one protocol version to another.

MBA Boot Protocol Configuration

To configure the boot protocol:

1. Restart the system.
2. From the PXE banner, select **CTRL+S**. The **MBA Configuration Menu** displays.

3. From the **MBA Configuration Menu**, use the **UP ARROW** or **DOWN ARROW** to move to the Boot Protocol option. Use the **LEFT ARROW** or **RIGHT ARROW** to change the Boot Protocol option for iSCSI.
4. Select **iSCSI Boot Configuration** from **Main Menu**.

iSCSI Boot Configuration

There are two ways to configure iSCSI boot:

- Static iSCSI Boot Configuration.
- Dynamic iSCSI Boot Configuration.

Static iSCSI Boot Configuration

In a static configuration, you must enter data for the system's IP address, the system's initiator IQN, and the target parameters obtained in ["Configuring the iSCSI Target" on page 49](#). For information on configuration options, see [Table 31 on page 50](#).

To configure the iSCSI boot parameters using static configuration:

1. From the **General Parameters** menu, set the following:
 - TCP/IP parameters via DHCP – Disabled. (For IPv4.)
 - IP Autoconfiguration – Disabled. (For IPv6, non-offload.)
 - iSCSI parameters via DHCP – Disabled
 - CHAP Authentication – Disabled
 - DHCP Vendor ID – BRCM ISAN
 - Link Up Delay Time – 0
 - Use TCP Timestamp – Enabled (for some targets such as the Dell/EMC AX100i, it is necessary to enable Use TCP Timestamp)
 - Target as First HDD – Disabled
 - LUN Busy Retry Count – 0
 - IP Version – IPv6. (For IPv6, non-offload.)
2. Select **ESC** to return to the **Main** menu.
3. From the **Main** menu, select **Initiator Parameters**.
4. From the **Initiator Parameters** screen, enter values for the following:
 - IP Address (unspecified IPv4 and IPv6 addresses should be "0.0.0.0" and ":::", respectively)
 - Subnet Mask Prefix
 - Default Gateway
 - Primary DNS
 - Secondary DNS
 - iSCSI Name (corresponds to the iSCSI initiator name to be used by the client system)



Note: Enter the IP address. There is no error-checking performed against the IP address to check for duplicates or incorrect segment/network assignment.

5. Select **ESC** to return to the **Main** menu.
6. From the **Main** menu, select **1st Target Parameters**.



Note: For the initial setup, configuring a second target is not supported.

7. From the **1st Target Parameters** screen, enable **Connect** to connect to the iSCSI target. Type values for the following using the values used when configuring the iSCSI target:
 - IP Address
 - TCP Port
 - Boot LUN
 - iSCSI Name
8. Select **ESC** to return to the **Main** menu.
9. Select **ESC** and select **Exit** and **Save Configuration**.
10. Select **F4** to save the MBA configuration.

Dynamic iSCSI Boot Configuration

In a dynamic configuration, specify that the system's IP address and target/initiator information are provided by a DHCP server (see IPv4 and IPv6 configurations in [“Configuring the DHCP Server to Support iSCSI Boot” on page 54](#). For IPv4, with the exception of the initiator iSCSI name, any settings on the Initiator Parameters, 1st Target Parameters, or 2nd Target Parameters screens are ignored and do not need to be cleared. For IPv6, with the exception of the CHAP ID and Secret, any settings on the Initiator Parameters, 1st Target Parameters, or 2nd Target Parameters screens are ignored and do not need to be cleared. For information on configuration options, see [Table 31 on page 50](#).



Note: When using a DHCP server, the DNS server entries are overwritten by the values provided by the DHCP server. This occurs even if the locally provided values are valid and the DHCP server provides no DNS server information. When the DHCP server provides no DNS server information, both the primary and secondary DNS server values are set to 0.0.0.0. When the Windows OS takes over, the Microsoft iSCSI initiator retrieves the iSCSI Initiator parameters and configures the appropriate registries statically. It will overwrite whatever is configured. Since the DHCP daemon runs in the Windows environment as a user process, all TCP/IP parameters have to be statically configured before the stack comes up in the iSCSI Boot environment.

- If DHCP Option 17 is used, the target information is provided by the DHCP server, and the initiator iSCSI name is retrieved from the value programmed from the Initiator Parameters screen. If no value was selected, then the controller defaults to the name:

`iqn.1995-05.com.broadcom.<11.22.33.44.55.66>.iscsiboot`
where the string **11.22.33.44.55.66** corresponds to the controller's MAC address.
- If DHCP option 43 (IPv4 only) is used, then any settings on the Initiator Parameters, 1st Target Parameters, or 2nd Target Parameters screens are ignored and do not need to be cleared.

To configure the iSCSI boot parameters using a dynamic configuration:

1. From the **General Parameters** menu screen, set the following parameters:
 - TCP/IP parameters via DHCP – Enabled. (For IPv4.)
 - IP Autoconfiguration – Enabled. (For IPv6, non-offload.)
 - iSCSI parameters via DHCP – Enabled
 - CHAP Authentication – Disabled
 - DHCP Vendor ID – BCM ISAN
 - Link Up Delay Time – 0
 - Use TCP Timestamp – Enabled (for some targets such as the Dell/EMC AX100i, it is necessary to enable Use TCP Timestamp)
 - Target as First HDD – Disabled
 - LUN Busy Retry Count – 0
 - IP Version – IPv6. (For IPv6, non-offload.)
2. Select **ESC** to return to the **Main** menu.



Note: Information on the Initiator Parameters, and 1st Target Parameters screens are ignored and do not need to be cleared.

3. Select **Exit** and **Save Configurations**.

Enabling CHAP Authentication

Ensure that CHAP authentication is enabled on the target.

To enable CHAP authentication:

1. From the **General Parameters** screen, set **CHAP Authentication** to **Enabled**.
2. From the **Initiator Parameters** screen, enter the parameters for the following:
 - CHAP ID (up to 128 bytes)
 - CHAP Secret (if authentication is required, and must be 12 characters in length or longer)
3. Select **ESC** to return to the **Main** menu.
4. From the **Main** menu, select the **1st Target Parameters**.
5. From the **1st Target Parameters** screen, type values for the following using the values used when configuring the iSCSI target:
 - CHAP ID (optional if two-way CHAP)
 - CHAP Secret (optional if two-way CHAP, and must be 12 characters in length or longer)
6. Select **ESC** to return to the **Main** menu.
7. Select **ESC** and select **Exit** and **Save Configuration**.

Configuring the DHCP Server to Support iSCSI Boot

The DHCP server is an optional component and it is only necessary for dynamic iSCSI Boot configuration setup (see [“Dynamic iSCSI Boot Configuration” on page 52](#)).

Configuring the DHCP server to support iSCSI boot is different for IPv4 and IPv6. Refer to the following sections:

DHCP iSCSI Boot Configurations for IPv4

The DHCP protocol includes a number of options that provide configuration information to the DHCP client. For iSCSI boot, Broadcom adapters support the following DHCP configurations:

DHCP Option 17, Root Path

Option 17 is used to pass the iSCSI target information to the iSCSI client. The format of the root path as defined in IETF RFC 4173 is:

```
iscsi:"<servername>":"<protocol>":"<port>":"<LUN>":"<targetname>
```

The parameters are defined in [Table 32](#).

Table 32: DHCP Option 17 Parameter Definition

Parameter	Definition
"iscsi:"	A literal string.
<servername>	The IP address or FQDN of the iSCSI target
":"	Separator.
<protocol>	The IP protocol used to access the iSCSI target. Currently, only TCP is supported so the protocol is 6.
<port>	The port number associated with the protocol. The standard port number for iSCSI is 3260.
<LUN>	The Logical Unit Number to use on the iSCSI target. The value of the LUN must be represented in hexadecimal format. A LUN with an ID OF 64 would have to be configured as 40 within the option 17 parameter on the DHCP server.
<targetname>	The target name in either IQN or EUI format (refer to RFC 3720 for details on both IQN and EUI formats). An example IQN name would be "iqn.1995-05.com.broadcom:iscsi-target".

DHCP Option 43, Vendor-Specific Information

DHCP option 43 (vendor-specific information) provides more configuration options to the iSCSI client than DHCP option 17. In this configuration, three additional suboptions are provided that assign the initiator IQN to the iSCSI boot client along with two iSCSI target IQNs that can be used for booting. The format for the iSCSI target IQN is the same as that of DHCP option 17, while the iSCSI initiator IQN is simply the initiator's IQN.



Note: DHCP Option 43 is supported on IPv4 only.

The suboptions are listed below.

Table 33: DHCP Option 43 Suboption Definition

Suboption	Definition
201	First iSCSI target information in the standard root path format iscsi:"<servername>":"<protocol>":"<port>":"<LUN>":"<targetname>
203	iSCSI initiator IQN

Using DHCP option 43 requires more configuration than DHCP option 17, but it provides a richer environment and provides more configuration options. Broadcom recommends that customers use DHCP option 43 when performing dynamic iSCSI boot configuration.

Configuring the DHCP Server

Configure the DHCP server to support option 17 or option 43.



Note: If Option 43 is used, configure Option 60. The value of Option 60 should match the DHCP Vendor ID value. The DHCP Vendor ID value is BCM ISAN, as shown in General Parameters of the iSCSI Boot Configuration menu.

DHCP iSCSI Boot Configuration for IPv6

The DHCPv6 server can provide a number of options, including stateless or stateful IP configuration, as well as information to the DHCPv6 client. For iSCSI boot, Broadcom adapters support the following DHCP configurations:



Note: The DHCPv6 standard Root Path option is not yet available. Broadcom suggests using Option 16 or Option 17 for dynamic iSCSI Boot IPv6 support.

DHCPv6 Option 16, Vendor Class Option

DHCPv6 Option 16 (vendor class option) must be present and must contain a string that matches the configured DHCP Vendor ID parameter. The DHCP Vendor ID value is BCM ISAN, as shown in General Parameters of the iSCSI Boot Configuration menu.

The content of Option 16 should be <2-byte length> <DHCP Vendor ID>.

DHCPv6 Option 17, Vendor-Specific Information

DHCPv6 Option 17 (vendor-specific information) provides more configuration options to the iSCSI client. In this configuration, three additional suboptions are provided that assign the initiator IQN to the iSCSI boot client along with two iSCSI target IQNs that can be used for booting.

The suboptions are listed in [Table 34 on page 56](#).

Table 34: DHCP Option 17 Suboption Definition

Suboption	Definition
201	First iSCSI target information in the standard root path format "iscsi:[<servername>]":"<protocol>":"<port>": "<LUN>":"<targetname>"
203	iSCSI initiator IQN



Note: In [Table 34](#), the brackets [] are required for the IPv6 addresses.

The content of option 17 should be <2-byte Option Number 201|202|203> <2-byte length> <data>.

Configuring the DHCP Server

Configure the DHCP server to support Option 16 and Option 17.



Note: The format of DHCPv6 Option 16 and Option 17 are fully defined in RFC 3315.

VXLAN: Configuration and Use Case Examples

VXLAN encapsulation permits many layer 3 hosts residing on one server to send and receive frames by encapsulating on to single IP address associated with the NIC card installed on the same server.

This example discusses basic VxLan connectivity between two RHEL servers. Each server has one physical NIC enabled with outer IP address set to 1.1.1.4 and 1.1.1.2.

A VXLAN10 interface with VXLAN ID 10 is created with multicast group 239.0.0.10 and is associated with physical network port pxp1 on each server.

An IP address for the host is created on each server and associated that to VXLAN interface. Once the VXLAN interface is brought up, the host present in system 1 can communicate with host present in system 2. The VLXAN format is shown in [Table 35](#).

Table 35: VXLAN Frame Format

MAC header	Outer IP header with proto = UDP	UDP header with Destination port= VXLAN	VXLAN header (Flags, VNI)	Original L2 Frame	FCS
------------	----------------------------------	---	---------------------------	-------------------	-----

[Table 36](#) provides VXLAN command and configuration examples.

Table 36: VXLAN Command and Configuration Examples

System 1	System 2
PxPy: ifconfig PxPy 1.1.1.4/24	PxPy: ifconfig PxPy 1.1.1.2/24
ip link add vxlan10 type vxlan id 10 group 239.0.0.10 dev PxPy dstport 4789	ip link add vxlan10 type vxlan id 10 group 239.0.0.10 dev PxPy dstport 4789
ip addr add 192.168.1.5/24 broadcast 192.168.1.255 dev vxlan10	ip addr add 192.168.1.10/24 broadcast 192.168.1.255 dev vxlan10
ip link set vxlan10 up	ip link set vxlan10 up
ip -d link show vxlan10	
Ping 192.168.1.10	ifconfig vxlan10 (MTU 1450) (SUSE and RHEL)

Note: x represents the PCIe bus number of the physical adapter found in the system. y represents the port number on the physical adapter.

SR-IOV: Configuration and Use Case Examples

SR-IOV can be configured, enabled, and used on 10Gb and 25Gb Broadcom NetExtreme-E NICs.

Linux Use Case

1. Enable SR-IOV in the NIC cards:
 - a. SR-IOV in the NIC card can be enabled using the **HII** menu. During system boot, access the system **BIOS -> Device Settings -> NetXtreme-E NIC -> Device Level Configuration**.
 - b. Set the Virtualization mode to SR-IOV.
 - c. Set the number of virtual functions per physical function.
 - d. Set the number of MSI-X vectors per the VF and Max number of physical function MSI-X vectors. If the VF is running out of resources, balance the number of MSI-X vectors per VM using CCM.
2. Enable virtualization in the BIOS:
 - a. During system boot, enter the system **BIOS -> Processor settings -> Virtualization Technologies** and set it to **Enabled**.
 - b. During system boot, enter the system **BIOS -> Integrated Devices -> SR-IOV Global** and set it to **Enabled**.
3. Install the desired Linux version with Virtualization enabled (libvirt and Qemu).
4. Enable the iommu kernel parameter.
 - a. The IOMMU kernel parameter is enabled by editing `/etc/default/grub.cfg` and running `grub2-mkconfig -o /boot/grub2/grub.cfg` for legacy mode. For UEFI mode, edit `/etc/default/grub.cfg` and run `grub2-mkconfig -o /etc/grub2-efi.cfg`. Refer to the following example:

```
Linuxefi /vmlinuz-3.10.0-229.el7.x86_64 root=/dev/mapper/rhel-root no rd.lvm.lv=rhel/swap crashkernel=auto rd.lvm.lv=rhel/root rhgb intel_iommu=on quiet LANG=en_US.UTF.8
```
5. Install `bnxt_en` driver:
 - a. Copy the `bnxt_en` driver on to the OS and run `make; make install; modprobe bnxt_en`.



Note: Use `netxtreme-bnxt_en<version>.tar.gz` to install both `bnxt_re` and `bnxt_en` for RDMA functionality on SRIOV VFs.

6. Enable Virtual Functions through Kernel parameters:

- a. Once the driver is installed, `lspci` will display the NetXtreme-E NICs present in the system. Bus, device, and Function are needed for activating Virtual functions.
- b. To activate Virtual functions, enter the command shown below:

```
echo X >/sys/bus/pci/device/0000\:Bus\:Dev.Function/sriov_numvfs
```



Note: Ensure that the PF interfaces are up. VFs are only created if PFs are up. X is the number of VFs that will be exported to the OS.

A typical example would be:

```
echo 4 > /sys/bus/pci/devices/0000\:04\:00.0/sriov_numvfs
```

7. Check the PCI-E virtual functions:

- a. The `lspci` command will display the virtual functions with DID set to 16D3 for BCM57402/BCM57404/BCM57406, 16DC for non-RDMA BCM57412/BCM57414/BCM57416, and 16C1 or RDMA enabled BCM57412/BCM57414/BCM57416.

8. Use the Virtual Manager to install a Virtualized Client system (VMs).

Refer to the Linux documentation for Virtual Manager installation. Ensure that the hypervisor's built in driver is removed. An example would be `NIC:d7:73:a7 rt18139`. Remove this driver.

9. Assign a virtual function to the guest VMs.

- a. Assign this adapter to a guest VM as a physical PCI Device. Refer to the Linux documentation for information on assigning virtual functions to a VM guest.

10. Install `bnxt_en` drivers on VMs:

- a. On the guest VMs, copy the `netxtreme-bnxt_en-<version>.tar.gz` source file and extract the `tar.gz` file. Change directory to each driver and run `make`; `make install`; `modprobe bnxt_en` (and `bnxt_re` if enabling RDMA). Make sure that the driver loads properly by checking the interface using `modinfo` command. The user may need to run `modprobe -r bnxt_en` to unload existing or inbox `bnxt_en` module prior to loading the latest built module.

11. Test the guest VM connectivity to external world:

- a. Assign proper IP address to the adapter and test the network connectivity.

Windows Case

1. Enable SR-IOV in the NIC cards:

- a. SR-IOV in the NIC card can be enabled using the HII menu. During the system boot, access the system **BIOS -> Device Settings -> NetXtreme-E NIC -> Device Level Configuration**.
- b. Set the Virtualization mode to SR-IOV.

- c. Set the number of virtual functions per physical function.
 - d. Set the number of MSI-X vectors per the VF and Max number of physical function MSI-X vectors. If the VF is running out of resources, balance the number of MSI-X vectors per VM using CCM.
2. Enable virtualization in the BIOS:
 - a. During system boot, enter the system **BIOS -> Processor settings -> Virtualization Technologies** and set it to **Enabled**.
 - b. During system boot, enter the system **BIOS -> Integrated Devices -> SR-IOV Global** and set it to **Enabled**.
 3. Install the latest KB update for your Windows 2012 R2 or Windows 2016 OS.
 4. Install the appropriate Virtualization (Hyper-V) options. For more detail requirements and steps on setting up Hyper-V, Virtual Switch, and Virtual Machine, please visit Microsoft.com:
<https://technet.microsoft.com/en-us/windows-server-docs/compute/hyper-v/system-requirements-for-hyper-v-on-windows>
<https://technet.microsoft.com/en-us/windows-server-docs/compute/hyper-v/get-started/install-the-hyper-v-role-on-windows-server>
 5. Install the latest NetXtreme-E driver on the Hyper-V.
 6. Enable SR-IOV in the NDIS miniport driver advanced properties.
 7. In Hyper-V Manager, create your Virtual Switch with the selected NetXtreme-E interface.
 8. Check the **Enable Single-Root I/O Virtualization (SR-IOV)** box while creating the Hyper-V Virtual Adapter.
 9. Create a Virtual Machine (VM) and add the desired number of Virtual Adapters
 10. Under the Virtual Machine's Network Adapter settings for each Virtual Adapter, check **Enable SR-IOV** under the **Hardware Acceleration** section.
 11. Launch your VM and install the desired guest OS.
 12. Install the corresponding NetXtreme-E driver for each guest OS.



Note: The Virtual Function (VF) driver for NetXtreme-E is same driver as the base driver. For example, if the guest OS is Windows 2012 R2, the user needs to install Bnxtnd64.sys in VM. The user can do this by running the NetXtreme-E Driver Installer executable. Once the driver has been installed in the guest OS, the user can see the VF driver interface(s) appearing in Device Manager of the guest OS in the VM.

VMWare SRIOV Case

1. Enable SR-IOV in the NIC cards:
 - a. SR-IOV in the NIC card can be enabled using the HII menu. During the system boot, access the system **BIOS -> Device Settings -> NetXtreme-E NIC -> Device Level Configuration**.
 - b. Set the Virtualization mode to SR-IOV.
 - c. Set the number of virtual functions per physical function.
 - d. Set the number of MSI-X vectors per the VF and Max number of physical function MSI-X vectors. If the VF is running out of resources, balance the number of MSI-X vectors per VM using CCM.

2. Enable virtualization in the BIOS:
 - a. During system boot, enter the system **BIOS -> Processor settings -> Virtualization Technologies** and set it to **Enabled**.
 - b. During system boot, enter the system **BIOS -> Integrated Devices -> SR-IOV Global** and set it to **Enabled**.

3. On ESXi, install the Bnxtnet driver using the following steps:

- a. Copy the <bnxtnet>-<driver version>.vib file in /var/log/vmware.

```
$ cd /var/log/vmware.
$ esxcli software vib install --no-sig-check -v <bnxtnet>-<driver version>.vib.
```

- b. Reboot the machine.
- c. Verify that whether drivers are correctly installed:

```
$ esxcli software vib list | grep bnxtnet
```

4. Install the Broadcom provided BNXTNETCLI (esxcli bnxtnet) utility to set/view the miscellaneous driver parameters that are not natively supported in esxcli, such as: link speed to 25G, show driver/firmware chip information, show NIC configuration (NPAR, SRIOV). For more information, please see the bnxtnet driver README.txt.

To install this utility:

- a. Copy BCM-ESX-bnxtnetcli-<version>.vib in /var/log/vmware.

```
$ cd /var/log/vmware
$ esxcli software vib install --no-sig-check -v /BCM-ESX-bnxtnetcli-<version>.vib
```

- b. Reboot the system.
- c. Verify whether vib is installed correctly:

```
$ esxcli software vib list | grep bcm-esx-bnxtnetcli
```

- d. Set speed 10/20/25

```
$ esxcli bnxtnet link set -S <speed> -D <full> -n <iface>
```

This will return OK message if the speed is correctly set.

Example:

```
$ esxcli bnxtnet link set -S 25000 -D full -n vmnic5
```

- e. Show the link stats

```
$ esxcli bnxtnet link get -n vmnic6
```

- f. Show the driver/firmware/chip information

```
$ esxcli bnxtnet drvinfo get -n vmnic4
```

- g. Show the NIC information (e.g. BDF; NPAR, SRIOV configuration)

```
$ esxcli bnxtnet nic get -n vmnic4
```

5. Enabling SRIOV VFs:

Only the PFs are automatically enabled. If a PF supports SR-IOV, the PF(vmknix) is part of the output of the command shown below.

```
esxcli network sriovnic list
```

To enable one or more VFs, the driver uses the module parameter "max_vfs" to enable the desired number of VFs for PFs. For example, to enable four VFs on PF1:

```
esxcfg-module -s 'max_vfs=4' bnxtnet (reboot required)
```

To enable VFs on a set of PFs, use the command format shown below. For example, to enable four VFs on PF 0 and 2 VFs on PF 2:

```
esxcfg-module -s 'max_vfs=4,2' bnxtnet (reboot required)
```

The required VFs of each supported PF are enabled in order during the PF bring up. Refer to the VMware documentation for information on how to map a VF to a VM:

<https://pubs.vmware.com/vsphere-60/index.jsp#com.vmware.vsphere.networking.doc/GUID-EE03DC6F-32CA-42EF-98FC-12FDE06C0BE0.html>

<https://pubs.vmware.com/vsphere-60/index.jsp#com.vmware.vsphere.networking.doc/GUID-CC021803-30EA-444D-BCBE-618E0D836B9F.html>

NPAR – Configuration and Use Case Example

Features and Requirements

- OS/BIOS Agnostic – The partitions are presented to the operating system as "real" network interfaces so no special BIOS or OS support is required like SR-IOV.
- Additional NICs without requiring additional switch ports, cabling, PCIe expansion slots.
- Traffic Shaping – The allocation of bandwidth per partition can be controlled so as to limit or reserve as needed.
- Can be used in a Switch Independent manner – The switch does not need any special configuration or knowledge of the NPAR enablement.
- Can be used in conjunction with RoCE and SR-IOV.
- Supports stateless offloads such as, LSO, TPA, RSS/TSS, and RoCE (two PFs per port only).
- Alternative Routing-ID support for greater than eight functions per physical device.



Note: In the **Dell UEFI HII Menu-> Main Configuration -> Device Level Configuration** page, the user can enable **NParEP** to allow the NXE adapter to support up to 16 PFs per device on an ARI capable system. For a 2 port device, this means up to 8 PFs for each port.

Limitations

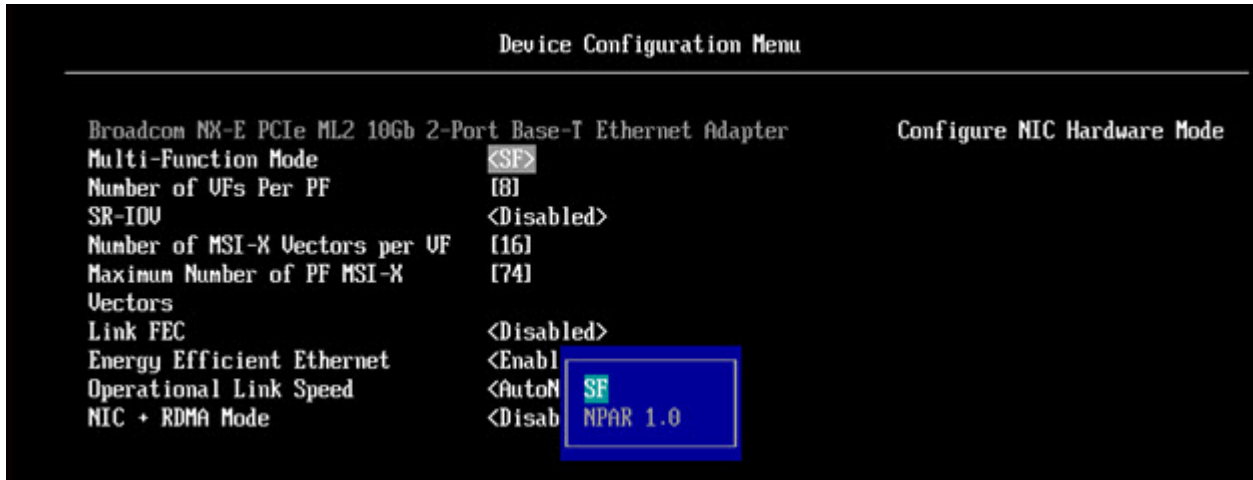
- Shared settings must be suppressed to avoid contention. For example: Speed, Duplex, Flow Control, and similar physical settings are hidden by the device driver to avoid contention.
- Non-ARI systems enable only eight partitions per physical device.
- RoCE is only supported on the first two partitions of each physical port, or a total of four partitions per physical device. In NPAR + SRIOV mode, only two VFs from each parent physical port can enable RDMA support, or total of four VFs + RDMA per physical device.

Configuration

NPAR can be configured using BIOS configuration HII menus or by using the Broadcom CCM utility on legacy boot systems. Some vendors also expose the configuration via additional proprietary interfaces.

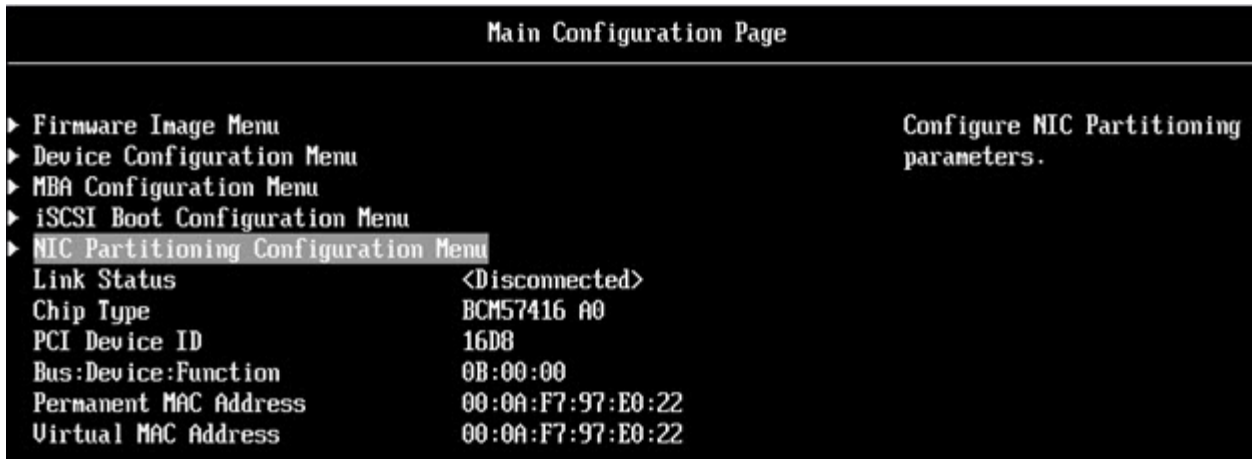
To enable NPAR:

1. Select the target NIC from the BIOS HII Menu or CCM interface and set the Multi-Function Mode or Virtualization Mode option.



NPAR is enabled in combination with SR-IOV. For some ARI capable OEM systems, the **NParEP** button is available to explicitly allow the BCM5741X to support up to 16 partitions. Switching from Single Function mode to Multifunction mode, the device needs to be re-enumerated, therefore changes will not take effect until a system reboot occurs.

- Once NPAR is enabled, the NIC Partitioning Main Configuration Menu option is available from the main NIC Configuration Menu associated with each physical port.



- The NIC Partition Configuration Menu (shown below) allows the user to choose the number of partitions that should be allocated from the selected physical port. Each BCM5741X NIC can support a maximum of 16 partitions on an ARI capable server. By default, dual-port adapters are configured for eight partitions per physical port. Configuration options for each partition are also accessible from this menu. For some OEM systems, the HII menu also includes a Global Bandwidth Allocation page where the minimum (reserved) and maximum (limit) TX Bandwidth for all partitions can be configured simultaneously.

```

NIC Partitioning Configuration Menu
-----
Broadcom BCM57416 NetXtreme-E 10GBASE-T RDMA Ethernet Controller   Configure Number of
Number of Partitions           [8]                               Partitions Per Port
▶ Partition 1 Configuration
▶ Partition 2 Configuration
▶ Partition 3 Configuration
▶ Partition 4 Configuration
▶ Partition 5 Configuration
▶ Partition 6 Configuration
▶ Partition 7 Configuration
▶ Partition 8 Configuration
    
```

4. Set the NIC Partition Configuration parameters (see [Table 37 on page 64](#)).

```

Partition 1 Configuration
-----
Broadcom NX-E PCIe ML2 10Gb 2-Port Base-T Ethernet Adapter       Configure RDMA Support
BW Reservation              [0]
BW Limit                    [100]
BW Reservation Valid       <False>
BW Limit Valid             <False>
Support RDMA               <Disabled>
MAC Address                 00:10:18:99:98:E0
    
```

Table 37: NPAR Parameters

Parameter	Description	Valid Options
BW Reservation	Percentage of total available bandwidth that should be reserved for this partition. 0 indicates equal division of bandwidth between all partitions.	Value 0-100
BW Limit	Maximum percentage of available bandwidth this partition is allowed.	Value 0-100
BW Reservation Valid	Functions as an on/off switch for the BW Reservation setting.	True/False
BW Limit Valid	Functions as an on/off switch for the BW Limit setting.	True/False
Support RDMA	Functions as an on/off switch for RDMA support on this partition. Please note, only two partitions per physical port can support RDMA. For a dual port device, up to 4 NPAR partitions can support RDMA.	Enabled/Disabled
MAC Address	MAC Address for this partition.	–

Notes on Reducing NIC Memory Consumption

Because of the faster link speeds supported in this NIC, the default number of receive buffers is larger. More packets can arrive within a given time interval when the link speed is higher, and if the host system is delayed in processing the receive interrupts, the NIC must drop packets if all available receive buffers are in use.

The default value of receive buffers was selected to work well for typical configurations. But if you have many NIC's in a system, have enabled NPAR on multiple NIC's, or if you have only a small amount of RAM, you may see a Code 12 "yellow bang" in the Device Manager for some of the NIC's. Code 12 means that the driver failed to load because there were not enough resources. In this case, the resource is a specific type of kernel memory called Non-Paged Pool (NPP) memory.

If you are getting a Code 12, or for other reasons wish to reduce the amount of NPP memory consumed by the NIC:

- Reduce the number of RSS queues from the default of 8 to 4 or 2. Each RSS queue has its own set of receive buffers allocated, so reducing the number of RSS queues reduces the allocated NPP memory. There can be performance implications from reducing the number of RSS queues, as fewer cores participate in processing receive packets from that NIC. Per processor CPU utilization should be monitored to ensure that there are no "hot" processors after this change.
- Reduce memory allocation by reducing the number of receive buffers allocated. The default value of 0 means the driver should automatically determine the number of receive buffers. For typical configurations, a setting of 0 (=auto) will map to XXXX receive buffers per queue. You can choose a smaller value such as 1500, 1000, or 500. (The value needs to be in multiples of 500 between the range of 500 and 15000.) As mentioned above, a smaller number of receive buffers increases the risk of packet drop and a corresponding impact to packet retransmissions and decreased throughput.

The parameters "Maximum Number of RSS Queues" and "Receive Buffers (0=Auto)" can be modified using the **Advanced** properties tab for each NIC in the **Device Manager**. If you want to modify multiple NIC's at the same time, it is faster to use the *Set-NetAdapterAdvancedProperty PowerShell cmdlet*. For example, to assign two RSS queues for all NIC's in a system whose NIC name starts with "SI", run the following command:

```
Set-NetAdapterAdvancedProperty SI* -RegistryKeyword *NumRSSQueues -RegistryValue 2
```

Similarly, to set the number of Receive buffers to 1500, run the following command:

```
Set-NetAdapterAdvancedProperty SI* -RegistryKeyword *ReceiveBuffers -RegistryValue 1500
```

See <https://blogs.technet.microsoft.com/wincat/2012/08/27/using-powershell-for-nic-configuration-tasks/> for an overview of using PowerShell to modify NIC properties.

RoCE – Configuration and Use Case Examples

This section provides configuration and use case examples for RoCE.

To enable RoCE for PFs or VFs, the user must enable the RDMA selection in the HII menu in the BIOS before the RDMA option takes effect in the host or guest OS.

To enable RDMA in single function mode (if **Virtualization Mode** is **None** or **SR-IOV**):

1. During the system boot, access the **System Setup -> Device Settings -> NetXtreme-E NIC -> Main Configuration Page** and set **NIC+ RMDA Mode** to **Enabled**.

To enable RDMA if Virtualization Mode is NPAR or NPAR+SR-IOV:

1. During the system boot, access the **System Setup -> Device Settings -> NetXtreme-E NIC ->NIC Partitioning Configuration-> Partition 1 (or 2) Configuration** and set **NIC+ RMDA Mode** to **Enabled**.



Note: If using NPAR+SRIOV mode, only two VFs from each parent physical port can enable RDMA support, or a total of four VFs + RDMA per physical device.

Linux Configuration

Requirements

To configure RoCE in Linux, the following items are required:

- `bnxt_en-roce` (RoCE supported `bnxt_en` driver which is part of released gzip compressed tar archive)
- `bnxt_re` (RoCE driver)
- `libbnxtre` (User mode RoCE library module)

BNXT_RE Driver Dependencies

The `Bnxt_re` driver requires a special RoCE enabled version of `bnxt_en` which is included in the `netxtreme-bnxt_en-1.7.9.tar.gz` (or newer) package. The `bnxt_re` driver compilation depends on whether IB stack is available along with the OS distribution or an external OFED is required.



Note: It is necessary to load the correct `bnxt_en` version that is included in the same `netxtreme-bnxt_en-1.7.x.tar.gz` package. `Bnxt_re` and `Bnxt_en` function as a pair to enable RoCE traffic. Using mismatching versions of these two drivers produces unreliable or unpredictable results.

- Distros that have an IB Stack available along with OS distribution:

RH7.1/7.2/7.3/6.7/6.8, SLES12SP2 and Ubuntu 16.04

If not already installed, the IB stack and useful utils can be installed in Redhat with the following commands prior to compiling `bnxt_re`:

```
yum -y install libibverbs* infiniband-diag perftest qperf librdmacm utils
```

To compile `bnxt_re`:

```
$make
```

- Distros that need external OFED to be installed:

SLES11SP4

Please refer OFED release notes from the following link and install OFED before compiling `bnxt_re` driver.

http://downloads.openfabrics.org/downloads/OFED/release_notes/OFED_3.18-2_release_notes

To compile `bnxt_re`:

```
$export OFED_VERSION=OFED-3.18-2
$make
```

Installation

To install RoCE in Linux:

1. Upgrade the NIC NVRAM using the RoCE supported firmware packages from Software Release 20.06.04.01 or newer.
2. In the OS, uncompress, build, and install the BCM5741X Linux L2 and RoCE drivers.
 - a. # `tar -xzf netxtreme-bnxt_en-1.7.9.tar.gz`
 - b. # `cd netxtreme-bnxt_en-bnxt_re`
 - c. # `make build && make install`
3. Uncompress, build, and install the NetXtreme-E Linux RoCE User Library.
 - a. # `tar xzf libbnxtre-0.0.18.tar.gz`
 - b. # `cd libbnxtre-0.0.18`
 - c. # `configure && make && make install.`
 - d. # `cp bnxtre.driver /etc/libibverbs.d/`
 - e. # `echo "/usr/local/lib" >> /etc/ld.so.conf`
 - f. # `ldconfig -v`

Please refer to the `bnxt_re` README.txt for more details on configurable options and recommendations.

Limitations

In dual port NICs, if both ports are on same subnet, `rdma perftest` commands may fail. The possible cause is due to an arp flux issue in the Linux OS. To workaround this limitation, use multiple subnets for testing or bring the second port/interface down.

Known Issues

`Bnxt_en` and `Bnxt_re` are designed to function in pair. Older `Bnxt_en` drivers prior to version 1.7.x do not support RDMA and cannot be loaded at the same time as the `Bnxt_re` (RDMA) driver. The user may experience a system crash and reboot if `Bnxt_re` is loaded with older `Bnxt_en` drivers. It is recommend that the user load the `Bnxt_en` and `Bnxt_re` module from the same `netxtreme-bnxt_en-<1.7.x>.tar.gz` bundle.

To prevent mismatching a combination of `bnxt_en` and `bnxt_re` from being loaded, the following is required:

- If RedHat/CentOS 7.2 OS was installed to the target system using PXEboot with `bnxt_en` DUD or a kernel module RPM, delete the file `bnxt_en.ko` found in `/lib/modules/$(uname -r)/extra/bnxt_en/bnxt_en.ko` or `edit /etc/depmod.d/`.
- `bnxt_en.conf` to override to use updated version. Users can also erase the current BCM5741X Linux kernel driver using the `rpm -e kmod-bnxt_en` command. RHEL 7.3/SLES 12 Sp2 has `bnxt_en` inbox driver (older than v1.7.x). This driver must be removed and the latest `bnxt_en` be added before applying the `bnxt_re` (RoCE drivers).

Windows

Kernel Mode

Windows Server 2012 and beyond invokes the RDMA capability in the NIC for SMB file traffic if both ends are enabled for RDMA. Broadcom NDIS miniport `bnxtnd.sys` v20.6.2 and beyond support RoCEv1 and RoCEv2 via the NDKPI interface. The default setting is RoCEv1.

To enable RDMA:

1. Upgrade the NIC NVRAM using the appropriate board packages. In CCM or in UEFI HII, enable support for RDMA.
2. Go to the adapter **Advanced Properties** page and set **NetworkDirect Functionality** to **Enabled** for each BCM5741X miniport, or using PowerShell window, run the following command:

```
Set-NetAdapterAdvancedProperty -RegistryKeyword *NetworkDirect -RegistryValue 1
```

3. The following PowerShell commands returns true if **NetworkDirect** is enabled.
 - a. `Get-NetOffLoadGlobalSetting`
 - b. `Get-NetAdapterRDMA`

Verifying RDMA

To verify RDMA:

1. Create a file share on the remote system and open that share using Windows Explorer or "net use". To avoid hard disk read/write speed bottleneck, a RAM disk is recommended as the network share under test.
2. From PowerShell run the following commands:

```
Get-SmbMultichannelConnection | fl *RDMA*
ClientRdmaCapable : True
ServerRdmaCapable : True
```

If both Client and Server show True, then any file transfers over this SMB connection use SMB.

3. The following commands can be used to enable/disable SMB Multichannel:

Server Side:

- Enable: `Set-SmbServerConfiguration -EnableMultiChannel $true`
- Disable: `Set-SmbServerConfiguration -EnableMultiChannel $false`

Client Side:

- Enable: Set-SmbClientConfiguration -EnableMultiChannel \$true
- Disable: Set-SmbClientConfiguration -EnableMultiChannel \$false



Note: By default, the driver sets up two RDMA connections for each network share per IP address (on a unique subnet). The user can scale up the number of RDMA connections by adding multiple IP addresses, each with different a subnet, for the same physical port under test. Multiple network shares can be created and mapped to each link partner using the unique IP addresses created.

For example:

On Server 1, create the following IP addresses for Network Port1.

```
172.1.10.1
172.2.10.2
172.3.10.3
```

On the same Server 1, create 3 shares.

```
Share1
Share2
Share3
```

On the network link partners,

```
Connect to \\172.1.10.1\share1
Connect to \\172.2.10.2\share2
Connect to \\172.3.10.3\share3
...etc
```

User Mode

Before you can run a user mode application written to NDSPI, copy and install the `bxndspi.dll` user mode driver. To copy and install the user mode driver:

1. Copy `bxndspi.dll` to `C:\Windows\System32`.
2. Install the driver by running the following command:

```
rundll32.exe .\bxndspi.dll,Config install|more
```

VMware ESX

Limitations

The current version of the RoCE supported driver requires ESXi-6.5.0 GA build 4564106 or above.

BNXT RoCE Driver Requirements

The BNXTNET L2 driver must be installed with the `disable_roce=0` module parameter before installing the driver.

To set the module parameter, run the following command:

```
esxcfg-module -s "disable_roce=0" bnxtnet
```

Please use the ESX6.5 L2 driver version 20.6.9.0 (RoCE supported L2 driver) or above.

Installation

To install the RoCE driver:

1. Copy the `<bnxtroce>-<driver version>.vib` file in `/var/log/vmware` using the following commands:

```
$ cd /var/log/vmware
$ esxcli software vib install --no-sig-check -v <bnxtroce>-<driver version>.vib
```

2. Reboot the machine.
3. Verify that the drivers are correctly installed using the following command:

```
esxcli software vib list | grep bnxtroce
```

4. To disable ECN (enabled by default) for RoCE traffic use the `"tos_ecn=0"` module parameter for `bnxtroce`.

Configuring Paravirtualized RDMA Network Adapters

Please refer to the vmware link below for additional information on setting up and using Paravirtualized RDMA (PVRDMA) network adapters.

<https://pubs.vmware.com/vsphere-65/index.jsp#com.vmware.vsphere.networking.doc/GUID-4A5EBD44-FB1E-4A83-BB47-BBC65181E1C2.html>

Configuring a Virtual Center for PVRDMA

To configure a Virtual Center for PVRDMA:

1. Create DVS (requires a Distributed Virtual Switch for PVRDMA)
2. Add the host to the DVS.

Tagging vmknic for PVRDMA on ESX Hosts

To tag a vmknic for PVRDMA to use on ESX hosts:

1. Select the host and right-click on **Settings** to switch to the settings page of the **Manage** tabs.
2. In the **Settings** page, expand **System** and click **Advanced System Settings** to show the Advanced System Settings key-pair value and its summary.
3. Click **Edit** to bring up the **Edit Advanced System Settings**.
Filter on **PVRDMA** to narrow all the settings to just `Net.PVRDMAvmknic`.
4. Set the `Net.PVRDMAvmknic` value to `vmknic`, as in example `vmk0`

Setting the Firewall Rule for PVRDMA

To set the firewall rule for PVRDMA:

1. Select the host and right-click on **Settings** to switch to the settings page of the **Manage** tabs.
2. In the **Settings** page, expand **System** and click **Security Profile** to show the firewall summary.
3. Click **Edit** to bring up the **Edit Security Profile**.
4. Scroll down to find `pvrDMA` and check the box to set the firewall.

Adding a PVRDMA Device to the VM

To add a PVRDMA device to the VM:

1. Select the VM and right-click on **Edit Settings**.
2. Add a new Network Adapter.
3. Select the network as a **Distributed Virtual Switch** and **Port Group**.
4. For the **Adapter Type**, select **PVRDMA** and click **OK**.

Configuring the VM on Linux Guest OS



Note: The user must install the appropriate development tools including git before proceeding with the configuration steps below.

1. Download the PVRDMA driver and library using the following commands:

```
git clone git://git.openfabrics.org/~aditr/pvrdma_driver.git
git clone git://git.openfabrics.org/~aditr/libpvrdma.git
```

2. Compile and install the PVRDMA guest driver and library.
3. To install the driver, execute `make && sudo insmod pvrdma.ko` in the directory of the driver.
The driver must be loaded after the paired `vmxnet3` driver is loaded.



The installed RDMA kernel modules may not be compatible with the PVRDMA driver. If so, remove the current installation and restart. Then follow the installation instructions. Please read the README in the driver's directory for more information about the different RDMA stacks.

4. To install the library, execute `./autogen.sh && ./configure --sysconfdir=/etc && make && sudo make install` in the directory of the library.



Note: The installation path of the library needs to be in the shared library cache. Follow the instructions in the INSTALL file in the library's directory.



Note: The firewall settings may need to be modified to allow RDMA traffic. Please ensure the proper firewall settings are in place

5. Add the `/usr/lib` in the `/etc/ld.so.conf` file and reload the `ldconf` by running `ldconfig`
6. Load `ib` modules using `modprobe rdma_ucm`.
7. Load the PVRDMA kernel module using `insmod pvrdma.ko`.
8. Assign an IP address to the PVRDMA interface.
9. Verify whether the IB device is created by running the `ibv_devinfo -v` command.

DCBX – Data Center Bridging

Broadcom NetXtreme-E controllers support IEEE802.1Qaz DCBX as well as the older CEE DCBX specification. DCB configuration is obtained by exchanging the locally configured settings with the link peer. Since the two ends of a link may be configured differently, DCBX uses a concept of 'willing' to indicate which end of the link is ready to accept parameters from the other end. This is indicated in the DCBX protocol using a single bit in the ETS Configuration and PFC TLV, this bit is not used with ETS Recommendation and Application Priority TLV. By Default, the NetXtreme-E NIC is in 'willing' mode while the link partner network switch is in 'non-willing' mode. This ensures the same DCBX setting on the Switch propagates to the entire network.

Users can manually set NetXtreme-E NIC to non-willing mode and perform various PFC, Strict Priority, ETS, APP configurations from the host side. Please refer to the driver readme.txt for more details on available configurations. This document will provide an example of how such setting can be done in Windows with Windows PowerShell. Additional information on DCBX, QoS, and associated use cases are described in more details in a separate white paper, beyond the scope of this user manual.

The following settings in the UEFI HII menu are required to enable DCBX support:

System Setup->Device Settings->NetXtreme-E NIC->Device Level Configuration

QoS Profile – Default QoS Queue Profile

Quality of Server (QoS) resources configuration is necessary to support various PFC and ETS requirements where finer tuning beyond bandwidth allocation is needed. The NetXtreme-E allows the administrator to select between devoting NIC hardware resources to support Jumbo Frames and/or combinations of lossy and lossless Class of Service queues (CoS queues). Many combinations of configuration are possible and therefore can be complicated to compute. This option allows a user to select from a list of precomputed QoS Queue Profiles. These precomputed profiles are design to optimize support for PFC and ETS requirements in typical customer deployments.

The following is a summary description for each QoS Profile.

Table 38: QoS Profiles

Profile No.	Jumbo Frame Support	No. of Lossy CoS Queues/Port	No. of Lossless CoS Queues/Port	Support for 2-Port SKU
Profile #1	Yes	0	1 (PFC Supported)	Yes (25 Gbps)
Profile #2	Yes	4	2 (PFC Supported)	No
Profile #3	No. (MTU <= 2 KB)	6	2 (PFC Supported)	Yes (25 Gbps)
Profile #4	Yes	1	2 (PFC Supported)	Yes (25 Gbps)
Profile #5	Yes	1	0 (No PFC Support)	Yes (25 Gbps)
Profile #6	Yes	8	0 (No PFC Support)	Yes (25 Gbps)
Profile #7	This Configuration maximizes packet-buffer allocations to two lossless CoS Queues to maximize RoCE performance while trading off flexibility.			
	Yes	0	2	Yes (25 Gbps)
Default	Yes	Same as Profile #4		Yes

DCBX Mode = Enable (IEEE only)

This option allows a user to enable/disable DCBX with the indicated specification. IEEE only indicates that IEEE802.1Qaz DCBX is selected.

Windows Driver setting:

After enabling the indicated options in the UEFI HII menu to set firmware level settings, perform the follow selection in the Windows driver advanced properties.

Open **Windows Device Manager -> Broadcom NetXtreme E Series adapter -> Advanced Properties -> Advanced tab**

- Quality of Service = Enabled
- Priority & VLAN = Priority& VLAN enabled
- VLAN = <ID>
- Set desired VLAN id

To exercise the DCB related command in Windows PowerShell, install the appropriate DCB Windows feature.

1. In the **Task Bar**, right-click the Windows PowerShell icon and then click **Run as Administrator**. Windows PowerShell opens in elevated mode.
2. In the Windows PowerShell console, type:

```
Install-WindowsFeature "data-center-bridging"
```

DCBX Willing Bit

The DCBX willing bit is specified in the DCB specification. If the Willing bit on a device is true, the device is willing to accept configurations from a remote device through DCBX. If the Willing bit on a device is false, the device rejects any configuration from a remote device and enforces only the local configurations.

Use the following to set the willing bit to True or False. 1 for enabled, 0 for disabled.

Example `set-netQoSdcbxSetting -willing 1`

Use the following to create a Traffic Class.

```
C:\> New-NetQoSTrafficClass -name "SMB class" -priority 4 -bandwidthPercentage 30 -Algorithm ETS
```



Note: By default, all 802.1p values are mapped to a default traffic class, which has 100% of the bandwidth of the physical link. The command shown above creates a new traffic class to which any packet tagged with eight IEEE 802.1p value 4 is mapped, and its Transmission Selection Algorithm (TSA) is ETS and has 30% of the bandwidth.

It is possible to create up to seven new traffic classes. In addition to the default traffic class, there is at most eight traffic classes in the system.

Use the following in displaying the created Traffic Class:

```
C:\> Get-NetQoSTrafficClass
Name           Algorithm Bandwidth(%) Priority
-----
[Default]     ETS         70         0-3,5-7
SMB class     ETS         30         4
```

Use the following in modifying the Traffic Class:

```
PS C:\> Set-NetQoSTrafficClass -Name "SMB class" -BandwidthPercentage 40
PS C:\> get-NetQoSTrafficClass
Name Algorithm Bandwidth(%) Priority
-----
[Default] ETS         60 0-3,5-7
SMB class ETS         40    4
```

Use the following to remove the Traffic Class:

```
PS C:\> Remove-NetQoSTrafficClass -Name "SMB class"
PS C:\> Get-NetQoSTrafficClass
Name Algorithm Bandwidth(%) Priority
-----
[Default] ETS         100        0-7
```

Use the following to create Traffic Class (Strict Priority):

```
C:\> New-NetQoSTrafficClass -name "SMB class" -priority 4 -bandwidthPercentage 30-Algorithm
Strict
```

Enabling PFC:

```
PS C:\> Enable-NetQoSFlowControl -priority 4
PS C:\> Get-NetQoSFlowControl -priority 4
Priority Enabled
-----
4 True

PS C:\> Get-NetQoSFlowControl
```

Disabling PFC:

```
PS C:\> disable-NetQoSflowControl -priority 4
PS C:\> get-NetQoSFlowControl -priority 4
Priority Enabled
-----
4 False
```

Use the following to create QoS Policy:

```
PS C:\> New-NetQoSPolicy -Name "SMB policy" -SMB -PriorityValue8021Action 4

Name : SMB policy
Owner : Group Policy (Machine)
NetworkProfile : All
Precedence : 127
```



Note: The above command creates a new policy for SMB. –SMB is an inbox filter that matches TCP port 445 (reserved for SMB). If a packet is sent to TCP port 445 it will be tagged by the operating system with 802.1p value of 4 before the packet is passed to a network miniport driver.

In addition to –SMB, other default filters include –iSCSI (matching TCP port 3260), -NFS (matching TCP port 2049), -LiveMigration (matching TCP port 6600), -FCOE (matching EtherType 0x8906) and –NetworkDirect.

NetworkDirect is an abstract layer we create on top of any RDMA implementation on a network adapter. –NetworkDirect must be followed by a Network Direct port.

In addition to the default filters, a user can classify traffic by application's executable name (as in the first example below), or by IP address, port, or protocol.

Use the following to create QoS Policy based on the Source/Destination Address:

```
PS C:\> New-NetQoSPolicy "Network Management" -IPDstPrefixMatchCondition 10.240.1.0/24 -
IPProtocolMatchCondition both -NetworkProfile all -PriorityValue8021Action 7
Name : Network Management
Owner : Group Policy (Machine)
Network Profile : All
Precedence : 127
IPProtocol : Both
IPDstPrefix : 10.240.1.0/24
PriorityValue : 7
```

Use the following to display QoS Policy:

```
PS C:\> Get-NetQoSPolicy
Name : Network Management
Owner : (382ACFAD-1E73-46BD-A0A-6-4EE0E587B95)
NetworkProfile : All
Precedence : 127
IPProtocol : Both
IPDstPrefix : 10.240.1.0/24
PriorityValue : 7
Name : SMB policy
Owner : (382AFAD-1E73-46BD-A0A-6-4EE0E587B95)
NetworkProfile : All
Precedence : 127
Template : SMB
PriorityValue : 4
```

Use the following to modify the QoS Policy:

```
PS C:\> Set-NetqoSPolicy -Name "Network Management" -IPSrcPrefixMatchCondition 10.235.2.0/24 -
IPProtocolMatchCondition both -PriorityValue802.1Action 7
PS C:\> Get-NetQoSPolicy -name "network management"
Name : Network Management
Owner : {382ACFD-1E73-46BD-A0A0-4EE0E587B95}
NetworkProfile : All
Precedence : 127
IPProtocol : Both
IPSrcPrefix : 10.235.2.0/24
IPDstPrefix : 10.240.1.0/24
PriorityValue : 7
```

Use the following to remove QoS Policy:

```
PS C:\> Remove-NetQosPolicy -Name "Network Management"
```

Frequently Asked Questions

- Do you support AutoNeg at 25G speed?

Yes. Please refer to ["Auto-Negotiation Configuration" on page 40](#) for more details.

- How would I connect SFP28 cable to QSFP ports?

Breakout cables are available from QSFP to 4xSFP28 ports.

- What are the compatible port speeds?

For BCM57404AXXX/BCM57414 dual-port devices, the port speed of each port must be compatible with the port speed of the other port. 10 Gbps and 25 Gbps are not compatible speed. If one port is set to 10 Gbps the other port can not be set to 25 Gbps. If user attempt to set incompatible port speeds, the second port to be brought up will not link. Please refer to ["Auto-Negotiation Configuration" on page 40](#) for more details.

- Can I use 10 Gbps for PXE connectivity on a 25 Gbps port?

Currently only 25 Gbps PXE speed is supported. It is not recommended to use 10 Gbps PXE connectivity on a 25 Gbps adapter. This is due to the lack of support for auto negotiation on existing 10 Gbps switches and possible complications caused by incompatible port link speed settings.