



User's Guide

Converged Network Adapters
41xxx Series



AH0054602-00 M October 16, 2019



User's Guide Ethernet iSCSI Adapters and Ethernet FCoE Adapters

For more information, visit our website at: http://www.marvell.com

Notice

THIS DOCUMENT AND THE INFORMATION FURNISHED IN THIS DOCUMENT ARE PROVIDED "AS IS" WITHOUT ANY WARRANTY. MARVELL EXPRESSLY DISCLAIMS AND MAKES NO WARRANTIES OR GUARANTEES REGARDING THE PRODUCT, WHETHER EXPRESS, ORAL, IMPLIED, STATUTORY, ARISING BY OPERATION OF LAW, OR AS A RESULT OF USAGE OF TRADE, COURSE OF DEALING, OR COURSE OF PERFORMANCE, INCLUDING THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR PARTICULAR PURPOSE AND NON-INFRINGEMENT.

Information furnished is believed to be accurate and reliable. However, Marvell assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license, express or implied, to any Marvell intellectual property rights is granted by this document. Marvell products are not authorized for use as critical components in medical devices, military systems, life or critical support devices, or related systems. Marvell retains the right to make changes to this document at any time, without notice.

Export Control

The user or recipient of this document acknowledges that the information included in this document may be subject to laws including, but not limited to, U.S. export control laws and regulations regarding the export, re-export, transfer, diversion, or release of such information. The user or recipient must comply with all applicable laws and regulations at all times. These laws and regulations include restrictions on prohibited destinations, end users, and end uses.

Patents/Trademarks

Products identified in this document may be covered by one or more Marvell patents and/or patent applications. You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning the Marvell products disclosed herein. Marvell and the Marvell logo are registered trademarks of Marvell or its affiliates. Please visit www.marvell.com for a complete list of Marvell trademarks and any guidelines for use of such trademarks. Other names and brands may be claimed as the property of others.

Copyright

Copyright © 2015–2019. Marvell International Ltd. All rights reserved.

Table of Contents

	Preface
	Supported Products xv Intended Audience xv What Is in This Guide xv Documentation Conventions xi Downloading Updates and Documentation xv Legal Notices xx Laser Safety—FDA Notice xx Agency Certification xx EMI and EMC Requirements xx KCC: Class A xxi VCCI: Class A xxi
	Product Safety Compliance
ı	Product Overview Functional Description Features Adapter Specifications Physical Characteristics Standards Specifications
2	Hardware Installation System Requirements Safety Precautions Preinstallation Checklist Installing the Adapter
3	Driver Installation Installing Linux Driver Software Installing the Linux Drivers Without RDMA
	Removing the Linux Drivers
	Installing Linux Drivers Using the kmp/kmod RPM Package 13 Installing Linux Drivers Using the TAR File

	Installing the Linux Drivers with RDMA14
	Linux Driver Optional Parameters
	Linux Driver Operation Defaults
	Linux Driver Messages
	Statistics
	Importing a Public Key for Secure Boot
	Installing Windows Driver Software
	Installing the Windows Drivers
	Running the DUP in the GUI
	DUP Installation Options
	DUP Installation Examples
	Removing the Windows Drivers
	Managing Adapter Properties
	Setting Power Management Options
	Configuring the Communication Protocol to Use with QCC GUI, QCC
	PowerKit, and QCS CLI
	Link Configuration in Windows
	Link Control Mode
	Link Speed and Duplex
	FEC Mode
	Installing VMware Driver Software
	VMware Drivers and Driver Packages
	Installing VMware Drivers
	VMware NIC Driver Optional Parameters
	VMware Driver Parameter Defaults
	Removing the VMware Driver
	FCoE Support
	iSCSI Support
4	Upgrading the Firmware
•	
	Running the DUP by Double-Clicking
	Running the DUP Using the .bin File
	Rulling the DOF Osing the .bill File
5	Adapter Preboot Configuration
	Getting Started
	Displaying Firmware Image Properties
	Configuring Device-level Parameters4
	Configuring NIC Parameters
	Configuring Data Center Bridging

	Configuring FCoE Boot Configuring iSCSI Boot Configuring Partitions Partitioning for VMware ESXi 6.5 and ESXi 6.7	55 56 61 65
6	Boot from SAN Configuration	
	iSCSI Boot from SAN	67
	iSCSI Out-of-Box and Inbox Support	68
	iSCSI Preboot Configuration	68
	Setting the BIOS Boot Mode to UEFI	69
	Enabling NPAR and the iSCSI HBA	71
	Configuring the Storage Target	71
	Selecting the iSCSI UEFI Boot Protocol	72
	Configuring iSCSI Boot Options	73
	Configuring the DHCP Server to Support iSCSI Boot	85
	Configuring iSCSI Boot from SAN on Windows	89
	Before You Begin	90
	Selecting the Preferred iSCSI Boot Mode	90
	Configuring iSCSI General Parameters	90
	Configuring the iSCSI Initiator	91
	Configuring the iSCSI Targets	92
	Detecting the iSCSI LUN and Injecting the Marvell Drivers	92
	Configuring iSCSI Boot from SAN on Linux	94
	Configuring iSCSI Boot from SAN for RHEL 7.5 and Later	95
	Configuring iSCSI Boot from SAN for SLES 12 SP3 and Later .	97
	Configuring iSCSI Boot from SAN for Other Linux Distributions.	97
	Configuring iSCSI Boot from SAN on VMware	108
	Setting the UEFI Main Configuration	108
	Configuring the System BIOS for iSCSI Boot (L2)	110
	Mapping the CD or DVD for OS Installation	112
	FCoE Boot from SAN	113
	FCoE Out-of-Box and Inbox Support	114
	FCoE Preboot Configuration	114
	Specifying the BIOS Boot Protocol	114
	Configuring Adapter UEFI Boot Mode	115
	Configuring FCoE Boot from SAN on Windows	120
	Windows Server 2012 R2 and 2016 FCoE Boot Installation	120
	Configuring FCoE on Windows	121
	FCoE Crash Dump on Windows	121

	Injecting (Slipstreaming) Adapter Drivers into Windows Image Files	122
	Configuring FCoE Boot from SAN on Linux	122
	Prerequisites for Linux FCoE Boot from SAN	122
	Configuring Linux FCoE Boot from SAN	123
	Configuring FCoE Boot from SAN on VMware	123
	Injecting (Slipstreaming) ESXi Adapter Drivers into Image Files	123
	Installing the Customized ESXi ISO	124
7	RoCE Configuration	
•	_	107
	Supported Operating Systems and OFED	127
	Planning for RoCE	128
	Preparing the Adapter	129
	Preparing the Ethernet Switch	129
	Configuring the Cisco Nexus 6000 Ethernet Switch	129
	Configuring the Dell Z9100 Ethernet Switch for RoCE	131
	Configuring RoCE on the Adapter for Windows Server	133
	Viewing RDMA Counters	136
	Configuring RoCE for SR-IOV VF Devices (VF RDMA)	141
	Configuration Instructions	141
	Limitations	149
	Configuring RoCE on the Adapter for Linux	150
	RoCE Configuration for RHEL	150
	RoCE Configuration for SLES	151
	Verifying the RoCE Configuration on Linux	152
	vLAN Interfaces and GID Index Values	154
	RoCE v2 Configuration for Linux	155
	Identifying the RoCE v2 GID Index or Address	155
	Verifying the RoCE v1 or RoCE v2 GID Index and Address	
	from sys and class Parameters	156
	Verifying the RoCE v1 or RoCE v2 Function Through	457
	perftest Applications	157
	Configuring RoCE for SR-IOV VF Devices (VF RDMA)	160
	Enumerating VFs for L2 and RDMA	161
	Number of VFs Supported for RDMA	163
	Limitations	164
	Configuring RoCE on the Adapter for VMware ESX	164
	Configuring RDMA Interfaces	165
	Configuring MTU	166

	RoCE Mode and Statistics	166
	Configuring a Paravirtual RDMA Device (PVRDMA)	167
	Configuring DCQCN	171
	DCQCN Terminology	171
	DCQCN Overview	172
	DCB-related Parameters	173
	Global Settings on RDMA Traffic	173
	Setting vLAN Priority on RDMA Traffic	173
	Setting ECN on RDMA Traffic	173
	Setting DSCP on RDMA Traffic	173
	Configuring DSCP-PFC	173
	Enabling DCQCN	174
	Configuring CNP	174
	DCQCN Algorithm Parameters	174
	MAC Statistics	175
	Script Example	175
	Limitations	176
8	iWARP Configuration	
	Preparing the Adapter for iWARP	177
	Configuring iWARP on Windows	178
	Configuring iWARP on Linux	182
	Installing the Driver	182
	Configuring iWARP and RoCE	182
	Detecting the Device	183
	Supported iWARP Applications	184
	Running Perftest for iWARP	185
	Configuring NFS-RDMA	186
9	iSER Configuration	
	_	188
	Before You Begin	189
	Configuring iSER for RHEL	192
	Using iSER with iWARP on RHEL and SLES	192
	Optimizing Linux Performance	195
	Configuring CPUs to Maximum Performance Mode	195
	Configuring Cros to Maximum Performance Mode	195
		196
	Configuring Block Dovice Staging	
	Configuring Block Device Staging	196

Configuring iSER on ESXi 6.7	196 196 197
iSCSI Configuration	
iSCSI Boot iSCSI Offload in Windows Server. Installing Marvell Drivers Installing the Microsoft iSCSI Initiator Configuring Microsoft Initiator to Use Marvell's iSCSI Offload iSCSI Offload FAQs. Windows Server 2012 R2, 2016, and 2019 iSCSI Boot Installation iSCSI Offload in Linux Environments	200 200 201 201 201 207 208 209 209
Differences from bnx2i	210 210 210
FCoE Configuration	
Configuring Linux FCoE Offload	213 214 214 215
SR-IOV Configuration	
Configuring SR-IOV on Windows	217 224 229 230
NVMe-oF Configuration with RDMA	
Installing Device Drivers on Both Servers Configuring the Target Server Configuring the Initiator Server. Preconditioning the Target Server Testing the NVMe-oF Devices Optimizing Performance. IRQ Affinity (multi_rss-affin.sh)	237 238 240 241 242 243 244 245
	Before You Begin. Configuring iSER for ESXi 6.7. iSCSI Configuration iSCSI Boot iSCSI Offload in Windows Server. Installing Marvell Drivers Installing the Microsoft iSCSI Initiator Configuring Microsoft Initiator to Use Marvell's iSCSI Offload iSCSI Offload FAQs. Windows Server 2012 R2, 2016, and 2019 iSCSI Boot Installation iSCSI Crash Dump iSCSI Offload in Linux Environments Differences from bnx2i. Configuring qedi.ko. Verifying iSCSI Interfaces in Linux FCOE Configuration Configuring Linux FCOE Offload. Differences Between qedf and bnx2fc. Configuring qedf.ko. Verifying FCOE Devices in Linux SR-IOV Configuration Configuring SR-IOV on Windows Configuring SR-IOV on Linux Enabling IOMMU for SR-IOV in UEFI-based Linux OS Installations. Configuring SR-IOV on VMware NVMe-oF Configuration with RDMA Installing Device Drivers on Both Servers Configuring the Target Server Configuring the Initiator Server. Preconditioning the Target Server Testing the NVMe-oF Devices Optimizing Performance

14	VXLAN Configuration	
	Configuring VXLAN in Linux	247
	Configuring VXLAN in VMware	249
	Configuring VXLAN in Windows Server 2016	250
	Enabling VXLAN Offload on the Adapter	250
	Deploying a Software Defined Network	251
15	Windows Server 2016	
	Configuring RoCE Interfaces with Hyper-V	252
	Creating a Hyper-V Virtual Switch with an RDMA NIC	253
	Adding a vLAN ID to Host Virtual NIC	254
	Verifying If RoCE is Enabled	255
	Adding Host Virtual NICs (Virtual Ports)	256
	Mapping the SMB Drive and Running RoCE Traffic	256
	RoCE over Switch Embedded Teaming	258
	Creating a Hyper-V Virtual Switch with SET and RDMA Virtual NICs	258
	Enabling RDMA on SET	258
	Assigning a vLAN ID on SET	259
	Running RDMA Traffic on SET	259
	Configuring QoS for RoCE	259
	Configuring QoS by Disabling DCBX on the Adapter	260
	Configuring QoS by Enabling DCBX on the Adapter	264
	Configuring VMMQ	268
	Enabling VMMQ on the Adapter	269
	Creating a Virtual Machine Switch with or Without SR-IOV	269
	Enabling VMMQ on the Virtual Machine Switch	270
	Getting the Virtual Machine Switch Capability	271
	Creating a VM and Enabling VMMQ on VMNetworkAdapters	271
	in the VM	272
	Monitoring Traffic Statistics	272
	Configuring Storage Spaces Direct	272
	Configuring the Hardware	273
	Deploying a Hyper-Converged System	273
	Deploying the Operating System	274
	Configuring the Network	274
	Configuring Storage Spaces Direct	276
	Configurity Storage Spaces Direct	210

16	Windows Server 2019	
	RSSv2 for Hyper-V. RSSv2 Description Known Event Log Errors Windows Server 2019 Behaviors VMMQ Is Enabled by Default Inbox Driver Network Direct (RDMA) Is Disabled by Default. New Adapter Properties Max Queue Pairs (L2) Per VPort. Network Direct Technology Virtualization Resources VMQ and VMMQ Default Accelerations Single VPort Pool	280 281 281 281 281 282 282 282 283 284 284
17	Troubleshooting	
	Troubleshooting Checklist Verifying that Current Drivers Are Loaded Verifying Drivers in Windows Verifying Drivers in Linux Verifying Drivers in VMware Testing Network Connectivity Testing Network Connectivity for Windows Testing Network Connectivity for Linux Microsoft Virtualization with Hyper-V Linux-specific Issues Miscellaneous Issues Collecting Debug Data	286 287 287 288 288 288 289 289 289 290
A	Adapter LEDS	
В	Cables and Optical Modules Supported Specifications	292 293 297
С	Dell Z9100 Switch Configuration	
D	Feature Constraints	
E	Revision History	
Glossary		

List of Figures

Figure		Page
3-1	Dell Update Package Window	19
3-2	QLogic InstallShield Wizard: Welcome Window	19
3-3	QLogic InstallShield Wizard: License Agreement Window	20
3-4	InstallShield Wizard: Setup Type Window	21
3-5	InstallShield Wizard: Custom Setup Window	22
3-6	InstallShield Wizard: Ready to Install the Program Window	22
3-7	InstallShield Wizard: Completed Window	23
3-8	Dell Update Package Window	24
3-9	Setting Advanced Adapter Properties	26
3-10	Power Management Options	27
3-11	Setting Driver Controlled Mode	29
3-12	Setting the Link Speed and Duplex Property	30
3-13	Setting the FEC Mode Property	31
4-1	Dell Update Package: Splash Screen	39
4-2	Dell Update Package: Loading New Firmware	40
4-3	Dell Update Package: Installation Results	40
4-4	Dell Update Package: Finish Installation	41
4-5	DUP Command Line Options	42
5-1	System Setup	45
5-2	System Setup: Device Settings	45
5-3	Main Configuration Page	46
5-4	Main Configuration Page, Setting Partitioning Mode to NPAR	46
5-5	Firmware Image Properties	48
5-6	Device Level Configuration	49
5-7	NIC Configuration	50
5-8	System Setup: Data Center Bridging (DCB) Settings	54
5-9	FCoE General Parameters	55
5-10	FCoE Target Configuration	56
5-11	iSCSI General Parameters	58
5-12	iSCSI Initiator Configuration Parameters	59
5-13	iSCSI First Target Parameters	59
5-14	iSCSI Second Target Parameters	60
5-15	NIC Partitioning Configuration, Global Bandwidth Allocation	61
5-16	Global Bandwidth Allocation Page	62
5-17	Partition 1 Configuration	63
5-18	Partition 2 Configuration: FCoE Offload	64
5-19	Partition 3 Configuration: iSCSI Offload	65
5-20	Partition 4 Configuration	65
6-1	System Setup: Boot Settings	70
6-2	System Setup: Device Settings	71
6-3	System Setup: NIC Configuration	72
6-4	System Setup: NIC Configuration, Boot Protocol	73
6-5	System Setup: iSCSI Configuration	74

6-6	System Setup: Selecting General Parameters	74
6-7	System Setup: iSCSI General Parameters	75
6-8	System Setup: Selecting iSCSI Initiator Parameters	77
6-9	System Setup: iSCSI Initiator Parameters	78
6-10	System Setup: Selecting iSCSI First Target Parameters	79
6-11	System Setup: iSCSI First Target Parameters	80
6-12	System Setup: iSCSI Second Target Parameters	81
6-13	System Setup: Saving iSCSI Changes	82
6-14	System Setup: iSCSI General Parameters	84
6-15	System Setup: iSCSI General Parameters, VLAN ID	89
6-16	Detecting the iSCSI LUN Using UEFI Shell (Version 2)	93
6-17	Windows Setup: Selecting Installation Destination	93
6-18	Windows Setup: Selecting Driver to Install	94
6-19	Integrated NIC: Device Level Configuration for VMware	108
6-20	Integrated NIC: Partition 2 Configuration for VMware	109
6-21	Integrated NIC: System BIOS, Boot Settings for VMware	110
6-22	Integrated NIC: System BIOS, Connection 1 Settings for VMware	111
6-23	Integrated NIC: System BIOS, Connection 1 Settings (Target) for VMware	111
6-24	VMware iSCSI BFS: Selecting a Disk to Install	112
6-25	VMware iSCSI Boot from SAN Successful	113
6-26	System Setup: Selecting Device Settings	115
6-27	System Setup: Device Settings, Port Selection	116
6-28	System Setup: NIC Configuration	117
6-29	System Setup: FCoE Mode Enabled	118
6-30	, ,	119
6-31		120
6-32		124
6-33		125
6-34	· ·	126
7-1		134
7-2	5	136
7-3	!	138
7-4	5	142
7-5	5	143
7-6		144
7-7	3	145
7-8	j i	146
7-9		147
7-10	0	148
7-11		149
7-12	5 ,	159
7-13	\mathbf{o}	159
7-14	<u> </u>	160
7-15	0 0 = 11	160
7-16	Configuring a New Distributed Switch	168

7-17	Assigning a vmknic for PVRDMA	169
7-18	Setting the Firewall Rule	170
8-1	Windows PowerShell Command: Get-NetAdapterRdma	179
8-2	Windows PowerShell Command: Get-NetOffloadGlobalSetting	179
8-3	Perfmon: Add Counters	180
8-4	Perfmon: Verifying iWARP Traffic	181
9-1	RDMA Ping Successful	190
9-2	iSER Portal Instances	190
9-3	Iface Transport Confirmed	191
9-4	Checking for New iSCSI Device	192
9-5	LIO Target Configuration	194
10-1	iSCSI Initiator Properties, Configuration Page	202
10-2	iSCSI Initiator Node Name Change	202
10-3	iSCSI Initiator—Discover Target Portal	203
10-4	Target Portal IP Address	204
10-5	Selecting the Initiator IP Address	205
10-6	Connecting to the iSCSI Target	206
10-7	Connect To Target Dialog Box	207
12-1	System Setup for SR-IOV: Integrated Devices	218
12-2	System Setup for SR-IOV: Device Level Configuration	218
12-3	Adapter Properties, Advanced: Enabling SR-IOV	219
12-4	Virtual Switch Manager: Enabling SR-IOV	220
12-5	Settings for VM: Enabling SR-IOV	222
12-6	Device Manager: VM with QLogic Adapter	223
12-7	Windows PowerShell Command: Get-NetadapterSriovVf	223
12-8	System Setup: Processor Settings for SR-IOV	224
12-9	System Setup for SR-IOV: Integrated Devices	225
12-10		226
12-11	Command Output for sriov_numvfs	227
	Command Output for ip link show Command	227
	RHEL68 Virtual Machine	228
	Add New Virtual Hardware	229
	VMware Host Edit Settings	233
13-1	NVMe-oF Network	236
13-2	Subsystem NQN	240
13-3	Confirm NVMe-oF Connection	241
13-4	FIO Utility Installation	242
14-1	Advanced Properties: Enabling VXLAN	250
15-1	Enabling RDMA in Host Virtual NIC	253
15-2	Hyper-V Virtual Ethernet Adapter Properties	254
15-3	Windows PowerShell Command: Get-VMNetworkAdapter	255
15-4	Windows PowerShell Command: Get-NetAdapterRdma	255
15-5	Add Counters Dialog Box	257
15-6	Performance Monitor Shows RoCE Traffic	257
15-7	Windows PowerShell Command: New-VMSwitch	258

User's Guide—Converged Network Adapters 41xxx Series

15-8	Windows PowerShell Command: Get-NetAdapter	259
15-9	Advanced Properties: Enable QoS	261
15-10	Advanced Properties: Setting VLAN ID	262
15-11	Advanced Properties: Enabling QoS	266
15-12	Advanced Properties: Setting VLAN ID	267
15-13	Advanced Properties: Enabling Virtual Switch RSS	269
15-14	Virtual Switch Manager	270
15-15	Windows PowerShell Command: Get-VMSwitch	271
15-16	Example Hardware Configuration	273
16-1	RSSv2 Event Log Error	281

List of Tables

Table		Page
2-1	Host Hardware Requirements	4
2-2	Minimum Host Operating System Requirements	5
3-1	41xxx Series Adapters Linux Drivers	9
3-2	qede Driver Optional Parameters	15
3-3	Linux Driver Operation Defaults	16
3-4	VMware Drivers	32
3-5	VMware NIC Driver Optional Parameters	34
3-6	VMware Driver Parameter Defaults	36
5-1	Adapter Properties	47
6-1	iSCSI Out-of-Box and Inbox Boot from SAN Support	68
6-2	iSCSI General Parameters	76
6-3	DHCP Option 17 Parameter Definitions	85
6-4	DHCP Option 43 Sub-option Definitions	86
6-5	DHCP Option 17 Sub-option Definitions	88
6-6	FCoE Out-of-Box and Inbox Boot from SAN Support	114
7-1	OS Support for RoCE v1, RoCE v2, iWARP, iSER, and OFED	127
7-2	Advanced Properties for RoCE	134
7-3	Marvell FastLinQ RDMA Error Counters	138
7-4	Supported Linux OSs for VF RDMA	161
7-5	DCQCN Algorithm Parameters	174
13-1	Target Parameters	238
16-1	Windows 2019 Virtualization Resources for Dell 41xxx Series Adapters	283
16-2	Windows 2019 VMQ and VMMQ Accelerations	284
17-1	Collecting Debug Data Commands	290
A-1	Adapter Port Link and Activity LEDs	291
B-1	Tested Cables and Optical Modules	293
B-2	Switches Tested for Interoperability	297

Preface

This preface lists the supported products, specifies the intended audience, explains the typographic conventions used in this guide, and describes legal notices.

Supported Products

NOTE

QConvergeConsole® (QCC) GUI is the *only* GUI management tool across all Marvell® FastLinQ® adapters. QLogic Control Suite™ (QCS) GUI is no longer supported for the FastLinQ 45000 Series Adapters and adapters based on 57xx/57xxx controllers, and has been replaced by the QCC GUI management tool. The QCC GUI provides single-pane-of-glass GUI management for all Marvell adapters.

In Windows® environments, when you run QCS CLI and Management Agents Installer, it will uninstall the QCS GUI (if installed on the system) and any related components from your system. To obtain the new GUI, download QCC GUI for your adapter from the Marvell Web site (see "Downloading Updates and Documentation" on page xxi).

This user's guide describes the following Marvell products:

- QL41112HFCU-DE 10Gb Converged Network Adapter, full-height bracket
- QL41112HLCU-DE 10Gb Converged Network Adapter, low-profile bracket
- QL41132HFRJ-DE 10Gb NIC Adapter, full-height bracket
- QL41132HLRJ-DE 10Gb NIC Adapter, low-profile bracket
- QL41132HQCU-DE 10Gb NIC Adapter
- QL41132HQRJ-DE 10Gb NIC Adapter
- QL41154HQRJ-DE 10Gb Converged Network Adapter
- QL41154HQCU-DE 10Gb Converged Network Adapter
- QL41162HFRJ-DE 10Gb Converged Network Adapter, full-height bracket
- QL41162HLRJ-DE 10Gb Converged Network Adapter, low-profile bracket

- QL41162HMRJ-DE 10Gb Converged Network Adapter
- QL41164HMCU-DE 10Gb Converged Network Adapter
- QL41164HMRJ-DE 10Gb Converged Network Adapter
- QL41164HFRJ-DE 10Gb Converged Network Adapter, full-height bracket
- QL41164HFRJ-DE 10Gb Converged Network Adapter, low-profile bracket
- QL41164HFCU-DE 10Gb Converged Network Adapter, full-height bracket
- QL41232HFCU-DE 10/25Gb NIC Adapter, full-height bracket
- QL41232HLCU-DE 10/25Gb NIC Adapter, low-profile bracket
- QL41232HMKR-DE 10/25Gb NIC Adapter
- QL41232HQCU-DE 10/25Gb NIC Adapter
- QL41262HFCU-DE 10/25Gb Converged Network Adapter, full-height bracket
- QL41262HLCU-DE 10/25Gb Converged Network Adapter, low-profile bracket
- QL41262HMCU-DE 10/25Gb Converged Network
- QL41262HMKR-DE 10/25Gb Converged Network Adapter
- QL41264HMCU-DE 10/25Gb Converged Network Adapter

Intended Audience

This guide is intended for system administrators and other technical staff members responsible for configuring and managing adapters installed on Dell[®] PowerEdge[®] servers in Windows[®], Linux[®], or VMware[®] environments.

What Is in This Guide

Following this preface, the remainder of this guide is organized into the following chapters and appendices:

- Chapter 1 Product Overview provides a product functional description, a list of features, and the adapter specifications.
- Chapter 2 Hardware Installation describes how to install the adapter, including the list of system requirements and a preinstallation checklist.
- Chapter 3 Driver Installation describes the installation of the adapter drivers on Windows, Linux, and VMware.
- Chapter 4 Upgrading the Firmware describes how to use the Dell Update Package (DUP) to upgrade the adapter firmware.

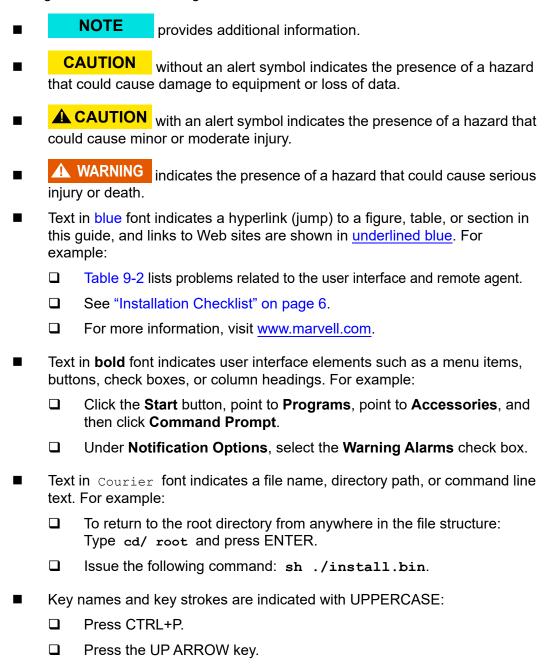
- Chapter 5 Adapter Preboot Configuration describes the preboot adapter configuration tasks using the Human Infrastructure Interface (HII) application.
- Chapter 6 Boot from SAN Configuration covers boot from SAN configuration for both iSCSI and FCoE.
- Chapter 7 RoCE Configuration describes how to configure the adapter, the Ethernet switch, and the host to use RDMA over converged Ethernet (RoCE).
- Chapter 8 iWARP Configuration provides procedures for configuring Internet wide area RDMA protocol (iWARP) on Windows, Linux, and VMware ESXi 6.7 systems.
- Chapter 9 iSER Configuration describes how to configure iSCSI Extensions for RDMA (iSER) for Linux RHEL, SLES, Ubuntu, and ESXi 6.7.
- Chapter 10 iSCSI Configuration describes iSCSI boot and iSCSI offload for Windows and Linux.
- Chapter 11 FCoE Configuration covers configuring Linux FCoE offload.
- Chapter 12 SR-IOV Configuration provides procedures for configuring single root input/output virtualization (SR-IOV) on Windows, Linux, and VMware systems.
- Chapter 13 NVMe-oF Configuration with RDMA demonstrates how to configure NVMe-oF on a simple network for 41xxx Series Adapters.
- Chapter 14 VXLAN Configuration describes how to configure VXLAN for Linux, VMware, and Windows Server 2016.
- Chapter 15 Windows Server 2016 describes the Windows Server 2016 features.
- Chapter 16 Windows Server 2019 describes the Windows Server 2019 features.
- Chapter 17 Troubleshooting describes a variety of troubleshooting methods and resources.
- Appendix A Adapter LEDS lists the adapter LEDs and their significance.
- Appendix B Cables and Optical Modules lists the cables, optical modules, and switches that the 41xxx Series Adapters support.
- Appendix C Dell Z9100 Switch Configuration describes how to configure the Dell Z9100 switch port for 25Gbps.
- Appendix D Feature Constraints provides information about feature constraints implemented in the current release.

Appendix E Revision History describes the changes made in this revision of the guide.

At the end of this guide is a glossary of terms.

Documentation Conventions

This guide uses the following documentation conventions:



•	Text in <i>italics</i> indicates terms, emphasis, variables, or document titles. For example:				
		What are shortcut keys?			
		To enter the date type $mm/dd/yyyy$ (where mm is the month, dd is the day, and $yyyy$ is the year).			
•	Topic titles between quotation marks identify related topics either within thi manual or in the online help, which is also referred to as <i>the help system</i> throughout this document.				
Command line interface (CLI) command syntax conventions include following:					
		Plain text indicates items that you must type as shown. For example:			
		■ qaucli -pr nic -ei			
		< > (angle brackets) indicate a variable whose value you must specify. For example:			
		<pre><serial_number></serial_number></pre>			
		NOTE			
		For CLI commands only, variable names are always indicated using angle brackets instead of <i>italics</i> .			
		[] (square brackets) indicate an optional parameter. For example:			
		[<file_name>] means specify a file name, or omit it to select the default file name.</file_name>			
		(vertical bar) indicates mutually exclusive options; select one option only. For example:			
		■ on off			
		1 2 3 4			
		\dots (ellipsis) indicates that the preceding item may be repeated. For example:			
		■ x means <i>one</i> or more instances of x.			
		x means one of more instances of x.			
		■ [x] means <i>zero</i> or more instances of x.			

- () (parentheses) and { } (braces) are used to avoid logical ambiguity. For example:
 - a|b c is ambiguous
 {(a|b) c} means a or b, followed by c
 {a|(b c)} means either a, or b c

Downloading Updates and Documentation

The Marvell Web site provides periodic updates to product firmware, software, and documentation.

To download Marvell firmware, software, and documentation:

- 1. Go to www.marvell.com.
- 2. Point to **Support**, and then under **Driver Downloads**, click **Marvell QLogic/FastLinQ Drivers**.
- 3. On the Drivers and Documentation page, click **Adapters**.
- 4. Click the corresponding button to search by Model or by Operating System.
- 5. To define a search, click an item in each selection column, and then click **Go**.
- 6. Locate the firmware, software, or document you need, and then click the item's name or icon to download or open the item.

Legal Notices

Legal notices covered in this section include laser safety (FDA notice), agency certification, and product safety compliance.

Laser Safety—FDA Notice

This product complies with DHHS Rules 21CFR Chapter I, Subchapter J. This product has been designed and manufactured according to IEC60825-1 on the safety label of laser product.

CLASS I LASER

Class 1 Caution—Class 1 laser radiation when open Laser Product Do not view directly with optical instruments

Appareil laser Attention—Radiation laser de classe 1

de classe 1 Ne pas regarder directement avec des instruments optiques

Produkt der Vorsicht—Laserstrahlung der Klasse 1 bei geöffneter Abdeckung

Laser Klasse 1 Direktes Ansehen mit optischen Instrumenten vermeiden

Luokan 1 Laserlaite Varoitus—Luokan 1 lasersäteilyä, kun laite on auki

Älä katso suoraan laitteeseen käyttämällä optisia instrumenttej

Agency Certification

The following sections summarize the EMC and EMI test specifications performed on the 41xxx Series Adapters to comply with emission, immunity, and product safety standards.

EMI and EMC Requirements

FCC Part 15 compliance: Class A

FCC compliance information statement: This device complies with Part 15 of the FCC Rules. Operation is subject to the following two conditions: (1) this device may not cause harmful interference, and (2) this device must accept any interference received, including interference that may cause undesired operation.

ICES-003 Compliance: Class A

This Class A digital apparatus complies with Canadian ICES-003. Cet appareil numériqué de la classe A est conformé à la norme NMB-003 du Canada.

CE Mark 2014/30/EU, 2014/35/EU EMC Directive Compliance:

EN55032:2012/ CISPR 32:2015 Class A

EN55024: 2010

EN61000-3-2: Harmonic Current Emission EN61000-3-3: Voltage Fluctuation and Flicker

Immunity Standards

EN61000-4-2: ESD

EN61000-4-3: RF Electro Magnetic Field

EN61000-4-4: Fast Transient/Burst

EN61000-4-5: Fast Surge Common/ Differential EN61000-4-6: RF Conducted Susceptibility EN61000-4-8: Power Frequency Magnetic Field EN61000-4-11: Voltage Dips and Interrupt

VCCI: 2015-04; Class A

AS/NZS; CISPR 32: 2015 Class A

CNS 13438: 2006 Class A

KCC: Class A

Korea RRA Class A Certified



Product Name/Model: Converged Network Adapters and

Intelligent Ethernet Adapters

Certification holder: QLogic Corporation

Manufactured date: Refer to date code listed on product Manufacturer/Country of origin: QLogic Corporation/USA

A class equipment

(Business purpose info/telecommunications

equipment)

As this equipment has undergone EMC registration for business purpose, the seller and/or the buyer is asked to beware of this point and in case a wrongful sale or purchase has been made, it is asked that a change to household use be

made.

Korean Language Format—Class A

A급 기기 (업무용 정보통신기기)

이 기기는 업무용으로 전자파적합등록을 한 기기이오니 판매자 또는 사용자는 이 점을 주의하시기 바라며, 만약 잘못판매 또는 구입하였을 때에는 가정용으로 교환하시기 바랍니다.

VCCI: Class A

This is a Class A product based on the standard of the Voluntary Control Council for Interference (VCCI). If this equipment is used in a domestic environment, radio interference may occur, in which case the user may be required to take corrective actions.

この装置は、クラスA情報技術装置です。この装置を家庭環境で使用すると電波妨害を引き起こすことがあります。この場合には使用者が適切な対策を講ずるよう要求されることがあります。 VCCI-A

Product Safety Compliance

UL, cUL product safety:

UL 60950-1 (2nd Edition) A1 + A2 2014-10-14 CSA C22.2 No.60950-1-07 (2nd Edition) A1 +A2 2014-10

Use only with listed ITE or equivalent.

Complies with 21 CFR 1040.10 and 1040.11, 2014/30/EU, 2014/35/EU.

2006/95/EC low voltage directive:

TUV EN60950-1:2006+A11+A1+A12+A2 2nd Edition TUV IEC 60950-1: 2005 2nd Edition Am1: 2009 + Am2: 2013 CB

CB Certified to IEC 60950-1 2nd Edition

1 Product Overview

This chapter provides the following information for the 41xxx Series Adapters:

- Functional Description
- Features
- "Adapter Specifications" on page 3

Functional Description

The Marvell FastLinQ 41000 Series Adapters include 10 and 25Gb Converged Network Adapters and Intelligent Ethernet Adapters that are designed to perform accelerated data networking for server systems. The 41000 Series Adapter includes a 10/25Gb Ethernet MAC with full-duplex capability.

Using the operating system's teaming feature, you can split your network into virtual LANs (vLANs), as well as group multiple network adapters together into teams to provide network load balancing and fault tolerance. For more information about teaming, see your operating system documentation.

Features

The 41xxx Series Adapters provide the following features. Some features may not be available on all adapters:

- NIC partitioning (NPAR)
- Single-chip solution:
 - □ 10/25Gb MAC
 - □ SerDes interface for direct attach copper (DAC) transceiver connection
 - ☐ PCI Express® (PCle®) 3.0 x8
 - Zero copy capable hardware
- Performance features:
 - ☐ TCP, IP, UDP checksum offloads
 - ☐ TCP segmentation offload (TSO)
 - Large segment offload (LSO)

	Generic segment offload (GSO)		
	Large receive offload (LRO)		
	Receive segment coalescing (RSC)		
	$Microsoft^{@}$ dynamic virtual machine queue (VMQ), and Linux Multiqueue		
Adap	otive interrupts:		
	Transmit/receive side scaling (TSS/RSS)		
	Stateless offloads for Network Virtualization using Generic Routing Encapsulation (NVGRE) and virtual LAN (VXLAN) L2/L3 GRE tunneled traffic ¹		
Manageability:			
	System management bus (SMB) controller		
	Advanced Configuration and Power Interface (ACPI) 1.1a compliant (multiple power modes)		
	Network controller-sideband interface (NC-SI) support		
Adva	Advanced network features:		
	Jumbo frames (up to 9,600 bytes). The OS and the link partner must support jumbo frames.		
	Virtual LANs (VLANs)		
	Flow control (IEEE Std 802.3x)		
Logi	cal link control (IEEE Std 802.2)		
 High-speed on-chip reduced instruction set computer (RISC) proces 			
Integ	rated 96KB frame buffer memory (not applicable to all models)		
1,02	4 classification filters (not applicable to all models)		
Supp	port for multicast addresses through 128-bit hashing hardware function		
Support for VMDirectPath I/O			
FastLinQ 41xxx Series Adapters support VMDirectPath I/O in Linux and ESX environments. VMDirectPath I/O is not supported in Windows environments.			
PCI all P virtua	LinQ 41xxx Series Adapters can be assigned to virtual machines for pass-through operation. However, due to function level dependencies, Cle functions associated with an adapter must be assigned to the same all machine. Sharing PCIe functions across the hypervisor and/or one or evirtual machines is not supported.		

2

¹ This feature requires OS or Hypervisor support to use the offloads.

- Serial flash NVRAM memory
- PCI Power Management Interface (v1.1)
- 64-bit base address register (BAR) support
- EM64T processor support
- iSCSI and FCoE boot support²

Adapter Specifications

The 41xxx Series Adapter specifications include the adapter's physical characteristics and standards-compliance references.

Physical Characteristics

The 41xxx Series Adapters are standard PCle cards and ship with either a full-height or a low-profile bracket for use in a standard PCle slot.

Standards Specifications

Supported standards specifications include:

- PCI Express Base Specification, rev. 3.1
- PCI Express Card Electromechanical Specification, rev. 3.0
- PCI Bus Power Management Interface Specification, rev. 1.2
- IEEE Specifications:
 - □ 802.1ad (QinQ)
 - 802.1AX (Link Aggregation)
 - □ 802.1p (Priority Encoding)
 - □ 802.1q (VLAN)
 - 802.3-2015 IEEE Standard for Ethernet (flow control)
 - □ 802.3-2015 Clause 78 Energy Efficient Ethernet (EEE)
 - ☐ 1588-2002 PTPv1 (Precision Time Protocol)
 - □ 1588-2008 PTPv2
- IPv4 (RFQ 791)
- IPv6 (RFQ 2460)

² Hardware support limit of SR-IOV VFs varies. The limit may be lower in some OS environments; refer to the appropriate section for your OS.

2 Hardware Installation

This chapter provides the following hardware installation information:

- System Requirements
- "Safety Precautions" on page 5
- "Preinstallation Checklist" on page 6
- "Installing the Adapter" on page 6

System Requirements

Before you install a Marvell 41xxx Series Adapter, verify that your system meets the hardware and operating system requirements shown in Table 2-1 and Table 2-2. For a complete list of supported operating systems, visit the Marvell Web site.

Table 2-1. Host Hardware Requirements

Hardware	Requirement
Architecture	IA-32 or EMT64 that meets operating system requirements
PCle	PCIe Gen 2 x8 (2x10G NIC) PCIe Gen 3 x8 (2x25G NIC) Full dual-port 25Gb bandwidth is supported on PCIe Gen 3 x8 or faster slots.
Memory	8GB RAM (minimum)
Cables and Optical Modules	The 41xxx Series Adapters have been tested for interoperability with a variety of 1G, 10G, and 25G cables and optical modules. See "Tested Cables and Optical Modules" on page 293.

Table 2-2. Minimum Host Operating System Requirements

Operating System	Requirement
Windows Server	2012 R2, 2019
Linux	RHEL [®] 7.6, 7.7, 8.0, 8.1 SLES [®] 12 SP4, SLES 15, SLES 15 SP1 CentOS 7.6
VMware	vSphere [®] ESXi 6.5 U3 and vSphere ESXi 6.7 U3
XenServer	Citrix Hypervisor 8.0 7.0, 7.1

NOTE

Table 2-2 denotes minimum host OS requirements. For a complete list of supported operating systems, visit the Marvell Web site.

Safety Precautions

MARNING

The adapter is being installed in a system that operates with voltages that can be lethal. Before you open the case of your 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 is 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.

Preinstallation Checklist

Before installing the adapter, complete the following:

- 1. Verify that the system meets the hardware and software requirements listed under "System Requirements" on page 4.
- 2. Verify that the system is using the latest BIOS.

NOTE

If you acquired the adapter software from the Marvell Web site, verify the path to the adapter driver files.

- 3. If the system is active, shut it down.
- 4. When system shutdown is complete, turn off the power and unplug the power cord.
- 5. Remove the adapter from its shipping package and place it on an anti-static surface.
- 6. Check the adapter for visible signs of damage, particularly on the edge connector. Never attempt to install a damaged adapter.

Installing the Adapter

The following instructions apply to installing the Marvell 41xxx Series Adapters in most systems. For details about performing these tasks, refer to the manuals that were supplied with the system.

To install the adapter:

- 1. Review "Safety Precautions" on page 5 and "Preinstallation Checklist" on page 6. Before you install the adapter, ensure that the system power is OFF, the power cord is unplugged from the power outlet, and that you are following proper electrical grounding procedures.
- 2. Open the system case, and select the slot that matches the adapter size, which can be PCle Gen 2 x8 or PCle Gen 3 x8. A lesser-width adapter can be seated into a greater-width slot (x8 in an x16), but a greater-width adapter cannot be seated into a lesser-width slot (x8 in an x4). If you do not know how to identify a PCle slot, refer to your system documentation.
- 3. Remove the blank cover-plate from the slot that you selected.
- 4. Align the adapter connector edge with the PCIe connector slot in the system.

5. Applying even pressure at both corners of the card, push the adapter card into the slot until it is firmly seated. When the adapter is properly seated, the adapter port connectors are aligned with the slot opening, and the adapter faceplate is flush against the system chassis.

CAUTION

Do not use excessive force when seating the card, because this may damage the system or the adapter. If you have difficulty seating the adapter, remove it, realign it, and try again.

- 6. Secure the adapter with the adapter clip or screw.
- 7. Close the system case and disconnect any personal anti-static devices.

7

3 Driver Installation

This chapter provides the following information about driver installation:

- Installing Linux Driver Software
- "Installing Windows Driver Software" on page 18
- "Installing VMware Driver Software" on page 31

Installing Linux Driver Software

This section describes how to install Linux drivers with or without remote direct memory access (RDMA). It also describes the Linux driver optional parameters, default values, messages, statistics, and public key for Secure Boot.

- Installing the Linux Drivers Without RDMA
- Installing the Linux Drivers with RDMA
- Linux Driver Optional Parameters
- Linux Driver Operation Defaults
- Linux Driver Messages
- Statistics
- Importing a Public Key for Secure Boot

The 41xxx Series Adapter Linux drivers and supporting documentation are available on the Dell Support page:

dell.support.com

Table 3-1 describes the 41xxx Series Adapter Linux drivers.

Table 3-1. 41xxx Series Adapters Linux Drivers

Linux Driver	Description
qed	The qed core driver module directly controls the firmware, handles interrupts, and provides the low-level API for the protocol specific driver set. The qed interfaces with the qede, qedr, qedi, and qedf drivers. The Linux core module manages all PCI device resources (registers, host interface queues, and so on). The qed core module requires Linux kernel version 2.6.32 or later. Testing was concentrated on the x86_64 architecture.
qede	Linux Ethernet driver for the 41xxx Series Adapter. This driver directly controls the hardware and is responsible for sending and receiving Ethernet packets on behalf of the Linux host networking stack. This driver also receives and processes device interrupts on behalf of itself (for L2 networking). The qede driver requires Linux kernel version 2.6.32 or later. Testing was concentrated on the x86_64 architecture.
qedr	Linux RoCE driver that works in the OpenFabrics Enterprise Distribution (OFED™) environment in conjunction with the qed core module and the qede Ethernet driverRDMA user space applications also require that the libqedr user library is installed on the server
qedi	Linux iSCSI-Offload driver for the 41xxx Series Adapters. This driver works with the Open iSCSI library.y
qedf	Linux FCoE-Offload driver for the 41xxx Series Adapters. This driver works with Open FCoE library.

Install the Linux drivers using either a source Red Hat[®] Package Manager (RPM) package or a kmod RPM package. The RHEL RPM packages are as follows:

- qlgc-fastlinq-<version>.<OS>.src.rpm
- qlgc-fastlinq-kmp-default-<version>.<arch>.rpm

The SLES source and kmp RPM packages are as follows:

- qlgc-fastlinq-<version>.<OS>.src.rpm
- qlgc-fastlinq-kmp-default-<version>.<OS>.<arch>.rpm

The following kernel module (kmod) RPM installs Linux drivers on SLES hosts running the Xen Hypervisor:

■ qlgc-fastlinq-kmp-xen-<version>.<OS>.<arch>.rpm

The following source RPM installs the RDMA library code on RHEL and SLES hosts:

qlgc-libqedr-<version>.<OS>.<arch>.src.rpm

The following source code TAR BZip2 (BZ2) compressed file installs Linux drivers on RHEL and SLES hosts:

fastling-<version>.tar.bz2

NOTE

For network installations through NFS, FTP, or HTTP (using a network boot disk), you may require a driver disk that contains the qede driver. Compile the Linux boot drivers by modifying the makefile and the make environment.

Installing the Linux Drivers Without RDMA

To install the Linux drivers without RDMA:

- Download the 41xxx Series Adapter Linux drivers from Dell: dell.support.com
- 2. Remove the existing Linux drivers, as described in "Removing the Linux Drivers" on page 10.
- 3. Install the new Linux drivers using one of the following methods:
 - ☐ Installing Linux Drivers Using the src RPM Package
 - ☐ Installing Linux Drivers Using the kmp/kmod RPM Package
 - ☐ Installing Linux Drivers Using the TAR File

Removing the Linux Drivers

There are two procedures for removing Linux drivers: one for a non-RDMA environment and another for an RDMA environment. Choose the procedure that matches your environment.

To remove Linux drivers in a non-RDMA environment, unload and remove the drivers:

Follow the procedure that relates to the original installation method and the OS.

■ If the Linux drivers were installed using an RPM package, issue the following commands:

```
rmmod qede
rmmod qed
depmod -a
rpm -e glqc-fastling-kmp-default-<version>.<arch>
```

■ If the Linux drivers were installed using a TAR file, issue the following commands:

```
rmmod qede
```

```
rmmod qed
depmod -a

    For RHEL:
    cd /lib/modules/<version>/extra/qlgc-fastlinq
    rm -rf qed.ko qede.ko qedr.ko

    For SLES:
    cd /lib/modules/<version>/updates/qlgc-fastlinq
    rm -rf qed.ko qede.ko qedr.ko
```

To remove Linux drivers in a non-RDMA environment:

1. To get the path to the currently installed drivers, issue the following command:

```
modinfo <driver name>
```

- 2. Unload and remove the Linux drivers.
 - ☐ If the Linux drivers were installed using an RPM package, issue the following commands:

```
modprobe -r qede
depmod -a
rpm -e qlgc-fastlinq-kmp-default-<version>.<arch>
```

☐ If the Linux drivers were installed using a TAR file, issue the following commands:

```
modprobe -r qede
depmod -a
```

NOTE

If the qedr is present, issue the <code>modprobe -r qedr command</code> instead.

3. Delete the qed.ko, qede.ko, and qedr.ko files from the directory in which they reside. For example, in SLES, issue the following commands:

```
cd /lib/modules/<version>/updates/qlgc-fastlinq
rm -rf qed.ko
rm -rf qede.ko
rm -rf qedr.ko
depmod -a
```

To remove Linux drivers in an RDMA environment:

1. To get the path to the installed drivers, issue the following command:

```
modinfo <driver name>
```

2. Unload and remove the Linux drivers.

```
modprobe -r qedr
modprobe -r qede
modprobe -r qed
depmod -a
```

- 3. Remove the driver module files:
 - If the drivers were installed using an RPM package, issue the following command:

```
rpm -e qlgc-fastlinq-kmp-default-<version>.<arch>
```

☐ If the drivers were installed using a TAR file, issue the following commands for your operating system:

```
For RHEL:
```

```
cd /lib/modules/<version>/extra/qlgc-fastlinq
rm -rf qed.ko qede.ko qedr.ko
For SLES:
```

cd /lib/modules/<version>/updates/qlgc-fastlinq
rm -rf qed.ko qede.ko qedr.ko

Installing Linux Drivers Using the src RPM Package

To install Linux drivers using the src RPM package:

1. Issue the following at a command prompt:

```
rpm -ivh RPMS/<arch>/qlgc-fastlinq-<version>.src.rpm
```

2. Change the directory to the RPM path and build the binary RPM for the kernel.

NOTE

For RHEL 8, install the kernel-rpm-nacros and kernel-abi-whitelists packages before building the binary RPM package.

For RHEL:

cd /root/rpmbuild
rpmbuild -bb SPECS/fastling-<version>.spec

For SLES:

cd /usr/src/packages
rpmbuild -bb SPECS/fastling-<version>.spec

3. Install the newly compiled RPM:

rpm -ivh RPMS/<arch>/qlgc-fastlinq-<version>.<arch>.rpm

NOTE

The --force option may be needed on some Linux distributions if conflicts are reported.

The drivers will be installed in the following paths.

For SLES:

/lib/modules/<version>/updates/qlgc-fastlinq

For RHEL:

/lib/modules/<version>/extra/qlgc-fastling

4. Turn on all ethX interfaces as follows:

ifconfig <ethX> up

5. For SLES, use YaST to configure the Ethernet interfaces to automatically start at boot by setting a static IP address or enabling DHCP on the interface.

Installing Linux Drivers Using the kmp/kmod RPM Package

To install kmod RPM package:

1. Issue the following command at a command prompt:

```
rpm -ivh qlgc-fastlinq-<version>.<arch>.rpm
```

2. Reload the driver:

```
modprobe -r qede
modprobe qede
```

Installing Linux Drivers Using the TAR File

To install Linux drivers using the TAR file:

1. Create a directory and extract the TAR files to the directory:

```
tar xjvf fastlinq-<version>.tar.bz2
```

2. Change to the recently created directory, and then install the drivers:

```
cd fastlinq-<version>
make clean; make install
```

The qed and qede drivers will be installed in the following paths.

For SLES:

/lib/modules/<version>/updates/qlgc-fastling

For RHEL:

/lib/modules/<version>/extra/qlgc-fastling

3. Test the drivers by loading them (unload the existing drivers first, if necessary):

```
rmmod qede
rmmod qed
modprobe qed
modprobe qede
```

Installing the Linux Drivers with RDMA

For information on iWARP, see Chapter 8 iWARP Configuration.

To install Linux drivers in an inbox OFED environment:

- Download the 41xxx Series Adapter Linux drivers from the Dell: dell.support.com
- 2. Configure RoCE on the adapter, as described in "Configuring RoCE on the Adapter for Linux" on page 150.
- 3. Remove existing Linux drivers, as described in "Removing the Linux Drivers" on page 10.
- 4. Install the new Linux drivers using one of the following methods:
 - ☐ Installing Linux Drivers Using the kmp/kmod RPM Package
 - ☐ Installing Linux Drivers Using the TAR File

5. Install libqedr libraries to work with RDMA user space applications. The libqedr RPM is available only for inbox OFED. You must select which RDMA (RoCE, RoCEv2, or iWARP) is used in UEFI until concurrent RoCE+iWARP capability is supported in the firmware). None is enabled by default. Issue the following command:

```
rpm -ivh qlgc-libqedr-<version>.<arch>.rpm
```

6. To build and install the libqedr user space library, issue the following command:

```
'make libqedr_install'
```

7. Test the drivers by loading them as follows:

```
modprobe qedr
make install libeqdr
```

Linux Driver Optional Parameters

Table 3-2 describes the optional parameters for the gede driver.

Table 3-2. qede Driver Optional Parameters

Parameter	Description
debug	Controls driver verbosity level similar to ethtool -s <dev> msglvl.</dev>
int_mode	Controls interrupt mode other than MSI-X.
gro_enable	Enables or disables the hardware generic receive offload (GRO) feature. This feature is similar to the kernel's software GRO, but is only performed by the device hardware.
err_flags_override	A bitmap for disabling or forcing the actions taken in case of a hardware error: bit #31 – An enable bit for this bitmask bit #0 – Prevent hardware attentions from being reasserted bit #1 – Collect debug data bit #2 – Trigger a recovery process bit #3 – Call WARN to get a call trace of the flow that led to the error

Linux Driver Operation Defaults

Table 3-3 lists the ged and gede Linux driver operation defaults.

Table 3-3. Linux Driver Operation Defaults

Operation	qed Driver Default	qede Driver Default
Speed	Auto-negotiation with speed advertised	Auto-negotiation with speed advertised
MSI/MSI-X	Enabled	Enabled
Flow Control	_	Auto-negotiation with RX and TX advertised
MTU	_	1500 (range is 46–9600)
Rx Ring Size	_	1000
Tx Ring Size	_	4078 (range is 128–8191)
Coalesce Rx Microseconds	_	24 (range is 0–255)
Coalesce Tx Microseconds	_	48
TSO	_	Enabled

Linux Driver Messages

To set the Linux driver message detail level, issue one of the following commands:

- ethtool -s <interface> msglvl <value>
- modprobe qede debug=<value>

Where <value> represents bits 0–15, which are standard Linux networking values, and bits 16 and greater are driver-specific.

Statistics

To view detailed statistics and configuration information, use the ethtool utility. See the ethtool man page for more information.

Importing a Public Key for Secure Boot

Linux drivers require that you import and enroll the QLogic public key to load the drivers in a Secure Boot environment. Before you begin, ensure that your server supports Secure Boot. This section provides two methods for importing and enrolling the public key.

16

To import and enroll the QLogic public key:

- 1. Download the public key from the following Web page:
 - http://ldriver.glogic.com/Module-public-key/
- 2. To install the public key, issue the following command:

```
# mokutil --root-pw --import cert.der
```

Where the --root-pw option enables direct use of the root user.

- 3. Reboot the system.
- 4. Review the list of certificates that are prepared to be enrolled:
 - # mokutil --list-new
- 5. Reboot the system again.
- 6. When the shim launches MokManager, enter the root password to confirm the certificate importation to the Machine Owner Key (MOK) list.
- 7. To determine if the newly imported key was enrolled:
 - # mokutil --list-enrolled

To launch MOK manually and enroll the QLogic public key:

- 1. Issue the following command:
 - # reboot
- 2. In the **GRUB 2** menu, press the C key.
- 3. Issue the following commands:

```
chainloader $efibootdir/MokManager.efi
```

- boot
- 4. Select Enroll key from disk.
- 5. Navigate to the cert.der file and then press ENTER.
- 6. Follow the instructions to enroll the key. Generally this includes pressing the 0 (zero) key and then pressing the Y key to confirm.

NOTE

The firmware menu may provide more methods to add a new key to the Signature Database.

For additional information about Secure Boot, refer to the following Web page:

https://www.suse.com/documentation/sled-12/book_sle_admin/data/sec_uefi_secboot.html

Installing Windows Driver Software

For information on iWARP, see Chapter 8 iWARP Configuration.

- Installing the Windows Drivers
- Removing the Windows Drivers
- Managing Adapter Properties
- Setting Power Management Options
- Link Configuration in Windows

Installing the Windows Drivers

Install Windows driver software using the Dell Update Package (DUP):

- Running the DUP in the GUI
- DUP Installation Options
- DUP Installation Examples

Running the DUP in the GUI

To run the DUP in the GUI:

1. Double-click the icon representing the Dell Update Package file.

NOTE

The actual file name of the Dell Update Package varies.

2. In the Dell Update Package window (Figure 3-1), click Install.

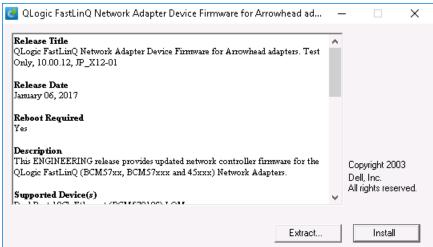


Figure 3-1. Dell Update Package Window

3. In the QLogic Super Installer—InstallShield® Wizard's Welcome window (Figure 3-2), click **Next**.

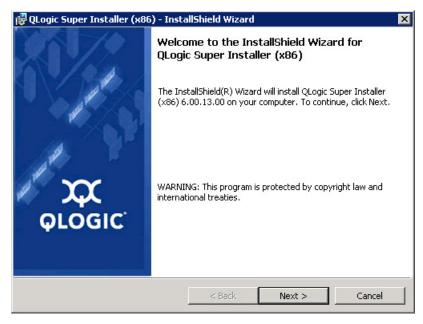


Figure 3-2. QLogic InstallShield Wizard: Welcome Window

- 4. Complete the following in the wizard's License Agreement window (Figure 3-3):
 - a. Read the End User Software License Agreement.
 - b. To continue, select I accept the terms in the license agreement.
 - c. Click Next.

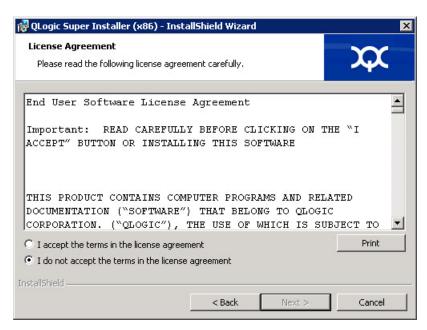


Figure 3-3. QLogic InstallShield Wizard: License Agreement Window

- 5. Complete the wizard's Setup Type window (Figure 3-4) as follows:
 - a. Select one of the following setup types:
 - Click Complete to install all program features.
 - Click Custom to manually select the features to be installed.
 - b. To continue, click Next.

If you clicked **Complete**, proceed directly to Step 6b.

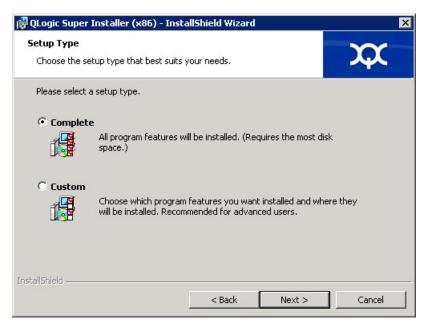


Figure 3-4. InstallShield Wizard: Setup Type Window

- 6. If you selected **Custom** in Step 5, complete the Custom Setup window (Figure 3-5) as follows:
 - a. Select the features to install. By default, all features are selected. To change a feature's install setting, click the icon next to it, and then select one of the following options:
 - This feature will be installed on the local hard drive—Marks the feature for installation without affecting any of its subfeatures.
 - This feature, and all subfeatures, will be installed on the local hard drive—Marks the feature and all of its subfeatures for installation.
 - This feature will not be available—Prevents the feature from being installed.
 - b. Click Next to continue.

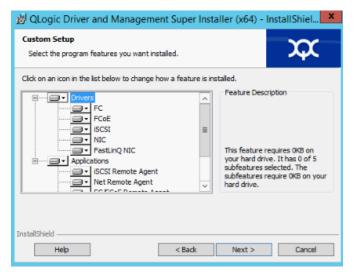


Figure 3-5. InstallShield Wizard: Custom Setup Window

7. In the InstallShield Wizard's Ready To Install window (Figure 3-6), click **Install**. The InstallShield Wizard installs the QLogic Adapter drivers and Management Software Installer.

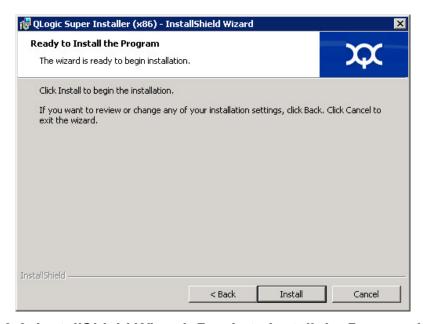


Figure 3-6. InstallShield Wizard: Ready to Install the Program Window

8. When the installation is complete, the InstallShield Wizard Completed window appears (Figure 3-7). Click **Finish** to dismiss the installer.

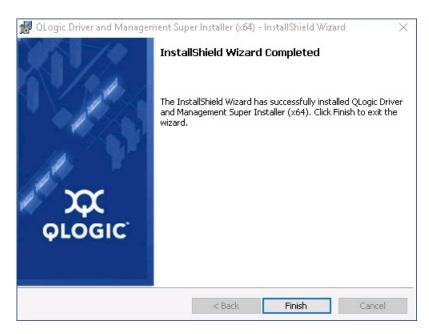


Figure 3-7. InstallShield Wizard: Completed Window

- 9. In the Dell Update Package window (Figure 3-8), "Update installer operation was successful" indicates completion.
 - ☐ (Optional) To open the log file, click **View Installation Log**. The log file shows the progress of the DUP installation, any previous installed versions, any error messages, and other information about the installation.
 - ☐ To close the Update Package window, click CLOSE.

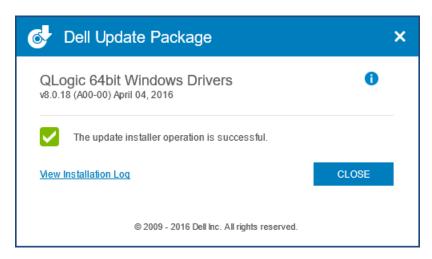


Figure 3-8. Dell Update Package Window

DUP Installation Options

To customize the DUP installation behavior, use the following command line options.

■ To extract only the driver components to a directory:

/drivers=<path>

NOTE

This command requires the /s option.

■ To install or update only the driver components:

/driveronly

NOTE

This command requires the /s option.

■ (Advanced) Use the /passthrough option to send all text following /passthrough directly to the QLogic installation software of the DUP. This mode suppresses any provided GUIs, but not necessarily those of the QLogic software.

/passthrough

(Advanced) To return a coded description of this DUP's supported features:

/capabilities

NOTE

This command requires the /s option.

DUP Installation Examples

The following examples show how to use the installation options.

To update the system silently:

```
<DUP file name>.exe /s
```

To extract the update contents to the C:\mydir\ directory:

```
<DUP file name>.exe /s /e=C:\mydir
```

To extract the driver components to the C:\mydir\ directory:

```
<DUP_file_name>.exe /s /drivers=C:\mydir
```

To install only the driver components:

```
<DUP file name>.exe /s /driveronly
```

To change from the default log location to C:\my path with spaces\log.txt:

<DUP file name>.exe /l="C:\my path with spaces\log.txt"

Removing the Windows Drivers

To remove the Windows drivers:

- 1. In the Control Panel, click **Programs**, and then click **Programs and Features**.
- 2. In the list of programs, select **QLogic FastLinQ Driver Installer**, and then click **Uninstall**.
- 3. Follow the instructions to remove the drivers.

Managing Adapter Properties

To view or change the 41xxx Series Adapter properties:

- 1. In the Control Panel, click **Device Manager**.
- 2. On the properties of the selected adapter, click the **Advanced** tab.
- 3. On the Advanced page (Figure 3-9), select an item under **Property** and then change the **Value** for that item as needed.

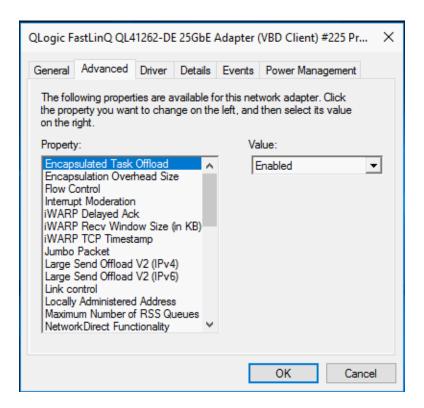


Figure 3-9. Setting Advanced Adapter Properties

Setting Power Management Options

You can set power management options to allow the operating system to turn off the controller to save power or to allow the controller to wake up the computer. If the device is busy (servicing a call, for example), the operating system will not shut down the device. The operating system attempts to shut down every possible device only when the computer attempts to go into hibernation. To have the controller remain on at all times, do not select the **Allow the computer to turn off the device to save power** check box (Figure 3-10).

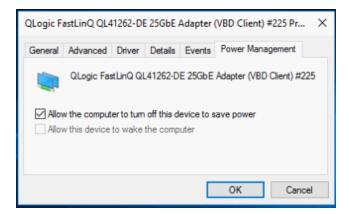


Figure 3-10. Power Management Options

NOTE

- The Power Management page is available only for servers that support power management.
- Do not select the Allow the computer to turn off the device to save power check box for any adapter that is a member of a team.

Configuring the Communication Protocol to Use with QCC GUI, QCC PowerKit, and QCS CLI

There are two main components of the QCC GUI, QCC PowerKit, and QCS CLI management applications: the RPC agent and the client software. An RPC agent is installed on a server, or managed host, that contains one or more Converged Network Adapters. The RPC agent collects information on the Converged Network Adapters and makes it available for retrieval from a management PC on which the client software is installed. The client software enables viewing information from the RPC agent and configuring the Converged Network Adapters. The management software includes QCC GUI and QCS CLI.

A communication protocol enables communication between the RPC agent and the client software. Depending on the mix of operating systems (Linux, Windows, or both) on the clients and managed hosts in your network, you can choose an appropriate utility.

For installation instructions for these management applications, refer to the following documents on the Marvell Web site:

- User's Guide, QLogic Control Suite CLI (part number BC0054511-00)
- User's Guide, PowerShell (part number BC0054518-00)
- Installation Guide, QConvergeConsole GUI (part number SN0051105-00)

Link Configuration in Windows

Link configuration can be done in the Windows OS with three different parameters, which are available for configuration on the Advanced tab of the Device Manager page.

Link Control Mode

There are two modes for controlling link configuration:

- **Preboot Controlled** is the default mode. In this mode, the driver uses the link configuration from the device, which is configurable from preboot components. This mode ignores the link parameters on the Advanced tab.
- **Driver Controlled** mode should be set when you want to configure the link settings from Advanced tab of the Device Manager (as shown in Figure 3-11).

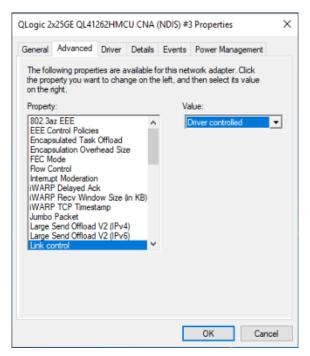


Figure 3-11. Setting Driver Controlled Mode

Link Speed and Duplex

The Speed & Duplex property (on the Advanced tab of the Device Manager) can be configured to any selection in the Value menu (see Figure 3-12).

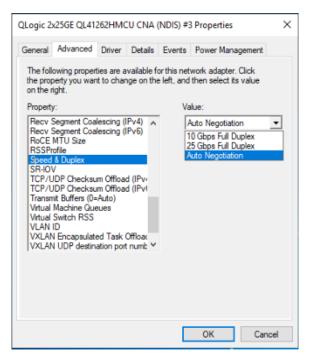


Figure 3-12. Setting the Link Speed and Duplex Property

This configuration is effective only when the link control property is set to Driver controlled (see Figure 3-11).

FEC Mode

FEC mode configuration at the OS level involves three driver advanced properties.

To set FEC mode:

- 1. Set Link Control. On the Advanced tab of the Device Manager:
 - a. In the Property menu, select **Link control**.
 - b. In the Value menu, select **Driver controlled**.

See Figure 3-11 for an example.

- 2. Set Speed & Duplex. On the Advanced tab of the Device Manager:
 - a. In the Property menu, select **Speed & Duplex**.
 - b. In the Value menu, select a fixed speed.

FEC mode configuration is active only when Speed & Duplex is set to a fixed speed. Setting this property to Auto Negotiation disables FEC configuration.

- 3. Set FEC Mode. On the Advanced tab of the Device Manager:
 - a. In the Property menu, select **FEC Mode**.
 - b. In the Value menu, select a valid value (see Figure 3-13).

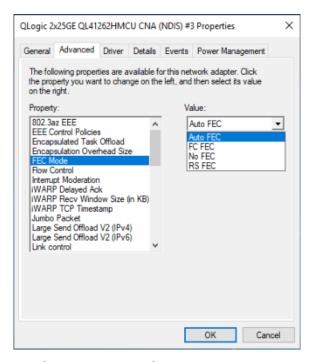


Figure 3-13. Setting the FEC Mode Property

This property is in effect only when Step 1 and Step 2 have been completed.

All FEC modes are not valid for each media; you must know the valid modes for your specific media. If the wrong FEC mode value is set, the link goes down.

Installing VMware Driver Software

This section describes the qedenty VMware ESXi driver for the 41xxx Series Adapters:

- VMware Drivers and Driver Packages
- Installing VMware Drivers
- VMware NIC Driver Optional Parameters

- VMware Driver Parameter Defaults
- Removing the VMware Driver
- FCoE Support
- iSCSI Support

VMware Drivers and Driver Packages

Table 3-4 lists the VMware ESXi drivers for the protocols.

Table 3-4. VMware Drivers

VMware Driver	Description
qedentv	Native networking driver
qedrntv	Native RDMA-Offload (RoCE and RoCEv2) driver ^a
qedf	Native FCoE-Offload driver
qedil	Legacy iSCSI-Offload driver
qedi	Native iSCSI-Offload driver (ESXi 6.7 and later) ^b

^a For ESXi 6.5, the NIC and RoCE drivers have been packaged together and can be installed as a single offline zip bundle using standard ESXi installation commands. The recommended installation sequence is the NIC/RoCE driver package, followed by the FCoE and iSCSI driver packages (as required).

The ESXi drivers are included as individual driver packages and are not bundled together, except as noted.

The VMware drivers are available for download only from the VMware web site:

https://www.vmware.com/resources/compatibility/search.php?deviceCategory=io &details=1&keyword=QL41&page=1&display_interval=10&sortColumn=Partner&sortOrder=Asc

Install individual drivers using either:

- Standard ESXi package installation commands (see Installing VMware Drivers)
- Procedures in the individual driver Read Me files

^b For ESXi 6.7, the NIC, RoCE, and iSCSI drivers have been packaged together and can be installed as a single offline zip bundle using standard ESXi installation commands. The recommended installation sequence is the NIC/RoCE/iSCSI driver package, followed by the FCoE driver package (as required).

Procedures in the following VMware KB article:

https://kb.vmware.com/selfservice/microsites/search.do?language=en_US&cmd=displayKC&externalId=2137853

You should install the NIC driver first, followed by the storage drivers.

Installing VMware Drivers

You can use the driver ZIP file to install a new driver or update an existing driver. Be sure to install the entire driver set from the same driver ZIP file. Mixing drivers from different ZIP files will cause problems.

To install the VMware driver:

1. Download the VMware driver for the 41xxx Series Adapter from the VMware support page:

www.vmware.com/support.html

- 2. Power up the ESX host, and then log into an account with administrator authority.
- 3. Use the Linux scp utility to copy the driver bundle from a local system into the /tmp directory on an ESX server with IP address 10.10.10.10. For example, issue the following command:
 - # scp qedentv-bundle-2.0.3.zip root@10.10.10.10:/tmp

You can place the file anywhere that is accessible to the ESX console shell.

- 4. Place the host in maintenance mode by issuing the following command:
 - # esxcli --maintenance-mode

NOTE

The maximum number of supported qedenty Ethernet interfaces on an ESXi host is 32 because the vmkernel allows only 32 interfaces to register for management callback.

- 5. Select one of the following installation options:
 - Option 1: Install the driver bundle (which will install all of the driver VIBs at one time) by issuing the following command:
 - # esxcli software vib install -d /tmp/qedentv-2.0.3.zip
 - Option 2: Install the .vib directly on an ESX server using either the CLI or the VMware Update Manager (VUM). To do this, unzip the driver ZIP file, and then extract the .vib file.
 - To install the .vib file using the CLI, issue the following command. Be sure to specify the full .vib file path:
- # esxcli software vib install -v /tmp/qedentv-1.0.3.11-10EM.550.0.0.1331820.x86 64.vib
 - To install the .vib file using the VUM, see the knowledge base article here:

Updating an ESXi/ESX host using VMware vCenter Update Manager 4.x and 5.x (1019545)

To upgrade the existing driver bundle:

- Issue the following command:
 - # esxcli software vib update -d /tmp/qedentv-bundle-2.0.3.zip

To upgrade an individual driver:

Follow the steps for a new installation (see **To install the VMware driver**), except replace the command in Option 1 with the following:

esxcli software vib update -v /tmp/qedentv-1.0.3.11-10EM.550.0.0.1331820.x86 64.vib

VMware NIC Driver Optional Parameters

Table 3-5 describes the optional parameters that can be supplied as command line arguments to the <code>esxcfg-module</code> command.

Table 3-5. VMware NIC Driver Optional Parameters

Parameter	Description
hw_vlan	Globally enables (1) or disables (0) hardware vLAN insertion and removal. Disable this parameter when the upper layer needs to send or receive fully formed packets. hw_vlan=1 is the default.

Table 3-5. VMware NIC Driver Optional Parameters (Continued)

Parameter	Description
num_queues	Specifies the number of TX/RX queue pairs. num_queues can be 1-11 or one of the following: -1 allows the driver to determine the optimal number of queue pairs (default). 0 uses the default queue. You can specify multiple values delimited by commas for multiport or multi-
multi_rx_filters	function configurations. Specifies the number of RX filters per RX queue, excluding the default queue. multi_rx_filters can be 1-4 or one of the following values: -1 uses the default number of RX filters per queue. 0 disables RX filters.
disable_tpa	Enables (0) or disables (1) the TPA (LRO) feature. disable_tpa=0 is the default.
max_vfs	Specifies the number of virtual functions (VFs) per physical function (PF). max_vfs can be 0 (disabled) or 64 VFs on a single port (enabled). The 64 VF maximum support for ESXi is an OS resource allocation constraint.
RSS	Specifies the number of receive side scaling queues used by the host or virtual extensible LAN (VXLAN) tunneled traffic for a PF. RSS can be 2, 3, 4, or one of the following values: ■ -1 uses the default number of queues. ■ 0 or 1 disables RSS queues. You can specify multiple values delimited by commas for multiport or multifunction configurations.
debug	Specifies the level of data that the driver records in the vmkernel log file. debug can have the following values, shown in increasing amounts of data: □ 0x80000000 indicates Notice level. □ 0x40000000 indicates Information level (includes the Notice level). □ 0x3FFFFFFF indicates Verbose level for all driver submodules (includes the Information and Notice levels).
auto_fw_reset	Enables (1) or disables (0) the driver automatic firmware recovery capability. When this parameter is enabled, the driver attempts to recover from events such as transmit timeouts, firmware asserts, and adapter parity errors. The default is auto_fw_reset=1 .

Table 3-5. VMware NIC Driver Optional Parameters (Continued)

Parameter	Description
vxlan_filter_en	Enables (1) or disables (0) the VXLAN filtering based on the outer MAC, the inner MAC, and the VXLAN network (VNI), directly matching traffic to a specific queue. The default is <pre>vxlan_filter_en=1</pre> . You can specify multiple values delimited by commas for multiport or multifunction configurations.
enable_vxlan_offld	Enables (1) or disables (0) the VXLAN tunneled traffic checksum offload and TCP segmentation offload (TSO) capability. The default is enable_vxlan_offld=1. You can specify multiple values delimited by commas for multiport or multifunction configurations.

VMware Driver Parameter Defaults

Table 3-6 lists the VMware driver parameter default values.

Table 3-6. VMware Driver Parameter Defaults

Parameter	Default
Speed	Autonegotiation with all speeds advertised. The speed parameter must be the same on all ports. If autonegotiation is enabled on the device, all of the device ports will use autonegotiation.
Flow Control	Autonegotiation with RX and TX advertised
MTU	1,500 (range 46-9,600)
Rx Ring Size	8,192 (range 128-8,192)
Tx Ring Size	8,192 (range 128-8,192)
MSI-X	Enabled
Transmit Send Offload (TSO)	Enabled
Large Receive Offload (LRO)	Enabled
RSS	Enabled (four RX queues)
HW VLAN	Enabled
Number of Queues	Enabled (eight RX/TX queue pairs)
Wake on LAN (WoL)	Disabled

Removing the VMware Driver

To remove the .vib file (gedenty), issue the following command:

esxcli software vib remove --vibname qedentv

To remove the driver, issue the following command:

vmkload_mod -u qedentv

FCoE Support

The Marvell VMware FCoE qedf driver included in the VMware software package supports Marvell FastLinQ FCoE converged network interface controllers (C-NICs). The driver is a kernel-mode driver that provides a translation layer between the VMware SCSI stack and the Marvell FCoE firmware and hardware. The FCoE and DCB feature set is supported on VMware ESXi 5.0 and later.

To enable FCoE-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

iSCSI Support

The Marvell VMware iSCSI qedil Host Bus Adapter (HBA) driver, similar to qedf, is a kernel mode driver that provides a translation layer between the VMware SCSI stack and the Marvell iSCSI firmware and hardware. The qedil driver leverages the services provided by the VMware iscsid infrastructure for session management and IP services.

To enable iSCSI-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

NOTE

The iSCSI interface supported by the QL41xxx Adapters is a dependent hardware interface that relies on networking services, iSCSI configuration, and management interfaces provided by VMware. The iSCSI interface includes two components: a network adapter and an iSCSI engine on the same interface. The iSCSI engine appears on the list of storage adapters as an iSCSI adapter (vmhba). For services such as ARP and DHCP needed by iSCSI, the iSCSI vmhba uses the services of the vmnic device created by the qedil driver. The vmnic is a thin dummy implementation intended to provide L2 functionality for iSCSI to operate. Do not configure, assign to a virtual switch, or use the vmnic in any manner for carrying regular networking traffic. The actual NIC interfaces on the adapter will be claimed by the gedenty driver, which is a fully-functional NIC driver.

4 Upgrading the Firmware

This chapter provides information about upgrading the firmware using the Dell Update Package (DUP).

The firmware DUP is a Flash update utility only; it is not used for adapter configuration. You can run the firmware DUP by double-clicking the executable file. Alternatively, you can run the firmware DUP from the command line with several supported command line options.

- Running the DUP by Double-Clicking
- "Running the DUP from a Command Line" on page 41
- "Running the DUP Using the .bin File" on page 42 (Linux only)

Running the DUP by Double-Clicking

To run the firmware DUP by double-clicking the executable file:

- Double-click the icon representing the firmware Dell Update Package file.
- 2. The Dell Update Package splash screen appears, as shown in Figure 4-1. Click **Install** to continue.

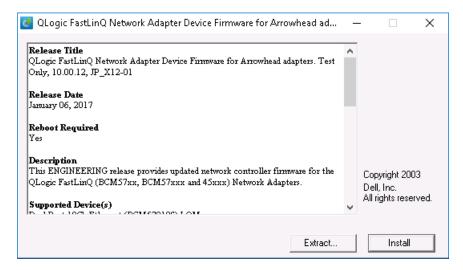


Figure 4-1. Dell Update Package: Splash Screen

3. Follow the on-screen instructions. In the Warning dialog box, click **Yes** to continue the installation.

The installer indicates that it is loading the new firmware, as shown in Figure 4-2.

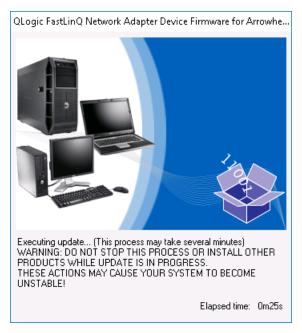


Figure 4-2. Dell Update Package: Loading New Firmware

When complete, the installer indicates the result of the installation, as shown in Figure 4-3.

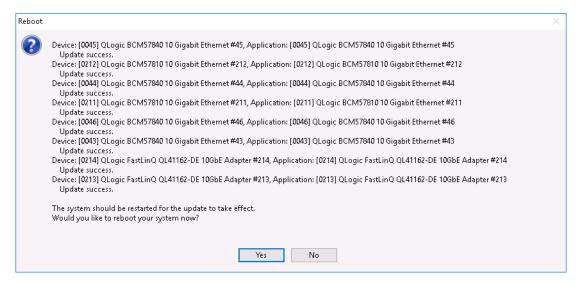


Figure 4-3. Dell Update Package: Installation Results

- 4. Click **Yes** to reboot the system.
- 5. Click **Finish** to complete the installation, as shown in Figure 4-4.

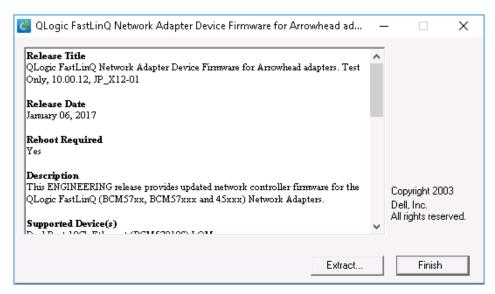


Figure 4-4. Dell Update Package: Finish Installation

Running the DUP from a Command Line

Running the firmware DUP from the command line, with no options specified, results in the same behavior as double-clicking the DUP icon. Note that the actual file name of the DUP will vary.

To run the firmware DUP from a command line:

Issue the following command:

C:\> Network_Firmware_2T12N_WN32_<version>_X16.EXE

Figure 4-5 shows the options that you can use to customize the Dell Update Package installation.

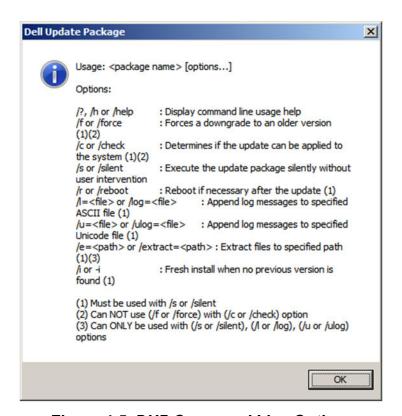


Figure 4-5. DUP Command Line Options

Running the DUP Using the .bin File

The following procedure is supported only on Linux OS.

To update the DUP using the .bin file:

- 1. Copy the Network_Firmware_NJCX1_LN_X.Y.Z.BIN file to the system or server.
- 2. Change the file type into an executable file as follows:

```
chmod 777 Network Firmware NJCX1 LN X.Y.Z.BIN
```

3. To start the update process, issue the following command:

```
./Network Firmware NJCX1 LN X.Y.Z.BIN
```

4. After the firmware is updated, reboot the system.

Example output from the SUT during the DUP update:

```
./Network Firmware NJCX1 LN 08.07.26.BIN
Collecting inventory...
Running validation...
BCM57810 10 Gigabit Ethernet rev 10 (p2p1)
The version of this Update Package is the same as the currently installed
version.
Software application name: BCM57810 10 Gigabit Ethernet rev 10 (p2p1)
Package version: 08.07.26
Installed version: 08.07.26
BCM57810 10 Gigabit Ethernet rev 10 (p2p2)
The version of this Update Package is the same as the currently installed
version.
Software application name: BCM57810 10 Gigabit Ethernet rev 10 (p2p2)
Package version: 08.07.26
Installed version: 08.07.26
Continue? Y/N:Y
Y entered; update was forced by user
Executing update...
WARNING: DO NOT STOP THIS PROCESS OR INSTALL OTHER DELL PRODUCTS WHILE UPDATE
IS IN PROGRESS.
THESE ACTIONS MAY CAUSE YOUR SYSTEM TO BECOME UNSTABLE!
Device: BCM57810 10 Gigabit Ethernet rev 10 (p2p1)
 Application: BCM57810 10 Gigabit Ethernet rev 10 (p2p1)
 Update success.
Device: BCM57810 10 Gigabit Ethernet rev 10 (p2p2)
 Application: BCM57810 10 Gigabit Ethernet rev 10 (p2p2)
 Update success.
Would you like to reboot your system now?
Continue? Y/N:Y
```

43

5 Adapter Preboot Configuration

During the host boot process, you have the opportunity to pause and perform adapter management tasks using the Human Infrastructure Interface (HII) application. These tasks include the following:

- "Getting Started" on page 45
- "Displaying Firmware Image Properties" on page 48
- "Configuring Device-level Parameters" on page 49
- "Configuring NIC Parameters" on page 50
- "Configuring Data Center Bridging" on page 53
- "Configuring FCoE Boot" on page 55
- "Configuring iSCSI Boot" on page 56
- "Configuring Partitions" on page 61

NOTE

The HII screen shots in this chapter are representative and may not match the screens that you see on your system.

Getting Started

To start the HII application:

- 1. Open the System Setup window for your platform. For information about launching the System Setup, consult the user guide for your system.
- 2. In the System Setup window (Figure 5-1), select **Device Settings**, and then press ENTER.



Figure 5-1. System Setup

3. In the Device Settings window (Figure 5-2), select the 41xxx Series Adapter port that you want to configure, and then press ENTER.

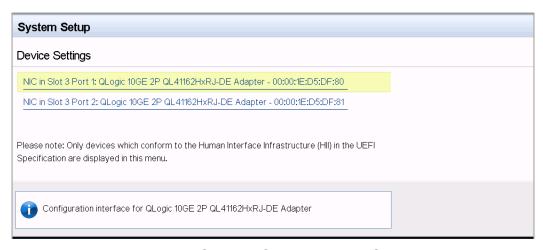


Figure 5-2. System Setup: Device Settings

Main Configuration Page Firmware Image Properties Device Level Configuration NIC Configuration Data Center Bridging (DCB) Settings QLogic 10GE 2P QL41162HxRJ-DE Adapter Chip Type -BCM57940S A2 PCI Device ID -----8070 86:00 PCI Address Blink LEDs -----0 Link Status ----Connected

The Main Configuration Page (Figure 5-3) presents the adapter management options where you can set the partitioning mode.

Figure 5-3. Main Configuration Page

00:00:00:00:00

4. Under **Device Level Configuration**, set the **Partitioning Mode** to **NPAR** to add the **NIC Partitioning Configuration** option to the Main Configuration Page, as shown in Figure 5-4.

----- 00:0E:1E:D5:F8:76



MAC Address ----

Virtual MAC Address

NPAR is not available on ports with a maximum speed of 1G.

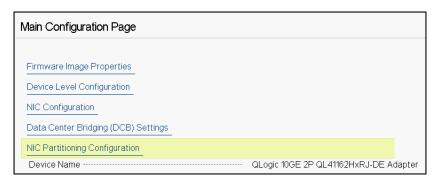


Figure 5-4. Main Configuration Page, Setting Partitioning Mode to NPAR

In Figure 5-3 and Figure 5-4, the Main Configuration Page shows the following:

■ Firmware Image Properties (see "Displaying Firmware Image Properties" on page 48)

- **Device Level Configuration** (see "Configuring Device-level Parameters" on page 49)
- NIC Configuration (see "Configuring NIC Parameters" on page 50)
- **iSCSI Configuration** (if iSCSI remote boot is allowed by enabling iSCSI offload in NPAR mode on the port's third partition) (see "Configuring iSCSI Boot" on page 56)
- FCoE Configuration (if FCoE boot from SAN is allowed by enabling FCoE offload in NPAR mode on the port's second partition) (see "Configuring FCoE Boot" on page 55)
- Data Center Bridging (DCB) Settings (see "Configuring Data Center Bridging" on page 53)
- NIC Partitioning Configuration (if NPAR is selected on the Device Level Configuration page) (see "Configuring Partitions" on page 61)

In addition, the Main Configuration Page presents the adapter properties listed in Table 5-1.

Table 5-1. Adapter Properties

Adapter Property	Description
Device Name	Factory-assigned device name
Chip Type	ASIC version
PCI Device ID	Unique vendor-specific PCI device ID
PCI Address	PCI device address in bus-device function format
Blink LEDs	User-defined blink count for the port LED
Link Status	External link status
MAC Address	Manufacturer-assigned permanent device MAC address
Virtual MAC Address	User-defined device MAC address
iSCSI MAC Address ^a	Manufacturer-assigned permanent device iSCSI Offload MAC address
iSCSI Virtual MAC Address ^a	User-defined device iSCSI Offload MAC address
FCoE MAC Address ^b	Manufacturer-assigned permanent device FCoE Offload MAC address
FCoE Virtual MAC Address ^b	User-defined device FCoE Offload MAC address

Adapter Property	Description
FCoE WWPN b	Manufacturer-assigned permanent device FCoE Offload WWPN (world wide port name)
FCoE Virtual WWPN b	User-defined device FCoE Offload WWPN
FCoE WWNN b	Manufacturer-assigned permanent device FCoE Offload WWNN (world wide node name)
FCoE Virtual WWNN b	User-defined device FCoE Offload WWNN

^a This property is visible only if **iSCSI Offload** is enabled on the NIC Partitioning Configuration page.

Displaying Firmware Image Properties

To view the properties for the firmware image, select **Firmware Image Properties** on the Main Configuration Page, and then press ENTER. The Firmware Image Properties page (Figure 5-5) specifies the following view-only data:

- Family Firmware Version is the multiboot image version, which comprises several firmware component images.
- **MBI Version** is the Marvell FastLinQ bundle image version that is active on the device.
- Controller BIOS Version is the management firmware version.
- **EFI Driver Version** is the extensible firmware interface (EFI) driver version.
- **L2B Firmware Version** is the NIC offload firmware version for boot.



Figure 5-5. Firmware Image Properties

^b This property is visible only if **FCoE Offload** is enabled on the NIC Partitioning Configuration page.

Configuring Device-level Parameters

NOTE

The iSCSI physical functions (PFs) are listed when the iSCSI Offload feature is enabled in NPAR mode only. The FCoE PFs are listed when the FCoE Offload feature is enabled in NPAR mode only. Not all adapter models support iSCSI Offload and FCoE Offload. Only one offload can be enabled per port, and only in NPAR mode.

Device-level configuration includes the following parameters:

- Virtualization Mode
- NPAREP Mode

To configure device-level parameters:

- 1. On the Main Configuration Page, select **Device Level Configuration** (see Figure 5-3 on page 46), and then press ENTER.
- 2. On the **Device Level Configuration** page, select values for the device-level parameters, as shown in Figure 5-6.

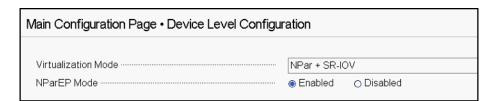


Figure 5-6. Device Level Configuration

NOTE

QL41264HMCU-DE (part number 5V6Y4) and QL41264HMRJ-DE (part number 0D1WT) adapters show support for NPAR, SR-IOV and NPAR-EP in the Device Level Configuration, though these features are not supported on 1Gbps ports 3 and 4.

- 3. For **Virtualization Mode**, select one of the following modes to apply to all adapter ports:
 - **None** (default) specifies that no virtualization mode is enabled.
 - NPAR sets the adapter to switch-independent NIC partitioning mode.
 - □ **SR-IOV** sets the adapter to SR-IOV mode.
 - □ NPar + SR-IOV sets the adapter to SR-IOV over NPAR mode.

- 4. **NParEP Mode** configures the maximum quantity of partitions per adapter. This parameter is visible when you select either **NPAR** or **NPar + SR-IOV** as the **Virtualization Mode** in **Step 2**.
 - ☐ **Enabled** allows you to configure up to 16 partitions per adapter.
 - ☐ **Disabled** allows you to configures up to 8 partitions per adapter.
- 5. Click Back.
- 6. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

Configuring NIC Parameters

NIC configuration includes setting the following parameters:

- Link Speed
- NIC + RDMA Mode
- RDMA Protocol Support
- Boot Mode
- FEC Mode
- **■** Energy Efficient Ethernet
- Virtual LAN Mode
- Virtual LAN ID

To configure NIC parameters:

1. On the Main Configuration Page, select **NIC Configuration** (Figure 5-3 on page 46), and then click **Finish**.

Figure 5-7 shows the NIC Configuration page.

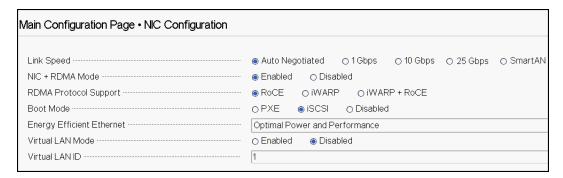


Figure 5-7. NIC Configuration

2.		ct one of the following Link Speed options for the selected port. Not all d selections are available on all adapters.
		Auto Negotiated enables Auto Negotiation mode on the port. FEC mode selection is not available for this speed mode.
		1 Gbps enables 1GbE fixed speed mode on the port. This mode is intended only for 1GbE interfaces and should not be configured for adapter interfaces that operate at other speeds. FEC mode selection is not available for this speed mode. This mode is not available on all adapters.
		10 Gbps enables 10GbE fixed speed mode on the port. This mode is not available on all adapters.
		25 Gbps enables 25GbE fixed speed mode on the port. This mode is not available on all adapters.
		SmartAN (Default) enables FastLinQ SmartAN [™] link speed mode on the port. No FEC mode selection is available for this speed mode. The SmartAN setting cycles through all possible link speeds and FEC modes until a link is established. This mode is intended for use only with 25G interfaces. This mode is not available on all adapters.
3.		NIC + RDMA Mode, select either Enabled or Disabled for RDMA on ort. This setting applies to all partitions of the port, if in NPAR mode.
4.	Link	Mode is visible when 25 Gbps fixed speed mode is selected as the Speed in Step 2 . For FEC Mode , select one of the following options. all FEC modes are available on all adapters.
		None disables all FEC modes.
		Fire Code enables Fire Code (BASE-R) FEC mode.
		Reed Solomon enables Reed Solomon FEC mode.
		Auto enables the port to cycle through None , Fire Code , and Reed Solomon FEC modes (at that link speed) in a round-robin fashion, until a link is established.
5.	NPA	RDMA Protocol Support setting applies to all partitions of the port, if in R mode. This setting appears if the NIC + RDMA Mode in Step 3 is set nabled. RDMA Protocol Support options include the following:
		RoCE enables RoCE mode on this port.
		iWARP enables iWARP mode on this port.
		iWARP + RoCE enables iWARP and RoCE modes on this port. This is the default. Additional configuration for Linux is required for this option

6.	For Boot Mode , select one of the following values:			
		PXE enables PXE boot.		
		FCoE enables FCoE boot from SAN over the hardware offload pathway. The FCoE mode is available only if FCoE Offload is enabled on the second partition in NPAR mode (see "Configuring Partitions" on page 61).		
		iSCSI enables iSCSI remote boot over the hardware offload pathway. The iSCSI mode is available only if iSCSI Offload is enabled on the third partition in NPAR mode (see "Configuring Partitions" on page 61)		
		Disabled prevents this port from being used as a remote boot source.		
7.	100E	The Energy Efficient Ethernet (EEE) parameter is visible only on 100BASE-T or 10GBASE-T RJ45 interfaced adapters. Select from the following EEE options:		
		☐ Disabled disables EEE on this port.		
		Optimal Power and Performance enables EEE in optimal power and performance mode on this port.		
		Maximum Power Savings enables EEE in maximum power savings mode on this port.		
		Maximum Performance enables EEE in maximum performance mode on this port.		
8.	The Virtual LAN Mode parameter applies to the entire port when in PXE remote install mode. It is not persistent after a PXE remote install finishes. Select from the following vLAN options:			
		Enabled enables vLAN mode on this port for PXE remote install mode		
		Disabled disables vLAN mode on this port.		
9.	The Virtual LAN ID parameter specifies the vLAN tag ID to be used on this port for PXE remote install mode. This setting applies only when Virtual LAN Mode is enabled in the previous step.			
10.	Click Back.			
11.	When prompted, click Yes to save the changes. Changes take effect after system reset.			

To configure the port to use RDMA:

NOTE

Follow these steps to enable RDMA on all partitions of an NPAR mode port.

- 1. Set NIC + RDMA Mode to Enabled.
- Click Back.
- 3. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

To configure the port's boot mode:

- 1. For a UEFI PXE remote installation, select **PXE** as the **Boot Mode**.
- Click Back.
- 3. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

To configure the port's PXE remote install to use a vLAN:

NOTE

This vLAN is not persistent after the PXE remote install is finished.

- 1. Set the **Virtual LAN Mode** to **Enabled**.
- 2. In the **Virtual LAN ID** box, enter the number to be used.
- 3. Click Back.
- 4. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

Configuring Data Center Bridging

The data center bridging (DCB) settings comprise the DCBX protocol and the RoCE priority.

To configure the DCB settings:

- On the Main Configuration Page (Figure 5-3 on page 46), select Data Center Bridging (DCB) Settings, and then click Finish.
- 2. On the Data Center Bridging (DCB) Settings page (Figure 5-8), select the appropriate **DCBX Protocol** option:
 - Disabled disables DCBX on this port.

- □ CEE enables the legacy Converged Enhanced Ethernet (CEE) protocol DCBX mode on this port.
- ☐ IEEE enables the IEEE DCBX protocol on this port.
- **Dynamic** enables dynamic application of either the CEE or IEEE protocol to match the attached link partner.
- 3. On the Data Center Bridging (DCB) Settings page, enter the **RoCE v1 Priority** as a value from **0–7**. This setting indicates the DCB traffic class priority number used for RoCE traffic and should match the number used by the DCB-enabled switching network for RoCE traffic. Typically, 0 is used for the default lossy traffic class, 3 is used for the FCoE traffic class, and 4 is used for the lossless iSCSI-TLV over DCB traffic class.

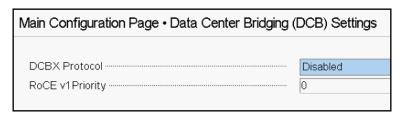


Figure 5-8. System Setup: Data Center Bridging (DCB) Settings

- 4. Click Back.
- 5. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

NOTE

When DCBX is enabled, the adapter periodically sends link layer discovery protocol (LLDP) packets with a dedicated unicast address that serves as the source MAC address. This LLDP MAC address is different from the factory-assigned adapter Ethernet MAC address. If you examine the MAC address table for the switch port that is connected to the adapter, you will see two MAC addresses: one for LLDP packets and one for the adapter Ethernet interface.

Configuring FCoE Boot

NOTE

The FCoE Boot Configuration Menu is only visible if **FCoE Offload Mode** is enabled on the second partition in NPAR mode (see Figure 5-18 on page 64). It is not visible in non-NPAR mode.

To enable FCoE-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

To configure the FCoE boot configuration parameters:

- 1. On the Main Configuration Page, select **FCoE Configuration**, and then select the following as needed:
 - **☐** FCoE General Parameters (Figure 5-9)
 - **☐** FCoE Target Configuration (Figure 5-10)
- 2. Press ENTER.
- 3. Choose values for the FCoE General or FCoE Target Configuration parameters.

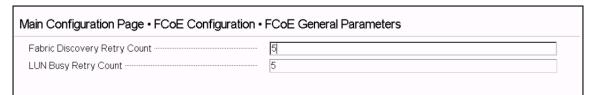


Figure 5-9. FCoE General Parameters

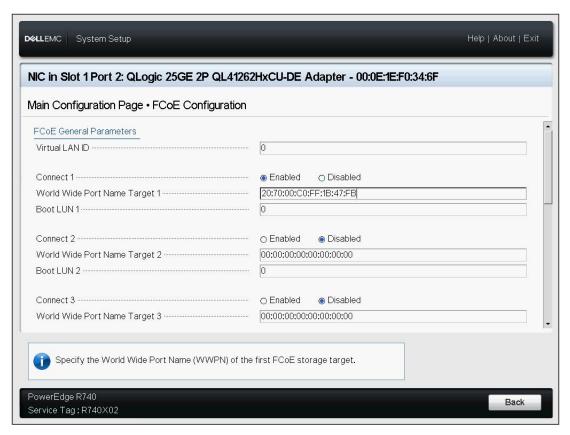


Figure 5-10. FCoE Target Configuration

- 4. Click Back.
- 5. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

Configuring iSCSI Boot

NOTE

The iSCSI Boot Configuration Menu is only visible if **iSCSI Offload Mode** is enabled on the third partition in NPAR mode (see Figure 5-19 on page 65). It is not visible in non-NPAR mode.

To enable FCoE-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

To configure the iSCSI boot configuration parameters:

1.		he Main Configuration Page, select iSCSI Boot Configuration Menu then select one of the following options:
		iSCSI General Configuration iSCSI Initiator Configuration iSCSI First Target Configuration iSCSI Second Target Configuration
2.	Pres	ss ENTER.
3.	Cho	ose values for the appropriate iSCSI configuration parameters:
		iSCSI General Parameters (Figure 5-11 on page 58)
		 TCP/IP Parameters Via DHCP iSCSI Parameters Via DHCP CHAP Authentication CHAP Mutual Authentication IP Version ARP Redirect DHCP Request Timeout Target Login Timeout DHCP Vendor ID
		iSCSI Initiator Parameters (Figure 5-12 on page 59)
		 IPv4 Address IPv4 Subnet Mask IPv4 Default Gateway IPv4 Primary DNS IPv4 Secondary DNS VLAN ID iSCSI Name CHAP ID CHAP Secret
		iSCSI First Target Parameters (Figure 5-13 on page 59)
		 Connect IPv4 Address TCP Port Boot LUN iSCSI Name CHAP ID

CHAP Secret

- □ iSCSI Second Target Parameters (Figure 5-14 on page 60)
 - Connect
 - IPv4 Address
 - TCP Port
 - Boot LUN
 - iSCSI Name
 - CHAP ID
 - CHAP Secret
- 4. Click Back.
- 5. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

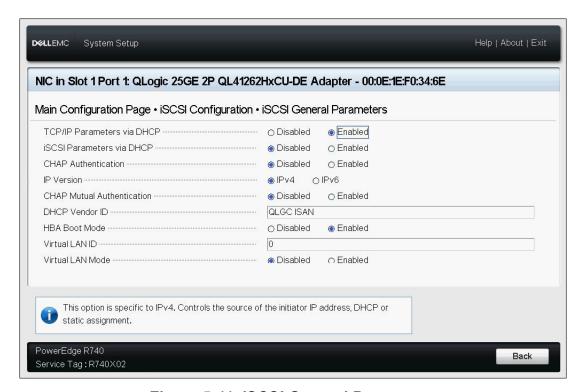


Figure 5-11. iSCSI General Parameters

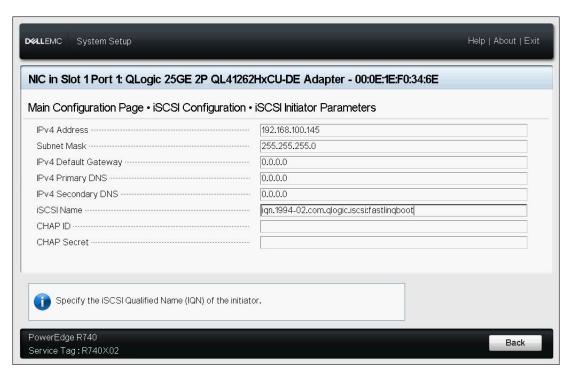


Figure 5-12. iSCSI Initiator Configuration Parameters

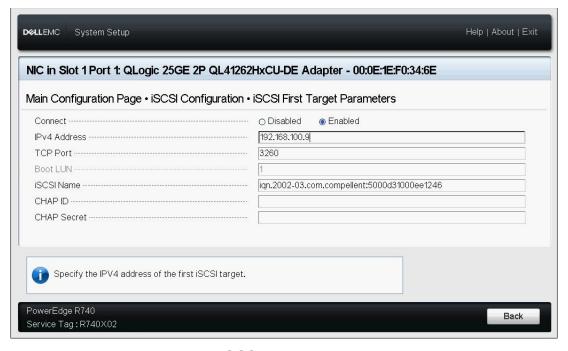


Figure 5-13. iSCSI First Target Parameters

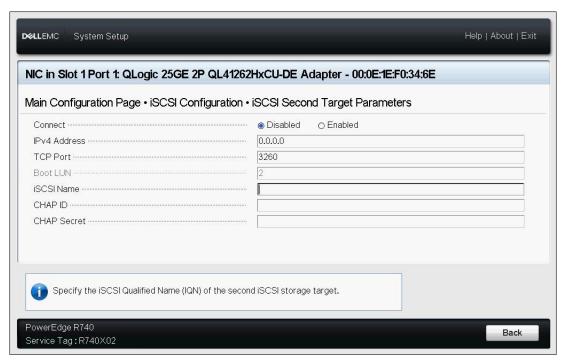


Figure 5-14. iSCSI Second Target Parameters

Configuring Partitions

You can configure bandwidth ranges for each partition on the adapter. For information specific to partition configuration on VMware ESXi 6.5, see Partitioning for VMware ESXi 6.5 and ESXi 6.7.

To configure the maximum and minimum bandwidth allocations:

- 1. On the Main Configuration Page, select **NIC Partitioning Configuration**, and then press ENTER.
- 2. On the Partitions Configuration page (Figure 5-15), select **Global Bandwidth Allocation**.



Figure 5-15. NIC Partitioning Configuration, Global Bandwidth Allocation

3. On the Global Bandwidth Allocation page (Figure 5-16), click each partition minimum and maximum TX bandwidth field for which you want to allocate bandwidth. There are eight partitions per port in dual-port mode.

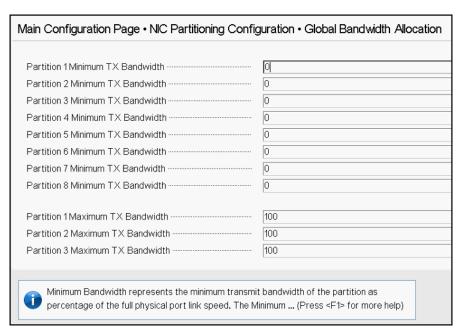


Figure 5-16. Global Bandwidth Allocation Page

□ Partition n Minimum TX Bandwidth is the minimum transmit bandwidth of the selected partition expressed as a percentage of the maximum physical port link speed. Values can be 0-100. When DCBX ETS mode is enabled, the per-traffic class DCBX ETS minimum bandwidth value is used simultaneously with the per-partition minimum TX bandwidth value. The total of the minimum TX bandwidth values of all partitions on a single port must equal 100 or be all zeros.

Setting the TX minimum bandwidth to all zeros is similar to equally dividing the available bandwidth over every active partition; however, the bandwidth is dynamically allocated over all actively sending partitions. A zero value (when one or more of the other values are set to a non-zero value) allocates a minimum of one percent to that partition, when congestion (from all of the partitions) is restricting TX bandwidth.

□ Partition n Maximum TX Bandwidth is the maximum transmit bandwidth of the selected partition expressed as a percentage of the maximum physical port link speed. Values can be 1-100. The per-partition maximum TX bandwidth value applies regardless of the DCBX ETS mode setting.

Type a value in each selected field, and then click **Back**.

4. When prompted, click **Yes** to save the changes. Changes take effect after a system reset.

To configure partitions:

- To examine a specific partition configuration, on the NIC Partitioning Configuration page (Figure 5-15 on page 61), select **Partition** n Configuration. If NParEP is not enabled, only four partitions exist per port.
- 2. To configure the first partition, select **Partition 1 Configuration** to open the Partition 1 Configuration page (Figure 5-17), which shows the following parameters:

NIC Mode (always enabled)
PCI Device ID
PCI (bus) Address
MAC Address
Virtual MAC Address

If NParEP is not enabled, only four partitions per port are available. On non-offload-capable adapters, the **FCoE Mode** and **iSCSI Mode** options and information are not displayed.

Main Configuration Page • NIC Partitioning Config	guration • Partition 1 Configuration
NIC Mode	Enabled
PCI Device ID	8070
PCI Address	86:00
MAC Address	00:0E:1E:D5:F8:76
Virtual MAC Address	00:00:00:00:00

Figure 5-17. Partition 1 Configuration

- 3. To configure the second partition, select **Partition 2 Configuration** to open the Partition 2 Configuration page. If FCoE Offload is present, the Partition 2 Configuration (Figure 5-18) shows the following parameters:
 - NIC Mode enables or disables the L2 Ethernet NIC personality on Partitions 2 and greater. To disable any of the remaining partitions, set the NIC Mode to Disabled. To disable offload-capable partitions, disable both the NIC Mode and respective offload mode.
 - FCoE Mode enables or disables the FCoE-Offload personality on the second partition. If you enable this mode on the second partition, you should disable NIC Mode. Because only one offload is available per port, if FCoE-Offload is enabled on the port's second partition, iSCSI-Offload cannot be enabled on the third partition of that same NPAR mode port. Not all adapters support FCoE Mode.

iSCSI Mode enables or disables the iSCSI-Offload personality on the third partition. If you enable this mode on the third partition, you should disable NIC Mode. Because only one offload is available per port, if iSCSI-Offload is enabled on the port's third partition, FCoE-Offload cannot be enabled on the second partition of that same NPAR mode port. Not all adapters support iSCSI Mode.
FIP MAC Address ¹
Virtual FIP MAC Address ¹
World Wide Port Name ¹
Virtual World Wide Port Name ¹
World Wide Node Name ¹
Virtual World Wide Node Name ¹
PCI Device ID
PCI (bus) Address

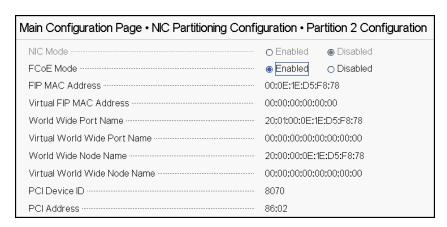


Figure 5-18. Partition 2 Configuration: FCoE Offload

□ NIC Mode (Disabled)□ iSCSI Offload Mode (Enabled)	4.
☐ iSCSI Offload MAC Address ²	

Virtual iSCSI Offload MAC Address²

¹ This parameter is only present on the second partition of an NPAR mode port of FCoE offload-capable adapters.

² This parameter is only present on the third partition of an NPAR mode port of iSCSI offload-capable adapters.



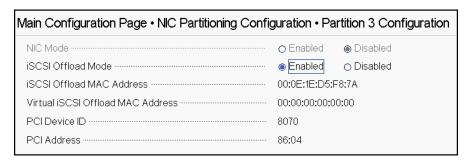


Figure 5-19. Partition 3 Configuration: iSCSI Offload

- 5. To configure the remaining Ethernet partitions, including the previous (if not offload-enabled), open the page for a partition 2 or greater partition (see Figure 5-20).
 - NIC Mode (Enabled or Disabled). When disabled, the partition is hidden such that it does not appear to the OS if fewer than the maximum quantity of partitions (or PCI PFs) are detected.
 - □ PCI Device ID
 - □ PCI Address
 - MAC Address
 - □ Virtual MAC Address

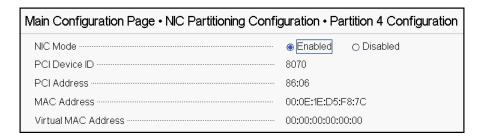


Figure 5-20. Partition 4 Configuration

Partitioning for VMware ESXi 6.5 and ESXi 6.7

If the following conditions exist on a system running either VMware ESXi 6.5 or ESXi 6.7, you must uninstall and reinstall the drivers:

- The adapter is configured to enable NPAR with all NIC partitions.
- The adapter is in Single Function mode.
- The configuration is saved and the system is rebooted.

- Storage partitions are enabled (by converting one of the NIC partitions as storage) while drivers are already installed on the system.
- Partition 2 is changed to FCoE.
- The configuration is saved and the system is rebooted again.

Driver re-installation is required because the storage functions may keep the <code>vmnicX</code> enumeration rather than <code>vmhbaX</code>, as shown when you issue the following command on the system:

esxcfg-scsidevs -a

```
vmnic4 gedf
                        link-up fc.2000000e1ed6fa2a:2001000e1ed6fa2a
(0000:19:00.2) QLogic Corp. QLogic FastLinQ QL41xxx Series 10/25 GbE
Controller (FCoE)
vmhba0 lsi mr3
                        link-n/a sas.51866da071fa9100
(0000:18:00.0) Avago (LSI) PERC H330 Mini
vmnic10 gedf
                        link-up fc.2000000e1ef249f8:2001000e1ef249f8
(0000:d8:00.2) QLogic Corp. QLogic FastLinQ QL41xxx Series 10/25 GbE
Controller (FCoE)
vmhbal vmw ahci
                        link-n/a sata.vmhba1
(0000:00:11.5) Intel Corporation Lewisburg SSATA Controller [AHCI mode]
vmhba2 vmw ahci
                        link-n/a sata.vmhba2
(0000:00:17.0) Intel Corporation Lewisburg SATA Controller [AHCI mode]
vmhba32 qedil
                      online iscsi.vmhba32
                                                                    QLogic
FastLinQ QL41xxx Series 10/25 GbE Controller (iSCSI)
vmhba33 gedil
                      online iscsi.vmhba33
                                                                    QLogic
FastLinQ QL41xxx Series 10/25 GbE Controller (iSCSI)
```

In the preceding command output, notice that <code>vmnic4</code> and <code>vmnic10</code> are actually storage adapter ports. To prevent this behavior, you should enable storage functions at the same time that you configure the adapter for NPAR mode.

For example, assuming that the adapter is in Single Function mode by default, you should:

- 1. Enable NPAR mode.
- 2. Change Partition 2 to FCoE.
- 3. Save and reboot.

6 Boot from SAN Configuration

SAN boot enables deployment of diskless servers in an environment where the boot disk is located on storage connected to the SAN. The server (initiator) communicates with the storage device (target) through the SAN using the Marvell Converged Network Adapter (CNA) Host Bus Adapter (HBA).

To enable FCoE-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

This chapter covers boot from SAN configuration for both iSCSI and FCoE:

- iSCSI Boot from SAN
- "FCoE Boot from SAN" on page 113

iSCSI Boot from SAN

Marvell 41xxx Series gigabit Ethernet (GbE) adapters support iSCSI boot to enable 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.

This section provides the following configuration information about iSCSI boot from SAN:

- iSCSI Out-of-Box and Inbox Support
- iSCSI Preboot Configuration
- Configuring iSCSI Boot from SAN on Windows
- Configuring iSCSI Boot from SAN on Linux
- Configuring iSCSI Boot from SAN on VMware

iSCSI Out-of-Box and Inbox Support

Table 6-1 lists the operating systems' inbox and out-of-box support for iSCSI boot from SAN (BFS).

Table 6-1. iSCSI Out-of-Box and Inbox Boot from SAN Support

	Out-of-Box		<u>Inbox</u>	
OS Version	SW iSCSI BFS Support	Hardware Offload iSCSI BFS Support	SW iSCSI BFS Support	Hardware Offload iSCSI BFS Support
Windows 2012 ^a	Yes	Yes	No	No
Windows 2012 R2 ^a	Yes	Yes	No	No
Windows 2016 ^b	Yes	Yes	Yes	No
Windows 2019	Yes	Yes	Yes	Yes
RHEL 7.5	Yes	Yes	Yes	Yes
RHEL 7.6	Yes	Yes	Yes	Yes
RHEL 8.0	Yes	Yes	Yes	Yes
SLES 12 SP3	Yes	Yes	Yes	Yes
SLES 15/15 SP1	Yes	Yes	Yes	Yes
vSphere ESXi 6.5 U3 ^c	Yes	No	Yes	No
vSphere ESXi 6.7 U2 ^c	Yes	No	Yes	No

^a Windows Server 2012 and 2012 R2 do not support the inbox iSCSI driver for SW or hardware offload.

iSCSI Preboot Configuration

For both Windows and Linux operating systems, configure iSCSI boot with **UEFI iSCSI HBA** (offload path with the Marvell offload iSCSI driver). Set this option using Boot Protocol, under **Port Level Configuration**. To support iSCSI boot, first enable the iSCSI HBA in the UEFI HII and then set the boot protocol accordingly.

^b Windows Server 2016 does not support the inbox iSCSI driver for hardware offload.

^c ESXi out-of-box and inbox do not support native hardware offload iSCSI boot. The system will perform a SW boot and connection and then will transition to hardware offload.

For both Windows and Linux operating systems, iSCSI boot can be configured to boot with two distinctive methods:

■ **iSCSI SW** (also known as non-offload path with Microsoft/Open-iSCSI initiator)

Follow the Dell BIOS guide for iSCSI software installation.

■ **ISCSI HW** (offload path with the Marvell FastLinQ offload iSCSI driver). This option can be set using **Boot Mode**.

iSCSI hardware installation instructions start in "Enabling NPAR and the iSCSI HBA" on page 71.

For VMware ESXi operating systems, only the iSCSI SW method is supported.

iSCSI preboot information in this section includes:

- Setting the BIOS Boot Mode to UEFI
- Enabling NPAR and the iSCSI HBA
- Selecting the iSCSI UEFI Boot Protocol
- Configuring the Storage Target
- Configuring iSCSI Boot Options
- Configuring the DHCP Server to Support iSCSI Boot

Setting the BIOS Boot Mode to UEFI

To configure the boot mode:

- 1. Restart the system.
- 2. Access the System BIOS menu.
- 3. For the **Boot Mode** setting, select **UEFI** (see Figure 6-1).

NOTE

SAN boot is supported in UEFI environment only. Make sure the system boot option is UEFI, and not legacy.

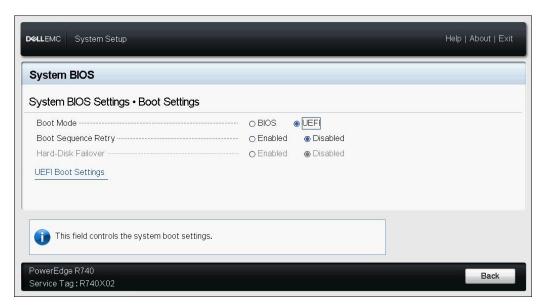


Figure 6-1. System Setup: Boot Settings

Enabling NPAR and the iSCSI HBA

To enable NPAR and the iSCSI HBA:

In the System Setup, Device Settings, select the QLogic device (Figure 6-2).
 Refer to the OEM user guide on accessing the PCI device configuration menu.

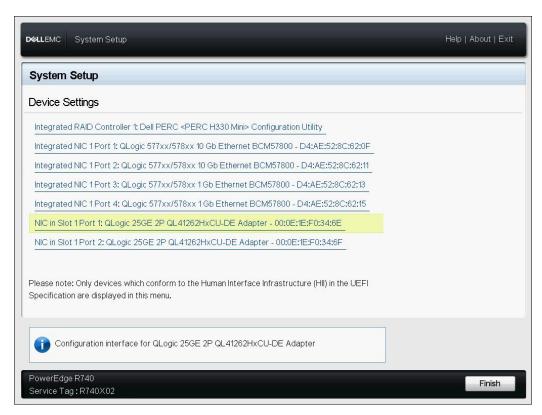


Figure 6-2. System Setup: Device Settings

2. Enable NPAR.

Configuring the Storage Target

Configuring the storage target varies by target vendors. For information on configuring the storage target, refer to the documentation provided by the vendor.

To configure the storage target:

- 1. Select the appropriate procedure based on your storage target, either:
 - □ Create an storage target using software such as SANBlaze[®] or Linux-IO (LIO[™]) Target.
 - ☐ Create a vdisk or volume using a target array such as EqualLogic® or EMC®.

Create a virtual disk.

Selecting the iSCSI UEFI Boot Protocol

Before selecting the preferred boot mode, ensure that the **Device Level Configuration** menu setting is **Enable NPAR** and that the **NIC Partitioning Configuration** menu setting is **Enable iSCSI HBA**.

The **Boot Mode** option is listed under **NIC Configuration** (Figure 6-3) for the adapter, and the setting is port specific. Refer to the OEM user manual for direction on accessing the device-level configuration menu under UEFI HII.

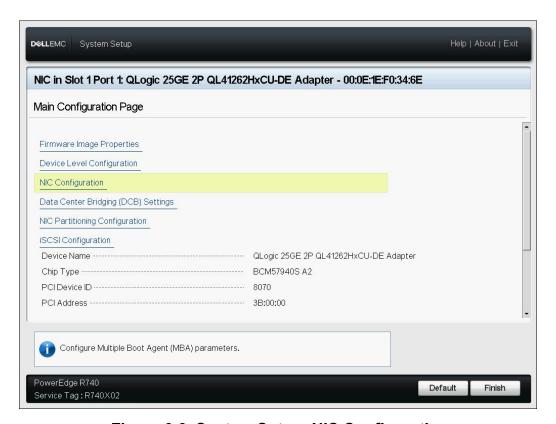


Figure 6-3. System Setup: NIC Configuration

NOTE

Boot from SAN boot is supported only in NPAR mode and is configured in UEFI, and not in legacy BIOS.

1. On the NIC Configuration page (Figure 6-4), for the **Boot Protocol** option, select **UEFI iSCSI HBA** (requires NPAR mode).

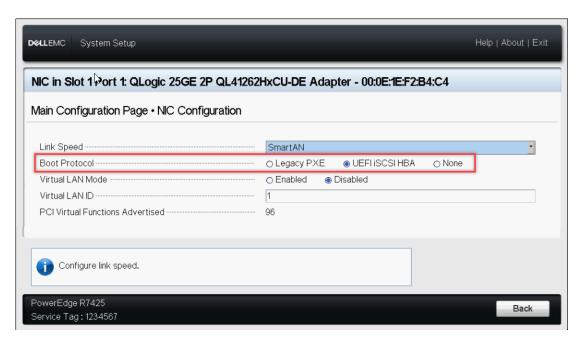


Figure 6-4. System Setup: NIC Configuration, Boot Protocol

NOTE

Use the **Virtual LAN Mode** and **Virtual LAN ID** options on this page only for PXE boot. If a vLAN is needed for UEFI iSCSI HBA boot mode, see Step 3 of Static iSCSI Boot Configuration.

Configuring iSCSI Boot Options

iSCSI boot configuration options include:

- Static iSCSI Boot Configuration
- Dynamic iSCSI Boot Configuration
- Enabling CHAP Authentication

Static iSCSI Boot Configuration

In a static configuration, you must enter data for the following:

- Initiator IP address
- Initiator IQN
- Target parameters (obtained in "Configuring the Storage Target" on page 71)

For information on configuration options, see Table 6-2 on page 76.

To configure the iSCSI boot parameters using static configuration:

1. In the Device HII **Main Configuration Page**, select **iSCSI Configuration** (Figure 6-5), and then press ENTER.

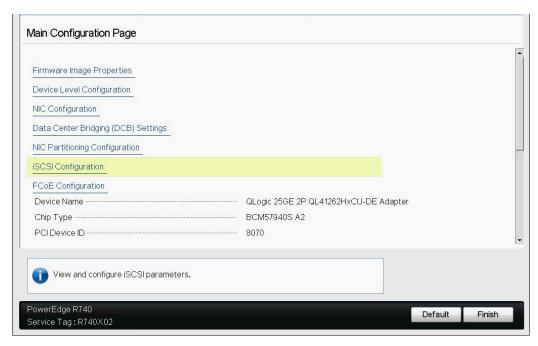


Figure 6-5. System Setup: iSCSI Configuration

2. On the **iSCSI Configuration** page, select **iSCSI General Parameters** (Figure 6-6), and then press ENTER.



Figure 6-6. System Setup: Selecting General Parameters

 On the iSCSI General Parameters page (Figure 6-7), press the DOWN ARROW key to select a parameter, and then press the ENTER key to input the following values (Table 6-2 on page 76 provides descriptions of these parameters):

TCP/IP Parameters via DHCP: Disabled iSCSI Parameters via DHCP: Disabled CHAP Authentication: As required **IP Version**: As required (IPv4 or IPv6) **CHAP Mutual Authentication**: As required **DHCP Vendor ID**: Not applicable for static configuration **HBA Boot Mode**: As required Virtual LAN ID: Default value or as required Virtual LAN Mode: As required



Figure 6-7. System Setup: iSCSI General Parameters

Table 6-2. iSCSI General Parameters

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 using 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 on the iSCSI Initiator Parameters Configuration page.
CHAP Authentication	Controls whether the iSCSI boot host software uses CHAP authentication when connecting to the iSCSI target. If CHAP Authentication is enabled, configure the CHAP ID and CHAP Secret on the iSCSI Initiator Parameters Configuration page.
IP Version	This option is specific to IPv6. Toggles between $\ \ \text{IPv4}$ and $\ \ \text{IPv6}.$ All IP settings are lost if you switch from one protocol version to another.
CHAP Mutual Authentication	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 on the iSCSI Initiator Parameters Configuration page.
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.
HBA Boot Mode	Controls whether SW or Offload is enabled or disabled. For Offload, this option is unavailable (grayed out). For information about SW (non-offload), refer to the Dell BIOS configuration.
Virtual LAN ID	vLAN ID range is 1–4094.
Virtual LAN Mode	Enables or disables vLAN.

4. Return to the iSCSI Configuration page, and then press the ESC key.



5. Select **iSCSI Initiator Parameters** (Figure 6-8), and then press ENTER.

Figure 6-8. System Setup: Selecting iSCSI Initiator Parameters

6.	he iSCSI Initiator Parameters page (Figure 6-9), select the following meters, and then type a value for each:
	IPv4* Address
	Subnet Mask
	IPv4* Default Gateway
	IPv4* Primary DNS
	IPv4* Secondary DNS
	iSCSI Name . Corresponds to the iSCSI initiator name to be used by the client system.
	CHAP ID

CHAP Secret

NOTE

For the preceding items with asterisks (*), note the following:

- The label will change to **IPv6** or **IPv4** (default) based on the IP version set on the iSCSI General Parameters page (Figure 6-7 on page 75).
- Carefully enter the IP address. There is no error-checking performed against the IP address to check for duplicates, incorrect segment, or network assignment.

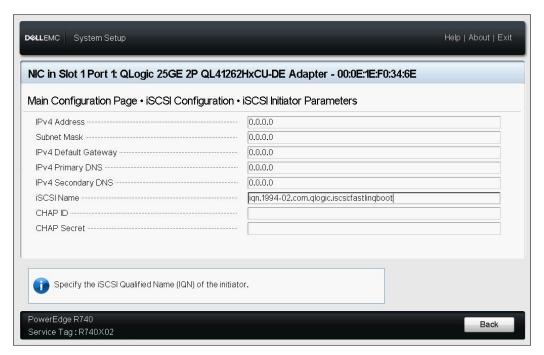


Figure 6-9. System Setup: iSCSI Initiator Parameters

7. Return to the iSCSI Configuration page, and then press ESC.

8. Select **iSCSI First Target Parameters** (Figure 6-10), and then press ENTER.

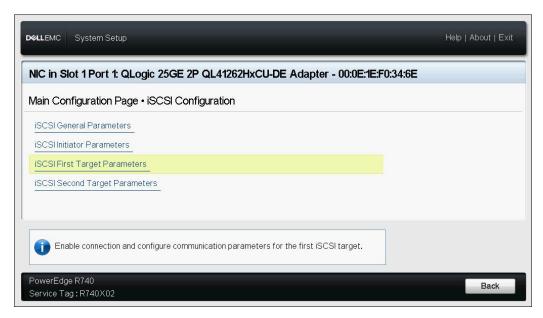


Figure 6-10. System Setup: Selecting iSCSI First Target Parameters

- 9. On the iSCSI First Target Parameters page, set the **Connect** option to **Enabled** for the iSCSI target.
- 10. Type values for the following parameters for the iSCSI target, and then press ENTER:
 - ☐ IPv4* Address
 - □ TCP Port
 - □ Boot LUN
 - ☐ iSCSI Name
 - □ CHAP ID
 - ☐ CHAP Secret

NOTE

For the preceding parameters with an asterisk (*), the label will change to **IPv6** or **IPv4** (default) based on IP version set on the iSCSI General Parameters page, as shown in Figure 6-11.

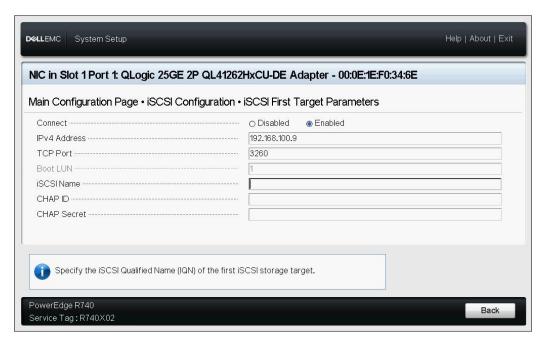


Figure 6-11. System Setup: iSCSI First Target Parameters

11. Return to the iSCSI Boot Configuration page, and then press ESC.

12. If you want to configure a second iSCSI target device, select **iSCSI Second Target Parameters** (Figure 6-12), and enter the parameter values as you did in Step 10. This second target is used if the first target cannot be connected to.Otherwise, proceed to Step 13.



Figure 6-12. System Setup: iSCSI Second Target Parameters

- 13. Press ESC once, and a second time to exit.
- 14. Click **Yes** to save changes, or follow the OEM guidelines to save the device-level configuration. For example, click **Yes** to confirm the setting change (Figure 6-13).

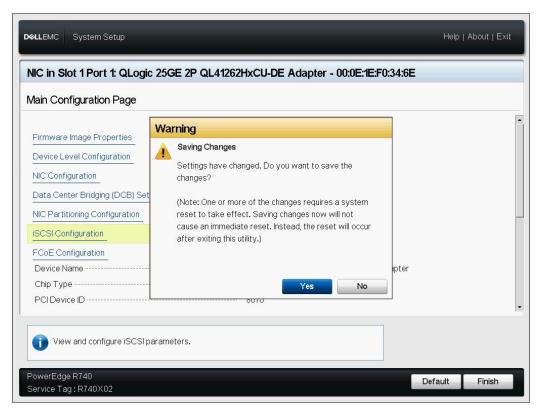


Figure 6-13. System Setup: Saving iSCSI Changes

15. After all changes have been made, reboot the system to apply the changes to the adapter's running configuration.

Dynamic iSCSI Boot Configuration

In a dynamic configuration, ensure that the system's IP address and target (or initiator) information are provided by a DHCP server (see IPv4 and IPv6 configurations in "Configuring the DHCP Server to Support iSCSI Boot" on page 85).

Any settings for the following parameters are ignored and do not need to be cleared (with the exception of the initiator iSCSI name for IPv4, CHAP ID, and CHAP secret for IPv6):

- Initiator Parameters
- First Target Parameters or Second Target Parameters

For information on configuration options, see Table 6-2 on page 76.

NOTE

When using a DHCP server, the DNS server entries are overwritten by the values provided by the DHCP server. This override 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 statically configures the appropriate registries. It will overwrite whatever is configured. Because the DHCP daemon runs in the Windows environment as a user process, all TCP/IP parameters must 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 window. If no value was selected, the controller defaults to the following name:

```
iqn.1995-05.com.qlogic.<11.22.33.44.55.66>.iscsiboot
```

The string 11.22.33.44.55.66 corresponds to the controller's MAC address. If DHCP Option 43 (IPv4 only) is used, any settings on the following windows are ignored and do not need to be cleared:

Initiator Parameters

■ First Target Parameters, or Second Target Parameters

To configure the iSCSI boot parameters using dynamic configuration:

	he iSCSI General Parameters page, set the following options, as shown gure 6-14:
	TCP/IP Parameters via DHCP: Enabled
	iSCSI Parameters via DHCP: Enabled
	CHAP Authentication: As required
	IP Version: As required (IPv4 or IPv6)
	CHAP Mutual Authentication: As required
	DHCP Vendor ID: As required
	HBA Boot Mode: As required
	Virtual LAN ID: As required

Virtual LAN Mode: As required¹

¹ **Virtual LAN Mode** is not necessarily required when using a dynamic (externally provided) configuration from the DHCP server.



Figure 6-14. System Setup: iSCSI General Parameters

Enabling CHAP Authentication

Ensure that the CHAP authentication is enabled on the target.

To enable CHAP authentication:

- 1. Go to the iSCSI General Parameters page.
- 2. Set CHAP Authentication to Enabled.
- 3. In the Initiator Parameters window, type values for the following:
 - ☐ CHAP ID (up to 255 characters)
 - ☐ CHAP Secret (if authentication is required; must be 12 to 16 characters in length)
- 4. Press ESC to return to the iSCSI Boot Configuration page.
- 5. On the iSCSI Boot Configuration Menu, select iSCSI First Target Parameters.
- 6. In the iSCSI First Target Parameters window, type values used when configuring the iSCSI target:
 - ☐ CHAP ID (optional if two-way CHAP)
 - ☐ CHAP Secret (optional if two-way CHAP; must be 12 to 16 characters in length or longer)

- 7. Press ESC to return to the iSCSI Boot Configuration Menu.
- 8. Press ESC, and then select confirm **Save Configuration**.

Configuring the DHCP Server to Support iSCSI Boot

The DHCP server is an optional component, and is only necessary if you will be doing a dynamic iSCSI boot configuration setup (see "Dynamic iSCSI Boot Configuration" on page 82).

Configuring the DHCP server to support iSCSI boot differs for IPv4 and IPv6:

- DHCP iSCSI Boot Configurations for IPv4
- Configuring the DHCP Server
- Configuring DHCP iSCSI Boot for IPv6
- Configuring vLANs for iSCSI Boot

DHCP iSCSI Boot Configurations for IPv4

DHCP includes several options that provide configuration information to the DHCP client. For iSCSI boot, Marvell FastLinQ adapters support the following DHCP configurations:

- DHCP Option 17, Root Path
- DHCP Option 43, Vendor-specific Information

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 IETC RFC 4173, is:

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

Table 6-3 lists the DHCP Option 17 parameters.

Table 6-3. DHCP Option 17 Parameter Definitions

Parameter	Definition		
"iscsi:"	A literal string		
<servername></servername>	IP address or fully qualified domain name (FQDN) of the iSCS target		
":"	Separator		
<pre><pre><pre><pre><pre><pre><pre><pre></pre></pre></pre></pre></pre></pre></pre></pre>	IP protocol used to access the iSCSI target. Because only TCP currently supported, the protocol is 6.		
<port></port>	Port number associated with the protocol. The standard port number for iSCSI is 3260.		

Table 6-3. DHCP Option 17 Parameter Definitions (Continued)

Parameter	Definition		
<lun></lun>	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 must be configured as 40 within the Option 17 parameter on the DHCP server.		
<targetname></targetname>	Target name in either IQN or EUI format. For details on both IQN and EUI formats, refer to RFC 3720. An example IQN name is iqn.1995-05.com.QLogic:iscsi-target.		

DHCP Option 43, Vendor-specific Information

DHCP Option 43 (vendor-specific information) provides more configuration options to the iSCSI client than does DHCP Option 17. In this configuration, three additional sub-options 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.

Table 6-4 lists the DHCP Option 43 sub-options.

Table 6-4. DHCP Option 43 Sub-option Definitions

Sub-option	Definition	
201	First iSCSI target information in the standard root path format:	
"iscsi:" <servername>":"<protocol>":"<port>":"<lun>": "<targetname>"</targetname></lun></port></protocol></servername>		
202	Second iSCSI target information in the standard root path format:	
	"iscsi:" <servername>":"<protocol>":"<port>":"<lun>": "<targetname>"</targetname></lun></port></protocol></servername>	
203	203 iSCSI initiator IQN	

Using DHCP Option 43 requires more configuration than DHCP Option 17, but it provides a richer environment and more configuration options. You should use DHCP Option 43 when performing dynamic iSCSI boot configuration.

Configuring the DHCP Server

Configure the DHCP server to support either Option 16, 17, or 43.

NOTE

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

If you use Option 43, you must also configure Option 60. The value of Option 60 must match the DHCP Vendor ID value, QLGC ISAN, as shown in the **iSCSI General Parameters** of the iSCSI Boot Configuration page.

Configuring DHCP iSCSI Boot for IPv6

The DHCPv6 server can provide several options, including stateless or stateful IP configuration, as well as information for the DHCPv6 client. For iSCSI boot, Marvell FastLinQ adapters support the following DHCP configurations:

- DHCPv6 Option 16, Vendor Class Option
- DHCPv6 Option 17, Vendor-Specific Information

NOTE

The DHCPv6 standard Root Path option is not yet available. Marvell 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 your configured DHCP Vendor ID parameter. The DHCP Vendor ID value is QLGC ISAN, as shown in the **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 sub-options are provided that assign the initiator IQN to the iSCSI boot client, along with two iSCSI target IQNs that can be used for booting.

Table 6-5 lists the DHCP Option 17 sub-options.

Table 6-5. DHCP Option 17 Sub-option Definitions

Sub-option	Definition		
201	First iSCSI target information in the standard root path format: "iscsi:"[<servername>]":"<protocol>":"<put>":"<lun> ": "<targetname>"</targetname></lun></put></protocol></servername>		
202	Second iSCSI target information in the standard root path format: "iscsi:"[<servername>]":"<protocol>":"<port>":"<lun> ": "<targetname>"</targetname></lun></port></protocol></servername>		
203	iSCSI initiator IQN		

Brackets [] are required for the IPv6 addresses.

The format of Option 17 should be:

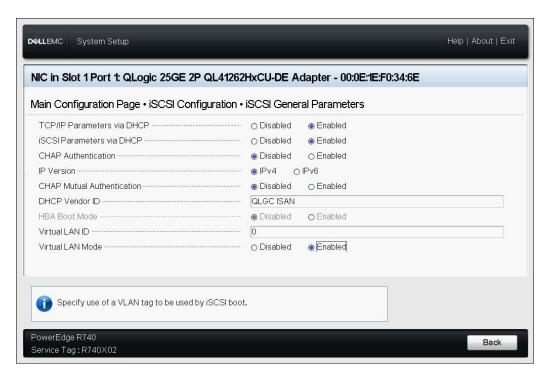
<2-byte Option Number 201|202|203> <2-byte length> <data>

Configuring vLANs for iSCSI Boot

iSCSI traffic on the network may be isolated in a Layer 2 vLAN to segregate it from general traffic. If this is the case, make the iSCSI interface on the adapter a member of that vLAN.

To configure vLAN for iSCSI boot:

- 1. Go to the **iSCSI Configuration Page** for the port.
- 2. Select iSCSI General Parameters.



3. Select **VLAN ID** to enter and set the vLAN value, as shown in Figure 6-15.

Figure 6-15. System Setup: iSCSI General Parameters, VLAN ID

Configuring iSCSI Boot from SAN on Windows

Adapters support iSCSI boot to enable network boot of operating systems to diskless systems. iSCSI boot allows a Windows operating system to boot from an iSCSI target machine located remotely over a standard IP network. You can set the L4 iSCSI option (offload path with Marvell offload iSCSI driver) by opening the **NIC Configuration** menu and setting the **Boot Protocol** to **UEFI iSCSI**.

iSCSI boot from SAN for Windows information includes the following:

- Before You Begin
- Selecting the Preferred iSCSI Boot Mode
- Configuring iSCSI General Parameters
- Configuring the iSCSI Initiator
- Configuring the iSCSI Targets
- Detecting the iSCSI LUN and Injecting the Marvell Drivers

Before You Begin

Before you begin configuring iSCSI boot from SAN on a Windows machine, note the following:

- iSCSI boot is only supported for NPAR with **NParEP Mode**. Before configuring iSCSI boot:
 - 1. Access the Device Level Configuration page.
 - 2. Set the **Virtualization Mode** to **Npar** (NPAR).
 - 3. Set the NParEP Mode to Enabled.
- The server boot mode must be UEFI.
- iSCSI boot on 41xxx Series Adapters is not supported in legacy BIOS.
- Marvell recommends that you disable the Integrated RAID Controller.

Selecting the Preferred iSCSI Boot Mode

To select the iSCSI boot mode on Windows:

- 1. On the NIC Partitioning Configuration page for a selected partition, set the iSCSI Offload Mode to Enabled.
- 2. On the NIC Configuration page, set the **Boot Protocol** option to **UEFI iSCSI HBA**.

Configuring iSCSI General Parameters

Configure the Marvell iSCSI boot software for either static or dynamic configuration. For configuration options available from the General Parameters window, see Table 6-2 on page 76, which lists parameters for both IPv4 and IPv6.

To set the iSCSI general parameters on Windows:

- 1. From the Main Configuration page, select **iSCSI Configuration**, and then select **iSCSI General Parameters**.
- On the iSCSI General Parameters page (see Figure 6-7 on page 75), press
 the DOWN ARROW key to select a parameter, and then press the ENTER
 key to input the following values (see Table 6-2 on page 76 provides
 descriptions of these parameters):

TCP/IP Parameters via DHCP: Disabled (for static iSCSI boot), or Enabled (for dynamic iSCSI boot)
iSCSI Parameters via DHCP: Disabled
CHAP Authentication: As required
IP Version: As required (IPv4 or IPv6)

□ Virtual LAN ID: (Optional) You can isolate iSCSI traffic on the network in a Layer 2 vLAN to segregate it from general traffic. To segregate traffic, make the iSCSI interface on the adapter a member of the Layer 2 vLAN by setting this value.

Configuring the iSCSI Initiator

To set the iSCSI initiator parameters on Windows:

- 1. From the Main Configuration page, select **iSCSI Configuration**, and then select **iSCSI Initiator Parameters**.
- 2. On the iSCSI Initiator Parameters page (see Figure 6-9 on page 78), select the following parameters, and then type a value for each:

IPv4* Address
Subnet Mask
IPv4* Default Gateway
IPv4* Primary DNS
IPv4* Secondary DNS
Virtual LAN ID : (Optional) You can isolate iSCSI traffic on the network in a Layer 2 vLAN to segregate it from general traffic. To segregate traffic, make the iSCSI interface on the adapter a member of the Layer 2 vLAN by setting this value.
iSCSI Name . Corresponds to the iSCSI initiator name to be used by the client system.
CHAP ID
CHAP Secret

NOTE

For the preceding items with asterisks (*), note the following:

- The label will change to **IPv6** or **IPv4** (default) based on the IP version set on the iSCSI General Parameters page (see Figure 6-7 on page 75).
- Carefully enter the IP address. There is no error-checking performed against the IP address to check for duplicates, incorrect segment, or network assignment.
- 3. Select **iSCSI First Target Parameters** (Figure 6-10 on page 79), and then press ENTER.

Configuring the iSCSI Targets

You can set up the iSCSI first target, second target, or both at once.

To set the iSCSI target parameters on Windows:

- 1. From the Main Configuration page, select **iSCSI Configuration**, and then select **iSCSI First Target Parameters**.
- 2. On the iSCSI First Target Parameters page, set the **Connect** option to **Enabled** for the iSCSI target.

3.	Type values for the following parameters for the iSCSI target, and then
	press ENTER:

IPv4* Address
TCP Port
Boot LUN
iSCSI Name
CHAP ID
CHAP Secret

NOTE

For the preceding parameters with an asterisk (*), the label will change to **IPv6** or **IPv4** (default) based on IP version set on the iSCSI General Parameters page, as shown in Figure 6-7 on page 75.

- 4. If you want to configure a second iSCSI target device, select **iSCSI Second Target Parameters** (Figure 6-12 on page 81), and enter the parameter values as you did in Step 3. This second target is used if the first target cannot be connected to.Otherwise, proceed to Step 5.
- 5. In the Warning dialog box, click **Yes** to save the changes, or follow the OEM guidelines to save the device-level configuration.

Detecting the iSCSI LUN and Injecting the Marvell Drivers

1. Reboot the system, access the HII, and determine if the iSCSI LUN is detected. Issue the following UEFI Shell (version 2) script command:

The output from the preceding command shown in Figure 6-16 indicates that the iSCSI LUN was detected successfully at the preboot level.

Figure 6-16. Detecting the iSCSI LUN Using UEFI Shell (Version 2)

- 2. On the newly detected iSCSI LUN, select an installation source such as using a WDS server, mounting the .iso with an integrated Dell Remote Access Controller (iDRAC), or using a CD/DVD.
- 3. In the Windows Setup window (Figure 6-17), select the drive name on which to install the driver.

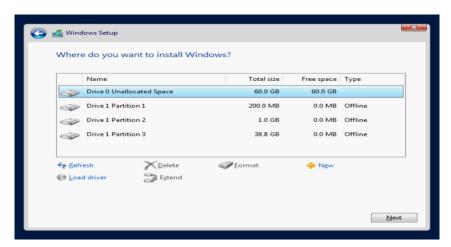


Figure 6-17. Windows Setup: Selecting Installation Destination

- 4. Inject the latest Marvell drivers by mounting drivers in the virtual media:
 - a. Click **Load driver**, and then click **Browse** (see Figure 6-18).

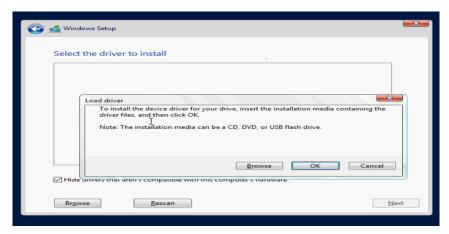


Figure 6-18. Windows Setup: Selecting Driver to Install

- b. Navigate to the driver location and choose the gevbd driver.
- c. Choose the adapter on which to install the driver, and then click **Next** to continue.
- 5. Repeat Step 4 to load the geios driver (Marvell L4 iSCSI driver).
- 6. After injecting the qevbd and qeios drivers, click **Next** to begin installation on the iSCSI LUN. Then follow the on-screen instructions.
 - The server will undergo a reboot multiple times as part of the installation process, and then will boot up from the iSCSI boot from SAN LUN.
- 7. If it does not automatically boot, access the **Boot Menu** and select the specific port boot entry to boot from the iSCSI LUN.

Configuring iSCSI Boot from SAN on Linux

This section provides iSCSI boot from SAN procedures for the following Linux distributions:

- Configuring iSCSI Boot from SAN for RHEL 7.5 and Later
- Configuring iSCSI Boot from SAN for SLES 12 SP3 and Later
- Configuring iSCSI Boot from SAN for Other Linux Distributions

Configuring iSCSI Boot from SAN for RHEL 7.5 and Later

To install RHEL 7.5 and later:

1. Boot from the RHEL 7.x installation media with the iSCSI target already connected in UEFI.

```
Install Red Hat Enterprise Linux 7.x
Test this media & install Red Hat Enterprise 7.x
Troubleshooting -->
Use the UP and DOWN keys to change the selection
Press 'e' to edit the selected item or 'c' for a command prompt
```

- 2. To install an out-of-box driver, press the E key. Otherwise, proceed to Step 6.
- 3. Select the kernel line, and then press the E key.
- 4. Issue the following command, and then press ENTER.

```
inst.dd modprobe.blacklist=qed,qede,qedr,qedi,qedf
```

The installation process prompts you to install the out-of-box driver.

- 5. If required for your setup, load the FastLinQ driver update disk when prompted for additional driver disks. Otherwise, if you have no other driver update disks to install, press the C key.
- 6. Continue with the installation. You can skip the media test. Click **Next** to continue with the installation.
- 7. In the Configuration window, select the language to use during the installation process, and then click **Continue**.
- 8. In the Installation Summary window, click **Installation Destination**. The disk label is *sda*, indicating a single-path installation. If you configured multipath, the disk has a device mapper label.
- 9. In the Specialized & Network Disks section, select the iSCSI LUN.
- 10. Type the root user's password, and then click **Next** to complete the installation.
- 11. Reboot the server, and then add the following parameters in the command line:

```
rd.iscsi.firmware
rd.driver.pre=qed,qedi (to load all drivers pre=qed,qedi,qede,qedf)
selinux=0
```

12. After a successful system boot, edit the

/etc/modprobe.d/anaconda-blacklist.conf file to remove the blacklist entry for the selected driver.

- 13. Edit the /etc/default/grub file as follows:
 - a. Locate the string in double quotes as shown in the following example. The command line is a specific reference to help find the string.

```
GRUB CMDLINE LINUX="iscsi firmware"
```

b. The command line may contain other parameters that can remain. Change only the iscsi firmware string as follows:

```
GRUB CMDLINE LINUX="rd.iscsi.firmware selinux=0"
```

- 14. Create a new grub.cfg file by issuing the following command:
 - # grub2-mkconfig -o /boot/efi/EFI/redhat/grub.cfg
- 15. Rebuild the ramdisk by issuing the **dracut** -f command, and then reboot.

NOTE

When installing iSCSI BFS in Linux with a multipath I/O (MPIO) configuration and a single path active, use the following settings in the multipath.conf file:

```
defaults {
    find_multipaths yes
    user_friendly_names yes
    polling_interval 5
    fast_io_fail_tmo 5
    dev_loss_tmo 10
    checker_timeout 15
    no_path_retry queue
}
```

These suggested settings are tunable and provided as guidance for iSCSI BFS to be operational.

For more information, contact the appropriate OS vendor.

Configuring iSCSI Boot from SAN for SLES 12 SP3 and Later

To install SLES 12 SP3 and later:

- 1. Boot from the SLES 12 SP3 installation media with the iSCSI target pre-configured and connected in UEFI.
- 2. Update the latest driver package by adding the <code>dud=1</code> parameter in the installer command parameter. The driver update disk is required because the necessary iSCSI drivers are not inbox.

NOTE

For SLES 12 SP3 only: If the server is configured for Multi-Function mode (NPAR), you must provide the following additional parameters as part of this step:

```
dud=1 brokenmodules=qed,qedi,qedf,qede withiscsi=1
[BOOT_IMAGE=/boot/x86_64/loader/linux dud=1
brokenmodules=qed,qedi,qedf,qede withiscsi=1]
```

3. Complete the installation steps specified by the SLES 12 SP3 OS.

Known Issue in DHCP Configuration

In DHCP configuration for SLES 12 SP3 and later, the first boot after an OS installation may fail if the initiator IP address acquired from the DHCP server is in a different range than the target IP address. To resolve this issue, boot the OS using static configuration, update the latest iscsiuio out-of-box RPM, rebuild the initrd, and then reboot the OS using DHCP configuration. The OS should now boot successfully.

Configuring iSCSI Boot from SAN for Other Linux Distributions

For distributions such as RHEL 6.9/6.10/7.2/7.3, SLES 11 SP4, and SLES 12 SP1/2, the inbox iSCSI user space utility (Open-iSCSI tools) lacks support for qedi iSCSI transport and cannot perform user space-initiated iSCSI functionality. During boot from SAN installation, you can update the qedi driver using a driver update disk (DUD). However, no interface or process exists to update userspace inbox utilities, which causes the iSCSI target login and boot from SAN installation to fail.

To overcome this limitation, perform the initial boot from SAN with the pure L2 interface (do not use hardware-offloaded iSCSI) using the following procedure during the boot from SAN.

iSCSI offload for other distributions of Linux includes the following information:

- Booting from SAN Using a Software Initiator
- Migrating from Software iSCSI Installation to Offload iSCSI
- Linux Multipath Considerations

Booting from SAN Using a Software Initiator

To boot from SAN using a software initiator with Dell OEM Solutions:

NOTE

The preceding step is required because DUDs contain qedi, which binds to the iSCSI PF. After it is bound, Open-iSCSI infrastructure fails due to unknown transport driver.

- Access the Dell EMC System BIOS settings.
- 2. Configure the initiator and target entries. For more information, refer to the Dell BIOS user guide.
- 3. At the beginning of the installation, pass the following boot parameter with the DUD option:
 - \Box For RHEL 6.x, 7.x, and older:

```
rd.iscsi.ibft dd
```

No separate options are required for older distributions of RHEL.

☐ For SLES 11 SP4 and SLES 12 SP1/SP2:

ip=ibft dud=1

☐ For the FastLinQ DUD package (for example, on RHEL 7):

fastling-8.18.10.0-dd-rhel7u3-3.10.0 514.el7-x86 64.iso

Where the DUD parameter is dd for RHEL 7.x and dud=1 for SLES 12.x.

4. Install the OS on the target LUN.

Migrating from Software iSCSI Installation to Offload iSCSI

Migrate from the non-offload interface to an offload interface by following the instructions for your operating system.

- Migrating to Offload iSCSI for RHEL 6.9/6.10
- Migrating to Offload iSCSI for SLES 11 SP4
- Migrating to Offload iSCSI for SLES 12 SP1/SP2

Migrating to Offload iSCSI for RHEL 6.9/6.10

To migrate from a software iSCSI installation to an offload iSCSI for RHEL 6.9 or 6.10:

1. Boot into the iSCSI non-offload/L2 boot from SAN operating system. Issue the following commands to install the Open-iSCSI and iscsiuio RPMs:

```
# rpm -ivh --force qlgc-open-iscsi-2.0_873.111-1.x86_64.rpm
# rpm -ivh --force iscsiuio-2.11.5.2-1.rhel6u9.x86 64.rpm
```

NOTE

To retain the original contents of the inbox RPM during installation, you must use the <code>-ivh</code> option (instead of the <code>-Uvh</code> option), followed by the <code>--force</code> option.

2. Edit the /etc/init.d/iscsid file, add the following command, and then save the file:

```
modprobe -q qedi
```

For example:

```
echo -n $"Starting $prog: "
modprobe -q iscsi_tcp
modprobe -q ib_iser
modprobe -q cxgb3i
modprobe -q cxgb4i
modprobe -q bnx2i
modprobe -q be2iscsi
modprobe -q qedi
daemon iscsiuio
```

3. Open the /boot/efi/EFI/redhat/grub.conf file, make the following changes, and save the file:

```
Remove ifname=eth5:14:02:ec:ce:dc:6d
Remove ip=ibft
Add selinux=0
```

For example:

```
kernel /vmlinuz-2.6.32-696.el6.x86_64 ro
root=/dev/mapper/vg_prebooteit-lv_root rd_NO_LUKS
iscsi_firmware LANG=en_US.UTF-8 ifname=eth5:14:02:ec:ce:dc:6d
rd_NO_MD SYSFONT=latarcyrheb-sun16 crashkernel=auto rd_NO_DM
rd_LVM_LV=vg_prebooteit/lv_swap ip=ibft KEYBOARDTYPE=pc
KEYTABLE=us rd LVM LV=vg prebooteit/lv root rhgb quiet
```

initrd /initramfs-2.6.32-696.el6.x86 64.img

```
kernel /vmlinuz-2.6.32-696.el6.x86_64 ro
root=/dev/mapper/vg_prebooteit-lv_root rd_NO_LUKS
iscsi_firmware LANG=en_US.UTF-8 rd_NO_MD
SYSFONT=latarcyrheb-sun16 crashkernel=auto rd_NO_DM
rd_LVM_LV=vg_prebooteit/lv_swap KEYBOARDTYPE=pc KEYTABLE=us
rd_LVM_LV=vg_prebooteit/lv_root selinux=0
    initrd /initramfs-2.6.32-696.el6.x86 64.img
```

- 4. Build the initramfs file by issuing the following command:
 - # dracut -f
- 5. Reboot the server, and then open the UEFI HII.
- 6. In the HII, disable iSCSI boot from BIOS, and then enable iSCSI HBA or boot for the adapter as follows:
 - a. Select the adapter port, and then select **Device Level Configuration**.
 - b. On the Device Level Configuration page, for the **Virtualization Mode**, select **NPAR**.
 - Open the NIC Partitioning Configuration page and set the iSCSI
 Offload Mode to Enabled. (iSCSI HBA support is on partition 3 for a two--port adapter and on partition 2 for a four-port adapter.)
 - d. Open the **NIC Configuration** menu and set the **Boot Protocol** to **UEFI iSCSI**.
 - e. Open the iSCSI Configuration page and configure iSCSI settings.
- 7. Save the configuration and reboot the server.

The OS can now boot through the offload interface.

Migrating to Offload iSCSI for SLES 11 SP4

To migrate from a software iSCSI installation to an offload iSCSI for SLES 11 SP4:

1. Update Open-iSCSI tools and iscsiuio to the latest available versions by issuing the following commands:

```
# rpm -ivh qlgc-open-iscsi-2.0_873.111.sles11sp4-3.x86_64.rpm --force
# rpm -ivh iscsiuio-2.11.5.3-2.sles11sp4.x86 64.rpm --force
```

NOTE

To retain the original contents of the inbox RPM during installation, you must use the -ivh option (instead of the -uvh option), followed by the --force option.

2.	Edit the /etc/elilo.conf file, make the following changes, and then s the file:		
		Remove the ip=ibft parameter (if present) Add iscsi_firmware	
3.	the /etc/sysconfig/kernel file as follows:		
	a.	Locate the line that begins with <code>INITRD_MODULES</code> . This line will look similar to the following, but may contain different parameters:	
		<pre>INITRD_MODULES="ata_piix ata_generic"</pre>	
		or	
		INITRD_MODULES="ahci"	
	b.	Edit the line by adding ${\tt qedi}$ to the end of the existing line (inside the quotation marks). For example:	
		<pre>INITRD_MODULES="ata_piix ata_generic qedi"</pre>	
		or	
		INITRD_MODULES="ahci qedi"	
	C.	Save the file.	
4.	Edit the /etc/modprobe.d/unsupported-modules file, change the value for allow_unsupported_modules to 1, and then save the file:		
	allo	w_unsupported_modules 1	
5. Locate and delete the		te and delete the following files:	
		<pre>/etc/init.d/boot.d/K01boot.open-iscsi /etc/init.d/boot.open-iscsi</pre>	
6.	Create a backup of initrd, and then build a new initrd by issuing the following commands.		

7. Reboot the server, and then open the UEFI HII.

cd /boot/
mkinitrd

- 8. In the UEFI HII, disable iSCSI boot from BIOS, and then enable iSCSI HBA or boot for the adapter as follows:
 - a. Select the adapter port, and then select **Device Level Configuration**.
 - b. On the Device Level Configuration page, for the **Virtualization Mode**, select **NPAR**.

- Open the NIC Partitioning Configuration page and set the iSCSI
 Offload Mode to Enabled. (iSCSI HBA support is on partition 3 for a two--port adapter and on partition 2 for a four-port adapter.)
- d. Open the **NIC Configuration** menu and set the **Boot Protocol** to **UEFI iSCSI**.
- e. Open the iSCSI Configuration page and configure iSCSI settings.
- 9. Save the configuration and reboot the server.

The OS can now boot through the offload interface.

Migrating to Offload iSCSI for SLES 12 SP1/SP2

To migrate from a software iSCSI installation to an offload iSCSI for SLES 12 SP1/SP2:

1. Boot into the iSCSI non-offload/L2 boot from SAN operating system. Issue the following commands to install the Open-iSCSI and iscsiuio RPMs:

```
# rpm -ivh qlgc-open-iscsi-2.0_873.111.sles12p2-3.x86_64.rpm --force
```

rpm -ivh iscsiuio-2.11.5.3-2.sles12sp2.x86 64.rpm --force

NOTE

To retain the original contents of the inbox RPM during installation, you must use the -ivh option (instead of the -uvh option), followed by the --force option.

- 2. Reload all the daemon services by issuing the following command:
 - # systemctl daemon-reload
- 3. Enable iscsid and iscsiulo services (if they are not already enabled) by issuing the following commands:
 - # systemctl enable iscsid
 - # systemctl enable iscsiuio
- 4. Issue the following command:

cat /proc/cmdline

- 5. Check if the OS has preserved any boot options, such as ip=ibft or rd.iscsi.ibft.
 - ☐ If there are preserved boot options, continue with Step 6.
 - ☐ If there are no preserved boot options, skip to Step 6 c.

- 6. Edit the /etc/default/grub file and modify the GRUB_CMDLINE_LINUX value:
 - a. Remove rd.iscsi.ibft (if present).
 - b. Remove any ip=<value> boot options (if present).
 - c. Add rd.iscsi.firmware. For older distros, add iscsi firmware.
- 7. Create a backup of the original <code>grub.cfg</code> file. The file is in the following locations:
 - ☐ Legacy boot: /boot/grub2/grub.cfg
 - ☐ UEFI boot: /boot/efi/EFI/sles/grub.cfg for SLES
- 8. Create a new grub.cfg file by issuing the following command:
 - # grub2-mkconfig -o <new file name>
- 9. Compare the old <code>grub.cfg</code> file with the new <code>grub.cfg</code> file to verify your changes.
- 10. Replace the original grub.cfg file with the new grub.cfg file.
- 11. Build the initramfs file by issuing the following command:
 - # dracut -f
- 12. Reboot the server, and then open the UEFI HII.
- 13. In the UEFI HII, disable iSCSI boot from BIOS, and then enable iSCSI HBA or boot for the adapter as follows:
 - a. Select the adapter port, and then select **Device Level Configuration**.
 - b. On the Device Level Configuration page, for the **Virtualization Mode**, select **NPAR**.
 - c. Open the NIC Partitioning Configuration page and set the **iSCSI Offload Mode** to **Enabled**. (iSCSI HBA support is on partition 3 for a two--port adapter and on partition 2 for a four-port adapter.)
 - d. Open the **NIC Configuration** menu and set the **Boot Protocol** to **UEFI iSCSI**.
 - e. Open the iSCSI Configuration page and configure iSCSI settings.
- 14. Save the configuration and reboot the server.

The OS can now boot through the offload interface.

Linux Multipath Considerations

iSCSI boot from SAN installations on Linux operating systems are currently supported only in a single-path configuration. To enable multipath configurations, perform the installation in a single path, using either an L2 or L4 path. After the server boots into the installed operating system, perform the required configurations for enabling multipath I/O (MPIO).

See the appropriate procedure in this section to migrate from L2 to L4 and configure MPIO for your OS:

- Migrating and Configuring MPIO to Offloaded Interface for RHEL 6.9/6.10
- Migrating and Configuring MPIO to Offloaded Interface for SLES 11 SP4
- Migrating and Configuring MPIO to Offloaded Interface for SLES 12 SP1/SP2

Migrating and Configuring MPIO to Offloaded Interface for RHEL 6.9/6.10

To migrate from L2 to L4 and configure MPIO to boot the OS over an offloaded interface for RHEL 6.9/6.10:

- Configure the iSCSI boot settings for L2 BFS on both ports of the adapter.
 The boot will log in through only one port, but will create an iSCSI Boot Firmware Table (iBFT) interface for both ports.
- 2. While booting to the CD, ensure that you specify the following kernel parameters:

```
ip=ibft
linux dd
```

- 3. Provide the DUD and complete the installation.
- 4. Boot to the OS with L2.
- 5. Update Open-iSCSI tools and iscsiulo by issuing the following commands:

```
 \begin{tabular}{ll} $\#$ rpm -ivh qlgc-open-iscsi-2.0_873.111.rhel6u9-3.x86_64.rpm --force \\ \end{tabular}
```

- # rpm -ivh iscsiuio-2.11.5.5-6.rhel6u9.x86_64.rpm --force
 - 6. Edit the /boot/efi/EFI/redhat/grub.conf file, make the following changes, and then save the file:
 - a. Remove ifname=eth5:14:02:ec:ce:dc:6d.
 - b. Remove ip=ibft.
 - c. Add selinux=0.
 - 7. Build the initramfs file by issuing the following command:
 - # dracut -f

- 8. Reboot and change the adapter boot settings to use L4 or **iSCSI (HW)** for both ports and to boot through L4.
- 9. After booting into the OS, set up the multipath daemon multipathd.conf:
 - # mpathconf --enable --with_multipathd y
 - # mpathconf -enable
- 10. Start the multipathd service:
 - # service multipathd start
- 11. Rebuild initramfs with multipath support.
 - # dracut --force --add multipath --include /etc/multipath
- 12. Reboot the server and boot into OS with multipath.

NOTE

For any additional changes in the /etc/multipath.conf file to take effect, you must rebuild the initrd image and reboot the server.

1.

2.

Migrating and Configuring MPIO to Offloaded Interface for SLES 11 SP4

To migrate from L2 to L4 and configure MPIO to boot the OS over an offloaded interface for SLES 11 SP4:

- Follow all the steps necessary to migrate a non-offload (L2) interface to an offload (L4) interface, over a single path. See Migrating to Offload iSCSI for SLES 11 SP4.
- 2. After you boot into the OS using the L4 interface, prepare multipathing as follows:
 - a. Reboot the server and open the HII by going to **System Setup/Utilities**.
 - b. In the HII, select **System Configuration**, and select the second adapter port to be used for multipath.
 - c. On the Main Configuration Page under **Port Level Configuration**, set **Boot Mode** to **iSCSI (HW)** and enable **iSCSI Offload**.
 - d. On the Main Configuration Page under **iSCSI Configuration**, perform the necessary iSCSI configurations.
 - e. Reboot the server and boot into OS.

- 3. Set MPIO services to remain persistent on re-boot as follows:
 - # chkconfig multipathd on
 - # chkconfig boot.multipath on
 - # chkconfig boot.udev on
- 4. Start multipath services as follows:
 - # /etc/init.d/boot.multipath start
 - # /etc/init.d/multipathd start
- 5. Run multipath -v2 -d to display multipath configuration with a dry run.
- 6. Locate the multipath.conf file under /etc/multipath.conf.

NOTE

If the file is not present, copy multipath.conf from the
/usr/share/doc/packages/multipath-tools folder:
cp /usr/share/doc/packages/multipath-tools/multipath.
conf.defaults /etc/multipath.conf

- 7. Edit the multipath.conf to enable the default section.
- 8. Rebuild initrd image to include MPIO support:
 - # mkinitrd -f multipath
- 9. Reboot the server and boot the OS with multipath support.

NOTE

For any additional changes in the /etc/multipath.conf file to take effect, you must rebuild the initrd image and reboot the server.

Migrating and Configuring MPIO to Offloaded Interface for SLES 12 SP1/SP2

To migrate from L2 to L4 and configure MPIO to boot the OS over an offloaded interface for SLES 12 SP1/SP2:

- Configure the iSCSI boot settings for L2 BFS on both ports of the adapter.
 The boot will log in through only one port, but will create an iSCSI Boot Firmware Table (iBFT) interface for both ports.
- 2. While booting to the CD, ensure that you specify the following kernel parameters:

```
dud=1
rd.iscsi.ibft
```

3. Provide the DUD and complete the installation.

- 4. Boot to the OS with L2.
- 5. Update the Open-iSCSI tools by issuing the following commands:
- # rpm -ivh qlgc-open-iscsi-2.0 873.111.sles12sp1-3.x86 64.rpm --force
- # rpm -ivh iscsiuio-2.11.5.5-6.sles12sp1.x86_64.rpm --force
 - 6. Edit the /etc/default/grub file by changing the rd.iscsi.ibft parameter to rd.iscsi.firmware, and then save the file.
 - 7. Issue the following command:
 - # grub2-mkconfig -o /boot/efi/EFI/suse/grub.cfg
 - 8. To load the multipath module, issue the following command:
 - # modprobe dm multipath
 - 9. To enable the multipath daemon, issue the following commands:
 - # systemctl start multipathd.service
 - # systemctl enable multipathd.service
 - # systemctl start multipathd.socket
 - 10. To run the multipath utility, issue the following commands:
 - # multipath (may not show the multipath devices because it is booted with a single path on L2)
 - # multipath -11
 - 11. To inject the multipath module in initrd, issue the following command:
 - # dracut --force --add multipath --include /etc/multipath
 - 12. Reboot the server and enter system settings by pressing the F9 key during the POST menu.
 - 13. Change the UEFI configuration to use L4 iSCSI boot:
 - a. Open System Configuration, select the adapter port, and then select **Port Level Configuration**.
 - b. On the Port Level Configuration page, set the Boot Mode to iSCSI (HW) and set iSCSI Offload to Enabled.
 - c. Save the settings, and then exit the System Configuration Menu.
 - 14. Reboot the system. The OS should now boot through the offload interface.

NOTE

For any additional changes in the /etc/multipath.conf file to take effect, you must rebuild the initrd image and reboot the server.

Configuring iSCSI Boot from SAN on VMware

Because VMware does not natively support iSCSI boot from SAN offload, you must configure BFS through the software in preboot, and then transition to offload upon OS driver loads. For more information, see "Enabling NPAR and the iSCSI HBA" on page 71.

In VMware ESXi, iSCSI BFS configuration procedures include:

- Setting the UEFI Main Configuration
- Configuring the System BIOS for iSCSI Boot (L2)
- Mapping the CD or DVD for OS Installation

Setting the UEFI Main Configuration

To configure iSCSI boot from SAN on VMware:

- 1. Plug the 41xxx Series Adapter into a Dell 14G server. For example, plug a PCIE and LOM (four ports or two ports) into an R740 server.
- 2. In the HII, go to **System Setup**, select **Device Settings**, and then select a an integrated NIC port to configure. Click **Finish**.
- 3. On the Main Configuration Page, select NIC Partitioning Configuration, and then click Finish.
- 4. On the **Main Configuration Page**, select **Firmware Image Properties**, view the non-configurable information, and then click **Finish**.
- 5. On the Main Configuration Page, select Device Level Configuration.
- 6. Complete the Device Level Configuration page (see Figure 6-19) as follows:
 - For Virtualization Mode, select either None, NPar, or NPar_EP for IBFT installation through the NIC interface.
 - b. For **NParEP Mode**, select **Disabled**.
 - c. For **UEFI Driver Debug Leve**l, select 10.

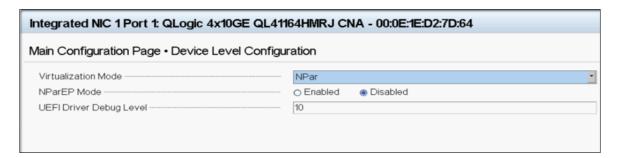


Figure 6-19. Integrated NIC: Device Level Configuration for VMware

- 7. Go to the **Main Configuration Page** and select **NIC Partitioning Configuration**.
- 8. On the NIC Partitioning Configuration page, select **Partition 1 Configuration**.
- 9. Complete the Partition 1 Configuration page as follows:
 - a. For Link Speed, select either Auto Neg, 10Gbps, or 1Gbps.
 - b. Ensure that the link is up.
 - c. For **Boot Protocol**, select **None**.
 - d. For Virtual LAN Mode, select Disabled.
- 10. On the NIC Partitioning Configuration page, select **Partition 2 Configuration**.
- 11. Complete the Partition 2 Configuration page (see Figure 6-20) as follows:
 - a. For FCoE Mode, select Disabled.
 - b. For iSCSI Offload Mode, select Disabled.

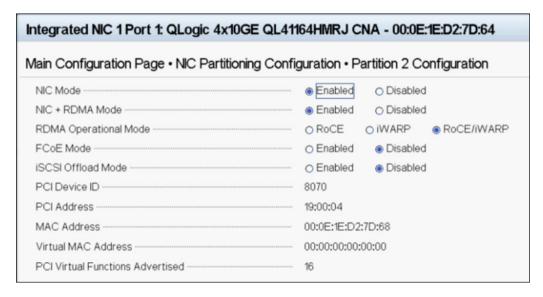


Figure 6-20. Integrated NIC: Partition 2 Configuration for VMware

Configuring the System BIOS for iSCSI Boot (L2)

To configure the System BIOS on VMware:

- 1. On the System BIOS Settings page, select **Boot Settings**.
- 2. Complete the Boot Settings page as shown in Figure 6-21.



Figure 6-21. Integrated NIC: System BIOS, Boot Settings for VMware

- 3. On the System BIOS Settings page, select **Network Settings**.
- 4. On the Network Settings page under **UEFI iSCSI Settings**:
 - a. For iSCSI Device1, select Enabled.
 - b. Select **UEFI Boot Settings**.
- 5. On the iSCSI Device1 Settings page:
 - a. For Connection 1, select Enabled.
 - b. Select Connection 1 Settings.
- 6. On the Connection 1 Settings page (see Figure 6-22):
 - a. For **Interface**, select the adapter port on which to test the iSCSI boot firmware table (IBFT) boot from SAN.
 - b. For **Protocol**, select either **IPv4** or **IPv6**.
 - c. For **VLAN**, select either **Disabled** (the default) or **Enabled** (if you want to test IBFT with a vLAN).
 - d. For **DHCP**, select **Enabled** (if the IP address is from the DHCP server) or **Disabled** (to use static IP configuration).
 - e. For Target info via DHCP, select Disabled.

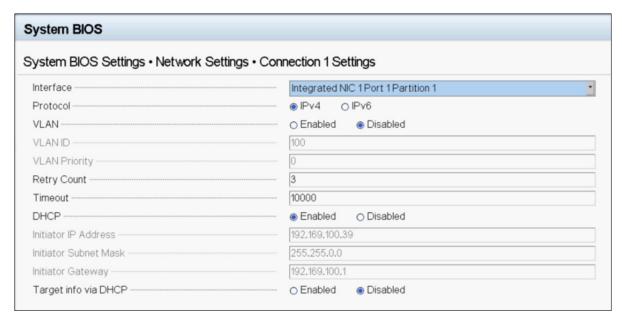


Figure 6-22. Integrated NIC: System BIOS, Connection 1 Settings for VMware

 Complete the target details, and for Authentication Type, select either CHAP (to set CHAP details) or None (the default). Figure 6-23 shows an example.

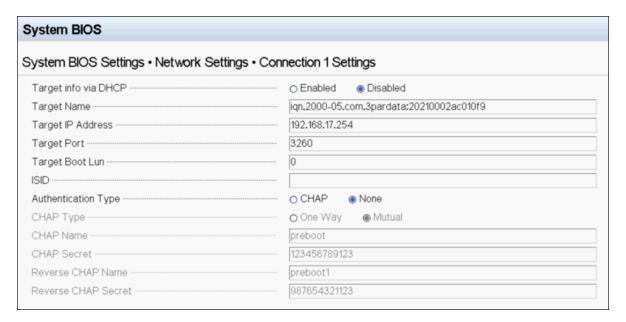


Figure 6-23. Integrated NIC: System BIOS, Connection 1 Settings (Target) for VMware

- 8. Save all configuration changes, and then reboot the server.
- 9. During system boot up, press the F11 key to start the Boot Manager.
- 10. In the Boot Manager under **Boot Menu**, **Select UEFI Boot Option**, select the **Embedded SATA Port AHCI Controller**.

Mapping the CD or DVD for OS Installation

To map the CD or DVD:

- 1. Create a customized ISO image using the ESXi-Customizer and inject the latest bundle or VIB.
- 2. Map the ISO to the server virtual console's virtual media.
- 3. On the virtual optical drive, load the ISO file.
- 4. After the ISO is loaded successfully, press the F11 key.
- 5. On the Select a Disk To Install Or Upgrade window, under **Storage Device**, select the **3PARdata W** disk, and then press the ENTER key. Figure 6-24 shows an example.



Figure 6-24. VMware iSCSI BFS: Selecting a Disk to Install

6. Start installation of the ESXi OS on the remote iSCSI LUN.

7. After the ESXi OS installation completes successfully, the system boots to the OS, as shown in Figure 6-25.

```
VMware ESXi 6.7.8 (VMKernel Release Build 18382608)

Dell Inc. PowerEdge R740

2 x Intel(R) Xeon(R) Gold 5118 CPU 9 2.30GHz
63.7 GiB Menory

To manage this host go to:
http://192.168.28.234/ (Haiting for DHCP...)
http://Ife88::28e:leff:fed2:7d641/ (STATIC)

Harning: DHCP lookup failed. You may be unable to access this system until you customize its
network configuration.

(F2) Customize System/View Logs
```

Figure 6-25. VMware iSCSI Boot from SAN Successful

FCoE Boot from SAN

Marvell 41xxx Series Adapters support FCoE boot to enable network boot of operating systems to diskless systems. FCoE boot allows a Windows, Linux, or VMware operating system to boot from a Fibre Channel or FCoE target machine located remotely over an FCoE supporting network. You can set the FCoE option (offload path with Marvell offload FCoE driver) by opening the **NIC Configuration** menu and setting the **Boot Protocol** option to **FCoE**.

This section provides the following configuration information about FCoE boot from SAN:

- FCoE Out-of-Box and Inbox Support
- **■** FCoE Preboot Configuration
- Configuring FCoE Boot from SAN on Windows
- Configuring FCoE Boot from SAN on Linux
- Configuring FCoE Boot from SAN on VMware

FCoE Out-of-Box and Inbox Support

Table 6-6 lists the operating systems' inbox and out-of-box support for FCoE boot from SAN (BFS).

Table 6-6. FCoE Out-of-Box and Inbox Boot from SAN Support

OS Version	Out-of-Box Hardware Offload FCoE BFS Support	Inbox Hardware Offload FCoE BFS Support
Windows 2012	Yes	No
Windows 2012 R2	Yes	No
Windows 2016	Yes	No
Windows 2019	Yes	Yes
RHEL 7.5	Yes	Yes
RHEL 7.6	Yes	Yes
RHEL 8.0	Yes	Yes
SLES 15/15 SP1	Yes	Yes
vSphere ESXi 6.5 U3	Yes	No
vSphere ESXi 6.7 U2	Yes	No

FCoE Preboot Configuration

This section describes the installation and boot procedures for the Windows, Linux, and ESXi operating systems. To prepare the system BIOS, modify the system boot order and specify the BIOS boot protocol, if required.

NOTE

FCoE boot from SAN is supported on ESXi 5.5 and later. Not all adapter versions support FCoE and FCoE boot from SAN.

Specifying the BIOS Boot Protocol

FCoE boot from SAN is supported in UEFI mode only. Set the platform in boot mode (protocol) using the system BIOS configuration to UEFI.

NOTE

FCoE BFS is not supported in legacy BIOS mode.

114

Configuring Adapter UEFI Boot Mode

To configure the boot mode to FCOE:

- 1. Restart the system.
- 2. Press the OEM hot key to enter System Setup (Figure 6-26). This is also known as UEFI HII.



Figure 6-26. System Setup: Selecting Device Settings

NOTE

SAN boot is supported in the UEFI environment only. Make sure the system boot option is UEFI, and not legacy.

3. On the Device Settings page, select the Marvell FastLinQ adapter (Figure 6-27).

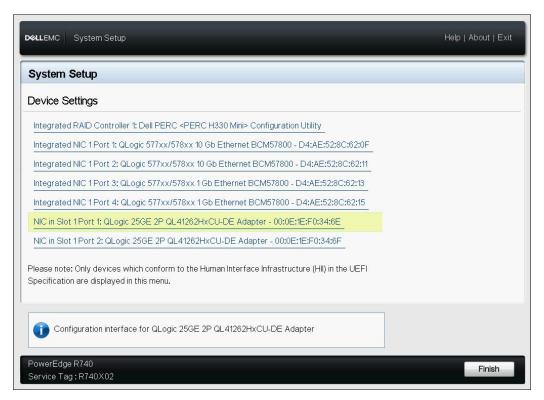


Figure 6-27. System Setup: Device Settings, Port Selection

4. On the Main Configuration Page, select **NIC Configuration** (Figure 6-28), and then press ENTER.

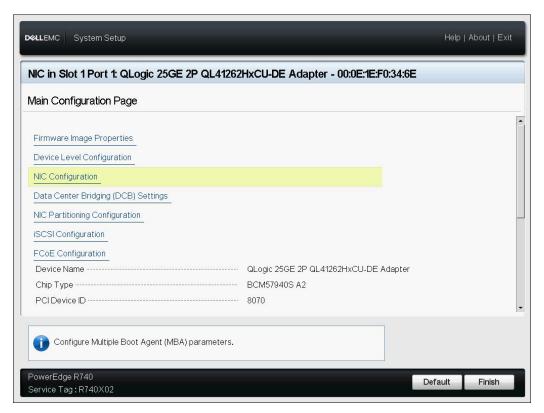


Figure 6-28. System Setup: NIC Configuration

5. On the NIC Configuration page, select **Boot Mode**, press ENTER, and then select **FCoE** as a preferred boot mode.

NOTE

FCoE is not listed as a boot option if the **FCoE Mode** feature is disabled at the port level. If the **Boot Mode** preferred is **FCoE**, make sure the **FCoE Mode** feature is enabled as shown in Figure 6-29. Not all adapter versions support FCoE.

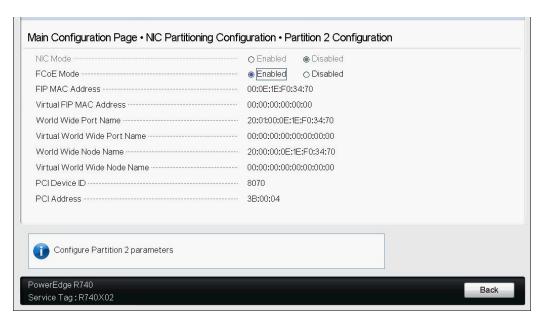


Figure 6-29. System Setup: FCoE Mode Enabled

To configure the FCoE boot parameters:

- 1. On the Device UEFI HII Main Configuration Page, select **FCoE Configuration**, and then press ENTER.
- 2. On the FCoE Configuration Page, select **FCoE General Parameters**, and then press ENTER.
- 3. On the FCoE General Parameters page (Figure 6-30), press the UP ARROW and DOWN ARROW keys to select a parameter, and then press ENTER to select and input the following values:
 - ☐ Fabric Discovery Retry Count: Default value or as required
 - □ LUN Busy Retry Count: Default value or as required

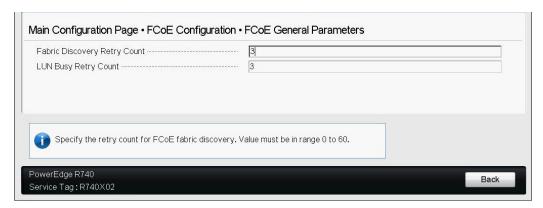


Figure 6-30. System Setup: FCoE General Parameters

- 4. Return to the FCoE Configuration page.
- 5. Press ESC, and then select **FCoE Target Parameters**.
- 6. Press ENTER.
- 7. In the **FCoE General Parameters Menu**, enable **Connect** to the preferred FCoE target.
- 8. Type values for the following parameters (Figure 6-31) for the FCoE target, and then press ENTER:
 - □ World Wide Port Name Target *n*
 - \Box Boot LUN n

Where the value of *n* is between 1 and 8, enabling you to configure 8 FCoE targets.

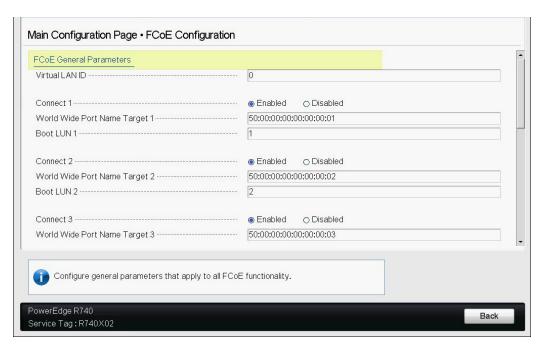


Figure 6-31. System Setup: FCoE General Parameters

Configuring FCoE Boot from SAN on Windows

FCoE boot from SAN information for Windows includes:

- Windows Server 2012 R2 and 2016 FCoE Boot Installation
- Configuring FCoE on Windows
- FCoE Crash Dump on Windows
- Injecting (Slipstreaming) Adapter Drivers into Windows Image Files

Windows Server 2012 R2 and 2016 FCoE Boot Installation

For Windows Server 2012R2/2016 boot from SAN installation, Marvell requires the use of a "slipstream" DVD, or ISO image, with the latest Marvell drivers injected. See "Injecting (Slipstreaming) Adapter Drivers into Windows Image Files" on page 122.

The following procedure prepares the image for installation and booting in FCoE mode.

To set up Windows Server 2012R2/2016 FCoE boot:

- 1. Remove any local hard drives on the system to be booted (remote system).
- 2. Prepare the Windows OS installation media by following the slipstreaming steps in "Injecting (Slipstreaming) Adapter Drivers into Windows Image Files" on page 122.

- Load the latest Marvell FCoE boot images into the adapter NVRAM.
- 4. Configure the FCoE target to allow a connection from the remote device. Ensure that the target has sufficient disk space to hold the new OS installation.
- 5. Configure the UEFI HII to set the FCoE boot type on the required adapter port, correct initiator, and target parameters for FCoE boot.
- 6. Save the settings and reboot the system. The remote system should connect to the FCoE target, and then boot from the DVD-ROM device.
- 7. Boot from DVD and begin installation.
- 8. Follow the on-screen instructions.
 - On the window that shows the list of disks available for the installation, the FCoE target disk should be visible. This target is a disk connected through the FCoE boot protocol, located in the remote FCoE target.
- 9. To proceed with Windows Server 2012R2/2016 installation, select **Next**, and then follow the on-screen instructions. The server will undergo a reboot multiple times as part of the installation process.
- 10. After the server boots to the OS, you should run the driver installer to complete the Marvell drivers and application installation.

Configuring FCoE on Windows

By default, DCB is enabled on Marvell FastLinQ 41xxx FCoE- and DCB-compatible C-NICs. Marvell 41xxx FCoE requires a DCB-enabled interface. For Windows operating systems, use QConvergeConsole GUI or a command line utility to configure the DCB parameters.

FCoE Crash Dump on Windows

Crash dump functionality is currently supported for FCoE boot for the FastLinQ 41xxx Series Adapters.

No additional configuration is required for FCoE crash-dump generation when in FCoE boot mode.

Injecting (Slipstreaming) Adapter Drivers into Windows Image Files

To inject adapter drivers into the Windows image files:

- 1. Obtain the latest driver package for the applicable Windows Server version (2012, 2012 R2, 2016, or 2019).
- 2. Extract the driver package to a working directory:
 - a. Open a command line session and navigate to the folder that contains the driver package.
 - b. To extract the driver Dell Update Package (DUP), issue the following command:
 - start /wait NameOfDup.exe /s /drivers=<folder path>
- 3. Download the Windows Assessment and Deployment Kit (ADK) version 10 from Microsoft:
 - https://developer.microsoft.com/en-us/windows/hardware/windows-assessment-deployment-kit
- 4. Follow the Microsoft "Add and Remove Drivers to an offline Windows Image article" and inject the OOB driver extracted on Step 2, part b. See https://docs.microsoft.com/en-us/windows-hardware/manufacture/desktop/add-and-remove-drivers-to-an-offline-windows-image

Configuring FCoE Boot from SAN on Linux

FCoE boot from SAN configuration for Linux covers the following:

- Prerequisites for Linux FCoE Boot from SAN
- Configuring Linux FCoE Boot from SAN

Prerequisites for Linux FCoE Boot from SAN

The following are required for Linux FCoE boot from SAN to function correctly with the Marvell FastLinQ 41xxx 10/25GbE Controller.

General

You no longer need to use the FCoE disk tabs in the Red Hat and SUSE installers because the FCoE interfaces are not exposed from the network interface and are automatically activated by the qedf driver.

SLES 12 and SLES 15

- Driver update disk is recommended for SLES 12 SP 3 and later.
- The installer parameter dud=1 is required to ensure that the installer will ask for the driver update disk.

■ Do not use the installer parameter withfcoe=1 because the software FCoE will conflict with the hardware offload if network interfaces from qede are exposed.

Configuring Linux FCoE Boot from SAN

This section provides FCoE boot from SAN procedures for the following Linux distributions:

- Configuring FCoE Boot from SAN for SLES 12 SP3 and Later
- Using an FCoE Boot Device as a kdump Target

Configuring FCoE Boot from SAN for SLES 12 SP3 and Later

No additional steps, other than injecting DUD for out-of-box driver, are necessary to perform boot from SAN installations when using SLES 12 SP3.

Using an FCoE Boot Device as a kdump Target

When using a device exposed using the qedf driver as a kdump target for crash dumps, Marvell recommends that you specify the kdump <code>crashkernel</code> memory parameter at the kernel command line to be a minimum of 512MB. Otherwise, the kernel crash dump may not succeed.

For details on how to set the kdump <code>crashkernel</code> size, refer to your Linux distribution documentation.

Configuring FCoE Boot from SAN on VMware

For VMware ESXi 6.5/6.7 boot from SAN installation, Marvell requires that you use a customized ESXi ISO image that is built with the latest Marvell Converged Network Adapter bundle injected. This section covers the following VMware FCoE boot from SAN procedures.

- Injecting (Slipstreaming) ESXi Adapter Drivers into Image Files
- Installing the Customized ESXi ISO

Injecting (Slipstreaming) ESXi Adapter Drivers into Image Files

This procedure uses ESXi-Customizer tool v2.7.2 as an example, but you can use any ESXi customizer.

To inject adapter drivers into an ESXi image file:

- 1. Download ESXi-Customizer v2.7.2 or later.
- 2. Go to the ESXi customizer directory.
- 3. Issue the ESXi-Customizer.cmd command.

- 4. In the ESXi-Customizer dialog box, click **Browse** to complete the following.
 - a. Select the original VMware ESXi ISO file.
 - b. Select either the Marvell FCoE driver VIB file or the Marvell offline gedenty bundle ZIP file.
 - c. For the working directory, select the folder in which the customized ISO needs to be created.
 - d. Click Run.

Figure 6-32 shows an example.

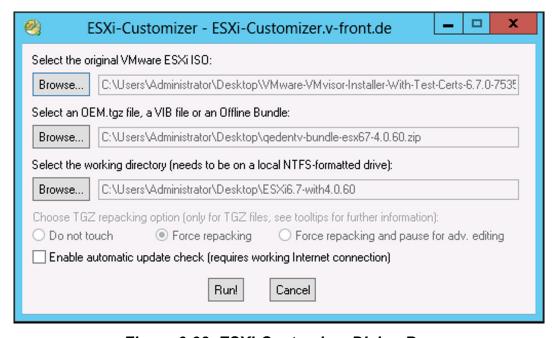


Figure 6-32. ESXi-Customizer Dialog Box

- 5. Burn a DVD that contains the customized ISO build located in the working directory specified in Step 4c.
- 6. Use the new DVD to install the ESXi OS.

Installing the Customized ESXi ISO

- 1. Load the latest Marvell FCOE boot images into the adapter NVRAM.
- 2. Configure the FCOE target to allow a valid connection with the remote machine. Ensure that the target has sufficient free disk space to hold the new OS installation.
- 3. Configure the UEFI HII to set the FCOE boot type on the required adapter port, the correct initiator, and the target parameters for FCOE boot.

4. Save the settings and reboot the system.

The initiator should connect to an FCOE target and then boot the system from the DVD-ROM device.

- 5. Boot from the DVD and begin installation.
- 6. Follow the on-screen instructions.

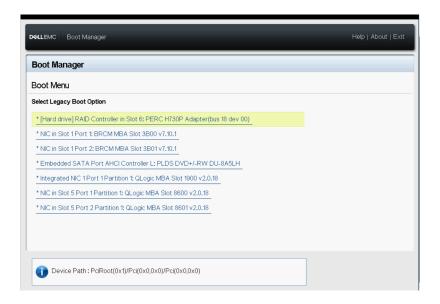
On the window that shows the list of disks available for the installation, the FCOE target disk should be visible because the injected Converged Network Adapter bundle is inside the customized ESXi ISO. Figure 6-33 shows an example.

```
Select a Disk to Install or Upgrade
     (any existing VMFS-3 will be automatically upgraded to VMFS-5)
* Contains a VMFS partition
# Claimed by VMware vSAN
  DGC
            RAID 5
                              (naa.600601602d9036008a096...)
                                                                2.00 GiB
                              (naa.600601602d9036008b096...)
  DGC
            RAID 5
                                                                2.00 GiB
            RAID 5
                              (naa.600601602d9036008c096...)
  DGC
                                                                2.00 GiB
                              (naa.600601602d9036008d096...)
  DGC
            RAID 5
                                                                2.00 GiB
  DGC
            RAID 5
                              (naa.600601602d9036008e096...)
                                                                2.00 GiB
            RAID 5
                              (naa.600601602d9036008f096...)
  DGC
                                                                2.00 GiB
                              (naa.600601602d903600d6449...)
  DGC
            RAID 5
                                                                2.00 GiB
  DGC
            RAID 5
                              (naa.600601602d903600d7449...)
                                                                2.00 GiB
```

Figure 6-33. Select a VMware Disk to Install

- 7. Select the LUN on which ESXi can install, and then press ENTER.
- 8. On the next window, click **Next**, and then follow the on-screen instructions.
- 9. When installation completes, reboot the server and eject the DVD.
- 10. During the server boot, press the F9 key to access the **One-Time Boot Menu**, and then select **Boot media to QLogic adapter port**.
- 11. Under **Boot Menu**, select the newly installed ESXi to load through boot from SAN.

Figure 6-34 provides two examples.



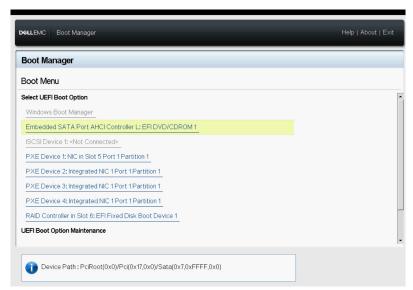


Figure 6-34. VMware USB Boot Options

7 RoCE Configuration

This chapter describes RDMA over converged Ethernet (RoCE v1 and v2) configuration on the 41xxx Series Adapter, the Ethernet switch, and the Windows, Linux, or VMware host, including:

- Supported Operating Systems and OFED
- "Planning for RoCE" on page 128
- "Preparing the Adapter" on page 129
- "Preparing the Ethernet Switch" on page 129
- "Configuring RoCE on the Adapter for Windows Server" on page 133
- "Configuring RoCE on the Adapter for Linux" on page 150
- "Configuring RoCE on the Adapter for VMware ESX" on page 164
- "Configuring DCQCN" on page 171

NOTE

Some RoCE features may not be fully enabled in the current release.

Supported Operating Systems and OFED

Table 7-1 shows the operating system support for RoCE v1, RoCE v2, iWARP, and OpenFabrics Enterprise Distribution (OFED). OFED is not supported on Windows or VMware ESXi.

Table 7-1. OS Support for RoCE v1, RoCE v2, iWARP, iSER, and OFED

Operating System	Inbox	OFED-4.17-1 GA		
Windows Server 2012	N/A	N/A		
Windows Server 2012 R2	No	N/A		
Windows Server 2016	No	N/A		
Windows Server 2019	RoCE v1, RoCE v2, iWARP	N/A		
RHEL 7.6	RoCE v1, RoCE v2, iWARP, iSER	RoCE v1, RoCE v2, iWARP		

Table 7-1. OS Support for RoCE v1, RoCE v2, iWARP, iSER, and OFED (Continued)

Operating System	Inbox	OFED-4.17-1 GA
RHEL 7.7	RoCE v1, RoCE v2, iWARP, iSER	No
RHEL 8.0	RoCE v1, RoCE v2, iWARP, iSER	No
RHEL 8.1	RoCE v1, RoCE v2, iWARP, iSER	No
SLES 12 SP4	RoCE v1, RoCE v2, iWARP, iSER	RoCE v1, RoCE v2, iWARP
SLES 15 SP0	RoCE v1, RoCE v2, iWARP, iSER	RoCE v1, RoCE v2, iWARP
SLES 15 SP1	RoCE v1, RoCE v2, iWARP, iSER	No
CentOS 7.6	RoCE v1, RoCE v2, iWARP, iSER	RoCE v1, RoCE v2, iWARP
VMware ESXi 6.5 U3	RoCE v1, RoCE v2	N/A
VMware ESXi 6.7 U2	RoCE v1, RoCE v2	N/A

Planning for RoCE

As you prepare to implement RoCE, consider the following limitations:

- If you are using the inbox OFED, the operating system should be the same on the server and client systems. Some of the applications may work between different operating systems, but there is no guarantee. This is an OFED limitation.
- For OFED applications (most often perftest applications), server and client applications should use the same options and values. Problems can arise if the operating system and the perftest application have different versions. To confirm the perftest version, issue the following command:
 - # ib_send_bw --version
- Building libqedr in inbox OFED requires installing libibverbs-devel.
- Running user space applications in inbox OFED requires installing the InfiniBand[®] Support group, by yum groupinstall "InfiniBand Support" that contains libibom, libibverbs, and more.
- OFED and RDMA applications that depend on libibverbs also require the Marvell RDMA user space library, libqedr. Install libqedr using the libqedr RPM or source packages.
- RoCE supports only little endian.

Preparing the Adapter

Follow these steps to enable DCBX and specify the RoCE priority using the HII management application. For information about the HII application, see Chapter 5 Adapter Preboot Configuration.

To prepare the adapter:

- 1. On the Main Configuration Page, select **Data Center Bridging (DCB) Settings**, and then click **Finish**.
- 2. In the Data Center Bridging (DCB) Settings window, click the **DCBX Protocol** option. The 41xxx Series Adapter supports both CEE and IEEE protocols. This value should match the corresponding value on the DCB switch. In this example, select **CEE** or **Dynamic**.
- 3. In the **RoCE Priority** box, type a priority value. This value should match the corresponding value on the DCB switch. In this example, type 5. Typically, 0 is used for the default lossy traffic class, 3 is used for the FCoE traffic class, and 4 is used for lossless iSCSI-TLV over DCB traffic class.
- 4. Click Back.
- 5. When prompted, click **Yes** to save the changes. Changes will take effect after a system reset.

For Windows, you can configure DCBX using the HII or QoS method. The configuration shown in this section is through HII. For QoS, refer to "Configuring QoS for RoCE" on page 259.

Preparing the Ethernet Switch

This section describes how to configure a Cisco® Nexus® 6000 Ethernet Switch and a Dell Z9100 Ethernet Switch for RoCE.

- Configuring the Cisco Nexus 6000 Ethernet Switch
- Configuring the Dell Z9100 Ethernet Switch for RoCE

Configuring the Cisco Nexus 6000 Ethernet Switch

Steps for configuring the Cisco Nexus 6000 Ethernet Switch for RoCE include configuring class maps, configuring policy maps, applying the policy, and assigning a vLAN ID to the switch port.

To configure the Cisco switch:

1. Open a config terminal session as follows:

```
Switch# config terminal
switch(config)#
```

2. Configure the quality of service (QoS) class map and set the RoCE priority (cos) to match the adapter (5) as follows:

```
switch(config) # class-map type qos class-roce
switch(config) # match cos 5
```

3. Configure queuing class maps as follows:

```
switch(config) # class-map type queuing class-roce
switch(config) # match qos-group 3
```

4. Configure network QoS class maps as follows:

```
switch(config) # class-map type network-qos class-roce
switch(config) # match qos-group 3
```

5. Configure QoS policy maps as follows:

```
switch(config) # policy-map type qos roce
switch(config) # class type qos class-roce
switch(config) # set qos-group 3
```

6. Configure queuing policy maps to assign network bandwidth. In this example, use a value of 50 percent:

```
switch(config) # policy-map type queuing roce
switch(config) # class type queuing class-roce
switch(config) # bandwidth percent 50
```

7. Configure network QoS policy maps to set priority flow control for no-drop traffic class as follows:

```
switch(config) # policy-map type network-qos roce
switch(config) # class type network-qos class-roce
switch(config) # pause no-drop
```

8. Apply the new policy at the system level as follows:

```
switch(config)# system qos
switch(config)# service-policy type qos input roce
switch(config)# service-policy type queuing output roce
switch(config)# service-policy type queuing input roce
switch(config)# service-policy type network-qos roce
```

9. Assign a vLAN ID to the switch port to match the vLAN ID assigned to the adapter (5).

```
switch(config)# interface ethernet x/x
switch(config)# switchport mode trunk
switch(config)# switchport trunk allowed vlan 1,5
```

Configuring the Dell Z9100 Ethernet Switch for RoCE

Configuring the Dell Z9100 Ethernet Switch for RoCE comprises configuring a DCB map for RoCE, configuring priority-based flow control (PFC) and enhanced transmission selection (ETS), verifying the DCB map, applying the DCB map to the port, verifying PFC and ETS on the port, specifying the DCB protocol, and assigning a VLAN ID to the switch port.

NOTE

For instructions on configuring a Dell Z91000 switch port to connect to the 41xxx Series Adapter at 25Gbps, see "Dell Z9100 Switch Configuration" on page 298.

To configure the Dell switch:

1. Create a DCB map.

```
Dell# configure
Dell(conf)# dcb-map roce
Dell(conf-dcbmap-roce)#
```

2. Configure two ETS traffic classes in the DCB map with 50 percent bandwidth assigned for RoCE (group 1).

```
Dell(conf-dcbmap-roce) # priority-group 0 bandwidth 50 pfc off
Dell(conf-dcbmap-roce) # priority-group 1 bandwidth 50 pfc on
```

3. Configure the PFC priority to match the adapter traffic class priority number (5).

```
Dell(conf-dcbmap-roce) # priority-pgid 0 0 0 0 1 0 0
```

4. Verify the DCB map configuration priority group.

5. Apply the DCB map to the port.

```
Dell(conf) # interface twentyFiveGigE 1/8/1
Dell(conf-if-tf-1/8/1) # dcb-map roce
```

6. Verify the ETS and PFC configuration on the port. The following examples show summarized interface information for ETS and detailed interface information for PFC.

Dell(conf-if-tf-1/8/1) # do show interfaces twentyFiveGigE 1/8/1 ets summary
Interface twentyFiveGigE 1/8/1
Max Supported TC is 4
Number of Traffic Classes is 8
Admin mode is on
Admin Parameters:

Admin is enabled

PG-grp	Priority#	BW-%		BW-COMMITTED		BW-PEAK	TSA	
	90	Rate	(Mbps)	Burst(KB)	Rate(Mbps)	Burst(KB)		
0	0,1,2,3,4,6,7	40	_	_	_		ETS	
1	5	60	_	-	-	-	ETS	
2		-	_	-	-	-	-	
3		_	_	_	_	_	_	

Dell(Conf) # do show interfaces twentyFiveGigE 1/8/1 pfc detail

Interface twentyFiveGigE 1/8/1

Admin mode is on
Admin is enabled, Priority list is 5
Remote is enabled, Priority list is 5
Remote Willing Status is enabled
Local is enabled, Priority list is 5
Oper status is init
PFC DCBX Oper status is Up
State Machine Type is Feature
TLV Tx Status is enabled
PFC Link Delay 65535 pause quntams
Application Priority TLV Parameters:

FCOE TLV Tx Status is disabled

```
ISCSI TLV Tx Status is enabled
Local FCOE PriorityMap is 0x0
Local ISCSI PriorityMap is 0x20
Remote ISCSI PriorityMap is 0x200
```

- 66 Input TLV pkts, 99 Output TLV pkts, 0 Error pkts, 0 Pause Tx pkts, 0 Pause Rx pkts
- 66 Input Appln Priority TLV pkts, 99 Output Appln Priority TLV pkts, 0 Error Appln Priority TLV Pkts
 - 7. Configure the DCBX protocol (CEE in this example).

```
Dell(conf)# interface twentyFiveGigE 1/8/1
Dell(conf-if-tf-1/8/1)# protocol lldp
Dell(conf-if-tf-1/8/1-lldp)# dcbx version cee
```

8. Assign a VLAN ID to the switch port to match the VLAN ID assigned to the adapter (5).

```
Dell(conf) # interface vlan 5
Dell(conf-if-vl-5) # tagged twentyFiveGigE 1/8/1
```

NOTE

The VLAN ID need not be the same as the RoCE traffic class priority number. However, using the same number makes configurations easier to understand.

Configuring RoCE on the Adapter for Windows Server

Configuring RoCE on the adapter for Windows Server host comprises enabling RoCE on the adapter and verifying the Network Direct MTU size.

To configure RoCE on a Windows Server host:

- 1. Enable RoCE on the adapter.
 - a. Open the Windows Device Manager, and then open the 41xxx Series Adapters NDIS Miniport Properties.
 - b. On the QLogic FastLinQ Adapter Properties, click the **Advanced** tab.
 - c. On the Advanced page, configure the properties listed in Table 7-2 by selecting each item under **Property** and choosing an appropriate **Value** for that item. Then click **OK**.

Table 7-2. Advanced Properties for RoCE

Property	Value or Description			
NetworkDirect Functionality	Enabled			
Network Direct Mtu Size	The network direct MTU size must be less than the jumbo packet size.			
Quality of Service	For RoCE v1/v2, always select Enabled to allow Windows DCB-QoS service to control and monitor DCB. For more information, see "Configuring QoS by Disabling DCBX on the Adapter" on page 260 and "Configuring QoS by Enabling DCBX on the Adapter" on page 264.			
NetworkDirect Technology	RoCE or RoCE v2.			
VLAN ID	Assign any vLAN ID to the interface. The value must be the same as is assigned on the switch.			

Figure 7-1 shows an example of configuring a property value.

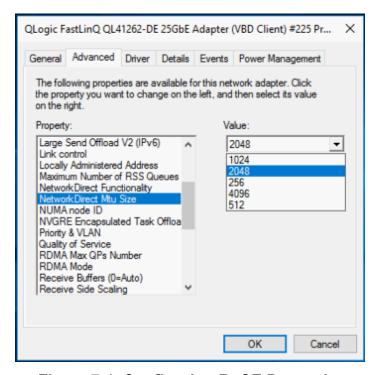


Figure 7-1. Configuring RoCE Properties

2. Using Windows PowerShell, verify that RDMA is enabled on the adapter. The Get-NetAdapterRdma command lists the adapters that support RDMA—both ports are enabled.

NOTE

If you are configuring RoCE over Hyper-V, do not assign a vLAN ID to the physical interface.

PS C:\Users\Admini	strator> Get-NetAdapterRdma	
Name	InterfaceDescription	Enabled
SLOT 4 3 Port 1	QLogic FastLinQ QL41262	True
SLOT 4 3 Port 2	QLogic FastLinQ QL41262	True

3. Using Windows PowerShell, verify that NetworkDirect is enabled on the host operating system. The Get-NetOffloadGlobalSetting command shows NetworkDirect is enabled.

```
PS C:\Users\Administrators> Get-NetOffloadGlobalSetting
ReceiveSideScaling : Enabled
ReceiveSegmentCoalescing : Enabled
Chimney : Disabled
```

TaskOffload : Enabled
NetworkDirect : Enabled
NetworkDirectAcrossIPSubnets : Blocked
PacketCoalescingFilter : Disabled

4. Connect a server message block (SMB) drive, run RoCE traffic, and verify the results.

To set up and connect to an SMB drive, view the information available online from Microsoft:

https://technet.microsoft.com/en-us/library/hh831795(v=ws.11).aspx

5. By default, Microsoft's SMB Direct establishes two RDMA connections per port, which provides good performance, including line rate at a higher block size (for example, 64KB). To optimize performance, you can change the quantity of RDMA connections per RDMA interface to four (or greater).

To increase the quantity of RDMA connections to four (or more), issue the following command in Windows PowerShell:

```
PS C:\Users\Administrator> Set-ItemProperty -Path
"HKLM:\SYSTEM\CurrentControlSet\Services\LanmanWorkstation\
Parameters" ConnectionCountPerRdmaNetworkInterface -Type
DWORD -Value 4 -Force
```

Viewing RDMA Counters

The following procedure also applies to iWARP.

To view RDMA counters for RoCE:

- 1. Launch Performance Monitor.
- 2. Open the Add Counters dialog box. Figure 7-2 shows an example.

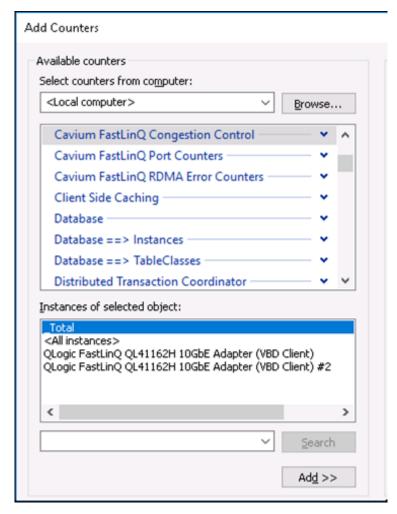


Figure 7-2. Add Counters Dialog Box

NOTE

If Marvell RDMA counters are not listed in the Performance Monitor Add Counters dialog box, manually add them by issuing the following command from the driver location:

Lodctr /M:qend.man

3. Select one of the following counter types:

Cavium FastLinQ Congestion Control:

- Increment when there is congestion in the network and ECN is enabled on the switch.
- Describe RoCE v2 ECN Marked Packets and Congestion Notification Packets (CNPs) sent and received successfully.
- Apply only to RoCE v2.

□ Cavium FastLinQ Port Counters:

- Increment when there is congestion in the network.
- Pause counters increment when flow control or global pause is configured and there is a congestion in the network.
- PFC counters increment when priority flow control is configured and there is a congestion in the network.

□ Cavium FastLinQ RDMA Error Counters:

- Increment if any error occurs in transport operations.
- For details, see Table 7-3.
- 4. Under Instances of selected object, select Total, and then click Add.

Figure 7-3 shows three examples of the counter monitoring output.

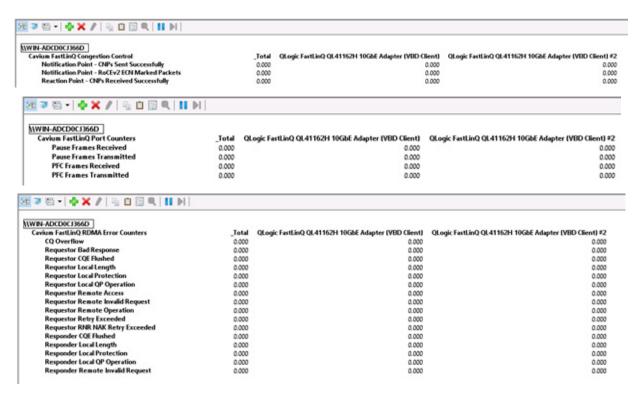


Figure 7-3. Performance Monitor: 41xxx Series Adapters' Counters

Table 7-3 provides details about error counters.

Table 7-3. Marvell FastLinQ RDMA Error Counters

RDMA Error Counter	Description	Applies to RoCE?	Applies to iWARP?	Troubleshooting
CQ overflow	A completion queue on which an RDMA work request is posted. This counter specifies the quantity of instances where there was a completion for a work request on the send or receive queue, but no space on the associated completion queue.	Yes	Yes	Indicates a software design issue causing an insufficient completion queue size.
Requestor Bad response	A malformed response was returned by the responder.	Yes	Yes	_

Table 7-3. Marvell FastLinQ RDMA Error Counters (Continued)

RDMA Error Counter	Description	Applies to RoCE?	Applies to iWARP?	Troubleshooting
Requestor CQEs flushed with error	Posted work requests may be flushed by sending completions with a flush status to the CQ (without completing the actual execution of the work request) if the QP moves to an error state for any reason and pending work requests exist. If a work request completed with error status, all other pending work requests for that QP are flushed.	Yes	Yes	Occurs when the RDMA connection is down.
Requestor Local length	The RDMA Read response message contained too much or too little payload data.	Yes	Yes	Usually indicates an issue with the host software components.
Requestor local protection	The locally posted work request's data segment does not reference a memory region that is valid for the requested operation.	Yes	Yes	Usually indicates an issue with the host software components.
Requestor local QP operation	An internal QP consistency error was detected while processing this work request.	Yes	Yes	_
Requestor Remote access	A protection error occurred on a remote data buffer to be read by an RDMA Read, written by an RDMA Write, or accessed by an atomic operation.	Yes	Yes	_
Requestor Remote Invalid request	The remote side received an invalid message on the channel. The invalid request may have been a Send message or an RDMA request.	Yes	Yes	Possible causes include the operation is not supported by this receive queue, insufficient buffering to receive a new RDMA or atomic operation request, or the length specified in an RDMA request is greater than 231 bytes.

Table 7-3. Marvell FastLinQ RDMA Error Counters (Continued)

RDMA Error Counter			Applies to iWARP?	Troubleshooting
Requestor remote operation	Remote side could not complete the operation requested due to a local issue.	Yes	Yes	A software issue at the remote side (for example, one that caused a QP error or a malformed WQE on the RQ) prevented operation completion.
Requestor retry exceeded	Transport retries have exceeded the maximum limit.	Yes	Yes	The remote peer may have stopped responding, or a network issue is preventing messages acknowledgment.
Requestor RNR Retries exceeded	Retry due to RNR NAK received have been tried the maximum number of times without success.	Yes	No	The remote peer may have stopped responding, or a network issue is preventing messages acknowledgment.
Responder CQE flushed	Posted work requests (receive buffers on RQ) may be flushed by sending completions with a flush status to the CQ if the QP moves to an error state for any reason, and pending receive buffers exist on the RQ. If a work request completed with an error status, all other pending work requests for that QP are flushed.	Yes	Yes	
Responder local length	Invalid length in inbound messages.	Yes	Yes	Misbehaving remote peer. For example, the inbound send messages have lengths greater than the receive buffer size.
Responder local protection	The locally posted work request's data segment does not reference a memory region that is valid for the requested operation.	Yes	Yes	Indicates a software issue with memory management.

Table 7-3. Marvell FastLinQ RDMA Error Counters (Continued)

RDMA Error Counter	Description	Applies to RoCE?	Applies to iWARP?	Troubleshooting
Responder Local QP Operation error	An internal QP consistency error was detected while processing this work request.	Yes	Yes	Indicates a software issue.
Responder remote invalid request	The responder detected an invalid inbound message on the channel.	Yes	Yes	Indicates possible mis- behavior by a remote peer. Possible causes include: the operation is not supported by this receive queue, insuffi- cient buffering to receive a new RDMA request, or the length specified in an RDMA request is greater than 2 ³¹ bytes.

Configuring RoCE for SR-IOV VF Devices (VF RDMA)

The following sections describe how to configure RoCE for SR-IOV VF devices (also referred to as *VF RDMA*). Associated information and limitations are also provided.

Configuration Instructions

To configure VF RDMA:

- 1. Install the VF RDMA capable components (drivers, firmware, multiboot image (MBI)).
- 2. Configure QoS for VF RDMA.

QoS configuration is needed to configure priority flow control (PFC) for RDMA. Configure QoS in the host as documented in "Configuring QoS for RoCE" on page 259. (QoS configuration must be done in the host, not in the VM).

- 3. Configure Windows Hyper-V for VF RDMA:
 - a. Enable SR-IOV in HII and on the Advanced tab in Windows Device Manager.
 - b. Open the Windows Hyper-V Manager on the host.

- c. Open the Virtual Switch Manager from the right pane.
- d. Select New Virtual Network switch with type External.

Figure 7-4 shows an example.

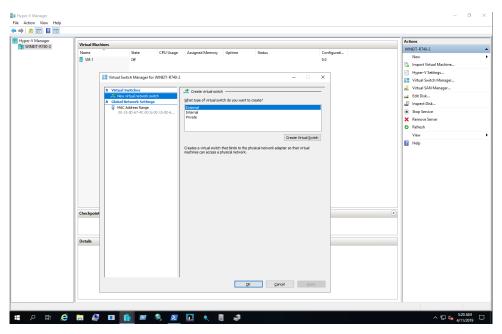


Figure 7-4. Setting an External New Virtual Network Switch

e. Click the **External network** button, and then select the appropriate adapter. Click **Enable single-root I/O virtualization (SR-IOV)**.

Figure 7-5 shows an example.

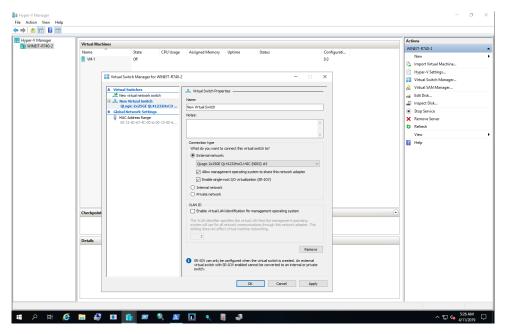


Figure 7-5. Setting SR-IOV for New Virtual Switch

f. Create a VM and open the VM settings.

Figure 7-6 shows an example.

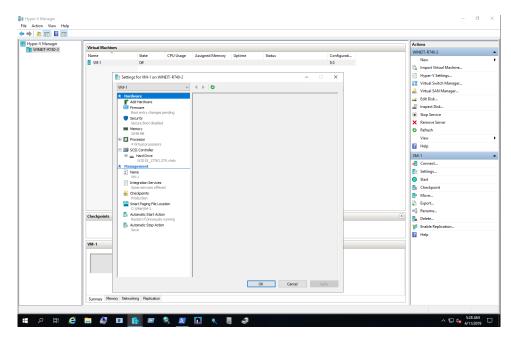


Figure 7-6. VM Settings

- g. Select **Add Hardware**, and then select **Network Adapter** to assign the virtual network adapters (VMNICs) to the VM.
- h. Select the newly created virtual switch.

i. Enable VLAN to the network adapter.

Figure 7-7 shows an example.

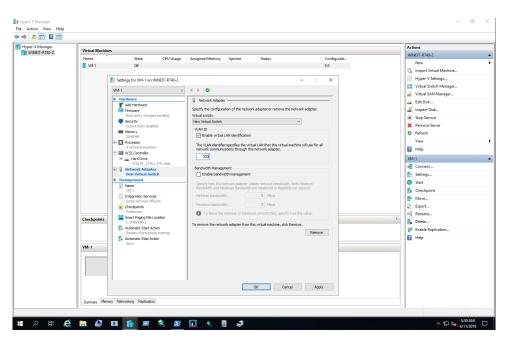


Figure 7-7. Enabling VLAN to the Network Adapter

j. Expand the network adapter settings. Under Single-root I/O virtualization, select **Enable SR-IOV** to enable SR-IOV capabilities for the VMNIC.

Figure 7-8 shows an example.

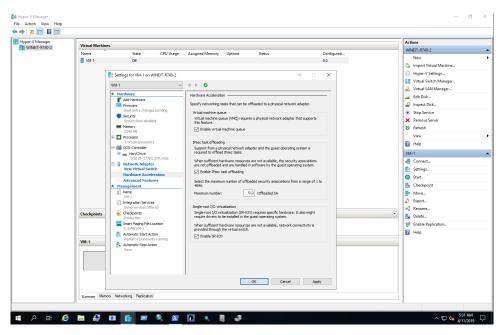


Figure 7-8. Enabling SR-IOV for the Network Adapter

4. Issue the following PowerShell command on the host to enable RDMA capabilities for the VMNIC (SR-IOV VF).

Set-VMNetworkAdapterRdma -VMName <VM_NAME> -VMNetworkAdapterName <VM NIC NAME> -RdmaWeight 100

NOTE

The VM must be powered off before issuing the PowerShell command.

5. Upgrade the Marvell drivers in the VM by booting the VM and installing the latest drivers using the Windows Super Installer on the Marvell CD.

Figure 7-9 shows an example.

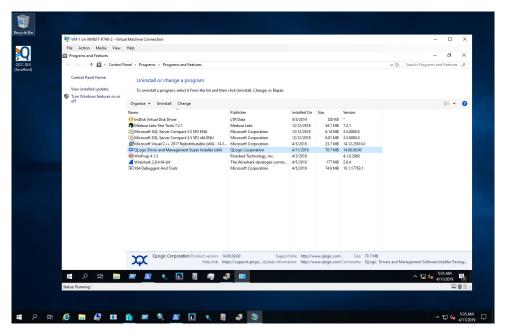
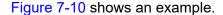


Figure 7-9. Upgrading Drivers in VM

6. Enable RMDA on the Microsoft network device associated with the VF inside the VM.



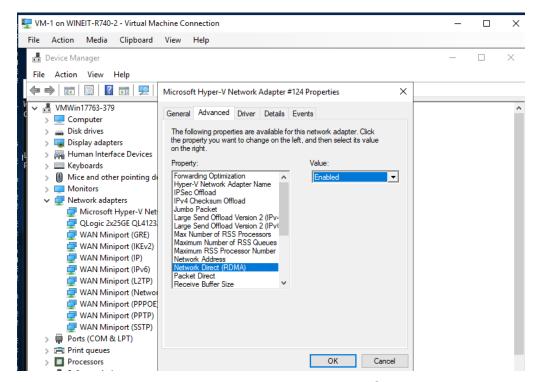


Figure 7-10. Enabling RDMA on the VMNIC

- 7. Start the VM RMDA traffic:
 - a. Connect a server message block (SMB) drive, run RoCE traffic, and verify the results.
 - b. Open the Performance monitor in the VM, and then add **RDMA Activity counter**.

c. Verify that RDMA traffic is running.

Figure 7-11 provides an example.

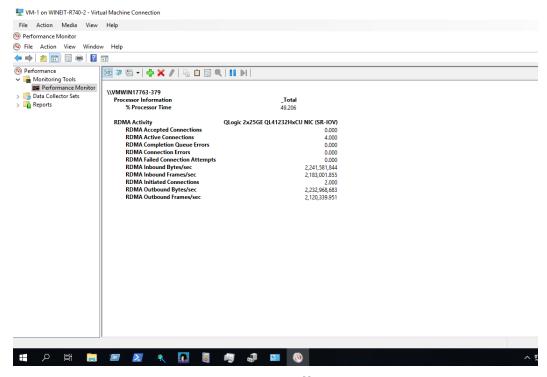


Figure 7-11. RDMA Traffic

Limitations

VF RDMA has the following limitations:

- VF RDMA is supported only for 41xxx-based devices.
- At the time of publication, only RoCEv2 is supported for VF RDMA. The same network direct technology must be configured in physical functions on both the host and SR-IOV VFs in the VM.
- A maximum of 16 VFs per PF can support VF RDMA. For quad-port adapters, the maximum is 8 VFs per PF.
- VF RDMA is supported only in Windows Server 2019 (for both the host and VM OSs).
- VF RDMA is not supported for Linux VMs on Windows Hypervisor.
- VF RDMA is not supported in NPAR mode.
- A maximum of 128 queue pairs (QPs)/connections are supported per VF.
- RDMA traffic between PF and its VFs, and among VFs of same PF, is supported. This traffic pattern is referred to as *loopback traffic*.

- On some older server platforms, VF devices may not be enumerated for one of the NIC PCI functions (PF). This limitation is because of the increased PCI base address register (BAR) requirements to support VF RDMA, meaning that the OS/BIOS cannot assign the required BAR for each VF.
- To support the maximum number of QPs in a VM, approximately 8GB of RAM must be available, assuming that only one VF is assigned to the VM. If less than 8GB of RAM is assigned to VM, there can be a sudden drop in the number of active connections due to insufficient memory and memory allocation failures.

Configuring RoCE on the Adapter for Linux

This section describes the RoCE configuration procedure for RHEL and SLES. It also describes how to verify the RoCE configuration and provides some guidance about using group IDs (GIDs) with vLAN interfaces.

- RoCE Configuration for RHEL
- RoCE Configuration for SLES
- Verifying the RoCE Configuration on Linux
- vLAN Interfaces and GID Index Values
- RoCE v2 Configuration for Linux
- Configuring RoCE for SR-IOV VF Devices (VF RDMA)

RoCE Configuration for RHEL

To configure RoCE on the adapter, the Open Fabrics Enterprise Distribution (OFED) must be installed and configured on the RHEL host.

To prepare inbox OFED for RHEL:

- 1. While installing or upgrading the operating system, select the InfiniBand and OFED support packages.
- 2. Install the following RPMs from the RHEL ISO image:

```
libibverbs-devel-x.x.x.x86_64.rpm
(required for libqedr library)
perftest-x.x.x.x86_64.rpm
(required for InfiniBand bandwidth and latency applications)

or, using Yum, install the inbox OFED:
yum groupinstall "Infiniband Support"
yum install perftest
yum install tcl tcl-devel tk zlib-devel libibverbs
libibverbs-devel
```

NOTE

During installation, if you already selected the previously mentioned packages, you need not reinstall them. The inbox OFED and support packages may vary depending on the operating system version.

3. Install the new Linux drivers as described in "Installing the Linux Drivers with RDMA" on page 14.

RoCE Configuration for SLES

To configure RoCE on the adapter for a SLES host, OFED must be installed and configured on the SLES host.

To install inbox OFED for SLES:

- 1. While installing or upgrading the operating system, select the InfiniBand support packages.
- 2. (SLES 12.x) Install the following RPMs from the corresponding SLES SDK kit image.

```
libibverbs-devel-x.x.x.x86_64.rpm (required for libqedr installation)
perftest-x.x.x86_64.rpm (required for bandwidth and latency applications)
```

3. (SLES 15/15 SP1) Install the following RPMs.

After installation, the rdma-core*, libibverbs*, libibumad*, libibmad*, libibmad*, librdmacm*, and perftest RPMs may be missing (all are required for RDMA). Install these packages using one of the following methods:

- □ Load the Package DVD and install the missing RPMs.
- ☐ Use the zypper command to install the missing RPMs. For example:

```
#zypper install rdma*
#zypper install libib*
#zypper install librdma*
#zypper install perftest
```

4. Install the Linux drivers, as described in "Installing the Linux Drivers with RDMA" on page 14.

Verifying the RoCE Configuration on Linux

After installing OFED, installing the Linux driver, and loading the RoCE drivers, verify that the RoCE devices were detected on all Linux operating systems.

To verify RoCE configuration on Linux:

- 1. Stop firewall tables using service/systematl commands.
- 2. For RHEL only: If the RDMA service is installed (yum install rdma), verify that the RDMA service has started.

NOTE

For RHEL 7.x and SLES 12 SPx and later, RDMA service starts itself after reboot.

On RHEL or CentOS: Use the service rdma status command to start service:

- ☐ If RDMA has not started, issue the following command:
 - # service rdma start
- ☐ If RDMA does not start, issue either of the following alternative commands:
 - # /etc/init.d/rdma start

or

- # systemctl start rdma.service
- 3. Verify that the RoCE devices were detected by examining the dmesg logs:
 - # dmesg|grep qedr

```
[87910.988411] qedr: discovered and registered 2 RoCE funcs
```

- 4. Verify that all of the modules have been loaded. For example:
 - # lsmod|grep qedr

5. Configure the IP address and enable the port using a configuration method such as ifconfig. For example:

```
# ifconfig ethX 192.168.10.10/24 up
```

6. Issue the <code>ibv_devinfo</code> command. For each PCI function, you should see a separate <code>hca id</code>, as shown in the following example:

```
root@captain:~# ibv devinfo
hca id: qedr0
       transport:
                                        InfiniBand (0)
                                        8.3.9.0
       fw ver:
       node guid:
                                        020e:1eff:fe50:c7c0
       sys image guid:
                                        020e:1eff:fe50:c7c0
       vendor id:
                                        0x1077
       vendor part id:
                                        5684
       hw ver:
                                        0 \times 0
       phys port cnt:
                                        1
               port: 1
                                              PORT ACTIVE (1)
                       state:
                       max mtu:
                                                4096 (5)
                       active mtu:
                                                1024 (3)
                        sm lid:
                                                0
                        port lid:
                        port lmc:
                                                0x00
                        link layer:
                                                Ethernet
```

- 7. Verify the L2 and RoCE connectivity between all servers: one server acts as a server, another acts as a client.
 - ☐ Verify the L2 connection using a simple ping command.
 - ☐ Verify the RoCE connection by performing an RDMA ping on the server or client:

On the server, issue the following command:

```
ibv rc pingpong -d <ib-dev> -g 0
```

On the client, issue the following command:

```
ibv rc pingpong -d <ib-dev> -g 0 <server L2 IP address>
```

The following are examples of successful ping pong tests on the server and the client.

Server Ping:

```
root@captain:~# ibv_rc_pingpong -d qedr0 -g 0
local address: LID 0x0000, QPN 0xff0000, PSN 0xb3e07e, GID
fe80::20e:leff:fe50:c7c0
remote address: LID 0x0000, QPN 0xff0000, PSN 0x934d28, GID
fe80::20e:leff:fe50:c570
8192000 bytes in 0.05 seconds = 1436.97 Mbit/sec
1000 iters in 0.05 seconds = 45.61 usec/iter
```

Client Ping:

```
root@lambodar:~# ibv_rc_pingpong -d qedr0 -g 0 192.168.10.165
local address: LID 0x0000, QPN 0xff0000, PSN 0x934d28, GID
fe80::20e:leff:fe50:c570
remote address: LID 0x0000, QPN 0xff0000, PSN 0xb3e07e, GID
fe80::20e:leff:fe50:c7c0
8192000 bytes in 0.02 seconds = 4211.28 Mbit/sec
1000 iters in 0.02 seconds = 15.56 usec/iter
```

- To display RoCE statistics, issue the following commands, where *x* is the device number:
 - > mount -t debugfs nodev /sys/kernel/debug
 - > cat /sys/kernel/debug/gedr/gedrX/stats

vLAN Interfaces and GID Index Values

If you are using vLAN interfaces on both the server and the client, you must also configure the same vLAN ID on the switch. If you are running traffic through a switch, the InfiniBand applications must use the correct GID value, which is based on the vLAN ID and vLAN IP address.

Based on the following results, the GID value (-x 4 / -x 5) should be used for any perftest applications.

ibv devinfo -d qedr0 -v|grep GID

```
GID[ 0]: fe80:0000:0000:020e:1eff:fe50:c5b0

GID[ 1]: 0000:0000:0000:0000:ffff:c0a8:0103

GID[ 2]: 2001:0db1:0000:020e:1eff:fe50:c5b0

GID[ 3]: 2001:0db2:0000:0000:020e:1eff:fe50:c5b0

GID[ 4]: 0000:0000:0000:0000:ffff:c0a8:0b03 IP address for vLAN interface

GID[ 5]: fe80:0000:0000:0000:020e:1e00:0350:c5b0 vLAN ID 3
```

NOTE

The default GID value is zero (0) for back-to-back or pause settings. For server and switch configurations, you must identify the proper GID value. If you are using a switch, refer to the corresponding switch configuration documents for the correct settings.

RoCE v2 Configuration for Linux

To verify RoCE v2 functionality, you must use RoCE v2 supported kernels.

To configure RoCE v2 for Linux:

- 1. Ensure that you are using one of the following supported kernels:
 - ☐ SLES 15/15 SP1
 - □ SLES 12 SP4 and later
 - ☐ RHEL 7.6, 7.7, and 8.0
- 2. Configure RoCE v2 as follows:
 - a. Identify the GID index for RoCE v2.
 - b. Configure the routing address for the server and client.
 - c. Enable L3 routing on the switch.

NOTE

You can configure RoCE v1 and RoCE v2 by using RoCE v2-supported kernels. These kernels allow you to run RoCE traffic over the same subnet, as well as over different subnets such as RoCE v2 and any routable environment. Only a few settings are required for RoCE v2, and all other switch and adapter settings are common for RoCE v1 and RoCE v2.

Identifying the RoCE v2 GID Index or Address

To find RoCE v1- and RoCE v2-specific GIDs, use either sys or class parameters, or run RoCE scripts from the 41xxx FastLinQ source package. To check the default **RoCE GID Index** and address, issue the <code>ibv_devinfo</code> command and compare it with the sys or class parameters. For example:

#ibv devinfo -d qedr0 -v|grep GID

```
GID[ 0]: fe80:0000:0000:020e:1eff:fec4:1b20
GID[ 1]: fe80:0000:0000:020e:1eff:fec4:1b20
GID[ 2]: 0000:0000:0000:0000:ffff:1e01:010a
GID[ 3]: 0000:0000:0000:0000:ffff:1e01:010a
GID[ 4]: 3ffe:ffff:0000:0f21:0000:0000:0004
```

GID[5]:	3ffe:ffff:0000:0f21:0000:0000:0000:0004
GID[6] :	0000:0000:0000:0000:0000:ffff:c0a8:6403
GID[7]:	0000:0000:0000:0000:0000:ffff:c0a8:6403

Verifying the RoCE v1 or RoCE v2 GID Index and Address from sys and class Parameters

Use one of the following options to verify the RoCE v1 or RoCE v2 GID Index and address from the sys and class parameters:

Option 1:

```
# cat /sys/class/infiniband/qedr0/ports/1/gid_attrs/types/0
IB/RoCE v1
# cat /sys/class/infiniband/qedr0/ports/1/gid_attrs/types/1
RoCE v2
# cat /sys/class/infiniband/qedr0/ports/1/gids/0
fe80:0000:0000:0000:020e:1eff:fec4:1b20
# cat /sys/class/infiniband/qedr0/ports/1/gids/1
```

fe80:0000:0000:0000:020e:1eff:fec4:1b20

■ Option 2:

Use the scripts from the FastLinQ source package.

#/../fastlinq-8.x.x.x/add-ons/roce/show_gids.sh

DEV	PORT	INDEX	GID	IPv4	VER	DEV
qedr0	1	0	fe80:0000:0000:0000:020e:1eff:fec4:1b20		v1	p4p1
qedr0	1	1	fe80:0000:0000:0000:020e:1eff:fec4:1b20		v2	p4p1
qedr0	1	2	0000:0000:0000:0000:0000:ffff:1e01:010a	30.1.1.10	v1	p4p1
qedr0	1	3	0000:0000:0000:0000:0000:ffff:1e01:010a	30.1.1.10	v2	p4p1
qedr0	1	4	3ffe:ffff:0000:0f21:0000:0000:0000:0004		v1	p4p1
qedr0	1	5	3ffe:ffff:0000:0f21:0000:0000:0000:0004		v2	p4p1
qedr0	1	6	0000:0000:0000:0000:0000:ffff:c0a8:6403	192.168.100.3	v1	p4p1.100
qedr0	1	7	0000:0000:0000:0000:0000:ffff:c0a8:6403	192.168.100.3	v2	p4p1.100
qedr1	1	0	fe80:0000:0000:0000:020e:1eff:fec4:1b21		v1	p4p2
qedr1	1	1	fe80:0000:0000:0000:020e:1eff:fec4:1b21		v2	p4p2

NOTE

You must specify the GID index values for RoCE v1- or RoCE v2-based server or switch configuration (Pause/PFC). Use the GID index for the link local IPv6 address, IPv4 address, or IPv6 address. To use vLAN tagged frames for RoCE traffic, you must specify GID index values that are derived from the vLAN IPv4 or IPv6 address.

Verifying the RoCE v1 or RoCE v2 Function Through perftest Applications

This section shows how to verify the RoCE v1 or RoCE v2 function through perftest applications. In this example, the following server IP and client IP are used:

Server IP: 192.168.100.3Client IP: 192.168.100.4

Verifying RoCE v1

Run over the same subnet and use the RoCE v1 GID Index.

```
Server# ib_send_bw -d qedr0 -F -x 0
Client# ib send_bw -d qedr0 -F -x 0 192.168.100.3
```

Verifying RoCE v2

Run over the same subnet and use the RoCE v2 GID Index.

```
Server# ib_send_bw -d qedr0 -F -x 1
Client# ib_send_bw -d qedr0 -F -x 1 192.168.100.3
```

NOTE

If you are running through a switch PFC configuration, use vLAN GIDs for RoCE v1 or v2 through the same subnet.

Verifying RoCE v2 Through Different Subnets

NOTE

You must first configure the route settings for the switch and servers. On the adapter, set the RoCE priority and DCBX mode using either the HII, UEFI user interface, or one of the Marvell management utilities.

To verify RoCE v2 through different subnets:

1. Set the route configuration for the server and client using the DCBX-PFC configuration.

```
□ System Settings:
```

```
Server VLAN IP: 192.168.100.3 and Gateway: 192.168.100.1 Client VLAN IP: 192.168.101.3 and Gateway: 192.168.101.1
```

□ Server Configuration:

```
#/sbin/ip link add link p4p1 name p4p1.100 type vlan id 100
#ifconfig p4p1.100 192.168.100.3/24 up
#ip route add 192.168.101.0/24 via 192.168.100.1 dev p4p1.100
```

□ Client Configuration:

```
#/sbin/ip link add link p4p1 name p4p1.101 type vlan id 101
#ifconfig p4p1.101 192.168.101.3/24 up
#ip route add 192.168.100.0/24 via 192.168.101.1 dev p4p1.101
```

- 2. Set the switch settings using the following procedure.
 - ☐ Use any flow control method (Pause, DCBX-CEE, or DCBX-IEEE), and enable IP routing for RoCE v2. See "Preparing the Ethernet Switch" on page 129 for RoCE v2 configuration, or refer to the vendor switch documents.
 - ☐ If you are using PFC configuration and L3 routing, run RoCE v2 traffic over the vLAN using a different subnet, and use the RoCE v2 vLAN GID index.

```
Server# ib_send_bw -d qedr0 -F -x 5
Client# ib send bw -d qedr0 -F -x 5 192.168.100.3
```

Server Switch Settings:

```
[root@RoCE-Auto-2 /] # ib_send_bw -d qedr0 -F -x 5 -q 2 --report_gbits
***********
* Waiting for client to connect... *
                         Send BW Test
                     : OFF Device
: 2 Transpor
Dual-port
Number of qps : 2
Connection type : RC
                                         Transport type : IB
                                       Using SRQ
RX depth
CQ Moderation : 100
                      : Ethernet
Gid index
rdma_cm QPs
Data ex. method : Ethernet
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:100:03
local address: LID 0000 QPN 0xff0002 PSN 0xa2b8f1
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:100:03
remote address: LID 0000 QPN 0xff0000 PSN 0x40473a
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:101:03
remote address: LID 0000 QFN 0xff0002 PSN 0x124cd3
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:101:03
                                  BW peak[Gb/sec] BW average[Gb/sec]
                                                                                         MsgRate[Mpps]
#bytes
```

Figure 7-12. Switch Settings, Server

Client Switch Settings:

```
oot@roce-auto-1 ~]# ib_send_bw -d qedr0 -F -x 5 192.168.100.3 -q 2 --report_gbits
                       Send BW Test
Dual-port
Number of qps
                                      Using SRQ
TX depth
CQ Moderation
Link type
Gid index
rdma_cm QPs : OFF
Data ex. method : Ethernet
local address: LID 0000 QPN 0xff0000 PSN 0x40473a
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:101:03
local address: LID 0000 QPN 0xff0002 PSN 0x124cd3
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:101:03
remote address: LID 0000 QFN 0xff0000 FSN 0xf0b2c3
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:100:03
BW peak[Gb/sec]
                                                        BW average[Gb/sec]
                                                                                  MsgRate[Mpps]
```

Figure 7-13. Switch Settings, Client

Configuring RoCE v1 or RoCE v2 Settings for RDMA_CM Applications

To configure RoCE, use the following scripts from the FastLinQ source package:

```
# ./show_rdma_cm_roce_ver.sh
qedr0 is configured to IB/RoCE v1
qedr1 is configured to IB/RoCE v1
# ./config_rdma_cm_roce_ver.sh v2
configured rdma_cm for qedr0 to RoCE v2
configured rdma cm for qedr1 to RoCE v2
```

Server Settings:

```
[root@Roce-Auto-2 /] # rping -s -v -c 10
server ping data: rdma-ping-0: ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqr
server ping data: rdma-ping-1: BCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrs
server ping data: rdma-ping-2: CDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrst
server ping data: rdma-ping-3: DEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstu
server ping data: rdma-ping-4: EFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuv
server ping data: rdma-ping-5: FGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvw
server ping data: rdma-ping-6: GHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwx
server ping data: rdma-ping-7: HIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxy
server ping data: rdma-ping-8: IJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-9: JKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server DISCONNECT EVENT...
wait for RDMA_READ_ADV state 10
[root@Roce-Auto-2 /]# []
```

Figure 7-14. Configuring RDMA_CM Applications: Server Client Settings:

```
[root@roce-auto-1 ~] # rping -c -v -C 10 -a 192.168.100.3
ping data: rdma-ping-0: ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqr
ping data: rdma-ping-1: BCDEFGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrs
ping data: rdma-ping-2: CDEFGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrst
ping data: rdma-ping-3: DEFGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstu
ping data: rdma-ping-4: EFGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstuv
ping data: rdma-ping-5: FGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstuvw
ping data: rdma-ping-6: GHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstuvwx
ping data: rdma-ping-7: HIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstuvwxy
ping data: rdma-ping-8: IJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-9: JKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmnopqrstuvwxyz
client DISCONNECT EVENT...
[root@roce-auto-1 ~]#
```

Figure 7-15. Configuring RDMA_CM Applications: Client

Configuring RoCE for SR-IOV VF Devices (VF RDMA)

The following sections describe how to configure RoCE for SR-IOV VF devices (also referred to as VFRMDA) on Linux. Associated information and limitations are also provided.

Table 7-4 lists the supported Linux OS combinations.

Table 7-4. Supported Linux OSs for VF RDMA

		Guest OS					
Hypervisor	RHEL 7.6	RHEL 7.7	RHEL 8.0	SLES12 SP4	SLES15 SP0	SLES15 SP1	
	Yes	Yes	s Yes Yes	Yes	Yes	Yes	
RHEL 7.7	Yes	Yes	Yes	Yes	Yes	Yes	
RHEL 8.0	Yes	Yes	Yes	Yes	Yes	Yes	
SLES12 SP4	Yes	Yes	Yes	Yes	Yes	Yes	
SLES15 SP0	Yes	Yes	Yes	Yes	Yes	Yes	
SLES15 SP1	Yes	Yes	Yes	Yes	Yes	Yes	

If you are using the inbox OFED, use the same OFED distribution between the hypervisor host OS and the guest (VM) OS. Check the out-of-box OFED distribution release notes for their specific supported host OS-to-VM OS distribution matrix.

Enumerating VFs for L2 and RDMA

There are two ways to enumerate the VFs:

- User Defined VF MAC Allocation
- Dynamic or Random VF MAC Allocation

User Defined VF MAC Allocation

When defining the VF MAC allocation, there are no changes in the default VF enumeration method. After creating the number of VFs, assign the static MAC address.

To create a user defined VF MAC allocation:

1. Enumerate the default VF.

```
# modprobe -v qede
# echo 2 > /sys/class/net/p6p1/device/sriov_numvfs
# ip link show
14: p6p1: <NO-CARRIER, BROADCAST, MULTICAST, UP> mtu 1500 qdisc mq state DOWN
mode DEFAULT group default qlen 1000
    link/ether 14:02:ec:ce:d0:e4 brd ff:ff:ff:ff
    vf 0 MAC 00:00:00:00:00:00, spoof checking off, link-state auto
```

vf 1 MAC 00:00:00:00:00:00, spoof checking off, link-state auto

2. Assign the static MAC address:

```
# ip link set dev p6p1 vf 0 mac 3c:33:44:55:66:77
```

ip link set dev p6p1 vf 1 mac 3c:33:44:55:66:89

#ip link show

14: p6p1: <BROADCAST, MULTICAST, UP, LOWER_UP> mtu 1500 qdisc mq state UP mode DEFAULT group default qlen 1000

link/ether 14:02:ec:ce:d0:e4 brd ff:ff:ff:ff:ff

vf 0 MAC 3c:33:44:55:66:77, tx rate 25000 (Mbps), $max_tx_rate 25000Mbps$, spoof checking off, link-state auto

 $\label{local_max_tx_rate} $$ vf 1 MAC 3c:33:44:55:66:89, tx rate 25000 (Mbps), max_tx_rate 25000Mbps, spoof checking off, link-state auto$

3. To reflect for RDMA, load/unload the gedr driver if it is already loaded.

#rmmod qedr		
#modprobe	qedr	
#ibv_devices		
device		node GUID
qedr0		1602ecfffeced0e4
qedr1		1602ecfffeced0e5
qedr_vf0		3e3344fffe556677
gedr vf1		3e3344fffe556689

Dynamic or Random VF MAC Allocation

To dynamically allocate a VF MAC:

```
# modprobe -r qedr
```

modprobe -v qed vf_mac_origin=3 [Use this module parameter for dynamic MAC allocation]

```
# modprobe -v qede
```

- # echo 2 > /sys/class/net/p6p1/device/sriov numvfs
- # modprobe qedr (This is an optional, mostly qedr driver loads itself)
- # ip link show|grep vf

vf 0 MAC ba:1a:ad:08:89:00, tx rate 25000 (Mbps), max_tx_rate
25000Mbps, spoof checking off, link-state auto

vf 1 MAC 96:40:61:49:cd:68, tx rate 25000 (Mbps), max_tx_rate
25000Mbps, spoof checking off, link-state auto

- # lsmod |grep gedr
- # ibv devices

device	node GUID
qedr0	1602ecfffececfa0
qedr1	1602ecfffececfa1
qedr_vf0	b81aadfffe088900
qedr vf1	944061fffe49cd68

Number of VFs Supported for RDMA

For the 41xxx Series Adapters, the number of VFs for L2 and RDMA are shared based on resources availability.

Dual Port Adapters

Each PF supports a maximum of 40 VFs for RDMA. If the number of VFs exceeds 56, it will be subtracted by the total number of VFs (96).

In the following example, PF0 is

```
/sys/class/net/<PF-interface>/device/sriov_numvfs

Echo 40 > PF0 (VFs for L2+RDMA=40+40 (40 VFs can use for both L2 and RDMA))

Echo 56 > PF0 (VFs for L2+RDMA=56+40)
```

After crossing 56 VFs, this number is subtracted by the total number of VFs. For example:

```
echo 57 > PF0 then 96-57=39 VFs for RDMA (57 VFs for L2 + 39VFs for RDMA) echo 96 > PF0 then 96-96=0 VFs for RDMA (all 96 VFs can use only for L2)
```

To view the available VFs for L2 and RDMA:

```
L2 : # ip link show RDMA: # ibv devices
```

Quad Port Adapters

Each PF supports a maximum of 20 VFs for RDMA; until 48 VFs, there are 20 VFs for RDMA. When exceeding 28 VFs, that number is subtracted by the total VFs (48).

For example, in a 4x10G:

```
Echo 20 > PF0 (VFs for L2+RDMA=20+20)
Echo 28 > PF0 (VFs for L2+RDMA=28+20)
```

When exceeding 28 VFs, this number is subtracted by the total number of VFs. For example:

```
echo 29 > PF0 (48–29=19VFs for RDMA; 29 VFs for L2 + 19 VFs for RDMA) echo 48 > PF0 (48-48=0 VFs for RDMA; all 48 VFs can use only for L2)
```

Limitations

VF RMDA has the following limitations:

- No iWARP support
- No NPAR support
- Cross OS is not supported; for example, a Linux hypervisor cannot use a Windows guest OS (VM)
- Perftest latency test on VF interfaces can be run only with the inline size zero -I 0 option. Neither the default nor more than one inline size works.
- To allow RDMA_CM applications to run on different MTU sizes (512–9000) other than the default (1500), follow these steps:
 - 1. Unload the gedr driver:
 - #rmmod qedr
 - 2. Set MTU on the VF interface:
 - #ifconfig <VF interface> mtu 9000
 - 3. Load the gedr driver:
 - #modprobe gedr
- The rdma_server/rdma_xserver does not support VF interfaces.
- No RDMA bonding support on VFs.

Configuring RoCE on the Adapter for VMware ESX

This section provides the following procedures and information for RoCE configuration:

- Configuring RDMA Interfaces
- Configuring MTU
- RoCE Mode and Statistics
- Configuring a Paravirtual RDMA Device (PVRDMA)

NOTE

Mapping Ethernet speeds to RDMA speeds is not always accurate, because values that can be specified by the RoCE driver are aligned with Infiniband[®]. For example, if RoCE is configured on an Ethernet interface operating at 1Gbps, the RDMA speed is shown as 2.5Gbps. There are no other suitable values available in the header files provided by ESXi that can be used by the RoCE driver to indicate 1Gbps speed.

Configuring RDMA Interfaces

To configure the RDMA interfaces:

- 1. Install both Marvell NIC and RoCE drivers.
- 2. Using the module parameter, enable the RoCE function from the NIC driver by issuing the following command:

```
esxcfg-module -s 'enable roce=1' qedentv
```

To apply the change, reload the NIC driver or reboot the system.

3. To view a list of the NIC interfaces, issue the <code>esxcfg-nics -l command</code>. For example:

esxcfg-nics -1

Name	PCI	Driver	Link	Speed	Duplex	MAC Address	MTU	Descript	tion
Vmnic0	0000:01:00.2	qedentv	Up	25000Mbps	Full	a4:5d:36:2b:6c:92	1500	QLogic (Corp.
QLogic	FastLinQ QL41:	xxx 1/10/25	GbE Et	hernet Adap	oter				
Vmnic1	0000:01:00.3	qedentv	Up	25000Mbps	Full	a4:5d:36:2b:6c:93	1500	QLogic (Corp.
QLogic	FastLinQ QL41:	xxx 1/10/25	GbE Et	hernet Adap	oter				

4. To view a list of the RDMA devices, issue the esxcli rdma device list command. For example:

esxcli rdma device list

Name	Driver	State	MTU	Speed	Paired Upli	ink D	escription			
vmrdma0	qedrntv	Active	1024	25 Gbps	vmnic0	QLogi	c FastLinQ	QL45xxx	RDMA	Interface
vmrdma1	qedrntv	Active	1024	25 Gbps	vmnic1	QLogi	c FastLinQ	QL45xxx	RDMA	Interface

5. To create a new virtual switch, issue the following command:

esxcli network vswitch standard add -v <new vswitch name>
For example:

esxcli network vswitch standard add -v roce_vs

This creates a new virtual switch named *roce vs*.

- 6. To associate the Marvell NIC port to the vSwitch, issue the following command:
 - # esxcli network vswitch standard uplink add -u <uplink
 device> -v <roce vswitch>

For example:

esxcli network vswitch standard uplink add -u vmnic0 -v roce vs

7. To create a new port group on this vSwitch, issue the following command:

esxcli network vswitch standard portgroup add -p roce_pg -v
roce vs

For example:

- # esxcli network vswitch standard portgroup add -p roce_pg -v
 roce vs
- 8. To create a vmknic interface on this port group and configure the IP, issue the following command:
 - # esxcfg-vmknic -a -i <IP address> -n <subnet mask> <roce port
 group name>

For example:

- # esxcfg-vmknic -a -i 192.168.10.20 -n 255.255.255.0 roce pg
- 9. To configure the vLAN ID, issue the following command:
 - # esxcfg-vswitch -v <VLAN ID> -p roce pg

To run RoCE traffic with a vLAN ID, configure the vLAN ID on the corresponding VMkernel port group.

Configuring MTU

To modify the MTU for an RoCE interface, change the MTU of the corresponding vSwitch. Set the MTU size of the RDMA interface based on the MTU of the vSwitch by issuing the following command:

esxcfg-vswitch -m <new MTU> <RoCE vswitch name>

For example:

esxcfg-vswitch -m 4000 roce vs

esxcli rdma device list

Name	Driver	State	MTU	Speed	Pai	red Uplink	Descr	iption				
vmrdma0	qedrntv	Active	204	48 25 0	Gbps	vmnic0	QLogic	FastLinQ	QL45xxx	RDMA	Interface	
vmrdma1	qedrntv	Active	102	24 25 0	Gbps	vmnic1	QLogic	FastLinQ	QL45xxx	RDMA	Interface	

RoCE Mode and Statistics

For the RoCE mode, ESXi requires concurrent support of both RoCE v1 and v2. The decision regarding which RoCE mode to use is made during queue pair creation. The ESXi driver advertises both modes during registration and initialization. To view RoCE statistics, issue the following command:

esxcli rdma device stats get -d vmrdma0

```
Packets received: 0
Packets sent: 0
Bytes received: 0
Bytes sent: 0
Error packets received: 0
Error packets sent: 0
Error length packets received: 0
Unicast packets received: 0
Multicast packets received: 0
Unicast bytes received: 0
Multicast bytes received: 0
Unicast packets sent: 0
Multicast packets sent: 0
Unicast bytes sent: 0
Multicast bytes sent: 0
Queue pairs allocated: 0
Queue pairs in RESET state: 0
Queue pairs in INIT state: 0
Queue pairs in RTR state: 0
Queue pairs in RTS state: 0
Queue pairs in SQD state: 0
Queue pairs in SQE state: 0
Queue pairs in ERR state: 0
Queue pair events: 0
Completion queues allocated: 1
Completion queue events: 0
Shared receive queues allocated: 0
Shared receive queue events: 0
Protection domains allocated: 1
Memory regions allocated: 3
Address handles allocated: 0
Memory windows allocated: 0
```

Configuring a Paravirtual RDMA Device (PVRDMA)

See VMware's documentation (for example, https://kb.vmware.com/articleview?docid=2147694) for details on configuring PVRDMA using a vCenter interface. The following instructions are only for reference.

To configure PVRDMA using a vCenter interface:

- 1. Create and configure a new distributed virtual switch as follows:
 - a. In the VMware vSphere[®] Web Client, right-click the **RoCE** node in the left pane of the Navigator window.
 - b. On the Actions menu, point to **Distributed Switch**, and then click **New Distributed Switch**.
 - c. Select version 6.5.0.
 - d. Under **New Distributed Switch**, click **Edit settings**, and then configure the following:
 - Number of uplinks. Select an appropriate value.
 - Network I/O Control. Select Disabled.
 - **Default port group**. Select the **Create a default port group** check box.
 - Port group name. Type a name for the port group.

Figure 7-16 shows an example.

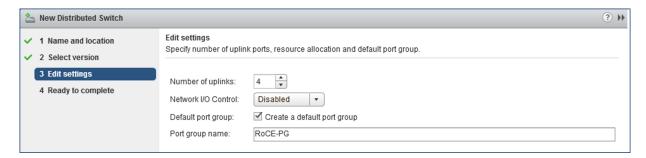


Figure 7-16. Configuring a New Distributed Switch

- 2. Configure a distributed virtual switch as follows:
 - a. In the VMware vSphere Web Client, expand the **RoCE** node in the left pane of the Navigator window.
 - b. Right-click **RoCE-VDS**, and then click **Add and Manage Hosts**.
 - c. Under **Add and Manage Hosts**, configure the following:
 - Assign uplinks. Select from the list of available uplinks.
 - Manage VMkernel network adapters. Accept the default, and then click Next.
 - Migrate VM networking. Assign the port group created in Step 1.

- 3. Assign a vmknic for PVRDMA to use on ESX hosts:
 - a. Right-click a host, and then click **Settings**.
 - b. On the Settings page, expand the **System** node, and then click **Advanced System Settings**.
 - c. The Advanced System Settings page shows the key-pair value and its summary. Click **Edit**.
 - d. On the Edit Advanced System Settings page, filter on **PVRDMA** to narrow all the settings to just Net.PVRDMAVmknic.
 - e. Set the **Net.PVRDMAVmknic** value to **vmknic**; for example, **vmk1**. Figure 7-17 shows an example.

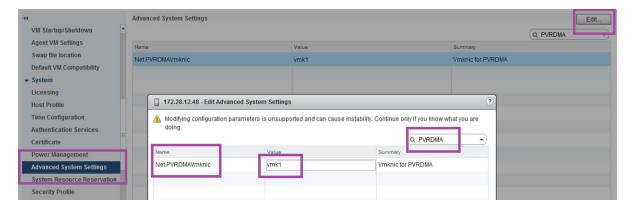


Figure 7-17. Assigning a vmknic for PVRDMA

- 4. Set the firewall rule for the PVRDMA:
 - a. Right-click a host, and then click **Settings**.
 - b. On the Settings page, expand the **System** node, and then click **Security Profile**.
 - c. On the Firewall Summary page, click Edit.
 - d. In the Edit Security Profile dialog box under **Name**, scroll down, select the **pvrdma** check box, and then select the **Set Firewall** check box.

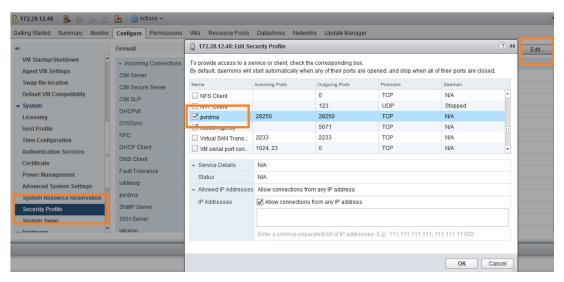


Figure 7-18 shows an example.

Figure 7-18. Setting the Firewall Rule

- 5. Set up the VM for PVRDMA as follows:
 - a. Install the following supported guest OS:
 - RHEL 7.5, 7.6, and 8.0
 - b. Install OFED4.17-1.
 - c. Compile and install the PVRDMA guest driver and library.
 - d. Add a new PVRDMA network adapter to the VM as follows:
 - Edit the VM settings.
 - Add a new network adapter.
 - Select the newly added DVS port group as **Network**.
 - Select PVRDMA as the adapter type.
 - After the VM is booted, ensure that the PVRDMA guest driver is loaded.

Configuring DCQCN

Data Center Quantized Congestion Notification (DCQCN) is a feature that determines how an RoCE receiver notifies a transmitter that a switch between them has provided an explicit congestion notification (notification point), and how a transmitter reacts to such notification (reaction point).

This section provides the following information about DCQCN configuration:

- DCQCN Terminology
- DCQCN Overview
- DCB-related Parameters
- Global Settings on RDMA Traffic
- Configuring DSCP-PFC
- Enabling DCQCN
- Configuring CNP
- DCQCN Algorithm Parameters
- MAC Statistics
- Script Example
- Limitations

DCQCN Terminology

The following terms describe DCQCN configuration:

- **ToS** (type of service) is a single-byte in the IPv4 header field. ToS comprises two ECN least significant bits (LSB) and six Differentiated Services Code Point (DSCP) most significant bits (MSB). For IPv6, traffic class is the equivalent of the IPv4 ToS.
- **ECN** (explicit congestion notification) is a mechanism where a switch adds to outgoing traffic an indication that congestion is imminent.
- CNP (congestion notification packet) is a packet used by the notification point to indicate that the ECN arrived from the switch back to the reaction point. CNP is defined in the Supplement to *InfiniBand Architecture Specification Volume 1 Release 1.2.1*, located here:
 - https://cw.infinibandta.org/document/dl/7781
- VLAN Priority is a field in the L2 vLAN header. The field is the three MSBs in the vLAN tag.
- **PFC** (priority-based flow control) is a flow control mechanism that applies to traffic carrying a specific vLAN priority.

- DSCP-PFC is a feature that allows a receiver to interpret the priority of an incoming packet for PFC purposes, rather than according to the vLAN priority or the DSCP field in the IPv4 header. You may use an indirection table to indicate a specified DSCP value to a vLAN priority value.

 DSCP-PFC can work across L2 networks because it is an L3 (IPv4) feature.
- Traffic classes, also known as priority groups, are groups of vLAN priorities (or DSCP values if DSCP-PFC is used) that can have properties such as being lossy or lossless. Generally, 0 is used for the default common lossy traffic group, 3 is used for the FCoE traffic group, and 4 is used for the iSCSI-TLV traffic group. You may encounter DCB mismatch issues if you attempt to reuse these numbers on networks that also support FCoE or iSCSI-TLV traffic. Marvell recommends that you use numbers 1–2 or 5–7 for RoCE-related traffic groups.
- **ETS** (enhanced transition services) is an allocation of maximum bandwidth per traffic class.

DCQCN Overview

Some networking protocols (RoCE, for example) require droplessness. PFC is a mechanism for achieving droplessness in an L2 network, and DSCP-PFC is a mechanism for achieving it across distinct L2 networks. However, PFC is deficient in the following regards:

- When activated, PFC completely halts the traffic of the specified priority on the port, as opposed to reducing transmission rate.
- All traffic of the specified priority is affected, even if there is a subset of specific connections that are causing the congestion.
- PFC is a single-hop mechanism. That is, if a receiver experiences congestion and indicates the congestion through a PFC packet, only the nearest neighbor will react. When the neighbor experiences congestion (likely because it can no longer transmit), it also generates its own PFC. This generation is known as *pause propagation*. Pause propagation may cause inferior route utilization, because all buffers must congest before the transmitter is made aware of the problem.

DCQCN addresses all of these disadvantages. The ECN delivers congestion indication to the reaction point. The reaction point sends a CNP packet to the transmitter, which reacts by reducing its transmission rate and avoiding the congestion. DCQCN also specifies how the transmitter attempts to increase its transmission rate and use bandwidth effectively after congestion ceases. DCQCN is described in the 2015 SIGCOMM paper, *Congestion Control for Large-Scale RDMA Deployments*, located here:

http://conferences.sigcomm.org/sigcomm/2015/pdf/papers/p523.pdf

DCB-related Parameters

Use DCB to map priorities to traffic classes (priority groups). DCB also controls which priority groups are subject to PFC (lossless traffic), and the related bandwidth allocation (ETS).

Global Settings on RDMA Traffic

Global settings on RDMA traffic include configuration of vLAN priority, ECN, and DSCP.

Setting vLAN Priority on RDMA Traffic

Use an application to set the vLAN priority used by a specified RDMA Queue Pair (QP) when creating a QP. For example, the <code>ib_write_bw</code> benchmark controls the priority using the <code>-sl</code> parameter. When RDMA-CM (RDMA Communication Manager) is present, you may be unable to set the priority.

Another method to control the vLAN priority is to use the <code>rdma_glob_vlan_pri</code> node. This method affects QPs that are created after setting the value. For example, to set the vLAN priority number to 5 for subsequently created QPs, issue the following command:

```
./debugfs.sh -n eth0 -t rdma glob vlan pri 5
```

Setting ECN on RDMA Traffic

Use the <code>rdma_glob_ecn</code> node to enable ECN for a specified RoCE priority. For example, to enable ECN on RoCE traffic using priority 5, issue the following command:

```
./debugfs.sh -n eth0 -t rdma glob ecn 1
```

This command is typically required when DCQCN is enabled.

Setting DSCP on RDMA Traffic

Use the <code>rdma_glob_dscp</code> node to control DSCP. For example, to set DSCP on RoCE traffic using priority 5, issue the following command:

```
./debugfs.sh -n eth0 -t rdma glob dscp 6
```

This command is typically required when DCQCN is enabled.

Configuring DSCP-PFC

Use <code>dscp_pfc</code> nodes to configure the <code>dscp->priority</code> association for PFC. You must enable the feature before you can add entries to the map. For example, to map DSCP value 6 to priority 5, issue the following commands:

```
./debugfs.sh -n eth0 -t dscp_pfc_enable 1 ./debugfs.sh -n eth0 -t dscp_pfc_set 6 5
```

Enabling DCQCN

To enable DCQCN for RoCE traffic, probe the qed driver with the <code>dcqcn_enable</code> module parameter. DCQCN requires enabled ECN indications (see "Setting ECN on RDMA Traffic" on page 173).

Configuring CNP

Congestion notification packets (CNPs) can have a separate configuration of vLAN priority and DSCP. Control these packets using the $dcqcn_cnp_dscp$ and $dcqcn_cnp_vlan_priority$ module parameters. For example:

modprobe qed dcqcn_cnp_dscp=10 dcqcn_cnp_vlan_priority=6

DCQCN Algorithm Parameters

Table 7-5 lists the algorithm parameters for DCQCN.

Table 7-5. DCQCN Algorithm Parameters

Parameter	Description and Values
dcqcn_cnp_send_timeout	Minimal difference of send time between CNPs. Units are in microseconds. Values range between 50500000.
dcqcn_cnp_dscp	DSCP value to be used on CNPs. Values range between 063.
dcqcn_cnp_vlan_priority	vLAN priority to be used on CNPs. Values range between 07. FCoE-Offload uses 3 and iSCSI-Offload-TLV generally uses 4. Marvell recommends that you specify a number from 1–2 or 5–7. Use this same value throughout the entire network.
dcqcn_notification_point	O – Disable DCQCN notification point. 1 – Enable DCQCN notification point.
dcqcn_reaction_point	0 – Disable DCQCN reaction point. 1 – Enable DCQCN reaction point.
dcqcn_rl_bc_rate	Byte counter limit
dcqcn_rl_max_rate	Maximum rate in Mbps
dcqcn_rl_r_ai	Active increase rate in Mbps
dcqcn_rl_r_hai	Hyperactive increase rate in Mbps.
dcqcn_gd	Alpha update gain denominator. Set to 32 for 1/32, and so on.

Table 7-5. DCQCN Algorithm Parameters (Continued)

Parameter	Description and Values
dcqcn_k_us	Alpha update interval
dcqcn_timeout_us	DCQCN timeout

MAC Statistics

To view MAC statistics, including per-priority PFC statistics, issue the phy_mac_stats command. For example, to view statistics on port 1 issue the following command:

```
./debugfs.sh -n eth0 -d phy_mac_stat -P 1
```

Script Example

The following example can be used as a script:

```
# probe the driver with both reaction point and notification point enabled
# with cnp dscp set to 10 and cnp vlan priority set to 6
modprobe qed dcqcn enable=1 dcqcn notification point=1 dcqcn reaction point=1
dcqcn cnp dscp=10 dcqcn cnp vlan priority=6
modprobe qede
# dscp-pfc configuration (associating dscp values to priorities)
# This example is using two DCBX traffic class priorities to better demonstrate
DCQCN in operation
debugfs.sh -n ens6f0 -t dscp pfc enable 1
debugfs.sh -n ens6f0 -t dscp pfc set 20 5
debugfs.sh -n ens6f0 -t dscp pfc set 22 6
# static DCB configurations. 0x10 is static mode. Mark priorities 5 and 6 as
# subject to pfc
debugfs.sh -n ens6f0 -t dcbx set mode 0x10
debugfs.sh -n ens6f0 -t dcbx set pfc 5 1
debugfs.sh -n ens6f0 -t dcbx set pfc 6 1
\# set roce global overrides for qp params. enable exn and open QPs with dscp 20
debugfs.sh -n ens6f0 -t rdma glob ecn 1
debugfs.sh -n ens6f0 -t rdma glob dscp 20
# open some QPs (DSCP 20)
ib write bw -d qedr0 -q 16 -F -x 1 --run infinitely
# change global dscp qp params
debugfs.sh -n ens6f0 -t rdma glob dscp 22
# open some more QPs (DSCP 22)
ib_write_bw -d qedr0 -q 16 -F -x 1 -p 8000 --run_infinitely
```

observe PFCs being generated on multiple priorities
debugfs.sh -n ens6f0 -d phy mac stat -P 0 | grep "Class Based Flow Control"

Limitations

DCQCN has the following limitations:

- DCQCN mode currently supports only up to 64 QPs.
- Marvell adapters can determine vLAN priority for PFC purposes from vLAN priority or from DSCP bits in the ToS field. However, in the presence of both, vLAN takes precedence.

8 iWARP Configuration

Internet wide area RDMA protocol (iWARP) is a computer networking protocol that implements RDMA for efficient data transfer over IP networks. iWARP is designed for multiple environments, including LANs, storage networks, data center networks, and WANs.

This chapter provides instructions for:

- Preparing the Adapter for iWARP
- "Configuring iWARP on Windows" on page 178
- "Configuring iWARP on Linux" on page 182

NOTE

Some iWARP features may not be fully enabled in the current release. For details, refer to Appendix D Feature Constraints.

Preparing the Adapter for iWARP

This section provides instructions for preboot adapter iWARP configuration using the HII. For more information about preboot adapter configuration, see Chapter 5 Adapter Preboot Configuration.

To configure iWARP through HII in Default mode:

- Access the server BIOS System Setup, and then click **Device Settings**.
- 2. On the Device Settings page, select a port for the 25G 41xxx Series Adapter.
- 3. On the Main Configuration Page for the selected adapter, click **NIC Configuration**.
- 4. On the NIC Configuration page:
 - a. Set the NIC + RDMA Mode to Enabled.
 - b. Set the **RDMA Protocol Support** to **RoCE/iWARP** or **iWARP**.
 - c. Click Back.
- 5. On the Main Configuration Page, click **Finish**.

- 6. In the Warning Saving Changes message box, click **Yes** to save the configuration.
- 7. In the Success Saving Changes message box, click **OK**.
- 8. Repeat Step 2 through Step 7 to configure the NIC and iWARP for the other ports.
- 9. To complete adapter preparation of both ports:
 - a. On the Device Settings page, click **Finish**.
 - b. On the main menu, click **Finish**.
 - c. Exit to reboot the system.

Proceed to "Configuring iWARP on Windows" on page 178 or "Configuring iWARP on Linux" on page 182.

Configuring iWARP on Windows

This section provides procedures for enabling iWARP, verifying RDMA, and verifying iWARP traffic on Windows. For a list of OSs that support iWARP, see Table 7-1 on page 127.

To enable iWARP on the Windows host and verify RDMA:

- Enable iWARP on the Windows host.
 - a. Open the Windows Device Manager, and then open the 41xxx Series Adapter NDIS Miniport Properties.
 - b. On the FastLinQ Adapter properties, click the **Advanced** tab.
 - c. On the Advanced page under **Property**, do the following:
 - Select Network Direct Functionality, and then select Enabled for the Value.
 - Select NetworkDirect Technology, and then select iWARP for the Value.
 - d. Click **OK** to save your changes and close the adapter properties.

2. Using Windows PowerShell, verify that RDMA is enabled. The Get-NetAdapterRdma command output (Figure 8-1) shows the adapters that support RDMA.

Figure 8-1. Windows PowerShell Command: Get-NetAdapterRdma

3. Using Windows PowerShell, verify that NetworkDirect is enabled. The Get-NetOffloadGlobalSetting command output (Figure 8-2) shows NetworkDirect as Enabled.

```
PS C:\Users\Administrator> Get-NetOffloadGlobalSetting

ReceiveSideScaling : Enabled
ReceiveSegmentCoalescing : Enabled
Chimney : Disabled
TaskOffload : Enabled
NetworkDirect : Enabled
NetworkDirectAcrossIPSubnets : Blocked
PacketCoalescingFilter : Disabled
```

Figure 8-2. Windows PowerShell Command: Get-NetOffloadGlobalSetting

To verify iWARP traffic:

- 1. Map SMB drives and run iWARP traffic.
- 2. Launch Performance Monitor (Perfmon).
- 3. In the Add Counters dialog box, click **RDMA Activity**, and then select the adapter instances.

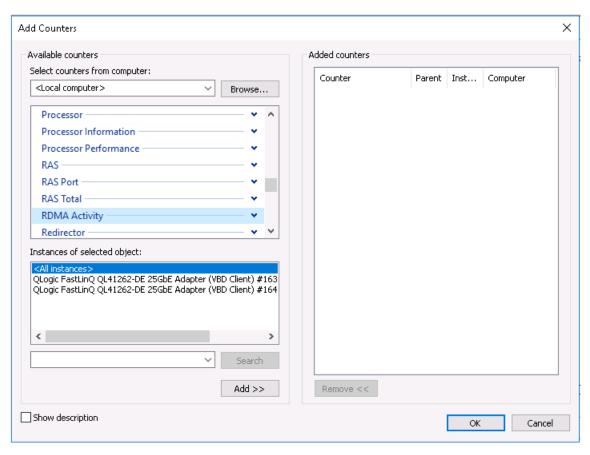


Figure 8-3 shows an example.

Figure 8-3. Perfmon: Add Counters

If iWARP traffic is running, counters appear as shown in the Figure 8-4 example.

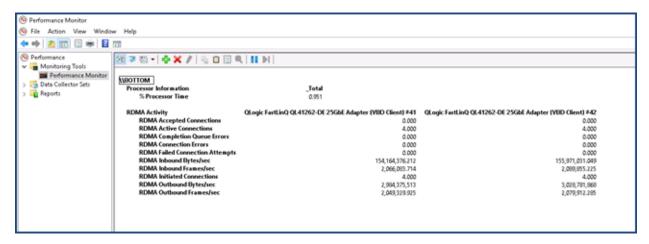


Figure 8-4. Perfmon: Verifying iWARP Traffic

NOTE

For more information on how to view Marvell RDMA counters in Windows, see "Viewing RDMA Counters" on page 136.

- 4. To verify the SMB connection:
 - a. At a command prompt, issue the net use command as follows:

C:\Users\Administrator> net use
New connections will be remembered.

Status Local Remote Network

OK F: \\192.168.10.10\Share1 Microsoft Windows Network

The command completed successfully.

b. Issue the netstat -xan command as follows, where Share1 is mapped as an SMB share:

C:\Users\Administrator> netstat -xan

Active NetworkDirect Connections, Listeners, ShareEndpoints

Mode	IfIndex	Type	Local Address	Foreign Address	PID
Kernel	56	Connection	192.168.11.20:16159	192.168.11.10:445	0
Kernel	56	Connection	192.168.11.20:15903	192.168.11.10:445	0
Kernel	56	Connection	192.168.11.20:16159	192.168.11.10:445	0

Kernel	56 Connection	192.168.11.20:15903 192.168.11.10:445	0
Kernel	60 Listener	[fe80::e11d:9ab5:a47d:4f0a%56]:445 NA	0
Kernel	60 Listener	192.168.11.20:445 NA	0
Kernel	60 Listener	[fe80::71ea:bdd2:ae41:b95f%60]:445 NA	0
Kernel	60 Listener	192 168 11 20.16159 192 168 11 10.445	Ο

Configuring iWARP on Linux

Marvell 41xxx Series Adapters support iWARP on the Linux Open Fabric Enterprise Distributions (OFEDs) listed in Table 7-1 on page 127.

iWARP configuration on a Linux system includes the following:

- Installing the Driver
- Configuring iWARP and RoCE
- Detecting the Device
- Supported iWARP Applications
- Running Perftest for iWARP
- Configuring NFS-RDMA

Installing the Driver

Install the RDMA drivers as shown in Chapter 3 Driver Installation.

Configuring iWARP and RoCE

NOTE

This procedure applies only if you previously selected **iWARP+RoCE** as the value for the RDMA Protocol Support parameter during preboot configuration using HII (see Configuring NIC Parameters, Step 5 on page 51).

To enable iWARP and RoCE:

- 1. Unload all FastLinQ drivers as follows:
 - # modprobe -r qedr or modprobe -r qede
- 2. Use the following command syntax to change the RDMA protocol by loading the qed driver with a port interface PCI ID (xx:xx.x) and an RDMA protocol value (p).
 - # modprobe -v qed rdma_protocol_map=<xx:xx.x-p>

The RDMA protocol (p) values are as follows:

- □ 0—Accept the default (RoCE)
- ☐ 1—No RDMA
- □ 2—**RoCE**
- □ 3—iWARP

For example, to change the interface on the port given by 04:00.0 from RoCE to iWARP, issue the following command:

```
# modprobe -v qed rdma_protocol_map=04:00.0-3
```

3. Load the RDMA driver by issuing the following command:

```
# modprobe -v qedr
```

The following example shows the command entries to change the RDMA protocol to iWARP on multiple NPAR interfaces:

```
# modprobe qed rdma protocol map=04:00.1-3,04:00.3-3,04:00.5-3,
04:00.7-3,04:01.1-3,04:01.3-3,04:01.5-3,04:01.7-3
# modprobe -v qedr
# ibv devinfo |grep iWARP
        transport:
                                         iWARP (1)
        transport:
                                         iWARP (1)
```

Detecting the Device

To detect the device:

1. To verify whether RDMA devices are detected, view the dmesg logs:

```
# dmesg | grep qedr
[10500.191047] qedr 0000:04:00.0: registered qedr0
[10500.221726] qedr 0000:04:00.1: registered qedr1
```

2. Issue the ibv devinfo command, and then verify the transport type.

If the command is successful, each PCI function will show a separate hca_id. For example (if checking the second port of the above dual-port adapter):

```
[root@localhost ~]# ibv_devinfo -d qedr1
hca id: qedr1
```

```
iWARP (1)
transport:
                                8.14.7.0
fw ver:
                                020e:1eff:fec4:c06e
node guid:
                                020e:1eff:fec4:c06e
sys image guid:
vendor id:
                                0x1077
                                5718
vendor part id:
                                0x0
hw ver:
phys port cnt:
       port: 1
               state:
                                      PORT ACTIVE (4)
                                        4096 (5)
               max mtu:
               active mtu:
                                        1024 (3)
                sm lid:
                port lid:
                port lmc:
                                       0x00
                link layer:
                                       Ethernet
```

Supported iWARP Applications

Linux-supported RDMA applications for iWARP include the following:

- ibv_devinfo, ib devices
- ib_send_bw/lat, ib_write_bw/lat, ib_read_bw/lat, ib_atomic_bw/lat For iWARP, all applications must use the RDMA communication manager (rdma_cm) using the ¬R option.
- rdma_server, rdma_client
- rdma xserver, rdma xclient
- rping
- NFS over RDMA (NFSoRDMA)
- iSER (for details, see Chapter 9 iSER Configuration)
- NVMe-oF (for details, see Chapter 13 NVMe-oF Configuration with RDMA)

Running Perftest for iWARP

All perftest tools are supported over the iWARP transport type. You must run the tools using the RDMA connection manager (with the -R option).

Example:

On one server, issue the following command (using the second port in this example):

```
# ib send bw -d qedr1 -F -R
```

2. On one client, issue the following command (using the second port in this example):

```
[root@localhost ~]# ib_send_bw -d qedr1 -F -R 192.168.11.3
```

```
Send BW Test
Dual-port : OFF Device : qedr1
Number of qps : 1
                     Transport type : IW
Connection type : RC
                     Using SRQ : OFF
TX depth : 128
CQ Moderation : 100
Mtu
         : 1024[B]
Link type
GID index
          : Ethernet
           : 0
Max inline data : 0[B]
rdma cm QPs : ON
Data ex. method : rdma cm
______
local address: LID 0000 QPN 0x0192 PSN 0xcde932
GID: 00:14:30:196:192:110:00:00:00:00:00:00:00:00:00:00
remote address: LID 0000 QPN 0x0098 PSN 0x46fffc
GID: 00:14:30:196:195:62:00:00:00:00:00:00:00:00:00:00
______
#bytes #iterations BW peak[MB/sec] BW average[MB/sec] MsgRate[Mpps]
                   2250.38
65536 1000
                                2250.36
                                              0.036006
```

NOTE

For latency applications (send/write), if the perftest version is the latest (for example, perftest-3.0-0.21.g21dc344.x86_64.rpm), use the supported inline size value: 0-128.

Configuring NFS-RDMA

NFS-RDMA for iWARP includes both server and client configuration steps.

To configure the NFS server:

- 1. Create an nfs-server directory and grant permission by issuing the following commands:
 - # mkdir /tmp/nfs-server
 - # chmod 777 /tmp/nfs-server
- 2. In the /etc/exports file for the directories that you must export using NFS-RDMA on the server, make the following entry:

```
/tmp/nfs-server *(rw,fsid=0,async,insecure,no root squash)
```

Ensure that you use a different file system identification (FSID) for each directory that you export.

- 3. Load the svcrdma module as follows:
 - # modprobe svcrdma
- 4. Load the service as follows:
 - ☐ For SLES, enable and start the NFS server alias:
 - # systemctl enable|start|status nfsserver
 - ☐ For RHEL, enable and start the NFS server and services:
 - # systemctl enable|start|status nfs
- 5. Include the default RDMA port 20049 into this file as follows:
 - # echo rdma 20049 > /proc/fs/nfsd/portlist
- 6. To make local directories available for NFS clients to mount, issue the exports command as follows:
 - # exportfs -v

To configure the NFS client:

NOTE

This procedure for NFS client configuration also applies to RoCE.

- 1. Create an nfs-client directory and grant permission by issuing the following commands:
 - # mkdir /tmp/nfs-client
 - # chmod 777 /tmp/nfs-client
- 2. Load the xprtrdma module as follows:
 - # modprobe xprtrdma
- 3. Mount the NFS file system as appropriate for your version:

For NFS Version 3:

mount -o rdma,port=20049 192.168.2.4:/tmp/nfs-server
/tmp/nfs-client

For NFS Version 4:

mount -t nfs4 -o rdma,port=20049 192.168.2.4:/tmp/nfs-server
/tmp/nfs-client

NOTE

The default port for NFSoRDMA is 20049. However, any other port that is aligned with the NFS client will also work.

- 4. Verify that the file system is mounted by issuing the mount command. Ensure that the RDMA port and file system versions are correct.
 - # mount |grep rdma

9 iSER Configuration

This chapter provides procedures for configuring iSCSI Extensions for RDMA (iSER) for Linux (RHEL and SLES) and VMware ESXi 6.7, including:

- Before You Begin
- "Configuring iSER for RHEL" on page 189
- "Configuring iSER for SLES 12 and Later" on page 192
- "Using iSER with iWARP on RHEL and SLES" on page 193
- "Optimizing Linux Performance" on page 195
- "Configuring iSER on ESXi 6.7" on page 196

Before You Begin

As you prepare to configure iSER, consider the following:

- iSER is supported only in inbox OFED for the following operating systems:
 - □ RHEL 8.x
 - □ RHEL 7.6 and later
 - □ SLES 12 SP4 and later
 - ☐ SLES 15 SP0 and later
 - VMware ESXi 6.7 U1
- After logging into the targets or while running I/O traffic, unloading the Linux RoCE gedr driver may crash the system.
- While running I/O, performing interface down/up tests or performing cable pull-tests can cause driver or iSER module errors that may crash the system. If this happens, reboot the system.

Configuring iSER for RHEL

To configure iSER for RHEL:

 Install inbox OFED as described in "RoCE Configuration for RHEL" on page 150.

NOTE

Out-of-box OFEDs are not supported for iSER because the ib_isert module is not available in the out-of-box OFED 3.18-2 GA/3.18-3 GA versions. The inbox ib_isert module does not work with any out-of-box OFED versions.

- 2. Unload any existing FastLinQ drivers as described in "Removing the Linux Drivers" on page 10.
- 3. Install the latest FastLinQ driver and libqedr packages as described in "Installing the Linux Drivers with RDMA" on page 14.
- 4. Load the RDMA services as follows:

```
systemctl start rdma
modprobe qedr
modprobe ib_iser
modprobe ib isert
```

- 5. Verify that all RDMA and iSER modules are loaded on the initiator and target devices by issuing the <code>lsmod | grep qed and lsmod | grep iser commands</code>.
- 6. Verify that there are separate hca_id instances by issuing the ibv devinfo command, as shown in Step 6 on page 153.
- 7. Check the RDMA connection on the initiator device and the target device.
 - a. On the initiator device, issue the following command:

```
rping -s -C 10 -v
```

b. On the target device, issue the following command:

```
rping -c -a 192.168.100.99 -C 10 -v
```

Figure 9-1 shows an example of a successful RDMA ping.

```
[root@localhost/home
[root@localhost home] # rping -c -a 192.168.100.99 -C 10 -v
ping data: rdma-ping-0: ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrr
ping data: rdma-ping-1: BCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrst
ping data: rdma-ping-2: CDEFGHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrst
ping data: rdma-ping-3: DEFGHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrst
ping data: rdma-ping-4: EFGHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrstuv
ping data: rdma-ping-5: FGHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrstuvw
ping data: rdma-ping-6: GHIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrstuvwx
ping data: rdma-ping-7: HIJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrstuvwxy
ping data: rdma-ping-8: IJKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-9: JKLMNOPQRSTUVWXYZ[\]^_'abcdefghijklmnopqrstuvwxyzA
client DISCONNECT EVENT...
[root@localhost home]#
```

Figure 9-1. RDMA Ping Successful

8. You can use a Linux TCM-LIO target to test iSER. The setup is the same for any iSCSI target, except that you issue the command enable_iser Boolean=true on the applicable portals. The portal instances are identified as iser in Figure 9-2.

Figure 9-2. iSER Portal Instances

- 9. Install Linux iSCSI Initiator Utilities using the yum install iscsi-initiator-utils commands.
 - a. To discover the iSER target, issue the <code>iscsiadm</code> command. For example:

iscsiadm -m discovery -t st -p 192.168.100.99:3260

b. To change the transport mode to iSER, issue the iscsiadm command. For example:

```
iscsiadm -m node -T iqn.2015-06.test.target1 -o update -n iface.transport name -v iser
```

c. To connect to or log in to the iSER target, issue the <code>iscsiadm</code> command. For example:

```
iscsiadm -m node -l -p 192.168.100.99:3260 -T
iqn.2015-06.test.target1
```

d. Confirm that the Iface Transport is iser in the target connection, as shown in Figure 9-3. Issue the iscsiadm command; for example:

```
iscsiadm -m session -P2
```

```
192.168.100.99:3260,1 iqn.2015-06.test.target1
192.168.100.99:3260,1 iqn.2015-06.test.target1
[root@localhost ~]# iscsiadm -m node -T iqn.2015-06.test.target1 -o update -n iface.transport name -v iser
[root@localhost ~]#
[root@localhost ~]# iscsiadm -m node -l -p 192.168.100.99:3260 -T iqn.2015-06.test.target1
ogging in to [iface: default, target: iqn.2015-06.test.target1, portal: 192.168.100.99,3260] (multiple)
ogin to [iface: default, target: iqn.2015-06.test.target1, portal: 192.168.100.99,3260] successful.
[root@localhost ~]#
root@localhost ~]# iscsiadm -m session -P2
Parget: iqn.2015-06.test.target1 (non-flash)
       Current Portal: 192.168.100.99:3260,1
               *******
               Iface Name: default

    ↓ Iface Transport: iser

                Iface Initiatorname: iqn.1994-05.com.redhat:c672dfb8b08f
               Iface IPaddress: <empty>
               Iface HWaddress: <empty>
               Iface Netdev: <empty>
               iSCSI Session State: LOGGED IN
               Internal iscsid Session State: NO CHANGE
               Timeouts:
                ******
               Recovery Timeout: 120
```

Figure 9-3. Iface Transport Confirmed

e. To check for a new iSCSI device, as shown in Figure 9-4, issue the lsscsi command.

```
root@localhost ~]# lsscsi
           disk
                             LOGICAL VOLUME
                                                    /dev/sdb
6:0:0:11
           disk
                   HP
                            LOGICAL VOLUME
                                                    /dev/sda
                             LOGICAL VOLUME
           storage HP
6:3:0:01
                             P440ar
           disk
                   LIO-ORG ram1
                                                    /dev/sdd
39:0:0:01
                                              4.0
coot@localhost ~]#
```

Figure 9-4. Checking for New iSCSI Device

Configuring iSER for SLES 12 and Later

Because the targetcli is not inbox on SLES 12 and later, you must complete the following procedure.

To configure iSER for SLES 12 and later:

1. Install targetcli.

For SLES 12:

Locate, copy and install the following RPMs from the ISO image (x86_64 and noarch location).

```
lio-utils-4.1-14.6.x86_64.rpm
python-configobj-4.7.2-18.10.noarch.rpm
python-PrettyTable-0.7.2-8.5.noarch.rpm
python-configshell-1.5-1.44.noarch.rpm
python-pyparsing-2.0.1-4.10.noarch.rpm
python-netifaces-0.8-6.55.x86_64.rpm
python-rtslib-2.2-6.6.noarch.rpm
python-urwid-1.1.1-6.144.x86_64.rpm
targetcli-2.1-3.8.x86 64.rpm
```

For SLES 15 and SLES 15 SP1:

Load the SLES Package DVD and install targetcli by issuing the following Zypper command, which installs all the dependency packages:

```
# zypper install python3-targetcli-fb
```

2. Before starting the targetcli, load all RoCE device drivers and iSER modules as follows:

```
# modprobe qed
# modprobe qede
# modprobe qedr
# modprobe ib_iser (initiator)
# modprobe ib isert (target)
```

- 3. Before configuring iSER targets, configure NIC interfaces and run L2 and RoCE traffic, as described in Step 7 on page 153.
- 4. For SLES 15 and SLES 15 SP1, insert the SLES Package DVD and install the targetcli utility. This command also installs all the dependency packages.
 - # zypper install python3-targetcli-fb
- 5. Start the targetcli utility, and configure your targets on the iSER target system.

NOTE

targetcli versions are different in RHEL and SLES. Be sure to use the proper backstores to configure your targets:

- RHEL uses ramdisk
- SLES uses rd mcp

Using iSER with iWARP on RHEL and SLES

Configure the iSER initiator and target similar to RoCE to work with iWARP. You can use different methods to create a Linux-IO (LIO^{TM}) target; one is listed in this section. You may encounter some difference in targetcli configuration in SLES 12 and RHEL 7.x because of the version.

To configure a target for LIO:

Create an LIO target using the targetcli utility. Issue the following command:

```
# targetcli
```

/> saveconfig

```
targetcli shell version 2.1.fb41
Copyright 2011-2013 by Datera, Inc and others.
For help on commands, type 'help'.
```

2. Issue the following commands:

```
/> /backstores/ramdisk create Ramdisk1-1 1g nullio=true
/> /iscsi create iqn.2017-04.com.org.iserport1.target1
/> /iscsi/iqn.2017-04.com.org.iserport1.target1/tpg1/luns create /backstores/ramdisk/Ramdisk1-1
/> /iscsi/iqn.2017-04.com.org.iserport1.target1/tpg1/portals/ create 192.168.21.4 ip_port=3261
/> /iscsi/iqn.2017-04.com.org.iserport1.target1/tpg1/portals/192.168.21.4:3261 enable_iser
boolean=true
/> /iscsi/iqn.2017-04.com.org.iserport1.target1/tpg1 set attribute authentication=0
demo_mode_write_protect=0 generate_node_acls=1 cache_dynamic_acls=1
```

Figure 9-5 shows the target configuration for LIO.

Figure 9-5. LIO Target Configuration

To configure an initiator for iWARP:

1. To discover the iSER LIO target using port 3261, issue the iscsiadm command as follows:

```
# iscsiadm -m discovery -t st -p 192.168.21.4:3261 -I iser 192.168.21.4:3261,1 iqn.2017-04.com.org.iserport1.target1
```

- 2. Change the transport mode to iser as follows:
- # iscsiadm -m node -o update -T iqn.2017-04.com.org.iserport1.target1 -n
 iface.transport name -v iser
 - 3. Log into the target using port 3261:
- # iscsiadm -m node -1 -p 192.168.21.4:3261 -T iqn.2017-04.com.org.iserport1.target1
 Logging in to [iface: iser, target: iqn.2017-04.com.org.iserport1.target1,
 portal: 192.168.21.4,3261] (multiple)
 Login to [iface: iser, target: iqn.2017-04.com.org.iserport1.target1, portal:
 192.168.21.4,3261] successful.
 - 4. Ensure that those LUNs are visible by issuing the following command:

lsscsi [1:0:0:0] P440ar 3.56 storage HP [1:1:0:0] disk LOGICAL VOLUME 3.56 /dev/sda ΗP UMD0 /dev/sr0 [6:0:0:0] cd/dvd hp DVD-ROM DUDON [7:0:0:0] disk LIO-ORG Ramdisk1-1 4.0 /dev/sdb

Optimizing Linux Performance

Consider the following Linux performance configuration enhancements described in this section.

- Configuring CPUs to Maximum Performance Mode
- Configuring Kernel sysctl Settings
- Configuring IRQ Affinity Settings
- Configuring Block Device Staging

Configuring CPUs to Maximum Performance Mode

Configure the CPU scaling governor to performance by using the following script to set all CPUs to maximum performance mode:

```
for CPUFREQ in
/sys/devices/system/cpu/cpu*/cpufreq/scaling_governor; do [ -f
$CPUFREQ ] || continue; echo -n performance > $CPUFREQ; done
```

Verify that all CPU cores are set to maximum performance mode by issuing the following command:

cat /sys/devices/system/cpu/cpu*/cpufreq/scaling_governor

Configuring Kernel sysctl Settings

Set the kernel sysctl settings as follows:

```
sysctl -w net.ipv4.tcp_mem="4194304 4194304 4194304"
sysctl -w net.ipv4.tcp_wmem="4096 65536 4194304"
sysctl -w net.ipv4.tcp_rmem="4096 87380 4194304"
sysctl -w net.core.wmem_max=4194304
sysctl -w net.core.rmem_max=4194304
sysctl -w net.core.wmem_default=4194304
sysctl -w net.core.rmem_default=4194304
sysctl -w net.core.netdev_max_backlog=250000
sysctl -w net.ipv4.tcp_timestamps=0
sysctl -w net.ipv4.tcp_timestamps=0
sysctl -w net.ipv4.tcp_low_latency=1
sysctl -w net.ipv4.tcp_low_latency=1
sysctl -w net.ipv4.tcp_adv_win_scale=1
echo 0 > /proc/sys/vm/nr_hugepages
```

Configuring IRQ Affinity Settings

The following example sets CPU core 0, 1, 2, and 3 to interrupt request (IRQ) XX, YY, ZZ, and XYZ respectively. Perform these steps for each IRQ assigned to a port (default is eight queues per port).

```
systemctl disable irqbalance
systemctl stop irqbalance
cat /proc/interrupts | grep qedr Shows IRQ assigned to each port queue
echo 1 > /proc/irq/XX/smp_affinity_list
echo 2 > /proc/irq/YY/smp_affinity_list
echo 4 > /proc/irq/ZZ/smp_affinity_list
echo 8 > /proc/irq/XYZ/smp affinity list
```

Configuring Block Device Staging

Set the block device staging settings for each iSCSI device or target as follows:

```
echo noop > /sys/block/sdd/queue/scheduler
echo 2 > /sys/block/sdd/queue/nomerges
echo 0 > /sys/block/sdd/queue/add_random
echo 1 > /sys/block/sdd/queue/rq affinity
```

Configuring iSER on ESXi 6.7

This section provides information for configuring iSER for VMware ESXi 6.7.

Before You Begin

Before you configure iSER for ESXi 6.7, ensure that the following is complete:

■ The CNA package with NIC and RoCE drivers is installed on the ESXi 6.7 system and the devices are listed. To view RDMA devices, issue the following command:

esxcli rdma device list

The iSER target is configured to communicate with the iSER initiator.

Configuring iSER for ESXi 6.7

esxcli rdma iser add

To configure iSER for ESXi 6.7:

1. Add iSER devices by issuing the following commands:

```
esxcli iscsi adapter list

Adapter Driver State UID Description

vmhba64 iser unbound iscsi.vmhba64 VMware iSCSI over RDMA (iSER) Adapter

vmhba65 iser unbound iscsi.vmhba65 VMware iSCSI over RDMA (iSER) Adapter
```

2. Disable the firewall as follows.

```
esxcli network firewall set --enabled=false
esxcli network firewall unload
vsish -e set /system/modules/iscsi_trans/loglevels/iscsitrans 0
vsish -e set /system/modules/iser/loglevels/debug 4
```

Create a standard vSwitch VMkernel port group and assign the IP:

esxcli network vswitch standard add -v vSwitch_iser1 esxcfg-nics -1

```
Name
     PCT
                Driver
                           Link Speed
                                          Duplex MAC Address
                                                                MTU
                                                                       Description
Broadcom
Corporation NetXtreme BCM5720 Gigabit Ethernet
vmnic1 0000:01:00.1 ntg3
                            Down OMbps Half e0:db:55:0c:5f:95 1500
                                                                      Broadcom
Corporation NetXtreme BCM5720 Gigabit Ethernet
                            Down OMbps Half e0:db:55:0c:5f:96 1500 Broadcom
vmnic2 0000:02:00.0 ntg3
Corporation NetXtreme BCM5720 Gigabit Ethernet
vmnic3 0000:02:00.1 ntg3
                           Down OMbps
                                         Half e0:db:55:0c:5f:97 1500 Broadcom
Corporation NetXtreme BCM5720 Gigabit Ethernet
vmnic4 0000:42:00.0 qedentv Up 40000Mbps Full 00:0e:1e:d5:f6:a2 1500 QLogic Corp.
QLogic FastLinQ QL41xxx 10/25/40/50/100 GbE Ethernet Adapter
vmnic5 0000:42:00.1 gedentv
                            Up 40000Mbps Full
                                                 00:0e:1e:d5:f6:a3 1500 QLogic Corp.
QLogic FastLinQ QL41xxx 10/25/40/50/100 GbE Ethernet Adapter
esxcli network vswitch standard uplink add -u vmnic5 -v vSwitch_iser1
esxcli network vswitch standard portgroup add -p "rdma_group1" -v vSwitch_iser1
esxcli network ip interface add -i vmk1 -p "rdma_group1"
esxcli network ip interface ipv4 set -i vmk1 -I 192.168.10.100 -N 255.255.255.0 -t static
esxcfg-vswitch -p "rdma group1" -v 4095 vSwitch iser1
```

```
esxcli iscsi networkportal add -A vmhba67 -n vmk1
esxcli iscsi networkportal list
esxcli iscsi adapter get -A vmhba65
vmhba65
  Name: iqn.1998-01.com.vmware:localhost.punelab.qlogic.com qlogic.org qlogic.com
mv.qlogic.com:1846573170:65
  Alias: iser-vmnic5
  Vendor: VMware
   Model: VMware iSCSI over RDMA (iSER) Adapter
   Description: VMware iSCSI over RDMA (iSER) Adapter
   Serial Number: vmnic5
   Hardware Version:
   Asic Version:
   Firmware Version:
   Option Rom Version:
   Driver Name: iser-vmnic5
   Driver Version:
   TCP Protocol Supported: false
   Bidirectional Transfers Supported: false
  Maximum Cdb Length: 64
   Can Be NIC: true
   Is NIC: true
   Is Initiator: true
   Is Target: false
   Using TCP Offload Engine: true
   Using ISCSI Offload Engine: true
```

4. Add the target to the iSER initiator as follows:

```
esxcli iscsi adapter target list
esxcli iscsi adapter discovery sendtarget add -A vmhba65 -a 192.168.10.11
esxcli iscsi adapter target list

Adapter Target Alias Discovery Method Last Error
----- ----- SENDTARGETS No Error
esxcli storage core adapter rescan --adapter vmhba65
```

List the attached target as follows:

```
esxcfg-scsidevs -1
mpx.vmhba0:C0:T4:L0
  Device Type: CD-ROM
  Size: 0 MB
  Display Name: Local TSSTcorp CD-ROM (mpx.vmhba0:C0:T4:L0)
  Multipath Plugin: NMP
```

```
Console Device: /vmfs/devices/cdrom/mpx.vmhba0:C0:T4:L0
  Devfs Path: /vmfs/devices/cdrom/mpx.vmhba0:C0:T4:L0
  Vendor: TSSTcorp Model: DVD-ROM SN-108BB Revis: D150
  SCSI Level: 5 Is Pseudo: false Status: on
  Is RDM Capable: false Is Removable: true
  Is Local: true Is SSD: false
  Other Names:
     vml.0005000000766d686261303a343a30
  VAAI Status: unsupported
naa.6001405e81ae36b771c418b89c85dae0
  Device Type: Direct-Access
  Size: 512 MB
  Display Name: LIO-ORG iSCSI Disk (naa.6001405e81ae36b771c418b89c85dae0)
  Multipath Plugin: NMP
  Console Device: /vmfs/devices/disks/naa.6001405e81ae36b771c418b89c85dae0
  Devfs Path: /vmfs/devices/disks/naa.6001405e81ae36b771c418b89c85dae0
  Vendor: LIO-ORG Model: ram1
                                            Revis: 4.0
  SCSI Level: 5 Is Pseudo: false Status: degraded
  Is RDM Capable: true Is Removable: false
  Is Local: false Is SSD: false
  Other Names:
     vml.02000000006001405e81ae36b771c418b89c85dae072616d312020
  VAAI Status: supported
naa.690b11c0159d050018255e2d1d59b612
```

10 iSCSI Configuration

This chapter provides the following iSCSI configuration information:

- iSCSI Boot
- "iSCSI Offload in Windows Server" on page 200
- "iSCSI Offload in Linux Environments" on page 209

NOTE

Some iSCSI features may not be fully enabled in the current release. For details, refer to Appendix D Feature Constraints.

To enable iSCSI-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

iSCSI Boot

Marvell 4xxxx Series gigabit Ethernet (GbE) adapters support iSCSI boot to enable 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.

Jumbo frames with iSCSI boot are supported only on Windows OSs, when the adapter is used as either an NDIS or HBA offload device.

For iSCSI boot from SAN information, see Chapter 6 Boot from SAN Configuration.

iSCSI Offload in Windows Server

iSCSI offload is a technology that offloads iSCSI protocol processing overhead from host processors to the iSCSI HBA. iSCSI offload increases network performance and throughput while helping to optimize server processor use. This section covers how to configure the Windows iSCSI offload feature for the Marvell 41xxx Series Adapters.

With the proper iSCSI offload licensing, you can configure your iSCSI-capable 41xxx Series Adapter to offload iSCSI processing from the host processor. The following sections describe how to enable the system to take advantage of Marvell's iSCSI offload feature:

- Installing Marvell Drivers
- Installing the Microsoft iSCSI Initiator
- Configuring Microsoft Initiator to Use Marvell's iSCSI Offload
- iSCSI Offload FAQs
- Windows Server 2012 R2, 2016, and 2019 iSCSI Boot Installation
- iSCSI Crash Dump

Installing Marvell Drivers

Install the Windows drivers as described in "Installing Windows Driver Software" on page 18.

Installing the Microsoft iSCSI Initiator

Launch the Microsoft iSCSI initiator applet. At the first launch, the system prompts for an automatic service start. Confirm the selection for the applet to launch.

Configuring Microsoft Initiator to Use Marvell's iSCSI Offload

After the IP address is configured for the iSCSI adapter, you must use Microsoft Initiator to configure and add a connection to the iSCSI target using the Marvell FastLinQ iSCSI adapter. For more details on Microsoft Initiator, see the Microsoft user guide.

To configure Microsoft Initiator:

- Open Microsoft Initiator.
- 2. To configure the initiator IQN name according to your setup, follow these steps:
 - a. On the iSCSI Initiator Properties, click the **Configuration** tab.
 - b. On the Configuration page (Figure 10-1), click **Change** next to **To** modify the initiator name.

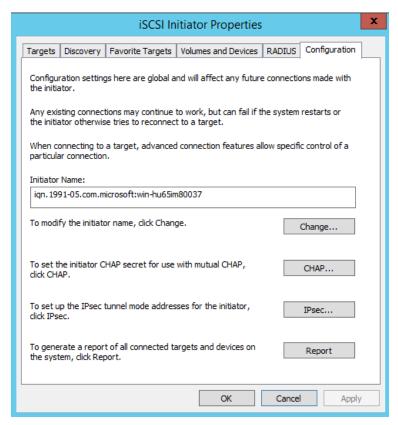


Figure 10-1. iSCSI Initiator Properties, Configuration Page

c. In the iSCSI Initiator Name dialog box, type the new initiator IQN name, and then click **OK**. (Figure 10-2)

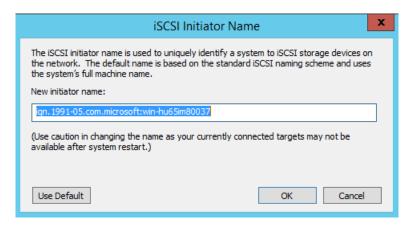


Figure 10-2. iSCSI Initiator Node Name Change

3. On the iSCSI Initiator Properties, click the **Discovery** tab.

iSCSI Initiator Properties Targets Discovery Favorite Targets Volumes and Devices RADIUS Configuration Target portals Refresh The system will look for \underline{T} argets on following portals: Port Address Adapter IP address To add a target portal, click Discover Portal. Discover Portal... To remove a target portal, select the address above and Remove . then click Remove. iSNS servers Refresh The system is registered on the following iSNS servers: Name 192.168.2.60 To add an iSNS server, click Add Server. Add Server... To remove an iSNS server, select the server above and Remove then dick Remove. OK Cancel **Apply**

4. On the Discovery page (Figure 10-3) under **Target portals**, click **Discover Portal**.

Figure 10-3. iSCSI Initiator—Discover Target Portal

- 5. In the Discover Target Portal dialog box (Figure 10-4):
 - a. In the **IP address or DNS name** box, type the IP address of the target.
 - b. Click Advanced.

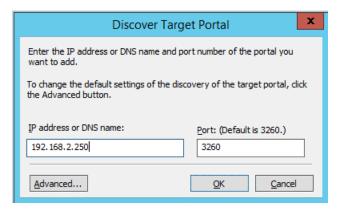


Figure 10-4. Target Portal IP Address

- 6. In the Advanced Settings dialog box (Figure 10-5), complete the following under **Connect using**:
 - a. For Local adapter, select the QLogic <name or model> Adapter.
 - b. For **Initiator IP**, select the adapter IP address.
 - c. Click OK.

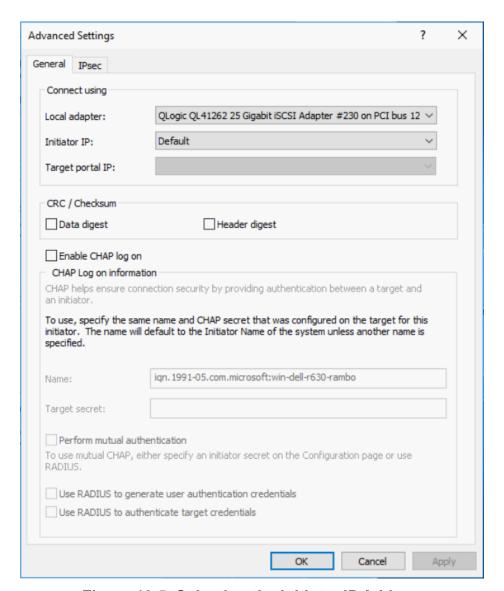


Figure 10-5. Selecting the Initiator IP Address

7. On the iSCSI Initiator Properties, Discovery page, click **OK**.

8. Click the **Targets** tab, and then on the Targets page (Figure 10-6), click **Connect**.

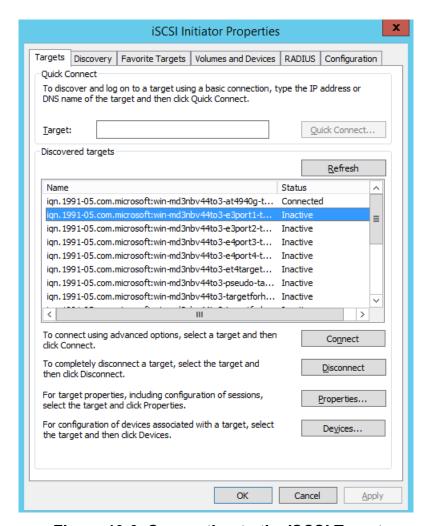


Figure 10-6. Connecting to the iSCSI Target

Connect To Target

Target name:

iqn. 1991-05.com.microsoft:win-md3nbv44to3-e3port1-target

✓ Add this connection to the list of Favorite Targets.

This will make the system automatically attempt to restore the connection every time this computer restarts.

□ Enable multi-path

Advanced...

OK Cancel

9. On the Connect To Target dialog box (Figure 10-7), click **Advanced**.

Figure 10-7. Connect To Target Dialog Box

- In the Local Adapter dialog box, select the QLogic <name or model>
 Adapter, and then click OK.
- 11. Click **OK** again to close Microsoft Initiator.
- 12. To format the iSCSI partition, use Disk Manager.

NOTE

Some limitations of the teaming functionality include:

- Teaming does not support iSCSI adapters.
- Teaming does not support NDIS adapters that are in the boot path.
- Teaming supports NDIS adapters that are not in the iSCSI boot path, but only for the switch-independent NIC team type.
- Switch dependent teaming (IEEE 802.3ad LACP and Generic/Static Link Aggregation (Trunking) cannot use a switch independent partitioned virtual adapter. IEEE standards require Switch Dependent Teaming (IEEE 802.3ad LACP and Generic/Static Link Aggregation (Trunking)) mode to work per the entire port instead of just the MAC address (fraction of a port) granularity.
- Microsoft recommends using their in-OS NIC teaming service instead of any adapter vendor-proprietary NIC teaming driver on Windows Server 2012 and later.

iSCSI Offload FAQs

Some of the frequently asked questions about iSCSI offload include:

Question: How do I assign an IP address for iSCSI offload?

Answer: Use the Configurations page in QConvergeConsole GUI.

Question: What tools should I use to create the connection to the target? **Answer:** Use Microsoft iSCSI Software Initiator (version 2.08 or later).

Question: How do I know that the connection is offloaded?

Answer: Use Microsoft iSCSI Software Initiator. From a command line, type

oiscsicli sessionlist. From Initiator Name, an iSCSI

offloaded connection will display an entry beginning with B06BDRV. A non-offloaded connection displays an entry beginning with Root.

Question: What configurations should be avoided?

Answer: The IP address should not be the same as the LAN.

Windows Server 2012 R2, 2016, and 2019 iSCSI Boot Installation

Windows Server 2012 R2, Windows Server 2016, and Windows Server 2019 support booting and installing in either the offload or non-offload paths. Marvell requires that you use a slipstream DVD with the latest Marvell drivers injected. See "Injecting (Slipstreaming) Adapter Drivers into Windows Image Files" on page 122.

The following procedure prepares the image for installation and booting in either the offload or non-offload path.

To set up Windows Server 2012 R2/2016/2019 iSCSI boot:

- 1. Remove any local hard drives on the system to be booted (remote system).
- 2. Prepare the Windows OS installation media by following the slipstreaming steps in "Injecting (Slipstreaming) Adapter Drivers into Windows Image Files" on page 122.
- 3. Load the latest Marvell iSCSI boot images into the NVRAM of the adapter.
- 4. Configure the iSCSI target to allow a connection from the remote device. Ensure that the target has sufficient disk space to hold the new OS installation.
- 5. Configure the UEFI HII to set the iSCSI boot type (offload or non-offload), correct initiator, and target parameters for iSCSI boot.
- 6. Save the settings and reboot the system. The remote system should connect to the iSCSI target and then boot from the DVD-ROM device.
- 7. Boot from DVD and begin installation.
- 8. Follow the on-screen instructions.

At the window that shows the list of disks available for the installation, the iSCSI target disk should be visible. This target is a disk connected through the iSCSI boot protocol and located in the remote iSCSI target.

- 9. To proceed with Windows Server 2012 R2/2016 installation, click **Next**, and then follow the on-screen instructions. The server will undergo a reboot multiple times as part of the installation process.
- 10. After the server boots to the OS, you should run the driver installer to complete the Marvell drivers and application installation.

iSCSI Crash Dump

Crash dump functionality is supported for both non-offload and offload iSCSI boot for the 41xxx Series Adapters. No additional configurations are required to configure iSCSI crash dump generation.

iSCSI Offload in Linux Environments

The Marvell FastLinQ 41xxx iSCSI software consists of a single kernel module called <code>qedi.ko</code> (qedi). The qedi module is dependent on additional parts of the Linux kernel for specific functionality:

- **qed.ko** is the Linux eCore kernel module used for common Marvell FastLinQ 41xxx hardware initialization routines.
- **scsi_transport_iscsi.ko** is the Linux iSCSI transport library used for upcall and downcall for session management.
- libiscsi.ko is the Linux iSCSI library function needed for protocol data unit (PDU) and task processing, as well as session memory management.
- iscsi_boot_sysfs.ko is the Linux iSCSI sysfs interface that provides helpers to export iSCSI boot information.
- uio.ko is the Linux Userspace I/O interface, used for light L2 memory mapping for iscsiuio.

These modules must be loaded before qedi can be functional. Otherwise, you might encounter an "unresolved symbol" error. If the qedi module is installed in the distribution update path, the requisite is automatically loaded by modprobe.

This section provides the following information about iSCSI offload in Linux:

- Differences from bnx2i
- Configuring qedi.ko
- Verifying iSCSI Interfaces in Linux

Differences from bnx2i

Some key differences exist between qedi—the driver for the Marvell FastLinQ 41xxx Series Adapter (iSCSI)—and the previous Marvell iSCSI offload driver—bnx2i for the Marvell 8400 Series Adapters. Some of these differences include:

- qedi directly binds to a PCI function exposed by the CNA.
- gedi does not sit on top of the net device.
- qedi is not dependent on a network driver such as bnx2x and cnic.
- qedi is not dependent on cnic, but it has dependency on qed.
- qedi is responsible for exporting boot information in sysfs using iscsi_boot_sysfs.ko, whereas bnx2i boot from SAN relies on the iscsi_ibft.ko module for exporting boot information.

Configuring qedi.ko

The qedi driver automatically binds to the exposed iSCSI functions of the CNA, and the target discovery and binding is done through the Open-iSCSI tools. This functionality and operation is similar to that of the bnx2i driver.

To load the <code>qedi.ko</code> kernel module, issue the following commands:

```
# modprobe qed
# modprobe libiscsi
# modprobe uio
# modprobe iscsi_boot_sysfs
# modprobe qedi
```

Verifying iSCSI Interfaces in Linux

After installing and loading the qedi kernel module, you must verify that the iSCSI interfaces were detected correctly.

To verify iSCSI interfaces in Linux:

1. To verify that the qedi and associated kernel modules are actively loaded, issue the following command:

```
# lsmod | grep qedi
qedi 114578 2
qed 697989 1 qedi
uio 19259 4 cnic,qedi
libiscsi 57233 2 qedi,bnx2i
scsi transport iscsi 99909 5 qedi,bnx2i,libiscsi
```

iscsi boot sysfs 16000 1 qedi

2. To verify that the iSCSI interfaces were detected properly, issue the following command. In this example, two iSCSI CNA devices are detected with SCSI host numbers 4 and 5.

dmesg | grep qedi

```
[0000:00:00.0]:[qedi_init:3696]: QLogic iSCSI Offload Driver v8.15.6.0.
....
[0000:42:00.4]:[__qedi_probe:3563]:59: QLogic FastLinQ iSCSI Module qedi 8.15.6.0, FW 8.15.3.0
....
[0000:42:00.4]:[qedi_link_update:928]:59: Link Up event.
....
[0000:42:00.5]:[__qedi_probe:3563]:60: QLogic FastLinQ iSCSI Module qedi 8.15.6.0, FW 8.15.3.0
....
[0000:42:00.5]:[qedi_link_update:928]:59: Link Up event
```

Use Open-iSCSI tools to verify that the IP is configured properly. Issue the following command:

iscsiadm -m iface | grep qedi

```
qedi.00:0e:le:c4:e1:6d
qedi,00:0e:le:c4:e1:6d,192.168.101.227,<empty>,iqn.1994-05.com.redhat:534ca9b6adf
qedi.00:0e:le:c4:e1:6c
qedi,00:0e:le:c4:e1:6c,192.168.25.91,<empty>,iqn.1994-05.com.redhat:534ca9b6adf
```

4. To ensure that the iscsiulo service is running, issue the following command:

systemctl status iscsiuio.service

5. To discover the iSCSI target, issue the iscsiadm command:

```
#iscsiadm -m discovery -t st -p 192.168.25.100 -I qedi.00:0e:1e:c4:e1:6c
192.168.25.100:3260,1 iqn.2003-
04.com.sanblaze:virtualun.virtualun.target-05000007
192.168.25.100:3260,1 iqn.2003-04.com.sanblaze:virtualun.virtualun.target-05000012
```

```
192.168.25.100:3260,1 iqn.2003-04.com.sanblaze:virtualun.virtualun.target-0500000c
192.168.25.100:3260,1 iqn.2003-
04.com.sanblaze:virtualun.virtualun.target-05000001
192.168.25.100:3260,1 iqn.2003-04.com.sanblaze:virtualun.virtualun.target-05000002
```

 Log into the iSCSI target using the IQN obtained in Step 5. To initiate the login procedure, issue the following command (where the last character in the command is a lowercase letter "L"):

```
#iscsiadm -m node -p 192.168.25.100 -T
iqn.2003-04.com.sanblaze:virtualun.virtualun.target-0000007 -1
Logging in to [iface: qedi.00:0e:le:c4:e1:6c,
target:iqn.2003-04.com.sanblaze:virtualun.virtualun.target-05000007, portal:192.168.25.100,3260]
(multiple)
Login to [iface: qedi.00:0e:le:c4:e1:6c, target:iqn.2003-
04.com.sanblaze:virtualun.virtualun.target-05000007, portal:192.168.25.100,3260] successful.
```

7. To verify that the iSCSI session was created, issue the following command:

```
# iscsiadm -m session
```

```
qedi: [297] 192.168.25.100:3260,1
iqn.2003-04.com.sanblaze:virtualun.virtualun.target-05000007 (non-flash)
```

8. To check for iSCSI devices, issue the iscsiadm command:

For advanced target configurations, refer to the Open-iSCSI README at:

https://github.com/open-iscsi/open-iscsi/blob/master/README

11 FCoE Configuration

This chapter provides the following Fibre Channel over Ethernet (FCoE) configuration information:

"Configuring Linux FCoE Offload" on page 213

NOTE

FCoE offload is supported on all 41xxx Series Adapters. Some FCoE features may not be fully enabled in the current release. For details, refer to Appendix D Feature Constraints.

To enable iSCSI-Offload mode, see the *Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters* at https://www.marvell.com/documents/5aa5otcbkr0im3ynera3/.

For FCoE boot from SAN information, see Chapter 6 Boot from SAN Configuration.

Configuring Linux FCoE Offload

The Marvell FastLinQ 41xxx Series Adapter FCoE software consists of a single kernel module called qedf.ko (qedf). The qedf module is dependent on additional parts of the Linux kernel for specific functionality:

- **qed.ko** is the Linux eCore kernel module used for common Marvell FastLinQ 41xxx hardware initialization routines.
- libfcoe.ko is the Linux FCoE kernel library needed to conduct FCoE forwarder (FCF) solicitation and FCoE initialization protocol (FIP) fabric login (FLOGI).
- libfc.ko is the Linux FC kernel library needed for several functions, including:
 - □ Name server login and registration
 - □ rport session management
- scsi_transport_fc.ko is the Linux FC SCSI transport library used for remote port and SCSI target management.

These modules must be loaded before qedf can be functional, otherwise errors such as "unresolved symbol" can result. If the qedf module is installed in the distribution update path, the requisite modules are automatically loaded by modprobe. Marvell FastLinQ 41xxx Series Adapters support FCoE offload.

This section provides the following information about FCoE offload in Linux:

- Differences Between gedf and bnx2fc
- Configuring gedf.ko
- Verifying FCoE Devices in Linux

Differences Between qedf and bnx2fc

Significant differences exist between qedf—the driver for the Marvell FastLinQ 41xxx 10/25GbE Controller (FCoE)—and the previous Marvell FCoE offload driver, bnx2fc. Differences include:

- qedf directly binds to a PCI function exposed by the CNA.
- qedf does not need the open-fcoe user space tools (fipvlan, fcoemon, fcoeadm) to initiate discovery.
- qedf issues FIP vLAN requests directly and does not need the fipvlan utility.
- gedf does not need an FCoE interface created by fipvlan for fcoemon.
- qedf does not sit on top of the net_device.
- gedf is not dependent on network drivers (such as bnx2x and cnic).
- qedf will automatically initiate FCoE discovery on link up (because it is not dependent on fipvlan or fcoemon for FCoE interface creation).

NOTE

FCoE interfaces no longer sit on top of the network interface. The qedf driver automatically creates FCoE interfaces that are separate from the network interface. Thus, FCoE interfaces do not show up in the FCoE interface dialog box in the installer. Instead, the disks show up automatically as SCSI disks, similar to the way Fibre Channel drivers work.

Configuring qedf.ko

No explicit configuration is required for qedf.ko. The driver automatically binds to the exposed FCoE functions of the CNA and begins discovery. This functionality is similar to the functionality and operation of the Marvell FC driver, qla2xx, as opposed to the older bnx2fc driver.

NOTE

For more information on FastLinQ driver installation, see Chapter 3 Driver Installation.

The load qedf.ko kernel module performs the following:

```
# modprobe qed
# modprobe libfcoe
# modprobe qedf
```

Verifying FCoE Devices in Linux

Follow these steps to verify that the FCoE devices were detected correctly after installing and loading the qedf kernel module.

To verify FCoE devices in Linux:

1. Check Ismod to verify that the qedf and associated kernel modules were loaded:

lsmod | grep qedf

```
69632 1 qedf libfc

143360 2 qedf,libfcoe scsi_transport_fc

65536 2 qedf,libfc qed

806912 1 qedf scsi_mod

262144 14 sg,hpsa,qedf,scsi_dh_alua,scsi_dh_rdac,dm_multipath,
scsi_transport_fc,scsi_transport_sas,libfc,scsi_transport_iscsi,scsi_dh_emc,
libata,sd_mod,sr_mod
```

2. Check dmesg to verify that the FCoE devices were detected properly. In this example, the two detected FCoE CNA devices are SCSI host numbers 4 and 5.

dmesg | grep qedf

```
[ 235.321185] [0000:00:00.0]: [qedf_init:3728]: QLogic FCoE Offload Driver
v8.18.8.0.
....
[ 235.322253] [0000:21:00.2]: [__qedf_probe:3142]:4: QLogic FastLinQ FCoE
Module qedf 8.18.8.0, FW 8.18.10.0
[ 235.606443] scsi host4: qedf
....
[ 235.624337] [0000:21:00.3]: [__qedf_probe:3142]:5: QLogic FastLinQ FCoE
Module qedf 8.18.8.0, FW 8.18.10.0
[ 235.886681] scsi host5: qedf
```

[243.991851] [0000:21:00.3]: [qedf link update:489]:5: LINK UP (40 GB/s).

3. Check for discovered FCoE devices using the lsscsi or lsblk -s commands. An example of each command follows.

# lsscsi							
[0:2:0:0]		disk	DELL	PERC H700		2.10	/dev/sda
[2:0:	:0:0]	cd/dvd	TEAC	DVD-ROM	DV-28SW	R.2A	/dev/sr0
[151:	:0:0:0]	disk	HP	P2000G3	FC/iSCSI	T252	/dev/sdb
[151:	:0:0:1]	disk	HP	P2000G3	FC/iSCSI	T252	/dev/sdc
[151:	:0:0:2]	disk	HP	P2000G3	FC/iSCSI	T252	/dev/sdd
[151:	:0:0:3]	disk	HP	P2000G3	FC/iSCSI	T252	/dev/sde
[151:	:0:0:4]	disk	HP	P2000G3	FC/iSCSI	T252	/dev/sdf
# lsblk -S							
NAME	HCTL	TYPE	VENDOR	MODEI	L	REV	TRAN
sdb	5:0:0:0	disk	SANBlaz	e VLUN	P2T1L0	V7.	3 fc
sdc	5:0:0:1	disk	SANBlaz	e VLUN	P2T1L1	V7.	3 fc
sdd	5:0:0:2	disk	SANBlaz	e VLUN	P2T1L2	V7.	3 fc
sde	5:0:0:3	disk	SANBlaz	e VLUN	P2T1L3	V7.	3 fc
sdf	5:0:0:4	disk	SANBlaz	e VLUN	P2T1L4	V7.	3 fc
sdg	5:0:0:5	disk	SANBlaz	e VLUN	P2T1L5	V7.	3 fc
sdh	5:0:0:6	disk	SANBlaz	e VLUN	P2T1L6	V7.	3 fc
sdi	5:0:0:7	disk	SANBlaz	e VLUN	P2T1L7	V7.	3 fc
sdj	5:0:0:8	disk	SANBlaz	e VLUN	P2T1L8	V7.	3 fc

Configuration information for the host is located in $/sys/class/fc_host/hostX$, where x is the number of the SCSI host. In the preceding example, x is 4. The hostX file contains attributes for the FCoE function, such as worldwide port name and fabric ID.

disk SANBlaze VLUN P2T1L9

sdk 5:0:0:9

V7.3 fc

12 SR-IOV Configuration

Single root input/output virtualization (SR-IOV) is a specification by the PCI SIG that enables a single PCI Express (PCIe) device to appear as multiple, separate physical PCIe devices. SR-IOV permits isolation of PCIe resources for performance, interoperability, and manageability.

NOTE

Some SR-IOV features may not be fully enabled in the current release.

This chapter provides instructions for:

- Configuring SR-IOV on Windows
- "Configuring SR-IOV on Linux" on page 224
- "Configuring SR-IOV on VMware" on page 230

Configuring SR-IOV on Windows

To configure SR-IOV on Windows:

- Access the server BIOS System Setup, and then click System BIOS Settings.
- 2. On the System BIOS Settings page, click Integrated Devices.
- 3. On the Integrated Devices page (Figure 12-1):
 - a. Set the SR-IOV Global Enable option to Enabled.
 - b. Click Back.

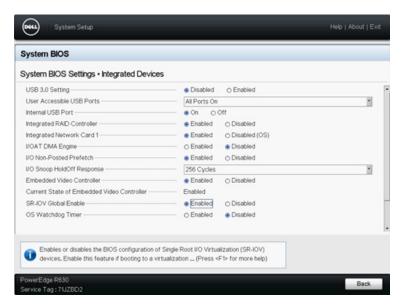


Figure 12-1. System Setup for SR-IOV: Integrated Devices

- 4. On the Main Configuration Page for the selected adapter, click **Device Level Configuration**.
- 5. On the Main Configuration Page Device Level Configuration (Figure 12-2):
 - Set the Virtualization Mode to SR-IOV, or NPAR+SR-IOV if you are using NPAR mode.
 - b. Click Back.



Figure 12-2. System Setup for SR-IOV: Device Level Configuration

- 6. On the Main Configuration Page, click **Finish**.
- 7. In the Warning Saving Changes message box, click **Yes** to save the configuration.
- 8. In the Success Saving Changes message box, click **OK**.

- 9. To enable SR-IOV on the miniport adapter:
 - a. Access Device Manager.
 - b. Open the miniport adapter properties, and then click the **Advanced** tab.
 - c. On the Advanced properties page (Figure 12-3) under **Property**, select **SR-IOV**, and then set the value to **Enabled**.
 - d. Click OK.

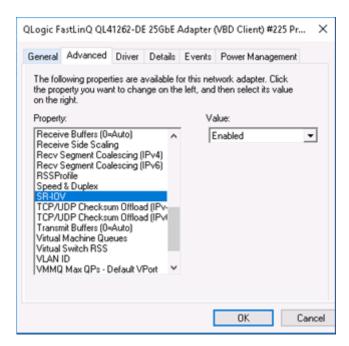


Figure 12-3. Adapter Properties, Advanced: Enabling SR-IOV

- 10. To create a Virtual Machine Switch (vSwitch) with SR-IOV (Figure 12-4 on page 220):
 - a. Launch the Hyper-V Manager.
 - b. Select Virtual Switch Manager.
 - c. In the **Name** box, type a name for the virtual switch.
 - d. Under Connection type, select External network.
 - e. Select the **Enable single-root I/O virtualization (SR-IOV)** check box, and then click **Apply**.

NOTE

Be sure to enable SR-IOV when you create the vSwitch. This option is unavailable after the vSwitch is created.

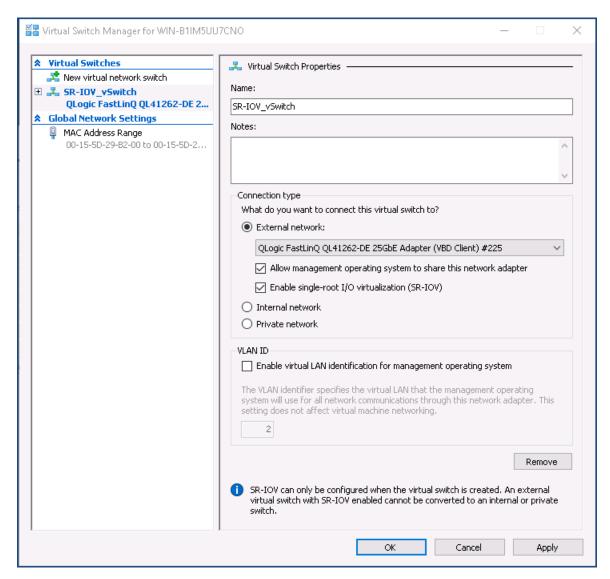


Figure 12-4. Virtual Switch Manager: Enabling SR-IOV

f. The Apply Networking Changes message box advises you that **Pending changes may disrupt network connectivity**. To save your changes and continue, click **Yes**.

11. To get the virtual machine switch capability, issue the following Windows PowerShell command:

PS C:\Users\Administrator> Get-VMSwitch -Name SR-IOV vSwitch | fl

Output of the Get-VMSwitch command includes the following SR-IOV capabilities:

IovVirtualFunctionCount : 80
IovVirtualFunctionsInUse : 1

- 12. To create a virtual machine (VM) and export the virtual function (VF) in the VM:
 - a. Create a virtual machine.
 - b. Add the VMNetworkadapter to the virtual machine.
 - c. Assign a virtual switch to the VMNetworkadapter.
 - d. In the Settings for VM <VM_Name> dialog box (Figure 12-5), Hardware Acceleration page, under Single-root I/O virtualization, select the Enable SR-IOV check box, and then click OK.

NOTE

After the virtual adapter connection is created, the SR-IOV setting can be enabled or disabled at any time (even while traffic is running).

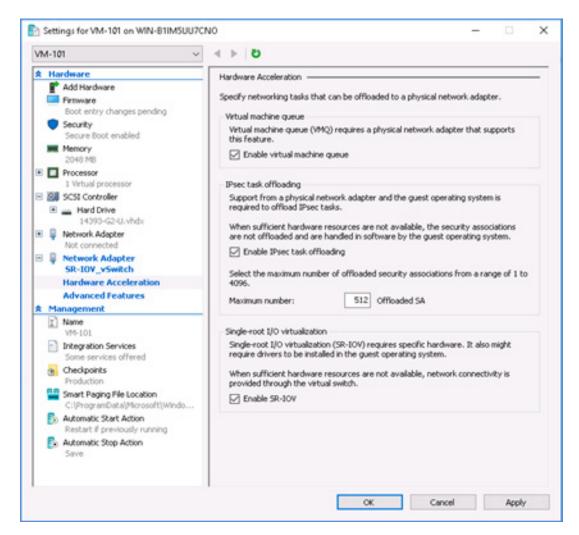
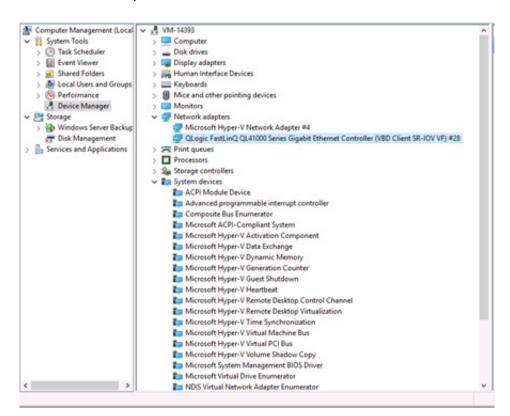


Figure 12-5. Settings for VM: Enabling SR-IOV

13. Install the Marvell drivers for the adapters detected in the VM. Use the latest drivers available from your vendor for your host OS (do not use inbox drivers).

NOTE

Be sure to use the same driver package on both the VM and the host system. For example, use the same qeVBD and qeND driver version on the Windows VM and in the Windows Hyper-V host.



After installing the drivers, the adapter is listed in the VM. Figure 12-6 shows an example.

Figure 12-6. Device Manager: VM with QLogic Adapter

14. To view the SR-IOV VF details, issue the following Windows PowerShell command:

PS C:\Users\Administrator> Get-NetadapterSriovVf

Figure 12-7 shows example output.



Figure 12-7. Windows PowerShell Command: Get-NetadapterSriovVf

Configuring SR-IOV on Linux

To configure SR-IOV on Linux:

- Access the server BIOS System Setup, and then click System BIOS Settings.
- 2. On the System BIOS Settings page, click **Integrated Devices**.
- 3. On the System Integrated Devices page (see Figure 12-1 on page 218):
 - a. Set the SR-IOV Global Enable option to Enabled.
 - b. Click **Back**.
- 4. On the System BIOS Settings page, click **Processor Settings**.
- 5. On the Processor Settings (Figure 12-8) page:
 - a. Set the Virtualization Technology option to Enabled.
 - b. Click Back.

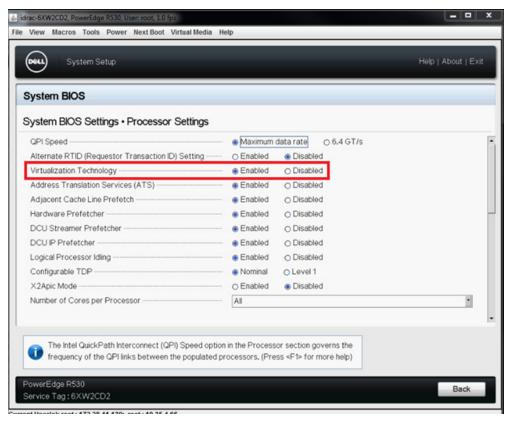


Figure 12-8. System Setup: Processor Settings for SR-IOV

6. On the System Setup page, select **Device Settings**.

- 7. On the Device Settings page, select **Port 1** for the Marvell adapter.
- 8. On the Device Level Configuration page (Figure 12-9):
 - a. Set the Virtualization Mode to SR-IOV.
 - b. Click Back.



Figure 12-9. System Setup for SR-IOV: Integrated Devices

- 9. On the Main Configuration Page, click **Finish**, save your settings, and then reboot the system.
- 10. To enable and verify virtualization:
 - a. Open the <code>grub.conf</code> file and configure the <code>iommu</code> parameter as shown in Figure 12-10. (For details, see "Enabling IOMMU for SR-IOV in UEFI-based Linux OS Installations" on page 229.)
 - For Intel-based systems, add intel_iommu=on.
 - For AMD-based systems, add amd iommu=on.

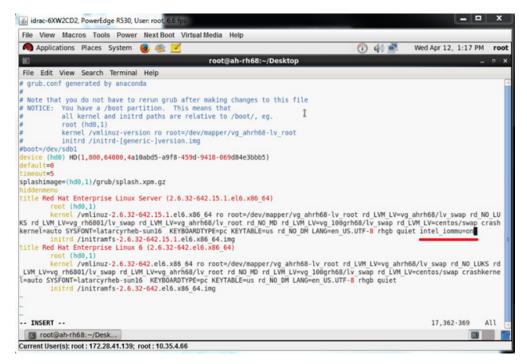


Figure 12-10. Editing the grub.conf File for SR-IOV

- b. Save the grub.conf file and then reboot the system.
- c. To verify that the changes are in effect, issue the following command:

```
dmesg | grep -i iommu
```

A successful input—output memory management unit (IOMMU) command output should show, for example:

```
Intel-IOMMU: enabled
```

d. To view VF details (number of VFs and total VFs), issue the following command:

```
find /sys/|grep -i sriov
```

- 11. For a specific port, enable a quantity of VFs.
 - a. Issue the following command to enable, for example, 8 VFs on PCI instance 04:00.0 (bus 4, device 0, function 0):

```
[root@ah-rh68 ~] # echo 8 >
/sys/devices/pci0000:00/0000:02.0/0000:04:00.0/
sriov_numvfs
```

b. Review the command output (Figure 12-11) to confirm that actual VFs were created on bus 4, device 2 (from the 0000:00:02.0 parameter), functions 0 through 7. Note that the actual device ID is different on the PFs (8070 in this example) versus the VFs (8090 in this example).

```
[root@ah-rh68 Desktop]#
[root@ah-rh68 Desktop]# echo 8 > /sys/devices/pci0000:00/0000:00:02.0/0000:04:00.0/sriov_numvfs
 [root@ah-rh68 Desktop]#
[root@ah-rh68 Desktop]# lspci -vv|grep -i Qlogic
04:00.0 Ethernet controller: QLogic Corp. Device 8070 (rev 02)
        Subsystem: QLogic Corp. Device 000b
                 Product Name: QLogic 25GE 2P QL41262HxCU-DE Adapter
[V4] Vendor specific: NMVQLogic
04:00.1 Ethernet controller: QLogic Corp. Device 8070 (rev 02)
        Subsystem: QLogic Corp. Device 000b
Product Name: QLogic 25GE 2P QL41262HxCU-DE Adapter
[V4] Vendor specific: NMVQLogic
04:02.0 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
         Subsystem: OLogic Corp. Device 000b
04:02.1 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
Subsystem: TQLogic Corp. Device 000b
04:02.2 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
         Subsystem: QLogic Corp. Device 000b
04:02.3 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
         Subsystem: QLogic Corp. Device 000b
04:02.4 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
         Subsystem: QLogic Corp. Device 000b
04:02.5 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
         Subsystem: QLogic Corp. Device 000b
04:02.6 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
         Subsystem: QLogic Corp. Device 000b
04:02.7 Ethernet controller: QLogic Corp. Device 8090 (rev 02)
        Subsystem: QLogic Corp. Device 000b
[root@ah-rh68 Desktop]#

☐ root@ah-rh68:~/Desk...
```

Figure 12-11. Command Output for sriov_numvfs

12. To view a list of all PF and VF interfaces, issue the following command:

```
# ip link show | grep -i vf -b2
```

Figure 12-12 shows example output.

```
[root@localhost ~]# ip link show | grep -i vf -b2
163-2: em1_1: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc mq state UP mode DEFAULT group default qlen 1000
271- link/ether f4:e9:d4:ee:54:c2 brd ff:ff:ff:ff:ff
326: vf 0 MAC 00:00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
439: vf 1 MAC 00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
552: vf 2 MAC 00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
665: vf 3 MAC 00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
778: vf 4 MAC 00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
891: vf 5 MAC 00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
1004: vf 6 MAC 00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
1117: vf 7 MAC 00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto
```

Figure 12-12. Command Output for ip link show Command

- 13. Assign and verify MAC addresses:
 - a. To assign a MAC address to the VF, issue the following command:
 - ip link set <pf device> vf <vf index> mac <mac address>

- b. Ensure that the VF interface is up and running with the assigned MAC address.
- 14. Power off the VM and attach the VF. (Some OSs support hot-plugging of VFs to the VM.)
 - a. In the Virtual Machine dialog box (Figure 12-13), click Add Hardware.

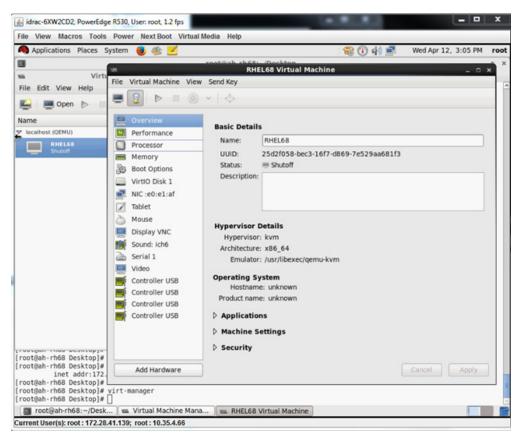


Figure 12-13. RHEL68 Virtual Machine

- b. In the left pane of the Add New Virtual Hardware dialog box (Figure 12-14), click **PCI Host Device**.
- c. In the right pane, select a host device.
- d. Click Finish.

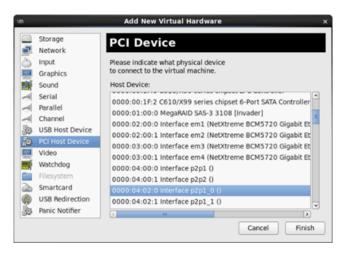


Figure 12-14. Add New Virtual Hardware

15. Power on the VM, and then issue the following command:

```
check lspci -vv|grep -I ether
```

- 16. Install the drivers for the adapters detected in the VM. Use the latest drivers available from your vendor for your host OS (do not use inbox drivers). The same driver version must be installed on the host and the VM.
- 17. As needed, add more VFs in the VM.

Enabling IOMMU for SR-IOV in UEFI-based Linux OS Installations

Follow the appropriate procedure for your Linux OS.

NOTE

For AMD systems, replace intel iommu=on with amd iommu=on.

To enable IOMMU for SR-IOV on RHEL 6.x:

■ In the /boot/efi/EFI/redhat/grub.conf file, locate the kernel line, and then append the intel iommu=on boot parameter.

To enable IOMMU for SR-IOV on RHEL 7.x and later:

- 1. In the /etc/default/grub file, locate GRUB_CMDLINE_LINUX, and then append the intel iommu=on boot parameter.
- 2. To update the grub configuration file, issue the following command:

```
grub2-mkconfig -o /boot/efi/EFI/redhat/grub.cfg
```

To enable IOMMU for SR-IOV on SLES 12.x:

- 1. In the /etc/default/grub file, locate GRUB_CMDLINE_LINUX_DEFAULT, and then append the intel iommu=on boot parameter.
- 2. To update the grub configuration file, issue the following command:

```
grub2-mkconfig -o /boot/grub2/grub.cfg
```

To enable IOMMU for SR-IOV on SLES 15.x and later:

- 1. In the /etc/default/grub file, locate GRUB_CMDLINE_LINUX_DEFAULT, and then append the intel iommu=on boot parameter.
- 2. To update the grub configuration file, issue the following command: grub2-mkconfig -o /boot/efi/EFI/sles/grub.cfg

Configuring SR-IOV on VMware

To configure SR-IOV on VMware:

- Access the server BIOS System Setup, and then click System BIOS Settings.
- 2. On the System BIOS Settings page, click Integrated Devices.
- 3. On the Integrated Devices page (see Figure 12-1 on page 218):
 - a. Set the SR-IOV Global Enable option to Enabled.
 - b. Click Back.
- 4. In the System Setup window, click **Device Settings**.
- 5. On the Device Settings page, select a port for the 25G 41xxx Series Adapter.
- 6. On the Device Level Configuration page (see Figure 12-2 on page 218):
 - a. Set the Virtualization Mode to SR-IOV.
 - b. Click Back.
- 7. On the Main Configuration Page, click **Finish**.
- 8. Save the configuration settings and reboot the system.

9. To enable the needed quantity of VFs per port (in this example, 16 on each port of a dual-port adapter), issue the following command:

```
"esxcfg-module -s "max_vfs=16,16" qedentv"
```

NOTE

Each Ethernet function of the 41xxx Series Adapter must have its own entry.

- 10. Reboot the host.
- 11. To verify that the changes are complete at the module level, issue the following command:

```
"esxcfg-module -g qedentv"
[root@localhost:~] esxcfg-module -g qedentv
qedentv enabled = 1 options = 'max vfs=16,16'
```

12. To verify if actual VFs were created, issue the <code>lspci</code> command as follows:

```
[root@localhost:~] lspci | grep -i QLogic | grep -i 'ethernet\|network' | more
0000:05:00.0 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx 10/25
GbE Ethernet Adapter [vmnic6]
0000:05:00.1 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx 10/25
GbE Ethernet Adapter [vmnic7]
0000:05:02.0 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.0 VF 0]
0000:05:02.1 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.0 VF 1]
0000:05:02.2 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.0 VF 2]
0000:05:02.3 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.0 VF 3]
0000:05:03.7 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.0 VF 15]
0000:05:0e.0 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.1 VF 0]
0000:05:0e.1 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.1 VF 1]
0000:05:0e.2 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF 0.5.1 VF 2]
```

231

```
0000:05:0e.3 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF_0.5.1_VF_3]
.
.
.
0000:05:0f.6 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF_0.5.1_VF_14]
0000:05:0f.7 Network controller: QLogic Corp. QLogic FastLinQ QL41xxx Series
10/25 GbE Controller (SR-IOV VF) [PF_0.5.1_VF_15]
```

- 13. Attach VFs to the VM as follows:
 - a. Power off the VM and attach the VF. (Some OSs support hot-plugging of VFs to the VM.)
 - b. Add a host to a VMware vCenter Server Virtual Appliance (vCSA).
 - c. Click **Edit Settings** of the VM.
- 14. Complete the Edit Settings dialog box (Figure 12-15) as follows:
 - a. In the **New Device** box, select **Network**, and then click **Add**.
 - b. For Adapter Type, select SR-IOV Passthrough.
 - c. For **Physical Function**, select the Marvell VF.
 - d. To save your configuration changes and close this dialog box, click **OK**.

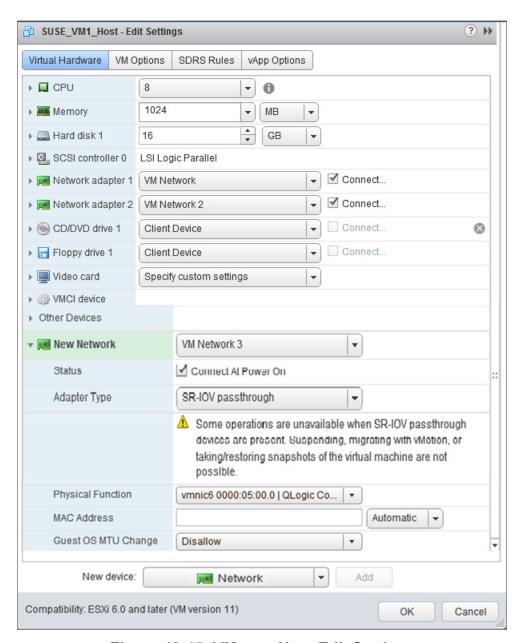


Figure 12-15. VMware Host Edit Settings

15. To validate the VFs per port, issue the <code>esxcli</code> command as follows:

```
2
    false 005:02.2
 3
    false 005:02.3
 4
    false 005:02.4
 5
    false 005:02.5
    false 005:02.6
 6
7
    false 005:02.7
    false 005:03.0
9
    false 005:03.1
10
    false 005:03.2
    false 005:03.3
11
12
    false 005:03.4
   false 005:03.5
13
    false 005:03.6
14
15
    false 005:03.7
```

- 16. Install the Marvell drivers for the adapters detected in the VM. Use the latest drivers available from your vendor for your host OS (do not use inbox drivers). The same driver version must be installed on the host and the VM.
- 17. Power on the VM, and then issue the ifconfig -a command to verify that the added network interface is listed.
- 18. As needed, add more VFs in the VM.

13 NVMe-oF Configuration with RDMA

Non-Volatile Memory Express over Fabrics (NVMe-oF) enables the use of alternate transports to PCIe to extend the distance over which an NVMe host device and an NVMe storage drive or subsystem can connect. NVMe-oF defines a common architecture that supports a range of storage networking fabrics for the NVMe block storage protocol over a storage networking fabric. This architecture includes enabling a front-side interface into storage systems, scaling out to large quantities of NVMe devices, and extending the distance within a data center over which NVMe devices and NVMe subsystems can be accessed.

The NVMe-oF configuration procedures and options described in this chapter apply to Ethernet-based RDMA protocols, including RoCE and iWARP. The development of NVMe-oF with RDMA is defined by a technical sub-group of the NVMe organization.

This chapter demonstrates how to configure NVMe-oF for a simple network. The example network comprises the following:

- Two servers: an initiator and a target. The target server is equipped with a PCIe SSD drive.
- Operating system: RHEL 7.6 and later, , RHEL 8.x and later, SLES 15.x and later
- Two adapters: One 41xxx Series Adapter installed in each server. Each port can be independently configured to use RoCE, RoCEv2, or iWARP as the RDMA protocol over which NVMe-oF runs.
- For RoCE and RoCEv2, an optional switch configured for data center bridging (DCB), relevant quality of service (QoS) policy, and vLANs to carry the NVMe-oF's RoCE/RoCEv2 DCB traffic class priority. The switch is not needed when NVMe-oF is using iWARP.

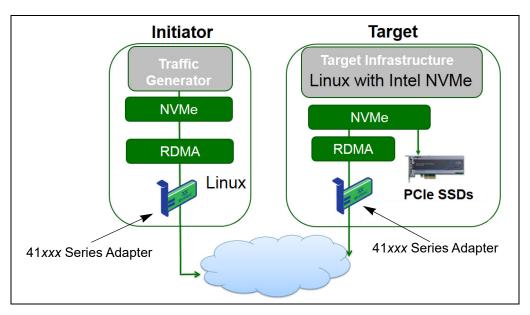


Figure 13-1 illustrates an example network.

Figure 13-1. NVMe-oF Network

The NVMe-oF configuration process covers the following procedures:

- Installing Device Drivers on Both Servers
- Configuring the Target Server
- Configuring the Initiator Server
- Preconditioning the Target Server
- Testing the NVMe-oF Devices
- Optimizing Performance

Installing Device Drivers on Both Servers

After installing your operating system (SLES 12 SP3), install device drivers on both servers. To upgrade the kernel to the latest Linux upstream kernel, go to:

https://www.kernel.org/pub/linux/kernel/v4.x/

- 1. Install and load the latest FastLinQ drivers (qed, qede, libqedr/qedr) following all installation instructions in the README.
- 2. (Optional) If you upgraded the OS kernel, you must reinstall and load the latest driver as follows:
 - a. Install the latest FastLinQ firmware following all installation instructions in the README.
 - b. Install the OS RDMA support applications and libraries by issuing the following commands:

```
# yum groupinstall "Infiniband Support"
# yum install tcl-devel libibverbs-devel libnl-devel
glib2-devel libudev-devel lsscsi perftest
```

- # yum install gcc make git ctags ncurses ncurses-devel
 openssl* openssl-devel elfutils-libelf-devel*
- c. To ensure that NVMe OFED support is in the selected OS kernel, issue the following command:

```
make menuconfig
```

d. Under **Device Drivers**, ensure that the following are enabled (set to **M**):

```
NVM Express block devices

NVM Express over Fabrics RDMA host driver

NVMe Target support

NVMe over Fabrics RDMA target support
```

e. (Optional) If the **Device Drivers** options are not already present, rebuild the kernel by issuing the following commands:

```
# make
# make modules
# make modules_install
# make install
```

f. If changes were made to the kernel, reboot to that new OS kernel. For instructions on how to set the default boot kernel, go to:

https://wiki.centos.org/HowTos/Grub2

- 3. Enable and start the RDMA service as follows:
 - # systemctl enable rdma.service
 - # systemctl start rdma.service

Disregard the RDMA Service Failed error. All OFED modules required by qedr are already loaded.

Configuring the Target Server

Configure the target server after the reboot process. After the server is operating, you cannot change the configuration without rebooting. If you are using a startup script to configure the target server, consider pausing the script (using the wait command or something similar) as needed to ensure that each command finishes before executing the next command.

To configure the target service:

- 1. Load target modules. Issue the following commands after each server reboot:
 - # modprobe qedr
 - # modprobe nvmet; modprobe nvmet-rdma
 - # 1smod | grep nvme (confirm that the modules are loaded)
- 2. Create the target subsystem NVMe Qualified Name (NQN) with the name indicated by <nvme-subsystem-name>. Use the NVMe-oF specifications; for example, nqn.<YEAR>-<Month>.org.<your-company>.
- # mkdir /sys/kernel/config/nvmet/subsystems/<nvme-subsystem-name>
- # cd /sys/kernel/config/nvmet/subsystems/<nvme-subsystem-name>
- 3. Create multiple unique NQNs for additional NVMe devices as needed.
- 4. Set the target parameters, as listed in Table 13-1.

Table 13-1. Target Parameters

Command	Description
# echo 1 > attr_allow_any_host	Allows any host to connect.
# mkdir namespaces/1	Creates a namespace.

Table 13-1. Target Parameters (Continued)

Command	Description		
<pre># echo -n /dev/nvme0n1 >namespaces/ 1/device_path</pre>	Sets the NVMe device path. The NVMe device path can differ between systems. Check the device path using the <code>lsblk</code> command. This system has two NVMe devices: <code>nvme0n1</code> and <code>nvme1n1</code> .		
	[root@localhost home]# lsblk NAME		
# echo 1 > namespaces/1/enable	Enables the namespace.		
<pre># mkdir /sys/kernel/config/nvmet/ ports/1</pre>	Creates NVMe port 1.		
<pre># cd /sys/kernel/config/nvmet/ports/1</pre>			
# echo 1.1.1.1 > addr_traddr	Sets the same IP address. For example, 1.1.1.1 is the IP address for the target port of the 41xxx Series Adapter.		
# echo rdma > addr_trtype	Sets the transport type RDMA.		
# echo 4420 > addr_trsvcid	Sets the RDMA port number. The socket port number for NVMe-oF is typically 4420. However, any port number can be used if it is used consistently throughout the configuration.		
# echo ipv4 > addr_adrfam	Sets the IP address type.		

- 5. Create a symbolic link (symlink) to the newly created NQN subsystem:
 - # ln -s /sys/kernel/config/nvmet/subsystems/
 nvme-subsystem-name subsystems/nvme-subsystem-name
- 6. Confirm that the NVMe target is listening on the port as follows:
 - # dmesg | grep nvmet_rdma
 [8769.470043] nvmet rdma: enabling port 1 (1.1.1:4420)

Configuring the Initiator Server

You must configure the initiator server after the reboot process. After the server is operating, you cannot change the configuration without rebooting. If you are using a startup script to configure the initiator server, consider pausing the script (using the wait command or something similar) as needed to ensure that each command finishes before executing the next command.

To configure the initiator server:

- 1. Load the NVMe modules. Issue these commands after each server reboot:
 - # modprobe qedr
 # modprobe nvme-rdma
- 2. Download, compile and install the nvme-cli initiator utility. Issue these commands at the first configuration—you do not need to issue these commands after each reboot.

```
# git clone https://github.com/linux-nvme/nvme-cli.git
# cd nvme-cli
# make && make install
```

3. Verify the installation version as follows:

```
# nvme version
```

4. Discover the NVMe-oF target as follows:

```
# nvme discover -t rdma -a 1.1.1.1 -s 1023
```

Make note of the subsystem NQN (subnqn) of the discovered target (Figure 13-2) for use in Step 5.

```
[root@localhost home] # nvme discover -t rdma -a 1.1.1.1 -s 1023

Discovery Log Number of Records 1, Generation counter 1

====Discovery Log Entry 0=====

trtype: rdma
adrfam: ipv4
subtype: nvme subsystem
treq: not specified
portid: 1
trsvcid: 1023

subnqn: nvme-qlogic-tgt1
traddr: 1.1.1.1

rdma_prtype: not specified
rdma_qptype: connected
rdma_qptype: connected
rdma_cms: rdma-cm
rdma_pkey: 0x0000
```

Figure 13-2. Subsystem NQN

- 5. Connect to the discovered NVMe-oF target (nvme-qlogic-tgt1) using the NQN. Issue the following command after each server reboot. For example:
 - # nvme connect -t rdma -n nvme-qlogic-tgt1 -a 1.1.1.1 -s 1023
- 6. Confirm the NVMe-oF target connection with the NVMe-oF device as follows:

```
# dmesg | grep nvme
# lsblk
# list nvme
```

Figure 13-3 shows an example.

```
[root@localhost home] #dmesg | grep nvme
[ 233.645554] nvme nvmeO: new ctrl: NQN "nvme-qlogic-tgt1", addr 1.1.1.1:1023
[root@localhost home] # lsblk
NAME MAJ:MIN RM SIZ
                             SIZE
1.1T
                                     RO TYPE
                                               MOUNTPOINT
sdb
             8:0
                       0
                                       disk
 -sdb2
                                        part
_sdb3
                                        part
             8:3
                                8G
                                     0
                                                [SWAP]
∟sdb1
                                               /boot/efi
                       0
                                1G
                                     0
                                        part
disk
             8:1
           259:0
nvme0n1
                           372.6G
[root@localhost home] # nvme list
Node
                 SN
                                       Model
                                               Namespace
                                                                                         Format
                                                                                                            FW Rev
/dev/nvme0n1
                 7a591f3ec788a367
                                       Linux 1
                                                              1.60 TB / 1.60 TB 512
                                                                                                B + 0 B 4.13.8
```

Figure 13-3. Confirm NVMe-oF Connection

Preconditioning the Target Server

NVMe target servers that are tested out-of-the-box show a higher-than-expected performance. Before running a benchmark, the target server needs to be *prefilled* or *preconditioned*.

To precondition the target server:

Data-Center-Tool

- Secure-erase the target server with vendor-specific tools (similar to formatting). This test example uses an Intel NVMe SSD device, which requires the Intel Data Center Tool that is available at the following link: https://downloadcenter.intel.com/download/23931/Intel-Solid-State-Drive-
- 2. Precondition the target server (nvme0n1) with data, which guarantees that all available memory is filled. This example uses the "DD" disk utility:

```
# dd if=/dev/zero bs=1024k of=/dev/nvme0n1
```

Testing the NVMe-oF Devices

Compare the latency of the local NVMe device on the target server with that of the NVMe-oF device on the initiator server to show the latency that NVMe adds to the system.

To test the NVMe-oF device:

- 1. Update the Repository (Repo) source and install the Flexible Input/Output (FIO) benchmark utility on both the target and initiator servers by issuing the following commands:
 - # yum install epel-release
 - # yum install fio

Figure 13-4. FIO Utility Installation

2. Run the FIO utility to measure the latency of the initiator NVMe-oF device. Issue the following command:

```
# fio --filename=/dev/nvme0n1 --direct=1 --time_based
--rw=randread --refill_buffers --norandommap --randrepeat=0
--ioengine=libaio --bs=4k --iodepth=1 --numjobs=1
--runtime=60 --group_reporting --name=temp.out
```

FIO reports two latency types: submission and completion. Submission latency (slat) measures application-to-kernel latency. Completion latency (clat), measures end-to-end kernel latency. The industry-accepted method is to read *clat percentiles* in the 99.00th range.

In this example, the initiator device NVMe-oF latency is 30µsec.

3. Run FIO to measure the latency of the local NVMe device on the target server. Issue the following command:

```
# fio --filename=/dev/nvme0n1 --direct=1 --time_based
--rw=randread --refill_buffers --norandommap --randrepeat=0
--ioengine=libaio --bs=4k --iodepth=1 --numjobs=1
--runtime=60 --group_reporting --name=temp.out
```

In this example, the target NVMe device latency is 8µsec. The total latency that results from the use of NVMe-oF is the difference between the initiator device NVMe-oF latency (30µsec) and the target device NVMe-oF latency (8µsec), or 22µsec.

4. Run FIO to measure bandwidth of the local NVMe device on the target server. Issue the following command:

```
fio --verify=crc32 --do_verify=1 --bs=8k --numjobs=1
--iodepth=32 --loops=1 --ioengine=libaio --direct=1
--invalidate=1 --fsync_on_close=1 --randrepeat=1
--norandommap --time_based --runtime=60
--filename=/dev/nvme0n1 --name=Write-BW-to-NVMe-Device
--rw=randwrite
```

Where --rw can be randread for reads only, randwrite for writes only, or randrw for reads and writes.

Optimizing Performance

To optimize performance on both initiator and target servers:

- 1. Configure the following system BIOS settings:
 - ☐ Power Profiles = 'Max Performance' or equivalent
 - ☐ ALL C-States = Disabled
 - ☐ Hyperthreading = Disabled
- 2. Configure the Linux kernel parameters by editing the grub file (/etc/default/grub).
 - a. Add parameters to end of line GRUB CMDLINE LINUX:

```
GRUB_CMDLINE_LINUX="nosoftlockup intel_idle.max_cstate=0
processor.max cstate=1 mce=ignore ce idle=poll"
```

- b. Save the grub file.
- c. Rebuild the grub file.
 - To rebuild the grub file for a legacy BIOS boot, issue the following command:
 - # grub2-mkconfig -o /boot/grub2/grub.cfg (Legacy BIOS boot)
 - To rebuild the grub file for an EFI boot, issue the following command:
 - # grub2-mkconfig -o /boot/efi/EFI/<os>/grub.cfg (EFI boot)
- d. Reboot the server to implement the changes.

3. Set the IRQ affinity for all 41xxx Series Adapters. The multi_rss-affin.sh file is a script file that is listed in "IRQ Affinity (multi_rss-affin.sh)" on page 244.

```
# systemctl stop irqbalance
# ./multi rss-affin.sh eth1
```

NOTE

A different version of this script, $qedr_affin.sh$, is in the 41xxx Linux Source Code Package in the $\add-ons\performance\roce$ directory. For an explanation of the IRQ affinity settings, refer to the multiple irqs.txt file in that directory.

4. Set the CPU frequency. The <code>cpufreq.sh</code> file is a script that is listed in "CPU Frequency (cpufreq.sh)" on page 245.

```
# ./cpufreq.sh
```

The following sections list the scripts that are used in Steps 3 and 4.

IRQ Affinity (multi_rss-affin.sh)

The following script sets the IRQ affinity.

```
#!/bin/bash
#RSS affinity setup script
#input: the device name (ethX)
#OFFSET=0 0/1 0/1/2 0/1/2/3
#FACTOR=1
           2
                  3
                           4
OFFSET=0
FACTOR=1
LASTCPU='cat /proc/cpuinfo | grep processor | tail -n1 | cut -d":" -f2'
MAXCPUID='echo 2 $LASTCPU ^ p | dc'
OFFSET='echo 2 $OFFSET ^ p | dc'
FACTOR='echo 2 $FACTOR ^ p | dc'
CPUID=1
for eth in $*; do
NUM='grep $eth /proc/interrupts | wc -l'
NUM FP=$((${NUM}))
INT='grep -m 1 $eth /proc/interrupts | cut -d ":" -f 1'
echo "$eth: ${NUM} (${NUM FP} fast path) starting irq ${INT}"
```

```
CPUID=$((CPUID*OFFSET))
for ((A=1; A<=${NUM_FP}; A=${A}+1)) ; do
INT='grep -m $A $eth /proc/interrupts | tail -1 | cut -d ":" -f 1'
SMP='echo $CPUID 16 o p | dc'
echo ${INT} smp affinity set to ${SMP}
echo $((${SMP})) > /proc/irq/$((${INT}))/smp_affinity
CPUID=$((CPUID*FACTOR))
if [ ${CPUID} -gt ${MAXCPUID} ]; then
CPUID=1
CPUID=$((CPUID*OFFSET))
fi
done
done
```

CPU Frequency (cpufreq.sh)

The following script sets the CPU frequency.

```
#Usage "./nameofscript.sh"
grep -E '^model name|^cpu MHz' /proc/cpuinfo
cat /sys/devices/system/cpu/cpu0/cpufreq/scaling_governor
for CPUFREQ in /sys/devices/system/cpu/cpu*/cpufreq/scaling_governor; do [ -f
$CPUFREQ ] || continue; echo -n performance > $CPUFREQ; done
cat /sys/devices/system/cpu/cpu0/cpufreq/scaling governor
```

To configure the network or memory settings:

```
sysctl -w net.ipv4.tcp_mem="16777216 16777216"
sysctl -w net.ipv4.tcp_wmem="4096 65536 16777216"
sysctl -w net.ipv4.tcp_rmem="4096 87380 16777216"
sysctl -w net.core.wmem_max=16777216
sysctl -w net.core.rmem_max=16777216
sysctl -w net.core.wmem_default=16777216
sysctl -w net.core.rmem_default=16777216
sysctl -w net.core.optmem_max=16777216
sysctl -w net.ipv4.tcp_low_latency=1
sysctl -w net.ipv4.tcp_timestamps=0
sysctl -w net.ipv4.tcp_sack=1
sysctl -w net.ipv4.tcp_window_scaling=0
sysctl -w net.ipv4.tcp_window_scaling=0
sysctl -w net.ipv4.tcp_adv win scale=1
```

NOTE

The following commands apply only to the initiator server.

- # echo 0 > /sys/block/nvme0n1/queue/add_random
- # echo 2 > /sys/block/nvme0n1/queue/nomerges

14 VXLAN Configuration

This chapter provides instructions for:

- Configuring VXLAN in Linux
- "Configuring VXLAN in VMware" on page 249
- "Configuring VXLAN in Windows Server 2016" on page 250

Configuring VXLAN in Linux

To configure VXLAN in Linux:

- 1. Download, extract, and configure the openvswitch (OVS) tar ball.
 - a. Download the appropriate openvswitch release from the following location:
 - http://www.openvswitch.org/download/
 - b. Extract the tar ball by navigating to the directory where you downloaded the openvswitch release, and then issue the following command:
 - ./configure;make;make install (compilation)
 - c. Configure openvswitch by issuing the following commands:

```
modprobe -v openvswitch
export PATH=$PATH:/usr/local/share/openvswitch/scripts
ovs-ctl start
ovs-ctl status
```

When running ovs-ctl status, the ovsdb-server and ovs-vswitchd should be running with pid. For example:

```
[root@localhost openvswitch-2.11.1]# ovs-ctl status
ovsdb-server is running with pid 8479
ovs-vswitchd is running with pid 8496
```

2. Create the bridge.

a. To configure Host 1, issue the following commands:

```
ovs-vsctl add-br br0
ovs-vsctl add-br br1
ovs-vsctl add-port br0 eth0
ifconfig eth0 0 && ifconfig br0 192.168.1.10 netmask 255.255.255.0
route add default gw 192.168.1.1 br0
ifconfig br1 10.1.2.10 netmask 255.255.255.0
ovs-vsctl add-port br1 vx1 -- set interface vx1 type=vxlan
options:remote_ip=192.168.1.11 (peer IP address)
```

b. To configure Host 2, issue the following commands:

```
ovs-vsctl add-br br0
ovs-vsctl add-br br1
ovs-vsctl add-port br0 eth0
ifconfig eth0 0 && ifconfig br0 192.168.1.11 netmask 255.255.255.0
route add default gw 192.168.1.1 br0
ifconfig br1 10.1.2.11 netmask 255.255.255.0
ovs-vsctl add-port br1 vx1 -- set interface vx1 type=vxlan options:
remote ip=192.168.1.10
```

3. Verify the configuration.

Run traffic between the host and peer using iperf. Ensure that the firewall and iptables stop and clean, respectively.

- 4. Configure the bridge as a passthrough to the VMs, and then check connectivity from the VM to the Peer.
 - a. Create a VM through virt-manager.
 - b. As there is no option to attach bridge <code>br1</code> through virt-manager, change the xml file as follows

Issue the following command:

```
command: virsh edit vm1
```

Add the following code:

```
<interface type='bridge'>
<source bridge='br1'/>
<virtualport type='openvswitch'>
<parameters/>
</virtualport>
<model type='virtio'/>
</interface>
```

c. Power up the VM and check br1 interfaces.

Ensure that br1 is in the OS. The br1 interface is named eth0, ens7; manually configure the static IP through th network device file and assign the same subnet IP to the peer (Host 2 VM).

Run traffic from the peer to the VM.

NOTE

You can use this procedure to test other tunnels, such as Generic Network Virtualization Encapsulation (GENEVE) and generic routing encapsulation (GRE), with OVS.

If you do not want to use OVS, you can continue with the legacy bridge option brctl.

Configuring VXLAN in VMware

To configure VXLAN in VMware, follow the instructions in the following locations:

https://docs.vmware.com/en/VMware-NSX-Data-Center-for-vSphere/6.3/com.vmware.nsx.cross-vcenter-install.doc/GUID-49BAECC2-B800-4670-AD8C-A5292ED6BC19.html

https://www.vmware.com/content/dam/digitalmarketing/vmware/en/pdf/products/nsx/vmw-nsx-network-virtualization-design-guide.pdf

https://pubs.vmware.com/nsx-63/topic/com.vmware.nsx.troubleshooting.doc/GUI D-EA1DB524-DD2E-4157-956E-F36BDD20CDB2.html

https://communities.vmware.com/api/core/v3/attachments/124957/data

Configuring VXLAN in Windows Server 2016

VXLAN configuration in Windows Server 2016 includes:

- Enabling VXLAN Offload on the Adapter
- Deploying a Software Defined Network

Enabling VXLAN Offload on the Adapter

To enable VXLAN offload on the adapter:

- 1. Open the miniport properties, and then click the **Advanced** tab.
- On the adapter properties' Advanced page (Figure 14-1) under Property, select VXLAN Encapsulated Task Offload.

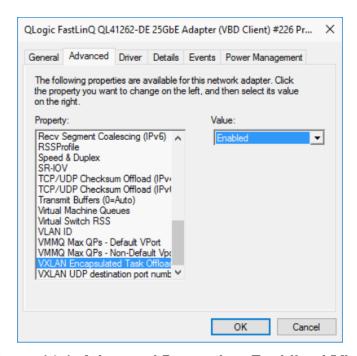


Figure 14-1. Advanced Properties: Enabling VXLAN

- 3. Set the Value to Enabled.
- 4. Click **OK**.

Deploying a Software Defined Network

To take advantage of VXLAN encapsulation task offload on virtual machines, you must deploy a Software Defined Networking (SDN) stack that utilizes a Microsoft Network Controller.

For more details, refer to the following Microsoft TechNet link on Software Defined Networking:

https://technet.microsoft.com/en-us/windows-server-docs/networking/sdn/software-defined-networking--sdn-

15 Windows Server 2016

This chapter provides the following information for Windows Server 2016:

- Configuring RoCE Interfaces with Hyper-V
- "RoCE over Switch Embedded Teaming" on page 258
- "Configuring QoS for RoCE" on page 259
- "Configuring VMMQ" on page 268
- "Configuring Storage Spaces Direct" on page 272

Configuring RoCE Interfaces with Hyper-V

In Windows Server 2016, Hyper-V with Network Direct Kernel Provider Interface (NDKPI) Mode-2, host virtual network adapters (host virtual NICs) support RDMA.

NOTE

DCBX is required for RoCE over Hyper-V. To configure DCBX, either:

- Configure through the HII (see "Preparing the Adapter" on page 129).
- Configure using QoS (see "Configuring QoS for RoCE" on page 259).

RoCE configuration procedures in this section include:

- Creating a Hyper-V Virtual Switch with an RDMA NIC
- Adding a vLAN ID to Host Virtual NIC
- Verifying If RoCE is Enabled
- Adding Host Virtual NICs (Virtual Ports)
- Mapping the SMB Drive and Running RoCE Traffic

Creating a Hyper-V Virtual Switch with an RDMA NIC

Follow the procedures in this section to create a Hyper-V virtual switch and then enable RDMA in the host VNIC.

To create a Hyper-V virtual switch with an RDMA virtual NIC:

- 1. On all physical interfaces, set the value of the **NetworkDirect Functionality** parameter to **Enabled**.
- 2. Launch Hyper-V Manager.
- 3. Click Virtual Switch Manager (see Figure 15-1).

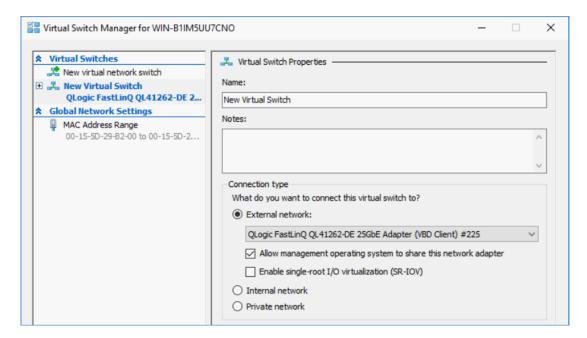


Figure 15-1. Enabling RDMA in Host Virtual NIC

- 4. Create a virtual switch.
- 5. Select the Allow management operating system to share this network adapter check box.

In Windows Server 2016, a new parameter—Network Direct (RDMA)—is added in the Host virtual NIC.

To enable RDMA in a host virtual NIC:

- 1. Open the Hyper-V Virtual Ethernet Adapter Properties window.
- 2. Click the Advanced tab.

- 3. On the Advanced page (Figure 15-2):
 - Under Property, select Network Direct (RDMA).
 - b. Under Value, select Enabled.
 - c. Click OK.

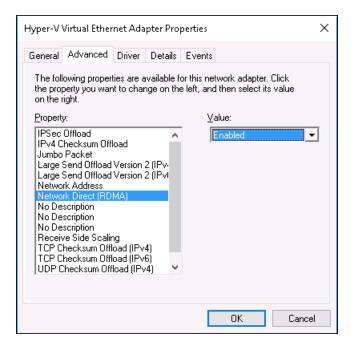


Figure 15-2. Hyper-V Virtual Ethernet Adapter Properties

4. To enable RDMA through PowerShell, issue the following Windows PowerShell command:

```
PS C:\Users\Administrator> Enable-NetAdapterRdma "vEthernet (New Virtual Switch)"
PS C:\Users\Administrator>
```

Adding a vLAN ID to Host Virtual NIC

To add a vLAN ID to a host virtual NIC:

1. To find the host virtual NIC name, issue the following Windows PowerShell command:

PS C:\Users\Administrator> Get-VMNetworkAdapter -ManagementOS

Figure 15-3 shows the command output.



Figure 15-3. Windows PowerShell Command: Get-VMNetworkAdapter

2. To set the vLAN ID to the host virtual NIC, issue the following Windows PowerShell command:

```
PS C:\Users\Administrator> Set-VMNetworkAdaptervlan
-VMNetworkAdapterName "New Virtual Switch" -VlanId 5 -Access
-ManagementOS
```

NOTE

Note the following about adding a vLAN ID to a host virtual NIC:

- A vLAN ID must be assigned to a host virtual NIC. The same vLAN ID must be assigned to ports on the switch.
- Make sure that the vLAN ID is not assigned to the physical interface when using a host virtual NIC for RoCE.
- If you are creating more than one host virtual NIC, you can assign a different vLAN to each host virtual NIC.

Verifying If RoCE is Enabled

To verify if the RoCE is enabled:

■ Issue the following Windows PowerShell command:

Get-NetAdapterRdma

Command output lists the RDMA supported adapters as shown in Figure 15-4.

```
PS C:\Users\Administrator> Get-NetAdapterRdma

Name InterfaceDescription Enabled
----
vEthernet (New Virtual... Hyper-V Virtual Ethernet Adapter True
```

Figure 15-4. Windows PowerShell Command: Get-NetAdapterRdma

Adding Host Virtual NICs (Virtual Ports)

To add host virtual NICs:

- 1. To add a host virtual NIC, issue the following command:
 - Add-VMNetworkAdapter -SwitchName "New Virtual Switch" -Name SMB ManagementOS
- 2. Enable RDMA on host virtual NICs as shown in "To enable RDMA in a host virtual NIC:" on page 253.
- 3. To assign a vLAN ID to the virtual port, issue the following command:

Set-VMNetworkAdapterVlan -VMNetworkAdapterName SMB -VlanId 5 -Access -ManagementOS

Mapping the SMB Drive and Running RoCE Traffic

To map the SMB drive and run the RoCE traffic:

- 1. Launch the Performance Monitor (Perfmon).
- 2. Complete the Add Counters dialog box (Figure 15-5) as follows:
 - a. Under Available counters, select RDMA Activity.
 - b. Under **Instances of selected object**, select the adapter.
 - c. Click Add.

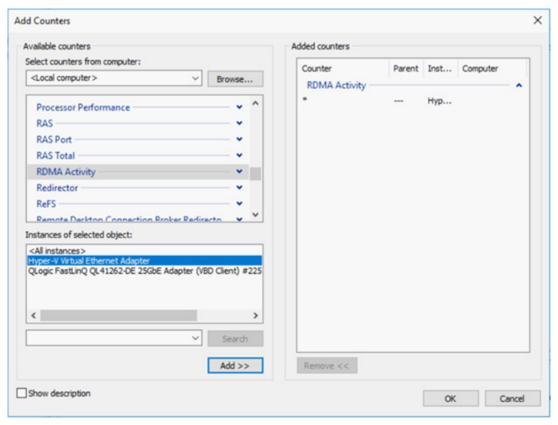


Figure 15-5. Add Counters Dialog Box

If the RoCE traffic is running, counters appear as shown in Figure 15-6.

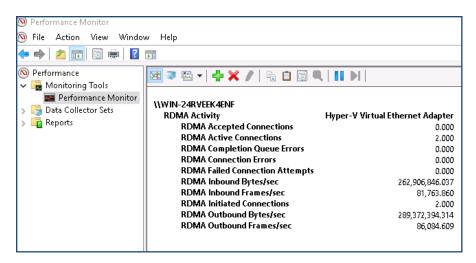


Figure 15-6. Performance Monitor Shows RoCE Traffic

RoCE over Switch Embedded Teaming

Switch Embedded Teaming (SET) is Microsoft's alternative NIC teaming solution available to use in environments that include Hyper-V and the Software Defined Networking (SDN) stack in Windows Server 2016 Technical Preview. SET integrates limited NIC Teaming functionality into the Hyper-V Virtual Switch.

Use SET to group between one and eight physical Ethernet network adapters into one or more software-based virtual network adapters. These adapters provide fast performance and fault tolerance if a network adapter failure occurs. To be placed on a team, SET member network adapters must all be installed in the same physical Hyper-V host.

RoCE over SET procedures included in this section:

- Creating a Hyper-V Virtual Switch with SET and RDMA Virtual NICs
- Enabling RDMA on SET
- Assigning a vLAN ID on SET
- Running RDMA Traffic on SET

Creating a Hyper-V Virtual Switch with SET and RDMA Virtual NICs

To create a Hyper-V virtual switch with SET and RDMA virtual NICs:

■ To create SET, issue the following Windows PowerShell command:

```
PS C:\Users\Administrator> New-VMSwitch -Name SET -NetAdapterName "Ethernet 2","Ethernet 3" -EnableEmbeddedTeaming $true
```

Figure 15-7 shows command output.

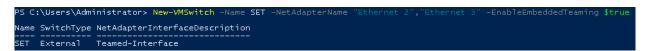


Figure 15-7. Windows PowerShell Command: New-VMSwitch

Enabling RDMA on SET

To enable RDMA on SET:

1. To view SET on the adapter, issue the following Windows PowerShell command:

```
PS C:\Users\Administrator> Get-NetAdapter "vEthernet (SET)"
```

Figure 15-8 shows command output.



Figure 15-8. Windows PowerShell Command: Get-NetAdapter

2. To enable RDMA on SET, issue the following Windows PowerShell command:

PS C:\Users\Administrator> Enable-NetAdapterRdma "vEthernet (SET)"

Assigning a vLAN ID on SET

To assign a vLAN ID on SET:

To assign a vLAN ID on SET, issue the following Windows PowerShell command:

```
PS C:\Users\Administrator> Set-VMNetworkAdapterVlan
-VMNetworkAdapterName "SET" -VlanId 5 -Access -ManagementOS
```

NOTE

Note the following when adding a vLAN ID to a host virtual NIC:

- Make sure that the vLAN ID is not assigned to the physical Interface when using a host virtual NIC for RoCE.
- If you are creating more than one host virtual NIC, a different vLAN can be assigned to each host virtual NIC.

Running RDMA Traffic on SET

For information about running RDMA traffic on SET, go to:

https://technet.microsoft.com/en-us/library/mt403349.aspx

Configuring QoS for RoCE

The two methods of configuring quality of service (QoS) include:

- Configuring QoS by Disabling DCBX on the Adapter
- Configuring QoS by Enabling DCBX on the Adapter

Configuring QoS by Disabling DCBX on the Adapter

All configuration must be completed on all of the systems in use before configuring QoS by disabling DCBX on the adapter. The priority-based flow control (PFC), enhanced transition services (ETS), and traffic classes configuration must be the same on the switch and server.

To configure QoS by disabling DCBX:

- 1. Disable DCBX on the adapter.
- 2. Using HII, set the RoCE Priority to 0.
- 3. To install the DCB role in the host, issue the following Windows PowerShell command:

```
PS C:\Users\Administrators> Install-WindowsFeature Data-Center-Bridging
```

4. To set the **DCBX Willing** mode to **False**, issue the following Windows PowerShell command:

```
PS C:\Users\Administrators> set-NetQosDcbxSetting -Willing 0
```

- 5. Enable QoS in the miniport as follows:
 - a. Open the miniport Properties, and then click the **Advanced** tab.
 - On the adapter properties' Advanced page (Figure 15-9) under Property, select Quality of Service, and then set the value to Enabled.
 - c. Click OK.

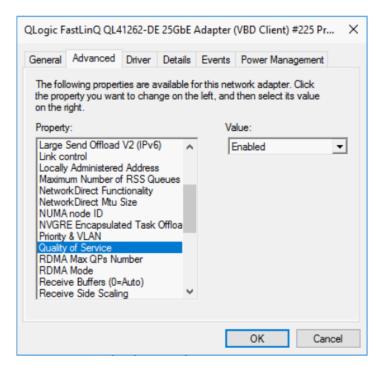


Figure 15-9. Advanced Properties: Enable QoS

- 6. Assign the vLAN ID to the interface as follows:
 - a. Open the miniport properties, and then click the **Advanced** tab.
 - On the adapter properties' Advanced page (Figure 15-10) under Property, select VLAN ID, and then set the value.
 - c. Click OK.

NOTE

The preceding step is required for priority flow control (PFC).

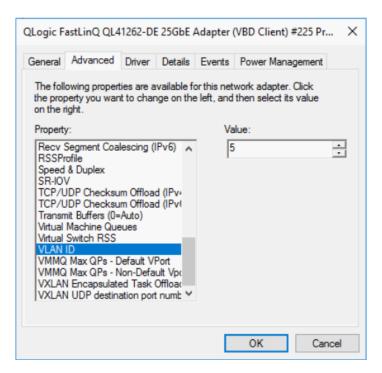


Figure 15-10. Advanced Properties: Setting VLAN ID

7. To enable PFC for RoCE on a specific priority, issue the following command:

PS C:\Users\Administrators> Enable-NetQoSFlowControl
-Priority 5

NOTE

If configuring RoCE over Hyper-V, do not assign a vLAN ID to the physical interface.

8. To disable priority flow control on any other priority, issue the following commands:

PS C:\Users\Administrator> Disable-NetQosFlowControl 0,1,2,3,4,6,7
PS C:\Users\Administrator> Get-NetQosFlowControl

Enabled	PolicySet	IfIndex	IfAlias
False	Global		
	False False False False	False Global False Global False Global False Global	False Global False Global False Global False Global

5 True Global 6 False Global 7 False Global

9. To configure QoS and assign relevant priority to each type of traffic, issue the following commands (where Priority 5 is tagged for RoCE and Priority 0 is tagged for TCP):

PS C:\Users\Administrators> New-NetQosPolicy "SMB"

-NetDirectPortMatchCondition 445 -PriorityValue8021Action 5 -PolicyStore ActiveStore

PS C:\Users\Administrators> New-NetQosPolicy "TCP" -IPProtocolMatchCondition TCP -PriorityValue8021Action 0 -Policystore ActiveStore

PS C:\Users\Administrator> Get-NetQosPolicy -PolicyStore activestore

Name : tcp

Owner : PowerShell / WMI

NetworkProfile : All Precedence : 127

JobObject :

IPProtocol : TCP
PriorityValue : 0

Name : smb

Owner : PowerShell / WMI

NetworkProfile : All
Precedence : 127
JobObject :
NetDirectPort : 445
PriorityValue : 5

10. To configure ETS for all traffic classes defined in the previous step, issue the following commands:

PS C:\Users\Administrators> New-NetQosTrafficClass -name "RDMA class" -priority 5 -bandwidthPercentage 50 -Algorithm ETS

PS C:\Users\Administrators> New-NetQosTrafficClass -name "TCP class" -priority 0 -bandwidthPercentage 30 -Algorithm ETS

PS C:\Users\Administrator> Get-NetQosTrafficClass

Name	Algorithm	Bandwidth(%)	Priority	PolicySet	IfIndex IfAlias
[Default]	ETS	20	1-4,6-7	Global	
RDMA class	ETS	50	5	Global	
TCP class	ETS	30	0	Global	

11. To see the network adapter QoS from the preceding configuration, issue the following Windows PowerShell command:

PS C:\Users\Administrator> Get-NetAdapterQos

Name : SLOT 4 Port 1

Enabled : True

Capabilities : Hardware Current

MacSecBypass : NotSupported NotSupported

DcbxSupport : None None NumTCs (Max/ETS/PFC) : 4/4/4 4/4/4

OperationalTrafficClasses : TC TSA Bandwidth Priorities

0 ETS 20% 1-4,6-7 1 ETS 50% 5 2 ETS 30% 0

OperationalFlowControl : Priority 5 Enabled

 ${\tt OperationalClassifications} \ : \ {\tt Protocol} \quad {\tt Port/Type} \ {\tt Priority}$

Default 0

NetDirect 445 5

- 12. Create a startup script to make the settings persistent across the system reboots.
- 13. Run RDMA traffic and verify as described in "RoCE Configuration" on page 127.

Configuring QoS by Enabling DCBX on the Adapter

All configuration must be completed on all of the systems in use. The PFC, ETS, and traffic classes configuration must be the same on the switch and server.

To configure QoS by enabling DCBX:

- 1. Enable DCBX (IEEE, CEE, or Dynamic).
- 2. Using HII, set the RoCE Priority to 0.

NOTE

If the switch does not have a way of designating the RoCE traffic, you may need to set the **RoCE Priority** to the number used by the switch. Arista[®] switches can do so, but some other switches cannot.

3. To install the DCB role in the host, issue the following Windows PowerShell command:

PS C:\Users\Administrators> Install-WindowsFeature
Data-Center-Bridging

NOTE

For this configuration, set the **DCBX Protocol** to **CEE**.

4. To set the **DCBX Willing** mode to **True**, issue the following command:

PS C:\Users\Administrators> set-NetQosDcbxSetting -Willing 1

- 5. Enable QoS in the miniport properties as follows:
 - On the adapter properties' Advanced page (Figure 15-11) under Property, select Quality of Service, and then set the value to Enabled.
 - b. Click **OK**.

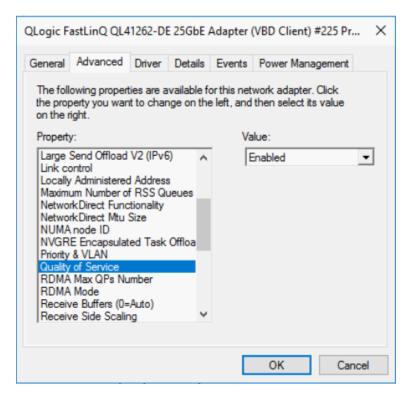


Figure 15-11. Advanced Properties: Enabling QoS

- 6. Assign the vLAN ID to the interface (required for PFC) as follows:
 - a. Open the miniport properties, and then click the **Advanced** tab.
 - b. On the adapter properties' Advanced page (Figure 15-12) under **Property**, select **VLAN ID**, and then set the value.
 - c. Click OK.

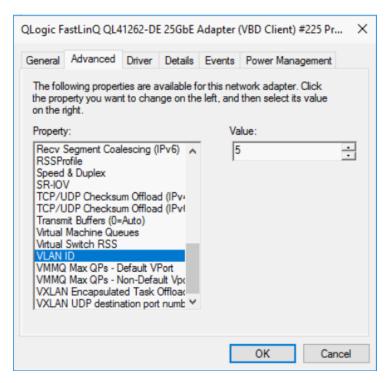


Figure 15-12. Advanced Properties: Setting VLAN ID

7. To configure the switch, issue the following Windows PowerShell command:

PS C:\Users\Administrators> Get-NetAdapterQoS

Name : Ethernet 5
Enabled : True

Capabilities : Hardware Current

MacSecBypass : NotSupported NotSupported

DcbxSupport : CEE CEE
NumTCs(Max/ETS/PFC) : 4/4/4 4/4/4

OperationalTrafficClasses : TC TSA Bandwidth Priorities

-- -- 0 ETS 5% 0-4,6-7 1 ETS 95% 5

OperationalFlowControl : Priority 5 Enabled

OperationalClassifications : Protocol Port/Type Priority

NetDirect 445 5

RemoteTrafficClasses : TC TSA Bandwidth Priorities

0 ETS 5% 0-4,6-7

1 ETS 95% 5

RemoteFlowControl : Priority 5 Enabled

RemoteClassifications : Protocol Port/Type Priority

NetDirect 445 5

NOTE

The preceding example is taken when the adapter port is connected to an Arista 7060X switch. In this example, the switch PFC is enabled on Priority 5. RoCE App TLVs are defined. The two traffic classes are defined as TC0 and TC1, where TC1 is defined for RoCE. **DCBX Protocol** mode is set to **CEE**. For Arista switch configuration, refer to "Preparing the Ethernet Switch" on page 129". When the adapter is in **Willing** mode, it accepts Remote Configuration and shows it as **Operational Parameters**.

Configuring VMMQ

Virtual machine multiqueue (VMMQ) configuration information includes:

- Enabling VMMQ on the Adapter
- Creating a Virtual Machine Switch with or Without SR-IOV
- Enabling VMMQ on the Virtual Machine Switch
- Getting the Virtual Machine Switch Capability
- Creating a VM and Enabling VMMQ on VMNetworkAdapters in the VM
- Enabling and Disabling VMMQ on a Management NIC
- Monitoring Traffic Statistics

Enabling VMMQ on the Adapter

To enable VMMQ on the adapter:

- 1. Open the miniport properties, and then click the **Advanced** tab.
- 2. On the adapter properties' Advanced page (Figure 15-13) under **Property**, select **Virtual Switch RSS**, and then set the value to **Enabled**.
- 3. Click OK.

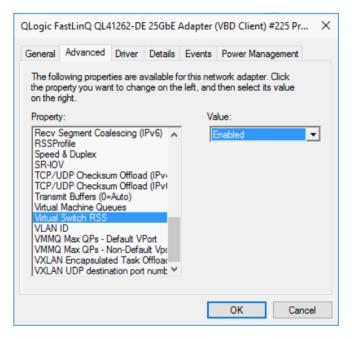


Figure 15-13. Advanced Properties: Enabling Virtual Switch RSS

Creating a Virtual Machine Switch with or Without SR-IOV

To create a virtual machine switch with or without SR-IOV:

- 1. Launch the Hyper-V Manager.
- 2. Select Virtual Switch Manager (see Figure 15-14).
- 3. In the **Name** box, type a name for the virtual switch.
- 4. Under Connection type:
 - Click External network.
 - b. Select the Allow management operating system to share this network adapter check box.

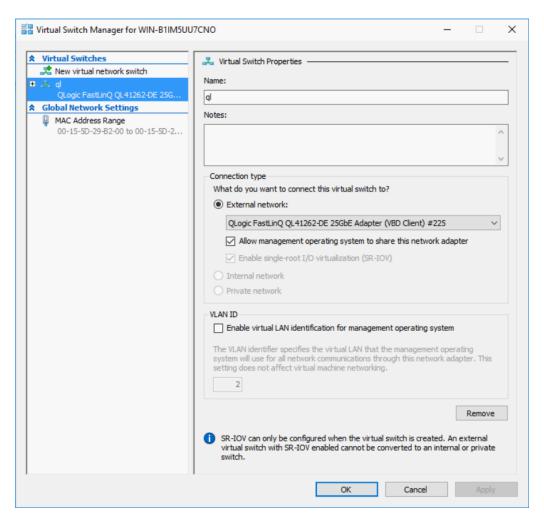


Figure 15-14. Virtual Switch Manager

5. Click OK.

Enabling VMMQ on the Virtual Machine Switch

To enable VMMQ on the virtual machine switch:

■ Issue the following Windows PowerShell command:

```
PS C:\Users\Administrators> Set-VMSwitch -name q1 -defaultqueuevmmqqnabled $true -defaultqueuevmmqqueuepairs 4
```

Getting the Virtual Machine Switch Capability

To get the virtual machine switch capability:

Issue the following Windows PowerShell command:

```
PS C:\Users\Administrator> Get-VMSwitch -Name ql | fl
```

Figure 15-15 shows example output.

```
PS C:\Users\Administrator> Get-VMSwitch -Name q1 | f1
Name
Id
                                                                                        q|
4dff5da3-f8bc-4146-a809-e1ddc6a04f7a
 Notes
                                                                                        {Microsoft Windows Filtering Platform, Microsoft Azure VFP Switch Extension, Microsoft NDIS Capture}
Extensions
                                                                                       None
False
False
True
BandwidthReservationMode
PacketDirectEnabled
 EmbeddedTeamingEnabled
IovEnabled
SwitchType
AllowManagementOS
NetAdapterInterfaceDescription
NetAdapterInterfaceDescription
TowSunport
                                                                                       External
True
                                                                                        OLogic FastLinQ QL41262-DE 25GbE Adapter (VBD Client) #225
{QLogic FastLinQ QL41262-DE 25GbE Adapter (VBD Client) #225}
True
NetAdapterInterfaceDescri
IovSupportReasons
AvailableIPSecSA
NumberIPSecSAAllocated
AvailableWQueues
NumberVmqAllocated
IovQueuePairCount
IovQueuePairFunde
IovVirtualFunctionCount
IovVirtualFunctionSInUse
PacketDirectInUse
                                                                                       0
103
                                                                                        127
                                                                                        96
                                                                                       96
0
False
True
True
False
16
16
 PacketDirectInUse
DefaultQueueVrssEnabledRequested
 DefaultQueueVrssEnabled :
DefaultQueueVrssEnabled :
DefaultQueueVmmqEnabledRequested :
DefaultQueueVmmqQueuePairsRequested :
DefaultQueueVmmqQueuePairs :
ReadwidtDansantes :
 BandwidthPercentage
DefaultFlowMinimumBandwidthAbsolute
DefaultFlowMinimumBandwidthWeight
                                                                                       CimSession: .
WIN-B1IM5UU7CNO
False
   imSession
  ComputerName
IsDeleted
```

Figure 15-15. Windows PowerShell Command: Get-VMSwitch

Creating a VM and Enabling VMMQ on VMNetworkAdapters in the VM

To create a virtual machine (VM) and enable VMMQ on VMNetworksadapters in the VM:

- 1. Create a VM.
- 2. Add the VMNetworkadapter to the VM.
- 3. Assign a virtual switch to the VMNetworkadapter.

4. To enable VMMQ on the VM, issue the following Windows PowerShell command:

PS C:\Users\Administrators> set-vmnetworkadapter -vmname vm1 -VMNetworkAdapterName "network adapter" -vmmqenabled \$true -vmmqqueuepairs 4

Enabling and Disabling VMMQ on a Management NIC

To enable or disable VMMQ on a management NIC:

- To enable VMMQ on a management NIC, issue the following command:
 - PS C:\Users\Administrator> Set-VMNetworkAdapter -ManagementOS -vmmqEnabled \$true
- To disable VMMQ on a management NIC, issue the following command:

PS C:\Users\Administrator> Set-VMNetworkAdapter -ManagementOS -vmmqEnabled \$false

A VMMQ will also be available for the multicast open shortest path first (MOSPF).

Monitoring Traffic Statistics

To monitor virtual function traffic in a virtual machine, issue the following Windows PowerShell command:

PS C:\Users\Administrator> Get-NetAdapterStatistics | fl

NOTE

Marvell supports the new parameter added for Windows Server 2016 and Windows Server 2019 to configure the maximum quantity of queue pairs on a virtual port. For details, see "Max Queue Pairs (L2) Per VPort" on page 282.

Configuring Storage Spaces Direct

Windows Server 2016 introduces Storage Spaces Direct, which allows you to build highly available and scalable storage systems with local storage. For more information, refer to the following Microsoft TechNet link:

https://technet.microsoft.com/en-us/windows-server-docs/storage/storage-spaces/st

Configuring the Hardware

Figure 15-16 shows an example of hardware configuration.

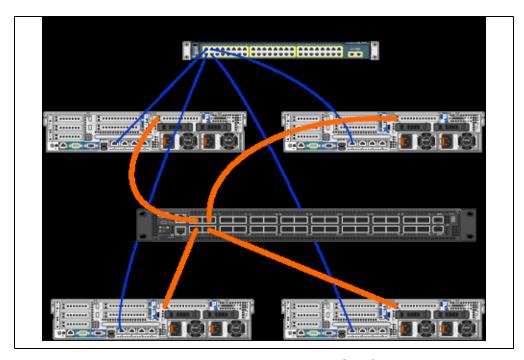


Figure 15-16. Example Hardware Configuration

NOTE

The disks used in this example are 4 × 400G NVMe[™], and 12 × 200G SSD disks.

Deploying a Hyper-Converged System

This section includes instructions to install and configure the components of a Hyper-Converged system using the Windows Server 2016. The act of deploying a Hyper-Converged system can be divided into the following three high-level phases:

- Deploying the Operating System
- Configuring the Network
- Configuring Storage Spaces Direct

Deploying the Operating System

To deploy the operating systems:

- Install the operating system.
- 2. Install the Windows Server roles (Hyper-V).
- 3. Install the following features:
 - □ Failover
 - □ Cluster
 - □ Data center bridging (DCB)
- 4. Connect the nodes to a domain and add domain accounts.

Configuring the Network

To deploy Storage Spaces Direct, the Hyper-V switch must be deployed with RDMA-enabled host virtual NICs.

NOTE

The following procedure assumes that there are four RDMA NIC ports.

To configure the network on each server:

- 1. Configure the physical network switch as follows:
 - Connect all adapter NICs to the switch port.

NOTE

If your test adapter has more than one NIC port, you must connect both ports to the same switch.

- b. Enable the switch port and make sure that:
 - The switch port supports switch-independent teaming mode.
 - The switch port is part of multiple vLAN networks.

Example Dell switch configuration:

```
no ip address
mtu 9416
portmode hybrid
switchport
dcb-map roce_S2D
protocol lldp
```

dcbx version cee no shutdown

2. Enable Network Quality of Service.

NOTE

Network Quality of Service is used to ensure that the Software Defined Storage system has enough bandwidth to communicate between the nodes to ensure resiliency and performance. To configure QoS on the adapter, see "Configuring QoS for RoCE" on page 259.

- 3. Create a Hyper-V virtual switch with Switch Embedded Teaming (SET) and RDMA virtual NIC as follows:
 - a. To identify the network adapters, issue the following command:

```
Get-NetAdapter | FT
Name,InterfaceDescription,Status,LinkSpeed
```

b. To create a virtual switch connected to all of the physical network adapters, and to then enable SET, issue the following command:

```
New-VMSwitch -Name SETswitch -NetAdapterName
"<port1>","<port2>","<port3>","<port4>"
-EnableEmbeddedTeaming $true
```

c. To add host virtual NICs to the virtual switch, issue the following commands:

Add-VMNetworkAdapter -SwitchName SETswitch -Name SMB_1 -managementOS

Add-VMNetworkAdapter -SwitchName SETswitch -Name SMB_2 -managementOS

NOTE

The preceding commands configure the virtual NIC from the virtual switch that you just configured for the management operating system to use.

d. To configure the host virtual NIC to use a vLAN, issue the following commands:

```
Set-VMNetworkAdapterVlan -VMNetworkAdapterName "SMB_1"
-VlanId 5 -Access -ManagementOS
Set-VMNetworkAdapterVlan -VMNetworkAdapterName "SMB_2"
-VlanId 5 -Access -ManagementOS
```

NOTE

These commands can be on the same or different vLANs.

e. To verify that the vLAN ID is set, issue the following command:

```
Get-VMNetworkAdapterVlan -ManagementOS
```

f. To disable and enable each host virtual NIC adapter so that the vLAN is active, issue the following commands:

```
Disable-NetAdapter "vEthernet (SMB_1)"
Enable-NetAdapter "vEthernet (SMB_1)"
Disable-NetAdapter "vEthernet (SMB_2)"
Enable-NetAdapter "vEthernet (SMB_2)"
```

g. To enable RDMA on the host virtual NIC adapters, issue the following command:

```
Enable-NetAdapterRdma "SMB1", "SMB2"
```

h. To verify RDMA capabilities, issue the following command:

```
Get-SmbClientNetworkInterface | where RdmaCapable -EQ
$true
```

Configuring Storage Spaces Direct

Configuring Storage Spaces Direct in Windows Server 2016 includes the following steps:

- Step 1. Running a Cluster Validation Tool
- Step 2. Creating a Cluster
- Step 3. Configuring a Cluster Witness
- Step 4. Cleaning Disks Used for Storage Spaces Direct
- Step 5. Enabling Storage Spaces Direct
- Step 6. Creating Virtual Disks
- Step 7. Creating or Deploying Virtual Machines

Step 1. Running a Cluster Validation Tool

Run the cluster validation tool to make sure server nodes are configured correctly to create a cluster using Storage Spaces Direct.

To validate a set of servers for use as a Storage Spaces Direct cluster, issue the following Windows PowerShell command:

Test-Cluster -Node <MachineName1, MachineName2, MachineName3,
MachineName4> -Include "Storage Spaces Direct", Inventory,
Network, "System Configuration"

Step 2. Creating a Cluster

Create a cluster with the four nodes (which was validated for cluster creation) in Step 1. Running a Cluster Validation Tool.

To create a cluster, issue the following Windows PowerShell command:

New-Cluster -Name <ClusterName> -Node <MachineName1, MachineName2,
MachineName3, MachineName4> -NoStorage

The -NoStorage parameter is required. If it is not included, the disks are automatically added to the cluster, and you must remove them before enabling Storage Spaces Direct. Otherwise, they will not be included in the Storage Spaces Direct storage pool.

Step 3. Configuring a Cluster Witness

You should configure a witness for the cluster, so that this four-node system can withstand two nodes failing or being offline. With these systems, you can configure file share witness or cloud witness.

For more information, go to:

https://docs.microsoft.com/en-us/windows-server/failover-clustering/manage-cluster-quorum

Step 4. Cleaning Disks Used for Storage Spaces Direct

The disks intended to be used for Storage Spaces Direct must be empty and without partitions or other data. If a disk has partitions or other data, it will not be included in the Storage Spaces Direct system.

The following Windows PowerShell command can be placed in a Windows PowerShell script (.PS1) file and executed from the management system in an open Windows PowerShell (or Windows PowerShell ISE) console with Administrator privileges.

NOTE

Running this script helps identify the disks on each node that can be used for Storage Spaces Direct. It also removes all data and partitions from those disks.

```
icm (Get-Cluster -Name HCNanoUSClu3 | Get-ClusterNode) {
Update-StorageProviderCache
```

```
Get-StoragePool |? IsPrimordial -eq $false | Set-StoragePool
-IsReadOnly:$false -ErrorAction SilentlyContinue
Get-StoragePool |? IsPrimordial -eq $false | Get-VirtualDisk |
Remove-VirtualDisk -Confirm: $false -ErrorAction SilentlyContinue
Get-StoragePool |? IsPrimordial -eq $false | Remove-StoragePool
-Confirm: $false -ErrorAction SilentlyContinue
Get-PhysicalDisk | Reset-PhysicalDisk -ErrorAction
SilentlyContinue
Get-Disk | ? Number -ne $null | ? IsBoot -ne $true | ? IsSystem -ne
$true |? PartitionStyle -ne RAW |% {
$ | Set-Disk -isoffline:$false
$ | Set-Disk -isreadonly:$false
$ | Clear-Disk -RemoveData -RemoveOEM -Confirm:$false
$ | Set-Disk -isreadonly:$true
$ | Set-Disk -isoffline:$true
Get-Disk | ? Number -ne $null | ? IsBoot -ne $true | ? IsSystem -ne
$true |? PartitionStyle -eq RAW | Group -NoElement -Property
FriendlyName
} | Sort -Property PsComputerName, Count
```

Step 5. Enabling Storage Spaces Direct

After creating the cluster, issue the <code>Enable-ClusterS2D</code> Windows PowerShell cmdlet. The cmdlet places the storage system into the Storage Spaces Direct mode and automatically does the following:

- Creates a single, large pool that has a name such as S2D on Cluster1.
- Configures a Storage Spaces Direct cache. If there is more than one media type available for Storage Spaces Direct use, it configures the most efficient type as cache devices (in most cases, read and write).
- Creates two tiers—Capacity and Performance—as default tiers. The cmdlet analyzes the devices and configures each tier with the mix of device types and resiliency.

Step 6. Creating Virtual Disks

If the Storage Spaces Direct was enabled, it creates a single pool using all of the disks. It also names the pool (for example S2D on Cluster1), with the name of the cluster that is specified in the name.

The following Windows PowerShell command creates a virtual disk with both mirror and parity resiliency on the storage pool:

New-Volume -StoragePoolFriendlyName "S2D*" -FriendlyName <VirtualDiskName> -FileSystem CSVFS_ReFS -StorageTierfriendlyNames Capacity,Performance -StorageTierSizes <Size of capacity tier in size units, example: 80GB>, <Size of Performance tier in size units, example: 80GB> -CimSession <ClusterName>

Step 7. Creating or Deploying Virtual Machines

You can provision the virtual machines onto the nodes of the hyper-converged S2D cluster. Store the virtual machine's files on the system's Cluster Shared Volume (CSV) namespace (for example, c:\ClusterStorage\Volume1), similar to clustered virtual machines on failover clusters.

16 Windows Server 2019

This chapter provides the following information for Windows Server 2019:

- RSSv2 for Hyper-V
- "Windows Server 2019 Behaviors" on page 281
- "New Adapter Properties" on page 282

RSSv2 for Hyper-V

In Windows Server 2019, Microsoft added support for Receive Side Scaling version 2 (RSSv2) with Hyper-V (RSSv2 per vPort).

RSSv2 Description

Compared to RSSv1, RSSv2 decreases the time between the CPU load measurement and the indirection table update. This feature prevents slowdown during high-traffic situations. RSSv2 can dynamically spread receive queues over multiple processors much more responsively than RSSv1. For more information, visit the following Web page:

https://docs.microsoft.com/en-us/windows-hardware/drivers/network/receive-side-scaling-version-2-rssv2-

RSSv2 is supported by default in the Windows Server 2019 driver when the **Virtual Switch RSS** option is also enabled. This option is enabled (the default), and the NIC is bound to the Hyper-V or vSwitch.

Known Event Log Errors

Under typical operation, the dynamic algorithm of RSSv2 may initiate an indirection table update that is incompatible with the driver and return an appropriate status code. In such cases, an event log error occurs, even though no functional operation issue exists. Figure 16-1 shows an example.



Figure 16-1. RSSv2 Event Log Error

Windows Server 2019 Behaviors

Windows Server 2019 introduced the following new behaviors affecting adapter configuration.

VMMQ Is Enabled by Default

In the inbox driver of Windows Server 2019, the **Virtual Switch RSS** (VMMQ) option is enabled by default in the NIC properties. In addition, Microsoft changed the default behavior of the **Virtual NICs** option to have VMMQ enabled with the 16 queue pairs. This behavior change impacts the quantity of available resources.

For example, suppose the NIC supports 32 VMQs and 64 queue pairs. In Windows Server 2016, when you add 32 virtual NICs (VNICs), they will have VMQ acceleration. However in Windows Server 2019, you will get 4 VNICs with VMMQ acceleration, each with 16 queue pairs and 30 VNICs with no acceleration.

Because of this functionality, Marvell introduced a new user property, **Max Queue Pairs (L2) Per VPort**. For more details, see New Adapter Properties.

Inbox Driver Network Direct (RDMA) Is Disabled by Default

In the inbox driver of Windows Server 2019, the **Network Direct** (RDMA) option is disabled by default in the NIC properties. However, when upgrading the driver to an out-of-box driver, **Network Direct** is enabled by default.

New Adapter Properties

New user-configurable properties available in Windows Server 2019 are described in the following sections:

- Max Queue Pairs (L2) Per VPort
- Network Direct Technology
- Virtualization Resources
- VMQ and VMMQ Default Accelerations
- Single VPort Pool

Max Queue Pairs (L2) Per VPort

As explained in VMMQ Is Enabled by Default, Windows 2019 (and Windows 2016) introduced a new user-configurable parameter, **Max Queue Pairs (L2) per VPort**. This parameter permits greater control of resource distribution by defining the maximum quantity of queue pairs that can be assigned to the following:

- VPort-Default VPort
- PF Non-Default VPort (VMQ/VMMQ)
- SR-IOV Non-Default VPort (VF)¹

The default value of the **Max Queue Pairs (L2) per VPort** parameter is set to **Auto**, which is one of the following:

- Max Queue Pairs for Default vPort = 8
- Max Queue Pairs for Non-Default vPort = 4

If you select a value less than 8, then:

- Max Queue Pairs for Default vPort = 8
- Max Queue Pairs for Non-Default vPort = value

If you select a value greater than 8, then:

- Max Queue Pairs for Default vPort = value
- Max Queue Pairs for Non-Default vPort = value

Network Direct Technology

Marvell supports the new **Network Direct Technology** parameter that allows you to select the underlying RDMA technology that adheres to the following Microsoft specification:

https://docs.microsoft.com/en-us/windows-hardware/drivers/network/inf-requirements-for-ndkpi

This option replaces the **RDMA Mode** parameter.

¹ This parameter also applies to Windows Server 2016.

Virtualization Resources

Table 16-1 lists the maximum quantities of virtualization resources in Windows 2019 for Dell 41xxx Series Adapters.

Table 16-1. Windows 2019 Virtualization Resources for Dell 41xxx Series
Adapters

Two-port NIC-only Single Function Non-CNA	Quantity
Maximum VMQs	102
Maximum VFs	80
Maximum QPs	112
Four-port NIC-only Single Function Non-CNA	Quantity
	Quantity 47
Function Non-CNA	

VMQ and VMMQ Default Accelerations

Table 16-2 lists the VMQ and VMMQ default and other values for accelerations in Windows Server 2019 for Dell 41xxx Series Adapters.

Table 16-2. Windows 2019 VMQ and VMMQ Accelerations

Two-port NIC-only Single Function Non-CNA	Default Value	Other Possible Values				
Maximum Queue Pairs (L2) per VPort a	Auto	1	2	4	8	16
Maximum VMQs	26	103	52	26	13	6
Default VPort Queue Pairs	8	8	8	8	8	16
PF Non-default VPort Queue Pairs	4	1	2	4	8	16
Four-port NIC-only Single Function Non-CNA	Default Value	Other Possible Values				
Maximum Queue Pairs (L2) per VPort ^a	Auto	1	2	4	8	16
Maximum VMQs	10	40	20	10	5	2
Default VPort Queue Pairs	8	8	8	8	8	16
PF Non-default VPort Queue Pairs	4	1	2	4	8	16

a Max Queue Pairs (L2) VPort is configurable parameter of NIC advanced properties.

Single VPort Pool

The 41xxx Series Adapter supports the **Single VPort Pool** parameter, which allows the system administrator to assign any available IOVQueuePair to either Default-VPort, PF Non-Default VPort, or VF Non-Default VPort. To assign the value, issue the following Windows PowerShell commands:

Default-VPort:

Set-VMSwitch -Name <vswitch name> -DefaultQueueVmmqEnabled:1
-DefaultQueueVmmqQueuePairs:<number>

NOTE

Marvell does not recommend that you disable VMMQ or decrease the quantity of queue pairs for the Default-VPort, because it may impact system performance.

- PF Non-Default VPort:
 - ☐ For the host:

```
Set-VMNetworkAdapter -ManagementOS -VmmqEnabled:1
-VmmqQueuePairs:<number>
```

☐ For the VM:

Set-VMNetworkAdapter -VMName <vm name> -VmmqEnabled:1
-VmmqQueuePairs:<number>

■ VF Non-Default VPort:

Set-VMNetworkAdapter -VMName <vm name> -IovWeight:100
-IovQueuePairsRequested:<number>

NOTE

The default quantity of QPs assigned for a VF (lovQueuePairsRequested) is still 1.

To apply multiple quantities of queue pairs to any vPort:

- The quantity of queue pairs must be less than or equal to the total number of CPU cores on the system.
- The quantity of queue pairs must be less than or equal to the value of **Max Queue Pairs (L2) Per VPort**. For more information, see Max Queue Pairs

 (L2) Per VPort.

17 Troubleshooting

This chapter provides the following troubleshooting information:

- Troubleshooting Checklist
- "Verifying that Current Drivers Are Loaded" on page 287
- "Testing Network Connectivity" on page 288
- "Microsoft Virtualization with Hyper-V" on page 289
- "Linux-specific Issues" on page 289
- "Miscellaneous Issues" on page 289
- "Collecting Debug Data" on page 290

Troubleshooting Checklist

CAUTION

Before you open the server cabinet to add or remove the adapter, review the "Safety Precautions" on page 5.

The following checklist provides recommended actions to resolve problems that may arise while installing the 41xxx Series Adapter or running it in your system.

- Inspect all cables and connections. Verify that the cable connections at the network adapter and the switch are attached properly.
- Verify the adapter installation by reviewing "Installing the Adapter" on page 6. Ensure that the adapter is properly seated in the slot. Check for specific hardware problems, such as obvious damage to board components or the PCI edge connector.
- Verify the configuration settings and change them if they are in conflict with another device.
- Verify that your server is using the latest BIOS.
- Try inserting the adapter in another slot. If the new position works, the original slot in your system may be defective.

- Replace the failed adapter with one that is known to work properly. If the second adapter works in the slot where the first one failed, the original adapter is probably defective.
- Install the adapter in another functioning system, and then run the tests again. If the adapter passes the tests in the new system, the original system may be defective.
- Remove all other adapters from the system, and then run the tests again. If the adapter passes the tests, the other adapters may be causing contention.

Verifying that Current Drivers Are Loaded

Ensure that the current drivers are loaded for your Windows, Linux, or VMware system.

Verifying Drivers in Windows

See the Device Manager to view vital information about the adapter, link status, and network connectivity.

Verifying Drivers in Linux

To verify that the qed.ko driver is loaded properly, issue the following command:

```
# lsmod | grep -i <module name>
```

If the driver is loaded, the output of this command shows the size of the driver in bytes. The following example shows the drivers loaded for the ged module:

```
# 1smod | grep -i qed
qed 199238 1
qede 1417947 0
```

If you reboot after loading a new driver, you can issue the following command to verify that the currently loaded driver is the correct version:

modinfo qede

Or, you can issue the following command:

```
[root@test1]# ethtool -i eth2
driver: qede
version: 8.4.7.0
firmware-version: mfw 8.4.7.0 storm 8.4.7.0
bus-info: 0000:04:00.2
```

If you loaded a new driver, but have not yet rebooted, the <code>modinfo</code> command will not show the updated driver information. Instead, issue the following <code>dmesg</code> command to view the logs. In this example, the last entry identifies the driver that will be active upon reboot.

```
# dmesg | grep -i "Cavium" | grep -i "qede"

[ 10.097526] QLogic FastLinQ 4xxxx Ethernet Driver qede x.x.x.x
[ 23.093526] QLogic FastLinQ 4xxxx Ethernet Driver qede x.x.x.x
[ 34.975396] QLogic FastLinQ 4xxxx Ethernet Driver qede x.x.x.x
[ 34.975896] QLogic FastLinQ 4xxxx Ethernet Driver qede x.x.x.x
[ 3334.975896] QLogic FastLinQ 4xxxx Ethernet Driver qede x.x.x.x
```

Verifying Drivers in VMware

To verify that the VMware ESXi drivers are loaded, issue the following command:

esxcli software vib list

Testing Network Connectivity

This section provides procedures for testing network connectivity in Windows and Linux environments.

NOTE

When using forced link speeds, verify that both the adapter and the switch are forced to the same speed.

Testing Network Connectivity for Windows

Test network connectivity using the ping command.

To determine if the network connection is working:

- 1. Click **Start**, and then click **Run**.
- 2. In the **Open** box, type cmd, and then click **OK**.
- To view the network connection to be tested, issue the following command: ipconfig /all
- 4. Issue the following command, and then press ENTER.

```
ping <ip address>
```

The displayed ping statistics indicate whether or not the network connection is working.

Testing Network Connectivity for Linux

To verify that the Ethernet interface is up and running:

- 1. To check the status of the Ethernet interface, issue the ifconfig command.
- 2. To check the statistics on the Ethernet interface, issue the netstat -i command.

To verify that the connection has been established:

1. Ping an IP host on the network. From the command line, issue the following command:

ping <ip_address>

2. Press ENTER.

The displayed ping statistics indicate whether or not the network connection is working.

The adapter link speed can be forced to 10Gbps or 25Gbps using either the operating system GUI tool or the ethtool command, ethtool -s ethx speed SSSS.

Microsoft Virtualization with Hyper-V

Microsoft Virtualization is a hypervisor virtualization system for Windows Server 2012 R2. For more information on Hyper-V, go to:

https://technet.microsoft.com/en-us/library/Dn282278.aspx

Linux-specific Issues

Problem: Errors appear when compiling driver source code.

Solution: Some installations of Linux distributions do not install the

development tools and kernel sources by default. Before compiling driver source code, ensure that the development tools for the Linux

distribution that you are using are installed.

Miscellaneous Issues

Problem: The 41xxx Series Adapter has shut down, and an error message

appears indicating that the fan on the adapter has failed.

Solution: The 41xxx Series Adapter may intentionally shut down to prevent

permanent damage. Contact Marvell Technical Support for

assistance.

Problem: In an ESXi environment, with the iSCSI driver (qedil) installed,

sometimes, the VI-client cannot access the host. This is due to the termination of the hostd daemon, which affects connectivity with

the VI-client.

Solution: Contact VMware technical support.

Collecting Debug Data

Use the commands in Table 17-1 to collect debug data.

Table 17-1. Collecting Debug Data Commands

Debug Data	Description	
demesg -T	Kernel logs	
ethtool -d	Register dump	
sys_info.sh	System information; available in the driver bundle	

A Adapter LEDS

Table A-1 lists the LED indicators for the state of the adapter port link and activity.

Table A-1. Adapter Port Link and Activity LEDs

Port LED	LED Appearance	Network State
	Off	No link (cable disconnected or port down)
Link LED	Continuously illuminated GREEN	Link at highest supported link speed
	Continuously illuminated AMBER	Link at lower supported link speed
Activity	Off	No port activity
LED	Blinking	Port activity

B Cables and Optical Modules

This appendix provides the following information for the supported cables and optical modules:

- Supported Specifications
- "Tested Cables and Optical Modules" on page 293
- "Tested Switches" on page 297

Supported Specifications

The 41xxx Series Adapters support a variety of cables and optical modules that comply with SFF8024. Specific form factor compliance is as follows:

- SFPs:
 - ☐ SFF8472 (for memory map)
 - □ SFF8419 or SFF8431 (low speed signals and power)
- Optical modules electrical input/output, active copper cables (ACC), and active optical cables (AOC):
 - □ 10G—SFF8431 limiting interface
 - □ 25G—IEEE 802.3by Annex 109B (25GAUI) (does not support RS-FEC)

Tested Cables and Optical Modules

Marvell does not guarantee that every cable or optical module that satisfies the compliance requirements will operate with the 41xxx Series Adapters. Marvell has tested the components listed in Table B-1 and presents this list for your convenience.

Table B-1. Tested Cables and Optical Modules

Speed/Form Factor	Manufac- turer	Part Number	Туре	Cable Length ^a	Gauge
		Cab	les		
		1539W	SFP+10G-to-SFP+10G	1	26
	Brocade [®]	V239T	SFP+10G-to-SFP+10G	3	26
		48V40	SFP+10G-to-SFP+10G	5	26
		H606N	SFP+10G-to-SFP+10G	1	26
	Cisco	K591N	SFP+10G-to-SFP+10G	3	26
100 D 10h		G849N	SFP+10G-to-SFP+10G	5	26
10G DAC ^b		V250M	SFP+10G-to-SFP+10G	1	26
		53HVN	SFP+10G-to-SFP+10G	3	26
	D. II	358VV	SFP+10G-to-SFP+10G	5	26
	Dell	407-BBBK	SFP+10G-to-SFP+10G	1	30
		407-BBBI	SFP+10G-to-SFP+10G	3	26
		407-BBBP	SFP+10G-to-SFP+10G	5	26
		NDCCGF0001	SFP28-25G-to-SFP28-25G	1	30
		NDCCGF0003	SFP28-25G-to-SFP28-25G	3	30
	Amphenol [®]	NDCCGJ0003	SFP28-25G-to-SFP28-25G	3	26
050 DA0		NDCCGJ0005	SFP28-25G-to-SFP28-25G	5	26
25G DAC		2JVDD	SFP28-25G-to-SFP28-25G	1	26
	Dell	D0R73	SFP28-25G-to-SFP28-25G	2	26
	Dell	OVXFJY	SFP28-25G-to-SFP28-25G	3	26
		9X8JP	SFP28-25G-to-SFP28-25G	5	26

Table B-1. Tested Cables and Optical Modules (Continued)

Speed/Form Factor	Manufac- turer	Part Number	Туре	Cable Length ^a	Gauge
40G Copper		TCPM2	QSFP+40G-to-4xSFP+10G	1	30
QSFP Splitter	Dell	27GG5	QSFP+40G-to-4xSFP+10G	3	30
(4 × 10G)		P8T4W	QSFP+40G-to-4xSFP+10G	5	26
1G Copper		8T47V	SFP+ to 1G RJ	1G RJ45	N/A
RJ45	Dell	XK1M7	SFP+ to 1G RJ	1G RJ45	N/A
Transceiver		XTY28	SFP+ to 1G RJ	1G RJ45	N/A
10G Copper RJ45 Transceiver	Dell	PGYJT	SFP+ to 10G RJ	10G RJ45	N/A
40G DAC		470-AAVO	QSFP+40G-to-4xSFP+10G	1	26
Splitter	Dell	470-AAXG	QSFP+40G-to-4xSFP+10G	3	26
(4 × 10G)		470-AAXH	QSFP+40G-to-4xSFP+10G	5	26
		NDAQGJ-0001	QSFP28-100G-to- 4xSFP28-25G	1	26
	Amphenol	NDAQGF-0002	QSFP28-100G-to- 4xSFP28-25G	2	30
	Amphenoi	NDAQGF-0003	QSFP28-100G-to- 4xSFP28-25G	3	30
100G DAC Splitter		NDAQGJ-0005	QSFP28-100G-to- 4xSFP28-25G	5	26
(4 × 25G)		026FN3 Rev A00	QSFP28-100G-to- 4XSFP28-25G	1	26
	Dell	0YFNDD Rev A00	QSFP28-100G-to- 4XSFP28-25G	2	26
		07R9N9 Rev A00	QSFP28-100G-to- 4XSFP28-25G	3	26
	FCI	10130795-4050LF	QSFP28-100G-to- 4XSFP28-25G	5	26

Table B-1. Tested Cables and Optical Modules (Continued)

Speed/Form Factor	Manufac- turer	Part Number	Туре	Cable Length ^a	Gauge		
	Optical Solutions						
	Avago [®]	AFBR-703SMZ	SFP+ SR	N/A	N/A		
	Avago	AFBR-701SDZ	SFP+ LR	N/A	N/A		
		Y3KJN	SFP+ SR	1G/10G	N/A		
	Dell	WTRD1	SFP+ SR	10G	N/A		
10G Optical Transceiver	Dell	3G84K	SFP+ SR	10G	N/A		
Transcerver		RN84N	SFP+ SR	10G-LR	N/A		
	Finisar®	FTLX8571D3BCL- QL	SFP+ SR	N/A	N/A		
	riilisai°	FTLX1471D3BCL- QL	SFP+ LR	N/A	N/A		
	Dell	P7D7R	SFP28 Optical Transceiver SR	25G SR	N/A		
25G Optical Transceiver	Finisar	FTLF8536P4BCL	SFP28 Optical Transceiver SR	N/A	N/A		
	i iilisai	FTLF8538P4BCL	SFP28 Optical Transceiver SR no FEC	N/A	N/A		
10/25G Dual Rate Transceiver	Dell	M14MK	SFP28	N/A	N/A		

Table B-1. Tested Cables and Optical Modules (Continued)

Speed/Form Factor	Manufac- turer	Part Number	Туре	Cable Length ^a	Gauge
		470-ABLV	SFP+ AOC	2	N/A
		470-ABLZ	SFP+ AOC	3	N/A
		470-ABLT	SFP+ AOC	5	N/A
		470-ABML	SFP+ AOC	7	N/A
		470-ABLU	SFP+ AOC	10	N/A
		470-ABMD	SFP+AOC	15	N/A
100 1000	Dell	470-ABMJ	SFP+ AOC	20	N/A
10G AOC ^c	Dell	YJF03	SFP+ AOC	2	N/A
		P9GND	SFP+ AOC	3	N/A
		T1KCN	SFP+ AOC	5	N/A
		1DXKP	SFP+AOC	7	N/A
		MT7R2	SFP+ AOC	10	N/A
		K0T7R	SFP+ AOC	15	N/A
	W5G04	SFP+AOC	20	N/A	
	Dell	X5DH4	SFP28 AOC	20	N/A
25G AOC	1 1 1 1 1	TF-PY003-N00	SFP28 AOC	3	N/A
	InnoLight [®]	TF-PY020-N00	SFP28 AOC	20	N/A

^a Cable length is indicated in meters.

^b DAC is direct attach cable.

^c AOC is active optical cable.

Tested Switches

Table B-2 lists the switches that have been tested for interoperability with the 41xxx Series Adapters. This list is based on switches that are available at the time of product release, and is subject to change over time as new switches enter the market or are discontinued.

Table B-2. Switches Tested for Interoperability

Manufacturer	Ethernet Switch Model
Arista	7060X 7160
Cisco	Nexus 3132 Nexus 3232C Nexus 5548 Nexus 5596T Nexus 6000
Dell EMC	S6100 Z9100
HPE	FlexFabric 5950
Mellanox	SN2410 SN2700

C Dell Z9100 Switch Configuration

The 41xxx Series Adapters support connections with the Dell Z9100 Ethernet Switch. However, until the auto-negotiation process is standardized, the switch must be explicitly configured to connect to the adapter at 25Gbps.

To configure a Dell Z9100 switch port to connect to the 41xxx Series Adapter at 25Gbps:

- 1. Establish a serial port connection between your management workstation and the switch.
- 2. Open a command line session, and then log in to the switch as follows:

```
Login: admin
Password: admin
```

3. Enable configuration of the switch port:

```
Dell> enable
Password: xxxxxx
Dell# config
```

4. Identify the module and port to be configured. The following example uses module 1, port 5:

```
Dell(conf) #stack-unit 1 port 5 ?

portmode Set portmode for a module

Dell(conf) #stack-unit 1 port 5 portmode ?

dual Enable dual mode

quad Enable quad mode

single Enable single mode

Dell(conf) #stack-unit 1 port 5 portmode quad ?

speed Each port speed in quad mode

Dell(conf) #stack-unit 1 port 5 portmode quad speed ?

Quad port mode with 10G speed
```

```
25G Quad port mode with 25G speed Dell(conf)#stack-unit 1 port 5 portmode quad speed 25G
```

For information about changing the adapter link speed, see "Testing Network Connectivity" on page 288.

5. Verify that the port is operating at 25Gbps:

```
Dell# Dell#show running-config | grep "port 5" stack-unit 1 port 5 portmode quad speed 25G
```

- 6. To disable auto-negotiation on switch port 5, follow these steps:
 - a. Identify the switch port interface (module 1, port 5, interface 1) and confirm the auto-negotiation status:

b. Disable auto-negotiation:

```
Dell(conf-if-tf-1/5/1) #no intf-type cr4 autoneg
```

c. Verify that auto-negotiation is disabled.

```
Dell(conf-if-tf-1/5/1)#do show run interface tw 1/5/1
!
interface twentyFiveGigE 1/5/1
no ip address
mtu 9416
switchport
flowcontrol rx on tx on
no shutdown
no intf-type cr4 autoneg
```

For more information about configuring the Dell Z9100 switch, refer to the *Dell Z9100 Switch Configuration Guide* on the Dell Support Web site:

support.dell.com

D Feature Constraints

This appendix provides information about feature constraints implemented in the current release.

These feature coexistence constraints may be removed in a future release. At that time, you should be able to use the feature combinations without any additional configuration steps beyond what would be usually required to enable the features.

Concurrent FCoE and iSCSI Is Not Supported on the Same Port in NPAR Mode

The device does not support configuration of both FCoE-Offload and iSCSI-Offload on the same port when in NPAR Mode. FCoE-Offload is supported on the second physical function (PF) and iSCSI-Offload is supported on the third PF in NPAR mode. The device does support configuration of both FCoE-Offload and iSCSI-Offload on the same port when in single Ethernet PF DEFAULT Mode. Not all devices support FCoE-Offload and iSCSI-Offload.

After a PF with either an iSCSI or FCoE personality has been configured on a port using either HII or Marvell management tools, configuration of the storage protocol on another PF is disallowed by those management tools.

Because storage personality is disabled by default, only the personality that has been configured using HII or Marvell management tools is written in NVRAM configuration. When this limitation is removed, users can configure additional PFs on the same port for storage in NPAR Mode.

Concurrent RoCE and iWARP Is Not Supported on the Same Physical Function

RoCE and iWARP are not supported on the same PF. The UEFI HII and Marvell management tools allow users to configure both concurrently, but the RoCE functionality takes precedence over the iWARP functionality in this case, unless overridden by the in-OS driver settings.

NIC and SAN Boot to Base Is Supported Only on Select PFs

Ethernet (such as software iSCSI remote boot) and PXE boot are currently supported only on the first Ethernet PF of a physical port. In NPAR Mode configuration, the first Ethernet PF (that is, not the other Ethernet PFs) supports Ethernet (such as software iSCSI remote boot) and PXE boot. Not all devices support FCoE-Offload and iSCSI-Offload.

- When the **Virtualization** or **Multi-Function Mode** is set to **NPAR**, FCoE-Offload boot is supported on the second PF of a physical port, iSCSI-Offload boot is supported on the third PF of a physical port, and Ethernet (such as software iSCSI) and PXE boot are supported on the first PF of a physical port.
- iSCSI and FCoE boot is limited to a single target per boot session.
- Only one boot mode is allowed per physical port.
- iSCSI-Offload and FCoE-Offload boot is only supported in NPAR mode.

E Revision History

Document Revision History	
Revision A, April 28, 2017	
Revision B, August 24, 2017	
Revision C, October 1, 2017	
Revision D, January 24, 2018	
Revision E, March 15, 2018	
Revision F, April 19, 2018	
Revision G, May 22, 2018	
Revision H, August 23, 2018	
Revision J, January 23, 2019	
Revision K, July 2, 2019	
Revision L, July 3, 2019	
Revision M, October 16, 2019	
Changes	Sections Affected
Added the following adapters to the list of Marvell products: QL41164HFRJ-DE, QL41164HFRJ-DE, QL41164HFCU-DE, QL41232HMKR-DE, QL41262HMKR-DE, QL41232HFCU-DE, QL41232HLCU-DE, QL41132HFRJ-DE, QL41132HLRJ-DE, QL41132HQRJ-DE, QL41232HQCU-DE, QL41132HQCU-DE, QL41154HQRJ-DE, QL41154HQCU-DE	"Supported Products" on page xvi
Added support for VMDirectPath I/O.	"Features" on page 1
In Table 2-2, updated the supported OSs for Windows Server, RHEL, SLES, XenServer.	"System Requirements" on page 4

In the bullets following the second paragraph, added text to further describe Dell iSCSI HW and SW installation.

Moved section to be closer to other relevant sections

In the first paragraph, corrected the first sentence to "The **Boot Mode** option is listed under **NIC Configuration...**"

Added instructions for setting UEFI iSCSI HBA.

Removed sections "Configuring iSCSI Boot Parameters" and "Configuring BIOS Boot Mode".

Added a reference to Application Note, Enabling Storage Offloads on Dell and Marvell FastLinQ 41000 Series Adapters.

In Step 3, clarified how the RoCE v1 Priority values are used.

At the end of the section, added a NOTE with an example of how to install iSCSI BFS in Linux with an MPIO configuration and a single path active.

Updated the steps for slipstreaming adapter drivers into Windows image files.

Removed the bullet stating that "RoCE does not work over a VF in an SR-IOV environment." VF RDMA is now supported.

In Table 7-1,

Removed RHEL 7.5; added RHEL 7.7. For RHEL7.6, added support for OFED-4.17-1 GA. Removed SLES 12 SP3; added SLES 12 SP4. Separated SLES 15 (SP0) and SLES 15 SP1; in SLES 15 SP1, added support for OFED-4.17-1 GA

CentOS 7.6: added support for OFED-4.17-1 GA.

Added information on configuring RoCE for VF RDMA for Windows and Linux.

"iSCSI Preboot Configuration" on page 68

"Configuring the Storage Target" on page 71

"Selecting the iSCSI UEFI Boot Protocol" on page 72

"Boot from SAN Configuration" on page 67

"FCoE Support" on page 37, "iSCSI Support" on page 37, "Configuring FCoE Boot" on page 55, "Configuring iSCSI Boot" on page 56, "Boot from SAN Configuration" on page 67, "iSCSI Configuration" on page 200, "FCoE Configuration" on page 213

"Configuring Data Center Bridging" on page 53

"Configuring iSCSI Boot from SAN for RHEL 7.5 and Later" on page 95

"Injecting (Slipstreaming) Adapter Drivers into Windows Image Files" on page 122

"Planning for RoCE" on page 128

"Supported Operating Systems and OFED" on page 127

"Configuring RoCE for SR-IOV VF Devices (VF RDMA)" on page 141 Configuring RoCE for SR-IOV VF Devices (VF RDMA)

In Step 1, updated the second and third bullets to the currently supported OSs for SLES 12 and RHEL, respectively.

Changed Step 4 part b to "Set the RDMA Protocol Support to RoCE/iWARP or iWARP."

Removed reference to appendix C; added configuration information.

Updated list of OSs that support inbox OFED.

Removed section "iWARP RDMA-Core Support on SLES 12 SP3 and OFED 4.8x".

In the bulleted list following the third paragraph, updated the list of supported OSs (second bullet).

Clarified Step 4: "To enable RDMA through PowerShell, issue the following Windows PowerShell command".

In Step 2, corrected the last word of the Power-Shell command to -ManagementOS.

Changed the commands in Step 1, part c. Added an example of ovsdb-server and ovs-vswitchd running with pid. In Step 4, part c, changed the second paragraph, second sentence to "The br1 interface is named eth0, ens7; manually configure the static IP through th network device file and assign the same subnet IP to the peer (Host 2 VM)."

Changed the PowerShell command for monitoring virtual function traffic on a virtual machine.

Changed the link for more information about configuring a cluster witness.

In the first paragraph, changed the cmdlet to Enable-ClusterS2D.

In Table A-1, in the Link LED section, updated the LED Appearance and Network State columns.

"RoCE v2 Configuration for Linux" on page 155

"Preparing the Adapter for iWARP" on page 177

"Configuring the Dell Z9100 Ethernet Switch for RoCE" on page 131

"Before You Begin" on page 188

"iWARP Configuration" on page 177

"NVMe-oF Configuration with RDMA" on page 235

"Creating a Hyper-V Virtual Switch with an RDMA NIC" on page 253

"Adding a vLAN ID to Host Virtual NIC" on page 254

"Configuring VXLAN in Linux" on page 247

"Monitoring Traffic Statistics" on page 272

"Step 3. Configuring a Cluster Witness" on page 277

"Step 5. Enabling Storage Spaces Direct" on page 278

"Adapter LEDS" on page 291

Glossary

ACPI

The Advanced Configuration and Power Interface (ACPI) specification provides an open standard for unified operating system-centric device configuration and power management. The ACPI defines platform-independent interfaces for hardware discovery, configuration, power management, and monitoring. The specification is central to operating system-directed configuration and Power Management (OSPM), a term used to describe a system implementing ACPI, which therefore removes device management responsibilities from legacy firmware interfaces.

adapter

The board that interfaces between the host system and the target devices. Adapter is synonymous with Host Bus Adapter, host adapter, and board.

adapter port

A port on the adapter board.

Advanced Configuration and Power Interface

See ACPL

bandwidth

A measure of the volume of data that can be transmitted at a specific transmission rateA 1Gbps or 2Gbps Fibre Channel port can transmit or receive at nominal rates of 1 or 2Gbps, depending on the device to which it is connected. This corresponds to actual bandwidth values of 106MB and 212MB, respectively.

BAR

Base address register. Used to hold memory addresses used by a device, or offsets for port addresses. Typically, memory address BARs must be located in physical RAM while I/O space BARs can reside at any memory address (even beyond physical memory).

base address register

See BAR.

basic input output system

See BIOS.

BIOS

Basic input output system. Typically in Flash PROM, the program (or utility) that serves as an interface between the hardware and the operating system and allows booting from the adapter at startup.

challenge-handshake authentication protocol

See CHAP.

CHAP

Challenge-handshake authentication protocol (CHAP) is used for remote logon, usually between a client and server or a Web browser and Web server. A challenge/response is a security mechanism for verifying the identity of a person or process without revealing a secret password that is shared by the two entities. Also referred to as a *three-way handshake*.

CNA

See Converged Network Adapter.

Converged Network Adapter

Marvell Converged Network Adapters support both data networking (TCP/IP) and storage networking (Fibre Channel) traffic on a single I/O adapter using two new technologies: Enhanced Ethernet and Fibre Channel over Ethernet (FCoE).

data center bridging

See DCB.

data center bridging exchange

See DCBX.

DCB

Data center bridging. Provides enhancements to existing 802.1 bridge specifications to satisfy the requirements of protocols and applications in the data center. Because existing high-performance data centers typically comprise multiple application-specific networks that run on different link layer technologies (Fibre Channel for storage and Ethernet for network management and LAN connectivity), DCB enables 802.1 bridges to be used for the deployment of a converged network where all applications can be run over a single physical infrastructure.

DCBX

Data center bridging exchange. A protocol used by DCB devices to exchange configuration information with directly connected peers. The protocol may also be used for misconfiguration detection and for configuration of the peer.

device

A target, typically a disk drive. Hardware such as a disk drive, tape drive, printer, or keyboard that is installed in or connected to a system. In Fibre Channel, a target device.

DHCP

Dynamic host configuration protocol. Enables computers on an IP network to extract their configuration from servers that have information about the computer only after it is requested.

driver

The software that interfaces between the file system and a physical data storage device or network media.

dynamic host configuration protocol See DHCP.

eCore

A layer between the OS and the hardware and firmware. It is device-specific and OS-agnostic. When eCore code requires OS services (for example, for memory allocation, PCI configuration space access, and so on) it calls an abstract OS function that is implemented in OS-specific layers. eCore flows may be driven by the hardware (for example, by an interrupt) or by the OS-specific portion of the driver (for example, loading and unloading the load and unload).

EEE

Energy-efficient Ethernet. A set of enhancements to the twisted-pair and backplane Ethernet family of computer networking standards that allows for less power consumption during periods of low data activity. The intention was to reduce power consumption by 50 percent or more, while retaining full compatibility with existing equipment. The Institute of Electrical and Electronics Engineers (IEEE), through the IEEE 802.3az task force, developed the standard.

EFI

Extensible firmware interface. A specification that defines a software interface between an operating system and platform firmware. EFI is a replacement for the older BIOS firmware interface present in all IBM PC-compatible personal computers.

energy-efficient Ethernet

See EEE.

enhanced transmission selection

See ETS.

Ethernet

The most widely used LAN technology that transmits information between computers, typically at speeds of 10 and 100 million bits per second (Mbps).

ETS

Enhanced transmission selection. A standard that specifies the enhancement of transmission selection to support the allocation of bandwidth among traffic classes. When the offered load in a traffic classe does not use its allocated bandwidth, enhanced transmission selection allows other traffic classes to use the available bandwidth. The bandwidth-allocation priorities coexist with strict priorities. ETS includes managed objects to support bandwidth allocation. For more information, refer to:

http://ieee802.org/1/pages/802.1az.html

extensible firmware interface

See EFI.

FCoE

Fibre Channel over Ethernet. A new technology defined by the T11 standards body that allows traditional Fibre Channel storage networking traffic to travel over an Ethernet link by encapsulating Fibre Channel frames inside Layer 2 Ethernet frames. For more information, visit www.fcoe.com.

Fibre Channel

A high-speed serial interface technology that supports other higher layer protocols such as SCSI and IP.

Fibre Channel over Ethernet

See FCoE.

file transfer protocol

See FTP.

FTP

File transfer protocol. A standard network protocol used to transfer files from one host to another host over a TCP-based network, such as the Internet. FTP is required for out-of-band firmware uploads that will complete faster than in-band firmware uploads.

HBA

See Host Bus Adapter.

HII

Human interface infrastructure. A specification (part of UEFI 2.1) for managing user input, localized strings, fonts, and forms, that allows OEMs to develop graphical interfaces for preboot configuration.

host

One or more adapters governed by a single memory or CPU complex.

Host Bus Adapter

An adapter that connects a host system (the computer) to other network and storage devices.

human interface infrastructure

See HII.

IEEE

Institute of Electrical and Electronics Engineers. An international nonprofit organization for the advancement of technology related to electricity.

Internet Protocol

See IP.

Internet small computer system interface

See iSCSI.

Internet wide area RDMA protocol

See iWARP.

IΡ

Internet protocol. A method by which data is sent from one computer to another over the Internet. IP specifies the format of packets, also called *datagrams*, and the addressing scheme.

IQN

iSCSI qualified name. iSCSI node name based on the initiator manufacturer and a unique device name section.

iSCSI

Internet small computer system interface. Protocol that encapsulates data into IP packets to send over Ethernet connections.

iSCSI qualified name

See IQN.

iWARP

Internet wide area RDMA protocol. A networking protocol that implements RDMA for efficient data transfer over IP networks. iWARP is designed for multiple environments, including LANs, storage networks, data center networks, and WANs.

jumbo frames

Large IP frames used in high-performance networks to increase performance over long distances. Jumbo frames generally means 9,000 bytes for Gigabit Ethernet, but can refer to anything over the IP MTU, which is 1,500 bytes on an Ethernet.

large send offload

See LSO.

Layer 2

Refers to the data link layer of the multilayered communication model, Open Systems Interconnection (OSI). The function of the data link layer is to move data across the physical links in a network, where a switch redirects data messages at the Layer 2 level using the destination MAC address to determine the message destination.

Link Layer Discovery Protocol

See LLDP.

LLDP

A vendor-neutral Layer 2 protocol that allows a network device to advertise its identity and capabilities on the local network. This protocol supersedes proprietary protocols like Cisco Discovery Protocol, Extreme Discovery Protocol, and Nortel Discovery Protocol (also known as SONMP).

Information gathered with LLDP is stored in the device and can be queried using SNMP. The topology of a LLDP-enabled network can be discovered by crawling the hosts and querying this database.

LSO

Large send offload. LSO Ethernet adapter feature that allows the TCP\IP network stack to build a large (up to 64KB) TCP message before sending it to the adapter. The adapter hardware segments the message into smaller data packets (frames) that can be sent over the wire: up to 1,500 bytes for standard Ethernet frames and up to 9,000 bytes for jumbo Ethernet frames. The segmentation process frees up the server CPU from having to segment large TCP messages into smaller packets that will fit inside the supported frame size.

maximum transmission unit

See MTU.

message signaled interrupts

See MSI, MSI-X.

MSI, MSI-X

Message signaled interrupts. One of two PCI-defined extensions to support message signaled interrupts (MSIs), in PCI 2.2 and later and PCI Express. MSIs are an alternative way of generating an interrupt through special messages that allow emulation of a pin assertion or deassertion.

MSI-X (defined in PCI 3.0) allows a device to allocate any number of interrupts between 1 and 2,048 and gives each interrupt separate data and address registers. Optional features in MSI (64-bit addressing and interrupt masking) are mandatory with MSI-X.

MTU

Maximum transmission unit. Refers to the size (in bytes) of the largest packet (IP datagram) that a specified layer of a communications protocol can transfer.

network interface card

See NIC.

NIC

Network interface card. Computer card installed to enable a dedicated network connection.

NIC partitioning

See NPAR.

non-volatile random access memory

See NVRAM.

non-volatile memory express

See NVMe.

NPAR

NIC partitioning. The division of a single NIC port into multiple physical functions or partitions, each with a user-configurable bandwidth and personality (interface type). Personalities include NIC, FCoE, and iSCSI.

NVRAM

Non-volatile random access memory. A type of memory that retains data (configuration settings) even when power is removed. You can manually configure NVRAM settings or restore them from a file.

NVMe

A storage access method designed for sold-state drives (SSDs).

OFED™

OpenFabrics Enterprise Distribution. An open source software for RDMA and kernel bypass applications.

PCI™

Peripheral component interface. A 32-bit local bus specification introduced by Intel®.

PCI Express (PCIe)

A third-generation I/O standard that allows enhanced Ethernet network performance beyond that of the older peripheral component interconnect (PCI) and PCI extended (PCI-X) desktop and server slots.

QoS

Quality of service. Refers to the methods used to prevent bottlenecks and ensure business continuity when transmitting data over virtual ports by setting priorities and allocating bandwidth.

quality of service

See QoS.

PF

Physical function.

RDMA

Remote direct memory access. The ability for one node to write directly to the memory of another (with address and size semantics) over a network. This capability is an important feature of VI networks.

reduced instruction set computer

See RISC.

remote direct memory access

See RDMA.

RISC

Reduced instruction set computer. A computer microprocessor that performs fewer types of computer instructions, thereby operating at higher speeds.

RDMA over Converged Ethernet

See RoCE.

RoCE

RDMA over Converged Ethernet. A network protocol that allows remote direct memory access (RDMA) over a converged or a non-converged Ethernet network. RoCE is a link layer protocol that allows communication between any two hosts in the same Ethernet broadcast domain.

SCSI

Small computer system interface. A high-speed interface used to connect devices, such as hard drives, CD drives, printers, and scanners, to a computer. The SCSI can connect many devices using a single controller. Each device is accessed by an individual identification number on the SCSI controller bus.

SerDes

Serializer/deserializer. A pair of functional blocks commonly used in high-speed communications to compensate for limited input/output. These blocks convert data between serial data and parallel interfaces in each direction.

serializer/deserializer

See SerDes.

single root input/output virtualization

See SR-IOV.

small computer system interface

See SCSI.

SR-IOV

Single root input/output virtualization. A specification by the PCI SIG that enables a single PCIe device to appear as multiple, separate physical PCIe devices. SR-IOV permits isolation of PCIe resources for performance, interoperability, and manageability.

target

The storage-device endpoint of a SCSI session. Initiators request data from targets. Targets are typically disk-drives, tape-drives, or other media devices. Typically a SCSI peripheral device is the target but an adapter may, in some cases, be a target. A target can contain many LUNs.

A target is a device that responds to a requested by an initiator (the host system). Peripherals are targets, but for some commands (for example, a SCSI COPY command), the peripheral may act as an initiator.

TCP

Transmission control protocol. A set of rules to send data in packets over the Internet protocol.

TCP/IP

Transmission control protocol/Internet protocol. Basic communication language of the Internet.

TLV

Type-length-value. Optional information that may be encoded as an element inside of the protocol. The type and length fields are fixed in size (typically 1–4 bytes), and the value field is of variable size. These fields are used as follows:

- Type—A numeric code that indicates the kind of field that this part of the message represents.
- Length—The size of the value field (typically in bytes).
- Value—Variable-sized set of bytes that contains data for this part of the message.

transmission control protocol

See TCP.

transmission control protocol/Internet protocol

See TCP/IP.

type-length-value

See TLV.

UDP

User datagram protocol. A connectionless transport protocol without any guarantee of packet sequence or delivery. It functions directly on top of IP.

UEFI

Unified extensible firmware interface. A specification detailing an interface that helps hand off control of the system for the preboot environment (that is, after the system is powered on, but before the operating system starts) to an operating system, such as Windows or Linux. UEFI provides a clean interface between operating systems and platform firmware at boot time, and supports an architecture-independent mechanism for initializing add-in cards.

unified extensible firmware interface

See UEFI.

user datagram protocol

See UDP.

VF

Virtual function.

VI

Virtual interface. An initiative for remote direct memory access across Fibre Channel and other communication protocols. Used in clustering and messaging.

virtual interface

See VI.

virtual logical area network

See vLAN.

virtual machine

See VM.

virtual port

See vPort.

vLAN

Virtual logical area network (LAN). A group of hosts with a common set of requirements that communicate as if they were attached to the same wire, regardless of their physical location. Although a vLAN has the same attributes as a physical LAN, it allows for end stations to be grouped together even if they are not located on the same LAN segment. vLANs enable network reconfiguration through software, instead of physically relocating devices.

VM

Virtual machine. A software implementation of a machine (computer) that executes programs like a real machine.

vPort

Virtual port. Port number or service name associated with one or more virtual servers. A virtual port number should be the same TCP or UDP port number to which client programs expect to connect.

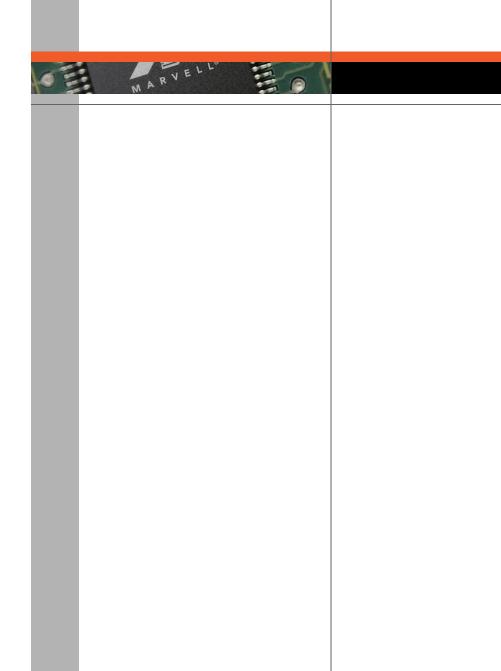
wake on LAN

See WoL.

WoL

Wake on LAN. An Ethernet computer networking standard that allows a computer to be remotely switched on or awakened by a network message sent usually by a simple program executed on another computer on the network.





Marvell Technology Group http://www.marvell.com