



# 用户指南

## 聚合网络适配器

41xxx 系列

文档修订历史	
修订版 A, 2017 年 4 月 28 日	
修订版 B, 2017 年 8 月 24 日	
修订版 C, 2017 年 10 月 1 日	
修订版 D, 2018 年 1 月 24 日	
修订版 E, 2018 年 3 月 15 日	
修订版 F, 2018 年 4 月 19 日	
更改	受影响的章节
<p>更新文件惯例示例。</p> <p>删除作废的 QLogic 许可协议和保证章节。</p> <p>在 <a href="#">表 3-5</a>, 添加脚注“公布本用户指南后, 可能提供其他 ESXi 驱动程序。有关更多信息, 请参考版本注释。”</p> <p>在 <a href="#">表 6-1</a> 中:</p> <ul style="list-style-type: none"> <li>■ 更新 Windows Server 和 VMware ESXi 的 OED 值。</li> <li>■ VMware ESXi 6.7 添加一行。</li> <li>■ 删除脚注: “此版本中未包括经认证的 RoCE 驱动程序。未经认证的驱动程序作为早期预览提供。”</li> </ul> <p>添加查看 Windows 上 RoCE 和 iWARP 的 Cavium RDMA 计数器的新步骤。</p> <p>遵循 <a href="#">图 7-4</a>, 添加注释, 交叉引用 <a href="#">第 67 页上的“查看 RDMA 计数器”</a> 步骤。</p> <p>添加配置 VMware ESXi 6.7 的 iSER 的信息。</p> <p>将以下主要章节更改为 Linux 环境中的 iSCSI 卸载下的小节:</p> <ul style="list-style-type: none"> <li>■ <a href="#">第 142 页上的“与 bnx2i 的差异”</a></li> <li>■ <a href="#">第 143 页上的“配置 qedi.ko”</a></li> <li>■ <a href="#">第 143 页上的“在 Linux 中验证 iSCSI 接口”</a></li> <li>■ <a href="#">第 145 页上的“Open-iSCSI 和从 SAN 引导注意事项”</a></li> </ul>	<p><a href="#">第 xviii 页上的“文档惯例”</a></p> <p><a href="#">前言</a></p> <p><a href="#">第 26 页上的“VMware 驱动程序和驱动程序包”</a></p> <p><a href="#">第 60 页上的“支持的操作系统和 OFED”</a></p> <p><a href="#">第 67 页上的“查看 RDMA 计数器”</a></p> <p><a href="#">第 89 页上的“在 Windows 上配置 iWARP”</a></p> <p><a href="#">第 108 页上的“在 ESXi 6.7 上配置 iSER”</a></p> <p><a href="#">第 142 页上的“Linux 环境中的 iSCSI 卸载”</a></p>

<p>在“<b>要从非卸载接口迁移到卸载接口</b>”步骤中：</p> <ul style="list-style-type: none"><li>■ 更改<b>步骤 1</b>为：“由 ..... 更新 open-iscsi 工具和 iscsiui0 到可用的最新版本”</li><li>■ 经编辑，<b>步骤 2</b>在第一个项目符号添加“（如果存在）”并删除最后一个项目符号（“附加 rd.driver.pre=qed rd.driver.pre=qedi）”</li></ul> <p>添加“及更高版本”到章节标题，并更新 <b>步骤 18</b>。</p> <p>添加注释：“在安装 SLES 11 或 SLES 12 时，不需要参数 withfcoe=1，因为 41000 系列适配器不再需要软件 FCoE 守护进程。”</p> <p>添加注释，描述 FCoE 接口当前功能。</p> <p>添加一个新的章节：FCoE 从 SAN 引导。</p> <p>删除作废的“从 SAN 引导注意事项”章节。</p> <p>在<b>要在 Linux 上配置 SR-IOV</b> 步骤中：</p> <ul style="list-style-type: none"><li>■ 在 <b>步骤 12</b>，将命令从 ip link show/ifconfig -a 更改为 ip link show   grep -i vf -b2。</li><li>■ 用新的截屏替换图 11-12。</li><li>■ 在 <b>步骤 15</b>，将命令从 check lspci -vv grep -I ether 更改为 lspci -vv grep -i ether。</li></ul> <p>在“<b>要在 VMware 上配置 SR-IOV</b>”步骤中，重排某些步骤：</p> <ul style="list-style-type: none"><li>■ 将步骤“要验证每个端口的 VF，请发出 esxcli 命令 .....”移到步骤“完成编辑设置对话框 .....”之后。</li><li>■ 将步骤“打开 VM 的电源.....”移到步骤“为检测到的适配器安装 QLogic 驱动程序”。</li></ul>	<p><a href="#">第 151 页上的“从 SAN 迁移的 SLES 11 SP4 iSCSI L4 引导”</a></p> <p><a href="#">第 156 页上的“从 RHEL 7.4 及更高版本的 SAN 配置 iSCSI 引导”</a></p> <p><a href="#">第 168 页上的“配置 Linux FCoE 卸载”</a></p> <p><a href="#">第 169 页上的“qedf 与 bnx2fc 之间的差异”</a></p> <p><a href="#">第 171 页上的“从 RHEL 7.4 及更高版本的 SAN 配置 FCoE 引导”</a></p> <p><a href="#">第 10 章 FCoE 配置</a></p> <p><a href="#">第 182 页上的“在 Linux 上配置 SR-IOV”</a></p> <p><a href="#">第 188 页上的“在 VMware 上配置 SR-IOV”</a></p>
--	--



# 目录

## 前言

支持的产品 .....	xvi
目标读者 .....	xvi
本指南的内容 .....	xvi
文档惯例 .....	xviii
法律声明 .....	xx
激光安全 - FDA 公告 .....	xx
机构认证 .....	xx
EMI 和 EMC 要求 .....	xx
KCC: A 级 .....	xxi
VCCI: A 类 .....	xxi
产品安全符合性 .....	xxii

## 1

### 产品概览

功能说明 .....	1
功能 .....	1
适配器规格 .....	3
物理特征 .....	3
标准规格 .....	3

## 2

### 硬件安装

系统要求 .....	4
安全预防措施 .....	5
预安装核查清单 .....	5
安装适配器 .....	6

## 3

### 驱动程序安装

安装 Linux 驱动程序软件 .....	7
安装不含 RDMA 的 Linux 驱动程序 .....	9
移除 Linux 驱动程序 .....	9
使用源代码 RPM 软件包安装 Linux 驱动程序 .....	11
使用 kmp/kmod RPM 软件包安装 Linux 驱动程序 .....	12
使用 TAR 文件安装 Linux 驱动程序 .....	12
安装具有 RDMA 的 Linux 驱动程序 .....	13

	Linux 驱动程序可选参数 .....	14
	Linux 驱动程序操作默认设置 .....	14
	Linux 驱动程序消息 .....	15
	统计信息 .....	15
安装	Windows 驱动程序软件 .....	15
	安装 Windows 驱动程序 .....	15
	在 GUI 中运行 DUP .....	16
	DUP 安装选项 .....	22
	DUP 安装示例 .....	23
	移除 Windows 驱动程序 .....	23
	管理适配器属性 .....	24
	设置电源管理选项 .....	25
安装	VMware 驱动程序软件 .....	26
	VMware 驱动程序和驱动程序包 .....	26
	安装 VMware 驱动程序 .....	27
	VMware 驱动程序可选参数 .....	28
	VMware 驱动程序参数默认值 .....	30
	移除 VMware 驱动程序 .....	30
	FCoE 支持 .....	31
	iSCSI 支持 .....	31
<b>4</b>	<b>升级固件</b>	
	通过双击运行 DUP .....	32
	从命令行运行 DUP .....	34
	使用 .bin 文件运行 DUP .....	35
<b>5</b>	<b>适配器预引导配置</b>	
	使用入门 .....	38
	显示固件映像属性 .....	41
	配置设备级参数 .....	42
	配置 NIC 参数 .....	43
	配置数据中心桥接 .....	47
	配置 FCoE 引导 .....	48
	配置 iSCSI 引导 .....	49
	配置分区 .....	54
	VMware ESXi 6.0 和 ESXi 6.5 分区 .....	59
<b>6</b>	<b>RoCE 配置</b>	
	支持的操作系统和 OFED .....	60
	规划 RoCE .....	61
	准备适配器 .....	62

准备以太网交换机.....	62
配置 Cisco Nexus 6000 以太网交换机.....	62
配置 Dell Z9100 以太网交换机.....	64
在 Windows Server 的适配器上配置 RoCE.....	64
查看 RDMA 计数器.....	67
在 Linux 的适配器上配置 RoCE.....	71
RHEL 的 RoCE 配置.....	72
SLES 的 RoCE 配置.....	72
验证 Linux 上的 RoCE 配置.....	73
VLAN 接口和 GID 索引值.....	75
Linux 的 RoCE v2 配置.....	76
识别 RoCE v2 GID 索引或地址.....	77
通过 sys 和 class 参数验证 RoCE v1 或 v2 GID 索引和地址...	77
通过 perfest 应用程序验证 RoCE v1 或 RoCE v2 功能.....	78
在 VMware ESX 的适配器上配置 RoCE.....	81
配置 RDMA 接口.....	82
配置 MTU.....	83
RoCE 模式和统计信息.....	83
配置半虚拟 RDMA 设备 (PVRDMA).....	85
<b>7</b>	
<b>iWARP 配置</b>	
为 iWARP 准备适配器.....	88
在 Windows 上配置 iWARP.....	89
在 Linux 上配置 iWARP.....	92
安装驱动程序.....	93
配置 iWARP 和 RoCE.....	93
检测设备.....	94
支持的 iWARP 应用程序.....	95
为 iWARP 运行 Perftest.....	95
配置 NFS-RDMA.....	96
SLES 12 SP3、RHEL 7.4 和 OFED 4.8x 上的 iWARP RDMA-Core 支持.....	97
<b>8</b>	
<b>iSER 配置</b>	
准备工作.....	100
为 RHEL 配置 iSER.....	100
为 SLES 12 配置 iSER.....	104
在 RHEL 和 SLES 上通过 iWARP 使用 iSER.....	105
优化 Linux 性能.....	106
将 CPU 配置为最高性能模式.....	106

配置内核 sysctl 设置 .....	107
配置 IRQ 关联设置 .....	107
配置块设备暂存 .....	107
在 ESXi 6.7 上配置 iSER .....	108
准备工作 .....	108
为 ESXi 6.7 配置 iSER .....	108
<b>9 iSCSI 配置</b>	
iSCSI 引导 .....	111
iSCSI 引导设置 .....	112
选择首选的 iSCSI 引导模式 .....	112
配置 iSCSI 目标 .....	113
配置 iSCSI 引导参数 .....	113
适配器 UEFI 引导模式配置 .....	115
配置 iSCSI 引导 .....	118
静态 iSCSI 引导配置 .....	118
动态 iSCSI 引导配置 .....	126
启用 CHAP 身份验证 .....	128
配置 DHCP 服务器以支持 iSCSI 引导 .....	129
IPv4 的 DHCP iSCSI 引导配置 .....	129
DHCP 选项 17, 根路径 .....	129
DHCP 选项 43, 供应商特定信息 .....	130
配置 DHCP 服务器 .....	131
为 IPv6 配置 DHCP iSCSI 引导 .....	132
DHCPv6 选项 16, 供应商类别选项 .....	132
DHCPv6 选项 17, 供应商特定信息 .....	132
配置 VLAN 用于 iSCSI 引导 .....	133
Windows Server 中的 iSCSI 卸载 .....	133
安装 QLogic 驱动程序 .....	134
安装 Microsoft iSCSI 启动器 .....	134
配置 Microsoft 启动器以使用 QLogic 的 iSCSI 卸载 .....	134
iSCSI 卸载常见问题 .....	140
Windows Server 2012 R2 和 2016 iSCSI 引导安装 .....	141
iSCSI 故障转储 .....	142
Linux 环境中的 iSCSI 卸载 .....	142
与 bnx2i 的差异 .....	142
配置 qedi.ko .....	143
在 Linux 中验证 iSCSI 接口 .....	143

	Open-iSCSI 和从 SAN 引导注意事项 .....	145
	从 SAN 迁移的 RHEL 6.9 iSCSI L4 引导 .....	146
	从 SAN 迁移的 RHEL 7.2/7.3 iSCSI L4 引导 .....	149
	从 SAN 迁移的 SLES 11 SP4 iSCSI L4 引导 .....	151
	从 SAN 迁移的 SLES 12 SP1/SP2 iSCSI L4 引导 .....	152
	使用 MPIO 从 SAN 迁移的 SLES 12 SP1/SP2 iSCSI L4 引导 ..	154
	从 RHEL 7.4 及更高版本的 SAN 配置 iSCSI 引导 .....	156
10	<b>FCoE 配置</b>	
	从 SAN 进行 FCoE 引导 .....	160
	为 FCoE 构建和引导来准备系统 BIOS .....	161
	指定 BIOS 引导协议 .....	161
	配置适配器 UEFI 引导模式 .....	161
	Windows 版从 SAN 进行的 FCoE 引导 .....	166
	Windows Server 2012 R2 和 2016 FCoE 引导安装 .....	166
	配置 FCoE .....	167
	FCoE 故障转储 .....	167
	将适配器驱动程序注入（滑流至）Windows 映像文件中 .....	167
	配置 Linux FCoE 卸载 .....	168
	qedf 与 bnx2fc 之间的差异 .....	169
	配置 qedf.ko .....	170
	在 Linux 中验证 FCoE 设备 .....	170
	从 RHEL 7.4 及更高版本的 SAN 配置 FCoE 引导 .....	171
11	<b>SR-IOV 配置</b>	
	在 Windows 上配置 SR-IOV .....	175
	在 Linux 上配置 SR-IOV .....	182
	在 VMware 上配置 SR-IOV .....	188
12	<b>使用 RDMA 进行 NVMe-oF 配置</b>	
	在两台服务器上安装设备驱动程序 .....	195
	配置目标服务器 .....	196
	配置启动器服务器 .....	197
	预处理目标服务器 .....	199
	测试 NVMe-oF 设备 .....	199
	优化性能 .....	201
	.IRQ 关联 (multi_rss-affin.sh) .....	202
	CPU 频率 (cpufreq.sh) .....	203

13	Windows Server 2016	
	使用 Hyper-V 配置 RoCE 接口	204
	创建带 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机	205
	将 VLAN ID 添加到主机虚拟 NIC	207
	验证 RoCE 是否启用	207
	添加主机虚拟 NIC（虚拟端口）	208
	映射 SMB 驱动器和运行 RoCE 流量	208
	Switch Embedded Teaming 上的 RoCE	210
	创建带 SET 和 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机	210
	在 SET 上启用 RDMA	211
	在 SET 上分配 VLAN ID	211
	在 SET 上运行 RDMA 流量	211
	为 RoCE 配置 QoS	212
	通过在适配器上禁用 DCBX 来配置 QoS	212
	通过在适配器上启用 DCBX 来配置 QoS	216
	配置 VMMQ	220
	在适配器上启用 VMMQ	220
	设置 VMMQ 最大 QP 默认和非默认 VPort	221
	创建带或不带 SR-IOV 的虚拟机交换机	222
	在虚拟机交换机上启用 VMMQ	224
	获取虚拟机交换机功能	224
	创建 VM 并在 VM 中的 VMNetworkadapter 上启用 VMMQ	225
	默认和最大 VMMQ 虚拟 NIC	226
	在管理 NIC 上启用和禁用 VMMQ	226
	监测流量统计信息	226
	配置 VXLAN	227
	在适配器上启用 VXLAN 卸载	227
	部署软件定义网络	228
	配置 Storage Spaces Direct	228
	配置硬件	228
	部署超聚合系统	229
	部署操作系统	229
	配置网络	229
	配置 Storage Spaces Direct	231
	部署和管理 Nano Server	234
	角色和功能	234
	在物理服务器上部署 Nano Server	236
	在虚拟机中部署 Nano Server	238

	远程管理 Nano Server . . . . .	240
	使用 Windows PowerShell 远程处理管理 Nano Server . . . . .	240
	将 Nano Server 添加到受信任的主机列表 . . . . .	240
	启动远程 Windows PowerShell 会话 . . . . .	241
	在 Windows Nano Server 上管理 QLogic 适配器 . . . . .	241
	RoCE 配置 . . . . .	241
14	<b>故障排除</b>	
	故障排除核查清单 . . . . .	245
	验证是否已加载最新驱动程序 . . . . .	246
	在 Windows 中验证驱动程序 . . . . .	246
	在 Linux 中验证驱动程序 . . . . .	246
	在 Vmware 中验证驱动程序 . . . . .	247
	测试网络连通性 . . . . .	247
	测试 Windows 的网络连通性 . . . . .	247
	测试 Linux 的网络连通性 . . . . .	248
	使用 Hyper-V 的 Microsoft Virtualization . . . . .	248
	Linux 特定问题 . . . . .	248
	其他问题 . . . . .	248
	收集调试数据 . . . . .	249
A	<b>适配器 LED</b>	
B	<b>电缆和光学模块</b>	
	支持的规格 . . . . .	251
	经测试的电缆和光学模块 . . . . .	252
	经测试的交换机 . . . . .	256
C	<b>Dell Z9100 交换机配置</b>	
D	<b>功能约束</b>	
	<b>词汇表</b>	

## 图列表

图		页
3-1	Dell Update Package 窗口	16
3-2	QLogic InstallShield 向导: Welcome (欢迎) 窗口	17
3-3	QLogic InstallShield 向导: License Agreement (许可协议) 窗口	18
3-4	InstallShield 向导: Setup Type (设置类型) 窗口	19
3-5	InstallShield 向导: 自定义设置窗口	20
3-6	InstallShield 向导: Ready to Install the Program (准备安装程序) 窗口	20
3-7	InstallShield 向导: Completed (完成) 窗口	21
3-8	Dell Update Package 窗口	22
3-9	设置高级适配器属性	24
3-10	电源管理选项	25
4-1	Dell Update Package: 初始屏幕	32
4-2	Dell Update Package: 加载新固件	33
4-3	Dell Update Package: 安装结果	33
4-4	Dell Update Package: 完成安装	34
4-5	DUP 命令行选项	35
5-1	系统设置	38
5-2	系统设置: 设备设置	38
5-3	主要配置页面	39
5-4	主要配置页面, 将分区模式设置为 NPAR	39
5-5	固件映像属性	42
5-6	设备级别配置	42
5-7	NIC 配置	44
5-8	系统设置: 数据中心桥接 (DCB) 设置	47
5-9	FCoE 常规参数	48
5-10	FCoE 目标配置	49
5-11	iSCSI 常规参数	51
5-12	iSCSI 启动器配置参数	52
5-13	iSCSI 第一目标参数	52
5-14	iSCSI 第二目标参数	53
5-15	NIC 分区配置, 全局带宽分配	54
5-16	全局带宽分配页面	55
5-17	分区 1 配置	56
5-18	分区 2 配置: FCoE 卸载	57
5-19	分区 3 配置: iSCSI 卸载	58
5-20	分区 4 配置: 以太网	58
6-1	配置 RoCE 属性	65
6-2	添加计数器对话框	67
6-3	性能监视器: Cavium FastLinQ 计数器	69
6-4	交换机设置, 服务器	80
6-5	交换机设置, 客户端	80
6-6	配置 RDMA_CM 应用程序: 服务器	81
6-7	配置 RDMA_CM 应用程序: 客户端	81
6-8	配置新的分布式交换机	85

6-9	为 PVRDMA 分配 vmknic	86
6-10	设置防火墙规则	87
7-1	Windows PowerShell 命令: Get-NetAdapterRdma	90
7-2	Windows PowerShell 命令: Get-NetOffloadGlobalSetting	90
7-3	Perfmon: 添加计数器	91
7-4	Perfmon: 验证 iWARP 流量	91
8-1	RDMA Ping 操作成功	101
8-2	iSER 门户实例	102
8-3	iface 传输确认	103
8-4	检查新 iSCSI 设备	103
8-5	LIO 目标配置	105
9-1	系统设置: NIC 配置	112
9-2	系统设置: 引导设置	115
9-3	系统设置: 设备设置配置实用程序	116
9-4	选择 NIC 配置	117
9-5	系统设置: NIC 配置、引导协议	118
9-6	系统设置: iSCSI 配置	119
9-7	系统设置: 选择常规参数	119
9-8	系统设置: iSCSI 常规参数	120
9-9	系统设置: 选择 iSCSI 启动器参数	121
9-10	系统设置: iSCSI 启动器参数	122
9-11	系统设置: 选择 iSCSI 第一个目标参数	123
9-12	系统设置: iSCSI 第一个目标参数	124
9-13	系统设置: iSCSI 第二个目标参数	125
9-14	系统设置: 保存 iSCSI 更改	126
9-15	系统设置: iSCSI 常规参数	128
9-16	系统设置: iSCSI 常规参数, VLANID	133
9-17	iSCSI 启动器属性, 配置页面	135
9-18	iSCSI 启动器节点名称更改	135
9-19	iSCSI 启动器 - 查找目标门户	136
9-20	目标门户 IP 地址	137
9-21	选择启动器 IP 地址	138
9-22	连接到 iSCSI 目标	139
9-23	连接到目标对话框	140
9-24	提示开箱即用安装	157
9-25	Red Hat Enterprise Linux 7.4 配置	158
10-1	系统设置: 选择设备设置	161
10-2	系统设置: 设备设置, 端口选择	162
10-3	系统设置: NIC 配置	163
10-4	系统设置: FCoE 模式已启用	164
10-5	系统设置: FCoE 常规参数	165
10-6	系统设置: FCoE 常规参数	166
10-7	提示开箱即用安装	173
10-8	Red Hat Enterprise Linux 7.4 配置	174
11-1	SR-IOV 的系统设置: 集成的设备	176

11-2	SR-IOV 的系统设置：设备级配置 .....	176
11-3	适配器属性，高级：启用 SR-IOV .....	177
11-4	虚拟交换机管理器：启用 SR-IOV .....	178
11-5	VM 的设置：启用 SR-IOV .....	180
11-6	设备管理器：带有 QLogic 适配器的 VM .....	181
11-7	Windows PowerShell 命令：Get-NetadapterSriovVf .....	181
11-8	系统设置：SR-IOV 的处理器设置 .....	183
11-9	SR-IOV 的系统设置：集成的设备 .....	184
11-10	为 SR-IOV 编辑 grub.conf 文件 .....	185
11-11	sriov_numvfs 的命令输出 .....	186
11-12	ip link show 命令的命令输出 .....	186
11-13	RHEL68 虚拟机 .....	187
11-14	添加新虚拟硬件 .....	188
11-15	VMware 主机编辑设置 .....	191
12-1	NVMe-oF 网络 .....	194
12-2	子系统 NQN .....	198
12-3	确认 NVMe-oF 连接 .....	199
12-4	FIO 公用程序安装 .....	200
13-1	在主机虚拟 NIC 中启用 RDMA .....	205
13-2	Hyper-V 虚拟以太网适配器属性 .....	206
13-3	Windows PowerShell 命令：Get-VMNetworkAdapter .....	207
13-4	Windows PowerShell 命令：Get-NetAdapterRdma .....	207
13-5	添加计数器对话框 .....	209
13-6	性能监视器显示 RoCE 流量 .....	209
13-7	Windows PowerShell 命令：New-VMSwitch .....	210
13-8	Windows PowerShell 命令：Get-NetAdapter .....	211
13-9	高级属性：启用 QoS .....	213
13-10	高级属性：设置 VLAN ID .....	214
13-11	高级属性：启用 QoS .....	217
13-12	高级属性：设置 VLAN ID .....	218
13-13	高级属性：启用虚拟交换机 RSS .....	221
13-14	高级属性：设置 VMMQ .....	222
13-15	虚拟交换机管理器 .....	223
13-16	Windows PowerShell 命令：Get-VMSwitch .....	224
13-17	高级属性：启用 VXLAN .....	227
13-18	示例硬件配置 .....	228
13-19	Windows PowerShell 命令：Get-NetAdapter .....	242
13-20	Windows PowerShell 命令：Get-NetAdapterRdma .....	242
13-21	Windows PowerShell 命令：New-Item .....	243
13-22	Windows PowerShell 命令：New-SMBShare .....	243
13-23	Windows PowerShell 命令：Get-NetAdapterStatistics .....	244

## 表格列表

表		页
2-1	主机硬件要求	4
2-2	最低主机操作系统要求	5
3-1	QLogic 41xxx 系列适配器 Linux 驱动程序	7
3-2	qed 驱动程序可选参数	14
3-3	Linux 驱动程序操作默认设置	14
3-4	VMware 驱动程序	26
3-5	按版本的 ESXi 驱动程序包	26
3-6	VMware 驱动程序可选参数	28
3-7	VMware 驱动程序参数默认值	30
3-8	QLogic 41xxx 系列适配器 VMware FCoE 驱动程序	31
3-9	QLogic 41xxx 系列适配器 iSCSI 驱动程序	31
5-1	适配器属性	40
6-1	RoCE v1、RoCE v2、iWARP 和 OFED 的操作系统支持	60
6-2	RoCE 的高级属性	64
6-3	Cavium FastLinQ RDMA 错误计数器	69
9-1	配置选项	114
9-2	DHCP 选项 17 参数定义	129
9-3	DHCP 选项 43 子选项定义	131
9-4	DHCP 选项 17 子选项定义	132
12-1	目标参数	196
13-1	Nano Server 的角色和功能	234
14-1	收集调试数据命令	249
A-1	适配器端口链路和活动 LED	250
B-1	经测试的电缆和光学模块	252
B-2	经测试互操作性的交换机	256

# 前言

本前言列出了支持的产品，指定了目标读者，说明了本指南中使用的排版惯例，并介绍了法律声明。

## 支持的产品

本用户指南介绍以下 Cavium™ 产品：

- QL41112HFCU-DE 10Gb 聚合网络适配器，全高支架
- QL41112HLCU-DE 10Gb 聚合网络适配器，薄型支架
- QL41162HFRJ-DE 10Gb 聚合网络适配器，全高支架
- QL41162HLRJ-DE 10Gb 聚合网络适配器，薄型支架
- QL41162HMRJ-DE 10Gb 聚合网络适配器
- QL41164HMCU-DE 10Gb 聚合网络适配器
- QL41164HMRJ-DE 10Gb 聚合网络适配器
- QL41262HFCU-DE 10/25Gb 聚合网络适配器，全高支架
- QL41262HLCU-DE 10/25Gb 聚合网络适配器，薄型支架
- QL41262HMCU-DE 10/25Gb 聚合网络适配器
- QL41264HMCU-DE 10/25Gb 聚合网络适配器

## 目标读者

本指南面向负责对安装在 Dell® PowerEdge® 服务器（在 Windows®、Linux® 或 VMware® 环境中）上的适配器进行配置和管理的系统管理员和其他技术人员。

## 本指南的内容

前言之后，本指南的其余部分划分为以下章节和附录：

- [第 1 章 产品概览](#) 提供产品的功能说明、功能列表以及适配器规格。
- [第 2 章 硬件安装](#) 介绍如何安装适配器，包括系统要求列表和预安装核查表。

- [第 3 章 驱动程序安装](#)介绍在 Windows、Linux 和 VMware 中安装适配器驱动程序。
- [第 4 章 升级固件](#)介绍如何使用 Dell Update Package (DUP) 升级适配器固件。
- [第 5 章 适配器预引导配置](#)介绍如何使用人机界面基础设施 (HII) 应用程序执行预引导适配器配置任务。
- [第 6 章 RoCE 配置](#)介绍如何配置适配器、以太网交换机和主机以使用基于聚合以太网的 RDMA (RoCE)。
- [第 7 章 iWARP 配置](#)提供在 Windows、Linux 和 VMware ESXi 6.7 系统中配置互联网广域 RDMA 协议 (iWARP) 的步骤。。
- [第 8 章 iSER 配置](#)介绍如何为 Linux RHEL 和 SLES 配置 RDMA 的 iSCSI 扩展 (iSER)。
- [第 9 章 iSCSI 配置](#)介绍适用于 Windows 和 Linux 的 iSCSI 引导、iSCSI 故障转储和 iSCSI 卸载。
- [第 10 章 FCoE 配置](#)介绍从 SAN 的以太网光纤信道 (FCoE) 引导以及安装后从 SAN 引导。
- [第 11 章 SR-IOV 配置](#)提供在 Windows、Linux 和 VMware 系统中配置单根输入 / 输出虚拟化 (SR-IOV) 的步骤。
- [第 12 章 使用 RDMA 进行 NVMe-oF 配置](#) 演示如何在简单的网络上配置 NVMe-oF。
- [第 13 章 Windows Server 2016](#) 介绍 Windows Server 2016 功能。
- [第 14 章 故障排除](#)介绍各种故障排除方法和资源。
- [附录 A 适配器 LED](#) 列出适配器 LED 及其意义。
- [附录 B 电缆和光学模块](#)列出 41xxx 系列适配器支持的电缆和光学模块。
- [附录 C Dell Z9100 交换机配置](#)介绍如何配置 Dell Z9100 交换机端口以建立 25Gbps 连接。
- [附录 D 功能约束](#)提供有关当前版本中所实现功能约束的信息。

本指南的最后部分是术语表。

## 文档惯例

本指南使用以下文档惯例：

- **注** 提供额外的信息。
- **小心** 不带警报符号，表示存在可能导致设备损坏或数据丢失的危险。
-  **小心** 带警报符号，表示存在可能造成轻度或中度伤害的危险。
-  **警告** 表示存在可能造成严重伤害或死亡的危险。
- 蓝色字体的文字表示至本指南中的插图、表格或章节的超链接（跳转），至网站的链接以下划线蓝色文字显示。例如：
  - 表 9-2 列出与用户界面和远程代理有关的问题。
  - 请参阅第 6 页上的“安装核查表”。
  - 有关更多信息，请访问 [www.cavium.com](http://www.cavium.com)。
- 黑体文字表示用户界面元素，如菜单项、按钮、复选框或列标题。例如：
  - 单击**开始按钮**，指向**程序**，指向**附件**，然后单击**命令提示符**。
  - 在**通知选项**下，选中**警报复选框**。
- Courier 字体文本表示文件名、目录路径或命令行文字。例如：
  - 要从文件结构的任何地方返回根目录：  
键入 `cd /root` 并按 ENTER 键。
  - 发出以下命令：`sh ./install.bin`。
- 键盘的键名和击键用大写字母表示：
  - 按 CTRL+P。
  - 按向上箭头键。
- 斜体文字表示术语、强调、变量或说明文件标题。例如：
  - 哪些是 *快捷键*？
  - 要输入日期键入 `mm/dd/yyyy`（其中 `mm` 是月，`dd` 是日，`yyyy` 是年）。
- 引号中的主题标题表示本手册中的联机帮助的相关主题；在本文档中，联机帮助又称为 *帮助系统*。

- 命令行界面 (CLI) 命令语法惯例包括以下内容：
  - 纯文本表示必须按照所示键入的项目。例如：
    - `qaucli -pr nic -ei`
  - `< >` (尖括号) 表示必须指定其值的变量。例如：
    - `<serial_number>`

---

**注**

仅适用于 CLI 命令，变量名称始终使用尖括号而不是斜体表示。

---

- `[ ]` (方括号) 表示可选的参数。例如：
  - `[<file_name>]` 表示指定文件名，或省略以选择默认文件名。
- `|` (竖线) 表示互斥的选项；只能选择一个选项。例如：
  - `on|off`
  - `1|2|3|4`
- `...` (省略号) 表示前面的项可重复。例如：
  - `x...` 表示 `x` 的一个或多个实例。
  - `[x...]` 表示 `x` 的零个或多个实例。
- 命令示例输出中的垂直省略号指示重复输出数据部分被有意省略的位置。
- `( )` (圆括号) 和 `{ }` (大括号) 用来避免逻辑模糊不清。例如：
  - `a|b c` 模糊不清
  - `{(a|b) c}` 表示 `a` 或 `b`，后跟 `c`
  - `{a|(b c)}` 表示 `a` 或 `b c`

## 法律声明

本节中涵盖的法律声明包括激光安全（FDA 公告）、机构认证和产品安全符合性。

### 激光安全 - FDA 公告

本产品符合 DHHS 规则 21CFR I 章 J 节的规定。本产品的设计和和生产符合 IEC60825-1 中有关激光产品安全标签的规定。

I 类激光产品

1 类 激光产品	<b>小心</b> — 在打开时存在 1 类激光辐射 请勿直视光学仪器
Appareil laser de classe 1	<b>Attention</b> —Radiation laser de classe 1 Ne pas regarder directement avec des instruments optiques
Produkt der Laser Klasse 1	<b>Vorsicht</b> —Laserstrahlung der Klasse 1 bei geöffneter Abdeckung Direktes Ansehen mit optischen Instrumenten vermeiden
Luokan 1 Laserlaite	<b>Varoitus</b> —Luokan 1 lasersäteilyä, kun laite on auki Älä katso suoraan laitteeseen käyttämällä optisia instrumenttejä

### 机构认证

以下章节概述对 41xxx 系列适配器执行的 EMC 和 EMI 检验规格，验证其是否符合辐射排放、抗辐射干扰性以及产品安全标准。

#### EMI 和 EMC 要求

##### FCC 第 15 部分符合性：A 级

**FCC 符合性信息声明：** 本设备符合 FCC 规则第 15 部分的规定。操作受限于以下两个条件：(1) 此设备不得造成有害干扰；(2) 此设备必须能够承受接收到的任何干扰，包括可能导致不良操作的干扰。

##### ICES-003 符合性：A 级

本 A 级数字装置符合加拿大 ICES-003 的规定。Cet appareil numérique de la classe A est conforme à la norme NMB-003 du Canada.

##### CE 标记 2014/30/EC、2014/35/EU EMC 指令符合性：

EN55032:2012/ CISPR 32:2015 A 类

EN55024:2010

EN61000-3-2: 谐波电流发射

EN61000-3-3: 电压波动和闪动

抗扰性标准  
EN61000-4-2: ESD  
EN61000-4-3: 射频电磁场  
EN61000-4-4: 快速瞬变 / 爆发  
EN61000-4-5: 快速电涌常见 / 差动  
EN61000-4-6: 射频传导敏感度  
EN61000-4-8: 电力频率电磁场  
EN61000-4-11: 电压骤降和中断

**VCCI: 2015-04 ; A 级**

**AS/NZS ; CISPR 32: 2015 A 类**

**CNS 13438: 2006 A 类**

## KCC: A 级

经过韩国 RRA A 级认证



产品名称 / 型号: 聚合网络适配器和智能以太网适配器  
证书持有者: QLogic Corporation  
制造日期: 请参阅产品上的日期代码  
制造商 / 原产国: QLogic Corporation/ 美国

A 级设备  
(商用信息 / 电讯设备)

由于本设备已就其商用性执行了 EMC 注册, 因此要求销售方和 / 或购买方对此引起充分注意, 倘若发生错误的销售或购买行为, 要求立即将其更换成家用型设备。

韩国语言格式 - A 级

### 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: A 类

根据 Voluntary Control Council for Interference (VCCI) 的标准, 该设备是 A 类产品。如果在家庭环境中使用该设备, 可能会发生无线电干扰, 在这种情况下, 用户可能需要采取纠正措施。

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

## 产品安全符合性

### UL、cUL 产品安全:

UL 60950-1 (第二版) A1 + A2 2014-10-14

CSA C22.2 No.60950-1-07 (第二版) A1 +A2 2014-10

只能用于所列 ITE 或等价对象。

符合 21 CFR 1040.10 和 1040.11、2014/30/EU、2014/35/EU。

### 2006/95/EC 低电压规程:

TUV EN60950-1:2006+A11+A1+A12+A2 第二版

TUV IEC 60950-1: 2005 第二版 Am1: 2009 + Am2: 2013 CB

根据 IEC 60950-1 第二版经过 CB 认证

# 1 产品概览

本章提供 41xxx 系列适配器的以下信息：

- 功能说明
- 功能
- 第 3 页上的“适配器规格”

## 功能说明

Cavium FastLinQ® 41000 系列适配器包括 10 和 25Gb 聚合网络适配器以及智能以太网适配器，旨在为服务器系统执行加速的数据网络。41000 系列适配器包括一个 10/25Gb 以太网 MAC（具有全双工能力）。

利用操作系统的组合功能，可将网络分割成虚拟 LAN (VLAN)，以及将多个网络适配器组合到各个组中，以便提供网络负载平衡和容错。有关组合的更多信息，请参阅操作系统说明文件。

## 功能

41xxx 系列适配器提供以下功能。并非所有适配器均具备所有功能：

- NIC 分区 (NPAR)
- 单芯片解决方案：
  - 10Gb/25Gb MAC
  - 用于直接附加铜缆 (DAC) 收发器连接的 SerDes 接口
  - PCIe® 3.0 x8
  - 具备零拷贝功能的硬件
- 性能特点：
  - TCP、IP、UDP 校验和卸载
  - TCP 分段卸载 (TSO)
    - 大段卸载 (LSO)

- 普通分段卸载 (GSO)
- 大量接收卸载 (LRO)
- 接收分段结合 (RSC)
- Microsoft® 动态虚拟机队列 (VMQ) 和 Linux 多队列
- 自适应中断:
  - 发送方 / 接收方缩放 (TSS/RSS)
  - 使用通用路由封装 (NVGRE) 和虚拟 LAN (VXLAN) L2/L3 GRE 隧道流量进行网络虚拟化无状态卸载<sup>1</sup>
- 可管理性:
  - 系统管理总线 (SMB) 控制器
  - 符合高级配置与电源接口 (ACPI) 1.1a 标准 (多电源模式)
  - 网络控制器边带接口 (NC-SI) 支持
- 高级网络特性:
  - 巨型帧 (多达 9,600 个字节)。操作系统和链路伙伴必须支持巨型帧。
  - 虚拟 LAN (VLAN)
  - 流控制 (IEEE Std 802.3x)
- 逻辑链路控制 (IEEE 标准 802.2)
- 高速芯片搭载精简指令集计算机 (RISC) 处理器
- 集成式 96KB 帧缓冲区存储区 (并非适用于所有型号)
- 1,024 分类过滤器 (并非适用于所有型号)
- 通过 128 位散列硬件功能支持多播地址
- 串行 NVRAM 闪存
- PCI 电源管理接口 (v1.1)
- 64 位基本地址寄存器 (BAR) 支持
- EM64T 处理器支持
- iSCSI 和 FCoE 引导支持<sup>2</sup>

---

<sup>1</sup> 此功能要求操作系统或虚拟机监控程序支持才能使用卸载。

<sup>2</sup> SR-IOV VF 的硬件支持限制有所不同。该限制在某些操作系统环境中可能较低；请参阅您的操作系统的相应章节。

## 适配器规格

41xxx 系列适配器规格包括适配器的物理特征和标准符合性引用。

### 物理特征

41xxx 系列适配器 是标准的 PCIe 卡，并附带一个全高或薄型支架以在标准 PCIe 插槽中使用。

### 标准规格

支持的标准规格包括：

- *PCI Express 基本规格*，修订版 3.1
- *PCI Express 卡机电规格*，修订版 3.0
- *PCI 总线电源管理接口规格*，修订版 1.2
- IEEE 规格：
  - 以太网的 802.3-2015 IEEE 标准*（流控制）
  - 802.1q (VLAN)*
  - 802.1AX*（链路聚合）
  - 802.1ad (QinQ)*
  - 802.1p*（优先级编码）
  - 1588-2002 PTPv1*（精确时间协议）
  - 1588-2008 PTPv2*
  - IEEE 802.3az 节能以太网 (EEE)*
- IPv4 (RFQ 791)
- IPv6 (RFC 2460)

# 2 硬件安装

本章提供以下硬件安装信息：

- [系统要求](#)
- [第 5 页上的“安全预防措施”](#)
- [第 5 页上的“预安装核查清单”](#)
- [第 6 页上的“安装适配器”](#)

## 系统要求

在您安装 Cavium 41xxx 系列适配器，之前，验证您的系统满足 [表 2-1](#) 和 [表 2-2](#) 中所示的硬件和操作系统的要求。支持的操作系统完整列表，请访问 [Cavium 网站](#)。

**表 2-1. 主机硬件要求**

硬件	要求
架构	可满足操作系统要求的 IA-32 或 EMT64
PCIe	PCIe Gen2 x8 (2x10G NIC) PCIe Gen3 x8 (2x25G NIC) PCIe Gen3 x8 或更快的插槽支持全双端口 25Gb 带宽。
内存	8GB RAM（最少）
电缆和光学模块	已经测试了 41xxx 系列适配器 与各种 1G、10G 和 25G 电缆和光学模块的互操作性。请参阅 <a href="#">第 252 页上的“经测试的电缆和光学模块”</a> 。

表 2-2. 最低主机操作系统要求

操作系统	要求
Windows Server	2012、2012 R2、2016（包括 Nano）
Linux	RHEL® 6.8、6.9、7.2、7.3、7.4 SLES® 11 SP4、SLES 12 SP2、SLES 12 SP3
VMware	用于 25G 适配器的 ESXi 6.0 u3 及更高版本

### 注

表 2-2 说明最低主机操作系统要求。支持的操作系统完整列表，请访问 Cavium 网站。

## 安全预防措施

### 警告

安装适配器的系统的操作电压可能会有致命危险。打开系统外壳之前，请遵从以下预防措施以保护您自己并避免损坏系统组件。

- 除去手上和手腕上的任何金属物体或首饰。
- 确保仅使用绝缘工具或非导电工具。
- 在触摸内部组件之前，确认系统电源已关闭并已拔掉电源线。
- 在不受静电干扰的环境中安装或卸下适配器。使用正确接地的腕带或其他人体防静电设备，强烈建议使用防静电地垫。

## 预安装核查清单

安装适配器之前，请完成以下操作：

1. 确认您的系统满足第 4 页上的“系统要求”中列出的硬件和软件要求。
2. 确认您的系统使用最新的 BIOS。

### 注

如果您从 Cavium 网站获取适配器软件，请验证适配器驱动程序文件的路径。

3. 如果系统正在运行，请将其关闭。

4. 系统关闭后，断开电源并拔下电源线。
5. 将适配器从其运输包装中取出并放在防静电表面上。
6. 检查适配器，特别是边缘连接器上是否有明显的损坏痕迹。切勿尝试安装损坏的适配器。

## 安装适配器

以下说明适用于在大多数系统中安装 Cavium 41xxx 系列适配器。有关执行这些任务的详细信息，请参阅随系统提供的手册。

### 要安装适配器：

1. 复查 [第 5 页上的“安全预防措施”](#) 和 [第 5 页上的“预安装核查清单”](#)。安装适配器前，确保系统电源已关闭而且电源线已从电源插座上拔下，并且遵守适当的电接地步骤。
2. 打开系统机箱，然后选择符合适配器大小的插槽，可以是 PCIe Gen 2 x8 或 PCIe Gen 3 x8。较窄的适配器可插入更宽的插槽中（x8 的适配器可插入 x16 的插槽中），但较宽的适配器不能插入更窄的插槽中（x8 的适配器不能插入 x4 的插槽中）。如果不知道如何识别 PCIe 插槽，请参考系统说明文件。
3. 从选择的插槽卸下空挡板。
4. 将适配器的连接器边缘与系统中的 PCIe 连接器插槽对齐。
5. 在适配器卡的两个边角均匀施压以推进插卡，直至其牢固就位在插槽中。当适配器正确就位时，其端口连接器将与插槽开口处对齐，适配器面板将与系统机箱齐平。

---

### 小心

将插卡推进到位时不要过度用力，否则可能损坏系统或适配器。如果无法固定适配器，将其卸下，重新对齐，并再次尝试。

---

6. 使用适配器夹或螺丝固定适配器。
7. 合上系统机箱，并断开任何个人防静电设备。

# 3 驱动程序安装

本章提供有关驱动程序安装的以下信息：

- 安装 Linux 驱动程序软件
- 第 15 页上的“安装 Windows 驱动程序软件”
- 第 26 页上的“安装 VMware 驱动程序软件”

## 安装 Linux 驱动程序软件

本节介绍如何安装包含或不包含远程直接内存访问 (RDMA) 的 Linux 驱动程序。本节还介绍 Linux 驱动程序的可选参数、默认值、消息和统计信息。

- 安装不含 RDMA 的 Linux 驱动程序
- 安装具有 RDMA 的 Linux 驱动程序
- Linux 驱动程序可选参数
- Linux 驱动程序操作默认设置
- Linux 驱动程序消息
- 统计信息

Dell 支持页面上提供了 41xxx 系列适配器 Linux 驱动程序和支持说明文件：

[dell.support.com](http://dell.support.com)

表 3-1 介绍了 41xxx 系列适配器 Linux 驱动程序。

**表 3-1. QLogic 41xxx 系列适配器 Linux 驱动程序**

Linux 驱动程序	说明
qed	qed 核心驱动程序模块直接控制固件、处理中断并为协议特定驱动程序集提供低级别 API。qed 与 qede、qedr、qedi 和 qedf 驱动程序接合。Linux 核心模块管理所有 PCI 设备资源（寄存器、主机接口队列，等等）。qed 核心模块需要使用 Linux 内核版本 2.6.32 或更高版本。测试集中于 x86_64 架构。

表 3-1. QLogic 41xxx 系列适配器 Linux 驱动程序 (续)

Linux 驱动程序	说明
qede	适用于 41xxx 系列适配器的 Linux 以太网驱动程序。该驱动程序直接控制硬件，并负责代表 Linux 主机网络堆栈发送和接收以太网数据包。该驱动程序还代表其自身接收和处理设备中断（适用于 L2 网络）。qede 驱动程序需要使用 Linux 内核版本 2.6.32 或更高版本。测试集中于 x86_64 架构。
qedr	Linux 基于聚合以太网的 RDMA (RoCE) 驱动程序。此驱动程序在开放结构企业分布 (OFED™) 环境中与 qed 核心模块和 qede 以太网驱动程序配合使用。RDMA 用户空间应用程序还需要在服务器上安装 libqedr 用户库。
qedi	适用于 41xxx 系列适配器的 Linux iSCSI 卸载驱动程序。此驱动程序与 Open iSCSI 库一起使用。
qedf	适用于 41xxx 系列适配器的 Linux FCoE 卸载驱动程序。此驱动程序与 Open FCoE 库一起使用。

可以使用源 Red Hat® Package Manager (RPM) 包或 kmod RPM 包安装 Linux 驱动程序。RHEL RPM 包内容如下：

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

SLES 源和 kmp RPM 包内容如下：

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

以下内核模块 (kmod) RPM 会在运行 Xen 虚拟机监控程序的 SLES 主机上安装 Linux 驱动程序：

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

以下源 RPM 会在 RHEL 和 SLES 主机上安装 RDMA 库代码：

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

以下源代码 TAR BZip2 (BZ2) 压缩文件会在 RHEL 和 SLES 主机上安装 Linux 驱动程序：

- fastlinq-<version>.tar.bz2

### 注

对于通过 NFS、FTP 或 HTTP（使用网络引导盘）执行的网络安装，可能需要使用其中包含 qede 驱动程序的驱动程序盘。可以通过修改 makefile 和 make 环境来编译 Linux 引导驱动程序。

---

## 安装不含 RDMA 的 Linux 驱动程序

要安装不含 RDMA 的 Linux 驱动程序：

1. 从 Dell 下载 41xxx 系列适配器 Linux 驱动程序：  
[dell.support.com](http://dell.support.com)
2. 如第 9 页上的“移除 Linux 驱动程序”中所述，移除现有 Linux 驱动程序。
3. 使用以下方法之一安装新 Linux 驱动程序：
  - 使用源代码 RPM 软件包安装 Linux 驱动程序
  - 使用 kmp/kmod RPM 软件包安装 Linux 驱动程序
  - 使用 TAR 文件安装 Linux 驱动程序

## 移除 Linux 驱动程序

移除 Linux 驱动程序有两种步骤：一种用于非 RDMA 环境，另一种用于 RDMA 环境。选择符合您的环境的步骤。

**要在非 RDMA 环境中移除 Linux 驱动程序，请卸载并移除驱动程序：**

按照与原始安装方法和操作系统相关的步骤进行操作。

- 如果以前是使用 RPM 软件包安装的 Linux 驱动程序，请发出以下命令：

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

- 如果以前是使用 TAR 文件安装的 Linux 驱动程序，请发出以下命令：

```
rmmod qede  
rmmod qed  
depmod -a
```

- ❑ 对于 RHEL:  

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

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

#### 要在非 RDMA 环境中移除 Linux 驱动程序:

1. 要获取当前安装的驱动程序的路径, 请发出以下命令:  

```
modinfo <driver name>
```
2. 卸载并移除 Linux 驱动程序。
  - ❑ 如果以前是使用 RPM 软件包安装的 Linux 驱动程序, 请发出以下命令:  

```
modprobe -r qede
depmod -a
rpm -e qlgc-fastlinq-kmp-default-<version>.<arch>
```
  - ❑ 如果以前是使用 TAR 文件安装的 Linux 驱动程序, 请发出以下命令:  

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

---

#### 注

如果存在 qedr, 请发出 `modprobe -r qedr` 命令代替。

---

3. 从 qed.ko、qede.ko 和 qedr.ko 文件所在的目录中删除它们。例如, 在 SLES 中发出以下命令:  

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

#### 要在 RDMA 环境中移除 Linux 驱动程序:

1. 要获取已安装的驱动程序的路径, 请发出以下命令:  

```
modinfo <driver name>
```

2. 卸载并移除 Linux 驱动程序。

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

3. 移除驱动程序模块文件:

- 如果以前是使用 RPM 软件包安装的驱动程序, 请发出以下命令:

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

- 如果以前是使用 TAR 文件安装的驱动程序, 则针对您的操作系统发出以下命令:

对于 RHEL:

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

对于 SLES:

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

## 使用源代码 RPM 软件包安装 Linux 驱动程序

要使用源代码 RPM 软件包安装 Linux 驱动程序:

1. 在命令提示符处发出以下命令:

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

2. 将目录更改至 RPM 路径并针对内核构建二进制 RPM。

对于 RHEL:

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

对于 SLES:

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

3. 安装新编译的 RPM:

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

### 注

如果报告了冲突，则可能需要对某些 Linux 分发版使用 `--force` 选项。

---

驱动程序将安装在以下路径中。

对于 SLES:

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

对于 RHEL:

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

- 按如下方式开启所有 ethX 接口:  

```
ifconfig <ethX> up
```
- 对于 SLES，使用 YaST 来配置以太网接口以在引导时自动启动（通过设置静态 IP 地址或启用接口上的 DHCP 实现）。

## 使用 kmp/kmod RPM 软件包安装 Linux 驱动程序

要安装 kmod RPM 软件包:

- 在命令提示符处发出以下命令:  

```
rpm -ivh qlgc-fastling-<version>.<arch>.rpm
```
- 重新加载驱动程序:  

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

## 使用 TAR 文件安装 Linux 驱动程序

要使用 TAR 文件安装 Linux 驱动程序:

- 创建目录并将 TAR 文件解压缩到该目录:  

```
tar xjvf fastling-<version>.tar.bz2
```
- 切换到最近创建的目录，然后安装驱动程序:  

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

qed 和 qede 驱动程序将安装在以下路径中。

对于 SLES:

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

对于 RHEL:

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

3. 加载驱动程序进行测试（如有必要，先卸载现有驱动程序）:

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

## 安装具有 RDMA 的 Linux 驱动程序

有关 iWARP 的信息，请参阅 [第 7 章 iWARP 配置](#)。

**要在内建 OFED 环境中安装 Linux 驱动程序:**

1. 从 Dell 下载 41xxx 系列适配器 Linux 驱动程序:
2. 如 [第 71 页](#)上的“[在 Linux 的适配器上配置 RoCE](#)”中所述，在适配器上配置 RoCE。
3. 如 [第 9 页](#)上的“[移除 Linux 驱动程序](#)”中所述，移除现有 Linux 驱动程序。
4. 使用以下方法之一安装新 Linux 驱动程序:
  - 使用 [kmp/kmod RPM 软件包安装 Linux 驱动程序](#)
  - 使用 [TAR 文件安装 Linux 驱动程序](#)
5. 安装 libqedr 库以使用 RDMA 用户空间应用程序。libqedr RPM 仅可用于内建 OFED。您必须选择在 UEFI 中使用哪个 RDMA（RoCE、RoCEv2 或 iWARP），直到固件中支持共存的 RoCE+iWARP 功能）。默认情况下均未启用。发出以下命令:

```
rpm -ivh qlgc-libqedr-<version>.<arch>.rpm
```
6. 要构建和安装 libqedr 用户空间库，请发出以下命令:

```
'make libqedr_install'
```
7. 通过加载驱动程序来测试它们，如下所示:

```
modprobe qedr  
make install_libqedr
```

## Linux 驱动程序可选参数

表 3-2 介绍 qede 驱动程序的可选参数。

表 3-2. qede 驱动程序可选参数

参数	说明
debug	控制驱动程序详细级别，与 <code>ethtool -s &lt;dev&gt; msglvl</code> 类似。
int_mode	控制除 MSI-X 以外的中断模式。
gro_enable	启用或禁用硬件通用接收卸载 (GRO) 功能。此功能与内核的软件 GRO 类似，但只能通过设备硬件执行。
err_flags_override	在发生硬件错误情况下用于禁用或强制采取措施的位图： <ul style="list-style-type: none"><li>■ 位 #31 - 此位掩码的启用位</li><li>■ 位 #0 - 阻止重新断言硬件关注</li><li>■ 位 #1 - 收集调试数据</li><li>■ 位 #2 - 触发恢复过程</li><li>■ 位 #3 - 调用 WARN 以获取导致错误的流的调用跟踪</li></ul>

## Linux 驱动程序操作默认设置

表 3-3 列出了 qed 和 qede Linux 驱动程序操作默认设置。

表 3-3. Linux 驱动程序操作默认设置

操作	qed 驱动程序默认值	qede 驱动程序默认值
Speed (速度)	自动协商并广告速度	自动协商并广告速度
MSI/MSI-X	Enabled (已启用)	Enabled (已启用)
Flow Control (流控制)	—	自动协商并广告 RX 和 TX
MTU	—	1500 (范围为 46 - 9600)
Rx Ring Size (Rx 环大小)	—	1000
Tx Ring Size (Tx 环大小)	—	4078 (范围为 128 - 8191)
Coalesce Rx Microseconds (合并 RX 微秒)	—	24 (范围为 0 - 255)

表 3-3. Linux 驱动程序操作默认设置 (续)

操作	qed 驱动程序默认值	qede 驱动程序默认值
Coalesce Tx Microseconds (合并 Tx 微秒)	—	48
TSO	—	Enabled (已启用)

## Linux 驱动程序消息

要设置 Linux 驱动程序消息详细级别，请发出以下命令之一：

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

其中 <value> 表示 0-15 位，这些是标准的 Linux 网络值，并且 16 及更高的位特定于驱动程序。

## 统计信息

要查看详细的统计信息和配置信息，请使用 `ethtool` 公用程序。参见 `ethtool` 手册页了解更多信息。

## 安装 Windows 驱动程序软件

有关 iWARP 的信息，请参阅 [第 7 章 iWARP 配置](#)。

- [安装 Windows 驱动程序](#)
- [移除 Windows 驱动程序](#)
- [管理适配器属性](#)
- [设置电源管理选项](#)

## 安装 Windows 驱动程序

使用 Dell Update Package (DUP) 安装 Windows 驱动程序软件：

- [在 GUI 中运行 DUP](#)
- [DUP 安装选项](#)
- [DUP 安装示例](#)

## 在 GUI 中运行 DUP

### 要在 GUI 中运行 DUP:

1. 双击代表 Dell Update Package 文件的图标。

### 注

Dell Update Package 的实际文件名称将会不同。

2. 在 Dell Update Package 窗口 (图 3-1) 中, 单击 **Install** (安装)。

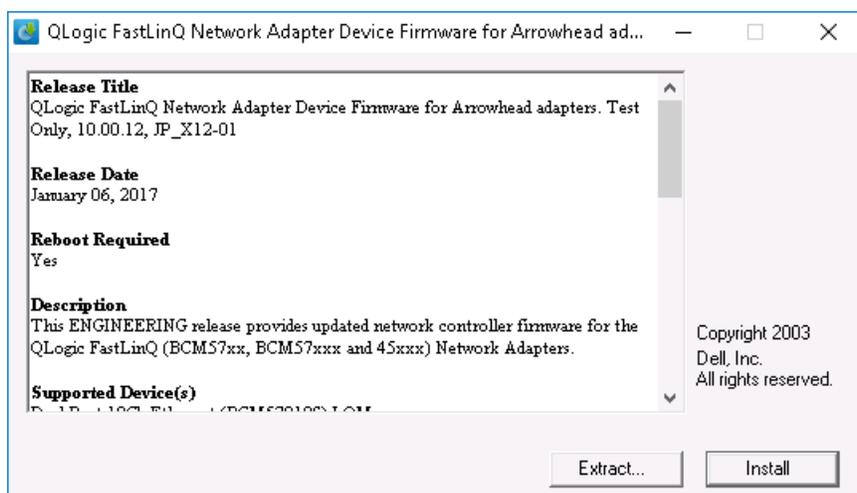
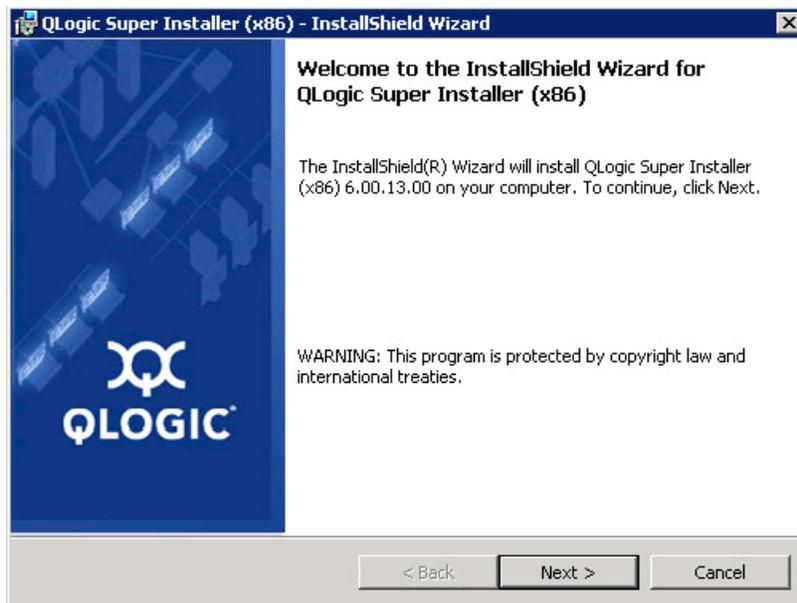


图 3-1. Dell Update Package 窗口

3. 在 QLogic Super Installer—InstallShield® 向导的 Welcome（欢迎）窗口（图 3-2）中，单击 **Next**（下一步）。



**图 3-2. QLogic InstallShield 向导：Welcome（欢迎）窗口**

4. 在向导的 License Agreement（许可协议）窗口（图 3-3）中完成以下操作：
  - a. 阅读 QLogic End User Software License Agreement（QLogic 最终用户许可协议）。
  - b. 选择 **I accept the terms in the license agreement**（我接受许可协议中的条款）继续。
  - c. 单击 **Next**（下一步）。

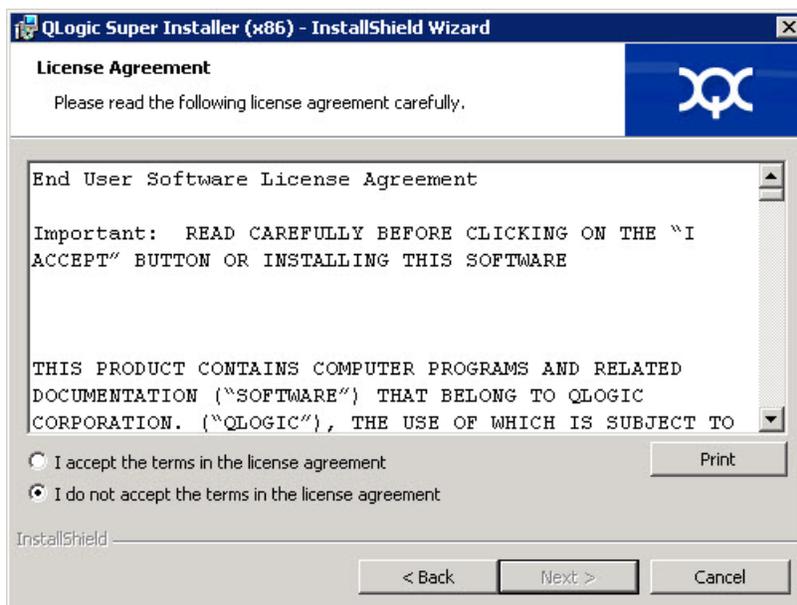


图 3-3. QLogic InstallShield 向导: License Agreement（许可协议）窗口

5. 如下完成向导的 Setup Type（设置类型）窗口（图 3-4）：
  - a. 选择以下设置类型之一：
    - 单击 **Complete**（完整）以安装所有程序功能。
    - 单击 **Custom**（自定义）以手动选择要安装的功能。
  - b. 单击 **Next**（下一步）继续。如果单击 **Complete**（完整），则直接继续步骤 6b。

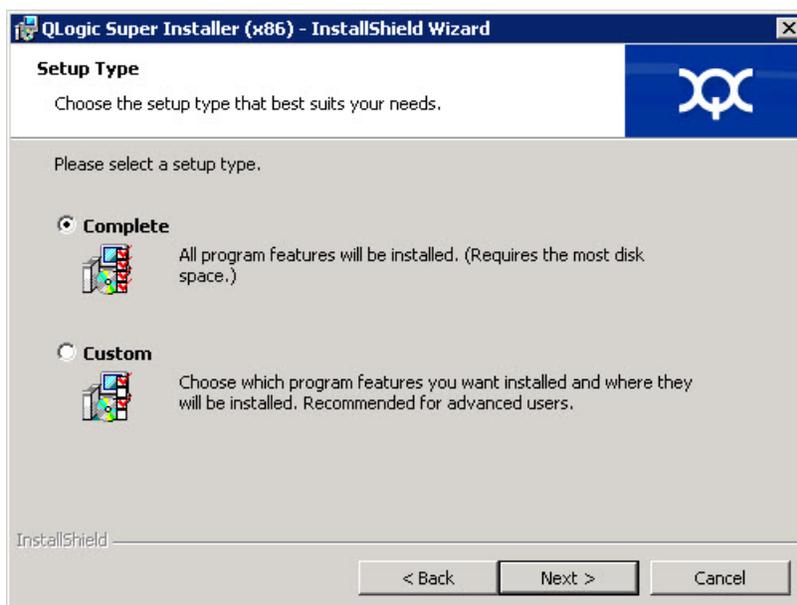
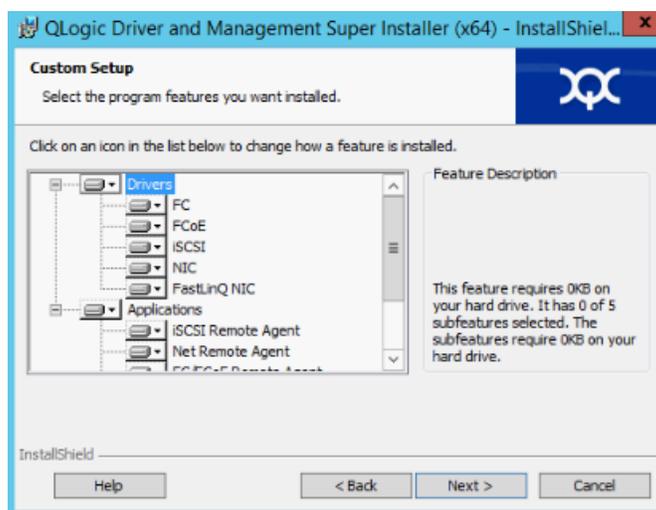


图 3-4. InstallShield 向导：Setup Type（设置类型）窗口

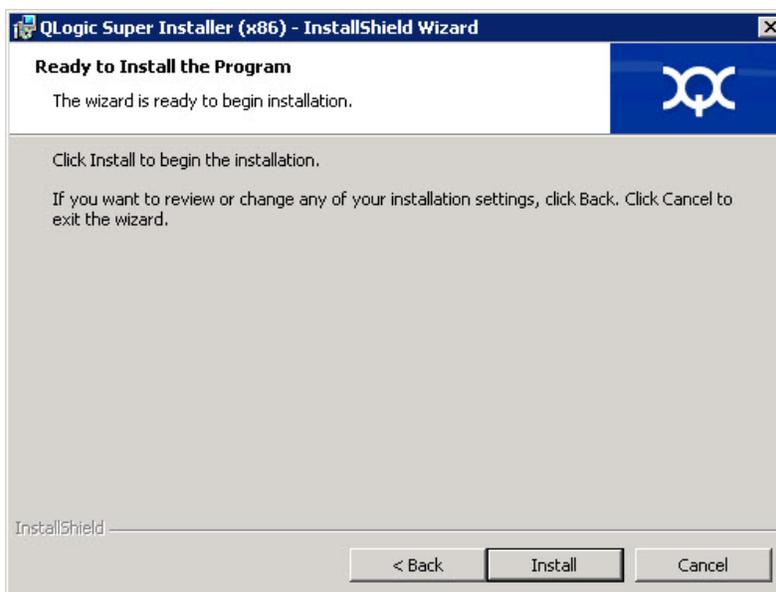
6. 如果您在步骤 5 中选择 **Custom**（自定义），则如下完成 Custom Setup（自定义设置）窗口（图 3-5）：
  - a. 选择要安装的功能。默认情况下，将选中所有功能。要更改某个功能的安装设置，请单击该功能旁边的图标，然后选择以下选项之一：
    - **This feature will be installed on the local hard drive**（此功能将安装在本地硬盘驱动器上）- 标记安装的功能，但不会影响其任何子功能。
    - **This feature, and all subfeatures, will be installed on the local hard drive**（此功能及所有子功能都将安装在本地硬盘驱动器上）- 标记安装的功能及其所有子功能。
    - **This feature will not be available**（此功能将不可用）- 阻止安装该功能。

- b. 单击 **Next**（下一步）继续。



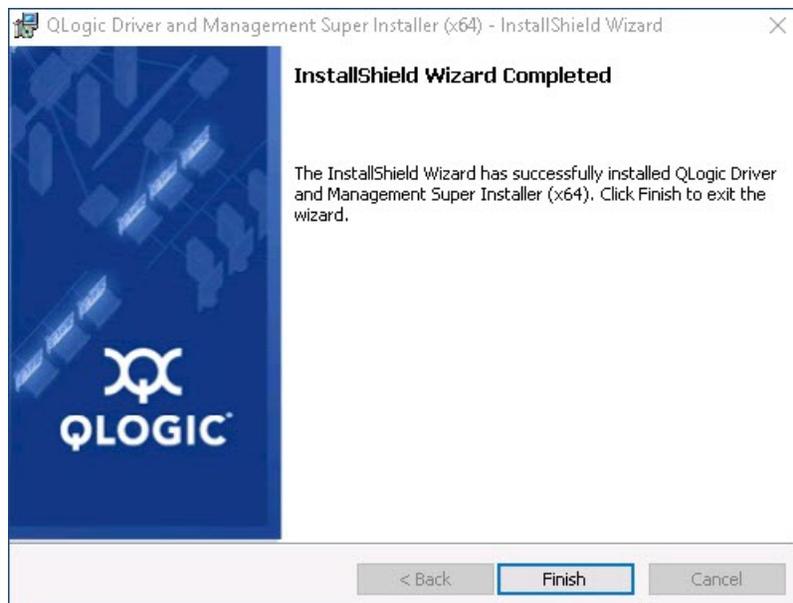
**图 3-5. InstallShield 向导：自定义设置窗口**

7. 在 InstallShield 向导的 Ready To Install（准备安装）窗口（图 3-6）中，单击 **Install**（安装）。InstallShield 向导将安装 QLogic 适配器驱动程序和管理软件安装程序。



**图 3-6. InstallShield 向导：Ready to Install the Program（准备安装程序）窗口**

8. 安装完成时，将显示 InstallShield Wizard Completed（InstallShield 向导完成）窗口（图 3-7）。单击 **Finish**（完成）关闭安装程序。



**图 3-7. InstallShield 向导：Completed（完成）窗口**

9. 在 Dell Update Package 窗口（图 3-8）中，“Update installer operation was successful（更新安装程序操作成功）”表示完成。
  - （可选）要打开日志文件，请单击 **View Installation Log**（查看安装日志）。日志文件显示 DUP 安装进度、任何之前安装的版本、任何错误消息，以及有关安装的其他信息。
  - 要关闭 Update Package（更新软件包）窗口，请单击 **CLOSE**（关闭）。

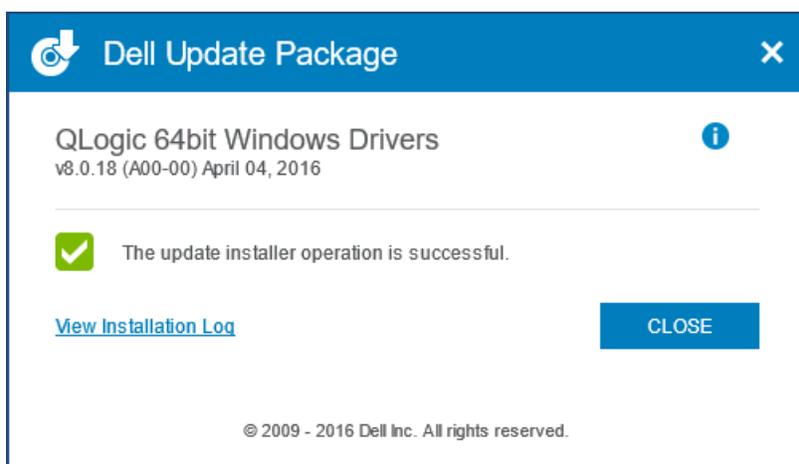


图 3-8. Dell Update Package 窗口

## DUP 安装选项

要自定义 DUP 安装行为，请使用以下命令行选项。

- 要仅将驱动程序组件提取到某个目录，请使用以下选项：

```
/drivers=<path>
```

### 注

此命令需要 **/s** 选项。

- 要仅安装或更新驱动程序组件，请使用以下选项：

```
/driveronly
```

### 注

此命令需要 **/s** 选项。

- (高级) 使用 `/passthrough` 选项可将 `passthrough` 后的所有文本直接发送至 DUP 的 QLogic 安装软件。此模式禁用提供的所有 GUI，但不一定禁用 QLogic 软件的 GUI。

`/passthrough`

- (高级) 要返回此 DUP 支持的功能的代码说明，请使用以下选项：

`/capabilities`

---

### 注

此命令需要 `/s` 选项。

---

## DUP 安装示例

以下示例显示如何使用安装选项。

要以静默方式更新系统，请使用以下命令：

```
<DUP_file_name>.exe /s
```

要将更新内容提取到 `C:\mydir\` 目录：

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

要将驱动程序组件提取到 `C:\mydir\` 目录：

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

要仅安装驱动程序组件，请使用以下命令：

```
<DUP_file_name>.exe /s /driveronly
```

要从默认日志位置更改为 `C:\my path with spaces\log.txt`，请使用以下命令：

```
<DUP_file_name>.exe /l="C:\my path with spaces\log.txt"
```

## 移除 Windows 驱动程序

要移除 Windows 驱动程序：

1. 在控制面板中，单击 **Programs**（程序），然后单击 **Programs and Features**（程序和功能）。
2. 在程序列表中，选择 **QLogic FastLinQ Driver Installer**（QLogic FastLinQ 驱动程序安装程序），然后单击 **Uninstall**（卸载）。
3. 按照说明操作以移除驱动程序。

## 管理适配器属性

要查看或更改 41xxx 系列适配器 属性：

1. 在控制面板中，单击 **Device Manager**（设备管理器）。
2. 在所选适配器的属性中，单击 **Advanced**（高级）选项卡。
3. 在 Advanced（高级）页面（图 3-9）中，选择 **Property**（属性）下的项目，然后根据需要更改该项目的 **Value**（值）。

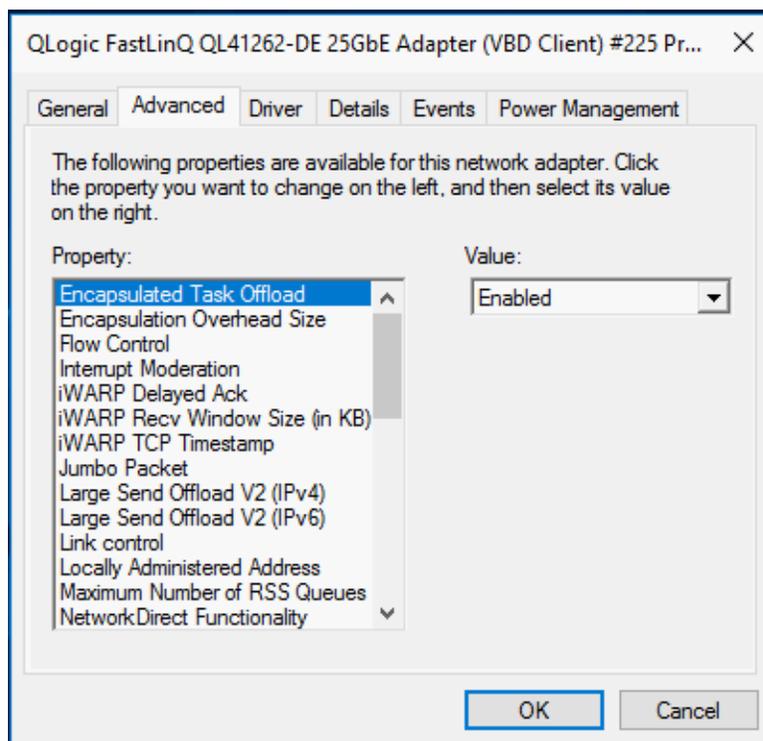


图 3-9. 设置高级适配器属性

## 设置电源管理选项

可以设置电源管理选项，以允许操作系统关闭该控制器以节约电源，或者允许该控制器唤醒计算机。如果设备正忙（例如，正在处理呼叫），操作系统将不会关闭设备。只有在计算机试图进入休眠状态时，操作系统才尝试尽可能关闭各个设备。要使控制器一直保持打开状态，不要选择 **Allow the computer to turn off the device to save power**（允许计算机关闭此设备以节约电源）复选框（图 3-10）。

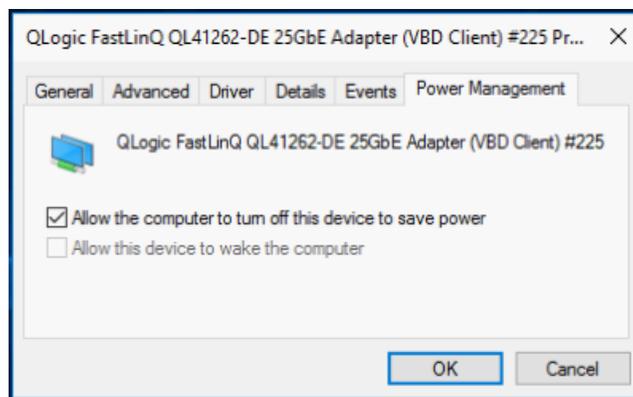


图 3-10. 电源管理选项

### 注

- 只有支持电源管理的服务器才有 Power Management（电源管理）页面。
  - 对于作为组成员的任何适配器，不要选中 **Allow the computer to turn off the device to save power**（允许计算机关闭设备以节约电源）复选框。
-

## 安装 VMware 驱动程序软件

本节介绍用于 41xxx 系列适配器的 qedentv VMware ESXi 驱动程序：

- [VMware 驱动程序和驱动程序包](#)
- [安装 VMware 驱动程序](#)
- [VMware 驱动程序可选参数](#)
- [VMware 驱动程序参数默认值](#)
- [移除 VMware 驱动程序](#)
- [FCoE 支持](#)
- [iSCSI 支持](#)

### VMware 驱动程序和驱动程序包

表 3-4 列出了协议的 VMware ESXi 驱动程序。

表 3-4. VMware 驱动程序

VMware 驱动程序	说明
qedentv	本机网络驱动程序
qedrntv	本机 RDMA 卸载（RoCE 和 RoCEv2）驱动程序 <sup>a</sup>
qedf	本机 FCoE 卸载驱动程序
qedil	旧版 iSCSI 卸载驱动程序

<sup>a</sup> 此版本中未包括经认证的 RoCE 驱动程序。未经认证的驱动程序可能会作为早期预览提供。

ESXi 驱动程序作为单独的驱动程序包随附，除非另有说明，否则不会捆绑在一起。表 3-5 列出了 ESXi 版本和适用的驱动程序版本。

表 3-5. 按版本的 ESXi 驱动程序包

ESXi 版本 <sup>a</sup>	协议	驱动程序名称	驱动程序版本
ESXi 6.5 <sup>b</sup>	NIC	qedentv	3.0.7.5
	FCoE	qedf	1.2.24.0
	iSCSI	qedil	1.0.19.0
	RoCE	qedrntv	3.0.7.5.1

表 3-5. 按版本的 ESXi 驱动程序包 (续)

ESXi 版本 <sup>a</sup>	协议	驱动程序名称	驱动程序版本
ESXi 6.0u3	NIC	qedentv	2.0.7.5
	FCoE	qedf	1.2.24.0
	iSCSI	qedil	1.0.19.0

<sup>a</sup> 公布本用户指南后, 可能提供其他 ESXi 驱动程序。有关更多信息, 请参考版本注释。

<sup>b</sup> 对于 ESXi 6.5, NIC 和 RoCE 驱动程序已打包在一起, 并且可使用标准 ESXi 安装命令作为单个脱机捆绑包进行安装。软件包名称为 `qedentv_3.0.7.5_qedrntv_3.0.7.5.1_signed_drivers.zip`。推荐的安装顺序是先安装 NIC 和 RoCE 驱动程序, 然后再安装 FCoE 和 iSCSI 驱动程序。

使用以下任一方法安装单独的驱动程序:

- 标准 ESXi 软件包安装命令 (请参阅[安装 VMware 驱动程序](#))
- 单独驱动程序自述文件中的步骤
- 以下 VMware KB 文章中的步骤:

[https://kb.vmware.com/selfservice/microsites/search.do?language=en\\_US&cmd=displayKC&externalId=2137853](https://kb.vmware.com/selfservice/microsites/search.do?language=en_US&cmd=displayKC&externalId=2137853)

您应该首先安装 NIC 驱动程序, 然后再安装存储驱动程序。

## 安装 VMware 驱动程序

可以使用驱动程序 ZIP 文件安装新驱动程序或更新现有驱动程序。请务必从同一驱动程序 ZIP 文件安装整个驱动程序集。混用来自不同 ZIP 文件的驱动程序会导致问题。

### 要安装 VMware 驱动程序:

1. 从 VMware 支持页面下载适用于 41xxx 系列适配器的 VMware 驱动程序:  
[www.vmware.com/support.html](http://www.vmware.com/support.html)
2. 启动 ESX 主机, 然后登录至具有管理员权限的帐户。
3. 解压缩驱动程序 ZIP 件, 然后提取 `.vib` 文件。
4. 使用 Linux `scp` 公用程序将 `.vib` 文件从本地系统复制到 IP 地址为 10.10.10.10 的 ESX 服务器上的 `/tmp` 目录中。例如, 发出以下命令:

```
#scp qedentv-1.0.3.11-10EM.550.0.0.1331820.x86_64.vib root@10.10.10.10:/tmp
```

您可以将该文件放在 ESX 控制台 shell 可以访问的任意位置。

### 注

如果没有 Linux 计算机，您可以使用 vSphere 数据存储文件浏览器将文件上传到服务器。

5. 通过发出以下命令，将主机置于维护模式：

```
#esxcli --maintenance-mode
```

6. 选择以下安装选项之一：

- 选项 1：**使用 CLI 或 VMware Update Manager (VUM) 直接在 ESX 服务器上安装 .vib。

- 要使用 CLI 安装 .vib，请发出以下命令。请务必指定完整的 .vib 文件路径。

```
# esxcli software vib install -v  
/tmp/qedentv-1.0.3.11-1OEM.550.0.0.1331820.x86_64.vib
```

- 要使用 VUM 安装 .vib 文件，请参阅知识库文章：

[使用 VMware vCenter Update Manager 4.x 和 5.x 更新 ESXi/ESX 主机 \(1019545\)](#)

- 选项 2：**发出以下命令，一次安装所有单独的 VIB：

```
# esxcli software vib install -d  
/tmp/qedentv-bundle-2.0.3.zip
```

#### 要升级现有的驱动程序：

按照全新安装的步骤进行操作，除了将之前选项 1 中的命令替换为以下命令：

```
#esxcli software vib update -v  
/tmp/qedentv-1.0.3.11-1OEM.550.0.0.1331820.x86_64.vib
```

## VMware 驱动程序可选参数

表 3-6 说明了可作为命令行参数提供给 esxcfg-module 命令的可选参数。

表 3-6. VMware 驱动程序可选参数

参数	说明
hw_vlan	全局启用 (1) 或禁用 (0) 硬件 VLAN 插入和移除。当上层需要发送或接收完整格式数据包时，禁用此参数。hw_vlan=1 为默认值。

**表 3-6. VMware 驱动程序可选参数 (续)**

参数	说明
num_queues	指定 TX/RX 队列对的数目。num_queues 可以是 1-11 或以下值之一： <ul style="list-style-type: none"> <li>■ -1 允许驱动程序决定队列对的最佳数目（默认值）。</li> <li>■ 0 使用默认队列。</li> </ul> 对于多端口或多功能配置，您可以指定用逗号分隔的多个值。
multi_rx_filters	指定每个 RX 队列的 RX 过滤器数目，默认队列除外。 multi_rx_filters 可以是 1-4 或以下值之一： <ul style="list-style-type: none"> <li>■ -1 使用默认的每队列 RX 过滤器数目。</li> <li>■ 0 禁用 RX 过滤器。</li> </ul>
disable_tpa	启用 (0) 或禁用 (1) TPA (LRO) 功能。disable_tpa=0 为默认值。
max_vfs	指定每个物理功能 (PF) 的虚拟功能 (VF) 数目。max_vfs 可以是一个端口上的 0 个 (已禁用) 或 64 个 VF (已启用)。ESXi 支持最大 64 个 VF 是操作系统资源分配约束。
RSS	指定主机或 PF 的虚拟可扩展 LAN (VXLAN) 隧道流量所使用的接收端伸缩队列数目。RSS 可以是 2、3、4 或以下值之一： <ul style="list-style-type: none"> <li>■ -1 使用默认队列数目。</li> <li>■ 0 或 1 禁用 RSS 队列。</li> </ul> 对于多端口或多功能配置，您可以指定用逗号分隔的多个值。
debug	指定驱动程序在 vmkernel 日志文件中记录的数据级别。debug 可具有以下值（按数据量升序显示）： <ul style="list-style-type: none"> <li>■ 0x80000000 表示通知级别。</li> <li>■ 0x40000000 表示信息级别（包括通知级别）。</li> <li>■ 0x3FFFFFFF 表示针对所有驱动程序子模块的详细级别（包括信息和通知级别）。</li> </ul>
auto_fw_reset	启用 (1) 或禁用 (0) 驱动程序自动固件恢复功能。启用此参数时，驱动程序会尝试从发送超时、固件断言和适配器奇偶校验错误等事件中恢复。默认值为 auto_fw_reset=1。
vxlan_filter_en	启用 (1) 或禁用 (0) 基于外层 MAC、内层 MAC 和 VXLAN 网络 (VNI) 的 VXLAN 过滤，将流量直接匹配至特定队列。默认值为 vxlan_filter_en=1。对于多端口或多功能配置，您可以指定用逗号分隔的多个值。
enable_vxlan_offld	启用 (1) 或禁用 (0) VXLAN 隧道流量校验和卸载和 TCP 分段卸载 (TSO) 功能。默认值为 enable_vxlan_offld=1。对于多端口或多功能配置，您可以指定用逗号分隔的多个值。

## VMware 驱动程序参数默认值

表 3-7 列出了 VMware 驱动程序参数默认值。

表 3-7. VMware 驱动程序参数默认值

参数	默认值
Speed (速度)	自动协商并广告所有速度。speed (速度) 参数在所有端口上必须相同。如果在设备上启用了自动协商, 则所有设备端口都将使用自动协商。
Flow Control (流控制)	自动协商并广告 RX 和 TX
MTU	1,500 (范围为 46 - 9,600)
Rx Ring Size (Rx 环大小)	8,192 (范围为 128 - 8,192)
Tx Ring Size (Tx 环大小)	8,192 (范围为 128 - 8,192)
MSI-X	Enabled (已启用)
Transmit Send Offload (传输发送卸载) (TSO)	Enabled (已启用)
Large Receive Offload (大量接收卸载) (LRO)	Enabled (已启用)
RSS	Enabled (已启用) (四个 RX 队列)
HW VLAN	Enabled (已启用)
Number of Queues (队列数)	Enabled (已启用) (八个 RX/TX 队列对)
Wake on LAN (局域网唤醒) (WoL)	Disabled (已禁用)

## 移除 VMware 驱动程序

要移除 .vib 文件 (qedentv), 请发出以下命令:

```
# esxcli software vib remove --vibname qedentv
```

要移除驱动程序, 请发出以下命令:

```
# vmkload_mod -u qedentv
```

## FCoE 支持

表 3-8 介绍 VMware 软件包中随附的驱动程序，用于支持 QLogic FCoE 聚合网络接口控制器 (C-NIC)。VMware ESXi 5.0 及以上版本支持 FCoE 和 DCB 功能。

**表 3-8. QLogic 41xxx 系列适配器 VMware FCoE 驱动程序**

驱动程序	说明
qedf	QLogic VMware FCoE 驱动程序是内核模式驱动程序，提供 VMware SCSI 堆栈与 QLogic FCoE 固件和硬件之间的转换层。

## iSCSI 支持

表 3-9 介绍 iSCSI 驱动程序。

**表 3-9. QLogic 41xxx 系列适配器 iSCSI 驱动程序**

驱动程序	说明
qedil	qedil 驱动程序是 QLogic VMware iSCSI HBA 驱动程序。qedil 与 qedf 类似，是内核模式驱动程序，提供 VMware SCSI 堆栈与 QLogic iSCSI 固件和硬件之间的转换层。qedil 利用 VMware iscsid 基础设施提供的服务进行会话管理和 IP 服务。

# 4 升级固件

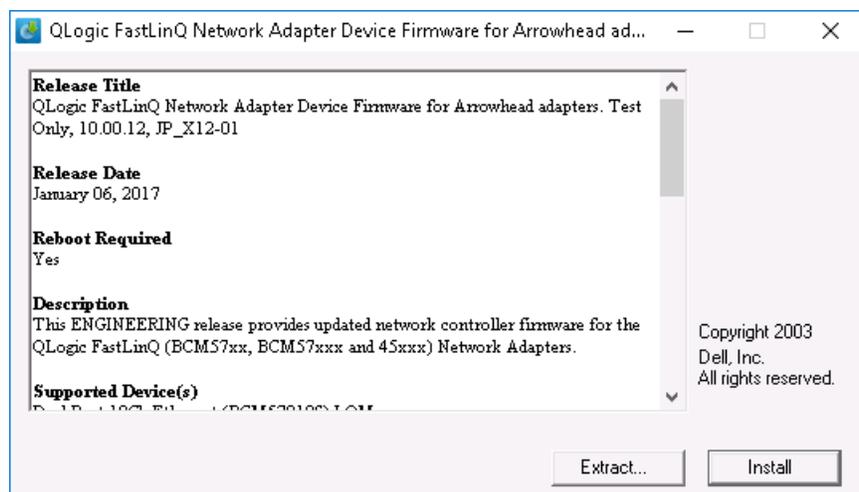
本章提供有关使用 Dell Update Package (DUP) 升级固件的信息。固件 DUP 是仅用于更新闪存的公用程序；它不用于适配器配置。您可以通过双击可执行文件来运行固件 DUP。或者，您可使用多个支持的命令行选项从命令行运行 DUP。

- [通过双击运行 DUP](#)
- [第 34 页上的“从命令行运行 DUP”](#)
- [第 35 页上的“使用 .bin 文件运行 DUP”](#)（仅适用于 Linux）

## 通过双击运行 DUP

要通过双击可执行文件来运行固件 DUP：

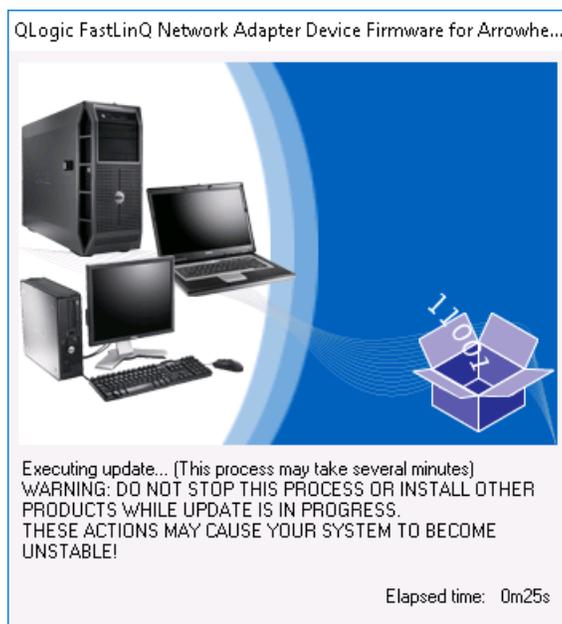
1. 双击代表固件 Dell Update Package 文件的图标。  
随即出现 Dell Update Package 初始屏幕，如 [图 4-1](#) 中所示。单击 **Install**（安装）以继续。



**图 4-1. Dell Update Package: 初始屏幕**

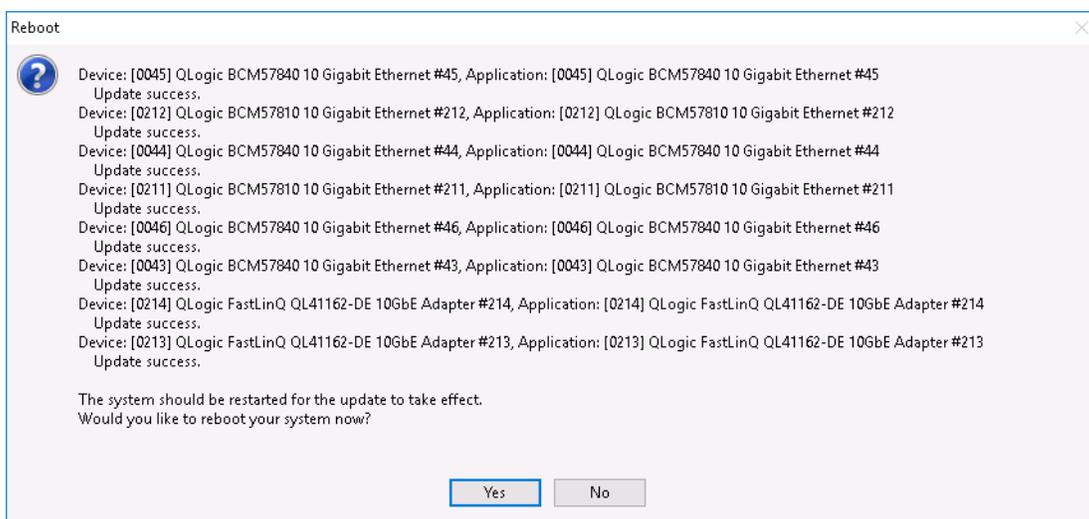
2. 请按照屏幕说明进行操作。在 Warning（警告）对话框中，单击 **Yes**（是）继续安装。

安装程序表示正在加载新固件，如 图 4-2 所示。



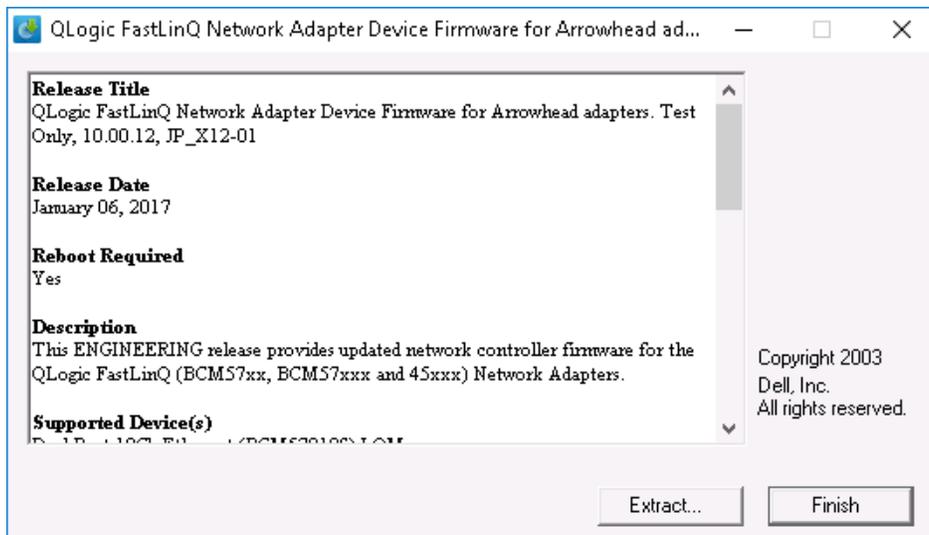
**图 4-2. Dell Update Package: 加载新固件**

完成时，安装程序表明安装结果，如 图 4-3 所示。



**图 4-3. Dell Update Package: 安装结果**

3. 单击 **Yes**（是）重新引导系统。
4. 单击 **Finish**（完成）以完成安装，如 [图 4-4](#) 所示。



**图 4-4. Dell Update Package: 完成安装**

## 从命令行运行 DUP

从命令行运行固件 DUP 而不指定选项时，会产生与双击 DUP 图标相同的行为。请注意，DUP 的实际文件名称各有不同。

### 要从命令行运行固件 DUP:

- 发出以下命令:

```
C:\> Network_Firmware_2T12N_WN32_<version>_X16.EXE
```

图 4-5 显示可用于自定义 Dell Update Package 安装的选项。

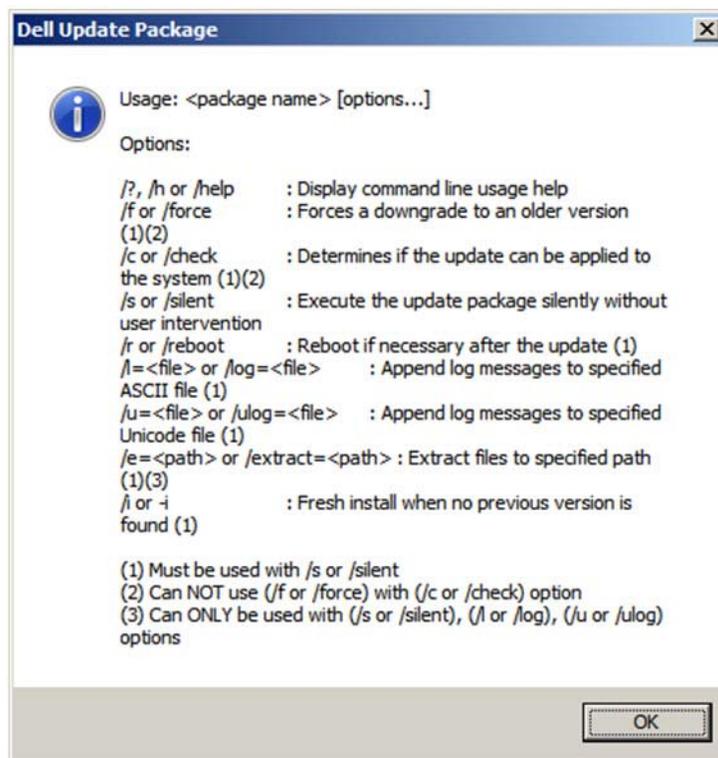


图 4-5. DUP 命令行选项

## 使用 .bin 文件运行 DUP

以下步骤仅在 Linux 操作系统上受支持。

### 要使用 .bin 文件更新 DUP:

1. 将 `Network_Firmware_NJCX1_LN_X.Y.Z.BIN` 文件复制到进行测试的系统 (SUT)。
2. 将文件类型更改为可执行文件，如下所示：  

```
chmod 777 Network_Firmware_NJCX1_LN_X.Y.Z.BIN
```
3. 要开始更新过程，请发出以下命令：  

```
./Network_Firmware_NJCX1_LN_X.Y.Z.BIN
```
4. 固件更新后，重新引导系统。

**DUP 更新期间 SUT 的示例输出:**

```
./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
```

# 5 适配器预引导配置

在主机引导过程中，您有机会暂停并使用人机界面基础设施 (HII) 应用程序执行适配器管理任务。包括以下任务：

- 第 38 页上的“使用入门”
- 第 41 页上的“显示固件映像属性”
- 第 42 页上的“配置设备级参数”
- 第 43 页上的“配置 NIC 参数”
- 第 47 页上的“配置数据中心桥接”
- 第 48 页上的“配置 FCoE 引导”
- 第 49 页上的“配置 iSCSI 引导”
- 第 54 页上的“配置分区”

---

## 注

本章中的 HII 屏幕截图具有代表性，可能不符合您在系统上看到的屏幕。

---

## 使用入门

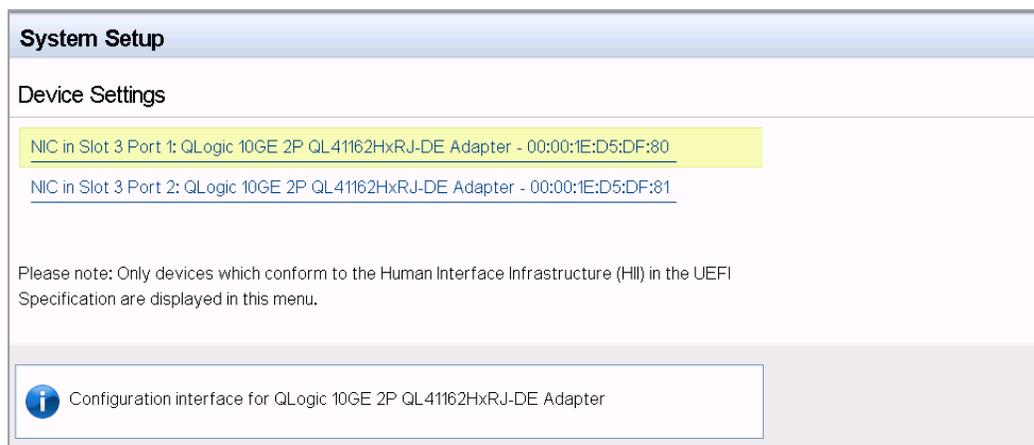
要启动 HII 应用程序：

1. 打开您的平台的 System Setup（系统设置）窗口。有关启动 System Setup（系统设置）的信息，请查询您的系统的用户指南。
2. 在 System Setup（系统设置）窗口（图 5-1）中，选择 **Device Settings**（设备设置），然后按 ENTER 键。



**图 5-1. 系统设置**

3. 在 Device Settings（设备设置）窗口（图 5-2）中，选择您要配置的 41xxx 系列适配器端口，然后按 ENTER 键。



**图 5-2. 系统设置：设备设置**

Main Configuration Page（主要配置页面，图 5-3）显示适配器管理选项，您可以在其中设置分区模式。

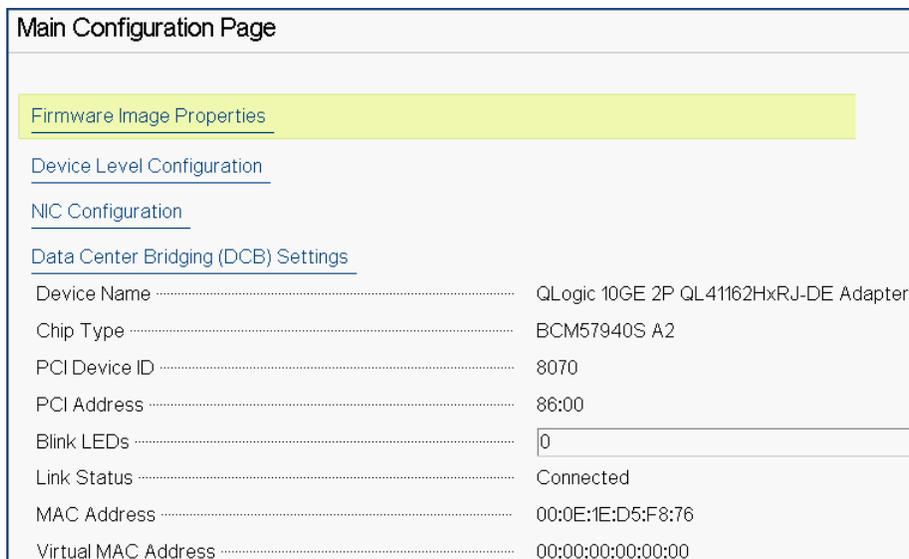


图 5-3. 主要配置页面

4. 在 **Device Level Configuration**（设备级配置）下，将 **Partitioning Mode**（分区模式）设置为 **NPAR**，以便将 **NIC Partitioning Configuration**（NIC 分区配置）选项添加到 Main Configuration Page（主要配置页面），如图 5-4 中所示。

**注**

NPAR 不可用于最高速度为 1G 的端口。

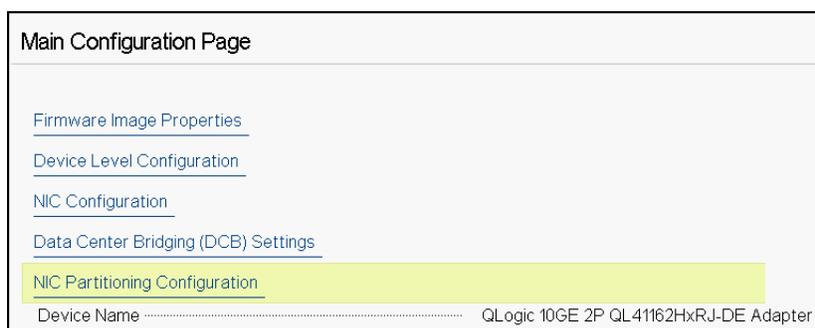


图 5-4. 主要配置页面，将分区模式设置为 NPAR

在图 5-3 和图 5-4 中，Main Configuration Page（主要配置页面）显示如下：

- **Firmware Image Properties**（固件映像属性）（请参阅第 41 页上的“显示固件映像属性”）
- **Device Level Configuration**（设备级配置）（请参阅第 42 页上的“配置设备级参数”）
- **NIC Configuration**（NIC 配置）（请参阅第 43 页上的“配置 NIC 参数”）
- **iSCSI Configuration**（iSCSI 配置）（如果通过 NPAR 模式启用端口第三个分区上的 iSCSI 卸载来允许 iSCSI 远程引导）（请参阅第 49 页上的“配置 iSCSI 引导”）
- **FCoE Configuration**（FCoE 配置）（如果通过 NPAR 模式启用端口第二个分区上的 FCoE 卸载来允许从 SAN 进行的 FCoE 引导）（请参阅第 48 页上的“配置 FCoE 引导”）
- **Data Center Bridging (DCB) Settings**（数据中心桥接 (DCB) 设置）（请参阅第 47 页上的“配置数据中心桥接”）
- **NIC Partitioning Configuration**（NIC 分区配置）（如果在 Device Level Configuration（设备级配置）页面上选择了 NPAR）（请参阅第 54 页上的“配置分区”）

此外，Main Configuration Page（主要配置页面）还提供表 5-1 中所列的适配器属性。

表 5-1. 适配器属性

适配器属性	说明
Device Name（设备名称）	工厂分配的设备名称
Chip Type（芯片类型）	ASIC 版本
PCI Device ID （PCI 设备 ID）	唯一的供应商特定 PCI 设备 ID
PCI Address（PCI 地址）	采用总线 - 设备功能格式的 PCI 设备地址
Blink LEDs（闪烁 LED）	用户定义的端口 LED 闪烁次数
Link Status（链路状态）	外部链路状态
MAC 地址	制造商分配的永久设备 MAC 地址
Virtual MAC Address （虚拟 MAC 地址）	用户定义的设备 MAC 地址

表 5-1. 适配器属性 (续)

适配器属性	说明
iSCSI MAC Address (iSCSI MAC 地址) <sup>a</sup>	制造商分配的永久设备 iSCSI 卸载 MAC 地址
iSCSI 虚拟 MAC 地址 <sup>a</sup>	用户定义的设备 iSCSI 卸载 MAC 地址
FCoE MAC Address (FCoE MAC 地址) <sup>b</sup>	制造商分配的永久设备 FCoE 卸载 MAC 地址
FCoE 虚拟 MAC 地址 <sup>b</sup>	用户定义的设备 FCoE 卸载 MAC 地址
FCoE WWPN <sup>b</sup>	制造商分配的永久设备 FCoE 卸载 WWPN (全局端口名称)
FCoE 虚拟 WWPN <sup>b</sup>	用户定义的设备 FCoE 卸载 WWPN
FCoE WWNN <sup>b</sup>	制造商分配的永久设备 FCoE 卸载 WWNN (全局节点名称)
FCoE 虚拟 WWNN <sup>b</sup>	用户定义的设备 FCoE 卸载 WWNN

<sup>a</sup> 只有在 NIC Partitioning Configuration (NIC 分区配置) 页面中启用了 **iSCSI Offload** (iSCSI 卸载) 时, 此属性才可见。

<sup>b</sup> 只有在 NIC Partitioning Configuration (NIC 分区配置) 页面中启用了 **FCoE Offload** (FCoE 卸载) 时, 此属性才可见。

## 显示固件映像属性

要查看固件映像的属性, 请选择 Main Configuration Page (主要配置页面) 中的 **Firmware Image Properties** (固件映像属性), 然后按 ENTER 键。Firmware Image Properties (固件映像属性) 页面 (图 5-5) 指定以下仅查看数据:

- **Family Firmware Version** (系列固件版本) 是多引导映像版本, 其包含多个固件组件映像。
- **MBI Version** (MBI 版本) 是设备上活动的 Cavium Qlogic 捆绑映像版本。
- **Controller BIOS Version** (控制器 BIOS 版本) 是管理固件版本。
- **EFI Driver Version** (EFI 驱动程序版本) 是可扩展固件接口 (EFI) 驱动程序版本。

- **L2B Firmware Version**（L2B 固件版本）是用于引导的 NIC 卸载固件版本。

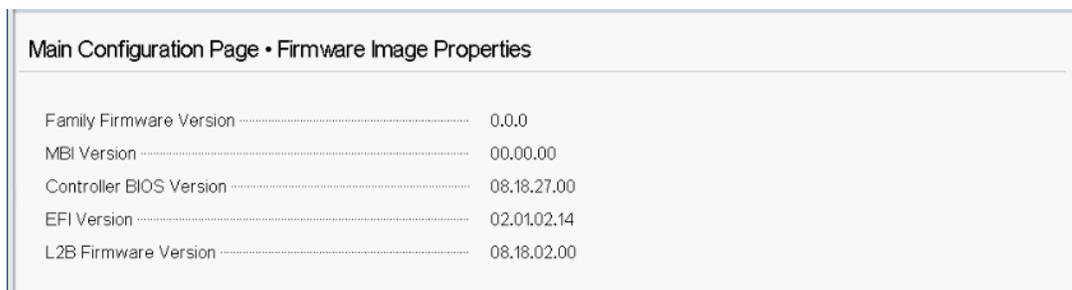


图 5-5. 固件映像属性

## 配置设备级参数

### 注

仅当在 NPAR 模式下启用了 iSCSI 卸载功能时，才会列出 iSCSI 物理功能 (PF)。仅当在 NPAR 模式下启用了 FCoE 卸载功能时，才会列出 FCoE PF。并非所有适配器型号都支持 iSCSI 卸载和 FCoE 卸载。每个端口只能启用一个卸载，并且只能在 NPAR 模式下启用。

设备级配置包括以下参数：

- 虚拟化模式
- NPAREP 模式

要配置设备级参数：

1. 在 Main Configuration Page（主要配置页面）中，选择 **Device Level Configuration**（设备级配置）（请参阅第 39 页上的图 5-3），然后按 ENTER 键。
2. 在 **Device Level Configuration**（设备级配置）页面中，选择设备级参数的值，如图 5-6 中所示。

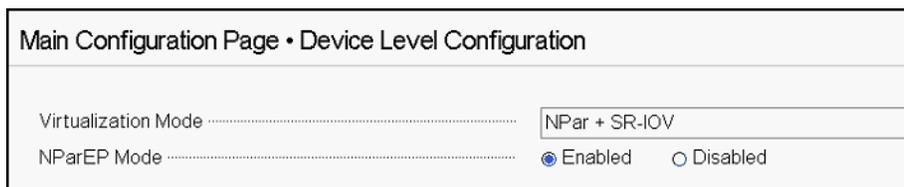


图 5-6. 设备级别配置

### 注

QL41264HMCU-DE（部件号 5V6Y4）和 QL41264HMRJ-DE（部件号 0D1WT）适配器在 Device Level Configuration（设备级配置）中显示对 NPAR、SR-IOV 和 NPAR-EP 的支持，不过 1Gbps 端口 3 和 4 上不支持这些功能。

---

- 对于 **Virtualization Mode**（虚拟化模式），选择将以下模式之一应用到所有适配器端口：
  - None**（无）（默认）指定未启用任何虚拟化模式。
  - NPAR** 将适配器设置为与交换机无关的 NIC 分区模式。
  - SR-IOV** 将适配器设置为 SR-IOV 模式。
  - NPar + SR-IOV** 将适配器设置为 SR-IOV over NPAR 模式。
- NParEP Mode**（NParEP 模式）配置每个适配器的最大分区数量。当在 [步骤 2](#) 中选择 **NPAR** 或 **NPar + SR-IOV** 作为 **Virtualization Mode**（虚拟化模式）时，此参数可见。
  - Enabled**（启用）允许为每个适配器配置最多 16 个分区。
  - Disabled**（禁用）允许为每个适配器配置最多 8 个分区。
- 单击 **Back**（后退）。
- 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

## 配置 NIC 参数

NIC 配置包括设置以下参数：

- 链路速度
- NIC + RDMA 模式
- RDMA 协议支持
- 引导模式
- FEC 模式
- 节能以太网
- 虚拟 LAN 模式
- 虚拟 LAN ID

**要配置 NIC 参数：**

1. 在 Main Configuration Page（主要配置页面）上，选择 **NIC Configuration**（NIC 配置，第 39 页上的图 5-3），然后单击 **Finish**（完成）。

图 5-7 显示 NIC Configuration（NIC 配置）页面。

Main Configuration Page • NIC Configuration	
Link Speed .....	<input checked="" type="radio"/> Auto Negotiated <input type="radio"/> 1 Gbps <input type="radio"/> 10 Gbps <input type="radio"/> 25 Gbps <input type="radio"/> SmartAN
NIC + RDMA Mode .....	<input checked="" type="radio"/> Enabled <input type="radio"/> Disabled
RDMA Protocol Support .....	<input checked="" type="radio"/> RoCE <input type="radio"/> iWARP <input type="radio"/> iWARP + RoCE
Boot Mode .....	<input type="radio"/> PXE <input checked="" type="radio"/> iSCSI <input type="radio"/> Disabled
Energy Efficient Ethernet .....	Optimal Power and Performance
Virtual LAN Mode .....	<input type="radio"/> Enabled <input checked="" type="radio"/> Disabled
Virtual LAN ID .....	1

**图 5-7. NIC 配置**

2. 为所选端口选择以下 **Link Speed**（链路速度）选项之一。并非所有适配器都提供所有速度选择。
  - Auto Negotiated**（自动协商）在端口上启用 Auto Negotiation（自动协商）模式。FEC 模式选择不可用于此速度模式。
  - 1 Gbps** 在端口上启用 1GbE 固定速度模式。此模式仅用于 1GbE 接口，不应配置用于以其他速度运行的适配器接口。FEC 模式选择不可用于此速度模式。此模式在所有适配器上都不可用。
  - 10 Gbps** 在端口上启用 10GbE 固定速度模式。并非所有适配器都提供此模式。
  - 25 Gbps** 在端口上启用 25GbE 固定速度模式。并非所有适配器都提供此模式。
  - SmartAN**（默认）在端口上启用 FastLinQ SmartAN™ 链路速度模式。FEC 模式选择不可用于此速度模式。**SmartAN** 设置在所有可能的链路速度和 FEC 模式之间循环，直到链路建立。该模式仅用于 25G 接口。如果您将 SmartAN 配置为 10Gb 接口，则系统将应用 10G 接口的设置。此模式在所有适配器上都不可用。
3. 对于 **NIC + RDMA Mode**（NIC + RDMA 模式），在端口上为 RDMA 选择 **Enabled**（启用）或 **Disabled**（禁用）。如果在 NPAR 模式下，此设置适用于端口的所有分区。

4. 在 [步骤 2](#) 中，当选择 **25 Gbps** 固定速度模式为 **Link Speed**（链路速度）时，**FEC Mode**（FEC 模式）可见。对于 **FEC Mode**（FEC 模式），选择以下选项之一：并非所有适配器都可使用 FEC 模式。
  - None**（无）禁用所有 FEC 模式。
  - Fire Code**（Fire 码）启用 Fire Code (BASE-R) FEC 模式。
  - Reed Solomon**（理德 - 所罗门）启用 Reed Solomon FEC 模式。
  - Auto**（自动）使端口能够采用循环法循环使用 **None**（无）、**Fire Code**（Fire 码）和 **Reed Solomon**（理德 - 所罗门）FEC 模式（以该链路速度），直到建立链路。
5. 如果在 NPAR 模式下，**RDMA Protocol Support**（RDMA 协议支持）设置适用于端口的所有分区。如果在 [步骤 3](#) 中将 **NIC + RDMA Mode**（NIC + RDMA 模式）设置为 **Enabled**（已启用），则会显示此设置。**RDMA Protocol Support**（RDMA 协议支持）选项包括以下内容：
  - RoCE** 在此端口上启用 RoCE 模式。
  - iWARP** 在此端口上启用 iWARP 模式。
  - iWARP + RoCE** 在此端口上启用 iWARP 和 RoCE 模式。这是默认值。此选项需要 Linux 的其他配置，如 [第 93 页上的“配置 iWARP 和 RoCE”](#) 中所述。
6. 对于 **Boot Mode**（引导模式），选择以下值之一：
  - PXE** 启用 PXE 引导。
  - FCoE** 启用硬件卸载路径上的从 SAN 的 FCoE 引导。只有在 NPAR 模式下启用了第二个分区上的 **FCoE Offload**（FCoE 卸载）时，**FCoE** 模式才可用（请参阅 [第 54 页上的“配置分区”](#)）。
  - iSCSI** 启用硬件卸载路径上的 iSCSI 远程引导。只有在 NPAR 模式下启用了第三个分区上的 **iSCSI Offload**（iSCSI 卸载）时，**iSCSI** 模式才可用（请参阅 [第 54 页上的“配置分区”](#)）。
  - Disabled**（禁用）阻止将此端口用作远程引导源。
7. **Energy Efficient Ethernet**（节能以太网，EEE）参数仅在 100BASE-T 或 10GBASE-T RJ45 接口的适配器上可见。从以下 EEE 选项中选择：
  - Disabled**（禁用）禁用此端口上的 EEE。
  - Optimal Power and Performance**（最佳电源和性能）在此端口上启用最佳电源和性能模式下的 EEE。
  - Maximum Power Savings**（最高节能）在此端口上启用最高节能模式下的 EEE。

- Maximum Performance**（最高性能）在此端口上启用最高性能模式下的 EEE。
- 8. 处于 PXE 远程安装模式时，**Virtual LAN Mode**（虚拟 LAN 模式）适用于整个端口。它在 PXE 远程安装完成后不会持续。从以下 VLAN 选项中选择：
  - 对于 PXE 远程安装模式，**Enabled**（启用）启用此端口上的 VLAN 模式。
  - Disabled**（禁用）禁用此端口上的 VLAN 模式。
- 9. **Virtual LAN ID**（虚拟 LAN ID）参数指定要在此端口上用于 PXE 远程安装模式的 VLAN 标记 ID。此设置仅当在上一步中启用了 **Virtual LAN Mode**（虚拟 LAN 模式）时适用。
- 10. 单击 **Back**（后退）。
- 11. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

**要将端口配置为使用 RDMA：**

---

**注**

按照以下步骤在 NPAR 模式端口的所有分区上启用 RDMA。

---

1. 将 **NIC + RDMA Mode**（NIC + RDMA 模式）设置为 **Enabled**（启用）。
2. 单击 **Back**（后退）。
3. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

**要配置端口的引导模式：**

1. 对于 UEFI PXE 远程安装，请选择 **PXE** 作为 **Boot Mode**（引导模式）。
2. 单击 **Back**（后退）。
3. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

**要配置端口的 PXE 远程安装以使用 VLAN：**

---

**注**

PXE 远程安装完成后，此 VLAN 不会持续。

---

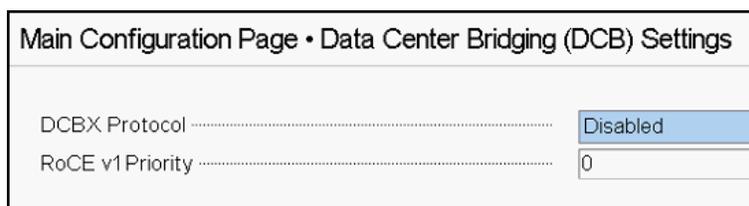
1. 将 **Virtual LAN Mode**（虚拟 LAN 模式）设置为 **Enabled**（启用）。
2. 在 **Virtual LAN ID**（虚拟 LAN ID）框中，输入要使用的编号。
3. 单击 **Back**（后退）。
4. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

## 配置数据中心桥接

数据中心桥接 (DCB) 设置包括 DCBX 协议和 RoCE 优先级。

### 要配置 DCB 设置：

1. 在 Main Configuration Page（主要配置页面）（第 39 页上的图 5-3）中，选择 **Data Center Bridging (DCB) Settings**（数据中心桥接 (DCB) 设置），然后单击 **Finish**（完成）。
2. 图 5-8 在 Data Center Bridging (DCB) Settings（数据中心桥接 (DCB) 设置）页面中，选择适当的 **DCBX Protocol**（DCBX 协议）选项：
  - Disabled**（禁用）禁用此端口上的 DCBX。
  - CEE** 在此端口上启用旧版聚合增强型以太网 (CEE) 协议 DCBX 模式。
  - IEEE** 在此端口上启用 IEEE DCBX 协议。
  - Dynamic**（动态）启用 CEE 或 IEEE 协议的动态应用，以匹配所连接的链路伙伴。
3. 在 Data Center Bridging (DCB) Settings（数据中心桥接 (DCB) 设置）页面上，输入 **RoCE v1 Priority**（RoCE v1 优先级）的值 **0-7**。此设置指明用作 RoCE 流量的 DCB 流量类优先级编号，并且应与启用了 DCB 的交换网络用于 RoCE 流量的编号相匹配。
  - 0** 指定有损默认或通用流量类使用的正常优先级编号。
  - 3** 指定无损 FCoE 流量使用的优先级编号。
  - 4** 指定通过 DCB 的无损 iSCSI-TLV 流量使用的优先级编号。
  - 1、2、5、6 和 7** 指定可用于 RoCE 用途的 DCB 流量类优先级编号。按照相应的操作系统 RoCE 设置说明来使用此 RoCE 控件。



Main Configuration Page • Data Center Bridging (DCB) Settings	
DCBX Protocol .....	Disabled
RoCE v1 Priority .....	0

图 5-8. 系统设置：数据中心桥接 (DCB) 设置

4. 单击 **Back**（后退）。
5. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

### 注

当启用 DCBX 时，适配器周期性地发送包含专用单播地址的链路层发现协议 (LLDP) 数据包，该专用单播地址充当源 MAC 地址。此 LLDP MAC 地址与工厂分配的适配器以太网 MAC 地址不同。如果您检查连接到适配器的交换机端口的 MAC 地址表，您将看到两个 MAC 地址：一个用于 LLDP 数据包，另一个用于适配器以太网接口。

## 配置 FCoE 引导

### 注

只有在 NPAR 模式下启用了第二个分区上的 **FCoE Offload Mode**（FCoE 卸载模式）时，FCoE Boot Configuration Menu（FCoE 引导配置菜单）才可见（请参阅第 57 页上的图 5-18）。它在非 NPAR 模式下不可见。

### 要配置 FCoE 引导配置参数：

1. 在 Main Configuration Page（主要配置页面）中，选择 **FCoE Configuration**（FCoE 配置），然后根据需要选择以下选项之一：
  - FCoE General Parameters**（FCoE 常规参数）（图 5-9）
  - FCoE Target Configuration**（FCoE 目标配置）（图 5-10）
2. 按 ENTER 键。
3. 选择 FCoE 常规参数或 FCoE 目标配置参数的值。

Main Configuration Page • FCoE Configuration • FCoE General Parameters	
Fabric Discovery Retry Count .....	<input type="text" value="5"/>
LUN Busy Retry Count .....	<input type="text" value="5"/>

图 5-9. FCoE 常规参数

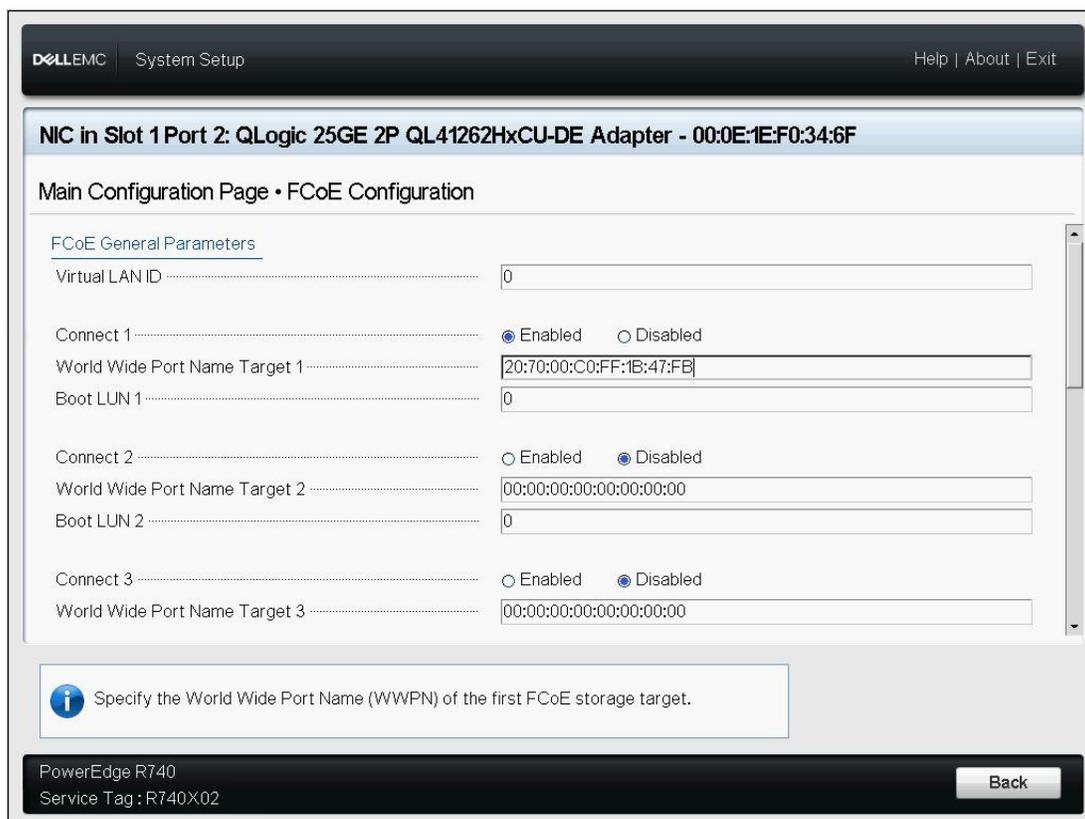


图 5-10. FCoE 目标配置

4. 单击 **Back**（后退）。
5. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

## 配置 iSCSI 引导

### 注

只有在 NPAR 模式下启用了第三个分区上的 **iSCSI Offload Mode**（iSCSI 卸载模式）时，iSCSI Boot Configuration Menu（iSCSI 引导配置菜单）才可见（请参阅第 58 页上的图 5-19）。它在非 NPAR 模式下不可见。

### 要配置 iSCSI 引导配置参数：

1. 在 Main Configuration Page（主要配置页面）中，选择 **iSCSI Boot Configuration**（iSCSI 引导配置）菜单，然后选择以下选项之一：
  - iSCSI General Configuration**（iSCSI 常规配置）
  - iSCSI Initiator Configuration**（iSCSI 启动器配置）

- iSCSI First Target Configuration** (iSCSI 第一目标配置)
- iSCSI Second Target Configuration** (iSCSI 第二目标配置)
- 2. 按 ENTER 键。
- 3. 选择相应 iSCSI 配置参数的值：
  - iSCSI General Parameters** (iSCSI 常规参数) (第 51 页上的图 5-11)
    - TCP/IP Parameters Via DHCP (通过 DHCP 获取 TCP/IP 参数)
    - iSCSI Parameters Via DHCP (通过 DHCP 获取 iSCSI 参数)
    - CHAP Authentication (CHAP 身份验证)
    - CHAP 相互身份验证
    - IP Version (IP 版本)
    - ARP Redirect (ARP 重定向)
    - DHCP Request Timeout (DHCP 请求超时)
    - Target Login Timeout (目标登录超时)
    - DHCP Vendor ID (DHCP 供应商 ID)
  - iSCSI Initiator Parameters** (iSCSI 启动器参数) (第 52 页上的图 5-12)
    - IPv4 地址
    - IPv4 子网掩码
    - IPv4 默认网关
    - IPv4 主要 DNS
    - IPv4 辅助 DNS
    - VLAN ID
    - iSCSI Name (iSCSI 名称)
    - CHAP ID
    - CHAP Secret (CHAP 机密)
  - iSCSI First Target Parameters** (iSCSI 第一目标参数) (第 52 页上的图 5-13)
    - Connect (连接)
    - IPv4 地址
    - TCP Port (TCP 端口)
    - Boot LUN (引导 LUN)
    - iSCSI Name (iSCSI 名称)
    - CHAP ID
    - CHAP Secret (CHAP 机密)
  - iSCSI Second Target Parameters** (iSCSI 第二目标参数) (第 53 页上的图 5-14)
    - Connect (连接)
    - IPv4 地址

- TCP Port (TCP 端口)
  - Boot LUN (引导 LUN)
  - iSCSI Name (iSCSI 名称)
  - CHAP ID
  - CHAP Secret (CHAP 机密)
4. 单击 **Back** (后退)。
  5. 看到提示时, 单击 **Yes** (是) 以保存更改。更改会在系统重设后生效。

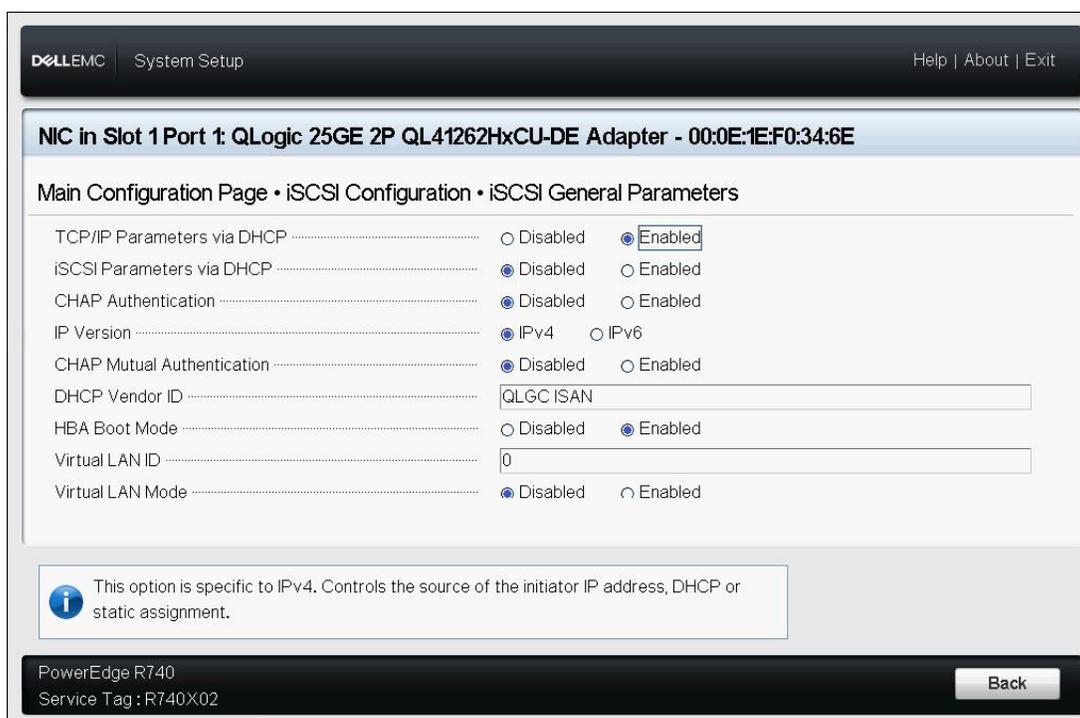


图 5-11. iSCSI 常规参数

The screenshot shows the 'iSCSI Initiator Parameters' configuration page. At the top, it identifies the network interface as 'NIC in Slot 1 Port 1: QLogic 25GE 2P QL41262HxCU-DE Adapter - 00:0E:1E:F0:34:6E'. The page title is 'Main Configuration Page • iSCSI Configuration • iSCSI Initiator Parameters'. The configuration fields are as follows:

IPv4 Address	192.168.100.145
Subnet Mask	255.255.255.0
IPv4 Default Gateway	0.0.0.0
IPv4 Primary DNS	0.0.0.0
IPv4 Secondary DNS	0.0.0.0
iSCSI Name	iqn.1994-02.com.qlogic.iscsi:fastlinqboot
CHAP ID	
CHAP Secret	

Below the fields is an information box: **i** Specify the iSCSI Qualified Name (IQN) of the initiator.

The footer shows 'PowerEdge R740' and 'Service Tag : R740X02' with a 'Back' button.

图 5-12. iSCSI 启动器配置参数

The screenshot shows the 'iSCSI First Target Parameters' configuration page. At the top, it identifies the network interface as 'NIC in Slot 1 Port 1: QLogic 25GE 2P QL41262HxCU-DE Adapter - 00:0E:1E:F0:34:6E'. The page title is 'Main Configuration Page • iSCSI Configuration • iSCSI First Target Parameters'. The configuration fields are as follows:

Connect	<input type="radio"/> Disabled <input checked="" type="radio"/> Enabled
IPv4 Address	192.168.100.9
TCP Port	3260
Boot LUN	1
iSCSI Name	iqn.2002-03.com.compellent:5000d31000ee1246
CHAP ID	
CHAP Secret	

Below the fields is an information box: **i** Specify the IPV4 address of the first iSCSI target.

The footer shows 'PowerEdge R740' and 'Service Tag : R740X02' with a 'Back' button.

图 5-13. iSCSI 第一目标参数

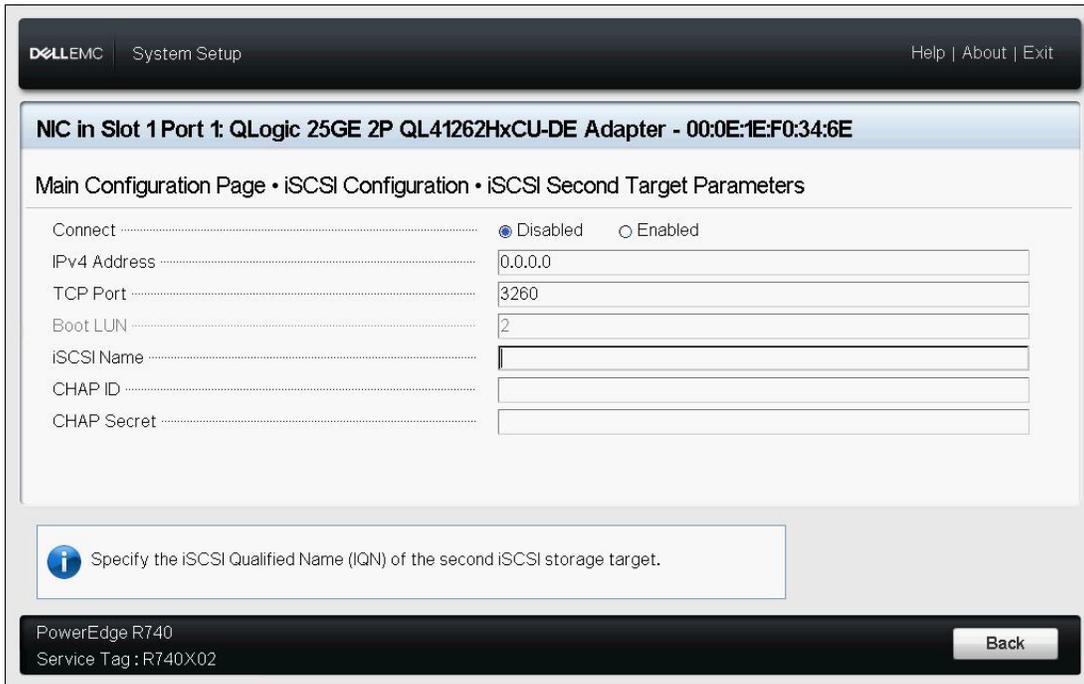


图 5-14. iSCSI 第二目标参数

## 配置分区

您可以为适配器上的每个分区配置带宽范围。有关 VMware ESXi 6.0/6.5 上分区配置的特定信息，请参阅第 59 页上的“VMware ESXi 6.0 和 ESXi 6.5 分区”。

要配置最大和最小带宽分配：

1. 在 Main Configuration Page（主要配置页面）中，选择 **NIC Partitioning Configuration**（NIC 分区配置），然后按 ENTER 键。
2. 在 Partitioning Configuration（分区配置）页面（图 5-15）中，选择 **Global Bandwidth Allocation**（全局带宽分配）。

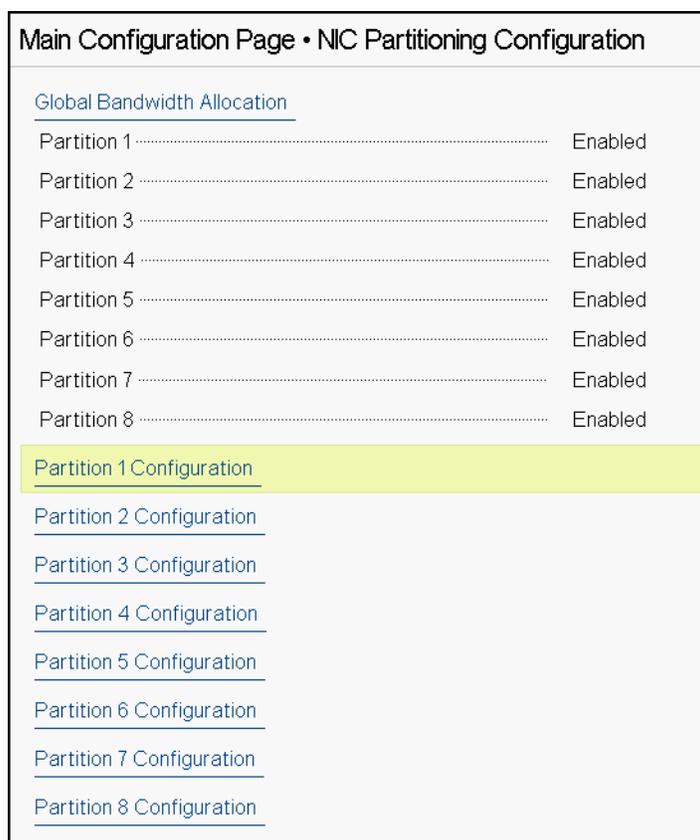


图 5-15. NIC 分区配置，全局带宽分配

3. 在 Global Bandwidth Allocation（全局带宽分配）页面（图 5-16）中，单击您要为其分配带宽的每个分区的最小和最大 TX 带宽字段。在双端口模式下，每个端口有 8 个分区。

Main Configuration Page • NIC Partitioning Configuration • Global Bandwidth Allocation	
Partition 1 Minimum TX Bandwidth .....	0
Partition 2 Minimum TX Bandwidth .....	0
Partition 3 Minimum TX Bandwidth .....	0
Partition 4 Minimum TX Bandwidth .....	0
Partition 5 Minimum TX Bandwidth .....	0
Partition 6 Minimum TX Bandwidth .....	0
Partition 7 Minimum TX Bandwidth .....	0
Partition 8 Minimum TX Bandwidth .....	0
Partition 1 Maximum TX Bandwidth .....	100
Partition 2 Maximum TX Bandwidth .....	100
Partition 3 Maximum TX Bandwidth .....	100

 Minimum Bandwidth represents the minimum transmit bandwidth of the partition as percentage of the full physical port link speed. The Minimum ... (Press <F1> for more help)

图 5-16. 全局带宽分配页面

- **Partition *n* Minimum TX Bandwidth**（分区 *n* 最小 TX 带宽）是所选分区的最小发送带宽，以物理端口最大链路速度的百分比形式表示。值的有效范围是 0 - 100。如果启用了 DCBX ETS 模式，则每流量类型 DCBX ETS 最小带宽值会与每分区最小 TX 带宽值同时使用。单个端口上所有分区的最小 TX 带宽值的总和必须等于 100 或全为零。
- 将 TX 带宽全设置为零类似于在每个活动分区上等分可用带宽；然而，带宽在所有主动发送分区上动态分配。当拥塞（来自所有分区）限制 TX 带宽时，零值（将一个或多个其他值设置为非零值时）将为该分区分配百分之一的最小值。
- **Partition *n* Maximum TX Bandwidth**（分区 *n* 最大 TX 带宽）是所选分区的最大发送带宽，以物理端口最大链路速度的百分比形式表示。值的有效范围是 1 - 100。无论 DCBX ETS 模式设置如何，每分区最大 TX 带宽值都适用。

在每个所选字段中键入值，然后单击 **Back**（后退）。

4. 看到提示时，单击 **Yes**（是）以保存更改。更改会在系统重设后生效。

**要配置分区：**

1. 要检查特定分区配置，在 NIC Partitioning Configuration （NIC 分区配置）页面（第 54 页上的图 5-15）中，选择 **Partition n Configuration** （分区 n 配置）。如果未启用 NParEP，则每个端口只能存在四个分区。
2. 要配置第一个分区，请选择 **Partition 1 Configuration** （分区 1 配置）以打开 Partition 1 Configuration （分区 1 配置）页面（图 5-17），其中显示以下参数：
  - NIC 模式** （始终启用）
  - PCI Device ID** （PCI 设备 ID）
  - PCI （总线）地址**
  - MAC 地址**
  - Virtual MAC Address** （虚拟 MAC 地址）

如果未启用 NParEP，则每个端口只有 4 个分区可用。在不具备卸载能力的适配器上，不会显示 **FCoE Mode** （FCoE 模式）及 **iSCSI Mode** （iSCSI 模式）的选项和信息。

Main Configuration Page • NIC Partitioning Configuration • 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 .....	00:00:00:00:00:00

**图 5-17. 分区 1 配置**

3. 要配置第二个分区，请选择 **Partition 2 Configuration** （分区 2 配置）以打开 Partition 2 Configuration （分区 2 配置）页面。如果存在 FCoE 卸载，则 Partition 2 Configuration （分区 2 配置）（图 5-18）将显示以下参数：
  - NIC Mode** （NIC 模式）可在分区 2 及更高分区上启用或禁用 L2 以太网 NIC 个性设置。要禁用所有剩余的分区，请将 **NIC Mode** （NIC 模式）设置为 **Disabled** （已禁用）。要禁用具备卸载能力分区，请同时禁用 **NIC Mode** （NIC 模式）和相应的卸载模式。
  - FCoE Mode** （FCoE 模式）启用或禁用第二个分区上的 FCoE 卸载个性设置。如果在第二个分区上启用此模式，则应禁用 **NIC Mode** （NIC 模式）。由于每个端口只有一个卸载可用，如果在端口的第二个分区上启用了 FCoE 卸载，就无法在同一 NPAR 模式端口的第三个分区上启用 iSCSI 卸载。并非所有适配器都支持 **FCoE Mode** （FCoE 模式）。

- iSCSI Mode** (iSCSI 模式) 启用或禁用第二个分区上的 iSCSI 卸载个性设置。如果在第三个分区上启用此模式, 则应禁用 **NIC Mode** (NIC 模式)。由于每个端口只有一个卸载可用, 如果在端口的第三个分区上启用了 iSCSI 卸载, 就无法在同一 NPAR 模式端口的第二个分区上启用 FCoE 卸载。并非所有适配器都支持 **iSCSI Mode** (iSCSI 模式)。
- FIP MAC Address** (FIP MAC 地址)<sup>1</sup>
- Virtual FIP MAC Address** (虚拟 FIP MAC 地址)<sup>1</sup>
- 全局端口名称**<sup>1</sup>
- Virtual World Wide Port Name** (虚拟全局端口名称)<sup>1</sup>
- 全局节点名称**<sup>1</sup>
- Virtual World Wide Node Name** (虚拟全局节点名称)<sup>1</sup>
- PCI Device ID** (PCI 设备 ID)
- PCI** (总线) 地址

Main Configuration Page • NIC Partitioning Configuration • Partition 2 Configuration		
NIC Mode .....	<input type="radio"/> Enabled	<input checked="" type="radio"/> Disabled
FCoE Mode .....	<input checked="" type="radio"/> Enabled	<input type="radio"/> Disabled
FIP MAC Address .....	00:0E:1E:D5:F8:78	
Virtual FIP MAC Address .....	00:00:00:00:00:00	
World Wide Port Name .....	20:01:00:0E:1E:D5:F8:78	
Virtual World Wide Port Name .....	00:00:00:00:00:00:00	
World Wide Node Name .....	20:00:00:0E:1E:D5:F8:78	
Virtual World Wide Node Name .....	00:00:00:00:00:00:00	
PCI Device ID .....	8070	
PCI Address .....	86:02	

**图 5-18. 分区 2 配置: FCoE 卸载**

4. 要配置第三个分区, 请选择 **Partition 3 Configuration** (分区 3 配置) 以打开 Partition 3 Configuration (分区 3 配置) 页面 (图 5-17)。如果存在 iSCSI 卸载, 则 Partition 3 Configuration (分区 3 配置) 将显示以下参数:
- NIC Mode** (NIC 模式) (**Disabled**) (禁用)
  - iSCSI Offload Mode** (iSCSI 卸载模式) (**Enabled**) (启用)
  - iSCSI Offload MAC Address** (iSCSI 卸载 MAC 地址)<sup>2</sup>
  - Virtual iSCSI Offload MAC Address** (虚拟 iSCSI 卸载 MAC 地址)<sup>2</sup>

<sup>1</sup> 此参数仅存在于具有 FCoE 卸载功能的适配器的 NPAR 模式端口的第二个分区上。

<sup>2</sup> 此参数仅存在于具有 iSCSI 卸载功能的适配器的 NPAR 模式端口的第三个分区上。

- PCI Device ID** (PCI 设备 ID)
- PCI Address** (PCI 地址)

Main Configuration Page • NIC Partitioning Configuration • Partition 3 Configuration	
NIC Mode .....	<input type="radio"/> Enabled <input checked="" type="radio"/> Disabled
iSCSI Offload Mode .....	<input checked="" type="radio"/> Enabled <input type="radio"/> Disabled
iSCSI Offload MAC Address .....	00:0E:1E:D5:F8:7A
Virtual iSCSI Offload MAC Address .....	00:00:00:00:00:00
PCI Device ID .....	8070
PCI Address .....	86:04

**图 5-19. 分区 3 配置: iSCSI 卸载**

5. 要配置剩余的以太网分区, 包括以前的分区 (如果尚未启用卸载), 请打开分区 2 或更大分区 (请参阅 [图 5-20](#)) 页面。
  - NIC Mode** (NIC 模式) (**Enabled** (启用或 **Disabled** (禁用)))。禁用时, 分区会隐藏, 以便其在检测到小于最大数量的分区 (或 PCI PF) 时不会向操作系统显示。
  - PCI Device ID** (PCI 设备 ID)
  - PCI Address** (PCI 地址)
  - MAC 地址**
  - Virtual MAC Address** (虚拟 MAC 地址)

Main Configuration Page • NIC Partitioning Configuration • Partition 4 Configuration	
NIC Mode .....	<input checked="" type="radio"/> Enabled <input type="radio"/> Disabled
PCI Device ID .....	8070
PCI Address .....	86:06
MAC Address .....	00:0E:1E:D5:F8:7C
Virtual MAC Address .....	00:00:00:00:00:00

**图 5-20. 分区 4 配置: 以太网**

## VMware ESXi 6.0 和 ESXi 6.5 分区

如果在运行 VMware ESXi 6.0 或 ESXi 6.5 的系统上存在以下情况，则必须卸载并重新安装驱动程序：

- 适配器配置为启用具有所有 NIC 分区的 NPAR。
- 适配器处于单功能模式。
- 保存配置并重新启动系统。
- 系统上已经安装了驱动程序时启用存储分区（通过将 NIC 分区之一转换为存储分区）。
- 分区 2 被更改为 FCoE。
- 保存配置并再次重新启动系统。

需要重新安装驱动程序是因为存储功能可能会保留 `vmnicX` 枚举而非 `vmhbaX`，如在系统上发出以下命令时所示：

```
# esxcfg-scsidevs -a
vmnic4 qedf                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 qedf                link-up    fc.2000000e1ef249f8:2001000e1ef249f8
(0000:d8:00.2) QLogic Corp. QLogic FastLinQ QL41xxx Series 10/25 GbE
Controller (FCoE)
vmhba1 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 qedil               online     iscsi.vmhba33                QLogic
FastLinQ QL41xxx Series 10/25 GbE Controller (iSCSI)
```

在上述命令输出中，请注意 `vmnic4` 和 `vmnic10` 实际上是存储适配器端口。为防止出现这种情况，应该在为 NPAR 模式配置适配器的同时启用存储功能。

例如，假设适配器默认处于单一功能模式，您应该：

1. 启用 NPAR 模式。
2. 将分区 2 更改为 FCoE。
3. 保存并重新启动。

# 6 RoCE 配置

本章介绍 41xxx 系列适配器、以太网交换机以及 Windows 或 Linux 主机上的基于聚合以太网的 RDMA（RoCE v1 和 v2）配置，包括以下内容：

- [支持的操作系统和 OFED](#)
- [第 61 页上的“规划 RoCE”](#)
- [第 62 页上的“准备适配器”](#)
- [第 62 页上的“准备以太网交换机”](#)
- [第 64 页上的“在 Windows Server 的适配器上配置 RoCE”](#)
- [第 71 页上的“在 Linux 的适配器上配置 RoCE”](#)
- [第 81 页上的“在 VMware ESX 的适配器上配置 RoCE”](#)

## 注

某些 RoCE 功能在当前版本中可能并未完全启用。

## 支持的操作系统和 OFED

表 6-1 显示 RoCE v1、RoCE v2、iWARP 和 OFED 的操作系统支持。

表 6-1. RoCE v1, RoCE v2, iWARP 和 OFED 的操作系统支持

操作系统	内建驱动程序	OFED 3.18-3 GA	OFED-4.8-1 GA
Windows Server 2012 R2	否	N/A	N/A
Windows Server 2016	否	N/A	N/A
RHEL 6.8	RoCE v1、iWARP	RoCE v1、iWARP	否
RHEL 6.9	RoCE v1、iWARP	否	否
RHEL 7.3	RoCE v1、RoCE v2、iWARP、iSER	否	RoCE v1、RoCE v2、iWARP

表 6-1. RoCE v1, RoCE v2, iWARP 和 OFED 的操作系统支持 (续)

操作系统	内建驱动程序	OFED 3.18-3 GA	OFED-4.8-1 GA
RHEL 7.4	RoCE v1、RoCE v2、iWARP、iSER	否	否
SLES 12 SP3	RoCE v1、RoCE v2、iWARP、iSER	否	否
CentOS 7.3	RoCE v1、RoCE v2、iWARP、iSER	否	RoCE v1、RoCE v2、iWARP
CentOS 7.4	RoCE v1、RoCE v2、iWARP、iSER	否	否
VMware ESXi 6.0 u3	否	N/A	N/A
VMware ESXi 6.5、6.5U1	RoCE v1、RoCE v2	N/A	N/A
VMware ESXi 6.7	RoCE v1、RoCE v2	N/A	N/A

## 规划 RoCE

在准备实施 RoCE 时，请考虑以下限制：

- 如果使用的是内建 OFED，则服务器和客户端系统上的操作系统应该相同。一些应用程序也许可以在不同的操作系统之间正常工作，但无法保证。这是 OFED 限制。
- 对于 OFED 应用程序（最常见的是 perftest 应用程序），服务器和客户端应用程序应使用相同的选项和值。如果操作系统和 perftest 应用程序具有不同的版本，则可能会出现问題。要确认 perftest 版本，请发出以下命令：  

```
# ib_send_bw --version
```
- 在内建 OFED 中生成 libqedr 需要安装 libibverbs-devel。
- 在内建 OFED 中运行用户空间应用程序需要安装 InfiniBand® 支持组，通过 yum groupinstall，“InfiniBand 支持”包含 libibcm、libibverbs 等等。
- 依赖于 libibverbs 的 OFED 和 RDMA 应用程序也需要 QLogic RDMA 用户空间库 libqedr。使用 libqedr RPM 或源文件包安装 libqedr。
- RoCE 仅支持小字节序。
- RoCE 无法在 SR-IOV 环境中的 VF 上工作。

## 准备适配器

按照以下步骤操作以启用 DCBX，并使用 HII 管理应用程序指定 RoCE 优先级。有关 HII 应用程序的信息，请参阅[第 5 章 适配器预引导配置](#)。

### 要准备适配器：

1. 在 Main Configuration Page（主要配置页面）中，选择 **Data Center Bridging (DCB) Settings**（数据中心桥接 (DCB) 设置），然后单击 **Finish**（完成）。
2. 在 Data Center Bridging (DCB) Settings（数据中心桥接 (DCB) 设置）窗口中，单击 **DCBX Protocol**（DCBX 协议）选项。41xxx 系列适配器同时支持 CEE 协议和 IEEE 协议。此值应该符合 DCB 交换机上的相应值。在本例中，选择 **CEE** 或 **Dynamic**（动态）。
3. 在 **RoCE Priority**（RoCE 优先级）框中，键入优先级值。此值应该符合 DCB 交换机上的相应值。在本示例中，键入 5。通常，0 用于默认有损流量类，3 用于 FCoE 流量类，4 用于通过 DCB 的无损 iSCSI-TLV 流量类。
4. 单击 **Back**（后退）。
5. 看到提示时，单击 **Yes**（是）以保存更改。更改将在系统重设后生效。

对于 Windows，您可以使用 HII 或 QoS 方法配置 DCBX。本节中所示的配置为通过 HII。对于 QoS，请参阅[第 212 页上的“为 RoCE 配置 QoS”](#)。

## 准备以太网交换机

本节介绍如何为 RoCE 配置 Cisco® Nexus® 6000 以太网交换机和 Dell® Z9100 以太网交换机。

- [配置 Cisco Nexus 6000 以太网交换机](#)
- [配置 Dell Z9100 以太网交换机](#)

### 配置 Cisco Nexus 6000 以太网交换机

为 RoCE 配置 Cisco Nexus 6000 以太网交换机的步骤包括配置类映射、配置策略映射、应用策略，以及为交换机端口分配 VLAN ID。

#### 要配置 Cisco 交换机：

1. 按如下方式打开配置终端会话：

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

2. 按如下方式配置服务质量 (QoS) 类映射, 并设置 RoCE 优先级以匹配适配器 (5):

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

3. 按如下方式配置排队类映射:

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

4. 按如下方式配置网络 QoS 类映射:

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

5. 按如下方式配置 QoS 策略映射:

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

6. 配置排队策略映射以分配网络带宽。在本示例中, 使用的值为 50%:

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

7. 按如下方式配置网络 QoS 策略映射以设置无丢弃流量类的优先级流控制:

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

8. 按如下方式在系统级别应用新策略:

```
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. 为交换机端口分配 VLAN ID, 以匹配分配给适配器 (5) 的 VLAN ID。

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

## 配置 Dell Z9100 以太网交换机

要为 RoCE 配置 Dell Z9100 以太网交换机，请参阅 [附录 C Dell Z9100 交换机配置](#) 中的步骤。

## 在 Windows Server 的适配器上配置 RoCE

在 Windows Server 的适配器上配置 RoCE 涉及在适配器上启用 RoCE 以及验证 Network Direct MTU 大小。

### 要在 Windows Server 主机上配置 RoCE：

1. 在适配器上启用 RoCE。
  - a. 打开 Windows 设备管理器，然后打开 41xxx 系列适配器 NDIS 微型端口属性。
  - b. 在 QLogic FastLinQ Adapter Properties（QLogic FastLinQ 适配器属性）中，单击 **Advanced**（高级）选项卡。
  - c. 在 Advanced（高级）页面中，通过选择 **Property**（属性）下的每个项目，然后选择该项目相应的 **Value**（值），配置 [表 6-2](#) 中列出的属性。然后单击 **OK**（确定）。

**表 6-2. RoCE 的高级属性**

属性	值或说明
Network Direct Functionality (Network Direct 功能)	Enabled（启用）
Network Direct Mtu Size (Network Direct Mtu 大小)	Network Direct MTU 大小必须小于巨型数据包大小。
RDMA Mode（RDMA 模式）	<b>RoCE v1</b> 或 <b>RoCE v2</b> 。 <b>iWARP</b> 值仅适用于为 iWARP 配置端口时，如 <a href="#">第 7 章 iWARP 配置</a> 中所述。
VLAN ID	将任何 VLAN ID 分配给接口。该值必须与交换机上分配的值相同。
服务质量	启用或禁用 QoS。 <ul style="list-style-type: none"> <li>■ 如果您通过 Windows DCB-QoS 服务控制 DCB，请选择 <b>启用</b>。有关更多信息，请参阅 <a href="#">第 212 页上的“通过在适配器上禁用 DCBX 来配置 QoS”</a>。</li> <li>■ 如果您通过连接的 DCB 配置开关控制 DCB，请选择 <b>Disabled</b>（已禁用）。有关更多信息，请参阅 <a href="#">第 216 页上的“通过在适配器上启用 DCBX 来配置 QoS”</a>。</li> </ul>

图 6-1 显示配置属性值的示例。

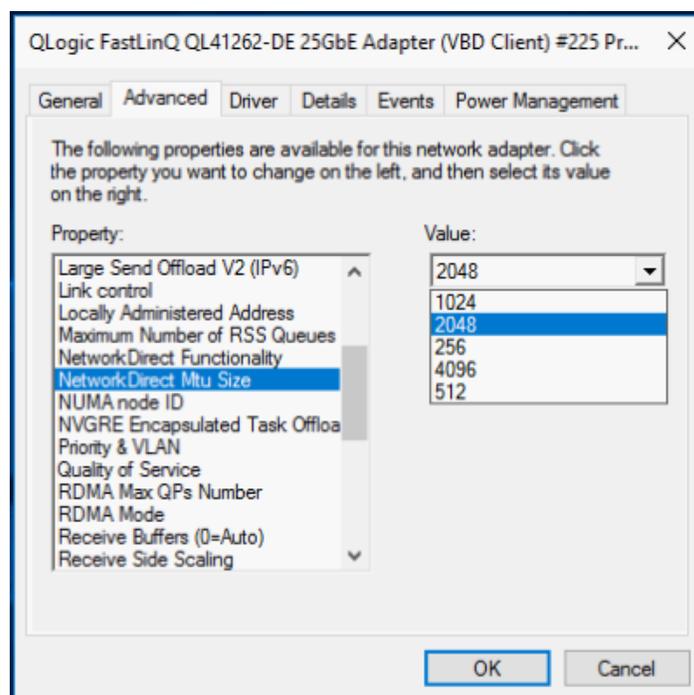


图 6-1. 配置 RoCE 属性

- 使用 Windows PowerShell，验证适配器上是否已启用 RDMA。  
Get-NetAdapterRdma 命令列出支持 RDMA 的适配器 - 两个端口同时启用。

### 注

如果您通过 Hyper-V 配置 RoCE，不要将 VLAN ID 分配给物理接口。

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

- 使用 Windows PowerShell，验证主机操作系统是否已启用 NetworkDirect。  
Get-NetOffloadGlobalSetting 命令显示 NetworkDirect 已启用。

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

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

4. 连接服务器消息块 (SMB) 驱动器，运行 RoCE 流量并验证结果。

要设置并连接到 SMB 驱动器，请查看 Microsoft 在线提供的信息：

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

5. 默认情况下，Microsoft 的 SMB Direct 为每个端口建立两个 RDMA 连接，以提供优质性能，包括较高块大小（例如，64KB）下的线速率。为优化性能，您可以将每个 RDMA 接口的 RDMA 连接数更改为四个（或更多）。

要将 RDMA 连接数增加到四个（或更多），请在 Windows PowerShell 中发出以下命令：

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

## 查看 RDMA 计数器

以下步骤也适用于 iWARP。

要查看 RoCE 的 RDMA 计数器：

1. 启动性能监视器。
2. 打开“添加计数器”对话框。图 6-2 显示一个示例。

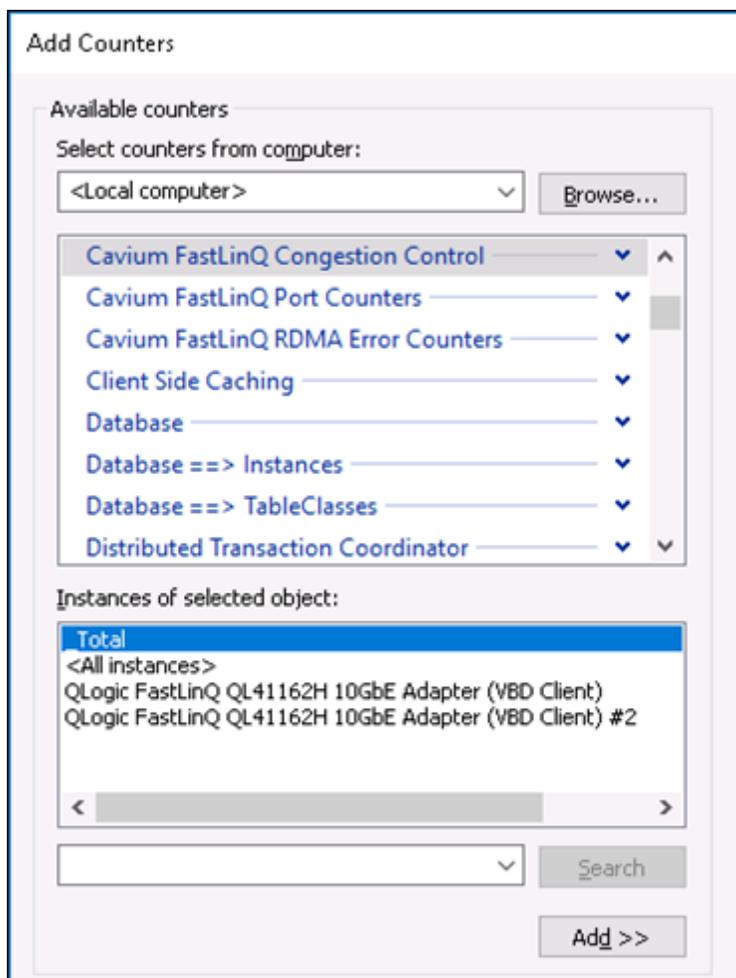


图 6-2. 添加计数器对话框

### 注

如果 Cavium RDMA 计数器没有列在性能监视器“添加计数器”对话框中，请通过从驱动程序位置发出以下命令手动添加：

```
Lodctr /M:qend.man
```

3. 选择以下计数器类型之一：
  - Cavium FastLinQ 拥塞控制：**
    - 如果网络拥塞且在开关上启用了 ECN，则增加。
    - 描述成功发送和接收的 RoCE v2 ECN 标记数据包和拥塞通知数据包 (CNP)。
    - 仅适用于 RoCE v2。
  - Cavium FastLinQ 端口计数器：**
    - 当网络拥塞时增加。
    - 如果配置了流量控制或全部暂停且网络拥塞，则暂停计数器增加。
    - 如果配置了优先流量控制且网络拥塞，则增加 PFC 计数器。
  - Cavium FastLinQ RDMA 错误计数器：**
    - 如果传输操作中出现任何错误，则增加。
    - 有关详细信息，请参阅 [表 6-3](#)。
4. 在**选定对象实例**中，选择**全部**，然后点击**添加**。

图 6-3 显示计数器监视输出的示例。

The figure consists of three screenshots of the Windows Performance Monitor tool, each displaying a different set of performance counters for the Cavium FastLinQ RDMA adapter. The screenshots are as follows:

- Top Screenshot:** Shows 'Cavium FastLinQ Congestion Control' counters. The table below summarizes the data shown.
- Middle Screenshot:** Shows 'Cavium FastLinQ Port Counters'. The table below summarizes the data shown.
- Bottom Screenshot:** Shows 'Cavium FastLinQ RDMA Error Counters'. The table below summarizes the data shown.

Counter Name	Total	QLogic FastLinQ QL41162H 10GbE Adapter (VBD Client)	QLogic FastLinQ QL41162H 10GbE Adapter (VBD Client) #2
Notification Point - CNPs Sent Successfully	0.000	0.000	0.000
Notification Point - RoCEv2 ECN Marked Packets	0.000	0.000	0.000
Reaction Point - CNPs Received Successfully	0.000	0.000	0.000

Counter Name	Total	QLogic FastLinQ QL41162H 10GbE Adapter (VBD Client)	QLogic FastLinQ QL41162H 10GbE Adapter (VBD Client) #2
Pause Frames Received	0.000	0.000	0.000
Pause Frames Transmitted	0.000	0.000	0.000
PFC Frames Received	0.000	0.000	0.000
PFC Frames Transmitted	0.000	0.000	0.000

Counter Name	Total	QLogic FastLinQ QL41162H 10GbE Adapter (VBD Client)	QLogic FastLinQ QL41162H 10GbE Adapter (VBD Client) #2
CQ Overflow	0.000	0.000	0.000
Requestor Bad Response	0.000	0.000	0.000
Requestor CQE Flushed	0.000	0.000	0.000
Requestor Local Length	0.000	0.000	0.000
Requestor Local Protection	0.000	0.000	0.000
Requestor Local QP Operation	0.000	0.000	0.000
Requestor Remote Access	0.000	0.000	0.000
Requestor Remote Invalid Request	0.000	0.000	0.000
Requestor Remote Operation	0.000	0.000	0.000
Requestor Retry Exceeded	0.000	0.000	0.000
Requestor RNR NAK Retry Exceeded	0.000	0.000	0.000
Responder CQE Flushed	0.000	0.000	0.000
Responder Local Length	0.000	0.000	0.000
Responder Local Protection	0.000	0.000	0.000
Responder Local QP Operation	0.000	0.000	0.000
Responder Remote Invalid Request	0.000	0.000	0.000

图 6-3. 性能监视器：Cavium FastLinQ 计数器

表 6-3 提供有关错误计数器的详细信息。

表 6-3. Cavium FastLinQ RDMA 错误计数器

RDMA 错误计数器	说明	适用于 RoCE?	适用于 iWARP?	故障排除
CQ 超值	一个完成队列，在上面发布 RDMA 工作请求。这个计数器规定实例的数量，存在发送或接收队列上工作请求的完成但在相关完成队列上没有空格。	是	是	指示导致完成队列不够大的软件设计问题。
请求程序错误回复	响应程序返回的错误响应。	是	是	—

表 6-3. Cavium FastLinQ RDMA 错误计数器 (续)

RDMA 错误计数器	说明	适用于 RoCE?	适用于 iWARP?	故障排除
请求程序 CQE 刷新错误	在 QP 由于任何原因转到错误状态且有待处理工作请求的情况下，发布的工作请求可以通过向 CQ 发送刷新状态的完成（未完成实际工作请求执行）来刷新。在有工作请求是在错误状态中完成的情况下，将刷新该 QP 的所有其他待处理工作请求。	是	是	在 RDMA 连接断开时发生。
请求程序本地长度	RDMA READ 响应消息包含太多或太少负载数据。	是	是	通常指示主机软件组件问题。
请求程序本地保护	本地发布的工作请求的数据段不引用对请求的操作有效的内存区域。	是	是	通常指示主机软件组件问题。
请求程序本地 QP 操作	在处理这个工作请求时检测到一个内部 QP 相容性错误。	是	是	—
请求程序远程访问	将由 RDMA Read 读取的、以 RDMA Write 写入的或通过原子操作访问的远程数据缓冲区上出现的保护错误。	是	是	—
请求程序远程无效请求	远程方在通道中收到无效消息。无效请求可以是已发送的消息或 RDMA 请求。	是	是	可能原因包括，这个接收队列不支持该操作、缓冲不足无法接收新的 RDMA 或原子操作请求、或 RDMA 请求中规定的长度超出 231 个字节。
请求程序远程操作	由于其本地问题，远程方无法完成请求的操作	是	是	远程方的软件问题（例如，导致 QP 错误或 RQ 上的错误 WQE 的问题）阻止完成操作。
请求程序重试超限	传输重试已超出最大限制	是	是	远程对方可能已停止响应，或网络问题阻止消息确认。
请求程序 RNR 重试超限	收到的 RNR NAK 导致的重试已经达到最多失败次数	是	否	远程对方可能已停止响应，或网络问题阻止消息确认。

表 6-3. Cavium FastLinQ RDMA 错误计数器 (续)

RDMA 错误计数器	说明	适用于 RoCE?	适用于 iWARP?	故障排除
请求程序 CQE 刷新	在 QP 由于任何原因转到错误状态且 RQ 上有待接收缓冲的情况下，发布的工作请求（在 RQ 上接收缓冲）可以通过向 CQ 发送刷新状态的完成来刷新。在有工作请求是在错误状态中完成的情况下，将刷新该 QP 的所有其他待处理工作请求。	是	是	—
响应程序本地长度	进站消息中的无效长度。	是	是	表现不良的远程对等方。例如，入站的发送消息的长度超出接收缓冲区大小。
响应程序本地保护	本地发布的工作请求的数据段不引用对请求的操作有效的内存区域。	是	是	指示内存管理方面的软件问题。
响应程序本地 QP 操作错误	在处理这个工作请求时检测到一个内部 QP 相容性错误。	是	是	指示软件问题。
响应程序远程无效请求	响应程序在通道中检测到无效进站消息。	是	是	指示远程对等方的可能违规行为。可能原因包括：这个接收队列不支持该操作、缓冲不足无法接收新的 RDMA 请求、或 RDMA 请求中规定的长度超出 $2^{31}$ 个字节。

## 在 Linux 的适配器上配置 RoCE

本节介绍 RHEL 和 SLES 的 RoCE 配置步骤。本节还介绍如何验证 RoCE 配置，并提供有关对 VLAN 接口使用组 ID (GID) 的一些指导。

- [RHEL 的 RoCE 配置](#)
- [SLES 的 RoCE 配置](#)
- [验证 Linux 上的 RoCE 配置](#)
- [VLAN 接口和 GID 索引值](#)
- [Linux 的 RoCE v2 配置](#)

## RHEL 的 RoCE 配置

要在适配器上配置 RoCE，必须在 RHEL 主机上安装并配置开放结构企业分发版 (OFED)。

### 要为 RHEL 准备内建 OFED：

1. 在安装或升级操作系统时，选择 InfiniBand 和 OFED 支持软件包。
2. 从 RHEL ISO 映像安装以下 RPM：

```
libibverbs-devel-x.x.x.x86_64.rpm  
(libqedr 库所需)  
perftest-x.x.x.x86_64.rpm  
(InfiniBand 带宽和延迟应用程序所需)
```

或使用 Yum，安装内建 OFED：

```
yum groupinstall "Infiniband Support"  
yum install perftest  
yum install tcl tcl-devel tk zlib-devel libibverbs  
libibverbs-devel
```

---

### 注

在安装过程中，如果您已经选择了前面提到的软件包，则不需要重新安装它们。内建 OFED 和支持软件包可能因操作系统版本而异。

---

3. 如第 13 页上的“[安装具有 RDMA 的 Linux 驱动程序](#)”中所述，安装新 Linux 驱动程序。

## SLES 的 RoCE 配置

要在 SLES 主机的适配器上配置 RoCE，必须在 SLES 主机上安装并配置 OFED。

### 要为 SLES Linux 安装内建 OFED：

1. 在安装或升级操作系统时，选择 InfiniBand 支持软件包。
2. 从相应的 SLES SDK 套件映像安装以下 RPM：

```
libibverbs-devel-x.x.x.x86_64.rpm  
(libqedr 安装所需)  
perftest-x.x.x.x86_64.rpm  
(带宽和延迟应用程序所需)
```

3. 如第 13 页上的“[安装具有 RDMA 的 Linux 驱动程序](#)”中所述，安装 Linux 驱动程序。

## 验证 Linux 上的 RoCE 配置

安装 OFED 后，安装 Linux 驱动程序，然后加载 RoCE 驱动程序，验证在所有 Linux 操作系统上是否检测到 RoCE 设备。

要在 Linux 上验证 RoCE 配置：

1. 使用 `service/systemctl` 命令停止防火墙表。
2. 仅适用于 RHEL：如果安装了 RDMA 服务 (`yum install rdma`)，请验证 RDMA 服务是否已启动。

### 注

对于 RHEL 6.x 和 SLES 11 SP4，重新引导后必须启动 RDMA 服务。  
对于 RHEL 7.x 和 SLES 12 SPX 及更高版本，RDMA 服务会在重新引导后自行启动。

在 RHEL 或 CentOS 上：使用 `service rdma` 状态命令以启动服务：

- ❑ 如果 RDMA 尚未启动，请发出以下命令：

```
# service rdma start
```

- ❑ 如果 RDMA 未启动，则发出以下备选命令之一：

```
# /etc/init.d/rdma start
```

或者

```
# systemctl start rdma.service
```

3. 通过检查 `dmesg` 日志验证是否检测到 RoCE 设备：

```
# dmesg|grep qedr
```

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

4. 验证是否已加载所有模块。例如：

```
# lsmod|grep qedr
```

```
qedr                89871  0
qede                96670  1 qedr
qed                 2075255  2 qede,qedr
ib_core             88311  16 qedr, rdma_cm, ib_cm,
                   ib_sa, iw_cm, xprtrdma, ib_mad, ib_srp,
                   ib_ucm, ib_iser, ib_srpt, ib_umad,
                   ib_uverbs, rdma_ucm, ib_ipoib, ib_isert
```

5. 使用配置方法（如 `ifconfig`）配置 IP 地址并启用端口：

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

6. 发出 `ibv_devinfo` 命令。对于每个 PCI 功能，您应该会看到一个单独的 `hca_id`，如下例中所示：

```
root@captain:~# ibv_devinfo
hca_id: qedr0
      transport:                InfiniBand (0)
      fw_ver:                    8.3.9.0
      node_guid:                 020e:1eff:fe50:c7c0
      sys_image_guid:           020e:1eff:fe50:c7c0
      vendor_id:                 0x1077
      vendor_part_id:           5684
      hw_ver:                    0x0
      phys_port_cnt:            1
      port: 1
      state:                     PORT_ACTIVE (1)
      max_mtu:                   4096 (5)
      active_mtu:               1024 (3)
      sm_lid:                    0
      port_lid:                 0
      port_lmc:                 0x00
      link_layer:               Ethernet
```

7. 验证所有服务器之间的 L2 和 RoCE 连通性：一个服务器充当服务器，另一个服务器充当客户端。

- ❑ 使用简单的 `ping` 命令验证 L2 连接。
- ❑ 通过在服务器或客户端上执行 RDMA ping 验证 RoCE 连接：

在服务器上发出以下命令：

```
ibv_rc_pingpong -d <ib-dev> -g 0
```

在客户端上发出以下命令：

```
ibv_rc_pingpong -d <ib-dev> -g 0 <server L2 IP address>
```

下面是服务器和客户端上的成功 ping pong 测试的示例。

### 服务器 Ping:

```
root@captain:~# ibv_rc_pingpong -d qedr0 -g 0
local address: LID 0x0000, QPN 0xff0000, PSN 0xb3e07e, GID
fe80::20e:1eff:fe50:c7c0
remote address: LID 0x0000, QPN 0xff0000, PSN 0x934d28, GID
fe80::20e:1eff:fe50:c570
8192000 bytes in 0.05 seconds = 1436.97 Mbit/sec
1000 iters in 0.05 seconds = 45.61 usec/iter
```

### 客户端 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:1eff:fe50:c570
remote address: LID 0x0000, QPN 0xff0000, PSN 0xb3e07e, GID
fe80::20e:1eff:fe50:c7c0
8192000 bytes in 0.02 seconds = 4211.28 Mbit/sec
1000 iters in 0.02 seconds = 15.56 usec/iter
```

■ 要显示 RoCE 统计信息，请发出以下命令，其中 *x* 是设备号：

```
> mount -t debugfs nodev /sys/kernel/debug
> cat /sys/kernel/debug/qedr/qedrX/stats
```

## VLAN 接口和 GID 索引值

如果在服务器和客户端上均使用 VLAN 接口，您还必须在交换机上配置相同的 VLAN ID。如果是通过交换机运行流量，则 InfiniBand 应用程序必须使用正确的 GID 值，该值基于 VLAN ID 和 VLAN IP 地址。

根据以下结果，应该为任意 perftest 应用程序使用 GID 值 (-x 4/-x 5)。

```
# ibv_devinfo -d qedr0 -v|grep GID
GID[ 0]: fe80:0000:0000:0000:020e:1eff:fe50:c5b0
GID[ 1]: 0000:0000:0000:0000:0000:ffff:c0a8:0103
GID[ 2]: 2001:0db1:0000:0000:020e:1eff:fe50:c5b0
GID[ 3]: 2001:0db2:0000:0000:020e:1eff:fe50:c5b0
GID[ 4]: 0000:0000:0000:0000:0000:ffff:c0a8:0b03 VLAN 接口的 IP 地址
GID[ 5]: fe80:0000:0000:0000:020e:1e00:0350:c5b0 VLAN ID 3
```

### 注

背靠背或暂停设置的默认 GID 值为零 (0)。对于服务器 / 交换机配置，您必须确定合适的 GID 值。如果使用的是交换机，请参阅相应的交换机配置说明文件以了解正确设置。

---

## Linux 的 RoCE v2 配置

要验证 RoCE v2 功能，您必须使用 RoCE v2 支持的内核。

### 要为 Linux 配置 RoCE v2：

1. 确保您使用以下支持的内核之一：
  - SLES 12 SP2 GA
  - RHEL 7.3 GA
2. 按如下方式配置 RoCE v2：
  - a. 识别 RoCE v2 的 GID 索引。
  - b. 配置服务器和客户端的路由地址。
  - c. 在交换机上启用 L3 路由。

### 注

您可以通过使用 RoCE v2 支持的内核来配置 RoCE v1 和 RoCE v2。这些内核允许您在同一子网以及不同子网诸如 RoCE v2 和任何可路由环境中运行 RoCE 流量。RoCE v2 只需少量设置，所有其他交换机和适配器设置均通用于 RoCE v1 和 RoCE v2。

---

## 识别 RoCE v2 GID 索引或地址

要查找 RoCE v1 和 RoCE v2 特定的 GID，请使用 `sys` 或 `class` 参数，或从 41xxx FastLinQ 源文件包运行 RoCE 脚本。要检查默认 **RoCE GID** 索引和地址，请发出 `ibv_devinfo` 命令并将其与 `sys` 或 `class` 参数进行比较。例如：

```
#ibv_devinfo -d qedr0 -v|grep GID
GID[ 0]:          fe80:0000:0000:0000:020e:1eff:fec4:1b20
GID[ 1]:          fe80:0000:0000:0000:020e:1eff:fec4:1b20
GID[ 2]:          0000:0000:0000:0000:0000:ffff:1e01:010a
GID[ 3]:          0000:0000:0000:0000:0000:ffff:1e01:010a
GID[ 4]:          3ffe:ffff:0000:0f21:0000:0000:0000:0004
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
```

## 通过 `sys` 和 `class` 参数验证 RoCE v1 或 v2 GID 索引和地址

使用以下选项之一，通过 `sys` 和 `class` 参数验证 RoCE v1 或 RoCE v2 GID 索引和地址：

### ■ 选项 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
```

### ■ 选项 2:

使用 FastLinQ 源文件包的脚本。

```
#/.../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
```

## 6-RoCE 配置

### 在 Linux 的适配器上配置 RoCE

---

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

---

#### 注

您必须基于服务器或交换机配置 (Pause/PFC) 指定 RoCE v1 或 RoCE v2 的 GID 索引值。为链路本地 IPv6 地址、IPv4 地址或 IPv6 地址使用 GID 索引。要将带 VLAN 标记的帧用于 RoCE 流量，您必须指定从 VLAN IPv4 或 IPv6 地址得出的 GID 索引值。

---

### 通过 perfest 应用程序验证 RoCE v1 或 RoCE v2 功能

本节展示如何通过 perfest 应用程序验证 RoCE v1 或 RoCE v2 功能。在本例中，使用以下服务器 IP 和客户端 IP：

- 服务器 IP: 192.168.100.3
- 客户端 IP: 192.168.100.4

#### 验证 RoCE v1

在同一子网上运行并使用 RoCE v1 GID 索引。

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

#### 验证 RoCE v2

在同一子网上运行并使用 RoCE v2 GID 索引。

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

---

#### 注

如果您通过交换机 PFC 配置来运行，请对通过同一子网的 RoCE v1 或 v2 使用 VLAN GID。

---

### 通过不同子网验证 RoCE v2

---

#### 注

您必须首先配置交换机和服务器的路由设置。在适配器上，使用 HII 或 UEFI 用户界面设置 RoCE 优先级和 DCBX 模式。

---

**要通过不同子网验证 RoCE v2:**

1. 使用 DCBX-PFC 配置设置服务器和客户端的路由配置。

 **系统设置:**

**服务器 VLAN IP:** 192.168.100.3 和**网关:** 192.168.100.1

**客户端 VLAN IP:** 192.168.101.3 和**网关:** 192.168.101.1

 **服务器配置:**

```
#!/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
```

 **客户端配置:**

```
#!/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. 使用以下程序设置交换机设置。

- 使用任何流控制方法（Pause、DCBX-CEE 或 DCBX-IEEE），并为 RoCE v2 启用 IP 路由。请参阅 [第 62 页上的“准备以太网交换机”](#) 以了解 RoCE v2 配置，或参阅供应商交换机文档。
- 如果您使用 PFC 配置和 L3 路由，则在使用不同子网的 VLAN 上运行 RoCE v2 流量，并使用 RoCE v2 VLAN GID 索引。

```
Server# ib_send_bw -d qedr0 -F -x 5
```

```
Client# ib_send_bw -d qedr0 -F -x 5 192.168.100.3
```

### 服务器交换机设置:

```
[root@RoCE-Auto-2 /]# ib_send_bw -d qedr0 -F -x 5 -q 2 --report_gbits
*****
* Waiting for client to connect... *
*****
-----
                Send BW Test
Dual-port      : OFF          Device       : qedr0
Number of qps  : 2           Transport type : IB
Connection type : RC         Using SRQ     : OFF
RX depth       : 512
CQ Moderation  : 100
MTU            : 1024[B]
Link type      : Ethernet
Gid index      : 5
Max inline data : 0[B]
rdma_cm QPs   : OFF
Data_ex. method : Ethernet
-----
local address: LID 0000 QPN 0xff0000 PSN 0xf0b2c3
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 QPN 0xff0002 PSN 0x124cd3
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:101:03
-----
#bytes      #iterations  BW peak[Gb/sec]  BW average[Gb/sec]  MsgRate[Mpps]
65536       1000          0.00             23.07                0.043995
-----
```

图 6-4. 交换机设置, 服务器

### 客户端交换机设置:

```
[root@roce-auto-1 ~]# ib_send_bw -d qedr0 -F -x 5 192.168.100.3 -q 2 --report_gbits
-----
                Send BW Test
Dual-port      : OFF          Device       : qedr0
Number of qps  : 2           Transport type : IB
Connection type : RC         Using SRQ     : OFF
TX depth       : 128
CQ Moderation  : 100
MTU            : 1024[B]
Link type      : Ethernet
Gid index      : 5
Max inline data : 0[B]
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 QPN 0xff0000 PSN 0xf0b2c3
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:100:03
remote address: LID 0000 QPN 0xff0002 PSN 0xa2b8f1
GID: 00:00:00:00:00:00:00:00:00:00:255:255:192:168:100:03
-----
#bytes      #iterations  BW peak[Gb/sec]  BW average[Gb/sec]  MsgRate[Mpps]
65536       1000          23.04            23.04                0.043936
-----
```

图 6-5. 交换机设置, 客户端

## 为 RDMA\_CM 应用程序配置 RoCE v1 或 RoCE v2 设置

要配置 RoCE，请使用 FastLinQ 源文件包的以下脚本：

```
# ./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
```

服务器设置：

```
[root@RoCE-Auto-2 /]# rping -s -v -C 10
server ping data: rdma-ping-0: ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-1: BCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-2: CDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-3: DEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-4: EFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-5: FGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-6: GHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
server ping data: rdma-ping-7: HIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
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 /]#
```

图 6-6. 配置 RDMA\_CM 应用程序：服务器

客户端设置：

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

图 6-7. 配置 RDMA\_CM 应用程序：客户端

## 在 VMware ESX 的适配器上配置 RoCE

本节提供 RoCE 配置的以下步骤和信息：

- 配置 RDMA 接口
- 配置 MTU

- [RoCE 模式和统计信息](#)
- [配置半虚拟 RDMA 设备 \(PVRDMA\)](#)

## 配置 RDMA 接口

要配置 RDMA 接口：

1. 安装 QLogic NIC 和 RoCE 驱动程序。
2. 使用模块参数，通过发出以下命令，从 NIC 驱动程序启用 RoCE 功能：

```
esxcfg-module -s 'enable_roce=1' qedentv
```

要应用该更改，请重新加载 NIC 驱动程序或重新启动系统。

3. 要查看 NIC 接口列表，请发出 `esxcfg-nics -l` 命令。例如：

```
esxcfg-nics -l
```

Name	PCI	Driver	Link	Speed	Duplex	MAC Address	MTU	Description
Vmnic0	0000:01:00.2	qedentv	Up	25000Mbps	Full	a4:5d:36:2b:6c:92	1500	QLogic Corp. QLogic FastLinQ QL41xxx 1/10/25 GbE Ethernet Adapter
Vmnic1	0000:01:00.3	qedentv	Up	25000Mbps	Full	a4:5d:36:2b:6c:93	1500	QLogic Corp. QLogic FastLinQ QL41xxx 1/10/25 GbE Ethernet Adapter

4. 要查看 RDMA 设备列表，请发出 `esxcli rdma device list` 命令。例如：

```
esxcli rdma device list
```

Name	Driver	State	MTU	Speed	Paired Uplink	Description
vmrdma0	qedrntv	Active	1024	25 Gbps	vmnic0	QLogic FastLinQ QL45xxx RDMA Interface
vmrdma1	qedrntv	Active	1024	25 Gbps	vmnic1	QLogic FastLinQ QL45xxx RDMA Interface

5. 要创建新的虚拟交换机，请发出以下命令：

```
esxcli network vswitch standard add -v <new vswitch name>
```

例如：

```
# esxcli network vswitch standard add -v roce_vs
```

这样将创建一个名为 `roce_vs` 的新虚拟交换机。

6. 要将 QLogic NIC 端口与 vSwitch 关联，请发出以下命令：

```
# esxcli network vswitch standard uplink add -u <uplink device> -v <roce vswitch>
```

例如：

```
# esxcli network vswitch standard uplink add -u vmnic0 -v roce_vs
```

7. 要在此 vSwitch 上创建新的端口组，请发出以下命令：

```
# esxcli network vswitch standard portgroup add -p roce_pg -v roce_vs
```

例如：

```
# esxcli network vswitch standard portgroup add -p roce_pg -v roce_vs
```

8. 要在此端口组上创建 vmknic 接口并配置 IP，请发出以下命令：

```
# esxcfg-vmknic -a -i <IP address> -n <subnet mask> <roce port group name>
```

例如：

```
# esxcfg-vmknic -a -i 192.168.10.20 -n 255.255.255.0 roce_pg
```

9. 要配置 VLAN ID，请发出以下命令：

```
# esxcfg-vswitch -v <VLAN ID> -p roce_pg
```

要使用 VLAN ID 运行 RoCE 流量，请在相应的 VMkernel 端口组上配置 VLAN ID。

## 配置 MTU

要修改 RoCE 接口的 MTU，请更改相应 vSwitch 的 MTU。通过发出以下命令，基于 vSwitch 的 MTU 设置 RDMA 接口的 MTU 大小：

```
# esxcfg-vswitch -m <new MTU> <RoCE vswitch name>
```

例如：

```
# esxcfg-vswitch -m 4000 roce_vs
```

```
# esxcli rdma device list
```

Name	Driver	State	MTU	Speed	Paired Uplink	Description
vmrdma0	qedrntv	Active	2048	25 Gbps	vmnic0	QLogic FastLinQ QL45xxx RDMA Interface
vmrdma1	qedrntv	Active	1024	25 Gbps	vmnic1	QLogic FastLinQ QL45xxx RDMA Interface

## RoCE 模式和统计信息

对于 RoCE 模式，ESXi 需要同时支持 RoCE v1 和 v2。在队列对创建期间决定使用哪种 RoCE 模式。在注册和初始化期间，ESXi 驱动程序将广告两种模式。要查看 RoCE 统计信息，请发出以下命令：

```
# 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
```

## 配置半虚拟 RDMA 设备 (PVRDMA)

要使用 vCenter 接口配置 PVRDMA:

1. 请按如下方式创建并配置新的分布式虚拟交换机:
  - a. 在 VMware vSphere Web Client 中, 右键单击 Navigator (导航) 窗口左侧窗格中的 **RoCE** 节点。
  - b. 在 Actions (操作) 菜单上, 指向 **Distributed Switch** (分布式交换机), 然后单击 **New Distributed Switch** (新的分布式交换机)。
  - c. 选择 6.5.0 版本。
  - d. 在 **New Distributed Switch** (新的分布式交换机) 下, 单击 **Edit settings** (编辑设置), 然后配置以下内容:
    - **Number of uplinks** (上行链路数)。选择适当的值。
    - **Network I/O Control** (网络 I/O 控制)。选择 **Disabled** (禁用)。
    - **Default port group** (默认端口组)。选中此复选框。
    - **Port group name** (端口组名称)。输入端口组的名称。

图 6-8 显示示例。

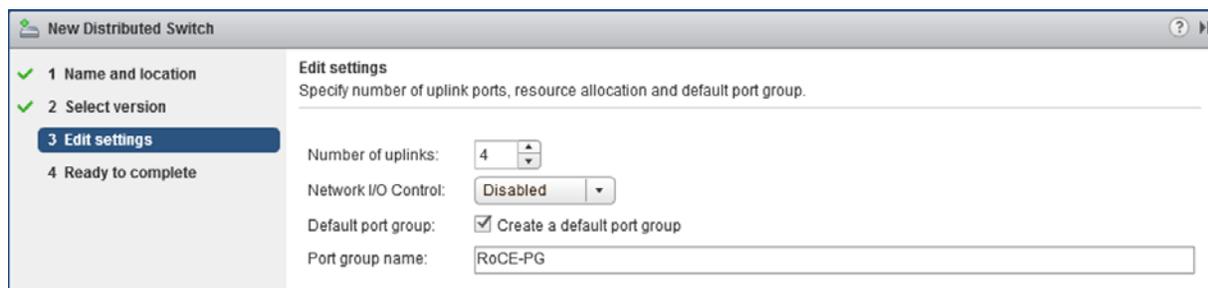


图 6-8. 配置新的分布式交换机

2. 请按如下方式配置分布式虚拟交换机:
  - a. 在 VMware vSphere Web Client 中, 展开 Navigator (导航) 窗口左侧窗格中的 **RoCE** 节点。
  - b. 右键单击 **RoCE-VDS**, 然后单击 **Add and Manage Hosts** (添加和管理主机)。
  - c. 在 **Add and Manage Hosts** (添加和管理主机) 下, 配置以下内容:
    - **Assign uplinks** (分配上行链路)。从可用上行链路列表中选择。
    - **Manage VMkernel network adapters** (管理 Vmkernel 网络适配器)。接受默认设置, 然后单击 **Next** (下一步)。

- **Migrate VM networking**（迁移 VM 网络）。分配在 [步骤 1](#) 中创建的端口组。
3. 为 PVRDMA 分配一个在 ESX 主机上使用的 vmknic:
    - a. 右键单击主机，然后选择 **Settings**（设置）。
    - b. 在 Settings（设置）页面中，展开 **System**（系统）节点，然后单击 **Advanced System Settings**（高级系统设置）。
    - c. Advanced System Settings（高级系统设置）页面显示密钥对的值及其摘要。单击 **Edit**（编辑）。
    - d. 在 Edit Advanced System Settings（编辑高级系统设置）页面中，在 **PVRDMA** 上筛选以将所有设置缩小为只有 Net.PVRDMAvmknic。
    - e. 将 **Net.PVRDMAvmknic** 值设置为 **vmknic**；例如，**vmk1**。

图 6-9 显示示例。

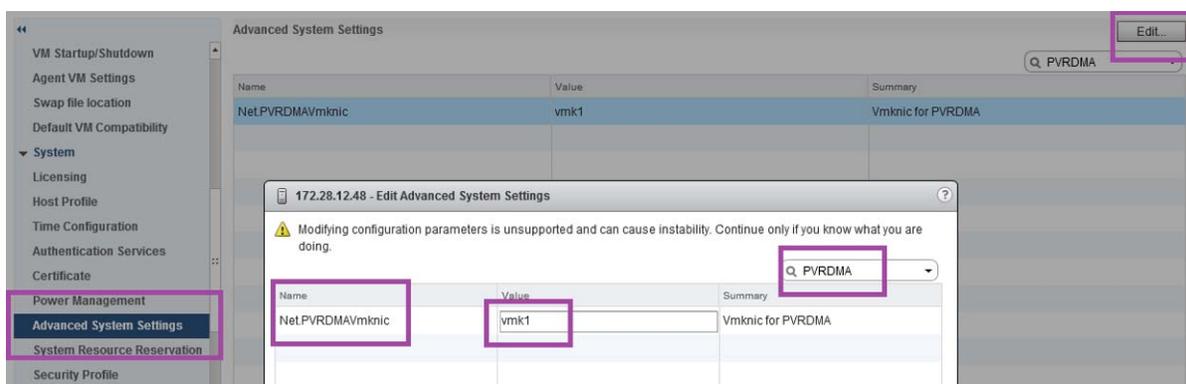


图 6-9. 为 PVRDMA 分配 vmknic

4. 设置 PVRDMA 的防火墙规则:
  - a. 右键单击主机，然后选择 **Settings**（设置）。
  - b. 在 Settings（设置）页面中，展开 **System**（系统）节点，然后单击 **Security Profile**（安全配置文件）。
  - c. 在 Firewall Summary（防火墙摘要）页面中，单击 **Edit**（编辑）。

- d. 在 **Name**（名称）下的 Edit Security Profile（编辑安全配置文件）对话框中，向下滚动，选中 **pvrDMA** 复选框，然后选中 **Set Firewall**（设置防火墙）复选框。

图 6-10 显示示例。

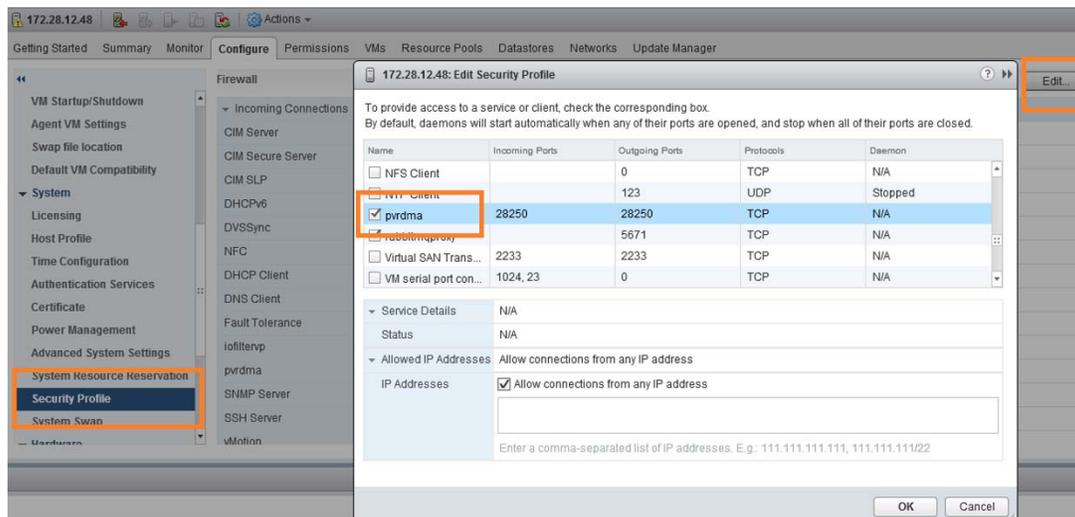


图 6-10. 设置防火墙规则

5. 请按如下方式设置 PVRDMA 的 VM：
  - a. 安装以下支持的客户操作系统之一：
    - RHEL 7.2
    - Ubuntu 14.04（内核版本 4.0）
  - b. 安装 OFED-3.18。
  - c. 编译并安装 PVRDMA 客户驱动程序和库。
  - d. 请按如下方式将新 PVRDMA 网络适配器添加至 VM：
    - 编辑 VM 设置。
    - 添加新的网络适配器。
    - 选择新添加的 DVS 端口组作为 **Network**（网络）。
    - 选择 **PVRDMA** 作为适配器类型。
  - e. 引导 VM 后，确保加载了 PVRDMA 客户驱动程序。

# 7 iWARP 配置

互联网广域 RDMA 协议 (iWARP) 是一种实现 RDMA 的计算机联网协议，用于通过 IP 网络进行有效的数据传输。iWARP 设计用于多种环境，包括 LAN、存储网络、数据中心网络和 WAN。

本章提供以下内容的说明：

- [为 iWARP 准备适配器](#)
- [第 89 页上的“在 Windows 上配置 iWARP”](#)
- [第 92 页上的“在 Linux 上配置 iWARP”](#)

## 注

某些 iWARP 功能在当前版本中可能并未完全启用。有关详细信息，请参阅 [附录 D 功能约束](#)。

## 为 iWARP 准备适配器

本节提供使用 HII 进行预引导适配器 iWARP 配置的说明。有关预引导适配器配置的更多信息，请参阅 [第 5 章 适配器预引导配置](#)。

**在默认模式下通过 HII 配置 iWARP：**

1. 访问服务器 BIOS System Setup（BIOS 系统设置），然后单击 **Device Settings**（设备设置）。
2. 在 Device Settings（设备设置）页面中，选择用于 25G 41xxx 系列适配器的端口。
3. 在所选适配器的 Main Configuration Page（主要配置页面）中，单击 **NIC Configuration**（NIC 配置）。
4. 在 NIC Configuration（NIC 配置）页面中：
  - a. 将 **NIC + RDMA Mode**（NIC + RDMA 模式）设置为 **Enabled**（已启用）。
  - b. 将 **RDMA Protocol Support**（RDMA 协议支持）设置为 **iWARP**。
  - c. 单击 **Back**（后退）。

5. 在 Main Configuration Page（主要配置页面）中，单击 **Finish**（完成）。
6. 在 Warning - Saving Changes（警告 - 保存更改）消息框中，单击 **Yes**（是）保存配置。
7. 在 Success - Saving Changes（成功 - 保存更改）消息框中，单击 **OK**（确定）。
8. 重复 [步骤 2](#) 至 [步骤 7](#) 配置其他端口的 NIC 和 iWARP。
9. 要完成两个端口的适配器准备：
  - a. 在 Device Settings（设备设置）页面中，单击 **Finish**（完成）。
  - b. 在主菜单中，单击 **Finish**（完成）。
  - c. 退出以重新引导系统。

继续 [第 89 页上的“在 Windows 上配置 iWARP”](#) 或 [第 92 页上的“在 Linux 上配置 iWARP”](#)。

## 在 Windows 上配置 iWARP

本节提供启用 iWARP、验证 RDMA 以及在 Windows 上验证 iWARP 流量的步骤。有关支持 iWARP 的操作系统列表，请参阅 [第 60 页上表 6-1](#)。

### 要在 Windows 主机上启用 iWARP 并验证 RDMA：

1. 在 Windows 主机上启用 iWARP。
  - a. 打开 Windows 设备管理器，然后打开 41xxx 系列适配器 NDIS 微型端口属性。
  - b. 在 FastLinQ 适配器属性中，单击 **Advanced**（高级）选项卡。
  - c. 在 Advanced（高级）页面的 **Property**（属性）下，进行以下操作：
    - 选择 **Network Direct Functionality**（Network Direct 功能），然后为 **Value**（值）选择 **Enabled**（已启用）。
    - 选择 **RDMA Mode**（RDMA 模式），然后为 **Value**（值）选择 **iWARP**。
  - d. 单击 **OK**（确定）保存您的更改并关闭适配器属性。

2. 使用 Windows PowerShell，验证是否已启用 RDMA。Get-NetAdapterRdma 命令输出（图 7-1）显示支持 RDMA 的适配器。

```
[172.28.41.178]: PS C:\Users\Administrator\Documents> Get-NetAdapterRdma
Name                               InterfaceDescription           Enabled
----                               -
SLOT 2 4 Port 2                    QLogic FastLinQ QL41262-DE 25GbE Adap... True
SLOT 2 3 Port 1                    QLogic FastLinQ QL41262-DE 25GbE Adap... True
```

图 7-1. Windows PowerShell 命令: Get-NetAdapterRdma

3. 使用 Windows PowerShell，验证是否已启用 NetworkDirect。Get-NetOffloadGlobalSetting 命令输出（图 7-2）显示 NetworkDirect 为 Enabled（已启用）。

```
PS C:\Users\Administrator> Get-NetOffloadGlobalSetting
ReceiveSideScaling           : Enabled
ReceiveSegmentCoalescing    : Enabled
Chimney                       : Disabled
TaskOffload                  : Enabled
NetworkDirect                 : Enabled
NetworkDirectAcrossIPSubnets : Blocked
PacketCoalescingFilter      : Disabled
```

图 7-2. Windows PowerShell 命令: Get-NetOffloadGlobalSetting

要验证 iWARP 流量:

1. 映射 SMB 驱动器，然后运行 iWARP 流量。
2. 启动性能监视器 (Perfmon)。
3. 在 Add Counters（添加计数器）对话框中，单击 **RDMA Activity**（RDMA 活动），然后选择适配器实例。

图 7-3 显示示例。

## 7-iWARP 配置 在 Windows 上配置 iWARP

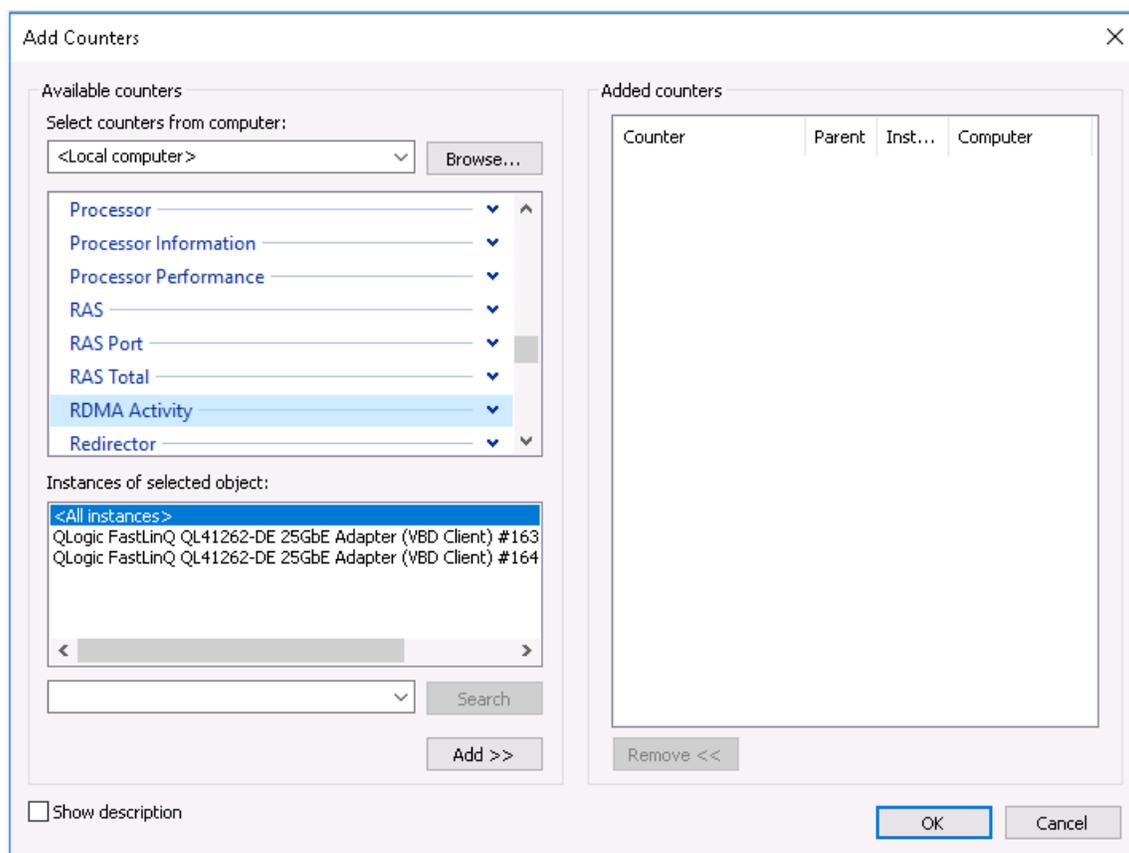


图 7-3. Perfmon: 添加计数器

如果 iWARP 流量正在运行，计数器将显示为如图 7-4 示例中所示。

Counter	QLogic FastLinQ QL41262-DE 25GbE Adapter (VBD Client) #41	QLogic FastLinQ QL41262-DE 25GbE Adapter (VBD Client) #42
RDMA Accepted Connections	0.000	0.000
RDMA Active Connections	4.000	4.000
RDMA Completion Queue Errors	0.000	0.000
RDMA Connection Errors	0.000	0.000
RDMA Failed Connection Attempts	0.000	0.000
RDMA Inbound Bytes/sec	154,164,376.212	155,971,031.049
RDMA Inbound Frames/sec	2,066,093.714	2,089,855.225
RDMA Initiated Connections	4.000	4.000
RDMA Outbound Bytes/sec	2,984,375.513	3,028,781,868
RDMA Outbound Frames/sec	2,049,329.925	2,079,912.285

图 7-4. Perfmon: 验证 iWARP 流量

**注**

有关如何在 Windows 中查看 Cavium RDMA 计数器的信息，请参阅 [第 67 页上的“查看 RDMA 计数器”](#)。

4. 要验证 SMB 连接:

- a. 在命令提示符处，请如下发出 `net use` 命令:

```
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. 请如下发出 `net -xan` 命令，其中 `Share1` 映射为 SMB 共享:

```
C:\Users\Administrator> net -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    0
```

## 在 Linux 上配置 iWARP

QLogic 41xxx 系列适配器在 [第 60 页上表 6-1](#) 中列出的 Linux 开放结构企业分布 (OFED) 上支持 iWARP。

Linux 系统中的 iWARP 配置包括以下各项:

- [安装驱动程序](#)
- [配置 iWARP 和 RoCE](#)
- [检测设备](#)

- 支持的 iWARP 应用程序
- 为 iWARP 运行 Perfest
- 配置 NFS-RDMA
- SLES 12 SP3、RHEL 7.4 和 OFED 4.8x 上的 iWARP RDMA-Core 支持

## 安装驱动程序

如 [第 3 章 驱动程序安装](#) 中所示，安装 RDMA 驱动程序。

## 配置 iWARP 和 RoCE

### 注

只有之前在使用 HII 进行预引导配置期间选择了 **iWARP + RoCE** 作为 RDMA 协议支持参数的值时，此步骤才适用（请参阅 [页 45](#) 上的 [配置 NIC 参数](#)，[步骤 5](#)）。

### 要启用 iWARP 和 RoCE：

1. 卸载所有 FastlinQ 驱动程序
2. 使用以下命令语法，通过加载具有端口接口 PCI ID (`xx:xx.x`) 和 RDMA 协议值 (`p`) 的 `qed` 驱动程序来更改 RDMA 协议。

```
#modprobe -r qedr or modprobe -r qede
```

RDMA 协议 (`p`) 值如下：

- 0 - 接受默认值 (RoCE)
- 1 - 无 RDMA
- 2 - RoCE
- 3 - iWARP

例如，发出以下命令，将 04:00.0 提供的端口从 RoCE 更改为 iWARP。

```
#modprobe -v qed rdma_protocol_map=04:00.0-3
```

3. 通过发出以下命令，加载 RDMA 驱动程序：

```
#modprobe -v qedr
```

以下示例显示用于在多个 NPAR 接口上将 RDMA 协议更改为 iWARP 的命令条目：

```
# 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)
```

## 检测设备

### 要检测设备:

1. 要验证是否检测到 RDMA 设备, 请查看 dmesg 日志:

```
# dmesg |grep qedr
[10500.191047] qedr 0000:04:00.0: registered qedr0
[10500.221726] qedr 0000:04:00.1: registered qedr1
```

2. 发出 `ibv_devinfo` 命令, 然后验证传输类型。

如果该命令成功, 则每个 PCI 功能将显示单独的 `hca_id`。例如 (如果检查上述双端口适配器的第二个端口):

```
[root@localhost ~]# ibv_devinfo -d qedr1
hca_id: qedr1
    transport:                iWARP (1)
    fw_ver:                   8.14.7.0
    node_guid:                 020e:1eff:fec4:c06e
    sys_image_guid:           020e:1eff:fec4:c06e
    vendor_id:                 0x1077
    vendor_part_id:           5718
    hw_ver:                    0x0
    phys_port_cnt:            1
        port: 1
            state:             PORT_ACTIVE (4)
            max_mtu:           4096 (5)
            active_mtu:        1024 (3)
            sm_lid:            0
            port_lid:          0
            port_lmc:          0x00
            link_layer:        Ethernet
```

## 支持的 iWARP 应用程序

适用于 iWARP 的 Linux 支持的 RDMA 应用程序包括以下各项：

- `ibv_devinfo`、`ib_devices`
- `ib_send_bw/lat`、`ib_write_bw/lat`、`ib_read_bw/lat`、`ib_atomic_bw/lat`  
对于 iWARP，所有应用程序必须通过 `-R` 选项使用 RDMA 通信管理器 (`rdma_cm`)。
- `rdma_server`、`rdma_client`
- `rdma_xserver`、`rdma_xclient`
- `rping`
- RDMA 上的 NFS (NFSoverRDMA)
- iSER（有关详细信息，请参阅 [第 8 章 iSER 配置](#)）。
- NVMe-oF（有关详细信息，请参阅 [第 12 章 使用 RDMA 进行 NVMe-oF 配置](#)）。

## 为 iWARP 运行 Perftest

所有 perftest 工具均通过 iWARP 传输类型受支持。您必须使用 RDMA 连接管理器（通过 `-R` 选项）运行工具。

### 示例：

1. 在一台服务器上，发出以下命令（本例中使用第二个端口）：

```
# ib_send_bw -d qedr1 -F -R
```

2. 在一个客户端上，发出以下命令（本例中使用第二个端口）：

```
[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      : Ethernet  
GID index      : 0  
Max inline data : 0[B]  
rdma_cm QPs    : ON  
Data ex. method : rdma_cm
```

```
-----  
local address: LID 0000 QPN 0x0192 PSN 0xcde932  
GID: 00:14:30:196:192:110:00:00:00:00:00:00:00:00:00:00:00  
remote address: LID 0000 QPN 0x0098 PSN 0x46fffc  
GID: 00:14:30:196:195:62:00:00:00:00:00:00:00:00:00:00:00  
-----
```

```
-----  
#bytes      #iterations    BW peak[MB/sec]    BW average[MB/sec]    MsgRate[Mpps]  
65536      1000           2250.38            2250.36                0.036006  
-----
```

### 注

对于延迟应用程序（发送 / 写入），如果 `perftest` 版本为最新（例如，`perftest-3.0-0.21.g21dc344.x86_64.rpm`），请使用支持的内嵌大小值：0 - 128。

---

## 配置 NFS-RDMA

适用于 iWARP 的 NFS-RDMA 包括服务器和客户端配置步骤。

### 要配置 NFS 服务器：

1. 在 `/etc/exports` 文件中，对您必须使用服务器上 NFS-RDMA 导出的目录，创建以下条目：

```
/tmp/nfs-server *(fsid=0,async,insecure,no_root_squash)
```

确保对您导出的每个目录使用不同的文件系统标识 (FSID)。

2. 按如下方式加载 `svcrdma` 模块：

```
# modprobe svcrdma
```

3. 启动 NFS 服务而不会发生任何错误：

```
# service nfs start
```

4. 请如下在此文件中包括默认 RDMA 端口 20049：

```
# echo rdma 20049 > /proc/fs/nfsd/portlist
```

5. 要使本地目录可供 NFS 客户端进行装载，请如下发出 `exportfs` 命令：

```
# exportfs -v
```

### 要配置 NFS 客户端：

#### 注

NFS 客户端配置的步骤也适用于 RoCE。

---

1. 请如下加载 `xprtrdma` 模块：

```
# modprobe xprtrdma
```

2. 根据您的版本，装载 NFS 文件系统：

对于 NFS 版本 3：

```
#mount -o rdma,port=20049 192.168.2.4:/tmp/nfs-server /tmp/nfs-client
```

对于 NFS 版本 4：

```
#mount -t nfs4 -o rdma,port=20049 192.168.2.4://tmp/nfs-client
```

#### 注

NFSv4 的默认端口为 20049。但是，与 NFS 客户端对齐的任何其他端口也将起作用。

---

3. 通过发出 `mount` 命令，验证文件系统已装载。确保 RDMA 端口和文件系统版本正确无误。

```
#mount |grep rdma
```

## SLES 12 SP3、RHEL 7.4 和 OFED 4.8x 上的 iWARP RDMA-Core 支持

用户空间库 `libqedr` 是 `rdma-core` 的一部分。但是，非内建 `libqedr` 不支持 SLES 12 SP3、RHEL 7.4、OFED 4.8x。因此，这些操作系统版本需要补丁才能支持 iWARP RDMA-Core。

### 要应用 iWARP RDMA-Core 补丁：

1. 要下载最新 RDMA-core 源，请发出以下命令：

```
# git clone https://github.com/linux-rdma/rdma-core.git
```

2. 按照 *RDMA-Core README* 中所述安装所有与操作系统有关的软件包 / 库。

对于 RHEL 和 CentOS，请发出以下命令：

```
# yum install cmake gcc libnl3-devel libudev-devel make  
pkgconfig valgrind-devel
```

对于 SLES 12 SP3 (ISO/SDK 工具包), 请安装以下 RPM:

```
cmake-3.5.2-18.3.x86_64.rpm (OS ISO)
libnl-1_1-devel-1.1.4-4.21.x86_64.rpm (SDK ISO)
libnl3-devel-3.2.23-2.21.x86_64.rpm (SDK ISO)
```

3. 要构建 RDMA-core, 请发出以下命令:

```
# cd <rdma-core-path>/rdma-core-master/
# ./build.sh
```

4. 要从当前 RDMA-core-master 位置运行所有 OFED 应用程序, 请发出以下命令:

```
# ls <rdma-core-master>/build/bin
cmpost  ib_acme          ibv_devinfo          ibv_uc_pingpong
iwpmc   rdma_client  rdma_xclient  rping          ucmatose
umad_compile_test  cmtime  ibv_asyncwatch  ibv_rc_pingpong
ibv_ud_pingpong  mckey  rdma-ndd    rdma_xserver  rstream
udaddy   umad_reg2  ibacm   ibv_devices   ibv_srq_pingpong
ibv_xsrq_pingpong  rcopy  rdma_server  riostream
srp_daemon  udpong   umad_register2
```

从当前 RDMA-core-master 位置运行应用程序。例如:

```
# ./rping -c -v -C 5 -a 192.168.21.3
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
client DISCONNECT EVENT...
```

5. 要运行内建 OFED 应用程序, 例如 perftest 和其他 InfiniBand 应用程序, 请发出以下命令以设置 iWARP 的库路径:

```
# export
LD_LIBRARY_PATH=/builds/rdma-core-path-iwarp/rdma-core-master/build/lib
```

例如:

```
# /usr/bin/rping -c -v -C 5 -a 192.168.22.3??rping -c -v -C 5 -a
192.168.22.3
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
client DISCONNECT EVENT...
```

# 8 iSER 配置

本章提供为 Linux（RHEL 和 SLES）和 ESXi 6.7 配置 RDMA 的 iSCSI 扩展 (iSER) 的步骤，包括：

- [准备工作](#)
- [为 RHEL 配置 iSER](#)
- [第 104 页上的“为 SLES 12 配置 iSER”](#)
- [第 105 页上的“在 RHEL 和 SLES 上通过 iWARP 使用 iSER”](#)
- [第 106 页上的“优化 Linux 性能”](#)
- [第 108 页上的“在 ESXi 6.7 上配置 iSER”](#)

## 准备工作

在准备配置 iSER 时，请考虑以下事项：

- 仅在以下操作系统的内建 OFED 中支持 iSER：
  - RHEL 7.1 和 7.2
  - SLES 12 和 12 SP1
- 在登录目标之后或运行 I/O 流量时，卸载 Linux RoCE qedr 驱动程序可能导致系统崩溃。
- 在运行 I/O 时，执行接口关闭 / 上线测试或执行电缆拉力测试可能会导致驱动程序或 iSER 模块错误，而这些错误可能导致系统崩溃。如果发生这种情况，请重新引导系统。

## 为 RHEL 配置 iSER

要为 RHEL 配置 iSER：

1. 如[第 72 页上的“RHEL 的 RoCE 配置”](#)中所述，安装内建 OFED。iSER 不支持非内建 OFED，因为 `ib_isert` 模块在非内建 OFED 3.18-2 GA/3.18-3 GA 版本中不可用。内建 `ib_isert` 模块不能与任何非内建 OFED 版本一起使用。

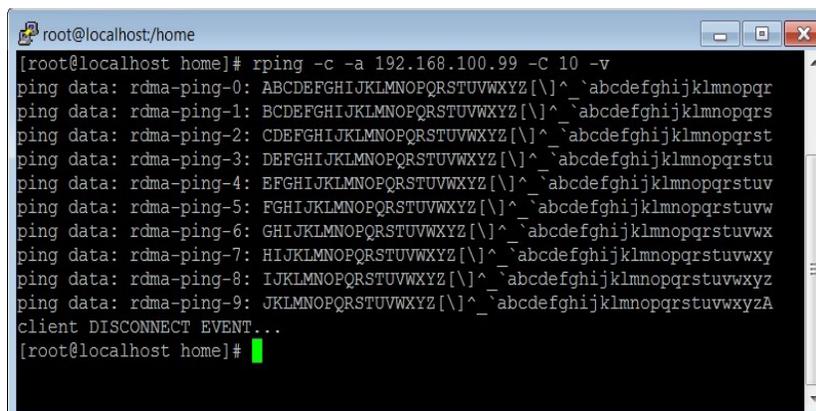
2. 如第 9 页上的“移除 Linux 驱动程序”中所述，卸载任何现有 FastLinQ 驱动程序。
3. 如第 13 页上的“安装具有 RDMA 的 Linux 驱动程序”中所述，安装最新 FastLinQ 驱动程序和 libqedr 软件包。
4. 如下加载 RDMA 服务：

```
systemctl start rdma
modprobe qedr
modprobe ib_iser
modprobe ib_isert
```
5. 通过发出 `lsmod | grep qed` 和 `lsmod | grep iser` 命令，验证在启动器和目标设备上加载的所有 RDMA 和 iSER 模块。
6. 如 74 页上的步骤 6 中所示，通过发出 `ibv_devinfo` 命令，验证存在独立的 `hca_id` 实例。
7. 检查启动器设备和目标设备上的 RDMA 连接。
  - a. 在启动器设备上，发出以下命令：

```
rping -s -C 10 -v
```
  - b. 在目标设备上，发出以下命令：

```
rping -c -a 192.168.100.99 -C 10 -v
```

图 8-1 显示了成功 RDMA ping 操作的示例。



```
root@localhost/home
[root@localhost home]# rping -c -a 192.168.100.99 -C 10 -v
ping data: rdma-ping-0: ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-1: BCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-2: CDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-3: DEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-4: EFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-5: FGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-6: GHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-7: HIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-8: IJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
ping data: rdma-ping-9: JKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyzA
Client DISCONNECT EVENT...
[root@localhost home]#
```

图 8-1. RDMA Ping 操作成功

8. 可以使用 Linux TCM-LIO 目标来测试 iSER。其设置与任何 iSCSI 目标相同，但要在适用门户上发出命令 `enable_iser Boolean=true`。门户实例在图 8-2 中标识为 `iser`。

```
/iscsi/iqn.20.../tpg1/portals> cd 192.168.100.99:3260
/iscsi/iqn.20...8.100.99:3260> enable_iser boolean=true
iSER enable now: True
/iscsi/iqn.20...8.100.99:3260>
/iscsi/iqn.20...8.100.99:3260> cd /
/> ls
o- / ..... [..]
  o- backstores ..... [..]
    | o- block ..... [Storage Objects: 0]
    | o- fileio ..... [Storage Objects: 0]
    | o- pscsi ..... [Storage Objects: 0]
    | o- ramdisk ..... [Storage Objects: 1]
    |   o- raml ..... [nullio (512.0MiB) activated]
  o- iscsi ..... [Targets: 1]
    | o- iqn.2015-06.test.target1 ..... [TPGs: 1]
    |   o- tpg1 ..... [gen-acls, no-auth]
    |     o- acls ..... [ACLs: 0]
    |     o- luns ..... [LUNs: 1]
    |       | o- lun0 ..... [ramdisk/raml]
    |     o- portals ..... [Portals: 1]
    |       o- 192.168.100.99:3260 ..... [iser]
  o- loopback ..... [Targets: 0]
  o- srpt ..... [Targets: 0]
/>
```

图 8-2. iSER 门户实例

9. 使用 `yum install iscsi-initiator-utils` 命令安装 Linux iSCSI 启动器公用程序。
  - a. 要查找 iSER 目标，请发出 `iscsiadm` 命令。例如：  
`iscsiadm -m discovery -t st -p 192.168.100.99:3260`
  - b. 要将传输模式更改为 iSER，请发出 `iscsiadm` 命令。例如：  
`iscsiadm -m node -T iqn.2015-06.test.target1 -o update -n iface.transport_name -v iser`
  - c. 要连接到或登录到 iSER 目标，请发出 `iscsiadm` 命令。例如：  
`iscsiadm -m node -l -p 192.168.100.99:3260 -T iqn.2015-06.test.target1`
  - d. 确认目标连接中的 `Iface Transport` 为 `iser`，如图 8-3 中所示。发出 `iscsiadm` 命令；例如：

```
iscsiadm -m session -P2
```

```
[root@localhost ~]# iscsiadm -m discovery -t st -p 192.168.100.99:3260
192.168.100.99:3260,1 iqn.2015-06.test.target1
192.168.100.99:3260,1 iqn.2015-06.test.target1
[root@localhost ~]#
[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
Logging in to [iface: default, target: iqn.2015-06.test.target1, portal: 192.168.100.99,3260] (multiple)
Login to [iface: default, target: iqn.2015-06.test.target1, portal: 192.168.100.99,3260] successful.
[root@localhost ~]#
[root@localhost ~]# iscsiadm -m session -P2
Target: iqn.2015-06.test.target1 (non-flash)
Current Portal: 192.168.100.99:3260,1
Persistent Portal: 192.168.100.99:3260,1
*****
Interface:
*****
Iface Name: default
Iface Transport: iser
Iface Initiatorname: iqn.1994-05.com.redhat:c672dfb8b08f
Iface IPaddress: <empty>
Iface HWaddress: <empty>
Iface Netdev: <empty>
SID: 33
iSCSI Connection State: LOGGED IN
iSCSI Session State: LOGGED_IN
Internal iscsid Session State: NO CHANGE
*****
Timeouts:
*****
Recovery Timeout: 120
```

图 8-3. Iface 传输确认

- e. 要检查新 iSCSI 设备，如图 8-4 中所示，发出 `lsscsi` 命令。

```
[root@localhost ~]# lsscsi
[6:0:0:0]   disk      HP          LOGICAL VOLUME  1.18  /dev/sdb
[6:0:0:1]   disk      HP          LOGICAL VOLUME  1.18  /dev/sda
[6:0:0:3]   disk      HP          LOGICAL VOLUME  1.18  /dev/sdc
[6:3:0:0]   storage  HP          P440ar          1.18  -
[39:0:0:0]  disk      LIO-ORG    ram1            4.0   /dev/sdd
[root@localhost ~]#
```

图 8-4. 检查新 iSCSI 设备

## 为 SLES 12 配置 iSER

由于 targetcli 没有在 SLES 12.x 上内建，您必须完成以下步骤。

### 要为 SLES 12 配置 iSER：

1. 要安装 targetcli，从 ISO 映像（x86\_64 和 noarch 位置）复制并安装以下 RPM：

```
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
```

2. 启动 targetcli 之前，请如下加载所有 RoCE 设备驱动程序和 iSER 模块：

```
# modprobe qed  
# modprobe qede  
# modprobe qedr  
# modprobe ib_iser （启动器）  
# modprobe ib_isert （目标）
```

3. 配置 iSER 目标之前，配置 NIC 接口并运行 L2 和 RoCE 流量，如 [74 页](#)上的 [步骤 7](#) 中所述。
4. 启动 targetcli 公用程序，然后在 iSER 目标系统上配置您的目标。

---

### 注

RHEL 和 SLES 中的 targetcli 版本是不同的。请务必使用合适的后备存储配置您的目标：

- RHEL 使用 *ramdisk*
  - SLES 使用 *rd\_mcp*
-

## 在 RHEL 和 SLES 上通过 iWARP 使用 iSER

配置与 RoCE 类似的 iSER 启动器和目标以与 iWARP 一起使用。您可以使用不同的方法创建 Linux-IO 目标 (LIO™)；本节中列出了一种方法。在 SLES 12 和 RHEL 7.x 中的 targetcli 配置可能会由于版本不同而存在一些差异。

### 要为 LIO 配置目标：

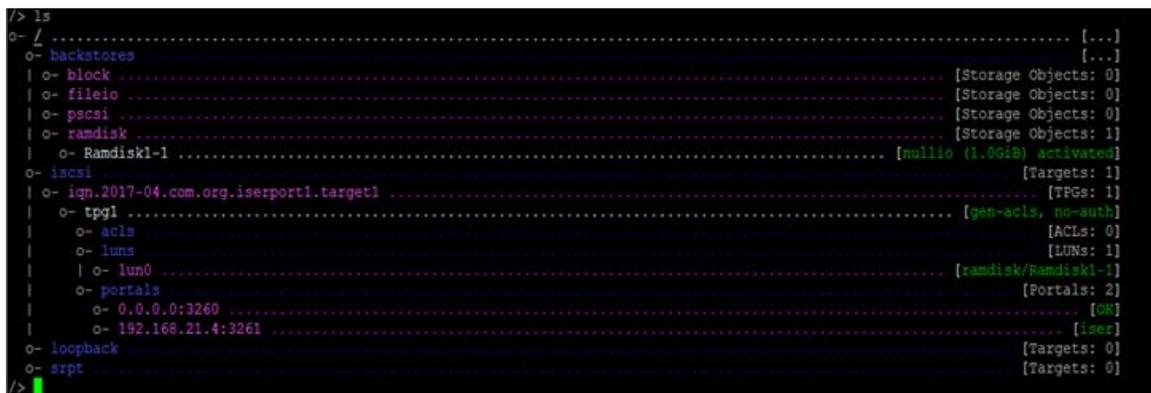
1. 使用 targetcli 公用程序创建 LIO 目标。发出以下命令：

```
# targetcli
targetcli shell version 2.1.fb41
Copyright 2011-2013 by Datera, Inc and others.
For help on commands, type 'help'.
```

2. 发出以下命令：

```
> /backstores/ramdisk create Ramdisk1-1 lg 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
> saveconfig
```

图 8-5 显示 LIO 的目标配置。



```
> ls
o- / ..... [..]
o- backstores ..... [..]
| o- block ..... [Storage Objects: 0]
| o- fileio ..... [Storage Objects: 0]
| o- pscsi ..... [Storage Objects: 0]
| o- ramdisk ..... [Storage Objects: 1]
|   o- Ramdisk1-1 ..... [nullio (1.0GiB) activated]
o- iscsi ..... [Targets: 1]
| o- iqn.2017-04.com.org.iserport1.target1 ..... [TPGs: 1]
|   o- tpg1 ..... [gen-acls, no-auth]
|     o- acls ..... [ACIs: 0]
|     o- luns ..... [LUNs: 1]
|       | o- lun0 ..... [ramdisk/Ramdisk1-1]
|     o- portals ..... [Portals: 2]
|       o- 0.0.0.0:3260 ..... [OH]
|       o- 192.168.21.4:3261 ..... [iser]
o- loopback ..... [Targets: 0]
o- srpt ..... [Targets: 0]
/>
```

图 8-5. LIO 目标配置

### 要为 iWARP 配置启动器:

1. 要使用端口 3261 查找 iSER LIO 目标, 请如下发出 `iscsiadm` 命令:

```
# 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. 将传输模式更改为 `iser`, 如下所示:

```
# iscsiadm -m node -o update -T iqn.2017-04.com.org.iserport1.target1 -n  
iface.transport_name -v iser
```

3. 使用端口 3261 登录到目标:

```
# iscsiadm -m node -l -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. 通过发出以下命令, 确保这些 LUN 可见:

```
# ls SCSI  
[1:0:0:0] storage HP P440ar 3.56 -  
[1:1:0:0] disk HP LOGICAL VOLUME 3.56 /dev/sda  
[6:0:0:0] cd/dvd hp DVD-ROM DUD0N UMD0 /dev/sr0  
[7:0:0:0] disk LIO-ORG Ramdisk1-1 4.0 /dev/sdb
```

## 优化 Linux 性能

考虑本节中介绍的以下 Linux 性能配置增强。

- [将 CPU 配置为最高性能模式](#)
- [配置内核 sysctl 设置](#)
- [配置 IRQ 关联设置](#)
- [配置块设备暂存](#)

### 将 CPU 配置为最高性能模式

使用以下脚本将 CPU scaling governor 配置为 `performance`, 从而将所有 CPU 设置为最高性能模式:

```
for CPUFREQ in  
/sys/devices/system/cpu/cpu*/cpufreq/scaling_governor; do [ -f  
$CPUFREQ ] || continue; echo -n performance > $CPUFREQ; done
```

通过发出以下命令，验证所有 CPU 核心是否设置为最高性能模式：

```
cat /sys/devices/system/cpu/cpu*/cpufreq/scaling_governor
```

## 配置内核 sysctl 设置

按如下所示设定内核 sysctl 设置：

```
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_sack=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
```

## 配置 IRQ 关联设置

以下示例将 CPU 核心 0、1、2 和 3 分别设置为中断请求 (IRQ) XX、YY、ZZ 和 XYZ。对分配给端口的每个 IRQ 执行以下步骤（默认为每端口八个队列）。

```
systemctl disable irqbalance  
systemctl stop irqbalance  
cat /proc/interrupts | grep qedr 显示分配给每个端口队列的 IRQ  
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
```

## 配置块设备暂存

请如下设定每个 iSCSI 设备或目标的块设备暂存设置：

```
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
```

## 在 ESXi 6.7 上配置 iSER

本节提供为 VMware ESXi 6.7 配置 iSER 的信息。

### 准备工作

在为 ESXi 6.7 配置 iSER 前，请确保完成以下步骤：

- 包含 NIC 和 RoCE 驱动程序的 CNA 程序包已在 ESXi 6.7 系统上安装，且已列出设备。要查看 RDMA 设备，请发出以下命令：

```
esxcli rdma device list
```

Name	Driver	State	MTU	Speed	Paired Uplink	Description
vmrdma0	qedrntv	Active	1024	40 Gbps	vmnic4	QLogic FastLinQ QL45xxx RDMA Interface
vmrdma1	qedrntv	Active	1024	40 Gbps	vmnic5	QLogic FastLinQ QL45xxx RDMA Interface

```
[root@localhost:~] esxcfg-vmknics -l
```

Interface	Port	Group/DVPort/Opaque	Network	IP Family	IP Address	Netmask	Broadcast	MAC Address	MTU	TSO	MSS	Enabled	Type
vmk0		Management	Network	IPv4	172.28.12.94	255.255.240.0	172.28.15.255	e0:db:55:0c:5f:94	1500	65535	true	DHCP	
vmk0		Management	Network	IPv6	fe80::e2db:55ff:fe0c:5f94	64		e0:db:55:0c:5f:94	1500	65535	true	STATIC, PREFERRED	

- 已配置 iSER 目标，以便与 iSER 启动程序通信。

### 为 ESXi 6.7 配置 iSER

要为 ESXi 6.7 配置 iSER：

1. 通过发出以下命令添加 iSER 设备：

```
esxcli rdma iser add
```

```
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. 按如下命令禁用防火墙。

```
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
```

3. 创建一个标准的 vSwitch Vmkernel 端口组并分配 IP：

```
esxcli network vswitch standard add -v vSwitch_iser1
```

## 8-iSER 配置

### 在 ESXi 6.7 上配置 iSER

---

#### esxcfg-nics -l

Name	PCI	Driver	Link	Speed	Duplex	MAC Address	MTU	Description
vmnic0	0000:01:00.0	ntg3	Up	1000Mbps	Full	e0:db:55:0c:5f:94	1500	Broadcom Corporation NetXtreme BCM5720 Gigabit Ethernet
vmnic1	0000:01:00.1	ntg3	Down	0Mbps	Half	e0:db:55:0c:5f:95	1500	Broadcom Corporation NetXtreme BCM5720 Gigabit Ethernet
vmnic2	0000:02:00.0	ntg3	Down	0Mbps	Half	e0:db:55:0c:5f:96	1500	Broadcom Corporation NetXtreme BCM5720 Gigabit Ethernet
vmnic3	0000:02:00.1	ntg3	Down	0Mbps	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	qedentv	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 -n vmk1 -A vmhba65
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
```

## 8-iSER 配置

### 在 ESXi 6.7 上配置 iSER

---

```
Is NIC: true
Is Initiator: true
Is Target: false
Using TCP Offload Engine: true
Using ISCSI Offload Engine: true
```

#### 4. 将目标添加到 iSER 启动器，如下所示：

```
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
-----
vmhba65 iqn.2015-06.test.target1 SENDTARGETS No Error
esxcli storage core adapter rescan --adapter vmhba65
```

#### 5. 列出附加目标如下：

```
esxcfg-scsidevs -l
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: raml 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
```

# 9 iSCSI 配置

本章提供以下 iSCSI 配置信息：

- [iSCSI 引导](#)
- [第 118 页上的“配置 iSCSI 引导”](#)
- [第 129 页上的“配置 DHCP 服务器以支持 iSCSI 引导”](#)
- [第 133 页上的“Windows Server 中的 iSCSI 卸载”](#)
- [第 142 页上的“Linux 环境中的 iSCSI 卸载”](#)
- [第 156 页上的“从 RHEL 7.4 及更高版本的 SAN 配置 iSCSI 引导”](#)

---

## 注

某些 iSCSI 功能在当前版本中可能并未完全启用。有关详细信息，请参阅 [附录 D 功能约束](#)。

---

## iSCSI 引导

Cavium 4xxxx 系列千兆位以太网 (GbE) 适配器支持 iSCSI 引导，从而实现无盘系统的操作系统网络引导。iSCSI 引导允许 Windows、Linux 或 VMware 操作系统通过标准 IP 网络从位于远程的 iSCSI 目标计算机引导。

本节中的 iSCSI 引导信息包括：

- [iSCSI 引导设置](#)
- [适配器 UEFI 引导模式配置](#)

对于 Windows 和 Linux 操作系统，可以使用 **UEFI iSCSI HBA**（使用 QLogic 卸载 iSCSI 驱动程序卸载路径）配置 iSCSI 引导。该选项使用引导协议在端口级配置下设置。

## iSCSI 引导设置

iSCSI 引导设置包括：

- 选择首选的 iSCSI 引导模式
- 配置 iSCSI 目标
- 配置 iSCSI 引导参数

### 选择首选的 iSCSI 引导模式

引导模式选项在适配器的 **iSCSI Configuration**（iSCSI 配置）（图 9-1）下列出，该设置为端口特定的。有关访问 UEFI HII 下设备级配置菜单的说明，请参阅 OEM 用户手册。

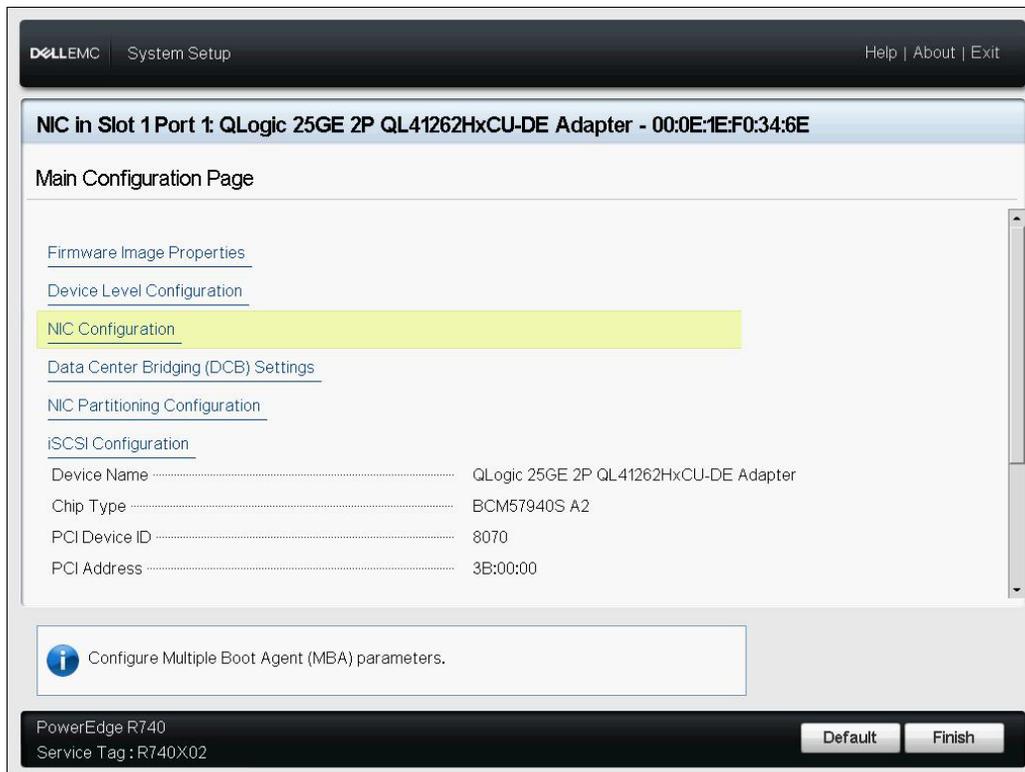


图 9-1. 系统设置：NIC 配置

### 注

从 SAN 引导仅在 NPAR 模式下受支持，并且在 UEFI 中配置，而不是在传统 BIOS 中配置。

## 配置 iSCSI 目标

配置 iSCSI 目标随目标供应商而异。有关配置 iSCSI 目标的信息，请参阅供应商提供的说明文件。

### 要配置 iSCSI 目标：

1. 基于您的 iSCSI 目标选择相应的步骤：
  - 为 SANBlaze® 或 IET® 等目标创建 iSCSI 目标。
  - 为 EqualLogic® 或 EMC® 等目标创建虚拟盘或卷。
2. 创建虚拟磁盘。
3. 将虚拟盘映射到 [步骤 1](#) 中创建的 iSCSI 目标。
4. 将 iSCSI 启动器与 iSCSI 目标关联。记录以下信息：
  - iSCSI 目标名称
  - TCP 端口号
  - iSCSI 逻辑单元号 (LUN)
  - 启动器 iSCSI 限定名称 (IQN)
  - CHAP 身份验证详细信息
5. 配置 iSCSI 目标之后，获取以下信息：
  - 目标 IQN
  - 目标 IP 地址
  - 目标 TCP 端口号
  - 目标 LUN
  - 启动器 IQN
  - CHAP ID 和机密

## 配置 iSCSI 引导参数

配置 QLogic iSCSI 引导软件以实现静态或动态配置。有关 General Parameters（常规参数）窗口提供的配置选项，请参阅[表 9-1](#)，其中列出了适用于 IPv4 和 IPv6 的参数。IPv4 或 IPv6 的特定参数将特别注明。

---

### 注

IPv6 iSCSI 引导的可用性取决于平台和设备。

---

表 9-1. 配置选项

选项	说明
TCP/IP parameters via DHCP 通过 DHCP 获取 TCP/IP 参数	此选项特定于 IPv4。控制 iSCSI 引导主机软件是使用 DHCP 获取 IP 地址信息 (Enabled [已启用]) 还是使用静态 IP 配置 (Disabled [已禁用])。
iSCSI parameters via DHCP 通过 DHCP 获取 iSCSI 参数	控制 iSCSI 引导主机软件是使用 DHCP 获取其 iSCSI 目标参数 (Enabled [已启用]) 还是通过静态配置 (Disabled [已禁用])。静态信息在 iSCSI Initiator Parameters Configuration (iSCSI 启动器参数配置) 页面中输入。
CHAP Authentication (CHAP 身份验证)	控制 iSCSI 引导主机软件在连接到 iSCSI 目标时是否使用 CHAP 身份验证。如果启用了 CHAP Authentication (CHAP 身份验证), 请在 iSCSI Initiator Parameters Configuration (iSCSI 启动器参数配置) 页面中配置 CHAP ID 和 CHAP Secret (CHAP 机密)。
IP Version (IP 版本)	此选项特定于 IPv6。在 IPv4 和 IPv6 之间切换。如果您从一个协议版本切换到另一个协议版本, 所有 IP 设置都将丢失。
DHCP Request Timeout (DHCP 请求超时)	允许您指定 DHCP 请求的最长等待时间 (以秒为单位), 然后响应完成。
Target Login Timeout (目标登录超时)	允许您指定启动器完成目标登录的最长等待时间 (以秒为单位)。
DHCP Vendor ID (DHCP 供应商 ID)	控制 iSCSI 引导主机软件如何解释在 DHCP 期间使用的 Vendor Class ID (供应商类别 ID) 字段。如果 DHCP 中的 Vendor Class ID (供应商类别 ID) 字段提供的数据包匹配字段中的值, 则 iSCSI 引导主机软件将查看所需 iSCSI 引导扩展中的 DHCP 选项 43 字段。如果 DHCP 被禁用, 不必设置此值。

## 适配器 UEFI 引导模式配置

要配置引导模式：

1. 重新启动系统。
2. 访问 System Utilities（系统公用程序）菜单（图 9-2）。

### 注

SAN 引导仅在 UEFI 环境中受支持。确保系统引导选项为 UEFI，而非旧版。

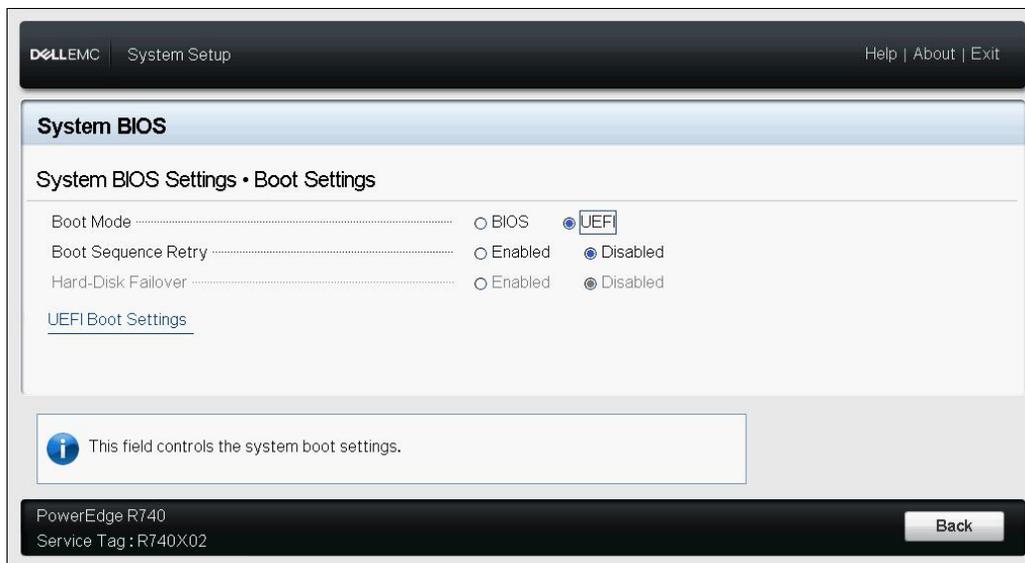


图 9-2. 系统设置：引导设置

3. 在 System Setup (系统设置)、Device Settings (设备设置) 中, 选择 QLogic 设备 (图 9-3)。有关访问 PCI 设备配置菜单, 请参阅 OEM 用户指南。

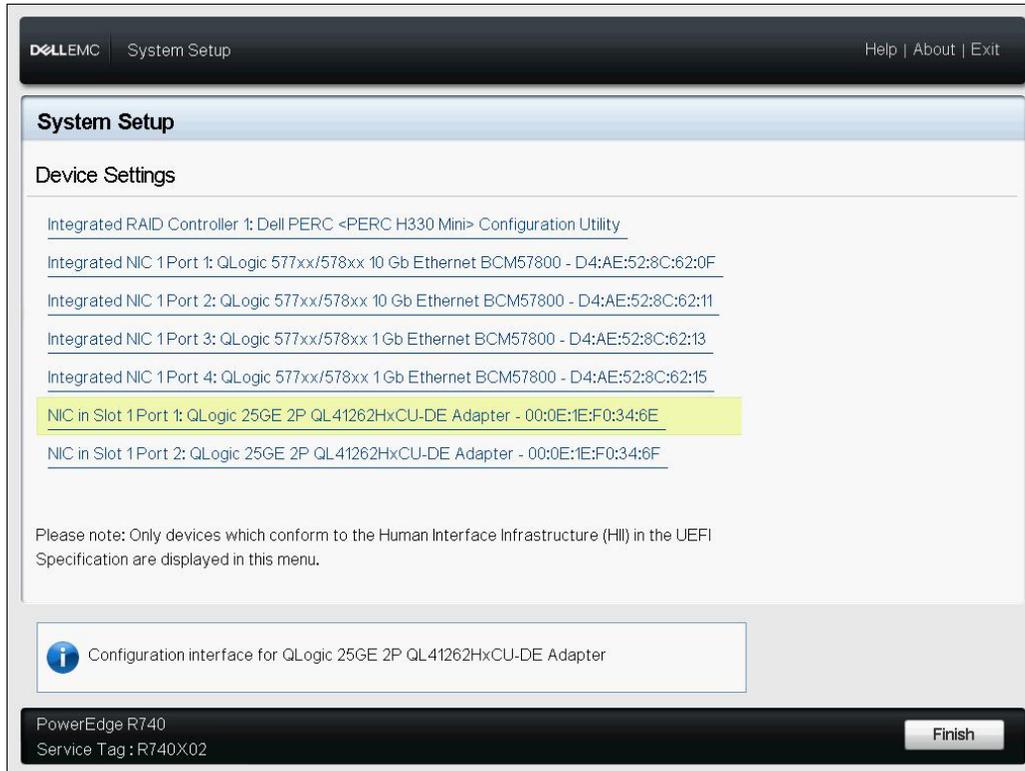


图 9-3. 系统设置：设备设置配置实用程序

4. 在 Main Configuration Page（主要配置页面）中，选择 **NIC Configuration**（NIC 配置）（图 9-4），然后按 ENTER 键。

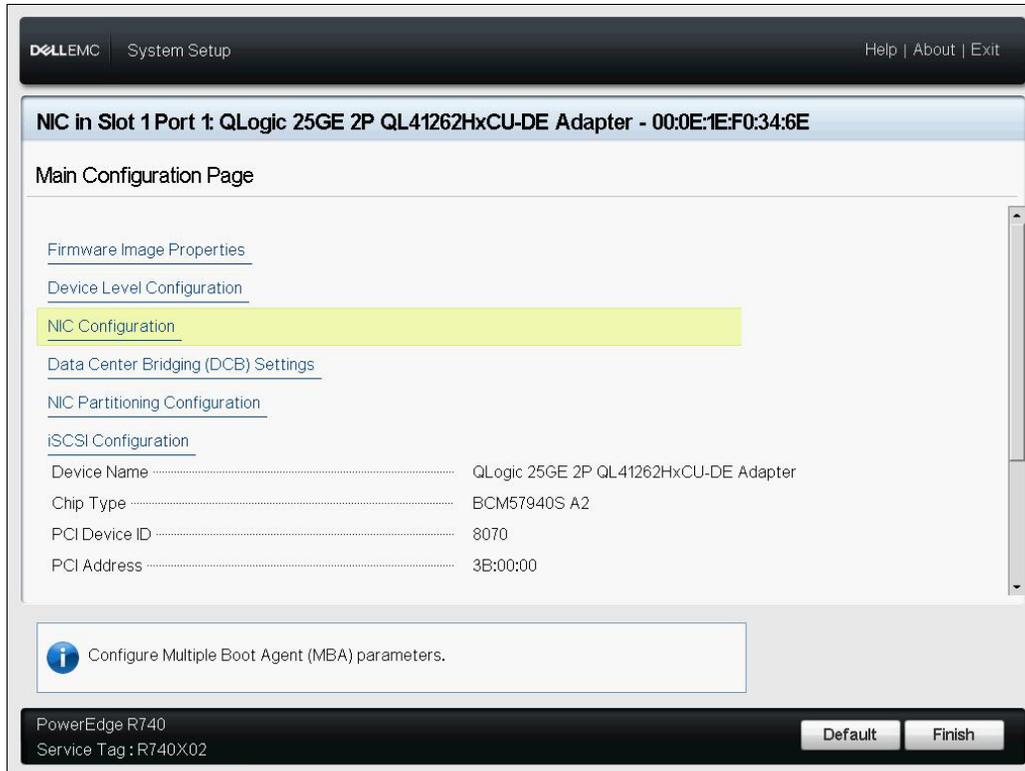


图 9-4. 选择 NIC 配置

5. 在 NIC 配置页面 (图 9-5) 上, 选择引导协议, 然后按 ENTER 键选择 **UEFI iSCSI HBA** (需要 NPAR 模式)。

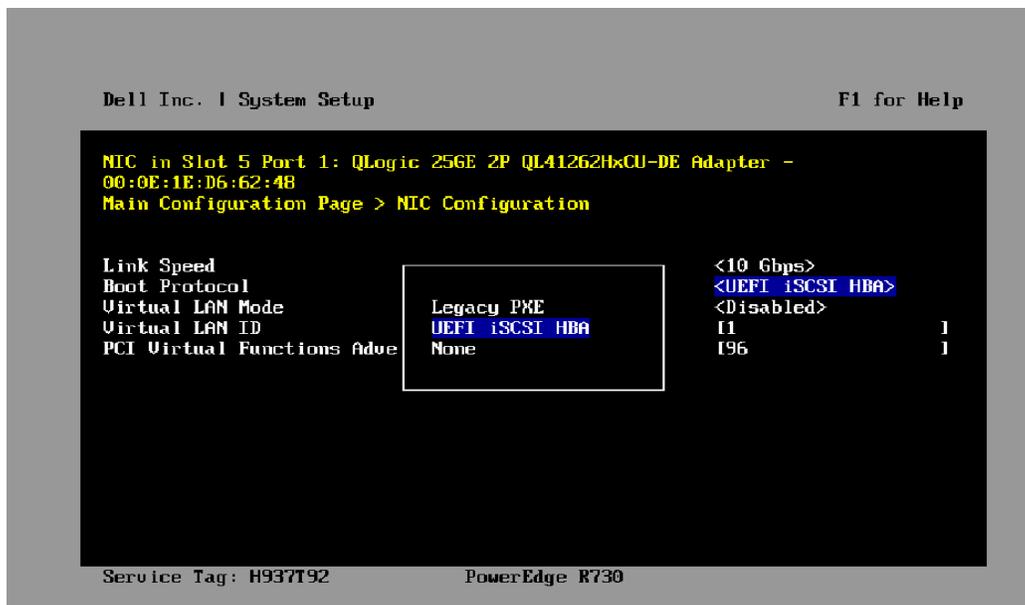


图 9-5. 系统设置: NIC 配置, 引导协议

6. 继续执行以下配置选项之一:
  - 第 118 页上的“静态 iSCSI 引导配置”
  - 第 126 页上的“动态 iSCSI 引导配置”

## 配置 iSCSI 引导

iSCSI 引导配置选项包括:

- 静态 iSCSI 引导配置
- 动态 iSCSI 引导配置
- 启用 CHAP 身份验证

### 静态 iSCSI 引导配置

在静态配置中, 您必须输入以下各项的数据:

- 系统的 IP 地址
- 系统的启动器 IQN
- 目标参数 (在第 113 页上的“配置 iSCSI 目标”中获得)

关于配置选项的信息，请参见 第 114 页上表 9-1。

要使用静态配置来配置 iSCSI 引导参数：

1. 在设备 HII 的 **Main Configuration Page**（主要配置页面）中，选择 **iSCSI Configuration**（iSCSI 配置）（图 9-6），然后按 ENTER 键。

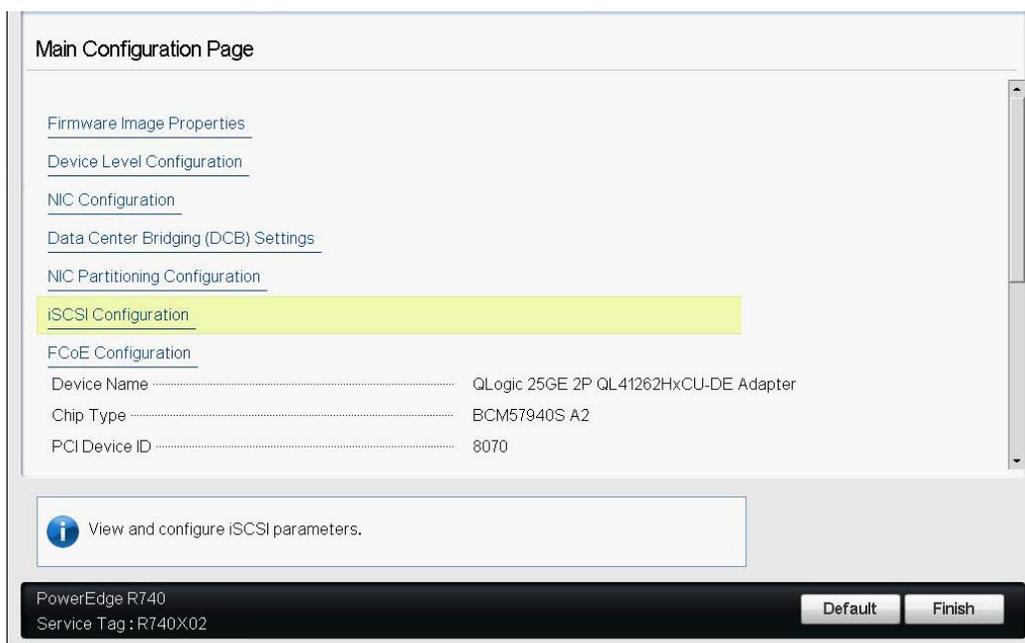


图 9-6. 系统设置：iSCSI 配置

2. 在 iSCSI Configuration（iSCSI 配置）页面中，选择 **iSCSI General Parameters**（iSCSI 常规参数）（图 9-7），然后按 ENTER 键。

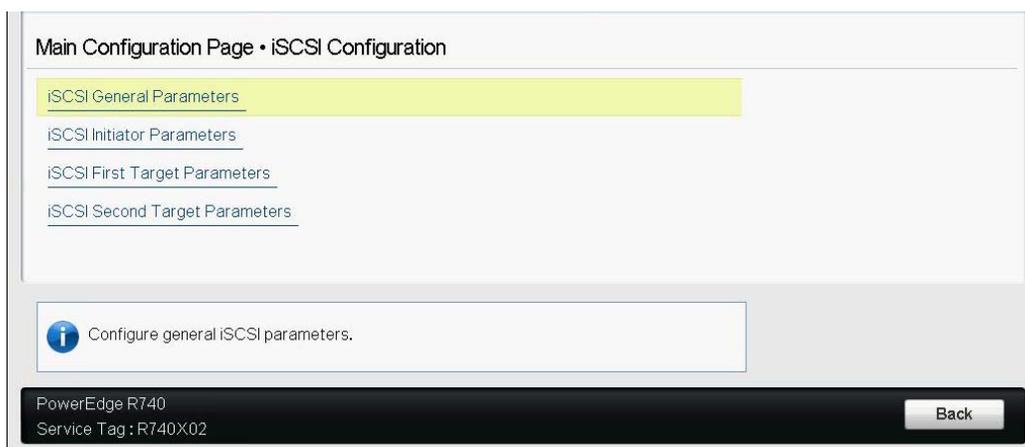


图 9-7. 系统设置：选择常规参数

3. 在 iSCSI General Parameters (iSCSI 常规参数) 页面 (图 9-8) 中, 按向上箭头键和向下箭头键选择参数, 然后按 ENTER 键选择或输入以下值:
  - TCP/IP Parameters via DHCP (通过 DHCP 获取 TCP/IP 参数): 已禁用
  - iSCSI Parameters via DHCP (通过 DHCP 获取 iSCSI 参数): 已禁用
  - CHAP Authentication (CHAP 身份验证): 根据需要
  - IP Version (IP 版本): 根据需要 (IPv4 或 IPv6)
  - CHAP Mutual Authentication (CHAP 相互认证): 根据需要
  - DHCP Vendor ID (DHCP 供应商 ID): 不适用于静态配置
  - HBA 引导模式: 已启用。
  - 虚拟 LAN ID: 默认值或根据需要设置
  - 虚拟 LAN 模式: 已禁用

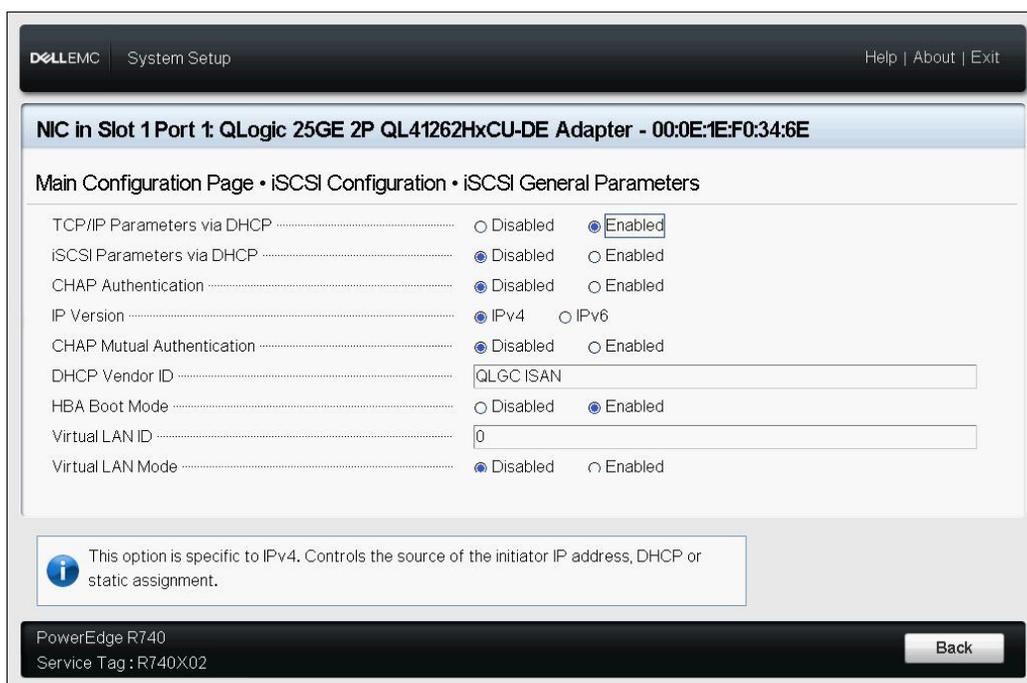


图 9-8. 系统设置: iSCSI 常规参数

4. 返回到 iSCSI Configuration (iSCSI 配置) 页面, 然后按 ESC 键。
5. 选择 **iSCSI Initiator Parameters** (iSCSI 启动器参数) (图 9-9), 然后按 ENTER 键。

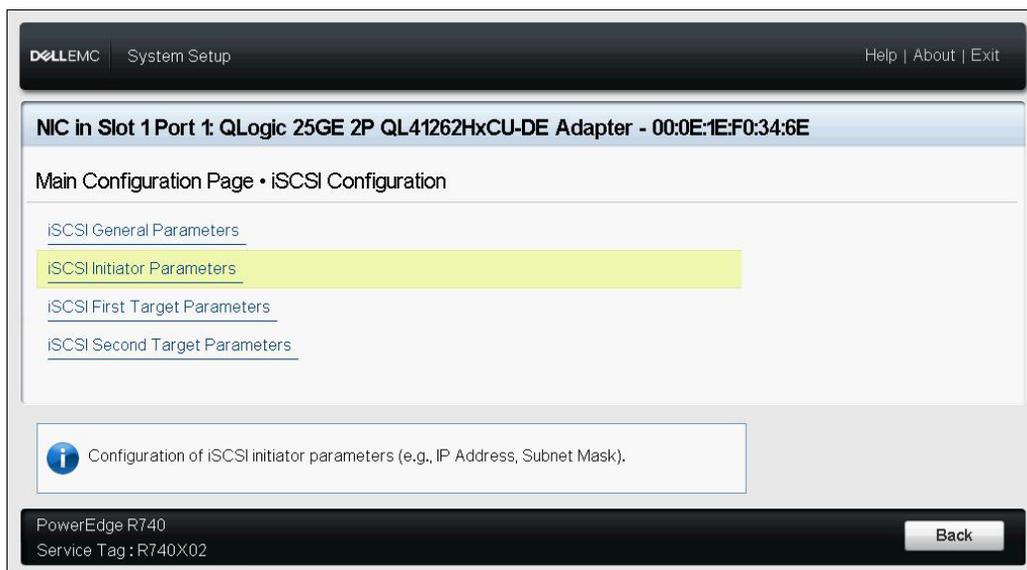


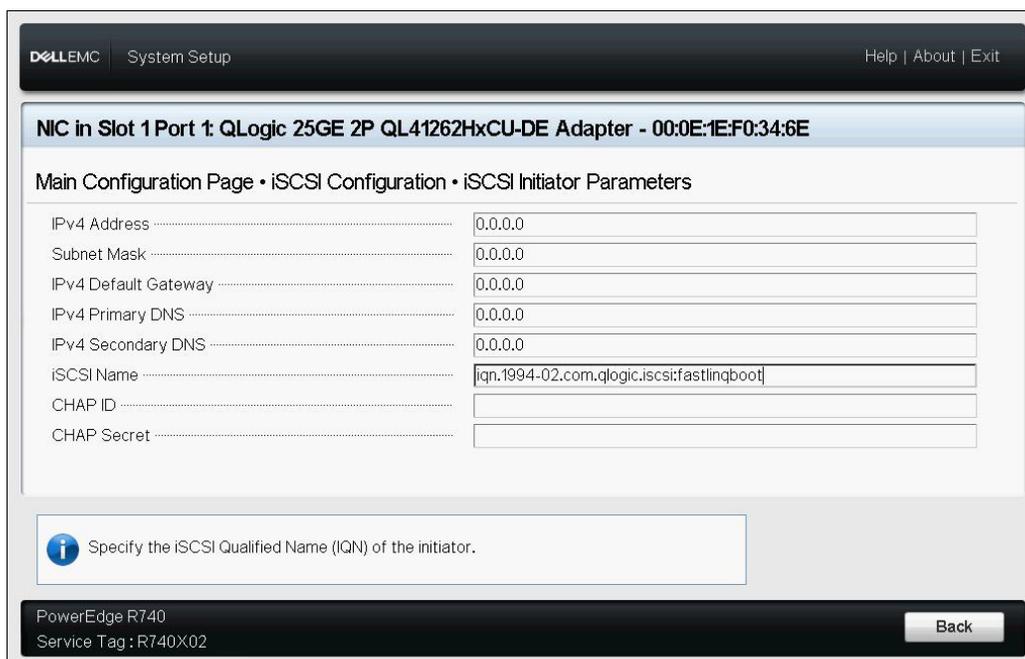
图 9-9. 系统设置：选择 iSCSI 启动器参数

6. 在 iSCSI Initiator Parameters (iSCSI 启动器参数) 页面 (图 9-10) 中, 选择以下参数, 然后键入各自的值:
  - IPv4\* Address** (IPv4 地址)
  - Subnet Mask** (子网掩码)
  - IPv4\* Default Gateway** (IPv4 默认网关)
  - IPv4\* Primary DNS** (IPv4 主 DNS)
  - IPv4\* Secondary DNS** (IPv4 辅助 DNS)
  - iSCSI Name** (iSCSI 名称)。与客户端系统将要使用的 iSCSI 启动器名称对应。
  - CHAP ID**
  - CHAP Secret** (CHAP 机密)

**注**

请注意，以下内容适用于上述带星号 (\*) 的项目：

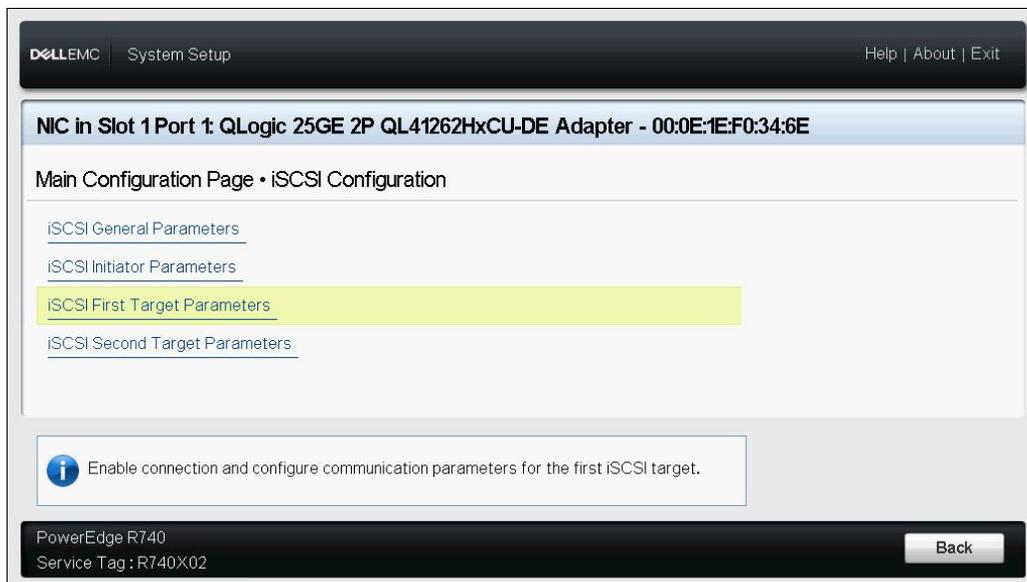
- 标签将基于 iSCSI General Parameters (iSCSI 常规参数) 页面 (第 120 页上的图 9-8) 中设置的 IP 版本更改为 IPv6 或 IPv4 (默认值)。
- 仔细输入 IP 地址。对 IP 地址不会检查是否有重复段、错误段或网络分配错误。



**图 9-10. 系统设置：iSCSI 启动器参数**

7. 返回到 iSCSI Configuration (iSCSI 配置) 页面，然后按 ESC 键。

8. 选择 **iSCSI First Target Parameters** (iSCSI 第一目标参数) (图 9-11), 然后按 ENTER 键。

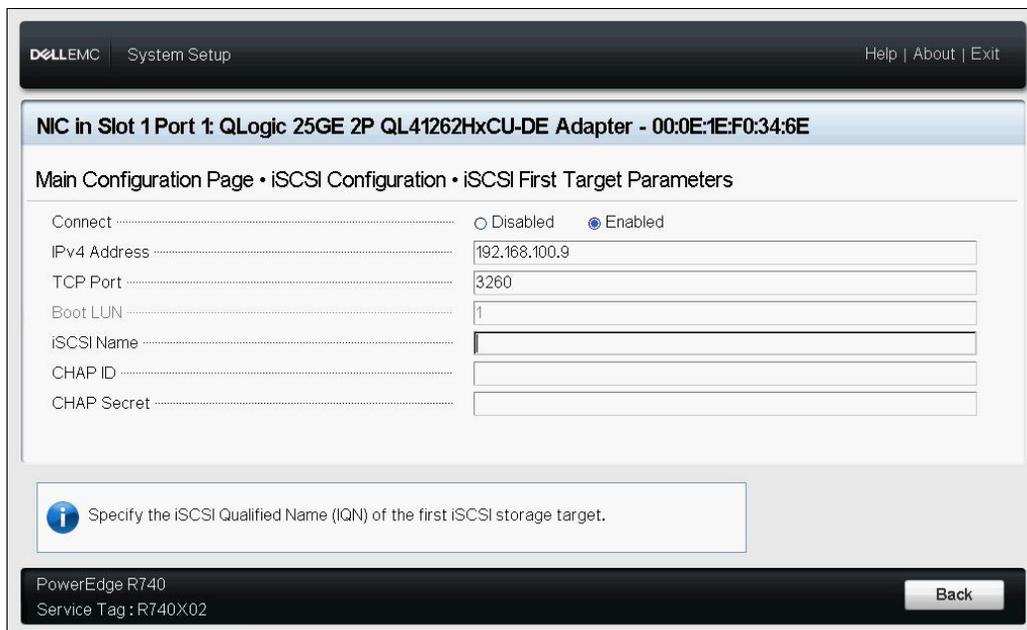


**图 9-11. 系统设置：选择 iSCSI 第一个目标参数**

9. 在 iSCSI First Target Parameters (iSCSI 第一目标参数) 页面中, 将 iSCSI 目标的 **Connect** (连接) 选项设置为 **Enabled** (已启用)。
10. 键入 iSCSI 目标以下参数的值, 然后按 ENTER 键:
  - IPv4\* Address** (IPv4 地址)
  - TCP Port** (TCP 端口)
  - Boot LUN** (引导 LUN)
  - iSCSI Name** (iSCSI 名称)
  - CHAP ID**
  - CHAP Secret** (CHAP 机密)

**注**

对于上述带有星号 (\*) 的参数，标签将基于 iSCSI General Parameters (iSCSI 常规参数) 页面中设置的 IP 版本更改为 **IPv6** 或 **IPv4** (默认值)，如图 9-12 中所示。



**图 9-12. 系统设置：iSCSI 第一个目标参数**

11. 返回到 iSCSI Boot Configuration (iSCSI 引导配置) 页面，然后按 ESC 键。

12. 如果要配置第二个 iSCSI 目标设备，选择 **iSCSI Second Target Parameters**（iSCSI 第二目标参数）（图 9-13），然后如同在步骤 10 中的操作一样输入参数值。否则，继续步骤 13。



**图 9-13. 系统设置：iSCSI 第二个目标参数**

13. 按 ESC 键一次，再次按下可退出。

14. 单击 **Yes**（是）保存更改，或按照 OEM 指导操作保存设备级配置。例如，单击 **Yes**（是）确认设置更改（图 9-14）。

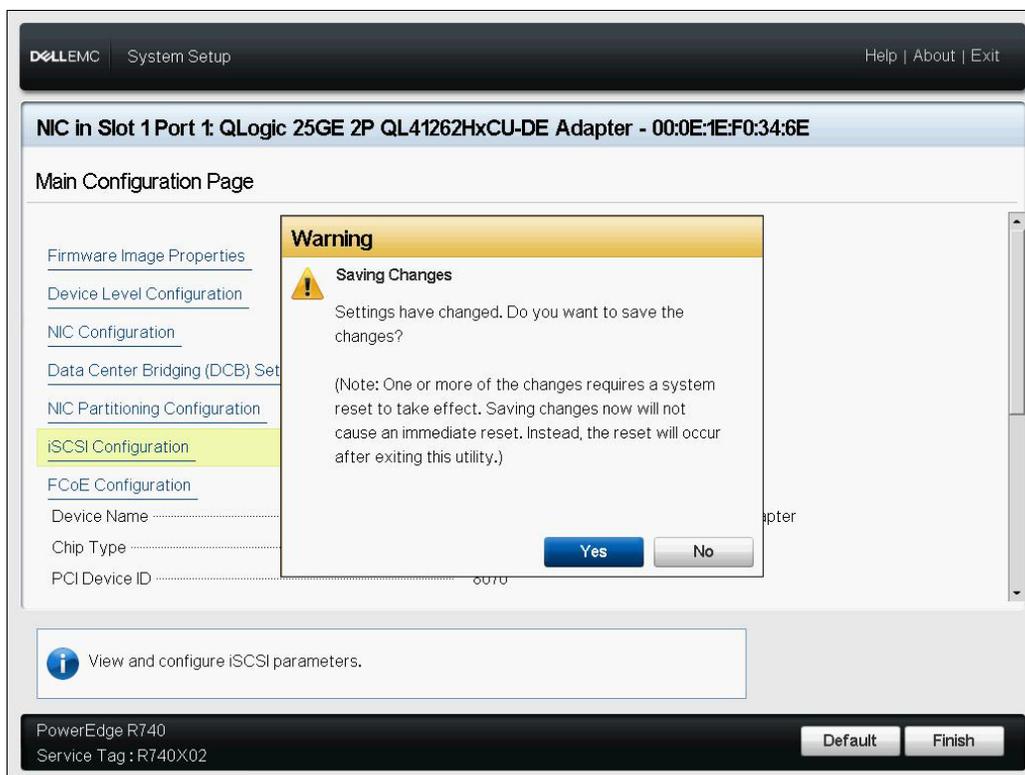


图 9-14. 系统设置：保存 iSCSI 更改

15. 进行所有更改后，重新引导系统以将更改应用到适配器正在运行的配置。

## 动态 iSCSI 引导配置

在动态配置中，确保系统的 IP 地址和目标（或启动器）信息由 DHCP 服务器提供（请参阅第 129 页上的“配置 DHCP 服务器以支持 iSCSI 引导”中的 IPv4 和 IPv6 配置）。

以下参数的任何设置都将被忽略并且无需清除（IPv4 的启动器 iSCSI 名称、IPv6 的 CHAP ID 和 CHAP 机密除外）：

- 启动器参数
- 第一目标参数或第二目标参数

关于配置选项的信息，请参见第 114 页上表 9-1。

## 注

使用 DHCP 服务器时，DNS 服务器条目将被 DHCP 服务器提供的值覆盖。即使本地提供的值有效并且 DHCP 服务器不提供 DNS 服务器信息，仍会发生这种覆盖。当 DHCP 服务器不提供 DNS 服务器信息时，主 DNS 服务器值和辅助 DNS 服务器值均设为 0.0.0.0。当 Windows 操作系统获得控制权时，Microsoft iSCSI 启动器将检索 iSCSI 启动器参数并静态配置相应的注册表。这将覆盖任何配置的参数。由于 DHCP 守护进程在 Windows 环境中作为一个用户进程运行，当堆栈在 iSCSI 引导环境中启动之前，所有 TCP/IP 参数都必须静态配置。

---

如果使用 DHCP 选项 17，则目标信息由 DHCP 服务器提供，且启动器 iSCSI 名称从 Initiator Parameters（启动器参数）窗口的编程值进行检索。如果未选择任何值，控制器默认名称为：

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

字符串 11.22.33.44.55.66 对应于控制器的 MAC 地址。如果使用 DHCP 选项 43（仅适用于 IPv4），则以下窗口中的任何设置都将被忽略并且无需清除：

- 启动器参数
- 第一目标参数或第二目标参数

**要使用动态配置来配置 iSCSI 引导参数：**

- 在 iSCSI General Parameters（iSCSI 常规参数）页面上，设置以下选项，如 [图 9-15](#) 中所示：
  - TCP/IP Parameters via DHCP**（通过 DHCP 获取 TCP/IP 参数）：已启用
  - iSCSI Parameters via DHCP**（通过 DHCP 获取 iSCSI 参数）：已启用
  - CHAP Authentication**（CHAP 身份验证）：根据需要
  - IP Version**（IP 版本）：根据需要（IPv4 或 IPv6）
  - CHAP Mutual Authentication**（CHAP 相互认证）：根据需要
  - DHCP Vendor ID**（DHCP 供应商 ID）：根据需要
  - HBA 引导模式**：已禁用
  - Virtual LAN ID**（虚拟 LAN ID）：根据需要
  - Virtual LAN Boot Mode**（虚拟 LAN 引导模式）：已启用

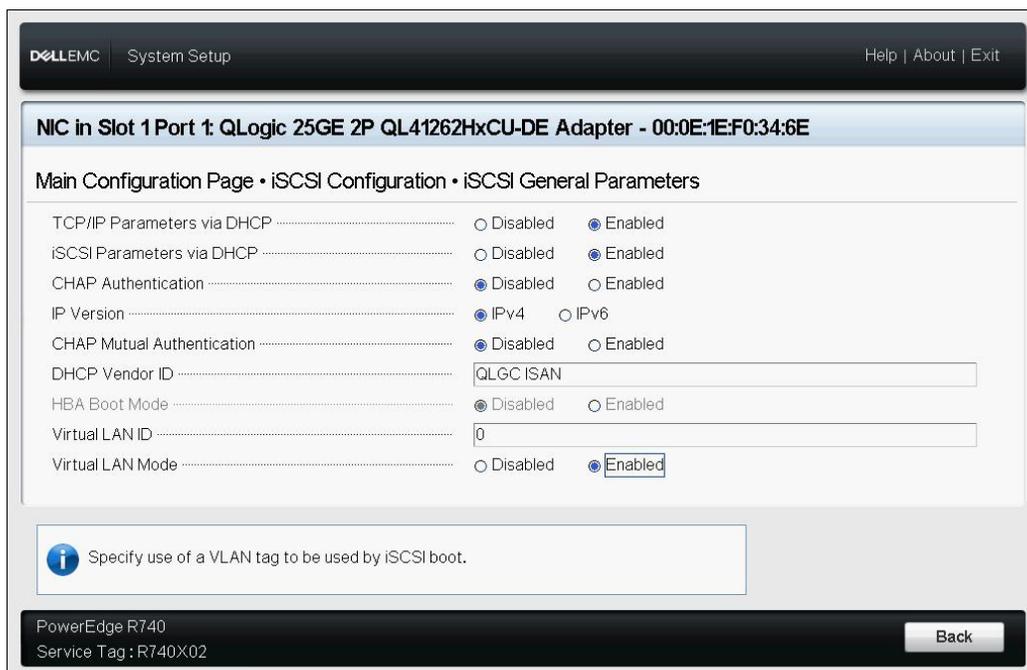


图 9-15. 系统设置：iSCSI 常规参数

## 启用 CHAP 身份验证

确保目标上启用 CHAP 身份验证。

要启用 CHAP 身份验证：

1. 转至 iSCSI General Parameters（iSCSI 常规参数）页面。
2. 将 **CHAP Authentication**（CHAP 身份验证）设置为 **Enabled**（已启用）。
3. 在 Initiator Parameters（启动器参数）窗口中，键入以下值：
  - CHAP ID**（最多 255 个字符）
  - CHAP Secret**（CHAP 机密）（如果身份验证需要；长度必须为 12 到 16 个字符）
4. 按 ESC 键返回到 iSCSI Boot Configuration（iSCSI 引导配置）页面。
5. 在 iSCSI Boot configuration（iSCSI 引导配置）页面中，选择 **iSCSI First Target Parameters**（iSCSI 第一目标参数）。
6. 在 iSCSI First Target Parameters（iSCSI 第一目标参数）窗口中，键入配置 iSCSI 目标时使用的值：
  - CHAP ID**（如果是双向 CHAP，可选填）

- **CHAP Secret** (CHAP 机密) (如果是双向 CHAP, 可选填; 长度必须为 12 到 16 个字符或更长)
7. 按 ESC 键返回到 iSCSI Boot Configuration (iSCSI 引导配置) 页面。
  8. 按 ESC 键, 然后选择确认 **Save Configuration** (保存配置)。

## 配置 DHCP 服务器以支持 iSCSI 引导

DHCP 服务器是一个可选组件, 只有在进行动态 iSCSI 引导配置设置时才必需 (请参阅第 126 页上的“动态 iSCSI 引导配置”)。

对 IPv4 和 IPv6 配置 DHCP 服务器以支持 iSCSI 引导的过程不同:

- [IPv4 的 DHCP iSCSI 引导配置](#)
- [配置 DHCP 服务器](#)
- [为 IPv6 配置 DHCP iSCSI 引导](#)
- [配置 VLAN 用于 iSCSI 引导](#)

### IPv4 的 DHCP iSCSI 引导配置

DHCP 包括向 DHCP 客户端提供配置信息的多个选项。对于 iSCSI 引导, QLogic 适配器支持以下 DHCP 配置:

- [DHCP 选项 17, 根路径](#)
- [DHCP 选项 43, 供应商特定信息](#)

#### DHCP 选项 17, 根路径

选项 17 用于将 iSCSI 目标信息传递到 iSCSI 客户端。

IETF RFC 4173 中定义的根路径的格式为:

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

表 9-2 列出 DHCP 选项 17 参数。

**表 9-2. DHCP 选项 17 参数定义**

参数	定义
"iscsi:"	字符串
<servername>	iSCSI 目标的 IP 地址或完全限定域名 (FQDN)
":"	分隔符

表 9-2. DHCP 选项 17 参数定义 (续)

参数	定义
<protocol>	用于访问 iSCSI 目标的 IP 协议。由于目前仅支持 TCP，因此协议为 6。
<port>	与协议关联的端口号。iSCSI 的标准端口号为 3260。
<LUN>	要在 iSCSI 目标上使用的逻辑单元号。LUN 的值必须以十六进制格式表示。ID 为 64 的 LUN 在 DHCP 服务器的选项 17 参数内必须配置为 40。
<targetname>	IQN 或 EUI 格式的目标名称。有关 IQN 和 EUI 格式的详细信息，请参阅 RFC 3720。IQN 名称示例： iqn.1995-05.com.QLogic:iscsi-target。

### DHCP 选项 43， 供应商特定信息

DHCP 选项 43（供应商特定信息）为 iSCSI 客户端提供比 DHCP 选项 17 更多的配置选项。在此配置中，还提供三个额外的子选项，将可用于引导的启动器 IQN 以及两个 iSCSI 目标 IQN 分配给 iSCSI 引导客户端。iSCSI 目标 IQN 的格式与 DHCP 选项 17 相同，而 iSCSI 启动器 IQN 仅仅是启动器的 IQN。

#### 注

DHCP 选项 43 仅在 IPv4 中受支持。

表 9-3 列出 DHCP 选项 43 子选项。

**表 9-3. DHCP 选项 43 子选项定义**

子选项	定义
201	标准根路径格式中的第一 iSCSI 目标信息： "iscsi:<servername>":"<protocol>":"<port>":"<LUN>": "<targetname>"
202	标准根路径格式中的第二 iSCSI 目标信息： "iscsi:<servername>":"<protocol>":"<port>":"<LUN>": "<targetname>"
203	iSCSI 启动器 IQN

使用 DHCP 选项 43 需要比 DHCP 选项 17 更多的配置，但它提供更丰富的环境和更多的配置选项。您应该在执行动态 iSCSI 引导配置时使用 DHCP 选项 43。

## 配置 DHCP 服务器

配置 DHCP 服务器以支持选项 16、17、43。

### 注

DHCPv6 选项 16 和选项 17 的格式在 RFC 3315 中全面定义。  
如果您使用选项 43，还必须配置选项 60。选项 60 的值必须匹配 DHCP Vendor ID（DHCP 供应商 ID）值 QLGC ISAN，如 iSCSI Boot Configuration（iSCSI 引导配置）页面的 **iSCSI General Parameters**（iSCSI 常规参数）中所示。

---

## 为 IPv6 配置 DHCP iSCSI 引导

DHCPv6 服务器可提供多个选项，包括无状态或有状态 IP 配置，以及 DHCPv6 客户端的信息。对于 iSCSI 引导，QLogic 适配器支持以下 DHCP 配置：

- DHCPv6 选项 16，供应商类别选项
- DHCPv6 选项 17，供应商特定信息

### 注

DHCPv6 标准根路径选项尚不可用。QLogic 建议对动态 iSCSI 引导 IPv6 支持使用选项 16 或选项 17。

### DHCPv6 选项 16，供应商类别选项

DHCPv6 选项 16（供应商类别选项）必须存在且必须包含匹配您配置的 DHCP Vendor ID（DHCP 供应商 ID）参数的字符串。DHCP Vendor ID（DHCP 供应商 ID）值为 QLGC ISAN，如 iSCSI Boot Configuration（iSCSI 引导配置）菜单的 **General Parameters**（常规参数）中所示。

选项 16 的内容应为 <2-byte length> <DHCP Vendor ID>。

### DHCPv6 选项 17，供应商特定信息

DHCPv6 选项 17（供应商特定信息）为 iSCSI 客户端提供更多的配置选项。在此配置中，还提供三个额外的子选项，将可用于引导的启动器 IQN 以及两个 iSCSI 目标 IQN 分配给 iSCSI 引导客户端。

表 9-4 列出 DHCP 选项 17 子选项。

**表 9-4. DHCP 选项 17 子选项定义**

子选项	定义
201	标准根路径格式中的第一 iSCSI 目标信息： "iscsi:"[<servername>]":"<protocol>":"<port>":"<LUN>": "<targetname>"
202	标准根路径格式中的第二 iSCSI 目标信息： "iscsi:"[<servername>]":"<protocol>":"<port>":"<LUN>": "<targetname>"
203	iSCSI 启动器 IQN

表格注释：

方括号 [] 是 IPv6 地址所必需。

选项 17 的内容应为：

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

## 配置 VLAN 用于 iSCSI 引导

网络上的 iSCSI 流量可以隔离在第二层 VLAN 中，以与常规流量隔离开来。如果是这种情况，请让适配器上的 iSCSI 接口成为该 VLAN 的成员。

**要配置 VLAN 用于 iSCSI 引导：**

1. 转至端口的 **iSCSI Configuration**（iSCSI 配置）页面。
2. 选择 **iSCSI General Parameters**（iSCSI 常规参数）。
3. 选择 **VLAN ID** 以输入和设置 VLAN 值，如图 9-16 中所示。

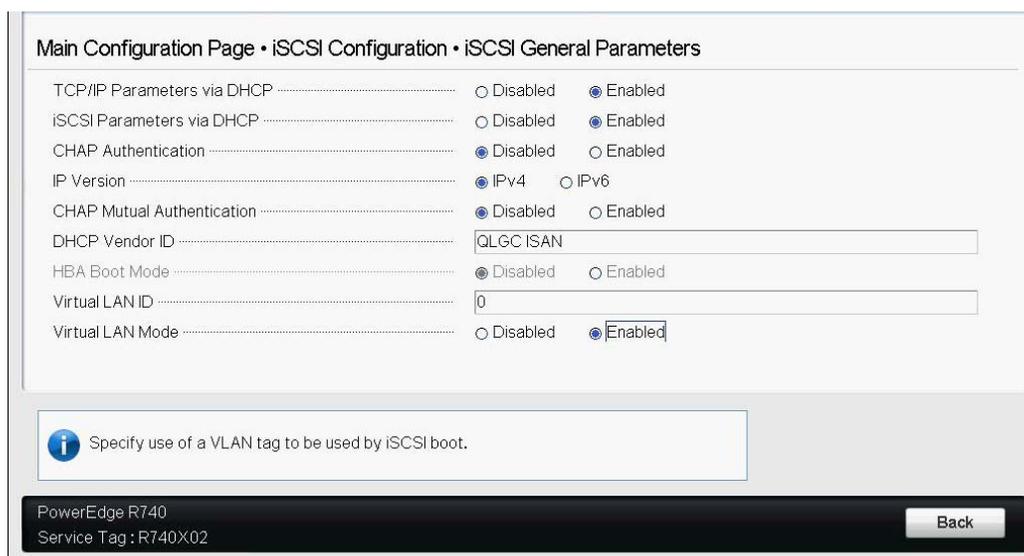


图 9-16. 系统设置：iSCSI 常规参数，VLANID

## Windows Server 中的 iSCSI 卸载

iSCSI 卸载是一种将 iSCSI 协议处理开销从主机处理器卸载到 iSCSI HBA 的技术。iSCSI 卸载可提高网络性能和吞吐量，同时帮助优化服务器处理器的使用。本节介绍如何为 QLogic 41xxx 系列适配器配置 Windows iSCSI 卸载功能。

通过适当的 iSCSI 卸载许可，可配置具有 iSCSI 功能的 41xxx 系列适配器，以便从主机处理器卸载 iSCSI 处理。以下各节介绍如何使系统能够充分利用 QLogic 的 iSCSI 卸载功能：

- [安装 QLogic 驱动程序](#)

- [安装 Microsoft iSCSI 启动器](#)
- [配置 Microsoft 启动器以使用 QLogic 的 iSCSI 卸载](#)
- [iSCSI 卸载常见问题](#)
- [Windows Server 2012 R2 和 2016 iSCSI 引导安装](#)
- [iSCSI 故障转储](#)

## 安装 QLogic 驱动程序

如第 15 页上的“安装 Windows 驱动程序软件”中所述，安装 Windows 驱动程序。

## 安装 Microsoft iSCSI 启动器

启动 Microsoft iSCSI 启动器小程序。初次启动时，系统会提示自动服务启动。确认要启动的小程序的选择。

## 配置 Microsoft 启动器以使用 QLogic 的 iSCSI 卸载

在为 iSCSI 适配器配置 IP 地址后，必须使用 Microsoft 启动器来配置和添加到使用 QLogic iSCSI 适配器的 iSCSI 目标的连接。有关 Microsoft 启动器的更多详细信息，请参阅 Microsoft 用户指南。

### 要配置 Microsoft 启动器：

1. 打开 Microsoft 启动器。
2. 要根据您的设置配置启动器 IQN 名称，请按照以下步骤操作：
  - a. 在 iSCSI Initiator Properties（iSCSI 启动器属性）中，单击 **Configuration**（配置）选项卡。
  - b. 在 Configuration（配置）页面（[图 9-17](#)）中，单击 **Change**（更改）

修改启动器名称。

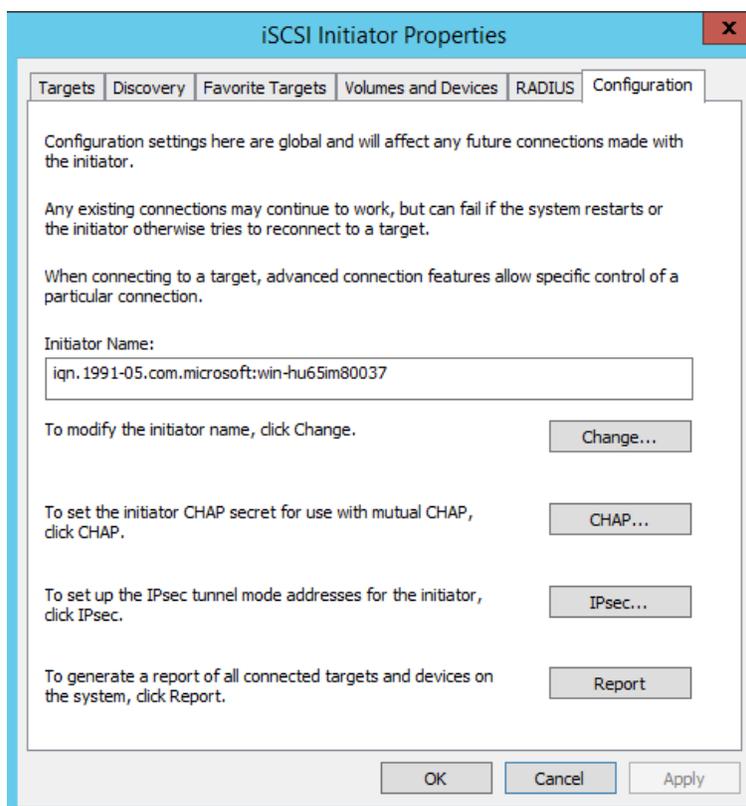


图 9-17. iSCSI 启动器属性, 配置页面

- c. 在 iSCSI Initiator Name (iSCSI 启动器名称) 对话框中, 键入新的启动器 IQN 名称, 然后单击 **OK** (确定)。(图 9-18)

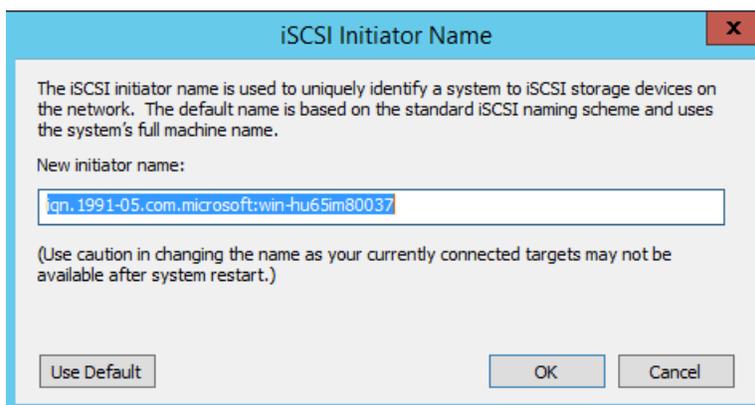


图 9-18. iSCSI 启动器节点名称更改

3. 在 iSCSI Initiator Properties (iSCSI 启动器属性) 中, 单击 **Discovery** (查找) 选项卡。
4. 在 Discovery (查找) 页面 (图 9-19) 的 **Target portals** (目标门户) 下, 单击 **Discover Portal** (查找门户)。

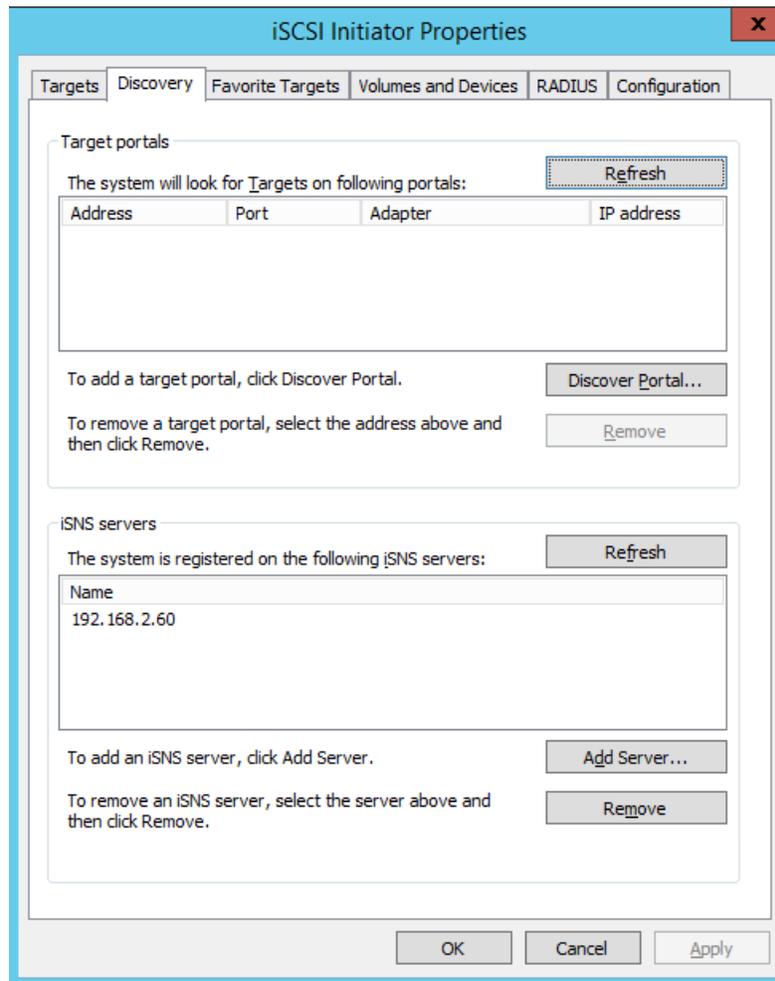


图 9-19. iSCSI 启动器 - 查找目标门户

5. 在 Discover Target Portal（查找目标门户）对话框（图 9-20）中：
  - a. 在 **IP address or DNS name**（IP 地址或 DNS 名称）框中，键入目标的 IP 地址。
  - b. 单击 **Advanced**（高级）。

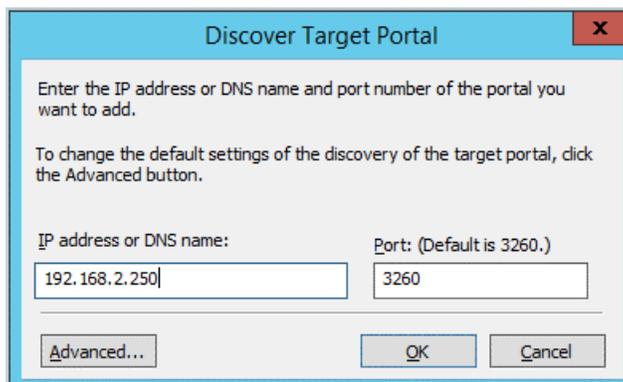
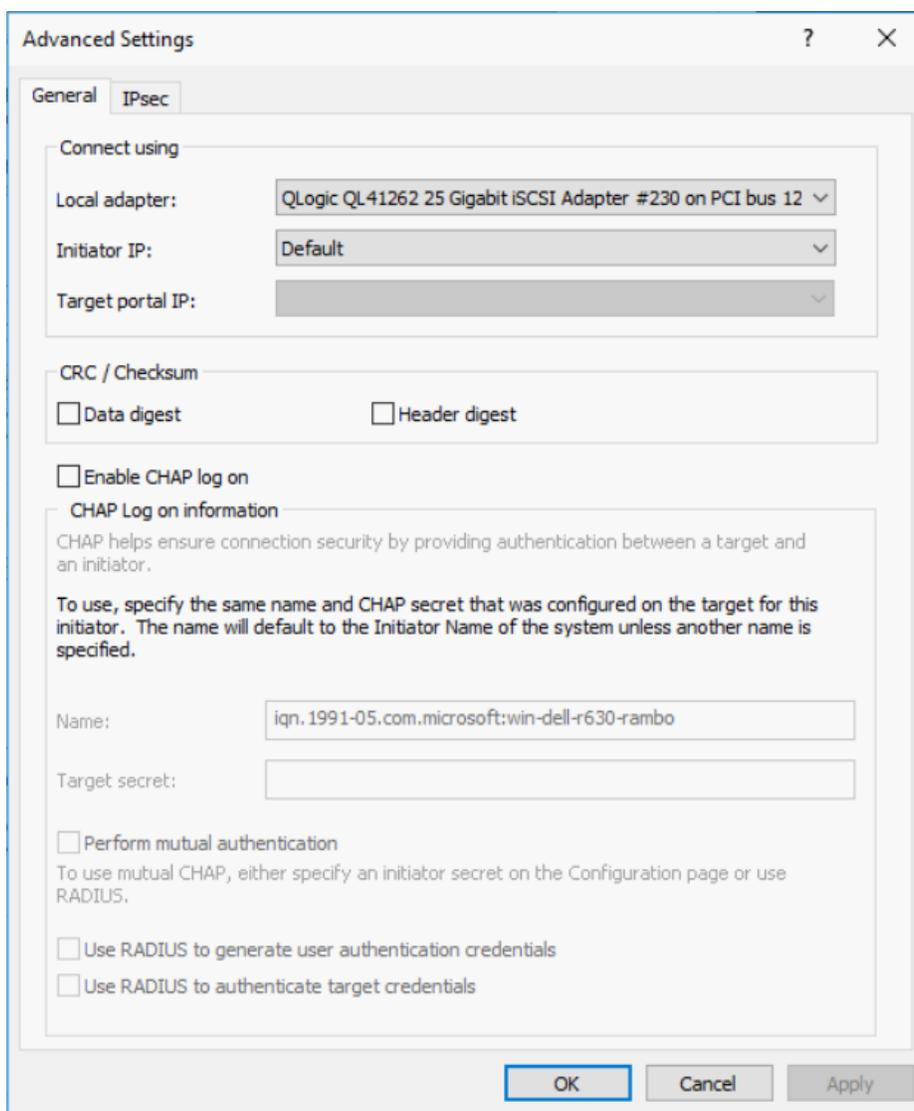


图 9-20. 目标门户 IP 地址

6. 在 Advanced Settings（高级设置）对话框（图 9-21）中，填写 **Connect using**（连接方式）下的以下各项：
  - a. 对于 **Local adapter**（本地适配器），选择 **QLogic <name or model> Adapter**（QLogic <名称或型号> 适配器）。
  - b. 对于 **Initiator IP**（启动器 IP），选择适配器 IP 地址。

- c. 单击 **OK**（确定）。



**图 9-21. 选择启动器 IP 地址**

7. 在 iSCSI Initiator Properties（iSCSI 启动器属性）、Discovery（查找）页面中，单击 **OK**（确定）。

- 单击 **Targets**（目标）选项卡，然后在 Targets（目标）页面（图 9-22）上，单击 **Connect**（连接）。

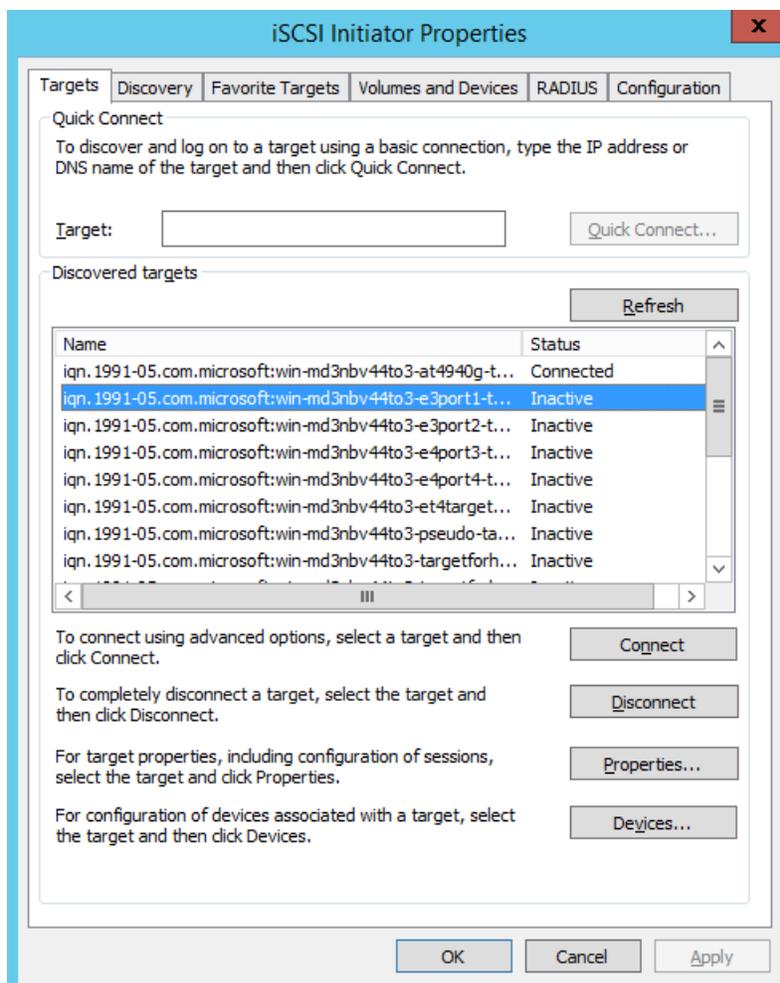


图 9-22. 连接到 iSCSI 目标

- 在 Connect To Target（连接到目标）对话框（图 9-23）中，单击 **Advanced**（高级）。

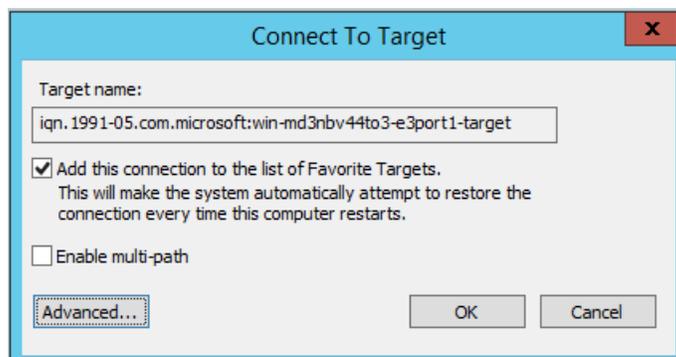


图 9-23. 连接到目标对话框

- 在 Local Adapter（本地适配器）对话框中，选择 **QLogic <name or model> Adapter**（QLogic <名称或型号> 适配器），然后单击 **OK**（确定）。
- 再次单击 **OK**（确定）关闭 Microsoft 启动器。
- 要格式化 iSCSI 分区，请使用磁盘管理器。

### 注

组合功能的一些限制包括：

- 组合不支持 iSCSI 适配器。
- 组合不支持位于引导路径中的 NDIS 适配器。
- 组合支持不位于 iSCSI 引导路径中的 NDIS 适配器，但仅用于 SLB 组类型。

## iSCSI 卸载常见问题

关于 iSCSI 卸载的一些常见问题包括：

问： 如何为 iSCSI 卸载分配 IP 地址？

答： 使用 QConvergeConsole GUI 中的 Configurations（配置）页面。

问： 创建到目标的连接时应使用哪些工具？

答： 使用 Microsoft iSCSI 软件启动器（版本 2.08 或以上）。

问： 怎样知道连接已卸载？

答： 使用 Microsoft iSCSI 软件启动器。在命令行中，键入 `oiscsicli sessionlist`。从 **Initiator Name**（启动器名称），iSCSI 卸载的连接将显示以 `B06BDRV` 开始的条目。非卸载的连接将显示以 `Root` 开始的条目。

问： 哪些配置应避免？

答： IP 地址不能与 LAN 相同。

## Windows Server 2012 R2 和 2016 iSCSI 引导安装

Windows Server 2012 R2 和 2016 支持在卸载或非卸载路径中引导和安装。QLogic 要求使用滑流 DVD，同时注入最新的 QLogic 驱动程序。请参阅 [第 167 页上的“将适配器驱动程序注入（滑流至）Windows 映像文件中”](#)。

以下过程准备通过卸载或非卸载路径安装和引导的映像。

### 要设置 Windows Server 2012R2/2016 iSCSI 引导：

1. 从要引导的系统（远程系统）上移除所有本地硬盘驱动器。
2. 通过按照 [第 167 页上的“将适配器驱动程序注入（滑流至）Windows 映像文件中”](#) 中的滑流步骤进行操作，准备 Windows 操作系统安装介质。
3. 将最新的 QLogic iSCSI 引导映像加载到适配器的 NVRAM 中。
4. 配置 iSCSI 目标以允许从远程设备连接。确保目标有足够磁盘空间安装新的操作系统。
5. 配置 UEFI HII 以设置 iSCSI 引导类型（卸载或非卸载）、正确的启动器和 iSCSI 引导的目标参数。
6. 保存设置并重新引导系统。远程系统应连接至 iSCSI 目标，然后从 DVD-ROM 设备引导。
7. 从 DVD 引导并开始安装。
8. 请按照屏幕说明进行操作。  
在显示可用于安装的磁盘列表的窗口中，应看到 iSCSI 目标磁盘。此目标是位于远程 iSCSI 目标中且通过 iSCSI 引导协议连接的磁盘。
9. 要继续 Windows Server 2012R2/2016 安装，请单击 **Next**（下一步），然后按照屏幕说明进行操作。作为安装过程的一部分，服务器将进行多次重新引导。
10. 您应该在服务器引导至操作系统后，运行驱动程序安装程序以完成 QLogic 驱动程序和应用程序安装。

## iSCSI 故障转储

41xxx 系列适配器的非卸载和卸载 iSCSI 引导均支持崩溃转储功能。配置 iSCSI 故障转储生成无需任何额外的配置。

## Linux 环境中的 iSCSI 卸载

QLogic FastLinQ 41xxx iSCSI 软件包含一个名为 `qed.ko` (`qed`) 的内核模块。`qed` 模块依赖于 Linux 内核的其他部分来实现特定功能：

- `qed.ko` 是用于常见 QLogic FastLinQ 41xxx 硬件初始化例程的 Linux eCore 内核模块。
- `scsi_transport_iscsi.ko` 是用于上行调用和下行调用会话管理的 Linux iSCSI 传输库。
- `libiscsi.ko` 是协议数据单元 (PDU) 和任务处理，以及会话内存管理所需的 Linux iSCSI 库功能。
- `iscsi_boot_sysfs.ko` 是为导出 iSCSI 引导信息提供帮助的 Linux iSCSI `sysfs` 接口。
- `uio.ko` 是 Linux 用户空间 I/O 接口，用于 `iscsiuio` 的轻型 L2 内存映射。

必须先加载这些模块，然后 `qed` 才能正常工作。否则，您可能会遇到“未解析的符号”错误。如果 `qed` 模块安装在分发版更新路径中，则必需模块通过 `modprobe` 自动加载。

本章节提供关于 iSCSI 在 Linux 中卸载的以下信息：

- [与 bnx2i 的差异](#)
- [配置 qed.ko](#)
- [在 Linux 中验证 iSCSI 接口](#)
- [Open-iSCSI 和从 SAN 引导注意事项](#)

## 与 bnx2i 的差异

QLogic FastLinQ 41xxx 系列适配器 (iSCSI) 的驱动程序 `qed` 与以前 QLogic 8400 系列适配器的 QLogic iSCSI 卸载驱动程序 `bnx2i` 之间存在一些重要的差异。其中一些差异包括：

- `qed` 直接绑定至 CNA 公开的 PCI 功能。
- `qed` 不会位于 `net_device` 的顶部。
- `qed` 不依赖于网络驱动程序（例如 `bnx2x` 和 `cnic`）。
- `qed` 不依赖于 `cnic`，但依赖于 `qed`。

- qedi 使用 `iscsi_boot_sysfs.ko` 负责导出 `sysfs` 中的引导信息，而从 SAN 的 `bnx2i` 引导依赖于 `iscsi_ibft.ko` 模块导出引导信息。

## 配置 qedi.ko

qedi 驱动程序自动绑定至 CNA 公开的 iSCSI 功能，且通过 `open-iscsi` 工具进行目标发现和绑定。此功能和操作与 `bnx2i` 驱动程序类似。

---

### 注

有关如何安装 FastLinQ 驱动程序的更多信息，请参阅 [第 3 章 驱动程序安装](#)。

---

要加载 `qedi.ko` 内核模块，请发出以下命令：

```
# modprobe qed
# modprobe libiscsi
# modprobe uio
# modprobe iscsi_boot_sysfs
# modprobe qedi
```

## 在 Linux 中验证 iSCSI 接口

安装和加载 `qedi` 内核模块后，必须验证是否正确检测到 iSCSI 接口。

要在 Linux 中验证 iSCSI 接口：

1. 要验证是否已主动加载 `qedi` 和关联的内核模块，请发出以下命令：

```
# 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. 要验证是否正确检测到 iSCSI 接口，请发出以下命令。在本例中，检测到的两个 iSCSI CNA 设备具有 SCSI 主机编号 4 和 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
```

#### 3. 使用 open-iscsi 工具验证 IP 是否正确配置。发出以下命令：

```
# iscsiadm -m iface | grep qedi
qedi.00:0e:1e:c4:e1:6d
qedi,00:0e:1e:c4:e1:6d,192.168.101.227,<empty>,iqn.1994-05.com.redhat:534ca9b6
adf
qedi.00:0e:1e:c4:e1:6c
qedi,00:0e:1e:c4:e1:6c,192.168.25.91,<empty>,iqn.1994-05.com.redhat:534ca9b6adf
```

#### 4. 要确保 iscsiui0 服务正在运行，请发出以下命令：

```
# systemctl status iscsiui0.service
iscsiui0.service - iSCSI UserSpace I/O driver
Loaded: loaded (/usr/lib/systemd/system/iscsiui0.service; disabled; vendor
preset: disabled)
Active: active (running) since Fri 2017-01-27 16:33:58 IST; 6 days ago
Docs: man:iscsiui0(8)
Process: 3745 ExecStart=/usr/sbin/iscsiui0 (code=exited, status=0/SUCCESS)
Main PID: 3747 (iscsiui0)
CGroup: /system.slice/iscsiui0.service !--3747 /usr/sbin/iscsiui0
Jan 27 16:33:58 localhost.localdomain systemd[1]: Starting iSCSI
UserSpace I/O driver...
Jan 27 16:33:58 localhost.localdomain systemd[1]: Started iSCSI UserSpace
I/O driver.
```

#### 5. 要查找 iSCSI 目标，请发出 iscsiadm 命令：

```
#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
```

6. 使用在 [步骤 5](#) 获取的 IQN 登录到 iSCSI 目标。要启动登录过程，请发出以下命令（其中命令中最后一个字符为小写字母“L”）：

```
#iscsiadm -m node -p 192.168.25.100 -T
iqn.2003-04.com.sanblaze:virtualun.virtualun.target-0)000007 -l
Logging in to [iface: qedi.00:0e:1e: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:1e:c4:e1:6c, target:iqn.2003-
04.com.sanblaze:virtualun.virtualun.target-05000007,
portal:192.168.25.100,3260] successful.
```

7. 要验证 iSCSI 会话已创建，请发出以下命令：

```
# iscsiadm -m session
qedi: [297] 192.168.25.100:3260,1
iqn.2003-04.com.sanblaze:virtualun.virtualun.target-05000007 (non-flash)
```

8. 要检查 iSCSI 设备，请发出 `iscsiadm` 命令：

```
# iscsiadm -m session -P3
...
*****
Attached SCSI devices:
*****
Host Number: 59 State: running
scsi59 Channel 00 Id 0 Lun: 0
Attached scsi disk sdb State: running scsi59 Channel 00 Id 0 Lun: 1
Attached scsi disk sdc State: running scsi59 Channel 00 Id 0 Lun: 2
Attached scsi disk sdd State: running scsi59 Channel 00 Id 0 Lun: 3
Attached scsi disk sde State: running scsi59 Channel 00 Id 0 Lun: 4
Attached scsi disk sdf State: running
```

有关高级目标配置，请参阅 Open-iSCSI 自述文件，网址为：

<https://github.com/open-iscsi/open-iscsi/blob/master/README>

## Open-iSCSI 和从 SAN 引导注意事项

在最新分发版（例如，RHEL 6/7 和 SLE 11/12）中，内建 iSCSI 用户空间公用程序（Open-iSCSI 工具）不支持 qedi iSCSI 传输，无法执行用户空间启动的 iSCSI 功能。从 SAN 引导安装期间，您可以使用驱动程序更新磁盘 (DUD) 来更新 qedi 驱动程序。不过，不存在用于更新用户空间内建公用程序的接口或进程，这会导致 iSCSI 目标登录和从 SAN 引导安装失败。

要克服这一限制，在从 SAN 引导期间，通过以下步骤使用纯 L2 接口（不使用硬件卸载 iSCSI）执行从 SAN 的初始引导。

#### 用 Dell OEM 解决方案使用软件启动器从 SAN 引导：

1. 在 NIC 配置页面上，选择 **Boot Protocol**( 引导协议 )，然后按 ENTER 键选择 **Legacy PXE**( 传统 PXE)。
2. 配置启动器和目标条目。
3. 在安装开始时，通过 DUD 选项传递以下引导参数：
  - 对于 RHEL 6.x 和 7.x：

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

旧分发版的 RHEL 无需单独的选项。
  - 对于 SLES 11 SP4 和 SLES 12 SP1/SP2/SP3：

```
ip=ibft dud=1
```
  - 对于 FastLinQ DUD 包（例如，在 RHEL 7 中）：

```
fastlinq-8.18.10.0-dd-rhel7u3-3.10.0_514.el7-x86_64.iso
```

其中 RHEL 7.x 的 DUD 参数是 `dd`，SLES 12.x 的是 `dud=1`。
4. 在目标 LUN 上安装操作系统。
5. 按照您操作系统的说明，从非卸载界面迁移到卸载界面：
  - [从 SAN 迁移的 RHEL 6.9 iSCSI L4 引导](#)
  - [从 SAN 迁移的 RHEL 7.2/7.3 iSCSI L4 引导](#)
  - [从 SAN 迁移的 SLES 11 SP4 iSCSI L4 引导](#)
  - [从 SAN 迁移的 SLES 12 SP1/SP2 iSCSI L4 引导](#)
  - [使用 MPIO 从 SAN 迁移的 SLES 12 SP1/SP2 iSCSI L4 引导](#)

### 从 SAN 迁移的 RHEL 6.9 iSCSI L4 引导

#### 要从非卸载接口迁移到卸载接口：

1. 引导到 iSCSI non-offload/L2 boot-from-SAN 操作系统。发出以下命令来安装 `open-iscsi` 和 `iscsiuio` RPM：

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

2. 编辑 `/etc/init.d/iscsid` 文件, 添加以下命令, 然后保存文件:

```
modprobe -q qedi
```

例如:

```
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 iscsiuiio
```

3. 编辑 `/etc/iscsi/iscsid.conf` 文件, 对以下行进行注释或取消注释, 然后保存文件:

注释:

```
iscsid.startup = /etc/rc.d/init.d/iscsid force-start
```

取消注释:

```
iscsid.startup = /sbin/iscsid
```

例如:

```
#####  
# iscsid daemon config  
#####  
# If you want iscsid to start the first time a iscsi tool  
# needs to access it, instead of starting it when the init  
# scripts run, set the iscsid startup command here.This  
# should normally only need to be done by distro package  
# maintainers.  
#  
# Default for Fedora and RHEL.(uncomment to activate).  
#iscsid.startup = /etc/rc.d/init.d/iscsid force-start  
#  
# Default for upstream open-iscsi scripts (uncomment to  
activate).  
iscsid.startup = /sbin/iscsid
```

4. Create an lface record for an L4 interface. 发出以下命令:

```
# iscsiadm -m iface -I qedi.14:02:ec:ce:dc:71 -o new  
New interface qedi.14:02:ec:ce:dc:71 added
```

iface 记录格式应该是 `qedi.<mac_address>`。在这种情况下，MAC 地址应与 iSCSI 会话处于活动状态的 L4 MAC 地址匹配。

5. 通过发出 `iscsiadm` 命令来更新 `iface` 记录中的 `iface` 字段。例如：

```
# iscsiadm -m iface -I qedi.14:02:ec:ce:dc:71 -n iface.hwaddress -v 14:02:ec:ce:dc:71 -o update
qedi.14:02:ec:ce:dc:71 updated.
# iscsiadm -m iface -I qedi.14:02:ec:ce:dc:71 -n iface.transport_name -v qedi -o update
qedi.14:02:ec:ce:dc:71 updated.
# iscsiadm -m iface -I qedi.14:02:ec:ce:dc:71 -n iface.bootproto -v dhcp -o update
qedi.14:02:ec:ce:dc:71 updated.
# iscsiadm -m iface -I qedi.14:02:ec:ce:dc:71 -n iface.ipaddress -v 0.0.0.0 -o update
qedi.14:02:ec:ce:dc:71 updated.
# iscsiadm -m node -T iqn.1986-03.com.hp:storage.p2000g3.13491b47fb -p 192.168.100.9:3260 -I
qedi.14:02:ec:ce:dc:71 -o new
New iSCSI node [qedi:[hw=14:02:ec:ce:dc:71,ip=0.0.0.0,net_if=,iscsi_if=qedi.14:02:ec:ce:dc:71]
192.168.100.9,3260,-1 iqn.1986-03.com.hp:storage.p2000g3.13491b47fb] added
```

6. 编辑 `/boot/efi/EFI/redhat/grub.conf` 文件，进行以下更改并保存该文件：

- 删除 `ifname=eth5:14:02:ec:ce:dc:6d`
- 删除 `ip=ibft`
- 添加 `selinux=0`

例如：

```
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=latacyrheb-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=latacyrheb-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
```

7. 通过发出以下命令建立 `initramfs` 文件：  

```
# dracut -f
```
8. 重新启动服务器，然后打开 HII。
9. 在 HII 中，启用 iSCSI 卸载模式：
  - a. 在主配置页面中，选择 **System Setup**（系统设置），**Device Settings**（设备设置）。
  - b. 在“设备设置”页面中，选择已配置 iSCSI 引导固件表 (iBFT) 的端口
  - c. 在“系统设置”页面中，选择 **NIC Partitioning Configuration**（NIC 分区配置），**Partition 3 Configuration**（分区 3 配置）。
  - d. 在分区 3 配置页面中，设置 **iSCSI Offload Mode**（iSCSI 卸载模式）为 **Enabled**（已启用）。
10. 在主配置页面上，选择 **iSCSI General Parameters**（iSCSI 常规参数），然后设置 **HBA Boot Mode**（HBA 引导模式）为 **Enabled**（已启用）。
11. 在“主配置”页面上，设置 **Boot Protocol**（引导协议）为 **UEFI iSCSI HBA**。
12. 保存配置并重新启动服务器。

---

### 注

现在，操作系统可通过卸载接口引导。

---

## 从 SAN 迁移的 RHEL 7.2/7.3 iSCSI L4 引导

要从非卸载接口迁移到卸载接口：

1. 通过发出以下命令更新 `open-iscsi` 工具和 `iscsiuio`：  

```
#rpm -ivh qlgc-open-iscsi-2.0_873.111.rhel7u3-3.x86_64.rpm --force  
#rpm -ivh iscsiuiio-2.11.5.3-2.rhel7u3.x86_64.rpm --force
```
2. 通过发出以下命令重新加载所有守护进程服务：  

```
#systemctl daemon-reload
```
3. 通过发出以下命令重新启动 `iscsid` 和 `iscsiuio` 服务：  

```
# systemctl restart iscsiuiio  
# systemctl restart iscsid
```
4. 通过发出以下命令为 L4 接口创建一个 `iface` 记录。  

```
# iscsiadm -m iface -I qedi.00:0e:1e:d6:7d:3a -o new
```

iface 记录格式应该是 `qedi<mac_address>`。在这种情况下，MAC 地址应与 iSCSI 会话处于活动状态的 L4 MAC 地址匹配。

5. 通过发出 `iscsiadm` 命令来更新 `iface` 记录中的 `iface` 字段。例如：

```
# iscsiadm -m iface -I qedi.00:0e:1e:d6:7d:3a -n iface.hwaddress -v 00:0e:1e:d6:7d:3a -o update
# iscsiadm -m iface -I qedi.00:0e:1e:d6:7d:3a -n iface.ipaddress -v 192.168.91.101 -o update
# iscsiadm -m iface -I qedi.00:0e:1e:d6:7d:3a -n iface.subnet_mask -v 255.255.0.0 -o update
# iscsiadm -m iface -I qedi.00:0e:1e:d6:7d:3a -n iface.transport_name -v qedi -o update
# iscsiadm -m iface -I qedi.00:0e:1e:d6:7d:3a -n iface.bootproto -v static -o update
```

6. 创建一个目标节点记录以使用 L4 接口，如下所示：

```
# iscsiadm -m node -T
iqn.2003-04.com.sanblaze:virtualun.virtualun.target-050123456
-p 192.168.25.100:3260 -I qedi.00:0e:1e:d6:7d:3a -o new
```

7. 编辑 `/usr/libexec/iscsi-mark-root-node` 文件并找到以下语句：

```
if [ "$transport" = bnx2i ]; then
start_iscsiuio = 1
```

添加 `|| [ "$transport" = qedi ]` 到 IF 表达式如下：

```
if [ "$transport" = bnx2i ] || [ "$transport" = qedi ]; then
start_iscsiuio = 1
```

8. 编辑 `/etc/default/grub` 文件并找到以下语句：

```
GRUB_CMDLINE_LINUX="iscsi_firmware ip=ibft"
```

将此语句更改为：

```
GRUB_CMDLINE_LINUX="rd.iscsi.firmware"
```

9. 通过发出以下命令创建新的 `grub.cfg` 文件：

```
# grub2-mkconfig -o /boot/efi/EFI/redhat/grub.cfg
```

10. 通过发出以下命令建立 `initramfs` 文件：

```
# dracut -f
```

11. 重新启动服务器，然后打开 HII。
12. 在 HII 中，启用 iSCSI 卸载模式：
  - a. 在主配置页面中，选择 **System Setup**（系统设置），**Device Settings**（设备设置）。
  - b. 在“设备设置”页面中，选择已配置 iSCSI 引导固件表 (iBFT) 的端口
  - c. 在“系统设置”页面中，选择 **NIC Partitioning Configuration**（NIC 分区配置），**Partition 3 Configuration**（分区 3 配置）。
  - d. 在分区 3 配置页面中，设置 **iSCSI Offload Mode**（iSCSI 卸载模式）为 **Enabled**（已启用）。
13. 在主配置页面上，选择 **iSCSI General Parameters**（iSCSI 常规参数），然后设置 **HBA Boot Mode**（HBA 引导模式）为 **Enabled**（已启用）。
14. 在“主配置”页面上，设置 **Boot Protocol**（引导协议）为 **UEFI iSCSI HBA**。
15. 保存配置并重新启动服务器。

---

### 注

现在，操作系统可通过卸载接口引导。

---

## 从 SAN 迁移的 SLES 11 SP4 iSCSI L4 引导

要从非卸载接口迁移到卸载接口：

1. 通过发出以下命令将 open-iscsi 工具和 iscsiui0 更新到最新可用版本：

```
# rpm -ivh qlgc-open-iscsi-2.0_873.111.sles11sp4-3.x86_64.rpm --force
# rpm -ivh iscsiui0-2.11.5.3-2.sles11sp4.x86_64.rpm --force
```

2. 编辑 `/etc/elilo.conf` 文件，进行以下更改，然后保存该文件：

- 删除 `ip=ibft` 参数（如果存在）
- 添加 `iscsi_firmware`

3. 编辑 `/etc/sysconfig/kernel` 文件并找到以下语句：

```
INITRD_MODULES="ata_piix ata_generic"
```

将语句更改为：

```
INITRD_MODULES="ata_piix ata_generic qedi"
```

保存文件。

4. 编辑 `/etc/modprobe.d/unsupported-modules` 文件，将 `allow_unsupported_modules` 的值更改为 1，并保存该文件：

```
allow_unsupported_modules 1
```
5. 找到并删除以下文件：
  - `/etc/init.d/boot.d/K01boot.open-iscsi`
  - `/etc/init.d/boot.open-iscsi`
6. 创建 `initrd` 的备份，然后通过发出以下命令来构建新的 `initrd`。

```
# cd /boot/  
# mkinitrd
```
7. 重新启动服务器，然后打开 HII。
8. 在 HII 中，启用 iSCSI 卸载模式：
  - a. 在主配置页面中，选择 **System Setup**（系统设置），**Device Settings**（设备设置）。
  - b. 在“设备设置”页面中，选择已配置 iSCSI 引导固件表 (iBFT) 的端口
  - c. 在“系统设置”页面中，选择 **NIC Partitioning Configuration**（NIC 分区配置），**Partition 3 Configuration**（分区 3 配置）。
  - d. 在分区 3 配置页面中，设置 **iSCSI Offload Mode**（iSCSI 卸载模式）为 **Enabled**（已启用）。
9. 在主配置页面上，选择 **iSCSI General Parameters**（iSCSI 常规参数），然后设置 **HBA Boot Mode**（HBA 引导模式）为 **Enabled**（已启用）。
10. 在“主配置”页面上，设置 **Boot Protocol**（引导协议）为 **UEFI iSCSI HBA**。
11. 保存配置并重新启动服务器。

---

### 注

现在，操作系统可通过卸载接口引导。

---

## 从 SAN 迁移的 SLES 12 SP1/SP2 iSCSI L4 引导

要从非卸载接口迁移到卸载接口：

1. 引导到 iSCSI non-offload/L2 boot-from-SAN 操作系统。发出以下命令来安装 `open-iscsi` 和 `iscsiuio` RPM：

```
# qlgc-open-iscsi-2.0_873.111.slessp2-3.x86_64.rpm  
# iscsiuio-2.11.5.3-2.sles12sp2.x86_64.rpm
```

2. 通过发出以下命令重新加载所有守护进程服务：  

```
# systemctl daemon-reload
```
3. 如果尚未通过发出以下命令启用 iscsid 和 iscsiui0 服务，请启用 iscsid 和 iscsiui0 服务：  

```
# systemctl enable iscsid  
# systemctl enable iscsiui0
```
4. 发出以下命令：  

```
cat /proc/cmdline
```
5. 检查操作系统是否保留了任何引导选项，诸如 ip=ibft 或 rd.iscsi.ibft。
  - 如果有保留的引导选项，继续步骤 6。
  - 如果没有保留的引导选项，跳至步骤 6c。
6. 编辑 /etc/default/grub 文件并修改 GRUB\_CMDLINE\_LINUX 值：
  - a. 移除 rd.iscsi.ibft（如果存在）。
  - b. 移除任何 ip=<value> 引导选项（如果存在）。
  - c. 添加 rd.iscsi.firmware。对于较老的发行版，添加 iscsi\_firmware。
7. 创建原文件 grub.cfg 的备份。该文件位于以下位置：
  - 传统引导：/boot/grub2/grub.cfg
  - UEFI 引导：SLES 的 /boot/efi/EFI/sles/grub.cfg
8. 通过发出以下命令创建新的 grub.cfg 文件：  

```
# grub2-mkconfig -o <new file name>
```
9. 比较旧的 grub.cfg 文件与新的 grub.cfg 文件来验证您的更改。
10. 用新的 grub.cfg 文件替换原 grub.cfg 文件。
11. 通过发出以下命令建立 initramfs 文件：  

```
# dracut -f
```
12. 重新启动服务器，然后打开 HII。
13. 在 HII 中，启用 iSCSI 卸载模式：
  - a. 在主配置页面中，选择 **System Setup**（系统设置），**Device Settings**（设备设置）。

- b. 在“设备设置”页面中，选择已配置 iSCSI 引导固件表 (iBFT) 的端口
  - c. 在“系统设置”页面中，选择 **NIC Partitioning Configuration** (NIC 分区配置)，**Partition 3 Configuration** (分区 3 配置)。
  - d. 在分区 3 配置页面中，设置 **iSCSI Offload Mode** (iSCSI 卸载模式) 为 **Enabled** (已启用)。
14. 在主配置页面上，选择 **iSCSI General Parameters** (iSCSI 常规参数)，然后设置 **HBA Boot Mode** (HBA 引导模式) 为 **Enabled** (已启用)。
  15. 在“主配置”页面上，设置 **Boot Protocol** (引导协议) 为 **UEFI iSCSI HBA**。
  16. 保存配置并重新启动服务器。

---

### 注

现在，操作系统可通过卸载接口引导。

---

## 使用 MPIO 从 SAN 迁移的 SLES 12 SP1/SP2 iSCSI L4 引导

要从 L2 迁移到 L4 并配置 Microsoft 多路径 I/O (MPIO) 设置以通过卸载界面引导操作系统：

1. 要更新 open-iscsi 工具，请发出以下命令：

```
# rpm -ivh --force qlgc-open-iscsi-2.0_873.111.sles12sp1-3.x86_64.rpm
# rpm -ivh --force iscsiuiio-2.11.5.3-2.sles12sp1.x86_64.rpm
```

2. 转到 `/etc/default/grub` 并将 `rd.iscsi.ibft` 参数更改为 `rd.iscsi.firmware`。
3. 发出以下命令：  

```
grub2-mkconfig -o /boot/efi/EFI/suse/grub.cfg
```
4. 要加载多路径模块，请发出以下命令：  

```
modprobe dm_multipath
```
5. 要启用多路径守护进程，请发出以下命令：  

```
systemctl start multipathd.service
systemctl enable multipathd.service
systemctl start multipathd.socket
```
6. 要将设备添加到多路径，请发出以下命令：  

```
multipath -a /dev/sda
multipath -a /dev/sdb
```

7. 要运行多路径实用程序，请发出以下命令：  

```
multipath （可能不会显示多路径设备，因为它使用 L2 上的单个路径引导）  
multipath -ll
```
8. 要在 initrd 中插入多路径模块，请发出以下命令：  

```
dracut --force --add multipath --include /etc/multipath
```
9. 重新启动服务器并通过在 POST 菜单中按 F9 键进入系统设置。
10. 更改 UEFI 配置以使用 L4 iSCSI 引导：
  - a. 打开系统设置窗口，然后选择 **Device Settings**（设备设置）。
  - b. 在“设备设置”窗口中，选择配置了 iSCSI 引导固件表 (iBFT) 的适配器端口，然后按 Enter。
  - c. 在 Main Configuration Page（主要配置页面）中，选择 **NIC Partitioning Configuration**（NIC 分区配置），然后按 ENTER 键。
  - d. 在“分区配置”页面上，选择 **Partition 3 Configuration**（分区 3 配置）。
  - e. 在分区 3 配置页面中，设置 **iSCSI Offload Mode**（iSCSI 卸载模式）为 **Enabled**（已启用）。
  - f. 转到主配置页面，然后选择 **iSCSI Configuration**（iSCSI 配置）。
  - g. 在“iSCSI 配置”页面上，选择 **iSCSI General Parameters**（iSCSI 常规参数）。
  - h. 在“iSCSI 常规参数”页面上，将 **HBA Boot Mode**（HBA 引导模式）设置为 **Enabled**（已启用）。
  - i. 进入主配置页面，然后选择 **NIC Configuration**（NIC 配置）。
  - j. 在“NIC 配置”页面上，设置 **Boot Protocol**（引导协议）为 **UEFI iSCSI HBA**。
  - k. 保存设置，然后退出系统配置菜单。
11. 为确保从驱动程序更新磁盘 (DUD) 正确安装开箱即用驱动程序并防止加载自带的驱动程序，请执行以下操作：
  - a. 编辑 `/etc/default/grub` 文件以包含以下命令：

```
BOOT_IMAGE=/boot/x86_64/loader/linux dud=1  
brokenmodules=qed,qedi,qedf linuxrc.debug=1
```
  - b. 编辑 DUD 上的 `dud.config` 文件并添加以下命令以清除损坏的模块列表：

```
brokenmodules=-qed,qedi,qedf  
brokenmodules=dummy_xxx
```

12. 重新启动系统。现在，操作系统应该通过卸载接口引导。

## 从 RHEL 7.4 及更高版本的 SAN 配置 iSCSI 引导

### 要安装 RHEL 7.4 及更高版本：

1. 从 RHEL 7.x 安装介质引导，并且 iSCSI 目标已经在 UEFI 中连接。

```
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. 要安装开箱即用驱动程序，请键入 `e`。否则，继续[步骤 7](#)。

3. 选择内核行，然后键入 `e`。

4. 发出以下命令，然后按 ENTER。

```
linux dd modprobe.blacklist=qed modprobe.blacklist=qede  
modprobe.blacklist=qedr modprobe.blacklist=qedi  
modprobe.blacklist=qedf
```

您可以使用 `inst.dd` 选项而不是 `linux dd`。

5. 安装过程会提示您安装开箱即用驱动程序，如示例 [图 9-24](#) 中所示。

```

Starting Driver Update Disk UI on tty1...
[ OK ] Started Show Plymouth Boot Screen.
[ OK ] Reached target Paths.
[ OK ] Reached target Basic System.
[ OK ] Started Device-Mapper Multipath Device Controller.
Starting Open-iSCSI...
[ OK ] Started Open-iSCSI.
Starting dracut initqueue hook...
[ OK ] Created slice system-driver\x2dupdates.slice.
Starting Driver Update Disk UI on tty1...
DD: starting interactive mode

(Page 1 of 1) Driver disk device selection
  /DEVICE  TYPE  LABEL  UUID
  1) sda1   ntfs   1A90FE4090FE2245
  2) sda2   ufat   A6FF-80A4
  3) sda4   ntfs   7490015F900128E6
  4) sr0    iso9660 2017-07-11-01-39-24-00
# to select, 'r'-refresh, or 'c'-continue: r

(Page 1 of 1) Driver disk device selection
  /DEVICE  TYPE  LABEL  UUID
  1) sda1   ntfs   Recovery 1A90FE4090FE2245
  2) sda2   ufat   A6FF-80A4
  3) sda4   ntfs   7490015F900128E6
  4) sr0    iso9660 CDROM    2017-07-11-13-08-37-00
# to select, 'r'-refresh, or 'c'-continue: 4
DD: Examining /dev/sr0
mount: /dev/sr0 is write-protected, mounting read-only

(Page 1 of 1) Select drivers to install
  1) [ ] /media/DD-1/rpms/x86_64/kmod-qlgc-fastlinq-8.22.0.0-1.rhel7u4.x86_64.rpm
# to toggle selection, or 'c'-continue: 1

(Page 1 of 1) Select drivers to install
  1) [x] /media/DD-1/rpms/x86_64/kmod-qlgc-fastlinq-8.22.0.0-1.rhel7u4.x86_64.rpm
# to toggle selection, or 'c'-continue: c
DD: Extracting: kmod-qlgc-fastlinq

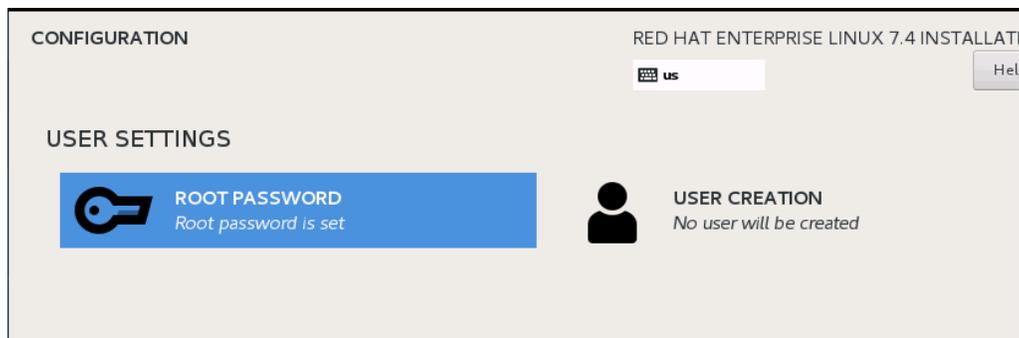
(Page 1 of 1) Driver disk device selection
  /DEVICE  TYPE  LABEL  UUID
  1) sda1   ntfs   Recovery 1A90FE4090FE2245
  2) sda2   ufat   A6FF-80A4
  3) sda4   ntfs   7490015F900128E6
  4) sr0    iso9660 CDROM    2017-07-11-13-08-37-00
# to select, 'r'-refresh, or 'c'-continue:

```

图 9-24. 提示开箱即用安装

6. 如果您的设置需要，请在提示其他驱动程序磁盘时加载 FastLinQ 驱动程序更新磁盘。否则，如果您没有其他驱动程序更新磁盘要安装，请键入 c。
7. 继续安装。您可以跳过介质测试。单击 **Next**（下一步）继续安装。

- 在配置窗口 (图 9-25) 中, 选择在安装过程中使用的语言, 然后单击 **Continue** (继续)。



**图 9-25. Red Hat Enterprise Linux 7.4 配置**

- 在 Installation Summary (安装摘要) 窗口, 单击 **Installation Destination** (安装目标)。磁盘标签是 `sda`, 指示单路径安装。如果您配置了多路径, 则该磁盘具有设备映射器标签。
- 在 **Specialized & Network Disks** 章节中, 选择 iSCSI LUN。
- 键入 root 用户的密码, 然后单击 **Next** (下一步) 完成安装。
- 在第一次引导期间, 添加以下内核命令行以落入 shell。

```
rd.iscsi.firmware rd.break=pre-pivot rd.driver.pre=qed,qede,
qedr,qedf,qedi selinux=0
```

- 发出以下命令:

```
# umount /sysroot/boot/efi
# umount /sysroot/boot/
# umount /sysroot/home/
# umount /sysroot
# mount /dev/mapper/rhel-root /sysroot
```

- 编辑 `/sysroot/usr/libexec/iscsi-mark-root-nodes` 文件, 并定位到以下语句:

```
if [ "$transport" = bnx2i ]; then
```

将语句更改为:

```
if [ "$transport" = bnx2i ] || [ "$transport" = qedi ]; then
```

- 通过发出以下命令卸载文件系统:

```
# umount /sysroot
```

16. 重新启动服务器，然后在命令行中添加以下参数：

```
rd.iscsi.firmware
rd.driver.pre=qed,qedi (以加载所有的驱动程序
pre=qed,qedi,qedi,qedf)
selinux=0
```

17. 成功引导系统后，编辑 `/etc/modprobe.d/anaconda-blacklist.conf` 文件，以删除所选驱动程序的黑名单条目。
18. 通过发出 `dracut -f` 命令重建虚拟磁盘，然后重新启动。

# 10 FCoE 配置

本章提供以下以太网光纤信道 (FCoE) 配置信息：

- [从 SAN 进行 FCoE 引导](#)
- [第 167 页上的“将适配器驱动程序注入（滑流至）Windows 映像文件中”](#)
- [第 168 页上的“配置 Linux FCoE 卸载”](#)
- [第 169 页上的“qedf 与 bnx2fc 之间的差异”](#)
- [第 170 页上的“配置 qedf.ko”](#)
- [第 170 页上的“在 Linux 中验证 FCoE 设备”](#)
- [第 171 页上的“从 RHEL 7.4 及更高版本的 SAN 配置 FCoE 引导”](#)

---

## 注

所有 41xxx 系列适配器 都支持 FCoE 卸载。某些 FCoE 功能在当前版本中可能并未完全启用。有关详细信息，请参阅 [附录 D 功能约束](#)。

---

## 从 SAN 进行 FCoE 引导

本节介绍 Windows、Linux 和 ESXi 操作系统的安装和引导步骤，包括：

- [为 FCoE 构建和引导来准备系统 BIOS](#)
- [Windows 版从 SAN 进行的 FCoE 引导](#)

---

## 注

ESXi 5.5 及更高版本支持从 SAN 进行 FCoE 引导。并非所有适配器版本均支持 FCoE 和从 SAN 进行的 FCoE 引导。

---

## 为 FCoE 构建和引导来准备系统 BIOS

要准备系统 BIOS，请修改系统引导顺序并指定 BIOS 引导协议（如果需要）。

### 指定 BIOS 引导协议

从 SAN 进行的 FCoE 引导仅在 UEFI 模式下受支持。使用系统 BIOS 配置将引导模式（协议）中的平台设置为 UEFI。

#### 注

FCoE BFS 在旧版 BIOS 模式下不受支持。

### 配置适配器 UEFI 引导模式

要将引导模式配置为 FCOE：

1. 重新启动系统。
2. 按 OEM 热键进入 System Setup（系统设置）（图 10-1）。这也称为 UEFI HII。

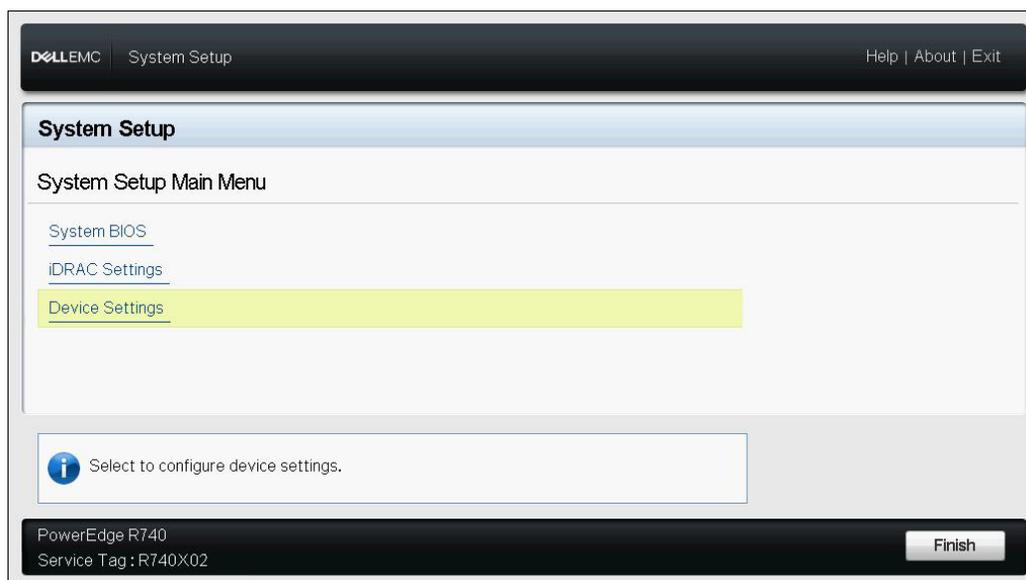


图 10-1. 系统设置：选择设备设置

#### 注

SAN 引导仅在 UEFI 环境中受支持。确保系统引导选项为 UEFI，而非旧版。

3. 在 Device Settings (设备设置) 页面中, 选择 QLogic 设备 (图 10-2)。

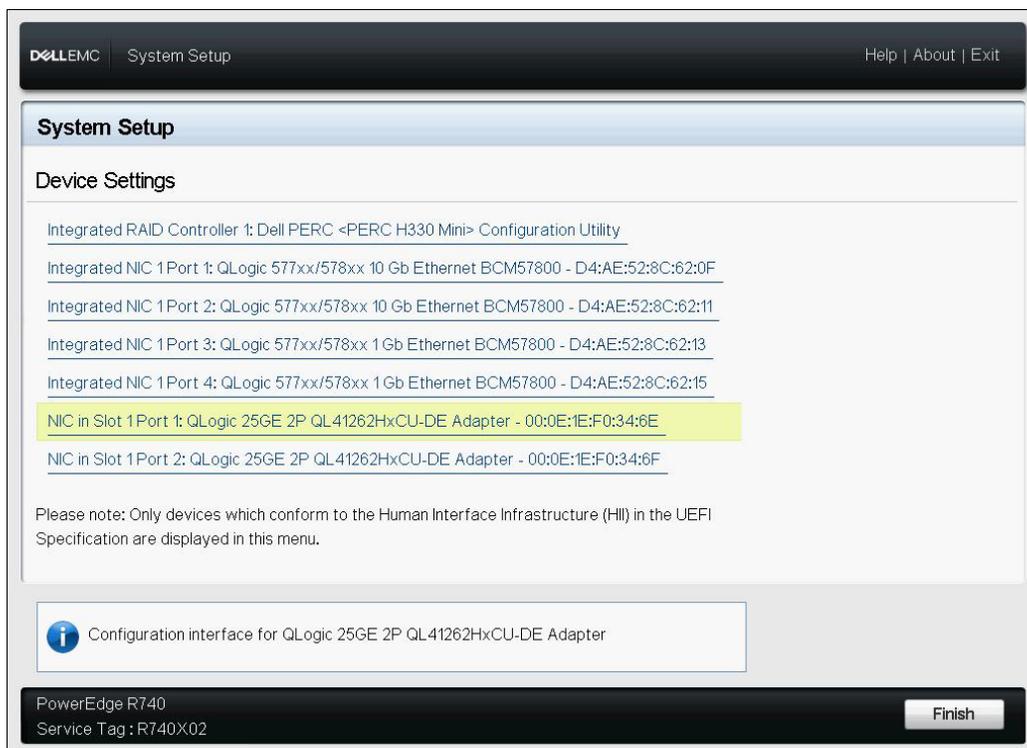


图 10-2. 系统设置: 设备设置, 端口选择

4. 在 Main Configuration Page（主要配置页面）中，选择 **NIC Configuration**（NIC 配置）（图 10-3），然后按 ENTER 键。

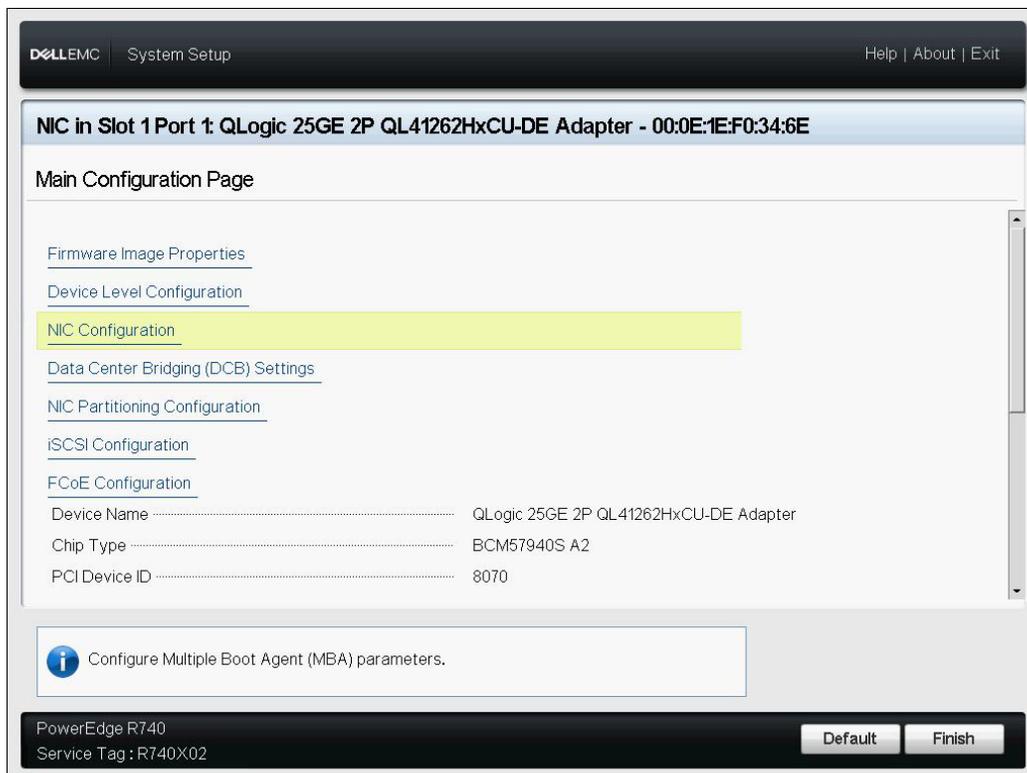


图 10-3. 系统设置：NIC 配置

5. 在 NIC Configuration（NIC 配置）页面上，选择 **Boot Mode**（引导模式），然后按 ENTER 键选择 **FCoE** 作为首选引导模式。

### 注

如果在端口级别禁用 **FCoE Mode**（FCoE 模式）功能，则 **FCoE** 不会作为引导选项列出。如果 **Boot Mode**（引导模式）首选为 **FCoE**，请确保 **FCoE Mode**（FCoE 模式）功能已启用，如图 10-4 中所示。并非所有适配器版本均支持 FCoE。

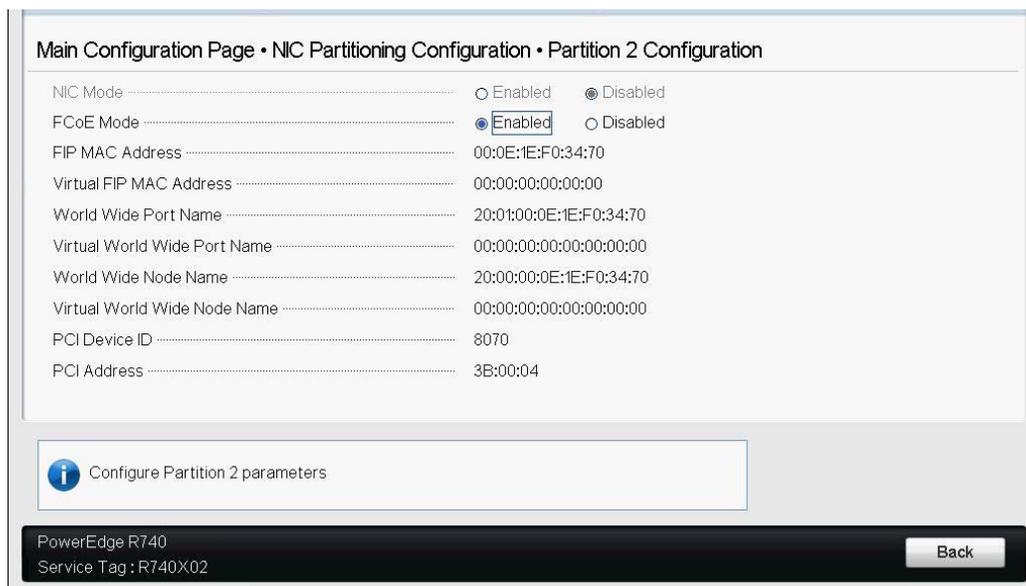
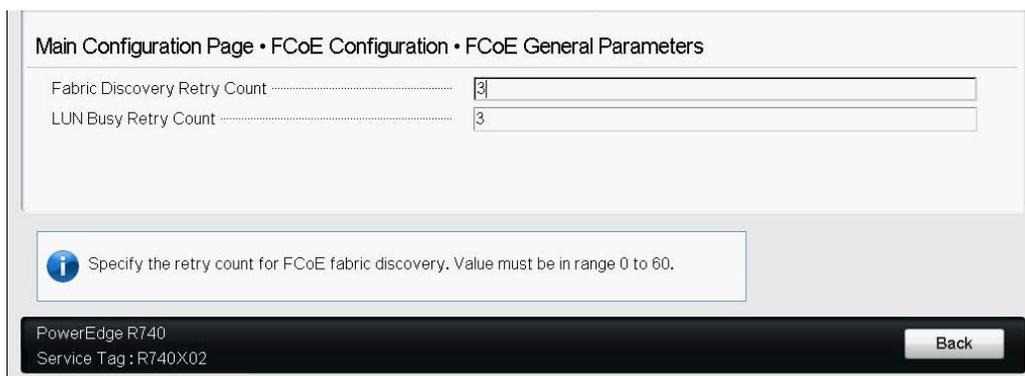


图 10-4. 系统设置：FCoE 模式已启用

**要配置 FCoE 引导参数：**

1. 在设备 HII 的 Main Configuration Page（主要配置页面）上，选择 **FCoE Configuration**（FCoE 配置），然后按 ENTER 键。
2. 在 FCoE Configuration（FCoE 配置）页面上，选择 **FCoE General Parameters**（FCoE 常规参数），然后按 ENTER 键。
3. 在 FCoE General Parameters（FCoE 常规参数）页面（图 10-5）上，按向上箭头键和向下箭头键选择参数，然后按 ENTER 键选择并输入以下值：
  - Fabric Discovery Retry Count**（结构查找重试计数）：默认值或根据需要
  - LUN Busy Retry Count**（LUN 繁忙重试次数）：默认值或根据需要



**图 10-5. 系统设置：FCoE 常规参数**

4. 返回 FCoE Configuration（FCoE 配置）页面。
5. 按 ESC 键，然后选择 **FCoE Target Parameters**（FCoE 目标参数）。
6. 按 ENTER 键。
7. 在 FCoE Target Parameters（FCoE 目标参数）菜单中，启用到首选 FCoE 目标的**连接**。
8. 键入 iSCSI 目标的以下参数的值（[图 10-6](#)），然后按 ENTER 键：
  - World Wide Port Name Target  $n$** （全局端口名称目标  $n$ ）
  - Boot LUN  $n$** （引导 LUN  $n$ ）

其中， $n$  的值介于 1 到 8 之间，使您能够配置 8 个 FCoE 目标。

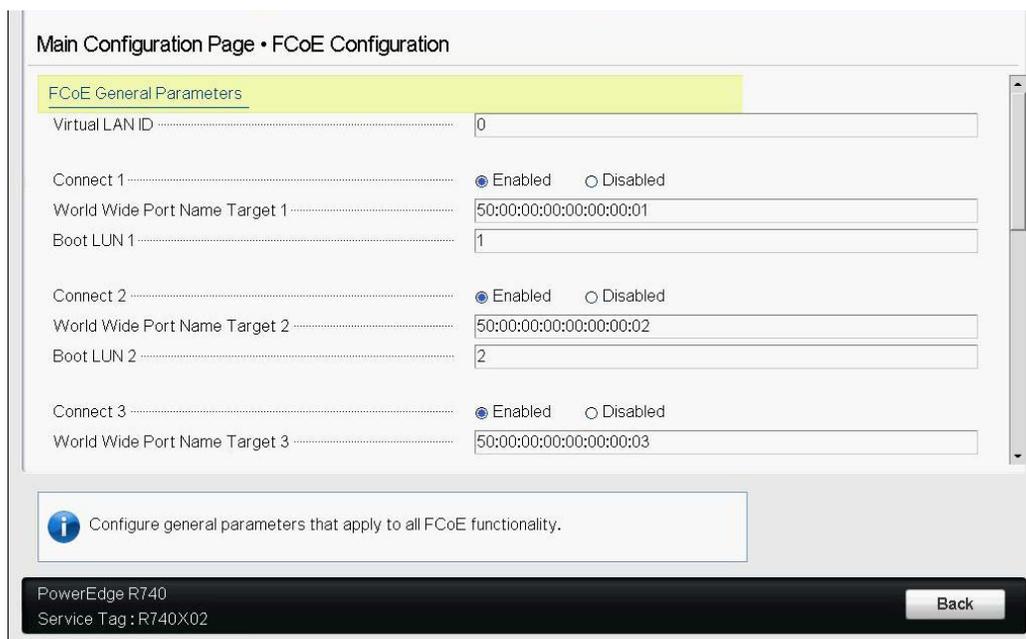


图 10-6. 系统设置：FCoE 常规参数

## Windows 版从 SAN 进行的 FCoE 引导

Windows 版从 SAN 进行 FCoE 引导信息包括：

- [Windows Server 2012 R2 和 2016 FCoE 引导安装](#)
- [配置 FCoE](#)
- [FCoE 故障转储](#)

### Windows Server 2012 R2 和 2016 FCoE 引导安装

对于 Windows Server 2012R2/2016 从 SAN 引导安装，QLogic 要求使用“滑流”DVD 或 ISO 映像，同时注入最新的 QLogic 驱动程序。请参阅第 167 页上的“将适配器驱动程序注入（滑流至）Windows 映像文件中”。

以下步骤准备在 FCoE 模式下安装和引导的映像。

#### 要设置 Windows Server 2012R2/2016 FCoE 引导：

1. 从要引导的系统（远程系统）上移除所有本地硬盘驱动器。
2. 通过按照第 167 页上的“将适配器驱动程序注入（滑流至）Windows 映像文件中”中的滑流步骤进行操作，准备 Windows 操作系统安装介质。
3. 将最新的 QLogic FCoE 引导映像加载到适配器 NVRAM 中。

4. 配置 FCoE 目标以允许从远程设备连接。确保目标有足够磁盘空间安装新的操作系统。
5. 配置 UEFI HII 以在所需的适配器端口上设置 FCoE 引导类型、正确的启动器和 FCoE 引导的目标参数。
6. 保存设置并重新引导系统。远程系统应连接至 FCoE 目标，然后从 DVD-ROM 设备引导。
7. 从 DVD 引导并开始安装。
8. 请按照屏幕说明进行操作。
9. 在显示可用于安装的磁盘列表的窗口上，应看到 FCoE 目标磁盘。此目标是通过 FCoE 引导协议连接的磁盘，位于远程 FCoE 目标中。
10. 要继续 Windows Server 2012R2/2016 安装，请选择 **Next**（下一步），然后按照屏幕说明进行操作。作为安装过程的一部分，服务器将进行多次重新引导。
11. 您应该在服务器引导至操作系统后，运行驱动程序安装程序以完成 QLogic 驱动程序和应用程序安装。

### 配置 FCoE

默认情况下，DCB 在 QLogic 41xxx FCoE 和 DCB 兼容 C-NIC 上已启用。QLogic 41xxx FCoE 需要启用 DCB 的接口。对于 Windows 操作系统，使用 QCC GUI 或命令行公用程序来配置 DCB 参数。

### FCoE 故障转储

故障转储功能目前支持 FastLinQ 41xxx 系列适配器的 FCoE 引导。

在 FCoE 引导模式下，FCoE 故障转储生成无需额外的配置。

## 将适配器驱动程序注入（滑流至）Windows 映像文件中

**要将适配器驱动程序注入 Windows 映像文件中：**

1. 获取适用 Windows Server 版本（2012、2012 R2 或 2016）的最新驱动程序包。
2. 将驱动程序包提取到工作目录：
  - a. 打开命令行会话，导航到包含驱动程序包的文件夹。
  - b. 要启动驱动程序安装程序，请发出以下命令：  

```
setup.exe /a
```

- c. 在 Network location (网络位置) 字段中, 输入驱动程序包要提取到的文件夹路径。例如, 键入 `c:\temp`。
- d. 按照驱动程序安装程序说明进行操作, 在指定的文件夹中安装驱动程序。在本例中, 驱动程序文件安装在以下位置:  
`c:\temp\Program File 64\QLogic Corporation\QDrivers`

3. 从 Microsoft 下载 Windows 评估和部署工具包 (ADK) 版本 10:

<https://developer.microsoft.com/en-us/windows/hardware/windows-assessment-deployment-kit>

4. 打开命令行会话 (以管理员权限), 浏览发行版本 CD 至 `Tools\Slipstream` 文件夹。
5. 找到 `slipstream.bat` 脚本文件, 然后发出以下命令:

```
slipstream.bat <path>
```

其中 `<path>` 是您在步骤 2 中指定的驱动器和子文件夹。例如:

```
slipstream.bat "c:\temp\Program Files 64\QLogic Corporation\QDrivers"
```

---

### 注

注意有关操作系统安装介质的以下信息:

- 操作系统安装介质预期为本地驱动器。不支持操作系统安装介质的网络路径。
- `slipstream.bat` 脚本将驱动程序组件注入到操作系统安装介质支持的所有 SKU 中。

- 
6. 刻录 DVD, 其中包含位于工作目录中的所生成驱动程序 ISO 映像文件。
  7. 使用新 DVD 安装 Windows Server 操作系统。

## 配置 Linux FCoE 卸载

Cavium FastLinQ 41xxx 系列适配器 FCoE 软件包含一个名为 `qedf.ko` (`qedf`) 的内核模块。`qedf` 模块依赖于 Linux 内核的其他部分来实现特定功能:

- `qed.ko` 是用于常见 Cavium FastLinQ 41xxx 硬件初始化例程的 Linux eCore 内核模块。
- `libfcoe.ko` 是执行 FCoE 转发器 (FCF) 请求和 FCoE 初始化协议 (FIP) 结构登录 (FLOGI) 所需的 Linux FCoE 内核库。

- **libfc.ko** 是多项功能所需的 Linux FC 内核库，这些功能包括：
  - 名称服务器登录和注册
  - rport 会话管理
- **scsi\_transport\_fc.ko** 是用于远程端口和 SCSI 目标管理的 Linux FC SCSI 传输库。

必须加载这些模块后 qedf 才能正常工作，否则可导致“未解析的符号”等错误。如果 qedf 模块安装在分发版更新路径中，则必要模块通过 modprobe 自动加载。Cavium 41xxx 系列适配器支持 FCoE 卸载。

---

**注**

使用 SLES 11 或 SLES 12 进行安装时，不需要参数 `withfcoe=1`，因为 41xxx 系列适配器不再需要软件 FCoE 守护进程。

---

## qedf 与 bnx2fc 之间的差异

Cavium FastLinQ 41xxx 10/25GbE 控制器 (FCoE) 的驱动程序 qedf 与以前的 Cavium FCoE 卸载驱动程序 bnx2fc 之间存在显著的差异。差异包括：

- qedf 直接绑定至 CNA 公开的 PCI 功能。
- qedf 无需 open-fcoe 用户空间工具（fipvlan、fcoemon、fcoeadm）即可启动查找。
- qedf 直接发出 FIP VLAN 请求而无需 fipvlan 公用程序。
- qedf 无需 fipvlan 为 fcoemon 创建的 FCoE 接口。
- qedf 不会位于 net\_device 的顶部。
- qedf 不依赖于网络驱动程序（例如 bnx2x 和 cnic）。
- qedf 将在链路正常工作时自动启动 FCoE 查找（因为它不依赖于 fipvlan 或 fcoemon 进行 FCoE 接口创建）。

---

**注**

FCoE 接口不再位于网络接口之上。Qedf 驱动程序自动创建独立于网络接口的 FCoE 接口。因此，FCoE 接口不会显示在安装程序的 FCoE 接口对话框中。相反，磁盘会自动显示为 SCSI 磁盘，与光纤信道驱动程序的工作方式类似。

---

## 配置 qedf.ko

qedf.ko 无需显式配置。驱动程序会自动绑定至 CNA 公开的 FCoE 功能并开始查找。相比于较旧的 bnx2fc 驱动程序，此功能与 QLogic 的 FC 驱动程序 qla2xx 的功能和运行类似。

### 注

有关 FastLinQ 驱动程序安装的更多信息，请参阅[第 3 章 驱动程序安装](#)。

---

加载 qedf.ko 内核模块会执行以下命令：

```
# modprobe qed  
# modprobe libfcoe  
# modprobe qedf
```

## 在 Linux 中验证 FCoE 设备

按照以下步骤操作，验证在安装和加载 qedf 内核模块后是否正确检测到 FCoE 设备。

**要在 Linux 中验证 FCoE 设备：**

1. 检查 lsmod，验证是否已加载 qedf 和关联的内核模块：

```
# 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. 检查 `dmesg`，验证是否正确检测到 FCoE 设备。在本示例中，检测到的两个 FCoE CNA 设备具有 SCSI 主机编号 4 和 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. 使用 `lsblk -s` 检查找到的 FCoE 设备:

```
# lsblk -s
NAME HCTL          TYPE VENDOR      MODEL          REV TRAN
sdb  5:0:0:0         disk SANBlaze    VLUN P2T1L0    V7.3 fc
sdc  5:0:0:1         disk SANBlaze    VLUN P2T1L1    V7.3 fc
sdd  5:0:0:2         disk SANBlaze    VLUN P2T1L2    V7.3 fc
sde  5:0:0:3         disk SANBlaze    VLUN P2T1L3    V7.3 fc
sdf  5:0:0:4         disk SANBlaze    VLUN P2T1L4    V7.3 fc
sdg  5:0:0:5         disk SANBlaze    VLUN P2T1L5    V7.3 fc
sdh  5:0:0:6         disk SANBlaze    VLUN P2T1L6    V7.3 fc
sdi  5:0:0:7         disk SANBlaze    VLUN P2T1L7    V7.3 fc
sdj  5:0:0:8         disk SANBlaze    VLUN P2T1L8    V7.3 fc
sdk  5:0:0:9         disk SANBlaze    VLUN P2T1L9    V7.3 fc
```

主机的配置信息位于 `/sys/class/fc_host/hostX`，其中 `x` 是 SCSI 主机的编号。在上例中，`x` 可能为 4 或 5。`hostX` 文件包含 FCoE 功能的属性，例如全局端口名称和结构 ID。

## 从 RHEL 7.4 及更高版本的 SAN 配置 FCoE 引导

### 要安装 RHEL 7.4 及更高版本:

1. 从 RHEL 7.x 安装介质引导，并且 FCoE 目标已经在 UEFI 中连接。

```
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. 要安装开箱即用驱动程序，请键入 **e**。否则，继续[步骤 7](#)。
3. 选择内核行，然后键入 **e**。
4. 发出以下命令，然后按 ENTER。

```
linux dd modprobe.blacklist=qed modprobe.blacklist=qede  
modprobe.blacklist=qedr modprobe.blacklist=qedi  
modprobe.blacklist=qedf
```

您可以使用 `inst.dd` 选项而不是 `linux dd`。

5. 安装过程会提示您安装开箱即用驱动程序，如示例图 10-7 中所示。

```

Starting Driver Update Disk UI on tty1...
[ OK ] Started Show Plymouth Boot Screen.
[ OK ] Reached target Paths.
[ OK ] Reached target Basic System.
[ OK ] Started Device-Mapper Multipath Device Controller.
Starting Open-iSCSI...
[ OK ] Started Open-iSCSI.
Starting dracut initqueue hook...
[ OK ] Created slice system-driver\x2dupdates.slice.
Starting Driver Update Disk UI on tty1...
DD: starting interactive mode

(Page 1 of 1) Driver disk device selection
  /DEVICE  TYPE    LABEL          UUID
  1) sda1   ntfs     Recovery       1A90FE4090FE2245
  2) sda2   ufat     A6FF-80A4      A6FF-80A4
  3) sda4   ntfs     7490015F900128E6 7490015F900128E6
  4) sr0    iso9660  2017-07-11-01-39-24-00 2017-07-11-01-39-24-00
# to select, 'r'-refresh, or 'c'-continue: r

(Page 1 of 1) Driver disk device selection
  /DEVICE  TYPE    LABEL          UUID
  1) sda1   ntfs     Recovery       1A90FE4090FE2245
  2) sda2   ufat     A6FF-80A4      A6FF-80A4
  3) sda4   ntfs     7490015F900128E6 7490015F900128E6
  4) sr0    iso9660  CDROM          2017-07-11-13-08-37-00 2017-07-11-13-08-37-00
# to select, 'r'-refresh, or 'c'-continue: 4
DD: Examining /dev/sr0
mount: /dev/sr0 is write-protected, mounting read-only

(Page 1 of 1) Select drivers to install
  1) [ ] /media/DD-1/rpms/x86_64/kmod-qlgc-fastlinq-8.22.0.0-1.rhel17u4.x86_64.rpm
# to toggle selection, or 'c'-continue: 1

(Page 1 of 1) Select drivers to install
  1) [x] /media/DD-1/rpms/x86_64/kmod-qlgc-fastlinq-8.22.0.0-1.rhel17u4.x86_64.rpm
# to toggle selection, or 'c'-continue: c
DD: Extracting: kmod-qlgc-fastlinq

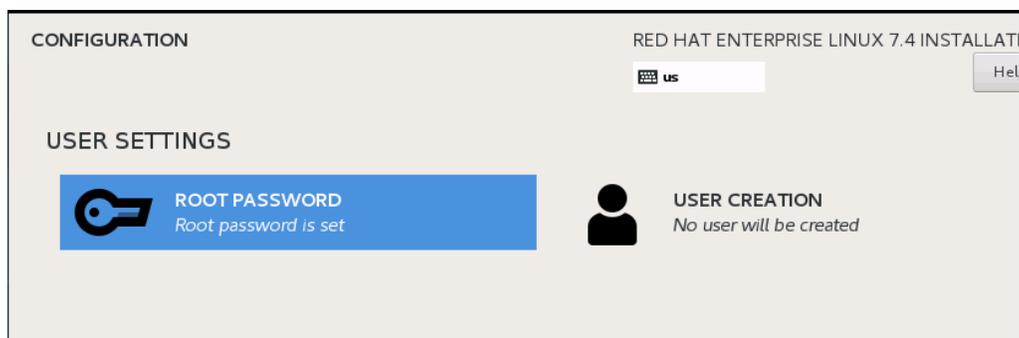
(Page 1 of 1) Driver disk device selection
  /DEVICE  TYPE    LABEL          UUID
  1) sda1   ntfs     Recovery       1A90FE4090FE2245
  2) sda2   ufat     A6FF-80A4      A6FF-80A4
  3) sda4   ntfs     7490015F900128E6 7490015F900128E6
  4) sr0    iso9660  CDROM          2017-07-11-13-08-37-00 2017-07-11-13-08-37-00
# to select, 'r'-refresh, or 'c'-continue:

```

图 10-7. 提示开箱即用安装

6. 如果您的设置需要，请在提示输入其他驱动程序磁盘时加载 FastLinQ 驱动程序更新磁盘。否则，如果您没有其他驱动程序更新磁盘要安装，请键入 c。
7. 继续安装。您可以跳过介质测试。单击 **Next**（下一步）继续安装。

- 在配置窗口 (图 10-8) 中, 选择在安装过程中使用的语言, 然后单击 **Continue** (继续)。



**图 10-8. Red Hat Enterprise Linux 7.4 配置**

- 在 Installation Summary (安装摘要) 窗口, 单击 **Installation Destination** (安装目标)。磁盘标签是 `sda`, 指示单路径安装。如果您配置了多路径, 则该磁盘具有设备映射器标签。
- 在 **Specialized & Network Disks** (专业化和网络磁盘) 章节, 选择 FCoE LUN。
- 键入 root 用户的密码, 然后单击 **Next** (下一步) 完成安装。
- 在第一次引导期间, 添加以下内核命令行以落入 shell。

```
rd.driver.pre=qed,qede,qedr,qedf,qedi
```
- 成功引导系统后, 编辑 `/etc/modprobe.d/anaconda-blacklist.conf` 文件, 以删除所选驱动程序的黑名单条目。
- 通过发出 `dracut -f` 命令重建虚拟磁盘, 然后重新启动。

### 注

关闭用于软件 FCoE 的 `lldpad` 和 `fcoe` 服务。如果这些服务处于活动状态, 则可能会干扰卸载 FCoE 的硬件的正常操作。

# 11 SR-IOV 配置

单根输入 / 输出虚拟化 (SR-IOV) 是一种 PCI SIG 规格，使单个 PCI Express (PCIe) 设备能够显示为多个单独的物理 PCIe 设备。SR-IOV 允许隔离 PCIe 资源以实现性能、互操作性和可管理性。

## 注

某些 SR-IOV 功能在当前版本中可能并未完全启用。

本章提供以下内容的说明：

- [在 Windows 上配置 SR-IOV](#)
- [第 182 页上的“在 Linux 上配置 SR-IOV”](#)
- [第 188 页上的“在 VMware 上配置 SR-IOV”](#)

## 在 Windows 上配置 SR-IOV

要在 Windows 上配置 SR-IOV：

1. 访问服务器 BIOS System Setup（BIOS 系统设置），然后单击 **System BIOS Settings**（系统 BIOS 设置）。
2. 在 System BIOS Settings（系统 BIOS 设置）页面中，单击 **Integrated Devices**（集成的设备）。
3. 在 Integrated Devices（集成的设备）页面（[图 11-1](#)）中：
  - a. 将 **SR-IOV Global Enable**（SR-IOV 全局启用）选项设置为 **Enabled**（启用）。
  - b. 单击 **Back**（后退）。

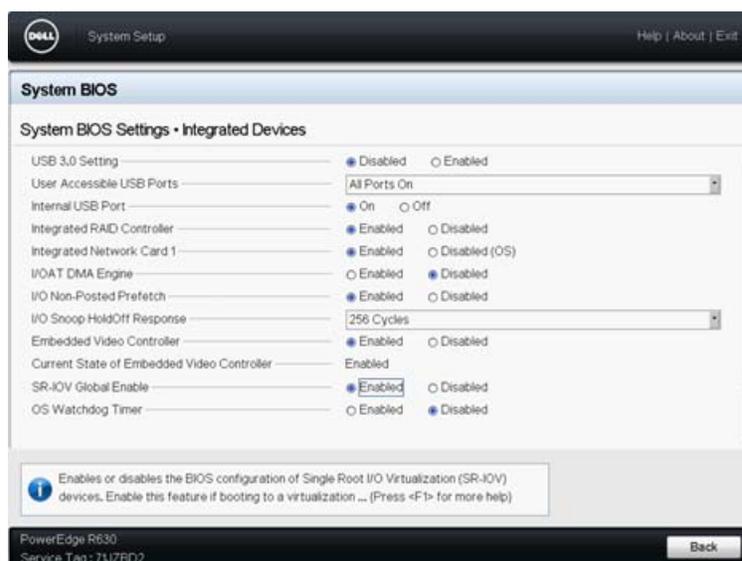


图 11-1. SR-IOV 的系统设置：集成的设备

4. 在所选适配器的 Main Configuration Page（主要配置页面）上，单击 **Device Level Configuration**（设备级配置）。
5. 在 Main Configuration Page - Device Level Configuration（主要配置页面 - 设备级配置）（图 11-2）中：
  - a. 如果您使用的是 NPAR 模式，请将 **Virtualization Mode**（虚拟化模式）设置为 **SR-IOV** 或 **NPAR+SR-IOV**。
  - b. 单击 **Back**（后退）。



图 11-2. SR-IOV 的系统设置：设备级配置

6. 在 Main Configuration Page（主要配置页面）中，单击 **Finish**（完成）。
7. 在 Warning - Saving Changes（警告 - 保存更改）消息框中，单击 **Yes**（是）保存配置。

8. 在 Success - Saving Changes（成功 - 保存更改）消息框中，单击 **OK**（确定）。
9. 要在微型端口适配器上启用 SR-IOV：
  - a. 访问设备管理器。
  - b. 打开微型端口适配器属性，然后单击 **Advanced**（高级）选项卡。
  - c. 在 Advanced（高级）属性页面（图 11-3）的 **Property**（属性）下，选择 **SR-IOV**，然后将值设置为 **Enabled**（已启用）。
  - d. 单击 **OK**（确定）。

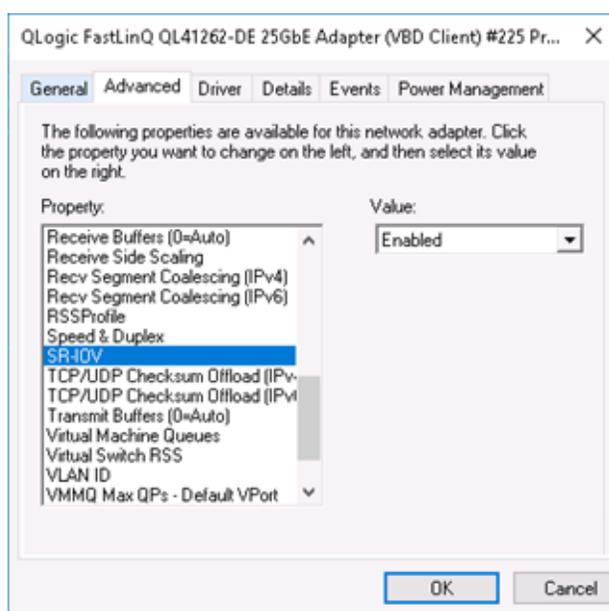


图 11-3. 适配器属性，高级：启用 SR-IOV

10. 要使用 SR-IOV 创建虚拟机交换机（第 178 页上的图 11-4）：
  - a. 启动 Hyper-V 管理器。
  - b. 选择 **Virtual Switch Manager**（虚拟交换机管理器）。
  - c. 在 **Name**（名称）框中，键入虚拟交换机的名称。
  - d. 在 **Connection type**（连接类型）下，选择 **External network**（外部网络）。
  - e. 选择 **Enable single-root I/O virtualization (SR-IOV)**（启用单根 I/O 虚拟化 (SR-IOV)）复选框，然后单击 **Apply**（应用）。

**注**

确保在创建 vSwitch 时启用 SR-IOV。创建 vSwitch 后，此选项将不可用。

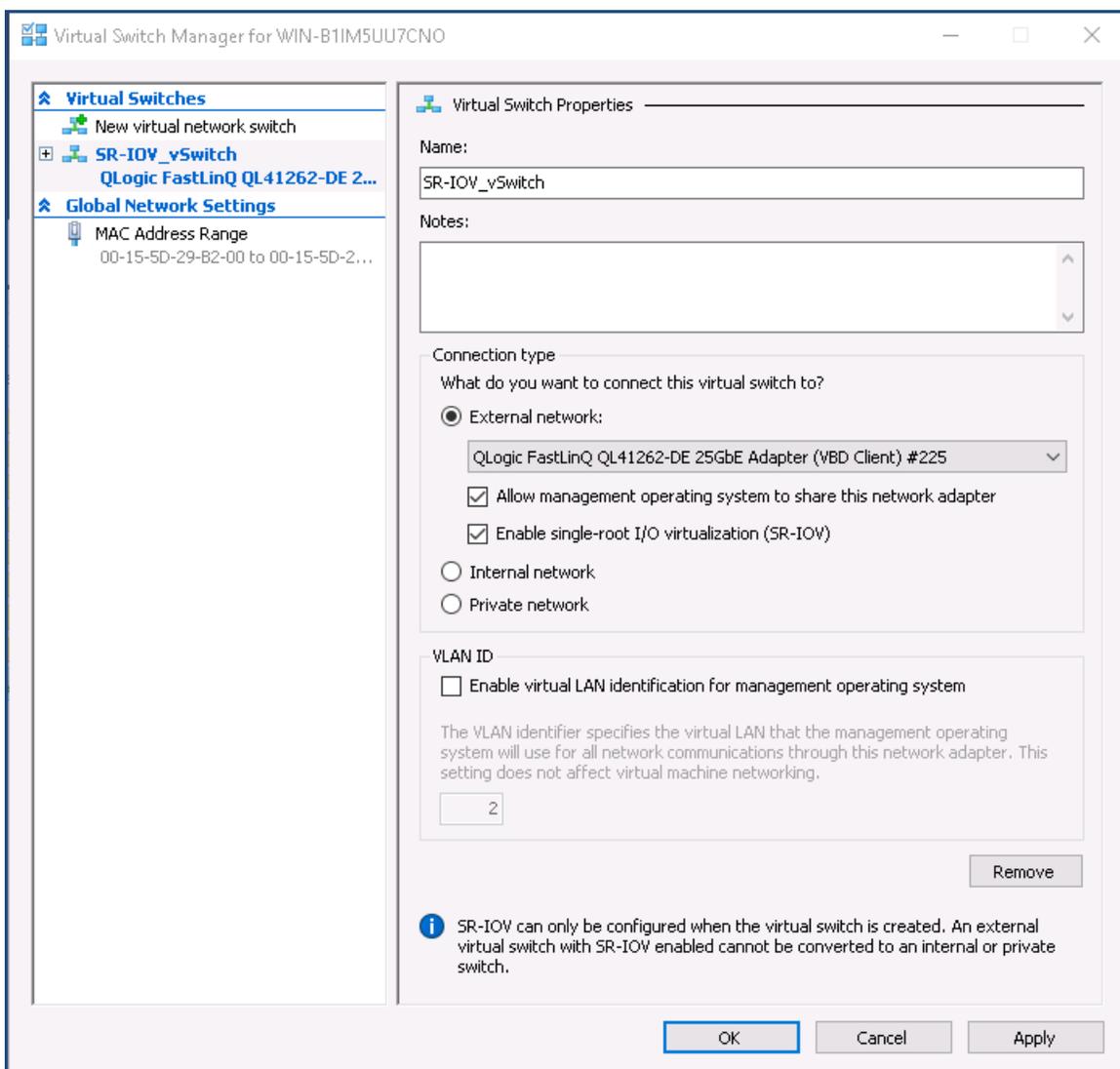


图 11-4. 虚拟交换机管理器：启用 SR-IOV

- f. Apply Networking Changes（应用网络更改）消息框告知您待处理更改可能会中断网络连接。要保存您的更改并继续，请单击 **Yes**（是）。

11. 要获取虚拟机交换机功能，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Get-VMSwitch -Name SR-IOV_vSwitch | fl
```

Get-VMSwitch 命令的输出包括以下 SR-IOV 功能：

```
IovVirtualFunctionCount           : 96  
IovVirtualFunctionsInUse          : 1
```

12. 要创建虚拟机 (VM) 并导出 VM 中的虚拟功能 (VF)：
- 创建一个虚拟机。
  - 将 VMNetworkadapter 添加到虚拟机。
  - 将虚拟交换机分配给 VMNetworkadapter。
  - 在 Settings for VM <VM\_Name> (VM <虚拟机名称> 的设置) 对话框 (图 11-5) 中，Hardware Acceleration (硬件加速) 页面的 **Single-root I/O virtualization** (单根 I/O 虚拟化) 下，选择 **Enable SR-IOV** (启用 SR-IOV) 复选框，然后单击 **OK** (确定)。

---

### 注

创建虚拟适配器连接后，可以随时启用或禁用 SR-IOV 设置 (即使流量正在运行)。

---

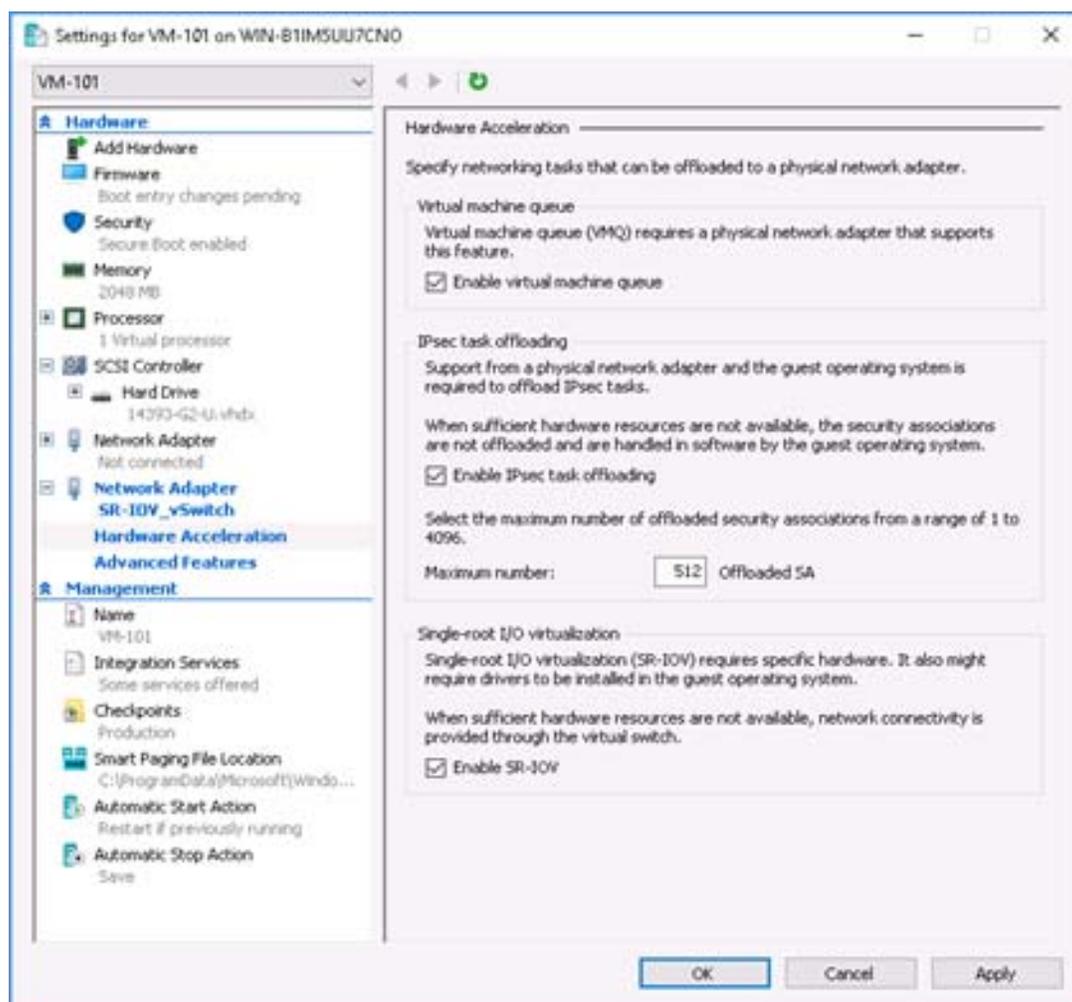


图 11-5. VM 的设置：启用 SR-IOV

13. 为虚拟机中检测到的适配器安装 QLogic 驱动程序。为您的主机操作系统使用可从供应商处获得的最新驱动程序（请勿使用自带的驱动程序）。

### 注

确保在 VM 和主机系统上使用相同的驱动程序包。例如，在 Windows VM 和 Windows Hyper-V 主机上使用相同的 qeVBD 和 qeND 驱动程序版本。

安装驱动程序后，QLogic 适配器在 VM 中列出。图 11-6 显示示例。

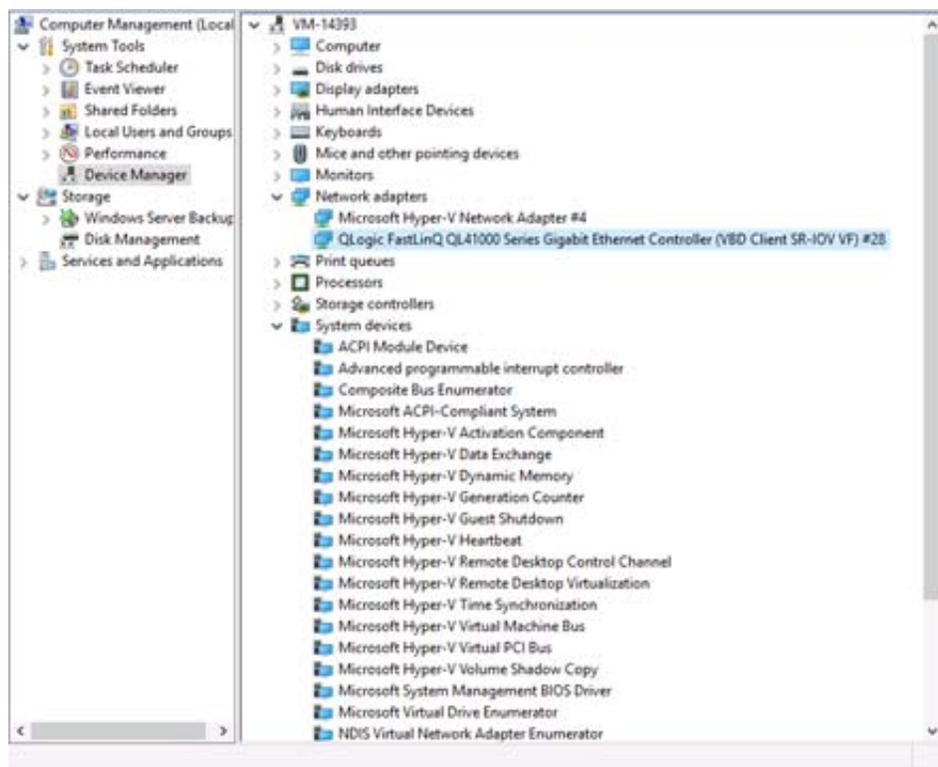


图 11-6. 设备管理器：带有 QLogic 适配器的 VM

14. 要查看 SR-IOV VF 详细信息，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Get-NetadapterSriovVF
```

图 11-7 显示示例输出。

```
PS C:\Users\Administrator> Get-NetAdapterSriovVF
Name                FunctionID VPortID  MacAddress          VmID                VmFriendlyName
-----
Ethernet 10         0         {2}          00-15-5D-29-B2-01  51F01C52-CDC6-4932-A95E-86D... VM-101
PS C:\Users\Administrator>
```

图 11-7. Windows PowerShell 命令：Get-NetadapterSriovVf

## 在 Linux 上配置 SR-IOV

要在 Linux 上配置 SR-IOV：

1. 访问服务器 BIOS System Setup（BIOS 系统设置），然后单击 **System BIOS Settings**（系统 BIOS 设置）。
2. 在 System BIOS Settings（系统 BIOS 设置）页面中，单击 **Integrated Devices**（集成的设备）。
3. 在 System Integrated Devices（系统集成设备）页面（请参阅第 176 页上的图 11-1）中：
  - a. 将 **SR-IOV Global Enable**（SR-IOV 全局启用）选项设置为 **Enabled**（启用）。
  - b. 单击 **Back**（后退）。
4. 在 System BIOS Settings（系统 BIOS 设置）页面中，单击 **Processor Settings**（处理器设置）。
5. 在 Processor Settings（处理器设置）（图 11-8）页面中：
  - a. 将 **Virtualization Technology**（虚拟化技术）选项设置为 **Enabled**（启用）。
  - b. 单击 **Back**（后退）。

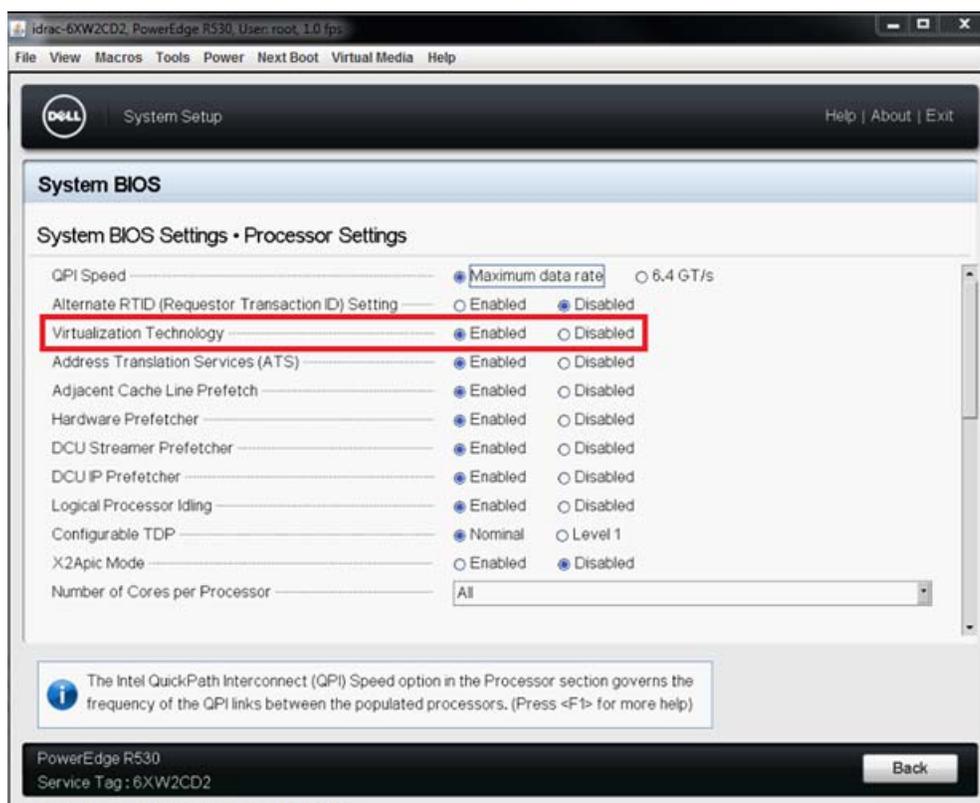


图 11-8. 系统设置：SR-IOV 的处理器设置

6. 在 System Setup（系统设置）页面中，选择 Device Settings（设备设置）。
7. 在 Device Settings（设备设置）页面中，为 QLogic 适配器选择 Port 1（端口 1）。

8. 在 Device Level Configuration（设备级配置）页面（图 11-9）中：
  - a. 将 **Virtualization Mode**（虚拟化模式）设置为 **SR-IOV**。
  - b. 单击 **Back**（后退）。



**图 11-9. SR-IOV 的系统设置：集成的设备**

9. 在 Main Configuration Page（主要配置页面）中，单击 **Finish**（完成），保存设置，然后重新引导系统。
10. 要启用并验证虚拟化：
  - a. 打开 `grub.conf` 文件并配置 `iommu` 参数，如图 11-10 中所示。
    - 对于基于 Intel 的系统，请添加 `intel_iommu=on`。
    - 对于基于 AMD 的系统，请添加 `amd_iommu=on`。

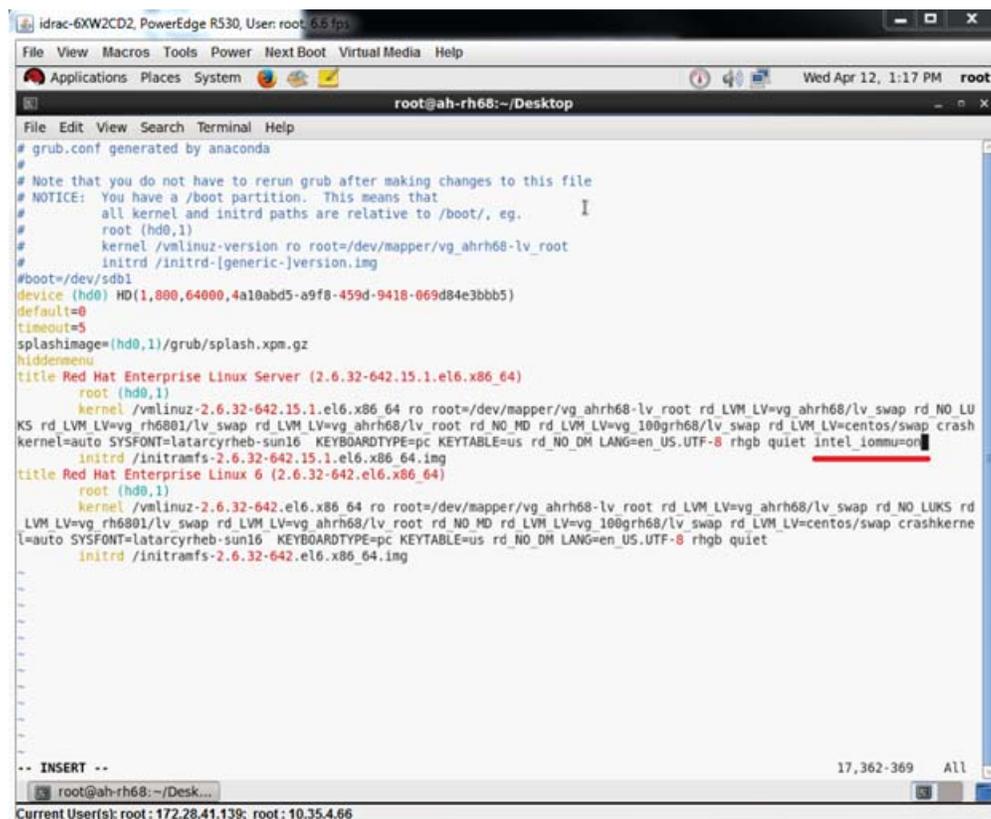


图 11-10. 为 SR-IOV 编辑 grub.conf 文件

b. 保存 grub.conf 文件，然后重新引导系统。

c. 要验证更改是否已生效，请发出以下命令：

```
dmesg | grep -I iommu
```

应显示成功的输入 - 输出内存管理单元 (IOMMU) 命令，例如：

```
Intel-IOMMU: enabled
```

d. 要查看 VF 详细信息 (VF 数和 VF 总计)，请发出以下命令：

```
find /sys/|grep -I sriov
```

11. 对于特定端口，启用一定数量的 VF。

a. 例如，发出以下命令以启用 PCI 实例 04:00.0 (总线 4，设备 0，功能 0) 上的 8 个 VF：

```
[root@ah-rh68 ~]# echo 8 >  
/sys/devices/pci0000:00/0000:00:02.0/0000:04:00.0/sriov_numvfs
```

- b. 查看命令输出（图 11-11）以确认在总线 4、设备 2（通过 0000:00:02.0 参数）、功能 0 至 7 上创建了实际 VF。请注意，实际的设备 ID 在 Pf（此示例中为 8070）与 Vf（此示例中为 9090）上不同。

```
[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: QLogic Corp. Device 000b  
04:02.1 Ethernet controller: QLogic Corp. Device 8090 (rev 02)  
    Subsystem: QLogic 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]#
```

图 11-11. sriov\_numvfs 的命令输出

12. 要查看所有 PF 和 VF 接口的列表，请发出以下命令：

```
# ip link show | grep -i vf -b2
```

图 11-12 显示示例输出。

```
[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:ff  
326: vf 0 MAC 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 00:00:00:00:00:00, tx rate 10000 (Mbps), max_tx_rate 10000Mbps, spoof checking off, link-state auto  
665: vf 3 MAC 00:00:00:00:00:00, 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
```

图 11-12. ip link show 命令的命令输出

13. 分配并验证 MAC 地址：

- a. 要将 MAC 地址分配给 VF，请发出以下命令：

```
ip link set <pf device> vf <vf index> mac <mac address>
```

- b. 确保 VF 接口正常工作并以分配的 MAC 地址运行。
14. 关闭 VM 的电源并连接 VF。（某些操作系统支持热插拔 VF 到 VM。）
- a. 在 Virtual Machine（虚拟机）对话框（图 11-13）中，单击 **Add Hardware**（添加硬件）。

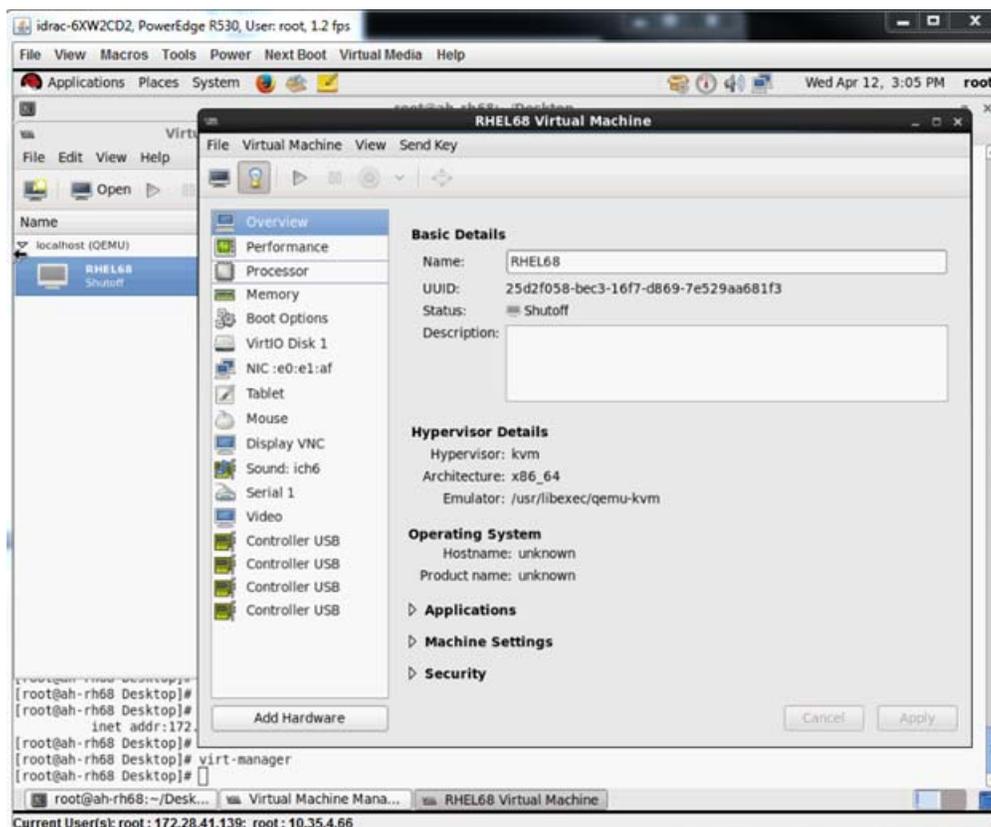


图 11-13. RHEL68 虚拟机

- b. 在 Add New Virtual Hardware（添加新虚拟硬件）对话框（图 11-14）的左侧窗格中，单击 **PCI Host Device**（PCI 主机设备）。
- c. 在右侧窗格中，选择一个主机设备。
- d. 单击 **Finish**（完成）。



图 11-14. 添加新虚拟硬件

15. 打开 VM 的电源，然后发出以下命令：  

```
check lspci -vv|grep -I ether
```
16. 为虚拟机中检测到的适配器安装驱动程序。为您的主机操作系统使用可从供应商处获得的最新驱动程序（请勿使用自带的驱动程序）。主机和虚拟机上必须安装相同的驱动程序版本。
17. 根据需要，在 VM 中添加更多 VF。

## 在 VMware 上配置 SR-IOV

要在 VMware 上配置 SR-IOV：

1. 访问服务器 BIOS System Setup（BIOS 系统设置），然后单击 **System BIOS Settings**（系统 BIOS 设置）。
2. 在 System BIOS Settings（系统 BIOS 设置）页面中，单击 **Integrated Devices**（集成的设备）。
3. 在 Integrated Devices（集成的设备）页面（请参阅第 176 页上的图 11-1）中：
  - a. 将 **SR-IOV Global Enable**（SR-IOV 全局启用）选项设置为 **Enabled**（启用）。
  - b. 单击 **Back**（后退）。
4. 在 System Setup（系统设置）窗口中，单击 **Device Settings**（设备设置）。
5. 在 Device Settings（设备设置）页面中，选择用于 25G 41xxx 系列适配器的端口。

6. 在 Device Level Configuration（设备级配置）（请参阅第 176 页上的图 11-2）中：
  - a. 将 **Virtualization Mode**（虚拟化模式）设置为 **SR-IOV**。
  - b. 单击 **Back**（后退）。
7. 在 Main Configuration Page（主要配置页面）中，单击 **Finish**（完成）。
8. 保存配置设置并重新引导系统。
9. 要启用每个端口所需数量的 VF（在此示例中，双端口适配器的每个端口上需要 16 个 VF），请发出以下命令：

```
"esxcfg-module -s "max_vfs=16,16" qedentv"
```

---

**注**

41xxx 系列适配器的每个以太网功能必须具备各自的条目。

---

10. 重新引导主机。
11. 要验证模块级别的更改是否已完成，请发出以下命令：

```
"esxcfg-module -g qedentv"
```

```
[root@localhost:~] esxcfg-module -g qedentv  
qedentv enabled = 1 options = 'max_vfs=16,16'
```

12. 要验证是否创建了实际 VF，请发出 `lspci` 命令如下：

```
[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 QL41xQL41xxxxx  
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]
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. 请如下将 VF 连接到 VM:
  - a. 关闭 VM 的电源并连接 VF。（某些操作系统支持热插拔 VF 到 VM。）
  - b. 将主机添加到 VMware vCenter Server 虚拟设备 (vCSA)。
  - c. 单击 VM 的 **Edit Settings**（编辑设置）。
14. 请如下填写 Edit Settings（编辑设置）对话框（[图 11-15](#)）：
  - a. 在 **New Device**（新设备）框中，选择 **Network**（网络），然后单击 **Add**（添加）。
  - b. 对于 **Adapter Type**（适配器类型），选择 **SR-IOV Passthrough**（SR-IOV 直通）。
  - c. 对于 **Physical Function**（物理功能），选择 QLogic VF。
  - d. 要保存您的配置更改并关闭此对话框，请单击 **OK**（确定）。

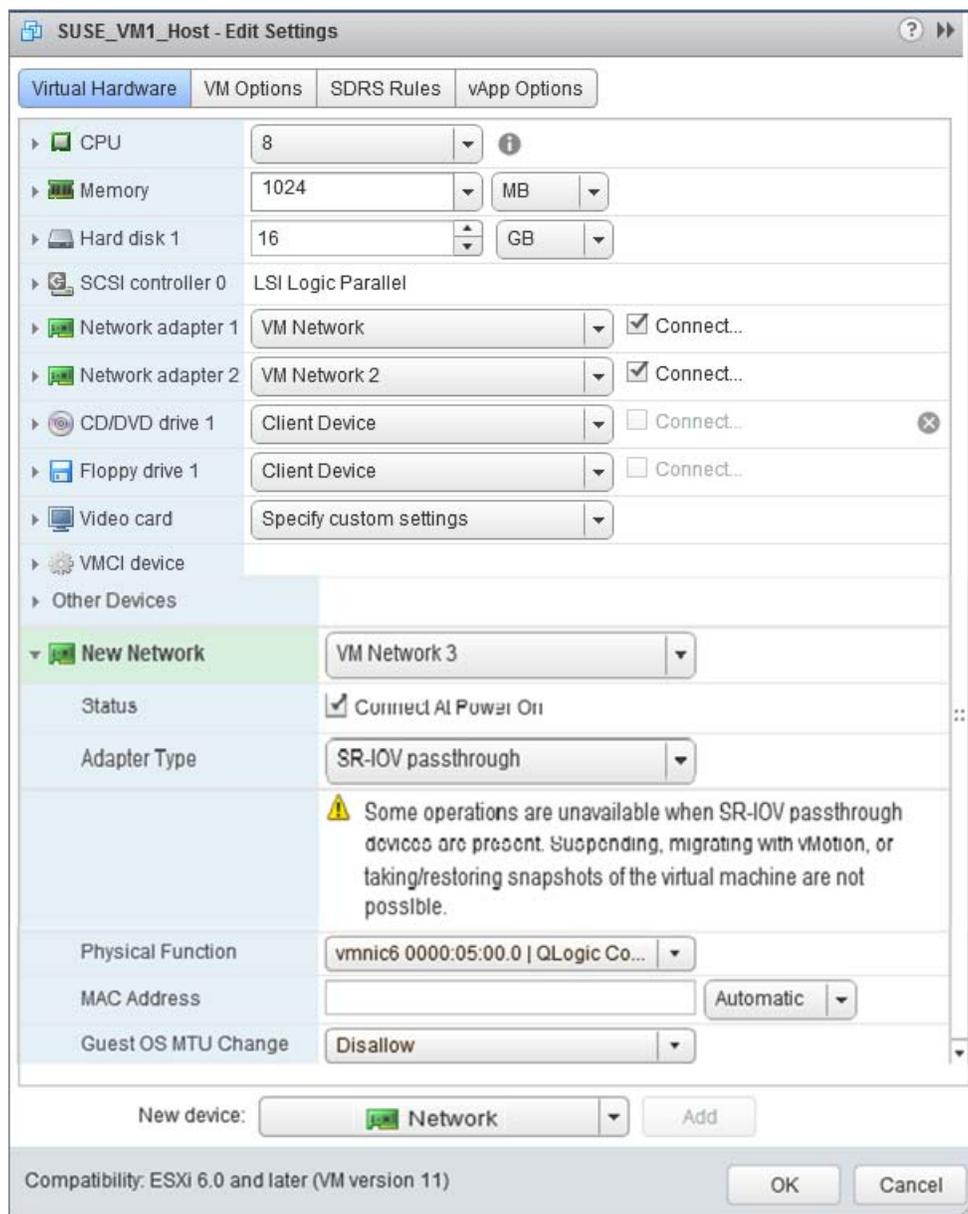


图 11-15. VMware 主机编辑设置

15. 要验证每个端口的 VF，请如下发出 `esxcli` 命令：

```
[root@localhost:~] esxcli network sriovnic vf list -n vmnic6
```

VF ID	Active	PCI Address	Owner World ID
0	true	005:02.0	60591
1	true	005:02.1	60591

2	false	005:02.2	-
3	false	005:02.3	-
4	false	005:02.4	-
5	false	005:02.5	-
6	false	005:02.6	-
7	false	005:02.7	-
8	false	005:03.0	-
9	false	005:03.1	-
10	false	005:03.2	-
11	false	005:03.3	-
12	false	005:03.4	-
13	false	005:03.5	-
14	false	005:03.6	-
15	false	005:03.7	-

16. 为 VM 中检测到的适配器安装 QLogic 驱动程序。为您的主机操作系统使用可从供应商处获得的最新驱动程序（请勿使用自带的驱动程序）。主机和虚拟机上必须安装相同的驱动程序版本。
17. 打开 VM 的电源，然后发出 `ifconfig -a` 命令以验证添加的网络接口是否已列出。
18. 根据需要，在 VM 中添加更多 VF。

# 12 使用 RDMA 进行 NVMe-oF 配置

基于结构的非易失性存储器标准 (NVMe-oF) 可实现向 PCIe 使用交替传输，从而延长 NVMe 主机设备和 NVMe 存储驱动器或子系统可连接的距离。NVMe-oF 定义了一种通用架构，其通过存储网络结构来支持 NVMe 块存储协议的一系列存储网络结构。此架构包括启用存储系统的前端接口、扩展到大量 NVMe 设备以及延长数据中心内可访问 NVMe 设备和 NVMe 子系统的距离。

本章介绍的 NVMe-oF 配置程序和选项适用于基于以太网的 RDMA 协议，包括 RoCE 和 iWARP。包含 RDMA 的 NVMe-oF 的开发由 NVMe 组织的技术小组定义。

本章演示如何在简单网络上配置 NVMe-oF。示例网络包括以下内容：

- 两台服务器：启动器和目标服务器。目标服务器配备了 PCIe SSD 驱动器。
- 操作系统：两台服务器上的 RHEL 7.4 或 SLES 12 SP3
- 两个适配器：每台服务器都安装一个 41xxx 系列适配器。每个端口都可以独立配置为使用 RoCE、RoCEv2 或 iWARP 作为运行 NVMe-oF 的 RDMA 协议。
- 对于 RoCE 和 RoCEv2，配置用于数据中心桥接 (DCB)、相关服务质量 (QoS) 策略和 vLAN 以承载 NVMe-oF 的 RoCE/RoCEv2 DCB 流量类别优先级的可选交换机。当 NVMe-oF 在使用 iWARP 时，不需要交换机。

图 12-1 展示示例网络。

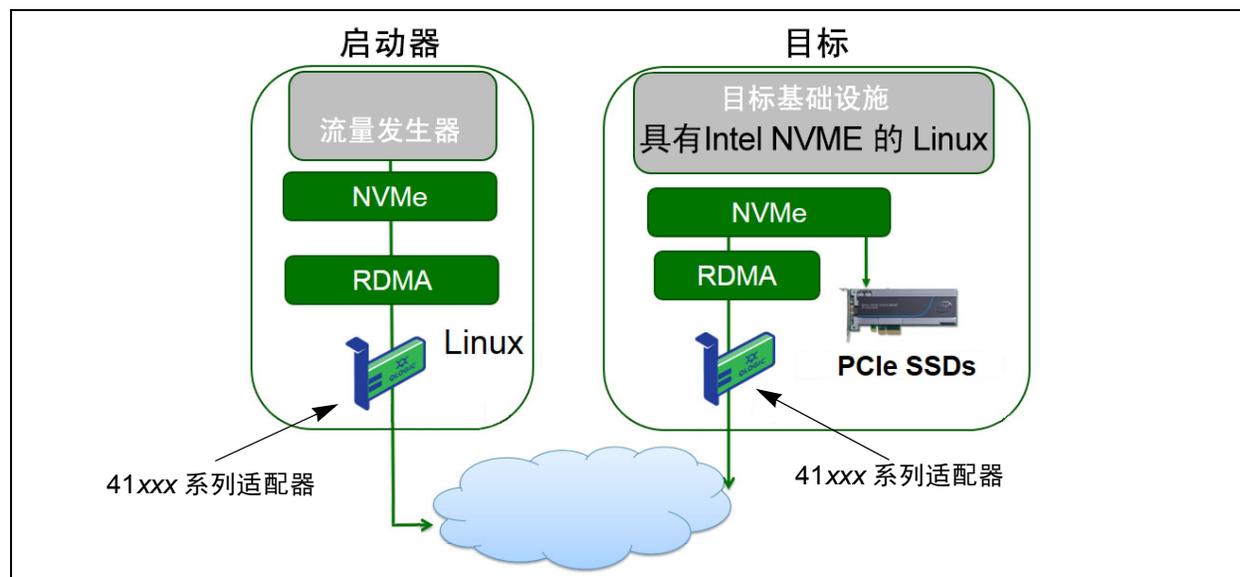


图 12-1. NVMe-oF 网络

NVMe-oF 配置过程包括以下步骤：

1. 在两台服务器上安装设备驱动程序
2. 配置目标服务器
3. 配置启动器服务器
4. 预处理目标服务器
5. 测试 NVMe-oF 设备
6. 优化性能

## 在两台服务器上安装设备驱动程序

安装操作系统（RHEL 7.4 或 SLES 12 SP3）后，请在两台服务器上安装设备驱动程序。要将内核升级到最新的 Linux 上游内核，请转至：

<https://www.kernel.org/pub/linux/kernel/v4.x/>

1. 按照 README 中的所有安装说明安装并加载最新的 FastLinQ 驱动程序（QED、QEDE、libqedr/QEDR）。
2. （可选）如果您升级了操作系统内核，则必须重新安装并加载最新的驱动程序，如下所示：
  - a. 按照 README 中的所有安装说明安装最新的 FastLinQ 固件
  - b. 通过发出以下命令来安装 OS RDMA 支持应用程序和库：

```
# 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. 要确保 NVMe OFED 支持位于所选的操作系统内核中，请发出以下命令：

```
make menuconfig
```
  - d. 在 **Device Drivers**（设备驱动程序）下，确保已启用以下设置（设置为 **M**）：

```
NVM Express block devices
NVM Express over Fabrics RDMA host driver
NVMe Target support
NVMe over Fabrics RDMA target support
```
  - e. （可选）如果 **Device Drivers**（设备驱动程序）选项尚不存在，请通过发出以下命令来重建内核：

```
# make
# make modules
# make modules_install
# make install
```
  - f. 如果对内核进行了更改，请重新引导到该新的操作系统内核。有关如何设置默认引导内核的说明，请转至：  
<https://wiki.centos.org/HowTos/Grub2>

3. 启用并启动 RDMA 服务，如下所示：

```
# systemctl enable rdma.service  
# systemctl start rdma.service
```

忽略 RDMA Service Failed (RDMA 服务失败) 错误。QEDR 所需的所有 OFED 模块已经加载。

## 配置目标服务器

重新启动过程之后配置目标服务器。服务器运行后，无需重新启动即可更改配置。如果您在使用启动脚本来配置目标服务器，请考虑根据需要暂停脚本（使用 wait 命令或类似的命令），以确保每个命令在执行下一个命令之前完成。

### 配置目标服务：

1. 加载目标模块。每次服务器重新启动后发出以下命令：

```
# modprobe qedr  
# modprobe nvmet; modprobe nvmet-rdma  
# lsmod | grep nvme (确认已加载模块)
```

2. 使用由 <nvme-subsystem-name> 指示的名称创建目标子系统 NVMe 限定名称 (NQN)。使用 NVMe-oF 规范；例如

```
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. 根据需要为其他 NVMe 设备创建多个唯一的 NQN。

4. 设置目标参数，如 [表 12-1](#) 中所列。

表 12-1. 目标参数

命令	说明
# echo 1 > attr_allow_any_host	允许任何主机连接。
# mkdir namespaces/1	创建名称空间

表 12-1. 目标参数 (续)

命令	说明
# echo -n /dev/nvme0n1 > namespaces/ 1/device_path	设置 NVMe 设备路径。系统之间的 NVMe 设备路径可能不同。使用 lsblk 命令检查设备路径。此系统有两个 NVMe 设备：nvme0n1 和 nvme1n1。  <pre>[root@localhost home]# lsblk NAME        MAJ:MIN RM  SIZE RO TYPE MOUNTPOINT nvme1n1    259:0   0 372.6G  0 disk sda         8:0     0  1.1T  0 disk ├─sda2      8:2     0   505G  0 part / ├─sda3      8:3     0     8G  0 part [SWAP] └─sda1      8:1     0     1G  0 part /boot/efi nvme0n1    259:1   0 372.6G  0 disk</pre>
# echo 1 > namespaces/1/enable	启用名称空间。
# mkdir /sys/kernel/config/nvmet/ ports/1 # cd /sys/kernel/config/nvmet/ports/1	创建 NVMe 端口 1。
# echo 1.1.1.1 > addr_traddr	设置相同的 IP 地址。例如，1.1.1.1 是 41xxx Series Adapter 的目标端口的 IP 地址。
# echo rdma > addr_trtype	设置传输类型 RDMA。
# echo 4420 > addr_trsvcid	设置 RDMA 端口号。NVMe-oF 的套接字端口号通常为 4420。但是，如果在整个配置中始终使用端口号，则可以使用任何端口号。
# echo ipv4 > addr_adrfam	设置 IP 地址类型。

5. 创建一个符号链接 (symlink) 到新创建的 NQN 子系统:

```
# ln -s /sys/kernel/config/nvmet/subsystems/  
nvme-subsystem-name subsystems/nvme-subsystem-name
```

6. 确认 NVMe 目标正在端口上监听，如下所示:

```
# dmesg | grep nvmet_rdma
[ 8769.470043] nvmet_rdma: enabling port 1 (1.1.1.1:4420)
```

## 配置启动器服务器

在重新启动过程之后配置启动器服务器。服务器运行后，不重新启动无法更改配置。如果您使用启动脚本来配置启动器服务器，请考虑根据需要暂停脚本（使用 wait 命令或类似的命令），以确保每一条命令在执行下一条命令之前完成。

### 要配置启动器服务器：

1. 加载 NVMe 模块。每次服务器重新启动后发出这些命令：

```
# modprobe qedr
# modprobe nvme-rdma
```
2. 下载、编译并安装 `nvme-cli` 启动器实用程序。在第一次配置时发出这些命令 — 每次重新启动后不需要发出这些命令。

```
# git clone https://github.com/linux-nvme/nvme-cli.git
# cd nvme-cli
# make && make install
```
3. 按如下方式验证安装版本：

```
# nvme version
```
4. 按如下方式发现 NVMe-oF 目标：

```
# nvme discover -t rdma -a 1.1.1.1 -s 1023
```

记下已发现目标 (图 12-2) 的子系统 NQN (`subnqn`) 用于步骤 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_qtype: connected
rdma_cms: rdma-cm
rdma_pkey: 0x0000
```

图 12-2. 子系统 NQN

5. 使用 NQN 连接到已发现的 NVMe-oF 目标 (`nvme-qlogic-tgt1`)。每次服务器重新启动后发出以下命令。例如：

```
# nvme connect -t rdma -n nvme-qlogic-tgt1 -a 1.1.1.1 -s 1023
```
6. 按如下方式确认 NVMe-oF 目标与 NVMe-oF 设备的连接：

```
# dmesg | grep nvme
# lsblk
# list nvme
```

图 12-3 显示示例。

```
[root@localhost home] #dmesg | grep nvme
[ 233.645554] nvme nvme0: new ctrl: NQN "nvme-qlogic-tgt1", addr 1.1.1.1:1023
[root@localhost home] # lsblk
NAME        MAJ:MIN RM   SIZE RO TYPE MOUNTPOINT
sdb         8:0     0    1.1T 0 disk
├─sdb2      8:2     0   493.2G 0 part /
├─sdb3      8:3     0     8G 0 part [SWAP]
└─sdb1      8:1     0     1G 0 part /boot/efi
nvme0n1    259:0     0   372.6G 0 disk
[root@localhost home] # nvme list
Node          SN                      Model      Namespace  Usage                Format          FW Rev
-----
/dev/nvme0n1  7a591f3ec788a367      Linux     1           1.60 TB / 1.60 TB   512 B + 0 B   4.13.8
```

图 12-3. 确认 NVMe-oF 连接

## 预处理目标服务器

经测试的、开箱即用的 NVMe 目标服务器显示出高于预期的性能。在运行基准测试之前，需要预装填或预处理目标服务器。

### 要预处理目标服务器：

1. 使用供应商特定的工具安全擦除目标服务器（类似于格式化）。此测试演示了使用 Intel NVMe SSD 设备，该设备需要 Intel Data Center Tool(Intel 数据中心工具)，可在以下链接获得：

<https://downloadcenter.intel.com/download/23931/Intel-Solid-State-Drive-Data-Center-Tool>

2. 使用数据来预处理目标服务器 (nvme0n1)，这样可以保证填满所有可用的存储器。此示例使用“DD”磁盘公用程序：

```
# dd if=/dev/zero bs=1024k of=/dev/nvme0n1
```

## 测试 NVMe-oF 设备

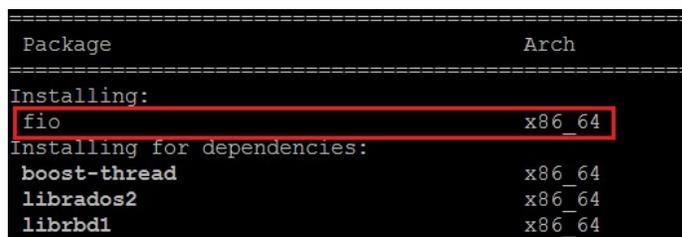
将目标服务器上本地 NVMe 设备的延迟与启动器服务器上 NVMe-oF 设备的延迟进行比较，以显示 NVMe 向系统添加的延迟。

### 要测试 NVMe-oF 设备：

1. 更新存储库 (Repo) 源并通过发出以下命令在目标服务器和启动器服务器上安装灵活输入 / 输出 (FIO) 基准测试实用程序：

```
# yum install epel-release
```

```
# yum install fio
```



```
=====  
Package                               Arch  
=====  
Installing:  
fio                                   x86_64  
Installing for dependencies:  
boost-thread                         x86_64  
librados2                             x86_64  
librbd1                               x86_64  
=====
```

图 12-4. FIO 公用程序安装

2. 运行 FIO 实用程序来测量启动器 NVMe-oF 设备的延迟。发出以下命令：

```
# 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 报告两种延迟类型：提交和完成。提交延迟 (slat) 测量应用程序到内核延迟。完成延迟 (clat) 测量端到端内核延迟。行业接受的方法是读取第 99.00 个范围中的 *clat* 百分位。

在此示例中，启动器设备 NVMe-oF 延迟为 30µsec。

3. 运行 FIO 以测量目标服务器上本地 NVMe 设备的延迟。发出以下命令：

```
# 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
```

在此示例中，目标 NVMe 设备延迟为 8µsec。使用 NVMe-oF 产生的总延迟为启动器设备 NVMe-oF 延迟 (30µsec) 与目标设备 NVMe-oF 延迟 (8µsec) 的差异，即 22 µsec。

4. 运行 FIO 以测量目标服务器上本地 NVMe 设备的带宽。发出以下命令：

```
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
```

其中 --rw 可以是只读的 *randread*、只写的 *randwrite*、可读取可写入的 *randrw*。

## 优化性能

要同时优化启动器和目标服务器上的性能：

1. 配置以下系统 BIOS 设置：
  - Power Profiles = 'Max Performance' or equivalent
  - ALL C-States = disabled
  - Hyperthreading = disabled
2. 通过编辑 grub 文件 (/etc/default/grub) 来配置 Linux 内核参数。
  - a. 将参数添加到 GRUB\_CMDLINE\_LINUX 行末尾：

```
GRUB_CMDLINE_LINUX="nosoftlockup intel_idle.max_cstate=0  
processor.max_cstate=1 mce=ignore_ce idle=poll"
```
  - b. 保存 grub 文件。
  - c. 重建 grub 文件。要重建旧版 BIOS 引导的 grub 文件，请发出以下命令：

```
# grub2-mkconfig -o /boot/grub2/grub.cfg
```

（旧版 BIOS 引导）  
要重建 EFI 引导的 grub 文件，请发出以下命令：

```
# grub2-mkconfig -o /boot/efi/EFI/<os>/grub.cfg
```

（EFI 引导）
  - d. 重新启动服务器以实施更改。
3. 设置所有 41xxx 系列适配器的 IRQ 关联。multi\_rss-affin.sh 文件为 [第 202 页上的“.IRQ 关联 \(multi\\_rss-affin.sh\)”](#) 中列出的脚本文件。

```
# systemctl stop irqbalance  
# ./multi_rss-affin.sh eth1
```

---

### 注

此脚本的不同版本 qedr\_affin.sh 位于  
\add-ons\performance\roce 目录中的 41xxx Linux 源代码包中。  
有关 IRQ 关联设置的说明，请参阅该目录中的  
multiple\_irqs.txt 文件。

---

4. 设置 CPU 频率。cpufreq.sh 文件为 [第 203 页上的“CPU 频率 \(cpufreq.sh\)”](#) 中列出的脚本。

```
# ./cpufreq.sh
```

以下部分列出 [步骤 3](#) 和 [步骤 4](#) 中使用的脚本。

## .IRQ 关联 (multi\_rss-affin.sh)

以下脚本设置 IRQ 关联。

```
#!/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 频率 (cpufreq.sh)

以下脚本设置 CPU 频率。

```
#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
```

**要配置网络或存储器设置：**

```
sysctl -w net.ipv4.tcp_mem="16777216 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_adv_win_scale=1
```

---

### 注

以下命令仅适用于启动器服务器。

---

```
# echo noop > /sys/block/nvme0n1/queue/scheduler
# echo 0 > /sys/block/nvme0n1/queue/add_random
# echo 2 > /sys/block/nvme0n1/queue/nomerges
```

# 13 Windows Server 2016

本章提供有关 Windows Server 2016 的以下信息：

- [使用 Hyper-V 配置 RoCE 接口](#)
- [第 210 页上的“Switch Embedded Teaming 上的 RoCE”](#)
- [第 212 页上的“为 RoCE 配置 QoS”](#)
- [第 220 页上的“配置 VMMQ”](#)
- [第 227 页上的“配置 VXLAN”](#)
- [第 228 页上的“配置 Storage Spaces Direct”](#)
- [第 234 页上的“部署和管理 Nano Server”](#)

## 使用 Hyper-V 配置 RoCE 接口

在 Windows Server 2016 中，具有网络直接内核提供程序接口 (NDKPI) 模式 2 的 Hyper-V、主机虚拟网络适配器（主机虚拟 NIC）均支持 RDMA。

### 注

Hyper-V 上的 RoCE 需要 DCBX。要配置 DCBX：

- [通过 HII 进行配置](#)（请参阅 [第 62 页上的“准备适配器”](#)）。
- [使用 QoS 进行配置](#)（请参阅 [第 212 页上的“为 RoCE 配置 QoS”](#)）。

本节中的 RoCE 配置步骤包括：

- [创建带 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机](#)
- [将 VLAN ID 添加到主机虚拟 NIC](#)
- [验证 RoCE 是否启用](#)
- [添加主机虚拟 NIC（虚拟端口）](#)
- [映射 SMB 驱动器和运行 RoCE 流量](#)

## 创建带 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机

按照本节中的步骤操作以创建 Hyper-V 虚拟交换机，然后在主机 VNIC 中启用 RDMA。

要创建带 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机：

1. 启动 Hyper-V 管理器。
2. 单击 **Virtual Switch Manager**（虚拟交换机管理器）（请参阅图 13-1）。

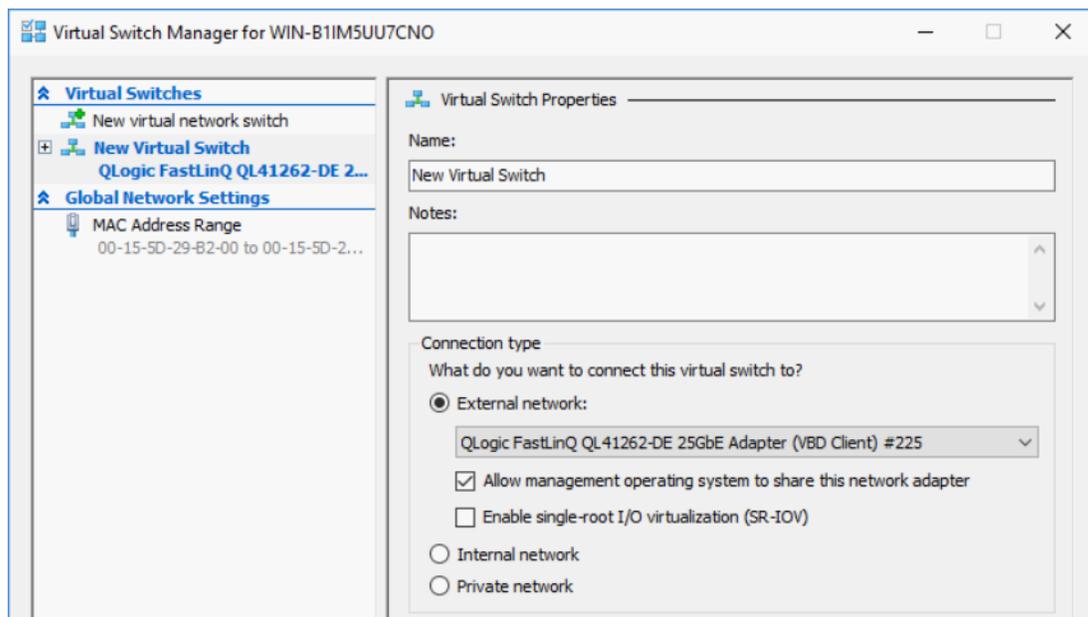


图 13-1. 在主机虚拟 NIC 中启用 RDMA

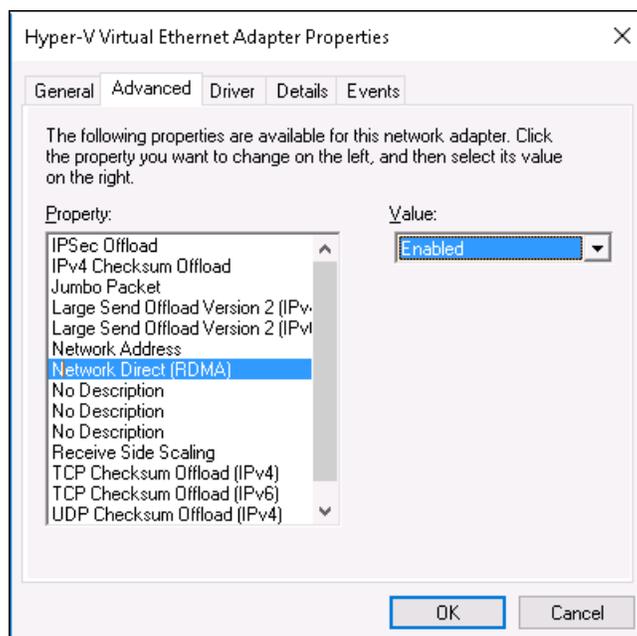
3. 创建虚拟交换机。
4. 选择 **Allow management operating system to share this network adapter**（允许管理操作系统以共享此网络适配器）复选框。

在 Windows Server 2016 中，将在主机虚拟 NIC 中添加一个新参数：Network Direct (RDMA)。

要在主机虚拟 NIC 中启用 RDMA：

1. 打开 Hyper-V Virtual Ethernet Adapter Properties（Hyper-V 虚拟以太网适配器属性）窗口。
2. 单击 **Advanced**（高级）选项卡。

3. 在 Advanced（高级）页面（图 13-2）中：
  - a. 在 **Property**（属性）下，选择 **Network Direct (RDMA)**。
  - b. 在 **Value**（值）下，选择 **Enabled**（启用）。
  - c. 单击 **OK**（确定）。



**图 13-2. Hyper-V 虚拟以太网适配器属性**

4. 要启用 RDMA，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Enable-NetAdapterRdma "vEthernet  
(New Virtual Switch)"  
PS C:\Users\Administrator>
```

## 将 VLAN ID 添加到主机虚拟 NIC

要将 VLAN ID 添加到主机虚拟 NIC：

1. 要查找主机虚拟 NIC 名称，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Get-VMNetworkAdapter -ManagementOS
```

图 13-3 显示命令输出。

```
PS C:\Users\Administrator> Get-VMNetworkAdapter -ManagementOS
Name                IsManagementOs VMName SwitchName          MacAddress          Status IPAddresses
-----
New Virtual Switch True                New Virtual Switch 000E1EC41F0B {Ok}
```

图 13-3. Windows PowerShell 命令：Get-VMNetworkAdapter

2. 要设置主机虚拟 NIC 的 VLAN ID，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Set-VMNetworkAdaptervlan
-VMNetworkAdapterName "New Virtual Switch" -VlanId 5 -Access
-ManagementOS
```

### 注

请注意关于将 VLAN ID 添加到主机虚拟 NIC 的以下事项：

- 必须为主机虚拟 NIC 分配一个 VLAN ID。必须为所有接口和交换机上分配相同的 VLAN ID。
- 将主机虚拟 NIC 用于 RoCE 时，确保 VLAN ID 没有分配给物理接口。
- 如果创建多个主机虚拟 NIC，可以为每个主机虚拟 NIC 分配不同的 VLAN。

## 验证 RoCE 是否启用

要验证 RoCE 是否启用：

- 发出以下 Windows PowerShell 命令：

```
Get-NetAdapterRdma
```

命令输出列出 RDMA 支持的适配器，如图 13-4 中所示。

```
PS C:\Users\Administrator> Get-NetAdapterRdma
Name                InterfaceDescription          Enabled
-----
vEthernet (New Virtual... Hyper-V Virtual Ethernet Adapter True
```

图 13-4. Windows PowerShell 命令：Get-NetAdapterRdma

## 添加主机虚拟 NIC（虚拟端口）

### 要添加主机虚拟 NIC：

1. 要添加主机虚拟 NIC，请发出以下命令：  

```
Add-VMNetworkAdapter -SwitchName "New Virtual Switch" -Name  
SMB - ManagementOS
```
2. 在主机虚拟 NIC 上启用 RDMA，如第 205 页上的“要在主机虚拟 NIC 中启用 RDMA：”中所示。
3. 要将 VLAN ID 分配给虚拟端口，请发出以下命令：  

```
Set-VMNetworkAdapterVlan -VMNetworkAdapterName SMB -VlanId 5  
-Access -ManagementOS
```

## 映射 SMB 驱动器和运行 RoCE 流量

### 要映射 SMB 驱动器和运行 RoCE 流量：

1. 启动性能监视器 (Perfmon)。
2. 按如下方式填写 Add Counters（添加计数器）对话框（图 13-5）：
  - a. 在 **Available counters**（可用计数器）下，选择 **RDMA Activity**（RDMA 活动）。
  - b. 在 **Instances of selected object**（所选对象的实例）下，选择适配器。

c. 单击 **ADD**（添加）。

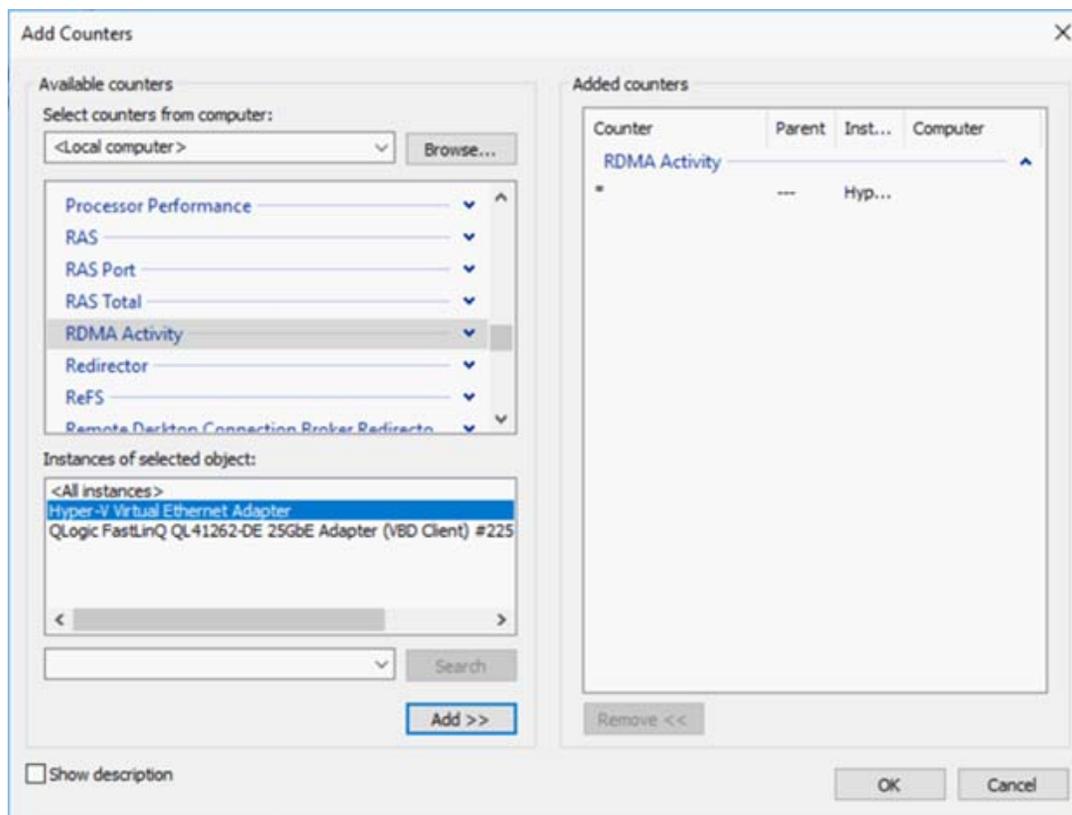


图 13-5. 添加计数器对话框

如果 RoCE 流量正在运行，计数器将显示为如图 13-6 中所示。

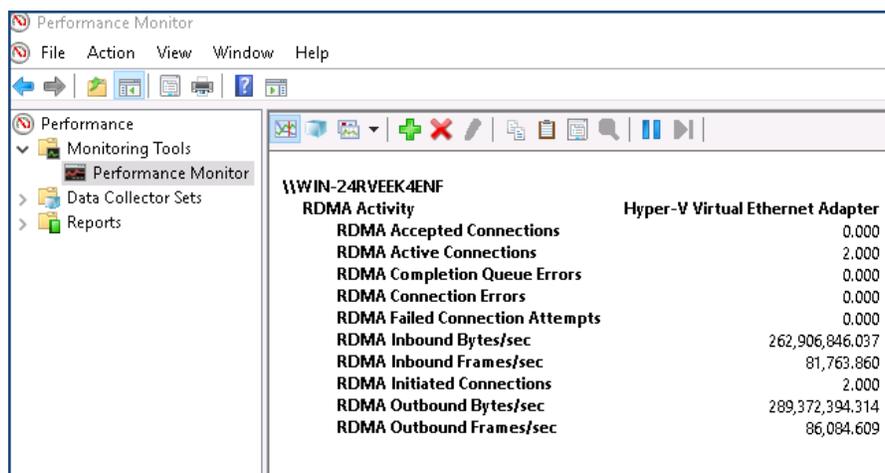


图 13-6. 性能监视器显示 RoCE 流量

## Switch Embedded Teaming 上的 RoCE

Switch Embedded Teaming (SET) 是 Microsoft 的备选 NIC 组合解决方案，可在 Windows Server 2016 Technical Preview 中包括 Hyper-V 和软件定义网络 (SDN) 堆栈的环境中使用。SET 将受限的 NIC 组合功能集成到 Hyper-V 虚拟交换机中。

使用 SET 将一到八个物理以太网网络适配器分组为一个或多个基于软件的虚拟网络适配器。如果网络适配器发生故障，这些适配器可提供快速性能和容错。SET 成员网络适配器必须均安装在同一物理 Hyper-V 主机中方可放在一个组中。

本节中包括以下 SET 上的 RoCE 步骤：

- 创建带 SET 和 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机
- 在 SET 上启用 RDMA
- 在 SET 上分配 VLAN ID
- 在 SET 上运行 RDMA 流量

### 创建带 SET 和 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机

要创建带 SET 和 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机：

- 要创建 SET，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> New-VMSwitch -Name SET  
-NetAdapterName "Ethernet 2","Ethernet 3"  
-EnableEmbeddedTeaming $true
```

图 13-7 显示命令输出。

```
PS C:\Users\Administrator> New-VMSwitch -Name SET -NetAdapterName "Ethernet 2","Ethernet 3" -EnableEmbeddedTeaming $true  
Name SwitchType NetAdapterInterfaceDescription  
-----  
SET External Teamed-Interface
```

图 13-7. Windows PowerShell 命令：New-VMSwitch

## 在 SET 上启用 RDMA

要在 SET 上启用 RDMA：

1. 要查看适配器上的 SET，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Get-NetAdapter "vEthernet (SET)"
```

图 13-8 显示命令输出。

```
PS C:\Users\Administrator> Get-NetAdapter "vEthernet (SET)"
```

Name	InterfaceDescription	ifIndex	Status	MacAddress	LinkSpeed
vEthernet (SET)	Hyper-V Virtual Ethernet Adapter	46	Up	00-0E-1E-C4-04-F8	50 Gbps

图 13-8. Windows PowerShell 命令：Get-NetAdapter

2. 要在 SET 上启用 RDMA，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Enable-NetAdapterRdma "vEthernet (SET)"
```

## 在 SET 上分配 VLAN ID

要在 SET 上分配 VLAN ID：

- 要在 SET 上分配 VLAN ID，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Set-VMNetworkAdapterVlan  
-VMNetworkAdapterName "SET" -VlanId 5 -Access -ManagementOS
```

### 注

请注意将 VLAN ID 添加到主机虚拟 NIC 时的以下事项：

- 将主机虚拟 NIC 用于 RoCE 时，确保 VLAN ID 没有分配给物理接口。
- 如果创建多个主机虚拟 NIC，可以为每个主机虚拟 NIC 分配不同的 VLAN。

## 在 SET 上运行 RDMA 流量

有关在 SET 上运行 RDMA 流量的信息，请访问：

<https://technet.microsoft.com/en-us/library/mt403349.aspx>

## 为 RoCE 配置 QoS

配置服务质量 (QoS) 的两种方法包括:

- [通过在适配器上禁用 DCBX 来配置 QoS](#)
- [通过在适配器上启用 DCBX 来配置 QoS](#)

### 通过在适配器上禁用 DCBX 来配置 QoS

通过在适配器上禁用 DCBX 来配置服务质量之前，必须完成所有使用中系统的所有配置。基于优先级的流控制 (PFC)、增强的转换服务 (ETS) 和流量类配置在交换机和服务器上必须相同。

**要通过禁用 DCBX 来配置 QoS:**

1. 在适配器上禁用 DCBX。
2. 使用 HII，将 **RoCE Priority** (RoCE 优先级) 设置为 0。
3. 要在主机中安装 DCB 角色，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrators> Install-WindowsFeature  
Data-Center-Bridging
```
4. 要将 **DCBX Willing** (DCBX 愿意) 模式设置为 **False** (假)，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrators> set-NetQosDcbxSetting -Willing 0
```
5. 请如下在微型端口中启用 QoS:
  - a. 打开微型端口窗口，然后单击 **Advanced** (高级) 选项卡。
  - b. 在适配器的 **Advanced Properties** (高级属性) 页面 (图 13-9) 的 **Property** (属性) 下，选择 **Quality of Service** (服务质量)，然后将值设置为 **Enabled** (启用)。
  - c. 单击 **OK** (确定)。

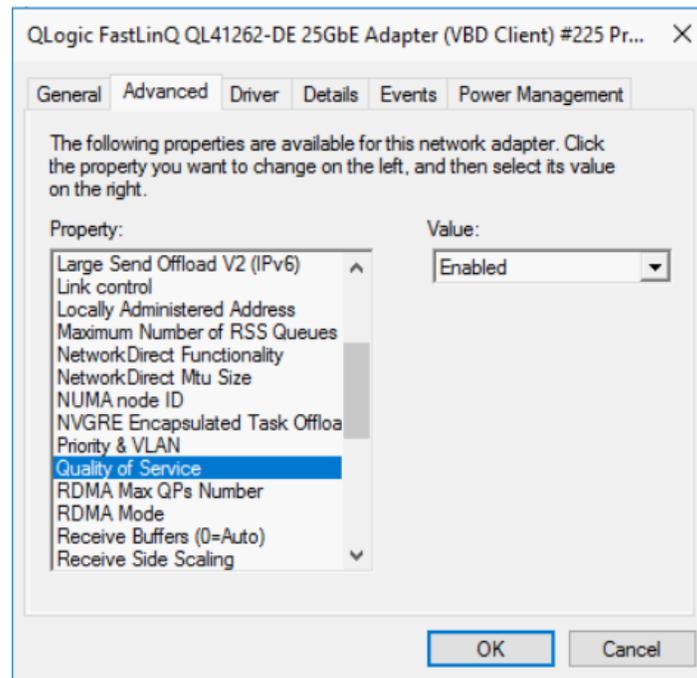


图 13-9. 高级属性：启用 QoS

6. 请如下将 VLAN ID 分配给接口：
  - a. 打开微型端口窗口，然后单击 **Advanced**（高级）选项卡。
  - b. 在适配器的 Advanced Properties（高级属性）页面（图 13-10）的 **Property**（属性）下，选择 **VLAN ID**，然后设置值。
  - c. 单击 **OK**（确定）。

#### 注

优先级流控制 (PFC) 需要上述步骤。

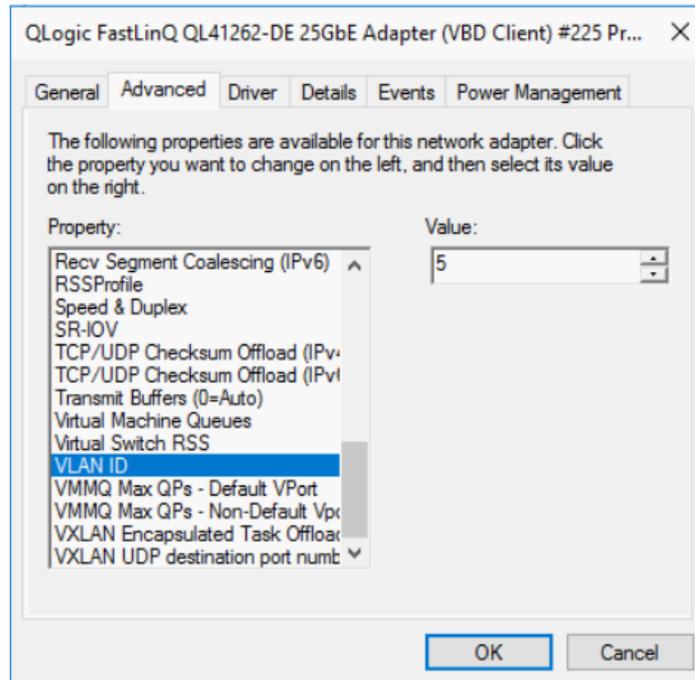


图 13-10. 高级属性：设置 VLAN ID

7. 要基于特定优先级启用 RoCE 的优先级流控制，请发出以下命令：

```
PS C:\Users\Administrators> Enable-NetQoSFlowControl  
-Priority 4
```

### 注

如果通过 Hyper-V 配置 RoCE，请勿将 VLAN ID 分配给物理接口。

8. 要基于任何其他优先级禁用优先级流控制，请发出以下命令：

```
PS C:\Users\Administrator> Disable-NetQoSFlowControl 0,1,2,3,5,6,7  
PS C:\Users\Administrator> Get-NetQoSFlowControl
```

Priority	Enabled	PolicySet	IfIndex	IfAlias
0	False	Global		
1	False	Global		
2	False	Global		
3	False	Global		
4	True	Global		
5	False	Global		
6	False	Global		
7	False	Global		

9. 要为每种类型的流量配置 QoS 和分配相关的优先级，请发出以下命令（其中，优先级 4 被标记用于 RoCE 而优先级 0 被标记用于 TCP）：

```
PS C:\Users\Administrators> New-NetQosPolicy "SMB"  
-NetDirectPortMatchCondition 445 -PriorityValue8021Action 4 -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  : 4
```

10. 要为之前步骤中定义的所有流量类配置 ETS，请发出以下命令：

```
PS C:\Users\Administrators> New-NetQosTrafficClass -name "RDMA class"  
-priority 4 -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	2-3,5-7	Global		
RDMA class	ETS	50	4	Global		
TCP class	ETS	30	0	Global		

11. 要查看之前配置的网络适配器 QoS，请发出以下 Windows PowerShell 命令：

```
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%      2-3,5-7
                    : 1 ETS    50%      4
                    : 2 ETS    30%      0

OperationalFlowControl   : Priority 4 Enabled
OperationalClassifications : Protocol  Port/Type Priority
                    : -----      -
                    : Default      0
                    : NetDirect 445    4
```

12. 创建启动脚本，以使设置在系统重新引导过程中保持不变。
13. 运行 RDMA 流量并验证，如第 60 页上的“RoCE 配置”中所述。

## 通过在适配器上启用 DCBX 来配置 QoS

必须在所有使用中的系统上完成所有配置。PFC、ETS 和流量类配置在交换机和服务器上必须相同。

### 要通过启用 DCBX 来配置 QoS：

1. 启用 DCBX（IEEE、CEE 或动态）。
2. 使用 HII，将 **RoCE Priority**（RoCE 优先级）设置为 0。

3. 要在主机中安装 DCB 角色，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrators> Install-WindowsFeature  
Data-Center-Bridging
```

### 注

对于此配置，将 **DCBX Protocol**（DCBX 协议）设置为 **CEE**。

4. 要将 **DCBX Willing**（DCBX 愿意）模式设置为 **True**（真），请发出以下命令：

```
PS C:\Users\Administrators> set-NetQosDcbxSetting -Willing 1
```

5. 请如下在微型端口中启用 QoS：

- a. 在适配器的 Advanced Properties（高级属性）页面（图 13-11）的 **Property**（属性）下，选择 **Quality of Service**（服务质量），然后将值设置为 **Enabled**（启用）。
- b. 单击 **OK**（确定）。

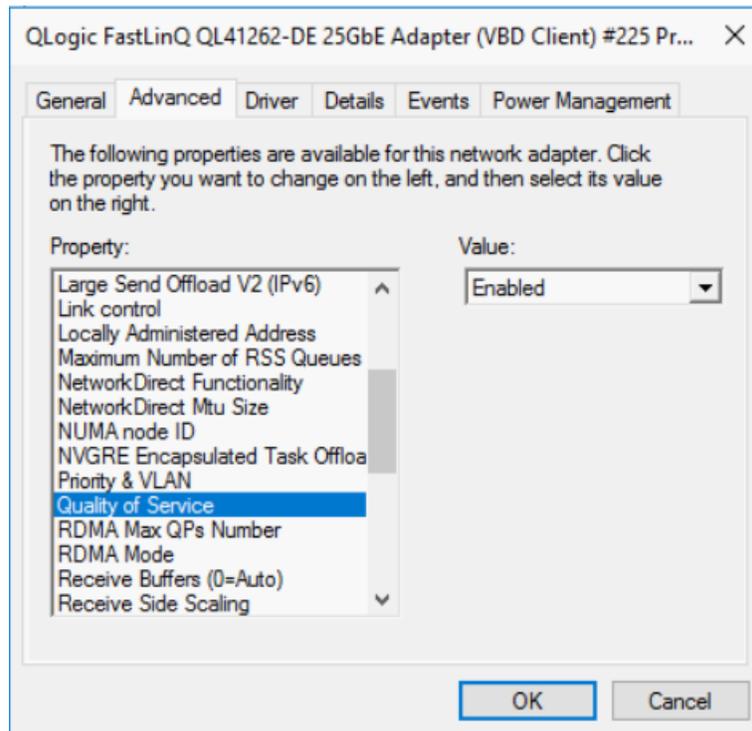


图 13-11. 高级属性：启用 QoS

6. 请如下将 VLAN ID 分配给接口（PFC 需要）：
  - a. 打开微型端口窗口，然后单击 **Advanced**（高级）选项卡。
  - b. 在适配器的 Advanced Properties（高级属性）页面（图 13-12）的 **Property**（属性）下，选择 **VLAN ID**，然后设置值。
  - c. 单击 **OK**（确定）。

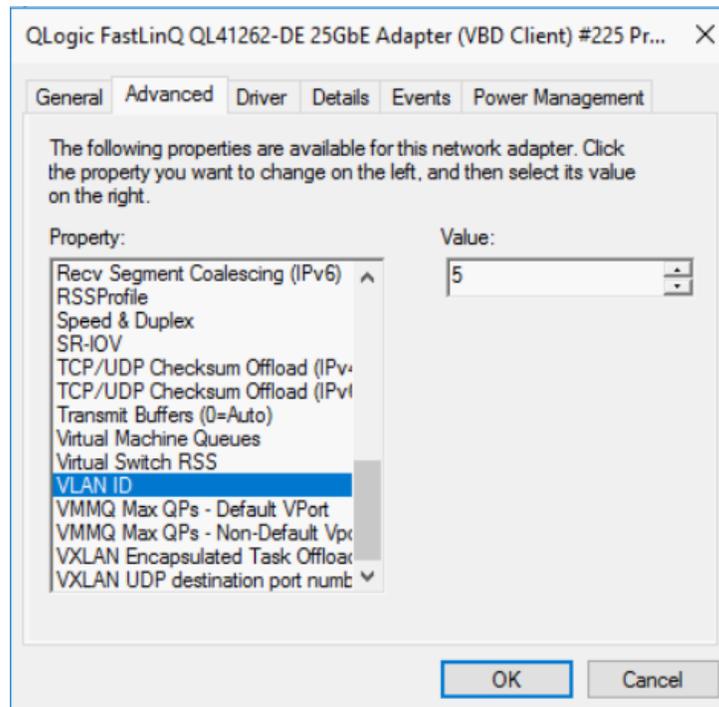


图 13-12. 高级属性：设置 VLAN ID

7. 要配置交换机，请发出以下 Windows PowerShell 命令：

```
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%      0-3,5-7
                          1 ETS      95%      4

OperationalFlowControl    : Priority 4 Enabled
OperationalClassifications : Protocol  Port/Type Priority
                          -----
                          NetDirect  445      4

RemoteTrafficClasses      : TC TSA      Bandwidth Priorities
                          --  ---      -
                          0 ETS      5%      0-3,5-7
                          1 ETS      95%      4

RemoteFlowControl         : Priority 4 Enabled
RemoteClassifications     : Protocol  Port/Type Priority
                          -----
                          NetDirect  445      4
```

---

### 注

上例为适配器端口连接到 Arista 7060X 交换机时采用。在此示例中，交换机 PFC 启用为优先级 4。RoCE App TLV 已定义。两个流量类定义为 TC0 和 TC1，其中 TC1 定义用于 RoCE。DCBX Protocol（DCBX 协议）模式设置为 CEE。有关 Arista 交换机配置，请参阅第 62 页上的“准备以太网交换机”。适配器处于 Willing（愿意）模式时，它会接受远程配置并将其显示为 Operational Parameters（操作参数）。

---

## 配置 VMMQ

虚拟机多队列 (VMMQ) 配置信息包括：

- 在适配器上启用 VMMQ
- 设置 VMMQ 最大 QP 默认和非默认 VPort
- 创建带或不带 SR-IOV 的虚拟机交换机
- 在虚拟机交换机上启用 VMMQ
- 获取虚拟机交换机功能
- 创建 VM 并在 VM 中的 VMNetworkadapter 上启用 VMMQ
- 默认和最大 VMMQ 虚拟 NIC
- 在管理 NIC 上启用和禁用 VMMQ
- 监测流量统计信息

### 在适配器上启用 VMMQ

要在适配器上启用 VMMQ：

1. 打开微型端口窗口，然后单击 **Advanced**（高级）选项卡。
2. 在 Advanced Properties（高级属性）页面（[图 13-13](#)）的 **Property**（属性）下，选择 **Virtual Switch RSS**（虚拟交换机 RSS），然后将值设置为 **Enabled**（启用）。

- 单击 **OK**（确定）。

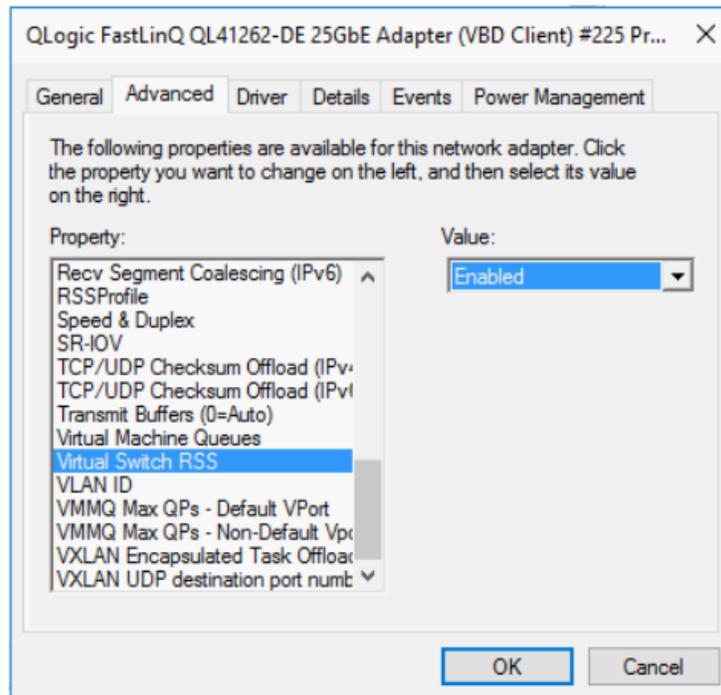


图 13-13. 高级属性：启用虚拟交换机 RSS

## 设置 VMMQ 最大 QP 默认和非默认 VPort

要设置 VMMQ 最大 QP 默认和非默认 VPort：

- 打开微型端口窗口，然后单击 **Advanced**（高级）选项卡。
- 在 Advanced Properties（高级属性）页面（图 13-14）的 **Property**（属性）下，选择以下项之一：
  - VMMQ Max QPs - Default VPort**（VMMQ 最大 QP - 默认 VPort）
  - VMMQ Max QPs - Non-Default VPort**（VMMQ 最大 QP - 非默认 VPort）
- 如果适用，调整所选属性的 **Value**（值）。

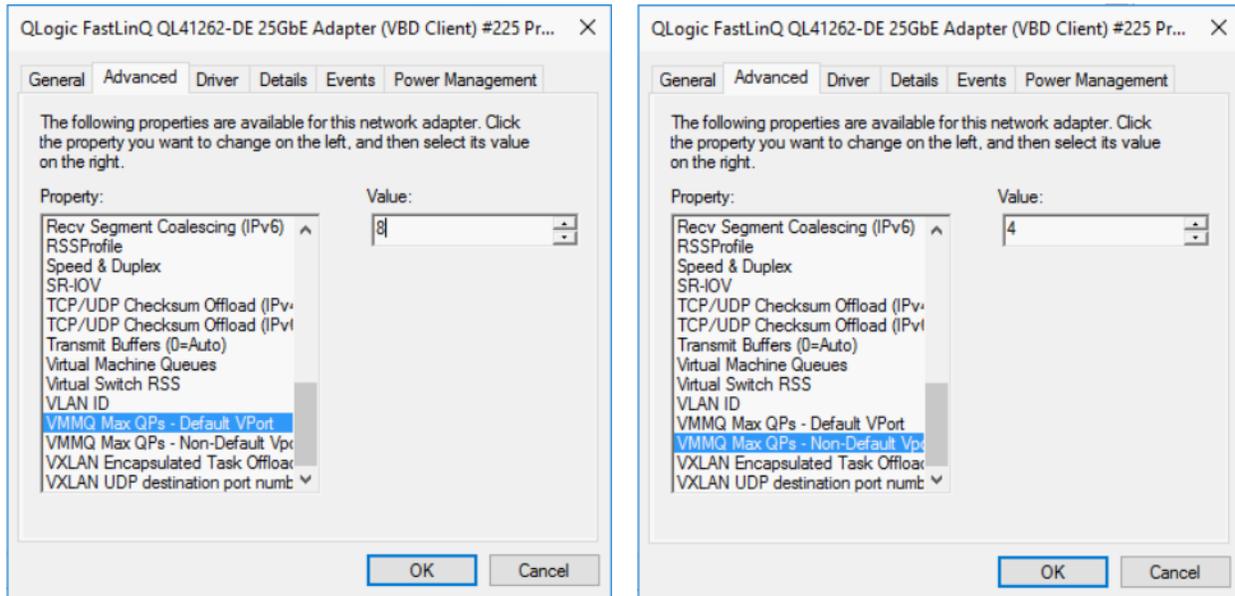


图 13-14. 高级属性：设置 VMMQ

4. 单击 **OK**（确定）。

## 创建带或不带 SR-IOV 的虚拟机交换机

要创建带或不带 SR-IOV 的虚拟机交换机：

1. 启动 Hyper-V 管理器。
2. 选择 **Virtual Switch Manager**（虚拟交换机管理器）（请参阅图 13-15）。
3. 在 **Name**（名称）框中，键入虚拟交换机的名称。

4. 在 **Connection type**（连接类型）下：
  - a. 单击 **External network**（外部网络）。
  - b. 选择 **Allow management operating system to share this network adapter**（允许管理操作系统共享此网络适配器）复选框。

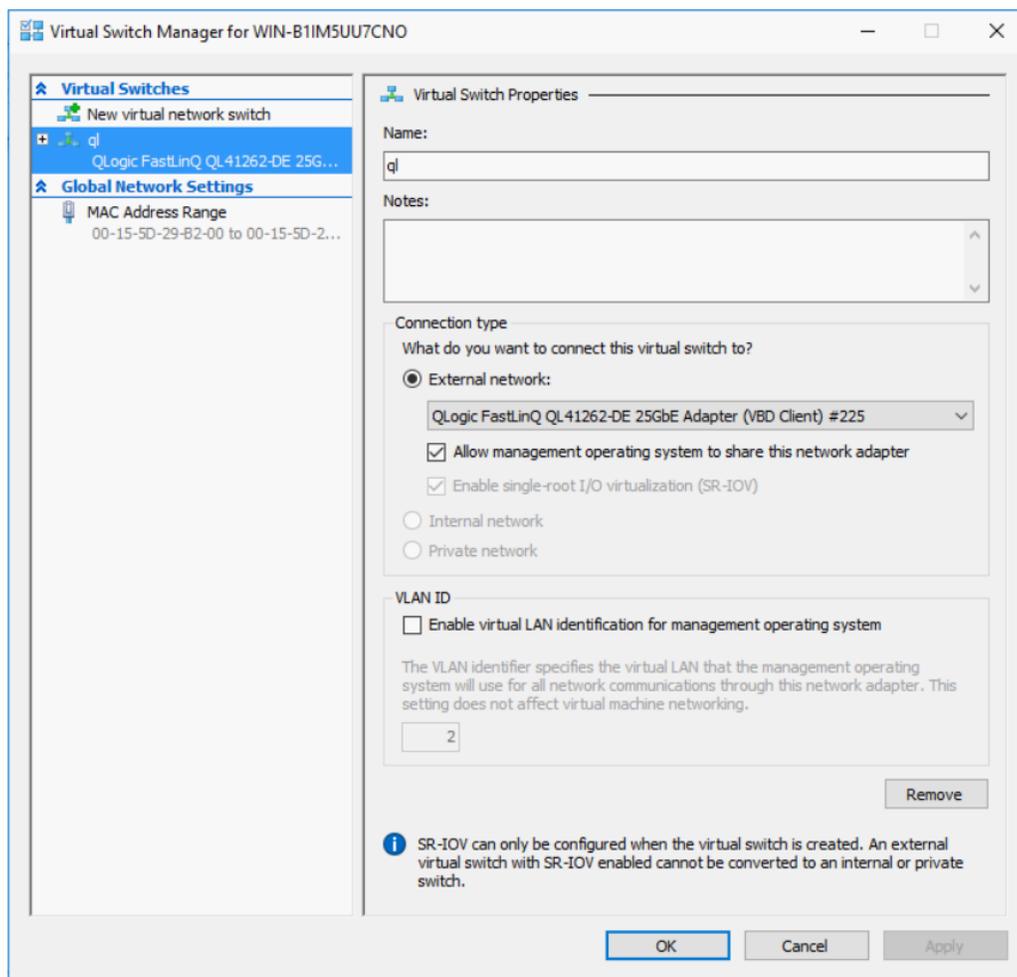


图 13-15. 虚拟交换机管理器

5. 单击 **OK**（确定）。

## 在虚拟机交换机上启用 VMMQ

要在虚拟机交换机上启用 VMMQ：

- 发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrators> Set-VMSwitch -name q1  
-defaultqueuevmmqenabled $true -defaultqueuevmmqqueuepairs 4
```

## 获取虚拟机交换机功能

要获取虚拟机交换机功能：

- 发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Get-VMSwitch -Name q1 | fl
```

图 13-16 显示示例输出。

```
PS C:\Users\Administrator> Get-VMSwitch -Name q1 | fl  
  
Name                : q1  
Id                  : 4dff5da3-f8bc-4146-a809-e1ddc6a04f7a  
Notes               :  
Extensions          : {Microsoft Windows Filtering Platform, Microsoft Azure VFP Switch Extension,  
Microsoft NDIS Capture}  
BandwidthReservationMode : None  
PacketDirectEnabled : False  
EmbeddedTeamingEnabled : False  
IovEnabled          : True  
SwitchType          : External  
AllowManagementOS  : True  
NetAdapterInterfaceDescription : QLogic FastLinQ QL41262-DE 25GbE Adapter (VBD Client) #225  
NetAdapterInterfaceDescriptions : {QLogic FastLinQ QL41262-DE 25GbE Adapter (VBD Client) #225}  
IovSupport           : True  
IovSupportReasons   :  
AvailableIPSecSA    : 0  
NumberIPSecSAAllocated : 0  
AvailableVMQueues   : 103  
NumberVmqAllocated  : 1  
IovQueuePairCount   : 127  
IovQueuePairsInUse  : 2  
IovVirtualFunctionCount : 96  
IovVirtualFunctionsInUse : 0  
PacketDirectInUse   : False  
DefaultQueueVrssEnabledRequested : True  
DefaultQueueVrssEnabled : True  
DefaultQueueVmmqEnabledRequested : False  
DefaultQueueVmmqEnabled : False  
DefaultQueueVmmqQueuePairsRequested : 16  
DefaultQueueVmmqQueuePairs : 16  
BandwidthPercentage : 0  
DefaultFlowMinimumBandwidthAbsolute : 0  
DefaultFlowMinimumBandwidthWeight : 0  
CimSession           : CimSession:  
ComputerName          : WIN-B1IM5U07CNO  
IsDeleted              : False
```

图 13-16. Windows PowerShell 命令：Get-VMSwitch

## 创建 VM 并在 VM 中的 VMNetworkadapter 上启用 VMMQ

要创建虚拟机 (VM) 并在该 VM 中的 VMNetworkadapter 上启用 VMMQ:

1. 创建 VM。
2. 将 VMNetworkadapter 添加到 VM。
3. 将虚拟交换机分配给 VMNetworkadapter。
4. 要在 VM 上启用 VMMQ，请发出以下 Windows PowerShell 命令:

```
PS C:\Users\Administrators> set-vmnetworkadapter -vmname vml
-VMNetworkAdapterName "network adapter" -vmmqenabled $true
-vmmqqueuepairs 4
```

### 注

对于支持 SR-IOV 的虚拟交换机：如果 VM 交换机和硬件加速启用了 SR-IOV，您必须创建 10 个 VM（每个带有 8 个虚拟 NIC）以利用 VMMQ。此要求是因为 SR-IOV 优先于 VMMQ。

64 个虚拟功能和 16 个 VMMQ 的示例输出如下所示:

```
PS C:\Users\Administrator> get-netadaptervport
```

Name	ID	MacAddress	VID	ProcMask	FID	State	ITR	QPairs
-----	--	-----	---	-----	---	-----	---	-----
Ethernet 3	0	00-15-5D-36-0A-FB		0:0	PF	Activated	Unknown	4
Ethernet 3	1	00-0E-1E-C4-C0-A4		0:8	PF	Activated	Adaptive	4
Ethernet 3	2			0:0	0	Activated	Unknown	1
Ethernet 3	3			0:0	1	Activated	Unknown	1
Ethernet 3	4			0:0	2	Activated	Unknown	1
Ethernet 3	5			0:0	3	Activated	Unknown	1
Ethernet 3	6			0:0	4	Activated	Unknown	1
Ethernet 3	7			0:0	5	Activated	Unknown	1
Ethernet 3	8			0:0	6	Activated	Unknown	1
Ethernet 3	9			0:0	7	Activated	Unknown	1
Ethernet 3	10			0:0	8	Activated	Unknown	1
Ethernet 3	11			0:0	9	Activated	Unknown	1
.								
.								
.								
Ethernet 3	64			0:0	62	Activated	Unknown	1
Ethernet 3	65			0:0	63	Activated	Unknown	1
Ethernet 3	66	00-15-5D-36-0A-04		0:16	PF	Activated	Adaptive	4
Ethernet 3	67	00-15-5D-36-0A-05		1:0	PF	Activated	Adaptive	4
Ethernet 3	68	00-15-5D-36-0A-06		0:0	PF	Activated	Adaptive	4
Name	ID	MacAddress	VID	ProcMask	FID	State	ITR	QPairs
-----	--	-----	---	-----	---	-----	---	-----
Ethernet 3	69	00-15-5D-36-0A-07		0:8	PF	Activated	Adaptive	4
Ethernet 3	70	00-15-5D-36-0A-08		0:16	PF	Activated	Adaptive	4

Ethernet 3	71	00-15-5D-36-0A-09	1:0	PF	Activated	Adaptive	4
Ethernet 3	72	00-15-5D-36-0A-0A	0:0	PF	Activated	Adaptive	4
Ethernet 3	73	00-15-5D-36-0A-0B	0:8	PF	Activated	Adaptive	4
Ethernet 3	74	00-15-5D-36-0A-F4	0:16	PF	Activated	Adaptive	4
Ethernet 3	75	00-15-5D-36-0A-F5	1:0	PF	Activated	Adaptive	4
Ethernet 3	76	00-15-5D-36-0A-F6	0:0	PF	Activated	Adaptive	4
Ethernet 3	77	00-15-5D-36-0A-F7	0:8	PF	Activated	Adaptive	4
Ethernet 3	78	00-15-5D-36-0A-F8	0:16	PF	Activated	Adaptive	4
Ethernet 3	79	00-15-5D-36-0A-F9	1:0	PF	Activated	Adaptive	4
Ethernet 3	80	00-15-5D-36-0A-FA	0:0	PF	Activated	Adaptive	4

```
PS C:\Users\Administrator> get-netadaptervmq
```

Name	InterfaceDescription	Enabled	BaseVmqProcessor	MaxProcessors	NumberOfReceive Queues
-----	-----	-----	-----	-----	-----
Ethernet 4	QLogic FastLinQ 41xxx	False	0:0	16	1

## 默认和最大 VMMQ 虚拟 NIC

根据当前实施，每个虚拟 NIC 最多 4 个可用的 VMMQ；也就是说，最多 16 个虚拟 NIC。

如之前使用 Windows PowerShell 命令所设置，有四个可用的默认队列。最大默认队列当前可设置为 8。要验证最大默认队列，请使用 VMswitch 功能。

## 在管理 NIC 上启用和禁用 VMMQ

**要在管理 NIC 上启用或禁用 VMMQ：**

- 要在管理 NIC 上启用 VMMQ，请发出以下命令：

```
PS C:\Users\Administrator> Set-VMNetworkAdapter -ManagementOS  
-vmmqEnabled $true
```

MOS vNIC 拥有四个 VMMQ。

- 要在管理 NIC 上禁用 VMMQ，请发出以下命令：

```
PS C:\Users\Administrator> Set-VMNetworkAdapter -ManagementOS  
-vmmqEnabled $false
```

VMMQ 还将可用于多播开放式最短路径优先 (MOSPF)。

## 监测流量统计信息

要监测虚拟机中的虚拟功能流量，请发出以下 Windows PowerShell 命令：

```
PS C:\Users\Administrator> Use get-netadapterstatistics | fl
```

## 配置 VXLAN

VXLAN 配置信息包括：

- 在适配器上启用 VXLAN 卸载
- 部署软件定义网络

### 在适配器上启用 VXLAN 卸载

要在适配器上启用 VXLAN 卸载：

1. 打开微型端口窗口，然后单击 **Advanced**（高级）选项卡。
2. 在 Advanced Properties（高级属性）页面（图 13-17）的 **Property**（属性）下，选择 **VXLAN Encapsulated Task Offload**（VXLAN 封装任务卸载）。

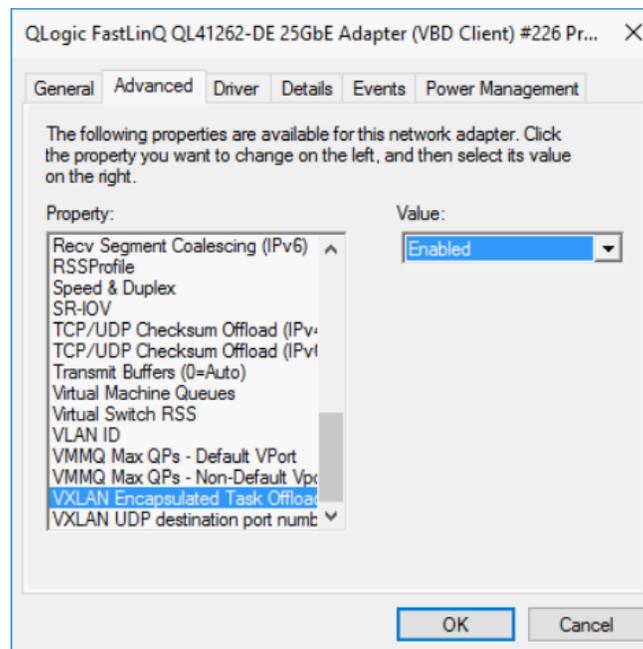


图 13-17. 高级属性：启用 VXLAN

3. 将 **Value**（值）设置为 **Enabled**（启用）。
4. 单击 **OK**（确定）。

## 部署软件定义网络

要利用虚拟机上的 VXLAN 封装任务卸载，必须部署使用 Microsoft 网络控制器的软件定义网络 (SDN) 堆栈。

有关更多详细信息，请参阅有关软件定义网络的以下 Microsoft TechNet 链接：

<https://technet.microsoft.com/en-us/windows-server-docs/networking/sdn/software-defined-networking--sdn->

## 配置 Storage Spaces Direct

Windows Server 2016 引入了 Storage Spaces Direct，这样便可通过本地存储构建高度可用和可扩展的存储系统。

有关更多信息，请参阅以下 Microsoft TechnNet 链接：

<https://technet.microsoft.com/en-us/windows-server-docs/storage/storage-spaces/storage-spaces-direct-windows-server-2016>

## 配置硬件

图 13-18 显示硬件配置的示例。

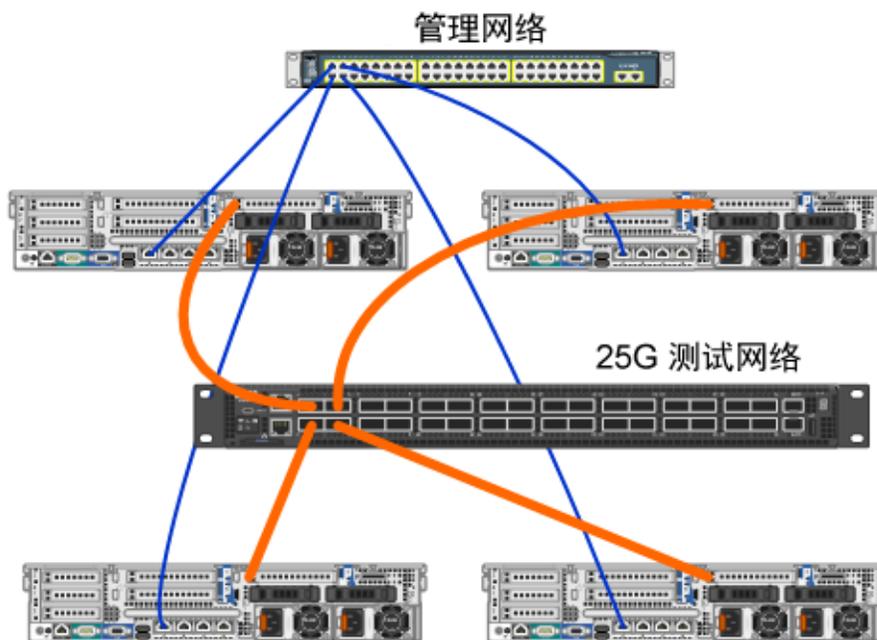


图 13-18. 示例硬件配置

### 注

本示例中使用的磁盘为 4 × 400G NVMe™ 和 12 × 200G SSD 磁盘。

---

## 部署超聚合系统

本节包括使用 Windows Server 2016 安装和配置超聚合系统组件的简介。部署超聚合系统的操作可以分为以下三个高级阶段：

- [部署操作系统](#)
- [配置网络](#)
- [配置 Storage Spaces Direct](#)

### 部署操作系统

#### 要部署操作系统：

1. 安装操作系统。
2. 安装 Windows 服务器角色 (Hyper-V)。
3. 安装以下功能：
  - 故障转移
  - 群集
  - 数据中心桥接 (DCB)
4. 将节点连接到域并添加域帐户。

### 配置网络

要部署 Storage Spaces Direct，Hyper-V 交换机必须部署有启用 RDMA 的主机虚拟 NIC。

### 注

以下步骤假设有四个 RDMA NIC 端口。

---

#### 要在每个服务器上配置网络：

1. 请按如下来配置物理网络交换机：
  - a. 将所有适配器 NIC 连接到交换机端口。

### 注

如果测试适配器有多个 NIC 端口，您必须将两个端口连接到同一交换机。

---

- b. 启用交换机端口，确保交换机端口支持与交换机无关的组合模式，同时是多个 VLAN 网络的组成部分。

示例 Dell 交换机配置：

```
no ip address
mtu 9416
portmode hybrid
switchport
dcb-map roce_S2D
protocol lldp
dcbx version cee
no shutdown
```

2. 启用 **Network Quality of Service**（网络服务质量）。

---

### 注

网络服务质量用于确保软件定义的存储系统有足够的带宽在节点之间通信，以确保弹性和性能。要在适配器上配置 QoS，请参阅[第 212 页上的“为 RoCE 配置 QoS”](#)。

---

3. 请如下创建带 SET 和 RDMA 虚拟 NIC 的 Hyper-V 虚拟交换机：

- a. 要标识网络适配器，请发出以下命令：

```
Get-NetAdapter | FT
Name,InterfaceDescription,Status,LinkSpeed
```

- b. 要创建连接到所有物理网络适配器的虚拟交换机，然后启用 Switch Embedded Teaming，请发出以下命令：

```
New-VMSwitch -Name SETswitch -NetAdapterName
"<port1>","<port2>","<port3>","<port4>"
-EnableEmbeddedTeaming $true
```

- c. 要将主机虚拟 NIC 添加到虚拟交换机，请发出以下命令：

```
Add-VMNetworkAdapter -SwitchName SETswitch -Name SMB_1
-managementOS
Add-VMNetworkAdapter -SwitchName SETswitch -Name SMB_2
-managementOS
```

---

### 注

上述命令从您刚刚为要使用的管理操作系统配置的虚拟交换机来配置虚拟 NIC。

---

- d. 要配置主机虚拟 NIC 以使用 VLAN，请发出以下命令：

```
Set-VMNetworkAdapterVlan -VMNetworkAdapterName "SMB_1"  
-VlanId 5 -Access -ManagementOS  
Set-VMNetworkAdapterVlan -VMNetworkAdapterName "SMB_2"  
-VlanId 5 -Access -ManagementOS
```

---

### 注

这些命令可位于相同或不同的 VLAN 上。

---

- e. 要验证 VLAN ID 是否已设置，请发出以下命令：

```
Get-VMNetworkAdapterVlan -ManagementOS
```

- f. 要禁用和启用每个主机虚拟 NIC 适配器以使 VLAN 处于活动状态，请发出以下命令：

```
Disable-NetAdapter "vEthernet (SMB_1)"  
Enable-NetAdapter "vEthernet (SMB_1)"  
Disable-NetAdapter "vEthernet (SMB_2)"  
Enable-NetAdapter "vEthernet (SMB_2)"
```

- g. 要在主机虚拟 NIC 适配器上启用 RDMA，请发出以下命令：

```
Enable-NetAdapterRdma "SMB1","SMB2"
```

- h. 要验证 RDMA 功能，请发出以下命令：

```
Get-SmbClientNetworkInterface | where RdmaCapable -EQ  
$true
```

## 配置 Storage Spaces Direct

在 Windows Server 2016 中配置 Storage Spaces Direct 包括以下步骤：

- [步骤 1. 运行群集验证工具](#)
- [步骤 2. 创建群集](#)
- [步骤 3. 配置群集见证](#)
- [步骤 4. 清理用于 Storage Spaces Direct 的磁盘](#)
- [步骤 5. 启用 Storage Spaces Direct](#)
- [步骤 6. 创建虚拟磁盘](#)
- [步骤 7. 创建或部署虚拟机](#)

### 步骤 1. 运行群集验证工具

运行群集验证工具以确保服务器节点正确配置，从而使用 Storage Spaces Direct 创建群集。

发出以下 Windows PowerShell 命令，验证用作 Storage Spaces Direct 群集的一组服务器：

```
Test-Cluster -Node <MachineName1, MachineName2, MachineName3,  
MachineName4> -Include "Storage Spaces Direct", Inventory,  
Network, "System Configuration"
```

### 步骤 2. 创建群集

创建 [步骤 1. 运行群集验证工具](#) 中具有四个节点的群集（群集创建已验证）。

要创建群集，请发出以下 Windows PowerShell 命令：

```
New-Cluster -Name <ClusterName> -Node <MachineName1, MachineName2,  
MachineName3, MachineName4> -NoStorage
```

-NoStorage 参数是必需的。如果未包括该参数，磁盘将自动添加到群集，您必须移除磁盘，然后才能启用 Storage Spaces Direct。否则，磁盘将不包括在 Storage Spaces Direct 存储池中。

### 步骤 3. 配置群集见证

您应该配置群集的见证，以使该四节点系统能够承受两个节点发生故障或脱机。对于这些系统，您可配置文件共享见证或云见证。

有关详细信息，请访问：

<https://blogs.msdn.microsoft.com/clustering/2014/03/31/configuring-a-file-share-witness-on-a-scale-out-file-server/>

### 步骤 4. 清理用于 Storage Spaces Direct 的磁盘

旨在用于 Storage Spaces Direct 的磁盘必须为空，且不带分区或其他数据。如果磁盘有分区或其他数据，该磁盘将不会包括在 Storage Spaces Direct 系统中。

可在 Windows PowerShell 脚本 (.PS1) 文件中放入以下 Windows PowerShell 命令，并从管理系统打开的 Windows PowerShell（或 Windows PowerShell ISE）控制台以管理员权限执行：

---

#### 注

运行此脚本可帮助识别每个节点上可用于 Storage Spaces Direct 的磁盘，并从这些磁盘移除所有数据和分区。

---

```
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
```

## 步骤 5. 启用 Storage Spaces Direct

创建群集后，发出 `Enable-ClusterStorageSpacesDirect Windows PowerShell cmdlet`。该 cmdlet 会将存储系统置于 Storage Spaces Direct 模式并自动执行以下操作：

- 创建名为诸如 *S2D on Cluster1* 的单个大型池。
- 配置 Storage Spaces Direct 高速缓存。如果有多种介质类型可供 Storage Spaces Direct 使用，它会配置最有效的类型作为高速缓存设备（在大多数情况下，读取和写入）。
- 创建两个层作为默认层：**Capacity**（容量）和 **Performance**（性能）。该 cmdlet 会分析设备并以混合的设备类型和弹性配置每个层。

## 步骤 6. 创建虚拟磁盘

如果 Storage Spaces Direct 已启用，它会使用所有磁盘创建一个池。它还会命名该池（例如，*S2D on Cluster1*），并在名称中指定群集名称。

以下 Windows PowerShell 命令在存储池上创建具有镜像和奇偶校验弹性的虚拟磁盘：

```
New-Volume -StoragePoolFriendlyName "S2D*" -FriendlyName
<VirtualDiskName> -FileSystem CSVFS_ReFS -StorageTierfriendlyNames
Capacity,Performance -StorageTierSizes <Size of capacity tier in
size units, example: 800GB>, <Size of Performance tier in size
units, example: 80GB> -CimSession <ClusterName>
```

### 步骤 7. 创建或部署虚拟机

您可以在超聚合 S2D 群集的节点上部署虚拟机。将虚拟机的文件存储在系统的 CSV 名称空间（例如，c:\ClusterStorage\Volume1）中，与故障转移群集上的群集虚拟机类似。

## 部署和管理 Nano Server

Windows Server 2016 提供 Nano Server 作为新的安装选项。Nano Server 是远程管理的服务器操作系统，优化用于私有云和数据中心。它类似于服务器核心模式下的 Windows Server，但明显更小，没有本地登录功能，并且仅支持 64 位应用程序、工具和代理。与 Windows Server 相比，Nano Server 占用较少的磁盘空间，设置更快，并且需要较少的更新和重新启动次数。重新启动时，它会重新启动得更快。

### 角色和功能

表 13-1 显示此版本的 Nano Server 中可用的角色和功能，以及将为其安装软件包的 Windows PowerShell 选项。某些软件包直接通过自己的 Windows PowerShell 选项（诸如 -Compute）进行安装。其他软件包作为 -Packages 选项的扩展进行安装，您可在逗号分隔的列表中进行组合。

**表 13-1. Nano Server 的角色和功能**

角色或功能	选项
Hyper-V 角色	-Compute
故障转移群集	-Clustering
用于将 Nano Server 托管为虚拟机的 Hyper-V 来宾驱动程序	-GuestDrivers
用于各种网络适配器和存储控制器的基本驱动程序。这与 Windows Server 2016 Technical Preview 服务器核心安装中包括的驱动程序集相同。	-OEMDrivers
文件服务器角色和其他存储组件	-Storage

**表 13-1. Nano Server 的角色和功能 (续)**

角色或功能	选项
Windows Defender 反恶意软件，包括默认的签名文件	-Defender
应用程序兼容性的反向转发器：例如，Ruby、Node.js 及其他常用应用程序框架。	-ReverseForwarders
DNS 服务器角色	-Packages Microsoft-NanoServer-DNS-Package
理想状态配置 (DSC)	-Packages Microsoft-NanoServer-DSC-Package
互联网信息服务器 (IIS)	-Packages Microsoft-NanoServer-IIS-Package
Windows 容器的主机支持	-Containers
系统中心虚拟机管理器代理 (System Center Virtual Machine Manager Agent)	-Packages Microsoft-Windows-Server-SCVMM-Package -Packages Microsoft-Windows-Server-SCVMM-Compute-Package 注：仅当您监测 Hyper-V 时使用此软件包。如果您安装此软件包，请勿对 Hyper-V 角色使用 -Compute 选项；而是使用 -Packages 选项来安装 -Packages Microsoft-NanoServer-Compute-Package、Microsoft-Windows-Server-SCVMM-Compute-Package。
网络性能诊断服务 (NPDS)	-Packages Microsoft-NanoServer-NPDS-Package
数据中心桥接	-Packages Microsoft-NanoServer-DCB-Package

接下来的各节介绍如何使用所需的软件包配置 Nano Server 映像，以及如何添加特定于 QLogic 设备的附加设备驱动程序。它们还说明如何使用 Nano Server Recovery Console，如何远程管理 Nano Server，以及如何从 Nano Server 运行 Ntttcp 流量。

## 在物理服务器上部署 Nano Server

按照以下步骤操作，使用预安装的设备驱动程序创建将在物理服务器上运行的 Nano Server 虚拟硬盘 (VHD)。

### 要部署 Nano Server：

1. 下载 Windows Server 2016 操作系统映像。
2. 装载 ISO。
3. 将以下文件从 NanoServer 文件夹复制到您的硬盘驱动器上的文件夹：
  - NanoServerImageGenerator.psml
  - Convert-WindowsImage.ps1
4. 以管理员身份启动 Windows PowerShell。
5. 将目录更改为 [步骤 3](#) 中您粘贴文件的文件夹。
6. 通过发出以下命令，导入 NanoServerImageGenerator 脚本：

```
Import-Module .\NanoServerImageGenerator.psml -Verbose
```
7. 要创建设置计算机名称并包括 OEM 驱动程序和 Hyper-V 的 VHD，请发出以下 Windows PowerShell 命令：

---

### 注

此命令将提示您为新 VHD 输入管理员密码。

---

```
New-NanoServerImage -DeploymentType Host -Edition  
<Standard/Datacenter> -MediaPath <path to root of media>  
-BasePath  
. \Base -TargetPath .\NanoServerPhysical\NanoServer.vhd  
-ComputerName  
<computer name> -Compute -Storage -Cluster -OEMDrivers  
-Compute  
-DriversPath "<Path to Qlogic Driver sets>"
```

### 示例：

```
New-NanoServerImage -DeploymentType Host -Edition Datacenter  
-MediaPath C:\tmp\TP4_iso\Bld_10586_iso  
-BasePath ".\Base" -TargetPath  
"C:\Nano\PhysicalSystem\Nano_phy_vhd.vhd" -ComputerName  
"Nano-server1" -Compute -Storage -Cluster -OEMDrivers  
-DriversPath  
"C:\Nano\Drivers"
```

在上例中，`C:\Nano\Drivers` 是 QLogic 驱动程序的路径。此命令需要大约 10 到 15 分钟来创建 VHD 文件。此命令的示例输出如下所示：

```
Windows(R) Image to Virtual Hard Disk Converter for Windows(R) 10
Copyright (C) Microsoft Corporation. All rights reserved.
Version 10.0.14300.1000.amd64fre.rs1_release_svc.160324-1723
INFO : Looking for the requested Windows image in the WIM file
INFO : Image 1 selected (ServerDatacenterNano)...
INFO : Creating sparse disk...
INFO : Mounting VHD...
INFO : Initializing disk...
INFO : Creating single partition...
INFO : Formatting windows volume...
INFO : Windows path (I:) has been assigned.
INFO : System volume location: I:
INFO : Applying image to VHD.This could take a while...
INFO : Image was applied successfully.
INFO : Making image bootable...
INFO : Fixing the Device ID in the BCD store on VHD...
INFO : Drive is bootable. Cleaning up...
INFO : Dismounting VHD...
INFO : Closing Windows image...
INFO : Done.
Done.The log is at:
C:\Users\ADMINI~1\AppData\Local\Temp\2\NanoServerImageGenerator.log
```

8. 在您要运行 Nano Server VHD 的物理服务器上以管理员身份登录。
9. 要将 VHD 复制到物理服务器并将其配置为从新 VHD 引导：
  - a. 转至 **Computer Management**（计算机管理）> **Storage**（存储）> **Disk Management**（磁盘管理）。
  - b. 右键单击 **Disk Management**（磁盘管理），然后选择 **Attach VHD**（附加 VHD）。
  - c. 提供 VHD 文件路径。
  - d. 单击 **OK**（确定）。
  - e. 运行 `bcdboot d:\windows`。

---

**注**

在本示例中，VHD 附加在 `D:\` 下。

---

- f. 右键单击 **Disk Management** (磁盘管理)，然后选择 **Detach VHD** (分离 VHD)。
10. 将物理服务器重新引导到 Nano Server VHD 中。
11. 使用在 [步骤 7](#) 中运行脚本时您提供的管理员和密码登录到 Recovery Console (恢复控制台)。
12. 获取 Nano Server 计算机的 IP 地址。
13. 使用 Windows PowerShell 远程处理 (或其他远程管理) 工具连接并远程管理服务器。

## 在虚拟机中部署 Nano Server

要创建 Nano Server 虚拟硬盘驱动器 (VHD) 以在虚拟机中运行：

1. 下载 Windows Server 2016 操作系统映像。
2. 转至 [步骤 1](#) 中下载的文件所在的 NanoServer 文件夹。
3. 将以下文件从 NanoServer 文件夹复制到您的硬盘驱动器上的文件夹：
  - NanoServerImageGenerator.psml
  - Convert-WindowsImage.ps1
4. 以管理员身份启动 Windows PowerShell。
5. 将目录更改为 [步骤 3](#) 中您粘贴文件的文件夹。
6. 通过发出以下命令，导入 NanoServerImageGenerator 脚本：

```
Import-Module .\NanoServerImageGenerator.psml -Verbose
```
7. 发出以下 Windows PowerShell 命令，创建设置计算机名称并包括 Hyper-V 来宾驱动程序程序的 VHD：

---

### 注

以下命令将提示您为新 VHD 输入管理员密码。

---

```
New-NanoServerImage -DeploymentType Guest -Edition  
<Standard/Datacenter> -MediaPath <path to root of media>  
-BasePath  
. \Base -TargetPath .\NanoServerPhysical\NanoServer.vhd  
-ComputerName  
<computer name> -GuestDrivers
```

示例:

```
New-NanoServerImage -DeploymentType Guest -Edition Datacenter
-MediaPath C:\tmp\TP4_iso\Bld_10586_iso
-BasePath .\Base -TargetPath .\Nano1\VM_NanoServer.vhd
-ComputerName
Nano-VM1 -GuestDrivers
```

上述命令需要大约 10 到 15 分钟来创建 VHD 文件。此命令的示例输出如下:

```
PS C:\Nano> New-NanoServerImage -DeploymentType Guest -Edition
Datacenter -MediaPath
C:\tmp\TP4_iso\Bld_10586_iso -BasePath .\Base -TargetPath
.\Nano1\VM_NanoServer.vhd -ComputerName Nano-VM1 -GuestDrivers
cmdlet New-NanoServerImage at command pipeline position 1
Supply values for the following parameters:
Windows(R) Image to Virtual Hard Disk Converter for Windows(R) 10
Copyright (C) Microsoft Corporation. All rights reserved.
Version 10.0.14300.1000.amd64fre.rs1_release_svc.160324-1723
INFO : Looking for the requested Windows image in the WIM file
INFO : Image 1 selected (ServerTuva)...
INFO : Creating sparse disk...
INFO : Attaching VHD...
INFO : Initializing disk...
INFO : Creating single partition...
INFO : Formatting windows volume...
INFO : Windows path (G:) has been assigned.
INFO : System volume location: G:
INFO : Applying image to VHD.This could take a while...
INFO : Image was applied successfully.
INFO : Making image bootable...
INFO : Fixing the Device ID in the BCD store on VHD...
INFO : Drive is bootable. Cleaning up...
INFO : Closing VHD...
INFO : Deleting pre-existing VHD : Base.vhd...
INFO : Closing Windows image...
INFO : Done.
Done.The log is at:
C:\Users\ADMINI~1\AppData\Local\Temp\2\NanoServerImageGenerator.log
```

8. 在 Hyper-V 管理器中创建新的虚拟机，并使用[步骤 7](#)中创建的 VHD。
9. 引导虚拟机。

10. 连接到 Hyper-V 管理器中的虚拟机。
11. 使用在 [步骤 7](#) 中运行脚本时您提供的管理员和密码登录到 Recovery Console（恢复控制台）。
12. 获取 Nano Server 计算机的 IP 地址。
13. 使用 Windows PowerShell 远程处理（或其他远程管理）工具连接并远程管理服务器。

## 远程管理 Nano Server

远程管理 Nano Server 的选项包括 Windows PowerShell、Windows Management Instrumentation (WMI)、Windows 远程管理和紧急管理服务 (EMS)。本节介绍如何使用 Windows PowerShell 远程处理访问 Nano Server。

### 使用 Windows PowerShell 远程处理管理 Nano Server

**要使用 Windows PowerShell 远程处理来管理 Nano Server：**

1. 将 Nano Server 的 IP 地址添加到管理计算机的受信任的主机列表。

---

#### 注

使用恢复控制台查找服务器 IP 地址。

---

2. 添加用于 Nano Server 管理员的帐户。
3. （可选）如果适用，启用 **CredSSP**。

### 将 Nano Server 添加到受信任的主机列表

在提升权限的 Windows PowerShell 提示符下，通过发出以下命令，将 Nano Server 添加到受信任的主机列表：

```
Set-Item WSMAN:\localhost\Client\TrustedHosts "<IP address of Nano Server>"
```

示例：

```
Set-Item WSMAN:\localhost\Client\TrustedHosts "172.28.41.152"  
Set-Item WSMAN:\localhost\Client\TrustedHosts "**"
```

---

#### 注

上述命令将所有主机服务器设置为受信任的主机。

---

## 启动远程 Windows PowerShell 会话

在提升权限的本地 Windows PowerShell 会话中，通过发出以下命令，启动远程 Windows PowerShell 会话：

```
$ip = "<IP address of Nano Server>"  
$user = "$ip\Administrator"  
Enter-PSSession -ComputerName $ip -Credential $user
```

现在，您可以在 Nano Server 上如常运行 Windows PowerShell 命令。不过，并非所有 Windows PowerShell 命令在此版本的 Nano Server 中均可用。要查看哪些命令可用，请发出命令 `Get-Command -CommandType Cmdlet`。要停止远程会话，请发出命令 `Exit-PSSession`。

有关 Nano Server 的更多详细信息，请访问：

<https://technet.microsoft.com/en-us/library/mt126167.aspx>

## 在 Windows Nano Server 上管理 QLogic 适配器

要在 Nano Server 环境中管理 QLogic 适配器，请参阅 Windows QConvergeConsole GUI 和 Windows QLogic Control Suite CLI 管理工具及关联的说明文件，这些内容可从 Cavium 网站获得：

## RoCE 配置

要使用 Windows PowerShell 远程处理管理 Nano Server，

1. 从另一台计算机通过 Windows PowerShell 远程处理连接到 Nano Server。  
例如：

```
PS C:\Windows\system32> $ip="172.28.41.152"  
PS C:\Windows\system32> $user="172.28.41.152\Administrator"  
PS C:\Windows\system32> Enter-PSSession -ComputerName $ip  
-Credential $user
```

---

### 注

在上例中，Nano Server IP 地址为 172.28.41.152，用户名为 Administrator。

---

如果 Nano Server 成功连接，将返回以下内容：

```
[172.28.41.152]: PS C:\Users\Administrator\Documents>
```

2. 要确定驱动程序是否已安装和链路是否正常工作，请发出以下 Windows PowerShell 命令：

```
[172.28.41.152]: PS C:\Users\Administrator\Documents>  
Get-NetAdapter
```

图 13-19 显示示例输出。

```
[172.28.41.178]: PS C:\Users\Administrator\Documents> Get-NetAdapter
```

Name	InterfaceDescription	ifIndex	Status	MacAddress	LinkSpeed
SLOT 2 4 Port 2	QLogic FastLinQ QL41262-DE 25GbE...#238	6	Up	00-0E-1E-FD-AB-C1	25 Gbps

图 13-19. Windows PowerShell 命令：Get-NetAdapter

3. 要验证在适配器上是否已启用 RDMA，请发出以下 Windows PowerShell 命令：

```
[172.28.41.152]: PS C:\Users\Administrator\Documents>  
Get-NetAdapterRdma
```

图 13-20 显示示例输出。

```
[172.28.41.178]: PS C:\Users\Administrator\Documents> Get-NetAdapterRdma
```

Name	InterfaceDescription	Enabled
SLOT 2 4 Port 2	QLogic FastLinQ QL41262-DE 25GbE Adap...	True
SLOT 2 3 Port 1	QLogic FastLinQ QL41262-DE 25GbE Adap...	True

图 13-20. Windows PowerShell 命令：Get-NetAdapterRdma

4. 要将 IP 地址和 VLAN ID 分配给适配器的所有接口，请发出以下 Windows PowerShell 命令：

```
[172.28.41.152]: PS C:\> Set-NetAdapterAdvancedProperty  
-InterfaceAlias "slot 1 port 1" -RegistryKeyword vlanid  
-RegistryValue 5
```

```
[172.28.41.152]: PS C:\> netsh interface ip set address  
name="SLOT 1 Port 1" static 192.168.10.10 255.255.255.0
```

5. 要在 Nano Server 上创建 SMBShare，请发出以下 Windows PowerShell 命令：

```
[172.28.41.152]: PS C:\Users\Administrator\Documents>  
New-Item -Path c:\-Type Directory -Name smbshare -Verbose
```

图 13-21 显示示例输出。

```
[172.28.41.152]: PS C:\Users\Administrator\Documents> New-Item -Path c:\ -Type Directory -Name smbshare -Verbose
VERBOSE: Performing the operation "Create Directory" on target "Destination: C:\smbshare".

Directory: C:\

Mode                LastWriteTime         Length Name
----                -
d-----          4/25/2016  1:34 AM             smbshare
```

图 13-21. Windows PowerShell 命令: New-Item

```
[172.28.41.152]: PS C:\> New-SMBShare -Name "smbshare" -Path c:\smbshare -FullAccess Everyone
```

图 13-22 显示示例输出。

```
[172.28.41.152]: PS C:\> New-SMBShare -Name "smbshare" -Path c:\smbshare -FullAccess Everyone

Name      ScopeName Path      Description
-----
smbshare *      c:\smbshare
```

图 13-22. Windows PowerShell 命令: New-SMBShare

6. 要将 SMBShare 映射为客户端计算机中的网络驱动器，请发出以下 Windows PowerShell 命令：

**注**

Nano Server 上接口的 IP 地址为 192.168.10.10。

```
PS C:\Windows\system32> net use z: \\192.168.10.10\smbshare
This command completed successfully.
```

7. 要在 SMBShare 上执行读取 / 写入操作，并检查 Nano Sever 上的 RDMA 统计信息，请发出以下 Windows PowerShell 命令：

```
[172.28.41.152]: PS C:\>
(Get-NetAdapterStatistics).RdmaStatistics
```

图 13-23 显示命令输出。

```
[172.28.41.152]: PS C:\> (Get-NetAdapterStatistics).RdmaStatistics  
  
AcceptedConnections      : 2  
ActiveConnections        : 2  
CompletionQueueErrors    : 0  
ConnectionErrors         : 0  
FailedConnectionAttempts : 0  
InboundBytes              : 403913290  
InboundFrames             : 4110373  
InitiatedConnections      : 0  
OutboundBytes             : 63902433706  
OutboundFrames            : 58728133  
PSComputerName           :
```

**图 13-23. Windows PowerShell 命令: Get-NetAdapterStatistics**

# 14 故障排除

本章提供了以下故障排除信息：

- [故障排除核查清单](#)
- [第 246 页上的“验证是否已加载最新驱动程序”](#)
- [第 247 页上的“测试网络连通性”](#)
- [第 248 页上的“使用 Hyper-V 的 Microsoft Virtualization”](#)
- [第 248 页上的“Linux 特定问题”](#)
- [第 248 页上的“其他问题”](#)
- [第 249 页上的“收集调试数据”](#)

## 故障排除核查清单

### 小心

在打开服务器机箱以添加或拆卸适配器之前，请先查阅[第 5 页上的“安全预防措施”](#)。

以下核查清单提供了一些建议的操作，以解决在系统中安装或运行 41xxx 系列适配器时可能遇到的问题。

- 检查所有电缆和连接。验证网络适配器和交换机上的电缆正确连接。
- 对照[第 6 页上的“安装适配器”](#)验证适配器安装。确保适配器正确插入插槽中。检查是否有特定的硬件问题，如插卡组件或 PCI 边缘连接器明显损坏。
- 验证配置设置，如果与其他设备冲突，则进行更改。
- 验证服务器在使用最新的 BIOS。
- 尝试将适配器插入另一插槽。如果在新位置没有问题，则系统中的原插槽可能有缺陷。
- 用已知工作正常的适配器替换故障适配器。如果第二个适配器在第一个适配器无法运行的插槽中可运行，则原适配器可能有缺陷。

- 将适配器安装在另一台运行正常的系统中，然后再次运行测试。如果适配器在新系统中通过测试，则原系统可能有缺陷。
- 从该系统卸下所有其他适配器，然后再次运行测试。如果适配器通过测试，则其他适配器可能导致了争用。

## 验证是否已加载最新驱动程序

确保已为您的 Windows、Linux 或 VMware 系统加载最新驱动程序。

### 在 Windows 中验证驱动程序

请参阅设备管理器，查看有关适配器、链路状态和网络连接的重要信息。

### 在 Linux 中验证驱动程序

要验证 qed.ko 驱动程序是否已正确加载，请发出以下命令：

```
# lsmod | grep -i <module name>
```

如果驱动程序已加载，此命令的输出将显示驱动程序大小（以字节为单位）。以下示例显示为 qed 模块加载的驱动程序：

```
# lsmod | grep -i qed
qed                199238  1
qede               1417947  0
```

如果加载新的驱动程序后重新引导，可以发出以下命令验证当前加载的驱动程序版本是否正确：

```
modinfo qede
```

也可以发出以下命令：

```
[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
```

如果已加载新驱动程序，但尚未重新引导，则 `modinfo` 命令不会显示更新的驱动程序信息。相反地，发出以下 `dmesg` 命令可查看日志。在此示例中，最后一个条目用于标识将在重新引导时激活的驱动程序。

```
# 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
```

## 在 VMware 中验证驱动程序

要验证 VMware ESXi 驱动程序已加载，请发出以下命令：

```
# esxcli software vib list
```

## 测试网络连通性

本节提供在 Windows 和 Linux 环境中测试网络连通性的步骤。

---

### 注

在使用强制链路速度时，验证适配器和交换机均被强制为同一速度。

---

## 测试 Windows 的网络连通性

使用 `ping` 命令测试网络连通性。

**要确定网络连通性是否工作：**

1. 单击 **Start**（开始），然后单击 **Run**（运行）。
2. 在 **Open**（打开）框中，键入 `cmd`，然后单击 **OK**（确定）。
3. 要查看被测试的网络连接，发出以下命令：

```
ipconfig /all
```

4. 发出以下命令，然后按 ENTER。

```
ping <ip_address>
```

显示的 ping 统计信息指示网络连接是否在工作。

## 测试 Linux 的网络连通性

要验证以太网接口正常工作并运行：

1. 要检查以太网接口的状态，请发出 `ifconfig` 命令。
2. 要检查以太网接口的统计信息，请发出 `netstat -i` 命令。

要验证连接是否已建立：

1. 对网络上的 IP 主机执行 Ping 操作。从命令行，发出以下命令：

```
ping <ip_address>
```

2. 按 ENTER 键。

显示的 ping 统计信息指示网络连接是否在工作。

使用操作系统 GUI 工具或 `ethtool` 命令 `ethtool -s ethX speed SSSS` 可以将适配器链路速度强制为 10 Gbps 或 25 Gbps。

## 使用 Hyper-V 的 Microsoft Virtualization

Microsoft Virtualization 是适用于 Windows Server 2012 R2 的虚拟机监控程序虚拟化系统。有关 Hyper-V 的更多信息，请访问：

<https://technet.microsoft.com/en-us/library/Dn282278.aspx>

## Linux 特定问题

- 问题：** 汇编驱动程序源码时出错。
- 解决方案：** Linux 分发版的有些安装没有默认安装开发工具和内核源。汇编驱动程序源码之前，确保所使用的 Linux 分发版的开发工具已安装。

## 其他问题

- 问题：** 41xxx 系列适配器已关闭并出现错误消息，提示适配器上的风扇发生故障。
- 解决方案：** 可有意关闭 41xxx 系列适配器 以防止永久性损坏。联系 Cavium 技术支持以获取帮助。
- 问题：** 在安装有 iSCSI 驱动程序 (qedil) 的 ESXi 环境中，有时 VI 客户端无法访问主机。这是由于 `hostd` 守护程序终止，影响了与 VI 客户端的连接性。
- 解决方案：** 联系 VMware 技术支持

## 收集调试数据

使用 [表 14-1](#) 中的命令收集调试数据。

**表 14-1. 收集调试数据命令**

调试数据	说明
dmesg-T	内核日志
ethtool-d	寄存器转储
sys_info.sh	系统信息；在驱动程序包中提供

# A 适配器 LED

表 A-1 列出用于适配器端口链路和活动状态的 LED 指示灯。

**表 A-1. 适配器端口链路和活动 LED**

端口 LED	LED 外观	网络状态
链路 LED	不亮 持续发亮	无链接（电缆断开） 链接
活动 LED	不亮 闪烁	无端口活动 端口活动

# B 电缆和光学模块

本附录提供有关支持的电缆和光学模块的以下信息：

- 支持的规格
- 第 252 页上的“经测试的电缆和光学模块”
- 第 256 页上的“经测试的交换机”

## 支持的规格

41xxx 系列适配器支持符合 SFF8024 的各种电缆和光学模块。具体的外形因素符合性如下：

- SFP:
  - SFF8472（用于内存映射）
  - SFF8419 或 SFF8431（低速信号和电源）
- 四通道小型可插拔接口 (QSFP):
  - SFF8636（用于内存映射）
  - SFF8679 或 SFF8436（低速信号和电源）
- 光学模块电输入 / 输出、有源铜缆 (ACC) 和有源光缆 (AOC):
  - 10G - SFF8431 限制接口
  - 25G - IEEE802.3by Annex 109B (25GAUI)

## 经测试的电缆和光学模块

Cavium 不保证满足符合性要求的每个电缆或光学模块均可与 41xxx 系列适配器一起使用。Cavium 已测试表 B-1 中列出的组件并提供此列表方便您使用。

表 B-1. 经测试的电缆和光学模块

速度 / 外形因素	制造商	部件号	类型	电缆长度 <sup>a</sup>	标准尺寸
<b>电缆</b>					
10G DAC <sup>b</sup>	Brocade®	1539W	SFP+10G-to-SFP+10G	1	26
		V239T	SFP+10G-to-SFP+10G	3	26
		48V40	SFP+10G-to-SFP+10G	5	26
	Cisco	H606N	SFP+10G-to-SFP+10G	1	26
		K591N	SFP+10G-to-SFP+10G	3	26
		G849N	SFP+10G-to-SFP+10G	5	26
	Dell	V250M	SFP+10G-to-SFP+10G	1	26
		53HVN	SFP+10G-to-SFP+10G	3	26
		358VV	SFP+10G-to-SFP+10G	5	26
		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
25G DAC	Amphenol®	NDCCGF0001	SFP28-25G-to-SFP28-25G	1	30
		NDCCGF0003	SFP28-25G-to-SFP28-25G	3	30
		NDCCGJ0003	SFP28-25G-to-SFP28-25G	3	26
		NDCCGJ0005	SFP28-25G-to-SFP28-25G	5	26
	Dell	2JVDD	SFP28-25G-to-SFP28-25G	1	26
		D0R73	SFP28-25G-to-SFP28-25G	2	26
		OVXFJY	SFP28-25G-to-SFP28-25G	3	26
		9X8JP	SFP28-25G-to-SFP28-25G	5	26

表 B-1. 经测试的电缆和光学模块 (续)

速度 / 外形因素	制造商	部件号	类型	电缆长度 <sup>a</sup>	标准尺寸
40G 铜缆 QSFP 分路器 (4×10G)	Dell	TCPM2	QSFP+40G-to-4xSFP+10G	1	30
		27GG5	QSFP+40G-to-4xSFP+10G	3	30
		P8T4W	QSFP+40G-to-4xSFP+10G	5	26
1G 铜缆 RJ45 收发器	Dell	8T47V	SFP+ to 1G RJ	1G RJ45	N/A
		XK1M7	SFP+ to 1G RJ	1G RJ45	N/A
		XTY28	SFP+ to 1G RJ	1G RJ45	N/A
10G 铜质 RJ45 收发器	Dell	PGYJT	SFP+ to 10G RJ	10G RJ45	N/A
40G DAC 分路器 (4×10G)	Dell	470-AAVO	QSFP+40G-to-4xSFP+10G	1	26
		470-AAHG	QSFP+40G-to-4xSFP+10G	3	26
		470-AAHX	QSFP+40G-to-4xSFP+10G	5	26
100G DAC 分路器 (4×25G)	Amphenol	NDAQGJ-0001	QSFP28-100G-to-4xSFP28-25G	1	26
		NDAQGF-0002	QSFP28-100G-to-4xSFP28-25G	2	30
		NDAQGF-0003	QSFP28-100G-to-4xSFP28-25G	3	30
		NDAQGJ-0005	QSFP28-100G-to-4xSFP28-25G	5	26
	Dell	026FN3 Rev A00	QSFP28-100G-to-4XSFP28-25G	1	26
		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

表 B-1. 经测试的电缆和光学模块 (续)

速度 / 外形因素	制造商	部件号	类型	电缆长度 <sup>a</sup>	标准尺寸
<b>光学解决方案</b>					
10G 光学收发器	Avago®	AFBR-703SMZ	SFP+ SR	N/A	N/A
		AFBR-701SDZ	SFP+ LR	N/A	N/A
	Dell	Y3KJN	SFP+ SR	1G/10G	N/A
		WTRD1	SFP+ SR	10G	N/A
		3G84K	SFP+ SR	10G	N/A
		RN84N	SFP+ SR	10G-LR	N/A
	Finisar®	FTLX8571D3BCL-QL	SFP+ SR	N/A	N/A
		FTLX1471D3BCL-QL	SFP+ LR	N/A	N/A
25G 光学收发器	Dell	P7D7R	SFP28 光学收发器 SR	25G SR	N/A
	Finisar	FTLF8536P4BCL	SFP28 光学收发器 SR	N/A	N/A
		FTLF8538P4BCL	SFP28 光学收发器 SR 无 FEC	N/A	N/A

表 B-1. 经测试的电缆和光学模块 (续)

速度 / 外形因素	制造商	部件号	类型	电缆长度 <sup>a</sup>	标准尺寸
10G AOC <sup>c</sup>	Dell	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
		470-ABMJ	SFP+ AOC	20	N/A
		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
25G AOC	Dell	X5DH4	SFP28 AOC	20	N/A
	InnoLight®	TF-PY003-N00	SFP28 AOC	3	N/A
		TF-PY020-N00	SFP28 AOC	20	N/A

<sup>a</sup> 电缆长度以米为单位。

<sup>b</sup> DAC 是直接连接电缆。

<sup>c</sup> AOC 是有源光学电缆。

## 经测试的交换机

表 B-2 列出经测试与 41xxx 系列适配器的互操作性的交换机。41xxx 系列适配器此列表基于产品发布时可购得的交换机，且随着新交换机进入市场或停产，将随时间更改。

**表 B-2. 经测试互操作性的交换机**

制造商	以太网交换机型号
Arista	
	7160
Cisco	Nexus 3132
	Nexus 3232C
	Nexus 5548
	Nexus 5596T
	Nexus 6000
Dell EMC	
	Z9100
HPE	FlexFabric 5950
Mellanox®	
	SN2700

# C

## Dell Z9100 交换机配置

41xxx 系列适配器支持与 Dell Z9100 以太网交换机连接。但是，在标准化自动协商过程之前，必须将交换机明确地配置为按 25 Gbps 速率连接到适配器。

### 要配置 Dell Z9100 交换机端口以便按 25 Gbps 速率连接 41xxx 系列适配器：

1. 在您的管理工作站与交换机之间建立串行端口连接。
2. 打开命令行会话，然后如下登录到交换机：

```
Login: admin
Password: admin
```

3. 启用交换机端口的配置：

```
Dell> enable
Password: xxxxxxx
Dell# config
```

4. 标识要配置的模块和端口。以下示例使用模块 1，端口 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 ?
10G               Quad port mode with 10G speed
25G               Quad port mode with 25G speed
Dell(conf)#stack-unit 1 port 5 portmode quad speed 25G
```

有关更改适配器链路速度的信息，请参阅第 247 页上的“测试网络连通性”。

5. 验证端口的运行速率为 25 Gbps:

```
Dell# Dell#show running-config | grep "port 5"
stack-unit 1 port 5 portmode quad speed 25G
```

6. 要在交换机端口 5 上禁用自动协商, 请按照以下步骤操作:

- a. 标识交换机端口接口 (模块 1, 端口 5, 接口 1), 并确认自动协商状态:

```
Dell(conf)#interface tw 1/5/1

Dell(conf-if-tf-1/5/1)#intf-type cr4 ?
autoneg                               Enable autoneg
```

- b. 禁用自动协商:

```
Dell(conf-if-tf-1/5/1)#no intf-type cr4 autoneg
```

- c. 验证已禁用自动协商。

```
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
```

有关配置 Dell Z9100 交换机的更多信息, 请参阅 Dell 支持网站上的 *Dell Z9100 Switch Configuration Guide* (Dell Z9100 交换机配置指南), 网址为:

[support.dell.com](http://support.dell.com)

# D 功能约束

本附录提供有关当前版本中实施的功能约束的信息。

这些功能共存约束可能会在未来版本中移除。届时，您应该能够使用功能组合，而无需超出启用功能通常所需的任何附加的配置步骤。

## 不支持 NPAR 模式下同一端口上共存的 FCoE 和 iSCSI

处于 NPAR 模式时，当前版本不支持在属于同一物理端口的 PF 上进行 FCoE 和 iSCSI 配置（仅在默认模式下支持同一端口上共存的 FCoE 和 iSCSI）允许在 NPAR 模式下的物理端口上配置 FCoE 或 iSCSI。

使用 HII 或 QLogic 管理工具在一个端口上配置具有 iSCSI 或 FCoE 个性设置的 PF 后，禁止通过这些管理工具在另一个 PF 上配置存储协议。

由于存储个性设置默认已禁用，因此只有使用 HII 或 QLogic 管理工具配置的个性设置才能写入 NVRAM 配置中。移除此限制后，用户可在 NPAR 模式下的同一端口上配置附加 PF 用于存储。

## 不支持同一端口上共存的 RoCE 和 iWARP

不支持同一端口上的 RoCE 和 iWARP。HII 和 QLogic 管理工具不允许用户同时配置两者。

## 仅在所选 PF 上支持 NIC 和 SAN 引导至库

以太网和 PXE 引导目前仅在 PF0 和 PF1 上受支持。在 NPAR 配置中，其他 PF 不支持以太网和 PXE 引导。

- 当 **Virtualization Mode**（虚拟化模式）设置为 **NPAR** 时，分区 2（PF2 和 PF3）上支持非卸载 FCoE 引导，而分区 3（PF4 和 PF5）上支持 iSCSI 引导。iSCSI 和 FCoE 引导限制为每个引导会话一个目标。iSCSI 引导目标 LUN 支持仅限于 LUN ID 0。
- 当 **Virtualization Mode**（虚拟化模式）设置为 **None**（无）或 **SR-IOV** 时，不支持从 SAN 引导。

# 词汇表

## 传输控制协议

请参阅 [TCP](#)。

## 传输控制协议 / 互联网协议

请参阅 [TCP/IP](#)。

## 串行化器 / 并行化器

请参阅 [SerDes](#)。

## 大量发送卸载

请参阅 [LSO](#)。

## 带宽

以特定传输速率可传输的数据量的度量。1Gbps 或 2Gbps 光纤信道端口可传输或接收的标称速率为 1 或 2Gbps，具体视所连接的设备而定。这分别对应于实际带宽值 106MB 和 212MB。

## 单根输入 / 输出虚拟化

请参阅 [SR-IOV](#)。

## 第二层

指多层通信模型开放系统互连 (OSI) 的数据链路层。数据链路层的功能是跨网络中的物理链路移动数据，其中交换机使用目标 MAC 地址在第二层级别重定向数据消息以确定消息目标。

## 动态主机配置协议

请参阅 [DHCP](#)。

## 非易失性存储器表示

请参阅 [NVMe](#)。

## 非易失性随机存取存储器

请参阅 [NVRAM](#)。

## 服务质量

请参阅 [QoS](#)。

## 高级配置和电源接口

请参阅 [ACPI](#)。

## 互联网广域 RDMA 协议

请参阅 [iWARP](#)。

## 互联网小型计算机系统接口

请参阅 [iSCSI](#)。

## 互联网协议

请参阅 [IP](#)。

## 基本地址寄存器

请参阅 [BAR](#)。

## 基本输入输出系统

请参阅 [BIOS](#)。

## 基于聚合以太网的 RDMA

请参阅 [RoCE](#)。

## 节能以太网

请参阅 [EEE](#)。

## 精简指令集计算机

请参阅 [RISC](#)。

## 巨型帧

在高性能网络中使用以提高长距离性能的大型 IP 帧。对于千兆位以太网，巨型帧通常意味着 9,000 字节，但可指超出 IP MTU（在以太网上为 1,500 字节）的任意大小。

## 可扩展固件接口

请参阅 [EFI](#)。

## 类型长度值

请参阅 [TLV](#)。

## 链路层发现协议

请参阅 [LLDP](#)。

## 目标

SCSI 会话的存储设备端点。启动器从目标请求数据。目标通常为磁盘驱动器、磁带驱动器或其他媒体设备。通常，SCSI 外围设备是目标，但某些情况下，适配器也可能是目标。目标可包含许多 LUN。

目标是响应启动器（主机系统）请求的设备。外围设备是目标，但对于某些命令（例如，SCSI COPY 命令）来说，外围设备可充当启动器。

## 驱动程序

用于充当文件系统与物理数据存储设备或网络介质之间的接口的软件。

## 人机界面基础设施

请参阅 [HII](#)。

## 设备

**目标**，通常为磁盘驱动器。安装或连接到系统的磁盘驱动器、磁带驱动器、打印机或键盘等硬件。在光纤信道中，为 *目标设备*。

## 适配器

用于联接主机系统与目标设备的板卡。适配器的其他表示形式包括主机总线适配器、主机适配器和板。

## 适配器端口

适配器板上的端口。

## 数据中心桥接

请参阅 [DCB](#)。

## 数据中心桥接交换

请参阅 [DCBX](#)。

## 统一可扩展固件接口

请参阅 [UEFI](#)。

## 网络接口卡

请参阅 [NIC](#)。

## 文件传输协议

请参阅 [FTP](#)。

## 消息信号中断

请参阅 [MSI](#)，[MSI-X](#)。

## 小型计算机系统接口

请参阅 [SCSI](#)。

## 虚拟机

请参阅 [VM](#)。

## 虚拟接口

请参阅 [VI](#)。

## 虚拟逻辑区域网络

请参阅 [VLAN](#)。

## 以太网

最为广泛使用的 LAN 技术，用于在计算机之间发送信息，通常速度为每秒 10 和 100 兆比特 (Mbps)。

## 以太网光纤信道

请参阅 [FCoE](#)。

## 用户数据报协议

请参阅 [UDP](#)。

## 远程直接内存访问

请参阅 [RDMA](#)。

## 增强的传输选择

请参阅 [ETS](#)。

## 最大传输单元

请参阅 [MTU](#)。

## ACPI

*高级配置和电源接口 (ACPI) 规格*提供了统一的、以操作系统为中心的设备配置和电源管理的开放标准。ACPI 定义了与平台无关的接口，用于硬件查找、配置、电源管理和监测。该规格以操作系统导向的配置和电源管理（OSPM，一个用于描述实施 ACPI 的系统的术语）为中心，从而移除了旧版固件接口的设备管理责任。

## BAR

基本地址寄存器。用于保留设备所使用的内存地址，或端口地址的偏移。通常，内存地址 BAR 必须位于物理 RAM 中，而 I/O 空间 BAR 可以驻留在任意内存地址（即使超出物理内存范围）。

## BIOS

基本输入输出系统。通常位于 Flash PROM 中，用于充当硬件和操作系统间接口的程序（或公用程序），并且允许在启动时从适配器引导。

## DCB

数据中心桥接。提供对现有 802.1 桥接规范的增强，以满足数据中心内的协议和应用程序的需求。由于现有高性能数据中心通常包含在不同链路层技术上运行的多个应用特定网络（光纤信道用于存储，以太网用于网络管理和 LAN 连接），DCB 允许使用 802.1 桥接来部署聚合网络，从而使所有应用可以在单个物理基础结构上运行。

## DCBX

数据中心桥接交换。DCB 设备用来与直连对等方交换配置信息的协议。该协议也可用于误配置检测和对等方配置。

## DHCP

动态主机配置协议。仅当请求后，允许 IP 网络上的计算机从具有该计算机相关信息的服务器提取其配置。

## eCore

操作系统与硬件和固件之间的层。它是设备特定的，与操作系统无关。当 eCore 代码需要 OS 服务时（例如，对于内存分配，PCI 配置空间访问等），它会调用在 OS 特定的层中实现的抽象 OS 函数。eCore 流可能由硬件驱动（例如，通过中断）或由驱动程序的操作系统特定部分驱动（例如，加载和卸载加载和卸载）。

## EEE

节能以太网计算机网络标准的双绞线和背板以太网家族的一系列增强功能，可在低数据活动期间降低功耗。目的是将功耗降低 50% 或更多，同时保持与现有设备的完全兼容性。电气和电子工程师协会 (IEEE) 通过 IEEE 802.3az 特别工作组制定了该标准。

## EFI

可扩展固件接口。此规范定义操作系统与平台固件之间的软件接口。EFI 取代了在所有 IBM PC 兼容个人计算机中提供的旧版 BIOS 固件接口。

## ETS

增强的传输选择。用于指定传输选择增强的标准，以支持流量类型之间的带宽分配。当某个流量类型中提供的负载不为其分配的带宽时，增强的传输选择允许其他流量类型使用可用带宽。带宽分配优先级与严格优先级共存。ETS 包含管理对象以支持带宽分配。有关更多信息，请参阅：<http://ieee802.org/1/pages/802.1az.html>

## FCoE

以太网光纤信道。一种由 T11 标准机构定义的新技术，通过在第二层以太网帧内封装光纤信道帧，允许传统的光纤信道存储网络流量在以太网链路上传输。有关更多信息，请访问 [www.fcoe.com](http://www.fcoe.com)。

## FTP

文件传输协议。一种标准网络协议，用于将文件通过基于 TCP 的网络（如互联网）从一台主机传送到另一台主机。进行带外固件上传（比带内固件上传速度更快）时需使用 FTP。

## HII

人机界面基础设施。此规范（UEFI 2.1 的一部分）用于管理用户输入、本地化字符串、字体和窗体，允许 OEM 开发用于预引导配置的图形界面。

## IEEE

电气和电子工程师协会。一个旨在促进电气相关技术进步的国际化非营利组织。

## IP

互联网协议。一种通过互联网将数据从一台计算机发送到另一台计算机的方法。IP 指定数据包的格式（也称为 *数据报*）和寻址方案。

## IQN

iSCSI 限定名称。iSCSI 节点名称基于启动器制造商和唯一的设备名称部分。

## iSCSI

互联网小型计算机系统接口。一种将数据封装到 IP 数据包以通过以太网连接发送的协议。

## iSCSI 限定名称

请参阅 [IQN](#)。

## iWARP

互联网广域 RDMA 协议。一种实现 RDMA 的网络协议，用于通过 IP 网络进行有效的数据传输。iWARP 设计用于多种环境，包括 LAN、存储网络、数据中心网络和 WAN。

## LAN 唤醒

请参阅 [WoL](#)。

## LLDP

供应商中立的第 2 层协议，允许网络设备在本地网络上公布其身份和功能。该协议取代了专有协议，如思科发现协议、极端发现协议和 Nortel 发现协议（也称为 SONMP）。

使用 LLDP 收集的信息存储在设备中，可以使用 SNMP 查询。通过抓取主机并查询该数据库，可以发现启用 LLDP 的网络的拓扑。

## LSO

大量发送卸载。LSO 以太网适配器功能，允许 TCP/IP 网络堆栈构建大型（最多 64KB）TCP 消息，然后将其发送给适配器。适配器硬件将消息分段为可通过线路发送的较小数据包（帧）：对于标准以太网帧最多为 1,500 字节，对于巨型以太网帧最多为 9,000 字节。分段过程可释放服务器 CPU，从而不必将大型 TCP 消息分段为装入所支持帧大小的较小数据包。

## MSI, MSI-X

消息信号中断。两个由 PCI 定义的扩展之一，用于支持 PCI 2.2 及更高版本和 PCI Express 中的消息信号中断 (MSI)。MSI 是另一种通过特殊消息来生成中断的方式，可模拟引脚有效或无效置位。

MSI-X（在 PCI 3.0 中定义）允许设备分配介于 1 到 2,048 之间任意数量的中断，并为每个中断提供单独的数据和地址寄存器。MSI 中可选的特征（64 位寻址和中断屏蔽）对 MSI-X 为强制的。

## MTU

最大传输单元。是指指定通信层协议可以传输的最大数据包（IP 数据报）大小（以字节为单位）。

## NIC

网络接口卡。安装以启用专用网络连接的计算机卡。

## NIC 分区

请参阅 [NPAR](#)。

## NPAR

**NIC 分区**。将一个 NIC 端口分成多个物理功能或分区，每个均有用户可配置的带宽和个性设置（接口类型）。个性设置包括 [NIC](#)、[FCoE](#) 和 [iSCSI](#)。

## NVMe

专为销售状态驱动程序 (SSD) 设计的存储访问方法。

## NVRAM

非易失性随机存取存储器。一种在即使断电时仍可保留数据（配置设置）的存储器。您可以手动配置或从文件还原 NVRAM 设置。

## OFED™

OpenFabrics 企业分布。RDMA 和内核旁路应用程序的开源软件。

## PCI™

外围组件接口。一种由 Intel® 引入的 32 位本地总线规格。

## PCI Express (PCIe)

第三代 I/O 标准，除支持旧式外围组件互连 (PCI) 和 PCI 扩展 (PCI-X) 台式机和服务器插槽外，还支持增强型以太网网络性能。

## PF

物理功能。

## QoS

服务质量。在通过虚拟端口传输数据时，此方法用于通过设置优先级和分配带宽来预防瓶颈和确保业务连续性。

## RDMA

远程直接内存访问。此功能允许一个节点通过网络直接写入另一个节点的内存（使用地址和大小语义）。此功能是 [VI](#) 网络的一项重要功能。

## RISC

精简指令集计算机。一种计算机微处理器，它执行较少类型的计算机指令，从而以更高的速度运行。

## RoCE

基于聚合以太网的 RDMA。一种网络协议，允许通过聚合或非聚合以太网网络进行远程直接内存访问 (RDMA)。RoCE 是链路层协议，允许位于同一以太网广播域中的任意两个主机之间进行通信。

## SCSI

小型计算机系统接口。一种用于将硬盘驱动器、CD 驱动器、打印机和扫描仪等设备连接到计算机的高速接口。SCSI 可使用一个控制器连接许多设备。每个设备均可通过 SCSI 控制器总线上单独的标识号进行访问。

## SerDes

串行化器 / 并行化器。一对功能块，通常在高速通信中用于弥补受限的输入 / 输出。这些功能块在串行数据接口和并行接口之间沿每个方向转换数据。

## SR-IOV

单根输入 / 输出虚拟化。一种 PCI SIG 的规格，使单个 PCIe 设备能够显示为多个单独的物理 PCIe 设备。SR-IOV 允许隔离 PCIe 资源以实现性能、互操作性和可管理性。

## TCP

传输控制协议。一组通过互联网协议发送数据包中数据的规则。

## TCP/IP

传输控制协议 / 互联网协议。互联网的基本通信语言。

## TLV

类型长度值。可编码为协议内元素的可选信息。类型和长度字段为固定大小（通常 1-4 字节），而值字段为可变大小。这些字段用法如下：

- 类型 - 表示此部分消息所代表字段类型的数字代码。
- 长度 - 值字段的大小（通常以字节为单位）。
- 值 - 一组包含此部分消息数据的、可变大小的字节。

## UDP

用户数据报协议。一种不保证数据包顺序或传递的无连接传输协议。它直接在 IP 之上工作。

## UEFI

统一可扩展固件接口。此规范详细说明了一个接口，可帮助将预引导环境（即从系统开机之后到操作系统启动之前）的系统控制移交给操作系统（如 Windows 或 Linux）。在引导时，UEFI 提供操作系统与平台固件之间的干净接口，并支持用于初始化添加式卡的独立于体系结构的机制。

## VF

虚拟功能。

## VI

虚拟接口。用于跨光纤信道和其他通信协议进行远程直接内存访问的方案。用于群集和消息传递。

## VLAN

虚拟逻辑区域网络 (LAN)。具有一组通用要求的主机组，如同连接到同一线路一样通信，无论其物理位置如何。尽管 VLAN 具有与物理 LAN 相同的属性，但它允许终端站组合在一起，即使其不在同一 LAN 网段上。VLAN 能够通过软件重新配置网络，而不必物理重新定位设备。

## VM

虚拟机。机器（计算机）的软件实现，如同真实机器一样执行程序。

## WoL

唤醒 LAN。一种以太网计算机联网标准，允许通过发送的网络消息远程开启或唤醒计算机，这些消息通常由在网络中另一台计算机上执行的简单程序发送。



**公司总部** Cavium, Inc. 2315 N. First Street San Jose, CA 95131 408-943-7100

**国际办事处** 英国 | 爱尔兰 | 德国 | 法国 | 印度 | 日本 | 中国 | 中国香港 | 新加坡 | 中国台湾 | 以色列

---

版权所有 © 2017, 2018 Cavium, Inc. 在世界范围内保留所有权利。QLogic Corporation 是 Cavium, Inc. 的全资子公司。Cavium、FastLinQ、QConvergeConsole、QLogic 和 SmartAN 是 Cavium, Inc. 的注册商标。所有其他品牌和产品名称均为其各自所有者的商标或注册商标。

本文档提供的信息仅供参考，并且可能包含错误。Cavium 保留在不另行通知的情况下对本文档或者产品设计或规格进行更改的权利。Cavium 不提供任何明示或暗示的任何担保，不保证您可获得本文档中描述的任何结果或性能。所有关于 Cavium 未来发展方向和意图的声明，均有可能在不另行通知的情况下进行更改或者取消，且仅表示目标与目的。

