

INFINIBAND

RADE ASSOCIATION

t Roce

Demartek RoCE Deployment Guide

Produced with the InfiniBand Trade Association in support of the RoCE Initiative.

Executive Summary

Data center managers and administrators face increasing pressure to accelerate application performance because data continues to grow, and demand for quick access to this data continues to intensify.

The computing industry continues to develop new technologies such as flash storage and more powerful processors to keep pace with these demands. Product design engineers are also incorporating performance acceleration techniques into network adapters and switches.

One of these networking adapter technologies, Remote Direct Memory Access (RDMA) has been built into special network adapters that use RDMA over Converged Ethernet, or RoCE technology. This performance acceleration technology is built into specialized network adapters that work with a streamlined software stack in the operating system to improve performance and lower the host CPU consumption required to process network traffic.

This guide is designed for managers and technical professionals within IT departments who are exploring the benefits of deploying RoCE technology or who are looking for practical guidance and deployment examples of RoCE solutions.

As RoCE includes server, adapter and switching technologies, this guide provides information and guidance in each area. A basic understanding of each of these areas is needed to successfully deploy RoCE technology.

This guide is intended to be used as a reference and includes screen shots and information from actual deployments of specific RoCE-enabled products.

All of the work was performed in the Demartek lab in Golden, Colorado, USA.



RoCE Basic Definitions

RDMA – Remote Direct Memory Access - the remote memory management capability that allows server-toserver data movement directly between the application memory of each without any CPU involvement.

RoCE – (pronounced "Rocky") RDMA over Converged Ethernet – a network protocol that allows remote direct memory access over a converged Ethernet network.

Converged Ethernet – A term associated with a collection of standards-based extensions to traditional Ethernet that provide lossless characteristics that enable the convergence of LAN and SAN technology onto a single unified fabric.

RoCE Initiative

The RoCE Initiative is an education and resource program of the InfiniBand Trade Association (IBTA) that is committed to increasing awareness of RoCE by providing technical education and reference solutions for high-performance Ethernet topologies in traditional and cloud-based data centers.

Additional information is available on the <u>RoCE Initiative</u> <u>website</u> including white papers, solution briefs and contact information.

Use of Bookmarks

This document uses a feature known as "bookmarks" that can be used for navigation within this document. Enable the display of bookmarks in your PDF reader so that you can navigate to any section of this document by clicking on the bookmark entry.

About Demartek

Demartek is a computer industry analyst organization with its own ISO 17025 accredited computer test lab. We provide real-world, hands-on research and analysis by focusing on lab validation testing and power efficiency testing of data center computer equipment such as servers, networking and storage systems. We produce white papers and videos that provide the results of our tests along with our commentary. Additional information regarding our services is available on our lab testing page.

To be notified when new Deployment Guides and lab validation reports become available, you can subscribe to our free monthly newsletter, <u>Demartek Lab Notes</u>, available on our website. We do not give out, rent or sell our email list.



RoCE and RDMA vs. Traditional TCP/IP

Remote Direct Memory Access (RDMA) enables more direct movement of data in and out of a server. RDMA bypasses the normal system software network stack components and the multiple buffer copy operations that they normally perform. This is displayed in the **Standard TCP/IP vs. RDMA Data Movement** diagram below that shows a standard network connection on the left and an RDMA connection on the right.

This elimination of buffer copy operations reduces overall CPU consumption and improves (lowers) the latency of the host software stack, since it uses fewer instructions to complete a data transfer.

RDMA can be used in Ethernet applications by using specialized adapters that support RDMA. The adapters are sometimes known as RNICs or RDMA network interface cards (NICs).

RDMA can be implemented for networking and/or storage applications. Some RDMA adapters may have offload functions built into them. The initiator and the target adapters must use the same type of RDMA technology, such as RoCE v1, RoCE v2, etc.

RDMA can be used for storage traffic such as block protocols like iSCSI and file protocols such as NFS and SMB (formerly known as CIFS). RDMA can improve the latency of storage access and reduce CPU consumption of the host server performing I/O requests, allowing a higher rate of I/O requests to be sustained, or perhaps a smaller server to perform the same rate of I/O requests.

RoCE v1 vs. RoCE v2

The original implementation of RoCE, known as "v1", provided the semantics to allow devices to perform DMA transfers that significantly reduce the CPU activity by eliminating the copy functions as shown in the **Standard TCP/IP vs. RDMA Data Movement** diagram. RoCE v1 does not span across IP subnets. RoCE v2 enables communication across IP subnets.

RoCE v1	RoCE v2	
Applications	Applications	
RDMA API	RDMA API	
IB Transport	IB Transport	
IB Network	UDP/IP	
Ethernet	Ethernet	

A more detailed technical explanation of the difference between RoCE v1 and v2 is available at <u>https://community.mellanox.com/docs/DOC-1451</u>

DCB Switch

To achieve low latency, RoCE adapters work best with Ethernet switches that support Data Center Bridging (DCB) that provide lossless characteristics and support priority flow control (PFC). DCB Ethernet switches have become more common over the last few years.

RoCE can also work without a DCB switch because the RoCE specification defines how to use standard IP explicit congestion notification (ECN) and a congestion notification packet response (CNP) which provides the mechanisms needed for congestion management without requiring PFC.



Oemartek



RoCE Applications

There are a variety of applications that can take advantage of RoCE technology, including the three described below. Typically, RoCE is available for 10Gbps and faster technologies. Increased throughput, reduced CPU consumption and low latency are among the benefits available.

Windows: SMB Direct & Storage Spaces Direct

Beginning with Windows Server 2012, Microsoft has included a feature known as SMB Direct that supports the use of RDMA-capable network adapters, such as RoCE network adapters. SMB Direct uses SMBv3 (SMB3) and these adapters to enable file transfer operations to operate at high speed and very low latency, improving the overall performance of a file server. SMB Direct is enabled automatically if RDMA-capable adapters are present.

Beginning with Windows Server 2016, Storage Spaces Direct takes advantage of RDMA-capable network adapters to provide highly scalable software-defined storage. Storage Spaces Direct can be deployed on physical machines or virtual machines with Hyper-V and can be deployed in converged or hyper-converged configurations.

Linux: iSER and NFS over RDMA

For Linux environments, iSCSI Extensions for RDMA (iSER) allows iSCSI (block storage) traffic to take advantage of RDMA-capable network adapters. Similarly, NFS over RDMA allows file traffic to take advantage of RDMA-capable network adapters.

NVMe over Fabrics

NVM Express over Fabrics (NVMe-oF[™]) defines a common architecture that supports a range of storage networking fabrics for NVMe block storage protocol over a storage networking fabric, including RDMA fabrics. RoCE is supported for NVMe-oF.

Products Tested in this Guide

We used several different products in various configurations to product this *RoCE Deployment Guide*. These include the following:

- > HPE 3PAR StoreServ File Controller v3 and HPE 3PAR 8450 all-flash storage array
- > Mellanox SN2410 25GbE/100GbE switch
- > Mellanox ConnectX-4 25GbE and 100GbE adapters
- > QLogic QLE45212 25GbE adapters

HPE

The HPE 3PAR StoreServ File Controller v3 provides file services from any model of HPE 3PAR StoreServ storage. This was connected to an HPE 3PAR 8450 all-flash storage array.

The HPE target consists of an Apollo 2600 chassis with two HPE ProLiant XL190r Gen9 server nodes. It is designed to provide a RoCE-capable file server front end for our HPE 3PAR StoreServ 8450 block storage system. SMB 3.0 over TCP or SMB Direct over RoCE can both be used; the setup is similar for both. Each server node is configured with Windows Storage Server 2012 R2 and contains an Emulex 16GbFC HBA and an HPE Ethernet 10/25Gb 2-port 640FLR-SFP28 Adapter. This adapter uses the Mellanox ConnectX-4 Lx controller and provides 2 ports of 25GbE for each server node.

Mellanox

The Mellanox SN2410 is an Open Ethernet top-of-rack switch with 48-ports of 25GbE and 8-ports of 100GbE. The Mellanox ConnectX-4 Lx EN network controller supports 1/10/25/40/50GbE. The Mellanox ConnectX-4 EN network controller supports 100GbE.

QLogic

The QLogic® FastLinQ[™] QL45212 dual-port Intelligent Ethernet Adapter supports 25GbE.



Deploying RoCE on a Network Switch

The user will need to appropriately configure the switch, configure the adapters according to vendor-specific instructions, and then follow the appropriate application setup.

General considerations for a RoCE deployment are as follows:

- > Standard MTU: While some concurrently running TCP/IP applications may benefit from Jumbo frames, and Jumbo Frames may be used with RoCE, it does increase the latency and is not necessary. Keeping the MTU at 1500 is better in most cases.
- > PFC (Priority Flow Control): Unlike TCP/IP that uses dropped packets as a form of implicit congestion notification; RoCE does not have its own congestion management and requires a lossless network to function properly when there is network congestion. While a lossless network may be achieved most simplistically by enabling basic flow control (global pause), standard practice is to implement Priority Flow Control and QoS with VLAN. These settings need to be implemented on the adapter, in the OS, and on the switch.
- Switch settings must match the server and driver settings. (We used VLAN 2 and priority 3 for RoCE Traffic in our RoCE examples because we are using a Mellanox switch. Any other switch would most likely have Priority 3 set to FCoE, and we would be unable to change it. In these cases, Priority 5 would be recommended.)

A typical way to achieve this setup on a Mellanox switch is as follows:

(Default username and password for a new Mellanox switch is admin,admin).

- Gain access to the administrator command set: #enable #config t
- Display the current configuration: #show running-config

- Display an interface:
 #show interfaces ethernet 1/1
- 4. Make sure flow control is disabled (often it is best to revert to the initial configuration):#config switch-to initial

🔇 Demai

January 201

- Enable PFC priority 3 on the switch:
 #dcb priority-flow-control enable force
 #dcb priority-flow-control priority 3 enable
- Enable PFC on each port: #interface ethernet 1/1 dcb priority-flow-control mode on force #interface ethernet 1/2 dcb priority-flow-control mode on force
- Configure VLAN on switch: #vlan 2
- Configure VLAN on each port: #interface ethernet 1/1 switchport trunk allowed-vlan 2 #interface ethernet 1/2 switchport trunk allowed-vlan 2
- 9. Configure lossless (RoCE) and lossy (all other traffic) buffers for each port: #interface 1/1 ingress-buffer iPort.pg0 map pool iPool0 type lossy reserved 20480 shared alpha 8 #interface ethernet 1/1 ingress-buffer iPort.pg3 map pool iPool0 type lossless reserved 90000 xoff 17000 xon 17000 shared alpha 2 #interface ethernet 1/1 ingress-buffer iPort.pg3 bind switch-priority 3 #interface ethernet 1/1 egress-buffer ePort.tc3 map pool ePool0 reserved 4096 shared alpha inf #interface ethernet 1/2 ingress-buffer iPort.pg0 map pool iPool0 type lossy reserved 20480 shared alpha 8 ...Each port gets 4 commands total to configure the buffers.
- 10. Write the configuration to come back to later: #config write to RoCEPFCwVLAN

- 11. Switch back to it at any time if you make changes:#config switch-to RoCEPFCwVLAN
- 12. Double-check PFC with the command: #show dcb priority-flow-control

Most of the time, this configuration change will require a reboot. If PFC is not enabled but all configuration changes are complete, a manual reboot may be needed after the configuration is written.

Windows: Deploying SMB Direct over RoCE with Storage Spaces

The user will need to appropriately configure the switch, configure the adapters according to vendor-specific instructions, and then follow the appropriate application setup.

TCP/IP connection setup should be completed before proceeding to adapter and application setup steps. The steps to do so are not outlined here in this guide.

Adapter and application setup will need to be completed twice, once on the client and once on the server.

*Please note that the HPE File Controller uses Mellanox adapters.

Mellanox

To, enable PFC and QoS within the Windows Server 2012 OS, the Data Center Bridging Feature is installed using the Add Roles and Features Wizard. This will enable the following PowerShell commands which put RoCE on PFC Priority 3:

>Remove-NetQosTrafficClass
>Remove-NetQosPolicy -Confirm:\$False
>Set-NetQosDcbxSetting -Willing 0
New-NetQosPolicy "SMB" -PolicyStore
ActiveStore -NetDirectPortMatchCondition 445 PriorityValue8021Action 3
>New-NetQosPolicy "DEFAULT" -PolicyStore
Activestore -Default -PriorityValue8021Action 3
New-NetQosPolicy "TCP" -PolicyStore
ActiveStore -IPProtocolMatchCondition TCP PriorityValue8021Action 0

>New-NetQosPolicy "UDP" -PolicyStore ActiveStore -IPProtocolMatchCondition UDP -PriorityValue8021Action 0

🗇 Demai

January 201

>Disable-NetQosFlowControl 0,1,2,4,5,6,7

>Enable-NetQosFlowControl -Priority 3

>Enable-NetAdapterQos -InterfaceAlias "<ethernet interface>"

The adapter settings were as follows:

Name	Value	
Encapsulated Task Offload	Enabled	
Flow Control	Rx & Tx Enabled	
Interrupt Moderation	Enabled	
IPV4 Checksum Offload	Rx & Tx Enabled	
Jumbo Packet	1500	
Large Send Offload V2 (IPv4)	Enabled	
Large Send Offload V2 (IPv6)	Enabled	
Maximum number of RSS Proce	20	
NetworkDirect Functionality Enabled		
Preferred NUMA Node	Default Settings	
Maximum Number of RSS Queues	20	
Priority & Vlan Tag	Priority & VLAN Enabled	
Quality Of Service	Enabled	
Receive Buffers	4096	
Recv Segment Coalescing (IPv4)	Enabled	
Recv Segment Coalescing (IPv6)	Enabled	
Receive Side Scaling	Enabled	
RSS Base Processor Number	0	
RSS Maximum Processor Number	40	
RSS Load Balancing Profile	ClosestProcessor	
SR-IOV	Disabled	
TCP/UDP Checksum Offload (l	Rx & Tx Enabled	
TCP/UDP Checksum Offload (l	Rx & Tx Enabled	
Send Buffers	4096	
Virtual Machine Queues	Disabled	



VMQ VLAN Filtering	Disabled
Network Address	
Receive Completion Method	Adaptive
R/RoCE Max Frame Size	Auto
Rx Interrupt Moderation Type	Adaptive
Rx Interrupt Moderation Pro	Aggressive
Tx Interrupt Moderation Pro	Aggressive
VLAN ID	2

QLogic

In order to use RoCE with QLogic cards, the DCBX and priority number are selected in the UEFI HII interface. Here is an example of this on a Dell server:

1. When the menu appears at boot in the top right, press F2 to enter System Setup.



2. Wait while normal boot procedures are completed.



3. Select Device Settings at the System Setup Main Menu.



Demart

January 2017

4. Choose the QLogic Device Port to configure (Here there are two ports to be configured separately, both shown as "NIC in Slot 4...").

System Setup		
Device Settings		
RAD Controller in Slot 1: Dell PERC <per< th=""><th>CH730P Adapter> Configuration Utility</th><th></th></per<>	CH730P Adapter> Configuration Utility	
Integrated NIC 1 Port 1: Intel(R) Gigabit 4P I	350-t rNDC - B8;CA:3A;6A:24;8C	
Integrated NIC 1 Port 2: Intel(R) Gigabit 4P	350-t rNDC - BRCA(3A)6A(24)8D	
Integrated NIC 1 Port 3: Intel(R) Gigabit 4P	350-t rNDC - B8:CA:3A:6A:24:8E	
Integrated NIC 1 Port 4: Intel(R) Gigabit 4P	350-t rNDC - B8:CA:3A:6A:24:8F	
NIC in Slot 4 Port 't QLogic FastLinQ QL48	212H 25GbE Adapter - 000E:1E:C4:CD:A6	
NIC in Slot 4 Port 2: QLogic FastLinQ QL4	212H 25GbE Adapter - 00:0E:1E:C4:CD:A7	
PCle Sold State Drive in Slot 0 in Bay 2: G	eneric NVMe PCIe SSD Configuration Data	
Unknown in Slot 3: Generic NVMe PCle SS	D Configuration Data	
Unknown in Slot 5: Generic NVMe PCIe SS	D Configuration Data	
Manage RAD Controller Configurati	115,	
•		

5. Select Device Level Configuration.

air Corniguration Page	
Firmware Image Properties	
Device Level Configuration	
Partitioning Mode	Befault O NPAR
Device Name	
Chip Type	BCM57980S B0
PCI Device ID	1056
PCI Address	2100:00
Link Status	Disconnected
Permanent MAC Address	00:0E:1E:C4:CD:A6
Virtual MAC Address	005E:1E:C4:CD:A6

6. Select Dynamic DCBX protocol and the appropriate RoCE Priority. Here Priority 3 is used, which is permitted with a Mellanox switch.



lain Configuration Page • Device	Level Configuration	
Link Speed	Auto Negroteted C Endeled @ Deabled O UNDI @ Deabled Dynamic 3	

Other manufacturers should have similar UEFI hook-ins for the QLogic cards allowing the DCBX option and traffic priority class to be set. If any difficulty is encountered using the HII interface on a server, contact QLogic technical support for an engineering solution.

The adapter settings were as follows:

Name	Value	
Encapsulated Task Offload	Enabled	
NVGRE Encapsulated Task	Enabled	
Off		
VXLAN Encapsulated Task Off	Enabled	
Flow Control	Rx & Tx Enabled	
Interrupt Moderation	Enabled	
Jumbo Packet	1514	
Large Send Offload V2 (IPv4)	Enabled	
Large Send Offload V2 (IPv6)	Enabled	
NetworkDirect	Enabled	
Functionality		
NUMA node ID	65535	
Maximum Number of RSS	16	
Queues	Priority & VI AN	
Priority & VLAN	enabled	
Receive Buffers (0=Auto)	3000	
Recv Segment Coalescing	Enabled	
(IPv4)		
Recv Segment Coalescing (IPv6)	Enabled	
Receive Side Scaling	Enabled	

NUMAScalingStatic	
25 Gbps Full Duplex	
Disabled	
Py & Ty Enabled	
Ry & Ty Enabled	
5000	
Disabled	
1780	
4709	
1024	
Preboot controlled	
RoCE v2	
2	

Demart

January 2017

Notice the option to specify the RoCE mode as RoCE v2 or RoCE v1. RoCE v2 was used for all our SMB Direct testing.



> Network Direct MTU size for QLogic may be increased to 4096 for possible performance gains, as long as the Ethernet MTU size is larger than 4096, and a 4096 MTU is configured on the switch, client, and server.

Application Setup

When deploying SMB Direct over RoCE, perform the following steps after switch configuration, OS install, cluster configuration (if needed, as in in the case of the HPE File Controller) and adapter configuration:

- Make sure the global setting for SMB Direct is enabled using the PowerShell command: \>get-NetOffloadGlobalSetting
- 2. Make sure multichannel is enabled. (While this feature was enabled but not leveraged in our testing due to our simple setup, it should be used when there are, physically, multiple channels available)

\> get-SMBClientConfiguration | Select Enable
Multichannel

OR

\> get-SMBServerConfiguration | Select Enable
Multichannel

 Verify the network interface is listed as being RDMA capable: \>get-SMBClientNetworkInterface

OR \>get-SMBServerNetworkInterface AND

\>get-NetAdapterRDMA

- 4. Create the storage server.
 - > HPE File Controller:
 - In the 3PAR management software, a small volume should be configured for Quorum and exported to both server nodes. For our testing, 1GB was sufficient. A larger volume should also be configured and exported to both nodes, so that the cluster configuration wizard will automatically configure a file server role on the cluster and make the larger volume available for shares.
 - > Following the ICT will configure the failover cluster with file server role for the user. Once complete, the user needs to go to Server Manager, File and Storage Services, Shares, and configure a share on the large volume.

- > If an active/active configuration is desired, add a second file server role to run on the second node. First create a new large LUN exported to both nodes. Then rescan the storage in Server Manager, File and Storage Services, Volumes, Disks. Create a new volume on the new LUN. Then go to the Server Manager Tools Menu to find the Failover Clustering. Go to roles and create a new file server role, with a new IP address and the newly created volume for storage. After the role is created, assign the role to the other node, if it hasn't been already. After this, a share can be created on the volume as was done with the previously created file server role.
- > Generic Servers: Install the desired block device(s) locally in the target server and use Storage Spaces to create Virtual Drives containing File Shares.
- If all is configured correctly, there should be two SMB Direct Listeners waiting for connections on both client and server: \>netstat -xan 1 Once a connection to the remote file share is initiated and traffic has started, netstat will show the connections.
- 6. To view RDMA connection attempts, active connections, and bytes or frames sent, open the Perfmon program, right click to add counters, and select RDMA.

While not employed in our own testing, the following may be explored to potentially enhance performance of SMB Direct:

> Increase the number of SMB Direct RDMA connections from the default (2) to 8 or 16 using the PowerShell command \>Set-ItemProperty -Path "HKLM:\SYSTEM\CurrentControlSet\Services\L anmanWorkstation\Parameters" ConnectionCountPerRdmaNetworkInterface -Type DWORD -Value <number of connections> -Force

It should be noted that there will not be a key present in the registry if the default number of connections are being used; it will only be present once the value is changed. It should also be noted that the PowerShell command sets the value immediately, while using regedit necessitates a reboot. This option will have to be set on both the target server and on the client.

> If RoCE v2 is not necessary for routing purposes, RoCE v1 may give better performance.

Linux: Deploying iSER over RoCE

The user will need to appropriately configure the switch, configure the adapters according to vendor-specific instructions, and then follow the appropriate application setup.

TCP/IP connection setup is should be completed before proceeding to adapter and application setup steps. The steps to do so are not outlined here in this guide.

Adapter and application setup will need to be completed twice, once on the initiator and once on the target.

Mellanox

Using iSER with RoCE on Red Hat Enterprise Linux requires the use of InfiniBand software and tools, called OFED. While Red Hat Enterprise Linux does provide inbox InfiniBand, Mellanox also provides its own OFED. While Mellanox can use the inbox OFED, if Mellanox is the only brand of NIC using RDMA on a server, it is best to download and install the Mellanox OFED. Please note that the Mellanox OFED should not be installed on a server alongside the generic OFED and cannot be used with non-Mellanox NICs.

- Download the Mellanox OFED. 3.3 is available from <u>http://www.mellanox.com/page/products_dyn?p</u> <u>roduct_family=26&mtag=linux_sw_drivers</u>
- Install prerequisite packages.
 #yum install cpp-4.8.5-4.el7.x86_64.rpm gcc-4.8.5-4.el7.x86_64.rpm gcc-gfortran-4.8.5-4.el7.x86_64.rpm glibc-2.17-

106.el7_2.8.x86_64.rpm glibc-common-2.17-106.el7_2.8.x86_64.rpm glibc-devel-2.17-106.el7_2.8.x86_64.rpm glibc-headers-2.17-106.el7_2.8.x86_64.rpm kernel-headers-3.10.0-327.36.1.el7.x86_64.rpm libgfortran-4.8.5-4.el7.x86_64.rpm libmpc-1.0.1-3.el7.x86_64.rpm libquadmath-4.8.5-4.el7.x86_64.rpm libquadmath-devel-4.8.5-4.el7.x86_64.rpm mpfr-3.1.1-4.el7.x86_64.rpm

Demar

January 201

- Install Mellanox OFED.
 #tar -xvf MLNX_OFED_LINUX-3.3-1.0.4.0-rhel7.2x86_64.tgz
 #cd MLNX_OFED_LINUX-3.3-1.0.4.0-rhel7.2x86_64
 #./mInxofedinstall
 It may warn of additional missing packages.
 These should be installed with yum and then the Mellanox OFED install should be attempted again.
- 4. Check if the drivers are loaded correctly.
 #ibdev2netdev
 If the drivers are loaded correctly, this should list the name of the Mellanox interface.
 If nothing was present:

#/etc/init.d/openibd restart (It may be necessary to unload several other modules in order to run this command:

#rmmod ib_isert
#rmmod xprtrdma
#rmmod ib_srpt
#/etc/init.d/openibd restart
#ibdev2netdev)

- 5. Configure PFC.
 - > Display the current priorities. This should show that there is no priority set. #mlnx_qos -i <interface name>
 - > Set the desired priority. Each number corresponds to a priority, 0 through 7. Setting the 4th number to 1 enables the 4th priority in the sequence, which is 3 (0,1,2,3,4,5,6,7). #mlnx_qos -i <interface name> 0,0,0,1,0,0,0,0
 - > Set the egress priority mappings.



There is supposed to be a way to set the egress mapping by adding a line to the network script. However, in our case it didn't work and we used alternate means. Unfortunately, this set of commands will have to be re-entered each time the server or network service is restarted. It may be useful to put them in a script.

#ip link set <VLAN interface> type vlan egress
 0:<RoCE PFC priority>

#ip link set <VLAN interface> type vlan egress
1:<RoCE PFC priority>

... repeat for all priorities up to 15. Only priorities 0 and 1 are shown above.

<pre>#ip link set p9p1.2 #ip link set p9p1.2</pre>	type vlan egress 0:3 type vlan egress 1:3 type vlan egress 2:3 type vlan egress 3:3 type vlan egress 4:3 type vlan egress 5:3 type vlan egress 6:3 type vlan egress 7:3 type vlan egress 9:3 type vlan egress 10:3 type vlan egress 11:3 type vlan egress 12:3 type vlan egress 13:3 type vlan egress 13:3 type vlan egress 14:3 type vlan egress 15:3
--	---

> Check the egress commands worked.

> Optional: Check the GIDs with the Mellanox provided script show-gids (below). This script can be obtained from the Mellanox website: <u>https://community.mellanox.com/docs/DOC-2421</u>

Index 7 is running RoCE v2 on VLAN 2 interface p9p1.2, which is what is needed.

#./sho	ow-gids				
DEV	PORT	INDEX	GID	IPv4	VER DEV
mlx5_() 1	0	fe80:0000:0000:0000:7efe:90ff:fe22:a00a		V1 p9p1
mlx5_() 1	1	fe80:0000:0000:0000:7efe:90ff:fe22:a00a		V2 p9p1
mlx5_0) 1	2	fe80:0000:0000:0000:7efe:90ff:fe22:a00a		V1.5 p9p1
mlx5_0) 1	3	0000:0000:0000:0000:0000:ffff:0a00:01fc	10.0.1.252	V1 p9p1
mlx5_0) 1	4	0000:0000:0000:0000:0000:ffff:0a00:01fc	10.0.1.252	V2 p9p1
mlx5_0) 1	5	0000:0000:0000:0000:0000:ffff:0a00:01fc	10.0.1.252	V1.5 p9p1
mlx5_0) 1	6	0000:0000:0000:0000:0000:ffff:0a00:01fc	10.0.1.252	V1 p9p1.2
m]x5_() 1	7	0000:0000:0000:0000:0000:ffff:0a00:01fc	10.0.1.252	V2 p9p1.2
mlx5_0) 1	8	0000:0000:0000:0000:0000:ffff:0a00:01fc	10.0.1.252	V1.5 p9p1.2
n_gids	s_found=9	9			



> Check everything with tc_wrap.py. Make sure everything for VLAN 2 is priority 3.

#tc_wrap.py -i p9p1	
skprio: 0	
skprio: 9	
skprio: 10	
skprio: 11 skprio: 12	
skprio: 13	
skprio: 14 skprio: 15	
priority 1	
skprio: 1	
skprio: 2 (tos: 8)	
priority 3	
skprio: 3 skprio: 0 (vlan 2)	
skprio: 1 (vlan 2)	
skprio: 2 (vlan 2 tos: 8)	
skprio: 4 (vlan 2 tos: 24)	
skprio: 5 (vlan 2)	
skprio: 6 (Vian 2 tos: 16) skprio: 7 (vian 2)	
skprio: 8 (vlan 2)	
skprio: 9 (vlan 2)	
skprio: 11 (vlan 2)	
skprio: 12 (vlan 2)	
skprio: 14 (vlan 2)	
skprio: 15 (vlan 2)	
priority 4 skprio: 4 (tos: 24)	
priority 5	
skprio: 5 priority 6	
skprio: 6 (tos: 16)	
priority 7	

Demartek

January 2017

QLogic

In order to use RoCE with QLogic cards, the DCBX and priority number are selected in the UEFI HII interface. Here is an example of this on a Dell server:

1. When the menu appears at boot in the top right, press F2 to enter System Setup.



2. Wait while normal boot procedures are completed.



3. Select Device Settings at the System Setup Main Menu.

System Soluc		Hip About Delt
System Setup		
System Setup Main Menu		
System BIOS		
DRAC Settings		
Device Settings		
Select to configure system BIOS settings		
PowerEdge RBP0	Arrow keys and Enter to select	The late
Service Tag : GHI6CZ1	Esc to exit page, Tab to change focus	Finish

4. Choose the QLogic Device Port to configure (Here there are two ports to be configured separately, both shown as "NIC in Slot 4...").



5. Select Device Level Configuration.

an comga acor Fage	
irmware Image Properties	
Device Level Configuration	
Partitioning Mode	Default NPAR
Device Name	QLogic FastLinQ QL45212H 25GbE Adapter
Chip Type	BCM57980S B0
PCI Device ID	1656
PCI Address	2100:00
Link Status	Disconnected
Permanent MAC Address	000E:E:C4:CD:A6
Virtual MAC Address	00:0E:1E:C4:CD:46

6. Select Dynamic DCBX protocol and the appropriate RoCE Priority. Here Priority 3 is used, which is permitted with a Mellanox switch.

Vain Configuration Page • Dev	vice Level Configuration	
Link Speed	Auto Negotiated	
SR40V	 Enabled	
Boot Mode	O UNDI	
DCBX Protocol	Dynamic	
RoCE Priority	3	

Other manufacturers should have similar UEFI hook-ins for the QLogic cards allowing the DCBX option and traffic priority class to be set. If any difficulty is encountered using the HII interface on a server, contact QLogic technical support for an engineering solution.

Demartek January 2017

While in system setup, configure the following in the System BIOS for optimal performance later:

- 1. Disable hyperthreading.
- 2. Set the power management to performance, and disable any Energy Performance Tuning or Energy Efficient Turbo settings.
- 3. Set Package C State Limit C0/C1 state.
- 4. Disable CPU C3 report and CPU C6 report and Enhanced Halt State (C1E).
- 5. Disable ACPI T-States.

Once the system has booted:

 Install an InfiniBand OFED. There are two choices: either use the Inbox OFED provided by Red Hat Enterprise Linux 7.2 or to compile and install the OFED provided at the OpenFabrics website. QLogic recommends using the inbox OFED:

yum -y groupinstall "InfiniBand Support"

Install extra packages for configuration and testing:

#yum -y install perftest infiniband-diags

- Install the RDMA drivers. The NIC drivers should already have been installed when TCP/IP was set-up. Go back to the same driver install directory. #cd fastling-minor-8.10.4.0/libgedr-8.10.3.0
- 4. The readme contains the appropriate configure line to use for Red Hat systems. Compile the RDMA driver.

#./configure -prefix=/usr libdir=\${exec_prefix}/lib64 -sysconfdir=/etc
#make install

- Check RDMA drivers. There should already be qed and qede. The RDMA driver is qedr, and may need to be loaded manually. #lsmod | grep qed lf necessary
 - #modprobe gedr
- 6. Check to make sure everything is working. ibstat should list all QLogic ports available for RDMA.

#ibstat #service rdma status

7. Make sure the CPUs are set to performance mode. Check each processor.

#cat

/sys/devices/system/cpu/cpu0/cpufreq/scaling_g overnor

#cat

...

/sys/devices/system/cpu/cpu1/cpufreq/scaling_g overnor

... It chould road by

It should read performance for all. If any processors are in powersave, manually change each of them to performance.

echo -n performance >

/sys/devices/system/cpu/cpu0/cpufreq/scaling_g overnor

echo -n performance >

/sys/devices/system/cpu/cpu1/cpufreq/scaling_g overnor

8. Check memory limits.#ulimit –lIf there is a limit, this must be rectified:

> Add the following lines to

/etc/security/limits.conf:

- * soft memlock unlimited
 * hard memlock unlimited
- > Then log off and log back in again for changes to take effect.
- > Check the status again. #ulimit –l Should give "unlimited"
- 9. Disable throttling.
 - > in /etc/default/grub, find the end of the line "GRUB_CMDLINE_LINUX=" and add "intel_idle.max_cstate=0 processor.max_cstate=1 mce=ignore_ce idle=poll"
 - > #grub2-mkconfig -o /boot/grub2/grub.cfg
- 10. Reboot. Do checks starting at step 5 again.
- 11. Disable firewall.

January 2017

#systemctl disable firewalld

12. Check InfiniBand.
On the target:
#ib_write_bw -p 12500 -x 0 -R -a
On the initiator:
#ib_write_bw <target address> -p 12500 -x 0 -R -a
-a

Application Setup

When deploying iSER over RoCE, perform the following steps after switch configuration, OS install, and adapter configuration:

- On the initiator system, install the iscsi-initiatorutils if it was not installed with the OS. #yum install iscsi-initiator-utils
- 2. Find out the initiator IQN for future reference. #cat /etc/iscsi/initiatorname.iscsi

#cat /etc/iscsi/initiatorname.iscsi
iqn.1994-05.com.Red Hat:2ffa3543d275

- 3. Set up the target system.
 - > install targetcli if it was not installed with the OS.

#yum install targetcli

- > List the block devices and take note for future reference. #lsblk
- > Run targetcli. At the present time, the ls command should give us an empty tree as pictured below. As the target is configured, the tree will fill out. (If the tree is not empty, "\>clearconfig confirm=true") will clear it. #targetcli

\>ls

> Create backstores out of the block devices. (In our case we used NVMe drives.)

\>cd backstores/block

- \>create name=<backstorename1>
 dev=/dev/<device name>
- \>create name=<backstorename2>
 dev=/dev/<device name>

... Repeat for additional backstores.

```
\>cd backstores/block
/backstores/block> create name=nvme0
dev=/dev/nvme0n1
 Created block storage object nvme0 using
/dev/nvme0n1
/backstores/block> create name=nvme1
dev=/dev/nvme1n1
Created block storage object nvme1 using
/dev/nvme1n1
/backstores/block> create name=nvme2
dev=/dev/nvme2n1
 Created block storage object nvme2 using
/dev/nvme2n1.
/backstores/block> create name=nvme3
dev=/dev/nvme3n1
 Created block storage object nvme3 using
/dev/nvme3n1
/backstores/block> create name=nvme4
dev=/dev/nvme4n1
  Created block storage object nvme4 using
/dev/nvme4n1.
```

> Create a target and give it LUNs created from the previously created backstores.

<>cd /iscsi

\>create

\>cd <new target created >/tpg1/luns
\>create /backstores/block/<backstorename1>
\>create /backstores/block/<backstorename2>

```
/backstores/block> cd /iscsi
/iscsi> create
/iscsi> cd iqn.2003-01.org.linux-
iscsi.dmrtk-srvr-
d.x8664:sn.cdf664e56d8c/tpg1/luns
/iscsi/iqn.20...d8c/tpg1/luns> create
/backstores/block/nvme0
/iscsi/iqn.20...d8c/tpg1/luns> create
/backstores/block/nvme1
Created LUN 1.
/iscsi/iqn.20...d8c/tpg1/luns> create
/backstores/block/nvme2
Created LUN 2.
/iscsi/iqn.20...d8c/tpg1/luns> create
/backstores/block/nvme3
Created LUN 3.
/iscsi/iqn.20...d8c/tpg1/luns> create
/backstores/block/nvme4
Created LUN 4.
```

> Designate an IP address from the target server as a portal for the initiators to reach the



target. Always designate port 3260 as this is the default port for iSCSI. >cd ../portals
>delete 0.0.0.0 3260
>create <target portal IP> 3260

/iscsi/iqn.20...d8c/tpg1/luns> cd ../portals /iscsi/iqn.20.../tpg1/portals> delete 0.0.0.0 3260 Deleted network portal 0.0.0.0:3260 /iscsi/iqn.20.../tpg1/portals> create 10.0.8.252 3260 Using default IP port 3260 Created network portal 10.0.8.252:3260. > Enable iSER for each portal.

\>cd <target portal IP>:3260

\>enable_iser true
/iscsi/iqn.20.../tpg1/portals> cd
10.0.8.252:3260
/iscsi/iqn.20...0.8.252:3260 >
enable_iser true
iSER enable now: True

> Designate which initiators can connect to this target, using the initiator name From Step 2. \>cd ../../acls

\>create <initiator name>

/iscsi/iqn.20....0.8.254:3260> cd ../../acls /iscsi/iqn.20...d8c/tpg1/acls> create iqn.1994-05.com.Red Hat:2ffa3543d275 Created Node ACL for iqn.1994-05.com.Red Hat:2ffa3543d275 Created mapped LUN 4. Created mapped LUN 3. Created mapped LUN 2. Created mapped LUN 1. Created mapped LUN 1. Created mapped LUN 0.



Check the configuration using ls to view the whole tree, then save and exit (as shown below).
 \>ls
 \>saveconfig
 \>exit

/> ls
· · · · · · · · · · · · · · · · · · ·
o- backstores
o- block
/dev/pyme0n1 (1.5TiB) write-thru activated
[/dev/nyme1n1 (1 5TiB) write-thru activated]
[/dev/nyme2n1 (1 5TiB) write-thru activated]
[/dev/nyme3n1 (1 5TiB) write-thru activated]
[/dev/nymeAn] $[/dev/nymeAn]$ $(15 TiB) write-thru activated]$
a fileo
o neci
o ramisk
a jessi
L a ign 2002-01 and linux-isosi dmrtk-styr-d v864/st cdf6405648c
$[0^{-} dc]_{5}$ [ALLS: I]
[] 0- 1(1.1994-05.Com.Rev Hat:211a534302/5[mapped LUNS: 5]
o mapped_luno
0- mapped_luni
o- mapped_lunz
o- mapped_lungling investigation of the second s
[] o- mapped_lun4[lun4 block/nvme4 (rw)]
0- luns[LUNS: 5]
o- lunolplock/nvmeO (/dev/nvmeOnl)]
o- [nui[block/nvmei (/dev/nvmeiui)]
o- lun2[block/nvme2 (/dev/nvme2n1)]
o- lun3
o- lun4[block/nvme4 (/dev/nvme4n1)]
o- portals[Portals: 1]
o- 10.0.8.252:3260[iser]
o- loopback[Targets: 0]
o- srpt[Targets: 0]
/> saveconfig
Last 10 configs saved in /etc/target/backup
Configuration saved to /etc/target/saveconfig.json
/> exit



4. Start a connection from the initiator and login to the target. Once this is done, new devices should show up with lsblk.

#iscsiadm -m discovery --op=new --op=delete
--type sendtargets --portal <target portal IP>
-l iser

#iscsiadm –m node –l #lsblk The new LUNs will be visible to the initiator.

For our testing, we ran vdbench against our LUNs raw, but in most cases, the user would want to employ parted and mkfs to implement file systems on the new devices.

iscsiadm -m discovery --op=new --op=delete --type sendtargets --portal 10.0.8.253 -I iser 10.0.8.252:3260,1 ign.2003-01.org.linux-iscsi.dmrtk-srvr-d.x8664:sn. cdf664e56d8c Logging in to [iface: iser, target: iqn.2003-01.org.linux-iscsi.dmrtk-srvr-d.x8664:sn.cdf664e56d8c, portal: 10.0.8.252,3260] #lsblk NAME MAJ:MIN RM SIZE RO TYPE MOUNTPOINT 0 sda 8:0 477G 0 disk ⊢sda1 8:1 0 500M 0 part /boot 476.5G 50G part lvm -sda2 8:2 0 0 -rhel-root 253:0 Õ 0 -rhel-swap -rhel-home 253:1 253:2 0 4G 0 l∨m [SWAP] Ó 422.4G 0 l∨m /home sdb 8:16 0 1.5T 0 disk sdc 8:32 0 1.5T 0 disk 1.5T 1.5T sdd 8:48 0 0 disk sde 8:64 0 0 disk 0 1.5T 0 disk sdf 8:80 sr0 11:0 1 1024M 0 rom



Test Results

Expected Performance from RoCE

When we replace TCP/IP with RoCE, we expect CPU consumption to decrease due to the Remote Direct Memory Access capabilities eliminating CPU cycles use for data movement and transport layer processing. In addition when using TCP/IP, despite configuring Receive Side Scaling (RSS) on NICs, often certain CPU cores are heavily utilized for transport processing while the rest of the cores are left idle. On a smaller server with less processing power, this can result in a throughput performance cap and large latencies when TCP/IP is used. If RoCE is instead deployed on the same small server, we would expect cores to utilized more evenly and less CPU to be required for the same or greater amounts of traffic. If RoCE is deployed on a large server, we may notice the same throughput performance increase due to better traffic distribution between the cores. In most cases this leads to a greater amount of throughput, and in some cases, the throughput is so much greater that the CPU consumption will increase.

Below is a typical example of what a server running SMB over TCP/IP might look like. Several cores are heavily utilized, while many are left idle. The server in our example is a larger, more powerful server.

For additional detail on CPU and core consumption levels, see the *Processor Consumption Comments* section on page **30**.







Below is a typical example of what a server running SMB Direct over RoCE might look like. Notice that the Ethernet is not registering any traffic, however a file is being copied across the network and RDMA traffic is showing up in Perfmon. To the right is netstat, showing an active SMB Direct Connection.





100GbE Windows Storage Spaces – Reads

Throughput limitations do not hit line rate for writes. It is suspected that we are nearing the performance limit for the storage utilized (Windows Storage Spaces with NVMe drives.) For the Windows Storage Spaces tests, the throughput was higher and the CPU consumption was lower for the RoCE technology compared to traditional TCP/IP.







100GbE Windows Storage Spaces – Writes

Throughput limitations do not hit line rate for writes. It is suspected that we are nearing the performance limit for the storage utilized (Windows Storage Spaces with NVMe drives.) For the Windows Storage Spaces tests, the throughput was higher and the CPU consumption was lower for the RoCE technology compared to traditional TCP/IP.



4K

Block Size and Queue Depth

16

StorageSpaces SMB 3.0 over RoCE (SMB Direct)

256K

64K

Storage Spaces SMB 3.0 over TCP/IP



25GbE Windows Storage Spaces - Reads

In this example, the client is capable of processing 4x's as much IO as our NIC allows. This results in the NIC causing much of the bottleneck, in the case of both TCP/IP and RoCE.

For the Windows Storage Spaces tests, the throughput was higher and the CPU consumption was lower for the RoCE technology compared to traditional TCP/IP.



StorageSpaces SMB 3.0 over RoCE (SMB Direct)

Storage Spaces SMB 3.0 over TCP/IP



25GbE Windows Storage Spaces – Writes

In this example, the client is capable of processing 4x's as much IO as our NIC allows. This results in the NIC causing much of the bottleneck, in the case of both TCP/IP and RoCE.

For the Windows Storage Spaces tests, the throughput was generally higher and the CPU consumption was generally lower for the RoCE technology compared to traditional TCP/IP.





100GbE Linux iSER/iSCSI - Reads

iSCSI over TCP/IP has a low performance cap, enabling RoCE to vastly outperform it in throughput in both the case of 100Gb and 25Gb connections.







100GbE Linux iSER/iSCSI - Writes

iSCSI over TCP/IP has a low performance cap, enabling RoCE to vastly outperform it in throughput in both the case of 100Gb and 25Gb connections.





100GbE Linux iSER/iSCSI CPU Consumption vs. Throughput

Taking a deeper look into CPU consumption and throughput, we took the 8K block size data and plotted the results on the chart below.

Using iSER (RoCE) provided higher throughput at various levels of CPU consumption compared to traditional TCP/IP. In addition, using the traditional TCP/IP technology the throughput levels did not reach the levels achieved using the RoCE technology.





25GbE Linux iSER/iSCSI - Reads

iSCSI over TCP/IP has a low performance cap, enabling RoCE to vastly outperform it in throughput in both the case of 100Gb and 25Gb connections.







25GbE Linux iSER/iSCSI - Writes

iSCSI over TCP/IP has a low performance cap, enabling RoCE to vastly outperform it in throughput in both the case of 100Gb and 25Gb connections.



January 2017

Processor Consumption Comments

It should be noted that the client in these examples has plenty of processor to spare. Had the client been equipped with less processor, we would expect a cap on the throughput performance shown earlier and processor be near 100%. Below we show what the utilization was of the individual processor cores on our client. We can see that certain cores are heavily utilized.

Processor Consumption for 4K Reads with TCP/IP					
Processor	16 Worker 1 Queue Depth	24 Worker 1 Queue Depth	32 Worker 1 Queue Depth	48 Worker 1 Queue Depth	
0	37%	54%	67%	83%	
1	0%	0%	0%	0%	
2	22%	30%	36%	47%	
3	1%	1%	1%	0%	
4	13%	19%	25%	35%	
5	0%	0%	0%	0%	
6	6%	9%	14%	19%	
7	0%	0%	0%	0%	
8	0%	1%	1%	2%	
9	0%	0%	0%	0%	
10	0%	1%	0%	0%	
11	0%	0%	0%	0%	
12	0%	0%	0%	0%	
13	0%	0%	0%	0%	
14	0%	0%	0%	0%	
15	0%	0%	0%	0%	
16	0%	0%	0%	0%	
17	0%	0%	0%	0%	
18	0%	0%	0%	0%	
19	1%	1%	2%	2%	
20	29%	37%	44%	52%	

21	0%	0%	0%	0%
22	50%	65%	78%	90%
23	0%	0%	0%	0%
24	21%	28%	35%	46%
25	0%	0%	0%	0%
26	19%	29%	35%	45%
27	0%	0%	0%	0%
28	39%	59%	72%	87%
29	0%	0%	0%	0%
30	11%	17%	22%	35%
31	0%	0%	0%	0%
32	35%	52%	65%	83%
33	0%	0%	0%	0%
34	5%	11%	15%	28%
35	0%	0%	0%	0%
36	2%	6%	10%	23%
37	0%	0%	0%	0%
38	1%	2%	6%	17%
39	0%	0%	0%	1%

Processor Consumption for 256 K Reads with TCP/IP					
Processor	16 Worker 1 Queue Depth	24 Worker 1 Queue Depth	32 Worker 1 Queue Depth	48 Worker 1 Queue Depth	
0	57%	70%	82%	94%	
1	0%	0%	0%	0%	
2	20%	24%	29%	33%	
3	1%	2%	1%	1%	
4	3%	6%	4%	6%	
5	0%	0%	0%	0%	
6	0%	1%	0%	1%	
7	0%	0%	0%	0%	
	070	070	0,0	• • •	



9	0%	0%	0%	0%
10	0%	0%	0%	1%
11	0%	0%	0%	0%
12	0%	0%	0%	0%
13	0%	0%	0%	0%
14	0%	0%	0%	0%
15	0%	0%	0%	0%
16	0%	0%	0%	0%
17	0%	0%	0%	0%
18	0%	0%	0%	0%
19	1%	1%	2%	2%
20	24%	32%	35%	36%
21	0%	0%	0%	0%
22	54%	66%	80%	94%
23	0%	0%	0%	0%
24	15%	18%	25%	31%
25	0%	0%	0%	0%
26	13%	17%	23%	31%
27	0%	0%	0%	0%
28	50%	62%	75%	92%
29	1%	0%	0%	0%
30	2%	4%	6%	13%
31	0%	0%	0%	0%
32	47%	60%	70%	90%
33	0%	0%	0%	0%
34	1%	1%	2%	7%
35	0%	0%	0%	0%
36	0%	0%	0%	2%
37	0%	0%	0%	0%
38	0%	0%	0%	2%
39	1%	1%	0%	1%



RoCE Configuration Diagrams





Servers

Supermicro Server

- > 2x Intel Xeon E5-2690 v2, 3.0GHz, 20 total cores, 40 total threads
- > 256 GB RAM

Dell PowerEdge R920 Server

- > 4x Intel Xeon E7-4880 v2, 2.5GHz, 60 total cores, 120 total threads
- > 416 GB RAM
- > 5x Samsung SM-1715 1.6 GB NVMe SSD

Switches

Mellanox SN2410 25/100GbE Switch

- > 48 ports 25GbE
- > 8 ports 100GbE

Adapters

Mellanox

- > ConnectX-4 Lx EN single-port 25GbE
- > ConnectX-4 EN single-port 100GbE

QLogic

> FastLinQ QL45212 dual-port 25GbE

Storage Systems

HPE Storage

- > HPE 3PAR StoreServ File Controller
- > HPE 3PAR StoreServ 8450 All-Flash Storage



Conclusion

Generally, we found that using RoCE technology provided higher throughput and lower latency for 25GbE and 100GbE applications.

Having knowledge of Ethernet switch commands is beneficial to deploying RoCE solutions, and these commands may vary by switch vendor.

Deploying RoCE technology can boost the performance or lower the CPU consumption of applications that use Ethernet networks.

The most current version of this report is available at <u>http://www.demartek.com/Demartek_RoCE_Deployment_Guide.html</u> on the Demartek website.

Mellanox, Mellanox logo and ConnectX are trademarks of Mellanox Technologies, Ltd.

QLogic, the QLogic logo and FastLinQ are trademarks or registered trademarks of QLogic Corporation.

Demartek is a registered trademark of Demartek, LLC.

All other trademarks are the property of their respective owners.