

WinOF VPI for Windows

User Manual

Rev 4.2

www.mellanox.com

NOTE:

THIS HARDWARE, SOFTWARE OR TEST SUITE PRODUCT ("PRODUCT(S)") AND ITS RELATED DOCUMENTATION ARE PROVIDED BY MELLANOX TECHNOLOGIES "AS-IS" WITH ALL FAULTS OF ANY KIND AND SOLELY FOR THE PURPOSE OF AIDING THE CUSTOMER IN TESTING APPLICATIONS THAT USE THE PRODUCTS IN DESIGNATED SOLUTIONS. THE CUSTOMER'S MANUFACTURING TEST ENVIRONMENT HAS NOT MET THE STANDARDS SET BY MELLANOX TECHNOLOGIES TO FULLY QUALIFY THE PRODUCTO(S) AND/OR THE SYSTEM USING IT. THEREFORE, MELLANOX TECHNOLOGIES CANNOT AND DOES NOT GUARANTEE OR WARRANT THAT THE PRODUCTS WILL OPERATE WITH THE HIGHEST QUALITY. ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT ARE DISCLAIMED. IN NO EVENT SHALL MELLANOX BE LIABLE TO CUSTOMER OR ANY THIRD PARTIES FOR ANY DIRECT, INDIRECT, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES OF ANY KIND (INCLUDING, BUT NOT LIMITED TO, PAYMENT FOR PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY FROM THE USE OF THE PRODUCT(S) AND RELATED DOCUMENTATION EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.



Mellanox Technologies 350 Oakmead Parkway, Suite 100 Sunnyvale, CA 94085 U.S.A. www.mellanox.com Tel: (408) 970-3400 Fax: (408) 970-3403 Mellanox Technologies, Ltd. Beit Mellanox PO Box 586 Yokneam 20692 Israel www.mellanox.com Tel: +972 (0)4 909 7200 ; +972 (0)74 723 7200 Fax: +972 (0)4 959 3245

© Copyright 2012. Mellanox Technologies. All rights reserved.

Mellanox®, Mellanox Logo®, BridgeX®, ConnectX®, CORE-Direct®, InfiniBridge®, InfiniHost®, InfiniScale®, PhyX®, SwitchX®, Virtual Protocol Interconnect® and Voltaire® are registered trademarks of Mellanox Technologies, Ltd.

Connect-IB[™], FabricIT[™], MLNX-OS[™], Unbreakable-Link[™], UFM[™] and Unified Fabric Manager[™] are trademarks of Mellanox Technologies, Ltd.

All other trademarks are property of their respective owners.

2

Table of Contents

Document l	Revision History	.7
About this]	Manual	.8
	Scope Intended Audience Supported Network Adapter Cards Documentation Conventions Common Abbreviations and Acronyms	. 8 . 8 . 9 . 9
Chapter 1	Introduction	11
-	 1.1 Mellanox VPI Package Contents	11 11 12 12 12
Chapter 2	Driver Features	13
	 2.1 Hyper-V with VMQ 2.2 Header Data Split 2.3 Receive Side Scaling (RSS)	13 13 14 14 14 14 15 17
Chapter 3	Deploying Windows Server 2012 with SMB Direct	18
-	 3.1 Overview 3.2 Hardware and Software Prerequisites- 3.3 SMB Configuration Verification - 3.3.1 Verifying SMB Configuration 3.3.2 Verifying SMB Connection 3.4 Verifying SMB Events that Confirm RDMA Connection - 	18 18 18 18 18 19 19
Chapter 4	Driver Configuration	20
	 4.1 Configuring the InfiniBand Driver 4.1.1 Modifying IPoIB Configuration 4.1.2 Displaying Adapter Related Information 4.2 Configuring the Ethernet Driver 	20 20 20 23
Chapter 5	Performance Tuning	25
	 5.1 General Performance Optimization and Tuning	25 25 25 25 26 30 30 30 31
Chanter 6	OpenSM - Subnet Manager	35
Chanter 7	InfiniBand Fabric	36
Suppor /	 7.1 Network Direct Interface	36 36 37

Rev 4.2

	7.3.1	Utilities Usage	7
	7.3.2	ibportstate	8
	7.3.3	ibroute	2
	7.3.4	smpquery	5
	7.3.5	perfquery	9
	7.3.6	ibping	3
	7.3.7	ibnetdiscover	4
	7.3.8	ibtracert	8
	7.3.9	sminfo	0
	7.3.10	ibclearerrors	1
	7.3.11	ibstat	2
	7.3.12	vstat	2
	7.3.13	osmtest	3
	7.4 InfiniBan	d Fabric Performance Utilities 60	5
	7.4.1	ib_read_bw	5
	7.4.2	ib_read_lat	/
	7.4.3	Ib_send_bw	5
	7.4.4	Ib_send_lat	۶ ۵
	7.4.5	ID_WIIIE_DW	1 0
	7.4.0	iby mod by 7	J 1
	7.4.7 7.4.8	iby read lat	1 2
	7.4.8	iby send by 7/) 1
	7.4.0	iby send lat 74	+ 5
	7.4.10	ibv_schu_iat	, 7
	7.4.11	ibv_write_lat	8
	7.4.12	ibaddr 70	9
	7.4.14	ibcacheedit 8	1
	7.4.15	iblinkinfo	2
	7.4.16	ibquerverrors	4
	7.4.17	ibsysstat	6
	7.4.18	perfquery	8
	7.4.19	saquery	1
	7.4.20	smpdump	3
Chapter 8	Software Do	evelopment Kit	5
Chapter 9	Troublesho	97	7
-	9.1 InfiniBan	d Troubleshooting 9'	7
	9.2 Ethernet	Froubleshooting9	7
Chapter 10	Document	ation)

4

List of Tables

Table 1:	Document Revision History	7
Table 2:	Documentation Conventions	9
Table 3:	Abbreviations and Acronyms	9
Table 4:	Registry Keys Setting	.14
Table 5:	ibportstate Flags and Options	. 39
Table 6:	ibroute Flags and Options	.42
Table 7:	smpquery Flags and Options	.45
Table 8:	perfquery Flags and Options	.49
Table 9:	ibping Flags and Options	.53
Table 10:	ibnetdiscover Flags and Options	. 54
Table 11:	ibtracert Flags and Options	. 59
Table 12:	sminfo Flags and Options	.60
Table 13:	ibclearerrors Flags and Options	.61
Table 14:	ibstat Flags and Options	.62
Table 15:	vstat Flags and Options	.63
Table 16:	osmtest Flags and Options	.64
Table 17:	ib_read_bw Flags and Options	.66
Table 18:	ib_read_lat Flags and Options	.67
Table 19:	ib_send_bw Flags and Options	.68
Table 20:	ib_send_lat Flags and Options	.69
Table 21:	ib_write_bw Flags and Options	.70
Table 22:	ib_write_lat Flags and Options	.71
Table 23:	ibv_read_bw Flags and Options	.72
Table 24:	ibv_read_lat Flags and Options	.73
Table 25:	ibv_send_bw Flags and Options	.74
Table 26:	ibv_send_lat Flags and Options	.76
Table 27:	ibv_write_bw Flags and Options	.77
Table 28:	ibv_write_lat Flags and Options	.78
Table 29:	ibaddr Flags and Options	.80
Table 30:	ibcacheedit Flags and Options	.82
Table 31:	iblinkinfo Flags and Options	.83
Table 32:	ibqueryerrors Flags and Options	.85
Table 33:	ibsysstat Flags and Options	.87
Table 34:	perfquery Flags and Options	. 89
Table 35:	saquery Flags and Options	.92

Table 36: sn	mpdump Flags and Options		1
--------------	--------------------------	--	---

Document Revision History

Table 1 - Document Revision History

Document Revision	Date	Changes
Rev 4.2	October 20, 2012	 Added Section 3, "Deploying Windows Server 2012 with SMB Direct," on page 18, and its subsections Updated Section 5, "Performance Tuning," on page 25 Added Section 2.2, "Header Data Split," on page 13 Added Section 7.2, "part_man - Virtual IPoIB Port Creation Utility," on page 36
Rev 3.2.0	July 23, 2012	No changes
Rev 3.1.0	May 21, 2012	 Added section Tuning the IPoIB Network Adapter Added section Tuning the Ethernet Network Adapter Added section Performance tuning tool application Removed section Tuning the Network Adapter Removed section part_man Removed section ibdiagnet
Rev 3.0.0	February 08, 2012	 Added section RDMA over Converged Ethernet (RoCE) and its subsections Added section Hyper-V with VMQ Added section Network Driver Interface Specification (NDIS) Added section Header Data Split Added section Auto Sensing Added section Adapter Teaming Added section Advanced Configuration Added section Advanced Configuration for InfiniBand Driver Added section Updated section Tunable Performance Parameters Added section Sockets Direct Protocol and its subsections Removed section Winsock Direct and Protocol and its subsections Removed section Added ConnectX®-3 support Removed section IPoIB Drivers Overview Removed section Booting Windows from an iSCSI Target
Rev 2.1.3	January 28. 2011	Complete restructure
Rev 2.1.2	October 10, 2010	 Removed section Debug Options. Updated Section 3, "Uninstalling Mellanox VPI Driver," on page 11 Added Section 6, "InfiniBand Fabric," on page 38 and its subsections Added Section 6.3, "InfiniBand Fabric Performance Utilities," on page 71 and its subsections
Rev 2.1.1.1	July 14, 2010	Removed all references of InfiniHost® adapter since it is not supported starting with WinOF VPI v2.1.1
Rev 2.1.1	May 2010	First release

About this Manual

Scope

The document describes WinOF Rev 4.2 features, performance, InfiniBand diagnostic, tools content and configuration. Additionally, this document provides information on various performance tools supplied with this version.

Intended Audience

This manual is intended for system administrators responsible for the installation, configuration, management and maintenance of the software and hardware of VPI (InfiniBand, Ethernet) adapter cards. It is also intended for application developers.

Supported Network Adapter Cards

Mellanox MLNX_WinOF_4.2_win8_x64.exe supports the following Mellanox network adapter cards:

- ConnectX®-2 SDR/DDR/QDR
- ConnectX®-3 FDR/SDR/QDR

Documentation Conventions

Table 2 -	Documentation	Conventions
-----------	---------------	-------------

Description	Convention	Example
File names	file.extension	
Directory names	directory	
Commands and their parameters	command param1	mts3610-1 > show hosts
Required item	<>	
Optional item	[]	
Mutually exclusive parameters	${p1, p2, p3}$ or ${p1 p2 p3}$	
Optional mutually exclusive parameters	[p1 p2 p3]	
Variables for which users supply specific values	Italic font	enable
Emphasized words	Italic font	These are emphasized words
Note	<text></text>	This is a note
Warning	<text></text>	May result in system instability.

Common Abbreviations and Acronyms

Table 3 - Abbreviations and Acronyms	(Sheet 1 of 2)
--------------------------------------	----------------

Abbreviation / Acronym	Whole Word / Description
В	(Capital) 'B' is used to indicate size in bytes or multiples of bytes (e.g., 1KB = 1024 bytes, and 1MB = 1048576 bytes)
b	(Small) 'b' is used to indicate size in bits or multiples of bits (e.g., 1Kb = 1024 bits)
FW	Firmware
НСА	Host Channel Adapter
HW	Hardware
IB	InfiniBand
LSB	Least significant byte
lsb	Least significant bit
MSB	Most significant byte

Rev 4.2

Abbreviation / Acronym	Whole Word / Description
msb	Most significant bit
NIC	Network Interface Card
SW	Software
VPI	Virtual Protocol Interconnect
IPoIB	IP over InfiniBand
PFC	Priority Flow Control
PR	Path Record
RDS	Reliable Datagram Sockets
RoCE	RDMA over Converged Ethernet
SL	Service Level
MPI	Message Passing Interface
EoIB	Ethernet over Infiniband
QoS	Quality of Service
ULP	Upper Level Protocol
VL	Virtual Lane

Table 3 - Abbreviations and Acronyms (Sheet 2 of 2)

1 Introduction

This User Manual addresses the Mellanox WinOF VPI driver Rev 4.2 package distributed for Windows Server 2012 (x64).

Mellanox WinOF VPI is composed of several software modules that contain an InfiniBand and Ethernet driver. The Mellanox WinOF VPI driver supports Infiniband and 10GB Ethernet ports. The port type is determined upon boot based on card's capability and user setting.

1.1 Mellanox VPI Package Contents

The Mellanox WinOF for Windows package contains the following components:

- Core and ULPs
 - IB network adapter cards low-level drivers (mlx4)
 - IB Access Layer (IBAL)
 - Ethernet driver (ETH)
 - IP over InfiniBand (IPoIB)
 - Upper Layer Protocols (ULPs):
 - NetworkDirect (ND)
- Utilities
- SW Development Kit (SDK)
- Documentation

1.2 Hardware and Software Requirements

- Administrator privileges on your machine(s)
- Disk Space for installation: 100MB

1.3 Supported Network Adapter Cards

Mellanox WinOF 4.2 for Windows 2012 supports the following Mellanox network adapter cards:

- ConnectX®-3/ConnectX®-3 EN QDR/FDR10/FDR
- ConnectX®-2/ConnectX®-2 EN SDR/DDR/QDR



Mellanox recommends upgrading ConnectX®-2 and ConnectX®-3 to the latest General Availability (GA) firmware version.

To enable improved functionality while using this WinOF release, it is recommended to upgrade ConnectX@-2 adapter cards to firmware v2.10.0720 or higher, and ConnectX@-3 adapter cards to firmware v2.11.0500 or higher.

The adapter card may not have been shipped with the latest firmware version. This section describes how to update firmware.

1.3.1 Downloading the Firmware Tools Package

1. Download Mellanox Firmware Tools

Please download the current firmware tools package (MFT) from http://www.mellanox.com > Products > Software/Drivers > InfiniBand & VPI SW/Drivers > Firmware Tools.

The tools package to download is "MFT Software for Windows x64" for x64 architecture.

2. Install and Run WinMFT

To install the WinMFT package, double click the MSI or run it from the command prompt.



Install the WinMFT package from the command line with administrator privileges.

Enter:

Rev 4.2

```
msiexec.exe /i WinMFT <arch> <version>.msi
```

- 3. Check the Device Status
 - start/stop mst is automatically done by the tools > C:\Users\herod\Desktop>mst start
 - To check device status run > mst status

If no card installation problems occur, the status command should produce the following output:

mt<device id> pciconf0 mt<device id>_pci_cr0

where device ID will be one of the supported PCI device IDs.

1.3.2 Downloading the Firmware Image of the Adapter Card

- To download the correct card firmware image, please visit http://www.mellanox.com > Support > Firmware Download
- To identify your adapter card, please visit http://www.mellanox.com > Support > Firmware Downloads > Identifying Adapter Cards

1.3.3 Updating Adapter Card Firmware

Using a card specific binary firmware image file, enter the following command:

> flint -d mt<device id> pci cr0 -i <image name.bin> burn

For additional details, please check the MFT user's manual under

http://www.mellanox.com > Products > Adapter IB/VPI SW

2 Driver Features

The Mellanox VPI WinOF driver release introduces the following capabilities:

- One or two ports
- Up to 16 Rx queues per port
- Rx steering mode (RSS)
- Hardware Tx/Rx checksum calculation
- Large Send Offload (i.e., TCP Segmentation Offload)
- Hardware multicast filtering
- Adaptive interrupt moderation
- MSI-X support
- Auto Sensing
- RoCE

Ethernet Only:

- HW VLAN filtering
- Header Data Split

For the complete list of Ethernet and InfiniBand Known Issues and Limitations, see MLNX_WinVPI_ReleaseNotes.txt.

2.1 Hyper-V with VMQ

Mellanox WinOF Rev 4.2 includes a virtual machine queue (VMQ) interface to support Microsoft Hyper-V network performance improvements and security enhancement.

VMQ interface supports:

- Classification of received packets by using the destination MAC address to route the packets to different receive queues
- NIC ability to use DMA to transfer packets directly to a Hyper-V child-partition's shared memory
- Scaling to multiple processors, by processing packets for different virtual machines on different processors.

2.2 Header Data Split

The header-data split feature improves network performance by splitting the headers and data in received Ethernet frames into separate buffers. The feature is disabled by default and can be enabled in the Advanced tab (Performance Options) from the Properties sheet.

For further information, please refer to the MSDN library:

http://msdn.microsoft.com/en-us/library/windows/hardware/ff553723(v=VS.85).aspx

2.3 Receive Side Scaling (RSS)

Mellanox WinOF Rev 4.2 IPoIB and Ethernet drivers use NDIS 6.30 new RSS capabilities. The main changes are:

- Supports unlimited number of processors (previously 64)
- Individual network adapter RSS configuration usage

To set the RSS capability for individual adapter instead of global setting, and to improve RSS on Windows 2012 server, set the registry keys listed in the table below:

Sub-key	Description
HKLM\SYSTEM\CurrentControlSet\Control\Class\{XXXX72- XXX}\ <network adapter="" number="">*MaxRSSProcessors</network>	Maximum number of CPUs allotted. Sets the desired maximum number of processors for each interface. The number can be different for each interface. Note: Restart the network adapter when you change this registry key.
HKLM\SYSTEM\CurrentControlSet\Control\Class\{XXXX72- XXX}\ <network adapter="" number="">*RssBaseProcNumber</network>	Base CPU number. Sets the desired base CPU number for each interface. The number can be different for each interface. This allows partition- ing of CPUs across network adapters. Note: Restart the network adapter when you change this registry key.
HKLM\SYSTEM\CurrentControlSet\Control\Class\{XXXX72- XXX}\ <network adapter="" number="">*NumaNodeID</network>	NUMA node affinitization
HKLM\SYSTEM\CurrentControlSet\Control\Class\{XXXX72- XXX}\ <network adapter="" number="">*RssBaseProcGroup</network>	Sets the RSS base processor group for systems with more than 64 processors.

Table 4 - Registry Keys Setting

2.4 Port Configuration

After MLNX_VPI installation, it is possible to modify the network protocol that runs on each port of VPI adapter cards. Each port can be set to run as InfiniBand, Ethernet or Auto Sensing.

2.4.1 Auto Sensing

Auto Sensing enables the NIC to automatically sense the link type (InfiniBand or Ethernet) based on the cable connected to the port and load the appropriate driver stack (InfiniBand or Ethernet).

Auto Sensing is performed only when rebooting the machine or after disabling/enabling the mlx4_bus interface from the Device Manager. Hence, if you replace cables during the runtime, the NIC will not perform Auto Sensing.

2.4.2 Port Protocol Configuration

Step 1 Display the Device Manager and expand "System devices".



Step 2. Right-click on the Mellanox ConnectX VPI network adapter and left-click Properties. Select the Port Protocol tab from the Properties sheet.



The figure below is an example of the displayed Port Protocol sheet for a dual port VPI adapter card.

Mellanox Connec	tX-3 VPI (MT04099) - PCIe 3.0 5GT/s, 🗴	
General Port Protoc	ol Driver Details Events Resources	
VECHNOLODIES	Current Setting Port1: IB Port2: Eth	
HCA Port Type Con	figuration	
HW Defaults	Port 1	
	CIB © ETH CAUTO	
Port Protocol Configuration This menu displays the adapter's port type and enables you to set the network protocols for the network adapter ports. The network protocol is determined according to the NIC's Hardware Defaults port type. You can choose the protocol explicitly by selecting the port type to InfiniBand (IB) or Ethernet (Eth). To enable Auto Sensing please choose ALITO. If the NIC		
	OK Cancel	

Step 3. In this step, you can perform the following functions:

- Choose HW Defaults option. If you choose the HW Defaults option, the port protocols will be determine according to the NIC's hardware default values.
- Choose the desired port protocol for the available port(s). If you choose IB or ETH, both ends of the connection must be of the same type (IB or ETH).
- Enable Auto Sensing by checking the AUTO checkbox. If the NIC does not support Auto Sensing, the AUTO option will be grayed out.



If you choose AUTO, Current Setting will indicate the actual port settings: IB or ETH.

WinOF VPI for Windows User Manual

2.5 Load Balancing, Fail-Over (LBFO) and VLAN

Windows 2012 Microsoft supports load balancing as part of the operating system. plese refer to Microsoft document "Windows 2012 Link Teaming".

3 Deploying Windows Server 2012 with SMB Direct

3.1 Overview

The Server Message Block (SMB) Protocol is a network file sharing protocol implemented in Microsoft Windows. The set of message packets that defines a particular version of the protocol is called a dialect.

The Microsoft SMB Protocol is a client-server implementation and consists of a set of data packets, each containing a request sent by the client or a response sent by the server.

SMB protocol is used on top of the TCP/IP protocol or other network protocols. Using the SMB protocol allows applications to access files or other resources on a remote server, to read, create, and update them. In addition, it enables communication with any server program that is set up to receive an SMB client request.

3.2 Hardware and Software Prerequisites

The following are Hardware and Software prerequisites:

- Two or more machines running Windows Server 2012
- One or more Mellanox ConnectX®-2 or ConnectX®-3 adapters for each server
- One or more Mellanox InfiniBand switches
- Two or more QSFP cables required for InfiniBand



SMB Direct in Windows Server 2012 does not support older Mellanox InfiniBand adapters (including ConnectX® and InfiniHost® III adapters).

3.3 SMB Configuration Verification

3.3.1 Verifying SMB Configuration

Use the following PowerShell cmdlets to verify SMB Multichannel is enabled, confirm the adapters are recognized by SMB and that their RDMA capability is properly identified.

• On the SMB client, run the following PowerShell cmdlets:

```
Get-SmbClientConfiguration | Select EnableMultichannel
Get-SmbClientNetworkInterface
```

18

• On SMB server, run the following PowerShell cmdlets:

```
Get-SmbServerConfiguration | Select EnableMultichannel
Get-SmbServerNetworkInterface
netstat.exe<sup>a</sup> -xan | ? {$_ -match "445"}
```

a. The NETSTAT command confirms if the File Server is listening on the RDMA interfaces.

3.3.2 Verifying SMB Connection

» To verify SMB connection on SMB client:

- 1. Start a long-running file copy to create a lasting session with the SMB Server.
- 2. Open a PowerShell window while the copy is ongoing, and run the following cmdlets to verify the connection is using the right SMB dialect and that SMB Direct is working:

```
Get-SmbConnection
Get-SmbMultichannelConnection
netstat.exe -xan | ? {$_ -match "445"}
```



If you have no activity while you run the commands above, you might get an empty list due to session expiration and no current connections.

3.4 Verifying SMB Events that Confirm RDMA Connection

On the SMB client, open a PowerShell window and run the following cmdlets to view the SMB events that confirm that you have an SMB Direct connection. Any RDMA-related connection errors will be displayed as well:

Get-WinEvent -LogName Microsoft-Windows-SMBClient/Operational | ? Message -match "RDMA"

4 Driver Configuration

Once you have installed Mellanox WinOF VPI package, you can perform various modifications to your driver to make it suitable for your system's needs



Rev 4.2

Changes made to the Windows registry happen immediately, and no backup is automatically made.

Do not edit the Windows registry unless you are confident regarding the changes.

4.1 Configuring the InfiniBand Driver

4.1.1 Modifying IPoIB Configuration

To modify the IPoIB configuration after installation, perform the following steps:

- Step 1 Open Device Manager and expand Network Adapters in the device display pane.
- Step 2. Right-click the Mellanox IPoIB Adapter entry and left-click Properties.
- Step 3. Click the Advanced tab and modify the desired properties.



The IPoIB network interface is automatically restarted once you finish modifying IPoIB parameters. Consequently, it might affect any running traffic.

4.1.2 Displaying Adapter Related Information

To display a summary of network adapter software, firmware- and hardware-related information such as driver version, firmware version, bus interface, adapter identity, and network port link information, perform the following steps:





Step 2. Right-click a Mellanox ConnectX VPI adapter (under "System devices" list) and left-click Properties.

Details	Events		Power Manage	ement	
General Advan	ced Inform	ation F	Performance	Driver	
Mellanox	Adapter Info	ormation			
Information		Value			
Driver Version		4.2.11	165.0		
Firmware Version Port Number		2.11.50	JU		
Bus Type		PCI-E S	5.0 Gbps x8		
Link Speed					
Part Number		MCX35	54A-FCBT		
Device Id		4099			
Revision Id Commit MAC Address		00.021			
Permanent MAC Address	; 1800	00-02-C3-30-3E-F0 00-02-C9-35-9E-E0			
Network Status		Discon	nected		
Adapter Friendly Name		Ethern	et 3		
IPv4 Address 169.254.27.228					
Adapter User Name					
			Save	To File	

Step 3. Select the Information tab from the Properties sheet.



To save this information for debug purposes, click **Save To File** and provide the output file name.

4.2 Configuring the Ethernet Driver

The following steps describe how to configure advanced features.

Step 1 Display the Device Manager.



Step 2. Right-click a Mellanox network adapter (under "Network adapters" list) and left-click Properties. Select the Advanced tab from the Properties sheet.

Mellano	x Conr	nectX-3	3 Ethern	et Ac	dapter Prope	erties	x
Detail:	3	Events			Power Management		nt
General	Adva	anced	Informa	tion	Performance Driv		Driver
The followin the property on the right <u>Property:</u> <u>Bus-maste</u> Flow Contr Header Da Interrupt M Interrupt M	r DMA 0 ol ata Split oderatio oderatio oderatio deratio deratio deratio deratio deratio deratio deratio deratio deratio	nt to cha peration n RX Par n RX Par n TX Par n TX Par fload d (LSO) d V2 (IPv d V2 (IPv d Version	available for nge on the cket Cr cket Ti cket Cr cket Ti cket Ti (IPv. ~	r this r left, a	network adapter. nd then select its ⊻alue: Enabled	Click value	Y
-			OK		Cancel	ŀ	lelp

Step 3. Modify configuration parameters to suit your system.

Please note the following:

- a. For help on a specific parameter/option, check the help window at the bottom of the dialog.
- b. If you select one of the entries Offload Options, Performance Options, or Flow Control Options, you'll need to click the Properties button to modify parameters via a pop-up dialog. See example in the two figures below.
- c. A "Use Default for All" button appears on the Advanced dialog. Click this button to set all entries (and their sub-entries) to the Mellanox Ethernet driver default values. You will be prompted to approve this action.

Rev 4.2

5 Performance Tuning



This document describes how to modify Windows registry parameters in order to improve performance.

Please note that modifying the registry incorrectly might lead to serious problems, including the loss of data, system hang, and you may need to reinstall Windows. As such it is recommended to back up the registry on your system before implementing recommendations included in this document. If the modifications you apply lead to serious problems, you will be able to restore the original registry state. For more details about backing up and restoring the registry, please visit www.microsoft.com.

5.1 General Performance Optimization and Tuning

To achieve the best performance for Windows, you may need to modify some of the Windows registries.

5.1.1 Registry Tuning

The registry entries that may be added/changed by this "General Tuning" procedure are:

Under HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Services\Tcpip\Parameters:

• Disable TCP selective acks option for better cpu utilization:

SackOpts, type REG_DWORD, value set to 0.

Under HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Services\AFD\Parameters:

• Enable fast datagram sending for UDP traffic:

FastSendDatagramThreshold, type REG_DWORD, value set to 64K.

Under HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Services\Ndis\Parameters:

• Set RSS parameters:

RssBaseCpu, type REG_DWORD, value set to 1.

5.1.2 Enable RSS

Enabling Receive Side Scaling (RSS) is performed by means of the following command:

"netsh int tcp set global rss = enabled"

5.1.3 Tuning the IPolB Network Adapter

The IPoIB Network Adapter tuning can be performed either during installation by modifying some of Windows registries as explained in Section 5.1.1, "Registry Tuning," on page 25. or can be set post-installation manually. To improve the network adapter performance, activate the performance tuning tool as follows:

- 1. Start the "Device Manager" (open a command line window and enter: devmgmt.msc).
- 2. Open "Network Adapters".
- 3. Select Mellanox IPoIB adapter, right click and select Properties.
- 4. Select the "Performance tab".
- 5. Choose one of the tuning scenarios:
 - a. Single port traffic Improves performance for running single port traffic each time.
 - b. Dual port traffic Improves performance for running traffic on both ports simultaneously.
 - c. Forwarding traffic Improves performance for running scenarios that involve both ports (for example: via IXIA)
 - d. Multicast traffic Improves performance when the main traffic runs on multicast.
 - e. Single stream traffic Optimizes tuning for applications with single connection
- 6. Click on "Run Tuning" button.

Clicking the "Run Tuning" button will change several registry entries (described below), and will check for system services that may decrease network performance. It will also generate a log including the applied changes.

Users can view this log to restore the previous values. The log path is:

%HOMEDRIVE%\Windows\System32\LogFiles\PerformanceTunning.log

This tuning is required to be performed only once after the installation is completed, and on one adapter only (as long as these entries are not changed directly in the registry, or by some other installation or script).

Please note that a reboot may be required for the changes to take effect.

5.1.4 Tuning the Ethernet Network Adapter

The Ethernet Network Adapter general tuning can be performed during installation by modifying some of Windows registries as explained in section "Registry Tuning" on page 32. Specific scenarios tuning can be set post-installation manually. To improve the network adapter performance, activate the performance tuning tool as follows:

- 1. Start the "Device Manager" (open a command line window and enter: devmgmt.msc).
- 2. Open "Network Adapters".
- 3. Select Mellanox Ethernet adapter, right click and select Properties.
- 4. Select the "Performance tab".
- 5. Choose one of the tuning scenarios:
 - a. Single port traffic Improves performance for running single port traffic each time.
 - b. Dual port traffic Improves performance for running traffic on both ports simultaneously.
 - Forwarding traffic Improves performance for running scenarios that involve both ports (for example: via IXIA)
 - d. Multicast traffic Improves performance when the main traffic runs on multicast.
 - e. Single stream traffic Optimizes tuning for applications with single connection

6. Click on "Run Tuning" button.



Clicking the "Run Tuning" button will activate the general tuning as explained above and change several driver registry entries for the current adapter and it's sibling device, if the sibling is an Ethernet device as well. It will also generate a log including the applied changes.

Users can view this log to restore the previous values. The log path is:

%HOMEDRIVE%\Windows\System32\LogFiles\PerformanceTunning.log

This tuning is required to be performed only once after the installation is completed, and on one adapter only (as long as these entries are not changed directly in the registry, or by some other installation or script).



Please note that a reboot may be required for the changes to take effect.

5.1.4.1 Performance Tuning Tool Application

You can also activate the performance tuning through a script called perf_tuning.exe. This script has 4 options, which include the 3 scenarios described above and an additional manual tuning through which you can set the RSS base and number of processors for each Ethernet adapter. The adapters you wish to tune are supplied to the script by their name according to the "Network Connections".

5.1.4.2 Synopsys

perf_tuning.exe -s -cl <first connection name> [-c2 <second connection name>]
perf_tuning.exe -d -cl <first connection name> -c2 <second connection name>
perf_tuning.exe -f -cl <first connection name> -c2 <second connection name>
perf_tuning.exe -m -cl <first connection name> -b <base RSS processor number> -n <number of RSS
processors>

5.1.4.3 Options

Flag	Description
-S	 Single port traffic scenario. This option can be followed by one or two connection names. The tuning will be performed on each one of them separately. This option chooses the best processors and values to assign to: DefaultRecvRingProcessor TxInterruptProcessor TxForwardingProcessor *RssBaseProcNumber *MaxRssProcessors This option also sets the following NDIS registry keys: HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\services\NDIS\Parameters\Rss-BaseCpu to 0 HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\services\NDIS\Parameters\MaxN-umRssCpus to the number of processors
-d	 Dual port traffic scenario. This option must be followed by two connection names. The tuning in this case is codependent. This option chooses the best processors and values to assign to: DefaultRecvRingProcessor TxInterruptProcessor TxForwardingProcessor *RssBaseProcNumber *MaxRssProcessors This option also sets the following NDIS registry keys: HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\services\NDIS\Parameters\Rss-BaseCpu to 0 HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\services\NDIS\Parameters\MaxN-umRssCpus to the number of processors

Flag	Description			
-f	Forwarding traffic scenario. This option must be followed by two connection names. The tuning in this case is codependent. This option automatically sets: • SendCompletionMethod = 1 • RecvCompletionMethod = 0 • *ReceiveBuffers = 4096 • UseRSSForRawIP = 0 Additionally, this option chooses the best processors to assign to: • DefaultRecvRingProcessor • TxInterruptProcessor • TxForwardingProcessor			
-m	 Manual configuration This option must be followed by one connection name. This option assigns the provided base and number of CPUs to: *RssBaseProcNumber *MaxRssProcessors Additionally, this option assigns the following with processors inside the range: DefaultRecvRingProcessor TxInterruptProcessor TxInterruptProcessor This option also sets the following NDIS registry keys: HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\services\NDIS\Parameters\Rss-BaseCpu to 0 HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\services\NDIS\Parameters\MaxN-umRssCpus to the number of processors 			
-r	Restore default settings. This option can be followed by one or two connection names. This option automatically sets the driver registry values back to their default values: • SendCompletionMethod = 0 • RecvCompletionMethod = 2 • *ReceiveBuffers = 1024 • UseRSSForRawIP = 1 • DefaultRecvRingProcessor = -1 • TxInterruptProcessor = -1			
-c1	Specifies first connection name. See examples			
-c2	Specifies second connection name. See examples			
-b	Specifies base RSS processor number. See examples. Used for manual option (-m) only.			
-n	Specifies number of RSS processors. See examples. Used for manual option (-m) only.			

Examples

For example, if the adapter is represented by "Local Area Connection 6" and "Local Area Connection 7"

```
For single port streams tuning type:
perf_tuning.exe -s -c1 "Local Area Connection 6" -c2 "Local Area Connection 7"
or to set one adapter only:
perf_tuning.exe -s -c1 "Local Area Connection 6"
For dual port streams tuning type:
perf_tuning.exe -d -c1 "Local Area Connection 6" -c2 "Local Area Connection 7"
For forwarding streams tuning type:
perf_tuning.exe -f -c1 "Local Area Connection 6" -c2 "Local Area Connection 7"
For manual tuning of the first adapter to use RSS on CPUs 0-3:
perf_tuning.exe -m -c1 "Local Area Connection 6" -b 0 -n 4
In order to restore defaults type:
perf tuning.exe -r -c1 "Local Area Connection 6" -c2 "Local Area Connection 7"
```

5.2 Application Specific Optimization and Tuning

5.2.1 Ethernet Performance Tuning

The user can configure the Ethernet adapter by setting some registry keys. The registry keys may affect Ethernet performance.

To improve performance, activate the performance tuning tool as follows:

- 1. Start the "Device Manager" (open a command line window and enter: devmgmt.msc).
- 2. Open "Network Adapters".
- 3. Right click the relevant Ethernet adapter and select Properties.
- 4. Select the "Advanced" tab and select Performance Options
- 5. Modify performance parameters (properties) as desired.

5.2.1.1 Performance Known Issues

- On Intel I/OAT supported systems, it is highly recommended to install and enable the latest I/OAT driver (download from www.intel.com).
- With I/OAT enabled, sending 256-byte messages or larger will activate I/OAT. This will cause a significant latency increase due to I/OAT algorithms. On the other hand, throughput will increase significantly when using I/OAT.

5.2.2 IPolB Performance Tuning

The user can configure the IPoIB adapter by setting some registry keys. The registry keys may affect IPoIB performance.

For the complete list of registry entries that may be added/changed by the performance tuning procedure, see the IPoIB_registry_values.pdf file.

To improve performance, activate the performance tuning tool as follows:

- 1. Start the "Device Manager" (open a command line window and enter: devmgmt.msc).
- 2. Open "Network Adapters".
- 3. Right click the relevant IPoIB adapter and select Properties.
- 4. Select the "Advanced" tab
- 5. Modify performance parameters (properties) as desired.

5.3 **Tunable Performance Parameters**

The following is a list of key parameters for performance tuning.

• Jumbo Packet

The maximum available size of the transfer unit, also known as the Maximum Transmission Unit (MTU). For IPoIB, the MTU should not include the size of the IPoIB header (=4B). For example, if the network adapter card supports a 4K MTU, the upper threshold for payload MTU is 4092B and not 4096B. The MTU of a network can have a substantial impact on performance. A 4K MTU size improves performance for short messages, since NDIS can coalesce a small message into a larger one.

Valid MTU values range for an Ethernet driver is between 600 and 9600.

Valid MTU values range for an IPoIB driver is between 1500 and 4092.



All devices on the same physical network, or on the same logical network, must have the same MTU.

• Receive Buffers

The number of receive buffers (default 1024).

Send Buffers

The number of sent buffers (default 2048).

Performance Options

Configures parameters that can improve adapter performance.

Interrupt Moderation

Moderates or delays the interrupts' generation. Hence, optimizes network throughput and CPU utilization (default Enabled).

- When the interrupt moderation is enabled, the system accumulates interrupts and sends a single interrupt rather than a series of interrupts. An interrupt is generated after receiving 5 packets or after 10ms from the first packet received. It improves performance and reduces CPU load however, it increases latency.
- When the interrupt moderation is disabled, the system generates an interrupt each time a packet is received or sent. In this mode, the CPU utilization data rates increase, as the system handles a larger number of interrupts. However, the latency decreases as the packet is handled faster.

Improves incoming packet processing performance. RSS enables the adapter port to utilize the multiple CPUs in a multi-core system for receiving incoming packets and steering them to the designated destination. RSS can significantly improve the number of transactions, the number of connections per second, and the network throughput.

This parameter can be set to one of the following values:

- Enabled (default): Set RSS Mode
- Disabled: The hardware is configured once to use the Toeplitz hash function, and the indirection table is never changed.



IOAT is not used while in RSS mode.

Receive Completion Method

Sets the completion methods of the received packets, and can affect network throughput and CPU utilization.

Polling Method

Increases the CPU utilization as the system polls the received rings for the incoming packets. However, it may increase the network performance as the incoming packet is handled faster.

Interrupt Method

Optimizes the CPU as it uses interrupts for handling incoming messages. However, in certain scenarios it can decrease the network throughput.

Adaptive (Default Settings)

A combination of the interrupt and polling methods dynamically, depending on traffic type and network usage. Choosing a different setting may improve network and/or system performance in certain configurations.

Interrupt Moderation RX Packet Count

Number of packets that need to be received before an interrupt is generated on the receive side (default 5).

Interrupt Moderation RX Packet Time

Maximum elapsed time (in usec) between the receiving of a packet and the generation of an interrupt, even if the moderation count has not been reached (default 10).

Rx Interrupt Moderation Type

Sets the rate at which the controller moderates or delays the generation of interrupts making it possible to optimize network throughput and CPU utilization. The default setting (Adaptive) adjusts the interrupt rates dynamically depending on the traffic type and network usage. Choosing a different setting may improve network and system performance in certain configurations.

Send completion method

Sets the completion methods of the Send packets and it may affect network throughput and CPU utilization.

Interrupt Moderation TX Packet Count

Number of packets that need to be sent before an interrupt is generated on the send side (default 0).

Interrupt Moderation TX Packet Time

Maximum elapsed time (in usec) between the sending of a packet and the generation of an interrupt even if the moderation count has not been reached (default 0).

Bus-master DMA Operations

Sets the addressing type: NDIS DMA addressing (UseDma=Enabled) or physical addressing (UseDma=Disabled) (default Disabled).

Offload Options

Allows you to specify which TCP/IP offload settings are handled by the adapter rather than the operating system.

Enabling offloading services increases transmission performance as the offload tasks are performed by the adapter hardware rather than the operating system. Thus, freeing CPU resources to work on other tasks.

IPv4 Checksums Offload

Enables the adapter to compute IPv4 checksum upon transmit and/or receive instead of the CPU (default Enabled).

TCP/UDP Checksum Offload for IPv4 packets

Enables the adapter to compute TCP/UDP checksum over IPv4 packets upon transmit and/or receive instead of the CPU (default Enabled).

TCP/UDP Checksum Offload for IPv6 packets

Enables the adapter to compute TCP/UDP checksum over IPv6 packets upon transmit and/or receive instead of the CPU (default Enabled).

Large Send Offload (LSO)

Allows the TCP stack to build a TCP message up to 64KB long and sends it in one call down the stack. The adapter then re-segments the message into multiple TCP packets for transmission on the wire with each pack sized according to the MTU. This option offloads a large amount of kernel processing time from the host CPU to the adapter.

• IB Options

Configures parameters related to InfiniBand functionality.

SA Query Retry Count

Sets the number of SA query retries once a query fails. The valid values are 1 - 64 (default 10).

SA Query Timeout

Sets the waiting timeout (in millisecond) of an SA query completion. The valid values are 500 - 60000 (default 1000 ms).



This document describes how to modify Windows registry parameters in order to improve performance.

Please note that modifying the registry incorrectly might lead to serious problems, including the loss of data, system hang, and you may need to reinstall Windows. As such it is recommended to back up the registry on your system before implementing recommendations included in this document. If the modifications you apply lead to serious problems, you will be able to restore the original registry state. For more details about backing up and restoring the registry, please visit www.microsoft.com.

6 OpenSM - Subnet Manager

OpenSM v3.3.11 is an InfiniBand Subnet Manager. For Mellanox WinOF VPI to operate, OpenSM must be running on at least one host machine in the InfiniBand cluster.



Please use the embedded OpenSM in the WinOF package for testing purpose and small cluster. Otherwise, we recommend using OpenSM from FabricIT EFMTM or UFMTM.

OpenSM can run as a Windows service which can be started manually from the following directory: <installation_directory>\tools. However, in WinOF 4.2 OpenSM is not defined as a service. OpenSM as a service will use the first port which is not in "down" state.

» To register it as a service run the following command:

```
sc create OpenSM1 binPath= "c:\Program Files\Mellanox\MLNX_VPI\IB\Tools\
opensm.exe --service" start=auto
```

» To start OpenSM as a service run:

sc start opensm

» To run OpenSM manually, enter on the command line:

opensm.exe

For additional run options, enter: opensm.exe -h

Notes

- For long term running, please avoid using the '-v' (verbosity) option to avoid exceeding disk quota.
- Running OpenSM on multiple servers may lead to incorrect OpenSM behavior. Please do not run more than a single instance of OpenSM in the subnet.

7 InfiniBand Fabric

7.1 Network Direct Interface

The Network Direct Interface (NDI) architecture provides application developers with a networking interface that enables zero-copy data transfers between applications, kernel-bypass I/O generation and completion processing, and one-sided data transfer operations.

NDI is supported by Microsoft and is the recommended method to write InfiniBand application. NDI exposes the advanced capabilities of the Mellanox networking devices and allows applications to leverage advances of InfiniBand.

For further information please refer to:

http://msdn.microsoft.com/en-us/library/cc904397(v=vs.85).aspx

7.2 part_man - Virtual IPoIB Port Creation Utility

part_man is used to create virtual ipoib ports. Currently, we allow to create only one virtual ipoib port per Mellanox IPoIB port. Each one of the virtual ports is created with PKey of 0xffff. part_man is used for adding/removing the virtual ports.

» Usage

part_man.exe [-v] <show|add|rem> ["Local area connection #"] [name]

- -v : increases verbosity level.
- Show: shows the currently configured virtual ipoib ports.
- Add : adds new virtual IPoIB port. Where add should be used with interface name as in Network connection in control panel.
- Name: any printable name without ':', ',' ',' '-' and ' ' and starting from i.
- Rem: removes existing virtual IPoIB port. Need to first run with Show, then copy the parameters.
- » Example

Adding and removing virtual port:

```
part_man add "Ethernet 4" ipoib_4_1
Done...
Part_man show
Ethernet 6 ipoib_4_1
part_man rem "Ethernet 6" ipoib_4_1
Done
```
7.3 InfiniBand Fabric Diagnostic Utilities

The diagnostic utilities described in this chapter provide means for debugging the connectivity and status of InfiniBand (IB) devices in a fabric. The tools are:

- Section 7.3.2, "ibportstate," on page 38
- Section 7.3.3, "ibroute," on page 42
- Section 7.3.4, "smpquery," on page 45
- Section 7.3.5, "perfquery," on page 49
- Section 7.3.6, "ibping," on page 53
- Section 7.3.7, "ibnetdiscover," on page 54
- Section 7.3.8, "ibtracert," on page 58
- Section 7.3.9, "sminfo," on page 60
- Section 7.3.10, "ibclearerrors," on page 61
- Section 7.3.11, "ibstat," on page 62
- Section 7.3.12, "vstat," on page 62
- Section 7.3.13, "osmtest," on page 63

7.3.1 Utilities Usage

This section first describes common configuration, interface, and addressing for all the tools in the package. Then it provides detailed descriptions of the tools themselves including: operation, synopsis and options descriptions, error codes, and examples.

7.3.1.1 Common Configuration, Interface and Addressing

Topology File (Optional)

An InfiniBand fabric is composed of switches and channel adapter (HCA/TCA) devices. To identify devices in a fabric (or even in one switch system), each device is given a GUID (a MAC equivalent). Since a GUID is a non-user-friendly string of characters, it is better to alias it to a meaningful, user-given name. For this objective, the IB Diagnostic Tools can be provided with a "topology file", which is an optional configuration file specifying the IB fabric topology in user-given names.

For diagnostic tools to fully support the topology file, the user may need to provide the local system name (if the local hostname is not used in the topology file).

To specify a topology file to a diagnostic tool use one of the following two options:

- 1. On the command line, specify the file name using the option '-t <topology file name>'
- 2. Define the environment variable IBDIAG_TOPO_FILE

To specify the local system name to an diagnostic tool use one of the following two options:

- 1. On the command line, specify the system name using the option '-s <local system name>'
- 2. Define the environment variable IBDIAG_SYS_NAME

7.3.1.2 IB Interface Definition

The diagnostic tools installed on a machine connect to the IB fabric by means of an HCA port through which they send MADs. To specify this port to an IB diagnostic tool use one of the following options:

- 1. On the command line, specify the port number using the option '-p <local port number>' (see below)
- 2. Define the environment variable IBDIAG_PORT_NUM

In case more than one HCA device is installed on the local machine, it is necessary to specify the device's index to the tool as well. For this use on of the following options:

- 1. On the command line, specify the index of the local device using the following option: '-i <index of local device>'
- 2. Define the environment variable IBDIAG DEV IDX

7.3.1.3 Addressing



This section applies to the ibdiagpath tool only. A tool command may require defining the destination device or port to which it applies.

The following addressing modes can be used to define the IB ports:

• Using a Directed Route to the destination: (Tool option '-d')

This option defines a directed route of output port numbers from the local port to the destination.

• Using port LIDs: (Tool option '-l'):

In this mode, the source and destination ports are defined by means of their LIDs. If the fabric is configured to allow multiple LIDs per port, then using any of them is valid for defining a port.

• Using port names defined in the topology file: (Tool option '-n')

This option refers to the source and destination ports by the names defined in the topology file. (Therefore, this option is relevant only if a topology file is specified to the tool.) In this mode, the tool uses the names to extract the port LIDs from the matched topology, then the tool operates as in the '-1' option.

7.3.2 ibportstate

Enables querying the logical (link) and physical port states of an InfiniBand port. It also allows adjusting the link speed that is enabled on any InfiniBand port.

If the queried port is a *swich* port, then ibportstate can be used to

- disable, enable or reset the port
- validate the port's link width and speed against the peer port

7.3.2.1 ibportstate Applicable Hardware

All InfiniBand devices.

7.3.2.2 ibportstate Synopsis

```
ibportstate [-d] [-e] [-v] [-V] [-D] [-L] [-G] [-s <smlid>] \ [-C
<ca_name>] [-P <ca_port>] [-u] [-t <timeout_ms>] \ [<dest
dr_path|lid|guid>] <portnum> [<op> [<value>]]
```

7.3.2.3 ibportstate Options

The table below lists the various flags of the command.

Table 5 -	ibportstate	Flags and	Options
-----------	-------------	-----------	---------

Flag	Description
-h/help	Print the help menu
-d/debug	Raise the IB debug level. May be used several times for higher debug levels (-ddd or -d -d -d)
-e/errors	Show send and receive errors (timeouts and others)
-v/verbose	Increase verbosity level. May be used several times for additional verbosity (-vvv or -v -v -v)
-V/version	Show version info
-D/Direct	Use directed path address arguments. The path is a comma separated list of out ports. Examples: '0' – self port '0,1,2,1,4' – out via port 1, then 2,
-L/Lid	Use Lid address argument
-G/Guid	Use GUID address argument. In most cases, it is the Port GUID. Example: '0x08f1040023'
-s/sm_port	Use <smlid> as the target lid for SM/SA queries</smlid>
-C/Ca	Use the specified channel adapter or router
-P/Port	Use the specified port
-u/usage	Usage message
-t/timeout	Override the default timeout for the solicited MADs [msec]
<dest dr_path="" guid="" lid="" =""></dest>	Destination's directed path, LID, or GUID.
<portnum></portnum>	Destination's port number
<op> [<value>]</value></op>	Define the allowed port operations: enable, disable, reset, speed, and query

In case of multiple channel adapters (CAs) or multiple ports without a CA/port being specified, a port is chosen by the utility according to the following criteria:

1. The first ACTIVE port that is found.

2. If not found, the first port that is UP (physical link state is LinkUp).

Examples

1. Query the status of Port 1 of CA mlx4_0 (using ibstatus) and use its output (the LID – 3 in this case) to obtain additional link information using ibportstate.

```
> ibstatus mlx4 0:1
Infiniband device 'mlx4 0' port 1 status:
      default gid:
                   fe80:0000:0000:0000:0000:0000:9289:3895
      base lid:
                   0x3
      sm lid:
                   0x3
      state:
                   2: INIT
      phys state: 5: LinkUp
              20 Gb/sec (4X DDR)
      rate:
> ibportstate -C mlx4 0 3 1 query
PortInfo:
# Port info: Lid 3 port 1
LinkState:....Initialize
PhysLinkState:....LinkUp
LinkWidthSupported:.....1X or 4X
LinkWidthEnabled:....1X or 4X
LinkWidthActive:.....4X
LinkSpeedSupported:.....2.5 Gbps or 5.0 Gbps
LinkSpeedEnabled:.....2.5 Gbps or 5.0 Gbps
LinkSpeedActive:.....5.0 Gbps
```

2. Query the status of two channel adapters using directed paths.

```
> ibportstate -C mlx4_0 -D 0 1
PortInfo:
# Port info: DR path slid 65535; dlid 65535; 0 port 1
LinkState:.....Initialize
PhysLinkState:....LinkUp
LinkWidthSupported:.....1X or 4X
LinkWidthEnabled:.....4X
LinkWidthActive:.....4X
LinkSpeedSupported:.....2.5 Gbps or 5.0 Gbps
```

```
Rev 4.2
```

LinkSpeedEnabled:.....2.5 Gbps or 5.0 Gbps LinkSpeedActive:....5.0 Gbps > ibportstate -C mthca0 -D 0 1 PortInfo: # Port info: DR path slid 65535; dlid 65535; 0 port 1 LinkState:....Down PhysLinkState:....Down PhysLinkState:....Polling LinkWidthSupported:.....1X or 4X LinkWidthEnabled:.....4X LinkWidthActive:.....4X LinkSpeedSupported:.....2.5 Gbps LinkSpeedActive:.....2.5 Gbps

3. Change the speed of a port.

First query for current configuration

```
> ibportstate -C mlx4_0 -D 0 1
PortInfo:
# Port info: DR path slid 65535; dlid 65535; 0 port 1
LinkState:.....Initialize
PhysLinkState:....LinkUp
LinkWidthSupported:.....1X or 4X
LinkWidthEnabled:.....4X
LinkWidthActive:....4X
LinkSpeedSupported:....2.5 Gbps or 5.0 Gbps
LinkSpeedEnabled:....5.0 Gbps
# Now change the enabled link speed
> ibportstate -C mlx4_0 -D 0 1 speed 2
```

Initial PortInfo:

Port info: DR path slid 65535; dlid 65535; 0 port 1

LinkSpeedEnabled:.....2.5 Gbps

ibportstate -C mlx4_0 -D 0 1 speed 2

After PortInfo set: # Port info: DR path slid 65535; dlid 65535; 0 port 1 LinkSpeedEnabled:.....5.0 Gbps (IBA extension) # Show the new configuration > ibportstate -C mlx4 0 -D 0 1 PortInfo: # Port info: DR path slid 65535; dlid 65535; 0 port 1 LinkState:....Initialize PhysLinkState:....LinkUp LinkWidthSupported:....1X or 4X LinkWidthEnabled:....1X or 4X LinkWidthActive:.....4X LinkSpeedSupported:.....2.5 Gbps or 5.0 Gbps LinkSpeedEnabled:.....5.0 Gbps (IBA extension) LinkSpeedActive:.....5.0 Gbps

7.3.3 ibroute

Uses SMPs to display the forwarding tables for unicast (LinearForwardingTable or LFT) or multicast (MulticastForwardingTable or MFT) for the specified switch LID and the optional lid (mlid) range. The default range is all valid entries in the range of 1 to FDBTop.

7.3.3.1 ibroute Applicable Hardware

InfiniBand switches.

7.3.3.2 ibroute Synopsis

```
ibroute [-h] [-d] [-V] [-U] [-n] [-D] [-G] [-M] [-L] [-e] [-u] [-s <smlid>] (-C < ca_name)
[-P <ca_port>] [ -t <timeout_ms>] \
                                                            [<dest dr path|lid|guid> [<star-
tlid> [<endlid>]]]
```

7.3.3.3 ibroute Options

The table below lists the various ibroute flags of the command.

Table 6 - ibroute Flags and Options

Flag	Description
-h/help	Print the help menu

Table 6 -	ibroute	Flags	and	Options
-----------	---------	-------	-----	---------

Flag	Description
-d/debug	Raise the IB debug level. May be used several times for higher debug levels (-ddd or -d -d -d)
-a/all	Show all LIDs in range, including invalid entries
-v/verbose	Increase verbosity level. May be used several times for additional verbosity (-vvv or -v -v -v)
-V/version	Show version info
-n/no_dests	Do not try to resolve destinations
-D/Direct	Use directed path address arguments. The path is a comma separated list of out ports. Examples: '0' - self port '0,1,2,1,4' - out via port 1, then 2,
-G/Guid	Use GUID address argument. In most cases, it is the Port GUID. Example: '0x08f1040023'
-M/Multicast	Show multicast forwarding tables. The parameters <startlid> and <endlid> spec- ify the MLID range.</endlid></startlid>
-L/Lid	Use Lid address argument
-u/usage	Usage message
-e/errors	Show send and receive errors (timeouts and others)
-s/sm_port <smlid></smlid>	Use <smlid> as the target LID for SM/SA queries</smlid>
-C/Ca <ca_name></ca_name>	Use the specified channel adapter or router
-P/Port <ca_port></ca_port>	Use the specified port
-t/timeout <timeout_ms></timeout_ms>	Override the default timeout for the solicited MADs [msec]
<dest dr_path="" guid="" lid="" =""></dest>	Destination's directed path, LID, or GUID
<startlid></startlid>	Starting LID in an MLID range
<endlid></endlid>	Ending LID in an MLID range

Examples

1. Dump all Lids with valid out ports of the switch with Lid 2.

```
> ibroute 2
Unicast lids [0x0-0x8] of switch Lid 2 guid 0x0002c902fffff00a (MT47396 Infiniscale-III Mellanox
Technologies):
Lid Out Destination
Port Info
0x0002 000 : (Switch portguid 0x0002c902fffff00a: 'MT47396 Infiniscale-III Mellanox Technolo-
gies')
```

```
Rev 4.2
```

```
0x0003 021 : (Switch portguid 0x000b8cffff004016: 'MT47396 Infiniscale-III Mellanox Technolo-
gies')
0x0006 007 : (Channel Adapter portguid 0x0002c90300001039: 'sw137 HCA-1')
0x0007 021 : (Channel Adapter portguid 0x0002c9020025874a: 'sw157 HCA-1')
0x0008 008 : (Channel Adapter portguid 0x0002c902002582cd: 'sw136 HCA-1')
5 valid lids dumped
```

2. Dump all Lids in the range 3 to 7 with valid out ports of the switch with Lid 2.

```
> ibroute 2 3 7
Unicast lids [0x3-0x7] of switch Lid 2 guid 0x0002c902fffff00a (MT47396 Infiniscale-III Mellanox
Technologies):
Lid Out Destination
Port Info
0x0003 021 : (Switch portguid 0x000b8cffff004016: 'MT47396 Infiniscale-III Mellanox Technolo-
gies')
0x0006 007 : (Channel Adapter portguid 0x0002c90300001039: 'sw137 HCA-1')
0x0007 021 : (Channel Adapter portguid 0x0002c9020025874a: 'sw157 HCA-1')
3 valid lids dumped
```

3. Dump all Lids with valid out ports of the switch with portguid 0x000b8cfff004016.

```
> ibroute -G 0x000b8cffff004016
Unicast lids [0x0-0x8] of switch Lid 3 guid 0x000b8cffff004016 (MT47396 Infiniscale-III Mellanox
Technologies):
Lid Out Destination
Port Info
0x0002 023 : (Switch portguid 0x0002c902fffff00a: 'MT47396 Infiniscale-III Mellanox Technolo-
gies')
0x0003 000 : (Switch portguid 0x000b8cffff004016: 'MT47396 Infiniscale-III Mellanox Technolo-
gies')
0x0006 023 : (Channel Adapter portguid 0x0002c90300001039: 'sw137 HCA-1')
0x0007 020 : (Channel Adapter portguid 0x0002c9020025874a: 'sw157 HCA-1')
0x0008 024 : (Channel Adapter portguid 0x0002c902002582cd: 'sw136 HCA-1')
5 valid lids dumped
```

4. Dump all non-empty mlids of switch with Lid 3.

Ports: 0 1 2 3 4	5 6 7 8 9 0 1 2 3 4 5 6 7 8 9	0 1 2 3 4
MLid		
0xc000		Х
0xc001		Х
0xc002		Х
0xc003		Х
0xc020	Х	
0xc021	Х	
0xc022	Х	
0xc023	Х	
0xc024	Х	
0xc040	Х	
0xc041	Х	
0xc042	X	
12 valid mlids dumped		

7.3.4 smpquery

Provides a basic subset of standard SMP queries to query Subnet management attributes such as node info, node description, switch info, and port info.

7.3.4.1 smpquery Applicable Hardware

All InfiniBand devices.

7.3.4.2 smpquery Synopsys

smpquery [-h] [-d] [-e] [-c] [-v] [-D] [-G] [-s <smlid>] [-L] [-u] [-V] [-C <ca_name>] [P <ca_port>] [-t <timeout_ms>] [--node-name-map <node-name-map>] <op> <dest
dr path|lid|guid> [op params]

7.3.4.3 smpquery Options

Table 7 - smpquery	Flags and	Options
--------------------	-----------	---------

Flag	Description
-h/help	Print the help menu
-d/debug	Raise the IB debug level. May be used several times for higher debug levels (-ddd or -d -d)
-e/errors	Show send and receive errors (timeouts and others)

Table 7 - smpquer	y Flags and	Options
-------------------	-------------	---------

Flag	Description
-v/verbose	Increase verbosity level. May be used several times for additional verbos- ity (-vvv or -v -v -v)
-D/Direct	Use directed path address arguments. The path is a comma separated list of out ports. Examples: '0' - self port '0,1,2,1,4' - out via port 1, then 2,
-G/Guid	Use GUID address argument. In most cases, it is the Port GUID. Example: '0x08f1040023'
-s/sm_port <smlid></smlid>	Use <smlid> as the target LID for SM/SA queries</smlid>
-V/version	Show version info
-L/Lid	Use Lid address argument
-c/combined	Use combined route address argument
-u/usage	Usage message
-C/Ca <ca_name></ca_name>	Use the specified channel adapter or router
-P/Port <ca_port></ca_port>	Use the specified port
-t/timeout <timeout_ms></timeout_ms>	Override the default timeout for the solicited MADs [msec]
<op></op>	Supported operations: NodeInfo (NI) <addr> NodeDesc (ND) <addr> PortInfo (PI) <addr> [<portnum>] SwitchInfo (SI) <addr> PKeyTable (PKeys) <addr> [<portnum>] SL2VLTable (SL2VL) <addr> [<portnum>] VLArbitration (VLArb) <addr> [<portnum>] GUIDInfo (GI) <addr></addr></portnum></addr></portnum></addr></portnum></addr></addr></portnum></addr></addr></addr>
<dest dr_path="" guid="" lid="" =""></dest>	Destination's directed path, LID, or GUID
node-name-map <file></file>	Node name map file
-x/extended	Use extended speeds

Examples

1. Query PortInfo by LID, with port modifier.

SMLid:0x0001
CapMask:0x251086a
IsSM
IsTrapSupported
IsAutomaticMigrationSupported
IsSLMappingSupported
IsSystemImageGUIDsupported
IsCommunicatonManagementSupported
IsVendorClassSupported
IsCapabilityMaskNoticeSupported
IsClientRegistrationSupported
DiagCode:0x0000
MkeyLeasePeriod:0
LocalPort:1
LinkWidthEnabled:1X or 4X
LinkWidthSupported:1X or 4X
LinkWidthActive:4X
LinkSpeedSupported:2.5 Gbps or 5.0 Gbps
LinkState:Active
PhysLinkState:LinkUp
LinkDownDefState:Polling
ProtectBits:0
LMC:0
LinkSpeedActive:5.0 Gbps
LinkSpeedEnabled:2.5 Gbps or 5.0 Gbps
NeighborMTU:2048
SMSL:0
VLCap:VL0-7
InitType:0x00
VLHighLimit:4
VLArbHighCap:
VLArbLowCap:
InitReply:0x00
MtuCap:

VLStallCount:....0 OperVLs:.....VL0-3 PartEnforceInb:....0 PartEnforceOutb:....0 FilterRawInb:....0 FilterRawOutb:....0 MkeyViolations:....0 PkeyViolations:....0 QkeyViolations:....0 GuidCap:.....128 ClientReregister:....0 SubnetTimeout:.....18 RespTimeVal:.....16 LocalPhysErr:.....8 OverrunErr:.....8 MaxCreditHint:....0 RoundTrip:....0

2. Query SwitchInfo by GUID.

> smpquery -G switchinfo 0x000b8cffff004016
Switch info: Lid 3
LinearFdbCap:......49152
RandomFdbCap:.....0
McastFdbCap:....0
McastFdbTop:....8
DefPort:...0
DefPort:...0
DefMcastPrimPort:...0
LifeTime:....18
StateChange:...0
LidsPerPort:...0
PartEnforceCap:....32
InboundPartEnf:....1
OutboundPartEnf:....1

FilterRawInbound:.....1 FilterRawOutbound:.....1 EnhancedPort0:....0

3. Query NodeInfo by direct route.

> smpquery -D nodeinfo 0
Node info: DR path slid 65535; dlid 65535; 0
BaseVers:1
ClassVers:1
NodeType:Channel Adapter
NumPorts:2
SystemGuid:0x0002c9030000103b
Guid:0x0002c90300001038
PortGuid:0x0002c90300001039
PartCap:128
DevId:0x634a
Revision:0x000000a0
LocalPort:1
VendorId:0x0002c9

7.3.5 perfquery

Queries InfiniBand ports' performance and error counters. Optionally, it displays aggregated counters for all ports of a node. It can also reset counters after reading them or simply reset them.

7.3.5.1 perfquery Applicable Hardware

All InfiniBand devices.

7.3.5.2 perfquery Synopsys

perfquery [-h] [-d] [-G] [--xmtsl, -X] [--xmtdisc, -D] [--rcvsl, -S] [--rcverr, -E] [--smplctl, c] [-a] [--Lid, -L] [--sm_port, -s <lid>] [-errors, -e] [--verbose, -v] [--usage, -u][-1] [-r] [C <ca_name>] [-P <ca_port>] [-R][-t <timeout_ms>] [-V] [<lid|guid> [[port][reset_mask]]]

The table below lists the various flags of the command.

Table 8 - perfquery Flags and Options

Flag	Description
help, -h	Print the help menu

Table 8 -	perfquery	Flags and	Options
-----------	-----------	-----------	---------

Flag	Description
debug, -d	Raise the IB debug level. May be used several times for higher debug levels (- ddd or -d -d)
Guid,-G	Use GUID address argument. In most cases, it is the Port GUID. Example: '0x08f1040023'
xmtsl, -X	Show Xmt SL port counters
revsl, -S	Show Rcv SL port counters
xmtdisc, -D	Show Xmt Discard Details
rcverr, -E	Show Rcv Error Details
smplctl, -c	Show samples control
all_ports, -a	Apply query to all ports
Lid, -L	Use LID address argument
sm_port, -s <lid></lid>	SM port lid
errors, -e	Show send and receive errors
verbose, -v	Increase verbosity level
usage, -u	Usage message
loop_ports, -l	Loop ports
reset_after_read, -r	Reset the counters after reading them
Ca, -C <ca_name></ca_name>	Use the specified channel adapter or router
Port, -P <ca_port></ca_port>	Use the specified port
Reset_only, -R	Reset the counters
timeout, -t <timeout_ms></timeout_ms>	Override the default timeout for the solicited MADs [msec]
version, -V	Show version info
<lid guid="" =""> [[port][reset_mask]]</lid>	LID or GUID
extended, -x	show extended port counters
extended_speeds, -T	show port extended speeds counters
oprevcounters	show Rev Counters per Op code
flowctlcounters	show flow control counters
vloppackets	show packets received per Op code per VL
vlopdata	show data received per Op code per VL
vlxmitflowctlerrors	show flow control update errors per VL
vlxmitcounters	show ticks waiting to transmit counters per VL
swportvlcong	show sw port VL congestion

Table 8 - perfquery	[,] Flags and	Options
---------------------	------------------------	---------

Flag	Description
rcvcc	show Rcv congestion control counters
slrcvfecn	show SL Rev FECN counters
slrcvbecn	show SL Rev BECN counters
xmitce	show Xmit congestion control counters
vlxmittimecc	show VL Xmit Time congestion control counters

Examples

```
perfquery -r 32 1  # read performance counters and reset
perfquery -e -r 32 1  # read extended performance counters and reset
perfquery -R 0x20 1  # reset performance counters of port 1 only
perfquery -e -R 0x20 1 # reset extended performance counters of port 1 only
perfquery -R -a 32  # reset performance counters of all ports
perfquery -R 32 2 0x0fff# reset only error counters of port 2
perfquery -R 32 2 0xf000# reset only non-error counters of port 2
```

1. Read local port's performance counters.



Rev 4.2

 XmtData:
 .55178210

 RcvData:
 .55174680

 XmtPkts:
 .766366

 RcvPkts:
 .766315

2. Read performance counters from LID 2, all ports.

> smpquery -a 2 # Port counters: Lid 2 port 255 CounterSelect:.....0x0100 SymbolErrors:.....65535 LinkRecovers:.....255 LinkDowned:.....16 RcvRemotePhysErrors:....0 RcvSwRelayErrors:.....70 XmtDiscards:.....488 XmtConstraintErrors:....0 RcvConstraintErrors:....0 LinkIntegrityErrors:....0 ExcBufOverrunErrors:.....0 VL15Dropped:.....0 XmtData:.....129840354 RcvData:.....129529906 XmtPkts:.....1803332 RcvPkts:.....1799018

3. Read then reset performance counters from LID 2, port 1.

```
> perfquery -r 2 1
# Port counters: Lid 2 port 1
PortSelect:.....0x0100
SymbolErrors:....0
LinkRecovers:....0
LinkDowned:....0
```

RcvErrors:0
RcvRemotePhysErrors:0
RcvSwRelayErrors:0
XmtDiscards:3
XmtConstraintErrors:0
RcvConstraintErrors:0
LinkIntegrityErrors:0
ExcBufOverrunErrors:0
VL15Dropped:0
XmtData:0
RcvData:0
XmtPkts:0
RcvPkts:0

7.3.6 ibping

ibping uses vendor MADs to validate connectivity between IB nodes. On exit, (IP) ping like output is shown. ibping is run as client/server, however the default is to run it as a client. Note also that in addition to ibping, a default server is implemented within the kernel.

7.3.6.1 ibping Synopsys

ibping [-d(ebug)] [-e(rr_show)] [-v(erbose)] [-G(uid)] [-C ca_name] [-P ca_port] [-s smlid] [t(imeout)timeout_ms] [-V(ersion)] [-L(id)][-u(sage)] [-c ping_count] [-f(lood)] [-o oui] [-S(erver)] [-h(elp)] <dest lid | guid>

7.3.6.2 ibping Options

Table 9 - ibping Flags and Options

Flag	Description
count, -c <num></num>	Stops after count packets
-f, (flood)	Floods destination: send packets back to back without delay
-o, (oui)	Uses specified OUI number to multiplex vendor mads
Server, -S	Starts in server mode (do not return)
debug, -d/-ddd/ -d -d -d	Raises the IB debugging level
errors, -e	Shows send and receive errors (timeouts and others)
help, -h	Shows the usage message

Flag	Description
verbose, -v/-vvv/-v -v -v	Increases the application verbosity level
version, -V	Shows the version info
Lid, -L	Use LID address argument
usage, -u	Usage message
Guid, -G	Uses GUID address argument. In most cases, it is the Port GUID. For example: "0x08f1040023"
sm_port, -s <smlid></smlid>	Uses 'smlid' as the target lid for SM/SA queries
Ca, -C <ca_name></ca_name>	Uses the specified ca_name
Port, -P <ca_port></ca_port>	Uses the specified ca_port
timeout, -t <timeout ms=""></timeout>	Overrides the default timeout for the solicited mads

Table 9 - ibping Flags and Options

7.3.7 ibnetdiscover

Rev 4.2

ibnetdiscover performs IB subnet discovery and outputs a readable topology file. GUIDs, node types, and port numbers are displayed as well as port LIDs and NodeDescriptions. All nodes (and links) are displayed (full topology). Optionally, this utility can be used to list the current connected nodes by node-type. The output is printed to standard output unless a topology file is specified.

7.3.7.1 ibnetdiscover Synopsys

ibnetdiscover [-d(ebug)] [-e(rr_show)] [-v(erbose)] [-s(how)] [-l(ist)] [-g(rouping)] [-H(ca_list)][-S(witch_list)] [-R(outer_list)] [-C ca_name] [-P ca_port] [-t(imeout) timeout_ms] [-V(ersion)] [--outstanding_smps -o <val>] [-u(sage)] [--node-name-map <node-name-map>] [-cache <filename>] [--load-cache <filename>] [-p(orts)] [-m(ax_hops)]

[-h(elp)] [<topology-file>]

7.3.7.2 ibnetdiscover Options

The table below lists the various flags of the command.

Most OpenIB diagnostics take the following common flags. The exact list of supported flags per utility can be found in the usage message and can be shown using the util name -h syntax.

Flag	Description
-l,list	List of connected nodes
-g,grouping	Show grouping. Grouping correlates IB nodes by different vendor specific schemes. It may also show the switch external ports correspondence.
-H,Hca_list	List of connected CAs

Table 10 - ibnetdiscover Flags and Options

Table 10 - ibnetdiscover Flags and Options

Flag	Description
-S,Switch_list	List of connected switches
-R,Router_list	List of connected routers
-s,show	Show progress information during discovery
node-name-map <node-name- map></node-name- 	Specify a node name map. The node name map file maps GUIDs to more user friendly names. See "Topology File Format" on page 56.
cache <filename></filename>	Cache the ibnetdiscover network data in the specified filename. This cache may be used by other tools for later analysis
load-cache <filename></filename>	Load and use the cached ibnetdiscover data stored in the specified filename. May be useful for outputting and learning about other fabrics or a previous state of a fabric
diff <filename></filename>	Load cached ibnetdiscover data and do a diff comparison to the current net- work or another cache. A special diff output for ibnetdiscover output will be displayed showing differences between the old and current fabric. By default, the following are compared for differences: switches, channel adapt- ers, routers, and port connections
diffcheck <key(s)></key(s)>	Specify what diff checks should be done in thediff option above. Comma separate multiple diff check key(s). The available diff checks are: sw = switches, ca = channel adapters, router = routers, port = port connections, lid = lids, nodedesc = node descriptions. Note that port, lid, and nodedesc are checked only for the node types that are specified (e.g. sw, ca, router). If port is specified alongside lid or nodedesc, remote port lids and node descriptions will also be com-pared
-p,ports	Obtain a ports report which is a list of connected ports with relevant informa- tion (like LID, port-num, GUID, width, speed, and NodeDescription)
-m,max_hops	Report max hops discovered
debug, -d/-ddd/ -d -d -d	Raise the IB debugging level
errors, -e	Show send and receive errors (timeouts and others)
help, -h	Show the usage message
verbose,-v/-vv/ -v -v -v	Increase the application verbosity level
version, -V	Show the version info
outstanding_smps -o <val></val>	Specify the number of outstanding SMPs which should be issued during the scan
-usage, -u	Usage message
Ca, -C <ca_name></ca_name>	Use the specified ca_name
Port, -P <ca_port></ca_port>	Use the specified ca_port
timeout, -t <timeout_ms></timeout_ms>	Override the default timeout for the solicited mads
full, -f	show full information (ports' speed and width)
show, -s	show more information

7.3.7.3 Topology File Format

The topology file format is largely intuitive. Most identifiers are given textual names like vendor ID (vendid), device ID (device ID), GUIDs of various types (sysimgguid, caguid, switchguid, etc.). PortGUIDs are shown in parentheses (). For switches, this is shown on the switchguid line. For CA and router ports, it is shown on the connectivity lines. The IB node is identified followed by the number of ports and the node GUID. On the right of this line is a comment (#) followed by the NodeDescription in quotes. If the node is a switch, this line also contains whether switch port 0 is base or enhanced, and the LID and LMC of port 0. Subsequent lines pertaining to this node show the connectivity. On the left is the port number of the current node. On the right is the peer node (node at other end of link). It is identified in quotes with nodetype followed by - followed by NodeGUID with the port number in square brackets. Further on the right is a comment (#). What follows the comment is dependent on the node type. If it it a switch node, it is followed by the NodeDescription in quotes and the LID of the peer node. If it is a CA or router node, it is followed by the local LID and LMC and then followed by the NodeDescription in quotes and the LID of the peer node. The active link width and speed are then appended to the end of this output line.

Example

```
# Topology file: generated on Tue Jun 5 14:15:10 2007
#
# Max of 3 hops discovered
# Initiated from node 0008f10403960558 port 0008f10403960559
```

Non-Chassis Nodes

When grouping is used, IB nodes are organized into chasses which are numbered. Nodes which cannot be determined to be in a chassis are displayed as "Non-Chassis Nodes". External ports are also shown on the connectivity lines.

```
vendid=0x8f1
devid=0x5a06
sysimgguid=0x5442ba00003000
switchguid=0x5442ba00003080(5442ba00003080)
Switch 24 "S-005442ba00003080"
                                        # "ISR9024 Voltaire" base port 0 lid 6 lmc 0
       "H-0008f10403961354"[1](8f10403961355)
                                                       # "MT23108 InfiniHost Mellanox Technolo-
gies" lid 4 4xSDR
[10]
        "S-0008f10400410015"[1]
                                        # "SW-6IB4 Voltaire" lid 3 4xSDR
        "H-0008f10403960558"[2](8f1040396055a)
                                                       # "MT23108 InfiniHost Mellanox Technolo-
[8]
gies" lid 14 4xSDR
[6]
        "S-0008f10400410015"[3]
                                        # "SW-6IB4 Voltaire" lid 3 4xSDR
        "H-0008f10403960558"[1](8f10403960559)
                                                       # "MT23108 InfiniHost Mellanox Technolo-
[12]
gies" lid 10 4xSDR
vendid=0x8f1
```

```
devid=0x5a05
switchguid=0x8f10400410015(8f10400410015)
Switch 8 "S-0008f10400410015" # "SW-6IB4 Voltaire" base port 0 lid 3 lmc 0
[6] "H-0008f10403960984"[1](8f10403960985) # "MT23108 InfiniHost Mellanox Technolo-
gies" lid 16 4xSDR
      "H-005442b100004900"[1](5442b100004901)
                                             # "MT23108 InfiniHost Mellanox Technolo-
[4]
gies" lid 12 4xSDR
[1] "S-005442ba00003080"[10] # "ISR9024 Voltaire" lid 6 1xSDR
[3] "S-005442ba00003080"[6] # "ISR9024 Voltaire" lid 6 4xSDR
vendid=0x2c9
devid=0x5a44
caguid=0x8f10403960984
Ca 2 "H-0008f10403960984" # "MT23108 InfiniHost Mellanox Technologies"
[1](8f10403960985) "S-0008f10400410015"[6] # lid 16 lmc 1 "SW-6IB4 Voltaire" lid 3
4xSDR
vendid=0x2c9
devid=0x5a44
caguid=0x5442b100004900
Ca 2 "H-005442b100004900"
                                # "MT23108 InfiniHost Mellanox Technologies"
[1](5442b100004901) "S-0008f10400410015"[4] # lid 12 lmc 1 "SW-6IB4 Voltaire" lid 3
4xSDR
vendid=0x2c9
devid=0x5a44
caguid=0x8f10403961354
Ca 2 "H-0008f10403961354"
                                # "MT23108 InfiniHost Mellanox Technologies"
[1](8f10403961355) "S-005442ba00003080"[22] # lid 4 lmc 1 "ISR9024 Voltaire"
lid 6 4xSDR
vendid=0x2c9
devid=0x5a44
caguid=0x8f10403960558
Ca 2 "H-0008f10403960558"
                                # "MT23108 InfiniHost Mellanox Technologies"
[2](8f1040396055a) "S-005442ba00003080"[8] # lid 14 lmc 1 "ISR9024 Voltaire" lid 6
4xSDR
[1](8f10403960559) "S-005442ba00003080"[12] # lid 10 lmc 1 "ISR9024 Voltaire"
lid 6 1xSDR
```

Node Name Map File Format

The node name map is used to specify user friendly names for nodes in the output. GUIDs are used to perform the lookup.

comment
<guid> "<name>"

Example

# IB1			
# Line cards			
0x0008f104003f125c	"IB1	(Rack 11 slot 1) ISR9288/ISR9096 Voltaire sLB-24D"	
0x0008f104003f125d	"IB1	(Rack 11 slot 1) ISR9288/ISR9096 Voltaire sLB-24D"	
0x0008f104003f10d2	"IB1	(Rack 11 slot 2) ISR9288/ISR9096 Voltaire sLB-24D"	
0x0008f104003f10d3	"IB1	(Rack 11 slot 2) ISR9288/ISR9096 Voltaire sLB-24D"	
0x0008f104003f10bf	"IB1	(Rack 11 slot 12) ISR9288/ISR9096 Voltaire sLB-24D"	
# Spines			
0x0008f10400400e2d	"IB1	(Rack 11 spine 1) ISR9288 Voltaire sFB-12D"	
0x0008f10400400e2e	"IB1	(Rack 11 spine 1) ISR9288 Voltaire sFB-12D"	
0x0008f10400400e2f	"IB1	(Rack 11 spine 1) ISR9288 Voltaire sFB-12D"	
0x0008f10400400e31	"IB1	(Rack 11 spine 2) ISR9288 Voltaire sFB-12D"	
0x0008f10400400e32	"IB1	(Rack 11 spine 2) ISR9288 Voltaire sFB-12D"	
# GUID Node Name			
0x0008f10400411a08	"SW1	(Rack 3) ISR9024 Voltaire 9024D"	
0x0008f10400411a28	"SW2	(Rack 3) ISR9024 Voltaire 9024D"	
0x0008f10400411a34	"SW3	(Rack 3) ISR9024 Voltaire 9024D"	
0x0008f104004119d0	"SW4	(Rack 3) ISR9024 Voltaire 9024D"	

7.3.8 ibtracert

58

ibtracert uses SMPs to trace the path from a source GID/LID to a destination GID/LID. Each hop along the path is displayed until the destination is reached or a hop does not respond. By using the -m option, multicast path tracing can be performed between source and destination nodes.

7.3.8.1 ibtracert Synopsys

```
ibtracert [-d(ebug)] [-v(erbose)] [-D(irect)] [-L(id)] [-e(rrors)] [-u(sage)] [-G(uids)] [-
f(orce)] [-n(o_info)] [-m mlid] [-s smlid] [-C ca_name][-P ca_port] [-t(imeout) timeout_ms] [-
V(ersion)] [--node-name--map <node-name-map>] [-h(elp)] [<dest dr_path|lid|guid> [<startlid>
[<endlid>]]
```

7.3.8.2 ibtracert Options

The table below lists the various flags of the command.

Most OpenIB diagnostics take the following common flags. The exact list of supported flags per utility can be found in the usage message and can be shown using the util_name -h syntax..

Flag	Description
force, -f	Force
-n,no_info	Simple format; do not show additional information
mlid, -m <mlid></mlid>	Show the multicast trace of the specified mlid
node-name-map <node-name- map></node-name- 	Specify a node name map. The node name map file maps GUIDs to more user friendly names. See "Topology File Format" on page 56.
debug, -d/-ddd/-d -d -d	Raise the IB debugging level
Lid, -L	Use LID address argument
errors, -e	Show send and receive errors
usage, -u	Usage message
Guid, -G	Use GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023"
sm_port, -s <smlid></smlid>	Use 'smlid' as the target lid for SM/SA queries
help, -h	Show the usage message
-verbose, -v/-vv/-v -v -v	Increase the application verbosity level
version, -V	Show the version info
Ca, -C <ca_name></ca_name>	Use the specified ca_name
Port, -P <ca_port></ca_port>	Use the specified ca_port
timeout, -t <timeout_ms></timeout_ms>	Override the default timeout for the solicited mads

Table 11 - ibtracert Flags and Options

Examples

• Unicast examples

ibtracert 4 16	# show path between lids 4 and 16
ibtracert -n 4 16	<pre># same, but using simple output format</pre>

ibtracert -G 0x8f1040396522d 0x002c9000100d051 # use guid addresses

• Multicast example

ibtracert -m 0xc000 4 16 # show multicast path of mlid 0xc000 between lids 4 and 16

7.3.9 sminfo

Optionally sets and displays the output of a sminfo query in a readable format. The target SM is the one listed in the local port info, or the SM specified by the optional SM lid or by the SM direct routed path.



Using sminfo for any purposes other then simple query may result in a malfunction of the target SM.

7.3.9.1 sminfo Synopsys

sminfo [-d(ebug)] [-e(rr_show)] [-s state] [-p prio] [-a activity] [-D(irect)] [-L(id)] [u(sage)] [-G(uid)] [-C ca_name] [-P ca_port] [-t(imeout) timeout_ms] [-V(ersion)] [-h(elp)] sm_lid | sm_dr_path [modifier]

7.3.9.2 sminfo Options

The table below lists the various flags of the command.

Most OpenIB diagnostics take the following common flags. The exact list of supported flags per utility can be found in the usage message and can be shown using the util_name -h syntax..

Flag	Description
state, -s	Set SM state • 0 - not active • 1 - discovering • 2 - standby • 3 - master
priority, -p	Set priority (0-15)
activity, -a	Set activity count
debug, -d/-ddd/-d -d -d	Raise the IB debugging level
Direct, -D	 Use directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0,1,2,1,4" # out via port 1, then 2,
Lid, -L	Use LID address argument

Table 12 - sminfo Flags and Options

Flag	Description
usage, -u	Usage message
errors, -e	Show send and receive errors (timeouts and others)
Guid, -G	Use GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023"
help, -h	Show the usage message
-verbose, -v/-vv/-v -v -v	Increase the application verbosity level
version, -V	Show the version info
Ca, -C <ca_name></ca_name>	Use the specified ca_name
Port, -P <ca_port></ca_port>	Use the specified ca_port
timeout, -t <timeout ms=""></timeout>	Override the default timeout for the solicited mads

Table 12 - sminfo Flags and Options

Examples

--timeout, -t <timeout ms>

sminfo	<pre># local ports sminfo</pre>
sminfo 32	<pre># show sminfo of lid 32</pre>
sminfo -G 0x8f1040023	<pre># same but using guid address</pre>

7.3.10 ibclearerrors

ibclearerrors is a script which clears the PMA error counters in PortCounters by either waking the IB subnet topology or using an already saved topology file.

7.3.10.1 ibclearerrors Synopsys

ibclearerrors [-h] [-N | -nocolor] [<topology-file> | -C ca_name -P ca_port -t(imeout) timeout_ms]

7.3.10.2 ibclearerrors Options

Table 13 - ibclearerrors Flags and Options

Flag	Description
-C <ca_name></ca_name>	Use the specified ca_name
-P <ca_port></ca_port>	Use the specified ca_port
-t <timeout_ms></timeout_ms>	Override the default timeout for the solicited mads

7.3.11 ibstat

Rev 4.2

ibstat is a binary which displays basic information obtained from the local IB driver. Output includes LID, SMLID, port state, link width active, and port physical state.

7.3.11.1 ibstat Synopsys

ibstat [-d(ebug)] [-l(ist_of_cas)] [-s(hort)] [-p(ort_list)] [-V(ersion)] [-h] <ca_name> [portnum]

7.3.11.2 ibstat Options

The table below lists the various flags of the command.

Most OpenIB diagnostics take the following common flags. The exact list of supported flags per utility can be found in the usage message and can be shown using the util_name -h syntax..

Table 14 - ibstat Flags and Options

Flag	Description
-l,list_of_cas	List all IB devices
-s,short	Short output
-p,port_list	Show port list
ca_name	InfiniBand device name
portnum	Port number of InfiniBand device
debug, -d/-ddd/-d -d -d	Raise the IB debugging level
help, -h	Show the usage message
-verbose, -v/-vv/-v -v -v	Increase the application verbosity level
version, -V	Show the version info
usage, -u	usage message

Examples

ibstat	<pre># display status of all ports on all IB devices</pre>
ibstat -l	# list all IB devices
ibstat -p	# show port guids
ibstat mthca0 2	<pre># show status of port 2 of 'mthca0'</pre>

7.3.12 vstat

vstat is a binary which displays information on the HCA attributes.

• vstat synopsys is

vstat [-v] [-c]

7.3.12.1 vstat Options

The table below lists the various flags of the command...

Table 15 -	vstat	Flags	and	Options
------------	-------	-------	-----	---------

Flag	Description
-V -	Verbose mode
-c	HCA error/statistic counters
-m	more verbose mode
-p N	repeat every N sec

7.3.13 osmtest

osmtest is a test program to validate InfiniBand subnet manager and administration (SM/SA). Default is to run all flows with the exception of the QoS flow. osmtest provides a test suite for opensm. osmtest has the following capabilities and testing flows:

- It creates an inventory file of all available Nodes, Ports, and PathRecords, including all their fields.
- It verifies the existing inventory, with all the object fields, and matches it to a presaved one.
- A Multicast Compliancy test.
- An Event Forwarding test.
- A Service Record registration test.
- An RMPP stress test.
- A Small SA Queries stress test.

It is recommended that after installing opensm, the user should run "osmtest -f c" to generate the inventory file, and immediately afterwards run "osmtest -f a" to test OpenSM.

Additionally, it is recommended to create the inventory when the IB fabric is stable, and occasionally run "osmtest -v" to verify that nothing has changed.

7.3.13.1 osmtest Synopsys

```
osmtest [-f(low) <c|a|v|s|e|f|m|q|t>] [-w(ait) <trap_wait_time>] [-d(ebug) <number>] [-
m(ax_lid) <LID in hex>] [-g(uid) [=]<GUID in hex>] [-p(ort)] [-i(nventory) <filename>] [-
s(tress)] [-M(ulticast_Mode)] [-t(imeout) <milliseconds>] [-1 | --log_file] [-v] [-vf <flags>]
[-h(elp)]
```

7.3.13.2 osmtest Options

Table 16 - osmtest Flags and Options

Flag	Description
-f,flow	 This option directs osmtest to run a specific flow. The following is the flow's description: c = create an inventory file with all nodes, ports and paths a = run all validation tests (expecting an input inventory) v = only validate the given inventory file s = run service registration, deregistration, and lease test e = run event forwarding test f = flood the SA with queries according to the stress mode m = multicast flow q = QoS info: dump VLArb and SLtoVL tables t = run trap 64/65 flow (this flow requires running of external tool, default is all flows except QoS)
-w,wait	This option specifies the wait time for trap $64/65$ in seconds It is used only when running -f t - the trap $64/65$ flow (default to 10 sec)
-d,debug	This option specifies a debug option. These options are not normally needed. The number following -d selects the debug option to enable as follows: OPT Description
-m,max_lid	This option specifies the maximal LID number to be searched for during inventory file build (default to 100)
-g,guid	This option specifies the local port GUID value with which OpenSM should bind. OpenSM may be bound to 1 port at a time. If GUID given is 0, OpenSM displays a list of possible port GUIDs and waits for user input. Without -g, OpenSM trys to use the default port
-p,port	This option displays a menu of possible local port GUID values with which osmtest could bind
-i,inventory	This option specifies the name of the inventory file Normally, osmtest expects to find an inventory file, which osmtest uses to validate real-time information received from the SA during testing If -i is not specified, osmtest defaults to the file osmtest.dat See -c option for related information

Table 16 - osmtest Flags and Options

Flag	Description
-s,stress	This option runs the specified stress test instead of the normal test suite Stress test options are as follows: OPT Description
	 -s1 - Single-MAD (RMPP) response SA queries -s2 - Multi-MAD (RMPP) response SA queries -s3 - Multi-MAD (RMPP) Path Record SA queries -s4 - Single-MAD (non RMPP) get Path Record SA queries Without -s, stress testing is not performed
-M,Multicast_Mode	 This option specify length of Multicast test: OPT Description
	• Multiple mode - Could be run with other apps using MC with OpenSM. Without -M, default flow testing is performed
-t	This option specifies the time in milliseconds used for transaction timeouts. Specifying -t 0 disables timeouts. Without -t, OpenSM defaults to a timeout value of 200 milliseconds.
-l,log_file	This option defines the log to be the given file. By default the log goes to std- out.
-V	This option increases the log verbosity level. The -v option may be speci- fied multiple times to further increase the verbosity level. See the -vf option for more information about. log verbosity.
-V	This option sets the maximum verbosity level and forces log flushing. The - V is equivalent to '-vf0xFF -d 2'. See the -vf option for more information about. log verbosity.
-vf	This option sets the log verbosity level. A flags field must follow the -D option. A bit set/clear in the flags enables/disables a specific log level as follows: BIT LOG LEVEL ENABLED
	 0x01 - ERROR (error messages) 0x02 - INFO (basic messages, low volume) 0x04 - VERBOSE (interesting stuff, moderate volume) 0x08 - DEBUG (diagnostic, high volume) 0x10 - FUNCS (function entry/exit, very high volume) 0x20 - FRAMES (dumps all SMP and GMP frames) 0x40 - ROUTING (dump FDB routing information) 0x80 - currently unused. Without -vf, osmtest defaults to ERROR + INFO (0x3) Specifying -vf 0 disables all messages Specifying -vf 0xFF enables all messages (see -V) High verbosity levels may require increasing the transaction timeout with the -t
-h,help	Display this usage info then exit.

)

7.4 InfiniBand Fabric Performance Utilities

The performance utilities described in this chapter are intended to be used as a performance micro-benchmark. The tools are:

- Section 7.4.1,"ib_read_bw," on page 66
- Section 7.4.2, "ib_read_lat," on page 67
- Section 7.4.3, "ib_send_bw," on page 68
- Section 7.4.4, "ib_send_lat," on page 69
- Section 7.4.5, "ib_write_bw," on page 69
- Section 7.4.6, "ib_write_lat," on page 70
- Section 7.4.7, "ibv_read_bw," on page 71
- Section 7.4.8, "ibv_read_lat," on page 73
- Section 7.4.9, "ibv_send_bw," on page 74
- Section 7.4.10, "ibv_send_lat," on page 75
- Section 7.4.11, "ibv_write_bw," on page 77
- Section 7.4.12, "ibv_write_lat," on page 78

7.4.1 ib_read_bw

Rev 4.2

ib_read_bw calculats the BW of RDMA read between a pair of machines. One acts as a server and the other as a client. The client RDMA reads the server memory and calculate the BW by sampling the CPU each time it receive a successfull completion. The test supports features such as Bidirectional, in which they both RDMA read from each other memory's at the same time,change of mtu size, tx size, number of iteration, message size and more. Read is availible only in RC connection mode (as specified in IB spec).

7.4.1.1 ib_read_bw Synopsys

ib_read_bw [-i(b_port) ib_port] [-m(tu) mtu_size] [-s(ize) message_size] [-n iteration_num] [p(ort) PDT_port] [-b(idirectional)] [-o(uts) outstanding reads] [-a(ll)] [-V(ersion)]

7.4.1.2 ib_read_bw Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-o,outs= <num></num>	The number of outstanding read/atom(default 4)
-s,size= <size></size>	The size of message to exchange (default 65536)

Table 17 - ib_read_bw Flags and Options

Flag	Description
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-b,bidirectional	Measures bidirectional bandwidth (default unidirectional)
-V,version	Displays version number
-g,grh	Use GRH with packets (mandatory for RoCE)

Table 17 - ib_read_bw Flags and Options

7.4.2 ib_read_lat

ib_read_lat calculats the latency of RDMA read operation of message_sizeB between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which one side RDMA reads the memory of the other side only after the other side have read his memory. Each of the sides samples the CPU clock each time they read the other side memory , in order to calculate latency. Read is available only in RC connection mode (as specified in IB spec).

7.4.2.1 ib_read_lat Synopsys

ib_read_lat [-i(b_port) ib_port] [-m(tu) mtu_size] [-s(ize) message_size] [-t(x-depth) tx_size] [-n iteration_num] [-p(ort) PDT_port] [-o(uts) outstanding reads] [-a(ll)] [-V(ersion)] [-C report cycles] [-H report histogram] [-U report unsorted]

7.4.2.2 ib_read_lat Options

Table 18 - ib_read_lat Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-o,outs= <num></num>	The number of outstanding read/atom(default 4)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-C,report-cycles	Reports times in cpu cycle units (default microseconds)
-H,report-histogram	Print out all results (default print summary only)

Table 18	· ib_	_read_	lat	Flags	and	Options
----------	-------	--------	-----	-------	-----	---------

Flag	Description
-U,report-unsorted (implies -H)	Print out unsorted results (default sorted)
-V,version	Displays version number
-g,grh	Use GRH with packets (mandatory for RoCE)

7.4.3 ib_send_bw

ib_send_bw calculats the BW of SEND between a pair of machines. One acts as a server and the other as a client. The server receive packets from the client and they both calculate the throughput of the operation. The test supports features such as Bidirectional, on which they both send and receive at the same time, change of mtu size, tx size, number of iteration, message size and more. Using the "-a" provides results for all message sizes.

7.4.3.1 ib_send_bw Synopsys

ib_send_bw [-i(b_port) ib_port] [-c(onnection_type) RC\UC\UD] [-m(tu) mtu_size] [-s(ize)
message_size] [-t(x-depth) tx_size] [-n iteration_num] [-p(ort)
PDT_port] [-b(idirectional)] [-a(ll)] [-V(ersion)]

7.4.3.2 ib_send_bw Options

Table 19 - ib_send_bw Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc=""></rc>	Connection type RC/UC/UD (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-b,bidirectional	Measures bidirectional bandwidth (default unidirectional)
-V,version	Displays version number
-g,grh	Use GRH with packets (mandatory for RoCE)

7.4.4

ib_send_lat calculats the latency of sending a packet in message_sizeB between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which you send packet only if you receive one. Each of the sides samples the CPU each time they receive a packet in order to calculate the latency.

7.4.4.1 ib_send_lat Synopsys

```
ib_send_lat [-i(b_port) ib_port] [-c(onnection_type) RC\UC\UD] [-m(tu) mtu_size] [-s(ize)
message_size] [-t(x-depth) tx_size] [-n iteration_num] [-p(ort)
PDT_port] [-a(ll)] [-V(ersion)] [-C report cycles] [-H report
histogram] [-U report unsorted]
```

7.4.4.2 ib_send_lat Options

The table below lists the various flags of the command.

Table 20 - ib_send_lat Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc=""></rc>	Connection type RC/UC/UD (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-l,signal	Signal completion on each msg
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-C,report-cycles	Reports times in cpu cycle units (default microseconds)
-H,report-histogram	Print out all results (default print summary only)
-U,report-unsorted (implies -H)	Print out unsorted results (default sorted)
-V,version	Displays version number
-g,grh	Use GRH with packets (mandatory for RoCE)

7.4.5 ib_write_bw

ib_write_bw calculats the BW of RDMA write between a pair of machines. One acts as a server and the other as a client. The client RDMA writes to the server memory and calculate the BW by sampling the CPU each time it receive a successful completion. The test supports features such as Bidirectional, in which they both RDMA write to each other at the same time, change of mtu size, tx size, number of iteration, message size and more. Using the "-a" flag provides results for all message sizes.

7.4.5.1 ib_write_bw Synopsys

```
ib_write_bw [-q num of qps] [-c(onnection_type) RC\UC\UD] [-i(b_port) ib_port] [-m(tu) mtu_size]
[-s(ize) message_size] [-t(x-depth) tx_size] [-n iteration_num] [-p(ort) PDT_port] [-b(idirec-
tional)]
                                                                 [-a(ll)] [-V(ersion)]
```

7.4.5.2 ib_write_bw Options

The table below lists the various flags of the command.

Table 21 - ib_write_bw Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc=""></rc>	Connection type RC/UC/UD (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-b,bidirectional	Measures bidirectional bandwidth (default unidirectional)
-V,version	Displays version number
-o,post= <num of="" posts=""></num>	The number of posts for each qp in the chain (default tx_depth)
-q,qp= <num of="" qp's=""></num>	The number of qp's(default 1)
-g,grh	Use GRH with packets (mandatory for RoCE)

7.4.6 ib write lat

ib write lat calculats the latency of RDMA write operation of message sizeB between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which one side RDMA writes to the other side memory only after the other side wrote on his memory. Each of the sides samples the CPU clock each time they write to the other side memory, in order to calculate latency.

70

7.4.6.1 ib_write_lat Synopsys

```
ib_write_lat [-i(b_port) ib_port] [-c(onnection_type) RC\UC\UD] [-m(tu) mtu_size] [-s(ize)
message_size] [-t(x-depth) tx_size] [-n iteration_num] [-p(ort)
PDT_port] [-a(ll)] [-V(ersion)] [-C report cycles] [-H report
histogram] [-U report unsorted]
```

7.4.6.2 ib_write_lat Options

The table below lists the various flags of the command.

Table 22 - ib_write_lat Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc=""></rc>	Connection type RC/UC/UD (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-f,freq= <dep></dep>	How often the time stamp is taken
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-C,report-cycles	Reports times in cpu cycle units (default microseconds)
-H,report-histogram	Print out all results (default print summary only)
-U,report-unsorted (implies -H)	Print out unsorted results (default sorted)
-V,version	Displays version number
-g,grh	Use GRH with packets (mandatory for RoCE)

7.4.7 ibv_read_bw

This is a more advanced version of ib_read_bw and contains more flags and featurs than the older version and also improved algorithms. ibv_read_bw Calculats the BW of RDMA read between a pair of machines. One acts as a server, and the other as a client. The client RDMA reads the server memory and calculate the BW by sampling the CPU each time it receive a successfull completion. The test supports a large variety of features as described below, and has better performance than ib_send_bw in Nahelem systems. Read is available only in RC connection mode (as specified in the InfiniBand spec).

7.4.7.1 ibv_read_bw Synopsys

```
ibv_read_bw [-i(b_port) ib_port] [-d ib device] [-o(uts) outstanding reads] [-m(tu) mtu_size] [-
s(ize) message_size] [-t(x-depth) tx_size] [-n iteration_num] [-p(ort)
PDT_port] [-u qp timeout] [-S(1) sl type] [-x gid index] [-e(vents) use
events] [-F CPU freq fail] [-b(idirectional)] [-a(ll)] [-V(ersion)]
```

7.4.7.2 ibv_read_bw Options

Table 23 - ibv_read_bw Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-o,outs= <num></num>	The number of outstanding read/atom(default for hermon 16 (others 4)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-u,qp-timeout= <timeout></timeout>	QP timeout. The timeout value is 4 usec * 2 ^(timeout), default 14
-S,sl= <sl></sl>	The service level (default 0)
-x,gid-index= <index></index>	Test uses GID with GID index taken from command line (for RDMAoE index should be 0)
-b,bidirectional	Measures bidirectional bandwidth (default unidirectional)
-V,version	Displays version number
-g,post= <num of="" posts=""></num>	The number of posts for each qp in the chain (default tx_depth)
-e,events	Inactive during CQ events (default poll)
-F,CPU-freq	The CPU frequency test. It is active even if the cpufreq_ondemand module is loaded
-R,rdma_cm	Connect QPs with rdma_cm and run test on those QPs
-z,com_rdma_cm	Communicate with rdma_cm module to exchange data - use regular QPs
-c,connection= <rc uc="" ud=""></rc>	Connection type RC/UC/UD (default RC)
-I,inline_size= <size></size>	Max size of message to be sent in inline (default 0)
-Q,cq-mod	Generate Cqe only after <cq-mod> completion</cq-mod>
-N,no peak-bw	Cancel peak-bw calculation (default with peak)
7.4.8 ibv_read_lat

This is a more advanced version of ib_read_lat ,and contains more flags and featurs than the older version and also improved algorithms. ibv_read_lat calculats the latency of RDMA read operation of message_sizeB between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which one side RDMA reads the memory of the other side only after the other side have read his memory. Each of the sides samples the CPU clock each time they read the other side memory, to calculate latency. Read is available only in RC connection mode (as specified in InfiniBand spec).

7.4.8.1 ibv_read_lat Synopsys

<pre>ibv_read_lat [-i(b_port) ib_port] [-m(tu) mt</pre>	u_size] [-s(ize) m	essage_size] [-t(x-depth) tx_size]
[-I(nline_size) inline size] [-u qp timeout]		[-S(L) sl type] [-d ib_device
name] [-x gid index]		[-n iteration_num] [-o(uts)
outstanding reads]		[-e(vents) use events] [-p(ort)
<pre>PDT_port] [-a(ll)] [-V(ersion)]</pre>		[-C report cycles] [-H report
histogram] [-U report unsorted]	[-F CPU freq fail]	

7.4.8.2 ibv_read_lat Options

Table 24 - ibv_read_lat Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-o,outs= <num></num>	The number of outstanding read/atom(default for hermon 16 (others 4)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-u,qp-timeout= <timeout></timeout>	QP timeout. The timeout value is 4 usec * 2 ^(timeout), default 14
-S,sl= <sl></sl>	The service level (default 0)
-x,gid-index= <index></index>	Test uses GID with GID index taken from command line (for RDMAoE index should be 0)
-C,report-cycles	Reports times in cpu cycle units (default microseconds)
-H,report-histogram	Print out all results (default print summary only)
-U,report-unsorted (implies -H)	Print out unsorted results (default sorted)
-V,version	Displays version number

Flag	Description
-e,events	Inactive during CQ events (default poll)
-F,CPU-freq	The CPU frequency test. It is active even if the cpufreq_ondemand module is loaded
-R,rdma_cm	Connect QPs with rdma_cm and run test on those QPs
-z,com_rdma_cm	Communicate with rdma_cm module to exchange data - use regular QPs
-c,connection= <rc uc="" ud=""></rc>	Connection type RC/UC/UD (default RC)
-I,inline_size= <size></size>	Max size of message to be sent in inline (default 400)

Table 24 - ibv_read_lat Flags and Options

7.4.9 ibv_send_bw

This is a more advanced version of ib_send_bw and contains more flags and featurs than the older version and also improved algorithms. ibv_send_bw calculats the BW of SEND between a pair of machines. One acts as a server and the other as a client. The server receive packets from the client and they both calculate the throughput of the operation. The test supports a large variety of features as described below, and has better performance than ib_send_bw in Nahelem systems.

7.4.9.1 ibv_send_bw Synopsys

ibv_send_bw [-i(b_port) ib_port] [-d ib devi	.ce] [-c(onnection_type) RC\UC\UD] [-m(tu) mtu_size]
[-s(ize) message_size] [-t(x-depth) tx_size]	[-r(x_dpeth) rx_size] [-n iteration_num] [-p(ort)
PDT_port]	[-I(nline_size) inline size] [-u qp timeout] [-S(1)
sl type]	[-x gid index] [-e(vents) use events] [-N(o_peak)
use peak calc]	[-F CPU freq fail] [-g num of
qps in mcast group] [-M mcast gid]	[-b(idirectional)] [-a(ll)] [-V(ersion)]

7.4.9.2 ibv_send_bw Options

Table 25 - ibv send	bw Flags and Options
---------------------	----------------------

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc="" ud=""></rc>	Connection type RC/UC/UD (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)

Flag	Description
-u,qp-timeout= <timeout></timeout>	QP timeout. The timeout value is 4 usec * 2 ^(timeout), default 14
-S,sl= <sl></sl>	The service level (default 0)
-x,gid-index= <index></index>	Test uses GID with GID index taken from command line (for RDMAoE index should be 0)
-b,bidirectional	Measures bidirectional bandwidth (default unidirectional)
-V,version	Displays version number
-g,post= <num of="" posts=""></num>	The number of posts for each qp in the chain (default tx_depth)
-e,events	Inactive during CQ events (default poll)
-F,CPU-freq	The CPU frequency test. It is active even if the cpufreq_ondemand module is loaded
-r,rx-depth= <dep></dep>	Makes rx queue bigger than tx (default 600)
-I,inline_size= <size></size>	The maximum size of message to be sent in "inline mode" (default 0)
-N,no peak-bw	Cancels peak-bw calculation (default with peak-bw)
-g,mcg= <num_of_qps></num_of_qps>	Sends messages to multicast group with <num_of_qps> qps attached to it.</num_of_qps>
-M,MGID= <multicast_gid></multicast_gid>	In case of multicast, uses <multicast_gid> as the group MGID. The format must be '255:1:X:X:X:X:X:X:X:X:X:X:X:X:X:X, where X is a vlaue within [0,255]</multicast_gid>
-R,rdma_cm	Connect QPs with rdma_cm and run test on those QPs
-z,com_rdma_cm	Communicate with rdma_cm module to exchange data - use regular QPs
-Q,cq-mod	Generate Cqe only after <cq-mod> completion</cq-mod>

Table 25 - ibv_send_bw Flags and Options

7.4.10 ibv_send_lat

This is a more advanced version of ib_send_lat and contains more flags and featurs than the older version and also improved algorithms. ibv_send_lat calculats the latency of sending a packet in message_sizeB between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which you send packet only after you receive one. Each of the sides samples the CPU clock each time they receive a send packet, in order to calculate the latency.

7.4.10.1 ibv_send_lat Synopsys

Rev 4.2

<pre>ibv_send_lat [-i(b_port) ib_port]</pre>	[-c(onnection_type)	$RC\setminus UC\setminus UD]$	[-d ib_device name] [-m(tu)
mtu_size] [-s(ize) message_size]			[-t(x-depth) tx_size] [-
I(nline_size) inline size]			[-u qp timeout] [-S(L) sl type]
[-x gid index]			[-e(events) use events] [-n
iteration_num]			[-g num of qps in mcast group]
[-p(ort) PDT_port] [-a(ll)]			[-V(ersion)] [-C report cycles]
[-H report histogram]			[-U report unsorted] [-F CPU
freq fail]			

7.4.10.2 ibv_send_lat Options

Table 26 - ibv_send_lat Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc="" ud=""></rc>	Connection type RC/UC/UD (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-u,qp-timeout= <timeout></timeout>	QP timeout. The timeout value is 4 usec * 2 ^(timeout), default 14
-S,sl= <sl></sl>	The service level (default 0)
-x,gid-index= <index></index>	Test uses GID with GID index taken from command line (for RDMAoE index should be 0)
-C,report-cycles	Reports times in cpu cycle units (default microseconds)
-H,report-histogram	Print out all results (default print summary only)
-U,report-unsorted (implies -H)	Print out unsorted results (default sorted)
-V,version	Displays version number
-F,CPU-freq	The CPU frequency test. It is active even if the cpufreq_ondemand module is loaded
-g,post= <num of="" posts=""></num>	The number of posts for each qp in the chain (default tx_depth)
-I,inline_size= <size></size>	The maximum size of message to be sent in "inline mode" (default 0)
-e,events	Inactive during CQ events (default poll)
-g,mcg= <num_of_qps></num_of_qps>	Sends messages to multicast group with <num_of_qps> qps attached to it.</num_of_qps>

Flag	Description
-M,MGID= <multicast_gid></multicast_gid>	In case of multicast, uses <multicast_gid> as the group MGID. The format must be '255:1:X:X:X:X:X:X:X:X:X:X:X:X:X:X', where X is a vlaue within [0,255]. You must specify a different MGID on both sides to avoid loopback.</multicast_gid>
-R,rdma_cm	Connect QPs with rdma_cm and run test on those QPs
-z,com_rdma_cm	Communicate with rdma_cm module to exchange data - use regular QPs

Table 26 - ibv_send_lat Flags and Options

7.4.11 ibv_write_bw

This is a more advanced version of ib_write_bw ,and contains more flags and features than the older version and also improved algorithms. ibv_write_bw calculates the BW of RDMA write between a pair of machines. One acts as a server and the other as a client. The client RDMA writes to the server memory and calculate the BW by sampling the CPU each time it receives a successful completion. The test supports a large variety of features as described below, and has better performance than ib_send_bw in Nahelem systems.

7.4.11.1 ibv_write_bw Synopsys

<pre>ibv_write_bw [-i(b_port) ib_port] [-d ib device]</pre>	[-c(onnection_type) RC\UC\UD] [-m(tu) mtu_size]
[-s(ize) message_size] [-t(x-depth) tx_size]	<pre>[-n iteration_num] [-p(ort)</pre>
PDT_port] [-I(nline_size) inline size]	[-u qp timeout] [-S(l) sl type]
[-x gid index] [-e(vents) use events] [-N(o_peak) use peak calc] [-
F CPU freq fail] [-g num of posts] [-q num of qps] [-b(idirectional)] [-a(ll)] [-
V(ersion)]	

7.4.11.2 ibv_write_bw Options

Table 27 - ibv_write_bw Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <port> (default 18515)</port>
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc=""></rc>	Connection type RC/UC(default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-u,qp-timeout= <timeout></timeout>	QP timeout. The timeout value is 4 usec * 2 ^(timeout), default 14

Flag	Description
-S,sl= <sl></sl>	The service level (default 0)
-x,gid-index= <index></index>	Test uses GID with GID index taken from command line (for RDMAoE index should be 0)
-b,bidirectional	Measures bidirectional bandwidth (default unidirectional)
-V,version	Displays version number
-g,post= <num of="" posts=""></num>	The number of posts for each qp in the chain (default tx_depth)
-F,CPU-freq	The CPU frequency test. It is active even if the cpufreq_ondemand module is loaded
-q,qp= <num of="" qp's=""></num>	The number of qp's (default 1)
-I,inline_size= <size></size>	The maximum size of message to be sent in "inline mode" (default 0)
-N,no peak-bw	Cancels peak-bw calculation (default with peak-bw)
-R,rdma_cm	Connect QPs with rdma_cm and run test on those QPs
-z,com_rdma_cm	Communicate with rdma_cm module to exchange data - use regular QPs
-Q,cq-mod	Generate Cqe only after <cq-mod> completion</cq-mod>

Table 27 - ibv_write_bw Flags and Options

7.4.12 ibv_write_lat

This is a more advanced version of ib_write_lat and contains more flags and featurs than the older version and also improved algorithms. ibv_write_lat calculats the latency of RDMA write operation of message_sizeB between a pair of machines. One acts as a server, and the other as a client. They perform a ping pong benchmark on which one side RDMA writes to the other side memory only after the other side wrote on his memory. Each of the sides samples the CPU clock each time they write to the other side memory to calculate latency.

7.4.12.1 ibv_write_lat Synopsis

<pre>ibv_write_lat [-i(b_port) ib_port]</pre>	[-c(onnection_type) RC\UC\UD	[-m(tu) mtu_size] [-s(ize)
<pre>message_size] [-t(x-depth) tx_size]</pre>		[-I(nline_size) inline size] [-
u qp timeout] [-S(L) sl type]	[-d ib_device name]	[-x gid index] [-n
iteration_num]		[-p(ort) PDT_port] [-a(ll)] [-
V(ersion)] [-C report cycles]		[-H report histogram] [-U report
unsorted]		

7.4.12.2 ibv_write_lat Options

Table 28 - ibv_write_lat Flags and Options

Flag	Description
-p,port= <port></port>	Listens on/connect to port <pre>port> (default 18515)</pre>

Table 20 - Iby Write lat rags and Options	Table	28 -	ibv	write	lat F	lags	and	Options
---	-------	------	-----	-------	-------	------	-----	---------

Flag	Description
-d,ib-dev= <dev></dev>	Uses IB device <device guid=""> (default first device found)</device>
-i,ib-port= <port></port>	Uses port <port> of IB device (default 1)</port>
-m,mtu= <mtu></mtu>	The mtu size (default 1024)
-c,connection= <rc uc=""></rc>	Connection type RC/UC (default RC)
-s,size= <size></size>	The size of message to exchange (default 65536)
-a,all	Runs sizes from 2 till 2^23
-t,tx-depth= <dep></dep>	The size of tx queue (default 100)
-n,iters= <iters></iters>	The number of exchanges (at least 2, default 1000)
-u,qp-timeout= <timeout></timeout>	QP timeout. The timeout value is 4 usec * 2 ^(timeout), default 14
-S,sl= <sl></sl>	The service level (default 0)
-x,gid-index= <index></index>	Test uses GID with GID index taken from command line (for RDMAoE index should be 0)
-C,report-cycles	Reports times in cpu cycle units (default microseconds)
-H,report-histogram	Print out all results (default print summary only)
-U,report-unsorted (implies -H)	Print out unsorted results (default sorted)
-V,version	Displays version number
-F,CPU-freq	The CPU frequency test. It is active even if the cpufreq_ondemand module is loaded
-I,inline_size= <size></size>	The maximum size of message to be sent in "inline mode" (default 0)
-R,rdma_cm	Connect QPs with rdma_cm and run test on those QPs
-z,com_rdma_cm	Communicate with rdma_cm module to exchange data - use regular QPs

7.4.13 ibaddr

Displays the lid (and range) as well as the GID address of the port specified (by DR path, lid, or GUID) or the local port by default.



This utility can be used as simple address resolver.

7.4.13.1 ibaddr Synopsis

```
ibaddr [-d(ebug)] [-D(irect)] [-G(uid)] [-l(id_show)] [-g(id_show)] [-C
ca_name] [-P ca_port] [-t(imeout) timeout_ms] [-V(ersion)] [-h(elp)]
[<lid | dr_path | guid>]
```

7.4.13.2 ibaddr Options

Table 29 - ibaddr Flags and Options

Flags	Description
-G,Guid	shows lid range and gid for GUID address
-l,lid_show	shows lid range only
-L,Lid_show	shows lid range (in decimal) only
-g,gid_show	shows gid address only
Debugging Flags	Description
NOTE: Most OpenIB diagnostics take the followin ity can be found in the usage message and can be s	g common flags. The exact list of supported flags per util- hown using the util_name -h syntax.
-d	Raises the IB debugging level. Can be used several times (-ddd or -d -d -d).
-e	shows send and receive errors (timeouts and others)
-h	shows the usage message
-V	Increases the application verbosity level. Can be used several times (-vv or -v -v -v)
-v	shows the version info.
Addressing Flags	Description
-D	Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0,1,2,1,4" # out via port 1, then 2,
-G	Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023"
-s <smlid></smlid>	Uses 'smlid' as the target lid for SM/SA queries.
Other Common Flags	Description
-C <ca_name></ca_name>	Uses the specified ca_name.
-C <ca_name> -P <ca_port></ca_port></ca_name>	Uses the specified ca_name. Uses the specified ca_port.

7.4.13.3 Multiple CA/Multiple Port Support

When no IB device or port is specified, the port to use is selected by the following criteria:

- 1. the first port that is ACTIVE.
- 2. if not found, the first port that is UP (physical link up).

If a port and/or CA name is specified, the user request is attempted to be fulfilled, and will fail if it is not possible.

Examples

ibaddr	<pre># local port's address</pre>	
ibaddr 32	<pre># show lid range and gid of lid 32</pre>	
ibaddr -G 0x8f1040023	<pre># same but using guid address</pre>	
ibaddr -1 32	<pre># show lid range only</pre>	
ibaddr -L 32	<pre># show decimal lid range only</pre>	
ibaddr -g 32	# show gid address only	

7.4.14 ibcacheedit

ibcacheedit allows users to edit an ibnetdiscover cache created through the --cache option in ibnetdiscover(8).

7.4.14.1 ibcacheedit Synopsis

ibcacheedit [--switchguid BEFOREGUID:AFTERGUID] [--caguid BEFORE:AFTER]
 [--sysimgguid BEFOREGUID:AFTERGUID] [-h(elp)] <orig.cache> <new.cache>

7.4.14.2 ibcacheedit Options

Table 30 - ibcacheedit Flags and Options

Flags	Description		
switchguid BEFOREGUID:AFTERGUID	Specifies a switchguid that should be changed. The before and after guid should be separated by a colon. On switches, port guids are identical to the switch guid, so port guids will be adjusted as well on switches.		
caguid BEFOREGUID:AFTERGUID	Specifies a caguid that should be changed. The before and after guid should be separated by a colon.		
sysimgguid BEFOREGUID:AFTERGUID	Specifies a sysimgguid that should be changed. The before and after guid should be spearated by a colon.		
portguid NODEGUID:BEFOREGUID:AFTERGUID	Specifies a portguid that should be changed. The nodeguid of the port (e.g. switchguid or caguid) should be specified first, followed by a colon, the before port guid, another colon, then the after port guid.On switches, port guids are identical to the switch guid, so the switch guid will be adjusted as well on switches.		
Debugging Flags	Description		
NOTE: Most OpenIB diagnostics take the following common flags. The exact list of supported flags per utility can be found in the usage message and can be shown using the util_name -h syntax.			
-h	shows the usage message		
-v	shows the version info.		

7.4.15 iblinkinfo

iblinkinfo reports link info for each port in an IB fabric, node by node. Optionally, iblinkinfo can do partial scans and limit its output to parts of a fabric.

7.4.15.1 iblinkinfo Synopsis

7.4.15.2 iblinkinfo Flags and Options

Table 31 - iblinkinfo Flags and Options

Flags	Description
-S <port_guid> -G <port_guid>port-guid</port_guid></port_guid>	start partial scan at the port specified by <port_guid> (hex for- mat)</port_guid>
-D <direct_route></direct_route>	start partial scan atthe port specified by the direct route path.
-1	Print all information for each link on one line. Defaultis to print a headerwith the node information and then a list for each port (useful for grep'ing output).
-d	Prints only nodes which have a port in the "Down" state.
-p	Prints additional port settings (<life- Time>,<hoqlife>,<vlstall-count>)</vlstall-count></hoqlife></life-
-C <ca_name></ca_name>	Uses the specified ca_name for the search.
-P <ca_port></ca_port>	Uses the specified ca_port for the search.
-R	(This option is obsolete and does nothing)
load-cache <filename></filename>	Load and use the cached ibnetdiscover data stored in the speci- fied filename. May be useful for outputting and learning about other fabrics or a previous state of a fabric.Cannot be used if user specifies a direct route path. See ibnetdiscover for information on caching ibnetdiscover output.
diff <filename></filename>	Load cached ibnetdiscover data and do a diff comparison to the current network or another cache. A special diff output for iblinkinfo output will be displayed showing differences between the old and current fabric links. Be default, the following are com- paredfor differences: port connections and port state. See ibnetdis- cover for information on caching ibnetdiscover output.
diffcheck <key(s)></key(s)>	Specify what diff checks should be done in thediffoption above. Comma separate multiple diff check key(s). The available diff checks are:port = port connections,state = port state, lid = lids, nodedesc = node descriptions. If port is specified alongside lid or nodedesc, remote port lids and node descriptions will also be compared.
filterdownports <filename></filename>	Filter downports indicated in a ibnetdiscover cache. If a port was previously indicated as down in the specified cache, and is still down, donot output it in the resulting output. This option may be particularly useful for environments where switchesare not fully populated, thus much of the default iblinkinfo info is considered un-useful. See ibnetdiscover for information on caching ibnetdis- cover output.

7.4.16 ibqueryerrors

The default behavior is to report the port error counters which exceed a threshold for each port in the fabric. The default threshold is zero (0). Error fields can also be suppressed entirely.

In addition to reporting errors onevery port. ibqueryerrors can report the port transmit and receive data as well as report full link information to the remote port if available.

7.4.16.1 ibqueryerrors Synopsis

ibqueryerrors [options]

7.4.16.2 ibqueryerrors Options

Table 32 - ibqueryerrors	Flags and Options
--------------------------	-------------------

Flags	Description
-s <err1,err2,></err1,err2,>	Suppresses the errors listed in the comma separated list provided.
-c	Suppresses some of the common "side effect" counters. These counters usually do not indicate an error condition and can be usually be safely ignored.
-G <port_guid> -S <port_guid>port-guid</port_guid></port_guid>	Report results for the port specified. For switches results are printed for all ports not just switch port 0.
-S same as "-G"	Provided only for backward compatibility
-D <direct_route></direct_route>	Reports results for the port specified. For switches results are printed for all ports not just switch port 0.
-r	Reports the port information. This includes LID, port, external port (if applicable), link speed setting, remote GUID, remote port, remote external port (if applicable), and remote node description information.
data	Includes the optional transmit and receive data counters.
threshold-file	Specifies an alternate threshold file. The default is: /opt/ufm/files/conf/infiniband-diags/error_thresholds
switch	Prints data for switches only.
ca	Prints data for CA's only.
router	Prints data for routers only
clear-errors-k	Clear error counters after readk and -K can be used together to clear both errors and counters.
clear-counts -K	Clear data counters after read.
	CAUTION: clearing data counters will occur regardless of if they are printed or not. This is because data counters are only printed on ports which have errors. This means if a port has 0 errors and the -K option is specified the data counters will be cleared without any printed output.
-details	Includes receive error and transmits discard details
load-cache <filename></filename>	Loads and uses the cached ibnetdiscover data stored in the specified filename. May be useful for outputting and learning about other fabrics or a previous state of a fabric. Cannot be used if user speci- fies a direct route path. See ibnetdiscover for information on caching ibnetdiscover output.
-R	This option is obsolete (and has no effect).
-d	Raises the IB debugging level. May be used several times (-ddd or - d -d -d).
-е	Shows send and receive errors (time-outs and others)

rapie 32 - inquei yerrors riags and Options	Table 32 -	ibqueryerrors	Flags and	Options
---	------------	---------------	-----------	---------

Flags	Description
-h	Shows the usage message
-v	Increases the application verbosity level. May be used several times (-vv or -v -v -v)
-C <ca_name></ca_name>	Uses the specified ca_name.
-P <ca_port></ca_port>	Uses the specified ca_port.
-t <timeout_ms></timeout_ms>	Overrides the default timeout for the solicited mads.

7.4.16.3 ibqueryerrors Exit Status

If a failure to scan the fabric occurs return -1. If the scan succeeds without errors beyond thresholds return 0. If errors found on ports beyond thresholds return 1.

7.4.16.4 ibqueryerrors Files

Rev 4.2

/opt/ufm/files/conf/infiniband-diags/error_thresholds

Define threshold values for errors. File format is simple "name=val".

Comments begin with '#'

Example:

Define thresholds for error counters SymbolErrorCounter=10 LinkErrorRecoveryCounter=10 VL15Dropped=100

7.4.17 ibsysstat

ibsysstat usesvendormads to validate connectivity between IB nodes and obtain other information about the IB node.ibsysstat is run as client/server. Default is to run as client.

7.4.17.1ibsysstat Synopsis

ibsysstat [-d(ebug)] [-e(rr_show)] [-v(erbose)] [-G(uid)] [-C ca_name] [-P ca_port] [-s smlid] [-t(imeout) timeout_ms] [-V(ersion)] [-o oui] [-S(erver)] [-h(elp)] <dest lid | guid> [<op>]

7.4.17.2 ibsysstat Options

Table 3	3 - ib	svsstat	Flags	and	Options
Tuble 0		Syssiai	i iugo	unu	opuons

Flags	Description
ping	Verifies connectivity to server (default)
host	Obtains host information from server
сри	Obtains cpu information from server
-o,oui	Uses specified OUI number to multiplex vendor mads
-S,Server	Starts in server mode (do not return)
Debugging Flags	Description
NOTE: Most OpenIB diagnostics take the fo per utility can be found in the usage message	llowing common flags. The exact list of supported flags e and can be shown using the util_name -h syntax.
-d	Raises the IB debugging level. Can be used several times (-ddd or -d -d -d).
-е	Shows send and receive errors (timeouts and others)
-h	Shows the usage message
-v	Increases the application verbosity level. Can be used several times (-vv or -v -v).
-V	Shows the version info.
Addressing Flags	Description
-G	Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023"
-s <smlid></smlid>	Uses 'smlid' as the target lid for SM/SA queries.
Other Common Flags	Description
-C <ca_name></ca_name>	Uses the specified ca_name.
-P <ca_port></ca_port>	Uses the specified ca_port.
-t <timeout_ms></timeout_ms>	Overrides the default timeout for the solicited mads.

7.4.17.3 Multiple CA/Multiple Port Support

When no IB device or port is specified, the port to use is selected by the following criteria:

- 1. The first port that is ACTIVE.
- 2. If not found, the first port that is UP (physical link up).

If a port and/or CA name is specified, the user request is attempted to be fulfilled, and will fail if it is not possible.

7.4.18 perfquery

Rev 4.2

perfquery usesPerfMgt GMPs to obtain the PortCounters (basic performance and error counters), PortExtendedCounters, PortXmitDataSL, PortR-cvDataSL, PortRcvErrorDetails,PortXmitDiscardDetails, PortExtended- SpeedsCounters, or PortSamplesControl from the PMA atthe node/port specified. Optionally shows aggregated counters for all ports of node. In addition, you may reset after read, or only reset counters.



• In PortCounters, PortCountersExtended, PortXmitDataSL, and PortR-cvDataSL, components that represent Data (e.g. PortXmitData and PortR- cvData) indicate octets divided by 4 rather than just octets.

Inputting a port of 255 indicates an operation be performed on all ports.

7.4.18.1 perfquery Synopsis

perfquery [-d(ebug)]	[-G(uid)] [-x extended] [-X xmtsl]
[-S rcvsl]	[-D]xmtdisc] [-E]rcverr]
[-T extended_speeds][-	oprcvcounters] [flowctlcounters] [vlop-
packets]	[vlopdata] [vlxmitflowctlerrors] [vlxmitcounters]
[swportvlcong]	[rcvcc] [slrcvfecn] [slrcvbecn] [xmitcc]
[vlxmittimecc]	[-c smplctl] [-a(ll_ports)] [-l(oop_ports)]
[-r(eset_after_read)]	[-R(eset_only)] [-C ca_name] [-P ca_port]
[-t(imeout) timeout_ms]	[-V(ersion)] [-h(elp)] [<lid guid> [[port]</lid guid>
[reset_mask]]]	

7.4.18.2 perfquery Options

Table 34 - perfquery Flags and Options

Flags	Description
-x,extended	Shows extended port counters rather than (basic) port counters. Note that extended port counters attribute is optional.
-X,xmtsl	Shows transmit data SL counter. This is an optional counter for QoS.
-S,revsl	Shows receive data SL counter. This is an optional counter for QoS.
-D,xmtdisc	Shows transmit discard details. This is an optional counter.
-E,reverr	Shows receive error details. This is an optional counter.
-D,xmtdisc	Shows transmit discard details. This is an optional counter.
-T,extended_speeds	Shows extended speeds port counters. This is an optional counter.
opreveounters	Shows Rcv Counters per Op code. This is an optional counter.
flowctlcounters	Shows flow control counters. This is an optional counter.
vloppackets	Shows packets received per Op code per VL. This is an optional counter.
vlopdata	Show data received per Op code per VL. This is an optional counter.
vlxmitflowctlerrors	Shows flow control update errors per VL. This is an optional counter.
vlxmitcounters	Shows ticks waiting to transmit countersper VL. Thisis an optional counter.
swportvlcong	Shows sw port VL congestion. This is an optional counter.

Table 34 -	perfquery	Flags and	Options
------------	-----------	-----------	---------

Flags	Description
revee	Shows Rcv congestion control counters. This is an optional counter.
slrcvfecn	Shows SL Rcv FECN counters. This is an optional counter.
slrcvbecn	Shows SL Rcv BECN counters. This is an optional counter.
xmitce	Shows Xmit congestion control counters. This is an optional counter.
vlxmittimecc	Shows VLXmit Time congestioncontrol counters. This is an optional counter.
-c,smplctl	Shows port samples control.
-a,all_ports	Shows aggregated counters for all ports of the destination lid or reset all counters for all ports. If the destination lid does not support the All-PortSelect flag, all ports will be iterated through to emulate AllPortSelect behavior.
-l,loop_ports	If all ports are selected by the user (either through the -a option or port 255) iterate through each port rather than doing than aggre- gate operation.
-r,reset_after_read	Resets counters after read
-R,Reset_only	Only reset counters
Debugging Flags	Description
-d	Raises the IB debugging level. Can be used several times (-ddd or -d -d -d).
-е	Shows send and receive errors (timeouts and others)
-h	Shows the usage message
-v	Increases the application verbosity level. Can be used several times (-vv or -v -v -v)
-V	Shows the version info.
Addressing Flags	Description
-G	Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023"
-s <smlid></smlid>	Uses 'smlid' as the target lid for SM/SA queries.
-C <ca_name></ca_name>	Uses the specified ca_name.

Table 34 - perfquery Flags and Options

Flags	Description
-P <ca_port></ca_port>	Uses the specified ca_port.
-t <timeout_ms></timeout_ms>	override the default timeout for the solicited mads.

7.4.18.3 Multiple CA/Multiple Port Support

When no IB device or port is specified, the port to use is selected by the following criteria:

- 1. The first port that is ACTIVE.
- 2. If not found, the first port that is UP (physical link up).

If a port and/or CA name is specified, the user request is attempted to be fulfilled, and will fail if it is not possible.

7.4.19 saquery

saqueryissuesthe selected SA query. Node records are queried by default.

7.4.19.1 saquery Synopsis

```
saquery [-h] [-d] [-p] [-N] [--list | -D] [-S] [-I] [-L] [-I] [-G] [-0]
[-U] [-c] [-s] [-g] [-m] [-x] [-C ca_name] [-P ca_port][--smkey val]
[-t(imeout) <msec>] [--src-to-dst <src:dst>] [--sgid-to-dgid
<sgid-dgid>] [--node-name-map <node-name-map>] [<name> | <lid>| |
<guid>]
```

7.4.19.2 saquery Options

Table 35 - saquery Flags and Options

Flags	Description
-р	Gets PathRecord info.
-N	Gets NodeRecord info.
list -D	Gets NodeDescriptions of CAs only.
-S	Gets ServiceRecord info.
-I	Gets InformInfoRecord (subscription) info.
-L	Returns the Lids of the name specified
-1	Returns the unique Lid of the name specified
-G	Returns the Guids of the name specified
-0	Returns the name for the Lid specified
-U	Returns the name for the Guid specified
-C	Gets the SA's class port info
-S	Returns the PortInfoRecords with isSM or isSMdisabled capability mask bit on.
-g	Gets multicast group info
-m	Gets multicast member info. If a group is speci- fied, limit the output to the group specified and print one line containing only the GUID and node description for each entry. Example: saquery -m 0xc000
-x	Gets LinkRecord info.
src-to-dst	Gets a PathRecord for <src:dst> where src and dst are either node names or LIDs.</src:dst>
sgid-to-dgid	Gets a PathRecord for sgid to dgid where both GIDs are in an IPv6 format acceptable to inet_pton(3).
-C <ca_name></ca_name>	Uses the specified ca_name.
-P <ca_port></ca_port>	Uses the specified ca_port.

Flags	Description
smkey <val></val>	Uses SM_Key value for the query. Will be used only with "trusted" queries. If non-numeric value (like'x') is specified then saquery will prompt for a value.
-t, -timeout <msec></msec>	Specifies SA query response timeout in milli- seconds. Default is 100 milliseconds. You may want to use this option if IB_TIMEOUT is indi- cated.
node-name-map <node-name-map></node-name-map>	Specifies a node name map. The node name map file maps GUIDs to more user friendly names. See ibnetdiscover(8) for node name map file format.Only used with the -O and -U options.
	 Supported query names (and aliases): ClassPortInfo (CPI) NodeRecord (NR) [lid] PortInfoRecord (PIR) [[lid]/[port]/ [options]] SL2VLTableRecord (SL2VL) [[lid]/ [in_port]/[out_port]] PKeyTableRecord (PKTR) [[lid]/[port]/ [block]] VLArbitrationTableRecord (VLAR) [[lid]/[port]/[block]] InformInfoRecord (IIR) LinkRecord (LR) [[from_lid]/ [from_port]] [[to_lid]/[to_port]] ServiceRecord (SR) PathRecord (PR) MCMemberRecord (MCMR) LFTRecord (LFTR) [[lid]/[block]] MFTRecord (MFTR) [[mlid]/ [position]/[block]] GUIDInfoRecord (GIR) [[lid]/[block]]
-d	enables debugging.
-h	Shows help.

Table 35 - saquery Flags and Options

7.4.20 smpdump

smpdump is a general purpose SMP utility which gets SM attributes from a specified SMA. The result is dumped in hex by default.

7.4.20.1 smpdump Synopsis

smpdump [-s(ring)] [-D(irect)] [-C ca_name] [-P ca_port] [-t(imeout)
timeout_ms] [-V(ersion)] [-h(elp)] <dlid|dr_path> <attr> [mod]

7.4.20.2 smpdump Options

Rev 4.2

Table 36 -	smpdum	p Flags	and O	ptions

Flags	Description
attr	IBA attribute ID for SM attribute
mod	IBA modifier for SM attribute
Debugging Flags	Description
NOTE: Most OpenIB diagnostics take the followin can be found in the usage message and can be show	g common flags. The exact list of supported flags per utility vn using the util_name -h syntax.
-d	Raises the IB debugging level. Can be used several times (-ddd or -d -d -d).
-е	Shows send and receive errors (timeouts and others)
-h	Shows the usage message
-v	Increases the application verbosity level. Can be used several times (-vv or -v -v)
-V	Shows the version info.
Addressing Flags	Description
Addressing Flags -D	Description Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0,1,2,1,4" # out via port 1, then 2,
Addressing Flags -D -G	Description Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0" # self port "0,1,2,1,4" # out via port 1, then 2, Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023"
Addressing Flags -D -G -s <smlid></smlid>	Description Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0" # self port "0,1,2,1,4" # out via port 1, then 2, Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023" Uses 'smlid' as the target lid for SM/SA queries.
Addressing Flags -D -G -s <smlid> Flags</smlid>	Description Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0,1,2,1,4" # out via port 1, then 2, Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023" Uses 'smlid' as the target lid for SM/SA queries.
Addressing Flags -D -G -s <smlid> Flags -C <ca_name></ca_name></smlid>	Description Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0" # self port "0,1,2,1,4" # out via port 1, then 2, Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023" Uses 'smlid' as the target lid for SM/SA queries. Description Uses the specified ca_name.
Addressing Flags -D -G -G -s <smlid> Flags -C <ca_name> -P <ca_port></ca_port></ca_name></smlid>	Description Uses directed path address arguments. The path is a comma separated list of out ports. Examples: "0" # self port "0,1,2,1,4" # out via port 1, then 2, Uses GUID address argument. In most cases, it is the Port GUID. Example: "0x08f1040023" Uses 'smlid' as the target lid for SM/SA queries. Uses the specified ca_name. Uses the specified ca_port.

7.4.20.3 Multiple CA/Multiple Port Support

When no IB device or port is specified, the port to use is selected by the following criteria:

- 1. The first port that is ACTIVE.
- 2. If not found, the first port that is UP (physical link up).

If a port and/or CA name is specified, the user request is attempted to be fulfilled, and will fail if it is not possible.

Examples

Direct Routed Examples:

smpdump -D 0,1,2,3,5 16 # NODE DESC
smpdump -D 0,1,2 0x15 2 # PORT INFO, port 2

LID Routed Examples:

smpdump 3 0x15 2	<pre># PORT INFO, lid 3 port 2</pre>
smpdump 0xa0 0x11	# NODE INFO, lid 0xa0

8 Software Development Kit

Software Development Kit (SDK) a set of development tools that allows the creation of Infini-Band applications for MLNX_VPI software package.

The SDK package contains, header files, libraries, and code examples. To open the SDK package you must run the sdk.exe file and get the complete list of files. SDK package can be found under <installation_directory>\IB\SDK

9 Troubleshooting

9.1 InfiniBand Troubleshooting

Issue # 1: The IB interfaces is not up after the first reboot after the installation process is completed.

Suggestion: To troubleshoot this issue, follow the steps bellow:

- 1. Check that the IB driver is running on all nodes by using 'vstat'. The vstat utility located at <installation_directory>\tools, displays the status and capabilities of the network adaptor card(s).
- 2. On the command line, enter "vstat" (use -h for options) to retrieve information about one or more adapter ports. The field port_state will be equal to:
 - PORT_DOWN when there is no InfiniBand cable ("no link");
 - PORT_INITIALIZED when the port is connected to some other port ("physical link");
 - PORT_ACTIVE when the port is connected and OpenSM is running ("logical link")
 - PORT_ARMED when the port is connected to some other port ("physical link");
- 3. Run OpenSM see OpenSM operation instructions in the OpenSM section above.
- **4.** Verify the status of ports by using vstat: All connected ports should report "PORT_ACTIVE" state.

9.2 Ethernet Troubleshooting

Issue #1: The installation of MLNX_VPI for Windows fails with the following (or a similar) error message:

This installation package is not supported by this processor type. Contact your product vendor."

Suggestion: This message is printed if you have downloaded and attempted to install an incorrect MSI -- for example, if you are trying to install a 64-bit MSI on a 32-bit machine (or vice versa).

Issue # 2: The performance is low.

Suggestion: This can be due to non-optimal system configuration. See the section "Performance Tuning" to take advantage of Mellanox 10 GBit NIC performance.

Issue # 3: The driver does no start.

Suggestion 1: This can happen due to an RSS configuration mismatch between the TCP stack and the Mellanox adapter. To confirm this scenario, open the event log and look under "System" for the "mlx4eth5" or "mlx4eth6" source. If found, enable RSS as follows:

1. Run the following command: "netsh int tcp set global rss = enabled".

Suggestion 2: This is a less recommended suggestion, and will cause low performance. Disable RSS on the adapter. To do this set RSS mode to "No Dynamic Rebalancing".

Issue #4: The Ethernet driver fails to start. In the Event log, under the mlx4_bus source, the following error message appears: RUN_FW command failed with error -22

Suggestion: The error message indicates that the wrong firmware image has been programmed on the adapter card.

See http://www.mellanox.com > Support > Firmware Download

Issue # 5: The Ethernet driver fails to start. A yellow sign appears near the "Mellanox ConnectX 10Gb Ethernet Adapter" in the Device Manager display.

Suggestion: This can happen due to a hardware error. Try to disable and re-enable "Mellanox ConnectX Adapter" from the Device Manager display.

Issue # 6: No connectivity to a Fault Tolerance bundle while using network capture tools (e.g., Wireshark).

Suggestion: This can happen if the network capture tool captures the network traffic of the nonactive adapter in the bundle. This is not allowed since the tool sets the packet filter to "promiscuous", thus causing traffic to be transferred on multiple interfaces. Close the network capture tool on the physical adapter card, and set it on the LBFO interface instead.

Issue #7: No Ethernet connectivity on 1Gb/100Mb adapters after activating Performance Tuning (part of the installation).

Suggestion: This can happen due to adding a TcpWindowSize registry value. To resolve this issue, remove the value key under HKEY_LOCAL_MACHINE\SYSTEM\CurrentControl-Set\Services\Tcpip\Parameters\TcpWindowSize or set its value to 0xFFFF.

Issue #8: System reboots on an I/OAT capable system on Windows Server 2008.

Suggestion: This may occur if you have an Intel I/OAT capable system with Direct Cache Access enabled, and 9K jumbo frames enabled. To resolve this issue, disable 9K jumbo frames.

Issue # 9: Packets are being lost.

Suggestion: This may occur if the port MTU has been set to a value higher than the maximum MTU supported by the switch.

Issue # 10: Issue(s) not listed above.

Suggestion: The MLNX_EN for Windows driver records events in the system log of the Windows event system. Using the event log you'll be able to identify, diagnose, and predict sources of system problems.

To see the log of events, open System Event Viewer as follows:

1. Right click on My Computer, click Manage, and then click Event Viewer.

OR

- 1. Click start-->Run and enter "eventvwr.exe".
- 2. In Event Viewer, select the system log.

The following events are recorded:

- Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> has been successfully initialized and enabled.
- Failed to initialize Mellanox ConnectX EN 10Gbit Ethernet Adapter.
- Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> has been successfully initialized and enabled. The port's network address is <MAC Address>
- The Mellanox ConnectX EN 10Gbit Ethernet was reset.
- Failed to reset the Mellanox ConnectX EN 10Gbit Ethernet NIC. Try disabling then re-enabling the "Mellanox Ethernet Bus Driver" device via the Windows device manager.
- Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> has been successfully stopped.
- Failed to initialize the Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> because it uses old firmware version (<old firmware version>). You need to burn firmware version <new firmware version> or higher, and to restart your computer.
- Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> device detected that the link connected to port <Y> is up, and has initiated normal operation.
- Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> device detected that the link connected to port <Y> is down. This can occur if the physical link is disconnected or damaged, or if the other endport is down.
- Mismatch in the configurations between the two ports may affect the performance. When Using MSI-X, both ports should use the same RSS mode. To fix the problem, configure the RSS mode of both ports to be the same in the driver GUI.
- Mellanox ConnectX EN 10Gbit Ethernet Adapter <X> device failed to create enough MSI-X vectors. The Network interface will not use MSI-X interrupts. This may affects the performance. To fix the problem, configure the number of MSI-X vectors in the registry to be at least <Y>.

10 Documentation

- Under <installation_directory>\Documentation:
 - License file
 - User Manual (this document)
 - MLNX_VPI_Installation Guide
 - MLNX_VPI_Release Notes