



Connect. Accelerate. Outperform.™

Mellanox OFED for Linux Release Notes

Rev 3.2-2.0.0.0

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 PRODUCT(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

© Copyright 2016. Mellanox Technologies LTD. All Rights Reserved.

Mellanox®, Mellanox logo, BridgeX®, CloudX logo, Connect-IB®, ConnectX®, CoolBox®, CORE-Direct®, EZchip®, EZchip logo, EZappliance®, EZdesign®, EZdriver®, EZsystem®, GPUDirect®, InfiniHost®, InfiniScale®, Kotura®, Kotura logo, Mellanox Federal Systems®, Mellanox Open Ethernet®, Mellanox ScalableHPC®, Mellanox Connect Accelerate Outperform logo, Mellanox Virtual Modular Switch®, MetroDX®, MetroX®, MLNX-OS®, NP-1c®, NP-2®, NP-3®, Open Ethernet logo, PhyX®, SwitchX®, Tiler®, Tiler logo, TestX®, The Generation of Open Ethernet logo, UFM®, Virtual Protocol Interconnect®, Voltaire® and Voltaire logo are registered trademarks of Mellanox Technologies, Ltd.

All other trademarks are property of their respective owners.

For the most updated list of Mellanox trademarks, visit <http://www.mellanox.com/page/trademarks>

Table of Contents

Table of Contents	1
List Of Tables	3
Release Update History	5
Chapter 1 Overview	6
1.1 Content of Mellanox OFED for Linux	6
1.2 Supported Platforms and Operating Systems	7
1.2.1 Supported Hypervisors	8
1.2.2 Supported Non-Linux Virtual Machines	8
1.2.3 Tested Hypervisors in Paravirtualized and SR-IOV Environments	9
1.3 Hardware and Software Requirements	9
1.4 Supported HCAs Firmware Versions	10
1.5 Compatibility Matrix	11
1.6 RoCE Modes Matrix	11
Chapter 2 Changes and New Features in Rev 3.2-2.0.0.0	12
2.1 API Changes in MLNX_OFED Rev 3.2-2.0.0.0	12
Chapter 3 Known Issues	13
3.1 Driver Installation/Loading/Unloading/Start Known Issues	13
3.1.1 Installation Known Issues	13
3.1.2 Driver Unload Known Issues	14
3.1.3 Driver Start Known Issues	14
3.1.4 System Time Known Issues	15
3.1.5 UEFI Secure Boot Known Issues	16
3.2 Performance Known Issues	16
3.3 HCAs Known Issues	18
3.3.1 ConnectX®-3 (mlx4 Driver) Known Issues	18
3.3.2 ConnectX®-4 (mlx5 Driver) Known Issues	18
3.4 Ethernet Network	19
3.4.1 Ethernet Known Issues	19
3.4.2 Port Type Management Known Issues	21
3.4.3 Flow Steering Known Issues	22
3.4.4 Quality of Service Known Issues	22
3.4.5 Ethernet Performance Counters Known Issues	22
3.5 InfiniBand Network	23
3.5.1 iPoIB Known Issues	23
3.5.2 eIPoIB Known Issues	26
3.5.3 XRC Known Issues	26
3.5.4 Verbs Known Issues	27
3.5.5 RoCE Known Issues	27
3.5.6 iSCSI over iPoIB Known Issues	28
3.6 Storage Protocols Known Issues	29
3.6.1 Storage Known Issues	29

3.6.2	SRP Known Issues	29
3.6.3	SRP Interop Known Issues	29
3.6.4	DDN Storage Fusion 10000 Target Known Issues	29
3.6.5	Oracle Sun ZFS Storage 7420 Known Issues	29
3.6.6	iSER Initiator Known Issues	30
3.6.7	iSER Target Known Issues	30
3.6.8	ZFS Appliance Known Issues	32
3.6.9	Erasure Coding Verbs Known Issues	32
3.7	Virtualization	32
3.7.1	SR-IOV Known Issues	32
3.8	Resiliency	35
3.8.1	Reset Flow Known Issues	35
3.9	Miscellaneous Known Issues	36
3.9.1	General Known Issues	36
3.9.2	ABI Compatibility Known Issues	36
3.9.3	Connection Manager (CM) Known Issues	36
3.9.4	Fork Support Known Issues	37
3.9.5	MLNX_OFED Sources Known Issues	37
3.9.6	Uplinks Known Issues	37
3.9.7	Resources Limitation Known Issues	37
3.9.8	Accelerated Verbs Known Issues	38
3.10	InfiniBand Fabric Utilities Known Issues	39
3.10.1	Performance Tools Known Issues	39
3.10.2	Diagnostic Utilities Known Issues	39
3.10.3	Tools Known Issues	39
Chapter 4	Bug Fixes History	40
Chapter 5	Change Log History	46
Chapter 6	API Change Log History	56

List Of Tables

Table 1:	Release Update History	5
Table 2:	Supported Uplinks to Servers	6
Table 3:	Mellanox OFED for Linux Software Components	6
Table 4:	Supported Platforms and Operating Systems	7
Table 5:	Tested Hypervisors in Paravirtualized and SR-IOV Environments	9
Table 6:	Additional Software Packages	10
Table 7:	Supported HCAs Firmware Versions	10
Table 8:	MLNX_OFED Rev 3.2-2.0.0.0 Compatibility Matrix	11
Table 9:	RoCE Modes Matrix	11
Table 10:	Changes in Rev 3.2-2.0.0.0	12
Table 11:	API Changes in MLNX_OFED Rev 3.2-2.0.0	12
Table 12:	Installation Known Issues	13
Table 13:	Driver Unload Known Issues	14
Table 14:	Driver Start Known Issues	14
Table 15:	System Time Known Issues	15
Table 16:	UEFI Secure Boot Known Issues	16
Table 17:	Performance Known Issues	16
Table 18:	ConnectX®-3 (mlx4 Driver) Known Issues	18
Table 19:	ConnectX-4 (mlx5 Driver) Known Issues	18
Table 20:	Ethernet Known Issues	19
Table 21:	Port Type Management Known Issues	21
Table 22:	Flow Steering Known Issues	22
Table 23:	Quality of Service Known Issues	22
Table 24:	Ethernet Performance Counters Known Issues	22
Table 25:	IPoIB Known Issues	23
Table 26:	eIPoIB Known Issues	26
Table 27:	XRC Known Issues	26
Table 28:	Verbs Known Issues	27
Table 29:	RoCE Known Issues	27
Table 30:	iSCSI over IPoIB Known Issues	28
Table 31:	Storage Known Issues	29
Table 32:	SRP Known Issues	29
Table 33:	SRP Interop Known Issues	29
Table 34:	DDN Storage Fusion 10000 Target Known Issues	29
Table 35:	Oracle Sun ZFS Storage 7420 Known Issues	29

Table 36: iSER Initiator Known Issues	30
Table 37: iSER Target Known Issues	30
Table 38: ZFS Appliance Known Issues	32
Table 39: Erasure Coding Verbs Known Issues	32
Table 40: SR-IOV Known Issues	32
Table 41: Resiliency Known Issues	35
Table 42: General Known Issues	36
Table 43: ABI Compatibility Known Issues	36
Table 44: Connection Manager (CM) Known Issues	36
Table 45: Fork Support Known Issues	37
Table 46: MLNX_OFED Sources Known Issues	37
Table 47: Uplinks Known Issues	37
Table 48: Resources Limitation Known Issues	37
Table 49: Accelerated Verbs Known Issues	38
Table 50: Performance Tools Known Issues	39
Table 51: Diagnostic Utilities Known Issues	39
Table 52: Tools Known Issues	39
Table 53: Fixed Bugs List	40
Table 54: Change Log History	46
Table 55: API Change Log History	56

Release Update History

Table 1 - Release Update History

Release	Date	Description
Rev 3.2-2.0.0.0	March 09, 2016	Added section Section 1.2.3, “Tested Hypervisors in Paravirtualized and SR-IOV Environments” , on page 9
	February 29, 2016	This is the initial release of this MLNX_OFED release.

1 Overview

These are the release notes of Mellanox OFED for Linux Driver, Rev 3.2-2.0.0.0. Mellanox OFED is a single Virtual Protocol Interconnect (VPI) software stack and operates across all Mellanox network adapter solutions supporting the following uplinks to servers:

Table 2 - Supported Uplinks to Servers

Uplink/HCAs	Uplink Speed
ConnectX®-4	<ul style="list-style-type: none"> InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GigE, 10GigE, 25GigE, 40GigE, 50GigE, 56GigE^a, and 100GigE
ConnectX®-4 Lx	<ul style="list-style-type: none"> Ethernet: 1GigE, 10GigE, 25GigE, 40GigE, and 50GigE
Connect-IB®	<ul style="list-style-type: none"> InfiniBand: SDR, QDR, FDR10, FDR
ConnectX®-3/ConnectX®-3 Pro	<ul style="list-style-type: none"> InfiniBand: SDR, QDR, FDR10, FDR Ethernet: 10GigE, 40GigE and 56GigE^a
ConnectX®-2	<ul style="list-style-type: none"> InfiniBand: SDR, DDR Ethernet: 10GigE, 20GigE
PCI Express 2.0	2.5 or 5.0 GT/s
PCI Express 3.0	8 GT/s

- a. 56 GbE is a Mellanox propriety link speed and can be achieved while connecting a Mellanox adapter cards to Mellanox SX10XX switch series or connecting a Mellanox adapter card to another Mellanox adapter card.

1.1 Content of Mellanox OFED for Linux

Mellanox OFED for Linux software contains the following components:

Table 3 - Mellanox OFED for Linux Software Components

Components	Description
OpenFabrics core and ULPs	<ul style="list-style-type: none"> InfiniBand and Ethernet HCA drivers (mlx4, mlx5) core Upper Layer Protocols: IPoIB, SRP, iSER and iSER Initiator and Target
OpenFabrics utilities	<ul style="list-style-type: none"> OpenSM: IB Subnet Manager with Mellanox proprietary Adaptive Routing Diagnostic tools Performance tests SSA (SLES12): libopensmssa plugin for OpenSM, ibssa, ibacm
MPI	<ul style="list-style-type: none"> OSU MPI (mvapich2-2.0) stack supporting the InfiniBand interface Open MPI stack 1.6.5 and later supporting the InfiniBand interface MPI benchmark tests (OSU benchmarks, Intel MPI benchmarks, Presta)
PGAS	<ul style="list-style-type: none"> HPC-X OpenSHMEM v2.2 supporting InfiniBand, MXM and FCA HPC-X UPC v2.2 supporting InfiniBand, MXM and FCA

Table 3 - Mellanox OFED for Linux Software Components

Components	Description
HPC Acceleration packages	<ul style="list-style-type: none"> Mellanox MXM v3.0 (p2p transport library acceleration over InfiniBand) Mellanox FCA v2.5 (MPI/PGAS collective operations acceleration library over InfiniBand) KNEM, Linux kernel module enabling high-performance intra-node MPI/PGAS communication for large messages
Extra packages	<ul style="list-style-type: none"> ibutils2 ibdump MFT
Sources of all software modules (under conditions mentioned in the modules' LICENSE files) except for MFT, OpenSM plugins, ibutils2, and ibdump	
HCA's	<ul style="list-style-type: none"> ConnectX-4 EN driver Rev 3.2-2.0.0.0 ConnectX-3 EN driver Rev 3.2-2.0.0.0
Documentation	

1.2 Supported Platforms and Operating Systems

The following are the supported OSs in MLNX_OFED Rev 3.2-2.0.0.0:

Table 4 - Supported Platforms and Operating Systems

Operating System	Platform
RHEL6.2/ CentOS6.2	x86_64
RHEL6.5/ CentOS6.5	x86_64
RHEL6.6/ CentOS6.6	x86_64/PPC64
RHEL6.7/ CentOS6.7	x86_64/PPC64
RHEL7.0/ CentOS7.0	x86_64/PPC64
RHEL7.1/ CentOS7.1	x86_64/PPC64/PPC64LE (Power8)
RHEL7.2/ CentOS7.2	x86_64/PPC64/PPC64LE (Power8)
Debian 6.0.10	x86_64
Debian 7.6	x86_64
Debian 8.2	x86_64
Debian 8.1	x86_64
Fedora 19	x86_64
Fedora 20	x86_64
Fedora 21	x86_64/PPC64LE (Power 8)
Fedora 22	x86_64/PPC64LE (Power 8)
Fedora 23	x86_64/PPC64LE (Power 8)/ ARM (AMD) (ARM is at beta level)
OEL 6.5	x86_64
OEL 6.6	x86_64
OEL 6.7	x86_64 (UEK)
OEL 7.1	x86_64 (UEK 3)
Sles10 SP3	x86_64
Sles11 SP1	x86_64

Table 4 - Supported Platforms and Operating Systems

Operating System	Platform
Sles11 SP2	x86_64
Sles11 SP3	x86_64/PPC64 (Power7)
Sles11 SP4	x86_64/PPC64
Sles12	x86_64/PPC64LE (Power 8)
Sles12SP1	x86_64/PPC64LE
Ubuntu 12.04.4	x86_64
Ubuntu 14.04	x86_64/PPC64LE (Power 8)/ARM (ARM is at beta level)
Ubuntu 14.10	x86_64/PPC64LE (Power8)
Ubuntu 15.04	x86_64/PPC64LE (Power 8)
Ubuntu 15.10	x86_64/PPC64LE (Power8)
Xen 4.2	x86_64
Wind River	x86_64
Kernels	3.10.28, 3.11.10, 3.14.3, 3.15, 3.16, 3.17, 3.18, 3.19, 4.0, 4.1, 4.2, 4.3, 4.4



For RPM based Distributions, if you wish to install OFED on a different kernel, you need to create a new ISO image, using `mlnx_add_kernel_support.sh` script. See the MLNX_OFED User Guide for instructions.



Upgrading MLNX_OFED on your cluster requires upgrading all of its nodes to the newest version as well.

1.2.1 Supported Hypervisors

The following are the supported Hypervisors in MLNX_OFED Rev 3.2-2.0.0.0:

- KVM:
 - RedHat/CentOS 6.6, 6.7, 7.1, 7.2
 - Ubuntu 14.10, 15.10
 - SLES11SP3 (InfiniBand only), SLES11SP4, SLES12, SLES12SP1
 - Debian 6.0.10

1.2.2 Supported Non-Linux Virtual Machines

The following are the supported Non-Linux (InfiniBand only) Virtual Machines in MLNX_OFED Rev 3.2-2.0.0.0:

- Windows Server 2016 (Beta)
- Windows Server 2012 R2
- Windows Server 2012
- Windows Server 2008 R2

1.2.3 Tested Hypervisors in Paravirtualized and SR-IOV Environments

Table 5 - Tested Hypervisors in Paravirtualized and SR-IOV Environments

Tested Hypervisors	HCA's	Operating System
SR-IOV	ConnectX-4 Lx	RHEL7.1 Ethernet
		XenServer6.5 PV (Partially tested)
	ConnectX-4	RHEL7.0
		RHEL6.7
		Ubuntu 15.10
		SLES11 SP4 KVM
		Xen4.2 Ethernet + RoCE native
		XenServer6.5 PV (Partially tested)
		Debian 6.0.10 + ConnectX-4 SRIOV (DDN)
		ConnectX-3 Pro
	ConnectX-3	RHEL6.7 InfiniBand
		SLES12 Ethernet
		SLES11 SP4 KVM
	Connect-IB	RHEL7.0
		RHEL 6.7
		Debian 6.0.10 SRIOV (DDN)
Paravirtualized	ConnectX-4 Lx	Ubuntu 15.10
		SLES11 SP4 KVM
	ConnectX-4	RH6.7 KVM
	ConnectX-3	RHEL6.5 eIPoIB
		RHEL6.5 Ethernet
		Ubuntu 15.04 IPoIB
		RH6.7 KVM
		Ubuntu 15.10
		SLES11.4 KVM

1.3 Hardware and Software Requirements

The following are the hardware and software requirements of MLNX_OFED Rev 3.2-2.0.0.0.

- Linux operating system
- Administrator privileges on your machine(s)
- Disk Space: 1GB

For the OFED Distribution to compile on your machine, some software packages of your operating system (OS) distribution are required.

To install the additional packages, run the following commands per OS:

Table 6 - Additional Software Packages

Operating System	Required Packages Installation Command
RHEL/OEL/Fedora	yum install perl pciutils python gcc-gfortran libxml2-python tesh libnl.i686 libnl expat glib2 tcl libstdc++ bc tk gtk2 atk cairo numactl pkgconfig
XenServer	yum install perl pciutils python libxml2-python libnl expat glib2 tcl bc libstdc++ tk pkgconfig
SLES 10 SP3	zypper install pkgconfig pciutils python libxml2-python libnl lsof expat glib2 tcl libstdc++ bc tk
SLES 11 SP2	zypper install perl pciutils python libnl-32bit libxml2-python tesh libnl libstdc++46 expat glib2 tcl bc tklibcurl4 gtk2 atk cairo pkg-config
SLES 11 SP3	zypper install perl pciutils python libnl-32bit libxml2-python tesh libstdc++43 libnl expat glib2 tcl bc tk libcurl4 gtk2 atk cairo pkg-config
SLES 12	zypper install pkg-config expat libstdc++6 libglib-2_0-0 libgtk-2_0-0 tcl libcairo2 tesh python bc pciutils libatk-1_0-0 tk python-libxml2 lsof libnl1
Ubuntu/Debian	apt-get install perl dpkg autotools-dev autoconf libtool automake1.10 automake m4 dkms debhelper tcl tcl8.4 chrpath swig graphviz tcl-dev tcl8.4-dev tk-dev tk8.4-dev bison flex dpatch zlib1g-dev curl libcurl4-gnutls-dev python-libxml2 libvirt-bin libvirt0 libnl-dev libglib2.0-dev libgfortran3 automake m4 pkg-config libnuma logrotate
Debian 8	apt-get install libnl-3-200 automake debhelper curl dkms logrotate libglib2.0-0 python-libxml2 graphviz tk tcl libvirt-bin coreutils pkg-config autotools-dev flex autoconf pciutils quilt module-init-tools libvirt0 libstdc++6 dpkg libgfortran3 procs lsof libltdl-dev gcc dpatch chrpath grep m4 gfortran bison libnl-route-3-200 swig perl make

1.4 Supported HCAs Firmware Versions

MLNX_OFED Rev 3.2-2.0.0.0 supports the following Mellanox network adapter cards firmware versions:

Table 7 - Supported HCAs Firmware Versions

HCA	Recommended Firmware Rev.	Additional Firmware Rev. Supported
Connect-IB®	10.14.2036	10.14.1100
ConnectX®-4 Lx	14.14.2036	14.14.1100
ConnectX®-4	12.14.2036	12.14.1100
ConnectX®-3 Pro	2.36.5000	2.35.5100
ConnectX®-3	2.36.5000	2.35.5100
ConnectX®-2	2.9.1000	2.9.1000

For official firmware versions please see:

http://www.mellanox.com/content/pages.php?pg=firmware_download

1.5 Compatibility Matrix

MLNX_OFED Rev 3.2-2.0.0.0 is compatible with the following:

Table 8 - MLNX_OFED Rev 3.2-2.0.0.0 Compatibility Matrix

Mellanox Product	Description/Version
MLNX-OS®	MSX6036 w/w MLNX-OS® version 3.4.3202 ^a
Grid Director™	4036 w/w Grid Director™ version 3.9.1-985
Unified Fabric Manager (UFM®)	v5.5
MXM	v3.4
HPC-X UPC	v2.22
HPC-X OpenSHMEM	v1.8.3
FCA	v2.5 and v3.4
OpenMPI	v1.10
MVAPICH	v2.0

- a. MLNX_OFED Rev 3.2-2.0.0.0 was tested with this switch however, additional switches might be supported as well.

1.6 RoCE Modes Matrix

The following is RoCE modes matrix:

Table 9 - RoCE Modes Matrix

Software Stack / Inbox Distribution	RoCEv1 (Layer 2) Supported as of Version	RoCEv2 (Layer 3) Supported as of Version	RoCEv1 & RoCEv2 (Layer 3) Supported as of Version
MLNX_OFED	2.1-x.x.x	2.3-x.x.x	3.0-x.x.x
Kernel.org	3.14	4.4	-
RHEL	6.6; 7.0	-	-
SLES	12	-	-
Ubuntu	14.04	-	-

2 Changes and New Features in Rev 3.2-2.0.0.0

Table 10 - Changes in Rev 3.2-2.0.0.0

Feature/Change	Description
API Changes	<ul style="list-style-type: none"> Support FCS scattering for Raw Packet QPs and WQs. Indication of L4 packet type on the receive side completions Support CVLAN insertion for WQs <p>For further information, please see Section 2.1, “API Changes in MLNX_OFED Rev 3.2-2.0.0.0”, on page 12.</p>
IPoIB	<ul style="list-style-type: none"> Added support for the following IPoIB UD QP offloads: <ul style="list-style-type: none"> RX check summing (AKA RX csu) Large Send Offloads (AKA LSO) <p>To see the new IPoIB UD mode, run: "ethtool -k <interface>"</p>

2.1 API Changes in MLNX_OFED Rev 3.2-2.0.0.0

The following are the API additions/changes in MLNX_OFED Rev 3.2-2.0.0:

Table 11 - API Changes in MLNX_OFED Rev 3.2-2.0.0

Release	Name	Description
Rev 3.2-2.0.0.0	libibverbs	<ul style="list-style-type: none"> Support FCS scattering for Raw Packet QPs and WQs. <ul style="list-style-type: none"> Query: <code>ibv_exp_query_device</code> reports <code>IBV_EXP_DEVICE_SCATTER_FCS</code> when it is supported. Enablement of this feature is done in the creation: <ol style="list-style-type: none"> For Raw Packet QPs: Set <code>IBV_EXP_QP_CREATE_SCATTER_FCS</code> in <code>exp_create_flags</code>. For WQs: Set <code>IBV_EXP_CREATE_WQ_FLAG_SCATTER_FCS</code> in flags of <code>ibv_exp_wq_init_attr</code>. Indication of L4 packet type on the receive side completions: <ul style="list-style-type: none"> Query: <code>ibv_exp_query_device</code> reports <code>IBV_EXP_DEVICE_RX_TCP_UDP_PKT_TYPE</code> when it is supported. <code>ibv_exp_cq_family_flags</code> was extended with two flags <code>IBV_EXP_CQ_RX_TCP_PACKET</code> and <code>IBV_EXP_CQ_RX_UDP_PACKET</code> to support L4 packet type when using <code>poll_length_flags()</code>. Support CVLAN insertion for WQs: <ul style="list-style-type: none"> Query: <code>IBV_EXP_RECEIVE_WQ_CVLAN_INSERTION</code> is set in <code>ibv_exp_vlan_offloads</code> when CVLAN insertion is supported. Enablement: The <code>ibv_exp_qp_burst_family</code> was extended to support CVLAN insertion: <ol style="list-style-type: none"> <code>send_pending_vlan</code>: Put one message in the provider send queue and insert <code>vlan_tci</code> to header. <code>send_pending_inline_vlan</code>: Put one inline message in the provider send queue and insert <code>vlan_tci</code> to header. <code>send_pending_sg_list_vlan</code>: Put one scatter-gather(sg) message in the provider send queue and insert <code>vlan_tci</code> to header.

3 Known Issues

The following is a list of general limitations and known issues of the various components of this Mellanox OFED for Linux release.

3.1 Driver Installation/Loading/Unloading/Start Known Issues

3.1.1 Installation Known Issues

Table 12 - Installation Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When upgrading from an earlier Mellanox OFED version, the installation script does not stop the earlier version prior to uninstalling it.	Stop the old OFED stack (/etc/init.d/openibd stop) before upgrading to this new version.
2.	Upgrading from the previous OFED installation to this release, does not unload the kernel module ipoib_helper.	Reboot after installing the driver.
3.	"--total-vfs <0-63>" installation parameter is no longer supported	Use '--enable-sriov' installation parameter to burn firmware with SR-IOV support. The number of virtual functions (VFs) will be set to 16. For further information, please refer to the User Manual.
4.	When using bonding on Ubuntu OS, the "ifenslave" package must be installed.	-
5.	On PPC systems, the ib_srp module is not installed by default since it breaks the ibmvscsi module.	If your system does not require the ibmvscsi module, run the mlnxofedinstall script with the "--with-srp" flag.
6.	#679801: Updating MLNX_OFED via Yum (e.g. running "yum update mlnx-ofed-all") can fail with the following error: --> Finished Dependency Resolution Error: Package: mpitests_openmpi_1_8_8-3.2.16-fe5387c.x86_64 (installed) Requires: liboshmem.so.3()(64bit) Removing: openmpi-1.8.8-1.x86_64 (installed) liboshmem.so.3()(64bit) Updated By: openmpi-1.10.2rc4-1.32008.x86_64 (mlnx_ofed) ~liboshmem.so.9()(64bit)	Remove the mpitests packages manually: # rpm -e --allmatches \$(rpm -qa grep mpitests_)
7.	#690799: OpenSM package removal fails with the following error on Ubuntu12.04: Removing opensm ... /sbin/insserv: No such file or directory	1. Create the missing link by running this command: # ln -s /usr/lib/insserv/insserv /sbin/insserv 2. Remove the package.

3.1.2 Driver Unload Known Issues

Table 13 - Driver Unload Known Issues

Index	Internal Reference Number: Description	Workaround
1.	"openibd stop" can sometime fail with the error: Unloading ib_cm [FAILED] ERROR: Module ib_cm is in use by ib_i- poib	Re-run "openibd stop"

3.1.3 Driver Start Known Issues

Table 14 - Driver Start Known Issues

Index	Internal Reference Number: Description	Workaround
1.	"Out-of-memory" issues may rise during drivers load depending on the values of the driver module parameters set (e.g. log_num_cq).	-
2.	When reloading/starting the driver using the /etc/init.d/openibd the following messages are displayed if there is a third party RPM or driver installed: "Module mlx4_core does not belong to MLNX_OFED" or "Module mlx4_core belong to <rpm name> which is not a part of MLNX_OFED"	Remove the third party RPM/non MLNX_OFED drivers directory, run: "depmod" and then rerun "/etc/init.d/openibd restart"
3.	Occasionally, when trying to repetitively reload the NES hardware driver on SLES11 SP2, a soft lockups occurs that required reboot.	-
4.	When downgrading from MLNX_OFED 3.0-x.x.x, driver reload might fail with the following errors in dmesg: [166271.886407] compat: exports duplicate symbol __ethtool_get_settings (owned by mlx_compat)	The issues will be resolved automatically after system reboot or by invoking the following commands: rmmod mlx_compat depmod -a /etc/init.d/openibd restart
5.	In ConnectX-2, (when the debug_level module parameter for module mlx4_core is non-zero), if the driver load succeeds, the message below is presented: "mlx4_core 0000:0d:00.0: command SET_PORT (0xc) failed: in_param=0x120064000, in_mod=0x2, op_mod=0x0, fw status = 0x40" This message is simply part of the learning process for setting the maximum port VLs compatible with a 4K port mtu, and should be ignored.	-
6.	"openibd start" unloads kernel modules that were loaded from initrd/initramfs upon boot. This affects only kernel modules which come with MLNX_OFED and are included in initrd/initramfs.	-

Table 14 - Driver Start Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
7.	If a Lustre storage is used, it must be fully unloaded before restarting the driver or rebooting the machine, otherwise machine might get stuck/panic.	1. Unmount any mounted Lustre storages: # umount<lustre_mount_point> 2. Unload all Lustre modules: # lustre_rmmod
8.	Driver unload fails with the following error message: Unloading rdma_cm [FAILED] rmmod: ERROR: Module rdma_cm is in use by: xprtrdma	Make sure that there are no mount points over NFS/RDMA prior to unloading the driver and run: # modprobe -r xprtrdma In case that the xprtrdma module keeps getting loaded automatically even though it is not used, add a pre-stop hook for the openibd service script to always unload it. Create an executable file "/etc/infiniband/pre-stop-hook.sh" with the following content: #!/bin/bash modprobe -r xprtrdma
9.	When loading or unloading the driver on HP ProLiant systems, you may see log messages like: dmar: DMAR:[DMA Write] Request device [07:00.0] fault addr 3df7f000 DMAR:[fault reason 05] PTE Write access is not set This is a known issue with ProLiant systems (see their support notice for Emulex adapters: http://h20564.www2.hp.com/hpsc/doc/public/display?docId=emr_na-c04446026&lang=en-us&cc=us) The messages are harmless, and may be ignored.	If you are <i>*not*</i> running SR-IOV on your system, you may eliminate these messages by removing the term "intel_iommu=on" from the boot line in file /boot/grub/menu.lst. For SR-IOV systems, this term must remain, you can ignore the log messages.
10.	#677998: False alarms errors may be printed to dmesg	-
11.	#610395: On RHEL 7.1, after updating to kernel version 3.10.0-229.14.1.el7 or later, driver load fails with unknown symbols errors in dmesg.	Use the <code>mlnx_add_kernel_support.sh</code> script to compile MLNX_OFED drivers against the new kernel.

3.1.4 System Time Known Issues

Table 15 - System Time Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Loading the driver using the openibd script when no InfiniBand vendor module is selected (for example <code>mlx4_ib</code>), may cause the execution of the <code>/sbin/start_udev</code> script. In RedHat 6.x and OEL6.x this may change the local system time.	-

3.1.5 UEFI Secure Boot Known Issues

Table 16 - UEFI Secure Boot Known Issues

Index	Internal Reference Number: Description	Workaround
1.	<p>On RHEL7 and SLES12, the following error is displayed in dmesg if the Mellanox's x.509 Public Key is not added to the system:</p> <pre>[4671958.383506] Request for unknown module key 'Mellanox Technologies signing key: 61feb074fc7292f958419386ffdd9d5-ca999e403' err -11</pre> <p>This error can be safely ignored as long as Secure Boot is disabled on the system.</p>	For further information, please refer to the User Manual section "Enrolling Mellanox's x.509 Public Key On your Systems".
2	Ubuntu12 requires update of user space open-iscsi to v2.0.873	-
3	The initiator does not respect interface parameter while logging in.	Configure each interface on a different subnet.

3.2 Performance Known Issues

Table 17 - Performance Known Issues

Index	Internal Reference Number: Description	Workaround
1.	On machines with irqbalancer daemon turned off, the default InfiniBand interrupts will be routed to a single core which may cause overload and software/hardware lockups.	Execute the following script as root: <pre>set_irq_affinity.sh <interface or IB device> [2nd interface or IB device]</pre>
2.	#414827: Out-of-the-box throughput performance in Ubuntu14.04 is not optimal and may achieve results below the line rate in 40GE link speed.	For additional performance tuning, please refer to Performance Tuning Guide.
3.	<p>UDP receiver throughput may be lower than expected, when running over mlx4_en Ethernet driver.</p> <p>This is caused by the adaptive interrupt moderation routine, which sets high values of interrupt coalescing, causing the driver to process large number of packets in the same interrupt, leading UDP to drop packets due to overflow in its buffers.</p>	<p>Disable adaptive interrupt moderation and set lower values for the interrupt coalescing manually.</p> <pre>ethtool -C <eth>X adaptive-rx off rx-usecs 64 rx-frames 24</pre> <p>Values above may need tuning, depending the system, configuration and link speed.</p>
4.	Performance degradation might occur when bonding Ethernet interfaces	-
5.	#656415: In RHEL7.0, when the irqbalance service is started or restarted, it incorrectly re-balances the IRQs, including the banned ones.	-

Table 17 - Performance Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
6.	#651322: In RH7.0/RH7.1, performance issue with ConnectX-4 cards over 100GbE link might occur when the process of forwarding the packets between the ports, which is done by the kernel, fib_table_lookup() function is called. For further information, please refer to: http://comments.gmane.org/gmane.linux.network/344243	Use RH7.2 to avoid such performance issues.

3.3 HCAs Known Issues

3.3.1 ConnectX®-3 (mlx4 Driver) Known Issues

Table 18 - ConnectX®-3 (mlx4 Driver) Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Using RDMA READ with a higher value than 30 SGEs in the WR might lead to "local length error".	Do not set the value of SGEs higher than 30 when RDMA READ is used.

3.3.2 ConnectX®-4 (mlx5 Driver) Known Issues

Table 19 - ConnectX-4 (mlx5 Driver) Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Atomic Operations in Connect-IB® are fully supported on big-endian machines (e.g. PPC). Their support is limited on little-endian machines (e.g. x86)	-
2.	#435583: EEH events that arrive while the mlx5 driver is loading may cause the driver to hang.	-
3.	#434570: The mlx5 driver can handle up to 5 EEH events per hour.	If more events are received, cold reboot the machine.
4.	#554120: When working with Connect-IB® firmware v10.10.5054, the following message would appear in driver start. command failed, status bad system state(0x4), syndrome 0x408b33 The message can be safely ignored.	Upgrade Connect-IB firmware to the latest available version.
5.	Changing the link speed is not supported in Ethernet driver when connected to a ConnectX-4 card.	-
6.	#538843: Bonding "active-backup" mode does not function properly.	-
7.	Rate, speed and width using IB sysfs/tools are available in RoCE mode in ConnectX-4 only after port physical speed configuration is done.	-
8.	#598092: Since MLNX_OFED's openibd does not unload modules while OpenSM is running, removing <code>mlx5_core</code> manually while OpenSM is running, may cause it to be out of sync when probed again.	Restart OpenSM
9.	#563022: ConnectX-4 port GIDs table shows a duplicated RoCE v2 default GID.	-

3.4 Ethernet Network

3.4.1 Ethernet Known Issues



Ethernet Known Issues are applicable to ConnectX-3/ConnectX-3 Pro only.

Table 20 - Ethernet Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When creating more than 125 VLANs and SR-IOV mode is enabled, a kernel warning message will be printed indicating that the native VLAN is created but will not work with RoCE traffic. kernel warning: mlx4_core 0000:07:00.0: vhcr command ALLOC_RES (0xf00) slave:0 in_param 0x7e in_mod=0x107, op_mod=0x1 failed with error:0, status -28	-
2.	Kernel panic might occur during FIO splice in kernels before 2.6.34-rc4.	Use kernel v2.6.34-rc4 which provides the following solution: baff42a net: Fix oops from tcp_collapse() when using splice()
3.	In PPC systems when QoS is enabled a harmless Kernel DMA mapping error messages might appear in kernel log (IOMMU related issue).	-
4.	Transmit timeout might occur on RH6.3 as a result of lost interrupt (OS issue). In this case, the following message will be shown in dmesg: do_IRQ: 0.203 No irq handler for vector (irq -1)	-
5.	Mixing ETS and strict QoS policies for TCs in 40GbE ports may cause inaccurate results in bandwidth division among TCs.	-
6.	Creating a VLAN with user priority >= 4 on ConnectX®-2 HCA is not supported.	-
7.	Affinity hints are not supported in Xen Hypervisor (an irqblancer issue). This causes a non-optimal IRQ affinity.	To overcome this issues, run: set_irq_affinity.sh eth<x>
8.	#433366: Reboot might hang in SR-IOV when using the "probe_vf" parameter with many Virtual Functions. The following message is logged in the kernel log: "waiting for eth to become free. Usage count =1"	-
9.	In ConnectX®-2, RoCE UD QP does not include VLAN tags in the Ethernet header	

Table 20 - Ethernet Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
10.	VXLAN may not be functional when configured over Linux bridge in RH7.0 or Ubuntu14.04. The issue is within the bridge modules in those kernels. In Vanilla kernels above 3.16 issues is fixed.	-
11.	In RH6.4, ping may not work over VLANs that are configured over Linux bridge when the bridge has a mlx4_en interface attached to it.	-
12.	The interfaces LRO needs to be set to "OFF" manually when there is a bond configured on Mellanox interfaces with a Bridge over that bond.	Run: <code>ethtool -K ethX lro off</code>
13.	#539117: On SLES12, the bonding interface over Mellanox Ethernet slave interfaces does not get IP address after reboot.	<ol style="list-style-type: none"> 1. Set "STARTMODE=hotplug" in the bonding slave's ifcfg files. More details can be found in the SUSE documentations page: https://www.suse.com/documentation/sles-12/book_sle_admin/?page=/documentation/sles-12/book_sle_admin/data/sec_bond.html 2. Enable the "nanny" service to support hotplugging: Open the "/etc/wicked/common.xml" file. Change: "<use-nanny>false</use-nanny>" to "<use-nanny>>true</use-nanny>" 3. Run: <pre># systemctl restart wickedd.service wicked</pre>
14.	ethtool -x command does not function in SLES OS.	-
15.	#516136: Ethertype proto 0x806 not supported by ethtool	-
16.	ETS configuration is not supported in the following kernels: <ul style="list-style-type: none"> • 3.7 • 3.8 • 3.9 • 3.10 • 3.2.37-94_fbk17_01925_g8e3b329 • 3.14 • 3.2.55-106_fbk22_00877_g6902630 • 3.2.28-76_fbk14_00230_g3c40d9e 	
17.	ETS is not supported in kernels that do not have MQPRIO as QDISC_KIND option in the tc tool.	-
18.	#592229: When NC-SI is ON, the port's MTU cannot be set to lower than 1500.	-
19.	#600242: GRO is not functional when using VXLAN in ConnectX-3 adapter cards.	-
20.	#596075: ethtool -X: The driver supports only the 'equal' mode and cannot be set by using weight flags.	-

Table 20 - Ethernet Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
21.	#600752: Q-in-Q infrastructure in the kernel is supported only in kernel version 3.10 and up.	-
22.	#596537: When SLES11 SP4 is used as a DHCP client over ConnectX-3 or ConnectX-3 adapters, it might fail to get an IP from the DHCP server.	-
23.	#560575: When using a hardware that has Time Stamping enabled, the system time might be higher than the expected variance.	-
24.	#597758: In Q-in-Q, ping failed when sending traffic with package size > 1468	-
25.	#665131: Call trace may occur when configuring VXLAN or under high traffic stress.	-
26.	HW LRO does not function in ConnectX-4 adapter cards.	-
27.	#685069: ethtool header does not currently support link speeds 25/50/100, therefore these speeds cannot be seen as advertised/supported.	-

3.4.2 Port Type Management Known Issues

Table 21 - Port Type Management Known Issues

Index	Internal Reference Number: Description	Workaround
1.	OpenSM must be stopped prior to changing the port protocol from InfiniBand to Ethernet.	-
2.	After changing port type using <code>connectx_port_config</code> interface ports' names can be changed. For example. <code>ib1 -> ib0</code> if port1 changed to be Ethernet port and port2 left IB.	Use udev rules for persistent naming configuration. For further information, please refer to the User Manual
3.	A working IP connectivity between the RoCE devices is required when creating an address handle or modifying a QP with an address vector.	-
4.	IPv4 multicast over RoCE requires the MGID format to be as follow <code>::ffff:<Multicast IPv4 Address></code>	-
5.	IP routable RoCE does not support Multicast Listener Discovery (MLD) therefore, multicast traffic over IPv6 may not work as expected.	-
6.	DIF: When running IO over FS over DM during unstable ports, block layer BIOS merges may cause false DIF error.	-
7.	<code>connectx_port_config</code> configurations is not saved after unbind/bind.	Re-run " <code>connectx_port_config</code> "

3.4.3 Flow Steering Known Issues

Table 22 - Flow Steering Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Flow Steering is disabled by default in firmware version < 2.32.5100.	To enable it, set the parameter below as follow: log_num_mgm_entry_size should set to -1
2.	IPv4 rule with source IP cannot be created in SLES 11.x OSES.	-
3.	RFS does not support UDP.	-
4.	When working in DMFS:A0 mode and VM/hypervisor is MLNX_OFED 2.3-x.x.x, the second side (hypervisor/VM respectively) should be MLNX_OFED 2.3-x.x.x as well.	-
5.	#516136: Setting ARP flow rules through ethtool is not allowed.	-

3.4.4 Quality of Service Known Issues

Table 23 - Quality of Service Known Issues

Index	Internal Reference Number: Description	Workaround
1.	QoS is not supported in XenServer, Debian 6.0 and 6.2 with uek kernel	-
2.	When QoS features are not supported by the kernel, mlnx_qos tool may crash.	-
3.	#448981: QoS default settings are not returned after configuring QoS.	-

3.4.5 Ethernet Performance Counters Known Issues

Table 24 - Ethernet Performance Counters Known Issues

Index	Internal Reference Number: Description	Workaround
1.	In ConnectX®-3, in a system with more than 61 VFs, the 62nd VF and onwards is assigned with the SINKQP counter, and as a result will have no statistics, and loopback prevention functionality for SINK counter.	-
2.	In ConnectX®-3, since each VF tries to allocate 2 more QP counter for its RoCE traffic statistics, in a system with less than 61 VFs, if there is free resources it receives new counter otherwise receives the default counter which is shared with Ethernet. In this case RoCE statistics is not available.	-
3.	In ConnectX®-3, when we enable function-based loopback prevention for Ethernet port by default (i.e., based on the QP counter index), the dropped self-loopback packets increase the IfRxErrorFrames/Octets counters.	-

3.5 InfiniBand Network

3.5.1 IPoIB Known Issues

Table 25 - IPoIB Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When user increases receive/send a buffer, it might consume all the memory when few child's interfaces are created.	-
2.	The size of send queue in Connect-IB® cards cannot exceed 1K.	-
3.	In 32 bit devices, the maximum number of child interfaces that can be created is 16. Creating more that, might cause out-of-memory issues.	-
4.	In RHEL7.0, the Network-Manager can detect when the carrier of one of the IPoIB interfaces is OFF and can decide to disable its IP address.	Set "ignore-carrier" for the corresponding device in NetworkManager.conf. For further information, please refer to " <i>man NetworkManager.conf</i> "
5.	IPoIB interface does not function properly if a third party application changes the PKey table. We recommend modifying PKey tables via OpenSM.	-
6.	Fallback to the primary slave of an IPoIB bond does not work with ARP monitoring. (https://bugs.openfabrics.org/show_bug.cgi?id=1990)	-
7.	Out-of memory issue might occur due to overload of interfaces created.	To calculate the allowed memory per each IPoIB interface check the following: <ul style="list-style-type: none"> • Num-rings = min(num-cores-on-that-device, 16) • Ring-size = 512 (by default, it is module parameter) • UD memory: 2 * num-rings * ring-size * 8K • CM memory: ring-size * 64k • Total memory = UD mem + CM mem
8.	Connect-IB does not reach the bidirectional line rate	Optimize the IPoIB performance in Connect-IB: <pre>cat /sys/class/net/<interface>/device/local_cpus > /sys/class/net/<interface>/queues/rx-0/rps_cpus</pre>
9.	If the CONNECTED_MODE=no parameter is set to "no" or missing from the ifcfg file for Connect-IB® IPoIB interface then the "service network restart" will hang.	Set the CONNECTED_MODE=yes parameter in the ifcfg file for Connect-IB® interface.
10.	Joining a multicast group in the SM using the RDMA_CM API requires IPoIB to first join the broadcast group.	-

Table 25 - IPoIB Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
11.	<p>Whenever the IOMMU parameter is enabled in the kernel it can decrease the number of child interfaces on the device according to resource limitation. The driver will stuck after unknown amount of child interfaces creation.</p> <p>For further information, please see: https://access.redhat.com/site/articles/66747 http://support.citrix.com/article/CTX136517 http://www.novell.com/support/kb/doc.php?id=7012337 https://bugzilla.redhat.com/show_bug.cgi?id=1044595</p>	<p>To avoid such issue:</p> <ul style="list-style-type: none"> • Decrease the amount of the RX receive buffers (module parameter, the default is 512) • Decrease the number of RX rings (sys/fs or ethtool in new kernels) • Avoid using IOMMU if not required <p>For KVM users: Run: <pre>echo 1 > /sys/module/kvm/parameters/allow_unsafe_assigned_interrupts</pre></p> <p>To make this change persist across reboots, add the following to the <code>/etc/modprobe.d/kvm.conf</code> file (or create this file, if it does not exist): <pre>options kvm allow_unsafe_assigned_interrupts=1 kernel parameters</pre></p>
12.	<p>System might crash in <code>skb_checksum_help()</code> while performing TCP retransmit involving packets with 64k packet size.</p> <p>A similar out to the below will be printed: kernel BUG at net/core/dev.c:1707! invalid opcode: 0000 [#1] SMP RIP: 0010: [<ffffffff81448988>] skb_checksum_help+0x148/0x160 Call Trace: <IRQ> [<ffffffff81448d83>] dev_hard_start_xmit+0x3e3/0x530 [<ffffffff8144c805>] dev_queue_xmit+0x205/0x550 [<ffffffff8145247d>] neigh_connected_output+0xbd/0x1 </p>	Use UD mode in IPoIB
13.	When InfiniBand ports are removed from the host (e.g when changing port type from IB to Eth or removing a card from the PCI bus) the remaining IPoIB interface might be renamed.	<p>To avoid it and have persistent IPoIB network devices names for ConnectX ports, add to the <code>/etc/udev/rules.d/70-persistent-net.rules</code> file:</p> <pre>SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?*", ATTR{address}=="*<Port GID>", NAME="ibN"</pre> <p>Where N is the IPoIB required interface index</p>
14.	After releasing a bond interface that contains IPoIB slaves, a call trace might be printed to the dmesg.	-

Table 25 - IPoIB Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
15.	IPoIB interfaces are loaded without an IP address on SLES 12.	<ol style="list-style-type: none"> 1. Open the "/etc/wicked/common.xml" file. 2. Change: <pre>"<use-nanny>>false</use-nanny>" to "<use-nanny>>true</use-nanny>"</pre> 3. Run: <pre># systemctl restart wickedd.ser- vice wicked # ifup all</pre>
16.	In RHEL7.0, running ifdown then ifup on an interface after changing CONNECTED_MODE in its ifcfg file, will cause the interface bring up to fail.	Reload the driver "/etc/init.d/openibd restart" or reboot the system.
17.	Clone interfaces receive a duplicated IPv6 address when a child interface with the same PKey (a.k.a clone interface) is used for all the clones.	-
18.	eth_ipoib module is not loaded after reloading the ib_ipoib module.	To restart the IPoIB driver, run "/etc/init.d/openibd restart". Do not restart it by manually restarting each module.
19.	In Ubuntu and Debian, the default of the recv_queue_size and send_queue_size is 128 according to the io_mmu issue.	-
20.	In RHEL6.7, when the Network Manager service is enabled and an IPoIB interface is configured using the "nm-connection-editor" tool, the generated ifcfg file is missing the "DEVICE=<interface name>" parameter. Hence, changing the CONNECTED_MODE in the ifcfg file, will cause a failure in the interface bring up.	Either disable the Network Manager, or add "DEVICE=<interface name>" to the interface's ifcfg file.
21.	#552840: ifdown command does not function in RH7.x	-
22.	#665143: Kernel Oops may occur after reboot.	-
23.	#555632Kernel panic may occur while re-assigning LIDs.	-
24.	#556352: ICMP traffic might be lost after Vnic restart	-
25.	#560575: Spikes may occur while running PTP protocol over ConnectX-3/ConnectX-3 Pro.	-
26.	#684720: ifdown fails on SLES12SP0/SP1 with the following errors <pre># ifdown ib0 wicked: ifdown: no matching interfaces</pre> The error indicates that there are active interfaces using the interface you are trying to bring down, and you must ifdown all dependent interfaces.	To see the list of all dependent interfaces, run: <pre># wicked --debug all ifdown ib0 wicked: skipping ib0 interface: unable to ifdown due to lowerdev dependency to: ib0.8001 wicked: ifdown: no matching inter- faces wicked: Exit with status: 0</pre>

3.5.2 eIPoIB Known Issues

Table 26 - eIPoIB Known Issues

Index	Internal Reference Number: Description	Workaround
1.	#383034: On rare occasions, upon driver restart the following message is shown in the dmesg: 'cannot create duplicate filename '/class/net/eth_ipoib_interfaces'	-
2.	No indication is received when eIPoIB is non functional.	Run 'ps -ef grep ipoibd' to verify its functionality.
3.	eIPoIB requires libvirtd, python	-
4.	eIPoIB supports only active-backup mode for bonding.	-
5.	eIPoIB supports only VLAN Switch Tagging (VST) mode on guests.	-
6.	IPv6 is currently not supported in eIPoIB	-
7.	#384279: eIPoIB cannot run when Flow Steering is enabled	-
8.	eIPoIB daemon requires the following libs in order to run: python-libxml2, libvirt-bin, libvirt0	-
9.	The eIPoIB driver in ConnectX®-3 and Connect-IB is currently at beta level.	-

3.5.3 XRC Known Issues

Table 27 - XRC Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Legacy API is deprecated, thus when recompiling applications over MLNX_OFED v2.0-3.x.x, warnings such as the below are displayed. rdma.c:1699: warning: 'ibv_open_xrc_domain' is deprecated (declared at /usr/include/infiniband/ofa_verbs.h:72) rdma.c:1706: warning: 'ibv_create_xrc_srq' is deprecated (declared at /usr/include/infiniband/ofa_verbs.h:89) These warnings can be safely ignored.	-
2.	XRC is not functional in heterogeneous clusters containing non Mellanox HCAs.	-
3.	XRC options do not work when using qperf tool.	Use perftest instead
4.	Out-of memory issue might occur due to overload of XRC receive QP with non zero receive queue size created. XRC QPs do not have receive queues.	-

3.5.4 Verbs Known Issues

Table 28 - Verbs Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Using libnl1_1_3~26 or earlier, requires <code>ibv_create_ah</code> protection by a lock for multi-threaded applications.	-
2.	In MLNX_OFED v2.4-1.0.0, if several CQEs are received on a CQ, they will be coalesced and a user-space event will be triggered only once.	When getting an event, poll the CQ until it is empty.
3.	#420847: <code>ibv_task_pingpong</code> over ConnectX-2 adapter cards in not supported.	-

3.5.5 RoCE Known Issues

Table 29 - RoCE Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Not configuring the Ethernet devices or independent VMs with a unique IP address in the physical port, may result in RoCE GID table corruption.	Restart the driver
2.	If <code>RDMA_CM</code> is not used for connection management, then the source and destination GIDs used to modify a QP or create AH should be of the same type - IPv4 or IPv6.	-
3.	#392592: On rare occasions, the driver reports a wrong GID table (read from <code>/sys/class/infiniband/mlx4_*/ports*/gids/*</code>). This may cause communication problems.	-
4.	MLNX_OFED v2.1-1.0.0 and onwards is not interoperable with older versions of MLNX_OFED.	-
5.	Since the number of GIDs per port is limited to 128, there cannot be more than the allowed IP addresses configured to Ethernet devices that are associated with the port. Allowed number is: <ul style="list-style-type: none"> "127" for a single function machine "15" for a hypervisor in a multifunction machine "(127-15)/n" for a guest in a multifunction machine (where n is the number of virtual functions) Note also that each IP address occupies 2 entries when RoCE mode is set to 4 (RoCEv1 + RoCE v2). This further reduces the number of allowed IP addresses.	-
6.	A working IP connectivity between the RoCE devices is required when creating an address handle or modifying a QP with an address vector.	-
7.	IPv4 multicast over RoCE requires the MGID format to be as follow : <code>ffff:<Multicast IPv4 Address></code>	-

Table 29 - RoCE Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
8.	IP RoCEv2 does not support Multicast Listener Discovery (MLD) therefore, multicast traffic over IPv6 may not work as expected.	-
9.	Using GID index 0 (the default GID) is possible only if the matching IPv6 link local address is configured on the net device of the port. This behavior is possible even though the default GID is configured regardless of the presence of the IPv6 address.	-
10.	Using IPv6 link local address (GID0) when VLANs are configured is not supported.	-
11.	Using GID index 0 (the default GID) on port 2 is currently not supported on kernel 3.14 and below.	-
12.	#559276/591244: Dynamically Connected (DC) in RoCE in ConnectX®-4 is currently not supported.	-
13.	Enslaving a Mellanox device to a bond with already configured IPs (or configured upper devices), prevents these IPs from being configured as GIDs.	1. Enslave the Mellanox device. 2. Configure IP devices.
14.	#517825: <code>ibv_create_ah_from_wc</code> is not supported for multicast messages.	-
15.	#592652: InfiniBand error counters found under: <code>/sys/class/infiniband/<mlx5_dev>/ports/<port>/</code> do not function properly in ConnectX-4 adapter cards.	-
16.	#609950/649407: Occasionally, when the Bonding Mode is set to other than active/backup mode (mode 1), the GID table is not populated correctly.	Add slave devices to the master before giving it an IP address.
17.	#667399: In ConnectX-4 adapter cards, when the port speed is lower than 10Gbps, the IB tools will present a higher rate.	-

3.5.6 iSCSI over IPoIB Known Issues

Table 30 - iSCSI over IPoIB Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When working with iSCSI over IPoIB, LRO must be disabled (even if IPoIB is set to connected mode) due to a bug in older kernels which causes a kernel panic.	-

3.6 Storage Protocols Known Issues

3.6.1 Storage Known Issues

Table 31 - Storage Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Older versions of <code>rescan_scsi_bus.sh</code> may not recognize some newly created LUNs.	If encountering such issues, it is recommended to use the '-c' flag.
2.	#664110: SDP is currently not supported in mlx5 driver (Connect-IB and Connect-X 4 adapter cards)	-

3.6.2 SRP Known Issues

Table 32 - SRP Known Issues

Index	Internal Reference Number: Description	Workaround
1.	MLNX_OFED SRP installation breaks the <code>ibmvstgt</code> and <code>ibmvscsi</code> symbol resolution in RHEL7.0	-

3.6.3 SRP Interop Known Issues

Table 33 - SRP Interop Known Issues

Index	Internal Reference Number: Description	Workaround
1.	The driver is tested with Storage target vendors recommendations for <code>multipath.conf</code> extensions (ZFS, DDN, TMS, Nimbus, NetApp).	-

3.6.4 DDN Storage Fusion 10000 Target Known Issues

Table 34 - DDN Storage Fusion 10000 Target Known Issues

Index	Internal Reference Number: Description	Workaround
1.	DDN does not accept non-default <code>P_Key</code> connection establishment.	-

3.6.5 Oracle Sun ZFS Storage 7420 Known Issues

Table 35 - Oracle Sun ZFS Storage 7420 Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Ungraceful power cycle of an initiator connected with Targets DDN, Nimbus, NetApp may result in temporary "stale connection" messages when initiator reconnects.	-

3.6.6 iSER Initiator Known Issues

Table 36 - iSER Initiator Known Issues

Index	Internal Reference Number: Description	Workaround
1.	On SLES OSs, the <code>ib_iser</code> module does not load on boot.	Add a dummy interface using <code>iscsiadm</code> : <ul style="list-style-type: none"> <code># iscsiadm -m iface -I ib_iser -o new</code> <code># iscsiadm -m iface -I ib_iser -o update -n iface.transport_name -v ib_iser</code>
2	Ubuntu12 requires update of user space <code>open-iscsi</code> to v2.0.873	-
3	The initiator does not respect interface parameter while logging in.	Configure each interface on a different subnet.
4	iSCSID v2.0.873 can enter an endless loop on bind error.	-
5	iSCSID may hang if target crashes during logout sequence (reproducible with TCP).	-
6	#440756: SLES12: Logging in with PI disabled followed by a log out and re-log in with PI enabled, without flushing multipath might cause the block layer to panic.	-
7	#489943: Rarely, in InfiniBand devices, when a catastrophic error scenario occurs, iSCSI/iSER initiator might not fully recover and result in system hang.	-
8	#453232: Ubuntu14.04: Stress login/logout might cause block layer to invoke a WARN trace.	-
9	#683370: iSER small read IO (< 8k) performance degrades compared to previous versions. iSER performs memory registration for each IO and avoids sending a global memory key to the target. Sending the global memory key to the wire should only be done in a trusted environment and is not recommended to use over the Internet protocol.	Set module param <code>always_register=N</code> <pre>\$ modprobe ib_iser always_register=N</pre>

3.6.7 iSER Target Known Issues

Table 37 - iSER Target Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Currently only the following OSs are supported: RHEL/CentOS 7.0, SLES12, Ubuntu14.04.	-
2	Stress login/logout from multiple initiators may cause iSER target to panic.	-
3	RHEL/CentOS 7.0: Discovery over RDMA is not supported.	-
4	<code>ib_isert</code> is unavailable on custom kernels after running the <code>mlnx_add_kernel_support.sh</code> script.	<ol style="list-style-type: none"> Add "<code>isert=y</code>" to the <code>mlnx_add_kernel_support.sh</code> script after "<code>cat << EOF > ofed.conf</code>". Use the updated script to build <code>MLNX_OFED</code> for the custom kernel.

Table 37 - iSER Target Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
5	#736410: iSER Target currently supports only the following OSs (distribution kernel): <ul style="list-style-type: none">• RHEL 7.0/7.1/7.2• SLES12/12.1• Ubuntu14.04, Ubuntu15.04	-

3.6.8 ZFS Appliance Known Issues

Table 38 - ZFS Appliance Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Connection establishment occurs twice which may cause iSER to log a stack trace.	-

3.6.9 Erasure Coding Verbs Known Issues

Table 39 - Erasure Coding Verbs Known Issues

Index	Internal Reference Number: Description	Workaround
1.	The Erasure-coding logical block size must be aligned to 64 bytes	-
2	Only w=1,2,3,4 are supported (w corresponds to the Galois symbol size - GF(2 ^w))	-
3	ibv_exp_ec_mem must pass with the following restrictions: <ul style="list-style-type: none"> num_data_sge must be equal to K (property of the EC calc) num_code_sge must be equal to M (property of the EC calc) 	-

3.7 Virtualization

3.7.1 SR-IOV Known Issues

Table 40 - SR-IOV Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When using legacy VMs with MLNX_OFED 2.x hypervisor, you may need to set the 'enable_64b_cqe_eqe' parameter to zero on the hypervisor. It should be set in the same way that other module parameters are set for mlx4_core at module load time. For example, add "options mlx4_core enable_64b_cqe_eqe=0" as a line in the file /etc/modprobe.d/mlx4_core.conf.	-
2.	#381754: mlx4_port1_mtu sysfs entry shows a wrong MTU number in the VM.	-
3.	#426988: When at least one port is configured as InfiniBand, and the num_vfs is provided but the probe_vf is not, HCA initialization fails.	Use both the num_vfs and the probe_vf in the modprobe line.
4.	#385750/378528: When working with a bonding device to enslave the Ethernet devices in active-backup mode and failover MAC policy in a Virtual Machine (VM), establishment of RoCE connections may fail.	Unload the module mlx4_ib and reload it in the VM.

Table 40 - SR-IOV Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
5.	Attaching or detaching a Virtual Function on SLES11 SP3 to a guest Virtual Machine while the <code>mlx4_core</code> driver is loaded in the Virtual Machine may cause a kernel panic in the hypervisor.	Unload the <code>mlx4_core</code> module in the hypervisor before attaching or detaching a function to or from the guest.
6.	#392172: When detaching a VF without shutting down the driver from a VM and reattaching it to another VM with the same IP address for the Mellanox NIC, RoCE connections will fail	Shut down the driver in the VM before detaching the VF.
7.	Enabling SR-IOV requires appending the <code>"intel_iommu=on"</code> option to the relevant OS in file <code>/boot/grub/grub.conf</code> or <code>/boot/grub2/grub.cfg</code> , depending on the OS installed. Without that SR-IOV cannot be loaded.	-
8.	On various combinations of Hypervisor/OSes and Guest/OSes, an issue might occur when attaching/detaching VFs to a guest while that guest is up and running.	Attach/detach VFs to/from a VM only while that VM is down.
9.	The known PCI BDFs for all VFs in kernel command line should be specified by adding <code>xen-pci-back.hide</code> For further information, please refer to http://wiki.xen.org/wiki/Xen_PCI_Passthrough	-
10.	The inbox qemu version (2.0) provided with Ubuntu 14.04 does not work properly when more than 2 VMs are run over an Ubuntu 14.04 Hypervisor.	-
11.	SR-IOV UD QPs are forced by the Hypervisor to use the base GID (i.e., the GID that the VF sees in its GID entry at its paravirtualized index 0). This is needed for security, since UD QPs use Address Vectors, and any GID index may be placed in such a vector, including <code>indices</code> not belonging to that VF.	-
12.	Attempting to attach a PF to a VM when SR-IOV is already enabled on that PF may result in a kernel panic.	-
13.	<code>osmtest</code> on the Hypervisor fails when SR-IOV is enabled. However, only the test fails, OpenSM will operate correctly with the host. The failure reason is that if an <code>mcg</code> is already joined by the host, a subsequent join request for that group succeeds automatically (even if the join parameters in the request are not correct). This success does no harm.	-
14.	If a VM does not support PCI hot plug, detaching an <code>mlx4</code> VF and probing it to the hypervisor may cause the hypervisor to crash.	-
15.	QPerf test is not supported on SR-IOV guests in Connect-IB cards.	-
16.	On ConnectX®-3 HCAs with firmware version 2.32.5000 and later, SR-IOV VPI mode works only with Port 1 = ETH and Port 2 = IB.	-

Table 40 - SR-IOV Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
17.	Occasionally, the <code>lspci grep Mellanox</code> command shows incorrect or partial information due to the current <code>pci.ids</code> file on the machine.	1. Locate the file: <code>\$locate pci.ids</code> 2. Manually update the file according to the latest version available online at: https://pci-ids.ucw.cz/v2.2/pci.ids This file can also be downloaded (using the following command: <code>update-pciids</code>).
18	SR-IOV is not supported in AMD architecture.	-
19	#506512: Setting 1 Mbit/s rate limit on Virtual Functions (Qos Per VF feature) may cause TX queue transmit timeout.	-
20	DC transport type is not supported on SR-IOV VMs.	-
21	#567908: Attaching a VF to a VM before unbinding it from the hypervisor and then attempting to destroy the VM, may cause the system to hang for a few minutes.	-
22	When using SR-IOV make sure to set interface to down and unbind BEFORE unloading driver/removing VF/restarting VM or kernel will lock. (reboot needed) Basically, clean-up might not work perfectly so user should do it manually.	-
23	#568602: Repeating change of the <code>m1x5_num_vfs</code> value from 0 to non-zero, might cause kernel panic in the PF driver.	-
24	#601749: Since the guest MAC addresses are configured to be all zeroes by default, in ConnectX-4 the administrator must explicitly set the VFs' MAC addresses. otherwise the Guest VM will see MAC zero and traffic is not passed.	-
25	#649366: Restarting the PF (Hypervisor) driver while Virtual Functions are assigned is not allowed in due to a <code>vfio-pci</code> bug.	-
26	#639046: Due to an issue with SR-IOV loopback, prevention "Duplicate IPv6 detected" are seen in the VF driver.	-
27	#648680/655070: When a VF uses <code>ethtool</code> facilities, error messages are shown in <code>dmesg</code> .	-
28	#655410: [ConnectX-4/Connect-IB] Failed to enable SR-IOV due to errors in PCI or BIOS.	1. Add <code>pci=realloc=on</code> to the grub command line. 2. Add more memory to the server. 3. Upgrade BIOS version.
29	#651119: Kernel panic may occur while running IPv6 UDP on SR-IOV ConnectX-4 environment	-
30	#669910: Bind/Unbind over ConnectX-4 Hypervisor may cause system lockup.	-
31	#650458: Occasionally, IPv6 might not function properly and cause lockup on SR-IOV ConnectX-4 environment.	-

Table 40 - SR-IOV Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
32	#688551: In ConnectX-3 adapter cards, the extended counter <code>port_rcv_data_64</code> on the VF may not be updated in some flows.	-
33	#690656/690674: When the physical link is down, any traffic from the PF to any VF on the same port will be dropped.	-
34	#691661: When in LAG mode and the Virtual Functions are present (VF LAG), the IP address given to the bonding interface (in the hypervisor) cannot be used for RoCE as well.	Probe one of the VFs in the hypervisor and use for RoCE.
35	#691661: Ethernet SR-IOV in ConnectX-4 requires firmware version 12.14.1100 and higher	-
36	#737434: VF vport statistics are not cleared upon <code>ifconfig up/down</code> .	-

3.8 Resiliency

3.8.1 Reset Flow Known Issues

Table 41 - Resiliency Known Issues

Index	Internal Reference Number: Description	Workaround
1.	SR-IOV non persistent configuration (such as VGT, VST, Host assigned GUIDs, and QP0-enabled VFs) may be lost upon Reset Flow.	Reset Admin configuration post Reset Flow
2.	Upon Reset Flow or after running restart driver, Ethernet VLANs are lost.	Reset the VLANs using the <code>ifup</code> command.
3.	Restarting the driver or running <code>connectx_port_config</code> when Reset Flow is running might result in a kernel panic	-
4.	Networking configuration (e.g. VLANs, IPv6) should be statically defined in order to have them set after Reset Flow as of after restart driver.	-
5.	After recovering from an EEH event, <code>mlx5_core/mlx4_core</code> unload may fail due to a bug in some kernel versions. The bug is fixed in Kernel 3.15	-

3.9 Miscellaneous Known Issues

3.9.1 General Known Issues

Table 42 - General Known Issues

Index	Internal Reference Number: Description	Workaround
1.	On ConnectX-2/ConnectX-3 Ethernet adapter cards, there is a mismatch between the GUID value returned by firmware management tools and that returned by fabric/driver utilities that read the GUID via device firmware (e.g., using <code>ibstat</code>). <code>Mlxburn/flint</code> return <code>0xffff</code> as GUID while the utilities return a value derived from the MAC address. For all driver/firmware/software purposes, the latter value should be used.	N/A. Please use the GUID value returned by the fabric/driver utilities (not <code>0xffff</code>).
2.	#552870/548518: On rare occasions, under extremely heavy MAD traffic, MAD (Management Datagram) storms might cause soft-lockups in the UMAC layer.	-
3.	Packets are dropped on the SM server on big clusters.	Increase the <code>recv_queue_size</code> of <code>ib_mad</code> module parameter for SM server to 8K. The <code>recv_queue_size</code> default size (4K)
4.	#663434: On ConnectX-4/ConnectX-4 Lx, when running " <code>lspci</code> " in RH7.0/7.1, the device information is displayed incorrect or the device is unnamed.	Run <code>update-pciids</code>

3.9.2 ABI Compatibility Known Issues

Table 43 - ABI Compatibility Known Issues

Index	Internal Reference Number: Description	Workaround
1.	MLNX_OFED v2.3-1.0.1 is not ABI compatible with previous MLNX_OFED/OFED versions.	Recompile the application over the new MLNX_OFED version

3.9.3 Connection Manager (CM) Known Issues

Table 44 - Connection Manager (CM) Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When 2 different ports have identical GIDs, the CM might send its packets on the wrong port.	All ports must have different GIDs.

3.9.4 Fork Support Known Issues

Table 45 - Fork Support Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Fork support from kernel 2.6.12 and above is available provided that applications do not use threads. <code>fork()</code> is supported as long as the parent process does not run before the child exits or calls <code>exec()</code> . The former can be achieved by calling <code>wait(childpid)</code> , and the latter can be achieved by application specific means. The Posix <code>system()</code> call is supported.	-

3.9.5 MLNX_OFED Sources Known Issues

Table 46 - MLNX_OFED Sources Known Issues

Index	Internal Reference Number: Description	Workaround
1.	MLNX_OFED includes the OFED source RPM packages used as a build platform for kernel code but does not include the sources of Mellanox proprietary packages.	-

3.9.6 Uplinks Known Issues

Table 47 - Uplinks Known Issues

Index	Internal Reference Number: Description	Workaround
1.	On rare occasions, ConnectX®-3 Pro adapter card may fail to link up when performing parallel detect to 40GbE.	Restart the driver

3.9.7 Resources Limitation Known Issues

Table 48 - Resources Limitation Known Issues

Index	Internal Reference Number: Description	Workaround
1.	The device capabilities reported may not be reached as it depends on the system on which the device is installed and whether the resource is allocated in the kernel or the userspace.	-
2.	#387061: <code>mlx4_core</code> can allocate up to 64 MSI-X vectors, an MSI-X vector per CPU.	-
3.	Setting more IP addresses than the available GID entries in the table results in failure and the "update_gid_table error message is displayed: GID table of port 1 is full. Can't add <address>" message.	-
4.	#553657: Registering a large amount of Memory Regions (MR) may fail because of DMA mapping issues on RHEL 7.0.	-

Table 48 - Resources Limitation Known Issues (Continued)

Index	Internal Reference Number: Description	Workaround
5.	Occasionally, a user process might experience some memory shortage and not function properly due to Linux kernel occupation of the system's free memory for its internal cache.	<p>To free memory to allow it to be allocated in a user process, run the <code>drop_caches</code> procedure below.</p> <p>Performing the following steps will cause the kernel to flush and free pages, dentries and inodes caches from memory, causing that memory to become free.</p> <p>Note: As this is a non-destructive operation and dirty objects are not freeable, run <code>`sync'</code> first.</p> <ul style="list-style-type: none"> • To free the pagecache: <code>echo 1 > /proc/sys/vm/drop_caches</code> • To free dentries and inodes: <code>echo 2 > /proc/sys/vm/drop_caches</code> • To free pagecache, dentries and inodes: <code>echo 3 > /proc/sys/vm/drop_caches</code>

3.9.8 Accelerated Verbs Known Issues

Table 49 - Accelerated Verbs Known Issues

Index	Internal Reference Number: Description	Workaround
1.	<p>On ConnectX®-4 Lx, the following may not be supported when using Multi-Packet WR flag (<code>IBV_EXP_QP_BURST_CREATE_ENABLE_MULTI_PACKET_SEND_WR</code>) on QP-burst family creation:</p> <ul style="list-style-type: none"> • ACLs • SR-IOV (eSwitch offloads) • priority and dscp forcing • Loopback decision. • VLAN insertion • encapsulation (encap/decap) • sniffer • Signature 	

3.10 InfiniBand Fabric Utilities Known Issues

3.10.1 Performance Tools Known Issues

Table 50 - Performance Tools Known Issues

Index	Internal Reference Number: Description	Workaround
1.	perftest package in MLNX_OFED v2.2-1.0.1 and onwards does not work with older versions of the driver.	-

3.10.2 Diagnostic Utilities Known Issues

Table 51 - Diagnostic Utilities Known Issues

Index	Internal Reference Number: Description	Workaround
1.	When running the ibdiagnet check nodes_info on the fabric, a warning specifying that the card does not support general info capabilities for all the HCAs in the fabric will be displayed.	Run <code>ibdiagnet --skip nodes_info</code>
2.	ibdumpp does not work when IPoIB device managed Flow Steering is OFF and at least one of the ports is configured as InfiniBand.	Enable IPoIB Flow Steering and restart the driver. For further information, please refer to MLNX_OFED User Manual section Enable/Disable Flow Steering.
3.	#736136: The maximum number of HCAs shown by ibstat is 32 HCAs.	-

3.10.3 Tools Known Issues

Table 52 - Tools Known Issues

Index	Internal Reference Number: Description	Workaround
1.	Running ibdump in InfiniBand mode with firmware older than v2.33.5000, may cause the server to hang due to a firmware issue.	Run ibdump with firmware v2.33.5000 and higher

4 Bug Fixes History

Table 53 lists the bugs fixed in this release.

Table 53 - Fixed Bugs List

#	Issue	Internal Reference Number: Description	Discovered in Release	Fixed in Release
1.	mlx5 driver	#708299: Fixed kernel's back-ports of XPS and affinity that did not have CONFIG_CPUMASK_OFFSTACK	3.2-1.0.1.1	3.2-2.0.0
2.		#685082: Added support for Rate Limit 0 to enable unlimited rate limiter and to prevent max rate zero traffic lose.	3.2-1.0.1.1	3.2-2.0.0
3.	SR-IOV	#667559: Fixed an issue which enabled SR-IOV on RHEL 6.7 although SR-IOV was already enabled. A check was added to make sure SR-IOV is not enabled before enabling it.	3.2-1.0.1.1	3.2-2.0.0
4.	eIPoIB	#682750: Fixed race between the udev that changes the interface name of eth_ipoib driver and the eIPoIB daemon that configured the same interface.	3.0-1.0.1	3.2-2.0.0
5.	Ethernet traffic	#692520: Fixed an issue which prevented ConnectX-4/ConnectX-4 Lx adapter cards from running Ethernet traffic on Big Endian arch machines.	3.2-1.0.1.1	3.2-2.0.0
6.	Performance	#668346: Set close NUMA node as default for RSS.	3.2-1.0.1.1	3.2-2.0.0
7.	mlx4_en	#696150: Fixed an issue where the ARP request packets destined for a proxy VXLAN interface were not handled correctly when GRO was enabled.	3.2-1.0.1.1	3.2-2.0.0
8.	Counters	#698795: Fixed an issue which prevented the calculated software counters (the correct ones) from being shown and provided the error counters that were previously inactive.	3.0-1.0.1	3.2-2.0.0
9.	Virtualization	#597110: Fixed an issue which prevented the driver from reaching VLAN when the VLAN was created over a Linux bridge.	3.1-1.0.3	3.2-1.0.1.1
10.	mlx5 driver	# 656298: Fixed an issue in the driver (in ConnectX-4) that discarded s-tag VLAN packets when in Promiscuous Mode.	3.1-1.0.3	3.2-1.0.1.1
11.		# 647865: Fixed an issue which prevented PORT_ERR event to be propagated to the user-space application when the port state was changed from Active to Initializing.	3.0-1.0.1	3.2-1.0.1.1
12.	HPC Acceleration packages	# 663975: Fixed a rare issue which allowed the knem package to run depmod on the wrong kernel version.	3.1-1.0.3	3.2-1.0.1.1
13.	IB/Core	# 666992: Fixed a race condition in the IB/umad layer that caused NULL pointer dereference.	3.0-2.0.1	3.2-1.0.1.1
14.	IPoIB	# 657718: Fixed an IPoIB issue that caused connectivity lost after server's restart in a cluster.	3.1-1.0.3	3.2-1.0.1.1

Table 53 - Fixed Bugs List

#	Issue	Internal Reference Number: Description	Discovered in Release	Fixed in Release
15.	Driver un-installation	# 619272: Fixed an issue causing MLNX_OFED to remove the “mutt” package upon driver uninstall.	3.1-1.0.3	3.2-1.0.1.1
16.	PFC	# 613514: Added a warning message in dmesg, notifying the user that the PFC RX/TX cannot be enabled simultaneously with Global Pauses. In this case Global Pauses will be disabled.	3.1-1.0.3	3.2-1.0.1.1
17.	IB MAD	#606916: Fixed an issue causing MADs to drop in large scale clusters.	3.1-1.0.0	3.1-1.0.3
18.	SR-IOV	#367410: Fixed InfiniBand counters which were unavailable in the VM.	2.1-1.0.0	3.1-1.0.0
19.	RoCE	#549687: Fixed InfiniBand traffic counters that are found under <code>/sys/class/infiniband/<mlx-5_dev>/ports/<port>/</code> which do not function properly in ConnectX-4 adapter cards.	3.0-1.0.1	3.1-1.0.0
20.	Virtualization	#589247/591877: Fixed VXLAN functionality issues.	3.0-2.0.1	3.1-1.0.0
21.	Performance	TCP/UDP latency on ConnectX®-4 was higher than expected.	3.0-2.0.1	3.1-1.0.0
22.		TCP throughput on ConnectX®-4 achieved full line rate.	3.0-2.0.1	3.1-1.0.0
23.		#568718: Fixed an issue causing inconsistent performance with ConnectX-3 and PowerKVM 2.1.1.	3.0-2.0.1	3.1-1.0.0
24.		#552658: Fixed ConnectX-4 traffic counters.	3.0-2.0.1	3.1-1.0.0
25.	num_entries	#572068: Updated the desired num_entries in each iteration, and accordingly updated the offset of the WC in the given WC array.	3.0-1.0.1	3.1-1.0.0
26.	mlx5 driver	#536981/554293: Fixed incorrect port rate and port speed values in RoCE mode in ConnectX-4.	3.0-2.0.1	3.1-1.0.0
27.	IPoIB	#551898: In RedHat7.1 kernel 3.10.0-299, when sending ICMP/TCP/UDP traffic over Connect-IB/ ConnectX-4 in UD mode, the packets were dropped with the following error: UDP: bad checksum...	3.0-2.0.1	3.1-1.0.0
28.	openibd	#596458: Fixed an issue which prevented openibd from starting correctly during boot.	3.0-2.0.1	3.1-1.0.0
29.	Ethernet	#589207: Added a new module parameter to control the number of IRQs allocated to the device.	3.0-2.0.1	3.1-1.0.0
30.	mlx5 driver	#576326: Fixed an issue on PPC servers which prevented PCI from reloading after EEH error recovery.	3.0-2.0.1	3.1-1.0.0
31.	OpenSM	#569369: Fixed an issue which prevented the OpenSM package from being fully removed when uninstalling MLNX_OFED v3.0-2.0.1	3.0-2.0.1	3.1-1.0.0

Table 53 - Fixed Bugs List

#	Issue	Internal Reference Number: Description	Discovered in Release	Fixed in Release
32.	mlx5_en	#568169: Added the option to toggle LRO ON/OFF using the "-κ" flags. The priv flag hw_lro will determine the type of LRO to be used, if the flag is ON, the hardware LRO will be used, otherwise the software LRO will be used.	3.0-2.0.1	3.1-1.0.0
33.		#568168: Added the option to toggle LRO ON/OFF using the "-κ" flags.	3.0-2.0.1	3.1-1.0.0
34.		#551075: Fixed race when updating counters.	3.0-2.0.1	3.1-1.0.0
35.		#550275: Fixed scheduling while sending atomic dmesg warning during bonding configuration.	3.0-2.0.1	3.1-1.0.0
36.		#550824: Added set_rx_csum callback implementation.	3.0-2.0.1	3.1-1.0.0
37.	mlx4_ib	#535884: Fixed mismatch between SL and VL in outgoing QP1 packets, which caused buffer overruns in attached switches at high MAD rates.	3.0-1.0.1	3.1-1.0.0
38.	SR-IOV/RoCE	#542722: Fixed a problem on VFs where the RoCE driver registered a zero MAC into the port's MAC table (during QP1 creation) because the ETH driver had not yet generated a non-zero random MAC for the ETH port.t	2.3-1.0.1	3.1-1.0.0
39.		#561866: Removed BUG_ON assert when checking if the ring is full.	3.0-1.0.1	3.1-1.0.0
40.	libvma	#541149: Added libvma support for Debian 8.0 x86_64 and Ubuntu 15.04	3.0-2.0.1	3.1-1.0.0
41.	IPoIB	Fixed an issue which prevented the failure to destroy QP upon IPoIB unload on debug kernel.	3.0-1.0.1	3.0-2.0.0
42.	Configuration	Fixed an issue which prevented the driver version to be reported to the Remote Access Controller tools (such as iDRAC)	3.0-1.0.1	3.0-2.0.0
43.	SR-IOV	Passed the correct port number in port-change event to single-port VFs, where the actual physical port used is port 2.	2.4-1.0.0	3.0-2.0.0
44.		Enabled OpenSM, running over a ConnectX-3 HCA, to manage a mixed ConnectX-3/ConnectX-4 network (by recognizing the "Well-known GID" in mad demux processing).	3.0-1.0.1	3.0-2.0.0
45.		Fixed double-free memory corruption in case where SR-IOV enabling failed (error flow).	3.0-1.0.1	3.0-2.0.0
46.	Start-up sequence	Fixed a crash in EQ's initialization error flow.	3.0-1.0.1	3.0-2.0.0
47.	Installation	#554253: MLNX_OFED v3.0-1.0.1 installation using yum fails on RH7.1	3.0-1.0.1	3.0-2.0.0

Table 53 - Fixed Bugs List

#	Issue	Internal Reference Number: Description	Discovered in Release	Fixed in Release
48.	mlx5 driver	#542686: In PPC systems, when working with ConnectX®-4 adapter card configured as Ethernet, driver load fails with BAD INPUT LENGTH. dmesg: command failed, status bad input length(0x50), syndrome 0x9074aa	3.0-1.0.1	3.0-2.0.0
49.		Error counters such as: CRC error counters, RX out range length error counter, are missing in the ConnectX-4 Ethernet driver.	3.0-1.0.1	3.0-2.0.0
50.		Changing the RX queues number is not supported in Ethernet driver when connected to a ConnectX-4 card.	3.0-1.0.1	3.0-2.0.0
51.	Ethernet	Hardware checksum call trace may appear when receiving IPV6 traffic on PPC systems that uses CHECKSUM COMPLETE method.	3.0-1.0.1	3.0-2.0.0
52.	mlx4_en	Fixed ping/traffic issue occurred when RXVLAN offload was disabled and CHECKSUM COMPLETE was used on ingress packets.	2.4-1.0.4	3.0-1.0.1
53.	Security	CVE-2014-8159 Fix: Prevented integer overflow in IB-core module during memory registration.	2.0-2.0.5	2.4-1.0.4
54.	mlx5_ib	Fixed the return value of max inline received size in the created QP.	2.3-2.0.1	2.4-1.0.0
55.		Resolved soft lock on massive amount of user memory registrations	2.3-2.0.1	2.4-1.0.0
56.	InfiniBand Counters	Occasionally, port_rcv_data and port_xmit_data counters may not function properly.	2.3-1.0.1	2.4-1.0.0
57.	mlx4_en	LRO fixes and improvements for jumbo MTU.	2.3-2.0.1	2.4-1.0.0
58.		Fixed a crash occurred when changing the number of rings (ethtool set-channels) when interface connected to netconsole.	2.2-1.0.1	2.4-1.0.0
59.		Fixed ping issues with IP fragmented datagrams in MTUs 1600-1700.	2.2-1.0.1	2.4-1.0.0
60.		The default priority to TC mapping assigns all priorities to TC0. This configuration achieves fairness in transmission between priorities but may cause undesirable PFC behavior where pause request for priority "n" affects all other priorities.	2.3-1.0.1	2.4-1.0.0
61.	mlx5_ib	Fixed an issue related to large memory regions registration. The problem mainly occurred on PPC systems due to the large page size, and on non PPC systems with large pages (contiguous pages).	2.3-2.0.1	2.3-2.0.5
62.		Fixed an issue in verbs API: fallback to glibc on contiguous memory allocation failure	2.3-2.0.1	2.3-2.0.5
63.	IPoIB	Fixed a memory corruption issue in multi-core system due to intensive IPoIB transmit operation.	2.3-2.0.1	2.3-2.0.5

Table 53 - Fixed Bugs List

#	Issue	Internal Reference Number: Description	Discovered in Release	Fixed in Release
64.	IB MAD	Fixed an issue to prevent process starvation due to MAD packet storm.	2.3-2.0.1	2.3-2.0.5
65.	IPoIB	#433348: Fixed an issue which prevented the spread of events among the closet NUMA CPU when only a single RX queue existed in the system.	2.3-1.0.1	2.3-2.0.0
66.		Returned the CQ to its original state (armed) to prevent traffic from stopping	2.3-1.0.1	2.3-2.0.0
67.		Fixed a TX timeout issue in CM mode, which occurred under heavy stress combined with ifup/ ifdown operation on the IPoIB interface.	2.1-1.0.0	2.3-2.0.0
68.	mlx4_core	Fixed "sleeping while atomic" error occurred when the driver ran many firmware commands simultaneously.	2.3-1.0.1	2.3-2.0.0
69.	mlx4_ib	Fixed an issue related to spreading of completion queues among multiple MSI-X vectors to allow better utilization of multiple cores.	2.1-1.0.0	2.3-2.0.0
70.		Fixed an issue that caused an application to fail when attaching Shared Memory.	2.3-1.0.1	2.3-2.0.0
71.	mlx4_en	Fixed dmesg warnings: "NOHZ: local_soft-irq_pending 08".	2.3-1.0.1	2.3-2.0.0
72.		Fixed erratic report of hardware clock which caused bad report of PTP hardware Time Stamping.	2.1-1.0.0	2.3-2.0.0
73.	mlx5_core	Fixed race when async events arrived during driver load.	2.3-1.0.1	2.3-2.0.0
74.		Fixed race in mlx5_eq_int when events arrived before eq->dev was set.	2.3-1.0.1	2.3-2.0.0
75.		Enabled all pending interrupt handlers completion before freeing EQ memory.	2.3-1.0.1	2.3-2.0.0
76.	mlnx.conf	Defined mlnx.conf as a configuration file in mlnx-ofa_kernel RPM	2.1-1.0.0	2.3-2.0.0
77.	SR-IOV	Fixed counter index allocation for VFs which enables Ethernet port statistics.	2.3-1.0.1	2.3-2.0.0
78.	iSER	Fixed iSER DIX sporadic false DIF errors caused in large transfers when block merges were enabled.	2.3-1.0.1	2.3-2.0.0
79.	RoCE v2	RoCE v2 was non-functional on big Endian machines.	2.3-1.0.1	2.3-2.0.0
80.	Verbs	Fixed registration memory failure when fork was enabled and contiguous pages or ODP were used.	2.3-1.0.1	2.3-2.0.0
81.	Installation	Using both '-c --config' and '--add-kernel-support' flags simultaneously when running the mlnxofedinstall.sh script caused installation failure with the following on screen message "--config does not exist".	2.2-1.0.1	2.3-2.0.0

Table 53 - Fixed Bugs List

#	Issue	Internal Reference Number: Description	Discovered in Release	Fixed in Release
82.	IPoIB	Changing the GUID of a specific SR-IOV guest after the driver has been started, causes the ping to fail. Hence, no traffic can go over that InfiniBand interface.	2.1-1.0.0	2.3-1.0.1
83.	XRC	XRC over ROCE in SR-IOV mode is not functional	2.0-3.1.0	2.2-1.0.1
84.	mlx4_en	Fixed wrong calculation of packet true-size reporting in LRO flow.	2.1-1.0.0	2.2-1.0.1
85.		Fixed kernel panic on Debian-6.0.7 which occurred when the number of TX channels was set above the default value.	2.1-1.0.0	2.2-1.0.1
86.		Fixed a crash incidence which occurred when enabling Ethernet Time-stamping and running VLAN traffic.	2.0-2.0.5	2.2-1.0.1
87.	IB Core	Fixed the QP attribute mask upon smac resolving	2.1-1.0.0	2.1-1.0.6
88.	mlx5_ib	Fixed a send WQE overhead issue	2.1-1.0.0	2.1-1.0.6
89.		Fixed a NULL pointer de-reference on the debug print	2.1-1.0.0	2.1-1.0.6
90.		Fixed arguments to kzalloc	2.1-1.0.0	2.1-1.0.6
91.	mlx4_core	Fixed the locks around completion handler	2.1-1.0.0	2.1-1.0.6
92.	mlx4_core	Restored port types as they were when recovering from an internal error.	2.0-2.0.5	2.1-1.0.0
93.		Added an N/A port type to support port_type_array module param in an HCA with a single port	2.0-2.0.5	2.1-1.0.0
94.	SR-IOV	Fixed memory leak in SR-IOV flow.	2.0-2.0.5	2.0-3.0.0
95.		Fixed communication channel being stuck	2.0-2.0.5	2.0-3.0.0
96.	mlx4_en	Fixed ALB bonding mode failure when enslaving Mellanox interfaces	2.0-3.0.0	2.1-1.0.0
97.		Fixed leak of mapped memory	2.0-3.0.0	2.1-1.0.0
98.		Fixed TX timeout in Ethernet driver.	2.0-2.0.5	2.0-3.0.0
99.		Fixed ethtool stats report for Virtual Functions.	2.0-2.0.5	2.0-3.0.0
100.		Fixed an issue of VLAN traffic over Virtual Machine in paravirtualized mode.	2.0-2.0.5	2.0-3.0.0
101.		Fixed ethtool operation crash while interface down.	2.0-2.0.5	2.0-3.0.0
102.	IPoIB	Fixed memory leak in Connected mode.	2.0-2.0.5	2.0-3.0.0
103.		Fixed an issue causing IPoIB to avoid pkey value 0 for child interfaces.	2.0-2.0.5	2.0-3.0.0

5 Change Log History

Table 54 - Change Log History

Release	Category	Description
3.2-1.0.1.1	VXLAN Hardware Stateless Offloads	[ConnectX-4 / ConnectX-4 Lx] Provides scalability and security challenges solutions.
	Priority Flow Control (PFC)	[ConnectX-4 / ConnectX-4 Lx] Applies pause functionality to specific classes of traffic on the Ethernet link.
	Offloaded Traffic Sniffer/TCP Dump	[ConnectX-4 / ConnectX-4 Lx] Allows bypass kernel traffic (such as, RoCE, VMA, DPDK) to be captured by existing packet analyzer such as tcpdump.
	Ethernet Time Stamping	[ConnectX-4 / ConnectX-4 Lx] Keeps track of the creation of a packet. A time-stamping service supports assertions of proof that a datum existed before a particular time.
	Custom RoCE Counters	[ConnectX-4 / ConnectX-4 Lx] Provide a clear indication on RDMA send/receive statistics and errors.
	LED Beaconing	[ConnectX-4 / ConnectX-4 Lx] Enables visual identification of the port by LED blinking.
	Enhanced Transmission Selection standard (ETS)	[ConnectX-4 / ConnectX-4 Lx] Exploits the time periods in which the offered load of a particular Traffic Class (TC) is less than its minimum allocated bandwidth.
	Strided WQE User Space	[ConnectX-4 / ConnectX-4 Lx] Striding RQ is a receive queue comprised by work queue elements (i.e. WQEs), where multiple packets of LRO segments (i.e. message) are written to the same WQE.
	VLAN Stripping in Linux Verbs	[ConnectX-4 / ConnectX-4 Lx] Adds access to the device's ability to offload the Customer VLAN (cVLAN) header stripping from an incoming packet.
	iSER: Remote invalidation support (target and initiator)	[ConnectX-4 / ConnectX-4 Lx] Improves performance by enabling the hardware to perform implicit memory region invalidation.
	iSER: Zero-Copy ImmediateData	[ConnectX-4 / ConnectX-4 Lx] Reduces the latency of small writes by avoiding an extra memory copy in the iSER target stack.
	iSER: Indirect Memory Registration	[ConnectX-4 / ConnectX-4 Lx] Uses ConnectX®-4 adapter card's Indirect Memory Registration capabilities to avoid bounce buffer strategy implementation and to reduce the latency of highly unaligned vectored IO operations, and also in cases of BIO merging.
Vector Calculation/ Erasure coding offload	[ConnectX-4 / ConnectX-4 Lx] Uses the HCA for offloading erasure coding calculations.	

Table 54 - Change Log History

Release	Category	Description
3.2-1.0.1.1 (cont.)	Virtual Guest Tagging (VGT+)	[ConnectX-3 / ConnectX-3 Pro] VGT+ is an advanced mode of Virtual Guest Tagging (VGT), in which a VF is allowed to tag its own packets as in VGT, but is still subject to an administrative VLAN trunk policy.
	Link Aggregation for Virtual Functions	[ConnectX-3 / ConnectX-3 Pro] Protects a VM with an attached ConnectX-3 VF from VF port failure, when VFs are present and RoCE Link Aggregation is configured in the Hypervisor.
3.1-1.0.3	User Access Region (UAR)	Allows the ConnectX-3 driver to operate on PPC machines without requiring a change to the MMIO area size.
	CQE Compression	Saves PCIe bandwidth by compressing a few CQEs into a smaller amount of bytes on PCIe
	Bug fixes	See Section 4, “Bug Fixes History”, on page 40
3.1-1.0.0	Wake-on-LAN (WOL)	Wake-on-LAN (WOL) is a technology that allows a network professional to remotely power on a computer or to wake it up from sleep mode.
	Hardware Accelerated 802.1ad VLAN (Q-in-Q Tunneling)	Q-in-Q tunneling allows the user to create a Layer 2 Ethernet connection between two servers. The user can segregate a different VLAN traffic on a link or bundle different VLANs into a single VLAN.
	ConnectX-4 ECN	ECN in ConnectX-4 enables end-to-end congestions notifications between two end-points when a congestion occurs, and works over Layer 3.
	RSS Verbs Support for ConnectX-4 HCAs	Receive Side Scaling (RSS) technology allows spreading incoming traffic between different receive descriptor queues. Assigning each queue to different CPU cores allows better load balancing of the incoming traffic and improve performance.
	Minimal Bandwidth Guarantee (ETS)	The amount of bandwidth (BW) left on the wire may be split among other TCs according to a minimal guarantee policy.
	SR-IOV Ethernet	SR-IOV Ethernet at Beta level
3.0-2.0.1	Virtualization	Added support for SR-IOV for ConnectX-4/Connect-IB adapter cards.

Table 54 - Change Log History

Release	Category	Description
3.0-1.0.1	HCA's	Added support for ConnectX®-4 Single/Dual-Port Adapter supporting up to 100Gb/s.
	RoCE per GID	RoCE per GID provides the ability to use different RoCE versions/modes simultaneously.
	RoCE Link Aggregation (RoCE LAG): ConnectX-3/ConnectX-3 Pro only	RoCE Link Aggregation (available in kernel 4.0 only) provides failover and link aggregation capabilities for mlx4 device physical ports. In this mode, only one IB port that represents the two physical ports, is exposed to the application layer.
	Resource Domain Experimental Verbs	Resource domain is a verb object which may be associated with QP and/or CQ objects on creation to enhance data-path performance.
	Alias GUID Support in InfiniBand	Enables the <code>query_gid</code> verb to return the admin desired value instead of the value that was approved by the SM, to prevent a case where the SM is unreachable or a response is delayed, or if the VF is probed into a VM before their GUID is registered with the SM.

Table 54 - Change Log History

Release	Category	Description
3.0-1.0.1 (cont.)	Denial Of Service (DOS) MAD Prevention	Denial Of Service MAD prevention is achieved by assigning a threshold for each agent's RX. Agent's RX threshold provides a protection mechanism to the host memory by limiting the agents' RX with a threshold.
	QoS per VF (Rate Limit per VF)	Virtualized QoS per VF, (supported in ConnectX-3/ConnectX-3 Pro adapter cards only with firmware v2.33.5100 and above), limits the chosen VFs' throughput rate limitations (Maximum throughput). The granularity of the rate limitation is 1Mbits.
	Ignore Frame Check Sequence (FCS) Errors	Upon receiving packets, the packets go through a checksum validation process for the FCS field. If the validation fails, the received packets are dropped. Using this feature, enables you to choose whether or not to drop the frames in case the FCS is wrong and use the FCS field for other info.
	Sockets Direct Protocol (SDP)	Sockets Direct Protocol (SDP) is a byte-stream transport protocol that provides TCP stream semantics. and utilizes InfiniBand's advanced protocol offload capabilities.
	Scalable Subnet Administration (SSA)	The Scalable Subnet Administration (SSA) solves Subnet Administrator (SA) scalability problems for Infiniband clusters. It distributes the needed data to perform the path-record-calculation needed for a node to connect to another node, and caches these locally in the compute (client) nodes. SSA ^a requires AF_IB address family support (3.12.28-4 kernel and later).
	SR-IOV in ConnectX-3 cards	Changed the Alias GUID support behavior in InfiniBand.
	LLR max retransmission rate	Added LLR max retransmission rate as specified in Vendor Specific MAD V1.1, Table 110 - PortLLRStatistics MAD Description ibdiagnet presents the LLR max_retransmission_rate counter as part of the PM_INFO in db_csv file.
	Experimental Verbs	Added the following verbs: <ul style="list-style-type: none"> • <code>ibv_exp_create_res_domain</code> • <code>ibv_exp_destroy_res_domain</code> • <code>ibv_exp_query_intf</code> • <code>ibv_exp_release_intf</code> Added the following interface families: <ul style="list-style-type: none"> • <code>ibv_exp_qp_burst_family</code> • <code>ibv_exp_cq_family</code>
2.4-1.0.4	Bug Fixes	See "Bug Fixes History" on page 40.

Table 54 - Change Log History

Release	Category	Description
2.4-1.0.0	mlx4_en net-device Ethtool	Added support for Ethtool speed control and advertised link mode.
		Added ethtool txvlan control for setting ON/OFF hardware TX VLAN insertion: <code>ethtool -k txvlan [on/off]</code>
		Ethtool report on port parameters improvements.
		Ethernet TX packet rate improvements.
	RoCE	RoCE uses now all available EQs and not only the 3 legacy EQs.
	InfiniBand	IRQ affinity hints are now set when working in InfiniBand mode.
	Virtualization	VXLAN fixes and performance improvements.
	libmlx4 & libmlx5	Improved message rate of short messages.
	libmlx5	Added ConnectX®-4 device (4114) to the list of supported devices (<code>hca_table</code>),
	Storage	Added iSER Target driver.
	Ethernet net-device	New adaptive interrupt moderation scheme to improve CPU utilization.
RSS support of fragmented IP datagram.		
Connect-IB Virtual Function	Added Connect-IB Virtual Function to the list of supported devices.	
2.3-2.0.5	mlx5_core	<p>Added the following files under <code>/sys/class/infiniband/mlx5_0/mr_cache/</code>:</p> <ul style="list-style-type: none"> <code>rel_timeout</code>: Defines the minimum allowed time between the last MR creation to the first MR released from the cache. When <code>rel_timeout = -1</code>, MRs are not released from the cache <code>rel_imm</code>: Triggers the immediate release of excess MRs from the cache when set to 1. When all excess MRs are released from the cache, <code>rel_imm</code> is reset back to 0.
	Bug Fixes	See “Bug Fixes History” on page 40.
2.3-2.0.1	Bug Fixes	See “Bug Fixes History” on page 40.
2.3-2.0.0	Connect-IB	Added Suspend to RAM (S3).
	Reset Flow	Added Enhanced Error Handling for PCI (EEH), a recovery strategy for I/O errors that occur on the PCI bus.
	Register Contiguous Pages	Added the option to ask for a specific address when the register memory is using contiguous page.
	mlx5_core	Moved the <code>mr_cache</code> subtree from <code>debugfs</code> to <code>mlx5_ib</code> while preserving all its semantics.
	InfiniBand Utilities	Updated the <code>ibutils</code> package. Added to the <code>ibdiagnet</code> tool the " <code>ibdiagnet2.mlx_cntrs</code> " option to enable reading of Mellanox diagnostic counters.
	Bug Fixes	See “Bug Fixes History” on page 40.

Table 54 - Change Log History

Release	Category	Description
2.3-1.0.1	OpenSM	Added Routing Chains support with Minhop/UPDN/FTree/DOR/Torus-2QoS
		Added double failover elimination. When the Master SM is turned down for some reason, the Standby SM takes ownership over the fabric and remains the Master SM even when the old Master SM is brought up, to avoid any unnecessary re-registrations in the fabric. To enable this feature, set the "master_sm_priority" parameter to be greater than the "sm_priority" parameter in all SMs in the fabric. Once the Standby SM becomes the Master SM, its priority becomes equal to the "master_sm_priority". So that additional SM handover is avoided. Default value of the master_sm_priority is 14. To disable this feature, set the "master_sm_priority" in opensm.conf to 0.
		Added credit-loop free unicast/multicast updn/ftree routing
		Added multithreaded Minhop/UPDN/DOR routing
	RoCE	Added IP routable RoCE modes. For further information, please refer to the MLNX_OFED User Manual.
	Installation	Added apt-get installation support.
	Ethernet	Added support for arbitrary UDP port for VXLAN. From upstream 3.15-rc1 and onward, it is possible to use arbitrary UDP port for VXLAN. This feature requires firmware version 2.32.5100 or higher. Additionally, the following kernel configuration option CONFIG_MLX4_EN_VXLAN=y must be enabled.
		MLNX_OFED no longer changes the OS sysctl TCP parameters.
		Added Explicit Congestion Notification (ECN) support
		Added Flow Steering: A0 simplified steering support
		Added RoCE v2 support

Table 54 - Change Log History

Release	Category	Description
2.3-1.0.1 (cont.)	InfiniBand Network	Added Secure host to enable the device to protect itself and the subnet from malicious software.
		Added User-Mode Memory Registration (UMR) to enable the usage of RDMA operations and to scatter the data at the remote side through the definition of appropriate memory keys on the remote side.
		Added On-Demand-Paging (ODP), a technique to alleviate much of the shortcomings of memory registration.
		Added Masked Atomics operation support
		Added Checksum offload for packets without L4 header support
		Added Memory re-registration to allow the user to change attributes of the memory region.
	Resiliency	Added Reset Flow for ConnectX@-3 (+SR-IOV) support.
	SR-IOV	Added Virtual Guest Tagging (VGT+), an advanced mode of Virtual Guest Tagging (VGT), in which a VF is allowed to tag its own packets as in VGT, but is still subject to an administrative VLAN trunk policy.
	Ethtool	Added Cable EEPROM reporting support
		Disable/Enable ethernet RX VLAN tag striping offload via ethtool
128 Byte Completion Queue Entry (CQE)		
Non-Linux Virtual Machines	Added Windows Virtual Machine over Linux KVM Hypervisor (SR-IOV with InfiniBand only) support	
Rev 2.2-1.0.1	mlnxofedinstall	32-bit libraries are no longer installed by default on 64-bit OS. To install 32-bit libraries use the ' <code>--with-32bit</code> ' installation parameter.
	openibd	Added pre/post start/stop scripts support. For further information, please refer to section " <i>openibd Script</i> " in the MLNX_OFED User Manual.
	Reset Flow	Reset Flow is not activated by default. It is controlled by the <code>mlx-4_core'internal_err_reset'</code> module parameter.

Table 54 - Change Log History

Release	Category	Description
Rev 2.2-1.0.1	InfiniBand Core	Asymmetric MSI-X vectors allocation for the SR-IOV hypervisor and guest instead of allocating 4 default MSI-X vectors. The maximum number of MSI-X vectors is <code>num_cpu</code> for port ConnectX®-3 has 1024 MSI-X vectors, 28 MSI-X vectors are reserved. <ul style="list-style-type: none"> Physical Function - gets the number of MSI-X vectors according to the <code>pf_msix_table_size</code> (multiple of 4 - 1) INI parameter Virtual Functions – the remaining MSI-X vectors are spread equally between all VFs, according to the <code>num_vfs mlx-4_core</code> module parameter
	Ethernet	Ethernet VXLAN support for kernels 3.12.10 or higher
		Power Management Quality of Service: when the traffic is active, the Power Management QoS is enabled by disabling the CPU states for maximum performance.
		Ethernet PTP Hardware Clock support on kernels/OSes that support it
	Verbs	Added additional experimental verbs interface. This interface exposes new features which are not integrated yet in to the upstream libibverbs. The Experimental API is an extended API therefore, it is backward compatible, meaning old application are not required to be recompiled to use MLNX-OFED v2.2-1.0.1.
Performance	Out of the box performance improvements: <ul style="list-style-type: none"> Use of affinity hints (based on NUMA node of the device) to indicate the IRQ balancer daemon on the optimal IRQ affinity Improvement in buffers allocation schema (based on the hint above) Improvement in the adaptive interrupt moderation algorithm 	
Rev 2.1-1.0.6	IB Core	Added allocation success verification process to <code>ib_alloc_device</code> .
	dapl	dapl is recompiled with no FCA support.
	openibd	Added the ability to bring up child interfaces even if the parent's <code>ifcfg</code> file is not configured.
	libmlx4	Unmapped the <code>hca_clock_page</code> parameter from <code>mlx4_uninit_context</code> .
	scsi_transport_srp	<code>scsi_transport_srp</code> cannot be cleared up when <code>rport</code> reconnecting fails.
	mlnxofedinstall	Added support for the following parameters: <ul style="list-style-type: none"> '--umad-dev-na' '--without-<package>'
Rev 2.1-1.0.6	Content Packages Updates	The following packages were updated: <ul style="list-style-type: none"> bupc to v2.2-407 mstflint to v3.5.0-1.1.g76e4acf perftest to v2.0-0.76.gb9a463 hcoll to v2.0.472-1 Openmpi to v1.6.5-440ad47 dapl to v2.0.40

Table 54 - Change Log History

Release	Category	Description
Rev 2.1-1.0.0	EoIB	EoIB is supported only in SLES11SP2 and RHEL6.4.
	eIPoIB	eIPoIB is currently at GA level.
	Connect-IB®	Added the ability to resize CQs.
	IPoIB	Reusing DMA mapped SKB buffers: Performance improvements when IOMMU is enabled.
	mlx_en	Added reporting autonegotiation support.
		Added Transmit Packet Steering (XPS) support.
		Added reporting 56Gbit/s link speed support.
		Added Low Latency Socket (LLS) support.
		Added check for dma_mapping errors.
eIPoIB	Added non-virtual environment support.	
Rev 2.0-3.0.0	Operating Systems	Additional OS support: <ul style="list-style-type: none"> • SLES11SP3 • Fedora16, Fedora17
	Drivers	Added Connect-IB™ support
	Installation	Added ability to install MLNX_OFED with SR-IOV support.
		Added Yum installation support
	EoIB	EoIB (at beta level) is supported only in SLES11SP2 and RHEL6.4
	mlx4_core	Modified module parameters to associate configuration values with specific PCI devices identified by their bus/device/function value format
	mlx4_en	Reusing DMA mapped buffers: major performance improvements when IOMMU is enabled
		Added Port level QoS support
	IPoIB	Reduced memory consumption
Limited the number TX and RX queues to 16		
Default IPoIB mode is set to work in Datagram, except for Connect-IB™ adapter card which uses IPoIB with Connected mode as default.		
Rev 2.0-3.0.0	Storage	iSER (at GA level)

Table 54 - Change Log History

Release	Category	Description
Rev 2.0-2.0.5 ^b	Virtualization	SR-IOV for both Ethernet and InfiniBand (at Beta level)
	Ethernet Network	RoCE over SR-IOV (at Beta level)
		eIPoIB to enable IPoIB in a Para-Virtualized environment (at Alpha level)
		Ethernet Performance Enhancements (NUMA related and others) for 10G and 40G
		Ethernet Time Stamping (at Beta level)
		Flow Steering for Ethernet and InfiniBand. (at Beta level)
		Raw Eth QPs: <ul style="list-style-type: none"> • Checksum TX/RX • Flow Steering
		InfiniBand Network
	Installation	YUM update support
	VMA	OFED_VMA integration to a single branch
	Storage	iSER (at Beta level) and SRP
	Operating Systems	Errata Kernel upgrade support
	API	VERSION query API: library and headers
Counters	64bit wide counters (port xmit/recv data/packets unicast/mcast)	

- a. SSA is tested on SLES 12 only (x86-64 architecture).
b. SR-IOV, Ethernet Time Stamping and Flow Steering are ConnectX®-3 HCA capability.

6 API Change Log History

Table 55 - API Change Log History

Release	Name	Description
3.2-1.0.1.1	libibverbs	<ul style="list-style-type: none"> Added API and primitives for Erasure Coding calculations. <ul style="list-style-type: none"> Verbs: <ul style="list-style-type: none"> ibv_exp_alloc_ec_calc ibv_exp_dealloc_ec_calc ibv_exp_ec_encode_sync ibv_exp_ec_decode_sync ibv_exp_ec_encode_async ibv_exp_ec_decode_async ibv_exp_ec_encode_send For further information, please refer to the manual page of the verbs. <ul style="list-style-type: none"> Structs: <ul style="list-style-type: none"> ibv_exp_ec_calc ibv_exp_ec_mem ibv_exp_ec_stripe ibv_exp_ec_comp Added version 1 of the CQ family with support for: <ul style="list-style-type: none"> Multi-Packet RQ (also called striding RQ) Cvlan stripping offload Added enhanced masked-atomic device capability Added a flag to the create QP/WQ option to enable end of RX message padding
Rev 3.1-1.0.3	libibverbs	<ul style="list-style-type: none"> Added <code>ibv_exp_wq_family</code> interface family (Supported only by ConnectX®-4) Added flag to the QP-burst family to enable Multi-Packet WR Added return error statuses to the <code>ibv_exp_query_intf</code> to notify that common-flags/family-flags are not supported. Added <code>ibv_exp_query_gid_attr</code> verb. For further information, please refer to the manual page of the verb.
Rev 3.0-1.0.0	libibverbs	<ul style="list-style-type: none"> Added the following new APIs: <ul style="list-style-type: none"> <code>ibv_exp_create_res_domain</code> - create resource domain <code>ibv_exp_destroy_res_domain</code> - destroy resource domain <code>ibv_exp_query_intf</code> - query for family of verbs interface for specific QP/CQ <code>ibv_exp_release_intf</code> - release the queried interface Updated the following APIs: <ul style="list-style-type: none"> <code>ibv_exp_create_qp</code> - Add resource-domain to the verb parameters <code>ibv_exp_create_cq</code> - Add resource-domain to the verb parameters

Table 55 - API Change Log History

Release	Name	Description
Rev 2.4-1.0.0	libibverbs	<p>Added the following verbs interfaces:</p> <ul style="list-style-type: none"> • <code>ibv_create_flow</code> • <code>ibv_destroy_flow</code> • <code>ibv_exp_use_priv_env</code> • <code>ibv_exp_setenv</code>
Rev 2.3-1.0.1	libibverbs	<ul style="list-style-type: none"> • <code>ibv_exp_rereg_mr</code> - Added new API for memory region re-integration (For further information, please refer to MLNX_OFED User Manual) • Added to the experimental API <code>ibv_exp_post_send</code> the following opcodes: <ul style="list-style-type: none"> • <code>IBV_EXP_WR_EXT_MASKED_ATOMIC_CMP_AND_SWP</code> • <code>IBV_EXP_WR_EXT_MASKED_ATOMIC_FETCH_AND_ADD</code> • <code>IBV_EXP_WR_NOP</code> <p>and these completion opcodes:</p> <ul style="list-style-type: none"> • <code>IBV_EXP_WC_MASKED_COMP_SWAP</code> • <code>IBV_EXP_WC_MASKED_FETCH_ADD</code>
Rev 2.2-1.0.1	libibverbs	<p>The following verbs changed to align with upstream libibverbs:</p> <ul style="list-style-type: none"> • <code>ibv_reg_mr</code> - <code>ibv_access_flags</code> changed. • <code>ibv_post_send</code> - opcodes and send flags changed and <code>wr</code> fields removed (<code>task</code>, <code>op</code>, <code>dc</code> and <code>bind_mw</code>) • <code>ibv_query_device</code> - capability flags changed. • <code>ibv_poll_cq</code> - opcodes and <code>wc</code> flags changed. • <code>ibv_modify_qp</code> - mask bits changed • <code>ibv_create_qp_ex</code> - <code>create_flags</code> field removed. <p>The following verbs removed to align with upstream libibverbs:</p> <ul style="list-style-type: none"> • <code>ibv_bind_mw</code> • <code>ibv_post_task</code> • <code>ibv_query_values_ex</code> • <code>ibv_query_device_ex</code> • <code>ibv_poll_cq_ex</code> • <code>ibv_reg_shared_mr_ex</code> • <code>ibv_reg_shared_mr</code> • <code>ibv_modify_cq</code> • <code>ibv_create_cq_ex</code> • <code>ibv_modify_qp_ex</code>

Table 55 - API Change Log History

Release	Name	Description
Rev 2.2-1.0.1	Verbs Experimental API	<p>The following experimental verbs added (replacing the removed extended verbs):</p> <ul style="list-style-type: none"> • <code>ibv_exp_bind_mw</code> • <code>ibv_exp_post_task</code> • <code>ibv_exp_query_values</code> • <code>ibv_exp_query_device</code> • <code>ibv_exp_poll_cq</code> • <code>ibv_exp_reg_shared_mr</code> • <code>ibv_exp_modify_cq</code> • <code>ibv_exp_create_cq</code> • <code>ibv_exp_modify_qp</code> <p>New experimental verbs:</p> <ul style="list-style-type: none"> • <code>ibv_exp_arm_dct</code> • <code>ibv_exp_query_port</code> • <code>ibv_exp_create_flow</code> • <code>ibv_exp_destroy_flow</code> • <code>ibv_exp_post_send</code> • <code>ibv_exp_reg_mr</code> • <code>ibv_exp_get_provider_func</code>
Rev 2.1-1.0.0	Dynamically Connected (DC)	<p>The following verbs were added:</p> <ul style="list-style-type: none"> • <code>struct ibv_dct *ibv_exp_create_dct(struct ibv_context *context, struct ibv_exp_dct_init_attr *attr)</code> • <code>int ibv_exp_destroy_dct(struct ibv_dct *dct)</code> • <code>int ibv_exp_query_dct(struct ibv_dct *dct, struct ibv_exp_dct_attr *attr)</code>
	<p>Verbs Extension API: Verbs extension API defines OFA APIs extension scheme to detect ABI compatibility and enable backward and forward compatibility support.</p>	<ul style="list-style-type: none"> • <code>ibv_post_task</code> • <code>ibv_query_values_ex</code> • <code>ibv_query_device_ex</code> • <code>ibv_create_flow</code> • <code>ibv_destroy_flow</code> • <code>ibv_poll_cq_ex</code> • <code>ibv_reg_shared_mr_ex</code> • <code>ibv_open_xrcd</code> • <code>ibv_close_xrcd</code> • <code>ibv_modify_cq</code> • <code>ibv_create_srq_ex</code> • <code>ibv_get_srq_num</code> • <code>ibv_create_qp_ex</code> • <code>ibv_create_cq_ex</code> • <code>ibv_open_qp</code> • <code>ibv_modify_qp_ex</code>

Table 55 - API Change Log History

Release	Name	Description
Rev 2.1-1.0.0	Verbs Experimental API: Verbs experimental API defines MLNX-OFED APIs extension scheme which is similar to the “Verbs extension API”. This extension provides a way to introduce new features before they are integrated into the formal OFA API and to the upstream kernel and libs.	<ul style="list-style-type: none"> • <code>ibv_exp_create_qp</code> • <code>ibv_exp_query_device</code> • <code>ibv_exp_create_dct</code> • <code>ibv_exp_destroy_dct</code> • <code>ibv_exp_query_dct</code>
Rev 2.0-3.0.0	XRC	<p>The following verbs have become deprecated:</p> <ul style="list-style-type: none"> • <code>struct ibv_xrc_domain *ibv_open_xrc_domain</code> • <code>struct ibv_srq *ibv_create_xrc_srq</code> • <code>int ibv_close_xrc_domain</code> • <code>int ibv_create_xrc_rcv_qp</code> • <code>int ibv_modify_xrc_rcv_qp</code> • <code>int ibv_query_xrc_rcv_qp</code> • <code>int ibv_reg_xrc_rcv_qp</code> • <code>int ibv_unreg_xrc_rcv_qp</code>
Rev 2.0-2.0.5	Libibverbs - Extended speeds	<ul style="list-style-type: none"> • Missing the <code>ext_active_speed</code> attribute from the <code>struct ibv_port_attr</code> • Removed function <code>ibv_ext_rate_to_int</code> • Added functions <code>ibv_rate_to_mbps</code> and <code>mbps_to_ibv_rate</code>
	Libibverbs - Raw QPs	QP types <code>IBV_QPT_RAW_PACKET</code> and <code>IBV_QPT_RAW_ETH</code> are not supported
	Libibverbs - Contiguous pages	<ul style="list-style-type: none"> • Added Contiguous pages support • Added function <code>ibv_reg_shared_mr</code>
	Libmverbs	<ul style="list-style-type: none"> • The enumeration <code>IBV_M_WR_CALC</code> was renamed to <code>IBV_M_WR_CALC_SEND</code> • The enumeration <code>IBV_M_WR_WRITE_WITH_IMM</code> was added • In the structure <code>ibv_m_send_wr</code>, the union <code>wr.send</code> was renamed to <code>wr.calc_send</code> and <code>wr.rdma</code> was added • The enumerations <code>IBV_M_WQE_CAP_CALC_RDMA_WRITE_WITH_IMM</code> was added • The following enumerations were renamed: <ul style="list-style-type: none"> • From <code>IBV_M_WQE_SQ_ENABLE_CAP</code> to <code>IBV_M_WQE_CAP_SQ_ENABLE</code> • From <code>IBV_M_WQE_RQ_ENABLE_CAP</code> to <code>IBV_M_WQE_CAP_RQ_ENABLE</code> • From <code>IBV_M_WQE_CQE_WAIT_CAP</code> to <code>IBV_M_WQE_CAP_CQE_WAIT</code> • From <code>IBV_M_WQE_CALC_CAP</code> to <code>IBV_M_WQE_CAP_CALC_SEND</code>