



# **NVIDIA MLNX\_OFED Documentation v4.9-7.1.0.0 LTS**

# Table of contents

<b>Release Notes</b>	7
General Support in MLNX_OFED	11
Changes and New Features	21
Bug Fixes	22
Known Issues	54
<b>User Manual</b>	92
Introduction	92
Installation	108
Features Overview and Configuration	108
Programming	109
InfiniBand Fabric Utilities	112
Troubleshooting	124
Common Abbreviations and Related Documents	124
<b>Documentation History</b>	128
Release Notes History	
User Manual Revision History	

# List of Figures

Figure 0. Image2019 1 10 10 55 56

---

Figure 1. Procedure Heading Icon

---

Figure 2. Worddave635fed9c99097774044df72a47e9130

---

Figure 3. Image2019 2 12 10 26 53

---

Figure 4. Procedure Heading Icon

---

Figure 5. Procedure Heading Icon

---

Figure 6. Procedure Heading Icon

---

Figure 7. Procedure Heading Icon

---

Figure 8. Procedure Heading Icon

---

Figure 9. Image2019 2 20 17 49 3

---

Figure 10. Procedure Heading Icon

---

Figure 11. Procedure Heading Icon

---

Figure 12. Procedure Heading Icon

---

Figure 13. Procedure Heading Icon

---

Figure 14. Procedure Heading Icon

---

Figure 15. Procedure Heading Icon

---

Figure 16. Image2022 7 26 10 29 13

---

Figure 17. Image2022 7 28 14 58 53

---

Figure 18. Image2022 7 28 15 0 34

---

Figure 19. Image2022 7 28 14 59 57

---

Figure 20. Image2022 7 28 15 1 11

---

Figure 21. Image2022 7 28 15 1 42

---

Figure 22. Image2022 7 28 15 2 14

---

Figure 23. Image2022 7 28 15 4 11

---

Figure 24. Image2022 7 28 15 6 40

---

Figure 25. Procedure Heading Icon

---

Figure 26. Procedure Heading Icon

---

Figure 27. Procedure Heading Icon

---

Figure 28. Procedure Heading Icon

---

Figure 29. Procedure Heading Icon

---

Figure 30. Procedure Heading Icon

---

Figure 31. Procedure Heading Icon

---

Figure 32. Procedure Heading Icon

---

Figure 33. Procedure Heading Icon

---

Figure 34. Procedure Heading Icon

---

Figure 35. Procedure Heading Icon

---

Figure 36. Procedure Heading Icon

---

Figure 37. Procedure Heading Icon

---

Figure 38. Procedure Heading Icon

---

Figure 39. Procedure Heading Icon

---

Figure 40. Procedure Heading Icon

---

Figure 41. Procedure Heading Icon

---

Figure 42. Procedure Heading Icon

---

Figure 43. Procedure Heading Icon

---

Figure 44. Worddavb2ee67a7eb9aae5c536610e39a37dcc5

---

Figure 45. Worddav6931c32564b3b0c166f4a26788219144

---

Figure 46. Procedure Heading Icon

---

Figure 47. Image2019 3 8 12 50 6

---

Figure 48. Procedure Heading Icon

---

Figure 49. Procedure Heading Icon

---

Figure 50. Procedure Heading Icon

---

Figure 51. Procedure Heading Icon

---

Figure 52. Procedure Heading Icon

---

Figure 53. Procedure Heading Icon

---

Figure 54. Procedure Heading Icon

---

Figure 55. Procedure Heading Icon

---

Figure 56. Procedure Heading Icon

---

Figure 57. Procedure Heading Icon

---

Figure 58. Procedure Heading Icon

---

Figure 59. Procedure Heading Icon

---

Figure 60. Procedure Heading Icon

---

Figure 61. Procedure Heading Icon

---

Figure 62. Procedure Heading Icon

---

Figure 63. Procedure Heading Icon

---

Figure 64. Procedure Heading Icon

---

Figure 65. Procedure Heading Icon

---

Figure 66. Procedure Heading Icon

---

Figure 67. Procedure Heading Icon

---

Figure 68. Procedure Heading Icon

---

Figure 69. Procedure Heading Icon

---

Figure 70. Worddav336f9b6791fd85e08c8e6897697cd75b

---

Figure 71. Procedure Heading Icon

---

Figure 72. Procedure Heading Icon

---

Figure 73. Procedure Heading Icon

---

Figure 74. Procedure Heading Icon

---

## Overview

NVIDIA OpenFabrics Enterprise Distribution for Linux (MLNX\_OFED) is a single Virtual Protocol Interconnect (VPI) software stack that operates across all NVIDIA network adapter solutions.

NVIDIA OFED (MLNX\_OFED) is an NVIDIA tested and packaged version of OFED and supports two interconnect types using the same RDMA (remote DMA) and kernel bypass APIs called OFED verbs—InfiniBand and Ethernet. Up to 200Gb/s InfiniBand and RoCE (based on the RDMA over Converged Ethernet standard) over 10/25/40/50/100/200GbE are supported with OFED by NVIDIA to enable OEMs and System Integrators to meet the needs end users in the said markets.

Further information on this product can be found in the following MLNX\_OFED documents:

- [Release Notes](#)
- [User Manual](#)

## Software Download

Please visit [nvidia.com/en-us/networking](https://nvidia.com/en-us/networking) Products Software InfiniBand/VPI Drivers [NVIDIA MLNX\\_OFED](#)

## Document Revision History

For the list of changes made to the User Manual, refer to [User Manual Revision History](#).

For the list of older release notes, refer to [Release Notes Revision History](#).

---

# Release Notes

## Release Notes Update History

These are the release notes for MLNX\_OFED LTS. This version provides long-term support (LTS) for customers who wish to utilize the following:

- ConnectX-3
- ConnectX-3 Pro
- Connect-IB
- RDMA experimental verbs library (mlx\_lib)

For other use-cases, it is recommended to use the latest MLNX\_OFED version 5.x.

Version	Date	Description
4.9-7.1.0.0	June 29, 2023	Initial release of this LTS document version.

## Supported NICs Speeds

This document provides instructions on how to install the driver on NVIDIA ConnectX® network adapter solutions supporting the following uplinks to servers.

Uplink/NICs	Driver Name	Uplink Speed
ConnectX®-3/Connect X-3 Pro	mlx4	<ul style="list-style-type: none"><li>• InfiniBand: SDR, QDR, FDR10, FDR</li><li>• Ethernet: 10GbE, 40GbE 56GbE<sup>1</sup></li></ul>
ConnectX-4	mlx5	<ul style="list-style-type: none"><li>• InfiniBand: SDR, QDR, FDR, FDR10, EDR</li><li>• Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 56GbE<sup>1</sup>, 100GbE</li></ul>



Uplink/NICs	Driver Name	Uplink Speed
ConnectX-4 Lx		<ul style="list-style-type: none"> <li>Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE</li> </ul>
ConnectX-5/ConnectX-5 Ex		<ul style="list-style-type: none"> <li>InfiniBand: SDR, QDR, FDR, FDR10, EDR</li> <li>Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 100GbE</li> </ul>
ConnectX-6		<ul style="list-style-type: none"> <li>InfiniBand: SDR, FDR, EDR, HDR</li> <li>Ethernet: 10GbE, 25GbE, 40GbE, 50GbE<sup>2</sup>, 100GbE<sup>2</sup>, 200GbE<sup>2</sup></li> </ul>
ConnectX-6 Dx		<ul style="list-style-type: none"> <li>Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE<sup>1</sup>, 100GbE<sup>1</sup>, 200GbE<sup>2</sup></li> </ul>
Innova™ IPsec EN		<ul style="list-style-type: none"> <li>Ethernet: 10GbE, 40GbE</li> </ul>
Connect-IB®		<ul style="list-style-type: none"> <li>InfiniBand: SDR, QDR, FDR10, FDR</li> </ul>

1. 56GbE is a NVIDIA propriety link speed can be achieved while connecting a NVIDIA adapter card to NVIDIA SX10XX switch series, or connecting a NVIDIA adapter card to another NVIDIA adapter card.

2. Supports both NRZ and PAM4 modes.

## Package Contents

Package	Revision	Licenses
ar_mgr	1.0-0.2.MLNX20201014.g8577618.49710	Mellanox Confidential and Proprietary
dapl	2.1.10.1.mlnx-OFED.4.9.0.1.4.49710	Dual GPL/BSD/CPL
dump_pr	1.0-0.2.MLNX20201014.g8577618.49710	GPLv2 or BSD
fabric-collector	1.1.0.MLNX20170103.89bb2aa-0.1.49710	GPLv2 or BSD

Package	Revision	Licenses
gpio-mlxbf	1.0-0.g6d44a8a	GPLv2
hcoll	4.4.2970-1.49710	Proprietary
i2c-mlx	1.0-0.g422740c	GPLv2
ibacm	41mlnx1-OFED.4.3.3.0.0.49710	GPLv2 or BSD
ibdump	6.0.0-1.49710	BSD2+GPL2
ibsim	0.10-1.49710	GPLv2 or BSD
ibutils	1.5.7.1-0.12.gdcaee2.49710	GPL/BSD
ibutils2	2.1.1-0.121.MLNX20200324.g061a520.49710	<a href="https://www.openib.org/">OpenIB.org</a> BSD license.
infiniband-diags	5.6.0.MLNX20200211.354e4b7-0.1.49710	GPLv2 or BSD
iser	4.9-OFED.4.9.7.1.0.1	GPLv2
isert	4.9-OFED.4.9.7.1.0.1	GPLv2
kernel-mft	4.15.1-9	Dual BSD/GPL
knem	1.1.4.90mlnx2-OFED.23.04.0.5.2.1	BSD and GPLv2
libibcm	41mlnx1-OFED.4.1.0.1.0.49710	GPL/BSD
libibmad	5.4.0.MLNX20190423.1d917ae-0.1.49710	GPLv2 or BSD
libibumad	43.1.1.MLNX20200211.078947f-0.1.49710	GPLv2 or BSD
libibverbs	41mlnx1-OFED.4.9.3.0.0.49710	GPLv2 or BSD
libmlx4	41mlnx1-OFED.4.7.3.0.3.49710	GPLv2 or BSD
libmlx5	41mlnx1-OFED.4.9.0.1.2.49710	GPLv2 or BSD
libpka	1.0-1.g6cc68a2.49710	BSD
librdmacm	41mlnx1-OFED.4.7.3.0.6.49710	GPLv2 or BSD
libvma	9.0.2-1	GPLv2
mlnx-en	4.9-7.1.0.0.g382c630	GPLv2

Package	Revision	Licenses
mlnx-ethtool	5.4-1.49710	GPL
mlnx-iproute2	5.4.0-1.49710	GPL
mlnx-nfsrdma	4.9-OFED.4.9.7.1.0.1	GPLv2
mlnx-nvme	4.9-OFED.4.9.7.1.0.1	GPLv2
mlnx-ofa_kernel	4.9-OFED.4.9.7.1.0.1	GPLv2
mlxbf-livfish	1.0-0.gec08328	GPLv2
mlx-bootctl	1.3-0.g2aa74b7	GPLv2
mlx-l3cache	0.1-1.gebb0728	GPLv2
mlx-pmc	1.1-0.g1141c2e	GPLv2
mlx-trio	0.1-1.g9d13513	GPLv2
mpi-selector	1.0.3-1.49710	BSD
mpitests	3.2.20-e1a0676.49710	BSD
mstflint	4.14.0-3.49710	GPL/BSD
multiperf	3.0-0.14.g5f0fd0e.49710	BSD 3-Clause, GPL v2 or later
multiperf	3.0.0.mlnxlibs-0.13.gcdaa426.49017.49417.49710	BSD 3-Clause, GPL v2 or later
mxm	3.7.3112-1.49710	Proprietary
nvme-snap	2.1.0-126.mlnx	Proprietary
ofed-docs	4.9-OFED.4.9.7.1.0	GPL/BSD
ofed-scripts	4.9-OFED.4.9.7.1.0	GPL/BSD
openmpi	4.0.3rc4-1.49710	BSD
opensm	5.7.2.MLNX20201014.9378048-0.1.49710	GPLv2 or BSD
openvswitch	2.12.1-1.49710	ASL 2.0 and LGPLv2+ and SISSL
perftest	4.5-0.1.g23b8f9c.49710	BSD 3-Clause, GPL v2 or later

Package	Revision	Licenses
perftest	4.5.0.mlnxlibs-0.3.g1121951.49417.49710	BSD 3-Clause, GPL v2 or later
pka-mlxbf	1.0-0.g963f663	GPLv2
qperf	0.4.11-1.49710	BSD 3-Clause, GPL v2
rdma-core	50mlnx1-1.49710	GPLv2 or BSD
rshim	1.18-0.gb99e894	GPLv2
sharp	2.1.2.MLNX20200428.ddda184-1.49710	Proprietary
sockperf	3.7-0.gita1e8e835a689.49710	BSD
srp	4.9-OFED.4.9.7.1.0.1	GPLv2
srptools	41mlnx1-5.49710	GPL/BSD
tmfifo	1.5-0.g31e8a6e	GPLv2
ucx	1.8.0-1.49710	BSD

Release Notes contain the following sections:

- [General Support in MLNX\\_OFED](#)
- [Changes and New Features](#)
- [Bug Fixes](#)
- [Known Issues](#)

## General Support in MLNX\_OFED

### MLNX\_OFED Supported Operating Systems

Operating System	Platform	Default Kernel Version
ALIOS7.2	AArch64	4.19.48-006.ali4000.alios7.aarch64
BCLINUX7.3	x86_64	3.10.0-514.el7.x86_64

Operating System	Platform	Default Kernel Version
BCLINUX7.4	x86_64	3.10.0-693.el7.x86_64
BCLINUX7.5	x86_64	3.10.0-862.el7.x86_64
BCLINUX7.6	x86_64	3.10.0-957.el7.x86_64
BCLINUX7.7	x86_64	3.10.0-1062.el7.bclinux.x86_64
BCLINUX8.1	x86_64	4.19.0-193.1.3.el8.bclinux.x86_64
Debian10.0	x86_64	4.19.0-5-arm64
	AArch64	4.19.0-5-amd64
Debian8.11	x86_64	3.16.0-6-amd64
Debian8.9	x86_64	3.16.0-4-amd64
Debian9.11	x86_64	4.9.0-11-amd64
Debian9.6	x86_64	4.9.0-8-amd64
Debian9.9	x86_64	4.9.0-9-amd64
EulerOS2.0sp9	AArch64	4.19.90- vhulk2006.2.0.h171.eulerosv2r9.aarch64
	x86_64	4.18.0-147.5.1.0.h269.eulerosv2r9.x86_64
Fedora30	x86_64	5.0.9-301.fc30.x86_64
Oracle Linux 6.10	x86_64	4.1.12-124.16.4.el6uek.x86_64
Oracle Linux 7.4	x86_64	4.1.12-94.3.9.el7uek.x86_64
Oracle Linux 7.7	x86_64	4.14.35-1902.3.2.el7uek.x86_64
Oracle Linux 7.8	x86_64	4.14.35-1902.300.11.el7uek.x86_64
Oracle Linux 7.9	x86_64	5.4.17-2011.6.2.el7uek.x86_64
Oracle Linux 8.0	x86_64	4.18.0-80.7.2.el8_0.x86_64
Oracle Linux 8.1	x86_64	4.18.0-147.el8.x86_64
Oracle Linux 8.2	x86_64	5.4.17-2011.1.2.el8uek.x86_64
Oracle Linux 8.3	x86_64	5.4.17-2011.7.4.el8uek.x86_64
RHEL/CentOS6.10	x86_64	2.6.32-754.el6.x86_64

Operating System	Platform	Default Kernel Version
RHEL/CentOS6.3	x86_64	2.6.32-279.el6.x86_64
RHEL/CentOS7.2	ppc64	3.10.0-327.el7.ppc64
	ppc64le	3.10.0-327.el7.ppc64le
	x86_64	3.10.0-327.el7.x86_64
RHEL/CentOS7.3	x86_64	3.10.0-514.el7.x86_64
RHEL/CentOS7.4	ppc64	3.10.0-693.el7.ppc64
	ppc64le	3.10.0-693.el7.ppc64le
	x86_64	3.10.0-693.el7.x86_64
RHEL/CentOS7.4al ternate	AArch64	4.11.0-44.el7a.aarch64
RHEL/CentOS7.5	ppc64	3.10.0-862.el7.ppc64
	ppc64le	3.10.0-862.el7.ppc64le
	x86_64	3.10.0-862.el7.x86_64
RHEL/CentOS7.5al ternate	AArch64	4.14.0-49.el7a.aarch64
RHEL/CentOS7.6	x86_64	3.10.0-957.el7.ppc64
	ppc64le	3.10.0-957.el7.ppc64le
	ppc64	3.10.0-957.el7.x86_64
RHEL/CentOS7.6al ternate	AArch64	4.14.0-115.el7a.aarch64
	ppc64le	4.14.0-115.el7a.ppc64le
RHEL/CentOS7.7	ppc64	3.10.0-1062.el7.ppc64
	ppc64le	3.10.0-1062.el7.ppc64le
	x86_64	3.10.0-1062.el7.x86_64
RHEL/CentOS7.8	ppc64	3.10.0-1127.el7.ppc64
	ppc64le	3.10.0-1127.el7.ppc64le
	x86_64	3.10.0-1127.el7.x86_64
RHEL/CentOS7.9	ppc64	3.10.0-1160.el7.ppc64

Operating System	Platform	Default Kernel Version
	ppc64le	3.10.0-1160.el7.ppc64le
	x86_64	3.10.0-1160.el7.x86_64
RHEL/CentOS8.0	AArch64	4.18.0-80.el8.aarch64
	ppc64le	4.18.0-80.el8.ppc64le
	x86_64	4.18.0-80.el8.x86_64
RHEL/CentOS8.1	AArch64	4.18.0-147.el8.aarch64
	ppc64le	4.18.0-147.el8.ppc64le
	x86_64	4.18.0-147.el8.x86_64
RHEL/CentOS8.2	AArch64	4.18.0-193.el8.aarch64
	ppc64le	4.18.0-193.el8.ppc64le
	x86_64	4.18.0-193.el8.x86_64
RHEL/CentOS8.3	AArch64	4.18.0-240.el8.aarch64
	ppc64le	4.18.0-240.el8.ppc64le
	x86_64	4.18.0-240.el8.x86_64
RHEL/CentOS8.4	AArch64	4.18.0-305.el8.aarch64
	ppc64le	4.18.0-305.el8.ppc64le
	x86_64	4.18.0-305.el8.x86_64
RHEL/CentOS8.5	AArch64	4.18.0-348.el8.aarch64
	ppc64le	4.18.0-348.el8.ppc64le
	x86_64	4.18.0-348.el8.x86_64
RHEL/CentOS8.6	AArch64	4.18.0-372.9.1.el8.aarch64
	ppc64le	4.18.0-372.9.1.el8.ppc64le
	x86_64	4.18.0-372.9.1.el8.x86_64
RHEL/CentOS8.7	AArch64	4.18.0-425.3.1.el8.aarch64
	ppc64le	4.18.0-425.3.1.el8.ppc64le
	x86_64	4.18.0-425.3.1.el8.x86_64

Operating System	Platform	Default Kernel Version
RHEL/Rocky 8.8	AArch64	4.18.0-477.10.1.el8_8.aarch64
	ppc64le	4.18.0-477.10.1.el8_8.ppc64le
	x86_64	4.18.0-477.10.1.el8_8.x86_64
SLES11SP3	x86_64	3.0.76-0.11-default
SLES11SP4	ppc64	3.0.101-63-ppc64
	x86_64	3.0.101-63-default
SLES12SP2	x86_64	4.4.21-69-default
SLES12SP3	x86_64	4.4.73-5-default
	ppc64le	4.4.73-5-default
SLES12SP4	x86_64	4.12.14-94.41-default
	ppc64le	4.12.14-94.41-default
	AArch64	4.12.14-94.41-default
SLES12SP5	x86_64	4.12.14-120-default
	ppc64le	4.12.14-120-default
	AArch64	4.12.14-120-default
SLES15SP0	x86_64	4.12.14-23-default
SLES15SP1	x86_64	4.12.14-195-default
	ppc64le	4.12.14-195-default
	AArch64	4.12.14-195-default
SLES15SP2	x86_64	5.3.18-22-default
	ppc64le	5.3.18-22-default
	AArch64	5.3.18-22-default
SLES15SP3	x86_64	5.3.18-57-default
	ppc64le	5.3.18-57-default
	AArch64	5.3.18-57-default
Ubuntu14.04	x86_64	3.13.0-27-generic



Operating System	Platform	Default Kernel Version
Ubuntu16.04	ppc64le	4.4.0-21-generic
	x86_64	4.4.0-21-generic
Ubuntu18.04	x86_64	4.15.0-20-generic
	ppc64le	4.15.0-20-generic
	AArch64	4.15.0-20-generic
Ubuntu19.04	x86_64	5.0.0-13-generic
Ubuntu19.10	x86_64	5.3.0-19-generic
Ubuntu20.04	x86_64	5.4.0-26-generic
	ppc64le	5.4.0-26-generic
	AArch64	5.4.0-26-generic
Kernel 5.5	x86_64	5.5

**Notes:**

- 32 bit platforms are no longer supported in MLNX\_OFED.
- 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 Manual for instructions.
- Upgrading MLNX\_OFED on your cluster requires upgrading all of its nodes to the newest version as well.
- All OSs listed above are fully supported in Paravirtualized and SR-IOV Environments with Linux KVM Hypervisor.

## Supported Non-Linux Virtual Machines

The following are the supported non-Linux Virtual Machines in this current MLNX\_OFED version:

NIC	Windows Virtual Machine Type	WinOF version	Protocol
ConnectX-3	Windows 2012 R2 DC	MLNX_VPI 5.50	IPoIB, ETH

NIC	Windows Virtual Machine Type	WinOF version	Protocol
ConnectX-3 Pro	Windows 2016 DC	MLNX_VPI 5.50	IPoIB, ETH
ConnectX-4	Windows 2012 R2 DC	MLNX_WinOF2 2.40	IB, IPoIB, ETH
ConnectX-4 Lx	Windows 2016 DC	MLNX_WinOF2 2.40	IB, IPoIB, ETH

## Support in ASAP2™

### Warning

Accelerated Switch and Packet Processing (ASAP<sup>2</sup>) is not supported in this MLNX\_OFED version.

## NFS over RDMA (NFSoverRDMA) Supported Operating Systems

Below is a list of all the OSs on which NFSoverRDMA is supported.

- SLES12 SP4
- SLES12 SP5
- SLES15 SP1
- Ubuntu 18.04.3
- RedHat 7.5
- RedHat 7.6
- RedHat 7.7
- RedHat 7.8
- RedHat 8.0

- RedHat 8.1
- RedHat 8.6

## Lustre Versions Supported by MLNX\_OFED

- Lustre 2.12.3
- Lustre 2.13.0

## Hardware and Software Requirements

The following are the hardware and software requirements of the current MLNX\_OFED version.

- 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:

Operating System	Required Packages Installation Command
RHEL/Oracle Linux/Fedora	<code>yum install perl pciutils python gcc-gfortran libxml2-python tcsh libnl.i686 libnl expat glib2 tcl libstdc++ bc tk gtk2 atk cairo numactl pkgconfig ethtool lsof</code>
XenServer	<code>yum install perl pciutils python libxml2-python libnl expat glib2 tcl bc libstdc++ tk pkgconfig ethtool</code>
SLES 12	<code>zypper install pkg-config expat libstdc++6 libglib-2_0-0 lib-gtk-2_0-0 tcl libcairo2 tcsh python bc pciutils libatk-1_0-0 tk python-libxml2 lsof libnl3-200 ethtool lsof</code>
SLES 15	<code>python ethtool libatk-1_0-0 python2-libxml2-python tcsh lib-stdc++6-devel-gcc7 libgtk-2_0-0 tcl libopenssl1_1 libnl3-200 make libcairo2 expat libmnl0 insserv-compat pciutils lsof lib-glib-2_0-0 pkg-config tk</code>

Operating System	Required Packages Installation Command
Ubuntu/Debian	apt-get install perl dpkg autotools-dev autoconf libtool auto- make1.10 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 ethtool lsof

## Supported NICs Firmware Versions

### Warning

This MLNX\_OFED version provides long term support (LTS) for customers who wish to utilize ConnectX-3, ConnectX-3 Pro and Connect-IB, as well as RDMA experimental verbs library (mlnx\_lib). Any MLNX\_OFED version starting from v5.1 and above does not support any of the adapter cards mentioned.

This current MLNX\_OFED version supports the following Mellanox network adapter cards firmware versions:

NIC	Recommended Firmware Rev.	Additional Firmware Rev. Supported
ConnectX®-3/ConnectX-3 Pro	2.42.5000	2.40.7000
ConnectX-4	12.28.2006	12.27.4000
ConnectX-4 Lx	14.28.2006	14.27.1016
ConnectX-5/ConnectX-5 Ex	16.28.2006	16.27.2008
ConnectX-6	20.28.2006	20.27.2008
ConnectX-6 Dx	22.28.2006	N/A
Innova IPsec EN	16.28.2006	16.27.2008

NIC	Recommended Firmware Rev.	Additional Firmware Rev. Supported
Connect-IB	10.16.1200	10.16.1020

For the official firmware versions, please visit the following site:  
<https://network.nvidia.com/support/firmware/firmware-downloads/>

## RDMA CM and RoCE Modes

### RoCE Modes Matrix

Software Stack / Inbox Distribution	RoCEv1 (IP Based GIDs) Supported as of Version		RoCEv2 Supported as of Version		RoCEv1 & RoCEv2 (RoCE per GID) Supported as of Version
	ConnectX-3/ ConnectX-3 Pro	ConnectX-4/ ConnectX-4 Lx/ ConnectX-5/ ConnectX-5 Ex	ConnectX-3 Pro	ConnectX-4/ ConnectX-4 Lx/ ConnectX-5/ ConnectX-5 Ex	
MLNX_OFED	2.1-x.x.x	3.0-x.x.x	2.3-x.x.x	3.0-x.x.x	3.0-x.x.x
<a href="http://kernel.org">Kernel.org</a>	3.14	4.4	4.4	4.4	4.4
RHEL	6.6, 7.0	-	-	-	-
SLES	12	-	-	-	-
Ubuntu	14.04.4, 16.04, 15.10	-	-	-	-

**Note:** Support for ConnectX-5 and ConnectX-5 Ex adapter cards in MLNX\_OFED starts from v4.0.

## RDMA CM Default RoCE Mode

The default RoCE mode on which RDMA CM runs is RoCEv2 instead of RoCEv1, starting from MLNX\_OFED v4.1. RDMA\_CM session requires both the client and server sides to support the same RoCE mode. Otherwise, the client will fail to connect to the server. For further information, refer to [RDMA CM and RoCE Version Defaults](#) Community post.

## MLNX\_OFED Unsupported Functionalities/Features/NICs

The following are the unsupported functionalities/features/NICs in MLNX\_OFED:

- ConnectX®-2 Adapter Card
- Relational Database Service (RDS)
- Ethernet over InfiniBand (EoIB) - mlx4\_vnic
- mthca InfiniBand driver
- Ethernet IPoIB (eIPoIB)
- Soft-RoCE

# Changes and New Features

## New Features

The following are the new features that have been added to this version of MLNX\_OFED.

Feature	Description
NFSoRDMA on RHEL8.6	Added support for <a href="#">NFSoRDMA</a> support on RHEL8.6 OS.
Bug Fixes	See <a href="#">Bug Fixes</a> section.

For additional information on the new features, please refer to MLNX\_OFED User Manual.

## Customer Affecting Changes

Change	Description
Installation, NVMe SNAP	Support for NVMe SNAP is discontinued on RHEL7.6-alternate.
Installation, ibutils	Support for ibutils is discontinued on RHEL7.6-alternate.

## Bug Fixes

This table lists the bugs fixed in this release.

For the list of old bug fixes, please refer to MLNX\_OFED Archived Bug Fixes file at: [http://www.mellanox.com/pdf/prod\\_software/MLNX\\_OFED\\_Archived\\_Bug\\_Fixes.pdf](http://www.mellanox.com/pdf/prod_software/MLNX_OFED_Archived_Bug_Fixes.pdf)

Internal Reference Number	Description
3419761	Memory allocation issue may lead to OOM.
	Memory, OOM
	<b>Discovered in Release:</b> 4.9-2.2.4.0
	<b>Fixed in Release:</b> 4.9-7.1.0.0
3383081	<b>Description:</b> When using SLES15SP3 with updated kernel, the OFED build fails.
	<b>Keywords:</b> Installation, SLES15SP3
	<b>Discovered in Release:</b> 4.9-6.0.6.0
	<b>Fixed in Release:</b> 4.9-7.1.0.0
3427527	<b>Description:</b> Set coalesce parameters was not supported over old LTS branches.
	<b>Keywords:</b> NetDev, Coalesce Parameters

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-7.1.0.0
3201693	<b>Description:</b> In the flow of MR destruction, use-after-free was found.
	<b>Keywords:</b> RDMA
	<b>Fixed in Release:</b> 4.9-6.0.6.0
3203796/3251496	<b>Description:</b> Using RHEL8.6 with an updated kernel, causes issues with installation.
	<b>Keywords:</b> Installation
	<b>Fixed in Release:</b> 4.9-6.0.6.0
2944030	<b>Description:</b> An issue with the Udev script caused non-NVIDIA devices to be renamed.
	<b>Keywords:</b> ASAP <sup>2</sup> , Udev, Naming
	<b>Fixed in Release:</b> 4.9-5.1.0.0
3037901	<b>Description:</b> RDMA traffic may fail due to incorrect tracking of outstanding work requests.
	<b>Keywords:</b> RDMA
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-5.1.0.0
2976200	<b>Description:</b> On the passive side, when the RDMACM disconnectReq event arrives, if the current state is MRA_REP_RCVD, it needs to cancel the MAD before entering the DREQ_RCVD and TIMEWAIT state, otherwise the destroy_id may block the request until this MAD reaches timeout.
	<b>Keywords:</b> RDMACM, MRA, destroy_id
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-5.1.0.0



Internal Reference Number	Description
2753944	<b>Description:</b> On rare occasion, registering a device (ib_register_device()) and loading modules in parallel (in this case (ib_cm)), ay cause a race condition to occur which would stop ib_cm from loading properly.
	<b>Keywords:</b> RDMA, ib_core, Racing Condition
	<b>Fixed in Release:</b> 4.9-5.1.0.0
2802401	<b>Description:</b> When removing the nvmet port from configs caused a use-after-free condition.
	<b>Keywords:</b> nvmet-rdma Module
	<b>Discovered in Release:</b> 4.9-0.1.7.0 <b>Fixed in Release:</b> 4.9-4.1.7.0
2162639	<b>Description:</b> MLNX_OFED includes several python tools, such as mlnx_qos, which rely on python modules included in the same package. On Ubuntu 20.04 OS, those are installed into a directory that is not in python modules search path.
	<b>Keywords:</b> mlnx_qos, Ubuntu
	<b>Discovered in Release:</b> 4.9-0.1.7.0 <b>Fixed in Release:</b> 4.9-4.0.8.0
2635628	<b>Description:</b> openibd does not load automatically after reboot on Suler2sp9 OS.
	<b>Keywords:</b> openibd, Suler2sp9
	<b>Discovered in Release:</b> 4.9-3.1.5.0 <b>Fixed in Release:</b> 4.9-4.0.8.0
2748862	<b>Description:</b> add-kernel-support flag was not supported on Oracle Linux 7.9 causing an installation issue.
	<b>Keywords:</b> openibd, Euleros2u0sp9
	<b>Discovered in Release:</b> 4.9-0.1.7.0 <b>Fixed in Release:</b> 4.9-4.0.8.0

Internal Reference Number	Description
2748862	<b>Description:</b> add-kernel-support flag was not supported on Oracle Linux 7.9 causing an installation issue.
	<b>Keywords:</b> Installation, Oracle Linux 7.9
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-4.0.8.0
2396956	<b>Description:</b> Fixed an issue were device under massive load may hit iommu allocation failures. For more information see " <a href="#">RX Page Cache Size Limit</a> " section in the user manual.
	<b>Keywords:</b> Legacy libibverbs
	<b>Discovered in Release:</b> 4.9-2.2.4.0
	<b>Fixed in Release:</b> 4.9-3.1.5.0
2434638	<b>Description:</b> Fixed an issue where "ibv_devinfo -v" command did not print some of the MEM_WINDOW capabilities, even though they were supported.
	<b>Keywords:</b> Legacy libibverbs
	<b>Discovered in Release:</b> 4.9-2.2.4.0
	<b>Fixed in Release:</b> 4.9-3.1.5.0
2292762	<b>Description:</b> Fixed a kernel panic scenario that may have taken place when using sysfs to cancel the probing of VFs and performing reboot while the VFs are still managed by the mlx5 driver.
	<b>Keywords:</b> Proved VFs
	<b>Discovered in Release:</b> 5.1-2.3.7.1
	<b>Fixed in Release:</b> 5.1-2.5.8.0
2265055	<b>Description:</b> Added missing release of the lock held in the traffic class error flow.
	<b>Keywords:</b> mutex_unlock
	<b>Discovered in Release:</b> 4.9-0.1.7.0

Internal Reference Number	Description
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2245228	<b>Description:</b> Fixed an issue of a crash when attempting to access roce_enable sysfs in unprobed VFs.
	<b>Keywords:</b> roce_enable, unprobed VFs
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2211311	<b>Description:</b> Fixed an issue where Rx port buffers cell size was wrong, leading to wrong buffers size reported by mlnx_qos/netdev qos/buffer_size sysfs.
	<b>Keywords:</b> mlx5e, RX buffers, mlnx_qos
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2143067	<b>Description:</b> If Openibd was configured to enable the SRP daemon, it now also enables srp_daemon from rdma-core.
	<b>Keywords:</b> Openibd, SRP daemon, srp_daemon, rdma-core
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2192791	<b>Description:</b> Fixed the issue where packages neohost-backend and neohost-sdk were not properly removed by the uninstallation procedure and may have required manual removal before re-installing or upgrading the MLNX_OFED driver.
	<b>Keywords:</b> NEO-Host, SDK
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2226715	<b>Description:</b> Fixed an issue where bringing up PF interface failed when using SR-IOV and configuring RoCE mode for v2 only.
	<b>Keywords:</b> PF, SR-IOV, RoCE v2

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2242041	<b>Description:</b> Fixed leak of memory pages when using ODP.
	<b>Keywords:</b> ODP
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2249090	<b>Description:</b> Fixed an issue where IBV_EXP_ACCESS_TUNNELED_ATOMIC capability did not work for ibv_exp_reg_mr experimental verb.
	<b>Keywords:</b> ibv_exp_reg_mr
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2178677	<b>Description:</b> Marked "ibacm-devel" as a package name explicitly to avoid accidentally including a symlink to it in the UPSTREAM_LIBS version of MLNX_OFED.
	<b>Keywords:</b> ibacm-devel, UPSTREAM_LIBS
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2255829	<b>Description:</b> Fixed an issue with metadata packages generation in the eth-only directory. This allows using the directory as a repository for package managers.
	<b>Keywords:</b> Metadata packages
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2328754	<b>Description:</b> Fixed an issue that may have caused panic during peer memory invalidation flow.
	<b>Keywords:</b> GPUDirect

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.9-0.1.7.0
	<b>Fixed in Release:</b> 4.9-2.2.4.0
2119017	<b>Description:</b> Fixed the issue where injecting EEH may cause extra Kernel prints, such as: "EEH: Might be infinite loop in mlx5_core driver".
	<b>Keywords:</b> EEH, kernel
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2076546	<b>Description:</b> Fixed the issue where in RPM-based OSs with non-default kernels, using repositories after re-creating the installer (using --add-kernel-support) would result in improper installation of the drivers.
	<b>Keywords:</b> Installation, OS
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2114957	<b>Description:</b> Fixed the issue where MLNX_OFED installation may have depended on python2 package even when attempting to install it on OSs whose default package is python3.
	<b>Keywords:</b> Installation, python
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2143258	<b>Description:</b> Fixed a typo in perftest package where help messages wrongly displayed the conversion result between Gb/s and MB/s (20^2 instead of 2^20).
	<b>Keywords:</b> perftest
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0

Internal Reference Number	Description
2094216	<b>Description:</b> Fixed the issue of when one of the LAG slaves went down, LAG deactivation failed, ultimately causing bandwidth degradation.
	<b>Keywords:</b> RoCE LAG
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2133778	<b>Description:</b> The mlx5 driver maintains a subdirectory for every open eth port in /sys/kernel/debug/. For the default network namespace, the sub-directory name is the name of the interface, like "eth8". The new convention for the network interfaces moved to the non-default network namespaces is the interfaces name followed by "@" and the port's PCI ID. For example: "eth8@0000:af:00.3".
	<b>Keywords:</b> Namespace
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2122684	<b>Description:</b> Fixed the issue where OFED uninstallation resulted in the removal of dependency packages, such as qemu-system-* (qemu-system-x86).
	<b>Keywords:</b> Uninstallation, dependency, qemu-system-x86
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2135476	<b>Description:</b> Added KMP ability to install MLNX_OFED Kernel modules on SLES12 SP5 and SLES15 kernel maintenance updates.
	<b>Keywords:</b> KMP, SLES, kernel
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2149577	<b>Description:</b> Fixed the issue where openibd script load used to fail when esp6_offload module did not load successfully.

Internal Reference Number	Description
	<p><b>Keywords:</b> openibd, esp6_offload</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p> <p><b>Fixed in Release:</b> 4.9-0.1.7.0</p>
2163879	<p><b>Description:</b> Added dependency of package mpi-selectors on perl-Getopt-Long system package. On minimal installs of RPM-based OSs, installing mpi-selectors will also install the required system package perl-Getopt-Long.</p> <p><b>Keywords:</b> Dependency, perl-Getopt-Long</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p> <p><b>Fixed in Release:</b> 4.9-0.1.7.0</p>
2107532	<p><b>Description:</b> Fixed the issue where in certain rare scenarios, due to Rx page not being replenished, the same page fragment mistakenly became assigned to two different Rx descriptors.</p> <p><b>Keywords:</b> Memory corruption, Rx page recycle</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p> <p><b>Fixed in Release:</b> 4.9-0.1.7.0</p>
2116234	<p><b>Description:</b> Fixed the issue where ibsim was missing after OFED installation.</p> <p><b>Keywords:</b> ibsim, installation</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p> <p><b>Fixed in Release:</b> 4.9-0.1.7.0</p>
2116233	<p><b>Description:</b> Fixed an issue where ucx-kmem was missing after OFED installation.</p> <p><b>Keywords:</b> ucx-kmem, installation</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p> <p><b>Fixed in Release:</b> 4.9-0.1.7.0</p>

Internal Reference Number	Description
2107776	<b>Description:</b> Fixed a driver load issue with Errata-kernel on SLES15 SP1.
	<b>Keywords:</b> Load, SLES, Errata
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2105536	<b>Description:</b> Fixed an issue in the Hairpin feature which prevented adding hairpin flows using TC tool.
	<b>Keywords:</b> Hairpin, TC
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2090321	<b>Description:</b> Fixed the issue where WQ queue flushing was not handled properly in the event of EEH.
	<b>Keywords:</b> WQ, EEH
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2076311	<b>Description:</b> Fixed a rare kernel crash scenario when exiting an application that uses RMPP mads intensively.
	<b>Keywords:</b> MAD RMPP
	<b>Discovered in Release:</b> 4.0-1.0.1.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2096998	<b>Description:</b> Fixed the issue where NEO-Host could not be installed from the MLNX_OFED package when working on Ubuntu and Debian OSs.
	<b>Keywords:</b> NEO-Host, Ubuntu, Debian
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0



Internal Reference Number	Description
2057076	<b>Description:</b> Fixed the issue where installing MLNX_OFED using --add-kernel-support option did not work over RHEL 8 OSs.
	<b>Keywords:</b> --add-kernel-support, installation, RHEL
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2090186	<b>Description:</b> Fixed a possible kernel crash scenario when AER/slot reset in done in parallel to user space commands execution.
	<b>Keywords:</b> mlx5_core, AER, slot reset
	<b>Discovered in Release:</b> 4.3-1.0.1.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2093410	<b>Description:</b> Added missing ECN configuration under sysfs for PFs in SwitchDev mode.
	<b>Keywords:</b> sysfs, ASAP, SwitchDev, ECN
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1916029	<b>Description:</b> Fixed the issue of when firmware response time to commands became very long, some commands failed upon timeout. The driver may have then triggered a timeout completion on the wrong entry, leading to a NULL pointer call trace.
	<b>Keywords:</b> Firmware, timeout, NULL
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2036394	<b>Description:</b> Added driver support for kernels with the old XDP_REDIRECT infrastructure that uses the following NetDev operations: .ndo_xdp_flush and .ndo_xdp_xmit.
	<b>Keywords:</b> XDP_REDIRECT, Soft lockup
	<b>Discovered in Release:</b> 4.7-3.2.9.0

Internal Reference Number	Description
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2072871	<b>Description:</b> Fixed an issue where the usage of --excludedocs Open MPI RPM option resulted in the removal of non-documentation related files.
	<b>Keywords:</b> --excludedocs, Open MPI, RPM
	<b>Discovered in Release:</b> 4.5-1.0.1.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2072884	<b>Description:</b> Removed all cases of automated loading of MLNX_OFED kernel modules outside of openibd to preserve the startup process of previous MLNX_OFED versions. These loads conflict with openibd, which has its own logic to overcome issues. Such issues can be inbox driver load instead of MLNX_OFED, or module load with wrong parameter value. They might also load modules while openibd is trying to unload the driver stack.
	<b>Keywords:</b> Installation, openibd
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2052037	<b>Description:</b> Disabled automated loading of some modules through udev triggers to preserve the startup process of previous MLNX_OFED versions.
	<b>Keywords:</b> Installation, udev
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2022634	<b>Description:</b> Fixed a typo in the packages build command line which could cause the installation of MLNX_OFED on SLES OSs to fail when using the option --without-depcheck.
	<b>Keywords:</b> Installation, SLES
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0

Internal Reference Number	Description
2022619	<b>Description:</b> Fixed the issue where uninstallation of MLNX_OFED would hang due to a bug in the package dependency check.
	<b>Keywords:</b> Uninstallation, dependency
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2047221	<b>Description:</b> Reference count (refcount) for RDMA connection ID (cm_id) was not incremented in rdma_resolve_addr() function, resulting in a cm_id use-after-free access. A fix was applied to increment the cm_id refcount.
	<b>Keywords:</b> rdma_resolve_addr(), cm_id
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2045181	<b>Description:</b> Fixed a race condition which caused kernel panic when moving two ports to SwitchDev mode at the same time.
	<b>Keywords:</b> ASAP, SwitchDev, race
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2004488	<b>Description:</b> Allowed accessing sysfs hardware counters in SwitchDev mode.
	<b>Keywords:</b> ASAP, hardware counters, sysfs, SwitchDev
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2030943	<b>Description:</b> Function smp_processor_id() is called in the RX page recycle flow to determine the core to run on. This is intended to run in NAPI context. However, due to a bug in backporting, the RX page recycle was mistakenly called also in the RQ close flow when not needed.
	<b>Keywords:</b> Rx page recycle, smp_processor_id

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2074487	<b>Description:</b> Fixed an issue where port link state was automatically changed (without admin state involvement) to "UP" after reboot.
	<b>Keywords:</b> Link state, UP
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2022618	<b>Description:</b> Fixed a hang with ConnectX-3 adapter cards when running over SLES 11 OS.
	<b>Keywords:</b> OS, SLES, ConnectX-3
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2064711	<b>Description:</b> Fixed an issue where RDMA CM connection failed when port space was small.
	<b>Keywords:</b> RDMA CM
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2076424	<b>Description:</b> Traffic mirroring with OVS offload and non-offload over VxLAN interface is now supported. <b>Note:</b> For kernel 4.9, make sure to use a dedicated OVS version.
	<b>Keywords:</b> VxLAN, OVS
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1828321	<b>Description:</b> Fixed the issue of when working with VF LAG while the bond device is in active-active mode, running fwreset would result in unequal traffic on both PFs, and PFs would not reach line rate.
	<b>Keywords:</b> VF LAG, bonding, PF

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1975293	<b>Description:</b> Installing OFED with --with-openswitch flag no longer requires manual removal of the existing Open vSwitch.
	<b>Keywords:</b> OVS, Open vSwitch, openswitch
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1939719	<b>Description:</b> Fixed an issue of when running openibd restart after the installation of MLNX_OFED on SLES12 SP5 and SLES15 SP1 OSs with the latest Kernel (v4.12.14) resulted in an error that the modules did not belong to that Kernel. This was due to the fact that the module installed by MLNX_OFED was incompatible with new Kernel's module.
	<b>Keywords:</b> SLES, operating system, OS, installation, Kernel, module
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
2001966	<b>Description:</b> Fixed an issue of when bond was created over VF netdevices in SwitchDev mode, the VF netdevice would be treated as representor netdevice. This caused the mlx5_core driver to crash in case it received netdevice events related to bond device.
	<b>Keywords:</b> PF, VF, SwitchDev, netdevice, bonding
	<b>Discovered in Release:</b> 4.7-3.2.9.0
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1816629	<b>Description:</b> Fixed an issue where following a bad affinity occurrence in VF LAG mode, traffic was sent after the port went up/down in the switch.
	<b>Keywords:</b> Traffic, VF LAG
	<b>Discovered in Release:</b> 4.6-1.0.1.1

Internal Reference Number	Description
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1718531	<b>Description:</b> Added support for VLAN header rewrite on CentOS 7.2 OS.
	<b>Keywords:</b> VLAN, ASAP, switchdev, CentOS 7.2
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1556337	<b>Description:</b> Fixed the issue where adding VxLAN decapsulation rule with enc_tos and enc_ttl failed.
	<b>Keywords:</b> VxLAN, decapsulation
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.9-0.1.7.0
1949260	<b>Description:</b> Fixed a race condition that resulted in kernel panic when running IPoIB traffic in Connected mode.
	<b>Keywords:</b> IPoIB
	<b>Discovered in Release:</b> 4.5-1.0.1.0
	<b>Fixed in Release:</b> 4.7-3.2.9.0
1973828	<b>Description:</b> Fixed wrong EEPROM length for small form factor (SFF) 8472 from 256 to 512 bytes.
	<b>Keywords:</b> EEPROM, SFF
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.7-3.2.9.0
1915553	<b>Description:</b> Fixed the issue where errno field was not sent in all error flows of ibv_reg_mr API.
	<b>Keywords:</b> ibv_reg_mr
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.7-3.2.9.0

Internal Reference Number	Description
1970901	<b>Description:</b> Fixed the issue where mlx5 IRQ name did not change to express the state of the interface.
	<b>Keywords:</b> Ethernet, PCIe, IRQ
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.7-3.2.9.0
1915587	<b>Description:</b> Udaddy application is now functional in Legacy mode.
	<b>Keywords:</b> Udaddy, MLNX_OFED legacy, RDMA-CM
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.7-3.2.9.0
1931421	<b>Description:</b> Added support for E-Switch (SR-IOV Legacy) mode in RHEL 7.7 OSs.
	<b>Keywords:</b> E-Switch, SR-IOV, RHEL, RedHat
	<b>Discovered in Release:</b> 4.7-1.0.0.1
	<b>Fixed in Release:</b> 4.7-3.2.9.0
1945411/1839353	<b>Description:</b> Fixed the issue of when XDP_REDIRECT fails, pages got double-freed due to a bug in the refcnt_bias feature.
	<b>Keywords:</b> XDP, XDP_REDIRECT, refcnt_bias
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.7-3.2.9.0
1547200	<b>Description:</b> Fixed an issue where IPoIB Tx queue may get stuck, leading to timeout warnings in dmesg.
	<b>Keywords:</b> IPoIB
	<b>Discovered in Release:</b> 4.5-1.0.1.0
	<b>Fixed in Release:</b> 4.7-1.0.0.1
1817636	<b>Description:</b> Fixed the issue of when disabling one port on the Server side, VF-LAG Tx Affinity would not work on the Client side.

Internal Reference Number	Description
	<p><b>Keywords:</b> VF-LAG, Tx Affinity</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p> <p><b>Fixed in Release:</b> 4.7-1.0.0.1</p>
1800525	<p><b>Description:</b> When configuring the Time-stamping feature, CQE compression will be disabled. This fix entails the removal of a warning message that appeared upon attempting to disable CQE compression when it has already been disabled.</p> <p><b>Keywords:</b> Time-stamping, CQE compression</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p> <p><b>Fixed in Release:</b> 4.7-1.0.0.1</p>
1431282	<p><b>Description:</b> Fixed the issue where software reset may have resulted in an order inversion of interface names.</p> <p><b>Keywords:</b> Software reset</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.7-1.0.0.1</p>
1843020	<p><b>Description:</b> Server reboot may result in a system crash.</p> <p><b>Keywords:</b> reboot, crash</p> <p><b>Discovered in Release:</b> 4.2-1.2.0.0</p> <p><b>Fixed in Release:</b> 4.7-1.0.0.1</p>
1734102	<p><b>Description:</b> Fixed the issue where Ubuntu v16.04.05 and v16.04.05 OSs could not be used with their native kernels.</p> <p><b>Keywords:</b> Ubuntu, Kernel, OS</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p> <p><b>Fixed in Release:</b> 4.7-1.0.0.1</p>
1811973	<p><b>Description:</b> VF mirroring offload is now supported.</p> <p><b>Keywords:</b> ASAP<sup>2</sup>, VF mirroring</p>



Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.7-1.0.0.1
1841634	<b>Description:</b> The number of guaranteed counters per VF is now calculated based on the number of ports mapped to that VF. This allows more VFs to have counters allocated.
	<b>Keywords:</b> Counters, VF
	<b>Discovered in Release:</b> 4.4-1.0.0.0
	<b>Fixed in Release:</b> 4.7-1.0.0.1
1758983	<b>Description:</b> Installing MLNX_OFED on RHEL 7.6 OSs platform x86_64 and RHEL 7.6 ALT OSs platform PPCLE using YUM is now supported.
	<b>Keywords:</b> RHEL, RedHat, YUM, OS, operating system
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Fixed in Release:</b> 4.7-1.0.0.1
1523548	<b>Description:</b> Fixed the issue where RDMA connection persisted even after dropping the network interface.
	<b>Keywords:</b> Network interface, RDMA
	<b>Discovered in Release:</b> 4.4-1.0.0.0
	<b>Fixed in Release:</b> 4.6-1.0.1.1
1712870	<b>Description:</b> Fixed the issue where small packets with non-zero padding were wrongly reported as "checksum complete" even though the padding was not covered by the csum calculation. These packets now report "checksum unnecessary". In addition, an ethtool private flag has been introduced to control the "checksum complete" feature: <code>ethtool --set-priv-flags eth1 rx_no_csum_complete on/off</code>
	<b>Keywords:</b> csum error, checksum, mlx5_core
	<b>Discovered in Release:</b> 4.5-1.0.1.0

Internal Reference Number	Description
	<b>Fixed in Release:</b> 4.6-1.0.1.1
1648597	<b>Description:</b> Fixed the wrong wording in the FW tracer ownership startup message (from "FW Tracer Owner" to "FWTracer: Ownership granted and active").
	<b>Keywords:</b> FW Tracer
	<b>Discovered in Release:</b> 4.5-1.0.1.0
	<b>Fixed in Release:</b> 4.6-1.0.1.1
1581631	<b>Description:</b> Fixed the issue where GID entries referenced to by a certain user application could not be deleted while that user application was running.
	<b>Keywords:</b> RoCE, GID
	<b>Discovered in Release:</b> 4.5-1.0.1.0
	<b>Fixed in Release:</b> 4.6-1.0.1.1
1368390	<b>Description:</b> Fixed the issue where MLNX_OFED could not be installed on RHEL 7.x Alt OSs using YUM repository.
	<b>Keywords:</b> Installation, YUM, RHEL
	<b>Discovered in Release:</b> 4.3-3.0.2.1
	<b>Fixed in Release:</b> 4.6-1.0.1.1
1531817	<b>Description:</b> Fixed an issue of when the number of channels configured was less than the number of CPUs available, part of the CPUs would not be used by Tx queues.
	<b>Keywords:</b> Performance, Tx, CPU
	<b>Discovered in Release:</b> 4.4-1.0.0.0
	<b>Fixed in Release:</b> 4.5-1.0.1.0
1571977	<b>Description:</b> Fixed an issue of when the same CQ is connected to some QPs with SRQ and some without, wrong <i>wr_id</i> might be reported by <i>ibv_poll_cq</i> .
	<b>Keywords:</b> libmlx5, wr_id

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.4-1.0.0.0
	<b>Fixed in Release:</b> 4.5-1.0.1.0
1380135	<b>Description:</b> Fixed the issue where IB port link used to flap due to MAD heartbeat response delay when using new CQ API.
	<b>Keywords:</b> IB port link, CQ API, MAD heartbeat
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.5-1.0.1.0
1498931	<b>Description:</b> Fixed the issue where establishing TCP connection took too long due to failure of SA PathRecord query callback handler.
	<b>Keywords:</b> TCP, SA PathRecord
	<b>Discovered in Release:</b> 4.4-1.0.0.0
	<b>Fixed in Release:</b> 4.5-1.0.1.0
1514096	<b>Description:</b> Fixed the issue where lack of high order allocations caused driver load failure. All high order allocations are now changed to order-0 allocations.
	<b>Keywords:</b> mlx5, high order allocation
	<b>Discovered in Release:</b> 4.0-2.0.2.0
	<b>Fixed in Release:</b> 4.5-1.0.1.0
1524932	<b>Description:</b> Fixed a backport issue on some OSs, such as RHEL v7.x, where mlx5 driver would support <i>ip link set DEVICE vf NUM rate TXRATE</i> old command, instead of <i>ip link set DEVICE vf NUM max_tx_rate TXRATE min_tx_rate TXRATE</i> new command.
	<b>Keywords:</b> mlx5 driver
	<b>Discovered in Release:</b> 4.0-2.0.2.0
	<b>Fixed in Release:</b> 4.5-1.0.1.0
1498585	<b>Description:</b> Fixed the issue of when performing configuration changes, mlx5e counters values were reset.

Internal Reference Number	Description
	<p><b>Keywords:</b> Ethernet counters</p> <p><b>Discovered in Release:</b> 4.0-2.0.2.0</p> <p><b>Fixed in Release:</b> 4.5-1.0.1.0</p>
1484603	<p><b>Description:</b> Fixed the issue of when using <code>ibv_exp_cqe_ts_to_ns</code> verb to convert a packet's hardware timestamp to UTC time in nanoseconds, the result may appear backwards compared to the converted time of a previous packet.</p> <p><b>Keywords:</b> libibverbs</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.5-1.0.1.0</p>
1425027	<p><b>Description:</b> Fixed the issue where attempting to establish a RoCE connection on the default GID or on IPv6 link-local address might have failed when two or more netdevices that belong to HCA ports were slaves under a bonding master. This might also have resulted in the following error message in the kernel log: "<code>__ib_cache_gid_add: unable to add gid fe80:0000:0000:0000:f652:14ff:fe46:7391 error=-28</code>".</p> <p><b>Keywords:</b> RoCE, bonding</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.5-1.0.1.0</p>
1480206	<p><b>Description:</b> Modified <code>mlx5_ib</code> SRQs behavior. Now the SRQs are allocated to "order 1" pages instead of contiguous ones to lower the probability of out-of-memory scenarios.</p> <p><b>Keywords:</b> SRQ, <code>mlx5_ib</code></p> <p><b>Discovered in Release:</b> 4.2-1.2.0.0</p> <p><b>Fixed in Release:</b> 4.4-2.0.7.0</p>
1363375	<p><b>Description:</b> Modified <code>mlx5_ib</code> QPs behavior. Now the QPs are allocated to "order 1" pages instead of contiguous ones to lower the probability of out-of-memory scenarios.</p>

Internal Reference Number	Description
	<b>Keywords:</b> IPoIB, mlx5_ib
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.4-2.0.7.0
1332080	<b>Description:</b> Modified mlx4_ib QPs behavior. Now the QPs are allocated to "order 1" pages instead of contiguous ones to lower the probability of out-of-memory scenarios.
	<b>Keywords:</b> IPoIB, mlx4_ib
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.4-1.0.0.0
	<b>Description:</b> Added support for multi-host connection on mstflint's mstfwreset.
	<b>Keywords:</b> mstfwreset, mstflint, MFT, multi-host
1412468	<b>Discovered in Release:</b> 4.3-1.0.1.0
	<b>Fixed in Release:</b> 4.4-1.0.0.0
	<b>Description:</b> Removed the following prints on server shutdown: mlx5_core 0005:81:00.1: mlx5_enter_error_state:96:(pid1): start mlx5_core 0005:81:00.1: mlx5_enter_error_state:109:(pid1): end
1423319	<b>Keywords:</b> mlx5, fast shutdown
	<b>Discovered in Release:</b> 4.3-1.0.1.0
	<b>Fixed in Release:</b> 4.4-1.0.0.0
1433092	<b>Description:</b> Fixed an issue of when querying for IBV_EXP_VALUES_HW_CLOCK_NS (using ibv_exp_query_values function) without querying for IBV_EXP_VALUES_HW_CLOCK, 0 value was returned.
	<b>Keywords:</b> mlx5, CQE time-stamping
	<b>Discovered in Release:</b> 4.3-1.0.1.0
	<b>Fixed in Release:</b> 4.4-1.0.0.0

Internal Reference Number	Description
1318251	<b>Description:</b> Fixed the issue of when bringing mlx4/mlx5 devices up or down, a call trace in <i>nvme_rdma_remove_one</i> or <i>nvmet_rdma_remove_one</i> may occur.
	<b>Keywords:</b> NVMeoF, mlx4, mlx5, call trace
	<b>Discovered in Release:</b> 4.3-1.0.1.0
	<b>Fixed in Release:</b> 4.4-1.0.0.0
1181815	<b>Description:</b> Fixed an issue where 4K UD packets were dropped when working with 4K MTU on mlx4 devices.
	<b>Keywords:</b> mlx4, 4K MTU, UD
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0
1247458	<b>Description:</b> Added support for VLAN Tag (VST) creation on RedHat v7.4 with new iproute2 packages (iptool).
	<b>Keywords:</b> SR-IOV, VST, RedHat
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0
1229554	<b>Description:</b> Enabled RDMA CM to honor incoming requests coming from ports of different devices.
	<b>Keywords:</b> RDMA CM
	<b>Discovered in Release:</b> 4.2-1.0.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0
1262257	<b>Description:</b> Fixed an issue where sending Work Requests (WRs) with multiple entries where the first entry is less than 18 bytes used to fail.
	<b>Keywords:</b> ConnectX-5; libibverbs; Raw QP
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0

Internal Reference Number	Description
1249358/1261023	<b>Description:</b> Fixed the issue of when the interface was down, ethtool counters ceased to increase. As a result, RoCE traffic counters were not always counted.
	<b>Keywords:</b> Ethtool counters, mlx5
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0
1244509	<b>Description:</b> Fixed compilation errors of MLNX_OFED over kernel when CONFIG_PTP_1588_CLOCK parameter was not set.
	<b>Keywords:</b> PTP, mlx5e
	<b>Discovered in Release:</b> 4.2-1.2.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0
1266802	<b>Description:</b> Fixed an issue where the system used to hang when trying to allocate multiple device memory buffers from different processes simultaneously.
	<b>Keywords:</b> Device memory programming
	<b>Discovered in Release:</b> 4.2-1.0.0.0
	<b>Fixed in Release:</b> 4.3-1.0.1.0
1120424	<b>Description:</b> Fixed incorrect SGE number of RSS QP.
	<b>Keywords:</b> RSS, SGE
	<b>Discovered in Release:</b> 4.1-1.0.2.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1078887	<b>Description:</b> Fixed an issue where post_list and CQ_mod features in perftest did not function when running the <i>--run_infinitely</i> flag.
	<b>Keywords:</b> perftest, --run_infinitely
	<b>Discovered in Release:</b> 4.2-1.0.1.0
	<b>Fixed in Release:</b> 4.2-1.2.0.0

Internal Reference Number	Description
1186260	<p><b>Description:</b> Fixed the issue where CNP counters exposed under <code>/sys/class/infiniband/mlx5_bond_0/ports/1/hw_counters/</code> did not aggregate both physical functions when working in RoCE LAG mode.</p>
	<p><b>Keywords:</b> RoCE, LAG, ECN, Congestion Counters</p>
	<p><b>Discovered in Release:</b> 4.2-1.0.1.0</p>
	<p><b>Fixed in Release:</b> 4.2-1.2.0.0</p>
1178129	<p><b>Description:</b> Fixed an issue that prevented Windows virtual machines running over MLNX_OFED Linux hypervisors from operating ConnectX-3 IB ports. When such failures occurred, the following message (or similar) appeared in the Linux HV message log when users attempted to start up a Windows VM running a ConnectX-3 VF: "mlx4_core 0000:81:00.0: vhcr command 0x1a slave:1 in_param 0x793000 in_mod=0x210 op_mod=0x0 failed with error:0, status -22"</p>
	<p><b>Keywords:</b> SR-IOV, RDMA, VM, KVM, Windows</p>
	<p><b>Discovered in Release:</b> 4.2-1.0.1.0</p>
	<p><b>Fixed in Release:</b> 4.2-1.2.0.0</p>
1192374	<p><b>Description:</b> Fixed wrong calculation of <code>max_device_ctx</code> capability in ConnectX-4, ConnectX-4 Lx, and ConnectX-5 HCAs.</p>
	<p><b>Keywords:</b> <code>ibv_exp_query_device</code>, <code>max_device_ctx</code> <code>mlx5</code></p>
	<p><b>Discovered in Release:</b> 4.2-1.0.1.0</p>
	<p><b>Fixed in Release:</b> 4.2-1.2.0.0</p>
1084791	<p><b>Description:</b> Fixed the issue where occasionally, after reboot, <code>rpm</code> commands used to fail and create a core file, with messages such as "Bus error (core dumped)", causing the <code>openibd</code> service to fail to start.</p>
	<p><b>Keywords:</b> <code>rpm</code>, <code>openibd</code></p>
	<p><b>Discovered in Release:</b> 3.4-2.0.0.0</p>
	<p><b>Fixed in Release:</b> 4.2-1.0.0.0</p>



Internal Reference Number	Description
960642/960653	<b>Description:</b> Added support for <i>min_tx_rate</i> and <i>max_tx_rate</i> limit per virtual function ConnectX-5 and ConnectX-5 Ex adapter cards.
	<b>Keywords:</b> SR-IOV, mlx5
	<b>Discovered in Release:</b> 4.0-1.0.1.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
866072/869183	<b>Description:</b> Fixed the issue where RoCE v2 multicast traffic using RDMA-CM with IPv4 address was not received.
	<b>Keywords:</b> RoCE
	<b>Discovered in Release:</b> 3.4-1.0.0.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1163835	<b>Description:</b> Fixed an issue where <i>ethtool -P</i> output was 00:00:00:00:00:00 when using old kernels.
	<b>Keywords:</b> ethtool, Permanent MAC address, mlx4, mlx5
	<b>Discovered in Release:</b> 4.0-2.0.0.1
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1067158	<b>Description:</b> Replaced a few "GPL only" legacy libibverbs functions with upstream implementation that conforms with libibverbs GPL/BSD dual license model.
	<b>Keywords:</b> libibverbs, license
	<b>Discovered in Release:</b> 4.1-1.0.2.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1119377	<b>Description:</b> Fixed an issue where ACCESS_REG command failure used to appear upon RoCE Multihost driver restart in dmesg. Such an error message looked as follows: <i>mlx5_core 0000:01:00.0: mlx5_cmd_check:705:(pid 20037): ACCESS_REG(0x805) op_mod(0x0) failed, status bad parameter(0x3), syndrome (0x15c356)</i>
	<b>Keywords:</b> RoCE, multihost, mlx5

Internal Reference Number	Description
	<b>Discovered in Release:</b> 4.1-1.0.2.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1122937	<b>Description:</b> Fixed an issue where concurrent client requests got corrupted when working in persistent server mode due to a race condition on the server side.
	<b>Keywords:</b> librdmacm, rping
	<b>Discovered in Release:</b> 4.1-1.0.2.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1102158	<b>Description:</b> Fixed an issue where client side did not exit gracefully in RTT mode when the server side was not reachable.
	<b>Keywords:</b> librdmacm, rping
	<b>Discovered in Release:</b> 4.1-1.0.2.0
	<b>Fixed in Release:</b> 4.2-1.0.0.0
1038933	<b>Description:</b> Fixed a backport issue where IPv6 procedures were called while they were not supported in the underlying kernel.
	<b>Keywords:</b> iw_cm
	<b>Discovered in Release:</b> 4.0-2.0.0.1
	<b>Fixed in Release:</b> 4.1-1.0.2.0
1064722	<b>Description:</b> Added log debug prints when changing HW configuration via DCB. To enable log debug prints, run: <i>ethtool -s &lt;devname&gt; msglvl hw on/off</i>
	<b>Keywords:</b> DCB, msglvl
	<b>Discovered in Release:</b> 4.0-2.0.0.1
	<b>Fixed in Release:</b> 4.1-1.0.2.0
1013076	<b>Description:</b> Fixed the issue where reassembly of packets larger than 64k might have failed when ipfrag threshold was low. This issue was present only on RHEL 6.3, 6.4, 6.5, and Ubuntu 12.04.

Internal Reference Number	Description
	<p>This packet drop could be seen from the netstat tool, indicated by the “packet reassembles failed” counter.</p> <p><b>Keywords:</b> IPoIB, Packet Fragmentation</p> <p><b>Discovered in Release:</b> 4.0-2.0.0.1</p> <p><b>Fixed in Release:</b> 4.1-1.0.2.0</p>
1022251	<p><b>Description:</b> Fixed SKB memory leak issue that was introduced in kernel 4.11, and added warning messages to the Soft RoCE driver for easy detection of future SKB leaks.</p> <p><b>Keywords:</b> Soft RoCE</p> <p><b>Discovered in Release:</b> 4.0-2.0.0.1</p> <p><b>Fixed in Release:</b> 4.1-1.0.2.0</p>
1044546	<p><b>Description:</b> Fixed the issue where a kernel crash used to occur when Rxe device was coupled with a virtual (dummy) device.</p> <p><b>Keywords:</b> Soft RoCE</p> <p><b>Discovered in Release:</b> 4.0-2.0.0.1</p> <p><b>Fixed in Release:</b> 4.1-1.0.2.0</p>
1047617	<p><b>Description:</b> Fixed the issue where a race condition in the RoCE GID cache used to cause for the loss of IP-based GIDs.</p> <p><b>Keywords:</b> RoCE, GID</p> <p><b>Discovered in Release:</b> 4.0-2.0.0.1</p> <p><b>Fixed in Release:</b> 4.1-1.0.2.0</p>
1006768	<p><b>Description:</b> Fixed the issue where an rdma_cm connection between a client and a server that were on the same host was not possible when working over VLAN interfaces.</p> <p><b>Keywords:</b> RDMACM</p> <p><b>Discovered in Release:</b> 4.0-2.0.0.1</p> <p><b>Fixed in Release:</b> 4.1-1.0.2.0</p>

Internal Reference Number	Description
801807	<b>Description:</b> Fixed an issue where RDMACM connection used to fail upon high connection rate accompanied with the error message: <i>RDMA_CM_EVENT_UNREACHABLE</i> .
	<b>Keywords:</b> RDMACM
	<b>Discovered in Release:</b> 3.0-2.0.1
	<b>Fixed in Release:</b> 4.1-1.0.2.0
869768	<b>Description:</b> Fixed the issue where SR-IOV was not supported in systems with a page size greater than 16KB.
	<b>Keywords:</b> SR-IOV, mlx5, PPC
	<b>Discovered in Release:</b> 4.0-2.0.0.1
	<b>Fixed in Release:</b> 4.1-1.0.2.0
1155972	<b>Description:</b> Fixed mlx4 kernel crash upon server shutdown due to NULL pointer dereference.
	<b>Keywords:</b> mlx4, shutdown
	<b>Discovered in Release:</b> 3.3-1.0.4.0
	<b>Fixed in Release:</b> 4.0-2.0.0.1
919545	<b>Description:</b> Fixed the issue of when the Kernel becomes out of memory upon driver start, it could crash on SLES 12 SP2.
	<b>Keywords:</b> mlx_5 Eth Driver
	<b>Discovered in Release:</b> 3.4-2.0.0.0
	<b>Fixed in Release:</b> 4.0-2.0.0.1
970668	<b>Description:</b> Fixed the issue where very high stress on DC QP transport might have triggered NMI messages on specific servers.
	<b>Keywords:</b> mlx5 Driver
	<b>Discovered in Release:</b> 4.0-1.0.1.0
	<b>Fixed in Release:</b> 4.0-2.0.0.1

Internal Reference Number	Description
966134	<b>Description:</b> Allowed Ethernet VFs to open Raw Ethernet QPs even if RoCE is not supported for the VF.
	<b>Keywords:</b> mlx4_ib
	<b>Discovered in Release:</b> 3.0-1.0.1
	<b>Fixed in Release:</b> 4.0-2.0.0.1
864063	<b>Description:</b> Fixed the issue of when Spoof-check may have been turned on for MAC address 00:00:00:00:00:00.
	<b>Keywords:</b> mlx4
	<b>Discovered in Release:</b> 3.4-1.0.0.0
	<b>Fixed in Release:</b> 4.0-2.0.0.1
869209	<b>Description:</b> Fixed an issue that caused TCP packets to be received in an out of order manner when Large Receive Offload (LRO) is on.
	<b>Keywords:</b> mlx5_en
	<b>Discovered in Release:</b> 3.3-1.0.0.0
	<b>Fixed in Release:</b> 4.0-2.0.0.1
913319	<b>Description:</b> Fixed the issue of low performance when creating many address handles.
	<b>Keywords:</b> libibverbs
	<b>Discovered in Release:</b> 3.4-1.0.0.0
	<b>Fixed in Release:</b> 4.0-1.0.1.0
912897	<b>Description:</b> Added debug prints to <i>ib_umem_get</i> function to fix lack of error indication when this function fails.
	<b>Keywords:</b> InfiniBand
	<b>Discovered in Release:</b> 3.4-1.0.0.0
	<b>Fixed in Release:</b> 4.0-1.0.1.0
945887	<b>Description:</b> [ConnectX-3] Fixed the issue where multicast traffic over Raw Ethernet QP on virtual functions were received on the

Internal Reference Number	Description
	<p>same QP (loopback).</p> <p><b>Keywords:</b> SR-IOV</p> <p><b>Discovered in Release:</b> 3.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.0-1.0.1.0</p>
920292	<p><b>Description:</b> Fixed three issues in libmlx5 that were found by NVIDIA in the patches that are part of MLNX_OFED v3.4:</p> <ol style="list-style-type: none"> <li>1. mlx5_exp_peer_commit_qp returns number of entries = 4 instead of 3.</li> <li>2. Peer capability check is wrong - should fail the check when there is neither NOR nor GEQ support.</li> <li>3. Missing break in mlx5_exp_peer_peek_cq. There is now fallthrough in the IBV_EXP_PEER_PEEK_ABSOLUTE case.</li> </ol> <p><b>Keywords:</b> libmlx5</p> <p><b>Discovered in Release:</b> 3.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.0-1.0.1.0</p>
890285	<p><b>Description:</b> Fixed the issue where memory allocation for CQ buffers used to fail when increasing the RX ring size.</p> <p><b>Keywords:</b> mlx5_core</p> <p><b>Discovered in Release:</b> 3.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.0-1.0.1.0</p>
867094	<p><b>Description:</b> Fixed the issue where MLNX_OFED used to fail to load on 4K page Arm architecture.</p> <p><b>Keywords:</b> Arm</p> <p><b>Discovered in Release:</b> 3.4-1.0.0.0</p> <p><b>Fixed in Release:</b> 4.0-1.0.1.0</p>

# Known Issues

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

For the list of old known issues, please refer to Mellanox OFED Archived Known Issues file at:

[http://www.mellanox.com/pdf/prod\\_software/MLNX\\_OFED\\_Archived\\_Known\\_Issues.pdf](http://www.mellanox.com/pdf/prod_software/MLNX_OFED_Archived_Known_Issues.pdf)

Internal Ref. Number	Issue
2894838	<b>Description:</b> Running 'ip link show' command over RHEL8.5 using ConnectX-3 with VFs will print "Truncated VFs" to the screen.
	<b>Workaround:</b> Use the following OFED IP link command: <code>/opt/mellanox/iproute2/sbin/ip link show</code>
	<b>Keywords:</b> IP Link, Virtual Functions, ConnectX-3
	<b>Discovered in Release:</b> 4.9-4.1.7.0
2793596	<b>Description:</b> On Sles15Sp3, MFT restart does not work.
	<b>Workaround:</b> Install MFT manually from <a href="https://www.mellanox.com/products/adapter-software/firmware-tools">https://www.mellanox.com/products/adapter-software/firmware-tools</a> .
	<b>Keywords:</b> MFT
	<b>Discovered in Release:</b> 4.9-4.0.8.0
2794326	<b>Description:</b> When upgrading MLNX_OFED from 4.9-4 to 5.4-2 GA using Yum installation, the installation fails due to ibutils.
	<b>Workaround:</b> Before the upgrade, remove ibutils manually (and the metapackage with it) using the following command: <code>yum remove ibutils</code>
	<b>Keywords:</b> Installation, ibutils
	<b>Discovered in Release:</b> 4.9-4.0.8.0
2753944	<b>Description:</b> On rare occasion, registering a device ( <code>ib_register_device()</code> ) and loading modules in parallel in this case ( <code>ib_cm</code> ), a racing condition may occur which would stop <code>ib_cm</code> from loading properly.

Internal Ref. Number	Issue
	<p><b>Workaround:</b> Add modprobe.d rules to force the ib_cm driver to load before the mlx4_ib and mlx5_ib drivers:  install mlx4_ib { /sbin/modprobe ib_cm; /sbin/modprobe -ignore-install mlx4_ib \$CMDLINE_OPTS; }  install mlx5_ib { /sbin/modprobe ib_cm; /sbin/modprobe —ignore-install mlx5_ib \$CMDLINE_OPTS; }</p> <p><b>Keywords:</b> ib_core, Racing Condition</p> <p><b>Discovered in Release:</b> 4.9-4.0.8.0</p>
2636998	<p><b>Description:</b> When using Debian or Ubuntu operating systems, installing MLNX_OFED with mlnxofedinstall and then proceeding to upgrade with a package manager (apt), the mlnx-rdma-core-dkms package remains installed and fails to rebuild.</p> <p><b>Workaround:</b> Before upgrade, remove mlnx-rdma-rxe-dkms: dpkg --purge mlnx-rdma-rxe-dkms</p> <p><b>Keywords:</b> Upgrade, Debian, Ubuntu, mlnx-rdma-core-dkms</p> <p><b>Discovered in Release:</b> 4.9-3.1.5.0</p>
2338121	<p><b>Description:</b> UCX will not work while running with upstream-libs if librdmacm is not installed.</p> <p><b>Workaround:</b> Install rdmacm or disable VMC (-x HCOLL_MCAST=^vmc).</p> <p><b>Keywords:</b> RDMA</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2440042	<p><b>Description:</b> Using ODP on specific hardware may cause intermittent failures (issue only reproduced on IBM POWER8 S822LC).</p> <p><b>Workaround:</b> If the program fail is seen, disable ODP. Or, to use ODP with ConnectX-4 and above, it is recommended to use MLNX_OFED version 5.2 and above.</p> <p><b>Keywords:</b> ar_mgr; dump_pr; upgrade; installation</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2432304	<p><b>Description:</b> ar_mgr and dump_pr plugin versions are not updated when upgrading the MLNX_OFED version.</p>



Internal Ref. Number	Issue
	<p><b>Workaround:</b> Prior to upgrading your MLNX_OFED version, make sure to uninstall ar_mgr and dump_pr subnet manager plugins. For example, on Ubuntu systems, run:  <pre>dpkg --remove mlnx-ofed-all ar-mgr dump-pr</pre></p> <p><b>Keywords:</b> ar_mgr; dump_pr; upgrade; installation</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2339456	<p><b>Description:</b> OFED installation requires the --add-kernel-support flag on some of the Errata kernels of RedHat 7.8.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Installation, Errata, RedHat, OS</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2328653	<p><b>Description:</b> Dependency between qemu and libibverbs may cause qemu failures after OFED installation on Ubuntu v20.04 or SLES 15.2 KVM systems.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> qemu, libibverbs, installation, OS, Ubuntu, SLES, SUSE</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2345669	<p><b>Description:</b> AliOS installation with add-kernel-support may require the installation of additional packages.</p> <p><b>Workaround:</b> Install the required packages.</p> <p><b>Keywords:</b> Installation</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2312063	<p><b>Description:</b> MKEY_BY_NAME is not supported.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> MKEY_BY_NAME</p> <p><b>Discovered in Release:</b> 4.9-2.2.4.0</p>
2046307	<p><b>Description:</b> Excessive toggling between modes (Connected and Datagram) and interface up and down may cause a crash.</p>

Internal Ref. Number	Issue
	<p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> System crash, mode change</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
1550266	<p><b>Description:</b> XDP is not supported over ConnectX-3 and ConnectX-3 Pro adapter cards.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> XDP, ConnectX-3</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2117822	<p><b>Description:</b> On ConnectX-3 and ConnectX-3 Pro adapter cards, no traffic runs between VLANs of any type over VLAN of type ctag (protocol 802.1Q).</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> ConnectX-3, VLAN</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2142218	<p><b>Description:</b> On ConnectX-3 and ConnectX-3 Pro adapter cards, driver might hang when found under the following conditions, collectively:</p> <ul style="list-style-type: none"> <li>• OS kernel is older than 4.10</li> <li>• Interface is down</li> <li>• CONFIG_NET_RX_BUSY_POLL parameter is set</li> <li>• netdev_ops.ndo_busy_poll is defined</li> </ul> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> ConnectX-3</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2156645	<p><b>Description:</b> MLNX_LIBS provider packages, such as libmlx5, cannot be installed simultaneously with ibverbs-providers distribution package when working with Ubuntu and Debian OSs.</p> <p><b>Workaround:</b> Before installing MLNX_OFED of type MLNX_LIBS, make sure that ibverbs-providers package is not installed.</p>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> MLNX_LIBS, libmlx5, ibverbs-providers, Debian, Ubuntu</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2105447	<p><b>Description:</b> hns_roce warning messages will appear in the dmesg after reboot on Euler2 SP3 OSs.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> hns_roce, dmesg, Euler</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2110321	<p><b>Description:</b> Multiple driver restarts may cause IPoIB soft lockup.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Driver restart, IPoIB</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2112251	<p><b>Description:</b> On kernels 4.10-4.14, when Geneve tunnel's remote endpoint is defined using IPv6, packets larger than MTU are not fragmented, resulting in no traffic sent.</p> <p><b>Workaround:</b> Define geneve tunnel's remote endpoint using IPv4.</p> <p><b>Keywords:</b> Kernel, Geneve, IPv4, IPv6, MTU, fragmentation</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2119210	<p><b>Description:</b> Multiple driver restarts may cause a stress and result in mlx5 commands check error message in the log.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Driver restart, syndrome, error message</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2111349	<p><b>Description:</b> Ethtool --show-fec/--get-fec are not supported over ConnectX-6 and ConnectX-6 Dx adapter cards.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Ethtool, ConnectX-6 Dx</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>

Internal Ref. Number	Issue
2119984	<b>Description:</b> IPsec crypto offloads does not work when ESN is enabled.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> IPsec, ESN
	<b>Discovered in Release:</b> 4.9-0.1.7.0
2102902	<b>Description:</b> A kernel panic may occur over RH8.0-4.18.0-80.el8.x86_64 OS when opening kTLS offload connection due to a bug in kernel TLS stack.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> TLS offload, mlx5e
	<b>Discovered in Release:</b> 4.9-0.1.7.0
2111534	<b>Description:</b> A Kernel panic may occur over Ubuntu19.04-5.0.0-38-generic OS when opening kTLS offload connection due to a bug in the Kernel TLS stack.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> TLS offload, mlx5e
	<b>Discovered in Release:</b> 4.9-0.1.7.0
2117845	<b>Description:</b> Relaxed ordering memory regions are not supported when working with CAPI. Registering memory region with relaxed ordering while CAPI enabled will result in a registration failure.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Relaxed ordering, memory region, MR, CAPI
	<b>Discovered in Release:</b> 4.9-0.1.7.0
2083942	<b>Description:</b> The content of file /sys/class/net/ /statistics/multicast may be out of date and may display values lower than the real values.
	<b>Workaround:</b> Run ethtool -S <NETIF> to show the actual multicast counters and to update the content of file /sys/class/net/ /statistics/multicast.
	<b>Keywords:</b> Multicast counters
	<b>Discovered in Release:</b> 4.9-0.1.7.0

Internal Ref. Number	Issue
2035950	<p><b>Description:</b> An internal error might take place in the firmware when performing any of the following in VF LAG mode, when at least one VF of either PF is still bound/attached to a VM.</p> <ol style="list-style-type: none"> <li>1. Removing PF from the bond (using ifdown, ip link or any other function)</li> <li>2. Attempting to disable SR-IOV</li> </ol>
	<p><b>Workaround:</b> N/A</p>
	<p><b>Keywords:</b> VF LAG, binding, firmware, FW, PF, SR-IOV</p>
	<p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2094176	<p><b>Description:</b> When running in a large scale in VF-LAG mode, bandwidth may be unstable.</p>
	<p><b>Workaround:</b> N/A</p>
	<p><b>Keywords:</b> VF LAG</p>
	<p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2044544	<p><b>Description:</b> When working with OSs with Kernel v4.10, bonding module does not allow setting MTUs larger than 1500 on a bonding interface.</p>
	<p><b>Workaround:</b> Upgrade your Kernel version to v4.11 or above.</p>
	<p><b>Keywords:</b> Bonding, MTU, Kernel</p>
	<p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
1882932	<p><b>Description:</b> Libibverbs dependencies are removed during OFED installation, requiring manual installation of libraries that OFED does not reinstall.</p>
	<p><b>Workaround:</b> Manually install missing packages.</p>
	<p><b>Keywords:</b> libibverbs, installation</p>
	<p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2093746	<p><b>Description:</b> Devlink health dumps are not supported on kernels lower than v5.3.</p>
	<p><b>Workaround:</b> N/A</p>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> Devlink, health report, dump</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2020260	<p><b>Description:</b> When changing the Trust mode to DSCP, there is an interval between the change taking effect in the hardware and updating the inline mode of the SQ in the driver. If any traffic is transmitted during this interval, the driver will not inline enough headers, resulting in a CQE error in the NIC.</p> <p><b>Workaround:</b> Set the interface down, change the trust mode, then bring the interface back up.</p> <pre data-bbox="383 730 1461 919">ip link set eth0 down mlnx_qos -i eth0 --trust dscp ip link set eth0 up</pre> <p><b>Keywords:</b> DSCP, inline, SQ, CQE</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2083427	<p><b>Description:</b> For kernels with connection tracking support, neigh update events are not supported, requiring users to have static ARPs to work with OVS and VxLAN.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> VxLAN, VF LAG, neigh, ARP</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2043739	<p><b>Description:</b> Userspace RoCE UD QPs are not supported over distributions such as SLES11 SP4 and RedHat 6.10 for which the netlink 3 libraries (libnl-3 and libnl-route3) are not available.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RoCE UD, QP, SLES, RedHat, RHEL, netlink</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
2067746	<p><b>Description:</b> When attaching a second slave to a bond, some bond interface GIDs might disappear.</p> <p><b>Workaround:</b> Re-create and re-configure the bond device.</p>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> Bond, GID</p> <p><b>Discovered in Release:</b> 4.9-0.1.7.0</p>
-	<p><b>Description:</b> The argparse module is installed by default in Python versions <math>\geq 2.7</math> and <math>\geq 3.2</math>. In case an older Python version is used, the argparse module is not installed by default.</p> <p><b>Workaround:</b> Install the argparse module manually.</p> <p><b>Keywords:</b> Python, MFT, argparse, installation</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p>
1979834	<p><b>Description:</b> When running MLNX_OFED on Kernel 4.10 with ConnectX-3/ConnectX-3 Pro NICs, deleting VxLAN may result in a crash.</p> <p><b>Workaround:</b> Upgrade the Kernel version to v4.14 to avoid the crash.</p> <p><b>Keywords:</b> Kernel, OS, ConnectX-3, VxLAN</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p>
1973238	<p><b>Description:</b> <code>ib_core</code> unload may fail on Ubuntu 18.04.2 OS with the following error message: "Module <code>ib_core</code> is in use"</p> <p><b>Workaround:</b> Stop <code>ibacm.socket</code> using the following commands: <code>systemctl stop ibacm.socket</code> <code>systemctl disable ibacm.socket</code></p> <p><b>Keywords:</b> <code>ib_core</code>, Ubuntu, <code>ibacm</code></p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p>
1970429	<p><b>Description:</b> With HW offloading in SR-IOV SwitchDev mode, the fragmented ICMP echo request/reply packets (with length larger than MTU) do not function properly. The correct behavior is for the fragments to miss the offloading flow and go to the slow path. However, the current behavior is as follows.</p> <ul style="list-style-type: none"> <li>• Ingress (to the VM): All echo request fragments miss the corresponding offloading flow, but all echo reply fragments hit the corresponding offloading flow</li> <li>• Egress (from the VM): The first fragment still hits the corresponding offloading flow, and the subsequent fragments miss the</li> </ul>

Internal Ref. Number	Issue
	<p>corresponding offloading flow</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> HW offloading, SR-IOV, SwitchDev, ICMP, VM, virtualization</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p>
1969580	<p><b>Description:</b> RHEL 6.10 OS is not supported in SR-IOV mode.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RHEL, RedHat, OS, operating system, SR-IOV, virtualization</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p>
1919335	<p><b>Description:</b> On SLES 11 SP4, RedHat 6.9 and 6.10 OSs, on hosts where OpenSM is running, the low-level driver's internal error reset flow will cause a kernel crash when OpenSM is killed (after the reset occurs). This is due to a bug in these kernels where opening the umad device (by OpenSM) does not take a reference count on the underlying device.</p> <p><b>Workaround:</b> Run OpenSM on a host with a more recent Kernel.</p> <p><b>Keywords:</b> SLES, RedHat, CR-Dump, OpenSM</p> <p><b>Discovered in Release:</b> 4.7-3.2.9.0</p>
1893464	<p><b>Description:</b> ibacm is not tested with MLNX_OFED or its components.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> ibacm, component</p> <p><b>Discovered in Release:</b> 4.7-1.0.0.1</p>
1921981	<p><b>Description:</b> On Ubuntu, Debian and RedHat 8 and above OSS, parsing the mfa2 file using the mstarchive might result in a segmentation fault.</p> <p><b>Workaround:</b> Use mlxarchive to parse the mfa2 file instead.</p> <p><b>Keywords:</b> MFT, mfa2, mstarchive, mlxarchive, Ubuntu, Debian, RedHat, operating system</p> <p><b>Discovered in Release:</b> 4.7-1.0.0.1</p>



Internal Ref. Number	Issue
1921799	<b>Description:</b> MLNX_OFED installation over SLES15 SP1 ARM OSs fails unless --add-kernel-support flag is added to the installation command.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> SLES, installation
	<b>Discovered in Release:</b> 4.7-1.0.0.1
1840288	<b>Description:</b> MLNX_OFED does not support XDP features on RedHat 7 OS, despite the declared support by RedHat.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> XDP, RedHat
	<b>Discovered in Release:</b> 4.7-1.0.0.1
1919335	<b>Description:</b> On SLES 11 SP4, RedHat 6.9 and 6.10 OSs, bringing the OpenSM down after CR-Dump results in a panic.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> SLES, RedHat, CR-Dump, OpenSM
	<b>Discovered in Release:</b> 4.7-1.0.0.1
1821235	<b>Description:</b> When using mlx5dv_dr API for flow creation, for flows which execute the "encapsulation" action or "push vlan" action, metadata C registers will be reset to zero.
	<b>Workaround:</b> Use the both actions at the end of the flow process.
	<b>Keywords:</b> Flow steering
	<b>Discovered in Release:</b> 4.7-1.0.0.1
1911130	<b>Description:</b> When Offloaded Traffic Sniffer feature is on, the usage of "all default" flow steering rule could cause a deadlock.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Offloaded Traffic Sniffer, steering, deadlock
	<b>Discovered in Release:</b> 4.7-1.0.0.1
1897199	<b>Description:</b> When using the RDMA statistics feature and attempting to unbind a QP from a counter, not including the counter-id as an argument

Internal Ref. Number	Issue
	<p>in the CLI will result in a segmentation fault.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RDMA, QP, segfault, unbinding</p> <p><b>Discovered in Release:</b> 4.7-1.0.0.1</p>
1869219	<p><b>Description:</b> On Fedora 27 OSs, reboot/shutdown operations may fail after uninstalling the MLNX_OFED package.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Fedora 27, uninstall, reboot, shutdown</p> <p><b>Discovered in Release:</b> 4.7-1.0.0.1</p>
1892663	<p><b>Description:</b> mlnx_tune script does not support python3 interpreter.</p> <p><b>Workaround:</b> Run mlnx_tune with python2 interpreter only.</p> <p><b>Keywords:</b> mlnx_tune, python3, python2</p> <p><b>Discovered in Release:</b> 4.7-1.0.0.1</p>
1341833	<p><b>Description:</b> On CoreOS, assigning a static IP address to PKeys using ifcfg configuration file option fails after restarting the driver.</p> <p><b>Workaround:</b> Manually run "ifdown" and then "ifup".</p> <p><b>Keywords:</b> CoresOS, PKey, restart_driver</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1504785	<p><b>Description:</b> A lost interrupt issue in pass-through virtual machines may prevent the driver from loading, followed by printing managed pages errors to the dmesg.</p> <p><b>Workaround:</b> Restart the driver.</p> <p><b>Keywords:</b> VM, virtual machine</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1630228	<p><b>Description:</b> Tunnel stateless offloads are wrongly forbidden for E-Switch manager function.</p> <p><b>Workaround:</b> Set the stateless offloads cap to be permanently '1'.</p>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> Stateless offloads cap</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1764415	<p><b>Description:</b> Unbinding PFs on LAG devices results in a "Failed to modify QP to RESET" error message.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RoCE LAG, unbind, PF, RDMA</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1769208	<p><b>Description:</b> Contrary to the standard DSCP mode setting procedure in SR-IOV mode, now, in order for this configuration to take effect, the DSCP trust mode has to be set before the VF is created, and not the other way around.</p> <p><b>Workaround:</b> Make sure to set the DSCP trust mode before creating the VF.</p> <p><b>Keywords:</b> DSCP, trust mode, VF</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1779150	<p><b>Description:</b> Upgrading the MLNX_OFED version over SLES 15 SP0 and SP1 OSs on PPCLE platforms might fail due to an insert-kmp-default issue.</p> <p><b>Workaround:</b> Remove the insert-kmp-default package manually</p> <p><b>Keywords:</b> Installation, SLES, PPCLE</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1806565	<p><b>Description:</b> RoCE default GIDs v1 and v2 are derived from the MAC address of the corresponding netdevice's PCI function, and they resemble the IPv6 address. However, in systems where the IPv6 link local address generated does not depend on the MAC address, RoCEv2 default GID should not be used.</p> <p><b>Workaround:</b> Use RoCEv2 default GID.</p> <p><b>Keywords:</b> RoCE</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>

Internal Ref. Number	Issue
1834997	<b>Description:</b> When working with VF Lag while the bond device is in active-active mode, traffic on both physical ports may not reach line rate.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> VF LAG, bonding, bandwidth degradation, fairness
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1839907	<b>Description:</b> In mlx4 devices, enabling RX-FCS offload does not disable LRO, and vice-versa.
	<b>Workaround:</b> Disable the RX-FCS or LRO separately.
	<b>Keywords:</b> Frame Check Sequence (FCS), Large Receive Offload (LRO)
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1735161	<b>Description:</b> Innova cards do not support InfiniBand mode.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Innova, IB, InfiniBand
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1787667	<b>Description:</b> NVMe-oF driver of MLNX OFED v4.6-x.x.x.x does not function on SLES12 SP4 and SLES15 SP1 OSs, as they have a built-in NVME driver in the Linux image. Therefore, Mellanox NVME and NVMe-oF drivers cannot be loaded. For tracking purposes of this bug, see <a href="#">Bugzilla issue #1150850</a> and <a href="#">Bugzilla issue #1150846</a> .
	<b>Workaround:</b> Change the kernel configuration of NVMe-oF driver to be "=m" and recompile the kernel.
	<b>Keywords:</b> NVMe-oF, NVME, SLES
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1759593	<b>Description:</b> OFED installation on XenServer OSs requires using the -u flag.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Installation, XenServer, OS, operating system

Internal Ref. Number	Issue
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1753629	<b>Description:</b> A bonding bug found in Kernels 4.12 and 4.13 may cause a slave to become permanently stuck in BOND_LINK_FAIL state. As a result, the following message may appear in dmesg: bond: link status down for interface eth1, disabling it in 100 ms
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Bonding, slave
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1734102	<b>Description:</b> Ubuntu v16.04.05 and v16.04.05 OSs can only be used with Kernels of version 4.4.0-143 or below.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Ubuntu, Kernel, OS
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1712068	<b>Description:</b> Uninstalling MLNX_OFED automatically results in the uninstallation of several libraries that are included in the MLNX_OFED package, such as InfiniBand-related libraries.
	<b>Workaround:</b> If these libraries are required, reinstall them using the local package manager (yum/dnf).
	<b>Keywords:</b> MLNX_OFED libraries
	<b>Discovered in Release:</b> 4.6-1.0.1.1
-	<b>Description:</b> Due to changes in libraries, MFT v4.11.0 and below are not forward compatible with MLNX_OFED v4.6-1.0.0.0 and above. Therefore, with MLNX_OFED v4.6-1.0.0.0 and above, it is recommended to use MFT v4.12.0 and above.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> MFT compatible
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1730840	<b>Description:</b> On ConnectX-4 HCAs, GID index for RoCE v2 is inconsistent when toggling between enabled and disabled interface modes.

Internal Ref. Number	Issue
	<b>Workaround:</b> N/A
	<b>Keywords:</b> RoCE v2, GID
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1731005	<b>Description:</b> MLNX_OFED v4.6 YUM and Zypper installations fail on RHEL8.0, SLES15.0 and PPCLE OSs.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> YUM, Zypper, installation, RHEL, RedHat, SLES, PPCLE
1717428	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Description:</b> On kernels 4.10-4.14, MTUs larger than 1500 cannot be set for a GRE interface with any driver (IPv4 or IPv6).
	<b>Workaround:</b> Upgrade your kernel to any version higher than v4.14.
1748343	<b>Keywords:</b> Fedora 27, gretap, ip_gre, ip_tunnel, ip6_gre, ip6_tunnel
	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Description:</b> Driver reload takes several minutes when a large number of VFs exists.
1748537	<b>Workaround:</b> N/A
	<b>Keywords:</b> VF, SR-IOV
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1748537	<b>Description:</b> Cannot set max Tx rate for VFs from the ARM.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Host control, max Tx rate
1732940	<b>Discovered in Release:</b> 4.6-1.0.1.1
	<b>Description:</b> Software counters not working for representor net devices.
	<b>Workaround:</b> N/A
1732940	<b>Keywords:</b> mlx5, counters, representors
	<b>Discovered in Release:</b> 4.6-1.0.1.1

Internal Ref. Number	Issue
1733974	<b>Description:</b> Running heavy traffic (such as 'ping flood') while bringing up and down other mlx5 interfaces may result in "INFO: rcu_preempt dectected stalls on CPUS/tasks:" call traces.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> mlx5
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1731939	<b>Description:</b> Get/Set Forward Error Correction FEC configuration is not supported on ConnectX-6 HCAs with 200Gbps speed rate.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> Forward Error Correction, FEC, 200Gbps
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1715789	<b>Description:</b> Mellanox Firmware Tools (MFT) package is missing from Ubuntu v18.04.2 OS.
	<b>Workaround:</b> Manually install MFT.
	<b>Keywords:</b> MFT, Ubuntu, operating system
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1652864	<b>Description:</b> On ConnectX-3 and ConnectX-3 Pro HCAs, CR-Dump poll is not supported using sysfs commands.
	<b>Workaround:</b> If supported in your Kernel, use the devlink tool as an alternative to sysfs to achieve CR-Dump support.
	<b>Keywords:</b> mlx4, devlink, CR-Dump
	<b>Discovered in Release:</b> 4.6-1.0.1.1
1699031	<b>Description:</b> When attempting to destroy IPoIB bonding interface on PPCLE setups, a leak of resources might occur.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> IPoIB, bonding, PPCLE
	<b>Discovered in Release:</b> 4.6-1.0.1.1

Internal Ref. Number	Issue
-	<p><b>Description:</b> On ConnectX-6 HCAs and above, an attempt to configure advertisement (any bitmap) will result in advertising the whole capabilities.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> 200Gmbps, advertisement, Ethtool</p> <p><b>Discovered in Release:</b> 4.6-1.0.1.1</p>
1699289	<p><b>Description:</b> HW LRO feature is disabled OOB, which results in increased CPU utilization on the Receive side. On ConnectX-5 adapter cards and above, this causes a bandwidth drop for a few streams.</p> <p><b>Workaround:</b> Make sure to enable HW LRO in the driver:  ethtool -k &lt;intf&gt; lro  ethtool --set-priv-flag &lt;intf&gt; hw_lro on</p> <p><b>Keywords:</b> HW LRO, ConnectX-5 and above</p> <p><b>Discovered in Release:</b> 4.5-1.0.1.0</p>
1583487	<p><b>Description:</b> MPI package is not part of MLNX_OFED package in Fedora 28 OS.</p> <p><b>Workaround:</b> Manually install MPI package.</p> <p><b>Keywords:</b> MPI package, Fedora</p> <p><b>Discovered in Release:</b> 4.5-1.0.1.0</p>
1403313	<p><b>Description:</b> Attempting to allocate an excessive number of VFs per PF in operating systems with kernel versions below v4.15 might fail due to a known issue in the Kernel.</p> <p><b>Workaround:</b> Make sure to update the Kernel version to v4.15 or above.</p> <p><b>Keywords:</b> VF, PF, IOMMU, Kernel, OS</p> <p><b>Discovered in Release:</b> 4.5-1.0.1.0</p>
-	<p><b>Description:</b> NEO-Host is not supported on the following OSs:</p> <ul style="list-style-type: none"> <li>• SLES12 SP3</li> <li>• SLES12 SP4</li> <li>• SLES15</li> <li>• Fedora 28</li> </ul>



Internal Ref. Number	Issue
	<ul style="list-style-type: none"> <li>• RHEL7.1</li> <li>• RHEL7.4 ALT (Pegas1.0)</li> <li>• REL 7.5</li> <li>• RHEL7.6</li> <li>• XenServer 4.9</li> </ul> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> NEO-Host, operating systems</p> <p><b>Discovered in Release:</b> 4.5-1.0.1.0</p>
1521877	<p><b>Description:</b> On SLES 12 SP1 OSs, a kernel tracepoint issue may cause undefined behavior when inserting a kernel module with a wrong parameter.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> mlx5 driver, SLES 12 SP1</p> <p><b>Discovered in Release:</b> 4.5-1.0.1.0</p>
1547200	<p><b>Description:</b> When running IPoIB connected traffic with multicasts in parallel, SKB crashes.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> IPoIB, SKB</p> <p><b>Discovered in Release:</b> 4.5-1.0.1.0</p>
1504073	<p><b>Description:</b> When using ConnectX-5 with LRO over PPC systems, the HCA might experience back pressure due to delayed PCI Write operations. In this case, bandwidth might drop from line-rate to ~35Gb/s. Packet loss or pause frames might also be observed.</p> <p><b>Workaround:</b> Look for an indication of PCI back pressure ("outbound_pci_stalled_wr" counter in ethtools advancing). Disabling LRO helps reduce the back pressure and its effects.</p> <p><b>Keywords:</b> Flow Control, LRO</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>

Internal Ref. Number	Issue
1424233	<p><b>Description:</b> On RHEL v7.3, 7.4 and 7.5 OSs, setting IPv4-IP-forwarding will turn off LRO on existing interfaces. Turning LRO back on manually using ethtool and adding a VLAN interface may cause a warning call trace.</p>
	<p><b>Workaround:</b> Make sure IPv4-IP-forwarding and LRO are not turned on at the same time.</p>
	<p><b>Keywords:</b> IPv4 forwarding, LRO</p>
	<p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>
1418447	<p><b>Description:</b> When working in IPoIB ULP (non-enhanced) mode, IPv6 may disappear in case ring size is changed dynamically (while the driver is running).</p>
	<p><b>Workaround:</b> There are three workarounds for this issue:</p> <ul style="list-style-type: none"> <li>• Perform static configuration of ring size instead of dynamic configuration</li> <li>• In case you have run dynamic configuration, run <i>ifdown ifup</i> afterwards</li> <li>• On supported kernels, enable <i>keep_addr_on_down</i> IPv6 sysfs parameter before configuring the ring size dynamically</li> </ul>
	<p><b>Keywords:</b> IPoIB, ULP mode, ring size</p>
	<p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>
1442507	<p><b>Description:</b> Retpoline support in GCC causes an increase in CPU utilization, which results in IP forwarding's 15% performance drop.</p>
	<p><b>Workaround:</b> N/A</p>
	<p><b>Keywords:</b> Retpoline, GCC, CPU, IP forwarding, Spectre attack</p>
	<p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>
1417414	<p><b>Description:</b> When working with old kernel versions that do not include the <code>unregister_netdevice_notifier</code> function fix (introduced in "net: In <code>unregister_netdevice_notifier</code> unregister the netdevices" commit), reloading <code>ib_ipoib</code> module using <code>modprobe</code> will fail with the following error message: "<i>Cannot allocate memory</i>".</p>

Internal Ref. Number	Issue
	<p><b>Workaround:</b> Reload the driver instead of modprobe by running: <i>/etc/init.d/openibd restart</i></p> <p><b>Keywords:</b> IPoIB</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>
1400381	<p><b>Description:</b> On SLES 11 SP3 PPC64 OSs, a memory allocation issue may prevent the interface from loading after reboot, resulting in a call trace in the message log.</p> <p><b>Workaround:</b> Restart the driver.</p> <p><b>Keywords:</b> SLES11 SP3</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>
1425129	<p><b>Description:</b> MLNX_OFED cannot be installed on SLES 15 OSs using Zypper repository.</p> <p><b>Workaround:</b> Install MLNX_OFED using the standard installation script instead of Zypper repository.</p> <p><b>Keywords:</b> Installation, SLES, Zypper</p> <p><b>Discovered in Release:</b> 4.4-1.0.0.0</p>
1241056	<p><b>Description:</b> When working with ConnectX-4/ConnectX-5 HCAs on PPC systems with Hardware LRO and Adaptive Rx support, bandwidth drops from full wire speed (FWS) to ~60Gb/s.</p> <p><b>Workaround:</b> Make sure to disable Adaptive Rx when enabling Hardware LRO: <i>ethtool -C &lt;interface&gt; adaptive-rx off</i> <i>ethtool -C &lt;interface&gt; rx-usecs 8 rx-frames 128</i></p> <p><b>Keywords:</b> Hardware LRO, Adaptive Rx, PPC</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1090612	<p><b>Description:</b> NVMeoF protocol does not support LBA format with non-zero metadata size. Therefore, NVMe namespace configured to LBA format with metadata size bigger than 0 will cause Enhanced Error Handling (EEH) in PowerPC systems.</p> <p><b>Workaround:</b> Configure the NVMe namespace to use LBA format with zero sized metadata.</p>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> NVMeoF, PowerPC, EEH</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1243581	<p><b>Description:</b> In switchdev mode, the IB device exposed does not support MADs. As a result, tools such as ibstat that work with MADs will not function properly.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> switchdev, IB representors, mlx5, MADs</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1309621	<p><b>Description:</b> In switchdev mode default configuration, stateless offloads/steering based on inner headers is not supported.</p> <p><b>Workaround:</b> To enable stateless offloads/steering based on inner headers, disable encap by running:  <code>devlink dev eswitch show pci/0000:83:00.1 encap disable</code>  Or, in case devlink is not supported by the kernel, run:  <code>echo none &gt; /sys/kernel/debug/mlx5/&lt;BDF&gt;/compat/encap</code></p> <p><b>Note:</b> This is a hardware-related limitation.</p> <p><b>Keywords:</b> switchdev, stateless offload, steering</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1268718	<p><b>Description:</b> ConnectX-5 supports up to 62 IB representors. When attempting to move to switchdev mode where more than 62 VFs are initialized, the call will fail with the following error message:  <i>" devlink answers: Invalid argument "</i></p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> ConnectX-5, IB representors</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1275082	<p><b>Description:</b> When setting a non-default IPv6 link local address or an address that is not based on the device MAC, connection establishments over RoCEv2 might fail.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> IPV6, RoCE, link local address</p>

Internal Ref. Number	Issue
	<b>Discovered in Release:</b> 4.3-1.0.1.0
1307336	<b>Description:</b> In RoCE LAG mode, when running <code>ibdev2netdev -v</code> , the port state of the second port of the <code>mlx4_0</code> IB device will read "NA" since this IB device does not have a second port.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> mlx4, RoCE LAG, ibdev2netdev, bonding
	<b>Discovered in Release:</b> 4.3-1.0.1.0
1316654	<b>Description:</b> PKEY interface receives PTP delay requests without a time-stamp.
	<b>Workaround:</b> Run <code>ptp4l</code> over the parent interface.
	<b>Keywords:</b> PKEY, PTP
	<b>Discovered in Release:</b> 4.3-1.0.1.0
1296355	<b>Description:</b> Number of MSI-X that can be allocated for VFs and PFs in total is limited to 2300 on Power9 platforms.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> MSI-X, VF, PF, PPC, SR-IOV
	<b>Discovered in Release:</b> 4.3-1.0.1.0
1294934	<b>Description:</b> Firmware reset might cause Enhanced Error Handling (EEH) on Power7 platforms.
	<b>Workaround:</b> N/A
	<b>Keywords:</b> EEH, PPC
	<b>Discovered in Release:</b> 4.3-1.0.1.0
1259293	<p><b>Description:</b> On Fedora 20 operating systems, driver load fails with an error message such as: "<code>[185.262460] kmem_cache_sanity_check (fs_ftes_0000:00:06.0): Cache name already exists.</code>"</p> <p>This is caused by SLUB allocators grouping multiple slab <code>kmem_cache_create</code> into one slab cache alias to save memory and increase cache hotness. This results in the slab name to be considered stale.</p>

Internal Ref. Number	Issue
	<p><b>Workaround:</b> Upgrade the kernel version to kernel-3.19.8-100.fc20.x86_64. Note that after rebooting to the new kernel, you will need to rebuild MLNX_OFED against the new kernel version.</p> <p><b>Keywords:</b> Fedora, driver load</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1264359	<p><b>Description:</b> When running perfctest (ib_send_bw, ib_write_bw, etc.) in rdma-cm mode, the resp_cqe_error counter under /sys/class/infiniband/mlx5_0/ports/1/hw_counters/resp_cqe_error might increase. This behavior is expected and it is a result of receive WQEs that were not consumed.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> perfctest, RDMA CM, mlx5</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1294575	<p><b>Description:</b> Traffic may hang while working in IPoIB SR-IOV environment.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> IPoIB, SR-IOV</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1227577	<p><b>Description:</b> Due to Enhanced IPoIB's lack of priority-based flow control, PTP accuracy may adversely be affected by heavy TCP traffic.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Enhanced IPoIB, PTP</p> <p><b>Discovered in Release:</b> 4.3-1.0.1.0</p>
1264956	<p><b>Description:</b> Configuring SR-IOV after disabling RoCE LAG using sysfs (/sys/bus/pci/drivers/mlx5_core/ /roce_lag_enable) might result in RoCE LAG being enabled again in case SR-IOV configuration fails.</p> <p><b>Workaround:</b> Make sure to disable RoCE LAG once again.</p> <p><b>Keywords:</b> RoCE LAG, SR-IOV</p>

Internal Ref. Number	Issue
	<b>Discovered in Release:</b> 4.3-1.0.1.0
1263043	<p><b>Description:</b> On RHEL7.4, due to an OS issue introduced in kmod package version 20-15.el7_4.6, parsing the depmod configuration files will fail, resulting in either of the following issues:</p> <ul style="list-style-type: none"> <li>• Driver restart failure prompting an error message, such as: "<i>ERROR: Module mlx5_core belong to kernel which is not a part of MLNX_OFED, skipping...</i>"</li> <li>• nvmet_rdma kernel module dysfunction, despite installing MLNX_OFED using the "--with-nvme" option. An error message, such as: "<i>nvmet_rdma: unknown parameter 'offload_mem_start' ignored</i>" will be seen in <i>dmesg</i> output</li> </ul> <p><b>Workaround:</b> Go to <a href="#">RedHat webpage</a> to upgrade the kmod package version.</p> <p><b>Keywords:</b> driver restart, kmod, kmp, nvme, nvmet_rdma</p> <p><b>Discovered in Release:</b> 4.2-1.2.0.0</p>
1229160	<p><b>Description:</b> Changing IPoIB Tx/Rx ring size dynamically using ethtool is not permitted.</p> <p><b>Workaround:</b> Use the <code>send_queue_size/recv_queue_size</code> module parameters to change the Tx/Rx ring size.</p> <p><b>Keywords:</b> IPoIB, queue size</p> <p><b>Discovered in Release:</b> 4.2-1.2.0.0</p>
1214477	<p><b>Description:</b> On vRedHat 7.2 operating systems, when Network Manager is enabled, IPoIB interfaces may not get an IPv6 address due to an issue in the Network Manager.</p> <p><b>Workaround:</b> Disable Network Manager or upgrade its version.</p> <p><b>Keywords:</b> Network Manager, IPoIB, IPv6</p> <p><b>Discovered in Release:</b> 4.2-1.2.0.0</p>
-	<p><b>Description:</b> Packet Size (Actual Packet MTU) limitation for IPsec offload on Innova IPsec adapter cards: The current offload implementation does not support IP fragmentation. The original packet size should be such</p>

Internal Ref. Number	Issue
	<p>that it does not exceed the interface's MTU size after the ESP transformation (encryption of the original IP packet which increases its length) and the headers (outer IP header) are added:</p> <ul style="list-style-type: none"> <li>• Inner IP packet size <math>\leq</math> I/F MTU - ESP additions (20) - outer_IP (20) - fragmentation issue reserved length (56)</li> <li>• Inner IP packet size <math>\leq</math> I/F MTU - 96</li> </ul> <p>This mostly affects forwarded traffic into smaller MTU, as well as UDP traffic. TCP does PMTU discovery by default and clamps the MSS accordingly.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, MTU</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
-	<p><b>Description:</b> No LLC/SNAP support on Innova IPsec adapter cards.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, LLC/SNAP</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
-	<p><b>Description:</b> No support for FEC on Innova IPsec adapter cards. When using switches, there may be a need to change its configuration.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, FEC</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
955929	<p><b>Description:</b> Heavy traffic may cause SYN flooding when using Innova IPsec adapter cards.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, SYN flooding</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
-	<p><b>Description:</b> Priority Based Flow Control is not supported on Innova IPsec adapter cards.</p>



Internal Ref. Number	Issue
	<p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, Priority Based Flow Control</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
-	<p><b>Description:</b> Pause configuration is not supported when using Innova IPsec adapter cards. Default pause is global pause (enabled).</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, Global pause</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1045097	<p><b>Description:</b> Connecting and disconnecting a cable several times may cause a link up failure when using Innova IPsec adapter cards.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Innova IPsec, Cable, link up</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
-	<p><b>Description:</b> On Innova IPsec adapter cards, supported MTU is between 512 and 2012 bytes. Setting MTU values outside this range might fail or might cause traffic loss.</p> <p><b>Workaround:</b> Set MTU between 512 and 2012 bytes.</p> <p><b>Keywords:</b> Innova IPsec, MTU</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1177196	<p><b>Description:</b> If OpenSM version is 4.8.1 and below, the IB interfaces link remains Down while the "SRIOV_IB_ROUTING_MODE_P1=1" and "SRIOV_IB_ROUTING_MODE_P2=1" flags are enabled in the HCA.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> OpenSM, SR-IOV, IB link</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1118530	<p><b>Description:</b> On kernel versions 4.10-4.13, when resetting sriov_numvfs to 0 on PowerPC systems, the following dmesg warning will appear: mlx5_core &lt;BDF&gt;: can't update enabled VF BAR0</p>

Internal Ref. Number	Issue
	<p><b>Workaround:</b> Reboot the system to reset <i>sriov_numvfs</i> value.</p> <p><b>Keywords:</b> SR-IOV, numvfs</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1125184	<p><b>Description:</b> In old kernel versions, such as Ubuntu 14.04 and RedHat 7.1, VXLAN interface does not reply to ARP requests for a MAC address that exists in its own ARP table. This issue was fixed in the following newer kernel versions: Ubuntu 16.04 and RedHat 7.3.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> ARP, VXLAN</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1171764	<p><b>Description:</b> Connecting multiple ports on the same server to the same subnet (IP/IB) will cause all interfaces connected to that subnet to respond to ARP requests. As a result, wrong ARP replies might be received when trying to resolve IP addresses.</p> <p><b>Workaround:</b> Run the following to make sure only the interface with the requested IP address responds to the ARP request:  <code>sysctl -w net.ipv4.conf.all.arp_ignore=1</code></p> <p><b>Keywords:</b> IPoIB, librdmacm, ARP</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1134323	<p><b>Description:</b> When using kernel versions older than version 4.7 with IOMMU enabled, performance degradations and logical issues (such as soft lockup) might occur upon high load of traffic. This is caused due to the fact that IOMMU IOVA allocations are centralized, requiring many synchronization operations and high locking overhead amongst CPUs.</p>

Internal Ref. Number	Issue
	<p><b>Workaround:</b> Use kernel v4.7 or above, or a backported kernel that includes the following patches:</p> <ul style="list-style-type: none"> <li>• 2aac630429d9 iommu/vt-d: change intel-iommu to use IOVA frame numbers</li> <li>• 9257b4a206fc iommu/iova: introduce per-cpu caching to iova allocation</li> <li>• 22e2f9fa63b0 iommu/vt-d: Use per-cpu IOVA caching</li> </ul> <p><b>Keywords:</b> IOMMU, soft lockup</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1135738	<p><b>Description:</b> On 64k page size setups, DMA memory might run out when trying to increase the ring size/number of channels.</p> <p><b>Workaround:</b> Reduce the ring size/number of channels.</p> <p><b>Keywords:</b> DMA, 64K page</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1159650	<p><b>Description:</b> When configuring VF VST, VLAN-tagged outgoing packets will be dropped in case of ConnectX-4 HCAs. In case of ConnectX-5 HCAs, VLAN-tagged outgoing packets will have another VLAN tag inserted.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> VST</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1157770	<p><b>Description:</b> On Passthrough/VM machines with relatively old QEMU and libvirtd, CMD timeout might occur upon driver load. After timeout, no other commands will be completed and all driver operations will be stuck.</p> <p><b>Workaround:</b> Upgrade the QEMU and libvirtd on the KVM server. Tested with (Ubuntu 16.10) are the following versions:</p> <ul style="list-style-type: none"> <li>• libvirt 2.1.0</li> <li>• QEMU 2.6.1</li> </ul>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> QEMU</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1147703	<p><b>Description:</b> Using dm-multipath for High Availability on top of NVMeoF block devices must be done with “directio” path checker.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> NVMeoF</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1152408	<p><b>Description:</b> RedHat v7.3 PPCLE and v7.4 PPCLE operating systems do not support KVM qemu out of the box. The following error message will appear when attempting to run <i>virt-install</i> to create new VMs: <i>Cant find qemu-kvm package to install</i></p> <p><b>Workaround:</b> Acquire the following rpms from the beta version of 7.4ALT to 7.3/7.4 PPCLE (in the same order):</p> <ul style="list-style-type: none"> <li>• qemu-img-.el7a.ppc64le.rpm</li> <li>• qemu-kvm-common-.el7a.ppc64le.rpm</li> <li>• qemu-kvm-.el7a.ppc64le.rpm</li> </ul> <p><b>Keywords:</b> Virtualization, PPC, Power8, KVM, RedHat, PPC64LE</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1012719	<p><b>Description:</b> A soft lockup in the CQ polling flow might occur when running very high stress on the GSI QP (RDMA-CM applications). This is a transient situation from which the driver will later recover.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RDMA-CM, GSI QP, CQ</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1062940	<p><b>Description:</b> When running Network Manger on devices on which Enhanced IPoIB is enabled, CONNECTED_MODE can only be set to NO/AUTO. Setting it to YES will prevent the interface from being configured.</p> <p><b>Workaround:</b> N/A</p>

Internal Ref. Number	Issue
	<p><b>Keywords:</b> Enhanced IPoIB, network manager, connected_mode</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1078630	<p><b>Description:</b> When working in RoCE LAG over kernel v3.10, a kernel crash might occur when unloading the driver as the Network Manager is running.</p> <p><b>Workaround:</b> Stop the Network Manager before unloading the driver and start it back once the driver unload is complete.</p> <p><b>Keywords:</b> RoCE LAG, network manager</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1149557	<p><b>Description:</b> When setting VGT+, the maximal number of allowed VLAN IDs presented in the sysfs is 813 (up to the first 813).</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> VGT+</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1122619	<p><b>Description:</b> On Arm setups, DMA memory resource is limited due to a default CMA limitation.</p> <p><b>Workaround:</b> Increase the CMA limitation or cancel its use, using the kernel's CMD line parameters:</p> <ul style="list-style-type: none"> <li>• Add the parameter cma=256M to increase the CMA limit to 256MB</li> <li>• Add the parameter cma=0 to disable the use of CMA</li> </ul> <p><b>Keywords:</b> IPoIB, CMA</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
1146837	<p><b>Description:</b> On SLES11 SP1 operating system, IPoIB interface renaming process may fail due to a broken udev rule, leaving interfaces with names like ib0_rename.</p> <p><b>Workaround:</b></p> <ol style="list-style-type: none"> <li>1. Open the udev conf file "/etc/udev/rules.d/70-persistent-net.rules", and remove such lines as SUBSYSTEM=="net", ACTION=="add",</li> </ol>

Internal Ref. Number	Issue
	<p>DRIVERS=="?*" , =="" , NAME="eth0". 2. Reload the driver stack.</p> <p><b>Keywords:</b> IPoIB</p> <p><b>Discovered in Release:</b> 4.2-1.0.0.0</p>
-	<p><b>Description:</b> NVMeoF support is available for the following:</p> <ul style="list-style-type: none"> <li>• SLES 12.3 and above</li> <li>• RHEL 7.2 and above (Host side only)</li> <li>• RHEL 7.4 and above (Host and Target side)</li> <li>• OS with distribution/custom kernel &gt;= 4.8.x</li> </ul> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> NVMeoF Host/Target</p>
995665/1165 919	<p><b>Description:</b> In kernels below v4.13, connection between NVMeoF host and target cannot be established in a hyper-threaded system with more than 1 socket.</p> <p><b>Workaround:</b> On the host side, connect to NVMeoF subsystem using <code>--nr-io-queues &lt;num_queues&gt;</code> flag. Note that <code>num_queues</code> must be lower or equal to <code>num_sockets</code> multiplied with <code>num_cores_per_socket</code>.</p> <p><b>Keywords:</b> NVMeoF</p>
1039346	<p><b>Description:</b> Enabling multiple namespaces per subsystem while using NVMeoF target offload is not supported on ConnectX-5 adapter cards.</p> <p><b>Workaround:</b> To enable more than one namespace, create a subsystem for each one.</p> <p><b>Keywords:</b> NVMeoF Target Offload, namespace</p>
1072347	<p><b>Description:</b> Ethtool -i displays incorrect driver name for devices with enhanced IPoIB support.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Enhanced IPoIB, Ethtool</p>

Internal Ref. Number	Issue
1071457	<p><b>Description:</b> PKEY-related limitations in enhanced IPoIB:</p> <ul style="list-style-type: none"> <li>• Since the parent interface ib&lt;x&gt; and the child interface ib&lt;x&gt;.yyyy share the same receive resources, the parent interface's MTU cannot be less than the child interface's MTU</li> <li>• Interface counters and Ethtool control are not supported on child interfaces</li> <li>• Parent interface should be in UP state to enable child interface to receive traffic</li> </ul>
	<p><b>Workaround:</b> N/A</p>
	<p><b>Keywords:</b> PKEY, Enhanced IPoIB, MTU, Ethtool, Interface Counters</p>
1059451	<p><b>Description:</b> When Enhanced IPoIB is enabled, the following module parameters will not be functional:</p> <ul style="list-style-type: none"> <li>• send_queue_size</li> <li>• recv_queue_size</li> <li>• max_nonsrq_conn_qp</li> </ul>
	<p><b>Workaround:</b> N/A</p>
	<p><b>Keywords:</b> Enhance IPoIB</p>
1030301	<p><b>Description:</b> Creating virtual functions on a device that is in LAG mode will destroy the LAG configuration. The bonding device over the Ethernet NICs will continue to work as expected.</p>
	<p><b>Workaround:</b> N/A</p>
	<p><b>Keywords:</b> LAG, SR-IOV</p>
1047616	<p><b>Description:</b> When node GUID of a device is set to zero (0000:0000:0000:0000), RDMA_CM user space application may crash.</p>
	<p><b>Workaround:</b> Set node GUID to a nonzero value.</p>
	<p><b>Keywords:</b> RDMA_CM</p>
1061298	<p><b>Description:</b> Since enhanced IPoIB does not support connected mode on RedHat operating systems, when using network manger and</p>

Internal Ref. Number	Issue
	<p>enhanced IPoIB capable devices, <i>CONNECTED_MODE</i> must be set to NO/AUTO. Setting <i>CONNECTED_MODE</i> to yes will cause the interface to not be configured.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Enhanced IPoIB</p>
1068215	<p><b>Description:</b> When enhanced IPoIB mode is enabled, ring size limit is 8k. When it is disabled, ring size limit is decreased to 4k.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Enhanced IPoIB</p>
1051701	<p><b>Description:</b> New versions of iproute which support new kernel features may misbehave on old kernels that do not support these new features.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> iproute</p>
1007830	<p><b>Description:</b> When working on Xenserver hypervisor with SR-IOV enabled on it, make sure the following instructions are applied:</p> <ol style="list-style-type: none"> <li>1. Right after enabling SR-IOV, unbind all driver instances of the virtual functions from their PCI slots.</li> <li>2. It is not allowed to unbind PF driver instance while having active VFs.</li> </ol> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> SR-IOV</p>
1008583	<p><b>Description:</b> A soft lockup in the CQ polling flow might occur when running very high stress on the GSI QP (RDMA-CM applications). This is a transient situation and the driver recovers from it after a while.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RDMA-CM</p>
1007356	<p><b>Description:</b> Creating a PKEY interface using "<i>ip link</i>" is not supported.</p>



Internal Ref. Number	Issue
	<p><b>Workaround:</b> Use sysfs to create a PKEY interface.</p> <p><b>Keywords:</b> IPoIB, PKEY</p>
1000197	<p><b>Description:</b> Displaying multicast groups using sysfs may not show all the entries on Fedora 23 OS.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> IPoIB</p>
1010148	<p><b>Description:</b> Upgrading from MLNX_OFED v3.x to v4.x using yum and apt-get repositories fails.</p> <p><b>Workaround:</b> Remove MLNX_OFED v3.x using the <i>ofed_uninstall.sh</i> script, and only then install MLNX_OFED v4.x as usual.</p> <p><b>Keywords:</b> Installation</p>
1005786	<p><b>Description:</b> When using ConnectX-5 adapter cards, the following error might be printed to dmesg, indicating temporary lack of DMA pages:  <i>“mlx5_core ... give_pages:289:(pid x): Y pages alloc time exceeded the max permitted duration</i>  <i>mlx5_core ... page_notify_fail:263:(pid x): Page allocation failure notification on func_id(z) sent to fw</i>  <i>mlx5_core ... pages_work_handler:471:(pid x): give fail -12”</i></p> <p><b>Example:</b> This might happen when trying to open more than 64 VFs per port.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> mlx5_core, DMA</p>
1008066/1009004	<p><b>Description:</b> Performing some operations on the user end during reboot might cause call trace/panic, due to bugs found in the Linux kernel.  For example: Running <i>get_vf_stats</i> (via iptool) during reboot.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> mlx5_core, reboot</p>
1009488	<p><b>Description:</b> Mounting MLNX_OFED to a path that contains special characters, such as parenthesis or spaces is not supported. For example,</p>

Internal Ref. Number	Issue
	<p>when mounting MLNX_OFED to <code>"/media/CDROM(vcd)"/</code>, installation will fail and the following error message will be displayed:</p> <pre># cd /media/CDROM(vcd)/ # ./mlnxofedinstall sh: 1: Syntax error: "(" unexpected</pre> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Installation</p>
982144	<p><b>Description:</b> When offload traffic sniffer is on, the bandwidth could decrease up to 50%.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Offload Traffic Sniffer</p>
981045	<p><b>Description:</b> On kernels below v4.2, when removing a bonding module with devices different from ARPHRD_ETHER, a call trace may be received.</p> <p><b>Workaround:</b> Remove the bond in the following order: Remove the slaves, delete the bond, and only then remove the bonding module.</p> <p><b>Keywords:</b> Bonding</p>
980066/981314	<p><b>Description:</b> Soft RoCE does not support Extended Reliable Connection (XRC).</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Soft RoCE, XRC</p>
982534	<p><b>Description:</b> In ConnectX-3, when using a server with page size of 64K, the UAR BAR will become too small. This may cause one of the following issues:</p> <ol style="list-style-type: none"> <li>1. mlx4_core driver does not load.</li> <li>2. The mlx4_core driver does load, but calls to <code>ibv_open_device</code> may return ENOMEM errors.</li> </ol> <p><b>Workaround:</b></p> <ol style="list-style-type: none"> <li>1. Add the following parameter in the firmware's ini file under [HCA] section: <code>log2_uar_bar_megabytes = 7</code></li> </ol>

Internal Ref. Number	Issue
	<p>2. Re-burn the firmware with the new ini file.</p> <p><b>Keywords:</b> PPC</p>
981362	<p><b>Description:</b> On several OSs, setting a number of TC is not supported via the tc tool.</p> <p><b>Workaround:</b> Set the number of TC via the <code>/sys/class/net/ /qos/tc_num sysfs</code> file.</p> <p><b>Keywords:</b> Ethernet, TC</p>
980257	<p><b>Description:</b> An issue in InfiniBand bond interfaces may cause memory corruption in Ubuntu v14.04 and v14.10 OSs. The memory corruption happens when attempting to reload the driver while the bond is up with InfiniBand slaves.</p> <p><b>Workaround:</b> Delete the bond before restarting the driver.</p> <p><b>Keywords:</b> Bonding, IPoIB</p>
980034/981311	<p><b>Description:</b> Soft RoCE counters located under <code>/sys/class/infiniband/ /ports/1/counters/</code> directory are not supported.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Soft RoCE</p>
979907	<p><b>Description:</b> Only the following two experimental verbs are supported for Soft RoCE:</p> <ul style="list-style-type: none"> <li>• <code>ibv_exp_query_device</code></li> <li>• <code>ibv_exp_poll_cq</code>.</li> </ul> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> Soft RoCE</p>
979457	<p><b>Description:</b> When setting IOMMU=ON, a severe performance degradation may occur due to a bug in IOMMU.</p> <p><b>Workaround:</b> Make sure the following patches are found in your kernel:</p> <ul style="list-style-type: none"> <li>• <code>iommu/vt-d: Fix PASID table allocation</code></li> <li>• <code>iommu/vt-d: Fix IOMMU lookup for SR-IOV Virtual Functions</code></li> </ul>

Internal Ref. Number	Issue
	<p><b>Note:</b> These patches are already available in Ubuntu 16.04.02 and 17.04 OSs.</p> <p><b>Keywords:</b> Performance, IOMMU</p>
977852	<p><b>Description:</b> <i>rdma_cm</i> running over IB ports does not support UD QPs on ConnectX-3 adapter cards.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> SR-IOV, RDMA CM</p>
955113/977990	<p><b>Description:</b> In RoCE LAG over ConnectX-4 adapter cards, the script <i>ibdev2netdev</i> may show a wrong port state for the bonded device. This means that although the IB device/port <i>mlx5_bond_0/1</i> is up (as seen in <i>ibstat</i>), <i>ibdev2netdev</i> may report that it is down.</p> <p><b>Workaround:</b> N/A</p> <p><b>Keywords:</b> RoCE, LAG, bonding</p>

---

# User Manual

- [Introduction](#)
- [Installation](#)
- [Features Overview and Configuration](#)
- [Programming](#)
- [InfiniBand Fabric Utilities](#)
- [Troubleshooting](#)
- [Common Abbreviations and Related Documents](#)

## Introduction

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

This document provides instructions on how to install the driver on NVIDIA ConnectX® network adapter solutions supporting the following uplinks to servers.

Uplink/NICs	Driver Name	Uplink Speed
ConnectX®-3/Connect X-3 Pro	mlx4	<ul style="list-style-type: none"><li>• InfiniBand: SDR, QDR, FDR10, FDR</li><li>• Ethernet: 10GbE, 40GbE 56GbE<sup>1</sup></li></ul>
ConnectX-4	mlx5	<ul style="list-style-type: none"><li>• InfiniBand: SDR, QDR, FDR, FDR10, EDR</li><li>• Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 56GbE<sup>1</sup>, 100GbE</li></ul>

Uplink/NICs	Driver Name	Uplink Speed
ConnectX-4 Lx		<ul style="list-style-type: none"> <li>Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE</li> </ul>
ConnectX-5/ConnectX-5 Ex		<ul style="list-style-type: none"> <li>InfiniBand: SDR, QDR, FDR, FDR10, EDR</li> <li>Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 100GbE</li> </ul>
ConnectX-6		<ul style="list-style-type: none"> <li>InfiniBand: SDR, FDR, EDR, HDR</li> <li>Ethernet: 10GbE, 25GbE, 40GbE, 50GbE<sup>2</sup>, 100GbE<sup>2</sup>, 200GbE<sup>2</sup></li> </ul>
ConnectX-6 Dx		<ul style="list-style-type: none"> <li>Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE<sup>1</sup>, 100GbE<sup>1</sup>, 200GbE<sup>2</sup></li> </ul>
Innova™ IPsec EN		<ul style="list-style-type: none"> <li>Ethernet: 10GbE, 40GbE</li> </ul>
Connect-IB®		<ul style="list-style-type: none"> <li>InfiniBand: SDR, QDR, FDR10, FDR</li> </ul>

1. 56GbE is a NVIDIA propriety link speed can be achieved while connecting a NVIDIA adapter card to NVIDIA SX10XX switch series, or connecting a NVIDIA adapter card to another NVIDIA adapter card.

2. Supports both NRZ and PAM4 modes.

All NVIDIA network adapter cards are compatible with OpenFabrics-based RDMA protocols and software and are supported by major operating system distributions.

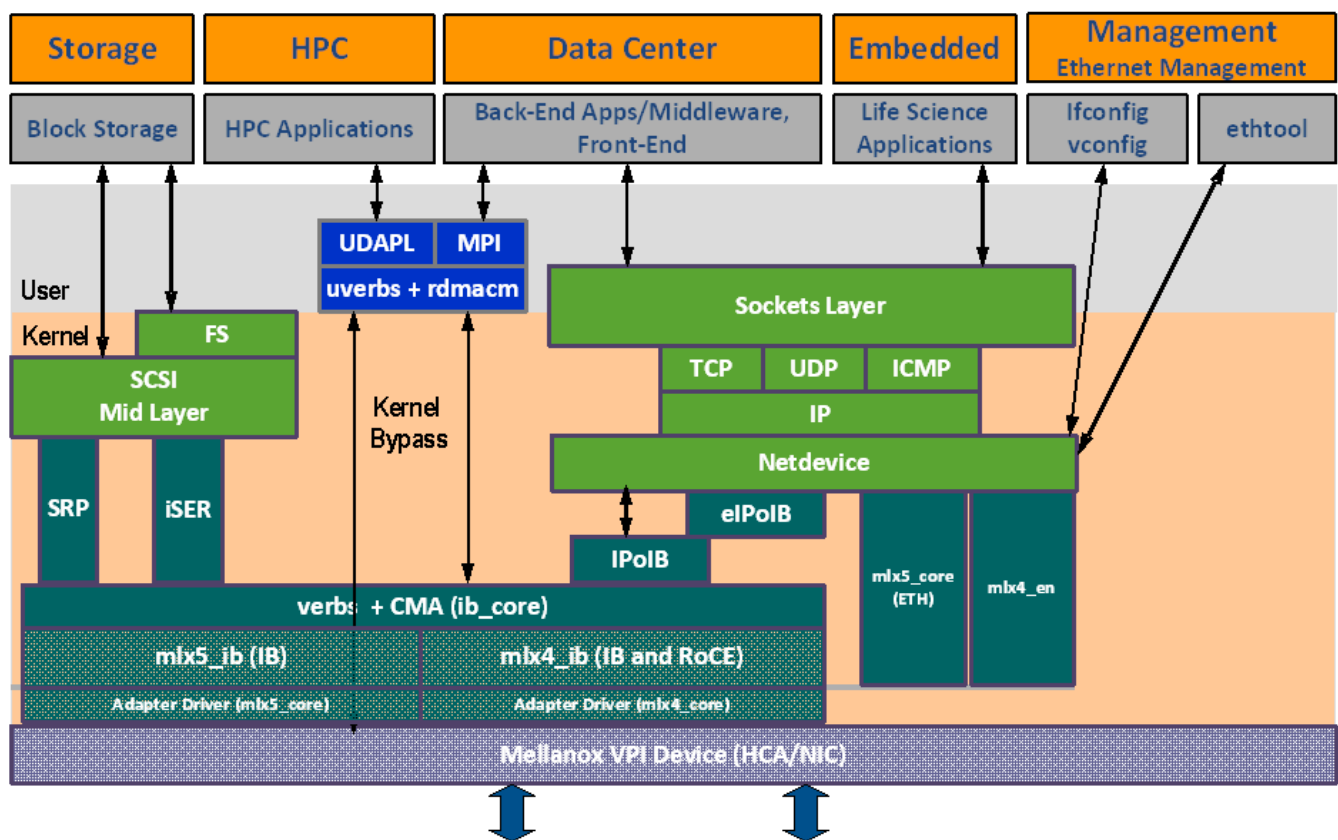
NVIDIA OFED is certified with the following products:

- NVIDIA Messaging Accelerator (VMA™) software: Socket acceleration library that performs OS bypass for standard socket-based applications. Please note, VMA support is provided separately from NVIDIA OFED support. For further information, please refer to the VMA documentation (<https://docs.nvidia.com/networking/category/vma>).

- NVIDIA Unified Fabric Manager (UFM®) software: Powerful platform for managing demanding scale-out computing fabric environments, built on top of the OpenSM industry standard routing engine.
- Fabric Collective Accelerator (FCA)—FCA is a NVIDIA MPI-integrated software package that utilizes CORE-Direct technology for implementing the MPI collectives communications.

## Stack Architecture

The figure below shows a diagram of the NVIDIA OFED stack, and how upper layer protocols (ULPs) interface with the hardware and with the kernel and userspace. The application level also shows the versatility of markets that NVIDIA OFED applies to.



The following subsections briefly describe the various components of the NVIDIA OFED stack.

### mlx4 VPI Driver

mlx4 is the low-level driver implementation for the ConnectX® family adapters designed by NVIDIA. ConnectX-3 adapters can operate as an InfiniBand adapter, or as an Ethernet

NIC. The OFED driver supports InfiniBand and Ethernet NIC configurations. To accommodate the supported configurations, the driver is split into the following modules:

### **mlx4\_core**

Handles low-level functions like device initialization and firmware commands processing. Also controls resource allocation so that the InfiniBand and Ethernet functions can share the device without interfering with each other.

### **mlx4\_ib**

Handles InfiniBand-specific functions and plugs into the InfiniBand mid layer.

### **mlx4\_en**

A 10/40GigE driver under drivers/net/ethernet/mellanox/mlx4 that handles Ethernet specific functions and plugs into the netdev mid layer.

## **mlx5 Driver**

mlx5 is the low-level driver implementation for the Connect-IB® and ConnectX-4 adapters designed by NVIDIA. Connect-IB® operates as an InfiniBand adapter whereas ConnectX®-4 operates as a VPI adapter (Infiniband and Ethernet). The mlx5 driver is comprised of the following kernel modules:

### **mlx5\_core**

Acts as a library of common functions (e.g. initializing the device after reset) required by Connect-IB® and ConnectX®-4 adapter cards. mlx5\_core driver also implements the Ethernet interfaces for ConnectX®-4. Unlike mlx4\_en/core, mlx5 drivers do not require the mlx5\_en module as the Ethernet functionalities are built-in in the mlx5\_core module.

### **mlx5\_ib**

Handles InfiniBand-specific functions and plugs into the InfiniBand mid layer.

## **libmlx5**

libmlx5 is the provider library that implements hardware specific user-space functionality. If there is no compatibility between the firmware and the driver, the driver will not load and a message will be printed in the dmesg.

The following are the libmlx5 environment variables:

- MLX5\_FREEZE\_ON\_ERROR\_CQE
  - Causes the process to hang in a loop of completion with error, which is not flushed with error or retry exceeded occurs/
  - Otherwise disabled



- MLX5\_POST\_SEND\_PREFER\_BF
  - Configures every work request that can use blue flame will use blue flame
- Otherwise - blue flame depends on the size of the message and inline indication in the packet
- MLX5\_SHUT\_UP\_BF
  - Disables blue flame feature
  - Otherwise - do not disable
- MLX5\_SINGLE\_THREADED
  - All spinlocks are disabled
  - Otherwise - spinlocks enabled
  - Used by applications that are single threaded and would like to save the overhead of taking spinlocks.
- MLX5\_CQE\_SIZE
  - 64 - completion queue entry size is 64 bytes (default)
  - 128 - completion queue entry size is 128 bytes
- MLX5\_SCATTER\_TO\_CQE
  - Small buffers are scattered to the completion queue entry and manipulated by the driver. Valid for RC transport.
  - Default is 1, otherwise disabled
- MLX5\_ENABLE\_CQE\_COMPRESSION
  - Saves PCIe bandwidth by compressing a few CQEs into a smaller amount of bytes on PCIe. Setting this variable 1 enables CQE compression.
  - Default value 0 (disabled)

- `MLX5_RELAXED_PACKET_ORDERING_ON`  
See [“Out-of-Order \(OOO\) Data Placement Experimental Verbs”](#) section.

## Mid-layer Core

Core services include management interface (MAD), connection manager (CM) interface, and Subnet Administrator (SA) interface. The stack includes components for both user-mode and kernel applications. The core services run in the kernel and expose an interface to user-mode for verbs, CM and management.

## Upper Layer Protocols (ULPs)

### IP over IB (IPoIB)

The IP over IB (IPoIB) driver is a network interface implementation over InfiniBand. IPoIB encapsulates IP datagrams over an InfiniBand connected or datagram transport service. IPoIB pre-appends the IP datagrams with an encapsulation header and sends the outcome over the InfiniBand transport service. The transport service is Unreliable Datagram (UD) by default, but it may also be configured to be Reliable Connected (RC), in case RC is supported. The interface supports unicast, multicast and broadcast. For details, see [“IP over InfiniBand \(IPoIB\)”](#) section.

### iSCSI Extensions for RDMA (iSER)

iSCSI Extensions for RDMA (iSER) extends the iSCSI protocol to RDMA. It permits data to be transferred directly into and out of SCSI buffers without intermediate data copies. For further information, please refer to [“iSCSI Extensions for RDMA \(iSER\)”](#) section.

### SCSI RDMA Protocol (SRP)

SCSI RDMA Protocol (SRP) is designed to take full advantage of the protocol offload and RDMA features provided by the InfiniBand architecture. SRP allows a large body of SCSI software to be readily used on InfiniBand architecture. The SRP driver—known as the SRP Initiator—differs from traditional low-level SCSI drivers in Linux. The SRP Initiator does not control a local HBA; instead, it controls a connection to an I/O controller—known as the SRP Target—to provide access to remote storage devices across an InfiniBand fabric. The SRP Target resides in an I/O unit and provides storage services. See [“SRP - SCSI RDMA Protocol”](#) section.

### User Direct Access Programming Library (uDAPL)

User Direct Access Programming Library (uDAPL) is a standard API that promotes data center application data messaging performance, scalability, and reliability over RDMA interconnects InfiniBand and RoCE. The uDAPL interface is defined by the DAT collaborative. This release of the uDAPL reference implementation package for both DAT 1.2 and 2.0 specification is timed to coincide with OFED release of the Open Fabrics ([www.openfabrics.org](http://www.openfabrics.org)) software stack.

## **MPI**

Message Passing Interface (MPI) is a library specification that enables the development of parallel software libraries to utilize parallel computers, clusters, and heterogeneous networks. OFED includes the following MPI implementation over InfiniBand:

- Open MPI – an open source MPI-2 implementation by the Open MPI Project

OFED also includes MPI benchmark tests such as OSU BW/LAT, Intel MPI BeBenchmark and Presta.

## **InfiniBand Subnet Manager**

All InfiniBand-compliant ULPs require a proper operation of a Subnet Manager (SM) running on the InfiniBand fabric, at all times. An SM can run on any node or on an IB switch. OpenSM is an InfiniBand-compliant Subnet Manager, and it is installed as part of OFED<sup>1</sup>.

1. OpenSM is disabled by default. See "[OpenSM](#)" section for details on enabling it.

## **Diagnostic Utilities**

OFED includes the following two diagnostic packages for use by network and data center managers:

- ibutils – NVIDIA diagnostic utilities
- infiniband-diags – OpenFabrics Alliance InfiniBand diagnostic tools

## **NVIDIA Firmware Tools**

The NVIDIA Firmware Tools (MFT) package is a set of firmware management tools for a single InfiniBand node. MFT can be used for:

- Generating a standard or customized NVIDIA firmware image

- Burning a firmware image to a single InfiniBand node

MFT includes a set of tools used for performing firmware update and configuration, as well as debug and diagnostics, and provides MST service. For the full list of available tools within MFT, please refer to MFT documentation ([docs.nvidia.com/networking/category/mft](https://docs.nvidia.com/networking/category/mft)).

## NVIDIA OFED Package

### ISO Image

OFED for Linux (MLNX\_OFED\_LINUX) is provided as ISO images or as a tarball, one per supported Linux distribution and CPU architecture, that includes *source code* and *binary* RPMs, firmware, utilities, and documentation. The ISO image contains an installation script (called `mlnxofedinstall`) that performs the necessary steps to accomplish the following:

- Discover the currently installed kernel
- Uninstall any InfiniBand stacks that are part of the standard operating system distribution or another vendor's commercial stack
- Install the MLNX\_OFED\_LINUX binary RPMs (if they are available for the current kernel)
- Identify the currently installed InfiniBand HCAs and perform the required firmware updates

### Software Components

MLNX\_OFED\_LINUX contains the following software components:

- NVIDIA Host Channel Adapter Drivers
  - `mlx5`, `mlx4` (VPI), which is split into multiple modules:
    - `mlx4_core` (low-level helper)
    - `mlx4_ib` (IB)
    - `mlx5_ib`

- mlx5\_core (includes Ethernet)
  - mlx4\_en (Ethernet)
- Mid-layer core
  - Verbs, MADs, SA, CM, CMA, uVerbs, uMADs
- Upper Layer Protocols (ULPs)
  - IPoIB, SRP Initiator and SRP
- MPI
  - Open MPI stack supporting the InfiniBand, RoCE and Ethernet interfaces
  - MPI benchmark tests (OSU BW/LAT, Intel MPI Benchmark, Presta)
- OpenSM: InfiniBand Subnet Manager
- Utilities
  - Diagnostic tools
  - Performance tests
  - Sysinfo (see [Sysinfo User Manual](#))
- Firmware tools (MFT)
- Source code for all the OFED software modules (for use under the conditions mentioned in the modules' LICENSE files)
- Documentation

## Firmware

The ISO image includes the following firmware item:

- mlnx-fw-updater RPM/DEB package, which contains firmware binaries for supported devices (using mlxfwmanager tool).

## Directory Structure

The ISO image of MLNX\_OFED\_LINUX contains the following files and directories:

- mlnxofedinstall - This is the MLNX\_OFED\_LINUX installation script.
- ofed\_uninstall.sh - This is the MLNX\_OFED\_LINUX un-installation script.
- <RPMS folders> - Directory of binary RPMs for a specific CPU architecture.
- src/ - Directory of the OFED source tarball.

### **Warning**

MLNX\_OFED includes the OFED source RPM packages used as a build platform for kernel code but does not include the sources of NVIDIA proprietary packages.

- mlnx\_add\_kernel\_support.sh - Script required to rebuild MLNX\_OFED\_LINUX for customized kernel version on supported Linux Distribution
- RPM based - A script required to rebuild MLNX\_OFED\_LINUX for customized kernel version on supported RPM-based Linux Distribution
- docs/ - Directory of NVIDIA OFED related documentation

## Module Parameters

### mlx4 Module Parameters

In order to set **mlx4** parameters, add the following line(s) to **/etc/modprobe.d/mlx4.conf**:

```
options mlx4_core parameter=<value>
```

and/or

```
options mlx4_ib parameter=<value>
```

and/or

```
options mlx4_en parameter=<value>
```

The following sections list the available mlx4 parameters.

### mlx4\_ib Parameters

sm_guid_assign:	Enable SM alias_GUID assignment if sm_guid_assign > 0 (Default: 0) (int)
dev_assign_str: <sup>1</sup>	Map device function numbers to IB device numbers (e.g.'0000:04:00.0-0,002b:1c:0b.a-1,...'). Hexadecimal digits for the device function (e.g. 002b:1c:0b.a) and decimal for IB device numbers (e.g. 1). Max supported devices - 32 (string)
en_ecn	Enable q/ecn [enable = 1, disable = 0 (default)] (bool)

1. In the current version, this parameter is using a decimal number to describe the InfiniBand device, and not a hexadecimal number as in previous versions. The purpose is to uniform the mapping of device function numbers with InfiniBand device numbers, as defined for other module parameters (e.g. num\_vfs and probe\_vf). For example, to map mlx4\_15 to device function number 04:00.0 in the current version, we use "options mlx4\_ib dev\_assign\_str=04:00.0-15", as opposed to the previous version where we used "options mlx4\_ib dev\_assign\_str=04:00.0-f"

### mlx4\_core Parameters

debug_level	Enable debug tracing if > 0 (int)
msi_x	0 - don't use MSI-X, 1 - use MSI-X, >1 - limit number of MSI-X irqs to msi_x (non-SRIOV only) (int)
enable_sys_tune	Tune the cpu's for better performance (default 0) (int)

block_loo pback	Block multicast loopback packets if > 0 (default: 1) (int)
num_vfs	Either a single value (e.g. '5') to define uniform num_vfs value for all devices functions or a string to map device function numbers to their num_vfs values (e.g. '0000:04:00.0-5,002b:1c:0b.a-15'). Hexadecimal digits for the device function (e.g. 002b:1c:0b.a) and decimal for num_vfs value (e.g. 15). (string)
probe_vf	Either a single value (e.g. '3') to indicate that the Hypervisor driver itself should activate this number of VFs for each HCA on the host, or a string to map device function numbers to their probe_vf values (e.g. '0000:04:00.0-3,002b:1c:0b.a-13'). Hexadecimal digits for the device function (e.g. 002b:1c:0b.a) and decimal for probe_vf value (e.g. 13). (string)
log_num _mgm_e ntry_size	log mgm size, that defines the num of qp per mcg, for example: 10 gives 248.range: 7 <= log_num_mgm_entry_size <= 12. To activate device managed flow steering when available, set to -1 (int)
high_rate _steer	Enable steering mode for higher packet rate (obsolete, set "Enable optimized steering" option in log_num_mgm_entry_size to use this mode). (int)
fast_drop	Enable fast packet drop when no recieve WQEs are posted (int)
enable_6 4b_cqe_e qe	Enable 64 byte CQEs/EQEs when the FW supports this if non-zero (default: 1) (int)
log_num _mac	Log2 max number of MACs per ETH port (1-7) (int)
log_num _vlan	(Obsolete) Log2 max number of VLANs per ETH port (0-7) (int)
log_mtts _per_seg	Log2 number of MTT entries per segment (0-7) (default: 0) (int)
port_typ e_array	Either pair of values (e.g. '1,2') to define uniform port1/port2 types configuration for all devices functions or a string to map device function numbers to their pair of port types values (e.g. '0000:04:00.0-1;2,002b:1c:0b.a-1;1'). Valid port types: 1-ib, 2-eth, 3-auto, 4-N/A If only a single port is available, use the N/A port type for port2 (e.g '1,4').



log_num_qp	log maximum number of QPs per HCA (default: 19) (int)
log_num_srq	log maximum number of SRQs per HCA (default: 16) (int)
log_rdma_rc_per_qp	log number of RDMARC buffers per QP (default: 4) (int)
log_num_cq	log maximum number of CQs per HCA (default: 16) (int)
log_num_mcg	log maximum number of multicast groups per HCA (default: 13) (int)
log_num_mpt	log maximum number of memory protection table entries per HCA (default: 19) (int)
log_num_mtt	log maximum number of memory translation table segments per HCA (default: max(20, 2*MTTs for register all of the host memory limited to 30)) (int)
enable_qos	Enable Quality of Service support in the HCA (default: off) (bool)
internal_err_reset	Reset device on internal errors if non-zero (default is 1) (int)
ingress_parser_mode	Mode of ingress parser for ConnectX3-Pro. 0 - standard. 1 - checksum for non TCP/UDP. (default: standard) (int)
roce_mode	Set RoCE modes supported by the port
ud_gid_type	Set gid type for UD QPs
log_num_mgm_entry_size	log mgm size, that defines the num of qp per mcg, for example: 10 gives 248.range: 7 <= log_num_mgm_entry_size <= 12 (default = -10).
use_prio	Enable steering by VLAN priority on ETH ports (deprecated) (bool)
enable_vfs_qos	Enable Virtual VFs QoS (default: off) (bool)

mlx4_en_only_mode	Load in Ethernet only mode (int)
enable_4k_uar	Enable using 4K UAR. Should not be enabled if have VFs which do not support 4K UARs (default: true) (bool)
mlx4_en_only_mode	Load in Ethernet only mode (int)
rr_proto	IP next protocol for RoCEv1.5 or destination port for RoCEv2. Setting 0 means using driver default values (deprecated) (int)

### mlx4\_en Parameters

inline_threshold	The threshold for using inline data (int) Default and max value is 104 bytes. Saves PCI read operation transaction, packet less than threshold size will be copied to hw buffer directly. (range: 17-104)
udp_rss:	Enable RSS for incoming UDP traffic (uint) On by default. Once disabled no RSS for incoming UDP traffic will be done.
pfctx	Priority-based Flow Control policy on TX[7:0]. Per priority bit mask (uint)
pfcrx	Priority-based Flow Control policy on RX[7:0]. Per priority bit mask (uint)
udev_dev_port_dev_id	Work with dev_id or dev_port when supported by the kernel. Range: 0 <= udev_dev_port_dev_id <= 2 (default = 0).
udev_dev_port_dev_id:	Work with dev_id or dev_port when supported by the kernel. Range: 0 <= udev_dev_port_dev_id <= 2 (default = 0). <ul style="list-style-type: none"> <li>• 0: Work with dev_port if supported by the kernel, otherwise work with dev_id.</li> <li>• 1: Work only with dev_id regardless of dev_port support.</li> <li>• 2: Work with both of dev_id and dev_port (if dev_port is supported by the kernel). (int)</li> </ul>

### mlx5\_core Module Parameters

The mlx5\_core module supports a single parameter used to select the profile which defines the number of resources supported.

<i>prof_sel</i>	The parameter name for selecting the profile. The supported values for profiles are: <ul style="list-style-type: none"> <li>• 0 - for medium resources, medium performance</li> <li>• 1 - for low resources</li> <li>• 2 - for high performance (int) (default)</li> </ul>
guids	charp
node_guid	guids configuration. This module parameter will be obsolete!
debug_mask	debug_mask: 1 = dump cmd data, 2 = dump cmd exec time, 3 = both. Default=0 (uint)
probe_vf	probe VFs or not, 0 = not probe, 1 = probe. Default = 1 (bool)
num_of_groups	Controls the number of large groups in the FDB flow table. Default=4; Range=1-1024

### ib\_core Parameters

send_queue_size	Size of send queue in number of work requests (int)
recv_queue_size	Size of receive queue in number of work requests (int)
force_mr	Force usage of MRs for RDMA READ/WRITE operations (bool)
roce_v1_noncompat_gid	Default GID auto configuration (Default: yes) (bool)

### ib\_ipoib Parameters

max_nonsrq_conn_qp	Max number of connected-mode QPs per interface (applied only if shared receive queue is not available) (int)
mcast_debug_level	Enable multicast debug tracing if > 0 (int)
send_queue_size	Number of descriptors in send queue (int)
recv_queue_size	Number of descriptors in receive queue (int)
debug_level	Enable debug tracing if > 0 (int)

ipoib_enhanced	Enable IPoIB enhanced for capable devices (default = 1) (0-1) (int)
----------------	---

## Devlink Parameters

The following parameters, supported in mlx4 driver only, can be changed using the Devlink user interface:

Parameter	Description	Parameter Type
internal_error_reset	Enables resetting the device on internal errors	Generic
max_macs	Max number of MACs per ETH port	Generic
region_snapshot_enable	Enables capturing region snapshots	Generic
enable_64b_cqe_eqe	Enables 64 byte CQEs/EQEs when supported by FW	Driver-specific
enable_4k_uar	Enables using 4K UAR	Driver-specific

## Device Capabilities

Normally, an application needs to query the device capabilities before attempting to create a resource. It is essential for the application to be able to operate over different devices with different capabilities.

Specifically, when creating a QP, the user needs to specify the maximum number of outstanding work requests that the QP supports. This value should not exceed the queried capabilities. However, even when you specify a number that does not exceed the queried capability, the verbs can still fail since some other factors such as the number of scatter/gather entries requested, or the size of the inline data required, affect the maximum possible work requests. Hence an application should try to decrease this size (halving is a good new value) and retry until it succeeds.

# Installation

This chapter describes how to install and test the Mellanox OFED for Linux package on a single host machine with Mellanox InfiniBand and/or Ethernet adapter hardware installed.

The chapter contains the following sections:

- [Hardware and Software Requirements](#)
- [Downloading MLNX\\_OFED](#)
- [Installing MLNX\\_OFED](#)
- [Uninstalling MLNX\\_OFED](#)
- [Updating Firmware After Installation](#)
- [UEFI Secure Boot](#)
- [Performance Tuning](#)

## Features Overview and Configuration

The chapter contains the following sections:

- [Ethernet Network](#)
- [InfiniBand Network](#)
- [Storage Protocols](#)
- [Virtualization](#)
- [Resiliency](#)
- [Docker Containers](#)
- [HPC-X™](#)
- [Fast Driver Unload](#)

- [OVS Offload Using ASAP<sup>2</sup> Direct](#)

# Programming

## Warning

This chapter is aimed for application developers and expert users that wish to develop applications over MLNX\_OFED.

## Raw Ethernet Programming

Raw Ethernet programming enables writing an application that bypasses the kernel stack. To achieve this, packet headers and offload options need to be provided by the application.

For a basic example on how to use Raw Ethernet programming, refer to the [Raw Ethernet Programming: Basic Introduction—Code Example](#) Community post.

## Packet Pacing

Packet pacing is a raw Ethernet sender feature that enables controlling the rate of each QP, per send queue.

For a basic example on how to use packet pacing per flow over libibverbs, refer to [Raw Ethernet Programming: Packet Pacing—Code Example](#) Community post.

## TCP Segmentation Offload (TSO)

TCP Segmentation Offload (TSO) enables the adapter cards to accept a large amount of data with a size greater than the MTU size. The TSO engine splits the data into separate packets and inserts the user-specified L2/L3/L4 headers automatically per packet. With the usage of TSO, CPU is offloaded from dealing with a large throughput of data.

To be able to program that on the sender side, refer to the [Raw Ethernet Programming: TSO—Code Example](#) Community post.

## ToS Based Steering

ToS/DSCP is an 8-bit field in the IP packet that enables different service levels to be assigned to network traffic. This is achieved by marking each packet in the network with a DSCP code and appropriating the corresponding level of service to it.

To be able to steer packets according to the ToS field on the receiver side, refer to the [Raw Ethernet Programming: ToS—Code Example](#) Community post.

## Flow ID Based Steering

Flow ID based steering enables developing a code that will steer packets using flow ID when developing Raw Ethernet over verbs. For more information on flow ID based steering, refer to the [Raw Ethernet Programming: Flow ID Steering—Code Example](#) Community post.

## VXLAN Based Steering

VXLAN based steering enables developing a code that will steer packets using the VXLAN tunnel ID when developing Raw Ethernet over verbs. For more information on VXLAN based steering, refer to the [Raw Ethernet Programming: VXLAN Steering—Code Example](#) Community post.

## Device Memory Programming

### Warning

This feature is supported on ConnectX-5/ConnectX-5 Ex adapter cards and above only.

Device Memory is an API that allows using on-chip memory located on the device as a data buffer for send/receive and RDMA operations. The device memory can be mapped and accessed directly by user and kernel applications, and can be allocated in various sizes, registered as memory regions with local and remote access keys for performing the send/receive and RDMA operations.

Using the device memory to store packets for transmission can significantly reduce transmission latency compared to the host memory.

## Device Memory Programming Model

The new API introduces a similar procedure to the host memory for sending packets from the buffer:

- `ibv_alloc_dm()/ibv_free_dm()` - to allocate/free device memory
- `ibv_reg_dm_mr` - to register the allocated device memory buffer as a memory region and get a memory key for local/remote access by the device
- `ibv_memcpy_to_dm` - to copy data to a device memory buffer
- `ibv_memcpy_from_dm` - to copy data from a device memory buffer
- `ibv_post_send/ibv_post_receive` - to request the device to perform a send/receive operation using the memory key

For examples, see [Device Memory](#).

## RDMA-CM QP Timeout Control

RDMA-CM QP Timeout Control feature enables users to control the QP timeout for QPs created with RDMA-CM.

A new option in 'rdma\_set\_option' function has been added to enable overriding calculated QP timeout, in order to provide QP attributes for QP modification. To achieve that, `rdma_set_option()` should be called with the new flag `RDMA_OPTION_ID_ACK_TIMEOUT`. Example:



```
rdma_set_option(cma_id, RDMA_OPTION_ID, RDMA_OPTION_ID_ACK_TIMEOUT,  
&timeout, sizeof(timeout));
```

## RDMA-CM Application Managed QP

Applications which do not create a QP through `rdma_create_qp()` may want to postpone the ESTABLISHED event on the passive side, to let the active side complete an application-specific connection establishment phase. For example, modifying the init state of the QP created by the application to RTR state, or make some preparations for receiving messages from the passive side. The feature returns a new event on the active side: `CONNECT_RESPONSE`, instead of `ESTABLISHED`, if `id->qp==NULL`. This gives the application a chance to perform the extra connection setup. Afterwards, the new `rdma_establish()` API should be called to complete the connection and generate an ESTABLISHED event on the passive side.

In addition, this feature exposes the 'rdma\_init\_qp\_attr' function in `librdmacm` API, which enables applications to get the parameters for creating Address Handler (AH) or control QP attributes after its creation.

## InfiniBand Fabric Utilities

This section first describes common configuration, interface, and addressing for all the tools in the package.

### Common Configuration, Interface and Addressing

#### Topology File (Optional)

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

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

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

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

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

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

## InfiniBand Interface Definition

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

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

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

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

## Addressing



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

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

- Using a Directed Route to the destination: (Tool option '-d')  
This option defines a directed route of output port numbers from the local port to the destination.
- Using port LIDs: (Tool option '-l'):  
In this mode, the source and destination ports are defined by means of their LIDs. If the fabric is configured to allow multiple LIDs per port, then using any of them is valid for defining a port.
- Using port names defined in the topology file: (Tool option '-n')  
This option refers to the source and destination ports by the names defined in the topology file. (Therefore, this option is relevant only if a topology file is specified to the tool.) In this mode, the tool uses the names to extract the port LIDs from the matched topology, then the tool operates as in the '-l' option.

## Diagnostic Utilities

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

### Diagnostic Utilities

Utility	Description
<b>dumpfts</b>	Dumps tables for every switch found in an <code>ibnetdiscover</code> scan of the subnet. The dump file format is compatible with loading into OpenSM using the <code>-R file -U /path/to/dump-file</code> syntax. For further information, please refer to the tool's man page.
<b>ibaddr</b>	Can be used to show the LID and GID addresses of the specified port or the local port by default. This utility can be used as simple address resolver. For further information, please refer to the tool's man page.

Utility	Description
<b>ibcachedit</b>	Allows users to edit an ibnetdiscover cache created through the --cache option in ibnetdiscover(8). For further information, please refer to the tool's man page.
<b>ibccconfig</b>	Supports the configuration of congestion control settings on switches and HCAs. For further information, please refer to the tool's man page.
<b>ibccquery</b>	Supports the querying of settings and other information related to congestion control. For further information, please refer to the tool's man page.
<b>ibcongest</b>	Provides static congestion analysis. It calculates routing for a given topology (topo-mode) or uses extracted lst/fdb files (lst-mode). Additionally, it analyzes congestion for a traffic schedule provided in a "schedule-file" or uses an automatically generated schedule of all-to-all-shift. To display a help message which details the tool's options, please run <code>"/opt/ibutils2/bin/ibcongest -h"</code> . For further information, please refer to the tool's man page.
<b>ibdev2netdev</b>	Enables association between IB devices and ports and the associated net device. Additionally it reports the state of the net device link. For further information, please refer to the tool's man page.
<b>ibdiagnet (of ibutils2)</b>	Scans the fabric using directed route packets and extracts all the available information regarding its connectivity and devices. An ibdiagnet run performs the following stages: <ul style="list-style-type: none"> <li>• Fabric discovery</li> <li>• Duplicated GUIDs detection</li> <li>• Links in INIT state and unresponsive links detection</li> <li>• Counters fetch</li> <li>• Error counters check</li> <li>• Routing checks</li> <li>• Link width and speed checks</li> <li>• Alias GUIDs check</li> <li>• Subnet Manager check</li> <li>• Partition keys check</li> <li>• Nodes information</li> </ul>

Utility	Description
	<p><b>Note:</b> This version of ibdiagnet is included in the ibutils2 package, and it is run by default after installing Mellanox OFED. To use this ibdiagnet version, run: ibdiagnet.</p> <p>For further information, either:</p> <ol style="list-style-type: none"> <li>1. Run ibdiagnet -H</li> </ol> <p>Or</p> <ol style="list-style-type: none"> <li>2. Refer to <a href="https://docs.nvidia.com/networking/display/ibdiagnetUserManualv10">docs.nvidia.com/networking/display/ibdiagnetUserManualv10</a></li> </ol>
<b>ibdi agp ath</b>	<p>Traces a path between two end-points and provides information regarding the nodes and ports traversed along the path. It utilizes device specific health queries for the different devices along the path.</p> <p>The way ibdiagpath operates depends on the addressing mode used in the command line. If directed route addressing is used (--dr_path flag), the local node is the source node and the route to the destination port is known apriori (for example: ibdiagpath --dr_path 0,1). On the other hand, if LID-route addressing is employed, --src_lid and --dest_lid, then the source and destination ports of a route are specified by their LIDs. In this case, the actual path from the local port to the source port, and from the source port to the destination port, is defined by means of Subnet Management Linear Forwarding Table queries of the switch nodes along that path. Therefore, the path cannot be predicted as it may change.</p> <p>Example: ibdiagpath --src_lid 1 --dest_lid 28</p> <p>For further information, please refer to the tool's -help flag.</p>
<b>ibd um p</b>	<p>Dump InfiniBand traffic that flows to and from Mellanox Technologies ConnectX® family adapters InfiniBand ports.</p> <p>Note the following:</p> <ul style="list-style-type: none"> <li>• ibdump is not supported for Virtual functions (SR-IOV).</li> <li>• Infiniband traffic sniffing is supported on all HCAs.</li> <li>• Ethernet and RoCE sniffing is supported only on Connect-X3 and Connect-X3 Pro cards.</li> </ul> <p>The dump file can be loaded by the Wireshark tool for graphical traffic analysis. The following describes a workflow for local HCA (adapter) sniffing:</p> <ol style="list-style-type: none"> <li>1. Run ibdump with the desired options</li> <li>2. Run the application that you wish its traffic to be analyzed</li> <li>3. Stop ibdump (CTRL-C) or wait for the data buffer to fill (in --mem-mode)</li> <li>4. Open Wireshark and load the generated file</li> </ol>

Utility	Description
	<p>To download Wireshark for a Linux or Windows environment go to <a href="http://www.wireshark.org">www.wireshark.org</a>.</p> <p><b>Note:</b> Although ibdump is a Linux application, the generated .pcap file may be analyzed on either operating system.</p> <p>[mlx4] In order for ibdump to function with RoCE, Flow Steering must be enabled. To do so:</p> <ol style="list-style-type: none"> <li>1. Add the following to /etc/modprobe.d/mlnx.conf file: options mlx4_core log_num_mgm_entry_size=-1</li> <li>2. Restart the drivers.</li> </ol> <p><b>Note:</b> If one of the HCA's port is configured as InfiniBand, ibdump requires IPoIB DMFS to be enabled. For further information, please refer to <a href="#">Flow Steering Configuration</a> section.</p> <p>For further information, please refer to the tool's man page.</p>
<b>iblinkinfo</b>	<p>Reports link info for each port in an InfiniBand fabric, node by node. Option- ally, iblinkinfo can do partial scans and limit its output to parts of a fabric.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibnetdiscover</b>	<p>Performs InfiniBand subnet discovery and outputs a human readable topology file. GUIDs, node types, and port numbers are displayed as well as port LIDs and node descriptions. All nodes (and links) are displayed (full topology).</p> <p>This utility can also be used to list the current connected nodes. The output is printed to the standard output unless a topology file is specified.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibnetsplit</b>	<p>Automatically groups hosts and creates scripts that can be run in order to split the network into sub-networks containing one group of hosts.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibnodes</b>	<p>Uses the current InfiniBand subnet topology or an already saved topology file and extracts the InfiniBand nodes (CAs and switches).</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibping</b>	<p>Uses vendor mads to validate connectivity between InfiniBand nodes. On exit, (IP) ping like output is show. ibping is run as client/server. The default is to run as client. Note also that a default ping server is implemented within the kernel.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibports</b>	<p>Enables querying the logical (link) and physical port states of an InfiniBand port. It also allows adjusting the link speed that is enabled on any InfiniBand port.</p>

Utility	Description
<b>ibportstate</b>	<p>If the queried port is a switch port, then <code>ibportstate</code> can be used to:</p> <ul style="list-style-type: none"> <li>• disable, enable or reset the port</li> <li>• validate the port's link width and speed against the peer port</li> </ul> <p>In case of multiple channel adapters (CAs) or multiple ports without a CA/ port being specified, a port is chosen by the utility according to the following criteria:</p> <ul style="list-style-type: none"> <li>• The first ACTIVE port that is found.</li> <li>• If not found, the first port that is UP (physical link state is LinkUp).</li> </ul> <p>For further information, please refer to the tool's man page.</p>
<b>ibqueryerrors</b>	<p>The default behavior is to report the port error counters which exceed a threshold for each port in the fabric. The default threshold is zero (0). Error fields can also be suppressed entirely.</p> <p>In addition to reporting errors on every port, <code>ibqueryerrors</code> can report the port transmit and receive data as well as report full link information to the remote port if available.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibroute</b>	<p>Uses SMPs to display the forwarding tables—unicast (<code>LinearForwardingTable</code> or <code>LFT</code>) or multicast (<code>MulticastForwardingTable</code> or <code>MFT</code>)—for the specified switch LID and the optional lid (mlid) range. The default range is all valid entries in the range 1 to <code>FDBTop</code>.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibstat</b>	<p><code>ibstat</code> is a binary which displays basic information obtained from the local IB driver. Output includes LID, SMLID, port state, link width active, and port physical state.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibstats</b>	<p>Displays basic information obtained from the local InfiniBand driver. Output includes LID, SMLID, port state, port physical state, port width and port rate. For further information, please refer to the tool's man page.</p>
<b>ibswitches</b>	<p>Traces the InfiniBand subnet topology or uses an already saved topology file to extract the InfiniBand switches.</p> <p>For further information, please refer to the tool's man page.</p>

Utility	Description
<b>ibsysstat</b>	<p>Uses vendor mads to validate connectivity between InfiniBand nodes and obtain other information about the InfiniBand node. ibsysstat is run as client/ server. The default is to run as client.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibtopdiff</b>	<p>Compares a topology file and a discovered listing of subnet.lst/ibdiagnet.lst and reports mismatches.</p> <p>Two different algorithms provided:</p> <ul style="list-style-type: none"> <li>• Using the -e option is more suitable for MANY mismatches it applies less heuristics and provide details about the match</li> <li>• Providing the -s, -p and -g starts a detailed heuristics that should be used when only small number of changes are expected</li> </ul> <p>For further information, please refer to the tool's man page.</p>
<b>ibtrace</b>	<p>Uses SMPs to trace the path from a source GID/LID to a destination GID/ LID. Each hop along the path is displayed until the destination is reached or a hop does not respond. By using the -m option, multicast path tracing can be performed between source and destination nodes.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibv_asyncwatch</b>	<p>Display asynchronous events forwarded to userspace for an InfiniBand device.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibv_devices</b>	<p>Lists InfiniBand devices available for use from userspace, including node GUIDs.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ibv_deviceinfo</b>	<p>Queries InfiniBand devices and prints about them information that is available for use from userspace.</p> <p>For further information, please refer to the tool's man page.</p>
<b>mstflint</b>	<p>Queries and burns a binary firmware-image file on non-volatile (Flash) memories of NVIDIA InfiniBand and Ethernet network adapters. The tool requires root privileges for Flash access.</p> <p>To run mstflint, you must know the device location on the PCI bus.</p>



Utility	Description
	<p><b>Note:</b> If you purchased a standard NVIDIA network adapter card, please download the firmware image from <a href="https://www.nvidia.com/en-us/networking/support/firmware-download">nvidia.com/en-us/networking/Support Support Firmware Download</a>. If you purchased a non-standard card from a vendor other than NVIDIA, please contact your vendor.</p> <p>For further information, please refer to the tool's man page.</p>
<b>per fqu ery</b>	<p>Queries InfiniBand ports' performance and error counters. Optionally, it displays aggregated counters for all ports of a node. It can also reset counters after reading them or simply reset them.</p> <p>For further information, please refer to the tool's man page.</p>
<b>saq uer y</b>	<p>Issues the selected SA query. Node records are queried by default. For further information, please refer to the tool's man page.</p>
<b>smi nfo</b>	<p>Issues and dumps the output of an sminfo query in human readable format. The target SM is the one listed in the local port info or the SM specified by the optional SM LID or by the SM direct routed path.</p> <p><b>Note:</b> Using sminfo for any purpose other than a simple query might result in a malfunction of the target SM.</p> <p>For further information, please refer to the tool's man page.</p>
<b>sm par que ry</b>	<p>Sends SMP query for adaptive routing and private LFT features. For further information, please refer to the tool's man page.</p>
<b>sm pdu mp</b>	<p>A general purpose SMP utility which gets SM attributes from a specified SMA. The result is dumped in hex by default.</p> <p>For further information, please refer to the tool's man page.</p>
<b>sm pqu ery</b>	<p>Provides a basic subset of standard SMP queries to query Subnet management attributes such as node info, node description, switch info, and port info.</p> <p>For further information, please refer to the tool's man page.</p>

## Link Level Retransmission (LLR) in FDR Links

With the introduction of FDR 56 Gbps technology, Mellanox enabled a proprietary technology called LLR (Link Level Retransmission) to improve the reliability of FDR links.

This proprietary LLR technology adds additional CRC checking to the data stream and retransmits portions of packets with CRC errors at the local link level. Customers should be aware of the following facts associated with LLR technology:

- Traditional methods of checking the link health can be masked because the LLR technology automatically fixes errors. The traditional IB symbol error counter will show no errors when LLR is active.
- Latency of the fabric can be impacted slightly due to LLR retransmissions. Traditional IB performance utilities can be used to monitor any latency impact.
- Bandwidth of links can be reduced if cable performance degrades and LLR retransmissions become too numerous. Traditional IB bandwidth performance utilities can be used to monitor any bandwidth impact.

Due to these factors, an LLR retransmission rate counter has been added to the `ibdiagnet` utility that can give end users an indication of the link health.

➤ *To monitor LLR retransmission rate:*

1. Run `ibdiagnet`, no special flags required.
2. If the LLR retransmission rate limit is exceeded it will print to the screen.
3. The default limit is set to 500 and requires further investigation if exceeded.
4. The LLR retransmission rate is reflected in the results file `/var/tmp/ibdiagnet2/ibdiagnet2.pm`.

The default value of 500 retransmissions/sec has been determined by Mellanox based on the extensive simulations and testing. Links exhibiting a lower LLR retransmission rate should not raise special concern.

## Performance Utilities

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

Utility	Description
<b>ib_atomic_bw</b>	<p>Calculates the BW of RDMA Atomic transactions between a pair of machines. One acts as a server and the other as a client. The client RDMA sends atomic operation to the server and calculate the BW by sampling the CPU each time it receive a successful completion. The test supports features such as Bidirectional, in which they both RDMA atomic to each other at the same time, change of MTU size, tx size, number of iteration, message size and more. Using the "-a" flag provides results for all message sizes.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ib_atomic_latency</b>	<p>Calculates the latency of RDMA Atomic transaction of message_size between a pair of machines. One acts as a server and the other as a client. The client sends RDMA atomic operation and sample the CPU clock when it receives a successful completion, in order to calculate latency.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ib_read_bw</b>	<p>Calculates the BW of RDMA read between a pair of machines. One acts as a server and the other as a client. The client RDMA reads the server memory and calculate the BW by sampling the CPU each time it receive a successful completion. The test supports features such as Bidirectional, in which they both RDMA read from each other memory's at the same time, change of MTU size, tx size, number of iteration, message size and more.</p> <p>Read is available only in RC connection mode (as specified in IB spec). For further information, please refer to the tool's man page.</p>
<b>ib_read_latency</b>	<p>Calculates the latency of RDMA read operation of message_size between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which one side RDMA reads the memory of the other side only after the other side have read his memory. Each of the sides samples the CPU clock each time they read the other side memory , in order to calculate latency. Read is available only in RC connection mode (as specified in IB spec).</p> <p>For further information, please refer to the tool's man page.</p>
<b>ib_send_bw</b>	<p>Calculates the BW of SEND between a pair of machines. One acts as a server and the other as a client. The server receive packets from the client and they both calculate the throughput of the operation. The test supports features such as Bidirectional, on which they both send and receive at the same time, change of MTU size, tx size, number of iteration, message size and more. Using the "-a" provides results for all message sizes.</p> <p>For further information, please refer to the tool's man page.</p>

Utility	Description
<b>ib_send_lat</b>	<p>Calculates the latency of sending a packet in message_size between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which you send packet only if you receive one. Each of the sides samples the CPU each time they receive a packet in order to calculate the latency. Using the "-a" provides results for all message sizes.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ib_write_bw</b>	<p>Calculates the BW of RDMA write between a pair of machines. One acts as a server and the other as a client. The client RDMA writes to the server memory and calculates the BW by sampling the CPU each time it receives a successful completion. The test supports features such as Bidirectional, in which they both RDMA write to each other at the same time, change of MTU size, tx size, number of iteration, message size and more. Using the "-a" flag provides results for all message sizes.</p> <p>For further information, please refer to the tool's man page.</p>
<b>ib_write_lat</b>	<p>Calculates the latency of RDMA write operation of message_size between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which one side RDMA writes to the other side memory only after the other side wrote on his memory. Each of the sides samples the CPU clock each time they write to the other side memory, in order to calculate latency.</p> <p>For further information, please refer to the tool's man page.</p>
<b>raw_ethernet_bw</b>	<p>Calculates the BW of SEND between a pair of machines. One acts as a server and the other as a client. The server receive packets from the client and they both calculate the throughput of the operation. The test supports features such as Bidirectional, on which they both send and receive at the same time, change of MTU size, tx size, number of iteration, message size and more. Using the "-a" provides results for all message sizes.</p> <p>For further information, please refer to the tool's man page.</p>
<b>raw_ethernet_lat</b>	<p>Calculates the latency of sending a packet in message_size between a pair of machines. One acts as a server and the other as a client. They perform a ping pong benchmark on which you send packet only if you receive one. Each of the sides samples the CPU each time they receive a packet in order to calculate the latency. Using the "-a" provides results for all message sizes.</p> <p>For further information, please refer to the tool's man page.</p>

# Troubleshooting

You may be able to easily resolve the issues described in this section. If a problem persists and you are unable to resolve it yourself, please contact your Mellanox representative or Mellanox Support at [support@mellanox.com](mailto:support@mellanox.com).

The chapter contains the following sections:

- [General Issues](#)
- [Ethernet Related Issues](#)
- [InfiniBand Related Issues](#)
- [Installation Related Issues](#)
- [Performance Related Issues](#)
- [SR-IOV Related Issues](#)
- [PXE \(FlexBoot\) Related Issues](#)
- [RDMA Related Issues](#)
- [Debugging Related Issues](#)
- [OVS Offload Using ASAP2 Direct Related Issues](#)

## Common Abbreviations and Related Documents

### Common Abbreviations and Acronyms

Abbreviation/ Acronym	Description
B	(Capital) 'B' is used to indicate size in bytes or multiples of bytes (e.g., 1KB = 1024 bytes, and 1MB = 1048576 bytes)

Abbreviation/ Acronym	Description
b	(Small) 'b' is used to indicate size in bits or multiples of bits (e.g., 1Kb = 1024 bits)
FW	Firmware
HCA	Host Channel Adapter
HW	Hardware
IB	InfiniBand
iSER	iSCSI RDMA Protocol
LSB	Least significant <i>byte</i>
lsb	Least significant <i>bit</i>
MSB	Most significant <i>byte</i>
msb	Most significant <i>bit</i>
NIC	Network Interface Card
SW	Software
VPI	Virtual Protocol Interconnect
IPoIB	IP over InfiniBand
PFC	Priority Flow Control
PR	Path Record
RoCE	RDMA over Converged Ethernet
SL	Service Level
SRP	SCSI RDMA Protocol
MPI	Message Passing Interface
QoS	Quality of Service
ULP	Upper Layer Protocol
VL	Virtual Lane
vHBA	Virtual SCSI Host Bus Adapter
uDAPL	User Direct Access Programming Library

## Glossary

The following is a list of concepts and terms related to InfiniBand in general and to Subnet Managers in particular. It is included here for ease of reference, but the main reference remains the *InfiniBand Architecture Specification*.

Term	Description
Channel Adapter (CA), Host Channel Adapter (HCA)	An IB device that terminates an IB link and executes transport functions. This may be an HCA (Host CA) or a TCA (Target CA)
HCA Card	A network adapter card based on an InfiniBand channel adapter device
IB Devices	An integrated circuit implementing InfiniBand compliant communication
IB Cluster/Fabric/Subnet	A set of IB devices connected by IB cables
In-Band	A term assigned to administration activities traversing the IB connectivity only
Local Identifier (ID)	An address assigned to a port (data sink or source point) by the Subnet Manager, unique within the subnet, used for directing packets within the subnet
Local Device/Node/System	The IB Host Channel Adapter (HCA) Card installed on the machine running IBDIAG tools
Local Port	The IB port of the HCA through which IBDIAG tools connect to the IB fabric
Master Subnet Manager	The Subnet Manager that is authoritative, that has the reference configuration information for the subnet
Multicast Forwarding Tables	A table that exists in every switch providing the list of ports to forward received multicast packet. The table is organized by MLID
Network Interface Card (NIC)	A network adapter card that plugs into the PCI Express slot and provides one or more ports to an Ethernet network
Standby Subnet Manager	A Subnet Manager that is currently quiescent, and not in the role of a Master Subnet Manager, by the agency of the master SM

Term	Description
Subnet Administrator (SA)	An application (normally part of the Subnet Manager) that implements the interface for querying and manipulating subnet management data
Subnet Manager (SM)	One of several entities involved in the configuration and control of the IB fabric
Unicast Linear Forwarding Tables (LFT)	A table that exists in every switch providing the port through which packets should be sent to each LID
Virtual Protocol Interconnect (VPI)	A Mellanox Technologies technology that allows Mellanox channel adapter devices (ConnectX®) to simultaneously connect to an InfiniBand subnet and a 10GigE subnet (each subnet connects to one of the adapter ports)

## Related Documentation

Document Name	Description
InfiniBand Architecture Specification, Vol. 1, Release 1.2.1	The InfiniBand Architecture Specification that is provided by IBTA
IEEE Std 802.3ae™-2002 (Amendment to IEEE Std 802.3-2002) Document # PDF: SS94996	Part 3: Carrier Sense Multiple Access with Collision Detection (CSMA/CD) Access Method and Physical Layer Specifications Amendment: Media Access Control (MAC) Parameters, Physical Layers, and Management Parameters for 10 Gb/s Operation
Firmware Release Notes for Mellanox adapter devices	See the <a href="#">Release Notes</a> relevant to your adapter device
MFT User Manual and Release Notes	Mellanox Firmware Tools (MFT) User Manual and Release Notes documents
WinOF User Manual	Mellanox WinOF User Manual describes the installation, configuration, and operation of Mellanox Windows driver
VMA User Manual	Mellanox VMA User Manual describes the installation, configuration, and operation of Mellanox VMA driver



---

# Documentation History

- [Release Notes History](#)
- [User Manual Revision History](#)

undefined

© Copyright 2023, NVIDIA. PDF Generated on 06/05/2024