

NVIDIA® MLNX_OFED Documentation Rev 5.2-1.0.4.0

Table of Contents

Release Notes	11
Supported NICs Speeds	11
General Support in MLNX_OFED	14
MLNX_OFED Supported Operating Systems	14
Supported Non-Linux Virtual Machines	16
Support in ASAP2™	17
ASAP2 Supported Operating Systems	17
ASAP2 Requirements	18
ASAP2 Supported Adapter Cards	18
NFS over RDMA (NFSoRDMA) Supported Operating Systems	18
NFSoRDMA Supported Operating Systems on Client and Target Sides	18
NVMEoF Supported Operating Systems	19
Lustre Versions Compatible with MLNX_OFED	19
NEO-Host Supported Operating Systems	19
GPUDirect Storage (GDS) Supported Operating Systems	19
Hardware and Software Requirements	19
Supported NICs Firmware Versions	20
MLNX_OFED Unsupported Functionalities/Features/NICs	21
Changes and New Features	21
MLNX_OFED New Features	21
API Changes in MLNX_OFED	24
MLNX_OFED Verbs API Migration	24
Known Issues	24
Bug Fixes	40
Introduction	65
Stack Architecture	66
Mellanox OFED Package	69
Module Parameters	70
Device Capabilities	71
Installation	72
Hardware and Software Requirements	72
Downloading Mellanox OFED	72

Installing Mellanox OFED	73
Installation Script	73
Installation Procedure	75
Installation Results	78
Installation Logging	78
Driver Load Upon System Boot	79
mlnxofedinstall Return Codes	79
Additional Installation Procedures	79
Installing MLNX_OFED Using YUM	79
Installing MLNX_OFED Using apt-get	82
Installing NEO-Host Using mlnxofedinstall Script	84
Uninstalling Mellanox OFED	84
Uninstalling Mellanox OFED Using the YUM and apt-get Tools	84
Updating Firmware After Installation	85
Updating the Device Online	85
Updating the Device Manually	85
Updating the Device Firmware Automatically upon System Boot	86
UEFI Secure Boot	86
Enrolling Mellanox's x.509 Public Key On your Systems	86
Removing Signature from Kernel Modules	87
Performance Tuning	88
Features Overview and Configuration	89
Ethernet Network	89
Ethernet Interface	89
Port Type Management/VPI Cards Configuration	89
Counters	90
Persistent Naming	91
Interrupt Request (IRQ) Naming	91
Quality of Service (QoS)	92
Mapping Traffic to Traffic Classes	93
Plain Ethernet Quality of Service Mapping	93
RoCE Quality of Service Mapping	94
Map Priorities with set_egress_map	94
Quality of Service Properties	94

Quality of Service Tools	96
Packet Pacing	99
Ethtool	101
Checksum Offload	105
Ignore Frame Check Sequence (FCS) Errors	105
RDMA over Converged Ethernet (RoCE)	106
RoCE Modes	106
GID Table Population	107
RoCE Lossless Ethernet Configuration	109
Installing and Loading the Driver	109
Type Of Service (ToS)	112
RoCE LAG	112
Disabling RoCE	113
Enabling/Disabling RoCE on VMs via VFs	113
Force DSCP	113
Force Time to Live (TTL)	114
Flow Control	114
Priority Flow Control (PFC)	114
Dropless Receive Queue (RQ)	118
Explicit Congestion Notification (ECN)	119
Enabling ECN	119
RSS Support	120
RSS Hash Function	120
Time-Stamping	121
Time-Stamping Service	121
RoCE Time-Stamping	124
One Pulse Per Second (1PPS)	125
Flow Steering	125
Flow Steering Support	125
Flow Domains and Priorities	125
Ethtool	126
Accelerated Receive Flow Steering (aRFS)	127
Flow Steering Dump Tool	128
Wake-on-LAN (WoL)	

	Hardware Accelerated 802.1ad VLAN (Q-in-Q Tunneling)	128
	VLAN Stripping in Linux Verbs	129
	Offloaded Traffic Sniffer	129
	Dump Configuration	129
	Dump Parameters (Bitmap Flag)	. 130
	Configuration	. 130
	Local Loopback Disable	131
	Kernel Transport Layer Security (kTLS) Offloads	132
	Overview	. 132
	Establishing a kTLS Connection	. 132
	Kernel Support	. 132
	Configuring kTLS Offloads	. 133
	OpenSSL with kTLS Offload	. 133
	IPsec Crypto Offload	133
	Overview and Configuration	. 133
	Configuring Security Associations for IPsec Offloads	. 134
ln [°]	finiBand Network	
	InfiniBand Interface	134
	InfiniBand Interface Port Type Management	
		. 134
	Port Type Management	. 134
	Port Type ManagementRDMA Counters	. 134 . 134 135
	Port Type Management RDMA Counters OpenSM	. 134 . 134 135 . 135
	Port Type Management RDMA Counters OpenSM opensm	. 134 . 134 135 . 135 . 136
	Port Type Management RDMA Counters OpenSM opensm osmtest	. 134 . 134 . 135 . 135 . 136
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms	. 134 . 135 . 135 . 136 . 136 . 139
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms Unicast Routing Cache	.134 .135 .135 .136 .136 .139 .139
	Port Type Management	.134 .135 .135 .136 .136 .139 .139 .156
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms Unicast Routing Cache Quality of Service Management in OpenSM Adaptive Routing Manager and SHIELD	.134 .135 .135 .136 .136 .139 .156 .156
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms Unicast Routing Cache Quality of Service Management in OpenSM Adaptive Routing Manager and SHIELD IB Router Support in OpenSM	.134 .135 .135 .136 .136 .139 .156 .156
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms Unicast Routing Cache Quality of Service Management in OpenSM Adaptive Routing Manager and SHIELD IB Router Support in OpenSM. OpenSM Activity Report	. 134 . 135 . 135 . 136 . 136 . 139 . 139 . 156 . 156 . 167 . 167
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms Unicast Routing Cache Quality of Service Management in OpenSM Adaptive Routing Manager and SHIELD IB Router Support in OpenSM OpenSM Activity Report Offsweep Balancing.	. 134 . 135 . 135 . 136 . 136 . 139 . 139 . 156 . 156 . 167 . 167 . 168
	Port Type Management RDMA Counters OpenSM opensm osmtest Partitions Effect of Topology Changes Routing Algorithms Unicast Routing Cache Quality of Service Management in OpenSM Adaptive Routing Manager and SHIELD IB Router Support in OpenSM. OpenSM Activity Report	. 134 . 135 . 135 . 136 . 136 . 139 . 139 . 156 . 156 . 167 . 167 . 168

Supported Policy	171
CMA Features	172
IP over InfiniBand (IPoIB)	172
Upper Layer Protocol (ULP)	172
Enhanced IPoIB	173
IPoIB Mode Setting	173
Port Configuration	174
IPoIB Configuration	174
Sub-interfaces	177
Verifying IPoIB Functionality	178
Bonding IPoIB	178
Dynamic PKey Change	179
Precision Time Protocol (PTP) over IPoIB	179
One Pulse Per Second (1PPS) over IPoIB	180
Advanced Transport	180
Atomic Operations	180
XRC - eXtended Reliable Connected Transport Service for InfiniBand	181
Dynamically Connected Transport (DCT)	181
MPI Tag Matching and Rendezvous Offloads	182
Optimized Memory Access	182
Memory Region Re-registration	182
Memory Window	183
User-Mode Memory Registration (UMR)	184
On-Demand-Paging (ODP)	184
Inline-Receive	185
Mellanox PeerDirect®	185
Mellanox PeerDirect Async	186
Mellanox Relaxed Ordering (RSYNC)	186
CPU Overhead Distribution	186
Out-of-Order (OOO) Data Placement	186
Overview	
IB Router	
MAD Congestion Control	188
Storage Protocols	189

SRP - SCSI RDMA Protocol	189
SRP Initiator	190
Shutting Down SRP	197
iSCSI Extensions for RDMA (iSER)	198
iSER Initiator	198
iSER Targets	199
Lustre	199
NVME-oF - NVM Express over Fabrics	200
NVME-oF	200
NVME-oF Target Offload	200
Virtualization	200
Single Root IO Virtualization (SR-IOV)	200
System Requirements	201
Setting Up SR-IOV	201
Configuring SR-IOV (Ethernet)	202
Configuring SR-IOV (InfiniBand)	202
Additional SR-IOV Configurations	205
Uninstalling the SR-IOV Driver	213
SR-IOV Live Migration	214
Enabling Paravirtualization	229
VXLAN Hardware Stateless Offloads	229
Enabling VXLAN Hardware Stateless Offloads	230
Important Note	231
Q-in-Q Encapsulation per VF in Linux (VST)	231
Setup	231
Prerequisites	231
Configuring Q-in-Q Encapsulation per Virtual Function for ConnectX-5/ConnectX-6	232
802.1Q Double-Tagging	232
Configuring 802.1Q Double-Tagging per Virtual Function	233
Resiliency	233
Reset Flow	233
Kernel ULPs	234
User Space Applications (IB/RoCE)	234
SR-I0V	234

Forcing the VF to Reset	234
Extended Error Handling (EEH)	235
CRDUMP	235
Firmware Tracer	235
Docker Containers	236
Docker Using SR-IOV	236
Kubernetes Using SR-IOV	236
Kubernetes with Shared HCA	237
Mediated Devices	237
Configuring Mediated Device	237
HPC-X™	238
Fast Driver Unload	238
OVS Offload Using ASAP ² Direct	238
Overview	238
Installing OVS-Kernel ASAP ² Packages	239
Installing OVS-DPDK ASAP ² Packages	239
Setting Up SR-IOV	239
OVS Hardware Offloads Configuration	240
OVS-Kernel Hardware Offloads	240
OVS-DPDK Hardware Offloads	258
VirtIO Acceleration through Hardware vDPA	265
Hardware vDPA Installation	265
Hardware vDPA Configuration	266
Running Hardware vDPA	267
Appendix: Mellanox Firmware Tools	268
Programming	270
Raw Ethernet Programming	270
Packet Pacing	270
TCP Segmentation Offload (TSO)	270
ToS Based Steering	270
Flow ID Based Steering	270
VXLAN Based Steering	271
Device Memory Programming	271
Device Memory Programming Model	271

RDMA-CM QP Timeout Control	271
RDMA-CM Application Managed QP	272
InfiniBand Fabric Utilities	273
Common Configuration, Interface and Addressing	273
Topology File (Optional)	273
InfiniBand Interface Definition	273
Addressing	273
Diagnostic Utilities	274
Link Level Retransmission (LLR) in FDR Links	279
Performance Utilities	279
Troubleshooting	282
General Issues	
Ethernet Related Issues	283
InfiniBand Related Issues	284
Installation Related Issues	285
Installation Issues	285
Fixing Application Binary Interface (ABI) Incompatibility with MLNX_OFI	ED Kernel
Modules	286
Overview	286
Performance Related Issues	288
SR-IOV Related Issues	288
PXE (FlexBoot) Related Issues	289
RDMA Related Issues	289
Debugging Related Issues	290
OVS Offload Using ASAP2 Direct Related Issues	290
Common Abbreviations and Related Documents	291
User Manual Revision History	294
Release Notes Change Log History	296

Overview

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

Mellanox OFED (MLNX_OFED) is a Mellanox 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 Mellanox 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 $\underline{\text{http://www.mellanox.com}} \rightarrow \text{Products} \rightarrow \text{Software} \rightarrow \text{InfiniBand/VPI Drivers} \rightarrow \text{Linux SW/Drivers}$

Document Revision History

For the list of changes made to the User Manual, refer to <u>User Manual Revision History</u>.

For the list of changes made to the Release Notes, refer to Release Notes Revision History.

Release Notes

Release Notes Update History

Revision	Date	Description
5.2-1.0.4.0	January 4, 2021	Initial release of this document version.

As of MLNX_OFED version v5.1-0.6.6.0, the following are no longer supported.

- ConnectX-3
- ConnectX-3 Pro
- Connect-IB
- RDMA experimental verbs libraries (mlnx_lib)

Users who wish to utilize the above devices/libraries are advised to refer to MLNX_OFED 4.9 long-term support (LTS) version.

Supported NICs Speeds

These are the release notes of MLNX_OFED for Linux Driver, which operates across all Mellanox network adapter solutions supporting the following uplinks to servers:

Uplink/NICs	Driver Name	Uplink Speed
ConnectX-4	mlx5	 InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 56GbE¹, 100GbE
ConnectX-4 Lx		• Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE
ConnectX-5/ConnectX-5 Ex		 InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 100GbE
ConnectX-6		InfiniBand: SDR, FDR, EDR, HDR
		 Ethernet: 10GbE, 25GbE, 40GbE, 50GbE², 100GbE², 200GbE²
ConnectX-6 Dx		• Ethernet: 10GbE, 25GbE, 40GbE, 50GbE ² , 100GbE ² , 200GbE ²
ConnectX-6 Lx		• Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE ²
BlueField		 InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 100GbE
BlueField-2		 InfiniBand: SDR, FDR, EDR, HDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE², 100GbE²

- 56GbE is a Mellanox propriety link speed and can be achieved while connecting a Mellanox adapter card to Mellanox SX10XX switch series, or connecting a Mellanox adapter card to another Mellanox adapter card.
- 2. Supports both NRZ and PAM4 modes.

Package Contents

Package Revision		Licenses	
ar_mgr	1.0-0.3.MLNX20201210.g3fdac47.5210	Mellanox Confidential and Proprietary	
clusterkit	1.0.36-1.52104	BSD	
dapl	2.1.10.1.mlnx-OFED.4.9.0.1.4.52104	Dual GPL/BSD/CPL	
dpcp	1.1.0-1.52104	Proprietary	
dump_pr	1.0-0.5.MLNX20201210.g3fdac47.5210	GPLv2 or BSD	
fabric-collector	1.1.0.MLNX20170103.89bb2aa-0.1.521 04	GPLv2 or BSD	
hcoll	4.7.3189-1.52104	Proprietary	
ibdump	6.0.0-1.52104	BSD2+GPL2	
ibsim	0.10-1.52104	GPLv2 or BSD	
2.1.1-0.128.MLNX20201210.g311e385. 52104		Mellanox Confidential and Proprietary	
iser	5.2-OFED.5.2.1.0.3.1	GPLv2	
isert	5.2-OFED.5.2.1.0.3.1	GPLv2	
kernel-mft	4.16.0-105	Dual BSD/GPL	
knem	1.1.4.90mlnx1-OFED.5.1.2.5.0.1	BSD and GPLv2	
libvma	9.2.2-1	GPLv2 or BSD	
mlnx-dpdk	20.11.0-1.52104	BSD-3-Clause AND GPL-2.0-only AND LGPL-2.1-only	
mlnx-en	5.2-1.0.4.0.g13e6dd1	GPLv2	
mlnx-ethtool	5.8-1.52104	GPL	
mlnx-iproute2	5.8.0-1.52104	GPL	
mlnx-nfsrdma	5.2-OFED.5.2.1.0.3.1	GPLv2	
mlnx-nvme	5.2-OFED.5.2.1.0.3.1	GPLv2	

Package	Revision	Licenses
mlnx-ofa_kernel	5.2-OFED.5.2.1.0.4.1	GPLv2
mpi-selector	1.0.3-1.52104	BSD
mpitests	3.2.20-5d20b49.52104	BSD
mstflint	4.15.0-1.52104	GPL/BSD
multiperf	3.0-0.14.g5f0fd0e.52104	BSD 3-Clause, GPL v2 or later
mxm	3.7.3112-1.52104	Proprietary
ofed-docs	5.2-OFED.5.2.1.0.4	GPL/BSD
ofed-scripts	5.2-OFED.5.2.1.0.4	GPL/BSD
openmpi	4.1.0rc5-1.52104	BSD
opensm	5.8.0.MLNX20201210.e0a89f5-0.1.521 04	GPLv2 or BSD
openvswitch	2.14.1-1.52104	ASL 2.0 and LGPLv2+ and SISSL
perftest	4.4-0.35.ga08c536.52104	BSD 3-Clause, GPL v2 or later
rdma-core	52mlnx1-1.52104	GPLv2 or BSD
rshim	2.0.5-8.g31b4357	GPLv2
sharp	2.4.1.MLNX20201210.efcfe48-1.52104	Proprietary
sockperf	3.7-0.gita1e8e835a689.52104	BSD
srp	5.2-OFED.5.2.1.0.3.1	GPLv2
исх	1.10.0-1.52104	BSD
xpmem	2.6.3-1.52104	GPLv2 and LGPLv2.1

Release Notes contain the following sections:

- General Support in MLNX_OFED
 Changes and New Features
 Known Issues
 Bug Fixes

General Support in MLNX_OFED

MLNX_OFED Supported Operating Systems

Operating System	Platform	Default Kernel Version
BCLinux 7.5	x86_64	3.10.0-862.el7.x86_64
BCLinux 7.6	x86_64	3.10.0-957.el7.x86_64
BCLinux 8.1	x86_64	4.19.0-193.1.3.el8.bclinux.x86_64
RHEL/CentOS 7.2	x86_64	3.10.0-327.el7.x86_64
RHEL/CentOS 7.3	x86_64	3.10.0-514.el7.x86_64
RHEL/CentOS 7.4	x86_64	3.10.0-693.el7.x86_64
	PPC64	3.10.0-693.el7.ppc64
	PPC64LE	3.10.0-693.el7.ppc64le
RHEL 7.4 ALT	Aarch64	4.11.0-44.el7a.aarch64
RHEL/CentOS 7.5	x86_64	3.10.0-862.el7.x86_64
	PPC64	3.10.0-862.el7.ppc64
	PPC64LE	3.10.0-862.el7.ppc64le
RHEL 7.5 ALT	Aarch64	4.14.0-49.el7a.aarch64
RHEL/CentOS 7.6	x86_64	3.10.0-957.el7.x86_64
	PPC64	3.10.0-957.el7.ppc64
	PPC64LE	3.10.0-957.el7.ppc64le
RHEL 7.6 ALT	Aarch64	4.14.0-115.el7a.0.1.aarch64
	PPC64LE	4.14.0-115.el7a.ppc64le
RHEL/CentOS 7.7	x86_64	3.10.0-1062.el7.x86_64
	PPC64	3.10.0-1062.el7.ppc64
	PPC64LE	3.10.0-1062.el7.ppc64le
RHEL/CentOS 7.8	x86_64	3.10.0-1127.el7.x86_64
	PPC64	3.10.0-1127.el7.ppc64
	PPC64LE	3.10.0-1127.el7.ppc64le
RHEL/CentOS 7.9	x86_64	3.10.0-1160.el7.x86_64
	PPC64	3.10.0-1160.el7.ppc64
	PPC64LE	3.10.0-1160.el7.ppc64le
RHEL/CentOS 8.0	x86_64	4.18.0-80.el8.x86_64
	Aarch64	4.18.0-80.el8.aarch64
	PPC64LE	4.18.0-80.el8.ppc64le
RHEL/CentOS 8.1	x86_64	4.18.0-147.el8.x86_64

Operating System	Platform	Default Kernel Version
	Aarch64	4.18.0-147.el8.aarch64
	PPC64LE	4.18.0-147.el8.ppc64le
RHEL/CentOS 8.2	x86_64	4.18.0-193.el8.x86_64
	Aarch64	4.18.0-193.el8.aarch64
	PPC64LE	4.18.0-193.el8.ppc64le
RHEL/CentOS 8.3	x86_64	4.18.0-240.el8.x86_64
	Aarch64	4.18.0-240.el8.aarch64
	PPC64LE	4.18.0-240.el8.ppc64le
Debian 9.9	x86_64	4.9.0-9-amd64
	Aarch64	4.9.0-9-arm64
Debian 9.11	x86_64	4.9.0-11-amd64
	Aarch64	4.9.0-11-arm64
Debian 10.0	x86_64	4.19.0-5-amd64
	Aarch64	4.19.0-5-arm64
Debian 10.3	x86_64	4.19.0-8-amd64
	Aarch64	4.19.0-8-arm64
Debian 10.5	x86_64	4.19.0-10-amd64
Fedora 32	x86_64	5.6.6-300.fc32.x86_64
OL 7.8	x86_64	4.14.35-1902.300.11.el7uek.x86_64
OL 7.9	x86_64	5.4.17-2011.6.2.el7uek.x86_64
OL 8.2	x86_64	5.4.17-2011.1.2.el8uek.x86_64
OL 8.3	x86_64	5.4.17-2011.7.4.el8uek.x86_64
SLES12 SP2	x86_64	4.4.21-69-default
SLES12 SP3	x86_64	4.4.73-5-default
	PPC64LE	4.4.73-5-default
SLES12 SP4	x86_64	4.12.14-94.41-default
	Aarch64	4.12.14-94.41-default
	PPC64LE	4.12.14-94.41-default
SLES12 SP5	x86_64	4.12.14-120-default
	Aarch64	4.12.14-120-default
	PPC64LE	4.12.14-120-default
SLES15 SP1	x86_64	4.12.14-195-default
	Aarch64	4.12.14-195-default
	PPC64LE	4.12.14-195-default
SLES15 SP2	x86_64	5.3.18-22-default
	Aarch64	5.3.18-22-default

Operating System	Platform	Default Kernel Version
	PPC64LE	5.3.18-22-default
Ubuntu 14.04	x86_64	3.13.0-27-generic
Ubuntu 16.04	x86_64	4.4.0-22-generic
	PPC64LE	4.4.0-21-generic
Ubuntu 18.04	x86_64	4.15.0-20-generic
	Aarch64	4.15.0-20-generic
	PPC64LE	4.15.0-20-generic
Ubuntu 20.04	x86_64	5.4.0-26-generic
	Aarch64	5.4.0-26-generic
	PPC64LE	5.4.0-26-generic
Ubuntu 20.10	x86_64	5.8.0-26-generic
Euler 2.0 SP5	x86_64	3.10.0-862.14.0.1.h43.eulerosv2r7.x86_64
Euler 2.0 SP8	Aarch64	4.19.36-vhulk1907.1.0.h748.eulerosv2r8.aarch64
Euler 2.0 SP9	x86_64	4.18.0-147.5.1.0.h269.eulerosv2r9.x86_64
	Aarch64	4.19.90-vhulk2006.2.0.h171.eulerosv2r9.aarch64
OpenEuler 20.03	x86_64	4.19.90-2003.4.0.0036.oe1.x86_64
	Aarch64	4.19.90-2003.4.0.0036.oe1.aarch64
AliOS 7.2	Aarch64	4.19.48-006.ali4000.alios7.aarch64
Kernel 5.10	x86_64	5.10

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-4	Windows 2012 R2 DC	MLNX_WinOF2 2.50	IB, IPoIB, ETH
ConnectX-4 Lx	Windows 2016 DC	MLNX_WinOF2 2.50	IB, IPoIB, ETH
ConnectX-5 family	All Windows server editions	MLNX_WinOF2 2.50	IPoIB, ETH

NIC	Windows Virtual Machine Type	WinOF version	Protocol
ConnectX-6 family		MLNX_WinOF2 2.50	IPoIB, ETH

Support in ASAP 2 TM

ASAP² Supported Operating Systems

OVS-Kernel SR-IOV Based Supported OSs

Below is a list of all the OSs that support OVS-Kernel ASAP² in the current MLNX_OFED package.

- BCLinux 7.4
- BCLinux 7.5
- BCLinux 7.6
- RHEL/CentOS 7.4
- RHEL/CentOS 7.5
- RHEL/CentOS 7.6
- RHEL/CentOS 7.7
- RHEL/CentOS 7.8
- RHEL/CentOS 8.0
- RHEL/CentOS 8.1
- RHEL/CentOS 8.2
- Fedora 31
- OL 7.4
- OL 7.6
- OL 7.7
- OL 7.8
- OL 8.1
- OL 8.2
- SLES12 SP4 • SLES12 SP5
- SLES15 SP1 • SLES15 SP2
- Ubuntu 16.04
- Ubuntu 18.04
- Ubuntu 19.04
- Ubuntu 20.04
- Kernel 5.7

OVS-DPDK SR-IOV Based Supported OSs

Below is a list of all the OSs that support OVS-DPDK ASAP² in the current MLNX_OFED package.

Adapter Card Type	Supported OSs
ConnectX	 RHEL/CentOS 7.4 RHEL/CentOS 7.5 RHEL/CentOS 7.6 RHEL/CentOS 7.7 Ubuntu 18.04 Ubuntu 20.04
BlueField	 RHEL/CentOS 7.4 RHEL/CentOS 7.5 RHEL/CentOS 7.6 RHEL/CentOS 7.7 Ubuntu 18.04 Ubuntu 20.04

ASAP² Requirements

- iproute >= 4.12 (for tc support)
- Upstream Open vSwitch >= 2.8 for CentOS 7.2 Mellanox openvswitch

ASAP² Supported Adapter Cards

- ConnectX-5
- ConnectX-6 Dx

NFS over RDMA (NFSoRDMA) Supported Operating Systems

NFSoRDMA Supported Operating Systems on Client and Target Sides

- SLES 12 SP4
- SLES 12 SP5
- SLES 15 SP1
- SLES 15 SP2
- Ubuntu 18.04.3
- Ubuntu 18.04.4
- Ubuntu 20.04
- Ubuntu 20.10
- RHEL 7.5
- RHEL 7.6
- RHEL 7.7
- RHEL 7.8
- RHEL 7.9
- RHEL 8.0
- RHEL 8.1RHEL 8.2
- RHEL 8.3

NVMEoF Supported Operating Systems

- SLES 12 SP3 and above
- RHEL 7.2 and above (Host side only)
- RHEL 7.4 and above (Host and Target sides)
- OS with distribution/custom kernel 4.8.x and above

Lustre Versions Compatible with MLNX_OFED

- Lustre 2.12.6
- Lustre 2.13.0

A Requires Upstream Lustre support for Kernel 5.8.

NEO-Host Supported Operating Systems

- RHEL 7.x
- OL 7.x
- SLES 12 SP3
- SLES 12 SP4
- SLES 12 SP5
- SLES 15 SP0
- SLES 15 SP1
- Ubuntu 16.04
- Ubuntu 18.04
- Ubuntu 19.04
- Debian 9.x
- Debian 10.x

GPUDirect Storage (GDS) Supported Operating Systems

• Ubuntu 18.04

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/OL/ Fedora	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
XenServer	yum install perl pciutils python libxml2-python libnl expat glib2 tcl bc libstdc++ tk pkgconfig ethtool
SLES 12	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
SLES 15	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
Ubuntu/Debian	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



As of MLNX_OFED v5.1, ConnectX-3, ConnectX-3 Pro or Connect-IB NICs are no longer supported. To work with a MLNX_OFED version that supports these adapter cards, please refer to MLNX_OFED 4.9 long-term support (LTS) version.

This current MLNX_OFED version supports the following Mellanox network adapter cards firmware versions:

NIC	Recommended Firmware Rev.	Additional Firmware Rev. Supported
ConnectX®-4	12.28.2006	12.28.1002
ConnectX®-4 Lx	14.29.1016	14.28.2006
ConnectX®-5/ConnectX®-5 Ex	16.29.1016	16.28.2006
ConnectX®-6	20.29.1016	20.28.2006
ConnectX®-6 Dx	22.29.1016	22.28.2006
ConnectX®-6 Lx	26.29.1016	26.28.2006
BlueField™	18.29.1016	18.28.2006

NIC	Recommended Firmware Rev.	Additional Firmware Rev. Supported
BlueField™-2	24.29.1016	N/A

For the official firmware versions, please see:

https://www.mellanox.com/support/firmware/firmware-downloads

MLNX_OFED Unsupported Functionalities/Features/NICs

The following are the unsupported functionalities/features/NICs in MLNX_OFED current version:

- · Innova IPsec EN
- ConnectX®-2 Adapter Card
- ConnectX®-3 Adapter Card
- ConnectX®-3 Pro Adapter Card
- Connect-IB® Adapter Card
- Relational Database Service (RDS)
- · mthca InfiniBand driver
- Ethernet IPolB (eIPolB)
- Soft-RoCE
- RDMA experimental verbs library (mlnx_lib)

Changes and New Features

MLNX_OFED New Features

The following are the changes and/or new features that have been added to this version of MLNX_OFED.

Feature/Change	Description
Supported NICs: ConnectX-5 and above	
Rx Multi-strides CQE Compression	Added CQE compression support for Rx multi-strides packets.
Multi-application QoS	Added support for configuring QoS on a single QP or on a group of QPs.
MPLS-over-UDP Hardware Offload Support	Added support for encap/decap hardware offload of IPv4 traffic over MPLS-over-UDP. This can be used in networks with MPLS routers to achieve more efficient routing.
Connection Tracking with Hairpin	Added support for adding connection tracking rules on VFs to forward traffic from one VF to the other.
sFlow Sampling Rules Offload	Added support for offloading sFlow sampling rules.
Supported NICs: BlueField/BlueField-2, ConnectX-5 & ConnectX-6 Dx	

mlx5dv_dr Software Steering Parallel Rules Insertion	Added support for a locking mechanism to enable parallel insertion of rules into the software steering using the mlx5dv_dr API. The parallel insertion improves the insertion rate and takes place when adding Rx and Tx rules via the FDB domain.		
mlx5dv_dr API Matching on Geneve Tunnel	Added support for the option to match mlx5dv_dr API on Geneve tunnel using a dynamic flex parser. The option header consists of class, type, length and data. The parser should be configured using devx command, after which a rule can be created to match on parser ID and data.		
OVS-DPDK Geneve Encap/Decap	Added support for Geneve tunneling offload, including matching on extension header.		
OVS-DPDK Parallel Offloads	Added support for parallel insertion and deletion of offloaded rules using multiple OVS threads.		
	Supported NICs: BlueField-2 & ConnectX-6 Dx		
GTP-U TEID Modification	[Beta] Added support to modify GTP-U TEID. This support requires flex parser configuration.		
OVS-DPDK E2E Cache Support	[Beta] Improved performance of OVS Connection Tracking flows by enabling the merge of the multi-table flow matches and actions into one joint flow.		
	Supported NICs: ConnectX-6 Dx and above		
Tx Port Time-Stamping	Transmitted packet timestamping accuracy can be improved when using a timestamp generated at the port level instead of a timestamp generated upon CQE creation. Tx port time-stamping better reflects the actual time of a packet's transmission.		
	This feature is disabled by default. The feature can be enabled or disabled using the following command.		
	ethtoolset-priv-flags <ifs-name> tx_port_ts on / off</ifs-name>		
	For further information on this feature, please see <u>Tx Port Time-Stamping</u> .		
Tunnel Rules Offload	Added support for offloading tunnel rules when the source interface is VF (in addition to uplink) in the Hypervisor.		
	Added support for offloading tunnel rules when the source interface is OpenvSwitch bridge (internal port).		
Connection Tracking Mirroring Offload	Added support for using Mirroring Offload with Connection Tracking.		
mlx5dv_dr API ASO Flow Meter	Added support for ASO flow meter using the mlx5dv_dr API, which allows for monitoring the packet rate for specific flows. When a packet hits a flow that is connected to a flow meter, the rate of packets through this meter is evaluated, and the packet is marked with a color copied into one of the C registers, according to the current rate compared to the reference rate.		
mlx5dv_dr API ASO First Hit	Added support for ASO first hit using the mlx5dv_dr API, which allows for tracking rule hits by packets. When a packet hits a rule with the ASO first hit action, a flag is set indicating this event, and the original value of the flag is copied to one of the C registers.		
mlx5dv_dr API GTP-U Extension Header	Added mlx5dv_dr API support for matching on a new field "gtpu_first_ext_dw_0". This field enables packet filtering based on the GTP-U first extension header (first dword only). To enable parsing of tunnel GTP-U extension header, run the following command.		
	./cloud_fw_reset.py FLEX_PARSER_PROFILE_ENABLE=3		
	Supported NICs: ConnectX-6 Lx and above		

IPsec Offload	Added IPsec full offload support for extended sequence number, replay protection window and lifetime packet limit.
Supported NICs: All	
Firmware Upgrade	Firmware upgrade during MLNX_OFED installation is now done on all supported devices simultaneously rather than consecutively.
RDMA-CM Disassociate Support	Added support for connecting kernel and RDMA-CM in a reliable way based on device index.
New Query GID API	Added support for a new query GID API that allows for querying a single GID entry by its port and GID index, or querying for all GID tables of a specific device. This API works over ioctl instead of sysfs, which accelerates the querying process.
Multi-Host Firmware Reset	Added support for performing multi-host firmware reset in order to upgrade the device firmware.
	Firmware reset loads the new firmware in case it was burnt on the flash and was pending activation, and reloads the current firmware image from the flash in case no new firmware was pending.
Firmware Live Patching	[Alpha] Added support for firmware live patching in the driver. Live patching updates the firmware without the need to perform firmware reset. However, it can only be applied in scenarios where the difference between the current and new firmware versions are minor, which is decided upon by the firmware itself.
Devlink Firmware Reset	Added support in the devlink tool for performing firmware reset in order to upgrade the device firmware.
	Firmware reset loads the new firmware in case it was burnt on the flash and was pending activation, and reloads the current firmware image from the flash in case no new firmware was pending.
	For further information, please refer to the the devlink man page.
	Note : In order for the firmware reset to run successfully, the following conditions should be met.
	 Each function should have the driver up and active with a version that supports this feature None of the functions has the devlink parameter enable_remote_dev_reset set to False.
	The current MLNX_OFED does not include the latest iproute2 version that provides support for this feature. Therefore, to be able to work with it, make sure to install the latest iproute2 version available on Github.
Command Interface Resiliency	Added a resiliency mechanism for the driver to manually poll the command event queue (EQ) in case of a command timeout. In case the resiliency mechanism finds unhandled event queue entry (EQE) due to a lost interrupt, the driver will handle it, after which the command interface returns to a healthy state.
Offloaded Traffic Sniffer	Setting a sniffer private flag is deprecated and no longer required. In order to capture offloaded/RoCE traffic, tcpdump can now be run on the RDMA device.
Devlink Port Health Reporters	Added per-port reporters to devlink health to manage per-port health activities. Users can now access the devlink port reporters by specifying the port index in addition to the device devlink name through the devlink health commands API. This update was first introduced in iproute2 v5.8. As part of this feature, mlx5e Tx and Rx reporters are now redefined as devlink port reporters. For examples, please see devlink-health manpage.

Memory Registration Optimization	Optimized memory consumption of memory registration in huge page systems. As an example, in a 2MB huge page system, 600 MB would be saved for 100 GB memory registration.		
mlx5dv API	Added support for mlx5dv API to modify the configured UDP source port for RoCE packets of a given RC/UC QP when QP is in RTS state.		
Enhanced Tx Multi- packet WQE (MPWQE)	Added support for accelerating Tx datapath by saving PCI bandwidth and CPU utilization. The savings are achieved by aggregating multiple packets into a single WQE. The feature is driven by xmit_more for certain traffic types, such as UDP.		
	Innova IPsec		
Innova IPsec NIC Support	Removed support for the network adapter Innova IPsec (EN).		
Supported NICs: All			
Bug Fixes	See <u>Bug Fixes</u> section.		

For additional information on the new features, please refer to MLNX_OFED User Manual.

API Changes in MLNX_OFED

MLNX_OFED Verbs API Migration

As of MLNX_OFED v5.0 release (Q1 of the year 2020), MLNX_OFED Verbs API have migrated from the legacy version of user space verbs libraries (libibervs, libmlx5, etc.) to the Upstream version rdmacore.

For the list of MLNX_OFED verbs APIs that have been migrated, refer to <u>Migration to RDMA-Core</u> document.

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

Internal Ref. Number	Issue
2302786	Description: On EulerOS 2.0 SP9 systems, the kernel ABI (kABI) between the base vhulk2006 kernel and the errata vhulk2008 kernel has been changed. It is now not possible to install MLNX_OFED compiled with KMP on vhulk2006 kernel on a vhulk2008 system.
	Workaround: Install MLNX_OFED withadd-kernel-support.
	Keywords: EulerOS; kABI; installation;add-kernel-support
	Discovered in Release: 5.2-1.0.4.0

Internal Ref. Number	Issue
2398281	Description: A crash in the TLS Rx socket cleanup flow may occur due to a kernel issue where a wrong extra call to tls_dev_del is made.
	Workaround: N/A
	Keywords: TLS RX device offload
	Discovered in Release: 5.2-1.0.4.0
2407415	Description: OpenEuler 20.03 Aarch64 with errata kernels 4.19.90-2011.6.0.0049.oe1.aarch64 and 4.19.90-2012.5.0.0054.oe1.aarch64 are incompatible with MLNX_OFED kmod-mlnx-ofa_kernel.
	Workaround: Install MLNX_OFED withadd-kernel-support.
	Keywords: OpenEuler; Aarch64; installation;add-kernel-support
	Discovered in Release: 5.2-1.0.4.0
2348077	Description: RDMA device name for VFs may change after resetting all VFs at once.
	Workaround: Either reset interfaces one by one with a delay in between, or use a network interface naming scheme with predictable interface names, such as NAME_PCI or NAME_GUID. Copy /lib/udev/rules.d/60-rdma-persistent-naming.rules to /etc/udev/rules.d/ and edit the last line accordingly. Note that this will change interface names.
	Keywords: RDMA; VF
	Discovered in Release: 5.2-1.0.4.0
2381713	Description: esp4_offload and esp6_offload modules are expected to be loaded according to the list determined by the default kernel. However, these modules cannot be loaded when working over Debian 10 with non-default custom kernel as they are not included in it.
	Workaround: Either install MLNX_OFED usingadd-kernel-support, or rebuild the non-default custom kernel to include these modules.
	Keywords: esp4_offload; esp6_offload; kernel, Debian
	Discovered in Release: 5.2-1.0.4.0
2382898	Description : On kernel 4.14, there is no traffic for UDP or TCP with payload size larger than 1398 on GENEVE IPv6 over VLAN tag interface.
	Workaround: N/A
	Keywords: GENEVE; stag; VLAN; UDP
	Discovered in Release: 5.2-1.0.4.0
2333971	Description: Changing the "other" channels count by "ethtool -L <interface> other <count>" command on Kernel 5.10 may cause a kernel panic.</count></interface>
	Workaround: N/A
	Keywords: Kernel 5.10, kernel panic; ethtool; "other" channels
	Discovered in Release: 5.2-1.0.4.0

Internal Ref. Number	Issue
2359599	Description: When an error occurs during adding a new flow rule (e.g. FW command failure in DMFS; unsupported sequence in SMFS; etc.), the driver might fail to release some resources.
	Workaround: N/A
	Keywords: Steering; SMFS; DMFS
	Discovered in Release: 5.2-1.0.4.0
2326155	Description: When toggling the link state while running RoCE traffic, the below warning may appear in the dmesg:
	ib_cache_gid_add: unable to add gid <gid> error=-28</gid>
	Workaround: N/A
	Keywords: RoCE;ib_cache_gid_add
	Discovered in Release: 5.2-1.0.4.0
2329654	Description: Running XDP over an IP tunnel may fail when working with kernels as old as version 4.14.
	Workaround: N/A
	Keywords: XDP, Kernel
	Discovered in Release: 5.2-1.0.4.0
2249156	Description: MLNX_OFED installation will remove qperf package in case it was done after qperf installation.
	Workaround: Make sure to install qperf package after installing MLNX_OFED, or reinstall qperf after installing MLNX_OFED.
	Keywords: Installation; qperf
	Discovered in Release: 5.2-1.0.4.0
2355956	Description: OFED installation requires kernel config CONFIG_DEBUG_INFO to be set.
	Workaround: N/A
	Keywords: Installation; CONFIG_DEBUG_INFO
	Discovered in Release: 5.2-1.0.4.0
2362781	Description: Openibd may fail to unload the Inbox driver mlx5_ib on Ubuntu 18.04 PPC Boston server due to a bug in the Inbox drivers.
	Workaround: N/A
	Keywords: Openibd; Inbox; Ubuntu; mlx5_ib
	Discovered in Release: 5.2-1.0.4.0

Internal Ref. Number	Issue
2367659	Description: Upgrading the MLNX_OFED version that is configured as a YUM repository may yield warning messages from depmod about unknown symbols, such as:
	depmod: WARNING: /lib/modules/4.18.0-240.el8.×8664/extra/iser/ib_iser.ko needs unknown symbol ib_fmr_pool_unmap depmod: WARNING: /lib/modules/4.18.0-240.el8.×8664/extra/srp/ib_srp.ko needs unknown symbol ib_create_qp_user
	These warnings appear since the RPM packages upgrade occurs sequentially, and there is an upgrade dependency between some of the modules, which would create a state of upgrade inconsistency. These warnings are temporary and can be ignored as eventually all modules will be upgraded, and the warnings will no longer appear.
	Workaround: N/A
	Keywords: YUM; RPM; symbol; depmod; ISER; SRP
	Discovered in Release: 5.2-1.0.4.0
2393352	Description: Using "with-openvswitch" flag during MLNX_OFED installation may not work on Debian 10 systems.
	Workaround: N/A
	Keywords:with-openvswitch; Debian
	Discovered in Release: 5.2-1.0.4.0
2383355	Description: eSwitch and eSwitch offloads are not supported for SR-IOV and its sub functions when installing MLNX_OFED over Upstream kernel v5.10 or higher.
	Workaround: N/A
	Keywords: eSwitch; Kernel; SR-IOV
	Discovered in Release: 5.2-1.0.4.0
2385269	Description: The number of connections offloaded is limited to 100K when working with Kernel v5.9.
	Workaround: N/A
	Keywords: ASAP2; Connection Tracking; Kernel
	Discovered in Release: 5.2-1.0.4.0
2393169	Description: Mirroring is not supported with Connection Tracking when the source port is a VxLAN device.
	Workaround: N/A
	Keywords: ASAP ² ; Connection Tracking; Mirroring
	Discovered in Release: 5.2-1.0.4.0
2395082	Description: A call trace may take place when moving from SwitchDev mode back to Legacy mode in Kernel v5.9 due to a kernel issue in tcf_block_unbind.
	Workaround: N/A
	Keywords: ASAP ² ;SwitchDev; call trace; kernel; tcf_block_unbind

Internal Ref. Number	Issue
	Discovered in Release: 5.2-1.0.4.0
2278352	Description: Under Containers environment, having "rdma-core" package installed on the Hypervisor may conflict with loading updated MLNX_OFED drivers from within a Container (using "/etc/init.d/openibd restart" for example). The "rdma-core" services may load the Inbox drivers on the Hypervisor, causing the loading of the MLNX_OFED drivers to fail, and errors about symbol incompatibility will appear in the dmesg log.
	Workaround: Delete the "rdma-core" package configuration files (*.conf) from the hypervisor under directory "/etc/rdma/modules". These files are responsible for defining which modules will be loaded by "rdma-core" services.
	Keywords: openibd; symbol; rdma-core; driver load
	Discovered in Release: 5.2-1.0.4.0
2328206	Description: Prior to stopping Open vSwitch, the following command should be run. ovs-appctl exitcleanup
	Workaround: N/A
	Keywords: Open vSwitch; ovs-appctl
	Discovered in Release: 5.2-1.0.4.0
2354899	Description: ODP is not supported on RHEL7.x systems when running over an ETH link layer with RoCE disabled.
	Workaround: N/A
	Keywords: ODP, RHEL, RoCE
	Discovered in Release: 5.1-2.5.8.0
2338150	Description: Scatter to CQE feature should be disabled for the GPUDirect tests to work.
	Workaround: Set the MLX5_SCATTER_TO_CQE environment variable to 0 before the ib_send_bw command. For example:
	MLX5_SCATTER_TO_CQE=0 ib_send_bw -d <>
	Keywords: CQE, GPUDirect
	Discovered in Release: 5.1-2.5.8.0
2295732	Description: Upgrading from legacy (mlnx-libs) to the current rdma-core based build using YUM (package manager) fails.
	Workaround: To perform this upgrade, either use the installer script or uninstall the old packages and install the new packages.
	Keywords: Legacy, mlnx-libs, rdma-core, installation
	Discovered in Release: 5.1-2.5.8.0
2295735	Description: Upgrading from legacy (mlnx-libs) to the current rdma-core based build using the apt-get (package manager) fails.
	Workaround: To perform this upgrade, either use the installer script or uninstall the old packages and install the new packages.

Internal Ref. Number	Issue
	Keywords: Legacy, mlnx-libs, rdma-core, apt, apt-get, installation
	Discovered in Release: 5.1-2.5.8.0
2293460	Description: In rare cases, mlx5_cmd_exec may lead to recount warning, such as:
	WARNING: CPU: 1 PID: 30811 at lib/refcount.c:28 refcount_warn_saturate+0xd9/0xe0
	This refcount warning can be ignored.
	Workaround: N/A
	Keywords: mlx5_cmd_exec
	Discovered in Release: 5.1-2.5.8.0
2278833	Description: Creating a bond via NetworkManager and restarting the driver (openibd restart) results in no pf0hpf and bond creation failure.
	Workaround: Before restarting driver, destroy bond configuration from NetworkManager. Or, instead of using NetworkManager to create bond, use "ip" commands.
	Keywords: Bond, LAG, network manager, driver reload
	Discovered in Release: 5.1-2.5.8.0
2209987	Description: aRFS feature (activated using "ethtool ntuple on") is disabled for kernel 4.1 or below.
	Workaround: N/A
	Keywords: aRFS
	Discovered in Release: 5.1-0.6.6.0
2200320	Description: In case MLNX_OFED is re-installed on a certain system without using force, the installation may fail requiring the removal of infiniband-diags package.
	Workaround: Remove the infiniband-diags package using rpm -e.
	Keywords: Installation, infiniband-diags
	Discovered in Release: 5.1-0.6.6.0
2248996	Description: Downgrading the firmware version for ConnectX-6 cards using "mlnx_ofed_installfw-update-onlyforce-fw-update" fails.
	Workaround: Manually downgrade the firmware version - please see <u>Firmware Update Instructions</u> .
	Keywords: Firmware, ConnectX-6
	Discovered in Release: 5.1-0.6.6.0
2175930	Description: When using OFED 5.1 on PPC architectures with kernels v5.5 or v5.6 and an old ethtool utility, a harmless warning call trace may appear in the dmesg due to mismatch between user space and kernel. The warning call trace mentions ethtool_notify.
	Workaround: Update the ethtool utility to version 5.6 on such systems in order to avoid the call trace.
	Keywords: PPC, ethtool_notify, kernel

Internal Ref. Number	Issue
	Discovered in Release: 5.1-0.6.6.0
2198764	Description: If MLNX_OFED is installed on a Debian or Ubuntu system that is run in chroot environment, the openibd service will not be enabled. If the chroot files are being used as a base of a full system, the openibd service is left disabled.
	Workaround: Currently, openibd is a sysv-init script that you can enable manually by running: update-rc.d openibd defaults
	Keywords: chroot, Debian , Ubuntu, openibd
	Discovered in Release: 5.1-0.6.6.0
2237134	Description: Running connection tracking (CT) with FW steering may cause CREATE_FLOW_TABLE command to fail with syndrome.
	Workaround: Configure OVS to use a single handler-thread:
	<pre>#ovs-vsctl set Open_vSwitch . other_config:n-handler-threads=1</pre>
	Keywords: Connection tracking, ASAP, OVS, FW steering
	Discovered in Release: 5.1-0.6.6.0
2239894	Description: Running OpenVSwitch offload with high traffic throughput can cause low insertion rate due to high CPU usage.
	Workaround: Reduce the number of combined channels of the uplink using "ethtool -L".
	Keywords: Insertion rate, ASAP2
	Discovered in Release: 5.1-0.6.6.0
2240671	Description: Header rewrite action is not supported over RHEL/CentOS 7.4.
	Workaround: N/A
	Keywords: ASAP, header rewrite, RHEL, RedHat, CentOS, OS
	Discovered in Release: 5.1-0.6.6.0
2242546	Description: Tunnel offload (encap/decap) may cause kernel panic if nf_tables module is not probed.
	Workaround: Make sure to probe the nf_tables module before inserting any rule.
	Keywords: Kernel v5.7, ASAP, kernel panic
	Discovered in Release: 5.1-0.6.6.0
2244416	Description: Configuring "other" channels over one representor is not supported and may cause a call trace.
	Workaround: N/A
	Keywords: ASAP, SwitchDev, ethtool, representor
	Discovered in Release: 5.1-0.6.6.0
2143007	Description: IPsec packets are dropped during heavy traffic due to a bug in net/xfrm Linux Kernel.
	Workaround: Make sure the Kernel is modified to apply the following patch: "xfrm: Fix double ESP trailer insertion in IPsec crypto offload".

Internal Ref. Number	Issue
	Keywords: IPsec, xfrm
	Discovered in Release: 5.1-0.6.6.0
2225952	Description: VF mirroring with TC policy skip_sw is not supported on RHEL/CentOS 7.4, 7.5 and 7.6 OSs.
	Workaround: N/A
	Keywords: ASAP ² , Mirroring, RHEL, RedHat, OS
	Discovered in Release: 5.1-0.6.6.0
2216521	Description: After upgrading MLNX_OFED from v5.0 or earlier, ibdev2netdev utility changes the installation prefix to /usr/sbin. Therefore, it cannot be found while found in the same SHELL environment.
	Workaround: After installing MLNX_OFED, log out and log in again to refresh the SHELL environment.
	Keywords: ibdev2netdev
	Discovered in Release: 5.1-0.6.6.0
2202520	Description: Rules with VLAN push/pop, encap/decap and header rewrite actions together are not supported.
	Workaround: N/A
	Keywords: ASAP ² , SwitchDev, VLAN push/pop, encap/decap, header rewrite
	Discovered in Release: 5.1-0.6.6.0
2210752	Description: Switching from Legacy mode to SwitchDev mode and vice-versa while TC rules exist on the NIC will result in failure.
	Workaround: Before attempting to switch mode, make sure to delete all TC rules on the NIC or stop OpenvSwitch.
	Keywords: ASAP ² , Devlink, Legacy SR-IOV
	Discovered in Release: 5.1-0.6.6.0
2125036/2125031	Description: Upgrading the MLNX_OFED from an UPSTREAM_LIBS based version to an MLNX_LIBS based version fails unless the driver is uninstalled and then re-installed.
	Workaround: Make sure to uninstall and re-install MLNX_OFED to complete the upgrade.
	Keywords: Installation, UPSTREAM_LIBS, MLNX_LIBS
	Discovered in Release: 5.0-2.1.8.0
2105447	Description: hns_roce warning messages will appear in the dmesg after reboot on Euler2 SP3 OSs.
	Workaround: N/A
	Keywords: hns_roce, dmesg, Euler
	Discovered in Release: 5.0-2.1.8.0

Internal Ref. Number	Issue
2110321	Description: Multiple driver restarts may cause IPoIB soft lockup.
	Workaround: N/A
	Keywords: Driver restart, IPoIB
	Discovered in Release: 5.0-2.1.8.0
2112251	Description: 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.
	Workaround: Define geneve tunnel's remote endpoint using IPv4.
	Keywords: Kernel, Geneve, IPv4, IPv6, MTU, fragmentation
	Discovered in Release: 5.0-2.1.8.0
2119210	Description: Multiple driver restarts may cause a stress and result in mlx5 commands check error message in the log.
	Workaround: N/A
	Keywords: Driver restart, syndrome, error message
	Discovered in Release: 5.0-2.1.8.0
2118956	Description: mlx5dv_dr API does not support sub functions (SFs) as destination actions.
	Workaround: Create the SFs only after domain creation.
	Keywords: mlx5dv_dr, sub functions, SF
	Discovered in Release: 5.0-2.1.8.0
2102902	Description: 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.
	Workaround: N/A
	Keywords: TLS offload, mlx5e
	Discovered in Release: 5.0-2.1.8.0
2111534	Description: 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.
	Workaround: N/A
	Keywords: TLS offload, mlx5e
	Discovered in Release: 5.0-2.1.8.0
2117845	Description: 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.
	Workaround: N/A
	Keywords: Relaxed ordering, memory region, MR, CAPI
	Discovered in Release: 5.0-2.1.8.0

Internal Ref. Number	Issue
2035950	Description: 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.
	 Removing PF from the bond (using ifdown, ip link or any other function) Attempting to disable SR-IOV
	Workaround: N/A
	Keywords: VF LAG, binding, firmware, FW, PF, SR-IOV
	Discovered in Release: 5.0-1.0.0.0
2094176	Description: When running in a large scale in VF-LAG mode, bandwidth may be unstable.
	Workaround: N/A
	Keywords: VF LAG
	Discovered in Release: 5.0-1.0.0.0
2044544	Description: When working with OSs with Kernel v4.10, bonding module does not allow setting MTUs larger than 1500 on a bonding interface.
	Workaround: Upgrade your Kernel version to v4.11 or above.
	Keywords: Bonding, MTU, Kernel
	Discovered in Release: 5.0-1.0.0.0
1882932	Description: Libibverbs dependencies are removed during OFED installation, requiring manual installation of libraries that OFED does not reinstall.
	Workaround: Manually install missing packages.
	Keywords: libibverbs, installation
	Discovered in Release: 5.0-1.0.0.0
2058535	Description: ibdev2netdev command returns duplicate devices with different ports in SwitchDev mode.
	Workaround: Use /opt/mellanox/iproute2/sbin/rdma link show command instead.
	Keywords: ibdev2netdev
	Discovered in Release: 5.0-1.0.0.0
2072568	Description: In RHEL/CentOS 7.2 OSs, adding drop rules when act_gact is not loaded may cause a kernel crash.
	Workaround: Preload all needed modules to avoid such a scenario (cls_flower, act_mirred, act_gact, act_tunnel_key and act_vlan).
	Keywords: RHEL/CentOS 7.2, Kernel 4.9, call trace, ASAP
	Discovered in Release: 5.0-1.0.0.0
2093698	Description: VF LAG configuration is not supported when the NUM_OF_VFS configured in mlxconfig is higher than 64.
	Workaround: N/A
	Keywords: VF LAG, SwitchDev mode, ASAP

Internal Ref. Number	Issue
	Discovered in Release: 5.0-1.0.0.0
2093746	Description: Devlink health dumps are not supported on kernels lower than v5.3.
	Workaround: N/A
	Keywords: Devlink, health report, dump
	Discovered in Release: 5.0-1.0.0.0
2000590	Description: Sending packets larger than MTU is not supported when working with OVS-DPDK.
	Workaround: N/A
	Keywords: MTU, OVS-DPDK
	Discovered in Release: 5.0-1.0.0.0
2062900	Description: Moving VF from SwitchDev mode to Legacy mode while the representor is being used by OVS-DPDK results in a segmentation fault.
	Workaround: To move VF to Legacy mode with no error, make sure to delete the ports from the OVS.
	Keywords: SwitchDev, Legacy, representor, OVS-DPDK
	Discovered in Release: 5.0-1.0.0.0
2075942	Description: Huge pages configuration is lost each time the server is configured.
	Workaround: Re-configure the huge pages after each reboot, or configure them as a kernel parameter.
	Keywords: Huge pages, reboot, OVS-DPDK
	Discovered in Release: 5.0-1.0.0.0
2083427	Description: For kernels with connection tracking support, neigh update events are not supported, requiring users to have static ARPs to work with OVS and VxLAN.
	Workaround: N/A
	Keywords: VxLAN, VF LAG, neigh, ARP
	Discovered in Release: 5.0-1.0.0.0
2067012	Description: MLNX_OFED cannot be installed on Debian 9.11 OS in SwitchDev mode.
	Workaround: Install OFED with the flagadd-kernel-support.
	Keywords: ASAP, SwitchDev, Debian, Kernel
	Discovered in Release: 5.0-1.0.0.0
2067746	Description : When attaching a second slave to a bond, some bond interface GIDs might disappear.
	Workaround: Re-create and re-configure the bond device.
	Keywords: Bond, GID
	Discovered in Release: 5.0-1.0.0.0

Internal Ref. Number	Issue
2036572	Description: When using a thread domain and the lockless rdma-core ibv_post_send path, there is an additional CPU penalty due to required barriers around the device MMIO buffer that were omitted in MLNX_OFED.
	Workaround: N/A
	Keywords: rdma-core, write-combining, MMIO buffer
	Discovered in Release: 5.0-1.0.0.0
-	Description: The argparse module is installed by default in Python versions =>2.7 and >=3.2. In case an older Python version is used, the argparse module is not installed by default.
	Workaround: Install the argparse module manually.
	Keywords: Python, MFT, argparse, installation
	Discovered in Release: 4.7-3.2.9.0
1997230	Description: Running mlxfwreset or unloading mlx5_core module while contrak flows are offloaded may cause a call trace in the kernel.
	Workaround: Stop OVS service before calling mlxfwreset or unloading mlx5_core module.
	Keywords: Contrak, ASAP, OVS, mlxfwrest, unload
	Discovered in Release: 4.7-3.2.9.0
1955352	Description : Moving 2 ports to SwitchDev mode in parallel is not supported.
	Workaround: N/A
	Keywords: ASAP, SwitchDev
	Discovered in Release: 4.7-3.2.9.0
1979958	Description: VxLAN IPv6 offload is not supported over CentOS/RHEL v7.2 OSs.
	Workaround: N/A
	Keywords: Tunnel, VXLAN, ASAP, IPv6
	Discovered in Release: 4.7-3.2.9.0
1980884	Description: Setting VF VLAN, state and spoofchk using ip link tool is not supported in SwitchDev mode.
	Workaround: N/A
	Keywords: ASAP, ip tool, VF, SwitchDev
	Discovered in Release: 4.7-3.2.9.0
1991710	Description: PRIO_TAG_REQUIRED_EN configuration is not supported and may cause call trace.
	Workaround: N/A
	Keywords: ASAP, PRIO_TAG, mstconfig
	Discovered in Release: 4.7-3.2.9.0

Internal Ref. Number	Issue
1970429	Description: 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.
	 Ingress (to the VM): All echo request fragments miss the corresponding offloading flow, but all echo reply fragments hit the corresponding offloading flow Egress (from the VM): The first fragment still hits the corresponding offloading flow, and the subsequent fragments miss the corresponding offloading flow
	Workaround: N/A
	Keywords: HW offloading, SR-IOV, SwitchDev, ICMP, VM, virtualization
	Discovered in Release: 4.7-3.2.9.0
1967866	Description: Enabling ECMP offload requires the VFs to be unbound and VMs to be shut down.
	Workaround: N/A
	Keywords: ECMP, Multipath, ASAP ²
	Discovered in Release: 4.7-3.2.9.0
1921981	Description: On Ubuntu, Debian and RedHat 8 and above OSS, parsing the mfa2 file using the mstarchive might result in a segmentation fault.
	Workaround: Use mlxarchive to parse the mfa2 file instead.
	Keywords: MFT, mfa2, mstarchive, mlxarchive, Ubuntu, Debian, RedHat, operating system
	Discovered in Release: 4.7-1.0.0.1
1840288	Description: MLNX_OFED does not support XDP features on RedHat 7 OS, despite the declared support by RedHat.
	Workaround: N/A
	Keywords: XDP, RedHat
	Discovered in Release: 4.7-1.0.0.1
1821235	Description: 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.
	Workaround: Use the both actions at the end of the flow process.
	Keywords: Flow steering
	Discovered in Release: 4.7-1.0.0.1
1888574	Description: Kernel support limitations in the current MLNX_OFED version:
	 SR-IOV SwitchDev is only supported on Kernel 4.14 and above, and on RedHat/CentOS 7.4, 7.5 and 7.6. SR-IOV Legacy is only supported on Kernel 4.3 and above, and on RedHat/CentOS 7.4, 7.5, 7.6 and 7.7.
	Workaround: N/A
	Keywords: SwitchDev, ASAP, Kernel , SR-IOV, RedHat, RHEL
	Discovered in Release: 4.7-1.0.0.1
1892663	Description: mlnx_tune script does not support python3 interpreter.

Internal Ref. Number	Issue
	Workaround: Run mlnx_tune with python2 interpreter only.
	Keywords: mlnx_tune, python3, python2
	Discovered in Release: 4.7-1.0.0.1
1504785	Description: A lost interrupt issue in pass-through virtual machines may prevent the driver from loading, followed by printing managed pages errors to the dmesg.
	Workaround: Restart the driver.
	Keywords: VM, virtual machine
	Discovered in Release: 4.6-1.0.1.1
1764415	Description: Unbinding PFs on LAG devices results in a "Failed to modify QP to RESET" error message.
	Workaround: N/A
	Keywords: RoCE LAG, unbind, PF, RDMA
	Discovered in Release: 4.6-1.0.1.1
1806565	Description: 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.
	Workaround: Use RoCEv2 default GID.
	Keywords: RoCE
	Discovered in Release: 4.6-1.0.1.1
1834997	Description: When working with VF Lag while the bond device is in active-active mode, traffic on both physical ports may not reach line rate.
	Workaround: N/A
	Keywords: VF LAG, bonding, bandwidth degradation, fairness
	Discovered in Release: 4.6-1.0.1.1
-	Description: Aging is not functional on bond device in RHEL 7.6.
	Workaround: N/A
	Keywords: VF LAG, ASAP ²
	Discovered in Release: 4.6-1.0.1.1
1747774	Description: In VF LAG mode, outgoing traffic in load balanced mode is according to the origin ring, thus, half of the rings will be coupled with port 1 and half with port 2. All the traffic on the same ring will be sent from the same port.
	Workaround: N/A
	Keywords: VF LAG, ASAP ²
	Discovered in Release: 4.6-1.0.1.1

Internal Ref. Number	Issue
1787667	Description: 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 <u>Bugzilla issue #1150850</u> and <u>Bugzilla issue #1150846</u> .
	Workaround: Change the kernel configuration of NVMe-oF driver to be "=m" and recompile the kernel.
	Keywords: NVME-oF, NVME, SLES
	Discovered in Release: 4.6-1.0.1.1
1753629	Description: 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
	Workaround: N/A
	Keywords: Bonding, slave
	Discovered in Release: 4.6-1.0.1.1
1712068	Description: Uninstalling MLNX_OFED automatically results in the uninstallation of several libraries that are included in the MLNX_OFED package, such as InfiniBand-related libraries.
	Workaround: If these libraries are required, reinstall them using the local package manager (yum/dnf).
	Keywords: MLNX_OFED libraries
	Discovered in Release: 4.6-1.0.1.1
-	Description: 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.
	Workaround: N/A
	Keywords: MFT compatible
	Discovered in Release: 4.6-1.0.1.1
1730840	Description: On ConnectX-4 HCAs, GID index for RoCE v2 is inconsistent when toggling between enabled and disabled interface modes.
	Workaround: N/A
	Keywords: RoCE v2, GID
	Discovered in Release: 4.6-1.0.1.1
1717428	Description : On kernels 4.10-4.14, MTUs larger than 1500 cannot be set for a GRE interface with any driver (IPv4 or IPv6).
	Workaround: Upgrade your kernel to any version higher than v4.14.
	Keywords: Fedora 27, gretap, ip_gre, ip_tunnel, ip6_gre, ip6_tunnel

Internal Ref. Number	Issue
1748343	Description : Driver reload takes several minutes when a large number of VFs exists.
	Workaround: N/A
	Keywords: VF, SR-IOV
	Discovered in Release: 4.6-1.0.1.1
1748537	Description : Cannot set max Tx rate for VFs from the ARM.
	Workaround: N/A
	Keywords: Host control, max Tx rate
	Discovered in Release: 4.6-1.0.1.1
1732940	Description : Software counters not working for representor net devices.
	Workaround: N/A
	Keywords: mlx5, counters, representors
	Discovered in Release: 4.6-1.0.1.1
1733974	Description : 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.
	Workaround: N/A
	Keywords: mlx5
	Discovered in Release: 4.6-1.0.1.1
1731939	Description: Get/Set Forward Error Correction FEC configuration is not supported on ConnectX-6 HCAs with 200Gbps speed rate.
	Workaround: N/A
	Keywords: Forward Error Correction, FEC, 200Gbps
	Discovered in Release: 4.6-1.0.1.1
-	Description : On ConnectX-6 HCAs and above, an attempt to configure advertisement (any bitmap) will result in advertising the whole capabilities.
	Workaround: N/A
	Keywords: 200Gmbps, advertisement, Ethtool
	Discovered in Release: 4.6-1.0.1.1
1699289	Description: 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.
	Workaround: Make sure to enable HW LRO in the driver:
	ethtool -k <intf> lro</intf>
	ethtoolset-priv-flag <intf> hw_lro on</intf>
	Keywords: HW LRO, ConnectX-5 and above
	Discovered in Release: 4.5-1.0.1.0

Internal Ref. Number	Issue
1403313	Description : 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.
	Workaround: Make sure to update the Kernel version to v4.15 or above.
	Keywords: VF, PF, IOMMU, Kernel, OS
	Discovered in Release: 4.5-1.0.1.0
-	Description : NEO-Host is not supported on the following OSs:
	 SLES12 SP3 SLES12 SP4 SLES15 Fedora 28 RHEL7.1 RHEL7.4 ALT (Pegas1.0) REL 7.5 RHEL7.6 XenServer 4.9
	Workaround: N/A
	Keywords: NEO-Host, operating systems
	Discovered in Release: 4.5-1.0.1.0
1521877	Description : On SLES 12 SP1 OSs, a kernel tracepoint issue may cause undefined behavior when inserting a kernel module with a wrong parameter.
	Workaround: N/A
	Keywords: mlx5 driver, SLES 12 SP1
	Discovered in Release: 4.5-1.0.1.0

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

Internal Reference Number	Description
2362665	Description: Fixed the issue where mdev module kernels dependency prevented their usage with MLNX_OFED. On kernel v4.18 through v5.1, mlnx-ofa kernel can now be rebuilt with mdev support by running the installer withadd-kernel-supportwithmdev. On RHEL 8.3 and kernel v5.2 and above, the internal mdev module can be used for mlnx-ofa kernel.
	Keywords: mlnx-ofa kernel; RHEL; Kernel;add-kernel-supportwith-mdev

Internal Reference Number	Description
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2395054	Description: Fixed the issue where installing MLNX_OFED usingadd-kernel-supportkmp used to fail. MLNX_OFED can now be installed using this option provided that kernel-rpm-macros and openEuler-rpm-config packages' minimal version is 30-9.h10.eulerosv2r9.
	Keywords: KMP;add-kernel-supportkmp; OpenEuler
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2312015	Description: Fixed the issue where clearing min_rate on all SR-IOV legacy VFs after setting min_rate to at least one of the VFs did not disable QoS min_rate.
	Keywords: SR-IOV; Legacy; QoS; VF; min_rate
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.2-1.0.4.0
2322044	Description: Allowed installation of KMP (kmod) on Lustre kernels usingadd-kernel-supportkmp on RHEL 7.x and above systems.
	Keywords: Installation; KMP; kmod; RHEL; Lustre; add-kernel-support
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2325883	Description: An empty /usr/src/mlnx-ofa_kernel/default is now no longer created.
	Keywords: mlnx-ofa_kernel
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.2-1.0.4.0
2326961	Description: Enabled installation of mlnx-ofa_kernel package to /bin/python3 instead of /usr/bin/python3 on RHEL 8.x systems.
	Keywords: /bin/python3; mlnx-ofa_kernel; RHEL; installation
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.2-1.0.4.0
2329152	Description: Fixed the issue where installing packages through a repository failed after generation of metapackages usingadd-kernel-support. The failure occurred due to excessive and incorrect Obsolete headers in those metapackages.
	Keywords: add-kernel-support; Installation; metapackage; repository
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0

Internal Reference Number	Description
2152217	Description: Fixed the issue where all IPoIB offload packets were wrongly counted as rx_csum_complete. These packets are now identified as rx_csum_unnecessary packets.
	Keywords: IPoIB; rx_csum_complete; rx_csum_unnecessary
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2392338	Description: Fixed an issue where mlx5dv_dr_rule_destroy segmentation fault could happen in rare cases with multiple rules on the same matcher with a different number of actions. This could happen after reusing an already deleted rule memory with less actions.
	Keywords: mlx5dv_dr; software steering; RMDA-Core
	Discovered in Release: 5.1-2.3.7.1
	Fixed in Release: 5.2-1.0.4.0
2297535	Description: Fixed the OFED udev script to not modify non-Mellanox NIC names.
	Keywords: udev; SwitchDev
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2333242	Description: Fixed an issue of moving to SwitchDev mode after configuring DSCP.
	Keywords: DSCP; SwitchDev
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2368330	Description: Fixed the issue where global traffic class configuration did not take effect on DC QPs.
	Keywords: DC; QP; TC
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2241983	Description: Fixed an issue where ib_send_bw traffic frequently dropped to zero when RDMA CM was used, because of incorrect min_rnr_timer setting on the responder side. The min_rnr_timer setting is now aligned with the setting in non-RDMA CM cases.
	Keywords: RDMA CM; ib_send_bw
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.2-1.0.4.0
2334009	Description: Fixed an issue where traffic did not pass over VFs with VST QinQ feature is enabled.
	Keywords: VST; QinQ

Internal Reference Number	Description
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2336091	Description: Fixed an issue of when traffic was sent over Geneve VLAN with Tx VLAN offload enabled and TSO or Tx csum enabled, traffic could be dropped and not sent to the wire.
	Keywords: Geneve; VLAN; Tx, TSO; offload
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2317257	Description: Fixed an issue that caused the firmware to restart upon installing mlnx-ofed-dpdk-upstream-libs package manually.
	Keywords: Firmware; mlnx-ofed-dpdk-upstream-libs; installation
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2365123	Description: Fixed bad utility paths in rdma-core "dracut" hooks on SLES systems, which used to result in the following errors when running "dracut" with the "add rdma" option.
	<pre>dracut-install: ERROR: installing '/usr/libexec/mlx4-setup.sh' dracut: /usr/lib/dracut/dracut-install -D /var/tmp/dracut.UdCOSJ/ initramfs /usr/libexec/mlx4-setup.sh dracut-install: ERROR: installing '/usr/libexec/rdma-set-sriov-vf' dracut: /usr/lib/dracut/dracut-install -D /var/tmp/dracut.UdCOSJ/ initramfs /usr/libexec/rdma-set-sriov-vf</pre>
	Keywords: rdma-core dracut
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2379620	Description: Fixed an issue where MLNX_OFED installation using 'yum' tool failed with the following errors.
	Requires: rdma-core(x86-64) = 50mlnx1-1.50218 Removing: rdma-core-22.4-1.el7.x86_64 (@anaconda) rdma-core(x86-64) = 22.4-1.el7
	Obsoleted By: mlnx-ofed-all-5.0-2.1.8.0.rhel7.8.noarch (MLNX_LOCAL_REPO) Not found
	Updated By: rdma-core-50mlnx1-1.50218.x86_64 (MLNX_LOCAL_REPO)
	rdma-core(x86-64) = 50mlnx1-1.50218 Available: rdma-core-22.4-5.el7.x86_64 (base)
	rdma-core(x86-64) = 22.4-5.el7
	Keywords: YUM; rdma-core conflict
	Discovered in Release: 5.0-2.1.8.0

Internal Reference Number	Description
	Fixed in Release: 5.2-1.0.4.0
2299982	Description: Fixed the issue where traffic class value was not updated in DCT when set via sysfs.
	Keywords: DCT; sysfs; TC
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.2-1.0.4.0
2335165	Description: Fixed a doorbell loss issue on AMD platforms with Secure Memory Encryption (SME).
	Keywords: AMD; SME
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2355878	Description: Fixed an issue with registering memory using mlx5dv_devx_umem_reg while forking. Without this fix, applications which use fork() or similar syscalls while using a memory registered with umem_reg could hang due to incorrect physical page mapping. This fix requires setting the IBV_FORK_SAFE environment variable.
	Keywords: mlx5dv_dr; SW steering; RDMA core
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2288599	Description: Fixed the issue where unbinding the device resulted in the following message being printed to the dmesg: "failed to disable DC tracer"
	Keywords: Unbind
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2083942	Description : Fixed the issue where the content of file /sys/class/net/ <netif>/ statistics/multicast may have been out of date and may have displayed values lower than the real values.</netif>
	Keywords: Multicast counters
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.2-1.0.4.0
2282316	Description: Fixed the issue where ERSPAN protocol was available only when turning off Tx checksum offload.
	Keywords: ERSPAN; TX checksum offload
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0

Internal Reference Number	Description
2310695	Description: Fixed a udev script issue which caused non-mellanox devices to be renamed.
	Keywords: udev; naming
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.2-1.0.4.0
2334518	Description: Fixed missing representor statistics when using ifconfig.
	Keywords: SwitchDev; representor; statistics; ifconfig
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2342348	Description: Fixed wrong value of skb mark of received packets on representors.
	Keywords: SwitchDev; skb mark
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2363982	Description: Fixed an issue which caused second port representors to be named as first port representors.
	Keywords: SwitchDev; udev; representor
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.2-1.0.4.0
2292762	Description: 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.
	Keywords: Proved VFs
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.1-2.5.8.0
2302010	Description: Fixed a GPUDirect locking bug that may have caused instability and communication loss in Peer-direct applications.
	Keywords: Peer-direct, GPUDirect
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.1-2.5.8.0
2298308	Description: Fixed the issue where pinned pages were not handled properly by peer flow, which resulted in ENOMEM error.
	Keywords: Peer flow
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.1-2.5.8.0

Internal Reference Number	Description
2323446	Description: Fixed the issue where num_free_callbacks counter was not functional, and querying the counters returned the value of 0 all the time.
	Keywords: GPUDirect
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.1-2.5.8.0
2288779	Description: Fixed a KMP support issue in EulerOS 2.0 SP9 with non-default kernel.
	Keywords: KMP, EulerOS, OS
	Discovered in Release: 5.1-2.5.8.0
	Fixed in Release: 5.1-2.5.8.0
2298285	 Description: Fixed the following VST issues. 1. Ingress VST traffic tagging issue caused by incorrect configuration settings over VF interface. 2. Enabled proper reset of ingress and egress legacy ACL.
	Keywords: VST
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2247363	Description: Added command interface resilience by manually polling the async commands EQE in case of command timeout.
	Keywords: Command interface, mlx5_core
	Discovered in Release: 4.2-1.2.0.0
	Fixed in Release: 5.1-2.5.8.0
2244729	Description: Enabled on-demand device memory sync to free cached memory.
	Keywords: Device memory, cache
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2229683	Description: Eliminated long mlx5 recovery-flow delay on a VF driver when PCI interface goes down.
	Keywords: mlx5, VF, PCI
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
	Description: ibv_get_device_list returns only accessible devices now.
	Keywords: ibv_get_device_list
	Discovered in Release: 5.1-0.6.6.0

Internal Reference Number	Description
	Fixed in Release: 5.1-2.5.8.0
2257363	Description: The install script will now also install the package rpm-build (that includes basic support for building RPM packages) on SLES15.x systems when trying to rebuild packages, as is the case with the rest of the distributions.
	Keywords: RPM, SLES, OS
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2244336	Description: AF_XDP is now functional.
	Keywords: AF_XDP
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2192791	Description: 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.
	Keywords: NEO-Host, SDK
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2247404	Description: srptools installation no longer fails in case the srp_daemon service fails to start on Debian and Ubuntu systems.
	Keywords: SRP, srptool, srp_daemon, Debian, Ubuntu, OS
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2248201	Description: Fixed the issue where during MLNX_OFED installation, warning messages related to modules iw_cxgb3 and iw_nes may have appeared in the log.
	Keywords: SLES, RHEL, KMP, weak updates, kmod
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2261475	Description : Simplified the sanity test that verifies which systems must use python3 and not python2, so that when rebuilding the kernel packages usingadd-kernel-support, the systems based on CentOS 8 are detected.
	Keywords: CentOS 8, RHEL 8, OS, sanity test, python
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0

Internal Reference Number	Description
2288099/2257445	Description: Simplified sanity tests for building MLNX_OFED/MLNX_EN with kernel v5.7.11 or newer.
	Keywords: Kernel, sanity test, installation
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2288783	Description: MLNX_OFED/MLNX_EN installer will no longer automatically remove the spdk package when installing the driver.
	Keywords: spdk, installation
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2292273	Description: Marked additional packages as obsolete so they would be removed when installing the HPE-roce mlnx-ofa-kernel package.
	Keywords: RPM, HPE-roce
	Discovered in Release: 5.1-1.0.7.1
	Fixed in Release: 5.1-2.5.8.0
2265094	Description: Fixed VM driver to avoid command timeouts when the Hypervisor disables the VF's PCI interface. This reduces the device removal time to 2 seconds.
	Keywords: VF disable
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2272539	Description: Added back rx_dct_connect counter that queries the number of received connection requests for the associated DCTs.
	Keywords: rx_dct_connect, DCT
	Discovered in Release: 5.1-0.6.6.0
	Fixed in Release: 5.1-2.5.8.0
2264237	Description: Fixed an issue with SLES and RHEL OS with Errata kernel where the driver failed to load as a result of missing dummy modules for cxgb3 and nes drivers.
	Keywords: SLES, RHEL, RedHat, OS
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-1.0.7.0
2020260	Description: Fixed the issue of when changing the Trust mode to DSCP, there was an interval between the change taking effect in the hardware and updating the inline mode of the SQ in the driver. If any traffic was transmitted during this interval, the driver would not inline enough headers, resulting in a CQE error in the NIC.
	Keywords: DSCP, inline, SQ, CQE

Internal Reference Number	Description
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0
2105631	Description: Removed IBV_FLOW_ATTR_FLAGS_ALLOW_LOOP_BACK flag as it is not used by the kernel.
	Keywords: IBV_FLOW_ATTR_FLAGS_ALLOW_LOOP_BACK
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.1-0.6.6.0
2099043	Description: Added QP isolation to improve SW steering performance under high packet load. This will allow SW steering RC QP to be executed on a separate scheduling queue without competing over hardware resources.
	Keywords: Software steering, ASAP, connection tracking, CT
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0
2097045	Description: Userspace Software Steering using mlx5dv_dr API support on ConnectX-6 Dx adapter cards is now at GA level.
	Keywords: Software Steering, SW, mlx5dv_dr, ConnectX-6 Dx
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2132332	Description: Fixed a sporadic reporting bandwidth issue in case of running with run_infinitely flag.
	Keywords: perftest, bandwidth
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2151658	Description: Optimized XRC target lookup by modifying the locking scheme to enable multiple readers and changing the linked list that holds the QPs to xarray.
	Keywords: XRC, QP, xarray
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2196118	Description: Fixed a driver issue that led to panic after DPDK application crashes.
	Keywords: DPDK, panic
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0

Internal Reference Number	Description
2245228	Description: Fixed an issue of a crash when attempting to access roce_enable sysfs in unprobed VFs.
	Keywords: roce_enable, unprobed VFs
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2061294	Description: Fixed a race of commands executed by command interface in parallel to AER recovery causing the kernel to crash.
	Keywords: mlx5e, AER
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.1-0.6.6.0
2131951	Description: Fixed an issue in MLNX_OFED build system that broke RPM sign process for random packages; all RPMs are now signed properly.
	Keywords: RPM, sign
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0
2143067	Description: If Openibd was configured to enable the SRP daemon, it now also enables srp_daemon from rdma-core.
	Keywords: Openibd, SRP daemon, srp_daemon, rdma-core
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0
2143094	Description: Regenerated package repository in the correct location after rebuilding the kernel using add-kernel-support. This allows for installing the newly generated packages with a package manager.
	Keywords: add-kernel-support, RPM, deb
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0
2172130	Description: Fixed an issue with metadata packages generation in the eth-only directory. This allows using the directory as a repository for package managers.
	Keywords: Metadata packages
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2214543	Description: Moved ibdev2netdev script from /usr/bin to /usr/sbin in the RPM package to avoid package conflict with RHEL 8 and consequent MLNX_OFED installation failure on some systems.
	Keywords: ibdev2netdev, RPM, RHEL, RedHat

Internal Reference Number	Description
	Discovered in Release:
	Fixed in Release: 5.1-0.6.6.0
2211311	Description: 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.
	Keywords: mlx5e, RX buffers, mlnx_qos
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2111349	Description: Fixed the issue where ethtoolshow-fec/get-fec were not supported over ConnectX-6 and ConnectX-6 Dx adapter cards.
	Keywords: Ethtool, ConnectX-6 Dx
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
2165668	Description: Fixed an issue related to mlx5 command interface that in some scenarios caused the driver to hang.
	Keywords: ConnectX-5, mlx5, panic
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.1-0.6.6.0
2119984	Description: Fixed the issue where IPsec crypto offloads did not work when ESN was enabled.
	Keywords: IPsec, ESN
	Discovered in Release: 5.0-2.1.8.0
	Fixed in Release: 5.1-0.6.6.0
1630228	Description: Fixed the issue where tunnel stateless offloads were wrongly forbidden for E-Switch manager function.
	Keywords: Stateless offloads cap
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.1-0.6.6.0
2089996	Description: Fixed the issue where dump flows were not supported and may have been corrupted when using tc tool with connection tracking rules.
	Keywords: ASAP, iproute2, tc, connection tracking
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.1-0.6.6.0

Internal Reference Number	Description
2094216	Description: Fixed the issue of when one of the LAG slaves went down, LAG deactivation failed, ultimately causing bandwidth degradation.
	Keywords: RoCE LAG
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2133778	Description: 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".
	Keywords: Namespace
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2076546	Description: Fixed the issue where in RPM-based OSs with non-default kernels, using repositories after re-creating the installer (usingadd-kernel-support) would result in improper installation of the drivers.
	Keywords: Installation, OS
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.1-0.6.6.0
2114957	Description: 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.
	Keywords: Installation, python
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2122684	Description: Fixed the issue where OFED uninstallation resulted in the removal of dependency packages, such as qemu-system-* (qemu-system-x86).
	Keywords: Uninstallation, dependency, qemu-system-x86
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2135476	Description: Added KMP ability to install MLNX_OFED Kernel modules on SLES12 SP5 and SLES15 kernel maintenance updates.
	Keywords: KMP, SLES, kernel
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2143258	Description: 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).

Internal Reference Number	Description
	Keywords: perftest
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2149577	Description: Fixed the issue where openibd script load used to fail when esp6_offload module did not load successfully.
	Keywords: openibd, esp6_offload
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2163879	Description: 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.
	Keywords: Dependency, perl-Getopt-Long
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2119017	Description: Fixed the issue where injecting EEH may cause extra Kernel prints, such as: "EEH: Might be infinite loop in mlx5_core driver".
	Keywords: EEH, kernel
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2107532	Description: 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.
	Keywords: Memory corruption, Rx page recycle
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-2.1.8.0
2116234	Description: Fixed the issue where ibsim was missing after OFED installation.
	Keywords: ibsim, installation
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2116233	Description: Fixed an issue where ucx-kmem was missing after OFED installation.
	Keywords: ucx-kmem, installation
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2109716	Description: Fixed a dependency issue between systemd and RDMA-Core.
	Keywords: Dependency, RDMA-Core

Internal Reference Number	Description
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2107776	Description: Fixed a driver load issue with Errata-kernel on SLES15 SP1.
	Keywords: Load, SLES, Errata
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2105536	Description: Fixed an issue in the Hairpin feature which prevented adding hairpin flows using TC tool.
	Keywords: Hairpin, TC
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2090321	Description: Fixed the issue where WQ queue flushing was not handled properly in the event of EEH.
	Keywords: WQ, EEH
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-2.1.8.0
2076311	Description: Fixed a rare kernel crash scenario when exiting an application that uses RMPP mads intensively.
	Keywords: MAD RMPP
	Discovered in Release: 4.0-1.0.1.0
	Fixed in Release: 5.0-2.1.8.0
2094545	Description: Fixed the issue where perftest applications (ib_read_*, ib_write_* and others) supplied with MLNX_OFED v5.0 and above did not work correctly if corresponding applications on another side of client-server communication were supplied with previous versions of MLNX_OFED due to an interoperability issue.
	Keywords: perftest, interoperability
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2096998	Description: Fixed the issue where NEO-Host could not be installed from the MLNX_OFED package when working on Ubuntu and Debian OSs.
	Keywords: NEO-Host, Ubuntu, Debian
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2094012	Description: Fixed the issue where MLNX_OFED installation failed to upgrade firmware version on ConnectX-6 Dx NICs with secure-fw.
	Keywords: ConnectX-6 Dx, installation, firmware, NIC

Internal Reference Number	Description
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2057076	Description: Added support for installing MLNX_OFED usingadd-kernel-support option over RHEL 8 OSs.
	Keywords:add-kernel-support, installation, RHEL
	Discovered in Release: 5.0-1.0.0.0
	Fixed in Release: 5.0-2.1.8.0
2090186	Description: Fixed a possible kernel crash scenario when AER/slot reset in done in parallel to user space commands execution.
	Keywords: mlx5_core, AER, slot reset
	Discovered in Release: 4.3-1.0.1.0
	Fixed in Release: 5.0-2.1.8.0
2093410	Description: Added missing ECN configuration under sysfs for PFs in SwitchDev mode.
	Keywords: sysfs, ASAP, SwitchDev, ECN
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-2.1.8.0
1731005	Description: Fixed the issue where MLNX_OFED v4.6 YUM and Zypper installations failed on RHEL8.0, SLES15.0 and PPCLE OSs.
	Keywords: YUM, Zypper, installation, RHEL, RedHat, SLES, PPCLE
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
1779150	Description: Fixed the issue of when upgrading the MLNX_OFED version over SLES 15 SPO and SP1 OSs on PPCLE platforms, it might have failed due to an isert-kmp-default issue.
	Keywords: Installation, SLES, PPCLE
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
1897199	Description: Fixed the issue of when using the RDMA statistics feature and attempting to unbind a QP from a counter, not including the counter-id as an argument in the CLI would have resulted in a segmentation fault.
	Keywords: RDMA, QP, segfault, unbinding
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
1916029	Description: 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.

Internal Reference Number	Description
	Keywords: Firmware, timeout, NULL
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2036394	Description: Added driver support for kernels with the old XDP_REDIRECT infrastructure that uses the following NetDev operations: .ndo_xdp_flush and .ndo_xdp_xmit.
	Keywords: XDP_REDIRECT, Soft lockup
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
1973238	Description: Fixed the issue where ib_core unload may fail on Ubuntu 18.04.2 OS with the following error message:
	"Module ib_core is in use"
	Keywords: ib_core, Ubuntu, ibacm
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.1-0.6.6.0
2072871	Description: Fixed an issue where the usage ofexcludedocs Open MPI RPM option resulted in the removal of non-documentation related files.
	Keywords:excludedocs, Open MPI, RPM
	Discovered in Release: 4.5-1.0.1.0
	Fixed in Release: 5.0-1.0.0.0
2060216	Description: Legacy mlnx-libs are now installed by default on SLES11 SP3 OS, as building MLNX_OFED on RDMA-Core based packages with this OS is not supported.
	Keywords: mlnx-libs, SLES, RDMA-Core
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2072884	Description: 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.
	Keywords: Installation, openibd, RDMA-Core
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2052037	Description: Disabled automated loading of some modules through udev triggers to preserve the startup process of previous MLNX_OFED versions.
	Keywords: Installation, udev, RDMA-Core

Internal Reference Number	Description
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2022634	Description: 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 optionwithout-depcheck.
	Keywords: Installation, SLES
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2022619	Description: Fixed the issue where uninstallation of MLNX_OFED would hang due to a bug in the package dependency check.
	Keywords: Uninstallation, dependency
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
1995843	Description: ibdump is now provided with the default rdma-core-based build.
	Keywords: ibdump, RDMA-Core
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
1995631	Description: Proper package dependencies are now set on Debian and Ubuntu libibverbs-dev package that is generated from RDMA-Core.
	Keywords: Dependency, libibverbs, RDMA-Core
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2047221	Description: 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.
	Keywords: rdma_resolve_addr(), cm_id
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
2045181	Description: Fixed a race condition which caused kernel panic when moving two ports to SwitchDev mode at the same time.
	Keywords: ASAP, SwitchDev, race
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
2004488	Description: Allowed accessing sysfs hardware counters in SwitchDev mode.

Internal Reference Number	Description
	Keywords: ASAP, hardware counters, sysfs, SwitchDev
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
2030943	Description: 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.
	Keywords: Rx page recycle, smp_processor_id
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
2074487	Description: Fixed an issue where port link state was automatically changed (without admin state involvement) to "UP" after reboot.
	Keywords: Link state, UP
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
2064711	Description: Fixed an issue where RDMA CM connection failed when port space was small.
	Keywords: RDMA CM
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
2076424	Description: Traffic mirroring with OVS offload and non-offload over VxLAN interface is now supported.
	Note : For kernel 4.9, make sure to use a dedicated OVS version.
	Keywords: VxLAN, OVS
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
1828321	Description: 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.
	Keywords: VF LAG, bonding, PF
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
1975293	Description: Installing OFED withwith-openvswitch flag no longer requires manual removal of the existing Open vSwitch.
	Keywords: OVS, Open vSwitch, openvswitch
	Discovered in Release: 4.7-3.2.9.0

Internal Reference Number	Description
	Fixed in Release: 5.0-1.0.0.0
1939719	Description: 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.
	Keywords: SLES, operating system, OS, installation, Kernel, module
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
2001966	Description: 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.
	Keywords: PF, VF, SwitchDev, netdevice, bonding
	Discovered in Release: 4.7-3.2.9.0
	Fixed in Release: 5.0-1.0.0.0
1816629	Description: 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.
	Keywords: Traffic, VF LAG
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
1718531	Description: Added support for VLAN header rewrite on CentOS 7.2 OS.
	Keywords: VLAN, ASAP, switchdev, CentOS 7.2
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 5.0-1.0.0.0
1556337	Description: Fixed the issue where adding VxLAN decapsulation rule with enc_tos and enc_ttl failed.
	Keywords: VxLAN, decapsulation
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 5.0-1.0.0.0
1921799	Description: Fixed the issue where MLNX_OFED installation over SLES15 SP1 ARM OSs failed unlessadd-kernel-support flag was added to the installation command.
	Keywords: SLES, installation
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 4.7-3.2.9.0
1949260	Description: Fixed a race condition that resulted in kernel panic when running IPoIB traffic in Connected mode.

Internal Reference Number	Description
	Keywords: IPoIB
	Discovered in Release: 4.5-1.0.1.0
	Fixed in Release: 4.7-3.2.9.0
1973828	Description: Fixed wrong EEPROM length for small form factor (SFF) 8472 from 256 to 512 bytes.
	Keywords: EEPROM, SFF
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 4.7-3.2.9.0
1915553	Description: Fixed the issue where errno field was not sent in all error flows of ibv_reg_mr API.
	Keywords: ibv_reg_mr
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-3.2.9.0
1970901	Description: Fixed the issue where mlx5 IRQ name did not change to express the state of the interface.
	Keywords: Ethernet, PCIe, IRQ
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 4.7-3.2.9.0
1915587	Description: Udaddy application is now functional in Legacy mode.
	Keywords: Udaddy, MLNX_OFED legacy, RDMA-CM
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 4.7-3.2.9.0
1931421	Description: Added support for E-Switch (SR-IOV Legacy) mode in RHEL 7.7 OSs.
	Keywords: E-Switch, SR-IOV, RHEL, RedHat
	Discovered in Release: 4.7-1.0.0.1
	Fixed in Release: 4.7-3.2.9.0
1945411/1839353	Description: Fixed the issue of when XDP_REDIRECT fails, pages got double-freed due to a bug in the refcnt_bias feature.
	Keywords: XDP, XDP_REDIRECT, refcnt_bias
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-3.2.9.0
1715789	Description: Fixed the issue where Mellanox Firmware Tools (MFT) package was missing from Ubuntu v18.04.2 OS.
	Keywords: MFT, Ubuntu, operating system
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-1.0.0.1

Internal Reference Number	Description
1547200	Description : Fixed an issue where IPoIB Tx queue may get stuck, leading to timeout warnings in dmesg.
	Keywords: IPoIB
	Discovered in Release: 4.5-1.0.1.0
	Fixed in Release: 4.7-1.0.0.1
1817636	Description: Fixed the issue of when disabling one port on the Server side, VF-LAG Tx Affinity would not work on the Client side.
	Keywords: VF-LAG, Tx Affinity
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-1.0.0.1
1800525	Description: 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.
	Keywords: Time-stamping, CQE compression
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-1.0.0.1
1431282	Description : Fixed the issue where software reset may have resulted in an order inversion of interface names.
	Keywords: Software reset
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.7-1.0.0.1
1843020	Description: Server reboot may result in a system crash.
	Keywords: reboot, crash
	Discovered in Release: 4.2-1.2.0.0
	Fixed in Release: 4.7-1.0.0.1
1734102	Description: Fixed the issue where Ubuntu v16.04.05 and v16.04.05 OSs could not be used with their native kernels.
	Keywords: Ubuntu, Kernel, OS
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-1.0.0.1
1811973	Description: VF mirroring offload is now supported.
	Keywords: ASAP ² , VF mirroring
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-1.0.0.1
1841634	Description: 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.
	Keywords: Counters, VF

Internal Reference Number	Description
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.7-1.0.0.1
1758983	Description: Installing MLNX_OFED on RHEL 7.6 OSs platform x86_64 and RHEL 7.6 ALT OSs platform PPCLE using YUM is now supported.
	Keywords: RHEL, RedHat, YUM, OS, operating system
	Discovered in Release: 4.6-1.0.1.1
	Fixed in Release: 4.7-1.0.0.1
1523548	Description: Fixed the issue where RDMA connection persisted even after dropping the network interface.
	Keywords: Network interface, RDMA
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.6-1.0.1.1
1712870	Description: 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: ethtoolset-priv-flags eth1 rx_no_csum_complete on/off
	Keywords: csum error, checksum, mlx5_core
	Discovered in Release: 4.5-1.0.1.0
	Fixed in Release: 4.6-1.0.1.1
1648597	Description : Fixed the wrong wording in the FW tracer ownership startup message (from "FW Tracer Owner" to "FWTracer: Ownership granted and active").
	Keywords: FW Tracer
	Discovered in Release: 4.5-1.0.1.0
	Fixed in Release: 4.6-1.0.1.1
1581631	Description : Fixed the issue where GID entries referenced to by a certain user application could not be deleted while that user application was running.
	Keywords: RoCE, GID
	Discovered in Release: 4.5-1.0.1.0
	Fixed in Release: 4.6-1.0.1.1
1368390	Description : Fixed the issue where MLNX_OFED could not be installed on RHEL 7.x Alt OSs using YUM repository.
	Keywords: Installation, YUM, RHEL
	Discovered in Release: 4.3-3.0.2.1
	Fixed in Release: 4.6-1.0.1.1
1531817	Description : 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.

Internal Reference Number	Description
	Keywords: Performance, Tx, CPU
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.5-1.0.1.0
1571977	Description : Fixed an issue of when the same CQ is connected to some QPs with SRQ and some without, wrong wr_id might be reported by ibv_poll_cq.
	Keywords: libmlx5, wr_id
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.5-1.0.1.0
1380135	Description : Fixed the issue where IB port link used to flap due to MAD heartbeat response delay when using new CQ API.
	Keywords: IB port link, CQ API, MAD heartbeat
	Discovered in Release: 4.2-1.2.0.0
	Fixed in Release: 4.5-1.0.1.0
1498931	Description : Fixed the issue where establishing TCP connection took too long due to failure of SA PathRecord query callback handler.
	Keywords: TCP, SA PathRecord
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.5-1.0.1.0
1514096	Description : Fixed the issue where lack of high order allocations caused driver load failure. All high order allocations are now changed to order-0 allocations.
	Keywords : mlx5, high order allocation
	Discovered in Release: 4.0-2.0.2.0
	Fixed in Release: 4.5-1.0.1.0
1524932	Description : 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.
	Keywords: mlx5 driver
	Discovered in Release: 4.0-2.0.2.0
	Fixed in Release: 4.5-1.0.1.0

Internal Reference Number	Description
1498585	Description : Fixed the issue of when performing configuration changes, mlx5e counters values were reset.
	Keywords: Ethernet counters
	Discovered in Release: 4.0-2.0.2.0
	Fixed in Release: 4.5-1.0.1.0
1425027	Description : 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: "ib_cache_gid_add: unable to add gid fe80:0000:0000:0000:f652:14ff:fe46:7391 error=-28".
	Keywords: RoCE, bonding
	Discovered in Release: 4.4-1.0.0.0
	Fixed in Release: 4.5-1.0.1.0

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.

Mellanox OFED is a single Virtual Protocol Interconnect (VPI) software stack which operates across all Mellanox network adapter solutions supporting the following uplinks to servers:

Uplink/NICs	Driver Name	Uplink Speed
ConnectX-4	mlx5	 InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 56GbE¹, 100GbE
ConnectX-4 Lx		Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE
ConnectX-5/ConnectX-5 Ex		 InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 100GbE
ConnectX-6		InfiniBand: SDR, FDR, EDR, HDR
		 Ethernet: 10GbE, 25GbE, 40GbE, 50GbE², 100GbE², 200GbE²
ConnectX-6 Dx		• Ethernet: 10GbE, 25GbE, 40GbE, 50GbE ² , 100GbE ² , 200GbE ²
ConnectX-6 Lx		• Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE ²
BlueField		 InfiniBand: SDR, QDR, FDR, FDR10, EDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE, 100GbE
BlueField-2		 InfiniBand: SDR, FDR, EDR, HDR Ethernet: 1GbE, 10GbE, 25GbE, 40GbE, 50GbE², 100GbE²

- 56GbE is a Mellanox propriety link speed and can be achieved while connecting a Mellanox adapter card to Mellanox SX10XX switch series, or connecting a Mellanox adapter card to another Mellanox
- 2. Supports both NRZ and PAM4 modes.

adapter card.

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

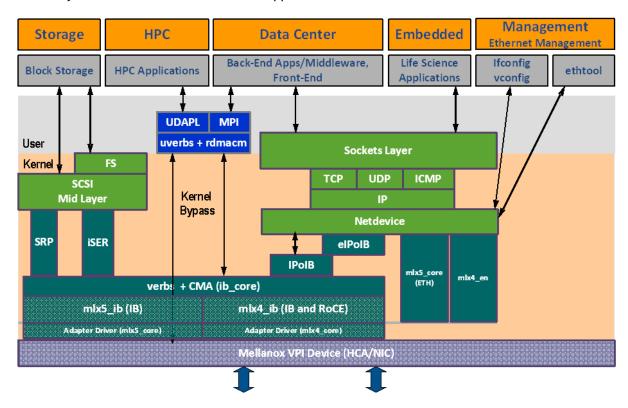
Mellanox OFED is certified with the following products:

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

- Mellanox 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 Mellanox 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 Mellanox 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 Mellanox OFED applies to.



The following subsections briefly describe the various components of the Mellanox OFED stack. mlx4 VPI Driver



⚠ This driver is no longer supported in MLNX_OFED. To work with ConnectX-3 and ConnectX-3 Pro NICs, please refer to MLNX_OFED LTS version available on the web.

mlx5 Driver

mlx5 is the low-level driver implementation for the Connect-IB® and ConnectX®-4 and above adapters designed by Mellanox Technologies. ConnectX®-4 and above adapter cards operate as a VPI adapter (Infiniband and Ethernet). The mlx5 driver is comprised of the following kernel modules:



Please note that Connect-IB card is no longer supported in MLNX_OFED. To work with this card, please refer to MLNX_OFED LTS version available on the web.

mlx5_core

Acts as a library of common functions (e.g. initializing the device after reset) required by ConnectX®-4 and above adapter cards. mlx5_core driver also implements the Ethernet interfaces for ConnectX®-4 and above. 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.

lihmlx5

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 Legacy and RDMA-Core 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 gueue 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

The following are libmlx5 Legacy only environment variables:

- MLX5_ENABLE_CQE_COMPRESSION
 - Saves PCIe bandwidth by compressing a few CQEs into a smaller amount of bytes on PCIe. Setting this variable enables CQE compression.
 - Default value 0 (disabled)
- MLX5_RELAXED_PACKET_ORDERING_ON

See "Out-of-Order (OOO) Data Placement" 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 preappends 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) 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. Mellanox OFED includes the following MPI implementation over InfiniBand:

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

Mellanox OFED also includes MPI benchmark tests such as OSU BW/LAT, Intel MPI BeBenchmarkand 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 Mellanox OFED¹.

1. OpenSM is disabled by default. See "OpenSM" section for details on enabling it.

Diagnostic Utilities

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

- ibutils Mellanox Technologies diagnostic utilities
- infiniband-diags OpenFabrics Alliance InfiniBand diagnostic tools

Mellanox Firmware Tools

The Mellanox 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 Mellanox 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 (https://docs.mellanox.com/category/mft).

Mellanox OFED Package

ISO Image

Mellanox 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:

- Mellanox Host Channel Adapter Drivers
 - mlx5
 - mlx5 ib
 - mlx5_core (includes 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.

⚠ MLNX_OFED includes the OFED source RPM packages used as a build platform for kernel code but does not include the sources of Mellanox proprietary packages.

- 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 Mellanox OFED related documentation

Module Parameters

mlx5_core Module Parameters

The mlx5_core module supports a single parameter used to select the profile which defines the number of resources supported.

prof_sel	The parameter name for selecting the profile. The supported values for profiles are: • 0 - for medium resources, medium performance • 1 - for low resources • 2 - for high performance (int) (default)
guids	charp
node_guid	guids configuration. This module parameter will be obsolete!
debug_mask	<pre>debug_mask: 1 = dump cmd data, 2 = dump cmd exec time, 3 = both. Default=0 (uint)</pre>
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)

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 Mellanox OFED
- Installing Mellanox OFED
- Uninstalling Mellanox OFED
- Updating Firmware After Installation
- UEFI Secure Boot
- · Performance Tuning

Hardware and Software Requirements

Requirements	Description
Platforms	A server platform with an adapter card based on one of the following Mellanox Technologies' VPI HCA devices listed in Supported NICs Speeds table.
	For the list of supported architecture platforms, please refer to the Mellanox OFED Release Notes file.
Required Disk Space for Installation	1GB
Device ID	For the latest list of device IDs, please visit Mellanox website.
Operating System	Linux operating system.
	For the list of supported operating system distributions and kernels, please refer to the Mellanox OFED Release Notes file.
Installer Privileges	The installation requires administrator (root) privileges on the target machine.

Downloading Mellanox OFED

1. Verify that the system has a Mellanox network adapter (HCA/NIC) installed.

The following example shows a system with an installed Mellanox HCA:

Note: For ConnectX-5 Socket Direct adapters, use ibdev2netdev to display the installed card and the mapping of logical ports to physical ports. Example:

```
[root@gen-1-vrt-203 ~]# ibdev2netdev -v | grep -i MCX556M-ECAT-S25 0000:84:00.0 mlx5_10 (MT4119 - MCX556M-ECAT-S25SN) CX556M - ConnectX-5 QSFP28 fw 16.22.0228 port 1 (DOWN ) ==> p2p1 (Down) 0000:84:00.1 mlx5_11 (MT4119 - MCX556M-ECAT-S25SN) CX556M - ConnectX-5 QSFP28 fw 16.22.0228 port 1 (DOWN ) ==> p2p2 (Down) 0000:05:00.0 mlx5_2 (MT4119 - MCX556M-ECAT-S25SN) CX556M - ConnectX-5 QSFP28 fw 16.22.0228 port 1 (DOWN ) ==> p5p1 (Down) 0000:05:00.1 mlx5_3 (MT4119 - MCX556M-ECAT-S25SN) CX556M - ConnectX-5 QSFP28 fw 16.22.0228 port 1 (DOWN ) ==> p5p2 (Down)
```

Notes:

- Each PCI card of ConnectX-5 Socket Direct has a different PCI address. In the output example above, the first two rows indicate that one card is installed in a PCI slot with PCI Bus address 84 (hexadecimal), and PCI Device Number 00, and PCI Function Number 0 and 1. RoCE assigned mlx5_10 as the logical port, which is the same as netdevice p2p1, and both are mapped to physical port of PCI function 0000:84:00.0.
- RoCE logical port mlx5_2 of the second PCI card (PCI Bus address 05) and netdevice p5p1 are mapped to physical port of PCI function 0000:05:00.0, which is the same physical port of PCI function 0000:84:00.0.

MT4119 is the PCI Device ID of the Mellanox ConnectX-5 adapters family.

For more details, please refer to ConnectX-5 Socket Direct Hardware User Manual, available at www.mellanox.com.

2. Download the ISO image to your host.

The image's name has the format MLNX_OFED_LINUX-<ver>--<OS label><CPU arch>.iso. You can download it from http://www.mellanox.com --> Products --> Software --> InfiniBand/VPI Drivers --> Mellanox OFED Linux (MLNX_OFED).

- a. Scroll down to the Download wizard, and click the Download tab.
- b. Choose your relevant package depending on your host operating system.
- c. Click the desired ISO/tgz package.
- d. To obtain the download link, accept the End User License Agreement (EULA).
- 3. Use the md5sum utility to confirm the file integrity of your ISO image. Run the following command and compare the result to the value provided on the download page.

host1\$ md5sum MLNX OFED LINUX-<ver>-<0S label>.iso

Installing Mellanox OFED

Installation Script

The installation script, mlnxofedinstall, performs the following:

- Discovers the currently installed kernel
- Uninstalls any software stacks that are part of the standard operating system distribution or another vendor's commercial stack
- Installs the MLNX_OFED_LINUX binary RPMs (if they are available for the current kernel)
- Identifies the currently installed InfiniBand and Ethernet network adapters and automatically upgrades the firmware

Note: If you wish to perform a firmware upgrade using customized FW binaries, you can provide a path to the folder that contains the FW binary files, by running --fw-image-dir. Using this option, the FW version embedded in the MLNX_OFED package will be ignored. Example:

./mlnxofedinstall --fw-image-dir /tmp/my_fw_bin_files



A If the driver detects unsupported cards on your system, it will abort the installation procedure. To avoid this, make sure to add --skip-unsupported-devices-check flag during installation.

Usage

./mnt/mlnxofedinstall [OPTIONS]

The installation script removes all previously installed Mellanox OFED packages and re-installs from scratch. You will be prompted to acknowledge the deletion of the old packages.



Pre-existing configuration files will be saved with the extension ".conf.rpmsave".

- If you need to install Mellanox OFED on an entire (homogeneous) cluster, a common strategy is to mount the ISO image on one of the cluster nodes and then copy it to a shared file system such as NFS. To install on all the cluster nodes, use cluster-aware tools (suchaspdsh).
- If your kernel version does not match with any of the offered pre-built RPMs, you can add your kernel version by using the "mlnx_add_kernel_support.sh" script located inside the MLNX_OFED package.



On Redhat and SLES distributions with errata kernel installed there is no need to use the mlnx_add_kernel_support.sh script. The regular installation can be performed and weak-updates mechanism will create symbolic links to the MLNX_OFED kernel modules.



▲ If you regenerate kernel modules for a custom kernel (using --add-kernelsupport), the packages installation will not involve automatic regeneration of the initramfs. In some cases, such as a system with a root filesystem mounted over a ConnectX card, not regenerating the initramfs may even cause the system to fail to reboot.

In such cases, the installer will recommend running the following command to update the initramfs:

dracut -f

On some OSs, dracut -f might result in the following error message which can be safely ignore.

libkmod: kmod_module_new_from_path: kmod_module 'mdev' already exists with different path

The "mlnx_add_kernel_support.sh" script can be executed directly from the mlnxofedinstall script. For further information, please see '--add-kernel-support' option below.

On Ubuntu and Debian distributions drivers installation use Dynamic Kernel Module Support (DKMS) framework. Thus, the drivers' compilation will take place on the host during MLNX_OFED installation.

Therefore, using "mlnx_add_kernel_support.sh" is irrelevant on Ubuntu and Debian distributions.

Example: The following command will create a MLNX_OFED_LINUX ISO image for RedHat 7.3 under the /tmp directory.

```
Building OFED RPMs. Please wait...
Removing OFED RPMs...
Created /tmp/MLNX_OFED_LINUX-x.x-x-rhel7.3-x86_64-ext.tgz
```

- The script adds the following lines to /etc/security/limits.conf for the userspace components such as MPI:
 - * soft memlock unlimited
 - * hard memlock unlimited
 - These settings set the amount of memory that can be pinned by a userspace application to unlimited. If desired, tune the value unlimited to a specific amount of RAM.

For your machine to be part of the InfiniBand/VPI fabric, a Subnet Manager must be running on one of the fabric nodes. At this point, Mellanox OFED for Linux has already installed the OpenSM Subnet Manager on your machine.

For the list of installation options, run: ./mlnxofedinstall --h



⚠ The DKMS (on Debian based OS) and the weak-modules (RedHat OS) mechanisms rebuild the initrd/initramfs for the respective kernel in order to add the MLNX_OFED drivers. When installing MLNX_OFED without DKMS support on Debian based OS, or without KMP support on RedHat or any other distribution, the initramfs will not be changed. Therefore, the inbox drivers may be loaded on boot. In this case, openibd service script will automatically unload them and load the new drivers that come with MLNX_OFED.

Installation Procedure

This section describes the installation procedure of MLNX_OFED on Mellanox adapter cards.

- 1. Log in to the installation machine as root.
- Mount the ISO image on your machine.

```
host1# mount -o ro,loop MLNX_OFED_LINUX-<ver>-<OS label>-<CPU arch>.iso /mnt
```

3. Run the installation script.

```
/mnt/mlnxofedinstall
/mmt/mlnxofedinstall
Logs dir: /tmp/MLNX_OFED_LINUX-x.x-x.logs
This program will install the MLNX_OFED_LINUX package on your machine.
Note that all other Mellanox, OEM, OFED, RDMA or Distribution IB packages will be removed.
Those packages are removed due to conflicts with MLNX_OFED_LINUX, do not reinstall them.
Starting MLNX_OFED_LINUX-x.x.x installation ...
Installation finished successfully.
Attempting to perform Firmware update...
Querying Mellanox devices firmware ...
```

For unattended installation, use the --force installation option while running the MLNX_OFED installation script:

/mnt/mlnxofedinstall --force

MLNX_OFED for Ubuntu should be installed with the following flags in chroot environment:

./mlnxofedinstall --without-dkms --add-kernel-support --kernel <kernel version in chroot> --without-fw-update --force

For example:

./mlnxofedinstall --without-dkms --add-kernel-support --kernel 3.13.0-85-generic -without-fw-update --force

Note that the path to kernel sources (--kernel-sources) should be added if the sources are not in their default location.

In case your machine has the latest firmware, no firmware update will occur and the installation script will print at the end of installation a message similar to the following:

Device #1:

Device Type: ConnectX4 Part Number: MCX456A-ECA

Description: ConnectX-4 VPI adapter card; EDR IB (100Gb/s) and 100GbE; dual-port

QSFP28; PCle3.0 x16; ROHS R6

PSID: MT_2190110032 PCI Device Name: 0b:00.0 Base MAC: 0000e41d2d5cf810 Versions: Current Available FW 12.14.0114 12.14.0114

Status: Up to date

In case your machine has an unsupported network adapter device, no firmware update will occur and one of the error messages below will be printed. Please contact your hardware vendor for help with firmware updates.

Error message #1:

Device #1:

Device Type: ConnectX4 Part Number: MCX456A-ECA

Description: ConnectX-4 VPI adapter card; EDR IB (100Gb/s) and 100GbE; dual-port

QSFP28; PCle3.0 x16; ROHS R6

PSID: MT_2190110032 PCI Device Name: 0b:00.0 Base MAC: 0000e41d2d5cf810 Versions: Current Available FW 12.14.0114 N/A

Status: No matching image found

Error message #2:

The firmware for this device is not distributed inside Mellanox driver: 0000:01:00.0

(PSID: IBM2150110033)

To obtain firmware for this device, please contact your HW vendor.

4. Case A: If the installation script has performed a firmware update on your network adapter, you need to either restart the driver or reboot your system before the firmware update can take effect. Refer to the table below to find the appropriate action for your specific card.

Action \ Adapter	Driver Restart	Standard Reboot (Soft Reset)	Cold Reboot (Hard Reset)
Standard ConnectX-4/ ConnectX-4 Lx or higher	-	+	-
Adapters with Multi- Host Support	-	-	+
Socket Direct Cards	-	-	+

Case B: If the installations script has not performed a firmware upgrade on your network adapter, restart the driver by running: "/etc/init.d/openibd restart".

- 5. (InfiniBand only) Run the hca_self_test.ofed utility to verify whether or not the InfiniBand link is up. The utility also checks for and displays additional information such as:
 - · HCA firmware version
 - · Kernel architecture
 - Driver version
 - Number of active HCA ports along with their states

For more details on hca_self_test.ofed, see the file docs/readme_and_user_manual/ hca_self_test.readme.

After installation completion, information about the Mellanox OFED installation, such as prefix, kernel version, and installation parameters can be retrieved by running the command /etc/ infiniband/info. Most of the Mellanox OFED components can be configured or reconfigured after the installation, by modifying the relevant configuration files. See the relevant chapters in this manual for details.

The list of the modules that will be loaded automatically upon boot can be found in the /etc/ infiniband/openib.conf file.



⚠ Installing Mellanox OFED will replace the RDMA stack and remove existing 3rd party RDMA connectors.

Installation Results

Software	 Most of MLNX_OFED packages are installed under the "/usr" directory except for the following packages which are installed under the "/opt" directory: fca and ibutils iproute2 (rdma tool) - installed under /opt/Mellanox/iproute2/sbin/rdma The kernel modules are installed under /lib/modules/`uname -r`/updates on SLES and Fedora Distributions /lib/modules/`uname -r`/extra/mlnx-ofa_kernel on RHEL and other RedHat like Distributions /lib/modules/`uname -r`/updates/dkms/ on Ubuntu
Firmware	 The firmware of existing network adapter devices will be updated if the following two conditions are fulfilled: The installation script is run in default mode; that is, without the option 'without- fw-update' The firmware version of the adapter device is older than the firmware version included with the Mellanox OFED ISO image Note: If an adapter's Flash was originally programmed with an Expansion ROM image, the automatic firmware update will also burn an Expansion ROM image. In case your machine has an unsupported network adapter device, no firmware update will occur and the error message below will be printed. The firmware for this device is not distributed inside Mellanox driver: 0000:01:00.0 (PSID: IBM2150110033) To obtain firmware for this device, please contact your HW vendor."

Installation Logging

While installing MLNX_OFED, the install log for each selected package will be saved in a separate log file.

The path to the directory containing the log files will be displayed after running the installation script in the following format:

Example:

Logs dir: /tmp/MLNX_OFED_LINUX-4.4-1.0.0.0.IBMM2150110033.logs

Driver Load Upon System Boot

Upon system boot, the Mellanox drivers will be loaded automatically.

- To prevent the automatic load of the Mellanox drivers upon system boot:
 - 1. Add the following lines to the "/etc/modprobe.d/mlnx.conf" file.

```
blacklist mlx5_core
blacklist mlx5_ib
```

- 2. Set "ONBOOT=no" in the "/etc/infiniband/openib.conf" file.
- If the modules exist in the initramfs file, they can automatically be loaded by the kernel. To
 prevent this behavior, update the initramfs using the operating systems' standard tools.
 Note: The process of updating the initramfs will add the blacklists from step 1, and will
 prevent the kernel from loading the modules automatically.

mlnxofedinstall Return Codes

The table below lists the mlnxofedinstall script return codes and their meanings.

Return Code	Meaning
0	The Installation ended successfully
1	The installation failed
2	No firmware was found for the adapter device
22	Invalid parameter
28	Not enough free space
171	Not applicable to this system configuration. This can occur when the required hardware is not present on the system
172	Prerequisites are not met. For example, missing the required software installed or the hardware is not configured correctly
173	Failed to start the mst driver

Additional Installation Procedures

Installing MLNX_OFED Using YUM

This type of installation is applicable to RedHat/OL and Fedora operating systems.

Setting up MLNX_OFED YUM Repository

- 1. Log into the installation machine as root.
- 2. Mount the ISO image on your machine and copy its content to a shared location in your network.

```
# mount -o ro,loop MLNX_OFED_LINUX-<ver>--<OS label>--<CPU arch>.iso /mnt
```

Download and install Mellanox Technologies GPG-KEY:
 The key can be downloaded via the following link:
 http://www.mellanox.com/downloads/ofed/RPM-GPG-KEY-Mellanox

4. Install the key.

```
# sudo rpm --import RPM-GPG-KEY-Mellanox
warning: rpmts_HdrFromFdno: Header V3 DSA/SHA1 Signature, key ID 6224c050: NOKEY
Retrieving key from file:///repos/MLNX_OFED/<MLNX_OFED file>/RPM-GPG-KEY-Mellanox
Importing GPG key 0x6224c050:
Userid: "Mellanox Technologies (Mellanox Technologies - Signing Key v2) <support@mellanox.com>"
From : /repos/MLNX_OFED/<MLNX_OFED file>/RPM-GPG-KEY-Mellanox
Is this ok [y/N]:
```

5. Check that the key was successfully imported.

```
# rpm -q gpg-pubkey --qf '%{NAME}-%{VERSION}-%{RELEASE}\t%{SUMMARY}\n' | grep Mellanox
gpg-pubkey-a9e4b643-520791ba gpg(Mellanox Technologies <support@mellanox.com>)
```

6. Create a yum repository configuration file called "/etc/yum.repos.d/mlnx_ofed.repo" with the following content:

```
[mlnx_ofed]
name=MLNX_OFED Repository
baseurl=file:///<path to extracted MLNX_OFED package>/RPMS
enabled=1
gpgkey=file:///<path to the downloaded key RPM-GPG-KEY-Mellanox>
gpgcheck=1
```

7. Check that the repository was successfully added.

```
# yum repolist
Loaded plugins: product-id, security, subscription-manager
This system is not registered to Red Hat Subscription Management. You can use subscription-manager to register.
repo id repo name status
mlnx_ofed MLNX_OFED Repository 108
rpmforge RHEL 6Server - RPMforge.net - dag 4,597
repolist: 8,351
```

Installing MLNX_OFED Using the YUM Tool

After setting up the YUM repository for MLNX_OFED package, perform the following:

1. View the available package groups by invoking:

```
# yum search mlnx-ofed-
mlnx-ofed-all.noarch: MLNX_OFED all installer package (with KMP support)
mlnx-ofed-all-user-only.noarch: MLNX_OFED all-user-only installer package (User Space packages only)
mlnx-ofed-basic-user-only.noarch: MLNX_OFED basic installer package (with KMP support)
mlnx-ofed-basic-user-only.noarch: MLNX_OFED basic-user-only installer package (User Space packages only)
mlnx-ofed-bluefield.noarch: MLNX_OFED bluefield installer package (with KMP support)
mlnx-ofed-bluefield-user-only.noarch: MLNX_OFED bluefield-user-only installer package (User Space packages only)
mlnx-ofed-dpdk.noarch: MLNX_OFED dpdk installer package (with KMP support)
mlnx-ofed-dpdk-upstream-libs.noarch: MLNX_OFED dpdk-upstream-libs installer package (with KMP support)
mlnx-ofed-dpdk-upstream-libs-user-only.noarch: MLNX_OFED dpdk-upstream-libs-user-only installer package (User Space packages only)
mlnx-ofed-dpdk-user-only.noarch: MLNX_OFED dpdk-user-only installer package (User Space packages only)
mlnx-ofed-dpdk-user-only.noarch: MLNX_OFED dpdk-user-only installer package (User Space packages only)
mlnx-ofed-guest.noarch: MLNX_OFED guest installer package (with KMP support)
mlnx-ofed-guest-user-only.noarch: MLNX_OFED guest-user-only installer package (User Space packages only)
mlnx-ofed-hpc-user-only.noarch: MLNX_OFED bpc installer package (with KMP support)
mlnx-ofed-hpc-user-only.noarch: MLNX_OFED hpc-user-only installer package (User Space packages only)
mlnx-ofed-hpc-user-only.noarch: MLNX_OFED hpc-user-only installer package (with KMP support)
mlnx-ofed-hpc-user-only.noarch: MLNX_OFED hpc-user-only installer package (with KMP support)
mlnx-ofed-wna-noarch: MLNX_OFED wma-th-user-only installer package (with KMP support)
mlnx-ofed-wna-noarch: MLNX_OFED wma-eth-installer package (With KMP support)
mlnx-ofed-vma-eth-noarch: MLNX_OFED vma-eth installer package (User Space packages only)
mlnx-ofed-vma-eth-user-only.noarch: MLNX_OFED vma-eth-user-only installer package (User Space packages only)
mlnx-ofed-vma-vpi-noarch: MLNX
```

where:

WHICH C.	
mlnx-ofed-all	Installs all available packages in MLNX_OFED
mlnx-ofed-basic	Installs basic packages required for running Mellanox cards
mlnx-ofed-guest	Installs packages required by guest OS
mlnx-ofed-hpc	Installs packages required for HPC
mlnx-ofed-hypervisor	Installs packages required by hypervisor OS
mlnx-ofed-vma	Installs packages required by VMA
mlnx-ofed-vma-eth	Installs packages required by VMA to work over Ethernet
mlnx-ofed-vma-vpi	Installs packages required by VMA to support VPI
bluefield	Installs packages required for BlueField
dpdk	Installs packages required for DPDK
dpdk-upstream-libs	Installs packages required for DPDK using RDMA-Core
kernel-only	Installs packages required for a non-default kernel

Note: MLNX_OFED provides kernel module RPM packages with KMP support for RHEL and SLES. For other operating systems, kernel module RPM packages are provided only for the operating system's default kernel. In this case, the group RPM packages have the supported kernel version in their package's name.

Example:

```
mlnx-ofed-all-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED all installer package for kernel 3.17.4-301.fc21.x8
mlnx-ofed-all-3.17.4-301.fc21.x80_64.noarch : MLNX_OFED basic installer package for kernel 3.17.4-301.fc2 1.x86_64 (without KMP support) mlnx-ofed-basic-3.17.4-301.fc21.x86_64 (without KMP support) mlnx-ofed-guest-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED guest installer package for kernel 3.17.4-301.fc2
1.886_64 (Without KMP Support)
mlnx-ofed-guest-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED guest installer package for kernel 3.17.4-301.fc21.x8
mlnx-ofed-hpc-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED hpc installer package for kernel 3.17.4-301.fc21.x8
mlnx-ofed-wpervisor-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED hypervisor installer package for kernel 3.17.4 01.fc21.x86_64 (without KMP support) mlnx-ofed-wma-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED vma installer package for kernel 3.17.4-301.fc21.x86_64 (without KMP support) mlnx-ofed-vma-eth-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED vma-eth installer package for kernel 3.17.4-301.
                                                           .4-301.fc21.x86_64.noarch : MLNX_OFED hypervisor installer package for kernel 3.17.4-
mlnx-ofed-vma-eth-3.17.4-301.fc21.x86_04.noarch : MLNX_OFED vma-vpi installer package for kernel 3.17.4-301. fc21.x86_64 (without KMP support) mlnx-ofed-vma-vpi-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED vma-vpi installer package for kernel 3.17.4-301. fc21.x86_64 (without KMP support) mlnx-ofed-hypervisor-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED hypervisor installer package for kernel 3.17.4-301.
mlnx-ofed-hypervisor-3.17.4-301.
01.fc21.x86_64 (without KMP supp
mlnx-ofed-vma-3.17.4-301.fc21.x86_64.noarch: MLNX_OFED vma installer package for kernel 3.17.4-301.fc21.x86_64 (without KMP support)
mlnx-ofed-vma-eth-3.17.4-301.fc21.x86_64.noarch: MLNX_OFED vma-eth installer package for kernel 3.17.4-301.
mlnx-ofed-vma-etn-5.17.4-301.c21.x86_fc21.x86_64 (without KMP support)
mlnx-ofed-vma-vpi-3.17.4-301.fc21.x86_64.noarch : MLNX_OFED vma-vpi installer package for kernel 3.17.4-301.
mlnx-ofed-vma-vpi-3.17.4-301.fc21
fc21.x86_64 (without KMP support)
```

If you have an operating system different than RHEL or SLES, or you have installed a kernel that is not supported by default in MLNX_OFED, you can use the mlnx_add_kernel_support.sh script to build MLNX_OFED for your kernel.

The script will automatically build the matching group RPM packages for your kernel so that you can still install MLNX_OFED via yum.

Please note that the resulting MLNX_OFED repository will contain unsigned RPMs, therefore, you should set 'gpgcheck=0' in the repository configuration file.

2. Install the desired group.

```
# vum install mlnx-ofed-all
Loaded plugins: langpacks, product-id, subscription-manager Resolving Dependencies
--> Running transaction check
---> Package mlnx-ofed-all.noarch 0:3.1-0.1.2 will be installed --> Processing Dependency: kmod-isert = 1.0-OFED.3.1.0.1.2.1.g832a737.rhel7u1 for package: mlnx-ofed-all-3.1-0.1.2.noarch
   qperf.x86_64 0:0.4.9-9
rds-devel.x86_64 0:2.0.7-1.12
rds-tools.x86_64 0:2.0.7-1.12
sdpnetstat.x86_64 0:1.60-26
Complete!
```

Installing MLNX_OFED using the "YUM" tool does not automatically update the firmware. To update the firmware to the version included in MLNX_OFED package, run:

```
# yum install mlnx-fw-updater
OR:
```

Update the firmware to the latest version available on Mellanox Technologies' Web site as described in "Updating Firmware After Installation" section.

Installing MLNX_OFED Using apt-get

This type of installation is applicable to Debian and Ubuntu operating systems.

Setting up MLNX OFED apt-get Repository

1. Log into the installation machine as root.

- Extract the MLNX_OFED package on a shared location in your network.
 You can download it from http://www.mellanox.com > Products > Software > InfiniBand Drivers.
- Create an apt-get repository configuration file called "/etc/apt/sources.list.d/mlnx_ofed.list" with the following content:

```
deb file:/<path to extracted MLNX_OFED package>/DEBS ./
```

4. Download and install Mellanox Technologies GPG-KEY.

```
# wget -qO - http://www.mellanox.com/downloads/ofed/RPM-GPG-KEY-Mellanox | sudo apt-key add -
```

5. Verify that the key was successfully imported.

```
# apt-key list
pub 1024D/A9E4B643 2013-08-11
uid Mellanox Technologies <support@mellanox.com>
sub 1024g/09FCC269 2013-08-11
```

6. Update the apt-get cache.

```
# sudo apt-get update
```

Installing MLNX_OFED Using the apt-get Tool

After setting up the apt-get repository for MLNX_OFED package, perform the following:

1. View the available package groups by invoking:

where:

mlnx-ofed-all	MLNX_OFED all installer package
mlnx-ofed-basic	MLNX_OFED basic installer package
mlnx-ofed-vma	MLNX_OFED vma installer package
mlnx-ofed-hpc	MLNX_OFED HPC installer package
mlnx-ofed-vma-eth	MLNX_OFED vma-eth installer package

mlnx-ofed-vma-vpi	MLNX_OFED vma-vpi installer package
knem-dkms	MLNX_OFED DKMS support for mlnx-ofed kernel modules
kernel-dkms	MLNX_OFED kernel-dkms installer package
kernel-only	MLNX_OFED kernel-only installer package
bluefield	MLNX_OFED bluefield installer package
mlnx-ofed-all-exact	MLNX_OFED mlnx-ofed-all-exact installer package
dpdk	MLNX_OFED dpdk installer package
mlnx-ofed-basic-exact	MLNX_OFED mlnx-ofed-basic-exact installer package
dpdk-upstream-libs	MLNX_OFED dpdk-upstream-libs installer package

2. Install the desired group.

```
apt-get install '<group name>'
```

Example:

apt-get install mlnx-ofed-all



⚠ Installing MLNX_OFED using the "apt-get" tool does not automatically update the

To update the firmware to the version included in MLNX_OFED package, run: # apt-get install mlnx-fw-updater

OR:

Update the firmware to the latest version available on Mellanox Technologies' Web site as described in "Updating Firmware After Installation" section.

Installing NEO-Host Using mlnxofedinstall Script

As of MLNX_OFED v4.5, NEO-Host users can opt to install the NEO-Host package embedded in MLNX_OFED package.

In order to install NEO-Host, add the --with-neohost-backend flag to the mlnxofedinstall script

Example:

/mnt/mlnxofedinstall --with-neohost-backend

Uninstalling Mellanox OFED

Uninstalling Mellanox OFED Using the YUM and apt-get Tools

Use the script /usr/sbin/ofed_uninstall.sh to uninstall the Mellanox OFED package. The script is part of the ofed-scripts RPM.

Updating Firmware After Installation

The firmware can be updated using one of the following methods:

Updating the Device Online

To update the device online on the machine from Mellanox site, use the following command line:

```
mlxfwmanager --online -u -d <device>
```

Example:

```
# mlxfwmanager --online -u -d 0000:01:00.0
Querying Mellanox devices firmware ...

Device #1:
--------

Device Type: ConnectX6
Part Number: MCX653106A-HDA_Ax
Description: ConnectX-6 VPI adapter card; HDR IB (200Gb/s) and 200GbE; dual-port QSFP56; PCIe4.0 x16; tall bracket; ROHS R6
PSID: MT_0000000225
PCI Device Name: 0000:01:00.0
Base MAC: 98039b970cc2
Versions: Current Available
FW 20.26.4012 20.27.1016
PXE 3.6.0101 3.5.0903
UEFI 14.21.0016 14.20.0025

Status: Up to date
```

Updating the Device Manually

In case you ran the mlnxofedinstall script with the '--without-fw-update' option or you are using an OEM card and now you wish to (manually) update firmware on your adapter card(s), you need to perform the steps below. The following steps are also appropriate in case you wish to burn newer firmware that you have downloaded from Mellanox Technologies' Web site (http://www.mellanox.com > Support > Firmware Download).

1. Get the device's PSID.

```
mlxfwmanager_pci | grep PSID
PSID: MT_1210110019
```

- 2. Download the firmware BIN file from the Mellanox website or the OEM website.
- 3. Burn the firmware.

```
mlxfwmanager_pci -i <fw_file.bin>
```

4. Reboot your machine once the firmware burning is completed.

Updating the Device Firmware Automatically upon System Boot

Firmware can be automatically updated upon system boot.

The firmware update package (mlnx-fw-updater) is installed in the "/opt/mellanox/mlnx-fw-updater" folder, and the openibd service script can invoke the firmware update process if requested on boot.

If the firmware is updated, the following message will be printed to the system's standard logging file:

fw_updater: Firmware was updated. Please reboot your system for the changes to take effect.

Otherwise, the following message will be printed:

fw_updater: Didn't detect new devices with old firmware.

Please note that this feature is disabled by default. To enable the automatic firmware update upon system boot, set the following parameter to "yes" "RUN_FW_UPDATER_ONBOOT=yes" in the openibd service configuration file "/etc/infiniband/openib.conf".

You can opt to exclude a list of devices from the automatic firmware update procedure. To do so, edit the configurations file "/opt/mellanox/mlnx-fw-updater/mlnx-fw-updater.conf" and provide a comma separated list of PCI devices to exclude from the firmware update. Example:

MLNX_EXCLUDE_DEVICES="00:05.0,00:07.0"

UEFI Secure Boot

All kernel modules included in MLNX_OFED for RHEL7 and SLES12 are signed with x.509 key to support loading the modules when Secure Boot is enabled.

Enrolling Mellanox's x.509 Public Key On your Systems

In order to support loading MLNX_OFED drivers when an OS supporting Secure Boot boots on a UEFI-based system with Secure Boot enabled, the Mellanox x.509 public key should be added to the UEFI Secure Boot key database and loaded onto the system key ring by the kernel.

Follow these steps below to add the Mellanox's x.509 public key to your system:

- A
- Prior to adding the Mellanox's x.509 public key to your system, please make sure:
 - The 'mokutil' package is installed on your system
 - The system is booted in UEFI mode
- 1. Download the x.509 public key.

wget http://www.mellanox.com/downloads/ofed/mlnx_signing_key_pub.der

Add the public key to the MOK list using the mokutil utility.You will be asked to enter and confirm a password for this MOK enrollment request.

```
# mokutil --import mlnx_signing_key_pub.der
```

3. Reboot the system.

The pending MOK key enrollment request will be noticed by shim.efi and it will launch MokManager.efi to allow you to complete the enrollment from the UEFI console. You will need to enter the password you previously associated with this request and confirm the enrollment. Once done, the public key is added to the MOK list, which is persistent. Once a key is in the MOK list, it will be automatically propagated to the system key ring and subsequent will be booted when the UEFI Secure Boot is enabled.



To see what keys have been added to the system key ring on the current boot, install the 'keyutils' package and run: #keyctl list %:.system_keyring

Removing Signature from Kernel Modules

The signature can be removed from a signed kernel module using the 'strip' utility which is provided by the 'binutils' package.

```
# strip -g my_module.ko
```

The strip utility will change the given file without saving a backup. The operation can be undone only by resigning the kernel module. Hence, we recommend backing up a copy prior to removing the signature.

To remove the signature from the MLNX_OFED kernel modules:

1. Remove the signature.

```
# rpm -qa | grep -E "kernel-ib|mlnx-ofa_kernel|iser|srp|knem|mlnx-rds|mlnx-nfsrdma|mlnx-nvme|mlnx-rdma-rxe" | xargs rpm -ql | grep "\.ko$" | xargs strip -g
```

Once the signature is removed, a message as the below will no longer be presented upon module loading:

```
"Request for unknown module key 'Mellanox Technologies signing key: 61feb074fc7292f958419386ffdd9d5ca999e403' err -11"
```

However, please note that a message similar to the following will be presented:

```
"my_module: module verification failed: signature and/or required key missing - tainting kernel"
```

This message is presented once, only upon first module boot that either has no signature or whose key is not in the kernel key ring. Therefore, this message may go unnoticed. Once the system is rebooted after unloading and reloading a kernel module, the message will appear (this message cannot be eliminated).

2. Update the initramfs on RHEL systems with the stripped modules.

```
mkinitrd /boot/initramfs-$(uname -r).img $(uname -r) --force
```

Performance Tuning

Depending on the application of the user's system, it may be necessary to modify the default configuration of network adapters based on the ConnectX® adapters. In case tuning is required, please refer to the Performance Tuning for Mellanox Adapters Community post.

Features Overview and Configuration

Unable to render include or excerpt-include. Could not retrieve page.

The chapter contains the following sections:

- Ethernet Network
- InfiniBand Network
- Storage Protocols
- Virtualization
- Resiliency
- Docker Containers
- HPC-X™
- Fast Driver Unload
- OVS Offload Using ASAP² Direct

Ethernet Network

The chapter contains the following sections:

- Ethernet Interface
- Quality of Service (QoS)
- Ethtool
- Checksum Offload
- Ignore Frame Check Sequence (FCS) Errors
- RDMA over Converged Ethernet (RoCE)
- Flow Control
- Explicit Congestion Notification (ECN)
- RSS Support
- Time-Stamping
- Flow Steering
- Wake-on-LAN (WoL)
- Hardware Accelerated 802.1ad VLAN (Q-in-Q Tunneling)
- VLAN Stripping in Linux Verbs
- Offloaded Traffic Sniffer
- Dump Configuration
- Local Loopback Disable
- Kernel Transport Layer Security (kTLS) Offloads
- IPsec Crypto Offload

Ethernet Interface

Port Type Management/VPI Cards Configuration

Ports of ConnectX-4 adapter cards and above can be individually configured to work as InfiniBand or Ethernet ports. By default, both VPI ports are initialized as InfiniBand ports. If you wish to change the port type, use the mlxconfig script after the driver is loaded.

For further information on how to set the port type, please refer to the MFT User Manual (www.mellanox.con --> Products --> Software --> InfiniBand/VPI Software --> MFT - Firmware Tools)

Counters

Counters are used to provide information about how well an operating system, an application, a service, or a driver is performing. The counter data help determine system bottlenecks and fine-tune the system and application performance. The operating system, network, and devices provide counter data that an application can consume to provide users with a graphical view of how well the system is performing.

The counter index is a Queue Pair (QP) attribute given in the QP context. Multiple QPs may be associated with the same counter set. If multiple QPs share the same counter, the counter value will represent the cumulative total.

RoCE Counters

- RoCE counters are available only through sysfs located under:
 - # /sys/class/infiniband/<device>/ports/*/counters/
 - # /sys/class/infiniband/<device>/ports/*/hw_counters/

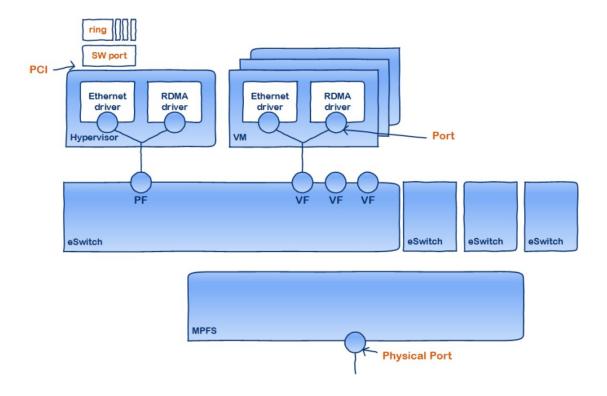
For mlx5 port and RoCE counters, refer to the <u>Understanding mlx5 Linux Counters</u> Community post.

SR-IOV Counters

Physical Function can also read Virtual Functions' port counters through sysfs located under # /sys/class/net/<interface name>/device/sriov/<index>/stats/

ethtool Counters

The ethtool counters are counted in different places, according to which they are divided into groups. Each counters group may also have different counter types.



For the full list of supported ethtool counters, refer to the <u>Understanding mlx5 ethtool Counters</u> Community post.

Persistent Naming

To avoid network interface renaming after boot or driver restart, set the desired constant interface name in the "/etc/udev/rules.d/70-persistent-net.rules" file.

• Example for Ethernet interfaces:

```
#PCI device 15b3:1019 (mlx5_core)
#NAME:="some name", := is used to make sure that device name will be persistent.
SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?", ATTR{address}=="00:02:c9:fa:c3:50", ATTR{dev_id}=="0x0",
ATTR(type)=="1", KERNEL=="eth", NAME:="eth1"
SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?", ATTR{address}=="00:02:c9:fa:c3:51", ATTR{dev_id}=="0x0",
ATTR{type}=="1", KERNEL=="eth", NAME:="eth2"
```

Example for IPoIB interfaces:

```
#1. ATTR{type}=="32" is IPoIB interfaces
#2. Find address value for each IB interfaces:
#cat /sys/class/net/ib0/address | head -1 | cut -d":" -f"13-"
#24:8a:07:03:00:a5:29:38
#Use above output to place into ATTR{address}.
#Address match must start with ?* and only reference the last 8 bytes of the address
SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?*", ATTR{type}=="32", ATTR{address}=="?
*24:8a:07:03:00:a5:29:38", NAME:="ib0"
SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?*", ATTR{type}=="32", ATTR{address}=="?
*24:8a:07:03:00:a5:29:38", NAME:="ib1"
```

Interrupt Request (IRQ) Naming

Once IRQs are allocated by the driver, they are named mlx5_comp<x>@pci:<pci_addr>. The IRQs corresponding to the channels in use are renamed to <interface>-<x>, while the rest maintain their default name.

The mlx5_core driver allocates all IRQs during loading time to support the maximum possible number of channels. Once the driver is up, no further IRQs are freed or allocated. Changing the number of working channels does not re-allocate or free the IRQs.

The following example demonstrates how reducing the number of channels affects the IRQs names.

```
$ ethtool -1 ens1
Channel parameters for ens1:
Pre-set maximums:
RX:
Current hardware settings:
Other:
Combined:
$ cat /proc/interrupts
98:
                                                                         7935
                                                                                       0
                                  mlx5_async@pci:0000:81:00.0
0 0 0
          0 IR-PCI-MSI-edge
         0 1
99:
                                                                          1
                                                                                       0
              IR-PCI-MSI-edge
                                  ens1-0
100:
                                                       0
                                                                  0
                                                                                       0
                                                                                                            0
     0
                                                                                                 0
                                  ens1-1
              IR-PCI-MSI-edge
101:
                                                                  0
                                  0 ens1-2 0 ens1-3 0 ens1-4
              IR-PCI-MSI-edge
          0 0
 102:
                                                     0
                                                                  0
                                                                                       0
                                                                                                 0
                                                                                                            0
              IR-PCI-MSI-edge
          0 0
 103:
                                                      0
                                                                  0
                                                                            1
                                                                                       0
                                                                                                 0
              IR-PCI-MSI-edge
                                           0
0
                                                                  0
                                                                            1
                                                                                       0
                                                                                                 0
                                                                                                            0
              IR-PCI-MSI-edge
                                  ens1-5
                                                                  0
                                                                                       0
                                                      0
                                                                                                 0
              IR-PCI-MSI-edge
                                  ens1-6
                                                      0
                                                                  0
                                                                            1
                                                                                       0
                                                                                                 0
                                                                                                            0
                                  ens1-7
              IR-PCI-MSI-edge
107:
          0
                                 0 ens1-8 0 ens1-9 0 ens1-10 0 ens1-11
                                                      0
                                                                  0
                                                                                       0
                                                                                                 0
              IR-PCI-MSI-edge
0
108:
                                                       0
                                                                 0
                                                                            1
                                                                                       0
                                                                                                 0
                                                                                                            0
              IR-PCI-MSI-edge
          0
        0 IR-PCI-MSI-eage
0 0
0 IR-PCI-MSI-edge
0 0 IR-PCI-MSI-edge
109:
                                                      0
                                                                0
                                                                            1
                                                                                       0
                                                                                                 0
                                                                                                            0
110:
                                                     0
                                                                0
                                                                           1
                                                                                     0
                                                                                                            0
$ ethtool -L ens1 combined 4
$ ethtool -l ens1
Channel parameters for ens1:
Current hardware settings:
RX:
TX:
Other:
Combined:
  cat /proc/interrupts 0 0
      mlx5_async@pci:0000:81:00.0
                                  mlx5_async@pci:0000:81:00.0
0 0 0 0
0 ens1-0
0 0 0 0 0
0 ens1-2
0 0 0 0 0 0
0 ens1-3
0 0 0 0 0
                                                                          1
99:
                                                                                       2
                                                                                                 0
                                                                                                            0
100:
                                                                            1
                                                                                       0
                                                                                                            0
0
101:
0
102:
                                                                                       0
                                                                                                 0
              IR-PCI-MSI-edge
          0
103:
          1 0
                                                                            1
                                                                                       0
                                                                                                 0
                                  mlx5_comp4@pci:0000:81:00.0
              IR-PCI-MSI-edge
                                  104:
          0
              IR-PCI-MSI-edge
                                  105:
                                                                                                 0
              IR-PCI-MSI-edge
106:
           0 IR-PCI-MSI-edge
           0
107:
                                                                                       0
                                                                                                 0
                                  mlx5_comp8@pci:0000:81:00.0
              IR-PCI-MSI-edge
          0
108:
                                                                            1
                                                                                       0
                                                                                                 0
                                  0 0 0 mlx5_comp9@pci:0000:81:00.0
              IR-PCI-MSI-edge
109:
                                                                 0
                                                                            1
                                                                                       0
                                                                                                 0
          1 IR-PCI-MSI-edge
0 0
0 IR-PCI-MSI-edge
                                  mlx5_comp10@pci:0000:81:00.0
110:
                                  0 0 0 0 mlx5_comp11@pci:0000:81:00.0
```

Quality of Service (QoS)

Quality of Service (QoS) is a mechanism of assigning a priority to a network flow (socket, rdma_cm connection) and manage its guarantees, limitations and its priority over other flows. This is accomplished by mapping the user's priority to a hardware TC (traffic class) through a 2/3 stage process. The TC is assigned with the QoS attributes and the different flows behave accordingly.

Mapping Traffic to Traffic Classes

Mapping traffic to TCs consists of several actions which are user controllable, some controlled by the application itself and others by the system/network administrators.

The following is the general mapping traffic to Traffic Classes flow:

- 1. The application sets the required Type of Service (ToS).
- 2. The ToS is translated into a Socket Priority (sk_prio).
- 3. The sk_prio is mapped to a User Priority (UP) by the system administrator (some applications set sk_prio directly).
- 4. The UP is mapped to TC by the network/system administrator.
- 5. TCs hold the actual QoS parameters

QoS can be applied on the following types of traffic. However, the general QoS flow may vary among them:

- Plain Ethernet Applications use regular inet sockets and the traffic passes via the kernel Ethernet driver
- RoCE Applications use the RDMA API to transmit using Queue Pairs (QPs)
- Raw Ethernet QP Application use VERBs API to transmit using a Raw Ethernet QP

Plain Ethernet Quality of Service Mapping

Applications use regular inet sockets and the traffic passes via the kernel Ethernet driver. The following is the Plain Ethernet QoS mapping flow:

- 1. The application sets the ToS of the socket using setsockopt (IP_TOS, value).
- 2. ToS is translated into the sk_prio using a fixed translation:

```
TOS 0 <=> sk_prio 0
TOS 8 <=> sk_prio 2
TOS 24 <=> sk_prio 4
TOS 16 <=> sk_prio 6
```

- 3. The Socket Priority is mapped to the UP:
 - If the underlying device is a VLAN device, egress_map is used controlled by the vconfig command. This is per VLAN mapping.
 - If the underlying device is not a VLAN device, the mapping is done in the driver.
 - 4. The UP is mapped to the TC as configured by the mlnx_qos tool or by the lldpad daemon if DCBX is used.
- ⚠ Socket applications can use setsockopt (SK_PRIO, value) to directly set the sk_prio of the socket. In this case, the ToS to sk_prio fixed mapping is not needed. This allows the application and the administrator to utilize more than the 4 values possible via ToS.
- ⚠ In the case of a VLAN interface, the UP obtained according to the above mapping is also used in the VLAN tag of the traffic.

RoCE Quality of Service Mapping

Applications use RDMA-CM API to create and use QPs. The following is the RoCE QoS mapping flow:

- 1. The application sets the ToS of the QP using the rdma_set_option option(RDMA OPTION ID TOS, value).
- 2. ToS is translated into the Socket Priority (sk_prio) using a fixed translation:

```
TOS 0 <=> sk_prio 0
TOS 8 <=> sk_prio 2
TOS 24 <=> sk_prio 4
TOS 16 <=> sk_prio 6
```

- 3. The Socket Priority is mapped to the User Priority (UP) using the tc command.
 - In the case of a VLAN device where the parent real device is used for the purpose of this mapping
 - If the underlying device is a VLAN device, and the parent real device was not used for the mapping, the VLAN device's egress map is used
 - 4. UP is mapped to the TC as configured by the mlnx_qos tool or by the lldpad daemon if DCBX is used.

With RoCE, there can only be 4 predefined ToS values for the purpose of QoS mapping.

Map Priorities with set_egress_map

For RoCE old kernels that do not support set_egress_map, use the tc_wrap script to map between sk_prio and UP. Use tc_wrap with option -u. For example:

```
tc_wrap -i <ethX> -u <skprio2up mapping>
```

Quality of Service Properties

The different QoS properties that can be assigned to a TC are:

- Strict Priority
- Enhanced Transmission Selection (ETS)
- Rate Limit
- Trust State
- Receive Buffer
- DCBX Control Mode

Strict Priority

When setting a TC's transmission algorithm to be 'strict', then this TC has absolute (strict) priority over other TC strict priorities coming before it (as determined by the TC number: TC 7 is the highest priority, TC 0 is lowest). It also has an absolute priority over nonstrict TCs (ETS).

This property needs to be used with care, as it may easily cause starvation of other TCs. A higher strict priority TC is always given the first chance to transmit. Only if the highest strict priority TC has nothing more to transmit, will the next highest TC be considered. Nonstrict priority TCs will be considered last to transmit.

This property is extremely useful for low latency low bandwidth traffic that needs to get immediate service when it exists, but is not of high volume to starve other transmitters in the system.

Enhanced Transmission Selection (ETS)

Enhanced Transmission Selection standard (ETS) exploits the time periods in which the offered load of a particular Traffic Class (TC) is less than its minimum allocated bandwidth by allowing the difference to be available to other traffic classes.

After servicing the strict priority TCs, the amount of bandwidth (BW) left on the wire may be split among other TCs according to a minimal guarantee policy.

If, for instance, TC0 is set to 80% guarantee and TC1 to 20% (the TCs sum must be 100), then the BW left after servicing all strict priority TCs will be split according to this ratio.

Since this is a minimum guarantee, there is no maximum enforcement. This means, in the same example, that if TC1 did not use its share of 20%, the reminder will be used by TC0.

ETS is configured using the mlnx_qos tool (mlnx_qos) which allows you to:

- Assign a transmission algorithm to each TC (strict or ETS)
- · Set minimal BW guarantee to ETS TCs Usage:

mlnx_qos -i \[options\]

Rate Limit

Rate limit defines a maximum bandwidth allowed for a TC. Please note that 10% deviation from the requested values is considered acceptable.

Trust State

Trust state enables prioritizing sent/received packets based on packet fields.

The default trust state is PCP. Ethernet packets are prioritized based on the value of the field (PCP/ DSCP).

For further information on how to configure Trust mode, please refer to HowTo Configure Trust State on Mellanox Adapters Community post.



⚠ Setting the Trust State mode shall be done before enabling SR-IOV in order to propagate the Trust State to the VFs.

Receive Buffer

By default, the receive buffer configuration is controlled automatically. Users can override the receive buffer size and receive buffer's xon and xoff thresholds using mlnx_gos tool. For further information, please refer to HowTo Tune the Receive buffers on Mellanox Adapters Community post.

DCBX Control Mode

DCBX settings, such as "ETS" and "strict priority" can be controlled by firmware or software. When DCBX is controlled by firmware, changes of QoS settings cannot be done by the software. The DCBX control mode is configured using the mlnx gos -d os/fw command.

For further information on how to configure the DCBX control mode, please refer to <u>mlnx qos</u> Community post.

Quality of Service Tools

mlnx_qos

mlnx_qos is a centralized tool used to configure QoS features of the local host. It communicates directly with the driver thus does not require setting up a DCBX daemon on the system. The mlnx_qos tool enables the administrator of the system to:

- Inspect the current QoS mappings and configuration
 The tool will also display maps configured by TC and vconfig set_egress_map tools, in order to give a centralized view of all QoS mappings.
- Set UP to TC mapping
- Assign a transmission algorithm to each TC (strict or ETS)
- Set minimal BW guarantee to ETS TCs
- Set rate limit to TCs
- Set DCBX control mode
- Set cable length
- Set trust state



⚠ For an unlimited ratelimit, set the ratelimit to 0.

Usage

mlnx_qos -i <interface> \[options\]

Options

version	Show the program's version number and exit
-h,help	Show this help message and exit
-f LIST,pfc=LIST	Set priority flow control for each priority. LIST is a comma separated value for each priority starting from 0 to 7. Example: 0,0,0,0,1,1,1,1 enable PFC on TC4-7
-p LIST,prio_tc=LIST	Maps UPs to TCs. LIST is 8 comma-separated TC numbers. Example: 0,0,0,0,1,1,1,1 maps UPs 0-3 to TC0, and UPs 4-7 to TC1
-s LIST,tsa=LIST	Transmission algorithm for each TC. LIST is comma separated algorithm names for each TC. Possible algorithms: strict, ets and vendor. Example: vendor, strict, ets, ets, ets, ets, ets sets TC0 to vendor, TC1 to strict, TC2-7 to ets
-t LIST,tcbw=LIST	Set the minimally guaranteed %BW for ETS TCs. LIST is commaseparated percents for each TC. Values set to TCs that are not configured to ETS algorithm are ignored but must be present. Example: if TC0,TC2 are set to ETS, then 10,0,90,0,0,0,0,0will set TC0 to 10% and TC2 to 90%. Percents must sum to 100

-r LIST,ratelimit=LIST	Rate limit for TCs (in Gbps). LIST is a comma-separated Gbps limit for each TC. Example: 1,8,8 will limit TC0 to 1Gbps, and TC1,TC2 to 8 Gbps each	
-d DCBX,dcbx=DCBX	Set dcbx mode to firmware controlled(fw) or OS controlled(os). Note, when in OS mode, mlnx_qos should not be used in parallel with other dcbx tools, such as lldptool	
trust=TRUST	set priority trust state to pcp or dscp	
dscp2prio=DSCP2PRIO	Set/del a (dscp,prio) mapping. Example 'set,30,2' maps dscp 30 to priority 2. 'del,30,2' resets the dscp 30 mapping back to the default setting priority 0	
cable_len=CABLE_LEN	Set cable_len for buffer's xoff and xon thresholds	
-i INTF,interface=INTF	Interface name	
-a	Show all interface's TCs	

Get Current Configuration

```
ofed_scripts/utils/mlnx_gos -i ens1f0
DCBX mode: OS controlled
Priority trust state: dscp
dscp2prio mapping:
    prio:0 dscp:07,06,05,04,03,02,01,00,
        prio:1 dscp:15,14,13,12,11,10,09,08,
        prio:2 dscp:23,22,21,20,19,18,17,16,
        prio:3 dscp:33,22,21,20,19,18,17,16,
        prio:4 dscp:39,38,37,36,35,34,33,32,
        prio:4 dscp:55,54,53,52,51,50,49,48,
        prio:6 dscp:55,54,53,52,51,50,49,48,
        prio:7 dscp:63,62,61,60,59,58,57,56,
Cable len: 7
PFC configuration:
    priority 0 1 2 3 4 5 6 7
        enabled 0 0 0 0 0 0 0 0
tc: 0 ratelimit: unlimited, tsa: vendor
        priority: 1
tc: 1 ratelimit: unlimited, tsa: vendor
        priority: 2
tc: 3 ratelimit: unlimited, tsa: vendor
        priority: 3
tc: 4 ratelimit: unlimited, tsa: vendor
        priority: 4
tc: 5 ratelimit: unlimited, tsa: vendor
        priority: 5
tc: 6 ratelimit: unlimited, tsa: vendor
        priority: 5
tc: 6 ratelimit: unlimited, tsa: vendor
        priority: 5
tc: 7 ratelimit: unlimited, tsa: vendor
        priority: 6
tc: 7 ratelimit: unlimited, tsa: vendor
```

Set ratelimit. 3Gbps for tc0 4Gbps for tc1 and 2Gbps for tc2

ConfigureQoS. Map UP0,7 to tc0,1,2,3 to tc1 and 4,5,6 to tc2. Set tc0,tc1 as ets and tc2 as strict. Divide ets 30% for tc0 and 70% for tc1

tc and tc_wrap.py

The tc tool is used to create 8 Traffic Classes (TCs).

The tool will either use the sysfs (/sys/class/net/<ethX>/qos/tc_num) or the tc tool to create the TCs.

Usage

```
tc_wrap.py -i <interface> \[options\]
```

Options

version	show program's version number and exit
-h,help	show this help message and exit

-u SKPRIO_UP,skprio_up=SKPRIO_UP	maps sk_prio to priority for RoCE. LIST is <=16 comma separated priority. index of element is sk_prio	
-i INTF,interface=INTF	Interface name	

Example

Run:

```
tc_wrap.py -i enp139s0
```

Output:

Additional Tools

tc tool compiled with the sch_mqprio module is required to support kernel v2.6.32 or higher. This is a part of iproute2 package v2.6.32-19 or higher. Otherwise, an alternative custom sysfs interface is available.

- mlnx_qos tool (package: ofed-scripts) requires python version 2.5 < = X
- tc_wrap.py (package: ofed-scripts) requires python version 2.5 < = X

Packet Pacing

ConnectX-4 and above devices allow packet pacing (traffic shaping) per flow. This capability is achieved by mapping a flow to a dedicated send queue and setting a rate limit on that Send queue. Note the following:

- Up to 512 send queues are supported
- 16 different rates are supported
- The rates can vary from 1 Mbps to line rate in 1 Mbps resolution
- Multiple queues can be mapped to the same rate (each queue is paced independently)
- It is possible to configure rate limit per CPU and per flow in parallel

System Requirements

- MLNX OFED, v3.3 or higher
- Linux kernel v4.1 or higher
- ConnectX-4 or ConnectX-4 Lx adapter cards with an official firmware version

Packet Pacing Configuration

A

This configuration is non-persistent and does not survive driver restart.

1. Firmware Activation:

To activate Packet Pacing in the firmware:

First, make sure Mellanox Firmware Tools service (mst) is started:

```
# mst start
```

Then run:

```
#echo "MLNX_RAW_TLV_FILE" > /tmp/mlxconfig_raw.txt
#echo "0x00000004 0x0000010c 0x000000000 0x00000001" >> /tmp/mlxconfig_raw.txt
#yes | mlxconfig -d <mst_dev> -f /tmp/mlxconfig_raw.txt set_raw > /dev/mull
#reboot /mlxfwreset
```

To deactivate Packet Pacing in the firmware, run:

```
#echo "MLNX_RAW_TLV_FILE" > /tmp/mlxconfig_raw.txt
#echo "0x00000004 0x0000010c 0x000000000 0x000000000" >> /tmp/mlxconfig_raw.txt
#yes | mlxconfig -d <mst_dev >-f /tmp/mlxconfig_raw.txt set_raw > /dev/null
#reboot /mlxfwreset
```

2. Driver Activation:

There are two operation modes for Packet Pacing:

a. Rate limit per CPU core:

When XPS is enabled, traffic from a CPU core will be sent using the corresponding send queue. By limiting the rate on that queue, the transmit rate on that CPU core will be limited. For example:

```
echo 300 > /sys/class/net/ens2f1/queues/tx-0/tx_maxrate
```

In this case, the rate on Core 0 (tx-0) is limited to 300Mbit/sec.

- b. Rate limit per flow:
 - i. The driver allows opening up to 512 additional send queues using the following command:

```
ethtool -L ens2f1 other 1200
```

In this case, 1200 additional queues are opened

ii. Create flow mapping.

Users can map a certain destination IP and/or destination layer 4 Port to a specific send queue. The match precedence is as follows:

- IP + L4 Port
- IP only
- · L4 Port only
- No match (the flow would be mapped to default queues) To create flow mapping:

Configure the destination IP. Write the IP address in hexadecimal representation to the relevant sysfs entry. For example, to map IP address 192.168.1.1 (0xc0a80101) to send queue 310, run the following command:

echo 0xc0a80101 > /sys/class/net/ens2f1/queues/tx-310/flow_map/dst_ip

To map Destination L4 3333 port (either TCP or UDP) to the same queue, run:

echo 3333 > /sys/class/net/ens2f1/queues/tx-310/flow_map/dst_port

From this point on, all traffic destined to the given IP address and L4 port will be sent using send queue 310. All other traffic will be sent using the original send queue.

iii. Limit the rate of this flow using the following command:

echo 100 > /sys/class/net/ens2f1/queues/tx-310/tx_maxrate

Note: Each queue supports only a single IP+Port combination.

Ethtool

Ethtool is a standard Linux utility for controlling network drivers and hardware, particularly for wired Ethernet devices. It can be used to:

- · Get identification and diagnostic information
- · Get extended device statistics
- · Control speed, duplex, auto-negotiation and flow control for Ethernet devices
- Control checksum offload and other hardware offload features
- Control DMA ring sizes and interrupt moderation
- Flash device firmware using a .mfa2 image

Ethtool Supported Options

Options	Description
ethtoolset-priv-flags eth <x> <priv flag=""> <on off=""></on></priv></x>	Enables/disables driver feature matching the given private flag.
ethtoolshow-priv-flags eth <x></x>	Shows driver private flags and their states (ON/OFF).
ethtool -a eth <x></x>	Queries the pause frame settings.
ethtool -A eth <x> [rx on off] [tx on off]</x>	Sets the pause frame settings.
ethtool -c eth <x></x>	Queries interrupt coalescing settings.
ethtool -C eth <x> [pkt-rate-low N] [pkt-rate-high N] [rx-usecs-low N] [rx-usecs-high N]</x>	Sets the values for packet rate limits and for moderation time high and low values.
ethtool -C eth <x> [rx-usecs N] [rx-frames N]</x>	Sets the interrupt coalescing setting.
	rx-frames will be enforced immediately, rx-usecs will be enforced only when adaptive moderation is disabled.
	Note : usec settings correspond to the time to wait after the *last* packet is sent/received before triggering an interrupt.

Options	Description
ethtool -C eth <x> adaptive-rx on off</x>	Enables/disables adaptive interrupt moderation.
	By default, the driver uses adaptive interrupt moderation for the receive path, which adjusts the moderation time to the traffic pattern.
ethtool -C eth <x> adaptive-tx on off</x>	Note : Supported by mlx5e for ConnectX-4 and above adapter cards.
	Enables/disables adaptive interrupt moderation.
	By default, the driver uses adaptive interrupt moderation for the transmit path, which adjusts the moderation parameters (time/frames) to the traffic pattern.
ethtool -g eth <x></x>	Queries the ring size values.
ethtool -G eth <x> [rx <n>] [tx <n>]</n></n></x>	Modifies the ring size.
ethtool -i eth <x></x>	Checks driver and device information.
	For example:
	driver: mlx5_core
	version: 5.1-0.4.0
	firmware-version: 4.6.4046 (MT_QEMU000000)
	expansion-rom-version:
	bus-info: 0000:07:00.0
	supports-statistics: yes
	supports-test: yes
	supports-eeprom-access: no
	supports-register-dump: no
	supports-priv-flags: yes
ethtool -k eth <x></x>	Queries the stateless offload status.

Options	Description		
ethtool -K eth <x> [rx on off] [tx on off] [sg on </x>	Sets the stateless offload status.		
off] [tso on off] [lro on off] [gro on off] [gso on off] [rxvlan on off] [txvlan on off] [ntuple on/off] [rxhash on/off] [rx-all on/off] [rx-fcs on/off]	TCP Segmentation Offload (TSO), Generic Segmentation Offload (GSO): increase outbound throughput by reducing CPU overhead. It works by queuing up large buffers and letting the network interface card split them into separate packets.		
	Large Receive Offload (LRO): increases inbound throughput of high-bandwidth network connections by reducing CPU overhead. It works by aggregating multiple incoming packets from a single stream into a larger buffer before they are passed higher up the networking stack, thus reducing the number of packets that have to be processed. LRO is available in kernel versions < 3.1 for untagged traffic.		
	Hardware VLAN insertion Offload (txvlan): When enabled, the sent VLAN tag will be inserted into the packet by the hardware.		
	Note : LRO will be done whenever possible. Otherwise GRO will be done. Generic Receive Offload (GRO) is available throughout all kernels.		
	Hardware VLAN Striping Offload (rxvlan): When enabled received VLAN traffic will be stripped from the VLAN tag by the hardware.		
	RX FCS (rx-fcs): Keeps FCS field in the received packets. Sets the stateless offload status.		
	RX FCS validation (rx-all): Ignores FCS validation on the received packets.		
ethtool -l eth <x></x>	Shows the number of channels.		
ethtool -L eth <x> [rx <n>] [tx <n>]</n></n></x>	Sets the number of channels.		
	 Notes: This also resets the RSS table to its default distribution, which is uniform across the cores on the NUMA (non-uniform memory access) node that is closer to the NIC. For ConnectX®-4 cards, use ethtool -L eth<x></x> 		
	combined <n> to set both RX and TX channels.</n>		
ethtool -m dump-module-eeprom eth <x> [raw on off] [hex on off] [offset N] [length N]</x>	Queries/decodes the cable module eeprom information.		
ethtool -p identify DEVNAME	Enables visual identification of the port by LED blinking [TIME-IN-SECONDS].		
ethtool -p identify eth <x> <led duration=""></led></x>	Allows users to identify interface's physical port by turning the ports LED on for a number of seconds.		
	Note : The limit for the LED duration is 65535 seconds.		
ethtool -S eth <x></x>	Obtains additional device statistics.		

Options	Options Description			
ethtool -s eth <x> advertise <n> autoneg on</n></x>	Changes the advertised link modes to requested link modes <n></n>			
	To check the link modes' hex values, run <man and="" check="" e="" eth<x="" ether="" link="" modes,="" run="" supported="" the="" to=""></man>			
	For advertising new link modes the entire bitmap as follows:	r advertising new link modes, make sure to configure e entire bitmap as follows:		
	200GAUI-4 / 200GBASE-CR4/ KR4	0x7c000000000000000		
	100GAUI-2 / 100GBASE-CR2 / KR2	0x3E00000000000000		
	CAUI-4 / 100GBASE-CR4 / KR4	0xF000000000		
	50GAUI-1 / LAUI-1/ 50GBASE-CR / KR	0x1F0000000000000		
	50GAUI-2 / LAUI-2/ 50GBASE-CR2/KR2	0x10C00000000		
	XLAUI-4/XLPPI-4 // 40G	0x7800000		
	25GAUI-1/ 25GBASE-CR / KR	0x380000000		
	XFI / XAUI-1 // 10G	0x7C0000181000		
	5GBASE-R	0x100000000000		
	2.5GBASE-X / 2.5GMII	0x82000000000		
	1000BASE-X / SGMII	0x20000020020		
	 Notes: Both previous and new link modes configurations are supported, however, they must be run separately. Any link mode configuration on Kernels below v5.1 and ConnectX-6 HCAs will result in the advertisement of the full capabilities. <autoneg on=""> only sends a hint to the driver that the user wants to modify advertised link modes and not speed.</autoneg> 			
ethtool -s eth <x> msglvl [N]</x>	Changes the current driver message level.			
ethtool -s eth <x> speed <speed> autoneg off</speed></x>	Changes the link speed to requested <speed>. To check the supported speeds, run ethtool eth<x>. Note: <autoneg off=""> does not set autoneg OFF, it only hints the driver to set a specific speed.</autoneg></x></speed>			
ethtool -t eth <x></x>	Performs a self-diagnostics test.			
ethtool -T eth <x></x>	Shows time stamping capabilities			
ethtool -x eth <x></x>	Retrieves the receive flow hash indirection table.			

Options	Description
ethtool -X eth <x> equal a b c</x>	Sets the receive flow hash indirection table.
	Note: The RSS table configuration is reset whenever the number of channels is modified (using ethtool -L command).
ethtoolshow-fec eth <x></x>	Queries current Forward Error Correction (FEC) encoding in case FEC is supported.
	Note : An output of "baser" implies Firecode encoding.
ethtoolset-fec eth <x> encoding auto off rs </x>	Configures Forward Error Correction (FEC).
baser	Note : 'baser' encoding applies to the Firecode encoding, and 'auto' regards the HCA's default.
ethtool -f flash <devname> FILE [N]</devname>	Flash firmware image on the device using the specified .mfa2 file (FILE). By default, the command flashes all the regions on the device unless a region number (N) is specified.

Checksum Offload

MLNX_OFED supports the following Receive IP/L4 Checksum Offload modes:

- CHECKSUM_UNNECESSARY: By setting this mode the driver indicates to the Linux Networking Stack that the hardware successfully validated the IP and L4 checksum so the Linux Networking Stack does not need to deal with IP/L4 Checksum validation. Checksum Unnecessary is passed to the OS when all of the following are true:
 - Ethtool -k <DEV> shows rx-checksumming: on
 - Received TCP/UDP packet and both IP checksum and L4 protocol checksum are correct.
- CHECKSUM_COMPLETE: When the checksum validation cannot be done or fails, the driver still
 reports to the OS the calculated by hardware checksum value. This allows accelerating
 checksum validation in Linux Networking Stack, since it does not have to calculate the whole
 checksum including payload by itself.

Checksum Complete is passed to OS when both of the following is true:

- Ethtool -k <DEV> shows rx-checksumming: on
- Received IPv4/IPv6 non-TCP/UDP packet
- CHECKSUM_NONE: By setting this mode the driver indicates to the Linux Networking Stack
 that the hardware failed to validate the IP or L4 checksum so the Linux Networking Stack
 must calculate and validate the IP/L4 Checksum.

Checksum None is passed to OS for all other cases.

Ignore Frame Check Sequence (FCS) Errors

Upon receiving packets, the packets go through a checksum validation process for the FCS field. If the validation fails, the received packets are dropped.

When FCS is enabled (disabled by default), the device does not validate the FCS field even if the field is invalid.

It is not recommended to enable FCS.

For further information on how to enable/disable FCS, please refer to ethtool option rx-fcs on/off.

RDMA over Converged Ethernet (RoCE)

Remote Direct Memory Access (RDMA) is the remote memory management capability that allows server-to-server data movement directly between application memory without any CPU involvement. RDMA over Converged Ethernet (RoCE) is a mechanism to provide this efficient data transfer with very low latencies on lossless Ethernet networks. With advances in data center convergence over reliable Ethernet, ConnectX® Ethernet adapter cards family with RoCE uses the proven and efficient RDMA transport to provide the platform for deploying RDMA technology in mainstream data center application at 10GigE and 40GigE link-speed. ConnectX® Ethernet adapter cards family with its hardware offload support takes advantage of this efficient RDMA transport (InfiniBand) services over Ethernet to deliver ultra-low latency for performance-critical and transaction-intensive applications such as financial, database, storage, and content delivery networks.

When working with RDMA applications over Ethernet link layer the following points should be noted:

- The presence of a Subnet Manager (SM) is not required in the fabric. Thus, operations that require communication with the SM are managed in a different way in RoCE. This does not affect the API but only the actions such as joining the multicast group, that need to be taken when using the API
- Since LID is a layer 2 attribute of the InfiniBand protocol stack, it is not set for a port and is displayed as zero when querying the port
- With RoCE, the alternate path is not set for RC QP. Therefore, APM (another type of High Availability and part of the InfiniBand protocol) is not supported
- Since the SM is not present, querying a path is impossible. Therefore, the path record structure must be filled with relevant values before establishing a connection. Hence, it is recommended working with RDMA-CM to establish a connection as it takes care of filling the path record structure
- VLAN tagged Ethernet frames carry a 3-bit priority field. The value of this field is derived from the IB SL field by taking the 3 least significant bits of the SL field
- RoCE traffic is not shown in the associated Ethernet device's counters since it is offloaded by
 the hardware and does not go through Ethernet network driver. RoCE traffic is counted in the
 same place where InfiniBand traffic is counted; /sys/class/infiniband/<device>/ports/<port
 number>/counters/

RoCE Modes

RoCE encapsulates IB transport in one of the following Ethernet packets:

- RoCEv1 dedicated ether type (0x8915)
- RoCEv2 UDP and dedicated UDP port (4791)

RoCEv1 and RoCEv2 Protocol Stack

? Unknown Attachment

RoCEv1

RoCE v1 protocol is defined as RDMA over Ethernet header (as shown in the figure above). It uses ethertype 0x8915 and can be used with or without the VLAN tag. The regular Ethernet MTU applies on the RoCE frame.

RoCEv2

A straightforward extension of the RoCE protocol enables traffic to operate in IP layer 3 environments. This capability is obtained via a simple modification of the RoCE packet format. Instead of the GRH used in RoCE, IP routable RoCE packets carry an IP header which allows traversal of IP L3 Routers and a UDP header (RoCEv2 only) that serves as a stateless encapsulation layer for the RDMA Transport Protocol Packets over IP.

The proposed RoCEv2 packets use a well-known UDP destination port value that unequivocally distinguishes the datagram. Similar to other protocols that use UDP encapsulation, the UDP source port field is used to carry an opaque flow-identifier that allows network devices to implement packet forwarding optimizations (e.g. ECMP) while staying agnostic to the specifics of the protocol header format.

Furthermore, since this change exclusively affects the packet format on the wire, and due to the fact that with RDMA semantics packets are generated and consumed below the AP, applications can seamlessly operate over any form of RDMA service, in a completely transparent way.



Both RoCEv1 and RoCEv2 are supported by default; the driver associates all GID indexes to RoCEv1 and RoCEv2, thus, a single entry for each RoCE version.

For further information, please refer to HowTo Configure RoCEv2 Community post.

GID Table Population

GID table entries are created whenever an IP address is configured on one of the Ethernet devices of the NIC's ports. Each entry in the GID table for RoCE ports has the following fields:

- GID value
- GID type
- · Network device

The GID table is occupied with two GIDs, both with the same GID value but with different types. The network device in an entry is the Ethernet device with the IP address that GID is associated with. The GID format can be of 2 types; IPv4 and IPv6. IPv4 GID is an IPv4-mapped IPv6 address, while IPv6 GID is the IPv6 address itself. Layer 3 header for packets associated with IPv4 GIDs will be IPv4 (for RoCEv2) and IPv6/GRH for packets associated with IPv6 GIDs and IPv4 GIDs for RoCEv1.

GID Table in sysfs

GID table is exposed to userspace via sysfs

• GID values can be read from:

/sys/class/infiniband/{device}/ports/{port}/gids/{index}

GID type can be read from:

/sys/class/infiniband/{device}/ports/{port}/gid_attrs/types/{index}

• GID net_device can be read from:

/sys/class/infiniband/{device}/ports/{port}/gid_attrs/ndevs/{index}

Setting the RoCE Mode for a Queue Pair (QP)

Setting RoCE mode for devices that support two RoCE modes is different for RC/UC QPs (connected QP types) and UD QP.

To modify an RC/UC QP (connected QP) from INIT to RTR, an Address Vector (AV) must be given. The AV, among other attributes, should specify the index of the port's GID table for the source GID of the QP. The GID type in that index will be used to set the RoCE type of the QP.

Setting RoCE Mode of RDMA_CM Applications

RDMA_CM interface requires only the active side of the peer to pass the IP address of the passive side. The RDMA_CM decides upon the source GID to be used and obtains it from the GID table. Since more than one instance of the GID value is possible, the lookup should be also according to the GID type. The type to use for the lookup is defined as a global value of the RDMA_CM module. Changing the value of the GID type for the GID table lookups is done using the cma_roce_mode script.

To print the current RoCE mod	de for a device port:	J	·
cma_roce_mode -d <dev> -p <port></port></dev>			
To set the RoCE mode for a de	evice port:		
cma_roce_mode -d <dev> -p <port> -m <</port></dev>	1 2>		

GID Table Example

The following is an example of the GID table.

DEV	PORT	INDEX	GID	IPv4	Туре	Netdev
mlx5_0	1	0	fe80:0000:0000:0000:ba59:9fff:fe1a:e3e a		v1	p4p1
mlx5_0	1	1	fe80:0000:0000:0000:ba59:9fff:fe1a:e3e a		v2	p4p1
mlx5_0	1	2	0000:0000:0000:0000:0000:ffff:0a0a:0a0	10.10.10.1	v1	p4p1
mlx5_0	1	3	0000:0000:0000:0000:0000:ffff:0a0a:0a0	10.10.10.1	v2	p4p1
mlx5_1	1	0	fe80:0000:0000:0000:ba59:9fff:fe1a:e3e b		v1	p4p2
mlx5_1	1	1	fe80:0000:0000:0000:ba59:9fff:fe1a:e3e b		v2	p4p2

where:

- Entries on port 1 index 0/1 are the default GIDs, one for each supported RoCE type
- Entries on port 1 index 2/3 belong to IP address 192.168.1.70 on eth1
- Entries on port 1 index 4/5 belong to IP address 193.168.1.70 on eth1.100

- Packets from a QP that is associated with these GID indexes will have a VLAN header (VID=100)
- Entries on port 1 index 6/7 are IPv6 GID. Packets from a QP that is associated with these GID indexes will have an IPv6 header

RoCE Lossless Ethernet Configuration

In order to function reliably, RoCE requires a form of flow control. While it is possible to use global flow control, this is normally undesirable, for performance reasons.

The normal and optimal way to use RoCE is to use Priority Flow Control (PFC). To use PFC, it must be enabled on all endpoints and switches in the flow path.

For further information, please refer to RoCE Over L2 Network Enabled with PFC User Guide:

http://www.mellanox.com/related-docs/prod_software/ RoCE_with_Priority_Flow_Control_Application_Guide.pdf

Configuring SwitchX® Based Switch System

? Unknown Attachment

To enable RoCE, the SwitchX should be configured as follows:

- · Ports facing the host should be configured as access ports, and either use global pause or Port Control Protocol (PCP) for priority flow control
- Ports facing the network should be configured as trunk ports, and use Port Control Protocol (PCP) for priority flow control
- For further information on how to configure SwitchX, please refer to SwitchX User Manual

Installing and Loading the Driver

? Unknown Attachment

To install and load the driver:

1. Install MLNX_OFED (See Installation section for further details). RoCE is installed as part of mlx5 and other modules upon driver's installation.

⚠ The list of the modules that will be loaded automatically upon boot can be found in the configuration file /etc/infiniband/openib.conf.

2. Query for the device's information. Example:

ofed_info -s MLNX_OFED_LINUX-5.0-2.1.8.0:

3. Display the existing MLNX_OFED version.

Output Notes:

The port's state is: Ethernet is in PORT_ACTIVE state	The port state can also be obtained by running the following command: # cat /sys/class/infiniband/mlx5_0/ports/ 1/state 4: ACTIVE
link_layer parameter shows that port 1 is Ethernet	The link_layer can also be obtained by running the following command: # cat /sys/class/infiniband/mlx5_0/ ports/1/link_layer Ethernet
The fw_ver parameter shows that the firmware version is 16.28.0578.	The firmware version can also be obtained by running the following command: # cat /sys/class/infiniband/mlx5_0/ fw_ver 16.28.0578

Associating InfiniBand Ports to Ethernet Ports

The mlx5_ib driver holds a reference to the net device for getting notifications about the state of the port, as well as using the mlx5_core driver to resolve IP addresses to MAC that are required for address vector creation. However, RoCE traffic does not go through the mlx5_core driver; it is completely offloaded by the hardware.

```
# ibdev2netdev
mlx5_0 port 1 <===> eth2
#
```

Configuring an IP Address to the netdev Interface

? Unknown Attachment To configure an IP address to the netdev interface:

1. Configure an IP address to the netdev interface on both sides of the link.

```
# ifconfig eth2 20.4.3.220
# ifconfig eth2
eth2    Link encap:Ethernet HWaddr 00:02:C9:08:E8:11
    inet addr:20.4.3.220 Bcast:20.255.255.255 Mask:255.0.0.0
UP BROADCAST MULTICAST MUITICAST MURTICAST M
```

2. Make sure that ping is working.

```
ping 20.4.3.219
PING 20.4.3.219 (20.4.3.219) 56(84) bytes of data.
64 bytes from 20.4.3.219: icmp_seq=1 ttl=64 time=0.873 ms
64 bytes from 20.4.3.219: icmp_seq=2 ttl=64 time=0.198 ms
64 bytes from 20.4.3.219: icmp_seq=3 ttl=64 time=0.167 ms
20.4.3.219 ping statistics -
3 packets transmitted, 3 received, 0% packet loss, time 2000ms rtt min/avg/max/mdev = 0.167/0.412/0.873/
0.326 ms
```

Adding VLANs

? Unknown Attachment

To add VLANs:

1. Make sure that the 8021.q module is loaded.

```
modprobe 8021q
```

2. Add VLAN.

```
# vconfig add eth2 7
Added VLAN with VID == 7 to IF -:eth2:-
#
```

3. Configure an IP address.

```
ifconfig eth2.7 7.4.3.220
```

Defining Ethernet Priority (PCP in 802.1q Headers)

1. Define Ethernet priority on the server.

```
# ibv_rc_pingpong -g 1 -i 2 -1 4
local address: LID 0x0000, QPN 0x1c004f, PSN 0x9daf6c, GID fe80::202:c900:708:e799
remote address: LID 0x0000, QPN 0x1c004f, PSN 0xb0a49b, GID fe80::202:c900:708:e811
8192000 bytes in 0.01 seconds = 4840.89 Mbit/sec
1000 iters in 0.01 seconds = 13.54 usec/iter
```

2. Define Ethernet priority on the client.

```
# ibv_rc_pingpong -g 1 -i 2 -1 4 sw419
local address: LID 0x0000, OPN 0x1c004f, PSN 0xb0a49b, GID fe80::202:c900:708:e811
remote address: LID 0x0000, QPN 0x1c004f, PSN 0x9daf6c, GID fe80::202:c900:708:e799
8192000 bytes in 0.01 seconds = 4855.96 Mbit/sec
1000 iters in 0.01 seconds = 13.50 usec/iter
```

Using rdma_cm Tests

1. Use rdma_cm test on the server.

```
# ucmatose
cmatose: starting server
initiating data transfers
completing sends
receiving data transfers
data transfers complete
cmatose: disconnecting
disconnected
test complete
return status 0
#
```

2. Use rdma_cm test on the client.

```
# ucmatose -s 20.4.3.219
cmatose: starting client
cmatose: connecting
receiving data transfers
sending replies
data transfers complete
test complete
return status 0
#
```

This server-client run is without PCP or VLAN because the IP address used does not belong to a VLAN interface. If you specify a VLAN IP address, then the traffic should go over VLAN.

Type Of Service (ToS)

Overview

The TOS field for rdma_cm sockets can be set using the rdma_set_option() API, just as it is set for regular sockets. If a TOS is not set, the default value (0) is used. Within the rdma_cm kernel driver, the TOS field is converted into an SL field. The conversion formula is as follows:

• SL = TOS >> 5 (e.g., take the 3 most significant bits of the TOS field)

In the hardware driver, the SL field is converted into PCP by the following formula:

• PCP = SL & 7 (take the 3 least significant bits of the TOS field)

⚠ SL affects the PCP only when the traffic goes over tagged VLAN frames.

DSCP

A new entry has been added to the RDMA-CM configfs that allows users to select default TOS for RDMA-CM QPs. This is useful for users that want to control the TOS field without changing their code. Other applications that set the TOS explicitly using the rdma_set_option API will continue to work as expected to override the configfs value.

For further information about DSCP marking, refer to <u>HowTo Set Egress ToS/DSCP on RDMA- CM QPs</u> Community post.

RoCE LAG

RoCE LAG is a feature meant for mimicking Ethernet bonding for IB devices and is available for dual port cards only.

This feature is supported on RHEL and SLES systems.

RoCE LAG mode is entered when both Ethernet interfaces are configured as a bond in one of the following modes:

- active-backup (mode 1)
- balance-xor (mode 2)
- 802.3ad (LACP) (mode 4)

Any change of bonding configuration that negates one of the above rules (i.e., bonding mode is not 1, 2 or 4, or both Ethernet interfaces that belong to the same card are not the only slaves of the bond interface), will result in exiting RoCE LAG mode and the return to normal IB device per port configuration.

Once RoCE LAG is enabled, instead of having two IB devices; mlx5_0 and mlx5_1, there will be one device named mlx5_bond_0.

For information on how to configure RoCE LAG, refer to HowTo Configure RoCE over LAG (ConnectX-4/ConnectX-5/ConnectX-6) Community post.

Disabling RoCE

By default, RoCE is enabled on all mlx5 devices. When RoCE is enabled, all traffic to UDP port 4791 is treated as RoCE traffic by the device.

In case you are only interested in Ethernet (no RDMA) and wish to enable forwarding of traffic to this port, you can disable RoCE through sysfs:

echo <0|1> > /sys/devices/{pci-bus-address}/roce_enable



• Once RoCE is disabled, only Ethernet traffic will be supported. Therefore, there will be no GID tables and only Raw Ethernet QPs will be supported.

The current RoCE state can be queried by sysfs:

cat /svs/devices/{pci-bus-address}/roce enable

Enabling/Disabling RoCE on VMs via VFs

By default, when configuring VFs on the hypervisor, all VFs will be enabled with RoCE. This means they require more OS memory (from the VM). In case you are only interested in Ethernet (no RDMA) on the VM, and you wish to save the VM memory, you can disable RoCE on the VF from the hypervisor. In addition, by disabling RoCE, a VM can have the capability of utilizing the RoCE UDP port (4791) for standard UDP traffic.

For details on how to enable/disable RoCE on a VF, refer to HowTo Enable/Disable RoCE on VMs via VFs Community post.

Force DSCP

This feature enables setting a global traffic_class value for all RC QPs, or setting a specific traffic class based on several matching criteria.

Usage

- To set a single global traffic class to be applied to all QPs, write the desired global traffic_class value to /sys/class/infiniband/<dev>/tc/<port>/traffic_class.
 Note the following:
 - Negative values indicate that the feature is disabled. traffic_class value can be set using ibv_modify_qp()
 - Valid values range between 0 255
- The ToS field is 8 bits, while the DSCP field is 6 bits. To set a DSCP value of X, you need to multiply this value by 4 (SHIFT 2). For example, to set DSCP value of 24, set the ToS bit to 96 (24x4=96).
 - To set multiple traffic class values based on source and/or destination IPs, write the desired rule to /sys/class/infiniband/<dev>/tc/<port>/traffic_class. For example:

```
echo "tclass=16,src_ip=1.1.1.2,dst_ip=1.1.1.0/24" > /sys/class/infiniband/mlx5_0/tc/1/traffic_class
```

Note: Adding "tclass" prefix to tclass value is optional.

In the example above, traffic class 16 will be set to any QP with source IP 1.1.1.2 and destination IP 1.1.1.0/24.

Note that when setting a specific traffic class, the following rule precedence will apply:

- If a global traffic class value is set, it will be applied to all QPs
- If no global traffic class value is set, and there is a rule with matching source and destination IPs applicable to at least one QP, it will be applied
- Rules only with matching source and/or destination IPs have no defined precedence over other rules with matching source and/or destination IPs

Notes:

- A mask can be provided when using destination IPv4 addresses
- The rule precedence is not affected by the order in which rules are inserted
- Overlapping rules are entirely up to the administrator.
- "tclass=-1" will remove the rule from the database

Force Time to Live (TTL)

This feature enables setting a global TTL value for all RC QPs. Write the desired TTL value to /sys/class/infiniband/<dev>/tc/<port>/ttl. Valid values range between 0 - 255

Flow Control

Priority Flow Control (PFC)

Priority Flow Control (PFC) IEEE 802.1Qbb applies pause functionality to specific classes of traffic on the Ethernet link. For example, PFC can provide lossless service for the RoCE traffic and best-effort service for the standard Ethernet traffic. PFC can provide different levels of service to specific classes of Ethernet traffic (using IEEE 802.1p traffic classes).

Configuring PFC on ConnectX-4 and above

1. Enable PFC on the desired priority:

```
mlnx_qos -i <ethX> --pfc <0/1>,<0/1>,<0/1>,<0/1>,<0/1>,<0/1>,<0/1>,<0/1>
```

Example (Priority=4):

```
mlnx_qos -i eth1 --pfc 0,0,0,0,1,0,0,0
```

2. Create a VLAN interface:

```
vconfig add <ethX> <VLAN_ID>
```

Example (VLAN_ID=5):

```
vconfig add eth1 5
```

- 3. Set egress mapping:
 - a. For Ethernet traffic:

```
vconfig set_egress_map <vlan_einterface> <skprio> <up>
```

Example (skprio=3, up=5):

```
vconfig set_egress_map eth1.5 3 5
```

4. Create 8 Traffic Classes (TCs):

```
tc_wrap.py -i <interface>
```

5. Enable PFC on the switch.

For information on how to enable PFC on your respective switch, please refer to Switch FC/ PFC Configuration sections in RDMA/RoCE Solutions Community page.

PFC Configuration Using LLDP DCBX

PFC Configuration on Hosts

PFC Auto-Configuration Using LLDP Tool in the OS

1. Start lldpad daemon on host.

```
lldpad -d Or
service lldpad start
```

2. Send Ildpad packets to the switch.

3. Set the PFC parameters.

• For the CEE protocol, use dcbtool:

```
dcbtool sc <ethX> pfc pfcup:<xxxxxxx>
```

Example:

```
dcbtool sc eth6 pfc pfcup:01110001
```

where:

[pfcup:x
xxxxxx]

Enables/disables priority flow control. From left to right (priorities 0-7) - x can be equal to either 0 or 1. 1 indicates that the priority is configured to transmit priority pause.

• For IEEE protocol, use lldptool:

```
lldptool -T -i <ethX> -V PFC enabled=x,x,x,x,x,x,x
```

Example:

```
lldptool -T -i eth2 -V PFC enabled=1,2,4
```

where:

enab
led

Displays or sets the priorities with PFC enabled. The set attribute takes a commaseparated list of priorities to enable, or the string none to disable all priorities.

PFC Auto-Configuration Using LLDP in the Firmware (for mlx5 driver)

There are two ways to configure PFC and ETS on the server:

- 1. Local Configuration Configuring each server manually.
- 2. Remote Configuration Configuring PFC and ETS on the switch, after which the switch will pass the configuration to the server using LLDP DCBX TLVs.

There are two ways to implement the remote configuration using mlx5 driver:

- a. Configuring the adapter firmware to enable DCBX.
- b. Configuring the host to enable DCBX.

For further information on how to auto-configure PFC using LLDP in the firmware, refer to the <u>HowTo Auto-Config PFC and ETS on ConnectX-4 via LLDP DCBX</u>Community post.

PFC Configuration on Switches

1. In order to enable DCBX, LLDP should first be enabled:

```
switch (config) # 11dp
show 11dp interfaces ethernet remote
```

Add DCBX to the list of supported TLVs per required interface. For IEEE DCBX:

```
switch (config) # interface 1/1
switch (config interface ethernet 1/1) # 11dp tlv-select dcbx
```

For CEE DCBX:

```
switch (config) # interface 1/1
switch (config interface ethernet 1/1) # 11dp tlv-select dcbx-cee
```

3. [Optional] Application Priority can be configured on the switch, with the required ethertype and priority. For example, IP packet, priority 1:

```
switch (config) # dcb application-priority 0x8100 1
```

- Make sure PFC is enabled on the host (for enabling PFC on the host, refer to <u>PFC</u> <u>Configuration on Hosts</u> section above). Once it is enabled, it will be passed in the LLDP TLVs.
- 5. Enable PFC with the desired priority on the Ethernet port.

```
dcb priority-flow-control enable force
dcb priority-flow-control priority <priority> enable
interface ethernet <port> dcb priority-flow-control mode on force
```

Example - Enabling PFC with priority 3 on port 1/1:

```
dcb priority-flow-control enable force
dcb priority-flow-control priority 3 enable
interface ethernet 1/1 dcb priority-flow-control mode on force
```

Priority Counters

MLNX_OFED driver supports several ingress and egress counters per priority. Run ethtool -S to get the full list of port counters.

ConnectX-4 Counters

- Rx and Tx Counters:
 - Packets
 - Bytes
 - Octets
 - Frames
 - Pause
 - · Pause frames
 - Pause Duration
 - Pause Transition

ConnectX-4 Example

```
# ethtool -S eth35 | grep prio4
    prio4_rx_octets: 62147780800
    prio4_rx_frames: 14885696
    prio4_tx_octets: 0
    prio4_tx_frames: 0
    prio4_rx_pause: 0
    prio4_rx_pause: 0
    prio4_rx_pause_duration: 0
    prio4_rx_pause_duration: 14508
    prio4_rx_pause_transition: 0
```

Note: The Pause counters in ConnectX-4 are visible via ethtool only for priorities on which PFC is enabled.

PFC Storm Prevention

PFC storm prevention enables toggling between default and auto modes.

The stall prevention timeout is configured to 8 seconds by default. Auto mode sets the stall prevention timeout to be 100 msec.

The feature can be controlled using sysfs in the following directory: /sys/class/net/eth*/settings/pfc_stall_prevention

• To query the PFC stall prevention mode:

```
cat /sys/class/net/eth*/settings/pfc_stall_prevention
```

Example

```
$ cat /sys/class/net/ens6/settings/pfc_stall_prevention
default
```

• To configure the PFC stall prevention mode:

```
Echo "auto"/"default" > /sys/class/net/eth*/settings/pfc_stall_prevention
```

The following two counters were added to the ethtool -S:

- tx_Pause_storm_warning_events when the device is stalled for a period longer than a preconfigured watermark, the counter increases, allowing the debug utility an insight into current device status.
- tx_pause_storm_error_events when the device is stalled for a period longer than a preconfigured timeout, the pause transmission is disabled, and the counter increase.

Dropless Receive Queue (RQ)

Dropless RQ feature enables the driver to notify the FW when SW receive queues are overloaded. This scenario takes place when the handling of SW receive queue is slower than the handling of the HW receive queues.

When this feature is enabled, a packet that is received while the receive queue is full will not be immediately dropped. The FW will accumulate these packets assuming posting of new WQEs will resume shortly. If received WQEs are not posted after a certain period of time, out_of_buffer counter will increase, indicating that the packet has been dropped.

This feature is disabled by default. In order to activate it, ensure that Flow Control feature is also enabled.

To enable the feature, run:

```
ethtool --set-priv-flags ens6 dropless_rq on
```

To get the feature state, run:

```
ethtool --show-priv-flags DEVNAME
```

Output example:

```
Private flags for DEVNAME:

rx_cqe_moder : on

rx_cqe_compress: off

sniffer : off
dropless_rq : off
hw_lro : off
```

To disable the feature, run:

Explicit Congestion Notification (ECN)

ECN is an extension to the IP protocol. It allows reliable communication by notifying all ends of communication when congestion occurs. This is done without dropping packets.

Please note that this feature requires all nodes in the path (nodes, routers etc) between the communicating nodes to support ECN to ensure reliable communication. ECN is marked as 2 bits in the traffic control IP header. This ECN implementation refers to RoCE v2.

Enabling ECN

- To enable ECN on the hosts:
 - 1. Enable ECN in sysfs.

```
/sys/class/net/<interface>/<protocol>/ecn_<protocol>_enable =1
```

2. Query the attribute.

cat /sys/class/net/<interface>/ecn/<protocol>/params/<requested attribute>

3. Modify the attribute.

echo <value> /sys/class/net/<interface>/ecn/<protocol>/params/<requested attribute>

ECN supports the following algorithms:

- r_roce_ecn_rp Reaction point
- r_roce_ecn_np Notification point

Each algorithm has a set of relevant parameters and statistics, which are defined per device, per port, per priority.

To query whether ECN is enabled per Priority X:

cat /sys/class/net/<interface>/ecn/<protocol>/enable/X

To read ECN configurable parameters:

cat /sys/class/net/<interface>/ecn/<protocol>/requested attributes

To enable ECN for each priority per protocol:

echo 1 > /sys/class/net/<interface>/ecn/<protocol>/enable/X

To modify ECN configurable parameters:

echo <value> > /sys/class/net/<interface>/ecn/<protocol>/requested attributes

where:

• X: priority {0..7}

- protocol: roce_rp / roce_np
- requested attributes: Next Slide for each protocol.

RSS Support

RSS Hash Function

The device has the ability to use XOR as the RSS distribution function, instead of the default Toplitz function.

The XOR function can be better distributed among driver's receive queues in a small number of streams, where it distributes each TCP/UDP stream to a different queue.

MLNX_OFED provides the following option to change the working RSS hash function from Toplitz to XOR, and vice-versa:

Through sysfs, located at: /sys/class/net/eth*/settings/hfunc.

To guery the operational and supported hash functions:

cat /sys/class/net/eth*/settings/hfunc

Example:

cat /sys/class/net/eth2/settings/hfunc Operational hfunc: toeplitz Supported hfuncs: xor toeplitz



To change the operational hash function:

echo xor > /sys/class/net/eth*/settings/hfunc

RSS Verbs Support

Receive Side Scaling (RSS) technology allows spreading incoming traffic between different receive descriptor queues. Assigning each queue to different CPU cores allows better load balancing of the incoming traffic and improves performance.

This technology was extended to user space by the verbs layer and can be used for RAW ETH QP.

RSS Flow Steering

Steering rules classify incoming packets and deliver a specific traffic type (e.g. TCP/UDP, IP only) or a specific flow to "RX Hash" QP. "RX Hash" QP is responsible for spreading the traffic it handles between the Receive Work Queues using RX hash and Indirection Table. The Receive Work Queue can point to different CQs that can be associated with different CPU cores.

Verbs

The below verbs should be used to achieve this task in both control and data path. Details per verb should be referenced from its man page.

- ibv_create_wq, ibv_modify_wq, ibv_destory_wq
- ibv_create_rwq_ind_table, ibv_destroy_rwq_ind_table
- ibv_create_qp_ex with specific RX configuration to create the "RX hash" QP

Time-Stamping

Time-Stamping Service

Time-stamping is the process of keeping track of the creation of a packet. A time-stamping service supports assertions of proof that a datum existed before a particular time. Incoming packets are time-stamped before they are distributed on the PCI depending on the congestion in the PCI buffers. Outgoing packets are time-stamped very close to placing them on the wire.

Enabling Time-Stamping

Time-stamping is off by default and should be enabled before use.



To enable time-stamping for a socket:

Call setsockopt() with SO_TIMESTAMPING and with the following flags:

SOF_TIMESTAMPING_TX_HARDWARE:	try to obtain send time-stamp in hardware
SOF_TIMESTAMPING_TX_SOFTWARE:	if SOF_TIMESTAMPING_TX_HARDWARE is off or fails, then do it in software
SOF_TIMESTAMPING_RX_HARDWARE:	return the original, unmodified time-stamp as generated by the hardware
SOF_TIMESTAMPING_RX_SOFTWARE:	if SOF_TIMESTAMPING_RX_HARDWARE is off or fails, then do it in software
SOF_TIMESTAMPING_RAW_HARDWARE:	return original raw hardware time-stamp
SOF_TIMESTAMPING_SYS_HARDWARE:	return hardware time-stamp transformed into the system time base
SOF_TIMESTAMPING_SOFTWARE:	return system time-stamp generated in software
SOF_TIMESTAMPING_TX/RX	determine how time-stamps are generated
SOF_TIMESTAMPING_RAW/SYS	determine how they are reported



To enable time-stamping for a net device:

Admin privileged user can enable/disable time stamping through calling ioctl (sock, SIOCSH-WTSTAMP, &ifreq) with the following values:

• Send side time sampling, enabled by ifreq.hwtstamp_config.tx_type when:

Receive side time sampling, enabled by ifreq.hwtstamp_config.rx_filter when:

```
* possible values for hwtstamp_config->rx_filter */
enum hwtstamp_rx_filters {
    /* time stamp no incoming packet at all */
    HWTSTAMP_FILTER_NONE,
               /* time stamp any incoming packet */
HWTSTAMP_FILTER_ALL,
* return value: time stamp all packets requested plus some others */
HWTSTAMP_FILTER_SOME,
                 /* PTP v1, UDP, any kind of event packet */
HWTSTAMP_FILTER_PTP_V1_L4_EVENT,
/* PTP v1, UDP. Sync packet */
                /* PTP v1, UDP, Sync packet */
HWTSTAMP_FILTER_PTP_V1_L4_SYNC,
                HWTSTAMP_FILTER_PTP_V1_L4_SINV.,
/* PTP v1, UDP, Delay_req packet */
HWTSTAMP_FILTER_PTP_V1_L4_DELAY_RRQ,
/* PTP v2, UDP, any kind of event pa
HWTSTAMP_FILTER_PTP_V2_L4_EVENT,
'* STD v2 UDD Sync_packet */
                                                                             event packet */
                 /* PTP v2, UDP, Sync packet */
HWTSTAMP_FILTER_PTP_V2_L4_SYNC,
                /* PTP v2, UDP, Delay_req packet */
HWTSTAMP_FILTER_PTP_V2_L4_DELAY_REQ,
                 /* 802.AS1, Ethernet, any kind of event packet */
HWTSTAMP_FILTER_PTP_V2_L2_EVENT,
                 /* 802.AS1, Ethernet, Sync packet */
HWTSTAMP_FILTER_PTP_V2_L2_SYNC,
                 /* 802.AS1, Ethernet, Delay_req packet */
HWTSTAMP_FILTER_PTP_V2_L2_DELAY_REQ,
                /* PTP v2/802.AS1, any layer, any kind of event packet */
HWTSTAMP_FILTER_PTP_V2_EVENT,
/* PTP v2/802.AS1, any layer, Sync packet */
HWTSTAMP_FILTER_PTP_V2_SYNC,
/* PTP v2/802.AS1, any layer, Delay_req packet */
                 /* PTP v2/802.AS1, any layer, Del.
HWTSTAMP_FILTER_PTP_V2_DELAY_REQ,
                                                                                Delay_req packet */
Note: for receive side time stamping currently only HWTSTAMP_FILTER_NONE and HWTSTAMP_FILTER_ALL are supported.
```

Getting Time-Stamping

Once time stamping is enabled time stamp is placed in the socket Ancillary data. recvmsg() can be used to get this control message for regular incoming packets. For send time stamps the outgoing packet is looped back to the socket's error queue with the send time-stamp(s) attached. It can be received with recvmsg (flags=MSG_ERRQUEUE). The call returns the original outgoing packet data including all headers prepended down to and including the link layer, the scm_time-stamping control message and a sock_extended_err control message with ee_errno==ENOMSG and ee_origin==SO_EE_ORIGIN_TIMESTAMPING. A socket with such a pending bounced packet is ready for reading as far as select() is concerned. If the outgoing packet has to be fragmented, then only the first fragment is time stamped and returned to the sending socket.

When time-stamping is enabled, VLAN stripping is disabled. For more info please refer to Documentation/networking/timestamping.txt in kernel.org.



On ConnectX-4 and above adapter cards, when time-stamping is enabled, RX CQE compression is disabled (features are mutually exclusive).

Time Stamping Capabilities via ethtool



To display Time Stamping capabilities via ethtool:

Show Time Stamping capabilities:

```
ethtool -T eth<x>
```

Example:

```
ethtool -T eth0
Time stamping parameters for p2p1:
Capabilities:
                                                                                                                                                                (SOF_TIMESTAMPING_TX_HARDWARE)
(SOF_TIMESTAMPING_TX_SOFTWARE)
(SOF_TIMESTAMPING_RX_HARDWARE)
(SOF_TIMESTAMPING_RX_SOFTWARE)
(SOF_TIMESTAMPING_RX_SOFTWARE)
(SOF_TIMESTAMPING_RAW_HARDWARE)
                                                                                          hardware-transmit
software-transmit
hardware-receive
software-receive
                                                                                          software-system-clock
hardware-raw-clock
PTP Hardware Clock: 1
Hardware Transmit Timestamp Modes:
off
                                                                                           (HWTSTAMP_TX_OFF)
(HWTSTAMP_TX_ON)
Hardware Receive Filter Modes:
                                                                                           (HWTSTAMP_FILTER_NONE)
(HWTSTAMP_FILTER_ALL)
```

For more details on PTP Hardware Clock, please refer to: https://www.kernel.org/doc/ Documentation/ptp/ptp.txt

Steering PTP Traffic to Single RX Ring

As a result of Receive Side Steering (RSS) PTP traffic coming to UDP ports 319 and 320, it may reach the user space application in an out of order manner. In order to prevent this, PTP traffic needs to be steered to single RX ring using ethtool.

Example:

```
# ethtool -u ens7
8 RX rings available
Total 0 rules
# ethtool -U ens7 flow-type udp4 dst-port 319 action 0 loc 1
# ethtool -U ens7 flow-type udp4 dst-port 320 action 0 loc 0
# ethtool -u ens7
8 RX rings available
Total 2 rules
Filter: 0
Rule Type: UDP over IPv4
Src IP addr: 0.0.0.0 mask: 255.255.255
Dest IP addr: 0.0.0.0 mask: 255.255.255
TOS: 0x0 mask: 0xff
Src port: 0 mask: 0xfff
Dest port: 320 mask: 0x0
Action: Direct to queue 0
Filter: 1
Rule Type: UDP over IPv4
Src IP addr: 0.0.0.0 mask: 255.255.255
TOS: 0x0 mask: 0xff
Src port: 0 mask: 0xfff
Src port: 0 mask: 0xffff
Src port addr: 0.0.0.0 mask: 255.255.255
TOS: 0x0 mask: 0xfff
Src port: 0 mask: 0xffff
Dest port: 319 mask: 0xffff
Dest port: 319 mask: 0x0
Action: Direct to queue 0
```

Tx Port Time-Stamping

Transmitted packet time-stamping accuracy can be improved when using a timestamp generated at the port level instead of a timestamp generated upon CQE creation. Tx port time-stamping better reflects the actual time of a packet's transmission.

Normal Send queues (SQs) are open with CQE time-stamp support. When this feature is enabled, the driver is expected to open extra Tx port time-stamped SQ per traffic class (TC).

The stream must meet the following conditions in order to be transmitted through a Tx port time-stamped SQ.

- SKBTX_HW_TSTAMP flag was set at tx_flag (SO_TIMESTAMPING was set via setsockopt() or similarly)
- 2. Packet type is:
 - a. Non-IP, with EtherType of PTP over IEEE 802.3 (0x88f7)
 - b. UDP over IPv4/IPv6

This feature is disabled by default in order to avoid extra SQ memory allocations. The feature can be enabled or disabled using the following command.

```
ethtool --set-priv-flags <ifs-name> tx_port_ts on / off
```

RoCE Time-Stamping

RoCE Time-Stamping allows you to stamp packets when they are sent to the wire/received from the wire. The time-stamp is given in raw hardware cycles but could be easily converted into hardware referenced nanoseconds based time. Additionally, it enables you to query the hardware for the hardware time, thus stamp other application's event and compare time.

Query Capabilities

Time-stamping is available if and only the hardware reports it is capable of reporting it. To verify whether RoCE Time-Stamping is available, run ibv_query_device_ex.

For further information, please see ibv_query_device_ex manual page.

Creating a Time-Stamping Completion Queue

To get time stamps, a suitable extended Completion Queue (CQ) must be created via a special call to ibv_create_cq_ex verb.

For further information, please see ibv_create_cq_ex manual page.

⚠ Time Stamping in not available when CQE zipping is used.

Querying the Hardware Time

Querying the hardware for time is done via the ibv_query_rt_values_ex verb. For example:

For further information, please see ibv_query_rt_values_ex manual page.

One Pulse Per Second (1PPS)

1PPS is a time synchronization feature that allows the adapter to be able to send or receive 1 pulse per second on a dedicated pin on the adapter card using an SMA connector (SubMiniature version A). Only one pin is supported and could be configured as 1PPS in or 1PPS out.

For further information, refer to HowTo Test 1PPS on Mellanox Adapters Community post.

Flow Steering

Flow steering is a new model which steers network flows based on flow specifications to specific QPs. Those flows can be either unicast or multicast network flows. In order to maintain flexibility, domains and priorities are used. Flow steering uses a methodology of flow attribute, which is a combination of L2-L4 flow specifications, a destination QP and a priority. Flow steering rules may be inserted either by using ethtool or by using InfiniBand verbs. The verbs abstraction uses different terminology from the flow attribute (ibv_flow_attr), defined by a combination of specifications (struct ibv_flow_spec_*).

Flow Steering Support

All flow steering features are enabled in the supported adapter cards.

Flow Steering support in InfiniBand is determined according to the MANAGED_FLOW_STEERING flag.

Flow Domains and Priorities

Flow steering defines the concept of domain and priority. Each domain represents a user agent that can attach a flow. The domains are prioritized. A higher priority domain will always supersede a lower priority domain when their flow specifications overlap. Setting a lower priority value will result in a higher priority.

In addition to the domain, there is a priority within each of the domains. Each domain can have at most 2¹² priorities in accordance with its needs.

The following are the domains at a descending order of priority:

- User Verbs allows a user application QP to be attached to a specified flow when using ibv_create_flow and ibv_destroy_flow verbs
 - ibv_create_flow

```
struct ibv_flow *ibv_create_flow(struct ibv_qp *qp, struct ibv_flow_attr
 *flow)
```

Input parameters:

- struct ibv_qp the attached QP.
- struct ibv_flow_attr attaches the QP to the flow specified. The flow contains mandatory control parameters and optional L2, L3 and L4 headers. The optional headers are detected by setting the size and num_of_specs fields: struct ibv_flow_attr can be followed by the optional flow headers structs:

```
struct ibv_flow_spec_eth
struct ibv_flow_spec_ipv4
struct ibv_flow_spec_tcp_udp
struct ibv_flow_spec_ipv6
```

For further information, please refer to the <code>ibv_create_flow</code> man page.

• ibv_destroy_flow

```
int ibv_destroy_flow(struct ibv_flow *flow_id)
```

Input parameters:

ibv_destroy_flow requires struct ibv_low which is the return value of ibv_create_flow in case of success.

Output parameters:

Returns 0 on success, or the value of errno on failure.

For further information, please refer to the <code>ibv_destroy_flow</code> man page.

Ethtool

Ethtool domain is used to attach an RX ring, specifically its QP to a specified flow. Please refer to the most recent ethtool man page for all the ways to specify a flow.

Examples:

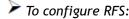
- ethtool -U eth5 flow-type ether dst 00:11:22:33:44:55 loc 5 action 2 All packets that contain the above destination MAC address are to be steered into rx-ring 2 (its underlying QP), with priority 5 (within the ethtool domain)
- ethtool -U eth5 flow-type tcp4 src-ip 1.2.3.4 dst-port 8888 loc 5 action 2
 All packets that contain the above destination IP address and source port are to be steered
 into rx- ring 2. When destination MAC is not given, the user's destination MAC is filled
 automatically.
- ethtool -U eth5 flow-type ether dst 00:11:22:33:44:55 vlan 45 m 0xf000 loc 5 action 2
 All packets that contain the above destination MAC address and specific VLAN are steered
 into ring 2. Please pay attention to the VLAN's mask 0xf000. It is required in order to add such
 a rule.
- ethtool -u eth5
 Shows all of ethtool's steering rule

When configuring two rules with the same priority, the second rule will overwrite the first one, so this ethtool interface is effectively a table. Inserting Flow Steering rules in the kernel requires support from both the ethtool in the user space and in kernel (v2.6.28).

Accelerated Receive Flow Steering (aRFS)

Receive Flow Steering (RFS) and Accelerated Receive Flow Steering (aRFS) are kernel features currently available in most distributions. For RFS, packets are forwarded based on the location of the application consuming the packet. aRFS boosts the speed of RFS by adding support for the hardware. By using aRFS (unlike RFS), the packets are directed to a CPU that is local to the thread running the application.

aRFSis an in-kernel-logic responsible for load balancing between CPUs by attaching flows to CPUs that are used by flow's owner applications. This domain allows the aRFS mechanism to use the flow steering infrastructure to support the aRFS logic by implementing the ndo_rx_flow_steer, which, in turn, calls the underlying flow steering mechanism with the aRFS domain.



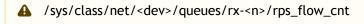
Configure the RFS flow table entries (globally and per core).

Note: The functionality remains disabled until explicitly configured (by default it is 0).

• The number of entries in the global flow table is set as follows:



• The number of entries in the per-queue flow table are set as follows:



Example:

```
# echo 32768 > /proc/sys/net/core/rps_sock_flow_entries
# for f in /sys/class/net/ens6/queues/rx-*/rps_flow_ent; do echo 32768 > $f; done
```

To Configure aRFS:

The aRFS feature requires explicit configuration in order to enable it. Enabling the aRFS requires enabling the 'ntuple' flag via the ethtool.

For example, to enable ntuple for eth0, run:

```
ethtool -K eth0 ntuple on
```

aRFS requires the kernel to be compiled with the CONFIG_RFS_ACCEL option. This option is available in kernels 2.6.39 and above. Furthermore, aRFS requires Device Managed Flow Steering support.



RFS cannot function if LRO is enabled. LRO can be disabled via ethtool.

Flow Steering Dump Tool

The mlx_fs_dump is a python tool that prints the steering rules in a readable manner. Python v2.7 or above, as well as pip, anytree and termcolor libraries are required to be installed on the host.

Running example:

```
./ofed_scripts/utils/mlx_fs_dump -d /dev/mst/mt4115_pciconf0
FT: 9 (level: 0x18, type: NIC_RX)
+-- FG: 0x15 (MISC)
|-- FTE: 0x0 (FWD) to (TIR:0x7e) out.ethtype:IPv4 out.ip_prot:UDP out.udp_dport:0x140
+-- FTE: 0x1 (FWD) to (TIR:0x7e) out.ethtype:IPv4 out.ip_prot:UDP out.udp_dport:0x13f
```

For further information on the mlx_fs_dump tool, please refer to mlx_fs_dump Community post.

Wake-on-LAN (WoL)

Wake-on-LAN (WoL) is a technology that allows a network professional to remotely power on a computer or to wake it up from sleep mode.

· To enable WoL:

```
ethtool -s <interface> wol q
```

To get WoL:

```
ethtool <interface> | grep Wake-on Wake-on: g
```

Where:

"g" is the magic packet activity.

Hardware Accelerated 802.1ad VLAN (Q-in-Q Tunneling)

Q-in-Q tunneling allows the user to create a Layer 2 Ethernet connection between two servers. The user can segregate a different VLAN traffic on a link or bundle different VLANs into a single VLAN. Q-in-Q tunneling adds a service VLAN tag before the user's 802.1Q VLAN tags. For Q-in-Q support in virtualized environments (SR-IOV), please refer to "Q-in-Q Encapsulation per

VF in Linux (VST)".



To enable device support for accelerated 802.1ad VLAN:

1. Turn on the new ethtool private flag "phy-bit" (disabled by default).

```
$ ethtool --set-priv-flags eth1 phv-bit on
```

Enabling this flag sets the phv_en port capability.

2. Change the interface device features by turning on the ethtool device feature "tx-vlan- staghw-insert" (disabled by default).

```
$ ethtool -K eth1 tx-vlan-stag-hw-insert on
```

Once the private flag and the ethtool device feature are set, the device will be ready for 802.1ad VLAN acceleration.



The "phv-bit" private flag setting is available for the Physical Function (PF) only. The Virtual Function (VF) can use the VLAN acceleration by setting the "tx-vlan-staghw-insert" parameter only if the private flag "phv-bit" is enabled by the PF. If the PF enables/disables the "phy-bit" flag after the VF driver is up, the configuration will take place only after the VF driver is restarted.

VLAN Stripping in Linux Verbs



A This capability is now accessible from userspace using the verbs.

VLAN stripping adds access to the device's ability to offload the Customer VLAN (cVLAN) header stripping from an incoming packet, thus achieving acceleration of VLAN handing in receive flow. It is configured per WQ option. You can either enable it upon creation or modify it later using the appropriate verbs (ibv create wg/ibv modify wg).

Offloaded Traffic Sniffer



To be able to activate this feature, make sure library v1.9 or above is installed on your setup.

To download libpcap, please visit https://www.tcpdump.org/.

Offloaded Traffic Sniffer allows bypass kernel traffic (such as RoCE, VMA, and DPDK) to be captured by existing packet analyzer, such as tcpdump.

To capture the interface's bypass kernel traffic, run tcpdump on the RDMA device.

For examples on how to dump RDMA traffic using the Inbox tcpdump tool for ConnectX-4 adapter cards and above, click here.



A Note that enabling Offloaded Traffic Sniffer can cause bypass kernel traffic speed degradation.



A In case you do not wish to install libpcap on your setup, you can use docker to run the tcpdump. For further information, please see https://hub.docker.com/r/mellanox/ tcpdump-rdma.

Dump Configuration

This feature helps dumping driver and firmware configuration using ethtool. It creates a backup of the configuration files into a specified dump file.

Dump Parameters (Bitmap Flag)

The following bitmap parameters are used to set the type of dump:

Bitmap Parameters

Value	Description	
1	MST dump	
2	Ring dump (Software context information for SQs, EQs, RQs, CQs)	
3	MST dump + Ring dump (1+2)	
4	Clear this parameter	

Configuration

In order to configure this feature, follow the steps below:

1. Set the dump bitmap parameter by running -W (uppercase) with the desired bitmap parameter value (see Bitmap Parameters table above). In the following example, the bitmap parameter value is 3.

```
ethtool -W ens1f0 3
```

2. Dump the file by running -w (lowercase) with the desired configuration file name.

```
ethtool -w ens1f0 data /tmp/dump.bin
```

3. [Optional] To get the bitmap parameter value, version and size of the dump, run the command above without the file name.

```
ethtool -w ens1f0 flag: 3, version: 1, length: 4312
```

4. To open the dump file, run:

```
mlnx_dump_parser -f /tmp/dump.bin -m mst_dump_demo.txt -r ring_dump_demo.txt
Version: 1 Flag: 3 Number of blocks: 123 Length 327584
MCION module number: 0 status: | present |
DRIVER VERSION: 1-23 (03 Mar 2015)
DEVICE NAME 0000:81:00.0:ens1f0
Parsing Complete!
```

where:

-f	For the file to be parsed (the file that was just created)	
-m	For the mst dump file	
-r	For the ring dump file	

For further information, refer to <u>HowTo Dump Driver Configuration (via ethtool)</u> Community post.

Output:

```
# mlnx_dump_parser -f /tmp/dump.bin -m mst_dump_demo.txt -r ring_dump_demo.txt
Version: 1 Flag: 3 Number of blocks: 123 Length 327584
MCION module number: 0 status: | present |
DRIVER VERSION: 1-23 (03 Mar 2015)
DEVICE NAME 0000:81:00.0:ens1f0
Parsing Complete!
```

- 5. Open the files.
 - a. The MST dump file will look as follows. In order to analyze it, contact Mellanox Support at support@mellanox.com.

b. The Ring dump file can help developers debug ring-related issues, and it looks as follows:

```
# cat ring_dump_demo.txt
SQ TYPE: 3, WQN: 102, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024...
SQ TYPE: 3, WQN: 102, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 65536, GROUP_IP: 0
CQ TYPE: 5, WQN: 20, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 1024, GROUP_IP: 0
RQ TYPE: 4, WQN: 103, PI: 15, CI: 0, STRIDE: 6, SIZE: 16, WQE_NUM: 512, GROUP_IP: 0
CQ TYPE: 5, WQN: 21, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 512, GROUP_IP: 0
EQ TYPE: 6, CI: 1, SIZE: 0, IRQN: 109, EQN: 19, NENT: 2048, MASK: 0, INDEX: 0, GROUP_ID: 0
SQ TYPE: 3, WQN: 106, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 65536, GROUP_IP: 1
CQ TYPE: 5, WQN: 23, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 65536, GROUP_IP: 1
CQ TYPE: 4, WQN: 107, PI: 15, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 512, GROUP_IP: 1
CQ TYPE: 5, WQN: 24, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 16384, GROUP_IP: 1
SQ TYPE: 3, WQN: 24, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 165536, GROUP_IP: 1
SQ TYPE: 3, WQN: 110, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 65536, GROUP_IP: 1
SQ TYPE: 3, WQN: 110, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 65536, GROUP_IP: 2
CQ TYPE: 5, WQN: 26, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 65536, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 1024, WQE_NUM: 1024, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 1024, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 16384, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 16384, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 16384, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 16384, GROUP_IP: 2
CQ TYPE: 5, WQN: 27, PI: 0, CI: 0, STRIDE: 6, SIZE: 16384, WQE_NUM: 16384, GROUP_IP: 2
```

Local Loopback Disable

Local Loopback Disable feature allows users to force the disablement of local loopback on the virtual port (vport). This disables both unicast and mutlicast loopback in the hardware.



To enable Local Loopback Disable, run the following command:

```
echo 1 > /sys/class/net/<ifname>/settings/force_local_lb_disable"
```



To disable Local Loopback Disable, run the following command:

```
echo 0 > /sys/class/net/<ifname>/settings/force_local_lb_disable"
```

Δ

When turned off, the driver configures the loopback mode according to its own logic.

Kernel Transport Layer Security (kTLS) Offloads

A This feature is supported on ConnectX-6 Dx crypto cards only.

Overview

Transport Layer Security (TLS) is a widely-deployed protocol used for securing TCP connections on the Internet. TLS is also a required feature for HTTP/2, the latest web standard. Kernel implementation of TLS (kTLS) provides new opportunities for offloading the protocol into the hardware.

TLS data-path offload allows the NIC to accelerate encryption, decryption and authentication of AES-GCM. TLS offload handles data as it goes through the device without storing any data, but only updating context. If the packet cannot be encrypted/decrypted by the device, then a software fallback handles the packet.

Establishing a kTLS Connection

To avoid unnecessary complexity in the kernel, the TLS handshake is kept in the user space. A full TLS connection using the socket is done using the following scheme:

- 1. Call connect() or accept() on a standard TCP file descriptor.
- 2. Use a user space TLS library to complete a handshake.
- 3. Create a new KTLS socket file descriptor.
- 4. Extract the TLS Initialization Vectors (IVs), session keys, and sequence IDs from the TLS library. Use the setsockopt function on the kTLS file descriptor (FD) to pass them to the
- 5. Use standard read(), write(), sendfile() and splice() system calls on the KTLS FD.

Drivers can offer Tx and Rx packet encryption/decryption offload from the kernel into the NIC hardware. Upon receipt of a non-data TLS message (a control message), the kTLS socket returns an error, and the message is left on the original TCP socket instead. The kTLS socket is automatically unattached. Transfer of control back to the original encrypted FD is done by calling getsockopt to receive the current sequence numbers, and inserting them into the TLS library.

Kernel Support

For support in the kernel, make sure the following flags are set as follows.

- CONFIG_TLS=y
- CONFIG_TLS_DEVICE=y | m

⚠ For kTLS Tx device offloads with OFED drivers, kernel TLS module (kernel/net/tls) must be aligned to kernel v5.3 and above.

For kTLS Rx device offloads with OFED drivers, kernel TLS module (kernel/net/tls) must be aligned to kernel v5.9 and above.

Configuring kTLS Offloads



To enable kTLS Tx offload, run:

ethtool -K <ifs> tls-hw-tx-offload on



To enable kTLS Rx offload, run:

ethtool -K <ifs> tls-hw-rx-offload on

For further information on TLS offloads, please visit the following kernel documentation:

- https://www.kernel.org/doc/html/latest/networking/tls-offload.html
- https://www.kernel.org/doc/html/latest/networking/tls.html#kernel-tls

OpenSSL with kTLS Offload

OpenSSL version 3.0.0 or above is required to support kTLS TX/RX offloads.

Supported OpenSSL version is available to download from distro packages, or can be downloaded and compiled from the OpenSSL github.

IPsec Crypto Offload

This feature is supported on ConnectX-6 Dx adapter cards (with crypto unit) only.

Overview and Configuration

IPsec crypto offload feature, also known as IPsec inline offload or IPsec aware offload feature enables the user to offload IPsec crypto encryption and decryption operations to the hardware.

Note that the hardware implementation only supports AES-GCM encryption scheme.

To enable the feature, support in both kernel and adapter firmware is required.

• For support in the kernel, make sure the following flags are set as follows.

CONFIG_XFRM_OFFLOAD=y CONFIG_INET_ESP_OFFLOAD=m CONFIG_INET6_ESP_OFFLOAD=m

Note: These flags are enabled by default in RedHat 8 and Ubuntu 18.04.

• For support in the firmware, make sure the below string is found in the dmesg.

mlx5e: IPSec ESP acceleration enabled

Configuring Security Associations for IPsec Offloads

To program the inline offload security associations (SA), add the option "offload dev <netdev interface> dir out/in" in the "ip xfrm state" command for transmitting and receiving SA.

Transmit inline offload SA xfrm command example:

```
sudo ip xfrm state add src 192.168.1.64/24 dst 192.168.1.65/24 proto esp spi 0x46dc6204 reqid 0x46dc6204 mode transport aead 'rfc4106(gcm(aes))' 0x60bd6c3eafba371a46411830fd56c53af93883261ed1fb26767820ff493f43ba35b0dcca 128 offload dev p4p1 dir out sel src 192.168.1.64 dst 192.168.1.65
```

Receive inline offload SA xfrm command example:

sudo ip xfrm state add src 192.168.1.65/24 dst 192.168.1.64/24 proto esp spi 0xaea0846c reqid 0xaea0846c mode transport aead 'rfc4106(gcm(aes))' 0x81d5c3167c912c1dd50dab0cb4b6d815b6ace8844304db362215a258cd19deda8f89deda 128 offload dev p4p1 dir in sel src 192.168.1.65 dst 192.168.1.64

InfiniBand Network

The chapter contains the following sections:

- · InfiniBand Interface
- OpenSM
- QoS Quality of Service
- IP over InfiniBand (IPoIB)
- Advanced Transport
- Optimized Memory Access
- Mellanox PeerDirect®
- CPU Overhead Distribution
- Out-of-Order (OOO) Data Placement
- IB Router
- MAD Congestion Control

InfiniBand Interface

Port Type Management

For information on port type management of ConnectX-4 and above adapter cards, please refer to Port Type Management/VPI Cards Configuration section.

RDMA Counters

- RDMA counters are available only through sysfs located under:
 - # /sys/class/infiniband/<device>/ports/*/counters/
 - # /sys/class/infiniband/<device>/ports/*/hw_counters/

For mlx5 port and RDMA counters, refer to the <u>Understanding mlx5 Linux Counters</u> Community post.

OpenSM

OpenSM is an InfiniBand compliant Subnet Manager (SM). It is provided as a fixed flow executable called <u>"opensm"</u>, accompanied by a testing application called <u>"osmtest"</u>. OpenSM implements an InfiniBand compliant SM according to the InfiniBand Architecture Specification chapters: Management Model, Subnet Management, and Subnet Administration.

opensm

opensm is an InfiniBand compliant Subnet Manager and Subnet Administrator that runs on top of the Mellanox OFED stack. opensm performs the InfiniBand specification's required tasks for initializing InfiniBand hardware. One SM must be running for each InfiniBand subnet.

opensm also provides an experimental version of a performance manager.

opensm defaults were designed to meet the common case usage on clusters with up to a few hundred nodes. Thus, in this default mode, opensm will scan the IB fabric, initialize it, and sweep occasionally for changes.

opensm attaches to a specific IB port on the local machine and configures only the fabric connected to it. (If the local machine has other IB ports, opensm will ignore the fabrics connected to those other ports). If no port is specified, opensm will select the first "best" available port. opensm can also present the available ports and prompt for a port number to attach to.

By default, the opensm run is logged to two files: /var/log/messages and /var/log/ opensm.log. The first file will register only general major events, whereas the second file will include details of reported errors. All errors reported in this second file should be treated as indicators of IB fabric health issues. (Note that when a fatal and non-recoverable error occurs, opensm will exit). Both log files should include the message "SUBNET UP" if opensm was able to set up the subnet correctly.

Syntax

opensm [OPTIONS]

For the complete list of opensm options, please run:

opensm --help / -h / -?

Environment Variables

The following environment variables control opensm behavior:

- OSM_TMP_DIR controls the directory in which the temporary files generated by opensm are created. These files are: opensm-subnet.lst, opensm.fdbs, and opensm.mcfdbs. By default, this directory is /var/log.
- OSM_CACHE_DIR opensm stores certain data to the disk such that subsequent runs are consistent. The default directory used is /var/cache/opensm. The following file is included in it:

guid2lid - stores the LID range assigned to each GUID

Signaling

When OpenSM receives a HUP signal, it starts a new heavy sweep as if a trap has been received or a topology change has been found.

Also, SIGUSR1 can be used to trigger a reopen of /var/log/opensm.log for logrotate purposes.

Running opensm

The defaults of opensm were designed to meet the common case usage on clusters with up to a few hundred nodes. Thus, in this default mode, opensm will scan the IB fabric, initialize it, and sweep occasionally for changes.

To run opensm in the default mode, simply enter:

host1# opensm

Note that opensm needs to be run on at least one machine in an IB subnet.

By default, an opensm run is logged to two files: /var/log/messages and /var/log/ opensm.log. The first file, message, registers only general major events; the second file, opensm.log, includes details of reported errors. All errors reported in opensm.log should be treated as indicators of IB fabric health. Both log files should include the message "SUBNET UP" if opensm was able to set up the subnet correctly.

⚠ If a fatal, non-recoverable error occurs, OpenSM will exit.

Running OpenSM As Daemon

OpenSM can also run as daemon. To run OpenSM in this mode, enter:

host1# /etc/init.d/opensmd start

osmtest

osmtest is a test program for validating the InfiniBand Subnet Manager and Subnet Administrator. osmtest provides a test suite for opensm. It can create an inventory file of all available nodes, ports, and PathRecords, including all their fields. It can also verify the existing inventory with all the object fields and matches it to a pre-saved one. osmtest has the following test flows:

- Multicast Compliancy test
- Event Forwarding test
- Service Record registration test
- RMPP stress test
- Small SA Queries stress test

For further information, please refer to the tool's man page.

Partitions

OpenSM enables the configuration of partitions (PKeys) in an InfiniBand fabric. By default, OpenSM searches for the partitions configuration file under the name /etc/opensm/partitions.conf. To change this filename, you can use opensm with the '--Pconfig' or '-P' flags.

The default partition is created by OpenSM unconditionally, even when a partition configuration file does not exist or cannot be accessed.

The default partition has a P_Key value of 0x7fff. The port out of which runs OpenSM is assigned full membership in the default partition. All other end-ports are assigned partial membership.

File Format



⚠ Line content followed after '#' character is comment and ignored by parser.

General File Format

<Partition Definition>:\[<newline>\]<Partition Properties>

<Partition Definition>:

[PartitionName\]\[=PKey\]\[,ipoib_bc_flags\]\[,defmember=full|limited\]

where:

PartitionName	String, will be used with logging. When omitted empty string will be used.
PKey	P_Key value for this partition. Only low 15 bits will be used. When omitted will be auto-generated.
ipoib_bc_flags	Used to indicate/specify IPoIB capability of this partition.
defmember=full limited both	Specifies default membership for port GUID list. Default is limited.

ipoib_bc_flags are:

ipoib	Indicates that this partition may be used for IPoIB, as a result the IPoIB broadcast group will be created with the flags given, if any.
rate= <val></val>	Specifies rate for this IPoIB MC group (default is 3 (10GBps))
mtu= <val></val>	Specifies MTU for this IPoIB MC group (default is 4 (2048))
sl= <val></val>	Specifies SL for this IPoIB MC group (default is 0)
scope= <val></val>	Specifies scope for this IPoIB MC group (default is 2 (link local))

• <Partition Properties>:

 $\[<Port list> | <MCast Group> \] * | <Port list> \]$

<Port List>:

<Port Specifier>[,<Port Specifier>]

• <Port Specifier>:

<PortGUID>[=[full|limited|both]]

where

PortGUID	GUID of partition member EndPort. Hexadecimal numbers should start from 0x, decimal numbers are accepted too.
full, limited	Indicates full and/or limited membership for this both port. When omitted (or unrecognized) limited membership is assumed. Both indicate full and limited membership for this port.

<MCast Group>:

mgid=gid[,mgroup_flag]*<newline>

where:

where.		
mgid=gi d	gid specified is verified to be a Multicast address IP groups are verified to match the rate and mtu of the broadcast group. The P_Key bits of the mgid for IP groups are verified to either match the P_Key specified in by "Partition Definition" or if they are 0x0000 the P_Key will be copied into those bits.	
mgroup_ flag	rate= <v al></v 	Specifies rate for this MC group (default is 3 (10GBps))
	mtu= <va< td=""><td>Specifies MTU for this MC group (default is 4 (2048))</td></va<>	Specifies MTU for this MC group (default is 4 (2048))
	sl= <val< td=""><td>Specifies SL for this MC group (default is 0)</td></val<>	Specifies SL for this MC group (default is 0)
	scope=< val>	Specifies scope for this MC group (default is 2 (link local)). Multiple scope settings are permitted for a partition. NOTE: This overwrites the scope nibble of the specified mgid. Furthermore specifying multiple scope settings will result in multiple MC groups being created.
	qkey= <v al></v 	Specifies the Q_Key for this MC group (default: 0x0b1b for IP groups, 0 for other groups)
	tclass= <val></val>	Specifies tclass for this MC group (default is 0)
	FlowLab el= <val ></val 	Specifies FlowLabel for this MC group (default is 0)

Note that values for rate, MTU, and scope should be specified as defined in the IBTA specification (for example, mtu=4 for 2048). To use 4K MTU, edit that entry to "mtu=5" (5 indicates 4K MTU to that specific partition).

PortGUIDs list:

PortGUID GUID of partition member EndPort. Hexadecimal numbers should start from 0x, decimal numbers are accepted too. full or limited indicates full or limited membership for this port. When omitted (or unrecognized) limited membership is assumed.

There are some useful keywords for PortGUID definition:

- 'ALL_CAS' means all Channel Adapter end ports in this subnet
- 'ALL_VCAS' means all virtual end ports in the subnet
- 'ALL_SWITCHES' means all Switch end ports in this subnet
- 'ALL_ROUTERS' means all Router end ports in this subnet
- 'SELF' means subnet manager's port. An empty list means that there are no ports in this partition

Notes:

- White space is permitted between delimiters ('=', ',',':',';').
- PartitionName does not need to be unique, PKey does need to be unique. If PKey is repeated then those partition configurations will be merged and the first PartitionName will be used (see the next note).
- It is possible to split partition configuration in more than one definition, but then PKey should be explicitly specified (otherwise different PKey values will be generated for those definitions).

Examples:

```
Default=0x7fff : ALL, SELF=full ;
Default=0x7fff : ALL, ALL_SWITCHES=full , SELF=full ;
NewPartition , ipoib : 0x123456=full , 0x3456789034=limi , 0x2134af2306 ;
YetAnotherOne = 0x300 : SELF=full ;
YetAnotherOne = 0x300 : ALL=limited ;
ShareIO = 0x80 , defmember=full : 0x123451 , 0x123452;
# 0x123453 , 0x123454 will be limited
ShareIO = 0x80 : 0x123453 , 0x123454 , 0x123455=full;
# 0x123456 , 0x123457 will be limited
ShareIO = 0x80 : defmember=limited : 0x123456 , 0x123457 , 0x123458=full;
ShareIO = 0x80 , defmember=full : 0x123459 , 0x12345a;
ShareIO = 0x80 , defmember=full : 0x123459 , 0x12345c=limited , 0x12345d;

# multicast groups added to default
Default=0x7fff,ipoib:
mgid=ff12:401b::16 # MLDv2-capable routers
mgid=ff12:61b::2 # All routers
mgid=ff12:1,sl=1,Q_Key=0xDEADBEEF,rate=3,mtu=2 # random group
ALL=full;
```

The following rule is equivalent to how OpenSM used to run prior to the partition manager:

```
Default=0x7fff,ipoib:ALL=full;
```

Effect of Topology Changes

If a link is added or removed, OpenSM may not recalculate the routes that do not have to change. A route has to change if the port is no longer UP or no longer the MinHop. When routing changes are performed, the same algorithm for balancing the routes is invoked.

In the case of using the file-based routing, any topology changes are currently ignored. The 'file' routing engine just loads the LFTs from the file specified, with no reaction to real topology. Obviously, this will not be able to recheck LIDs (by GUID) for disconnected nodes, and LFTs for non-existent switches will be skipped. Multicast is not affected by 'file' routing engine (this uses min hop tables).

Routing Algorithms

OpenSM offers the following routing engines:

1. Min Hop Algorithm

Based on the minimum hops to each node where the path length is optimized.

2. UPDN Algorithm

Based on the minimum hops to each node, but it is constrained to ranking rules. This algorithm should be chosen if the subnet is not a pure Fat Tree, and a deadlock may occur due to a loop in the subnet.

3. Fat-tree Routing Algorithm

This algorithm optimizes routing for a congestion-free "shift" communication pattern. It should be chosen if a subnet is a symmetrical Fat Tree of various types, not just a K-ary-N-Tree: non-constant K, not fully staffed, and for any CBB ratio. Similar to UPDN, Fat Tree routing is constrained to ranking rules.

4. LASH Routing Algorithm

Uses InfiniBand virtual layers (SL) to provide deadlock-free shortest-path routing while also distributing the paths between layers. LASH is an alternative deadlock-free, topology-agnostic routing algorithm to the non-minimal UPDN algorithm. It avoids the use of a potentially congested root node.

5. DOR Routing Algorithm

Based on the Min Hop algorithm, but avoids port equalization except for redundant links between the same two switches. This provides deadlock free routes for hypercubes when the fabric is cabled as a hypercube and for meshes when cabled as a mesh.

6. Torus-2QoS Routing Algorithm

Based on the DOR Unicast routing algorithm specialized for 2D/3D torus topologies. Torus-2QoS provides deadlock-free routing while supporting two quality of service (QoS) levels. Additionally, it can route around multiple failed fabric links or a single failed fabric switch without introducing deadlocks, and without changing path SL values granted before the failure.

7. Routing Chains

Allows routing configuration of different parts of a single InfiniBand subnet by different routing engines. In the current release, minhop/updn/ftree/dor/torus-2QoS/pqft can be combined.

MINHOP/UPDN/DOR routing algorithms are comprised of two stages:

- 1. MinHop matrix calculation. How many hops are required to get from each port to each LID. The algorithm to fill these tables is different if you run standard (min hop) or Up/Down. For standard routing, a "relaxation" algorithm is used to propagate min hop from every destination LID through neighbor switches. For Up/Down routing, a BFS from every target is used. The BFS tracks link direction (up or down) and avoid steps that will perform up after a down step was used.
- 2. Once MinHop matrices exist, each switch is visited and for each target LID a decision is made as to what port should be used to get to that LID. This step is common to standard and Up/ Down routing. Each port has a counter counting the number of target LIDs going through it. When there are multiple alternative ports with same MinHop to a LID, the one with less previously assigned ports is selected.

If LMC > 0, more checks are added. Within each group of LIDs assigned to same target port:

- a. Use only ports which have same MinHop
- b. First prefer the ones that go to different systemImageGuid (then the previous LID of the same LMC group)
- c. If none, prefer those which go through another NodeGuid
- d. Fall back to the number of paths method (if all go to same node).

Min Hop Algorithm

The Min Hop algorithm is invoked by default if no routing algorithm is specified. It can also be invoked by specifying '-R minhop'.

The Min Hop algorithm is divided into two stages: computation of min-hop tables on every switch

and LFT output port assignment. Link subscription is also equalized with the ability to override based on port GUID. The latter is supplied by:

```
-i <equalize-ignore-guids-file>
-ignore-guids <equalize-ignore-guids-file>
```

This option provides the means to define a set of ports (by GUIDs) that will be ignored by the link load equalization algorithm.

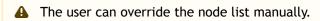
LMC awareness routes based on a (remote) system or on a switch basis.

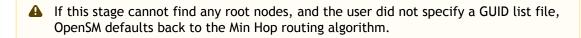
UPDN Algorithm

The UPDN algorithm is designed to prevent deadlocks from occurring in loops of the subnet. A loop-deadlock is a situation in which it is no longer possible to send data between any two hosts connected through the loop. As such, the UPDN routing algorithm should be sent if the subnet is not a pure Fat Tree, and one of its loops may experience a deadlock (due, for example, to high pressure).

The UPDN algorithm is based on the following main stages:

1. Auto-detect root nodes - based on the CA hop length from any switch in the subnet, a statistical histogram is built for each switch (hop num vs the number of occurrences). If the histogram reflects a specific column (higher than others) for a certain node, then it is marked as a root node. Since the algorithm is statistical, it may not find any root nodes. The list of the root nodes found by this auto-detect stage is used by the ranking process stage.





- 2. Ranking process All root switch nodes (found in stage 1) are assigned a rank of 0. Using the BFS algorithm, the rest of the switch nodes in the subnet are ranked incrementally. This ranking aids in the process of enforcing rules that ensure loop-free paths.
- 3. Min Hop Table setting after ranking is done, a BFS algorithm is run from each (CA or switch) node in the subnet. During the BFS process, the FDB table of each switch node traversed by BFS is updated, in reference to the starting node, based on the ranking rules and GUID values.

At the end of the process, the updated FDB tables ensure loop-free paths through the subnet.

UPDN Algorithm Usage

Activation through OpenSM:

- Use '-R updn' option (instead of old '-u') to activate the UPDN algorithm.
- Use '-a <root_guid_file>' for adding an UPDN GUID file that contains the root nodes for ranking. If the `-a' option is not used, OpenSM uses its auto-detect root nodes algorithm.

Notes on the GUID list file:

- A valid GUID file specifies one GUID in each line. Lines with an invalid format will be discarded
- The user should specify the root switch GUIDs

Fat-tree Routing Algorithm

The fat-tree algorithm optimizes routing for "shift" communication pattern. It should be chosen if a subnet is a symmetrical or almost symmetrical fat-tree of various types. It supports not just K- ary-N-Trees, by handling for non-constant K, cases where not all leafs (CAs) are present, any Constant Bisectional Ratio (CBB) ratio. As in UPDN, fat-tree also prevents credit-loop-dead-locks. If the root GUID file is not provided ('a' or '-root_guid_file' options), the topology has to be pure fat-tree that complies with the following rules:

- Tree rank should be between two and eight (inclusively)
- Switches of the same rank should have the same number of UP-going port groups, unless they are root switches, in which case the shouldn't have UP-going ports at all.

 Note: Ports that are connected to the same remote switch are referenced as 'port group'.
- Switches of the same rank should have the same number of DOWN-going port groups, unless they are leaf switches.
- Switches of the same rank should have the same number of ports in each UP-going port group.
- Switches of the same rank should have the same number of ports in each DOWN-going port group.
- All the CAs have to be at the same tree level (rank).

If the root GUID file is provided, the topology does not have to be pure fat-tree, and it should only comply with the following rules:

- Tree rank should be between two and eight (inclusively)
- All the Compute Nodes have to be at the same tree level (rank). Note that non-compute node
 CAs are allowed here to be at different tree ranks.
 Note: List of compute nodes (CNs) can be specified using '-u' or '--cn_guid_file' OpenSM
 options.

Topologies that do not comply cause a fallback to min-hop routing. Note that this can also occur on link failures which cause the topology to no longer be a "pure" fat-tree.

Note that although fat-tree algorithm supports trees with non-integer CBB ratio, the routing will not be as balanced as in case of integer CBB ratio. In addition to this, although the algorithm allows leaf switches to have any number of CAs, the closer the tree is to be fully populated, the more effective the "shift" communication pattern will be. In general, even if the root list is provided, the closer the topology to a pure and symmetrical fat-tree, the more optimal the routing will be.

The algorithm also dumps the compute node ordering file (opensm-ftree-ca-order.dump) in the same directory where the OpenSM log resides. This ordering file provides the CN order that may be used to create efficient communication pattern, that will match the routing tables.

Routing between non-CN Nodes

The use of the io_guid_file option allows non-CN nodes to be located on different levels in the fat tree. In such case, it is not guaranteed that the Fat Tree algorithm will route between two non-CN nodes. In the scheme below, N1, N2, and N3 are non-CN nodes. Although all the CN have routes to and from them, there will not necessarily be a route between N1,N2 and N3. Such routes would require to use at least one of the switches the wrong way around.

? Unknown Attachment

To solve this problem, a list of non-CN nodes can be specified by \'-G\' or \'--io_guid_file\' option. These nodes will be allowed to use switches the wrong way around a specific number of times (specified by \'-H\' or \'--max_reverse_hops\'. With the proper max_reverse_hops and io_guid_file values, you can ensure full connectivity in the Fat Tree. In the scheme above, with a max_reverse_hop of 1, routes will be instantiated between N1<-->N2 and N2<-->N3. With a max_reverse_hops value of 2, N1,N2 and N3 will all have routes between them.



Using max_reverse_hops creates routes that use the switch in a counter-stream way. This option should never be used to connect nodes with high bandwidth traffic between them! It should only be used to allow connectivity for HA purposes or similar. Also having routes the other way around can cause credit loops.

Activation through OpenSM

Use '-R ftree' option to activate the fat-tree algorithm.



LMC > 0 is not supported by fat-tree routing. If this is specified, the default routing algorithm is invoked instead.

LASH Routing Algorithm

LASH is an acronym for LAyered SHortest Path Routing. It is a deterministic shortest path routing algorithm that enables topology agnostic deadlock-free routing within communication networks. When computing the routing function, LASH analyzes the network topology for the shortest-path routes between all pairs of sources/destinations and groups these paths into virtual layers in such a way as to avoid deadlock.



▲ LASH analyzes routes and ensures deadlock freedom between switch pairs. The link from HCA between and switch does not need virtual layers as deadlock will not arise between switch and HCA.

Here is a detailed explanation of how this algorithm works:

- 1. LASH determines the shortest-path between all pairs of source/destination switches. Note, LASH ensures the same SL is used for all SRC/DST - DST/SRC pairs and there is no guarantee that the return path for a given DST/SRC will be the reverse of the route SRC/DST.
- 2. LASH then begins an SL assignment process where a route is assigned to a layer (SL) if the addition of that route does not cause deadlock within that layer. This is achieved by maintaining and analyzing a channel dependency graph for each layer. Once the potential addition of a path could lead to deadlock, LASH opens a new layer and continues the process.
- 3. Once this stage has been completed, it is highly likely that the first layers processed will contain more paths than the latter ones. To better balance the use of layers, LASH moves paths from one layer to another so that the number of paths in each layer averages out.

Note that the implementation of LASH in opensm attempts to use as few layers as possible. This number can be less than the number of actual layers available.

In general, LASH is a very flexible algorithm. It can, for example, reduce to Dimension Order Routing in certain topologies, it is topology agnostic and fares well in the face of faults.

It has been shown that for both regular and irregular topologies, LASH outperforms Up/Down. The reason for this is that LASH distributes the traffic more evenly through a network, avoiding the bottleneck issues related to a root node and always routes shortest-path.

The algorithm was developed by Simula Research Laboratory. Use '-R lash -Q' option to activate the LASH algorithm



A QoS support has to be turned on in order that SL/VL mappings are used.

▲ LMC > 0 is not supported by the LASH routing. If this is specified, the default routing algorithm is invoked instead.

For open regular cartesian meshes, the DOR algorithm is the ideal routing algorithm. For toroidal meshes, on the other hand, there are routing loops that can cause deadlocks. LASH can be used to route these cases. The performance of LASH can be improved by preconditioning the mesh in cases where there are multiple links connecting switches and also in cases where the switches are not cabled consistently. To invoke this, use '-R lash -Q --do_mesh_analysis'. This will add an additional phase that analyses the mesh to try to determine the dimension and size of a mesh. If it determines that the mesh looks like an open or closed cartesian mesh it reorders the ports in dimension order before the rest of the LASH algorithm runs.

DOR Routing Algorithm

The Dimension Order Routing algorithm is based on the Min Hop algorithm and so uses shortest paths. Instead of spreading traffic out across different paths with the same shortest distance, it chooses among the available shortest paths based on an ordering of dimensions. Each port must be consistently cabled to represent a hypercube dimension or a mesh dimension. Paths are grown from a destination back to a source using the lowest dimension (port) of available paths at each step. This provides the ordering necessary to avoid deadlock. When there are multiple links between any two switches, they still represent only one dimension and traffic is balanced across them unless port equalization is turned off. In the case of hypercubes, the same port must be used throughout the fabric to represent the hypercube dimension and match on both ends of the cable. In the case of meshes, the dimension should consistently use the same pair of ports, one port on one end of the cable, and the other port on the other end, continuing along the mesh dimension. Use '-R dor' option to activate the DOR algorithm.

Torus-2QoS Routing Algorithm

Torus-2QoS is a routing algorithm designed for large-scale 2D/3D torus fabrics. The torus-2QoS routing engine can provide the following functionality on a 2D/3D torus:

- Free of credit loops routing
- Two levels of QoS, assuming switches support 8 data VLs
- Ability to route around a single failed switch, and/or multiple failed links, without:
 - introducing credit loops
 - · changing path SL values
- Very short run times, with good scaling properties as fabric size increases

Unicast Routing

Torus-2 QoS is a DOR-based algorithm that avoids deadlocks that would otherwise occur in a torus using the concept of a dateline for each torus dimension. It encodes into a path SL which datelines the path crosses as follows:

```
s1 = 0;
for (d = 0; d < torus_dimensions; d++)</pre>
/* path_crosses_dateline(d) returns 0 or 1 */
sl |= path_crosses_dateline(d) << d;
```

For a 3D torus, that leaves one SL bit free, which torus-2 QoS uses to implement two QoS levels. Torus-2 QoS also makes use of the output port dependence of switch SL2VL maps to encode into one VL bit the information encoded in three SL bits. It computes in which torus coordinate direction each inter-switch link "points", and writes SL2VL maps for such ports as follows:

```
for (sl = 0; sl < 16; sl ++)
/* cdir(port) reports which torus coordinate direction a switch port
* "points" in, and returns 0, 1, or 2 */
sl2vl(iport,oport,sl) = 0xl & (sl >> cdir(oport));
```

Thus, on a pristine 3D torus, i.e., in the absence of failed fabric switches, torus-2 QoS consumes 8 SL values (SL bits 0-2) and 2 VL values (VL bit 0) per QoS level to provide deadlock-free routing on a 3D torus. Torus-2 QoS routes around link failure by "taking the long way around" any 1D ring interrupted by a link failure. For example, consider the 2D 6x5 torus below, where switches are denoted by [+a-zA-Z]:

? Unknown Attachment

For a pristine fabric the path from S to D would be S-n-T-r-D. In the event that either link S-n or n-T has failed, torus-2QoS would use the path S-m-p-o-T-r-D.

Note that it can do this without changing the path SL value; once the 1D ring m-S-n-T-o-p-m has been broken by failure, path segments using it cannot contribute to deadlock, and the x-direction dateline (between, say, x=5 and x=0) can be ignored for path segments on that ring. One result of this is that torus-2QoS can route around many simultaneous link failures, as long as no 1D ring is broken into disjoint segments. For example, if links n-T and T-o have both failed, that ring has been broken into two disjoint segments, T and o-p-m-S-n. Torus-2QoS checks for such issues, reports if they are found, and refuses to route such fabrics.

Note that in the case where there are multiple parallel links between a pair of switches, torus-2QoS will allocate routes across such links in a round-robin fashion, based on ports at the path destination switch that are active and not used for inter-switch links. Should a link that is one of severalsuch parallel links fail, routes are redistributed across the remaining links. When the last of such a set of parallel links fails, traffic is rerouted as described above.

Handling a failed switch under DOR requires introducing into a path at least one turn that would be otherwise "illegal", i.e. not allowed by DOR rules. Torus-2QoS will introduce such a turn as close as possible to the failed switch in order to route around it. n the above example, suppose switch T has failed, and consider the path from S to D. Torus-2QoS will produce the path S-n-I-r-D, rather than the S-n-T-r-D path for a pristine torus, by introducing an early turn at n. Normal DOR rules will cause traffic arriving at switch I to be forwarded to switch r; for traffic arriving from I due to the "early" turn at n, this will generate an "illegal" turn at I.

Torus-2QoS will also use the input port dependence of SL2VL maps to set VL bit 1 (which would be otherwise unused) for y-x, z-x, and z-y turns, i.e., those turns that are illegal under DOR. This causes the first hop after any such turn to use a separate set of VL values, and prevents deadlock in the presence of a single failed switch. For any given path, only the hops after a turn that is illegal under DOR can contribute to a credit loop that leads to deadlock. So in the example above with failed switch T, the location of the illegal turn at I in the path from S to D requires that any credit loop caused by that turn must encircle the failed switch at T. Thus the second and later hops after the illegal turn at I (i.e., hop r-D) cannot contribute to a credit loop because they cannot be used to construct a loop encircling T. The hop I-r uses a separate VL, so it cannot contribute to a credit loop encircling T. Extending this argument shows that in addition to being capable of routing around a single switch failure without introducing deadlock, torus-2QoS can also route around multiple failed switches on the condition they are adjacent in the last dimension routed by DOR. For example, consider the following case on a 6x6 2D torus:

? Unknown Attachment

Suppose switches T and R have failed, and consider the path from S to D. Torus-2QoS will generate the path S-n-q-I-u-D, with an illegal turn at switch I, and with hop I-u using a VL with bit 1 set. As a

further example, consider a case that torus-2QoS cannot route without deadlock: two failed switches adjacent in a dimension that is not the last dimension routed by DOR; here the failed switches are O and T:

? Unknown Attachment

In a pristine fabric, torus-2QoS would generate the path from S to D as S-n-O-T-r-D. With failed switches O and T, torus-2QoS will generate the path S-n-I-q-r-D, with an illegal turn at switch I, and with hop I-q using a VL with bit 1 set. In contrast to the earlier examples, the second hop after the illegal turn, q-r, can be used to construct a credit loop encircling the failed switches.

Multicast Routing

Since torus-2QoS uses all four available SL bits, and the three data VL bits that are typically available in current switches, there is no way to use SL/VL values to separate multicast traffic from unicast traffic. Thus, torus-2QoS must generate multicast routing such that credit loops cannot arise from a combination of multicast and unicast path segments. It turns out that it is possible to construct spanning trees for multicast routing that have that property. For the 2D 6x5 torus example above, here is the full-fabric spanning tree that torus-2QoS will construct, where "x" is the root switch and each "+" is a non-root switch:

? Unknown Attachment

For multicast traffic routed from root to tip, every turn in the above spanning tree is a legal DOR turn. For traffic routed from tip to root, and some traffic routed through the root, turns are not legal DOR turns. However, to construct a credit loop, the union of multicast routing on this spanning tree with DOR unicast routing can only provide 3 of the 4 turns needed for the loop. In addition, if none of the above spanning tree branches crosses a dateline used for unicast credit loop avoidance on a torus, and if multicast traffic is confined to SL 0 or SL 8 (recall that torus-2QoS uses SL bit 3 to differentiate QoS level), then multicast traffic also cannot contribute to the "ring" credit loops that are otherwise possible in a torus. Torus-2QoS uses these ideas to create a master spanning tree. Every multicast group spanning tree will be constructed as a subset of the master tree, with the same root as the master tree. Such multicast group spanning trees will in general not be optimal for groups which are a subset of the full fabric. However, this compromise must be made to enable support for two QoS levels on a torus while preventing credit loops. In the presence of link or switch failures that result in a fabric for which torus-2QoS can generate credit-loop-free unicast routes, it is also possible to generate a master spanning tree for multicast that retains the required properties. For example, consider that same 2D 6x5 torus, with the link from (2,2) to (3,2) failed. Torus-2QoS will generate the following master spanning tree:

? Unknown Attachment

Two things are notable about this master spanning tree. First, assuming the x dateline was between x=5 and x=0, this spanning tree has a branch that crosses the dateline. However, just as for unicast, crossing a dateline on a 1D ring (here, the ring for y=2) that is broken by a failure cannot contribute to a torus credit loop. Second, this spanning tree is no longer optimal even for multicast groups that encompass the entire fabric. That, unfortunately, is a compromise that must be made to retain the other desirable properties of torus-2QoS routing. In the event that a single switch fails, torus-2QoS will generate a master spanning tree that has no "extra" turns by appropriately selecting a root switch. In the 2D 6x5 torus example, assume now that the switch at (3,2), i.e. the root for a pristine fabric, fails. Torus-2QoS will generate the following master spanning tree for that case:

? Unknown Attachment

Assuming the dateline was between y=4 and y=0, this spanning tree has a branch that crosses a dateline. However, this cannot contribute to credit loops as it occurs on a 1D ring (the ring for x=3) that is broken by failure, as in the above example.

Torus Topology Discovery

The algorithm used by torus-2QoS to construct the torus topology from the undirected graph representing the fabric requires that the radix of each dimension be configured via torus-2QoS.conf. It also requires that the torus topology be "seeded"; for a 3D torus this requires configuring four switches that define the three coordinate directions of the torus. Given this starting information, the algorithm is to examine the cube formed by the eight switch locations bounded by the corners (x,y,z) and (x+1,y+1,z+1). Based on switches already placed into the torus topology at some of these locations, the algorithm examines 4-loops of inter-switch links to find the one that is consistent with a face of the cube of switch locations and adds its switches to the discovered topology in the correct locations.

Because the algorithm is based on examining the topology of 4-loops of links, a torus with one or more radix-4 dimensions requires extra initial seed configuration. See torus-2QoS.conf(5) for details. Torus-2QoS will detect and report when it has an insufficient configuration for a torus with radix-4 dimensions.

In the event the torus is significantly degraded, i.e., there are many missing switches or links, it may happen that torus-2QoS is unable to place into the torus some switches and/or links that were discovered in the fabric, and will generate a warning in that case. A similar condition occurs if torus-2QoS is misconfigured, i.e., the radix of a torus dimension as configured does not match the radix of that torus dimension as wired, and many switches/links in the fabric will not be placed into the torus.

Quality Of Service Configuration

OpenSM will not program switches and channel adapters with SL2VL maps or VL arbitration configuration unless it is invoked with -Q. Since torus-2QoS depends on such functionality for correct operation, always invoke OpenSM with -Q when torus-2QoS is in the list of routing engines. Any quality of service configuration method supported by OpenSM will work with torus-2QoS, subject to the following limitations and considerations. For all routing engines supported by OpenSM except torus-2QoS, there is a one-to-one correspondence between QoS level and SL. Torus-2QoS can only support two quality of service levels, so only the high-order bit of any SL value used for unicast QoS configuration will be honored by torus-2QoS. For multicast QoS configuration, only SL values 0 and 8 should be used with torus-2QoS.

Since SL to VL map configuration must be under the complete control of torus-2QoS, any configuration via qos_sl2vl, qos_swe_sl2vl, etc., must and will be ignored, and a warning will be generated. Torus-2QoS uses VL values 0-3 to implement one of its supported QoS levels, and VL values 4-7 to implement the other. Hard-to-diagnose application issues may arise if traffic is not delivered fairly across each of these two VL ranges. Torus-2QoS will detect and warn if VL arbitration is configured unfairly across VLs in the range 0-3, and also in the range 4-7. Note that the default OpenSM VL arbitration configuration does not meet this constraint, so all torus-2QoS users should configure VL arbitration via qos_vlarb_high, qos_vlarb_low, etc.

Operational Considerations

Any routing algorithm for a torus IB fabric must employ path SL values to avoid credit loops. As a result, all applications run over such fabrics must perform a path record query to obtain the correct path SL for connection setup. Applications that use rdma_cm for connection setup will automatically meet this requirement.

If a change in fabric topology causes changes in path SL values required to route without credit

loops, in general, all applications would need to repath to avoid message deadlock. Since torus-2QoS has the ability to reroute after a single switch failure without changing path SL values, repathing by running applications is not required when the fabric is routed with torus-2QoS. Torus-2QoS can provide unchanging path SL values in the presence of subnet manager failover provided that all OpenSM instances have the same idea of dateline location. See torus- 2QoS.conf(5) for details. Torus-2QoS will detect configurations of failed switches and links that prevent routing that is free of credit loops and will log warnings and refuse to route. If "no fall-back" was configured in the list of OpenSM routing engines, then no other routing engine will attempt to route the fabric. In that case, all paths that do not transit the failed components will continue to work, and the subset of paths that are still operational will continue to remain free of credit loops. OpenSM will continue to attempt to route the fabric after every sweep interval and after any change (such as a link up) in the fabric topology. When the fabric components are repaired, full functionality will be restored. In the event OpenSM was configured to allow some other engine to route the fabric if torus-2QoS fails, then credit loops and message deadlock are likely if torus-2QoS had previously routed the fabric successfully. Even if the other engine is capable of routing a torus without credit loops, applications that built connections with path SL values granted under torus-2QoS will likely experience message deadlock under routing generated by a different engine, unless they repath. To verify that a torus fabric is routed free of credit loops, use ibdmchk to analyze data collected via ibdiagnet -vlr.

Torus-2QoS Configuration File Syntax

The file torus-2QoS.conf contains configuration information that is specific to the OpenSM routing engine torus-2QoS. Blank lines and lines where the first non-whitespace character is "#" are ignored. A token is any contiguous group of non-whitespace characters. Any tokens on a line following the recognized configuration tokens described below are ignored.

```
[torus|mesh] x_radix[m|M|t|T] y_radix[m|M|t|T] z_radix[m|M|t|T]
```

Either torus or mesh must be the first keyword in the configuration and sets the topology that torus-2QoS will try to construct. A 2D topology can be configured by specifying one of x_radix, y_radix, or z_radix as 1. An individual dimension can be configured as mesh (open) or torus (looped) by suffixing its radix specification with one of m, M, t, or T. Thus, "mesh 3T 4 5" and "torus 3 4M 5M" both specify the same topology.

Note that although torus-2QoS can route mesh fabrics, its ability to route around failed components is severely compromised on such fabrics. A failed fabric components very likely to cause a disjoint ring; see UNICAST ROUTING in torus-2QoS(8).

```
xp_link sw0_GUID sw1_GUID
yp_link sw0_GUID sw1_GUID
zp_link sw0_GUID sw1_GUID
xm_link sw0_GUID sw1_GUID
ym_link sw0_GUID sw1_GUID
zm_link sw0_GUID sw1_GUID
```

These keywords are used to seed the torus/mesh topology. For example, "xp_link 0x2000 0x2001" specifies that a link from the switch with node GUID 0x2000 to the switch with node GUID 0x2001 would point in the positive x direction, while "xm_link 0x2000 0x2001" specifies that a link from the switch with node GUID 0x2000 to the switch with node GUID 0x2001 would point in the negative x direction. All the link keywords for a given seed must specify the same "from" switch. In general, it is not necessary to configure both the positive and negative directions for a given coordinate; either is sufficient. However, the algorithm used for topology discovery needs extra information for torus dimensions of radix four (see TOPOLOGY DISCOVERY in torus-2QoS(8)). For such cases, both the positive and negative coordinate directions must be specified. Based on the topology specified via the torus/mesh keyword, torus-2QoS will detect and log when it has insufficient seed configuration.

```
GUIDx_dateline position
y_dateline position
z_dateline position
```

In order for torus-2QoS to provide the guarantee that path SL values do not change under any conditions for which it can still route the fabric, its idea of dateline position must not change relative to physical switch locations. The dateline keywords provide the means to configure such behavior.

The dateline for a torus dimension is always between the switch with coordinate 0 and the switch with coordinate radix-1 for that dimension. By default, the common switch in a torus seed is taken as the origin of the coordinate system used to describe switch location. The position parameter for a dateline keyword moves the origin (and hence the dateline) the specified amount relative to the common switch in a torus seed.

```
next_seed
```

If any of the switches used to specify a seed were to fail torus-2QoS would be unable to complete topology discovery successfully. The next_seed keyword specifies that the following link and dateline keywords apply to a new seed specification.

For maximum resiliency, no seed specification should share a switch with any other seed specification. Multiple seed specifications should use dateline configuration to ensure that torus-2QoS can grant path SL values that are constant, regardless of which seed was used to initiate topology discovery.

portgroup_max_ports max_ports - This keyword specifies the maximum number of parallel interswitch links, and also the maximum number of host ports per switch, that torus-2QoS can accommodate. The default value is 16. Torus-2QoS will log an error message during topology discovery if this parameter needs to be increased. If this keyword appears multiple times, the last instance prevails.

port_order p1 p2 p3 ... - This keyword specifies the order in which CA ports on a destination switch are visited when computing routes. When the fabric contains switches connected with multiple parallel links, routes are distributed in a round-robin fashion across such links, and so changing the order that CA ports are visited changes the distribution of routes across such links. This may be advantageous for some specific traffic patterns.

The default is to visit CA ports in increasing port order on destination switches. Duplicate values in the list will be ignored.

Example:

```
# Look for a 2D (since x radix is one) 4x5 torus.
torus 1 4 5
# y is radix-4 torus dimension, need both
# ym_link and yp_link configuration.
yp_link 0x200000 0x200005 # sw @ y=0,z=0 -> sw @ y=1,z=0
ym_link 0x200000 0x200005 # sw @ y=0,z=0 -> sw @ y=3,z=0
# z is not radix-4 torus dimension, only need one of
# zm_link or zp_link configuration.
zp_link 0x200000 0x200001 # sw @ y=0,z=0 -> sw @ y=0,z=1
next_seed
yp_link 0x200000 0x200010 # sw @ y=2,z=1 -> sw @ y=3,z=1
ym_link 0x200000 0x200010 # sw @ y=2,z=1 -> sw @ y=1,z=1
zp_link 0x200000 0x200006 # sw @ y=2,z=1 -> sw @ y=2,z=2
y_dateline -2 # Move the dateline for this seed
z_dateline -1 # back to its original position.
# If OpenSM failover is configured, for maximum resiliency
# one instance should run on a host attached to a switch
# from the first seed, and another instance should run
# on a host attached to a switch from the second seed.
# Both instances should use this torus-200S.conf to ensure
# path SL values do not change in the event of SM failover.
# port_order defines the order on which the ports would be
# chosen for routing.
port_order 7 10 8 11 9 12 25 28 26 29 27 30
```

Routing Chains

The routing chains feature is offering a solution that enables one to configure different parts of the fabric and define a different routing engine to route each of them. The routings are done in a sequence (hence the name "chains") and any node in the fabric that is configured in more than one part is left with the routing updated by the last routing engine it was a part of.

Configuring Routing Chains



To configure routing chains:

- 1. Define the port groups.
- 2. Define topologies based on previously defined port groups.
- 3. Define configuration files for each routing engine.
- 4. Define routing engine chains over previously defined topologies and configuration files.

Defining Port Groups

The basic idea behind the port groups is the ability to divide the fabric into sub-groups and give each group an identifier that can be used to relate to all nodes in this group. The port groups is a separate feature from the routing chains but is a mandatory prerequisite for it. In addition, it is used to define the participants in each of the routing algorithms.

Defining a Port Group Policy File

In order to define a port group policy file, set the parameter 'pgrp_policy_file' in the opensm configuration file.

pgrp_policy_file /etc/opensm/conf/port_groups_policy_file

Configuring a Port Group Policy

The port groups policy file details the port groups in the fabric. The policy file should be composed of one or more paragraphs that define a group. Each paragraph should begin with the line 'port-group' and end with the line 'end-port-group'. For example:

port-group
...port group qualifiers...
end-port-group

Port Group Qualifiers



Unlike the port group's beginning and end which do not require a colon, all qualifiers must end with a colon (':'). Also - a colon is a predefined mark that must not be used inside qualifier values. The inclusion of a colon in the name or the use of a port group will result in the policy's failure.

Rule Qualifier

Parameter	Description	Example
name	Each group must have a name. Without a name qualifier, the policy fails.	name: grp1
use	'use' is an optional qualifier that one can define in order to describe the usage of this port group (if undefined, an empty string is used as a default).	use: first port group

There are several qualifiers used to describe a rule that determines which ports will be added to the group. Each port group may include one or more rules out of the rules described in the below table (at least one rule must be defined for each port group).

Param eter	Description	Example
guid list	Comma separated list of GUIDs to include in the group. If no specific physical ports were configured, all physical ports of the guid are chosen. However, for each guid, one can detail specific physical ports to be included in the group. This can be done using the following syntax:	port-guid: 0x283, 0x286, 0x289
	 Specify a specific port in a guid to be chosen port-guid: 0x283@3 Specify a specific list of ports in a guid to be chosen port-guid: 0x286@1/5/7 Specify a specific range of ports in a guid to be chosen port-guid: 0x289@2-5 Specify a list of specific ports and ports ranges in a guid to be chosen port-guid: 0x289@2-5/7/9-13/18 Complex rule port-guid: 0x283@5-8/12/14, 0x286, 0x289/6/ 8/12 	
port guid range	It is possible to configure a range of guids to be chosen to the group. However, while using the range qualifier, it is impossible to detail specific physical ports. Note: A list of ranges cannot be specified. The below example is invalid and will cause the policy to fail: port-guid-range: 0x283-0x289, 0x290- 0x295	port-guid- range: 0x283-0x289
port name	One can configure a list of hostnames as a rule. Hosts with a node description that is built out of these hostnames will be chosen. Since the node description contains the network card index as well, one might also specify a network card index and a physical port to be chosen. For example, the given configuration will cause only physical port 2 of a host with the node description 'kuku HCA-1' to be chosen. port and hca_idx parameters are optional. If the port is unspecified, all physical ports are chosen. If hca_idx is unspecified, all card numbers are chosen. Specifying a hostname is mandatory.	port-name: host- name=kuku; port=2; hca_idx=1
	One can configure a list of hostname/ port/hca_idx sets in the same qualifier as follows:	
	port-name: hostname=kuku; port=2; hca_idx=1 , hostname=host1; port=3, hostname=host2	
	Note: port-name qualifier is not relevant for switches, but for HCA's only.	

Param eter	Description	Example
port regexp	One can define a regular expression so that only nodes with a matching node description will be chosen to the group.	port-regexp: SW
	Note : This example shows how to choose nodes which their node description starts with 'SW'.	
	It is possible to specify one physical port to be chosen for matching nodes (there is no option to define a list or a range of ports). The given example will cause only nodes that match physical port 3 to be added to the group.	port-regexp: SW:3
union rule	It is possible to define a rule that unites two different port groups. This means that all ports from both groups will be included in the united group.	union-rule: grp1, grp2
subtrac t rule	One can define a rule that subtracts one port group from another. The given rule, for example, will cause all the ports which are a part of grp1, but not included in grp2, to be chosen. In subtraction (unlike union), the order does matter, since the purpose is to subtract the second group from the first one. There is no option to define more than two groups for union/subtraction. However, one can unite/subtract groups which are a union or a subtraction themselves, as shown in the port groups policy file example.	subtract-rule: grp1, grp2

Predefined Port Groups

There are 3 predefined, automatically created port groups that are available for use, yet cannot be defined in the policy file (if a group in the policy is configured with the name of one of these predefined groups, the policy fails) -

- ALL a group that includes all nodes in the fabric
- ALL_SWITCHES a group that includes all switches in the fabric
- ALL_CAS a group that includes all HCAs in the fabric
- ALL_ROUTERS a group that includes all routers in the fabric (supported in OpenSM starting from v4.9.0)

Port Groups Policy Examples

```
port-group
name: grp3
use: Subtract of groups grp1 and grp2
subtract-rule: grp1, grp2
end-port-group

port-group
name: grp1
port-grid: 0x281, 0x282, 0x283
end-port-group

port-group
name: grp2
port-group
name: grp2
port-group
port-group
port-group
port-group
port-group
port-group
port-group
port-group
name: grp4
port-name: hostname=kika port=1 hca_idx=1
end-port-group
port-group
name: grp3
union-rule: grp3, grp4
end-port-group
```

Defining a Topologies Policy File

In order to define a topology policy file, set the parameter 'topo_policy_file' in the opensm configuration file.

topo_policy_file /etc/opensm/conf/topo_policy_file.cfg

Configuring a Topology Policy

The topologies policy file details a list of topologies. The policy file should be composed of one or more paragraphs which define a topology. Each paragraph should begin with the line 'topol- ogy' and end with the line 'end-topology'.

For example:

topology ...topology qualifiers... end-topology

Topology Qualifiers

A

Unlike topology and end-topology which do not require a colon, all qualifiers must end with a colon (':'). Also - a colon is a predefined mark that must not be used inside qualifier values. An inclusion of a column in the qualifier values will result in the policy's failure.

All topology qualifiers are mandatory. Absence of any of the below qualifiers will cause the policy parsing to fail.

Topology Qualifiers

Parameter	Description	Example
id	Topology ID. Legal Values - any positive value. Must be unique.	id: 1
sw-grp	Name of the port group that includes all switches and switch ports to be used in this topology.	sw-grp: ys_switches
hca-grp	Name of the port group that includes all HCA's to be used in this topology.	hca-grp: ys_hosts

Configuration File per Routing Engine

Each engine in the routing chain can be provided by its own configuration file. Routing engine configuration file is the fraction of parameters defined in the main opensm configuration file. Some rules should be applied when defining a particular configuration file for a routing engine:

- Parameters that are not specified in specific routing engine configuration file are inherited from the main opensm configuration file.
- The following configuration parameters are taking effect only in the main opensm configuration file:
 - qos and qos_* settings like (vl_arb, sl2vl, etc.)
 - lmc

· routing_engine

Defining a Routing Chain Policy File

In order to define a port group policy file, set the parameter 'rch_policy_file' in the opensm configuration file.

rch_policy_file /etc/opensm/conf/chains_policy_file

First Routing Engine in the Chain

The first unicast engine in a routing chain must include all switches and HCAs in the fabric (topology id must be 0). The path-bit parameter value is path-bit 0 and it cannot be changed.

Configuring a Routing Chains Policy

The routing chains policy file details the routing engines (and their fallback engines) used for the fabric's routing. The policy file should be composed of one or more paragraphs which defines an engine (or a fallback engine). Each paragraph should begin with the line 'unicast-step' and end with the line 'end-unicast-step'.

For example:

unicast-step ...routing engine qualifiers... end-unicast-step

Routing Engine Qualifiers

A

Unlike unicast-step and end-unicast-step which do not require a colon, all qualifiers must end with a colon (':'). Also - a colon is a predefined mark that must not be used inside qualifier values. An inclusion of a colon in the qualifier values will result in the policy's failure.

Para mete r	Description	Example
id	 'id' is mandatory. Without an ID qualifier for each engine, the policy fails. Legal values - size_t value (0 is illegal). The engines in the policy chain are set according to an ascending id order, so it is highly crucial to verify that the id that is given to the engines match the order in which you would like the engines to be set. 	is: 1
engine	This is a mandatory qualifier that describes the routing algorithm used within this unicast step. Currently, on the first phase of routing chains, legal values are minhop/ftree/updn.	engine: minhop
use	This is an optional qualifier that enables one to describe the usage of this unicast step. If undefined, an empty string is used as a default.	use: ftree routing for for yellow stone nodes

Para mete r	Description	Example
config	This is an optional qualifier that enables one to define a separate opensm config file for a specific unicast step. If undefined, all parameters are taken from main opensm configuration file.	<pre>config: / etc/config/ opensm2.cfg</pre>
topolo gy	 Define the topology that this engine uses. Legal value - id of an existing topology that is defined in topologies policy (or zero that represents the entire fabric and not a specific topology). Default value - If unspecified, a routing engine will relate to the entire fabric (as if topology zero was defined). Notice: The first routing engine (the engine with the lowest id) MUST be configured with topology: 0 (entire fabric) or else, the routing chain parser will fail. 	topology: 1
fallba ck-to	 This is an optional qualifier that enables one to define the current unicast step as a fallback to another unicast step. This can be done by defining the id of the unicast step that this step is a fallback to. If undefined, the current unicast step is not a fallback. If the value of this qualifier is a non-existent engine id, this step will be ignored. A fallback step is meaningless if the step it is a fallback to did not fail. It is impossible to define a fallback to a fall- back step (such definition will be ignored) 	_
path- bit	This is an optional qualifier that enables one to define a specific lid offset to be used by the current unicast step. Setting lmc > 0 in main opensm configuration file is a prerequisite for assigning specific path-bit for the routing engine. Default value is 0 (if path-bit is not specified)	Path-bit: 1

Dump Files per Routing Engine

Each routing engine on the chain will dump its own data files if the appropriate log_flags is set (for instance 0x43).

The files that are dumped by each engine are:

- opensm-lid-matrix.dump
- opensm-lfts.dump
- opensm.fdbs
- opensm-subnet.lst

These files should contain the relevant data for each engine topology.



▲ sl2vl and mcfdbs files are dumped only once for the entire fabric and NOT by every routing engine.

- Each engine concatenates its ID and routing algorithm name in its dump files names, as follows:
 - opensm-lid-matrix.2.minhop.dump
 - opensm.fdbs.3.ftree
 - opensm-subnet.4.updn.lst

- In case that a fallback routing engine is used, both the routing engine that failed and the fallback engine that replaces it, dump their data.
 - If, for example, engine 2 runs ftree and it has a fallback engine with 3 as its id that runs minhop, one should expect to find 2 sets of dump files, one for each engine:
 - opensm-lid-matrix.2.ftree.dump
 - opensm-lid-matrix.3.minhop.dump
 - opensm.fdbs.2.ftree
 - opensm.fdbs.3.munhop

Unicast Routing Cache

Unicast routing cache prevents routing recalculation (which is a heavy task in a large cluster) when no topology change was detected during the heavy sweep, or when the topology change does not require new routing calculation (for example, when one or more CAs/RTRs/leaf switches going down, or one or more of these nodes coming back after being down).

Quality of Service Management in OpenSM

When Quality of Service (QoS) in OpenSM is enabled (using the '-Q' or '--qos' flags), OpenSM looks for a QoS Policy file. During fabric initialization and at every heavy sweep, OpenSM parses the QoS policy file, applies its settings to the discovered fabric elements, and enforces the provided policy on client requests. The overall flow for such requests is as follows:

- The request is matched against the defined matching rules such that the QoS Level definition is found
- Given the QoS Level, a path(s) search is performed with the given restrictions imposed by that level

? Unknown Attachment

There are two ways to define QoS policy:

- Advanced the advanced policy file syntax provides the administrator various ways to match a PathRecord/MultiPathRecord (PR/MPR) request, and to enforce various QoS constraints on the requested PR/MPR
- Simple the simple policy file syntax enables the administrator to match PR/MPR requests by various ULPs and applications running on top of these ULPs

Advanced QoS Policy File

The QoS policy file has the following sections:

- 1. Port Groups (denoted by port-groups) this section defines zero or more port groups that can be referred later by matching rules (see below). Port group lists ports by:
 - Port GUID
 - Port name, which is a combination of NodeDescription and IB port number
 - PKey, which means that all the ports in the subnet that belong to partition with a given PKey belong to this port group
 - Partition name, which means that all the ports in the subnet that belong to partition with a given name belong to this port group
 - Node type, where possible node types are: CA, SWITCH, ROUTER, ALL, and SELF (SM's port).

- QoS Setup (denoted by qos-setup) this section describes how to set up SL2VL and VL
 Arbitration tables on various nodes in the fabric. However, this is not supported in OFED.
 SL2VL and VLArb tables should be configured in the OpenSM options file (default location /
 var/cache/opensm/opensm.opts).
- QoS Levels (denoted by qos-levels) each QoS Level defines Service Level (SL) and a few optional fields:
 - MTU limit
 - Rate limit
 - PKev
 - Packet lifetime

When path(s) search is performed, it is done with regards to restriction that these QoS Level parameters impose. One QoS level that is mandatory to define is a DEFAULT QoS level. It is applied to a PR/MPR query that does not match any existing match rule. Similar to any other QoS Level, it can also be explicitly referred by any match rule.

- QoS Matching Rules (denoted by qos-match-rules) each PathRecord/MultiPathRecord query that OpenSM receives is matched against the set of matching rules. Rules are scanned in order of appearance in the QoS policy file such as the first match takes precedence.
 Each rule has a name of QoS level that will be applied to the matching query. A default QoS level is applied to a query that did not match any rule.
 Queries can be matched by:
 - Source port group (whether a source port is a member of a specified group)
 - Destination port group (same as above, only for destination port)
 - PKev
 - QoS class
 - Service ID

To match a certain matching rule, PR/MPR query has to match ALL the rule's criteria. However, not all the fields of the PR/MPR query have to appear in the matching rule. For instance, if the rule has a single criterion - Service ID, it will match any query that has this Service ID, disregarding rest of the query fields. However, if a certain query has only Service ID (which means that this is the only bit in the PR/MPR component mask that is on), it will not match any rule that has other matching criteria besides Service ID.

Simple QoS Policy Definition

Simple QoS policy definition comprises of a single section denoted by qos-ulps. Similar to the advanced QoS policy, it has a list of match rules and their QoS Level, but in this case a match rule has only one criterion - its goal is to match a certain ULP (or a certain application on top of this ULP) PR/MPR request, and QoS Level has only one constraint - Service Level (SL).

The simple policy section may appear in the policy file in combine with the advanced policy or as

The simple policy section may appear in the policy file in combine with the advanced policy, or as a stand-alone policy definition. See more details and list of match rule criteria below.

Policy File Syntax Guidelines

- Leading and trailing blanks, as well as empty lines, are ignored, so the indentation in the example is just for better readability.
- Comments are started with the pound sign (#) and terminated by EOL.
- Any keyword should be the first non-blank in the line, unless it's a comment.
- Keywords that denote section/subsection start have matching closing keywords.

- Having a QoS Level named "DEFAULT" is a must it is applied to PR/MPR requests that did not match any of the matching rules.
- Any section/subsection of the policy file is optional.

Examples of Advanced Policy Files

As mentioned earlier, any section of the policy file is optional, and the only mandatory part of the policy file is a default QoS Level.

Here is an example of the shortest policy file:

```
qos-levels
qos-level
name: DEFAULT
sl: 0
end-qos-level
end-qos-level
```

Port groups section is missing because there are no match rules, which means that port groups are not referred anywhere, and there is no need defining them. And since this policy file doesn't have any matching rules, PR/MPR query will not match any rule, and OpenSM will enforce default QoS level. Essentially, the above example is equivalent to not having a QoS policy file at all. The following example shows all the possible options and keywords in the policy file and their syntax:

```
See the comments in the following example. They explain different keywords and their meaning.
port-groups
        port-group # using port GUIDs
   name: Storage
   # "use" is just a description that is used for logging
   # Other than that, it is just a comment
   use: SRP Targets
   port-guid: 0x1000000000001, 0x100000000005-0x1000000000FFFA
   port-guid: 0x100000000FFFF
end-port-group
        port-group
  name: Virtual Servers
  # The syntax of the port name is as follows:
  # "node_description/Pnum".
  # node_description is compared to the NodeDescription of the node,
  # and "Pnum" is a port number on that node.
  port-name: "vsl HCA-1/Pl, vs2 HCA-1/Pl"
         # using partitions defined in the partition policy
port-group
   name: Partitions
   partition: Part1
   pkey: 0x1234
end-port-group
          # using node types: CA, ROUTER, SWITCH, SELF (for node that runs SM)
# or ALL (for all the nodes in the subnet)
         port-group
name: CAs and SM
node-type: CA, SELF
end-port-group
end-port-groups
qos-setup
# This section of the policy file describes how to set up SL2VL and VL
# Arbitration tables on various nodes in the fabric.
# However, this is not supported in OFED - the section is parsed
# and ignored. SL2VL and VLArb tables should be configured in the
# OpenSM options file (by default - /var/cache/opensm/opensm.opts).
end-qos-setup
gos-levels
         \# Having a QoS Level named "DEFAULT" is a must - it is applied to \# PR/MPR requests that didn't match any of the matching rules. qos-level _____
                   name: DEFAULT
use: default QoS Level
         sl: 0
end-qos-level
           # the whole set: SL, MTU-Limit, Rate-Limit, PKey, Packet Lifetime
         qos-level
name: WholeSet
                    sl: 1
mtu-limit: 4
                   rate-limit: 5
pkey: 0x1234
packet-life: 8
         end-qos-level
end-qos-levels
\# Match rules are scanned in order of their apperance in the policy file. \# First matched rule takes precedence. qos-match-rules
         # matching by single criteria: QoS class
qos-match-rule
   use: by QoS class
                    # Name of qos-level to apply to the matching PR/MPR qos-level-name: WholeSet
         end-qos-match-rule
          # show matching by destination group and service id
         # snow matching by destination group and service id

qos-match-rule

use: Storage targets

destination: Storage

service-id: 0x1000000000001, 0x100000000008-0x1000000000FFF

qos-level-name: WholeSet

end-qos-match-rule
        qos-match-rule
source: Storage
use: match by source group only
qos-level-name: DEFAULT
end-qos-match-rule
qos-match-rule
use: match by all parameters
qos-class: 7-9,11
source: Virtual Servers
destination: Storage
service-id: 0x000000000010000-0x0000000001FFFF
pkey: 0x0F00-0x0FFF
qos-level-name: WholeSet
```

Simple QoS Policy - Details and Examples

Simple QoS policy match rules are tailored for matching ULPs (or some application on top of a ULP) PR/MPR requests. This section has a list of per-ULP (or per-application) match rules and the SL that should be enforced on the matched PR/MPR query. Match rules include:

- Default match rule that is applied to PR/MPR query that didn't match any of the other match rules
- · IPoIB with a default PKey
- IPolB with a specific PKey
- Any ULP/application with a specific Service ID in the PR/MPR query
- Any ULP/application with a specific PKey in the PR/MPR query
- Any ULP/application with a specific target IB port GUID in the PR/MPR query

Since any section of the policy file is optional, as long as basic rules of the file are kept (such as no referring to nonexistent port group, having default QoS Level, etc), the simple policy section (qosulps) can serve as a complete QoS policy file.

The shortest policy file in this case would be as follows:

```
qos-ulps
default : 0 #default SL
end-qos-ulps
```

It is equivalent to the previous example of the shortest policy file, and it is also equivalent to not having policy file at all. Below is an example of simple QoS policy with all the possible keywords:

Similar to the advanced policy definition, matching of PR/MPR queries is done in order of appearance in the QoS policy file such as the first match takes precedence, except for the "default" rule, which is applied only if the query didn't match any other rule. All other sections of the QoS policy file take precedence over the qos-ulps section. That is, if a policy file has both qos-match-rules and qos-ulps sections, then any query is matched first against the rules in the qos-match-rules section, and only if there was no match, the query is matched against the rules in qos-ulps section. Note that some of these match rules may overlap, so in order to use the simple QoS definition effectively, it is important to understand how each of the ULPs is matched.

IPoIB

IPoIB query is matched by PKey or by destination GID, in which case this is the GID of the multicast group that OpenSM creates for each IPoIB partition.

Default PKey for IPoIB partition is 0x7fff, so the following three match rules are equivalent:

```
ipoib:<SL>ipoib, pkey 0x7fff : <SL>
any, pkey 0x7fff : <SL>
```

SRP

Service ID for SRP varies from storage vendor to vendor, thus SRP query is matched by the target IB port GUID. The following two match rules are equivalent:

```
srp, target-port-guid 0x1234 : <SL>
any, target-port-guid 0x1234 : <SL>
```

Note that any of the above ULPs might contain target port GUID in the PR query, so in order for these queries not to be recognized by the QoS manager as SRP, the SRP match rule (or any match rule that refers to the target port GUID only) should be placed at the end of the qos-ulps match rules.

MPI

SL for MPI is manually configured by an MPI admin. OpenSM is not forcing any SL on the MPI traffic, which explains why it is the only ULP that did not appear in the qos-ulps section.

SL2VL Mapping and VL Arbitration

OpenSM cached options file has a set of QoS related configuration parameters, that are used to configure SL2VL mapping and VL arbitration on IB ports. These parameters are:

- Max VLs: the maximum number of VLs that will be on the subnet
- High limit: the limit of High Priority component of VL Arbitration table (IBA 7.6.9)
- VLArb low table: Low priority VL Arbitration table (IBA 7.6.9) template
- VLArb high table: High priority VL Arbitration table (IBA 7.6.9) template
- SL2VL: SL2VL Mapping table (IBA 7.6.6) template. It is a list of VLs corresponding to SLs 0-15 (Note that VL15 used here means drop this SL).

There are separate QoS configuration parameters sets for various target types: CAs, routers, switch external ports, and switch's enhanced port 0. The names of such parameters are prefixed by "qos_<type>_" string. Here is a full list of the currently supported sets:

- qos_ca_ QoS configuration parameters set for CAs.
- qos_rtr_ parameters set for routers.
- gos_sw0_ parameters set for switches' port 0.
- qos_swe_ parameters set for switches' external ports.

Here's the example of typical default values for CAs and switches' external ports (hard-coded in OpenSM initialization):

```
qos_ca_max_vls 15
qos_ca_high_limit 0
qos_ca_vlarb_high 0:4,1:0,2:0,3:0,4:0,5:0,6:0,7:0,8:0,9:0,10:0,11:0,12:0,13:0,14:0
qos_ca_vlarb_low 0:0,1:4,2:4,3:4,4:4,5:4,6:4,7:4,8:4,9:4,10:4,11:4,12:4,13:4,14:4
qos_ca_s12vl 0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,7
qos_swe_max_vls 15
qos_swe_high_limit 0
qos_swe_vlarb_high 0:4,1:0,2:0,3:0,4:0,5:0,6:0,7:0,8:0,9:0,10:0,11:0,12:0,13:0,14:0
qos_swe_vlarb_low 0:0,1:4,2:4,3:4,4:4,5:4,6:4,7:4,8:4,9:4,10:4,11:4,12:4,13:4,14:4
qos_swe_sl2vl 0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,7
```

VL arbitration tables (both high and low) are lists of VL/Weight pairs. Each list entry contains a VL number (values from 0-14), and a weighting value (values 0-255), indicating the number of 64 byte units (credits) which may be transmitted from that VL when its turn in the arbitration occurs. A weight of 0 indicates that this entry should be skipped. If a list entry is programmed for VL15 or for a VL that is not supported or is not currently configured by the port, the port may either skip that entry or send from any supported VL for that entry.

Note, that the same VLs may be listed multiple times in the High or Low priority arbitration tables, and, further, it can be listed in both tables. The limit of high-priority VLArb table (qos_<type>_high_limit) indicates the number of high-priority packets that can be transmitted without an opportunity to send a low-priority packet. Specifically, the number of bytes that can be sent is high_limit times 4K bytes.

A high_limit value of 255 indicates that the byte limit is unbounded.

A

If the 255 value is used, the low priority VLs may be starved.

A value of 0 indicates that only a single packet from the high-priority table may be sent before an opportunity is given to the low-priority table.

Keep in mind that ports usually transmit packets of size equal to MTU. For instance, for 4KB MTU a single packet will require 64 credits, so in order to achieve effective VL arbitration for packets of 4KB MTU, the weighting values for each VL should be multiples of 64.

Below is an example of SL2VL and VL Arbitration configuration on subnet:

```
qos_ca_max_vls 15
qos_ca_high_limit 6
qos_ca_vlarb_high 0:4
qos_ca_vlarb_low 0:0,1:64,2:128,3:192,4:0,5:64,6:64,7:64
qos_ca_sl2vl 0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,7
qos_swe_max_vls 15
qos_swe_high_limit 6
qos_swe_vlarb_high 0:4
qos_swe_vlarb_low 0:0,1:64,2:128,3:192,4:0,5:64,6:64,7:64
qos_swe_sl2vl 0,1,2,3,4,5,6,7,8,9,10,11,12,13,14,7
```

In this example, there are 8 VLs configured on subnet: VL0 to VL7. VL0 is defined as a high priority VL, and it is limited to 6 x 4KB = 24KB in a single transmission burst. Such configuration would suilt VL that needs low latency and uses small MTU when transmitting packets. Rest of VLs are defined as low priority VLs with different weights, while VL4 is effectively turned off.

Deployment Example

The figure below shows an example of an InfiniBand subnet that has been configured by a QoS manager to provide different service levels for various ULPs.

QoS Deployment on InfiniBand Subnet Example

2

Unknown Attachment

QoS Configuration Examples

The following are examples of QoS configuration for different cluster deployments. Each example provides the QoS level assignment and their administration via OpenSM configuration files.

Typical HPC Example: MPI and Lustre

Assignment of QoS Levels

- MPI
- Separate from I/O load
- Min BW of 70%
- Storage Control (Lustre MDS)
 - Low latency
- Storage Data (Lustre OST)
 - Min BW 30%

Administration

• MPI is assigned an SL via the command line

```
host1# mpirun -sl 0
```

• OpenSM QoS policy file

```
qos-ulps

default

any, target-port-guid OST1,OST2,OST3,OST4

any, target-port-guid MDS1,MDS2

end-qos-ulps

:0 # default SL (for MPI)

:1 # SL for Lustre OST

:2 # SL for Lustre MDS
```

Note: In this policy file example, replace OST* and MDS* with the real port GUIDs.

• OpenSM options file

```
qos_max_vls 8
qos_high_limit 0
qos_vlarb_high 2:1
qos_vlarb_low 0:96,1:224
qos_sl2vl 0,1,2,3,4,5,6,7,15,15,15,15,15,15
```

EDC SOA (2-tier): IPoIB and SRP

The following is an example of QoS configuration for a typical enterprise data center (EDC) with service oriented architecture (SOA), with IPoIB carrying all application traffic and SRP used for storage.

QoS Levels

- Application traffic
 - · IPoIB (UD and CM) and SDP
 - Isolated from storage
 - Min BW of 50%
- SRP
- Min BW 50%
- · Bottleneck at storage nodes

Administration

OpenSM QoS policy file

```
qos-ulps
default :0
ipoib :1
sdp :1
srp, target-port-guid SRPT1,SRPT2,SRPT3 :2
end-qos-ulps :2
```

Note: In this policy file example, replace SRPT* with the real SRP Target port GUIDs.

• OpenSM options file

```
gos_max_vls 8
    qos_high_limit 0
    qos_vlarb_high 1:32,2:32
    qos_vlarb_low 0:1,
    qos_sl2vl 0,1,2,3,4,5,6,7,15,15,15,15,15,15
```

EDC (3-tier): IPolB, RDS, SRP

The following is an example of QoS configuration for an enterprise data center (EDC), with IPoIB carrying all application traffic, RDS for database traffic, and SRP used for storage.

QoS Levels

- Management traffic (ssh)
 - IPolB management VLAN (partition A)
 - Min BW 10%
- · Application traffic
 - IPolB application VLAN (partition B)
 - Isolated from storage and database
 - Min BW of 30%
- Database Cluster traffic
 - RDS
 - Min BW of 30%
- SRP
- Min BW 30%
- · Bottleneck at storage nodes

Administration

OpenSM QoS policy file

```
      qos-ulps
      :0

      ipoib, pkey 0x8001
      :1

      ipoib, pkey 0x8002
      :2

      rds
      :3

      srp, target-port-guid SRPT1, SRPT2, SRPT3
      :4

      end-qos-ulps
```

Note: In the following policy file example, replace SRPT* with the real SRP Initiator port GUIDs.

· OpenSM options file

```
qos_max_vls 8
qos_high_limit 0
qos_vlarb_high 1:32,2:96,3:96,4:96
qos_vlarb_low 0:1
qos_sl2v1 0,1,2,3,4,5,6,7,15,15,15,15,15,15
```

· Partition configuration file

```
Default=0x7fff,ipoib : ALL=full;PartA=0x8001, sl=1, ipoib : ALL=full;
```

Enhanced QoS

Enhanced QoS provides a higher resolution of QoS at the service level (SL). Users can configure rate limit values per SL for physical ports, virtual ports, and port groups, using enhanced_qos_policy_file configuration parameter.

Valid values of this parameter:

- Full path to the policy file through which Enhanced QoS Manager is configured
- "null" to disable the Enhanced QoS Manager (default value)



⚠ To enable Enhanced QoS Manager, QoS must be enabled in OpenSM.

Enhanced QoS Policy File

The policy file is comprised of three sections:

· BW_NAMES: Used to define bandwidth setting and name (currently, rate limit is the only setting). Bandwidth names can be used in BW_RULES and VPORT_BW_RULES sections. Bandwidth names are defined using the syntax:

```
<name> = <rate limit in 1Mbps units>
Example: My_bandwidth = 50
```

• BW RULES: Used to define the rules that map the bandwidth setting to a specific SL of a specific GUID.

Bandwidth rules are defined using the syntax:

```
<guid>|<port group name> = <sl id>:<bandwidth name>, <sl id>:<bandwidth name>...
```

```
0 \times 2 c 9 0 0 0 0 0 0 0 25 = 5:My\_bandwidth, 7:My\_bandwidth
Port_grp1 = 3:My_bandwidth, 9:My_bandwidth
```

 VPORT BW RULES: Used to define the rules that map the bandwidth setting to a specific SL of a specific virtual port GUID.

Bandwidth rules are defined using the syntax:

```
<guid>= <sl id>:<bandwidth name>, <sl id>:<bandwidth name>...
Examples:
0x2c9000000026= 5:My_bandwidth, 7:My_bandwidth
```

Special Keywords

 Keyword "all" allows setting a rate limit of all SLs to some BW for a specific physical or virtual port. It is possible to combine "all" with specific SL rate limits. Example:

```
0x2c90000000025 = all:BW1,SL3:BW2
```

In this case, SL3 will be assigned BW2 rate limit, while the rest of SLs get BW1 rate limit.

- "default" is a well-known name which can be used to define a default rule used for any GUID with no defined rule.
 - If no default rule is defined, any GUID without a specific rule will be configured with unlimited rate limit for all SLs.
 - Keyword "all" is also applicable to the default rule. Default rule is local to each section.

Special Subnet Manager Configuration Options

New SM configuration option enhanced qos_vport0_unlimit_default_rl was added to opensm.conf.

The possible values for this configuration option are:

- TRUE: For specific virtual port0 GUID, SLs not mentioned in bandwidth rule will be set to
 unlimited bandwidth (0) regardless of the default rule of the VPORT_BW_RULES section.
 Virtual port0 GUIDs not mentioned in VPORT_BW_SECTION will be set to unlimited BW on all
 SLs.
- FALSE: The GUID of virtual port0 is treated as any other virtual port in VPORT_BW_SECTION. SM should be signaled by HUP once the option is changed.

Default: TRUE

Notes

- When rate limit is set to 0, it means that the bandwidth is unlimited.
- Any unspecified SL in a rule will be set to 0 (unlimited) rate limit automatically if no default rule is specified.
- Failure to complete policy file parsing leads to an undefined behavior. User must confirm no relevant error messages in SM log in order to ensure Enhanced QoS Manager is configured properly.
- A file with only 'BW_NAMES' and 'BW_RULES' keywords configures the network with an unlimited rate limit.
- HCA physical port GUID can be specified in BW RULES and VPORT BW RULES sections.
- In BW_RULES section, the rate limit assigned to a specific SL will limit the total BW that can be sent through the PF on a given SL.
- In VPORT_BW_RULES section, the rate limit assigned to a specific SL will limit only the traffic sent from the IB interface corresponding to the physical port GUID (virtual port0 IB interface). The traffic sent from other virtual IB interfaces will not be limited if no specific rules are defined.

Policy File Example

All physical ports in the fabric are with a rate limit of 50Mbps on SL1, except for GUID 0x2c9000000025, which is configured with rate limit of 25Mbps on SL1. In this example, the traffic on SLs (other than SL1) is unlimited.

All virtual ports in the fabric (except virtual port0 of all physical ports) will be rate-limited to 15Mbps for all SLs because of the default rule of VPORT_BW_RULES section.

Virtual port GUID 0x2c90000000026 is configured with a rate limit of 10Mbps on SL3. The rest of the SLs on this virtual port will get a rate limit of 15 Mbps because of the default rule of VPORT_BW_RULES section.

```
BW_NAMES
bw1 = 50
bw2 = 25
bw3 = 15
bw4 = 10

BW_RULES
default= 1:bw1
0x2c90000000025= 1:bw2

VPORT_BW_RULES
default= all:bw3
0x2c90000000026= 3:bw4
```

Adaptive Routing Manager and SHIELD

Adaptive Routing Manager supports advanced InfiniBand features; Adaptive Routing (AR) and Self-Healing Interconnect Enhancement for InteLligent Datacenters (SHIELD).

For information on how to set up AR and SHIELD, please refer to <u>HowTo Configure Adaptive Routing</u> and <u>SHIELD</u> Community post.

DOS MAD Prevention

DOS MAD prevention is achieved by assigning a threshold for each agent's RX. Agent's RX threshold provides a protection mechanism to the host memory by limiting the agents' RX with a threshold. Incoming MADs above the threshold are dropped and are not queued to the agent's RX.

? Unknown Attachment

To enable DOS MAD Prevention:

- 1. Go to /etc/modprobe.d/mlnx.conf.
- 2. Add to the file the option below.

```
ib_umad enable_rx_threshold 1
```

The threshold value can be controlled from the user-space via libibumad.

To change the value, use the following API:

```
int umad_update_threshold(int fd, int threshold);

@fd: file descriptor, agent's RX associated to this fd.
@threshold: new threshold value
```

IB Router Support in OpenSM

In order to enable the IB router in OpenSM, the following parameters should be configured:

IB Router Parameters for OpenSM

Parameter	Description	Default Value
rtr_pr_flow_l abel	Defines whether the SM should create alias GUIDs required for router support for each port.	0 (Disabled)
	Defines flow label value to use in response for path records related to the router.	
rtr_pr_tclass	Defines TClass value to use in response for path records related to the router	0
rtr_pr_sl	Defines sl value to use in response for path records related to router.	0
rtr_p_mtu	Defines MTU value to use in response for path records related to the router.	4 (IB_MTU_LEN_2048)
rtr_pr_rate	Defines rate value to use in response for path records related to the router.	16 (IB_PATH_RE- CORD_RATE_100_GBS)

OpenSM Activity Report

OpenSM can produce an activity report in a form of a dump file which details the different activities done in the SM. Activities are divided into subjects. The OpenSM Supported Activities table below specifies the different activities currently supported in the SM activity report. Reporting of each subject can be enabled individually using the configuration parameter activity_report_subjects:

• Valid values:

Comma separated list of subjects to dump. The current supported subjects are:

- "mc" activity IDs 1, 2 and 8
- "prtn" activity IDs 3, 4, and 5
- "virt" activity IDs 6 and 7
- "routing" activity IDs 8-12

Two predefined values can be configured as well:

- "all" dump all subjects
- "none" disable the feature by dumping none of the subjects
- Default value: "none"

OpenSM Supported Activities

ACtivity ID	Activity Name	Additional Fields	Comments	Description
1	mcm_member	MLid MGid Port Guid Join State	Join state: 1 - Join -1 - Leave	Member joined/ left MC group
2	mcg_change	MLid MGid Change	Change: 0 - Create 1 - Delete	MC group created/ deleted
3	prtn_guid_add	Port GuidPKeyBlock indexPkey Index		Guid added to partition
4	prtn_create	-PKey • Prtn Name		Partition created
5	prtn_delete	PKey Delete Reason	Delete Reason: 0 - empty prtn 1 - duplicate prtn 2 - sm shutdown	Partition deleted
6	port_virt_discover	Port Guid Top Index		Port virtualization discovered

ACtivity ID	Activity Name	Additional Fields	Comments	Description
7	vport_state_change	Port GuidVPort GuidVPort IndexVNode GuidVPort State	VPort State: 1 - Down 2 - Init 3 - ARMED 4 - Active	Vport state changed
8	mcg_tree_calc	mlid		MCast group tree calculated
9	routing_succeed	routing engine name		Routing done successfully
10	routing_failed	routing engine name		Routing failed
11	ucast_cache_invali- dated			ucast cache invalidated
12	ucast_cache_rout- ing_done			ucast cache routing done

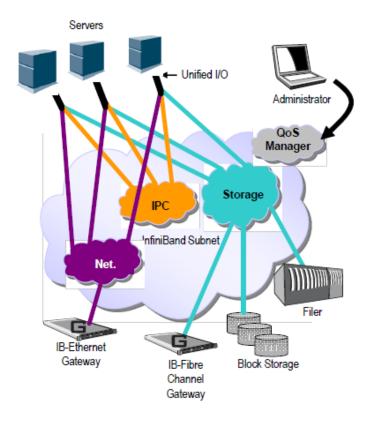
Offsweep Balancing

When working with minhop/dor/updn, subnet manager can re-balance routing during idle time (between sweeps).

- offsweep_balancing_enabled enables/disables the feature. Examples:
 - offsweep_balancing_enabled = TRUE
 - offsweep_balancing_enabled = FALSE (default)
- offsweep_balancing_window defines window of seconds to wait after sweep before starting the re-balance process. Applicable only if offsweep_balancing_enabled=TRUE. Example: offsweep_balancing_window = 180 (default)

QoS - Quality of Service

Quality of Service (QoS) requirements stem from the realization of I/O consolidation over an IB network. As multiple applications and ULPs share the same fabric, a means is needed to control their use of network resources.



The basic need is to differentiate the service levels provided to different traffic flows, such that a policy can be enforced and can control each flow utilization of fabric resources. The InfiniBand Architecture Specification defines several hardware features and management interfaces for supporting QoS:

Up to 15 Virtual Lanes (VL) carry traffic in a non-blocking manner

- Arbitration between traffic of different VLs is performed by a two-priority-level weighted round robin arbiter. The arbiter is programmable with a sequence of (VL, weight) pairs and a maximal number of high priority credits to be processed before low priority is served
- Packets carry class of service marking in the range 0 to 15 in their header SL field
- Each switch can map the incoming packet by its SL to a particular output VL, based on a programmable table VL=SL-to-VL-MAP(in-port, out-port, SL)
- The Subnet Administrator controls the parameters of each communication flow by providing them as a response to Path Record (PR) or MultiPathRecord (MPR) gueries

DiffServ architecture (IETF RFC 2474 & 2475) is widely used in highly dynamic fabrics. The following subsections provide the functional definition of the various software elements that enable a DiffServ-like architecture over the Mellanox OFED software stack.

QoS Architecture

QoS functionality is split between the SM/SA, CMA and the various ULPs. We take the "chronology approach" to describe how the overall system works.

1. The network manager (human) provides a set of rules (policy) that define how the network is being configured and how its resources are split to different QoS-Levels. The policy also define how to decide which QoS-Level each application or ULP or service use.

- 2. The SM analyzes the provided policy to see if it is realizable and performs the necessary fabric setup. Part of this policy defines the default QoS-Level of each partition. The SA is enhanced to match the requested Source, Destination, QoS-Class, Service-ID, PKey against the policy, so clients (ULPs, programs) can obtain a policy enforced QoS. The SM may also set up partitions with appropriate IPoIB broadcast group. This broadcast group carries its QoS attributes: SL, MTU, RATE, and Packet Lifetime.
- 3. IPOIB is being setup. IPOIB uses the SL, MTU, RATE and Packet Lifetime available on the multicast group which forms the broadcast group of this partition.
- 4. MPI which provides non IB based connection management should be configured to run using hard coded SLs. It uses these SLs for every QP being opened.
- 5. ULPs that use CM interface (like SRP) have their own pre-assigned Service-ID and use it while obtaining PathRecord/MultiPathRecord (PR/MPR) for establishing connections. The SA receiving the PR/MPR matches it against the policy and returns the appropriate PR/MPR including SL, MTU, RATE and Lifetime.
- 6. ULPs and programs (e.g. SDP) use CMA to establish RC connection provide the CMA the target IP and port number. ULPs might also provide QoS-Class. The CMA then creates Service-ID for the ULP and passes this ID and optional QoS-Class in the PR/MPR request. The resulting PR/ MPR is used for configuring the connection QP.

PathRecord and Multi Path Record Enhancement for QoS:

As mentioned above, the PathRecord and MultiPathRecord attributes are enhanced to carry the Service-ID which is a 64bit value. A new field QoS-Class is also provided.

A new capability bit describes the SM QoS support in the SA class port info. This approach provides an easy migration path for existing access layer and ULPs by not introducing new set of PR/MPR attributes.

Supported Policy

The QoS policy, which is specified in a stand-alone file, is divided into the following four subsections:

Port Group

A set of CAs, Routers or Switches that share the same settings. A port group might be a partition defined by the partition manager policy, list of GUIDs, or list of port names based on NodeDescription.

Fabric Setup

Defines how the SL2VL and VLArb tables should be set up.



⚠ In OFED this part of the policy is ignored. SL2VL and VLArb tables should be configured in the OpenSM options file (opensm.opts).

QoS-Levels Definition

This section defines the possible sets of parameters for QoS that a client might be mapped to. Each set holds SL and optionally: Max MTU, Max Rate, Packet Lifetime and Path Bits.



Path Bits are not implemented in OFED.

Matching Rules

A list of rules that match an incoming PR/MPR request to a QoS-Level. The rules are processed in order such as the first match is applied. Each rule is built out of a set of match expressions which should all match for the rule to apply. The matching expressions are defined for the following fields:

- SRC and DST to lists of port groups
- Service-ID to a list of Service-ID values or ranges
- QoS-Class to a list of QoS-Class values or ranges

CMA Features

The CMA interface supports Service-ID through the notion of port space as a prefix to the port number, which is part of the sockaddr provided to rdma_resolve_add(). The CMA also allows the ULP (like SDP) to propagate a request for a specific QoS-Class. The CMA uses the provided QoS-Class and Service-ID in the sent PR/MPR.

IPolB

IPoIB queries the SA for its broadcast group information and uses the SL, MTU, RATE and Packet Lifetime available on the multicast group which forms this broadcast group.

SRP

The current SRP implementation uses its own CM callbacks (not CMA). So SRP fills in the Service-ID in the PR/MPR by itself and use that information in setting up the QP. SRP Service-ID is defined by the SRP target I/O Controller (it also complies with IBTA Service-ID rules). The Service-ID is reported by the I/O Controller in the ServiceEntries DMA attribute and should be used in the PR/MPR if the SA reports its ability to handle QoS PR/MPRs.

IP over InfiniBand (IPoIB)

Upper Layer Protocol (ULP)

The IP over IB (IPoIB) ULP driver is a network interface implementation over InfiniBand. IPoIB encapsulates IP datagrams over an InfiniBand Connected or Datagram transport service. The IPoIB driver, ib_ipoib, exploits the following capabilities:

- VLAN simulation over an InfiniBand network via child interfaces
- · High Availability via Bonding
- Varies MTU values:
 - up to 4k in Datagram mode
 - up to 64k in Connected mode
- Uses any ConnectX® IB ports (one or two)
- Inserts IP/UDP/TCP checksum on outgoing packets
- Calculates checksum on received packets
- Support net device TSO through ConnectX® LSO capability to defragment large data- grams to MTU quantas.
- Dual operation mode datagram and connected
- Large MTU support through connected mode

IPoIB also supports the following software based enhancements:

- · Giant Receive Offload
- NAPI
- Ethtool support

Enhanced IPolB



Note that switching between Enhanced mode and ULP mode can be done by setting ib_ipoib module parameter "ipoib_enhanced" to 1 or 0.

Enhanced IPoIB feature enables offloading ULP basic capabilities to a lower vendor specific driver, in order to optimize IPoIB data path. This will allow IPoIB to support multiple stateless offloads, such as RSS/TSS, and better utilize the features supported, enabling IPoIB datagram to reach peak performance in both bandwidth and latency.

Enhanced IPoIB supports/performs the following:

- Stateless offloads (RSS, TSS)
- · Multi queues
- · Interrupt moderation
- · Multi partitions optimizations
- Sharing send/receive Work Queues
- · Vendor specific optimizations
- · UD mode only

In order to verify that the driver is using Enhanced IPolB, run:

ip link show ibX

Output example:

8: ib1: <BROADCAST,MULTICAST> mtu 4092 qdisc noop state DOWN mode DEFAULT qlen 1024 link/infiniband 00:00:00:67:fe:80:00:00:00:00:00:00:e4:1d:2d:03:00:a5:f0:2f brd 00:ff:ff:ff:ff:12:40:1b:ff:ff:00:00:00:00:00:ff:ff:ff:ff:

Note: The driver MAC should start with 00:xxxxxxxx (Upstream) or 20:xxxxxxx (OFED) in case Enhanced IPoIB is enabled.

IPolB Mode Setting

IPoIB ULP can run in two modes of operation: Connected mode and Datagram mode. By default, IPoIB ULP is set to work in Datagram mode.

For better scalability and performance, we recommend using the Datagram mode. However, the mode can be changed to Connected mode by editing the file /etc/infiniband/openib.conf and setting 'SET_IPOIB_CM=yes'.

After changing the mode, you need to restart the driver by running:

/etc/init.d/openibd restart



To check the current mode used for out-going connections, enter:

cat /sys/class/net/ib<n>/mode

Changing the IPoIB mode (CM vs UD) requires the interface to be in 'down' state.

Connected mode is not supported when using enhanced IPolB.

Port Configuration

The physical port MTU in Datagram mode (indicates the port capability) default value is 4k, whereas the IPoIB port MTU ("logical" MTU) default value is 2k as it is set by the OpenSM.

To change the IPoIB MTU to 4k, edit the OpenSM partition file in the section of IPoIB setting as follow:

Default=0xffff, ipoib, mtu=5 : ALL=full;

where:

"mtu=5" indicates that all IPoIB ports in the fabric are using 4k MTU, ("mtu=4" indi- cates 2k MTU)

IPolB Configuration

Unless you have run the installation script mlnxofedinstall with the flag '-n', then IPoIB has not been configured by the installation. The configuration of IPoIB requires assigning an IP address and a subnet mask to each HCA port, like any other network adapter card (i.e., you need to prepare a file called ifcfg-ib<n> for each port). The first port on the first HCA in the host is called interface ib0, the second port is called ib1, and so on.

IPoIB configuration can be based on DHCP or on a static configuration that you need to supply (see below). You can also apply a manual configuration that persists only until the next reboot or driver restart (see below).

IPolB Configuration Based on DHCP

Setting an IPoIB interface configuration based on DHCP is performed similarly to the configuration of Ethernet interfaces. In other words, you need to make sure that IPoIB configuration files include the following line:

· For RedHat:

BOOTPROTO=dhcp

· For SLES:

BOOTPROTO='dchp'

If IPoIB configuration files are included, ifcfg-ib<n> files will be installed under:

/etc/sysconfig/network-scripts/ on a RedHat machine /etc/sysconfig/network/ on a SuSE machine.

 A patch for DHCP may be required for supporting IPolB. For further information, please see the REAME file available under the docs/dhcp/ directory.

Standard DHCP fields holding MAC addresses are not large enough to contain an IPoIB hardware address. To overcome this problem, DHCP over InfiniBand messages convey a client identifier field used to identify the DHCP session. This client identifier field can be used to associate an IP address with a client identifier value, such that the DHCP server will grant the same IP address to any client that conveys this client identifier.

The length of the client identifier field is not fixed in the specification. For the Mellanox OFED for Linux package, it is recommended to have IPoIB use the same format that FlexBoot uses for this client identifier.

DHCP Server

In order for the DHCP server to provide configuration records for clients, an appropriate configuration file needs to be created. By default, the DHCP server looks for a configuration file called dhcpd.conf under /etc. You can either edit this file or create a new one and provide its full path to the DHCP server using the -cf flag (See a file example at docs/dhcpd.conf). The DHCP server must run on a machine which has loaded the IPoIB module. To run the DHCP server from the command line, enter:

```
dhcpd <IB network interface name> -d
```

Example:

```
host1# dhcpd ib0 -d
```

DHCP Client (Optional)



A DHCP client can be used if you need to prepare a diskless machine with an IB driver.

In order to use a DHCP client identifier, you need to first create a configuration file that defines the DHCP client identifier.

Then run the DHCP client with this file using the following command:

```
dhclient -cf <client conf file> <IB network interface name>
```

Example of a configuration file for the ConnectX (PCI Device ID 26428), called dhclient.conf:

```
The value indicates a hexadecimal number interface "ib1" {
send dhcp-client-identifier
ff:00:00:00:00:00:02:00:00:02:c9:00:00:02:c9:03:00:00:10:39;
```

Example of a configuration file for InfiniHost III Ex (PCI Device ID 25218), called dhclient.conf:

```
The value indicates a hexadecimal number interface "ib1" { send dhcp-client-identifier
20:00:55:04:01:fe:80:00:00:00:00:00:00:02:c9:02:00:23:13:92;
```



In order to use the configuration file, run:

```
host1# dhclient -cf dhclient.conf ib1
```

Static IPolB Configuration

If you wish to use an IPoIB configuration that is not based on DHCP, you need to supply the installation script with a configuration file (using the '-n' option) containing the full IP configuration. The IPoIB configuration file can specify either or both of the following data for an IPoIB interface:

- · A static IPoIB configuration
- An IPoIB configuration based on an Ethernet configuration
 See your Linux distribution documentation for additional information about configuring IP addresses.

The following code lines are an excerpt from a sample IPoIB configuration file:

```
# Static settings; all values provided by this file
IPADDR_ib0=11.4.3.175
NETMASK_ib0=255.255.0.0
NETWORK_ib0=11.4.0.0
BROADCAST_ib0=11.4.255.255
ONBOOT_ib0=1
# Based on eth0; each '*' will be replaced with a corresponding octet
# from eth0.
LAN_INTERPACE_ib0=eth0
IPADDR_ib0=11.4.'*'.'*'
NETMASK_ib0=255.255.0.0
NETWORK_ib0=11.4.0.0
BROADCAST_ib0=11.4.255.255
ONBOOT_ib0=1
# Based on the first eth<n> interface that is found (for n=0,1,...);
# each '*' will be replaced with a corresponding octet from eth<n>.
LAN_INTERPACE_ib0=
IPADDR_ib0=11.4.'*'.'*'
NETMASK_ib0=255.255.0.0
NETWORK_ib0=11.4.'*'.'*'
NETMASK_ib0=255.255.0.0
NETWORK_ib0=11.4.'*'.'*'
NETMASK_ib0=255.255.0.0
BROADCAST_ib0=11.4.0.0
BROADCAST_ib0=11.4.0.0
BROADCAST_ib0=11.4.0.0
BROADCAST_ib0=11.4.255.255
ONBOOT_ib0=1
```

Manually Configuring IPolB

A

This manual configuration persists only until the next reboot or driver restart.



To manually configure IPoIB for the default IB partition (VLAN), perform the following steps:

- 1. Configure the interface by entering the ifconfig command with the following items:
 - The appropriate IB interface (ib0, ib1, etc.)
 - The IP address that you want to assign to the interface
 - The netmask keyword
 - The subnet mask that you want to assign to the interface

The following example shows how to configure an IB interface:

```
host1$ ifconfig ib0 11.4.3.175 netmask 255.255.0.0
```

2. (Optional) Verify the configuration by entering the ifconfig command with the appropriate interface identifier *ib#* argument.

The following example shows how to verify the configuration:

```
host1$ ifconfig ib0
b0 Link encap:UNSPEC HWaddr 80-00-04-04-FE-80-00-00-00-00-00-00-00-00-00-00
inet addr:11.4.3.175 Bcast:11.4.255.255 Mask:255.255.0.0
UP BROADCAST MULTICAST MTU:65520 Metric:1
RX packets:0 errors:0 dropped:0 overruns:0 frame:0
TX packets:0 errors:0 dropped:0 overruns:0 carrier:0
collisions:0 txqueuelen:128
RX bytes:0 (0.0 b) TX bytes:0 (0.0 b)
```

3. Repeat the first two steps on the remaining interface(s).

Sub-interfaces

You can create sub-interfaces for a primary IPoIB interface to provide traffic isolation. Each such sub-interface (also called a child interface) has a different IP and network addresses from the primary (parent) interface. The default Partition Key (PKey), ff:ff, applies to the primary (parent) interface.

This section describes how to:

- · Create a subinterface
- Remove a subinterface

Creating a Subinterface

In the following procedure, ib0 is used as an example of an IB sub-interface.

To create a child interface (sub-interface), follow this procedure:

- 1. Decide on the PKey to be used in the subnet (valid values can be 0 or any 16-bit unsigned value). The actual PKey used is a 16-bit number with the most significant bit set. For example, a value of 1 will give a PKey with the value 0x8001.
- 2. Create a child interface by running:

```
host1$ echo <PKey> > /sys/class/net/<IB subinterface>/create_child
```

Example:

```
host1$ echo 1 > /sys/class/net/ib0/create_child
```

This will create the interface ib0.8001.

3. Verify the configuration of this interface by running:

```
host1$ ifconfig <subinterface>.<subinterface PKey>
```

Using the example of the previous step:

```
host1$ ifconfig ib0.8001
ib0.8001 Link encap:UNSPEC HWaddr 80-00-00-4A-FE-80-00-00-00-00-00-00-00-00-00-00
BROADCAST MULTICAST MTU:2044 Metric:1
RX packets:0 errors:0 dropped:0 overruns:0 frame:0
TX packets:0 errors:0 dropped:0 overruns:0 carrier:0
collisions:0 txqueuelen:128
RX bytes:0 (0.0 b) TX bytes:0 (0.0 b)
```

- 4. As can be seen, the interface does not have IP or network addresses. To configure those, you should follow the manual configuration procedure described in "Manually Configuring IPoIB" section above.
- 5. To be able to use this interface, a configuration of the Subnet Manager is needed so that the PKey chosen, which defines a broadcast address, be recognized.

Removing a Subinterface



To remove a child interface (subinterface), run:

```
echo <subinterface PKey> /sys/class/net/<ib_interface>/delete_child
```

Using the example of the second step from the previous chapter:

```
echo 0x8001 > /sys/class/net/ib0/delete_child
```

Note that when deleting the interface you must use the PKey value with the most significant bit set (e.g., 0x8000 in the example above).

Verifying IPolB Functionality



To verify your configuration and IPoIB functionality are successful, perform the following steps:

1. Verify the IPoIB functionality by using the ifconfig command. The following example shows how two IB nodes are used to verify IPoIB functionality. In the following example, IB node 1 is at 11.4.3.175, and IB node 2 is at 11.4.3.176:

```
host1# ifconfig ib0 11.4.3.175 netmask 255.255.0.0 host2# ifconfig ib0 11.4.3.176 netmask 255.255.0.0
```

- 2. Enter the ping command from 11.4.3.175 to 11.4.3.176.
- 3. The following example shows how to enter the ping command:

```
host1# ping -c 5 11.4.3.176
PING 11.4.3.176 (11.4.3.176) 56(84) bytes of data.
64 bytes from 11.4.3.176: icmp_seq=0 ttl=64 time=0.079 ms
64 bytes from 11.4.3.176: icmp_seq=1 ttl=64 time=0.044 ms
64 bytes from 11.4.3.176: icmp_seq=2 ttl=64 time=0.055 ms
64 bytes from 11.4.3.176: icmp_seq=3 ttl=64 time=0.049 ms
64 bytes from 11.4.3.176: icmp_seq=4 ttl=64 time=0.065 ms
--- 11.4.3.176 ping statistics ---
5 packets transmitted, 5 received, 0% packet loss, time 3999ms rtt min/avg/max/mdev = 0.044/0.058/0.079/
0.014 ms, pipe 2
 0.014 ms, pipe 2
```

Bonding IPolB

To create an interface configuration script for the ibX and bondX interfaces, you should use the standard syntax (depending on your OS).

Bonding of IPoIB interfaces is accomplished in the same manner as would bonding of Ethernet interfaces: via the Linux Bonding Driver.

- Network Script files for IPoIB slaves are named after the IPoIB interfaces (e.g. ifcfg- ib0)
- The only meaningful bonding policy in IPoIB is High-Availability (bonding mode number 1, or active-backup)
- · Bonding parameter "fail_over_mac" is meaningless in IPoIB interfaces, hence, the only supported value is the default: 0

For a persistent bonding IPoIB Network configuration, use the same Linux Network Scripts semantics, with the following exceptions/ additions:

• In the bonding master configuration file (e.g. ifcfg-bond0), in addition to Linux bonding semantics, use the following parameter: MTU=65520

♠ 65520 is a valid MTU value only if all IPoIB slaves operate in Connected mode (See "IPoIB Mode Setting") and are configured with the same value. For IPoIB slaves that work in datagram mode, use MTU=2044. If you do not set the correct MTU or do not set MTU at all, performance of the interface might decrease.

Dynamically Connected Transport (DCT)

- In the bonding slave configuration file (e.g. ifcfg-ib0), use the same Linux Network Scripts semantics. In particular: DEVICE=ib0
- In the bonding slave configuration file (e.g. ifcfg-ib0.8003), the line TYPE=InfiniBand is necessary when using bonding over devices configured with partitions (p_key)
- For RHEL users:

In /etc/modprobe.b/bond.conf add the following lines:

alias bond0 bonding

· For SLES users:

It is necessary to update the MANDATORY_DEVICES environment variable in /etc/sysconfig/ network/config with the names of the IPoIB slave devices (e.g. ib0, ib1, etc.). Otherwise, bonding master may be created before IPoIB slave interfaces at boot time. It is possible to have multiple IPoIB bonding masters and a mix of IPoIB bonding master and Ethernet bonding master. However, It is NOT possible to mix Ethernet and IPoIB slaves under the same bonding master.



Restarting openibd does no keep the bonding configuration via Network Scripts. You have to restart the network service in order to bring up the bonding master. After the configuration is saved, restart the network service by running: /etc/init.d/network restart.

Dynamic PKey Change

Dynamic PKey change means the PKey can be changed (add/removed) in the SM database and the interface that is attached to that PKey is updated immediately without the need to restart the

If the PKey is already configured in the port by the SM, the child-interface can be used immediately. If not, the interface will be ready to use only when SM adds the relevant PKey value to the port after the creation of the child interface. No additional configuration is required once the childinterface is created.

Precision Time Protocol (PTP) over IPolB

This feature allows for accurate synchronization between the distributed entities over the network. The synchronization is based on symmetric Round Trip Time (RTT) between the master and slave devices.

This feature is enabled by default, and is also supported over PKey interfaces.

For more on the PTP feature, refer to Running Linux PTP with ConnectX-4/ConnectX-5 Community post.

For further information on Time-Stamping, follow the steps in "Time-Stamping Service".

One Pulse Per Second (1PPS) over IPolB

1PPS is a time synchronization feature that allows the adapter to be able to send or receive 1 pulse per second on a dedicated pin on the adapter card using an SMA connector (SubMiniature version A). Only one pin is supported and could be configured as 1PPS in or 1PPS out. For further information, refer to HowTo Test 1PPS on Mellanox Adapters Community post.

Advanced Transport

Atomic Operations

Atomic Operations in mlx5 Driver

To enable atomic operation with this endianness contradiction, use the <code>ibv_create_qp</code> to create the QP and set the <code>IBV_QP_CREATE_ATOMIC_BE_REPLY</code> flag on <code>create_flags</code>.

Enhanced Atomic Operations

ConnectX® implements a set of Extended Atomic Operations beyond those defined by the IB spec. Atomicity guarantees, Atomic Ack generation, ordering rules and error behavior for this set of extended Atomic operations is the same as that for IB standard Atomic operations (as defined in section 9.4.5 of the IB spec).

Masked Compare and Swap (MskCmpSwap)

The MskCmpSwap atomic operation is an extension to the CmpSwap operation defined in the IB spec. MskCmpSwap allows the user to select a portion of the 64 bit target data for the "compare" check, as well as to restrict the swap to a (possibly different) portion. The pseudo-code below describes the operation:

The MFetchAdd Atomic operation extends the functionality of the standard IB FetchAdd by allowing the user to split the target into multiple fields of selectable length. The atomic add is done independently on each one of this fields. A bit set in the field_boundary parameter specifies the field boundaries. The pseudo-code below describes the operation:

The additional operands are carried in the Extended Transport Header. Atomic response generation and packet format for MskCmpSwap is as for standard IB Atomic operations.

```
| atomic_response = *va
| if (!((compare_add ^ *va) & compare_add_mask)) then
| *va = (*va & ~(swap_mask)) | (swap & swap_mask)
|
| return atomic_response
```

Masked Fetch and Add (MFetchAdd)

The MFetchAdd Atomic operation extends the functionality of the standard IB FetchAdd by allowing the user to split the target into multiple fields of selectable length. The atomic add is done independently on each one of this fields. A bit set in the field_boundary parameter specifies the field boundaries. The pseudocode below describes the operation:

```
| bit_adder(ci, b1, b2, *co)
| {
| value = ci + b1 + b2 |
| *co = !!(value & 2)
| return value & 1
| }
| #define MASK_IS_SET(mask, attr) | (!!((mask)&(attr)))
| bit_position = 1 |
| carry = 0 |
| atomic_response = 0
| for i = 0 to 63
| {
| if (i != 0) |
| bit_position = bit_position << 1
| bit_add_res = bit_adder(carry, MASK_IS_SET(*va, bit_position), MASK_IS_SET(compare_add, bit_position), &new_carry)
| if (bit_add_res) |
| atomic_response | = bit_position |
| carry = ((new_carry) && (!MASK_IS_SET(compare_add_mask, bit_position)))
| return atomic_response</pre>
```

XRC - eXtended Reliable Connected Transport Service for InfiniBand

XRC allows significant savings in the number of QPs and the associated memory resources required to establish all to all process connectivity in large clusters.

It significantly improves the scalability of the solution for large clusters of multicore end-nodes by reducing the required resources.

For further details, please refer to the "Annex A14 Supplement to InfiniBand Architecture Specification Volume 1.2.1"

A new API can be used by user space applications to work with the XRC transport. The legacy API is currently supported in both binary and source modes, however it is deprecated. Thus we recommend using the new API.

The new verbs to be used are:

- ibv_open_xrcd/ibv_close_xrcd
- ibv_create_srq_ex
- ibv_get_srq_num
- ibv_create_qp_ex
- ibv_open_qp

Please use <code>ibv_xsrq_pingpong</code> for basic tests and code reference. For detailed information regarding the various options for these verbs, please refer to their appropriate man pages.

Dynamically Connected Transport (DCT)

Dynamically Connected transport (DCT) service is an extension to transport services to enable a higher degree of scalability while maintaining high performance for sparse traffic. Utilization of DCT reduces the total number of QPs required system wide by having Reliable type QPs dynamically connect and disconnect from any remote node. DCT connections only stay connected while they are active. This results in smaller memory footprint, less overhead to set connections and higher on-chip cache utilization and hence increased performance. DCT is supported only in mlx5 driver.

MPI Tag Matching and Rendezvous Offloads



Supported in ConnectX®-5 and above adapter cards.

Tag Matching and Rendezvous Offloads is a technology employed by Mellanox to offload the processing of MPI messages from the host machine onto the network card. Employing this technology enables a zero copy of MPI messages, i.e. messages are scattered directly to the user's buffer without intermediate buffering and copies. It also provides a complete rendezvous progress by Mellanox devices. Such overlap capability enables the CPU to perform the application's computational tasks while the remote data is gathered by the adapter. For more information Tag Matching Offload, please refer to the Community post "Understanding MPI Tag Matching and Rendezvous Offloads (ConnectX-5)".

Optimized Memory Access

Memory Region Re-registration

Memory Region Re-registration allows the user to change attributes of the memory region. The user may change the PD, access flags or the address and length of the memory region. Memory region supports contagious pages allocation. Consequently, it de-registers memory region followed by register memory region. Where possible, resources are reused instead of de-allocated and reallocated.



Please note that the verb is implemented as an experimental verb.

Example:

int ibv_rereg_mr(struct ibv_mr *mr, int flags, struct ibv_pd *pd, void *addr, size_t length, uint64_t access,
struct ibv_rereg_mr_attr *attr);

@mr:	The memory region to modify.
@flags:	A bit-mask used to indicate which of the following properties of the memory region are being modified. Flags should be one of: IBV_REREG_MR_CHANGE_TRANSLATION /* Change translation (location and length) */ IBV_REREG_MR_CHANGE_PD/* Change protection domain*/ IBV_REREG_MR_CHANGE_ACCESS/* Change access flags*/
@pd:	If IBV_REREG_MR_CHANGE_PD is set in flags, this field specifies the new protection domain to associated with the memory region, otherwise, this parameter is ignored.
@addr:	If IBV_REREG_MR_CHANGE_TRANSLATION is set in flags, this field specifies the start of the virtual address to use in the new translation, otherwise, this parameter is ignored.
@length:	If IBV_REREG_MR_CHANGE_TRANSLATION is set in flags, this field specifies the length of the virtual address to use in the new translation, otherwise, this parameter is ignored.

@access:	If IBV_REREG_MR_CHANGE_ACCESS is set in flags, this field specifies the new memory access rights, otherwise, this parameter is ignored. Could be one of the following: IBV_ACCESS_LOCAL_WRITE IBV_ACCESS_REMOTE_WRITE IBV_ACCESS_REMOTE_READ IBV_ACCESS_ALLOCATE_MR /* Let the library allocate the memory for * the user, tries to get contiguous pages */
@attr:	Future extensions

ibv_rereg_mr returns 0 on success, or the value of an errno on failure (which indicates the error reason). In case of an error, the MR is in undefined state. The user needs to call ibv_dereg_mr in order to release it.

Please note that if the MR (Memory Region) is created as a Shared MR and a translation is requested, after the call, the MR is no longer a shared MR. Moreover, Re-registration of MRs that uses Mellanox PeerDirect[™] technology are not supported.

Memory Window

Memory Window allows the application to have a more flexible control over remote access to its memory. It is available only on physical functions/native machines The two types of Memory Windows supported are: type 1 and type 2B.

Memory Windows are intended for situations where the application wants to:

- · Grant and revoke remote access rights to a registered region in a dynamic fashion with less of a performance penalty
- Grant different remote access rights to different remote agents and/or grant those rights over different ranges within registered region

For further information, please refer to the InfiniBand specification document.



Memory Windows API cannot co-work with peer memory clients (Mellanox PeerDirect™).

Query Capabilities

Memory Windows are available if and only the hardware supports it. To verify whether Memory Windows are available, run ibv_query_device.

For example:

```
truct ibv_device_attr device_attr = {.comp_mask = IBV_DEVICE_ATTR_RESERVED - 1};
ibv_query_device(context, & device_attr);
if (device_attr.exp_device_cap_flags & IBV_DEVICE_MEM_WINDOW ||
    device_attr.exp_device_cap_flags & IBV_DEVICE_MW_TYPE_2B) {
    /* Memory window is supported */
```

Memory Window Allocation

Allocating memory window is done by calling the ibv_alloc_mw verb.

```
type_mw = IBV_MW_TYPE_2/ IBV_MW_TYPE_1
mw = ibv_alloc_mw(pd, type_mw);
```

Binding Memory Windows

After being allocated, memory window should be bound to a registered memory region. Memory Region should have been registered using the IBV_ACCESS_MW_BIND access flag.

For further information on how to bind memory windows, please see rdma-core man page.

Invalidating Memory Window

Before rebinding Memory Window type 2, it must be invalidated using ibv_post_send - see here.

Deallocating Memory Window

Deallocating memory window is done using the ibv_dealloc_mw verb.

ibv_dealloc_mw(mw);

User-Mode Memory Registration (UMR)

User-mode Memory Registration (UMR) is a fast registration mode which uses send queue. The UMR support enables the usage of RDMA operations and scatters the data at the remote side through the definition of appropriate memory keys on the remote side.

UMR enables the user to:

- Create indirect memory keys from previously registered memory regions, including creation of KLM's from previous KLM's. There are not data alignment or length restrictions associated with the memory regions used to define the new KLM's.
- Create memory regions, which support the definition of regular non-contiguous memory regions.

On-Demand-Paging (ODP)

On-Demand-Paging (ODP) is a technique to alleviate much of the shortcomings of memory registration. Applications no longer need to pin down the underlying physical pages of the address space, and track the validity of the mappings. Rather, the HCA requests the latest translations from the OS when pages are not present, and the OS invalidates translations which are no longer valid due to either non-present pages or mapping changes. ODP does not support contiguous pages. ODP can be further divided into 2 subclasses: Explicit and Implicit ODP.

- Explicit ODP
 In Explicit ODP, applications still register memory buffers for communication, but this operation is used to define access control for IO rather than pin-down the pages. ODP Memory Region (MR) does not need to have valid mappings at registration time.
- Implicit ODP In Implicit ODP, applications are provided with a special memory key that represents their complete address space. This all IO accesses referencing this key (subject to the access rights associated with the key) does not need to register any virtual address range.

Query Capabilities

On-Demand Paging is available if both the hardware and the kernel support it. To verify whether ODP is supported, run <code>ibv_query_device</code>.

For further information, please refer to the ibv_query_device manual page.

Registering ODP Explicit and Implicit MR

ODP Explicit MR is registered after allocating the necessary resources (e.g. PD, buffer), while ODP implicit MR registration provides an implicit lkey that represents the complete address space.

For further information, please refer to the <code>ibv_reg_mr</code> manual page.

De-registering ODP MR

ODP MR is deregistered the same way a regular MR is deregistered:

ibv_dereg_mr(mr);

Advice MR Verb

The driver can pre-fetch a given range of pages and map them for access from the HCA. The advice MR verb is applicable for ODP MRs only.

For further information, please refer to the ibv_advise_mr manual page.

ODP Statistics

To aid in debugging and performance measurements and tuning, ODP support includes an extensive set of statistics.

For further information, please refer to rdma-statistics manual page.

Inline-Receive

The HCA may write received data to the Receive CQE. Inline-Receive saves PCIe Read transaction since the HCA does not need to read the scatter list. Therefore, it improves performance in case of short receive-messages.

On poll CQ, the driver copies the received data from CQE to the user's buffers.

Inline-Receive is enabled by default and is transparent to the user application. To disable it globally, set MLX5_SCATTER_TO_CQE environment variable to the value of 0. Otherwise, disable it on a specific QP using mlx5dv_create_qp() with MLX5DV_QP_CREATE_DISABLE_SCATTER_TO_CQE.

For further information, please refer to the manual page of mlx5dv_create_qp().

Mellanox PeerDirect®

Mellanox PeerDirect® uses an API between IB CORE and peer memory clients, (e.g. GPU cards) to provide access to an HCA to read/write peer memory for data buffers. As a result, it allows RDMA-based (over InfiniBand/RoCE) application to use peer device computing power, and RDMA interconnect at the same time without copying the data between the P2P devices. For example, Mellanox PeerDirect™ is being used for GPUDirect RDMA. Detailed description for that API exists under MLNX OFED installation, please see docs/readme_and_user_manual/PEER_MEMORY_API.txt

Mellanox PeerDirect Async

Mellanox PeerDirect Async sub-system gives PeerDirect hardware devices, such as GPU cards, dedicated AS accelerators, and so on, the ability to take control over HCA in critical path offloading CPU. To achieve this, there is a set of verb calls and structures providing application with abstract description of operation sequences intended to be executed by peer device.

Mellanox Relaxed Ordering (RSYNC)



⚠ This feature is only supported for ConnectX-5 adapter cards family and above.

In GPU systems with relaxed ordering, RSYNC callback will be invoked to ensure memory consistency. The registration and implementation of the callback will be done using an external module provided by the system vendor. Loading the module will register the callback in MLNX_OFED to be used later to guarantee memory operations order.

CPU Overhead Distribution

When creating a CQ using the <code>ibv_create_cq()</code> API, a "comp_vector" argument is sent. If the value set for this argument is 0, while the CPU core executing this verb is not equal to zero, the driver assigns a completion EQ with the least CQs reporting to it. This method is used to distribute CQs amongst available completions EQ. To assign a CQ to a specific EQ, the EQ needs to be specified in the comp_vector argument.

Out-of-Order (OOO) Data Placement



A This feature is only supported on:

- ConnectX-5 adapter cards and above
- RC and XRC QPs
- DC transport

Overview

In certain fabric configurations, InfiniBand packets for a given QP may take up different paths in a network from source to destination. This results into packets being received in an out-of-order manner. These packets can now be handled instead of being dropped, in order to avoid retransmission, by:

- · Achieving better network utilization
- Decreasing latency

Data will be placed into host memory in an out-of-order manner when out-of-order messages are received.

For information on how to set up out-of-order processing by the QP, please refer to HowTo Configure Adaptive Routing and SHIELD Community post.

IB Router

IB router provides the ability to send traffic between two or more IB subnets thereby potentially expanding the size of the network to over 40k end-ports, enabling separation and fault resilience between islands and IB subnets, and enabling connection to different topologies used by different subnets.

The forwarding between the IB subnets is performed using GRH lookup. The IB router's basic functionality includes:

- Removal of current L2 LRH (local routing header)
- Routing
- table lookup using GID from GRH
- · Building new LRH according to the destination according to the routing table

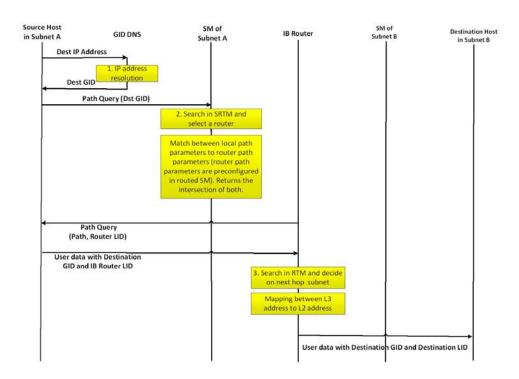
The DLID in the new LRH is built using simplified GID-to-LID mapping (where LID = 16 LSB bits of GID) thereby not requiring to send for ARP query/lookup.

Local Unicast GID Format



For this to work, the SM allocates an alias GID for each host in the fabric where the alias GID = {subnet prefix[127:64], reserved[63:16], LID[15:0}. Hosts should use alias GIDs in order to transmit traffic to peers on remote subnets.

Host-to-Host IB Router Unicast Flow



- For information on the architecture and functionality of IB Router, refer to IB Router Architecture and Functionality Community post.
- For information on IB Router configuration, refer to HowTo Configure IB Routers Community post.

MAD Congestion Control

The SA Management Datagrams (MAD) are General Management Packets (GMP) used to communicate with the SA entity within the InfiniBand subnet. SA is normally part of the subnet manager, and it is contained within a single active instance. Therefore, congestion on the SA communication level may occur.

Congestion control is done by allowing max_outstanding MADs only, where outstanding MAD means that is has no response yet. It also holds a FIFO queue that holds the SA MADs that their sending is delayed due to max_outstanding overflow.

The length of the queue is queue_size and meant to limit the FIFO growth beyond the machine memory capabilities. When the FIFO is full, SA MADs will be dropped, and the drops counter will increment accordingly.

When time expires (time_sa_mad) for a MAD in the queue, it will be removed from the queue and the user will be notified of the item expiration.

This features is implemented per CA port.

The SA MAD congestion control values are configurable using the following sysfs entries:

```
/sys/class/infiniband/mlx5_0/mad_sa_cc/
       drops
max_outstanding
       queue_size
time_sa_mad
     max outstanding
     queue_size
time_sa_mad
```

To print the current value:

```
cat /sys/class/infiniband/mlx5_0/mad_sa_cc/1/max_outstanding 16
```



To change the current value:

```
echo 32 > /sys/class/infiniband/mlx5_0/mad_sa_cc/1/max_outstanding cat /sys/class/infiniband/mlx5_0/mad_sa_cc/1/max_outstanding 32
```



To reset the drops counter:

```
echo 0 > /sys/class/infiniband/mlx5_0/mad_sa_cc/1/drops
```

Parameters' Valid Ranges

Parameter	Range		Default Values
	MIN	MAX	
max_oustanding	1	2^20	16
queue_size	16	2^20	16
time_sa_mad	1 milliseconds	10000	20 milliseconds

Storage Protocols

There are several storage protocols that use the advantage of InfiniBand and RDMA for performance reasons (high throughput, low latency and low CPU utilization). In this chapter we will discuss the following protocols:

- SCSI RDMA Protocol (SRP) is designed to take full advantage of the protocol off-load and RDMA features provided by the InfiniBand architecture.
- iSCSI Extensions for RDMA (iSER) is an extension of the data transfer model of iSCSI, a storage networking standard for TCP/IP. It uses the iSCSI components while taking the advantage of the RDMA protocol suite. ISER is implemented on various stor- age targets such as TGT, LIO, SCST and out of scope of this manual.
 - For various ISER targets configuration steps, troubleshooting and debugging, as well as other implementation of storage protocols over RDMA (such as Ceph over RDMA, nbdX and more) refer to Storage Solutions on Mellanox Community.
- Lustre is an open-source, parallel distributed file system, generally used for large-scale cluster computing that supports many requirements of leadership class HPC simulation environments.
- NVM Express[™] over Fabrics (NVME-oF)
 - NVME-oF is a technology specification for networking storage designed to enable NVMe message-based commands to transfer data between a host computer and a target solid-state storage device or system over a network such as Ethernet, Fibre Channel, and InfiniBand. Tunneling NVMe commands through an RDMA fabric provides a high throughput and a low latency. This is an alternative to the SCSi based storage networking protocols.
 - NVME-oF Target Offload is an implementation of the new NVME-oF standard Target (server) side in hardware. Starting from ConnectX-5 family cards, all regular IO requests can be processed by the HCA, with the HCA sending IO requests directly to a real NVMe PCI device, using peer-to-peer PCI communications. This means that excluding connection management and error flows, no CPU utilization will be observed during NVME-oF traffic.
 - For further information, please refer to Storage Solutions on Mellanox Community (https://community.mellanox.com).

SRP - SCSI RDMA Protocol

The 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 Initiator controls the connection to an SRP Target in order to provide access to remote storage devices across an InfiniBand fabric. The kSRP Target resides in an IO unit and provides storage services.

SRP Initiator

This SRP Initiator is based on open source from OpenFabrics (www.openfabrics.org) that implements the SCSI RDMA Protocol-2 (SRP-2). SRP-2 is described in Document # T10/1524-D available from http://www.t10.org.

The SRP Initiator supports

- Basic SCSI Primary Commands -3 (SPC-3) (www.t10.org/ftp/t10/drafts/spc3/spc3r21b.pdf)
- Basic SCSI Block Commands -2 (SBC-2) (www.t10.org/ftp/t10/drafts/sbc2/sbc2r16.pdf)
- · Basic functionality, task management and limited error handling



⚠ This package, however, does not include an SRP Target.

Loading SRP Initiator



To load the SRP module either:

- Execute the modprobe ib_srp command after the OFED driver is up.
- or
- 1. Change the value of SRP_LOAD in /etc/infiniband/openib.conf to "yes".
- 2. Run /etc/init.d/openibd restart for the changes to take effect.



⚠ When loading the ib_srp module, it is possible to set the module parameter srp_sg_tablesize. This is the maximum number of gather/scatter entries per I/O (default: 12).

SRP Module Parameters

When loading the SRP module, the following parameters can be set (viewable by the "modinfo ib_srp" command):

cmd_sg_entr	Default number of gather/scatter entries in the SRP command (default is 12, max 255)
allow_ext_s	Default behavior when there are more than cmd_sg_entries S/G entries after mapping; fails the request when false (default false)
topspin_wor karounds	Enable workarounds for Topspin/Cisco SRP target bugs
reconnect_d elay	Time between successive reconnect attempts. Time between successive reconnect attempts of SRP initiator to a disconnected target until dev_loss_tmo timer expires (if enabled), after that the SCSI target will be removed

fast_io_fai l_tmo	Number of seconds between the observation of a transport layer error and failing all I/O. Increasing this timeout allows more tolerance to transport errors, however, doing so increases the total failover time in case of serious transport failure. Note: fast_io_fail_tmo value must be smaller than the value of recon- nect_delay
dev_loss_tm o	Maximum number of seconds that the SRP transport should insulate transport layer errors. After this time has been exceeded the SCSI target is removed. Normally it is advised to set this to -1 (disabled) which will never remove the scsi_host. In deployments where different SRP targets are connected and disconnected frequently, it may be required to enable this timeout in order to clean old scsi_hosts representing targets that no longer exists

Constraints between parameters:

- dev_loss_tmo, fast_io_fail_tmo, reconnect_delay cannot be all disabled or negative values.
- reconnect_delay must be positive number.
- fast_io_fail_tmo must be smaller than SCSI block device timeout.
- fast_io_fail_tmo must be smaller than dev_loss_tmo.

SRP Remote Ports Parameters

Several SRP remote ports parameters are modifiable online on existing connection.



To modify dev loss tmo to 600 seconds:

echo 600 > /sys/class/srp_remote_ports/port-xxx/dev_loss_tmo



To modify fast_io_fail_tmo to 15 seconds:

echo 15 > /sys/class/srp_remote_ports/port-xxx/fast_io_fail_tmo



To modify reconnect_delay to 10 seconds:

echo 20 > /sys/class/srp_remote_ports/port-xxx/reconnect_delay

Manually Establishing an SRP Connection

The following steps describe how to manually load an SRP connection between the Initiator and an SRP Target. "Automatic Discovery and Connection to Targets" section explains how to do this automatically.

- Make sure that the ib_srp module is loaded, the SRP Initiator is reachable by the SRP Target, and that an SM is running.
- To establish a connection with an SRP Target and create an SRP (SCSI) device for that target under /dev, use the following command:

```
echo -n id_ext=[GUID value],ioc_guid=[GUID value],dgid=[port GID value],\
pkey=ffff,service_id=[service[0] value] > \
/sys/class/infiniband_srp/srp-mlx[hca number]-[port number]/add_target
```

See "SRP Tools - ibsrpdm, srp_daemon and srpd Service Script" section for instructions on how the parameters in this echo command may be obtained.

Notes:

- Execution of the above "echo" command may take some time
- The SM must be running while the command executes
- It is possible to include additional parameters in the echo command:
 - max_cmd_per_lun Default: 62
 - max_sect (short for max_sectors) sets the request size of a command
 - io_class Default: 0x100 as in rev 16A of the specification (In rev 10 the default was 0xff00)
 - tl_retry_count a number in the range 2..7 specifying the IB RC retry count. Default: 2
 - comp_vector, a number in the range 0..n-1 specifying the MSI-X completion vector.
 Some HCA's allocate multiple (n) MSI-X vectors per HCA port. If the IRQ affinity masks of these interrupts have been configured such that each MSI-X interrupt is handled by a different CPU then the comp_vector parameter can be used to spread the SRP completion workload over multiple CPU's.
 - cmd_sg_entries, a number in the range 1..255 that specifies the maximum number of data buffer descriptors stored in the SRP_CMD information unit itself. With allow_ext_sg=0 the parameter cmd_sg_entries defines the maximum S/G list length for a single SRP_CMD, and commands whose S/G list length exceeds this limit after S/G list collapsing will fail.
 - initiator_ext see "Multiple Connections from Initiator InfiniBand Port to the Target" section.
- To list the new SCSI devices that have been added by the echo command, you may use either of the following two methods:
 - Execute "fdisk -l". This command lists all devices; the new devices are included in this listing.
 - Execute "dmesg" or look at /var/log/messages to find messages with the names of the new devices.

SRP sysfs Parameters

Interface for making ib_srp connect to a new target. One can request ib_srp to connect to a new target by writing a comma-separated list of login parameters to this sysfs attribute. The supported parameters are:

id_ext	A 16-digit hexadecimal number specifying the eight byte identifier extension in the 16-byte SRP target port identifier. The target port identifier is sent by ib_srp to the target in the SRP_LOGIN_REQ request.
ioc_guid	A 16-digit hexadecimal number specifying the eight byte I/O controller GUID portion of the 16-byte target port identifier.
dgid	A 32-digit hexadecimal number specifying the destination GID.
pkey	A four-digit hexadecimal number specifying the InfiniBand partition key.
service_id	A 16-digit hexadecimal number specifying the InfiniBand service ID used to establish communication with the SRP target. How to find out the value of the service ID is specified in the documentation of the SRP target.
max_sect	A decimal number specifying the maximum number of 512-byte sectors to be transferred via a single SCSI command.
max_cmd_pe r_lun	A decimal number specifying the maximum number of outstanding commands for a single LUN.
io_class	A hexadecimal number specifying the SRP I/O class. Must be either 0xff00 (rev 10) or 0x0100 (rev 16a). The I/O class defines the format of the SRP initiator and target port identifiers.

initiator_ext	A 16-digit hexadecimal number specifying the identifier extension portion of the SRP initiator port identifier. This data is sent by the initiator to the target in the SRP_LOGIN_REQ request.
cmd_sg_entri es	A number in the range 1255 that specifies the maximum number of data buffer descriptors stored in the SRP_CMD information unit itself. With allow_ext_sg=0 the parameter cmd_sg_entries defines the maxi- mum S/G list length for a single SRP_CMD, and commands whose S/G list length exceeds this limit after S/G list collapsing will fail.
allow_ext_sg	Whether ib_srp is allowed to include a partial memory descriptor list in an SRP_CMD instead of the entire list. If a partial memory descriptor list has been included in an SRP_CMD the remaining memory descriptors are communicated from initiator to target via an additional RDMA transfer. Setting allow_ext_sg to 1 increases the maximum amount of data that can be transferred between initiator and target via a single SCSI command. Since not all SRP target implementations support partial memory descriptor lists the default value for this option is 0.
sg_tablesize	A number in the range 12048 specifying the maximum S/G list length the SCSI layer is allowed to pass to ib_srp. Specifying a value that exceeds cmd_sg_entries is only safe with partial memory descriptor list support enabled (allow_ext_sg=1).
comp_vector	A number in the range 0n-1 specifying the MSI-X completion vector. Some HCA's allocate multiple (n) MSI-X vectors per HCA port. If the IRQ affinity masks of these interrupts have been configured such that each MSI-X interrupt is handled by a different CPU then the comp_vector parameter can be used to spread the SRP completion workload over multiple CPU's.
tl_retry_coun t	A number in the range 27 specifying the IB RC retry count.

SRP Tools - ibsrpdm, srp_daemon and srpd Service Script

The OFED distribution provides two utilities: ibsrpdm and srp_daemon:

- They detect targets on the fabric reachable by the Initiator (Step 1)
- Output target attributes in a format suitable for use in the above "echo" command (Step 2)
- A service script srpd which may be started at stack startup

The utilities can be found under /usr/sbin/, and are part of the srptools RPM that may be installed using the Mellanox OFED installation. Detailed information regarding the various options for these utilities are provided by their man pages.

Below, several usage scenarios for these utilities are presented.

ibsrpdm

ibsrpdm has the following tasks:

- 1. Detecting reachable targets.
 - a. To detect all targets reachable by the SRP initiator via the default umad device (/sys/class/infiniband_mad/umad0), execute the following command:

ibsrpdm

This command will result into readable output information on each SRP Target detected. Sample:

b. To detect all the SRP Targets reachable by the SRP Initiator via another umad device, use the following command:

```
ibsrpdm -d <umad device>
```

- 2. Assisting in SRP connection creation.
 - a. To generate an output suitable for utilization in the "echo" command in "Manually Establishing an SRP Connection" section, add the '-c' option to ibsrpdm:

```
ibsrpdm -c
```

Sample output:

```
id_ext=200400A0B81146A1,ioc_guid=0002c90200402bd4,
dgid=fe80000000000000002c90200402bd5,pkey=ffff,service_id=200400a0b81146a1
```

b. To establish a connection with an SRP Target using the output from the 'ibsrpdm -c' example above, execute the following command:

```
echo -n id_ext=200400A0B81146A1,ioc_guid=0002c90200402bd4,
dgid=fe8000000000000000002c90200402bd5,pkey=ffff,service_id=200400a0b81146a1 > /sys/
class/infiniband_srp/srp-mlx5_0-1/add_target
```

The SRP connection should now be up; the newly created SCSI devices should appear in the listing obtained from the 'fdisk -l' command.

3. Discover reachable SRP Targets given an InfiniBand HCA name and port, rather than by just running /sys/class/infiniband_mad/umad<N> where <N> is a digit.

srpd

The srpd service script allows automatic activation and termination of the srp_daemon utility on all system live InfiniBand ports.

srp_daemon

srp_daemon utility is based on ibsrpdm and extends its functionality. In addition to the ibsrpdm functionality described above, srp_daemon can:

- Establish an SRP connection by itself (without the need to issue the "echo" command described in "Manually Establishing an SRP Connection" section)
- Continue running in background, detecting new targets and establishing SRP connections with them (daemon mode)
- Discover reachable SRP Targets given an infiniband HCA name and port, rather than just by / dev/umad<N> where <N> is a digit
- Enable High Availability operation (together with Device-Mapper Multipath)
- Have a configuration file that determines the targets to connect to:
- 1. srp_daemon commands equivalent to ibsrpdm:

```
"srp_daemon -a -o" is equivalent to "ibsrpdm"
"srp_daemon -c -a -o" is equivalent to "ibsrpdm -c"
```

Note: These srp_daemon commands can behave differently than the equivalent ibsrpdm command when /etc/srp_daemon.conf is not empty.

- 2. srp_daemon extensions to ibsrpdm.
 - To discover SRP Targets reachable from the HCA device < InfiniBand HCA name > and the port <port num>, (and to generate output suitable for 'echo'), execute:

```
host1# srp_daemon -c -a -o -i <InfiniBand HCA name> -p <port number>
```

Note: To obtain the list of InfiniBand HCA device names, you can either use the ibstat tool or run 'ls /sys/class/infiniband'.

- To both discover the SRP Targets and establish connections with them, just add the -e option to the above command.
- Executing srp_daemon over a port without the -a option will only display the reachable targets via the port and to which the initiator is not connected. If executing with the -e option it is better to omit -a.
- It is recommended to use the -n option. This option adds the initiator_ext to the connecting string (see "Multiple Connections from Initiator InfiniBand Port to the Target" section).
- srp_daemon has a configuration file that can be set, where the default is /etc/ srp_daemon.conf. Use the -f to supply a different configuration file that configures the targets srp_daemon is allowed to connect to. The configuration file can also be used to set values for additional parameters (e.g., max_cmd_per_lun, max_sect).
- A continuous background (daemon) operation, providing an automatic ongoing detection and connection capability. See "Automatic Discovery and Connection to Targets" section.

Automatic Discovery and Connection to Targets

- Make sure the ib_srp module is loaded, the SRP Initiator can reach an SRP Target, and that an SM is running.
- To connect to all the existing Targets in the fabric, run "srp_daemon -e -o". This utility will scan the fabric once, connect to every Target it detects, and then exit.



A srp_daemon will follow the configuration it finds in /etc/srp_daemon.conf. Thus, it will ignore a target that is disallowed in the configuration file.

- To connect to all the existing Targets in the fabric and to connect to new targets that will join the fabric, execute srp_daemon -e. This utility continues to execute until it is either killed by the user or encounters connection errors (such as no SM in the fabric).
- To execute SRP daemon as a daemon on all the ports:
 - srp_daemon.sh (found under /usr/sbin/). srp_daemon.sh sends its log to /var/log/ srp daemon.log.
 - Start the srpd service script, run service srpd start

For the changes in openib.conf to take effect, run:

/etc/init.d/openibd restart

Multiple Connections from Initiator InfiniBand Port to the Target

Some system configurations may need multiple SRP connections from the SRP Initiator to the same SRP Target: to the same Target IB port, or to different IB ports on the same Target HCA. In case of a single Target IB port, i.e., SRP connections use the same path, the configuration is enabled using a different initiator_ext value for each SRP connection. The initiator_ext value is a 16-hexadecimal-digit value specified in the connection command.

Also in case of two physical connections (i.e., network paths) from a single initiator IB port to two different IB ports on the same Target HCA, there is need for a different initiator_ext value on each path. The conventions is to use the Target port GUID as the initiator_ext value for the relevant path.

If you use srp_daemon with -n flag, it automatically assigns initiator_ext values according to this convention. For example:

Notes:

- It is recommended to use the -n flag for all srp_daemon invocations.
- ibsrpdm does not have a corresponding option.
- srp_daemon.sh always uses the -n option (whether invoked manually by the user, or automatically at startup by setting SRP_DAEMON_ENABLE to yes).

High Availability (HA)

High Availability works using the Device-Mapper (DM) multipath and the SRP daemon. Each initiator is connected to the same target from several ports/HCAs. The DM multipath is responsible for joining together different paths to the same target and for failover between paths when one of them goes offline. Multipath will be executed on newly joined SCSI devices.

Each initiator should execute several instances of the SRP daemon, one for each port. At startup, each SRP daemon detects the SRP Targets in the fabric and sends requests to the ib_srp module to connect to each of them. These SRP daemons also detect targets that subsequently join the fabric, and send the ib_srp module requests to connect to them as well.

Operation

When a path (from port1) to a target fails, the ib_srp module starts an error recovery process. If this process gets to the reset_host stage and there is no path to the target from this port, ib_srp will remove this scsi_host. After the scsi_host is removed, multipath switches to another path to this target (from another port/HCA).

When the failed path recovers, it will be detected by the SRP daemon. The SRP daemon will then request ib_srp to connect to this target. Once the connection is up, there will be a new scsi_host for this target. Multipath will be executed on the devices of this host, returning to the original state (prior to the failed path).

Manual Activation of High Availability

Initialization - execute after each boot of the driver:

- 1. Execute modprobe dm-multipath
- 2. Execute modprobe ib-srp
- 3. Make sure you have created file /etc/udev/rules.d/91-srp.rules as described above
- 4. Execute for each port and each HCA:

srp_daemon -c -e -R 300 -i <InfiniBand HCA name> -p <port number>

This step can be performed by executing srp_daemon.sh, which sends its log to /var/log/ srp_daemon.log.

Now it is possible to access the SRP LUNs on /dev/mapper/.



A It is possible for regular (non-SRP) LUNs to also be present; the SRP LUNs may be identified by their names. You can configure the /etc/multipath.conf file to change multipath behavior.



⚠ It is also possible that the SRP LUNs will not appear under /dev/mapper/. This can occur if the SRP LUNs are in the black-list of multipath. Edit the 'blacklist' section in /etc/ multipath.conf and make sure the SRP LUNs are not blacklisted.

Automatic Activation of High Availability

Start srpd service, run:

service srpd start

• From the next loading of the driver it will be possible to access the SRP LUNs on /dev/ mapper/



⚠ It is possible that regular (not SRP) LUNs are also present. SRP LUNs may be identified by their name.

• It is possible to see the output of the SRP daemon in /var/log/srp_daemon.log

Shutting Down SRP

SRP can be shutdown by using "rmmod ib srp", or by stopping the OFED driver ("/etc/init.d/ openibd stop"), or as a by-product of a complete system shutdown.

Prior to shutting down SRP, remove all references to it. The actions you need to take depend on the way SRP was loaded. There are three cases:

- 1. Without High Availability
 - When working without High Availability, you should unmount the SRP partitions that were mounted prior to shutting down SRP.
- 2. After Manual Activation of High Availability

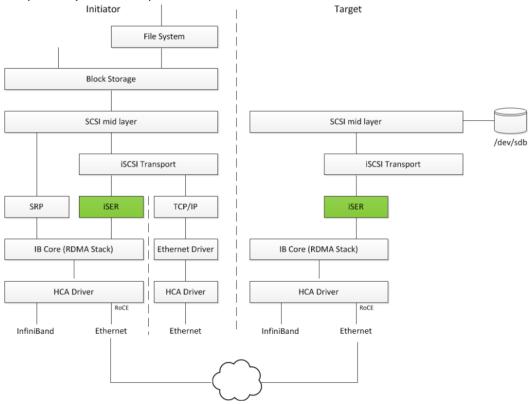
If you manually activated SRP High Availability, perform the following steps:

- a. Unmount all SRP partitions that were mounted.
- b. Stop service srpd (Kill the SRP daemon instances).
- c. Make sure there are no multipath instances running. If there are multiple instances, wait for them to end or kill them.
- d. Run: multipath -F
- 3. After Automatic Activation of High Availability

If SRP High Availability was automatically activated, SRP shutdown must be part of the driver shutdown ("/etc/init.d/openibd stop") which performs Steps 2-4 of case b above. However, you still have to unmount all SRP partitions that were mounted before driver shutdown.

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. iSER uses the RDMA protocol suite to supply higher bandwidth for block storage transfers (zero time copy behavior). To that fact, it eliminates the TCP/IP processing overhead while preserving the compatibility with iSCSI protocol.



There are three target implementation of ISER:

- Linux SCSI target framework (tgt)
- Linux-IO target (LIO)
- Generic SCSI target subsystem for Linux (SCST)

Each one of those targets can work in TCP or iSER transport modes.

iSER also supports RoCE without any additional configuration required. To bond the RoCE interfaces, set the fail_over_mac option in the bonding driver (see "Bonding IPoIB").

RDMA/RoCE is located below the iSER block on the network stack. In order to run iSER, the RDMA layer should be configured and validated (over Ethernet or InfiniBand). For troubleshooting RDMA, please refer to "HowTo Enable, Verify and Troubleshoot RDMA" on Mellanox Community.

iSER Initiator

The iSER initiator is controlled through the iSCSI interface available from the iscsi-initiator-utils package.

To discover and log into iSCSI targets, as well as access and manage the open-iscsi database use the iscasiadm utility, a command-line tool.

To enable iSER as a transport protocol use "I iser" as a parameter of the iscasiadm command.

Example for discovering and connecting targets over iSER:

```
iscsiadm -m discovery -o new -o old -t st -I iser -p <ip:port> -1
```

Note that the target implementation (e.g. LIO, SCST, TGT) does not affect he initiator process and configuration.

iSER Targets



Setting the iSER target is out of scope of this manual. For guidelines of how to do so, please refer to the relevant target documentation (e.g. stgt, targetcli).

Targets settings such as timeouts and retries are set the same as any other iSCSI targets.



⚠ If targets are set to auto connect on boot, and targets are unreachable, it may take a long time to continue the boot process if timeouts and max retries are set too high.

For various configuration, troubleshooting and debugging examples, refer to Storage Solutions on Mellanox Community.

Lustre

Lustre is an open-source, parallel distributed file system, generally used for large-scale cluster computing that supports many requirements of leadership class HPC simulation environments.

Lustre Compilation for MLNX OFED:



This procedure applies to RHEL/SLES OSs supported by Lustre. For further information, please refer to Lustre Release Notes.



To compile Lustre version 2.4.0 and higher:

```
./configure --with-o2ib=/usr/src/ofa_kernel/default/
```



To compile older Lustre versions:

```
$ EXTRA_LNET_INCLUDE="-I/usr/src/ofa_kernel/default/include/ -include /usr/src/ofa_kernel/default/include/linux/
compat-2.6.h" ./configure --with-o2ib=/usr/src/ofa_kernel/default/
$ EXTRA_LNET_INCLUDE="-I/usr/src/ofa_kernel/default/include/ -include /usr/src/ofa_kernel/default/include/linux/
compat-2.6.h" make rpms
```

For full installation example, refer to "HowTo Install Mellanox OFED driver for Lustre" Community post.

NVME-oF - NVM Express over Fabrics

NVME-oF

NVME-oF enables NVMe message-based commands to transfer data between a host computer and a target solid-state storage device or system over a network such as Ethernet, Fibre Channel, and InfiniBand. Tunneling NVMe commands through an RDMA fabric provides a high throughput and a low

For information on how to configure NVME-oF, please refer to the HowTo Configure NVMe over Fabrics Community post.

NVME-oF Target Offload



A This feature is only supported for ConnectX-5 adapter cards family and above.

NVME-oF Target Offload is an implementation of the new NVME-oF standard Target (server) side in hardware. Starting from ConnectX-5 family cards, all regular IO requests can be processed by the HCA, with the HCA sending IO requests directly to a real NVMe PCI device, using peer-to-peer PCI communications. This means that excluding connection management and error flows, no CPU utilization will be observed during NVME-oF traffic.

- For instructions on how to configure NVME-oF target offload, refer to HowTo Configure NVMEoF Target Offload Community post.
- For instructions on how to verify that NVME-oF target offload is working properly, refer to Simple NVMe-oF Target Offload Benchmark Community post.

Virtualization

The chapter contains the following sections:

- Single Root IO Virtualization (SR-IOV)
- Enabling Paravirtualization
- VXLAN Hardware Stateless Offloads
- Q-in-Q Encapsulation per VF in Linux (VST)
- 802.1Q Double-Tagging

Single Root IO Virtualization (SR-IOV)

Single Root IO Virtualization (SR-IOV) is a technology that allows a physical PCIe device to present itself multiple times through the PCIe bus. This technology enables multiple virtual instances of the device with separate resources. Mellanox adapters are capable of exposing up to 127 virtual instances (Virtual Functions (VFs) for each port in the Mellanox ConnectX® family cards. These virtual functions can then be provisioned separately. Each VF can be seen as an additional device connected to the Physical Function. It shares the same resources with the Physical Function, and its number of ports equals those of the Physical Function.

SR-IOV is commonly used in conjunction with an SR-IOV enabled hypervisor to provide virtual machines direct hardware access to network resources hence increasing its performance.

In this chapter we will demonstrate setup and configuration of SR-IOV in a Red Hat Linux environment using Mellanox ConnectX® VPI adapter cards family.

System Requirements

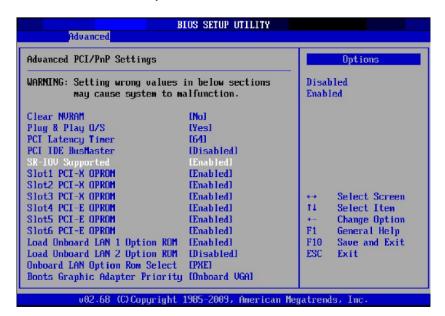
To set up an SR-IOV environment, the following is required:

- MLNX_OFED Driver
- A server/blade with an SR-IOV-capable motherboard BIOS
- Hypervisor that supports SR-IOV such as: Red Hat Enterprise Linux Server Version 6
- Mellanox ConnectX® VPI Adapter Card family with SR-IOV capability

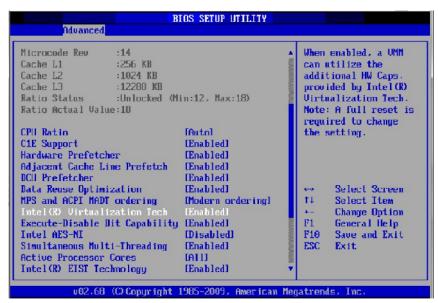
Setting Up SR-IOV

Depending on your system, perform the steps below to set up your BIOS. The figures used in this section are for illustration purposes only. For further information, please refer to the appropriate BIOS User Manual:

1. Enable "SR-IOV" in the system BIOS.



2. Enable "Intel Virtualization Technology".



- 3. Install a hypervisor that supports SR-IOV.
- 4. Depending on your system, update the /boot/grub/grub.conf file to include a similar command line load parameter for the Linux kernel. For example, to Intel systems, add:

Note: Please make sure the parameter "intel_iommu=on" exists when updating the /boot/grub/grub.conf file, otherwise SR-IOV cannot be loaded.

Some OSs use /boot/grub2/grub.cfg file. If your server uses such file, please edit this file instead (add "intel_iommu=on" for the relevant menu entry at the end of the line that starts with "linux16").

Configuring SR-IOV (Ethernet)

To set SR-IOV in Ethernet mode, refer to HowTo Configure SR-IOV for ConnectX-4/ConnectX-5/ ConnectX-6 with KVM (Ethernet) Community Post.

Configuring SR-IOV (InfiniBand)

- 1. Install the MLNX_OFED driver for Linux that supports SR-IOV.
- 2. Check if SR-IOV is enabled in the firmware.

```
mlxconfig -d /dev/mst/mt4115_pciconf0 q

Device #1:
------

Device type: Connect4
PCI device: /dev/mst/mt4115_pciconf0
Configurations: Current
SRIOV_EN 1
NUM_OF_VFS 8
```



If needed, use mlxconfig to set the relevant fields:

mlxconfig -d /dev/mst/mt4115_pciconf0 set SRIOV_EN=1 NUM_OF_VFS=16

- 3. Reboot the server.
- 4. Write to the sysfs file the number of Virtual Functions you need to create for the PF. You can use one of the following equivalent files:

You can use one of the following equivalent files:

A standard Linux kernel generated file that is available in the new kernels.

```
echo [num_vfs] > /sys/class/infiniband/mlx5_0/device/sriov_numvfs
```

Note: This file will be generated only if IOMMU is set in the grub.conf file (by adding intel_iommu=on, as seen in the fourth step under "Setting Up SR-IOV").

- A file generated by the mlx5_core driver with the same functionality as the kernel generated one.

```
echo [num_vfs] > /sys/class/infiniband/mlx5_0/device/mlx5_num_vfs
```

Note: This file is used by old kernels that do not support the standard file. In such kernels, using sriov_numvfs results in the following error: "bash: echo: write error: Function not implemented".

The following rules apply when writing to these files:

- If there are no VFs assigned, the number of VFs can be changed to any valid value (0 max **#VFs** as set during FW burning)
- If there are VFs assigned to a VM, it is not possible to change the number of VFs
- If the administrator unloads the driver on the PF while there are no VFs assigned, the driver will unload and SRI-OV will be disabled
- If there are VFs assigned while the driver of the PF is unloaded, SR-IOV will not be disabled. This means that VFs will be visible on the VM. However, they will not be operational. This is OSs with kernels that use pci_stub and not vfio.
 - The VF driver will discover this situation and will close its resources
- When the driver on the PF is reloaded, the VF becomes operational. The administrator of the VF will need to restart the driver in order to resume working with the VF.
- 5. Load the driver. To verify that the VFs were created. Run:

```
lspci | grep Mellanox
08:00.0 Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4]
08:00.1 Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4]
08:00.2 Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual Function]
08:00.3 Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual Function]
08:00.4 Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual Function]
08:00.5 Infiniband controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual Function]
```

6. Configure the VFs.

After VFs are created, 3 sysfs entries per VF are available under /sys/class/infiniband/ mlx5_<PF INDEX>/device/sriov (shown below for VFs 0 to 2):

```
+-- policy
+-- port
     node
     policy
port
```

For each Virtual Function, the following files are available:

- Node - Node's GUID:

The user can set the node GUID by writing to the /sys/class/infiniband/<PF>/device/sriov/ <index>/node file. The example below, shows how to set the node GUID for VF 0 of mlx5_0.

```
echo 00:11:22:33:44:55:1:0 > /sys/class/infiniband/mlx5_0/device/sriov/0/node
```

- Port - Port's GUID:

The user can set the port GUID by writing to the /sys/class/infiniband/<PF>/device/sriov/<index>/port file. The example below, shows how to set the port GUID for VF 0 of mlx5_0.

```
echo 00:11:22:33:44:55:2:0 > /sys/class/infiniband/mlx5_0/device/sriov/0/port
```

- Policy The vport's policy. The user can set the port GUID by writing to the /sys/class/infiniband/<PF>/device/sriov/<index>/port file. The policy can be one of:
 - Down the VPort PortState remains 'Down'
- Up if the current VPort PortState is 'Down', it is modified to 'Initialize'. In all other states, it is unmodified. The result is that the SM may bring the VPort up.
- Follow follows the PortState of the physical port. If the PortState of the physical port is 'Active', then the VPort implements the 'Up' policy. Otherwise, the VPort PortState is 'Down'. Notes:
- The policy of all the vports is initialized to "Down" after the PF driver is restarted except for VPort0 for which the policy is modified to 'Follow' by the PF driver.
- To see the VFs configuration, you must unbind and bind them or reboot the VMs if the VFs were assigned.
- 7. Make sure that OpenSM supports Virtualization (Virtualization must be enabled). The /etc/opensm/opensm.conf file should contain the following line:

```
virt_enabled 2
```

Note: OpenSM and any other utility that uses SMP MADs (ibnetdiscover, sminfo, iblink- info, smpdump, ibqueryerr, ibdiagnet and smpquery) should run on the PF and not on the VFs. In case of multi PFs (multi-host), OpenSM should run on HostO.

VFs Initialization Note

Since the same mlx5_core driver supports both Physical and Virtual Functions, once the Virtual Functions are created, the driver of the PF will attempt to initialize them so they will be available to the OS owning the PF. If you want to assign a Virtual Function to a VM, you need to make sure the VF is not used by the PF driver. If a VF is used, you should first unbind it before assigning to a VM.



To unbind a device use the following command:

1. Get the full PCI address of the device.

```
lspci -D

Example:
```

0000:09:00.2

2. Unbind the device.

echo 0000:09:00.2 > /sys/bus/pci/drivers/mlx5_core/unbind

3. Bind the unbound VF.

```
echo 0000:09:00.2 > /sys/bus/pci/drivers/mlx5_core/bind
```

PCI BDF Mapping of PFs and VFs

PCI addresses are sequential for both of the PF and their VFs. Assuming the card's PCI slot is 05:00 and it has 2 ports, the PFs PCI address will be 05:00.0 and 05:00.1. Given 3 VFs per PF, the VFs PCI addresses will be:

```
05:00.2-4 for VFs 0-2 of PF 0 (mlx5_0)
05:00.5-7 for VFs 0-2 of PF 1 (mlx5_1)
```

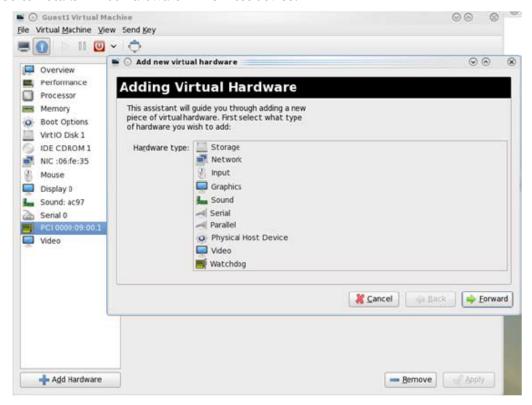
Additional SR-IOV Configurations

Assigning a Virtual Function to a Virtual Machine

This section describes a mechanism for adding a SR-IOV VF to a Virtual Machine.

Assigning the SR-IOV Virtual Function to the Red Hat KVM VM Server

- 1. Run the virt-manager.
- 2. Double click on the virtual machine and open its Properties.
- 3. Go to Details → Add hardware → PCI host device.



- 4. Choose a Mellanox virtual function according to its PCI device (e.g., 00:03.1)
- 5. If the Virtual Machine is up reboot it, otherwise start it.

6. Log into the virtual machine and verify that it recognizes the Mellanox card. Run:

```
lspci | grep Mellanox
```

Example:

```
lspci | grep Mellanox 01:00.0 Infiniband controller: Mellanox Technologies MT28800 Family [ConnectX-5 Ex]
```

7. Add the device to the /etc/sysconfig/network-scripts/ifcfg-ethX configuration file. The MAC address for every virtual function is configured randomly, therefore it is not necessary to add it.

Ethernet Virtual Function Configuration when Running SR-IOV

SR-IOV Virtual function configuration can be done through Hypervisor iprout2/netlink tool, if present. Otherwise, it can be done via sysfs.

```
ip link set { dev DEVICE | group DEVGROUP } [ { up | down } ]
...
[ vf NUM [ mac LLADDR ] [ vlan VLANID [ qos VLAN-QOS ] ]
...
[ spoofchk { on | off} ] ]
...
sysfs configuration (ConnectX-4):

/sys/class/net/enp8s0f0/device/sriov/[VF]

+-- [VF]
| +-- config
| +-- link_state
| +-- mac |
| +-- min_tx_rate
| +-- spoofcheck
| +-- stats
| +-- trunk
| +-- trust
```

VLAN Guest Tagging (VGT) and VLAN Switch Tagging (VST)

When running ETH ports on VGT, the ports may be configured to simply pass through packets as is from VFs (VLAN Guest Tagging), or the administrator may configure the Hypervisor to silently force packets to be associated with a VLAN/Qos (VLAN Switch Tagging).

In the latter case, untagged or priority-tagged outgoing packets from the guest will have the VLAN tag inserted, and incoming packets will have the VLAN tag removed. The default behavior is VGT.

To configure VF VST mode, run:

```
ip link set dev <PF device> vf <NUM> vlan <vlan_id> [qos <qos>]
```

where:

- NUM = 0..max-vf-num
- vlan_id = 0..4095
- qos = 0..7

For example:

• ip link set dev eth2 vf 2 vlan 10 qos 3 - sets VST mode for VF #2 belonging to PF eth2, with vlan_id = 10 and qos = 3

• ip link set dev eth2 vf 2 vlan 0 - sets mode for VF 2 back to VGT

Additional Ethernet VF Configuration Options

Guest MAC configuration - by default, guest MAC addresses are configured to be all zeroes. If
the administrator wishes the guest to always start up with the same MAC, he/she should
configure guest MACs before the guest driver comes up. The guest MAC may be configured by
using:

```
ip link set dev <PF device> vf <NUM> mac <LLADDR>
```

For legacy and ConnectX-4 guests, which do not generate random MACs, the administrator should always configure their MAC addresses via IP link, as above.

 Spoof checking - Spoof checking is currently available only on upstream kernels newer than 3.1.

```
ip link set dev <PF device> vf <NUM> spoofchk [on | off]
```

Guest Link State

```
ip link set dev <PF device> vf <UM> state [enable| disable| auto]
```

Virtual Function Statistics

Virtual function statistics can be queried via sysfs:

```
cat /sys/class/infiniband/mlx5_2/device/sriov/2/stats tx_packets : 5011
tx_bytes : 4450870
tx_dropped : 0
rx_packets : 5003
rx_bytes : 4450222
rx_broadcast : 0
rx_multicast : 0
tx_multicast : 0
tx_multicast : 8
rx_dropped : 0
```

Mapping VFs to Ports



To view the VFs mapping to ports:

Use the ip link tool v2.6.34~3 and above.

```
ip link
```

Output:

```
61: plp1: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT group default qlen 1000 link/ether 00:02:c9:f1:72:e0 brd ff:ff:ff:ff:ff:ff vf 0 MAC 00:00:00:00:00:00:00, vlan 4095, spoof checking off, link-state auto vf 37 MAC 00:00:00:00:00:00, vlan 4095, spoof checking off, link-state auto vf 38 MAC ff:ff:ff:ff:ff; vlan 65535, spoof checking off, link-state disable vf 39 MAC ff:ff:ff:ff:ff; vlan 65535, spoof checking off, link-state disable
```

When a MAC is ff:ff:ff:ff:ff; the VF is not assigned to the port of the net device it is listed under. In the example above, vf38 is not assigned to the same port as p1p1, in contrast to vf0.

However, even VFs that are not assigned to the net device, could be used to set and change its settings. For example, the following is a valid command to change the spoof check:

```
ip link set dev p1p1 vf 38 spoofchk on
```

This command will affect only the vf38. The changes can be seen in ip link on the net device that this device is assigned to.

Mapping VFs to Ports using the mlnx_get_vfs.pl tool



To map the PCI representation in BDF to the respective ports, run:

```
mlnx_get_vfs.pl
```

Output:

```
BDF 0000:04:00.0
Port 1:
                           0000:04:00.1
                  vf1
                          0000:04:00.2
                  vf2
vf3
         vf4 0000:04:00.5
```

RoCE Support

RoCE is supported on Virtual Functions and VLANs may be used with it. For RoCE, the hypervisor GID table size is of 16 entries while the VFs share the remaining 112 entries. When the number of VFs is larger than 56 entries, some of them will have GID table with only a single entry which is inadequate if VF's Ethernet device is assigned with an IP address.

Virtual Guest Tagging (VGT+)

VGT+ is an advanced mode of Virtual Guest Tagging (VGT), in which a VF is allowed to tag its own packets as in VGT, but is still subject to an administrative VLAN trunk policy. The policy determines which VLAN IDs are allowed to be transmitted or received. The policy does not determine the user priority, which is left unchanged.

Packets can be sent in one of the following modes: when the VF is allowed to send/receive untagged and priority tagged traffic and when it is not. No default VLAN is defined for VGT+ port. The send packets are passed to the eSwitch only if they match the set, and the received packets are forwarded to the VF only if they match the set.

Configuration



⚠ When working in SR-IOV, the default operating mode is VGT.



To enable VGT+ mode:

Set the corresponding port/VF (in the example below port eth5, VF0) range of allowed VLANs.

```
echo "<add> <start_vid> <end_vid>" > /sys/class/net/eth5/device/sriov/0/trunk
```

Examples:

• Adding VLAN ID range (4-15) to trunk:

```
echo add 4 15 > /sys/class/net/eth5/device/sriov/0/trunk
```

Adding a single VLAN ID to trunk:

```
echo add 17 17 > /sys/class/net/eth5/device/sriov/0/trunk
```

Note: When VLAN ID = 0, it indicates that untagged and priority-tagged traffics are allowed



To disable VGT+ mode, make sure to remove all VLANs.

echo rem 0 4095 > /sys/class/net/eth5/device/sriov/0/trunk



To remove selected VLANs.

• Remove VLAN ID range (4-15) from trunk:

```
echo rem 4 15 > /sys/class/net/eth5/device/sriov/0/trunk
```

• Remove a single VLAN ID from trunk:

echo rem 17 17 > /sys/class/net/eth5/device/sriov/0/trunk

SR-IOV Advanced Security Features

SR-IOV MAC Anti-Spoofing

Normally, MAC addresses are unique identifiers assigned to network interfaces, and they are fixed addresses that cannot be changed. MAC address spoofing is a technique for altering the MAC address to serve different purposes. Some of the cases in which a MAC address is altered can be legal, while others can be illegal and abuse security mechanisms or disguises a possible attacker.

The SR-IOV MAC address anti-spoofing feature, also known as MAC Spoof Check provides protection against malicious VM MAC address forging. If the network administrator assigns a MAC address to a VF (through the hypervisor) and enables spoof check on it, this will limit the end user to send traffic only from the assigned MAC address of that VF.

MAC Anti-Spoofing Configuration



▲ MAC anti-spoofing is disabled by default.

In the configuration example below, the VM is located on VF-0 and has the following MAC address: 11:22:33:44:55:66.

There are two ways to enable or disable MAC anti-spoofing:

- 1. Use the standard IP link commands available from Kernel 3.10 and above.
 - a. To enable MAC anti-spoofing, run:

ip link set ens785f1 vf 0 spoofchk on

b. To disable MAC anti-spoofing, run:

```
ip link set ens785f1 vf 0 spoofchk off
```

- 2. Specify echo "ON" or "OFF" to the file located under /sys/class/net/<ifname / device/sriov/ <VF index>/spoofcheck.
 - a. To enable MAC anti-spoofing, run:

```
echo "ON" > /sys/class/net/ens785f1/vf/0/spoofchk
```

b. To disable MAC anti-spoofing, run:

```
echo "OFF" > /sys/class/net/ens785f1/vf/0/spoofchk
```

A This configuration is non-persistent and does not survive driver restart.

Limit and Bandwidth Share Per VF

This feature enables rate limiting traffic per VF in SR-IOV mode. For details on how to configure rate limit per VF for ConnectX-4 and above adapter cards, please refer to HowTo Configure Rate Limit per VF for ConnectX-4/ConnectX-5/ConnectX-6 Community post.

Limit Bandwidth per Group of VFs

VFs Rate Limit for vSwitch (OVS) feature allows users to join available VFs into groups and set a rate limitation on each group. Rate limitation on a VF group ensures that the total Tx bandwidth that the VFs in this group get (altogether combined) will not exceed the given value.

With this feature, a VF can still be configured with an individual rate limit as in the past (under / sys/class/net/<ifname>/device/sriov/<vf_num>/max_tx_rate). However, the actual bandwidth limit on the VF will eventually be determined considering the VF group limitation and how many VFs are in the same group.

For example: 2 VFs (0 and 1) are attached to group 3.

Case 1: The rate limitation on the group is set to 20G. Rate limit of each VF is 15G Result: Each VF will have a rate limit of 10G

Case 2: Group's max rate limitation is still set to 20G. VF 0 is configured to 30G limit, while VF 1 is configured to 5G rate limit

Result: VF 0 will have 15G de-facto. VF 1 will have 5G

The rule of thumb is that the group's bandwidth is distributed evenly between the number of VFs in the group. If there are leftovers, they will be assigned to VFs whose individual rate limit has not been met yet.

VFs Rate Limit Feature Configuration

- 1. When VF rate group is supported by FW, the driver will create a new hierarchy in the SRI-OV sysfs named "groups" (/sys/class/net/<ifname>/device/sriov/groups/). It will contain all the info and the configurations allowed for VF groups.
- 2. All VFs are placed in group 0 by default since it is the only existing group following the initial driver start. It would be the only group available under /sys/class/net/<ifname>/device/
- 3. The VF can be moved to a different group by writing to the group file -> echo \$GROUP_ID > / sys/class/net/<ifname>/device/sriov/<vf_id>/group
- 4. The group IDs allowed are 0-255

- 5. Only when there is at least 1 VF in a group, there will be a group configuration available under /sys/class/net/<ifname>/device/sriov/groups/ (Except for group 0, which is always available even when it's empty).
- 6. Once the group is created (by moving at least 1 VF to that group), users can configure the group's rate limit. For example:
 - a. echo 10000 > /sys/class/net/<ifname>/device/sriov/5/max_tx_rate setting individual rate limitation of VF 5 to 10G (Optional)
 - b. echo 7 > /sys/class/net/<ifname>/device/sriov/5/group moving VF 5 to group 7
 - c. echo 5000 > /sys/class/net/<ifname>/device/sriov/groups/7/max_tx_rate setting group 7 with rate limitation of 5G
 - d. When running traffic via VF 5 now, it will be limited to 5G because of the group rate limit even though the VF itself is limited to 10G
 - e. echo 3 > /sys/class/net/<ifname>/device/sriov/5/group moving VF 5 to group 3
 - f. Group 7 will now disappear from /sys/class/net/<ifname>/device/sriov/groups since there are 0 VFs in it. Group 3 will now appear. Since there's no rate limit on group 3, VF 5 can transmit at 10G (thanks to its individual configuration)

Notes:

- You can see to which group the VF belongs to in the 'stats' sysfs (cat /sys/class/net/ <ifname>/device/sriov/<vf_num>/stats)
- You can see the current rate limit and number of attached VFs to a group in the group's 'config' sysfs (cat /sys/class/net/<ifname>/device/sriov/groups/<group_id>/config)

Bandwidth Guarantee per Group of VFs

Bandwidth guarantee (minimum BW) can be set on a group of VFs to ensure this group is able to transmit at least the amount of bandwidth specified on the wire.

Note the following:

- The minimum BW settings on VF groups determine how the groups share the total BW between themselves. It does not impact an individual VF's rate settings.
- The total minimum BW that is set on the VF groups should not exceed the total line rate.
 Otherwise, results are unexpected.
- It is still possible to set minimum BW on the individual VFs inside the group. This will determine how the VFs share the group's minimum BW between themselves. The total minimum BW of the VF member should not exceed the minimum BW of the group.

For instruction on how to create groups of VFs, see Limit Bandwidth per Group of VFs above.

Example

With a 40Gb link speed, assuming 4 groups and default group 0 have been created:

```
echo 20000 > /sys/class/net/<ifname>/device/sriov/group/1/min_tx_rate
echo 5000 > /sys/class/net/<ifname>/device/sriov/group/2/min_tx_rate
echo 15000 > /sys/class/net/<ifname>/device/sriov/group/3/min_tx_rate
```

```
Group 0(default): 0 - No BW guarantee is configured.
Group 1: 20000 - This is the maximum min rate among groups
Group 2: 5000 which is 25% of the maximum min rate
Group 3: 15000 which is 75% of the maximum min rate
Group 4: 0 - No BW guarantee is configured.
```

Assuming there are VFs attempting to transmit in full line rate in all groups, the results would look like: In which case, the minimum BW allocation would be:

```
Group0 - Will have no BW to use since no BW guarantee was set on it while other groups do have such settings.

Group1 - Will transmit at 20Gb/s

Group2 - Will transmit at 5Gb/s

Group3 - Will transmit at 15Gb/s

Group4 - Will have no BW to use since no BW guarantee was set on it while other groups do have such settings.
```

Privileged VFs

In case a malicious driver is running over one of the VFs, and in case that VF's permissions are not restricted, this may open security holes. However, VFs can be marked as trusted and can thus receive an exclusive subset of physical function privileges or permissions. For example, in case of allowing all VFs, rather than specific VFs, to enter a promiscuous mode as a privilege, this will enable malicious users to sniff and monitor the entire physical port for incoming traffic, including traffic targeting other VFs, which is considered a severe security hole.

Privileged VFs Configuration

In the configuration example below, the VM is located on VF-0 and has the following MAC address: 11:22:33:44:55:66.

There are two ways to enable or disable trust:

- 1. Use the standard IP link commands available from Kernel 4.5 and above.
 - a. To enable trust for a specific VF, run:

```
ip link set ens785f1 vf 0 trust on
```

b. To disable trust for a specific VF, run:

```
ip link set ens785f1 vf 0 trust off
```

- Specify echo "ON" or "OFF" to the file located under /sys/class/net/<ETH_IF_NAME> / device/ sriov/<VF index>/trust.
 - a. To enable trust for a specific VF, run:

```
echo "ON" > /sys/class/net/ens785f1/device/sriov/0/trust
```

b. To disable trust for a specific VF, run:

```
echo "OFF" > /sys/class/net/ens785f1/device/sriov/0/trust
```

Probed VFs

Probing Virtual Functions (VFs) after SR-IOV is enabled might consume the adapter cards' resources. Therefore, it is recommended not to enable probing of VFs when no monitoring of the VM is needed. VF probing can be disabled in two ways, depending on the kernel version installed on your server:

- 1. If the kernel version installed is v4.12 or above, it is recommended to use the PCI sysfs interface sriov_drivers_autoprobe. For more information, see linux-next branch.
- 2. If the kernel version installed is older than v4.12, it is recommended to use the mlx5_core module parameter probe_vf with MLNX_OFED v4.1 or above.

Example:

```
echo 0 > /sys/module/mlx5_core/parameters/probe_vf
```

For more information on how to probe VFs, see HowTo Configure and Probe VFs on mlx5 Drivers Community post.

VF Promiscuous Rx Modes

VF Promiscuous Mode

VFs can enter a promiscuous mode that enables receiving the unmatched traffic and all the multicast traffic that reaches the physical port in addition to the traffic originally targeted to the VF. The unmatched traffic is any traffic's DMAC that does not match any of the VFs' or PFs' MAC addresses.

Note: Only privileged/trusted VFs can enter the VF promiscuous mode.



To set the promiscuous mode on for a VF, run:

ifconfig eth2 promisc



To exit the promiscuous mode, run:

ifconfig eth2 -promisc

VF All-Multi Mode

VFs can enter an all-multi mode that enables receiving all the multicast traffic sent from/to the other functions on the same physical port in addition to the traffic originally targeted to the VF. Note: Only privileged/trusted VFs can enter the all-multi RX mode.



To set the all-multi mode on for a VF, run:

ifconfig eth2 allmulti



To exit the all-multi mode, run:

#ifconfig eth2 -allmulti

Uninstalling the SR-IOV Driver



To uninstall SR-IOV driver, perform the following:

- 1. For Hypervisors, detach all the Virtual Functions (VF) from all the Virtual Machines (VM) or stop the Virtual Machines that use the Virtual Functions. Please be aware that stopping the driver when there are VMs that use the VFs, will cause machine to hang.
- 2. Run the script below. Please be aware, uninstalling the driver deletes the entire driver's file, but does not unload the driver.

```
[root@sw1022 ~] # /usr/sbin/ofed_uninstall.sh
This program will uninstall all OFED packages on your machine.
Do you want to continue?[y/N]:y
Running /usr/sbin/vendor_pre_uninstall.sh
Removing OFED Software installations
Running /bin/rpm -e --allmatches kernel-ib kernel-ib-devel libibwerbs libibwerbs-devel libibwerbs-
devel-static libibwerbs-utils libmlx4 libmlx4-devel libibcm libibcm-devel libibumad libibumad-devel
libibumad-static libibmad libibmad-devel libibmad-static librdmacm-utils librdmacm-devel ibacm
opensm-libs opensm-devel perftest compat-dapl compat-dapl-devel dapl dapl-devel dapl-devel-static dapl-
utils srptools infiniband-diags-guest ofed-scripts opensm-devel
warning: /etc/infiniband/openib.conf saved as /etc/infiniband/openib.conf.rpmsave
Running /tmp/2818-ofed_vendor_post_uninstall.sh
```

3. Restart the server.

SR-IOV Live Migration



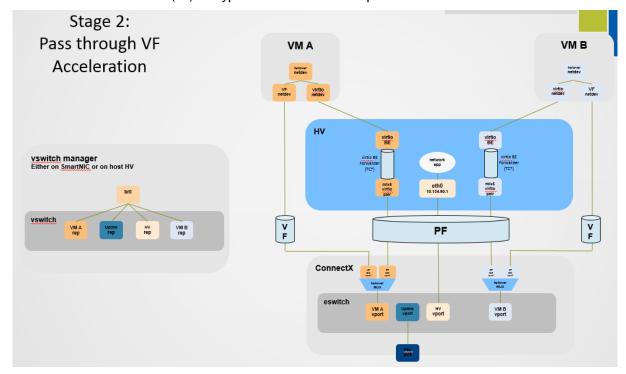
Support for this feature is at beta level.

Overview

This section describes how to set up and perform live migration on VMs with SR-IOV and with actively running traffic.

The below are the requirements for working with SR-IOV Live Migration.

- VM network persistency for applications VM's applications must survive the migration process
- No internal VM admin configuration
- Support for Mellanox ASAP² solution (kernel OVS hardware offload)
- No use of physical function (PF) network device in Paravirtual (PV) path for failover network
- Use of sub-function (SF) in HyperVisor as failover PV path



Prerequisites

- ConnectX-5 or higher adapter cards
- Hypervisor host with RedHat/CentOS minimal version of 8.0 with MLNX_OFED minimal version of 5.1
- VMs that run on CentOS v8.0 or higher
- Hypervisor host with latest libvirt from https://github.com/libvirt/libvirt.git
- Hypervisor host with atest gemu from https://github.com/gemu/gemu.git

This section consists of the following steps.

- 1. Host Servers Configuration.
- 2. VM OS Installation Using "virt-manager".
- 3. VFs to VMs Deployment.
- 4. Mellanox ASAP² with OVS Deployment.
- 5. Live Migration with Paravirtual Path and Traffic.

Host Servers Configuration

The following steps should be performed in both host servers.

- Install RedHat/CentOS v8.0 on the host server.
 The CentOS 8.0 ISO image can be downloaded via this link.
- 2. Connect the host servers to the Ethernet switch.

 Two host servers HV1 (eth0: 10.20.1.118) and HV2 (eth0: 10.20.1.119) connected via an Ethernet switch with switch ports configured/enabled VLAN. For example, vlan=100
- Install the latest MLNX_OFED version.
 Download and install Mellanox MLNX OFED driver for distribution RHEL/CentOS 8.0.

```
# mount -o loop MLNX_OFED_LINIx-4.7-3.2.0-rhel8.0-x86_64.iso /mnt
# cd /mnt
# ./mlnxofedinstall
# reboot
```

- 4. Configure the host server and Mellanox NIC with SR-IOV as instructed here.
- 5. Configure the host server and Mellanox NIC with sub-function by enabling it on BAR2.

```
# mlxconfig -d /dev/mst/mt4119_pciconf0 set PF_BAR2_ENABLE=1
# mlxconfig -d /dev/mst/mt4119_pciconf0 set PF_BAR2_SIZE=<0,1,2,3>

NOTES on PF_BAR2_SIZE:

0: 8 SFs support
1: 16 SFs support
2: 32 SFs support
Set the number of SFs same as that of VFs, as one SF-VF pair is used to attach to one VM.
```

6. Configure storage for VMs' images as shared.

The default location for VMs' images is /var/lib/libvirt/images, which is a shared location that is set up as an NFS directory in this PoC. For example:

```
# mount <nfs-server>:/opt/nfs/images /var/lib/libvirt/images
```

7. Set up the network bridge "installation" for VMs to enable external communication.

VMs must be able to download and install extra required packages for external sources.

```
# cd /etc/sysconfig/network-scripts
# vim ifcfg-installation

DEVICE=installation
    TYPE=Bridge
    BOOTPROTO=dhcp
    ONBOOT=yes
# vim ifcfg-eth0
    BRIDGE=installation
# systemctl network restart
```

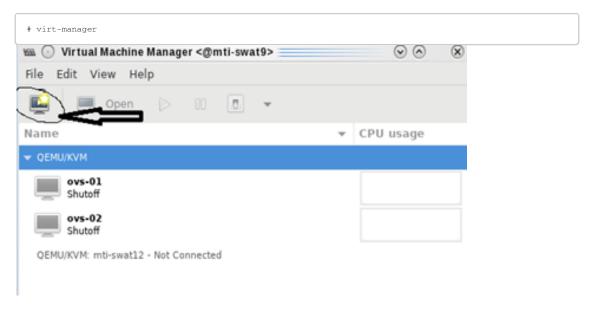
8. Download CentOS v8.0 ISO image for VM installation.

Download CentOS v8.0 ISO image from one of the mirror sites to the local host.

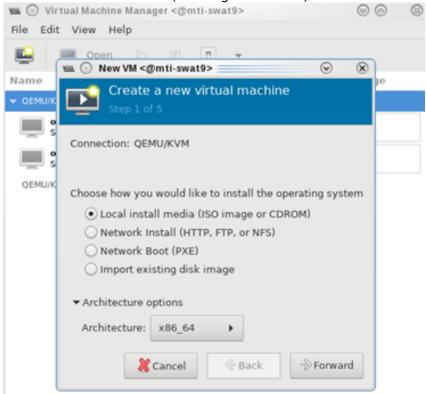
```
# wget
http://isoredirect.centos.org/centos/8/isos/x86_64/CentOS-8-x86_64-1905-dvd1.iso
```

VM OS installation Using "virt-manager"

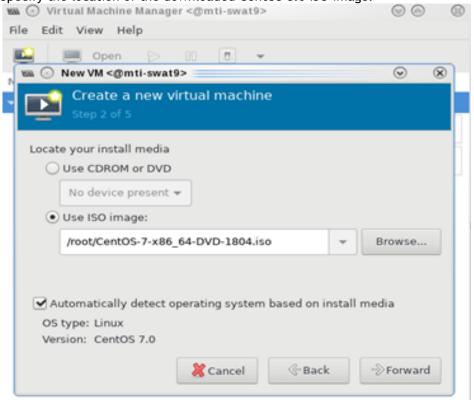
1. Launch "virt-manager" to create a new VM. Click the icon as shown below.



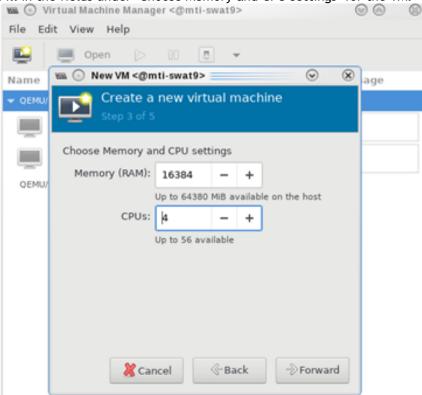
2. Choose "Local install media (ISO images or CDROM)".



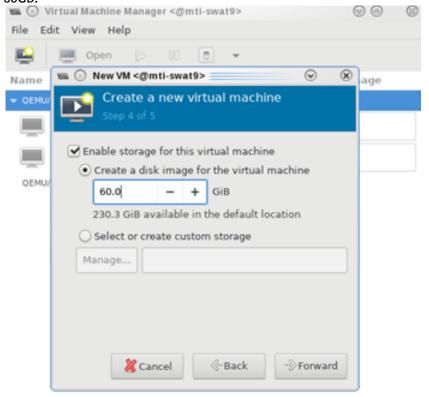
3. Specify the location of the downloaded CentOS 8.0 ISO image.



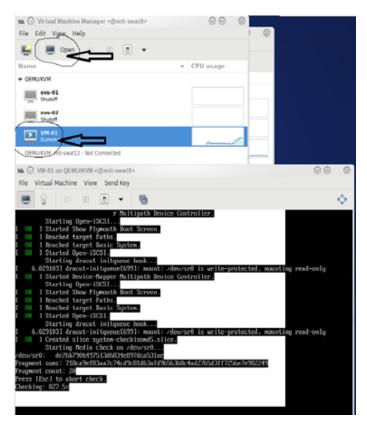
4. Fill in the fields under "Choose Memory and CPU settings" for the VM.



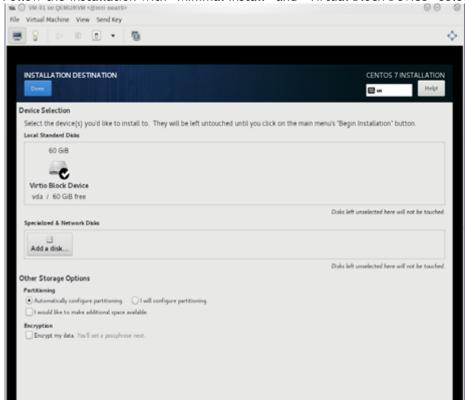
5. Create the disk image at the default root user location, for example: /var/lib/libvirt/images (NFS mount point). Make sure the storage is higher than 120GB and the virtual disk image is 60GB.



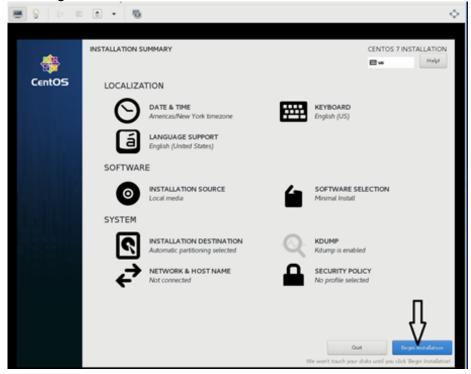
- 6. In the VM Name field, add "vm-01", and for the network selection, choose "Bridge installation: Host device eth0".
- 7. Click "vm-01", then "Open".



8. Follow the installation with "Minimal Install" and "Virtual Block Device" selection.



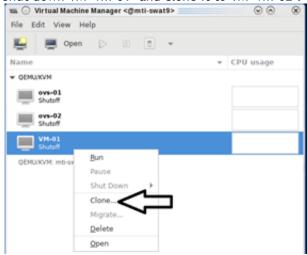
9. Click "Begin Installation".



- 10. Reboot VM "vm-01" after installation is completed.
- 11. Use the VM's console terminal to enable external communication.

[~]\$ vi /etc/sysconfig/network-scripts/ifcfg-eth0
 ONBOOT=yes
 BOOTPROTO=dhcp
systemctl network restart

12. Shut down VM "vm-01" and clone it to VM "vm-02".



13. Clone the virtual disk VM-01.qcow to VM-02.qcow.



- 14. Test the VM installation and migration without VFs.
 - a. Boot both VMs and run ping to each other.

```
[root@vm-01]# ping 10.20.4.8
[root@vm-02]# ping 10.20.4.20
```

b. Perform live migration using the "virsh" command line on HV1 where VM "vm-01" resides.

```
[root@HV1]# wirsh list --all
Id Name State
24 VM-01 running
```

c. Perform live migration.

```
[root@HV1]# virsh migrate --live --unsafe --persistent --verbose VM-01 qemu+ssh://HV2/system
```

d. Verify that vm-01 is migrated and resides at HV2.

```
[root@HV2]# wirsh list --all
Id Name State

1 VM-01 running
2 VM-02 running
```

VFs to VMs Deployment

1. Make sure SR-IOV is enabled and VFs are available.

```
[root@HV2]# lspci -D | grep Mellanox 0000:06:00.0 Ethernet controller: Mellanox Technologies MT27800 Family [ConnectX-5] 0000:06:00.1 Ethernet controller: Mellanox Technologies MT27800 Family [ConnectX-5 Virtual Function] 0000:06:00.2 Ethernet controller: Mellanox Technologies MT27800 Family [ConnectX-5 Virtual Function] 0000:06:00.3 Ethernet controller: Mellanox Technologies MT27800 Family [ConnectX-5 Virtual Function] 0000:06:00.4 Ethernet controller: Mellanox Technologies MT27800 Family [ConnectX-5 Virtual Function]
```

Enable the usage of VF2 "0000:06:00.3" and VF3 "0000:06:00.4" to assign to VM-01 and VM-02 respectively.

- 2. Attach VF to VM with XML file using "virsh" command line.
 - a. Create VM-01-vf.xml and VM-02-vf.xml files to assign VF1 to VM-01 and VF2 to VM-02 as "hostdev" with MAC-address assigned.

- b. Assign the VF to the VM by running the "virsh" command line.
 - Before attaching the VF, VM-01 and VM-02 should have a single network interface.

ii. On HV1 host, assign VF1 to VM-01.

```
[root@HV1]# virsh attach-device VM-01 VM-01-vf.xml Device attached successfully
```

iii. On HV2 host, assign VF2 to VM-02.

```
[root@HV2]# virsh attach-device VM-02 VM-02-vf.xml Device attached successfully
```

iv. After attaching the VFs, VM-01 and VM-02 should have two network interfaces. The second interface "ens12" is the VF with the MAC-address assigned.

c. Connect to VMs console, configure the IP address of the VF network interface and run traffic.

Install iperf, configure the IP address and run traffic between them on both VMs.

```
[root@VM-01]# yum -y install iperf3
...
Resolving Dependencies
--> Running transaction check
---> Package iperf3.x86_64 0:3.1.7-2.el7 will be installed
...
[root@VM-01]# ip addr add 11.11.11 dev ens12

[root@VM-02]# yum -y install iperf3
...
Resolving Dependencies
--> Running transaction check
---> Package iperf3.x86_64 0:3.1.7-2.el7 will be installed
...
[root@VM-02]# ip addr add 11.11.11.2 dev ens12
[root@VM-02]# ping 11.11.11.1
64 bytes from 11.11.11.1: icmp_seq=1 ttl=64 time=0.046 ms
64 bytes from 11.11.11.1: icmp_seq=2 ttl=64 time=0.039 ms
[root@VM-02]# iperf3 -s -f m
Server listening on 5201
```

```
[root@VM-01]# iperf3 -c 11.11.11.2 -P 4 -t 600 -i 1
```

Mellanox ASAP² with OVS Deployment

Perform the Mellanox ASAP² installation and configuration on both HV1 and HV2.

1. Download, build and install the latest IP-route2.

```
[root@HV1]# git clone git://git.kernel.org/pub/scm/linux/kernel/git/shemminger/iproute2.git
[root@HV1]# cd iproute2
[root@HV1]# ./configure
[root@HV1]# make -j32
[root@HV1]# make install
```

2. Download, build and install OpenvSwitch-2.12.0.

```
[root@HV1]# git clone http://github.com/openvswitch/ovs
[root@HV1]# cd ovs; git checkout -b v2.12.0 v2.12.0
[root@HV1]# ./boot.sh
[root@HV1]# ./configure
[root@HV1]# make -j32; make install
[root@HV1]# export PATH=$PATH:/usr/local/share/openvswitch/scripts
```

- Configure OVS with a single vSwitch "ovs-sriov" with hw-offload=true and tcpolicy=verbose.
 - $\textbf{a. Create OVS "ovs-sriov" and set } \verb|hw-offload=true| \textbf{and } \verb|tc-policy=verbose| \\$

```
[root@HV1]# export PATH=$PATH:/usr/local/share/openvswitch/scripts
[root@HV1]# ovs-ctl start

[root@HV1]# ovs-vsctl add-br ovs-sriov
[root@HV1]# ovs-vsctl set Open_vSwitch . other_config:hw-offload=true
[root@HV1]# ovs-vsctl set Open_vSwitch . other_config:tc-policy=verbose
[root@HV1]# ovs-ctl restart
[root@HV1]# ovs-vsctl get Open_vSwitch . other_config
{hw-offload="true", tc-policy=verbose}
```

Enable SR-IOV and SWITCHDEV mode by executing "asap_config.sh" script for PF port
 1.

c. Create a sub-function on PF port 1 with the script "create_sf.sh".

d. Rename the sub-function's netdevices.

```
Rename sub function's netdevices so that both source and destination systems can have same name.

[root@HV1]
$ ip link set eth0 down
$ ip link set eth0 name meth0
$ ip link set meth0 up

$ ip link set eth1 down
$ ip link set eth1 name meth1
$ ip link set meth1 up

$ ip link set eth2 down
$ ip link set eth2 name meth2
$ ip link set eth2 name meth2
$ ip link set meth2 up

$ ip link set meth3 down
$ ip link set eth3 name meth3
$ ip link set meth3 up
```

Live Migration with Paravirtual Path and Traffic

1. Create bonding devices of VF and sub-function (SF) representors.

```
NOTES:
=======
bond0 is bond device of sf1's representor (primary slave) and vf1's representor bond1 is bond device of sf2's representor (primary slave) and vf2's representor (protection of the primary slave) and vf2's representor (primary slave) and
```

2. Add Uplink "ens2" and bonding devices to "ovs-sriov" bridge.

```
[root@HV1]# ovs-vsctl add-port ovs-sriov ens2
[root@HV1]# ovs-vsctl add-port ovs-sriov bond0 tag=100
[root@HV1]# ovs-vsctl add-port ovs-sriov bond1 tag=100
```

- 3. Modify the VM-01's xml file to have the default SF's macvtap-based virtio netdevice.
 - a. Edit VM-01 configuration with the same MAC address assigned to VF-1.
 - b. Make sure the alias name has the prefix "ua-".

- c. Edit VM-02 configuration with the same MAC address assigned to VF-2.
- d. Make sure the alias name has the prefix "ua-".

4. Restart the VMs and verify that the PV path exists in both VMs and that it is accessible. Each VM should have two extra netdevs, such as: eth0, eth1 where eth0 is master and eth1 is automatically enslaved to eth0.

```
[root@VM-01] ip link show
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN mode DEFAULT qlen 1000
     link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00
2: ens3: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc pfifo_fast state UP mode DEFAULT qlen 1000
     link/ether 52:54:00:lc:91:0a brd ff:ff:fff:fff:ff
3: eth0: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT qlen 1000
     link/ether 52:54:00:53:53:53 brd fff:ff:fff:ff
4: eth1: <BROADCAST,MULTICAST> mtu 1500 qdisc noop master eth0 state DOWN mode DEFAULT qlen 1000
     link/ether 52:54:00:53:53:53:53 brd ff:ff:ff:ff:ff
```

a. Configure the IP address and run iperf on VMs over SF PF path.

```
[root@VM-01]# ip addr add 11.11.11.1 dev eth0
[root@VM-01]# iperf3 -s -f m
```

```
[root@VM-02]# ip addr add 11.11.11.2 dev eth0
[root@VM-02]# ping 11.11.11.1
[root@VM-02]# iperf3 -c 11.11.11.1 -P 4 -t 600 -i 1
```

b. Modify the XML file of the VFs to link to the persistent device of the SFs.

c. Attach VFs to VMs VM-01 and VM-02, and switch to the direct path in HyperVisors HV1 and HV2. I/O traffic should continue after the VFs have been successfully attached to the VMs.

```
[root@HV1] virsh attach-device VM-01 VM-01-vf.xml; echo ens2_0 > /sys/class/net/bond0/bonding/
active_slave
Device attached successfully
[root@HV2] virsh attach-device VM-02 VM-02-vf.xml; echo ens2_1 > /sys/class/net/bond1/bonding/
active_slave
Device attached successfully
```

Each VM should have one extra netdev from the attached VF that is automatically enslaved to eth0.

```
[root@VM-01] ip link show
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN mode DEFAULT qlen 1000
    link/loopback 00:00:00:00:00 brd 00:00:00:00:00
2: ens3: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc pfifo_fast state UP mode DEFAULT qlen
1000
    link/ether 52:54:00:1c:91:0a brd ff:ff:ff:ff:ff:
3: eth0: <BROADCAST,MULTICAST> mtu 1500 qdisc noop state DOWN mode DEFAULT qlen 1000
    link/ether 52:54:00:53:53:53 brd ff:ff:ff:ff:ff
4: eth1: <BROADCAST,MULTICAST> mtu 1500 qdisc noop master eth0 state DOWN mode DEFAULT qlen 1000
    link/ether 52:54:00:53:53:53 brd ff:ff:ff:ff:ff
5: ens12: <BROADCAST,MULTICAST> mtu 1500 qdisc noop master eth0 state DOWN mode DEFAULT qlen 1000
    link/ether 52:54:00:53:53:53 brd ff:ff:ff:ff:ff
```

5. Detach VF from VM, switch to SF PV path in HyperVisors HV1 and HV2. I/O traffic should pause 0.5s and then resume.

```
[root@HV1] virsh detach-device VM-01 VM-01-vf.xml; echo 06_00_0_32768 > /sys/class/net/bond0/bonding/active_slave
Device detached successfully
```

6. Perform Live Migration on VM-01. iperf traffic should run as usual.

```
[root@HV1] virsh migrate --live --unsafe --persistent --verbose VM-01 qemu+ssh://HV2/system
```

Attach VF to VM again and switch to the direct path in HyperVisor. I/O traffic should run as usual.

```
[root@HV2] \ virsh \ attach-device \ VM-01 \ VM-01-vf.xml; \ echo \ ens2\_0 \ > \ /sys/class/net/bond0/bonding/active\_slave \ Device \ attached \ successfully
```

Enabling Paravirtualization



⚠ The example below works on RHEL7.* without a Network Manager.

1. Create a bridge.

```
vim /etc/sysconfig/network-scripts/ifcfg-bridge0
DEVICE=bridge0
TYPE=Bridge
IPADDR=12.195.15.1
NETMASK=255.255.0.0
BOOTPROTO=static
ONBOOT-yes
NM_CONTROLLED=no
DELAY=0
```

2. Change the related interface (in the example below bridge0 is created over eth5).

```
DEVICE=eth5
BOOTPROTO-none
STARTMODE=on
HWADDR=00:02:c9:2e:66:52
TTYPE=Ethernet
NM_CONTROLLED=no
ONBOOT=yes
BRIDGE=bridge0
```

- 3. Restart the service network.
- 4. Attach a bridge to VM.

```
ifconfig -a
...
eth6
Link encap:Ethernet HWaddr 52:54:00:E7:77:99
inet addr:13.195.15.5 Bcast:13.195.255.255 Mask:255.255.0.0
inet6 addr: fe80::5054:fff:fee7:7799/64 Scope:Link
UP BROADCAST RUNNING MULTICAST MTU:1500 Metric:1
RX packets:481 errors:0 dropped:0 overruns:0 frame:0
TX packets:450 errors:0 dropped:0 overruns:0 carrier:0
collisions:0 txqueuelen:1000
RX bytes:22440 (21.9 KiB) TX bytes:19232 (18.7 KiB)
Interrupt:10 Base address:0xa000
```

VXLAN Hardware Stateless Offloads

VXLAN technology provides scalability and security challenges solutions. It requires extension of the traditional stateless offloads to avoid performance drop. ConnectX family cards offer the following

stateless offloads for a VXLAN packet, similar to the ones offered to non-encapsulated packets. VXLAN protocol encapsulates its packets using outer UDP header.

Available hardware stateless offloads:

- Checksum generation (Inner IP and Inner TCP/UDP)
- Checksum validation (Inner IP and Inner TCP/UDP)
- TSO support for inner TCP packets
- RSS distribution according to inner packets attributes
- Receive queue selection inner frames may be steered to specific QPs

Enabling VXLAN Hardware Stateless Offloads

VXLAN offload is enabled by default for ConnectX-4 family devices running the minimum required firmware version and a kernel version that includes VXLAN support.



To confirm if the current setup supports VXLAN, run:

```
ethtool -k $DEV | grep udp_tnl
```

Example:

```
ethtool -k ens1f0 | grep udp_tnl tx-udp_tnl-segmentation: on
```

ConnectX-4 family devices support configuring multiple UDP ports for VXLAN offload. Ports can be added to the device by configuring a VXLAN device from the OS command line using the "ip" command.

Note: If you configure multiple UDP ports for offload and exceed the total number of ports supported by hardware, then those additional ports will still function properly, but will not benefit from any of the stateless offloads.

Example:

```
ip link add vxlan0 type vxlan id 10 group 239.0.0.10 ttl 10 dev ens1f0 dstport 4789 ip addr add 192.168.4.7/24 dev vxlan0 ip link set up vxlan0
```

Note: dstport' parameters are not supported in Ubuntu 14.4.

The VXLAN ports can be removed by deleting the VXLAN interfaces.

Example:

ip link delete vxlan0



To verify that the VXLAN ports are offloaded, use debugfs (if supported):

1. Mount debugfs.

```
mount -t debugfs nodev /sys/kernel/debug
```

2. List the offloaded ports.

ls /sys/kernel/debug/mlx5/\$PCIDEV/VXLAN

Where \$PCIDEV is the PCI device number of the relevant ConnectX-4 family device. Example:

ls /sys/kernel/debug/mlx5/0000:81:00.0/VXLAN 4789

Important Note

VXLAN tunneling adds 50 bytes (14-eth + 20-ip + 8-udp + 8-vxlan) to the VM Ethernet frame. Please verify that either the MTU of the NIC who sends the packets, e.g. the VM virtio-net NIC or the host side veth device or the uplink takes into account the tunneling overhead. Meaning, the MTU of the sending NIC has to be decremented by 50 bytes (e.g 1450 instead of 1500), or the uplink NIC MTU has to be incremented by 50 bytes (e.g 1550 instead of 1500)

Q-in-Q Encapsulation per VF in Linux (VST)

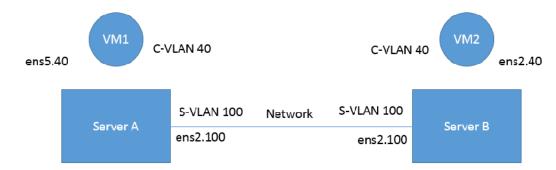
⚠ This feature is supported on ConnectX-5 and ConnectX-6 adapter cards only.

⚠ ConnectX-4 and ConnectX-4 Lx adapter cards support 802.1Q double-tagging (C-tag stacking on C-tag) - refer to "802.1Q Double-Tagging" section.

This section describes the configuration of IEEE 802.1ad QinQ VLAN tag (S-VLAN) to the hypervisor per Virtual Function (VF). The Virtual Machine (VM) attached to the VF (via SR- IOV) can send traffic with or without C-VLAN. Once a VF is configured to VST QinQ encapsulation (VST QinQ), the adapter's hardware will insert S-VLAN to any packet from the VF to the physical port. On the receive side, the adapter hardware will strip the S-VLAN from any packet coming from the wire to that VF.

Setup

The setup assumes there are two servers equipped with ConnectX-5/ConnectX-6 adapter cards.



Prerequisites

- Kernel must be of v3.10 or higher, or custom/inbox kernel must support vlan-stag
- Firmware version 16/20.21.0458 or higher must be installed for ConnectX-5/ConnectX-6 HCAs

- The server should be enabled in SR-IOV and the VF should be attached to a VM on the hypervisor.
 - In order to configure SR-IOV in Ethernet mode for ConnectX-5/ConnectX-6 adapter cards, please refer to "Configuring SR-IOV for ConnectX-4/ConnectX-5 (Ethernet)" section. In the following configuration example, the VM is attached to VFO.
- Network Considerations the network switches may require increasing the MTU (to support 1522 MTU size) on the relevant switch ports.

Configuring Q-in-Q Encapsulation per Virtual Function for ConnectX-5/ConnectX-6

- 1. Add the required S-VLAN (QinQ) tag (on the hypervisor) per port per VF. There are two ways to add the S-VLAN:
 - a. By using sysfs:

```
echo '100:0:802.1ad' > /sys/class/net/ens1f0/device/sriov/0/vlan
```

b. By using the ip link command (available only when using the latest Kernel version):

```
ip link set dev ens1f0 vf 0 vlan 100 proto 802.1ad
```

Check the configuration using the ip link show command:

2. Optional: Add S-VLAN priority. Use the gos parameter in the ip link command (or sysfs):

```
ip link set dev ens1f0 vf 0 vlan 100 qos 3 proto 802.1ad
```

Check the configuration using the ip link show command:

3. Create a VLAN interface on the VM and add an IP address.

```
ip link add link ens5 ens5.40 type vlan protocol 802.1q id 40 ip addr add 42.134.135.7/16 brd 42.134.255.255 dev ens5.40 ip link set dev ens5.40 up
```

4. To verify the setup, run ping between the two VMs and open Wireshark or tcpdump to capture the packet.

802.1Q Double-Tagging

This section describes the configuration of 802.1Q double-tagging support to the hypervisor per Virtual Function (VF). The Virtual Machine (VM) attached to the VF (via SR-IOV) can send traffic with

or without C-VLAN. Once a VF is configured to VST encapsulation, the adapter's hardware will insert C-VLAN to any packet from the VF to the physical port. On the receive side, the adapter hardware will strip the C-VLAN from any packet coming from the wire to that VF.

Configuring 802.1Q Double-Tagging per Virtual Function

- 1. Add the required C-VLAN tag (on the hypervisor) per port per VF. There are two ways to add the C-VLAN:
 - a. By using sysfs:

```
echo '100:0:802.1q' > /sys/class/net/ens1f0/device/sriov/0/vlan
```

b. By using the ip link command (available only when using the latest Kernel version):

```
ip link set dev ens1f0 vf 0 vlan 100
```

Check the configuration using the ip link show command:

```
# ip link show ens1f0
ens1f0: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 gdisc mq state UP mode DEFAULT glen 1000
link/ether ec:0d:9a:44:37:84 brd ff:ff:ff:ff:ff:ff
    vf 0 MAC 00:00:00:00:00:00:00, vlan 100, spoof checking off, link-state auto, trust off
    vf 1 MAC 00:00:00:00:00:00, spoof checking off, link-state auto, trust off
    vf 2 MAC 00:00:00:00:00:00, spoof checking off, link-state auto, trust off
    vf 3 MAC 00:00:00:00:00:00, spoof checking off, link-state auto, trust off
    vf 4 MAC 00:00:00:00:00:00:00, spoof checking off, link-state auto, trust off
```

2. Create a VLAN interface on the VM and add an IP address.

```
# ip link add link ens5 ens5.40 type vlan protocol 802.1q id 40
# ip addr add 42.134.135.7/16 brd 42.134.255.255 dev ens5.40
# ip link set dev ens5.40 up
```

3. To verify the setup, run ping between the two VMs and open Wireshark or tcpdump to capture the packet.

Resiliency

The chapter contains the following sections:

Reset Flow

Reset Flow

Reset Flow is activated by default. Once a "fatal device" error is recognized, both the HCA and the software are reset, the ULPs and user application are notified about it, and a recovery process is performed once the event is raised.

Currently, a reset flow can be triggered by a firmware assert with Recover Flow Request (RFR) only. Firmware RFR support should be enabled explicitly using mlxconfig commands.



To query the current value, run:

```
mlxconfig -d /dev/mst/mt4115_pciconf0 query | grep SW_RECOVERY_ON_ERRORS
```



To enable RFR bit support, run:

mlxconfig -d /dev/mst/mt4115_pciconf0 set SW_RECOVERY_ON_ERRORS=true

Kernel ULPs

Once a "fatal device" error is recognized, an IB_EVENT_DEVICE_FATAL event is created, ULPs are notified about the incident, and outstanding WQEs are simulated to be returned with "flush in error" message to enable each ULP to close its resources and not get stuck via calling its "remove_one" callback as part of "Reset Flow".

Once the unload part is terminated, each ULP is called with its "add one" callback, its resources are re-initialized and it is re-activated.

User Space Applications (IB/RoCE)

Once a "fatal device" error is recognized an IB_EVENT_DEVICE_FATAL event is created, applications are notified about the incident and relevant recovery actions are taken.

Applications that ignore this event enter a zombie state, where each command sent to the kernel is returned with an error, and no completion on outstanding WQEs is expected.

The expected behavior from the applications is to register to receive such events and recover once the above event is raised. Same behavior is expected in case the NIC is unbounded from the PCI and later is rebounded. Applications running over RDMA CM should behave in the same manner once the RDMA_CM_EVENT_DEVICE_REMOVAL event is raised.

The below is an example of using the unbind/bind for NIC defined by "0000:04:00.0"

echo 0000:04:00.0 > /sys/bus/pci/drivers/mlx5_core/unbind echo 0000:04:00.0 > /sys/bus/pci/drivers/mlx5_core/bind

SR-IOV

If the Physical Function recognizes the error, it notifies all the VFs about it by marking their communication channel with that information, consequently, all the VFs and the PF are reset. If the VF encounters an error, only that VF is reset, whereas the PF and other VFs continue to work unaffected.

Forcing the VF to Reset

If an outside "reset" is forced by using the PCI sysfs entry for a VF, a reset is executed on that VF once it runs any command over its communication channel.

For example, the below command can be used on a hypervisor to reset a VF defined by 0000:04:00.1:

echo 1 >/sys/bus/pci/devices/0000:04:00.1/reset

Extended Error Handling (EEH)

Extended Error Handling (EEH) is a PowerPC mechanism that encapsulates AER, thus exposing AER events to the operating system as EEH events.

The behavior of ULPs and user space applications is identical to the behavior of AER.

CRDUMP

CRDUMP feature allows for taking an automatic snapshot of the device CR-Space in case the device's FW/HW fails to function properly.

Snapshots Triggers:

The snapshot is triggered after firmware detects a critical issue, requiring a recovery flow.

This snapshot can later be investigated and analyzed to track the root cause of the failure. Currently, only the first snapshot is stored, and is exposed using a temporary virtual file. The virtual file is cleared upon driver reset.

When a critical event is detected, a message indicating CRDUMP collection will be printed to the Linux log. User should then back up the file pointed to in the printed message. The file location format is: /proc/driver/mlx5_core/crdump/<pci address>

Snapshot should be copied by Linux standard tool for future investigation.

Firmware Tracer

This mechanism allows for the device's FW/HW to log important events into the event tracing system (/sys/kernel/debug/tracing) without requiring any Mellanox tool.



⚠ To be able to use this feature, trace points must be enabled in the kernel.

This feature is enabled by default, and can be controlled using sysfs commands.



To disable the feature:

echo 0 > /sys/kernel/debug/tracing/events/mlx5/fw_tracer/enable



To enable the feature:

echo 1 > /sys/kernel/debug/tracing/events/mlx5/fw_tracer/enable



To view FW traces using vim text editor:

vim /sys/kernel/debug/tracing/trace

Docker Containers

On Linux, Docker uses resource isolation of the Linux kernel, to allow independent "containers" to run within a single Linux kernel instance.

Docker containers are supported on MLNX_OFED using Docker runtime. Virtual RoCE and InfiniBand devices are supported using SR-IOV mode.

Currently, RDMA/RoCE devices are supported in the modes listed in the following table.

Linux Containers Networking Modes

Orchestration and Clustering Tool	Version	Networking Mode	Link Layer	Virtualizati on Mode
Docker	Docker Engine 17.03 or higher	SR-IOV using sriov-plugin along with docker run wrapper tool	InfiniBand and Ethernet	SR-IOV
Kubernetes	Kubernetes 1.10.3 or higher	SR-IOV using device plugin, and using SR- IOV CNI plugin	InfiniBand and Ethernet	SR-IOV
		VXLAN using IPoIB bridge	InfiniBand	Shared HCA

Docker Using SR-IOV

In this mode, Docker engine is used to run containers along with SR-IOV networking plugin. To isolate the virtual devices, docker_rdma_sriov tool should be used. This mode is applicable to both InfiniBand and Ethernet link layers.

To obtain the plugin, visit: https://hub.docker.com/r/mellanox/sriov-plugin/

To install the docker_rdma_sriov tool, use the container tools installer available via https://hub.docker.com/r/mellanox/container_tools_installer/

For instructions on how to use Docker with SR-IOV, refer to <u>Docker RDMA SRIOV Networking with ConnectX4/ConnectX5/ConnectX6 Community post</u>.

Kubernetes Using SR-IOV

In order to use RDMA in Kubernetes environment with SR-IOV networking mode, two main components are required:

- 1. RDMA device plugin this plugin allows for exposing RDMA devices in a Pod
- 2. SR-IOV CNI plugin this plugin provisions VF net device in a Pod

When used in SR-IOV mode, this plugin enables SR-IOV and performs necessary configuration including setting GUID, MAC, privilege mode, and Trust mode.

The plugin also allocates the VF devices when Pods are scheduled and requested by Kubernetes framework.

Kubernetes with Shared HCA

One RDMA device (HCA) can be shared among multiple Pods running in a Kubernetes worker nodes. User defined networks are created using VXLAN or VETH networking devices. RDMA device (HCA) can be shared among multiple Pods running in a Kubernetes worker nodes.

Mediated Devices

The Mellanox mediated devices deliver flexibility in allowing to create accelerated devices without SR-IOV on the Bluefield® system. These mediated devices support NIC and RDMA, and offer the same level of ASAP² offloads as SR-IOV VFs. Mediates devices are supported using mlx5 sub-function acceleration technology.

Configuring Mediated Device

1. To support sub-functions, PCIe BAR2 must be enabled. Run:

```
$ mlxconfig -d /dev/mst/mst41682_pciconf0 s PF_BAR2_SIZE 4 PF_BAR2_ENABLE=True
```

Cold reboot the BlueField host system so that the above settings can be applied on subsequent reboot sequences.

2. By default, the firmware allows for a large number of maximum mdev devices. You must set the maximum number of mediated devices to 2 or 4 devices. Run:

```
\$ \  \, echo \  \, 4 \  \, > \  \, /sys/bus/pci/devices/0000:05:00.0/mdev\_supported\_types/mlx5\_core-local/max\_mdevs
```

3. Mediated devices are uniquely identified using UUID. To create one, run:

```
$ uuidgen
$ echo 49d0e9ac-61b8-4c91-957e-6f6dbc42557d > /sys/bus/pci/devices/0000:05:00.0/mdev_supported_types/
mlx5_core-local/create
```

4. By default, if the driver vfio_mdev is loaded, newly created mdev devices are bound to it. To make use of this newly created mdev device in order to create a netdevice and RDMA device, you must first unbind it from that driver. Run:

```
$ echo 49d0e9ac-61b8-4c91-957e-6f6dbc42557d > /sys/bus/mdev/drivers/vfio_mdev/unbind
```

5. Configure a MAC address for the mdev device. Run:

```
\label{localized} $$ echo 00:11:22:33:44:55 > /sys/bus/mdev/devices/49d0e9ac-61b8-4c91-957e-6f6dbc42557d/devlink-compatconfig/mac_addr
```

6. Query the representor netdevice of the mdev device. Run:

```
$ cat /sys/bus/mdev/devices/49d0e9ac-61b8-4c91-957e-6f6dbc42557d/devlink-compat-config/netdev
```

7. Bind the mediated device to mlx5_core driver. Run:

```
$ echo 49d0e9ac-61b8-4c91-957e-6f6dbc42557d > /sys/bus/mdev/drivers/mlx5_core/bind
```

When an mdev device is bound to the mlx5_core driver, its respective netdevice and/or RDMA device is also created.

8. To inspect the netdevice and RDMA device for the mdev, run:

```
$ ls /sys/bus/mdev/devices/49d0e9ac-61b8-4c91-957e-6f6dbc42557d/net/
$ ls /sys/bus/mdev/devices/49d0e9ac-61b8-4c91-957e-6f6dbc42557d/infiniband/
```

HPC-X™

For information on HPC-X, please refer to HPC-X $^{\text{m}}$ User Manual at www.docs.mellanox.com, under Software \rightarrow HPC-X.

Fast Driver Unload

This feature enables optimizing mlx5 driver teardown time in shutdown and kexec flows.

The fast driver unload is disabled by default. To enable it, the prof_sel module parameter of mlx5_core module should be set to 3.

OVS Offload Using ASAP² Direct

Overview



Supported on ConnectX-5 and above adapter cards.

Open vSwitch (OVS) allows Virtual Machines (VMs) to communicate with each other and with the outside world. OVS traditionally resides in the hypervisor and switching is based on twelve tuple matching on flows. The OVS software based solution is CPU intensive, affecting system performance and preventing full utilization of the available bandwidth.

Mellanox Accelerated Switching And Packet Processing (ASAP²) technology allows OVS offloading by handling OVS data-plane in Mellanox ConnectX-5 onwards NIC hardware (Mellanox Embedded Switch or eSwitch) while maintaining OVS control-plane unmodified. As a result, we observe significantly higher OVS performance without the associated CPU load.

As of v5.0, OVS-DPDK became part of MLNX_OFED package as well. OVS-DPDK supports ASAP² just as the OVS-Kernel (Traffic Control (TC) kernel-based solution) does, yet with a different set of features.

The traditional ASAP² hardware data plane is built over SR-IOV virtual functions (VFs), so that the VF is passed through directly to the VM, with the Mellanox driver running within the VM. An alternate approach that is also supported is vDPA (vhost Data Path Acceleration). vDPA allows the connection to the VM to be established using VirtIO, so that the data-plane is built between the SR-IOV VF and the standard VirtIO driver within the VM, while the control-plane is managed on the host by the vDPA application. Two flavors of vDPA are supported, Software vDPA; and Hardware vDPA. Software vDPA management functionality is embedded into OVS-DPDK, while Hardware vDPA uses a standalone application for management, and can be run with both OVS-Kernel and OVS-DPDK. For further information, please see sections VirtIO Acceleration through VF Relay (Software vDPA) and VirtIO Acceleration through Hardware vDPA.

Installing OVS-Kernel ASAP² Packages

Install the required packages. For the complete solution, you need to install supporting MLNX_OFED (v4.4 and above), iproute2, and openvswitch packages.

Installing OVS-DPDK ASAP² Packages

```
./mlnxofedinstall --ovs-dpdk -upstream-libs
```

Setting Up SR-IOV

Note that this section applies to both OVS-DPDK and OVS-Kernel similarly.

Unknown Attachment

To set up SR-IOV:

1. Choose the desired card.

The example below shows a dual-ported ConnectX-5 card (device ID 0x1017) and a single SR-IOV VF (Virtual Function, device ID 0x1018).

In SR-IOV terms, the card itself is referred to as the PF (Physical Function).

```
# lspci -nn | grep Mellanox
0a:00.0 Ethernet controller [0200]: Mellanox Technologies MT27800 Family [ConnectX-5] [15b3:1017] 0a:00.1 Ethernet controller [0200]: Mellanox Technologies MT27800 Family [ConnectX-5] [15b3:1017]
Oa:00.2 Ethernet controller [0200]: Mellanox Technologies MT27800 Family [ConnectX-5 Virtual Function] [15b3:1018]
```

Enabling SR-IOV and creating VFs is done by the firmware upon admin directive as explained in Step 5 below.

Identify the Mellanox NICs and locate net-devices which are on the NIC PCI BDF.

```
# 1s -1 /sys/class/net/ | grep 04:00
lrwxrwxrwx 1 root root 0 Mar 27 16:58 enp4s0f0 -> ../../devices/pci0000:00/0000:00:03.0/0000:04:00.0/net/
Irwxrwxrwx 1 root root 0 Mar 27 16:58 enp4s0f1 -> ../../devices/pci0000:00/0000:00:03.0/0000:04:00.1/net/enp4s0f1
ENDARGE INVESTIGATION OF TOOL 0 Mar 27 16:58 eth0 -> ../../devices/pci0000:00/0000:03.0/0000:04:00.2/net/eth0 lrwxrwxrwx 1 root root 0 Mar 27 16:58 eth1 -> ../../devices/pci0000:00/0000:03.0/0000:04:00.3/net/eth1
```

The PF NIC for port #1 is enp4s0f0, and the rest of the commands will be issued on it.

3. Check the firmware version.

Make sure the firmware versions installed are as state in the Release Notes document.

```
# ethtool -i enp4s0f0 | head -5
driver: mlx5_core
version: 5.0-5
firmware-version: 16.21.0338
expansion-rom-version:
bus-info: 0000:04:00.0
```

4. Make sure SR-IOV is enabled on the system (server, card).

Make sure SR-IOV is enabled by the server BIOS, and by the firmware with up to N VFs, where N is the number of VFs required for your environment. Refer to "Mellanox Firmware Tools" below for more details.

```
# cat /sys/class/net/enp4s0f0/device/sriov_totalvfs
```

5. Turn ON SR-IOV on the PF device.

```
# echo 2 > /sys/class/net/enp4s0f0/device/sriov_numvfs
```

6. Provision the VF MAC addresses using the IP tool.

```
# ip link set enp4s0f0 vf 0 mac e4:11:22:33:44:50
# ip link set enp4s0f0 vf 1 mac e4:11:22:33:44:51
```

Verify the VF MAC addresses were provisioned correctly and SR-IOV was turned ON.

```
# cat /sys/class/net/enp4s0f0/device/sriov_numvfs
# ip link show dev enp4s0f0
256: enp4s0f0: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc mq master ovs-system state UP mode DEFAULT
group default qlen 1000
    link/ether e4:1d:2d:60:95:a0 brd ff:ff:ff:ff:ff
    vf 0 MAC e4:11:22:33:44:50, spoof checking off, link-state auto
    vf 1 MAC e4:11:22:33:44:51, spoof checking off, link-state auto
```

In the example above, the maximum number of possible VFs supported by the firmware is 4 and only 2 are enabled.

8. Provision the PCI VF devices to VMs using PCI Pass-Through or any other preferred virt tool of choice, e.g virt-manager.

For further information on SR-IOV, refer to https://community.mellanox.com/docs/DOC-2386.

OVS Hardware Offloads Configuration

OVS-Kernel Hardware Offloads

SwitchDev Configuration

1. Unbind the VFs.

```
echo 0000:04:00.2 > /sys/bus/pci/drivers/mlx5_core/unbind echo 0000:04:00.3 > /sys/bus/pci/drivers/mlx5_core/unbind
```

VMs with attached VFs must be powered off to be able to unbind the VFs.

2. Change the eSwitch mode from Legacy to SwitchDev on the PF device. This will also create the VF representor netdevices in the host OS.

```
# devlink dev eswitch set pci/0000:3b:00.0 mode switchdev
```

Before changing the mode, make sure that all VFs are unbound.

⚠ To go back to SR-IOV legacy mode, run:

devlink dev eswitch set pci/0000:3b:00.0 mode legacy This will also remove the VF representor netdevices.

On old OSs or kernels that do not support Devlink, moving to SwitchDev mode can be done using sysfs.

```
echo switchdev > /sys/class/net/enp4s0f0/compat/devlink/mode
```

3. At this stage, VF representors have been created. To map representor to its VF, make sure to obtain the representor's switchid and portname from:

```
ip -d link show eth4
1: enp0s8f0_1: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc mq state UP mode DEFAULT group default
qlen 1000
qlen 1000 link/ether ba:e6:21:37:bc:d4 brd ff:ff:ff:ff:ff:ff promiscuity 0 addrgenmode eui64 numtxqueues 10 numrxqueues 10 gso_max_size 65536 gso_max_segs 65535 portname pf0vf1 switchid f4ab580003a1420c
```

switchid - used to map representor to device, both device PFs have the same switchid. portname - used to map representor to PF and VF, value returned is pfXvfY, where X is the PF number and Y is the number of VF.

On old kernels, switchid and portname can be acquired through sysfs:

```
# cat /sys/class/net/eth4/phys_switch_id
f4ab580003a1420c
# cat /sys/class/net/eth4/phys_port_name pf0vf1
```

During OFED installation, udev scripts are created. This will set the network VF representor device names to be in the form of \$PF_\$VFID, where \$PF is the PF netdev name, and \$VFID is the VF ID=0,1,[..],

SwitchDev Performance Tuning

SwitchDev performance can be further improved by tuning it.

Steering Mode

OVS-kernel supports two steering modes for rules insertion into hardware.

- 1. SMFS Software Managed Flow Steering (as of MLNX_OFED v5.1, this is the default mode) Rules are inserted directly to the hardrware by the software (driver). This mode is optimized for rules insertion.
- 2. DMFS Device Managed Flow Steering Rules insertion is done using firmware commands. This mode is optimized for throughput with

a small amount of rules in the system.

The mode can be controlled via sysfs or devlink API in kernels that support it:

```
Sysfs:
# echo smfs > /sys/class/net/<PF netdev>/compat/devlink/steering_mode

Devlink:
# devlink dev param set pci/0000:00:08.0 name flow_steering_mode value "smfs" cmode runtime

Replace smfs param with dmfs for device managed flow steering
```

Notes:

- The mode should be set prior to moving to SwitchDev, by echoing to the sysfs or invoking the devlink command.
- Only when moving to SwitchDev will the driver use the mode set by the previous step.
- Mode cannot be changed after moving to SwitchDev.
- The steering mode is applicable for SwitchDev mode only, meaning it does not affect legacy SR-IOV or other configurations.

vPort Match Mode

OVS-kernel support two modes that define how the rules on match on vport.

- Metadata rules match on metadata instead of vport number (default mode).
 This mode is needed in order to support SR-IOV Live migration and Dual port RoCE features.
 Matching on Metadata can have a performance impact.
- Legacy rules match on vport number.
 In this mode, performance can be higher in comparison to Metadata. It can still be used only if none of the above features (SR-IOV Live migration and Dual port RoCE) is enabled/used.
 The mode can be controlled via sysfs:

```
Set Legacy:
# echo legacy > /sys/class/net/<PF netdev>/compat/devlink/vport_match_mode

Set metadata:
Devlink:
# echo metadata > /sys/class/net/<PF netdev>/compat/devlink/vport_match_mode
```

Note: This mode should be set prior to moving to SwitchDev, by echoing to the sysfs.

Flow Table Large Group Number

Offloaded flows, including Connection Tracking, are added to Virtual Switch Forwarding Data Base (FDB) flow tables. FDB tables have a set of flow groups, where each flow group saves the same traffic pattern flows. E.g, for connection tracking offloaded flow, TCP and UDP are different traffic patterns which will end up in two different flow groups.

A flow group has a limited size to save flow entries. As default, the driver has 15 big FDB flow groups. Each of these big flow groups can save 4M / (15 + 1) = 256k different 5-tuple flow entries at most. For scenarios with more than 15 traffic patterns, the driver provides a module parameter (num_of_groups) to allow customization and performance tuning.

The mode can be controlled via module param or devlink API for kernels that support it:

Notes:

- In MLNX_OFED v5.1, the default value was changed from 4 to 15.
- The change takes effect immediately if there is no flow inside the FDB table (no traffic running and all offloaded flows are aged out). And it can be dynamically changed without reloading the driver.

If there are still offloaded flows residual when changing this parameter, it will only take effect after all flows have aged out.

Open vSwitch Configuration

Open vSwitch configuration is a simple OVS bridge configuration with SwitchDev.

1. Run the openvswitch service.

```
# systemctl start openvswitch
```

2. Create an OVS bridge (here it's named ovs-sriov).

```
# ovs-vsctl add-br ovs-sriov
```

3. Enable hardware offload (disabled by default).

```
# ovs-vsctl set Open_vSwitch . other_config:hw-offload=true
```

4. Restart the openvswitch service. This step is required for HW offload changes to take effect.

```
# systemctl restart openvswitch
```



HW offload policy can also be changed by setting the tc-policy using one on the following values:

- * none adds a TC rule to both the software and the hardware (default)
- * skip_sw adds a TC rule only to the hardware
- * skip_hw adds a TC rule only to the software

The above change is used for debug purposes.

5. Add the PF and the VF representor netdevices as OVS ports.

```
ovs-vsctl add-port ovs-sriov enp4s0f0
# ovs-vsctl add-port ovs-sriov enp4s0f0_0
# ovs-vsctl add-port ovs-sriov enp4s0f0_1
```

Make sure to bring up the PF and representor netdevices.

```
# ip link set dev enp4s0f0 up
# ip link set dev enp4s0f0_0 up
# ip link set dev enp4s0f0_1 up
```

The PF represents the uplink (wire).

```
# ovs-dpctl show
system@ovs-system:
    lookups: hit:0 missed:192 lost:1
    flows: 2
    masks: hit:384 total:2 hit/pkt:2.00
    port 0: ovs-system (internal)
    port 1: ovs-sriov (internal)
    port 2: enp4s0f0
    port 3: enp4s0f0
    port 3: enp4s0f0_0
    port 4: enp4s0f0_1
```

6. Run traffic from the VFs and observe the rules added to the OVS data-path.

```
# ovs-dpctl dump-flows
recirc_id(0),in_port(3),eth(src=e4:11:22:33:44:50,dst=e4:1d:2d:a5:f3:9d),
eth_type(0x0800),ipv4(frag=no), packets:33, bytes:3234, used:1.196s, actions:2
recirc_id(0),in_port(2),eth(src=e4:1d:2d:a5:f3:9d,dst=e4:11:22:33:44:50),
eth_type(0x0800),ipv4(frag=no), packets:34, bytes:3332, used:1.196s, actions:3
```

In the example above, the ping was initiated from VF0 (OVS port 3) to the outer node (OVS port 2), where the VF MAC is e4:11:22:33:44:50 and the outer node MAC is e4:1d:2d:a5:f3:9d As shown above, two OVS rules were added, one in each direction.

Note that you can also verify offloaded packets by adding type=offloaded to the command. For example:

```
# ovs-appctl dpctl/dump-flows type=offloaded
```

Open vSwitch Performance Tuning

Flow Aging

The aging timeout of OVS is given in ms and can be controlled using the following command.

```
# ovs-vsctl set Open_vSwitch . other_config:max-idle=30000
```

TC Policy

Specifies the policy used with HW offloading.

- none adds a TC rule to both the software and the hardware (default)
- skip_sw adds a TC rule only to the hardware
- skip_hw adds a TC rule only to the software

Example:

```
# ovs-vsctl set Open_vSwitch . other_config:tc-policy=skip_sw
```

Note: TC policy should only be used for debugging purposes.

Max-Revalidator

Specifies the maximum time (in ms) that revalidator threads will wait for kernel statistics before executing flow revalidation.

```
# ovs-vsctl set Open_vSwitch . other_config:max-revalidator=10000
```

n-handler-threads

Specifies the number of threads for software datapaths to use for handling new flows. The default value is the number of online CPU cores minus the number of revalidators.

```
# ovs-vsctl set Open_vSwitch . other_config:n-handler-threads=4
```

n-revalidator-threads

Specifies the number of threads for software datapaths to use for revalidating flows in the datapath.

```
# ovs-vsctl set Open_vSwitch . other_config:n-revalidator-threads=4
```

vlan-limit

Limits the number of VLAN headers that can be matched to the specified number.

```
# ovs-vsctl set Open_vSwitch . other_config:vlan-limit=2
```

Basic TC Rules Configuration

Offloading rules can also be added directly, and not only through OVS, using the tc utility.



 Create an ingress qdisc (queueing discipline) for each interface that you wish to add rules into.

```
# tc qdisc add dev enp4s0f0 ingress
# tc qdisc add dev enp4s0f0_0 ingress
# tc qdisc add dev enp4s0f0_1 ingress
```

2. Add TC rules using flower classifier in the following format.

```
# tc filter add dev NETDEVICE ingress protocol PROTOCOL prio PRIORITY \
[chain CHAIN] flower [ MATCH_LIST ] [ action ACTION_SPEC ]
```

Note: List of supported matches (specifications) and actions can be found in <u>Classification</u> <u>Fields (Matches)</u> section.

3. Dump the existing tc rules using flower classifier in the following format.

```
# tc [ -s ] filter show dev NETDEVICE ingress
```

SR-IOV VF LAG

SR-IOV VF LAG allows the NIC's physical functions (PFs) to get the rules that the OVS will try to offload to the bond net-device, and to offload them to the hardware e-switch. Bond modes supported are:

- · Active-Backup
- XOR
- LACP

SR-IOV VF LAG enables complete offload of the LAG functionality to the hardware. The bonding creates a single bonded PF port. Packets from up-link can arrive from any of the physical ports, and will be forwarded to the bond device.

When hardware offload is used, packets from both ports can be forwarded to any of the VFs. Traffic from the VF can be forwarded to both ports according to the bonding state. Meaning, when in active-backup mode, only one PF is up, and traffic from any VF will go through this PF. When in XOR or LACP mode, if both PFs are up, traffic from any VF will split between these two PFs.

SR-IOV VF LAG Configuration on ASAP²

To enable SR-IOV VF LAG, both physical functions of the NIC should first be configured to SR-IOV SwitchDev mode, and only afterwards bond the up-link representors.

The example below shows the creation of bond interface on two PFs:

1. Load bonding device and enslave the up-link representor (currently PF) net-device devices.

```
modprobe bonding mode=802.3ad

Ifup bond0 (make sure ifcfg file is present with desired bond configuration)
ip link set enp4s0f0 master bond0
ip link set enp4s0f1 master bond0
```

Add the VF representor net-devices as OVS ports. If tunneling is not used, add the bond device as well.

```
ovs-vsctl add-port ovs-sriov bond0
ovs-vsctl add-port ovs-sriov enp4s0f0_0
ovs-vsctl add-port ovs-sriov enp4s0f1_0
```

3. Make sure to bring up the PF and the representor netdevices.

```
ip link set dev bond0 up
ip link set dev enp4s0f0_0 up
ip link set dev enp4s0f1_0 up
```

wil

Once SR-IOV VF LAG is configured, all VFs of the two PFs will become part of the bond, and will behave as described above.

Limitations

- In VF LAG mode, outgoing traffic in load balanced mode is according to the origin ring, thus, half of the rings will be coupled with port 1 and half with port 2. All the traffic on the same ring will be sent from the same port.
- VF LAG configuration is not supported when the NUM_OF_VFS configured in mlxconfig is higher than 64.

Using TC with VF LAG

Both rules can be added using either of the following.

1. Shared block (supported from kernel 4.16 and RHEL/CentOS 7.7 and above).

```
# tc qdisc add dev bond0 ingress_block 22 ingress
# tc qdisc add dev ens4p0 ingress_block 22 ingress
# tc qdisc add dev ens4p1 ingress_block 22 ingress
```

a. Add drop rule.

```
# tc filter add block 22 protocol arp parent ffff: prio 3 \
    flower \
    dst_mac e4:11:22:11:4a:51 \
    action drop
```

b. Add redirect rule from bond to representor.

```
# tc filter add block 22 protocol arp parent ffff: prio 3 \
    flower \
    dst_mac e4:11:22:11:4a:50 \
    action mirred egress redirect dev ens4f0_0
```

c. Add redirect rule from representor to bond.

```
# tc filter add dev ens4f0_0 protocol arp parent ffff: prio 3 \
    flower \
    dst_mac ec:0d:9a:8a:28:42 \
    action mirred egress redirect dev bond0
```

- 2. Without shared block (supported from kernel 4.15 and below).
 - a. Add redirect rule from bond to representor.

```
# tc filter add dev bond0 protocol arp parent ffff: prio 1 \
    flower \
    dst_mac e4:11:22:11:4a:50 \
    action mirred egress redirect dev ens4f0_0
```

b. Add redirect rule from representor to bond.

```
# tc filter add dev ens4f0_0 protocol arp parent ffff: prio 3 \
    flower \
    dst_mac ec:0d:9a:8a:28:42 \
    action mirred egress redirect dev bond0
```

Classification Fields (Matches)

OVS-Kernel supports multiple classification fields which packets can fully or partially match.

Ethernet Layer 2

- · Destination MAC
- Source MAC
- Ethertype

Supported on all kernels.

In OVS dump flows:

```
skb_priority(0/0),skb_mark(0/0),in_port(eth6),eth(src=00:02:10:40:10:0d
,dst=68:54:ed:00:af:de),eth_type(0x8100), packets:1981, bytes:206024, used:0.440s, dp:tc, actions:eth7
```

Using TC rules:

```
tc filter add dev $rep parent ffff: protocol arp pref 1 \
flower \
dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
action mirred egress redirect dev $NIC
```

IPv4/IPv6

- · Source address
- · Destination address
- Protocol

- TCP/UDP/ICMP/ICMPv6
- TOS
- TTL (HLIMIT)

Supported on all kernels.

In OVS dump flows:

```
Ipv4:
    ipv4(src=0.0.0.0/0.0.0.0,dst=0.0.0.0/0.0.0.0,proto=17,tos=0/0,ttl=0/0,frag=no)
    Ipv6:
    ipv6(src=::/::,dst=1:1:1::3:1040:1008,label=0/0,proto=58,tclass=0/0x3,hlimit=64),
```

Using TC rules:

```
IPv4:
tc filter add dev $rep parent ffff: protocol ip pref 1 \
flower \
dst_ip 1.1.1.1 \
src_ip 1.1.1.2 \
ip_proto TCP \
ip_tso 0x3 \
ip_ttl 63 \
action mirred egress redirect dev $NIC

IPv6:
tc filter add dev $rep parent ffff: protocol ipv6 pref 1 \
flower \
dst_ip 1:1:1::3:1040:1009 \
src_ip 1:1:1::3:1040:1008 \
ip_proto TCP \
ip_tso 0x3 \
ip_ttl 63 \
action mirred egress redirect dev $NIC
```

TCP/UDP Source and Destination ports & TCP Flags

- TCP/UDP source and destinations ports
- TCP flags

Supported kernels are kernel > 4.13 and RHEL > 7.5

In OVS dump flows:

```
TCP: tcp(src=0/0,dst=32768/0x8000),

UDP: udp(src=0/0,dst=32768/0x8000),

TCP flags: tcp_flags(0/0)
```

Using TC rules:

```
tc filter add dev $rep parent ffff: protocol ip pref 1 \
flower \
ip_proto TCP \
dst_port 100 \
src_port 500 \
tcp_flags 0x4/0x7 \
action mirred egress redirect dev $NIC
```

VLAN

- ID
- Priority
- · Inner vlan ID and Priority

Supported kernels: All (QinQ: kernel 4.19 and higher, and RHEL 7.7 and higher)

In OVS dump flows:

```
eth_type(0x8100),vlan(vid=2347,pcp=0),
```

Using TC rules:

Tunnel

- ID (Key)
- Source IP address
- · Destination IP address
- Destination port
- TOS (supported from kernel 4.19 and above & RHEL 7.7 and above)
- TTL (support from kernel 4.19 and above & RHEL 7.7 and above)
- Tunnel options (Geneve)

Supported kernels:

- VXLAN: All
- GRE: Kernel > 5.0, RHEL 7.7 and above
- Geneve: Kernel > 5.0, RHEL 7.7 and above

In OVS dump flows:

```
tunnel(tun_id=0x5,src=121.9.1.1,dst=131.10.1.1,ttl=0/0,tp_dst=4789,flags(+key))
```

Using TC rules:

```
# tc filter add dev $rep protocol 802.10 parent ffff: pref 1
flower \
vlan_ethtype 0x800 \
vlan_prio 0 \
action mirred egress redirect dev $NIC
QinQ:
# tc filter add dev vxlan100 protocol ip parent ffff: \
flower \
skip_sw \
dst_mac e4:11:22:11:4a:51 \
src_mac e4*:11:22:11:4a:50 \
enc_src_ip 20.1.11.1 \
enc_dst_ip 20.1.12.1 \
enc_key_id 100 \
enc_dst_port 4789 \
action mirred egress redirect dev ens4f0_0
```

Supported Actions

Forward

Forward action allows for packet redirection:

- From VF to wire
- · Wire to VF
- VF to VF

Supported on all kernels.

In OVS dump flows:

```
skb_priority(0/0),skb_mark(0/0),in_port(eth6),eth(src=00:02:10:40:10:0d,dst=68:54:ed:00:af:de),eth_type(0x8100), packets:1981, bytes:206024, used:0.440s, dp:tc, actions:eth7
```

Using TC rules:

```
tc filter add dev $rep parent ffff: protocol arp pref 1 \
flower \
dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
action mirred egress redirect dev $NIC
```

Drop

Drop action allows to drop incoming packets.

Supported on all kernels.

In OVS dump flows:

```
skb_priority(0/0),skb_mark(0/0),in_port(eth6),eth(src=00:02:10:40:10:0d,dst=68:54:ed:00:af:de),eth_type(0x8100), packets:1981, bytes:206024, used:0.440s, dp:tc, actions:drop
```

Using TC rules:

```
tc filter add dev $rep parent ffff: protocol arp pref 1 \
flower \
dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
action drop
```

Statistics

By default, each flow collects the following statistics:

- Packets number of packets which hit the flow
- · Bytes total number of bytes which hit the flow
- Last used the amount of time passed since last packet hit the flow

Supported on all kernels.

In OVS dump flows:

```
skb_priority(0/0),skb_mark(0/0),in_port(eth6),eth(src=00:02:10:40:10:0d,dst=68:54:ed:00:af:de),eth_type(0x8100), packets:1981, bytes:206024, used:0.440s, dp:tc, actions:drop
```

Using TC rules:

```
#tc -s filter show dev $rep ingress

filter protocol ip pref 2 flower chain 0
filter protocol ip pref 2 flower chain 0 handle 0x2
eth_type ipv4
ip_proto tcp
src_jp 192.168.140.100
src_port 80
skip_sw
in_hw
    action order 1: mirred (Egress Redirect to device p0v11_r) stolen
    index 34 ref 1 bind 1 installed 144 sec used 0 sec
    Action statistics:
    Sent 388344 bytes 2942 pkt (dropped 0, overlimits 0 requeues 0)
    backlog 0b 0p requeues 0
```

Tunnels (Encapsulation/Decapsulation)

OVS-kernel supports offload of tunnels using encapsulation and decapsulation actions.

- Encapsulation pushing of tunnel header is supported on Tx
- Decapsulation popping of tunnel header is supported on Rx

Supported Tunnels:

- VXLAN (IPv4/IPv6) supported on all Kernels
- GRE (IPv4/IPv6) supported on kernel 5.0 and above & RHEL 7.6 and above
- Geneve (IPv4/IPv6) supported on kernel 5.0 and above & RHEL 7.6 and above

OVS configuration:

In case of offloading tunnel, the PF/bond should not be added as a port in the OVS datapath. It should rather be assigned with the IP address to be used for encapsulation.

The example below shows two hosts (PFs) with IPs 1.1.1.177 and 1.1.1.75, where the PF device on both hosts is enp4s0f0, and the VXLAN tunnel is set with VNID 98:

· On the first host:

```
# ip addr add 1.1.1.177/24 dev enp4s0f1
# ovs-vsctl add-port ovs-sriov vxlan0 -- set interface vxlan0 type=vxlan
options:local_ip=1.1.1.177 options:remote_ip=1.1.1.75 options:key=98
```

· On the second host:

```
# ip addr add 1.1.1.75/24 dev enp4s0f1

# ovs-vsctl add-port ovs-sriov vxlan0 -- set interface vxlan0 type=vxlan
options:local_ip=1.1.1.75 options:remote_ip=1.1.1.177 options:key=98

• for GRE IPv4 tunnel need use type=gre
• for GRE IPv6 tunnel need use type=ip6gre
• for GENEVE tunnel need use type=geneve
```

4

When encapsulating guest traffic, the VF's device MTU must be reduced to allow the host/ HW to add the encap headers without fragmenting the resulted packet. As such, the VF's MTU must be lowered by 50 bytes from the uplink MTU for IPv4 and 70 bytes for IPv6.

Tunnel offload using TC rules:

VLAN Push/Pop

OVS-kernel supports offload of vlan header push/pop actions.

- Push pushing of vlan header is supported on Tx
- Pop popping of tunnel header is supported on Rx

OVS Configuration

Add a tag=\$TAG section for the OVS command line that adds the representor ports. For example, VLAN ID 52 is being used here.

```
# ovs-vsctl add-port ovs-sriov enp4s0f0
# ovs-vsctl add-port ovs-sriov enp4s0f0_0 tag=52
# ovs-vsctl add-port ovs-sriov enp4s0f0_1 tag=52
```

The PF port should not have a VLAN attached. This will cause OVS to add VLAN push/pop actions when managing traffic for these VFs.

Dump Flow Example

```
recirc_id(0),in_port(3),eth(src=e4:11:22:33:44:50,dst=00:02:c9:e9:bb:b2),eth_type(0x0800),ipv4(frag=no), \
packets:0, bytes:0, used:never, actions:push_vlan(vid=52,pcp=0),2

recirc_id(0),in_port(2),eth(src=00:02:c9:e9:bb:b2,dst=e4:11:22:33:44:50),eth_type(0x8100), \
vlan(vid=52,pcp=0),encap(eth_type(0x0800),ipv4(frag=no)), packets:0, bytes:0, used:never, actions:pop_vlan,3
```

VLAN Offload using TC Rules Example

Header Rewrite

This action allows for modifying packet fields.

Ethernet Layer 2

- Destination MAC
- Source MAC

Supported kernels: Kernel 4.14 and above & RHEL 7.5 and above

In OVS dump flows:

```
skb_priority(0/0),skb_mark(0/0),in_port(eth6),eth(src=00:02:10:40:10:0d,dst=68:54:ed:00:af:de),eth_type(0x8100), packets:1981, bytes:206024, used:0.440s, dp:tc, actions: set(eth(src=68:54:ed:00:f4:ab,dst=fa:16:3e:dd:69:c4)),eth7
```

Using TC rules:

```
tc filter add dev $rep parent_fffff: protocol arp pref 1 \
                                                           rep parent ffff: protocol arp pret 1 \
flower \
dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
action pedit ex \
munge eth dst set 20:22:33:44:55:66 \
munge eth src set aa:ba:cc:dd:ee:fe \
action mirred egress redirect dev $NIC
```

IPv4/IPv6

- Source address
- · Destination address
- Protocol
- TOS
- TTL (HLIMIT)

Supported kernels: Kernel 4.14 and above & RHEL 7.5 and above

In OVS dump flows:

```
set(eth(src=de:e8:ef:27:5e:45,dst=00:00:01:01:01:01)),
set(ipv4(src=10.10.0.111,dst=10.20.0.122,tt1=63))
Ipv6:
      set(ipv6(dst=2001:1:6::92eb:fcbe:f1c8,hlimit=63)),
```

Using TC rules:

```
ip_ttl 63 \
pedit ex \
munge ip src set 2.2.2.1 \
munge ip dst set 2.2.2.2 \
munge ip tos set 0 \
munge ip ttl dec \
action mirred egress redirect dev $NIC
ip_ttl 63\
                               pedit ex \
munge ipv6 src set 2:2:2:3:1040:1009 \
munge ipv6 dst set 2:2:2:3:1040:1008 \
munge ipv6 hlimit dec \
action mirred egress redirect dev $NIC
```

▲ IPv4 and IPv6 header rewrite is only supported with match on UDP/TCP/ICMP protocols.

TCP/UDP Source and Destination Ports

• TCP/UDP source and destinations ports

Supported kernels: kernel > 4.16 & RHEL > 7.6

In OVS dump flows:

```
TCP:
        set(tcp(src= 32768/0xfffff,dst=32768/0xfffff)),
UDP:
        set(udp(src= 32768/0xfffff,dst=32768/0xfffff)),
```

Using TC rules:

```
tc filter add dev $rep parent ffff: protocol ip pref 1 \
flower \
dst_ip 1.1.1.1 \
src_ip 1.1.1.2 \
ip_proto tcp \
ip_tto 80x3 \
ip_ttl 63 \
pedit ex munge ip tcp sport set 200
pedit ex munge ip tcp dport set 200
action mirred egress redirect dev $NIC

UDP:

tc filter add dev $rep parent ffff: protocol ip pref 1 \
flower \
dst_ip 1.1.1.1 \
src_ip 1.1.1.2 \
ip_proto udp \
ip_tos 0x3 \
ip_ttl 63 \
pedit ex \
pedit ex munge ip udp sport set 200
pedit ex munge ip udp sport set 200
pedit ex munge ip udp dport set 200
pedit ex munge ip udp dport set 200
action mirred egress redirect dev $NIC
```

VLAN

ID

Supported on all kernels.

In OVS dump flows:

```
Set(vlan(vid=2347,pcp=0/0)),
```

Using TC rules:

```
tc filter add dev $rep parent ffff: protocol 802.10 pref 1 \
flower \
vlan_ethtype 0x800 \
vlan_jid 100 \
vlan_prio 0 \
action vlan modify id 11 pipe
action mirred egress redirect dev $NIC
```

Connection Tracking

The TC connection tracking action performs connection tracking lookup by sending the packet to netfilter conntrack module. Newly added connections may be associated, via the ct commit action, with a 32 bit mark, 128 bit label and source/destination NAT values.

The following example allows ingress tcp traffic from the uplink representor to vf1_rep, while assuring that egress traffic from vf1_rep is only allowed on established connections. In addition, mark and source IP NAT is applied.

In OVS dump flows:

```
ct(zone=2,nat)
ct_state(+est+trk)
actions:ct(commit,zone=2,mark=0x4/0xfffffffff,nat(src=5.5.5.5))
```

Using TC rules:

Forward to Chain (TC Only)

TC interface supports adding flows on different chains. Only chain 0 is accessed by default. Access to the other chains requires usafe of the goto action.

In this example, a flow is created on chain 1 without any match and redirect to wire.

The second flow is created on chain 0 and match on source MAC and action goto chain 1.

This example simulates simple MAC spoofing.

Port Mirroring (Flow Based VF Traffic Mirroring for ASAP²)

Unlike para-virtual configurations, when the VM traffic is offloaded to the hardware via SR-IOV VF, the host side Admin cannot snoop the traffic (e.g. for monitoring).

ASAP² uses the existing mirroring support in OVS and TC along with the enhancement to the offloading logic in the driver to allow mirroring the VF traffic to another VF.

The mirrored VF can be used to run traffic analyzer (tcpdump, wireshark, etc) and observe the traffic of the VF being mirrored.

The example below shows the creation of port mirror on the following configuration:

```
# ovs-vsctl show
09d8a574-9c39-465c-9f16-47d81c12f88a
Bridge br-vxlan
Port "enp4s0f0_1"
Interface "enp4s0f0_1"
Port "vxlan0"
Interface "vxlan0"
type: vxl
                           Interface c...
Port br-vxlan
Interface br-vxlan
type: internal
         ovs version: "2.14.1"
```

• To set enp4s0f0_0 as the mirror port, and mirror all of the traffic:

• To set enp4s0f0_0 as the mirror port, and only mirror the traffic, the destination is enp4s0f0_1:

• To set enp4s0f0_0 as the mirror port, and only mirror the traffic the source is enp4s0f0_1:

• To set enp4s0f0_0 as the mirror port and mirror, all the traffic on enp4s0f0_1:

```
# ovs-vsctl -- --id=@p1 get port enp4s0f0_0 \
    -- --id=@p2 get port enp4s0f0_1 \
    -- --id=@m create mirror name=m0 select-dst-port=@p2 select-src-port=@p2 output-port=@p1 \
    -- -- set bridge br-vxlan mirrors=@m
```

To clear the mirror port:

```
# ovs-vsctl clear bridge br-vxlan mirrors
```

Mirroring using TC:

```
tc filter add dev $rep parent ffff: protocol arp pref 1 \
                                        flower \
dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
                                         action mirred egress mirror dev $mirror_rep pipe \action mirred egress redirect dev $NIC
tc filter add dev $rep parent ffff: protocol arp pref 1 \
                                        dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
                           action tunnel_key set \
src_ip 1.1.1.1 \
dst_ip 1.1.1.2 \
dst_port 4789 \
id 768 \
                           pipe \
action mirred egress mirror dev vxlan100 pipe \
action mirred egress redirect dev $NIC
```

sFLOW

This feature allows for monitoring the traffic sent between two VMs on the same host using an sFlow collector.

The example below shows the creation of port mirroring on the following configuration.

To sample all traffic over the OVS bridge:

Parameter	Description	
SFLOW_AGENT	Indicates that the sFlow agent should send traffic from SFLOW_AGENT's IP address	
SFLOW_TARGET	Remote IP address of the sFLOW collector	
SFLOW_PORT	Size of packet header to sample (in bytes)	
SFLOW_SAMPLING	Sample rate	

To clear the sFLOW configuration:

```
# ovs-vsctl clear bridge br-vxlan mirrors
```

sFLOW using TC:

```
Sample to VF
tc filter add dev $rep parent ffff: protocol arp pref 1 \
flower \
dst_mac e4:1d:2d:5d:25:35 \
src_mac e4:1d:2d:5d:25:34 \
action sample rate 10 group 5 trunc 96 \
action mirred egress redirect dev $NIC
```

A

Userspace application is needed in order to process to sampled packet from the kernel. Example: https://github.com/Mellanox/libpsample

Rate Limit

OVS-kernel supports offload of VF rate limit using OVS configuration and TC.

The example below sets a rate limit to the VF related to representor eth0 to 10Mbps.

```
# ovs-vsctl set interface eth0 ingress policing rate=10000
tc:
# tc_filter add dev eth0 root prio 1 protocol ip matchall skip_sw action police rate 10mbit burst 20k
```

Kernel Requirements

This kernel config should be enabled in order to support switchdev offload.

- CONFIG_NET_ACT_CSUM needed for action csum
- CONFIG_NET_ACT_PEDIT needed for header rewrite
- CONFIG_NET_ACT_MIRRED needed for basic forward
- CONFIG_NET_ACT_CT needed for connection tracking (supported from kernel 5.6)
- CONFIG_NET_ACT_VLAN needed for action vlan push/pop
- CONFIG_NET_ACT_GACT
- CONFIG_NET_CLS_FLOWER
- CONFIG_NET_CLS_ACT
- CONFIG_NET_SWITCHDEV
- CONFIG_NET_TC_SKB_EXT needed for connection tracking (supported from kernel 5.6)
- CONFIG_NET_ACT_CT needed for connection tracking (supported from kernel 5.6)
- CONFIG_NFT_FLOW_OFFLOAD
- CONFIG_NET_ACT_TUNNEL_KEY
- CONFIG_NF_FLOW_TABLE needed for connection tracking (supported from kernel 5.6)
- CONFIG_SKB_EXTENSIONS needed for connection tracking (supported from kernel 5.6)
- CONFIG_NET_CLS_MATCHALL
- CONFIG_NET_ACT_POLICE
- CONFIG_MLX5_ESWITCH

OVS-DPDK Hardware Offloads

OVS-DPDK Hardware Offloads Configuration



To configure OVS-DPDK HW offloads:

1. Unbind the VFs.

```
echo 0000:04:00.2 > /sys/bus/pci/drivers/mlx5_core/unbind echo 0000:04:00.3 > /sys/bus/pci/drivers/mlx5_core/unbind
```

Note: VMs with attached VFs must be powered off to be able to unbind the VFs.

2. Change the e-switch mode from Legacy to SwitchDev on the PF device (make sure all VFs are unbound). This will also create the VF representor netdevices in the host OS.

```
echo switchdev > /sys/class/net/enp4s0f0/compat/devlink/mode
```

To revert to SR-IOV Legacy mode:

```
echo legacy > /sys/class/net/enp4s0f0/compat/devlink/mode
```

Note that running this command will also result in the removal of the VF representor netdevices.

3. Bind the VFs.

```
echo 0000:04:00.2 > /sys/bus/pci/drivers/mlx5_core/bind
echo 0000:04:00.3 > /sys/bus/pci/drivers/mlx5_core/bind
```

4. Run the Open vSwitch service.

```
systemctl start openvswitch
```

5. Enable hardware offload (disabled by default).

```
ovs-vsctl --no-wait set Open_vSwitch . other_config:dpdk-init=true ovs-vsctl set Open_vSwitch . other_config:hw-offload=true
```

6. Configure the DPDK white list.

```
ovs-vsctl --no-wait set Open_vSwitch . other_config:dpdk-extra="-w 0000:01:00.0,representor=[0],dv_flow_en=1,dv_esw_en=1,dv_xmeta_en=1"
```

7. Restart the Open vSwitch service. This step is required for HW offload changes to take effect.

```
systemctl restart openvswitch
```

8. Create OVS-DPDK bridge.

```
ovs-vsctl --no-wait add-br br0-ovs -- set bridge br0-ovs datapath_type=netdev
```

9. Add PF to OVS.

```
ovs-vsctl add-port br0-ovs pf -- set Interface pf type=dpdk options:dpdk-devargs=0000:88:00.0
```

10. Add representor to OVS.

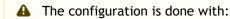
```
ovs-vsctl add-port br0-ovs representor -- set Interface representor type=dpdk options:dpdk-devargs=0000:88:0 0.0,representor=[$rep]
```

Offloading VXLAN Encapsulation/Decapsulation Actions

vSwitch in userspace rather than kernel-based Open vSwitch requires an additional bridge. The purpose of this bridge is to allow use of the kernel network stack for routing and ARP resolution.

The datapath needs to look-up the routing table and ARP table to prepare the tunnel header and transmit data to the output port.

Configuring VXLAN Encap/Decap Offloads



- PF on 0000:03:00.0 PCI and MAC 98:03:9b:cc:21:e8
- Local IP 56.56.67.1 br-phy interface will be configured to this IP
- Remote IP 56.56.68.1

To configure OVS-DPDK VXLAN:

1. Create a br-phy bridge.

ovs-vsctl add-br br-phy -- set Bridge br-phy datapath_type=netdev -- br-set-external-id br-phy bridge-id br-phy -- set bridge br-phy fail-mode=standalone other_config:hwaddr=98:03:9b:cc:21:e8

2. Attach PF interface to br-phy bridge.

ovs-vsctl add-port br-phy p0 -- set Interface p0 type=dpdk options:dpdk-devargs=0000:03:00.0

3. Configure IP to the bridge.

ip addr add 56.56.67.1/24 dev br-phy

4. Create a br-ovs bridge.

ovs-vsctl add-br br-ovs -- set Bridge br-ovs datapath_type=netdev -- br-set-external-id br-ovs bridge-id br-ovs -- set bridge br-ovs fail-mode=standalone

5. Attach representor to br-ovs.

 $ovs-vsct1 \ add-port \ br-ovs \ pf0vf0 \ -- \ set \ Interface \ pf0vf0 \ type=dpdk \ options: dpdk-devargs=0000:03:00.0, representor=[0]$

6. Add a port for the VXLAN tunnel.

ovs-vsctl add-port ovs-sriov vxlan0 -- set **interface** vxlan0 type=vxlan options:local_ip=56.56.67.1 options:remote_ip=56.56.68.1 options:key=45 options:dst_port=4789

Connection Tracking Offload

Connection tracking enables stateful packet processing by keeping a record of currently open connections.

OVS flows using connection tracking can be accelerated using advanced Network Interface Cards (NICs) by offloading established connections.

To view offloaded connections, run:

ovs-appctl dpctl/get-offload-stats

SR-IOV VF LAG

To configure OVS-DPDK SR-IOV VF LAG:

1. Enable SR-IOV on the NICs.

mlxconfig -d <PCI> set SRIOV_EN=1

2. Allocate the desired number of VFs per port.

echo \$n > /sys/class/net/<net name>/device/sriov_numvfs

3. Unbind all VFs.

echo <VF PCI> >/sys/bus/pci/drivers/mlx5_core/unbind

4. Change both NICs' mode to SwitchDev.

devlink dev eswitch set pci/<PCI> mode switchdev

5. Create Linux bonding using kernel modules.

modprobe bonding mode=<desired mode>

Note: Other bonding parameters can be added here. The supported Bond modes are: Active-Backup, XOR and LACP

6. Bring all PFs and VFs down.

ip link set <PF/VF> down

7. Attach both PFs to the bond.

ip link set <PF> master bond0

8. To work with VF-LAG with OVS-DPDK, add the bond master (PF) to the bridge.

ovs-vsctl add-port br-phy p0 -- set Interface p0 type=dpdk options:dpdk-devargs=0000:03:00.0 options:dpdk-lsc-interrupt=true

9. Add representor \$N of PF0 or PF1 to a bridge.

ovs-vsctl add-port br-phy rep\$N -- set Interface rep\$N type=dpdk options:dpdk-devargs=<PF0 PCI>,representor=pf0vf\$N

OR

ovs-vsctl add-port br-phy rep\$N -- set Interface rep\$N type=dpdk options:dpdk-devargs=<PF0 PCI>,representor=pf1vf\$N

VirtIO Acceleration through VF Relay (Software & Hardware vDPA)

▲ Hardware vDPA is supported on ConnectX-6 Dx & BlueField-2 cards and above only.

⚠ Hardware vDPA is enabled by default. In case your hardware does not support vDPA, the driver will fall back to Software vDPA.

To check which vDPA mode is activated on your driver, run: ovs-ofctl -0 OpenFlow14 dump-ports br0-ovs and look for hw-mode flag.

⚠ This feature has not been accepted to the OVS-DPDK Upstream yet, making its API subject to change.

In user space, there are two main approaches for communicating with a guest (VM), either through SR-IOV, or through virtIO.

Phy ports (SR-IOV) allow working with port representor, which is attached to the OVS and a matching VF is given with pass-through to the guest. HW rules can process packets from up-link and direct

them to the VF without going through SW (OVS). Therefore, using SR-IOV achieves the best performance.

However, SR-IOV architecture requires the guest to use a driver specific to the underlying HW. Specific HW driver has two main drawbacks:

- 1. Breaks virtualization in some sense (guest is aware of the HW). It can also limit the type of images supported.
- 2. Gives less natural support for live migration.

Using virtIO port solves both problems. However, it reduces performance and causes loss of some functionalities, such as, for some HW offloads, working directly with virtIO. To solve this conflict, a new netdev type- dpdkvdpa has been created. The new netdev is similar to the regular DPDK netdey, yet introduces several additional functionalities.

dpdkvdpa translates between phy port to virtIO port. It takes packets from the Rx queue and sends them to the suitable Tx queue, and allows transfer of packets from virtIO guest (VM) to a VF, and vice-versa, benefitting from both SR-IOV and virtIO.



To add vDPA port:

```
ovs-vsctl add-port br0 vdpa0 -- set Interface vdpa0 type=dpdkvdpa \
ovs-vsct1 add-port brU vdpaU -- set Interface vdpaU options:vdpa-socket-path-<br/>
options:vdpa-accelerator-devargs=<vf pci id> \<br/>
options:dpdk-devargs=<pf pci id>,representor=[id] \<br/>
options: vdpa-max-queues = <num queues> \<br/>
options: vdpa-sw=<true/false>
```

Note: vdpa-max-queues is an optional field. When the user wants to configure 32 vDPA ports, the maximum queues number is limited to 8.

vDPA Configuration in OVS-DPDK Mode

Prior to configuring vDPA in OVS-DPDK mode, follow the steps below.

1. Generate the VF.

```
echo 0 > /sys/class/net/enp175s0f0/device/sriov_numvfs
echo 4 > /sys/class/net/enp175s0f0/device/sriov_numvfs
```

Unbind each VF.

```
echo <pci> > /sys/bus/pci/drivers/mlx5_core/unbind
```

Switch to SwitchDev mode.

```
echo switchdev >> /svs/class/net/enp175s0f0/compat/devlink/mode
```

4. Bind each VF.

```
echo <pci> > /sys/bus/pci/drivers/mlx5 core/bind
```

5. Initialize OVS with:

```
ovs-vsctl --no-wait set Open_vSwitch . other_config:dpdk-init=true ovs-vsctl --no-wait set Open_vSwitch . other_config:hw-offload=true
```

To configure vDPA in OVS-DPDK mode on ConnectX-5 cards and above:

1. Open vSwitch configuration.

```
ovs-vsctl --no-wait set Open_vSwitch . other_config:dpdk-extra="-w
0000:01:00.0,representor=[0],dv_flow_en=1,dv_esw_en=1,dv_xmeta_en=
/usr/share/openvswitch/scripts/ovs-ctl restart
```

2. Create OVS-DPDK bridge.

```
ovs-vsctl add-br br0-ovs -- set bridge br0-ovs datapath_type=netdev ovs-vsctl add-port br0-ovs pf -- set Interface pf type=dpdk options:dpdk-devargs=0000:01:00.0
```

3. Create vDPA port as part of the OVS-DPDK bridge.

ovs-vsctl add-port br0-ovs vdpa0 -- set Interface vdpa0 type=dpdkvdpa options:vdpa-socket-path=/var/run/virtio-forwarder/sock0 options:vdpa-accelerator-devargs=0000:01:00.2 options:dpdk-devargs=0000:01:00.0,repr esentor=[0] options: vdpa-max-queues=8

To configure vDPA in OVS-DPDK mode on BlueField cards:

Set the bridge with the software or hardware vDPA port:

• On the ARM side: Create the OVS-DPDK bridge.

```
ovs-vsctl add-br br0-ovs -- set bridge br0-ovs datapath_type=netdev ovs-vsctl add-port br0-ovs pf -- set Interface pf type=dpdk options:dpdk-devargs=0000:af:00.0 ovs-vsctl add-port br0-ovs rep-- set Interface rep type=dpdk options:dpdk-devargs=0000:af:00.0,representor=[0]
```

• On the host side: Create the OVS-DPDK bridge.

```
ovs-vsctl add-br br1-ovs -- set bridge br1-ovs datapath_type=netdev protocols=OpenFlow14 ovs-vsctl add-port br0-ovs vdpa0 -- set Interface vdpa0 type=dpdkvdpa options:vdpa-socket-path=/var/run/virtio-forwarder/sock0 options:vdpa-accelerator-devargs=0000:af:00.2
```

Note: To configure SW vDPA, add "options:vdpa-sw=true" to the end of the command.

Software vDPA Configuration in OVS-Kernel Mode

SW vDPA can also be used in configurations where the HW offload is done through TC and not DPDK.

1. Open vSwitch configuration.

```
ovs-vsctl set Open_vSwitch . other_config:dpdk-extra="-w
0000:01:00.0,representor=[0],dv_flow_en=1,dv_esw_en=0,idv_xmeta_en=0,isolated_mode=1"
/usr/share/openvswitch/scripts/ovs-ctl restart
```

2. Create OVS-DPDK bridge.

```
ovs-vsctl add-br br0-ovs -- set bridge br0-ovs datapath_type=netdev
```

3. Create vDPA port as part of the OVS-DPDK bridge.

 $ovs-vsct1\ add-port\ br0-ovs\ vdpa0\ --\ set\ Interface\ vdpa0\ type=dpdkvdpa\ options: vdpa-socket-path=/var/run/virtio-forwarder/sock0\ options: vdpa-accelerator-devargs=0000:01:00.2\ options: dpdk-devargs=0000:01:00.0, representation of the control of the control options of the control option$ esentor=[0] options: vdpa-max-queues=8

4. Create Kernel bridge.

```
ovs-vsctl add-br br-kernel
```

5. Add representors to Kernel bridge.

```
ovs-vsctl add-port br-kernel enp1s0f0_0 ovs-vsctl add-port br-kernel enp1s0f0
```

E2E Cache



This feature is at beta level.

OVS offload rules are based on a multi-table architecture. E2E cache feature enables merging the multi-table flow matches and actions into one joint flow.

This improves connection tracking performance by using a single-table when exact match is detected.



To set the E2E cache size (default = 4k):

```
ovs-vsctl set open_vswitch . other_config:e2e-size=<size>systemctl restart openvswitch
```

Note: Make sure to restart the openvswitch service in order for the configuration to take effect.



To enable/disable E2E cache (default = disabled):

```
ovs-vsctl set open_vswitch . other_config:e2e-enable=< true/false > systemctl restart openvswitch
```

Note: Make sure to restart the openvswitch service in order for the configuration to take effect.



To run E2E cache statistics:

ovs-appctl dpctl/dump-e2e-stats



To run E2E cache flows:

ovs-appct1 dpct1/dump-e2e-flows

Geneve Encapsulation/Decapsulation

Geneve tunneling offload feature support includes matching on extension header.



To configure OVS-DPDK Geneve encap/decap:

1. Create a br-phy bridge.

```
ovs-vsctl --may-exist add-br br-phy -- set Bridge br-phy datapath_type=netdev -- br-set-external-id br-phy bridge-id br-phy -- set bridge br-phy fail-mode=standalone
```

2. Attach PF interface to br-phy bridge.

```
ovs-vsctl add-port br-phy pf -- set Interface pf type=dpdk options:dpdk-devargs=<PF PCI>
```

3. Configure IP to the bridge.

```
ifconfig br-phy <$local_ip_1> up
```

4. Create a br-int bridge.

```
ovs-vsctl --may-exist add-br br-int -- set Bridge br-int datapath_type=netdev -- br-set-external-id br-int bridge-id br-int -- set bridge br-int fail-mode=standalone
```

5. Attach representor to br-int.

```
ovs-vsctl add-port br-int rep$x -- set Interface rep$x type=dpdk options:dpdk-devargs=<PF PCI>,representor=[$x]
```

6. Add a port for the GENEVE tunnel.

```
ovs-vsctl add-port br-int geneve0 -- set interface geneve0 type=geneve options:key=<VNI>options:remote_ip=<$remote_ip_1> options:local_ip=<$local_ip_1>
```

Parallel Offloads

OVS-DPDK supports parallel insertion and deletion of offloads (flow & CT). While multiple threads are supported, by default only one is used.



To configure multiple threads:

```
ovs-vsctl set Open_vSwitch . other_config:hw-offload-thread-nb=3
```

Make sure to restart the openvswitch service in order for the configuration to take effect.

```
systemctl restart openvswitch
```

VirtIO Acceleration through Hardware vDPA

Hardware vDPA Installation

Hardware vDPA requires QEMU v2.12.1 and DPDK v20.11 as minimal versions.



To install QEMU:

1. Install prerequisites:

```
yum install -y libseccomp-devel
```

2. Clone the sources:

```
git clone https://git.qemu.org/git/qemu.git -b mlx_2.12.1 cd qemu
```

3. Build QEMU:

```
cd ~/qemu
mkdir bin
cd bin
../configure --target-list=x86_64-softmmu --enable-kvm --enable-seccomp
make
```

To install DPDK:

1. Clone the sources:

```
git clone git://dpdk.org/dpdk
cd dpdk
git checkout v20.11
```

2. Install dependencies (if needed):

```
yum install meson ninja-build
```

3. Configure and build DPDK:

```
export RTE_SDK=$PWD
sed -i 's/\(CONFIG_RTE_LIBRTE_MLX5_PMD=\)n/\ly/g' .config
sed -i 's/\(CONFIG_RTE_LIBRTE_MLX5_VDPA_PMD=\)n/\ly/g' .config
meson build -Dexamples=vdpa
ninja -C build
```

Hardware vDPA Configuration

To configure huge pages:

```
mkdir -p /hugepages
mount -t hugetlbfs hugetlbfs /hugepages
echo <more> > /sys/devices/system/node/node0/hugepages/hugepages-1048576kB/nr_hugepages
echo <more> > /sys/devices/system/node/node1/hugepages/hugepages-1048576kB/nr_hugepages
```

To configure a vDPA VirtIO interface in an existing VM's xml file (using libvirt):

1. Open the VM's configuration xml for editing:

```
virsh edit <domain name>
```

- 2. Modify/add the following:
 - a. Change the top line to:

```
<domain type='kvm' xmlns:qemu='http://libvirt.org/schemas/domain/qemu/1.0'>
```

b. Assign a memory amount and use 1GB page size for hugepages (size must be the same as used for the vDPA application), so that the memory configuration looks like the following.

c. Assign an amount of CPUs for the VM CPU configuration, so that the vcpu and cputune configuration looks like the following.

d. Set the memory access for the CPUs to be shared, so that the cpu configuration looks like the following.

```
<cpu mode='custom' match='exact' check='partial'>
    <model fallback='allow'>Skylake-Server-IBRS</model>
    <numa>
        <cell id='0' cpus='0-4' memory='8388608' unit='KiB' memAccess='shared'/>
        </numa>
        </cpu>
```

e. Set the emulator in use to be the one built in step "2. Build QEMU" above, so that the emulator configuration looks as follows.

```
<emulator><path to qemu executable></emulator>
```

f. Add a virtio interface using qemu command line argument entries, so that the new interface snippet looks as follows.

```
<qemu:commandline>
  <qemu:arg value='-chardev'/>
  <qemu:arg value='socket,id=charnet1,path=/tmp/sock-virtio0'/>
  <qemu:arg value='hetdev'/>
  <qemu:arg value='vhost-user,chardev=charnet1,queues=16,id=hostnet1'/>
  <qemu:arg value='device'/>
  <qemu:arg value='rdevice'/>
  <qemu:arg value='virtio-net-pci,mq=on,vectors=6,netdev=hostnet1,id=net1,mac=e4:11:c6:d3:45:f2,bus =pci.0,addr=0x6,
    page-per-vq=on,rx_queue_size=1024,tx_queue_size=1024'/>
  </qemu:commandline>
```

Note: In this snippet, the vhostuser socket file path, the amount of queues, the MAC and the PCI slot of the VirtIO device can be configured.

g. [Optional] In case you wish to configure jumbo frames, make sure the VM kernel version is 4.14 or above, and replace the virtio-net-pci line as follows.

```
<qemu:arg value='virtio-net-
pci,mq=on,vectors=34,netdev=hostnet1,id=net1,mac=00:21:21:24:02:01,bus=pci.0,addr=0xC,page-per-
vq=on,rx_queue_size=1024,tx_queue_size=1024,host_mtu=9216,csum=on,guest_csum=on,host_tso4=on,host_t
so6=on'/>
```

Running Hardware vDPA

- Hardware vDPA supports SwitchDev mode only.
- Create the ASAP² environment:
 - 1. Create the VFs.
 - 2. Enter switchdev mode.
 - 3. Set up OVS.
- **Run** the vDPA application.

```
cd $RTE_SDK/examples/vdpa/build
./vdpa -w <VF PCI BDF>,class=vdpa --log-level=pmd,info -- -i
```



Create a vDPA port via the vDPA application CLI.

```
create /tmp/sock-virtio0 <PCI DEVICE BDF>
```

Note: The vhostuser socket file path must be the one used when configuring the VM.



Start the VM.

```
virsh start <domain name>
```

For further information on the vDPA application, please visit: https://doc.dpdk.org/guides/ sample_app_ug/vdpa.html.

Appendix: Mellanox Firmware Tools

Download and install the MFT package corresponding to your computer's operating system. You would need the kernel-devel or kernel-headers RPM before the tools are built and installed. The package is available at http://www.mellanox.com => Products => Software => Firmware Tools.

1. Start the mst driver.

```
This Starting MST (Mellanox Software Tools) driver set Loading MST PCI module - Success Loading MST PCI configuration module - Success
```

2. Show the devices status.

```
ST modules:
    MST PCI module loaded
    MST PCI configuration module loaded
\# mlxconfig -d /dev/mst/mt4117_pciconf0 q | head -16
Device #1:
Device type: ConnectX41x
PCI device: /dev/mst/mt4117_pciconf0
Configurations:
        ations:
SRIOV_EN
NUM_OF_VFS
PF_LOG_BAR_SIZE
VF_LOG_BAR_SIZE
NUM_PF_MSIX
         NUM_VF_MSIX
LINK_TYPE_P1
         LINK_TYPE_P2
```

- 3. Make sure your configuration is as follows:
 - * SR-IOV is enabled (SRIOV_EN=1)
 - * The number of enabled VFs is enough for your environment (NUM_OF_VFS=N)
 - * The port's link type is Ethernet (LINK_TYPE_P1/2=2) when applicable If this is not the case, use mlxconfig to enable that, as follows:
 - a. Enable SR-IOV.

mlxconfig -d /dev/mst/mt4115_pciconf0 s SRIOV_EN=1

b. Set the number of required VFs.

mlxconfig -d /dev/mst/mt4115_pciconf0 s NUM_OF_VFS=8

c. Set the link type to Ethernet.

```
# mlxconfig -d /dev/mst/mt4115_pciconf0 s LINK_TYPE_P1=2
# mlxconfig -d /dev/mst/mt4115_pciconf0 s LINK_TYPE_P2=2
```

4. Reset the firmware.

mlxfwreset -d /dev/mst/mt4115_pciconf0 reset

5. Query the firmware to make sure everything is set correctly.

mlxconfig -d /dev/mst/mt4115_pciconf0 q

Programming



⚠ 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 <u>Raw Ethernet Programming: VXLAN Steering - Code Example</u> Community post.

Device Memory Programming

A

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, usergiven 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 on of the following options:

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

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	
dump_fts	Dumps tables for every switch found in an ibnetdiscover scan of the subnet The dump file format is compatible with loading into OpenSM using the -R file -U /path/to/dump-file syntax.	
	For further information, please refer to the tool's man page.	
ibaddr	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.	
ibcacheedit	Allows users to edit an ibnetdiscover cache created through thecache option in ibnetdiscover(8).	
	For further information, please refer to the tool's man page.	
ibccconfig	Supports the configuration of congestion control settings on switches and HCAs.	
	For further information, please refer to the tool's man page.	
ibccquery	Supports the querying of settings and other information related to congestion control.	
	For further information, please refer to the tool's man page.	
ibcongest	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 "/opt/ibutils2/bin/ibcongest -h".	
	For further information, please refer to the tool's man page.	

Utility	Description	
ibdev2netdev	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.	
ibdiagnet (of ibutils2)	For further information, please refer to the tool's man page. 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: • Fabric discovery • Duplicated GUIDs detection • Links in INIT state and unresponsive links detection • Counters fetch • Error counters check • Routing checks • Link width and speed checks • Alias GUIDs check • Subnet Manager check • Partition keys check • Nodes information Note: 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. For further information, either:	
	1. Run ibdiagnet -H Or 2. Refer to https://docs.mellanox.com/display/ibdiagnetUserManualv10	
ibdiagpath	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.	
	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: ibdiagpathdr_path 0,1). On the other hand, if LID-route addressing is employed,src_lid anddest_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.	
	Example: ibdiagpathsrc_lid 1dest_lid 28	
	For further information, please refer to the tool's -help flag.	

Utility	Description	
ibdump	Dump InfiniBand traffic that flows to and from Mellanox Technologies ConnectX® family adapters InfiniBand ports.	
	Note the following:	
	 ibdump is not supported for Virtual functions (SR-IOV) Infiniband traffic sniffing is supported on all HCAs The dump file can be loaded by the Wireshark tool for graphical traffic analysis. 	
	The following describes a workflow for local HCA (adapter) sniffing:	
	 Run ibdump with the desired options Run the application that you wish its traffic to be analyzed Stop ibdump (CTRL-C) or wait for the data buffer to fill (inmemmode) Open Wireshark and load the generated file To download Wireshark for a Linux or Windows environment go to www.wireshark.org. 	
	Notes:	
	 Although ibdump is a Linux application, the generated .pcap file may be analyzed on either operating system. If one of the HCA's ports is configured as InfiniBand, ibdump requires IPoIB DMFS to be enabled. For further information, please refer to Flow Steering Configuration section. For further information, please refer to the tool's man page. 	
iblinkinfo	Reports link info for each port in an InfiniBand fabric, node by node. Optionally, iblinkinfo can do partial scans and limit its output to parts of a fabric.	
	For further information, please refer to the tool's man page.	
ibnetdiscover	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).	
	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.	
	For further information, please refer to the tool's man page.	
ibnetsplit	Automatically groups hosts and creates scripts that can be run in order to split the network into sub-networks containing one group of hosts.	
	For further information, please refer to the tool's man page.	
ibnodes	Uses the current InfiniBand subnet topology or an already saved topology file and extracts the InfiniBand nodes (CAs and switches).	
	For further information, please refer to the tool's man page.	
ibping	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.	
	For further information, please refer to the tool's man page.	

Utility	Description		
ibportstate	Enables querying the logical (link) and physical port states of an InfiniBand port. It also allows adjusting the link speed that is enabled on any InfiniBand port.		
	If the queried port is a switch port, then ibportstate can be used to:		
	 disable, enable or reset the port validate the port's link width and speed against the peer port 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: 		
	 The first ACTIVE port that is found. If not found, the first port that is UP (physical link state is LinkUp). For further information, please refer to the tool's man page. 		
ibqueryerrors	The default behavior is to report the port error counters which exceed a threshold for each port in the fabric. The default threshold is zero (0). Error fields can also be suppressed entirely.		
	In addition to reporting errors on every port, ibqueryerrors can report the port transmit and receive data as well as report full link information to the remote port if available.		
	For further information, please refer to the tool's man page.		
ibroute	Uses SMPs to display the forwarding tables—unicast (LinearForwarding- Table or LFT) or multicast (MulticastForwardingTable or MFT)—for the specified switch LID and the optional lid (mlid) range. The default range is all valid entries in the range 1 to FDBTop.		
	For further information, please refer to the tool's man page.		
ibstat	ibstat is a binary which displays basic information obtained from the local IB driver. Output includes LID, SMLID, port state, link width active, and port physical state.		
	For further information, please refer to the tool's man page.		
ibstatus	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.		
ibswitches	Traces the InfiniBand subnet topology or uses an already saved topology file to extract the InfiniBand switches.		
	For further information, please refer to the tool's man page.		
ibsysstat	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.		
	For further information, please refer to the tool's man page.		

Utility	Description	
ibtopodiff	Compares a topology file and a discovered listing of subnet.lst/ibdiagnet.lst and reports mismatches.	
	Two different algorithms provided:	
	 Using the -e option is more suitable for MANY mismatches it applies less heuristics and provide details about the match Providing the -s, -p and -g starts a detailed heuristics that should be used when only small number of changes are expected For further information, please refer to the tool's man page. 	
ibtracert	Uses SMPs to trace the path from a source GID/LID to a destination GID/LID. Each hop along the path is displayed until the destination is reached or a hop does not respond. By using the -m option, multicast path tracing can be performed between source and destination nodes.	
	For further information, please refer to the tool's man page.	
ibv_asyncwatch	Display asynchronous events forwarded to userspace for an InfiniBand device.	
	For further information, please refer to the tool's man page.	
ibv_devices	Lists InfiniBand devices available for use from userspace, including node GUIDs.	
	For further information, please refer to the tool's man page.	
ibv_devinfo	Queries InfiniBand devices and prints about them information that is available for use from userspace.	
	For further information, please refer to the tool's man page.	
mstflint	Queries and burns a binary firmware-image file on non-volatile (Flash) memories of Mellanox InfiniBand and Ethernet network adapters. The tool requires root privileges for Flash access.	
To run mstflint, you must know the device location on the PCI bu		
	Note: If you purchased a standard Mellanox Technologies network adapter card, please download the firmware image from www.mellanox.com > Support > Firmware Download. If you purchased a non-standard card from a vendor other than Mellanox Technologies, please contact your vendor.	
	For further information, please refer to the tool's man page.	
perfquery	Queries InfiniBand ports' performance and error counters. Optionally, it displays aggregated counters for all ports of a node. It can also reset counters after reading them or simply reset them.	
	For further information, please refer to the tool's man page.	
saquery	Issues the selected SA query. Node records are queried by default. For further information, please refer to the tool's man page.	
sminfo	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.	
	Note: Using sminfo for any purpose other than a simple query might result in a malfunction of the target SM.	
	For further information, please refer to the tool's man page.	

Utility	Description	
smparquery	Sends SMP query for adaptive routing and private LFT features. For further information, please refer to the tool's man page.	
smpdump	A general purpose SMP utility which gets SM attributes from a specified SMA. The result is dumped in hex by default. For further information, please refer to the tool's man page.	
smpquery	Provides a basic subset of standard SMP queries to query Subnet management attributes such as node info, node description, switch info, and port info. For further information, please refer to the tool's man page.	

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 microbenchmark.

Utility	Description
ib_atomic_bw	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. For further information, please refer to the tool's man page.
ib_atomic_lat	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. For further information, please refer to the tool's man page.
ib_read_bw	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.
	Read is available only in RC connection mode (as specified in IB spec). For further information, please refer to the tool's man page.
ib_read_lat	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).
	For further information, please refer to the tool's man page.
ib_send_bw	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.
	For further information, please refer to the tool's man page.
ib_send_lat	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.
	For further information, please refer to the tool's man page.

Utility	Description	
ib_write_bw	Calculates the BW of RDMA write between a pair of machines. One acts as a server and the other as a client. The client RDMA writes to the server memory and 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.	
	For further information, please refer to the tool's man page.	
ib_write_lat	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.	
	For further information, please refer to the tool's man page.	
raw_ethernet_bw	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.	
	For further information, please refer to the tool's man page.	
raw_ethernet_lat	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.	
	For further information, please refer to the tool's man page.	

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.

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

General Issues

Issue	Cause	Solution
The system panics when it is booted with a failed adapter installed.	Malfunction hardware component	 Remove the failed adapter. Reboot the system.
Mellanox adapter is not identified as a PCI device.	PCI slot or adapter PCI connector dysfunctionality	 Run Ispci. Reseat the adapter in its PCI slot or insert the adapter to a different PCI slot. If the PCI slot confirmed to be functional, the adapter should be replaced.
Mellanox adapters are not installed in the system.	Misidentification of the Mellanox adapter installed	Run the command below and check Mellanox's MAC to identify the Mellanox adapter installed.
		lspci grep Mellanox' or 'lspci -d 15b3:
		Note: Mellanox MACs start with: 00:02:C9:xx:xx:xx, 00:25:8B:xx:xx:xx or F4:52:14:xx:xx:xx"
The default device may vary when invoking user apps (such as ibv_asyncwatch) which run using a specific device.	The default device for such apps is the first device in the device list generated by libibverbs. This first device in the list varies, depending on which and how many InfiniBand devices are installed on the host, which slot the devices are installed on, whether they use SR-IOV, and other factors.	Always specify the desired device explicitly when running userspace apps, by using the provided command line parameter (for example: ibv_asyncwatch -d <dev>).</dev>

Issue	Cause	Solution
Insufficient memory to be used by udev upon OS boot.	udev is designed to fork() new process for each event it receives so it could handle many events in parallel, and each udev instance consumes some RAM memory.	Limit the udev instances running simultaneously per boot by adding udev.children-max= <number> to the kernel command line in grub.</number>
Operating system running from root file system located on a remote storage (over Mellanox devices), hang during reboot/ shutdown (errors such as "No such file or directory" will appear).	The openibd service script is called using the 'stop' option by the operating system. This option unloads the driver stack. Therefore, the OS root file system disappears before the reboot/ shutdown procedure is completed, leaving the OS in a hang state.	Disable the openibd 'stop' option by setting 'ALLOW_STOP=no' in /etc/infiniband/openib.conf configuration file.
Mellanox adapter warning print to dmesg: Detected insufficient power on the PCIe slot (xxxW).	Insufficient PCI power.	Investigate the cause for lack of PCI power.

Ethernet Related Issues

Issue	Cause	Solution
Ethernet interfaces renaming fails leaving them with names such as renameXY.	Invalid udev rules.	Review the udev rules inside the "/etc/udev/rules.d/70-persistent-net.rules" file. Modify the rules such that every rule is unique to the target interface, by adding correct unique attribute values to each interface, such as dev_id, dev_port and KERNELS or address).
		Example of valid udev rules:
		SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?*",
		ATTR{dev_id}=="0x0", ATTR{type} =="1", KERNEL=="eth*",
		ATTR{dev_port}=="0", KER- NELS=="0000:08:00.0", NAME="eth4" SUBSYSTEM=="net", ACTION=="add", DRIVERS=="?*",
		ATTR{dev_id}=="0x0", ATTR{type} =="1", KERNEL=="eth*",
		ATTR{dev_port}=="1", KER- NELS=="0000:08:00.0", NAME="eth5"

Issue	Cause	Solution
No link.	Misconfiguration of the switch port or using a cable not supporting link rate.	 Ensure the switch port is not down Ensure the switch port rate is configured to the same rate as the adapter's port
Degraded performance is measured when having a mixed rate environment (10GbE, 40GbE and 56GbE).	Sending traffic from a node with a higher rate to a node with lower rate.	Enable Flow Control on both switch ports and nodes: • On the server side run: ethtool -A <interface> rx on tx on • On the switch side run the following command on the relevant interface: send on force and receive on force</interface>
No link with break-out cable.	Misuse of the break-out cable or misconfiguration of the switch's split ports	 Use supported ports on the switch with proper configuration. For further information, please refer to the MLNX_OS User Manual. Make sure the QSFP breakout cable side is connected to the SwitchX.

InfiniBand Related Issues

Issue	Cause	Solution
The following messages is logged after loading the driver: multicast join failed with status - 22	Trying to join a multicast group that does not exist or exceeding the number of multicast groups supported by the SM.	If this message is logged often, check for the multicast group's join requirements as the node might not meet them. Note: If this message is logged after driver load, it may safely be ignored.
Unable to stop the driver with the following on screen message: ERROR: Module <module> is in use</module>	An external application is using the reported module.	Manually unloading the module using the 'modprobe -r' command.
Logical link fails to come up while port logi- cal state is Initializing.	The logical port state is in the Initializing state while pending the SM for the LID assignment.	 Verify an SM is running in the fabric. Run 'sminfo' from any host connected to the fabric. If SM is not running, activate the SM on a node or on managed switch.

Issue	Cause	Solution
InfiniBand utilities commands fail to find devices on the system. For example, the 'ibv_devinfo' command fail with the following output:	The InfiniBand utilities commands are invoked when the driver is not loaded.	Load the driver: /etc/init.d/openibd start
Failed to get IB devices list: Function not implemented		

Installation Related Issues

Installation Issues

Issue	Cause	Solution
Driver installation fails.	The install script may fail for the following reasons: • Using an unsupported installation option • Failed to uninstall the previous installation due to dependencies being used • The operating system is not supported • The kernel is not supported. You can run mlnx_add_kernel_support.sh in order to to generate a MLNX-OFED package with drivers for the kernel • Required packages for installing the driver are missing • Missing kernel backport support for non supported kernel	 Use only supported installation options. The full list of installation options case be displayed on screen by using: mlnxofedinstallh Run 'rpm -e' to display a list of all RPMs and then manually uninstall them if the preliminary uninstallation failed due to dependencies being used. Use a supported operating system and kernel Manually install the missing packages listed on screen by the installation script if the installation failed due to missing prerequisites.
After driver installation, the openibd service fail to start. This message is logged by the driver: Unknown symbol	The driver was installed on top of an existing In-box driver.	 Uninstall the MLNX_OFED driver. ofed_uninstall.sh Reboot the server. Search for any remaining installed driver. If found, move them to the /tmp directory from the current directory. Re-install the MLNX_OFED driver. Restart the openibd service.

Fixing Application Binary Interface (ABI) Incompatibility with MLNX_OFED Kernel Modules

A

This section is relevant for RedHat and SLES distributions only.

Overview

MLNX_OFED package for RedHat comes with RPMs that support KMP (weak-modules), meaning that when a new errata kernel is installed, compatibility links will be created under the weak-updates directory for the new kernel. Those links allow using the existing MLNX_OFED kernel modules without the need for recompilation. However, at times, the ABI of the new kernel may not be compatible with the MLNX_OFED modules, which will prevent loading them. In this case, the MLNX_OFED modules must be rebuilt against the new kernel.

Detecting ABI Incompatibility with MLNX_OFED Modules

When MLNX_OFED modules are not compatible with a new kernel from a new OS or errata kernel, no links will be created under the weak-updates directory for the new kernel, causing the driver load to fail. Checking for the existence of needed module links under weak-updates directory can be done by reloading the MLNX_OFED modules. If one or more modules are missing, the driver reload will fail with an error message.

Example:

Resolving ABI Incompatibility with MLNX_OFED Modules

In order to fix ABI incompatibility with MLNX_OFED modules, the modules should be recompiled against the new kernel, using the mlnx_add_kernel_support.sh script, available in MLNX_OFED installation image.

There are two ways to recompile the MLNX_OFED modules:

Local recompilation and installation on one server.
 Run the mlnxofedinstall command to recompile the kernel modules and reinstall the whole MLNX_OFED on the server. Mount MLNX_OFED ISO image or extract the TGZ file:

```
# cd <MLNX_OFED dir>
# ./mlnxofedinstall --skip-distro-check --add-kernel-support --kmp --force
```

Notes:

- The --kmp flag will enable rebuilding RPMs with KMP (weak-updates) support for the new kernel. Therefore, in the next OS/kernel update, the same modules can be used with the new kernel (assuming that the ABI compatibility was not broken again).
- The command above will rebuild only the kernel RPMs (using mlnx_add_kernel_support.sh), and will save the resulting MLNX_OFED package under /tmp and start installing it

automatically. This package can be used for installation on other servers using regular mlnxofedinstall command or yum.

- 2. Preparing a new image on one server and deploying it on the cluster.
 - a. Use the mlnx_add_kernel_support.sh script directly only to rebuild the kernel RPMs (without running any installations) on one server. Mount MLNX_OFED ISO image or extract the TGZ file:

```
# cd <MLNX_OFED dir>
# ./mlnx_add_kernel_support.sh -m $PWD --kmp -y
```

Note: This command will save the resulting MLNX_OFED package under /tmp. Example:

- b. Install the newly created MLNX_OFED package on the cluster:
 - Option 1: Copy the package to the servers and install it using the mlnxofedinstall script.
 - Option 2: Deploy the MLNX_OFED package using YUM (for YUM installation instructions, refer to Installing MLNX_OFED Using YUM section):
 - i. Extract the resulting MLNX_OFED image and copy it to a shared NFS location.
 - ii. Create a YUM repository configuration.
 - iii. Install the new MLNX_OFED kernel RPMs on the servers: # yum update Example:

Note: The MLNX_OFED user-space packages will not change; only the kernel RPMs will be updated. However, "YUM update" can also update other inbox packages (not related to OFED). In order to install the MLNX_OFED kernel RPMs only, make sure to run:

yum install mlnx-ofed-kernel-only

Note: mlnx-ofed-kernel-only is a metadata RPM that requires the $MLNX_OFED$ kernel RPMs only.

c. Verify that the driver can be reloaded:

/etc/init.d/openibd restart

Performance Related Issues

Issue	Cause	Solution
The driver works but the transmit and/or receive data rates are not optimal.	-	These recommendations may assist with gaining immediate improvement:
		 Confirm PCI link negotiated uses its maximum capability Stop the IRQ Balancer service: /etc/init.d/irq_balancer stop Start mlnx_affinity service: mlnx_affinity start For best performance practices, please refer to the "Performance Tuning Guide for Mellanox Network Adapters".
Out of the box throughput performance in Ubuntu14.04 is not optimal and may achieve results below the line rate in 40GE link speed.	IRQ affinity is not set properly by the irq_balancer	For additional performance tuning, please refer to Performance Tuning Guide.

SR-IOV Related Issues

Issue	Cause	Solution
When assigning a VF to a VM the following message is reported on the screen: PCI-assgine: error: requires KVM support	SR-IOV and virtualization are not enabled in the BIOS.	 Verify they are both enabled in the BIOS Add to the GRUB configuration file to the following kernel parameter: "intel_immun=on" (see "Settin g Up SR-IOV" section).

PXE (FlexBoot) Related Issues

Issue	Cause	Solution
PXE boot timeout.	The 'always-broadcast' option is disabled.	Enable 'always-broadcast on'. For the complete procedure, please refer to Linux PXE User Guide.
PXE InfiniBand link fails with the following messages although the DHCP request was sent: Initializing and The socket is not connected.	Either the SM is not running in the fabric or the SM default multicast group was created with non-default settings.	 Activate the SM on a node or on managed switch. Check in the SM partitions.conf file that the default partition rate and MTU setting are SDR and 2K, respectively. The PXE is establishing by default an SDR link set with an MTU of 2K. If the default multicast group opened with different rate and/or MTU, the SM will deny the PXE request to join.
Mellanox adapter is not identified as a boot device.	The expansion ROM image is not installed on the adapter. or the server's BIOS is not configured to work on Legacy mode	 Run a flint query to display the expansion ROM information. For example: "flint -d /dev/mst/mt4099_pci_cr0 q" and look for the "Rom info:" line. For further information on how to burn the ROM, please refer to MFT User Manual. Make sure the BIOS is configured to work in Legacy mode if the adapter's firmware does not include a UEFI image.

RDMA Related Issues

Issue	Cause	Solution
Infiniband-diags tests, such as 'ib_write_bw', fail between systems with different driver releases.	When running a test between 2 systems in the fabric with different Infiniband-diags packages installed.	Run the test using the same perftest RPM on both systems.

Debugging Related Issues

Issue	Cause	Solution
False positive errors when running applications with valgrind.	Default MLNX_OFED libraries are compiled with- out valgrind support and several resources are managed by the kernel.	Libraries' files compiled with valgrind support are installed under "/usr/ lib64/mlnx_ofed/valgrind/" • To run an application over these libraries, thus prevent false positive errors: # env LD_LIBRARY_PATH=/usr/lib64/mlnx_ofed/valgrind/val-grind [valgrind options] <application cmd=""> • To suppress most of valgrind's false positive errors, generate the suppression file: #./generate_mlnx_ofed supp.sh > mlnx.supp</application>

OVS Offload Using ASAP2 Direct Related Issues

Issue	Cause	Solution(s)
Traffic is not offloaded	 OVS uses TC flower classifier to add offloading rules to both the software and the hardware. TC flower classifier fails to add a rule. A rule was added to the TC flower classifier but failed to be added to the firmware. 	 Check for system error in dmesg or the system logging facility like journalctl Check OVS logs for errors Dump the rules using the TC command line For example: Dump rules on a specific interface # tc filter show dev ens4f0 parent ffff:

Common Abbreviations and Related Documents

Common Abbreviations and Acronyms

Abbreviation/Acronym	Description
В	(Capital) 'B' is used to indicate size in bytes or multiples of bytes (e.g., 1KB = 1024 bytes, and 1MB = 1048576 bytes)
b	(Small) 'b' is used to indicate size in bits or multiples of bits (e.g., 1Kb = 1024 bits)
FW	Firmware
НСА	Host Channel Adapter
HW	Hardware
IB	InfiniBand
iSER	iSCSI RDMA Protocol
LSB	Least significant byte
lsb	Least significant <i>bit</i>
MSB	Most significant byte
msb	Most significant bit
NIC	Network Interface Card
SW	Software
VPI	Virtual Protocol Interconnect
IPolB	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

Abbreviation/Acronym	Description
VL	Virtual Lane
vНВА	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
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

Term	Description
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 Release Notes PDF file relevant to your adapter device on mellanox.com
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

User Manual Revision History

Rel ea se	Date	Description	
5.2	January 19, 2021	Added OpenSSL with kTLS Offload section.	
	January 4, 2021	 Added Offloaded Traffic Sniffer section. Added Tx Port Time-Stamping section. Added VLAN Push/Pop section. Added sFLOW section. Added E2E Cache section. Added Geneve Encapsulation/Decapsulation section. Added Parallel Offloads section. Updated SR-IOV VF LAG section. Removed Installing MLNX_OFED on Innova™ IPsec Adapter Cards section. Removed Updating Firmware and FPGA Image on Innova IPsec Cards section. 	
5.1- 2	September 17, 2020	Added <u>Packet Pacing for Hairpin Queues</u> section.	
5.1	July 28, 2020	Updated the content of the entire document following the removal of support for ConnectX-3, ConnectX-3 Pro and Connect-IB adapter cards, as well as the deprecation of RDMA experimental verbs library (mlnx_lib).	
		Added <u>SR-IOV Live Migration</u> section.	
		Added <u>SR-IOV VF LAG</u> section.	
5.0- 2	April 23, 2020	Added Interrupt Request (IRQ) Naming section.	
	April 6, 2020	Added Kernel Transport Layer Security (kTLS) Offloads section.	
5.0	March 3, 2020	 Added <u>IPSec Crypto Offload</u> section. Added <u>OVS-DPDK Hardware Offloads</u> section. Updated OVS Hardware Offloads Configuration section. 	
4.7	December 29, 2019	Added <u>Configuring Uplink Representor Mode</u> section.	
	December 13, 2019	 Added <u>Performance Tuning Based on Traffic Patterns</u> section. Added "num_of_groups" entry to table <u>mlx5_core Module Parameters</u>. Added <u>Mediated Devices</u> section. 	
	September 29, 2019	Updated Additional Installation Procedures section.	
4.6	May 13, 2019	 ethtool section updates: Added description of -f flashing option to <u>Ethtool Supported Options</u> table. 	
	April 30, 2019	 ethtool section updates: Updated the description of ethtool -s eth<x> advertise <n> autoneg on counter under Ethtool.</n></x> Added the following counters under Ethtool: ethtoolshow-fec eth<x></x> ethtoolset-fec eth<x> encoding auto off rs baser</x> Added Devlink Parameters section. Added Limit Bandwidth per Group of VFs section. Added Disabling ROCE section. Added RDMA-CM QP Timeout Control section. Added RDMA-CM Application Managed QP section. 	

Rel ea se	Date	Description
4.5	December 19, 2018	 Reorganized Chapter 2, "Installation": Consolidated the separate installation procedures under <u>Installing Mellanox OFED</u> and <u>Additional Installation Procedures</u> Added <u>Installing NEO-Host Using mlnxofedinstall Script</u>
	November 29, 2018	Added the following sections: • Local Loopback Disable • Offsweep Balancing

Release Notes Change Log History

Category	Description	
5.1-0.6.6.0		
	Adapters: ConnectX-4 and above	
IP-in-IP RSS Offload	Added support for receive side scaling (RSS) offload in IP-in-IP (IPv4 and IPv6).	
Devlink Port Support in Non- representor Mode	Added support for viewing the mlx5e physical devlink ports using the 'devlink port' command. This also may affect network interface names, if predictable naming scheme is configured. Suffix indicating a port number will be added to interface name.	
Devlink Health State Notifications	Added support for receiving notifications on devlink health state changes when an error is reported or recovered by one of the reporters. These notifications can be seen using the userspace 'devlink monitor' command.	
Legacy SR-IOV VF LAG Load Balancing	When VF LAG is in use, round-robin the Tx affinity of channels among the different ports, if supported by the firmware, enables all SQs of a channel to share the same port affinity. This allows the distribution of traffic sent from a VF between two ports, as well as round-robin the starting port among VFs to distribute traffic originating from single-core VMs.	
RDMA-CM DevX Support	Added support for DevX in RDMA-CM applications.	
RoCEv2 Flow Label and UDP Source Port Definition	This feature provides flow label and UDP source port definition in RoCE v2. Those fields are used to create entropy for network routes (ECMP), load balancers and 802.3ad link aggregation switching that are not aware of RoCE headers.	
RDMA Tx Steering	Enabled RDMA Tx steering flow table. Rules in this flow table will allow for steering transmitted RDMA traffic.	
Custom Parent- Domain Allocators for CQ	Enabled specific custom allocations for CQs.	
mlx5dv Helper APIs for Tx Affinity Port	Added support for the following mlx5dv helper APIs which enable the user application to query or set a RAW QP's Tx affinity port number in a LAG configuration.	
Selection	mlx5dv_query_qp_lag_portmlx5dv_modify_qp_lag_port	
RDMA-CM Path Alignment	Added support for RoCE network path alignment between RDMA-CM message and QP data. The drivers and network components in RoCE calculate the same hash results for egress port selection both on the NICs and the switches.	
IPoIB QP Number Creation	Enabled setting the QP number of an IPoIB PKey interface in Enhanced mode. This is done using the standard ip link add command while padding the hardware address of the newly created interface. The QP number is the 2nd-4th bytes. To enable the feature, the MKEY_BY_NAME configuration should firstly be enabled in the NvConfig.	
CQ and QP Context Exposure	Exposed QP, CQ and MR context in raw format via RDMA tool.	
In-Driver xmit_more	Enabled xmit_more feature by default in kernels that lack Rx bulking support (v4.19 and above) to ensure optimized IP forwarding performance when stress from Rx to Tx flow is insufficient.	
	In kernels with Rx bulking support, xmit_more is disabled in the driver by default, but can be enabled to achieve enhanced IP forwarding performance.	

Relaxed Ordering	Relaxed ordering is a PCIe feature which allows flexibility in the transaction order over the PCIe. This reduces the number of retransmissions on the lane, and increases performance up to 4 times.
	By default, mlx5e buffers are created with Relaxed Ordering support when firmware capabilities are on and the PCI subsystem reports that CPU is not on the kernel's blocklist.
	Note : Some CPUs which are not listed in the kernel's blocklist may suffer from buggy implementation of relaxed ordering, in which case the user may experience a degradation in performance and even unexpected behavior. To turn off relaxed ordering and restore previous behavior, run setpci command as instructed here . Example:
	"RlxdOrd-" : setpci -s82:00.0 CAP_EXP+8.w=294e
ODP Huge Pages Support	Enabled ODP Memory Region (MR) to work with huge pages by exposing IBV_ACCESS_HUGETLB access flag to indicate that the MR range is mapped by huge pages.
	The flag is applicable only in conjunction with IBV_ACCESS_ON_DEMAND.
Offloaded Traffic Sniffer	Removed support for Offloaded Traffic Sniffer feature and replaced its function with Upstream solution tcpdump tool.
	Adapters: ConnectX-5 and above
Connection Tracking Offload	Added support for offloading TC filters containing connection tracking matches and actions.
Dual-Port RoCE Support	Enabled simultaneous operation of dual-port RoCE and Ethernet in SwitchDev mode.
IP-in-IP Tunnel Offload for Checksum and TSO	Added support for the driver to offload checksum and TSO in IP-in-IP tunnels.
Packet Pacing DevX Support	Enabled RiverMax to work over DevX with packet pacing functionality by exposing a few DV APIs from rdma-core to enable allocating/destroying a packet pacing index. For further details on usage, see man page for: mlx5dv_pp_alloc() and mlx5dv_pp_free().
Software Steering Support for Memory Reclaiming	Added support for reclaiming device memory to the system when it is not in use. This feature is disabled by default and can be enabled using the command mlx5dv_dr_domain_set_reclaim_device_memory().
SR-IOV Live Migration	[Beta] Added support for performing a live migration for a VM with an SR-IOV NIC VF attached to it and with minimal to no traffic disruption. This feature is supported in SwitchDev mode; enabling users to fully leverage VF TC/OVS offloads, where the failover inbox driver is in the Guest VM, and the bonding driver is in the Hypervisor.
	Note that you must use the latest QEMU and libvirt from the Upstream github.com sources.
Uplink Representor Modes	Removed support for new_netdev mode in SwitchDev mode. The new default behaviour is to always keep the NIC netdev.
OVS-DPDK Offload Statistics	Added support for dumping connection tracking offloaded statistics.
OVS-DPDK Connection Tracking Labels Exact Matching	Added support for labels exact matching in OVS-DPDK CT openflow rules.
	Adapters: ConnectX-5 & ConnectX-6 Dx

Added support for LAG (modes 1,2,4) with OVS-DPDK.		
Adapters: ConnectX-6 and above		
Allowed configuration of Reed Solomon and Low Latency Reed Solomon over PAM4 link modes.		
Added support for allowing automatic enabling/disabling of vendor specific features during connection establishment between network nodes, which is performed over RDMA-CM messaging interface.		
This feature introduces a new QP retransmission mode in RoCE in which dropped packet recovery is done by re-sending the packet instead of re-sending the PSN window only (Go-Back-N protocol). This feature is enabled by default when RDMA-CM is being used and both connection nodes support it.		
Adapters: ConnectX-6 Dx & BlueField-2		
[Beta] Added support for IPsec full offload (VxLAN over ESP transport).		
Added support for configuring hardware vDPA on OVS-DPDK. This support includes the option to fall back to Software vDPA in case the NIC installed on the driver does not support hardware vDPA.		
Adapters: ConnectX-6 Dx		
Support for IPsec Crypto Offloads feature over ConnectX-6 Dx devices and up is now at GA level.		
Support for TLS Tx Hardware Offload feature over ConnectX-6 Dx devices and up is now at GA level.		
[Alpha] Added support for hardware offload decryption of TLS Rx traffic over cryptoenabled ConnectX-6 Dx NICs and above.		
Support for software steering on ConnectX-6 Dx adapter cards in the user-space RDMA-Core library through the mlx5dv_dr API is now at GA level.		
[Beta] Added support for kernel software steering on ConnectX-6 Dx adapter cards.		
Adapters: ConnectX-6 Lx		
Added support for ConnectX-6 Lx adapter cards.		
Adapters: All		
As of MLNX_OFED v5.1, Legacy verbs libraries have been fully replaced by RDMA-Core library.		
For the list of new APIs used for various MLNX_OFED features, please refer to the Migration to RDMA-Core document.		
Added support for safely inserting consecutive firmware images without the need to reset the NIC in between.		
Am Adr Trija II Aos SG Sa Le SC II A Ali FW A		

UCX-CUDA Support	UCX-CUDA is now supported on the following OS	s and platforms.
	OS	Platform
	RedHat 7.6 ALT	PPC64LE
	RedHat 7.7	x86_64
	RedHat 7.8	PPC64LE/x86_64
	RedHat 7.9	x86_64
	RedHat 8.1	x86_64
	RedHat 8.2	x86_64
HCOLL-CUDA	The hcoll package includes a CUDA plugin (hmcait is built on various platforms as the package has if the system has CUDA 10-2 installed. Notes: If you install MLNX_OFED from a package package hcoll-cuda explicitly to be able HCOLL-CUDA is supported on the same O (listed in the table above), except for Reference in the same of	e repository, you will need to install the to use it. Ss that include support for UCX-CUDA
GPUDirect Storage (GDS)	[Beta] Added support for the new technology of a direct data path between local or remote stor Fabric (NVMe-oF), and GPU memory. Both GPUD extra copies through a bounce buffer in the CPL memory access (DMA) engine near the NIC or sto or out of GPU memory, without burdening the C	age, such as NFS, NVMe or NVMe over irect RDMA and GPUDirect Storage avoid I's memory. They enable the direct prage to move data on a direct path into
	To enable the feature, run ./mlnxofedinstall enable-gdsadd-kernel-support	with-nfsrdmawith-nvmf
	To get access to GDS Beta, please reach out to t GPUDirectStorageExt@nvidia.com.	he GDS team at
	For the list of operating systems on which GDS i	s supported, see <u>here</u> .
Unable to render include or excerpt-include. Could not retrieve page.		
	Adapters: ConnectX-5 and abo	ove
Kernel Software Managed Flow Steering (SMFS) Performance	Improved the performance of Kernel software st consumption.	eering by reducing its memory
	Adapters: All	
NEO-Host SDK	Added support for NEO-Host SDK installation on	MLNX_OFED.
Bug Fixes	See <u>Bug Fixes</u> .	
Unable to render include or excerpt-include. Could not retrieve page.		
	Adapters: ConnectX-6 Dx	

Adapters	Added support for ConnectX-6 Dx adapter cards.
Userspace Software Steering ConnectX-6 Dx Support	[Beta] Added support for software steering on ConnectX-6 Dx adapter cards in the user-space RDMA-Core library through the mlx5dv_dr API.
	Adapters: ConnectX-6 Dx and above
Virtual Output Queuing (VoQ) Counters	Exposed rx_prio[p]_buf_discard, rx_prio[p]_wred_discard and rx_prio[p]_marked firmware counters that count the number of packets that were dropped due to insufficient resources.
IPsec Crypto Offloads	[Beta] IPsec crypto offloads are now supported on ConnectX-6 Dx devices and up. The offload functions use the existing ip xfrm tool to activate offloads on the device. It supports transport/tunnel mode with AES-GCM IPsec scheme.
TLS TX Hardware Offload	[Alpha] Added support for hardware offload encryption of TLS traffic.
VirtIO Acceleration through Datapath I/ O Processor (vDPA)	Added support to enable mapping the VirtIO access region (VAR) to be used for doorbells by vDPA applications. Specifically, the following DV APIs were introduced (see man page for more details):
	mlx5dv_alloc_var()mlx5dv_free_var()
	Adapters: ConnectX-5 and above
Resource Allocation on External Memory	Added support to enable overriding mlx5 internal allocations in order to let applications allocate some resources on external memory, such as that of the GPU.
	The above is achieved by extending the parent domain object with custom allocation callbacks. Currently supported verbs objects are: QP, DBR, RWQ, SRQ.
Hardware Clock Exposure	Added support for querying the adapter clock via mlx5dv_query_device.
ODP Diagnostic Counters	Added ODP diagnostics counters for the following items per MR (memory region) within IB/mlx5 driver:
	 Page faults: Total number of faulted pages. Page invalidations: Total number of pages invalidated by the OS during all invalidation events. The translations can no longer be valid due to either non-present pages or mapping changes. Prefetched pages: When prefetching a page, a page fault is generated in order to bring the page to the main memory.
Devlink Health CR- Space Dump	Added the option to dump configuration space via the devlink tool in order to improve debug capabilities.
Multi-packet TX WQE Support for XDP Transmit Flows	The conventional TX descriptor (WQE or Work Queue Element) describes a single packet for transmission. Added driver support for the HW feature of multi-packet TX WQEs in XDP transmit flows. With this, the HW becomes capable of working with a new and improved WQE layout that describes several packets. In effect, this feature saves PCI bandwidth and transactions, and improves transmit packet rate.
OVS-Kernel ToS Rewrite	Added support for Type of Service (ToS) rewrite in the OVS-Kernel.
OVS-Kernel Mirroring	Added support for mirroring output in SwitchDev mode in the OVS-Kernel. The mirroring port may either be a local or a remote VF, using VxLAN or GRE encapsulations.
GENEVE Encap/ Decap Rules Offload	Added support for GENEVE encapsulation/decapsulation rules offload.

GPRS Tunneling Protocol (GTP) Header	[Beta] Added support for matching (filtering) GTP header-based packets using mlx5dv_dr API over user-space RDMA-Core library.
Multi Packet Tx WQE Support for XDP Transmit Flows	Added driver support for the hardware feature of multi-packet Tx to work with a new and improved WQE layout that describes several packets instead of a single packet for XDP transmission flows. This saves PCI bandwidth and transactions, and improves transmit packet rate.
Userspace Software Steering Debugging API	[Beta] Added support for software steering to dump flows for debugging purposes in the user-space RDMA-Core library through the mlx5dv_dr API.
Kernel Software Steering for Connection Tracking (CT)	[Beta] Added support for updating CT rules using the software steering mechanism.
Kernel Software Steering Remote Mirroring	[Beta] Added support for updating remote mirroring rules using the software steering mechanism.
	Adapters: ConnectX-5 and BlueField
OVS-DPDK Support	Added OVS-DPDK component as part of the MLNX_OFED package with hardware offload capabilities.
OVS-DPDK Connection Tracking	[Beta] Added support for OvS-DPDK Connection Tracking hardware offload.
OVS-DPDK VirtIO Acceleration through VF Relay	Added support for OVS-DPDK VirtIO Acceleration through VF Relay (also known as Software vDPA) forwarding of traffic from VF to Virtio VM and vice-versa.
OVS-DPDK VXLAN Encap/Decap	Added support for OVS-DPDK VXLAN encapsulation and decapsulation hardware offload.
	Adapters: ConnectX-4 and above
Discard Counters	Exposed rx_prio[p]_discards discard counters per priority that count the number of received packets dropped due to lack of buffers on the physical port.
MPLS Traffic	Added support for reporting TSO and CSUM offload capabilities for MPLS tagged traffic and, allowed the kernel stack to use these offloads.
mlx5e Max Combined Channels	Increased the driver's maximal combined channels value from 64 to 128 (however, note that OOB value will not cross 64). 128 is the upper bound. Lower maximal value can be seen on the host, depending on the number of cores and MSIX's configured by the firmware.
RoCE Accelerator Counters	Added the following RoCE accelerator counters: • roce_adp_retrans - counts the number of adaptive retransmissions for RoCE traffic • roce_adp_retrans_to - counts the number of times RoCE traffic reached timeout due to adaptive retransmission • roce_slow_restart - counts the number of times RoCE slow restart was used • roce_slow_restart_cnps - counts the number of times RoCE slow restart generated CNP packets • roce_slow_restart_trans - counts the number of times RoCE slow restart changed state to slow restart
	All Adapters

The default installation of the userspace is now the RDMA-Core library instead of the legacy verbs. This achieves most of the legacy experimental verbs' functionalities, and more.	
For Mellanox VMA or Mellanox RiverMax, use experimental verbs (prefix "ibv_exp").	
For further information on the migration to RDMA-Core and the list of new APIs used for various MLNX_OFED features, please refer to the <u>Migration to RDMA-Core document</u> .	
ibdev2netdev tool output was changed such that the bonding device now points at the bond instead of the slave interface.	
Added support for the user to register memory regions with a relaxed ordering access flag. This can enhance performance, depending on architecture and scenario.	
Added support for monitoring and recovering from errors that occur on the RX queue, such as cookie errors and timeout.	
Improved GSO (Generic Segmentation Offload) workload performance by decreasing doorbells usage to the minimum required.	
Added support for TX CQE (Completion Queue Element) compression. Saves on outgoing PCIe bandwidth by compressing CQEs together. Disabled by default. Configurable via private flags of ethtool.	
Added the option to query for running and stored firmware versions using the devlink tool.	
Added the option to update the firmware image in the flash using the devlink tool.	
Usage: devlink dev flash <dev> file <file_name>.mfa2</file_name></dev>	
For further information on how to perform this update, see "Updating Firmware Using ethtool/devlink and .mfa2 File" section in MFT User Manual.	
Added support for WQE (Work Queue Element) dump, triggered by an error on Rx/Tx reporters. In addition, some dumps (not triggered by an error) can be retrieved by the user via devlink health reporters.	
Added support for GENEVE tunneled hardware offloads of TSO, CSUM and RSS.	
Added TCP segmentation and checksum offload support for MPLS-tagged traffic.	
4.7-3.2.9.0	
HCAs: ConnectX-5 and above	
Added support for new_netdev and nic_netdev uplink representor modes.	

mlx5_core	Added new mlx5_core module parameter "num_of_groups", which controls the number of large groups in the FDB flow table.		
	Note : The default value of num_of_groups may change per MLNX_OFED driver version. The following table lists the values that <u>must</u> be set when upgrading the MLNX_OFED version prior to driver load , in order to achieve the same OOB experience.		
	MLNX_OFED Version	num_of_groups Default Value	
	v4.7-3.2.9.0	4	
	v4.6-3.1.9.0.14	15	
	v4.6-3.1.9.0.15	15	
	v4.5-1.0.1.0.19	63	
	For further information, please refer to section in MLNX_OFED User Manual.	Performance Tuning Based on Traffic Patterns	
ConnectX-5			
VFs Groups Minimum Bandwidth Rate	Added support for setting a minimum bandwidth rate on a group of VFs (BW guarantee) to ensure this group is able to transmit at least the amount of bandwidth specified on the wire.		
Direct Verbs Support for Batch Counters on Root Table	Added support for mlx5dv_dr API to set batch counters for root tables.		
ConnectX-5 & BlueField			
Modify Header	Added support for mlx5dv_dr_actions to support up to 32 modify actions.		
mlx5dv_dr Memory Consumption	Reduced the mlx5dv_dr API memory cons	sumption by improving the memory allocator.	
mlx5dv_dr Memory Allocation	Reduced memory allocation time when using the mlx5dv_dr API. This is particularly significant for the first inserted rules on which memory is allocated.		
	BlueField		
Mediated Devices	Added support for mediated devices that without SR-IOV on the Bluefield® system	allows the creation of accelerated devices .	
	For further information on mediated dev to Mediated Devices section in MLNX_OF	ices and how to configure them, please refer ED User Manual.	
4.7-1.0.0.1			
	HCAs: ConnectX-4 and	d above	
Counters Monitoring	Added support for monitoring selected of (Monitor_Counter_Change event) upon of the counters to be monitored are selected command.		

Signature Offload Kernel Verbs Enhancements	Added a new API which enables posting a single WR that completes the Protection Information (PI) operation internally. This reduces CPU utilization for posting and processing multiple WRs and improves performance by choosing the optimal mkey for the hardware according to the buffer memory layout.
EEPROM Device Thresholds via Ethtool	Added support to read additional EEPROM information from high pages of modules such as SFF-8436 and SFF-8636. Such information can be: 1. Application Select table 2. User writable EEPROM 3. Thresholds and alarms - Ethtool dump works on active cables only (e.g. optic), but thresholds and alarms can be read with "offset" and "length" parameters in any cable by running: ethtool -m <devname> offset X length Y</devname>
Performance Improvements	 Updated Blueflame capability reporting to prevent redundant use of Blueflame when Write-combining is not supported. Added Blueflame capabilities over VFs.
RDMA_RX RoCE Steering Support	Added the ability to create rules to steer RDMA traffic, with two destinations supported: DevX object and QP. Multiple priorities are also supported.
SRQ and XRC Support on On Demand Paging (ODP) Memory Region (MR)	Added support for using ODP MR with SRQ WQEs and XRC transport.
Indirect Mkey ODP	Added the ability to create indirect Mkeys with ODP support over DevX interface.
DevX Asynchronous Query Commands	Added support for running QUERY commands over the DevX interface in an asynchronous mode. This enables applications to issue many commands in parallel while firmware processes the commands.
Implicit ODP	Added support for reporting implicit ODP support to user applications in order to allow better granularity over ODP creation.
Devlink Health Utility	Added support for real-time alerting of functionality issues that may be found in a system component (reporter). This utility helps detect and recover from a problem with a PCI device. It provides a centralize status of drivers' health activities in the generic Devlink instance and inter alia, supports the following:
	 Storing real-time error dumps Performing automatic (configurable) real-time reporter recovery Performing real-time reporter diagnosis Indicating real-time reporter's health status Providing admins with the ability to dump, diagnose and recover a reporter Providing admins with the ability to configure a reporter
User-Mode Memory Registration (UMR)	Enabled registration of memory patterns that can be used for future RDMA operations.
GENEVE Tunnel Stateless Offload	Added support for Generic Network Virtualization Encapsulation (GENEVE) tunneled hardware offload of TSO, CSUM and RSS.
ODP Pre-fetch	Added support for pre-fetching a range of an on-demand paging (ODP) memory region (MR), this way reducing latency by making pages present with RO/RW permissions before the actual IO is conducted.
Fragmented QPs Buffer	Added the ability to allocate a fragmented buffer to in-kernel QP creation requests, in cases of large QP size requests that used to fail due to low memory resources on the host.
Flow Counters Batch Query	Allowed flow counters created with the DevX interface to be attached to flows created with the raw flow creation API.
DevX Privilege Enforcement	Enforced DevX privilege by firmware. This enables future device functionality without the need to make driver changes unless a new privilege type is introduced.

DevX Interoperability APIs WQ, and IND_TBLAPIs) via the DevX interface. Added support for montifying and/or querying for a verb object (including CQ, QP, SRQ, Interoperability APIs) with the DevX interface. This enables interoperability between verbs and DevX. Counters Monitoring Added support for monitoring selected counters and generating a notification event (Monitor_Counter_Change event) upon changes made to these counters. The counters to be monitored are selected using the SET_MONITOR_COUNTER command. RX Hash Fields Configuration Added the ability to configure Rx hash fields used for traffic spreading into Rx queues using ETHTOOL_SRXFH and ETHTOOL_GRXFH ethtool commands. Built-in Receive Side Scaling (RSS) profiles can now be changed on the following traffic types: UDP4, UDP6, TCP4 and TCP6. This configuration affects both outer and inner headers. **HCAs: ConnectX-4 Lx and above** HCAs: ConnectX-4 Lx and above Equal Cost Multi-Pace and Pace and Pace and Pace and Pace and Inner headers. **HCAs: ConnectX-5 Lx and above** UPLIA Added support for lifty Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNiCs. **HCAs: ConnectX-5** **HCAs: ConnectX-5** Userspace Software Steering for eSwitch Added software steering capabilities to the SR-IOV eSwitch. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs Userspace Software Steering for NICs Userspace Software steering capabilities to NIC Rx-/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution w				
Counters Monitoring Added support for monitoring selected counters and generating a notification event (Monitor_Counter_Change event) upon changes made to these counters. The counters to be monitored are selected using the SET_MONITOR_COUNTER command. Rx Hash Fields Configuration Added the ability to configure Rx hash fields used for traffic spreading into Rx queues using ETHTOOL_SRXFH and ETHTOOL_GRXFH ethtool commands. Built-in Receive Side Scaling (RSS) profiles can now be changed on the following traffic types: UDP4, UDP6, TCP4 and TCP6. This configuration affects both outer and inner headers. HCAs: ConnectX-4 Lx and above Equal Cost Multi-Path (ECMP) Added support for offloading ECMP rules by tracking software multipath route and related next-hops, and reflecting this as port affinity to the hardware. VF LAG Added support for High Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNiCs. HCAs: ConnectX-5 Userspace Software Steering for eSwitch Steering for eSwitch Added software steering capabilities to the SR-IOV eSwitch. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NiCs Added software steering capabilities to NiC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above HCAs: ConnectX-5 and above HCAs: ConnectX-5 and above Incoparation and flow counters are now set per Process of the process				
(Monitor_Counter_Change event) upon changes made to these counters. The counters to be monitored are selected using the SET_MONITOR_COUNTER command. Rx Hash Fields Configuration Added the ability to configure Rx hash fields used for traffic spreading into Rx queues using ETHTOOL_SRXFH and ETHTOOL_GRXFH ethtool commands. Built-in Receive Side Scaling (RS5) profiles can now be changed on the following traffic types: UDP4, UDP6, TCP4 and TCP6. This configuration affects both outer and inner headers. **HCAs: ConnectX-4 Lx and above** HCAs: ConnectX-4 Lx and above		This enables interoperability between verbs and DevX.		
command. Rx Hash Fields Configuration Added the ability to configure Rx hash fields used for traffic spreading into Rx queues using ETHTOOL_SRXFH and ETHTOOL_GRXFH ethtool commands. Built-in Receive Side Scaling (RSS) profiles can now be changed on the following traffic types: UDP4, UDP6, TCP4 and TCP6. This configuration affects both outer and inner headers. HCAs: ConnectX-4 Lx and above Equal Cost Multi-Path (ECMP) Added support for offloading ECMP rules by tracking software multipath route and related next-hops, and reflecting this as port affinity to the hardware. VF LAG Added support for High Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNiCs. HCAs: ConnectX-5 Userspace Software Steering for eSwitch Added software steering capabilities to the SR-IOV eSwitch. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs Verspace Software Steering for Nics Note: Support will be enabled by default once the support for GID change is added. Verspace Software Note:	Counters Monitoring			
Configuration scaling (RSS) profiles can now be changed on the following traffic types: UDP4, UDP6, TCP4 and TCP6. This configuration affects both outer and inner headers. HCAs: ConnectX-4 Lx and above Equal Cost Multi-Path (ECMP) Feath (ECMP) Feath (ECMP) VF LAG Added support for offloading ECMP rules by tracking software multipath route and related next-hops, and reflecting this as port affinity to the hardware. VF LAG Added support for High Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNICs. HCAs: ConnectX-5 Weerspace Software Steering capabilities to the SR-IOV eSwitch. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Weerspace Software steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Note: Support will be enabled by default once the support for GID change is added. **ConnectX-5** and above** HCAs: ConnectX-5** and above** HCAs: ConnectX-5** and above** ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PiD QP counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be		_		
Equal Cost Multi-Path (ECMP) Redded support for offloading ECMP rules by tracking software multipath route and related next-hops, and reflecting this as port affinity to the hardware. VF LAG Added support for High Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNICs. HCAs: ConnectX-5 Added software steering capabilities to the SR-IOV eSwitch. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs Userspace Software Steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firm of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA > MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets of packets that are transmitted inline depends on CPUs load, where lower load leads to		using ETHTOOL_SRXFH and ETHTOOL_GRXFH ethtool commands. Built-in Receive Side Scaling (RSS) profiles can now be changed on the following traffic types: UDP4, UDP6,		
Path (ECMP) related next-hops, and reflecting this as port affinity to the hardware. VF LAG Added support for High Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNICs. HCAs: ConnectX-5 Userspace Software Steering Capabilities to the SR-IOV eSwitch. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above HCAs: ConnectX-5 and above ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID QP counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Address Translation Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		HCAs: ConnectX-4 Lx and above		
physical ports in SwitchDev SR-IOV mode. Uplink Representors Exposed PF (uplink) representors in SwitchDev mode, similarly to VF representors, as an infrastructure improvement for SmartNICs. HCAs: ConnectX-5 Userspace Software Steering for eSwitch achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs Added software steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Added software steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Smalt Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to				
Userspace Software Steering for eSwitch better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs achieved by performing calculations on the main CPU which allows for higher insertion rates. Added software steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above HCAs: ConnectX-5 and above Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID GP Counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of South Services that are transmitted inline depends on CPUs load, where lower load leads to	VF LAG			
Userspace Software Steering for eSwitch Steering for eSwitch better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering for NICs Added software steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above HCAs: ConnectX-5 and above ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID Firmware Errors per PID Firmware Errors per PID Firmware Errors per PID Added support for On-Demand Paging (ODP) over DC transport. Added support for On-Demand Paging (ODP) over DC transport. Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to	Uplink Representors			
Steering for eSwitch better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. Userspace Software Steering capabilities to NIC Rx/Tx. Software steering enables better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID QP counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Services would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		HCAs: ConnectX-5		
Steering for NICs rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates. This solution was designed to work with Virtio DPDK. Note: Support will be enabled by default once the support for GID change is added. HCAs: ConnectX-5 and above ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID QP counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		better rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion		
HCAs: ConnectX-5 and above Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID QP counters and GP counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		rules insertion rate compared to the current firmware-based solution. This is achieved by performing calculations on the main CPU which allows for higher insertion rates.		
ASAP ² Incorporated the documentation of Accelerated Switching And Packet Processing (ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per PID QP counters and flow counters are now set per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		Note : Support will be enabled by default once the support for GID change is added.		
(ASAP²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User Manual. QP Counters and Firmware Errors per Process ID (PID) to allow better visibility of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		HCAs: ConnectX-5 and above		
Firmware Errors per PID of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the QP type. ODP over DC Added support for On-Demand Paging (ODP) over DC transport. Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to	ASAP ²	(ASAP ²): Hardware Offloading for vSwitches into MLNX_OFED Release Notes and User		
Address Translation Services Added support for Address Translation Services (ATS) feature, which improves performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets Added support for when forwarding packets with XDP, a packet smaller than 256 bytes would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to	Firmware Errors per	of RDMA error states. Users will be able to manually tune the Q counter to monitor specific QPs, or automatically monitor QPs according to predefined criteria, such as the		
Services performance for virtualized PeerDirect applications by caching PA-> MA translations and preventing PCI transactions from going to the root complex. XDP Inline Transmission of Small Packets performance in line within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to	ODP over DC	Added support for On-Demand Paging (ODP) over DC transport.		
Transmission of Small Packets would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		performance for virtualized PeerDirect applications by caching PA-> MA translations and		
	Transmission of	would be sent inline within its WQE Tx descriptor for better performance. The number of packets that are transmitted inline depends on CPUs load, where lower load leads to		

VLAN Rewrite	Added support for offloading VLAN ID modify operation, allowing the user to replace the VLAN tag of the incoming frame with a user-specified VLAN tag value.	
CQE Padding	Added support for padding 64B CQEs to 128B cache lines to improve performance on 128B cache line systems, such as PPC.	
XDP Multi-Packet Tx Work Queue Element (WQE)	Added support for Multi-Packet Tx WQEs in XDP transmit flows to work with a new and improved WQE layout that describes several packets. This saves PCI bandwidth and transactions, and improves transmit packet rate.	
	HCAs: ConnectX-6	
ConnectX Device IDs	Added support for the following new device IDs:	
	 ConnectX-6 Dx (PF) ConnectX Family mlx5Gen Virtual Function (VF) Note that every new device (adapter) VF will be identified with this device ID. Different VF models will be distinguished by their revision ID. 	
	HCAs: ConnectX-6 and above	
Ethtool 200Gbps	ConnectX-6 hardware introduces support for 200Gbps and 50Gbps-per-lane link mode. The driver supports full backward compatibility with previous configurations. Note that in order to advertise newly added link-modes, the full bitmap related to the link modes must be advertised from ethtool man page.	
	NOTE : This feature is firmware-dependent. Currently, ConnectX-6 Ethernet firmware supports up to 100Gbps only. Thus, this capability may not function properly using the current driver and firmware versions.	
HDR Link Speed Exposure	Added support for HDR link speed in CapabilityMask2 field in port attributes.	
QP Packet Based Credit Mode	Added support for an alternative end-to-end credit mode for QP creation. Credits transported from the responder to the requester are now issued per packet. This is particularly useful for sending large RDMA messages from HCA to switches that are short in memory.	
	HCAs: BlueField	
Device Emulation Infrastructure	Added support for Device Emulation in BlueField. This mechanism allows function-A to perform operations on behalf of function-B. The emulation manager creates a channel (named VHCA_TUNNEL general object) that acts as the direct command interface between the emulated function host and the HCA hardware. The emulation software creates this tunnel for every managed function and issues commands via the DevX general command interface.	
HCAs: All		
Verbs Migration to RDMA-Core	Legacy verbs remain the default userspace installation option in the MLNX_OFED. However, as of MLNX_OFED v4.7, you can opt to install full RDMA-Core based userspace by adding theupstream-libs flag to the mlnxofedinstall script.	
MLNX_OFED Installation via	The repository providing legacy verbs has been moved from RPMS or DEBS folders to RPMS/MLNX_LIBS and DEBS/MLNX_LIBS.	
Repository	In addition, a new repository providing RDMA-Core based userspace has been added to RPMS/UPSTREAM_LIBS and DEBS/UPSTREAM_LIBS.	

NFSoRDMA	Added support for NFS over RDMA (NFSoRDMA) module over the OSs listed in NFSoRDMA Supported OSs section.
	As of MLNX_OFED v4.7, NFSoRDMA driver is no longer installed by default. In order to install it over a supported kernel, add the "with-nfsrdma" installation option to the "mlnxofedinstall" script.
RDMA-CM QP Timeout Control	Added a new option to rdma_set_option that allows applications to override the RDMA-CM's QP ACK timeout value.
Object IDs Exportation	Added a unique ID for each verbs object to allow direct query over rdma-tool and rdma-netlink for enhanced debuggability.
RDMA-CM Application Managed QP	Added support for the RDMA application to manage its own QPs and use RDMA-CM only for exchanging Address information.
Bug Fixes	See "Bug Fixes" section.
	4.6-1.0.1.0
	HCAs: ConnectX-3/ConnectX-3 Pro
Devlink Configuration Parameters Tool	Added support for a set of configuration parameters that can be changed by the user through the Devlink user interface.
	HCAs: ConnectX-4 and above
ODP Pre-fetch	Added support for pre-fetching a range of an on-demand paging (ODP) memory region (MR), this way reducing latency by making pages present with RO/RW permissions before the actual IO is conducted.
DevX Privilege Enforcement	Enforced DevX privilege by firmware. This enables future device functionality without the need to make driver changes unless a new privilege type is introduced.
DevX Interoperability APIs	Added support for modifying and/or querying for a verb object (including CQ, QP, SRQ, WQ, and IND_TBL APIs) via the DevX interface.
	This enables interoperability between verbs and DevX.
DevX Asynchronous Query Commands	Added support for running QUERY commands over the DevX interface in an asynchronous mode. This enables applications to issue many commands in parallel while firmware processes the commands.
DevX User-space PRM Handles	Exposed all PRM handles to user-space so DevX user application can mix verbs objects with DevX objects.
Exposure	For example: Take the cqn from the created ibv_cq and use it on a devx)create(QP).
Indirect Mkey ODP	Added the ability to create indirect Mkeys with ODP support over DevX interface.
XDP Redirect	Added support for XDP_REDIRECT feature for both ingress and egress sides. Using this feature, incoming packets on one interface can be redirected very quickly into the transmission queue of another capable interface. Typically used for load balancing.
RoCE Disablement	Added the option to disable RoCE traffic handling. This enables forwarding of traffic over UDP port 4791 that is handled as RoCE traffic when RoCE is enabled.
	When RoCE is disabled, there is no GID table, only Raw Ethernet QP type is supported and RoCE traffic is handled as regular Ethernet traffic.

Forward Error Correction (FEC) Encoding	Added the ability to query and modify Forward Error Correction (FEC) encoding, as well as disabling it via Ethtool.		
RAW Per-Lane Counters Exposure	Exposed RAW error counters per cable-module lane via ethtool stats. The counters show the number of errors before FEC correction (if enabled).		
	For further information, please see phy_raw_errors_lane[i] under Physical Port Counters section in <u>Understanding mlx5 ethtool Counters</u> Community post.		
	HCAs: ConnectX-4 Lx and above		
VF LAG	Added support for High Availability and load balancing for Virtual Functions of different physical ports in SwitchDev SR-IOV mode.		
	HCAs: ConnectX-5 and above		
ASAP ² Offloading VXLAN Decapsulation with HW LRO	Added support for performing hardware Large Receive Offload (HW LRO) on VFs with HW-decapsulated VXLAN. For further information on the VXLAN decapsulation feature, please refer to ASAP ² User Manual under www.mellanox.com -> Products -> Software -> ASAP ² .		
PCI Atomic Operations	Added the ability to run atomic operations on local memory without involving verbs API or compromising the operation's atomicity.		
Equal-Cost Multi- Path (ECMP) Routing Offloading	Enabled Equal-Cost Multi-Path (ECMP) Routing offloading. Equal-Cost Multi-Path (ECMP) is a forwarding mechanism for routing packets along multiple paths of equal cost with the goal to achieve almost equally distributed link load sharing.		
VXLAN over VLAN	VXLAN over VLAN enables the user to use VXLAN offloads' benefit to offload VLAN tagged tunnels thus boost system's performance.		
VLAN Rewrite	Rewriting VLAN tags allows the user to replace the VLAN tag of the incoming frame with a user-specified VLAN tag value.		
	HCAs: ConnectX-5		
Virtual Ethernet Port Aggregator (VEPA)	Added support for activating/deactivating Virtual Ethernet Port Aggregator (VEPA) mode on a single virtual function (VF). To turn on VEPA on the second VF, run: echo ON > /sys/class/net/enp59s0/device/sriov/1/vepa		
VFs Rate Limit	Added support for setting a rate limit on groups of Virtual Functions rather on an individual Virtual Function.		
	HCAs: ConnectX-6		
ConnectX-6 Support	[Beta] Added support for ConnectX-6 (VPI only) adapter cards.		
	NOTE : In HDR installations that are built with remotely managed Quantum-based switches, the switch's firmware must be upgraded to version 27.2000.1142 prior to upgrading the HCA's (ConnectX-6) firmware to version 20.25.1500. When using ConnectX-6 HCAs with firmware v20.25.1500 and connecting them to Quantum-based switches, make sure the Quantum firmware version is 27.2000.1142 in order to avoid any critical link issues.		

Ethtool 200Gbps	ConnectX-6 hardware introduces support for 200Gbps and 50Gbps-per-lane link mode. MLNX_OFED supports full backward compatibility with previous configurations.
	Note that in order to advertise newly added link-modes, the full bitmap related to the link modes must be advertised from ethtool man page. For the full bitmap list per link mode, please refer to MLNX_OFED User Manual.
	NOTE : This feature is firmware-dependent. Currently, ConnectX-6 Ethernet firmware supports up to 100Gbps only. Thus, this capability may not function properly using the current driver and firmware versions.
PCIe Power State	Added support for the following PCIe power state indications to be printed to dmesg:
	 Info message #1: PCle slot power capability was not advertised. Warning message: Detected insufficient power on the PCle slot (xxxW). Info message #2: PCle slot advertised sufficient power (xxxW). When indication #1 or #2 appear in dmesg, user should make sure to use a PCle slot that is capable of supplying the required power.
	HCAs: mlx5
Message Signaled Interrupts-X (MSI-X) Vectors	Added support for using a single MSI-X vector for all control event queues instead of one MSI-X vector per queue in a virtual function driver. This frees extra MSI-X vectors to be used for completion event queue, allowing for additional traffic channels in the network device.
Send APIs	Introduced a new set of QP Send operations (APIs) which allows extensibility for new Send opcodes.
DC Data-path	Added DC QP data-path support using new Send APIs introduced in Direct Verbs (DV).
	HCAs: BlueField
BlueField Support	BlueField is now fully supported as part of the Mellanox OFED mainstream version sharing the same code baseline with all the adapters product line.
Representor Name	In SwitchDev mode:
Change	 Uplink representors are now called p0/p1 Host PF representors are now called pf0hpf/pf1hpf VF representors are now called pf0vfN/pf1vfN
ECPF Net Devices	In SwitchDev mode, net devices enp3s0f0 and enp3s0f1 are no longer created.
Setting Host MAC and Tx Rate Limit from ECPF	Expanded to support VFs as well as the host PFs.
HCAs: All	
RDMA-CM Application Managed QP	Added support for the RDMA application to manage its own QPs and use RDMA-CM only for exchanging Address information.
RDMA-CM QP Timeout Control	Added a new option to rdma_set_option that allows applications to override the RDMA-CM's QP ACK timeout value.
MLNX_OFED Verbs API	As of MLNX_OFED v5.0 release (Q1 2020) onwards, MLNX_OFED Verbs API will be migrated from the legacy version of the user space verbs libraries (libibervs, libmlx5) to the upstream version rdma-core. More details are available in MLNX_OFED user manual under Installing Upstream rdma-core Libraries.
Bug Fixes	See <u>Bug Fixes</u> section.

4.5-1.0.1.0		
HCAs: ConnectX-5		
VFs per PF	Increased the amount of maximum virtual functions (VF) that can be allocated to a physical function (PF) to 127 VF.	
	HCAs: ConnectX-4/ ConnectX-4 Lx/ConnectX-5	
SW-Defined UDP Source Port for RoCE v2	UDP source port for RoCE v2 packets is now calculated by the driver rather than the firmware, achieving better distribution and less congestion. This mechanism works for RDMA- CM QPs only, and ensures that RDMA connection messages and data messages have the same UDP source port value.	
	HCAs: mlx5 Driver	
Local Loopback Disable	Added the ability to manually disable Local Loopback regardless of the number of open user-space transport domains.	
HCAs: ConnectX-6		
Adapter Cards	Added support for ConnectX-6 Ready. For further information, please contact <u>Mellanox Support</u> .	
HCAs: All		
NEO-Host	Integrated NEO-Host for orchestration and management of host networking into MLNX_OFED package.	
Bug Fixes	See <u>Bug Fixes</u> section.	
	4.4-2.0.7.0	
	HCAs: All	
Bug Fixes	See <u>Bug Fixes</u> section.	
	4.4-1.0.0.0	
	HCAs: ConnectX-4/ConnectX-4 Lx/ConnectX-5	
Adaptive Interrupt Moderation	Added support for adaptive Tx, which optimizes the moderation values of the Tx CQs on runtime for maximum throughput with minimum CPU overhead.	
	This mode is enabled by default.	
	Updated Adaptive Rx to ignore ACK packets so that queues that only handle ACK packets remain with the default moderation.	
Docker Containers [Beta]	Added support for Docker containers to run over Virtual RoCE and InfiniBand devices using SR-IOV mode.	
VF Statistics	Performed the following virtual function statistics changes:	
	Added tx_broadcast and tx_multicast countersIncluded RDMA statistics for existing counters	
Force TTL	Added support for setting a global TTL value for all RC QPs and rdma-cm QPs.	

Firmware Tracer	Added a new mechanism for the device's FW/HW to log important events into the event tracing system (/sys/kernel/debug/tracing) without requiring any Mellanox-specific tool.
	Note: This feature is enabled by default.
CR-Dump	Accelerated the original cr-dump by optimizing the reading process of the device's CR-Space snapshot.
RoCE ICRC Error Counter	Added support for a new counter that exposes the amount of corrupted RoCE packets that arrive with bad Invariant Cyclic Redundancy Code (ICRC).
	HCAs: ConnectX-4/ConnectX-4 Lx
VST Q-in-Q	Added support for C-tag (0x8100) VLAN insertion to tagged packets in VST mode.
	HCAs: ConnectX-4
Ethernet Tunneling Over IPoIB Driver (eIPoIB)	Re-added support for eth_ipoib driver, which provides a standard Ethernet interface to be used as a Physical Interface (PIF) into the Hypervisor virtual network, and serves one or more Virtual Interfaces (VIF).
	HCAs: ConnectX-4 Lx/ConnectX-5
OVS Offload using ASAP ²	Added support for Mellanox Accelerated Switching And Packet Processing (ASAP ²) technology, which allows OVS offloading by handling OVS data-plane, while maintaining OVS control-plane unmodified. OVS Offload using ASAP ² technology provides significantly higher OVS performance without the associated CPU load.
	For further information, refer to ASAP ² Release Notes under <u>www.mellanox.com</u> -> Products -> Software -> ASAP ²
	HCAs: All
Upstream Libraries	Added a repository repodata to support installing upstream libraries (based on upstream rdma-core), using he Operating System's standard package manager (yum, apt-get, etc.).
	For further information, please refer to "Installing Upstream rdma-core Libraries" section in MLNX_OFED User Manual
	Note: This is intended only for DPDK users.
Installation	Added support for new metadata packages that only install userspace packages at a time (without any kernel packages), using the Operating System's standard package manager (yum, apt-get, etc.). These metadata packages will have the suffix "-user-only". For example: "mlnx-ofed-all-user-only".
Bug Fixes	See <u>Bug Fixes</u> section.
	4.3-1.0.1.0
	HCAs: ConnectX-4/ConnectX-4 Lx/ConnectX-5
Adaptive Interrupt Moderation	Added support for adaptive Tx, which optimizes the moderation values of the Tx CQs on runtime for maximum throughput with minimum CPU overhead.
	This mode is enabled by default.

	Updated Adaptive Rx to ignore ACK packets so that queues that only handle ACK packets remain with the default moderation.
Docker Containers [Beta]	Added support for Docker containers to run over Virtual RoCE and InfiniBand devices using SR-IOV mode.
VF Statistics	Performed the following virtual function statistics changes:
	Added tx_broadcast and tx_multicast countersIncluded RDMA statistics for existing counters
Force TTL	Added support for setting a global TTL value for all RC QPs and rdma-cm QPs.
Firmware Tracer	Added a new mechanism for the device's FW/HW to log important events into the event tracing system (/sys/kernel/debug/tracing) without requiring any Mellanox-specific tool.
	Note: This feature is enabled by default.
CR-Dump	Accelerated the original cr-dump by optimizing the reading process of the device's CR-Space snapshot.
RoCE ICRC Error Counter	Added support for a new counter that exposes the amount of corrupted RoCE packets that arrive with bad Invariant Cyclic Redundancy Code (ICRC).
	HCAs: ConnectX-4/ConnectX-4 Lx
VST Q-in-Q	Added support for C-tag (0x8100) VLAN insertion to tagged packets in VST mode.
	HCAs: ConnectX-4
Ethernet Tunneling Over IPoIB Driver (eIPoIB)	Re-added support for eth_ipoib driver, which provides a standard Ethernet interface to be used as a Physical Interface (PIF) into the Hypervisor virtual network, and serves one or more Virtual Interfaces (VIF).
	HCAs: ConnectX-5/ConnectX-4 Lx
OVS Offload using ASAP2	Added support for Mellanox Accelerated Switching And Packet Processing (ASAP2) technology, which allows OVS offloading by handling OVS data-plane, while maintaining OVS control-plane unmodified. OVS Offload using ASAP2 technology provides significantly higher OVS performance without the associated CPU load.
	For further information, refer to ASAP ² Release Notes under <u>www.mellanox.com</u> -> Products -> Software -> ASAP2
	HCAs: All
Upstream Libraries	Added a repository repodata to support installing upstream libraries (based on upstream rdma-core), using the Operating System's standard package manager (yum, apt-get, etc.).
	For further information, please refer to "Installing Upstream rdma-core Libraries" section in MLNX_OFED User Manual
	Note: This is intended only for DPDK users.
Installation	Added support for new metadata packages that only install userspace packages at a time (without any kernel packages), using the Operating System's standard package manager (yum, apt-get, etc.). These metadata packages will have the suffix "-user-only". For example: "mlnx-ofed-all-user-only".

Bug Fixes	See <u>Bug Fixes</u> section.	
4.3-1.0.1.0		
	HCAs: ConnectX-5	
Multi-Packet Work Request (WR)	Added support for the following multi-packet WR related verbs for control path: • ibv_exp_query_device • ibv_exp_create_srq For further information on the use of these verbs, please refer to the Verbs man page.	
Coherent Accelerator Processor Interface (CAPI) [beta]	Added support for CAPI, an interface that enables ConnectX-5 adapter cards to provide the best performance for Power and OpenPower based platforms.	
Tunneled Atomic	Added support for RDMA atomic commands offload so that when an RDMA Write operation is issued, the payload indicates which atomic operation to perform, instead of being written to the Memory Region (MR).	
Packet Pacing	Added support for the following advanced burst control parameters: • max_burst_sz - for indicating the maximal burst size of packets • typical_pkt_sz - for improving the accuracy of the rate limiter	
Erasure Coding Offload verbs	Added support for erasure coding offload software verbs (encode/decode/update API) supporting a number of redundancy blocks (m) greater than 4.	
	HCAs: ConnectX-4/ConnectX-4 Lx/ConnectX-5	
Virtual MAC	Removed support for Virtual MAC feature.	
RoCE LAG	Added out of box RoCE LAG support for RHEL 7.2 and RHEL 6.9.	
Dropped Counters	Added a new counter <i>rx_steer_missed_packets</i> which provides the number of packets that were received by the NIC, yet were discarded/dropped since they did not match any flow in the NIC steering flow table.	
	Added the ability for SR-IOV counter <i>rx_dropped</i> to count the number of packets that were dropped while vport was down.	
Relaxed Ordering (RSYNC)	Added support for RSYNC feature to ensure correct ordering of memory operations between the GPU and HCA.	
HCAs: mlx5 Driver		
Reset Flow	Added support for triggering software reset for firmware/driver recovery. When fatal errors occur, firmware can be reset and driver reloaded.	
HCAs: ConnectX-4 Lx/ConnectX-5		
Striding RQ with HW Time-Stamping		
HCAs: ConnectX-4/ConnectX-4 Lx/ConnectX-5		
4.2-1.2.0.0		

DSCP Trust Mode	Added support for automatically setting the number of TC to 8 when the Trust state is changed to DSCP.
Receive Buffer	Added xon and xoff columns to the Receive Buffer configuration display.
	4.2-1.0.0.0
	HCAs: mlx5 Driver
Physical Address Memory Allocation	Added support to register a specific physical address range.
	HCAs: Innova IPsec EN
Innova IPsec Adapter Cards	Added support for Mellanox Innova IPsec EN adapter card, that provides security acceleration for IPsec-enabled networks.
HCAs: ConnectX-4/ConnectX-4 Lx/ConnectX-5	
Precision Time	Added support for PTP feature over PKEY interfaces.
Protocol (PTP)	This feature allows for accurate synchronization between the distributed entities over the network. The synchronization is based on symmetric Round Trip Time (RTT) between the master and slave devices, and is enabled by default.
1PPS Time Synchronization	Added support for One Pulse Per Second (1PPS) over IPoIB interfaces.
Virtual MAC	Added support for Virtual MAC feature, which allows users to add up to 4 virtual MACs (VMACs) per VF. All traffic that is destined to the VMAC will be forwarded to the relevant VF instead of PF. All traffic going out from the VF with source MAC equal to VMAC will go to the wire also when Spoof Check is enabled.
	For further information, please refer to "Virtual MAC" section in MLNX_OFED User Manual.
Receive Buffer	Added the option to change receive buffer size and cable length. Changing cable length will adjust the receive buffer's xon and xoff thresholds.
	For further information, please refer to "Receive Buffer" section in MLNX_OFED User Manual.
GRE Tunnel Offloads	Added support for the following GRE tunnel offloads:
	 TSO over GRE tunnels Checksum offloads over GRE tunnels RSS spread for GRE packets
NVMEoF	Added support for the host side (RDMA initiator) in RedHat 7.2 and above.
Dropless Receive Queue (RQ)	Added support for the driver to notify the FW when SW receive queues are overloaded.
PFC Storm Prevention	Added support for configuring PFC stall prevention in cases where the device unexpectedly becomes unresponsive for a long period of time. PFC stall prevention disables flow control mechanisms when the device is stalled for a period longer than the default pre-configured timeout. Users now have the ability to change the default timeout by moving to auto mode.
	For further information, please refer to "PFC Stall Prevention" section in MLNX_OFEDUser Manual.

Force DSCP	Added support for this feature that enables setting a	
	global traffic_class value for all RC QPs.	
	HCAs: ConnectX-5	
Q-in-Q	Added support for Q-in-Q VST feature in ConnectX-5 adapter cards family.	
Device Memory Programming [beta]	Added support for on-chip memory allocation and usage in send/receive and RDMA operations at beta level.	
Virtual Guest Tagging (VGT+)	Added support for VGT+ in ConnectX-4/ConnectX-5 HCAs. This feature is s an advanced mode of Virtual Guest Tagging (VGT), in which a VF is allowed to tag its own packets as in VGT, but is still subject to an administrative VLAN trunk policy. The policy determines which VLAN IDs are allowed to be transmitted or received. The policy does not determine the user priority, which is left unchanged.	
	For further information, please refer to "Virtual Guest Tagging (VGT+)" section in MLNX_OFED User Manual.	
Tag Matching Offload	Added support for hardware Tag Matching offload with Dynamically Connected Transport (DCT).	
	HCAs: ConnectX-3/ConnectX-3 Pro	
Shared Memory Region (MR)	Removed support for Shared MR feature on ConnectX-3/ConnectX-3 Pro adapter cards. As a result of this change, the following API/flags should not be used:	
	ibv_exp_reg_shared_mraccess shared flags for ibv_exp_reg_mr (IBV_EXP_ACCESS_SHARED_MR_XXX)	
	HCAs: All	
CR-DUMP	Added support for the driver to take an automatic snapshot of the device's CR-Space in cases of critical failures.	
	For further information, please refer to "CRDUMP" section in MLNX_OFED User Manual.	
Upstream Libraries	Added the option to install upstream libraries (based on upstream rdma-core) for DPDK users only.	
	For further information, please refer to "Installing Upstream rdma-core Libraries" section in MLNX_OFED User Manual.	
DiSNI	Added the option to install libdisni package as part of MLNX_OFED.	
	For further information, please refer to section "Installing libdisni Package" in MLNX_OFED User Manual.	
Service Scripts	Added the ability to disable the 'stop' option in the openibd service script, by setting ALLOW_STOP=no in /etc/infiniband/openib.conf.	
	Starting from the next release, 'stop' option will be disabled by default, and in order to enable it, ALLOW_STOP should be set to 'yes' in the conf file, or force-stop should be run.	
4.1-1.0.2.0		
HCAs: mlx5 Driver		

RoCE Diagnostics and ECN Counters	Added support for additional RoCE diagnostics and ECN congestion counters under /sys/class/infiniband/mlx5_0/ports/1/hw_counters/ directory.
	For further information, refer to the <u>Understanding mlx5 Linux Counters and Status Parameters</u> Community post.
rx-fcs Offload (ethtool)	Added support for rx-fcs ethtool offload configuration. Normally, the FCS of the packet will be truncated by the ASIC hardware before sending it to the application socket buffer (skb). Ethtool allows to set the rx-fcs not to be truncated, but to pass it to the application for analysis.
	For more information and usage, refer to <u>Understanding ethtool rx-fcs for mlx5</u> <u>Drivers</u> Community post.
DSCP Trust Mode	Added the option to enable PFC based on the DSCP value. Using this solution, VLAN headers will no longer be mandatory for use.
	For further information, refer to the <u>HowTo Configure Trust Mode on Mellanox</u> <u>Adapters</u> Community post.
RoCE ECN Parameters	ECN parameters have been moved to the following directory: /sys/kernel/debug/mlx5/ <pci bus="">/cc_params/</pci>
	For more information, refer to the <u>HowTo Configure DCQCN (RoCE CC) for ConnectX-4 (Linux)</u> Community post.
Flow Steering Dump Tool	Added support for mlx_fs_dump, which is a python tool that prints the steering rules in a readable manner.
Secure Firmware Updates	Firmware binaries embedded in MLNX_OFED package now support Secure Firmware Updates. This feature provides devices with the ability to verify digital signatures of new firmware binaries, in order to ensure that only officially approved versions are installed on the devices.
	For further information on this feature, refer to Mellanox Firmware Tools (MFT) User Manual.
Enhanced IPoIB	Added support for Enhanced IPoIB feature, which enables better utilization of features supported in ConnectX-4 adapter cards, by optimizing IPoIB data path and thus, reaching peak performance in both bandwidth and latency.
	Enhanced IPoIB is enabled by default.
PeerDirect	Added the ability to open a device and create a context while giving PCI peer attributes such as name and ID.
	For further details, refer to the <u>PeerDirect Programming</u> Community post.
Probed VFs	Added the ability to disable probed VFs on the hypervisor. For further information, see HowTo Configure and Probe VFs on mlx5 Drivers Community post.
Local Loopback	Improved performance by rendering Local loopback (unicast and multicast) disabled by mlx5 driver by default while local loopback is not in use. The mlx5 driver keeps track of the number of transport domains that are opened by user-space applications. If there is more than one user-space transport domain open, local loopback will automatically be enabled.
1PPS Time Synchronization (at alpha level)	Added support for One Pulse Per Second (1PPS), which is a time synchronization feature that allows the adapter to send or receive 1 pulse per second on a dedicated pin on the adapter card.
	For further information on this feature, refer to the <u>HowTo Test 1PPS on Mellanox</u> <u>Adapters</u> Community post.

Precision Time	Added support for PTP feature in IPoIB offloaded devices.
Protocol (PTP)	This feature allows for accurate synchronization between the distributed entities over the network.
	The synchronization is based on symmetric Round Trip Time (RTT) between the master and slave devices.
	The feature is enabled by default.
	For further information, refer to <u>Running Linux PTP with ConnectX-4</u> Community post.
Fast Driver Unload	Added support for fast driver teardown in shutdown and kexec flows.
	HCAs: ConnectX-5/ConnectX-5 Ex
NVMEoF Target Offload	Added support for NVMe over fabrics (NVMEoF) offload, an implementation of the new NVMEoF standard target (server) side in hardware.
	For further information on NVMEoF Target Offload, refer to <u>HowTo Configure NVMEoF Target Offload</u> .
MPI Tag Matching	Added support for offloading MPI tag matching to HCA.
	HCAs: All
RDMA CM	Changed the default RoCE mode on which RDMA CM runs to RoCEv2 instead of RoCEv1.
	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 <u>RDMA CM and RoCE Version Defaults</u> Community post.
Lustre	Added support for Lustre file system open-source project.
	4.0-2.0.2.0
Operating Systems	Added support for Ubuntu v17.04.
	4.0-2.0.0.1
PCIe Error Counting	[ConnectX-4/ConnectX-4 Lx] Added the ability to expose physical layer statistical counters to ethtool.
Multiprotocol Label Switching (MPLS) Tagged Packets Classification	[ConnectX-4/ConnectX-4 Lx] Enabled packet flow steering rules with IPv4/IPv6 classification (for raw packet QP (DPDK) only) to work on IPv4/IPv6 over MPLS (Ethertype 0x8847 and 0x8848) encapsulated packets.
RoCE VFs	[ConnectX-4/ConnectX-4 Lx] Added the ability to enable/disable RoCE on VFs.
RoCE LAG	[ConnectX-4/ConnectX-4 Lx] Added support for RoCE over LAG interface.
Standard ethtool	[ConnectX-4/ConnectX-4 Lx] Added support for flow steering and rx-all mode.
SR-IOV Bandwidth Share for Ethernet/ RoCE (beta)	[ConnectX-4/ConnectX-4 Lx] Added the ability to guarantee the minimum rate of a certain VF in SR-IOV mode.
Adapter Cards	Added support for ConnectX-5 and ConnectX-5 Ex HCAs.

DSCP ConfigFS Control for RDMA- CM QPs	Added the ability to configure ToS/DSCP for RDMA-CM QPs only.
Soft RoCE (beta)	Add software implementation of RoCE that allows RoCE to run on any Ethernet network adapter whether it offers hardware acceleration or not.
NVMe over Fabrics (NVMEoF)	NVMEoF related module installation has been disabled by default. In order to enable it, add the "with-nvmf" installation option to the "mlnxofedinstall" script.
NFS over RDMA (NFSoRDMA)	Removed support for NFSoRDMA drivers. These drivers are no longer provided along with the MLNX_OFED package.

Notice

This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. Neither NVIDIA Corporation nor any of its direct or indirect subsidiaries (collectively: "NVIDIA") make any representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code, or functionality.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.

Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer ("Terms of Sale"). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer's own risk. NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer's sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure the product is suitable and fit for the application planned by customer, and perform the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer's product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this document. NVIDIA accepts no liability related to any default, damage, costs, or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this document or (ii) customer product designs. No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA. Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices. THIS DOCUMENT AND ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's



aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.

Trademarks

NVIDIA, the NVIDIA logo, and Mellanox are trademarks and/or registered trademarks of Mellanox Technologies Ltd. and/or NVIDIA Corporation in the U.S. and in other countries. Other company and product names may be trademarks of the respective companies with which they are associated. For the complete and most updated list of Mellanox trademarks, visit http://www.mellanox.com/page/trademarks

Copyright

© 2021 Mellanox Technologies Ltd. All rights reserved.

