



# Mellanox ConnectX-4/ConnectX-5 NATIVE ESXi Driver for VMware vSphere 5.5/6.0 User Manual

---

Rev 4.15.8.8/4.5.8.8



NOTE:

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



Mellanox Technologies  
350 Oakmead Parkway Suite 100  
Sunnyvale, CA 94085  
U.S.A.  
[www.mellanox.com](http://www.mellanox.com)  
Tel: (408) 970-3400  
Fax: (408) 970-3403

© Copyright 2017. Mellanox Technologies Ltd. All Rights Reserved.

Mellanox®, Mellanox logo, Accelio®, BridgeX®, CloudX logo, CompustorX®, Connect-IB®, ConnectX®, CoolBox®, CORE-Direct®, EZchip®, EZchip logo, EZappliance®, EZdesign®, EZdriver®, EZsystem®, GPUDirect®, InfiniHost®, InfiniBridge®, InfiniScale®, Kotura®, Kotura logo, Mellanox CloudRack®, Mellanox CloudXMellanox®, Mellanox Federal Systems®, Mellanox HostDirect®, Mellanox Multi-Host®, Mellanox Open Ethernet®, Mellanox OpenCloud®, Mellanox OpenCloud Logo®, Mellanox PeerDirect®, Mellanox ScalableHPC®, Mellanox StorageX®, Mellanox TuneX®, Mellanox Connect Accelerate Outperform logo, Mellanox Virtual Modular Switch®, MetroDX®, MetroX®, MLNX-OS®, NP-1c®, NP-2®, NP-3®, Open Ethernet logo, PhyX®, PlatformX®, PSIPHY®, SiPhy®, StoreX®, SwitchX®, Tiler®, Tiler logo, TestX®, TuneX®, The Generation of Open Ethernet logo, UFM®, Unbreakable Link®, Virtual Protocol Interconnect®, Voltaire® and Voltaire logo are registered trademarks of Mellanox Technologies, Ltd.

All other trademarks are property of their respective owners.

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

# Table of Contents

<b>Table of Contents</b> .....	<b>3</b>
<b>List of Tables</b> .....	<b>4</b>
<b>Document Revision History</b> .....	<b>5</b>
<b>About this Manual</b> .....	<b>6</b>
<b>Chapter 1 Introduction</b> .....	<b>8</b>
1.1 nmlx5 Driver .....	8
1.2 Mellanox NATIVE ESXi Package for ConnectX-4/ConnectX-5 .....	8
1.2.1 Software Components .....	8
1.3 Module Parameters .....	8
1.3.1 nmlx5 Module Parameters .....	8
<b>Chapter 2 Installation</b> .....	<b>11</b>
2.1 Hardware and Software Requirements .....	11
2.2 Installing Mellanox NATIVE ESXi Driver for VMware vSphere .....	11
2.3 Removing the Previous Mellanox Driver .....	12
2.4 Firmware Programming .....	12
<b>Chapter 3 Features Overview and Configuration</b> .....	<b>13</b>
3.1 Ethernet Network .....	13
3.1.1 Wake-on-LAN (WoL) .....	13
3.1.2 Set Link Speed .....	14
3.1.3 Priority Flow Control (PFC) .....	14
3.1.4 Receive Side Scaling (RSS) .....	15
3.2 Virtualization .....	16
3.2.1 Single Root IO Virtualization (SR-IOV) .....	16
3.2.2 VXLAN Hardware Offload .....	19
<b>Chapter 4 Troubleshooting</b> .....	<b>21</b>
4.1 General Related Issues .....	21
4.2 Ethernet Related Issues .....	21
4.3 Installation Related Issues .....	22

## List of Tables

Table 1:	Document Revision History.....	5
Table 2:	Abbreviations and Acronyms.....	6
Table 3:	Reference Documents.....	7
Table 4:	nmlx5_core Module Parameters.....	9
Table 5:	Software and Hardware Requirements.....	11
Table 6:	General Related Issues.....	21
Table 7:	Ethernet Related Issues.....	21
Table 8:	Installation Related Issues.....	22

## Document Revision History

**Table 1 - Document Revision History**

Release	Date	Description
Rev 4.15.8.8/4.5.8.8	January 31, 2017	<ul style="list-style-type: none"> <li>Added support for ConnectX-5 adapter cards.</li> </ul>
Rev 4.15.6.22/ 4.5.6.22	September, 2016	<ul style="list-style-type: none"> <li>Added the following sections:               <ul style="list-style-type: none"> <li>Section 3.1.2, “Set Link Speed”, on page 14</li> <li>Section 3.1.3, “Priority Flow Control (PFC)”, on page 14</li> <li>Section 3.1.4, “Receive Side Scaling (RSS)”, on page 15</li> <li>Section 3.1.4.1, “Default Queue Receive Side Scaling (DRSS)”, on page 15</li> <li>Section 3.1.4.2, “NetQ RSS”, on page 15</li> <li>Section 3.1.4.3, “Important Notes”, on page 16</li> </ul> </li> </ul>
Rev 4.15.5.10/ 4.5.5.10	July, 2016	<ul style="list-style-type: none"> <li>Updated the following section:               <ul style="list-style-type: none"> <li>Section 1.3.1.1, “nmlx5_core Parameters”, on page 9</li> </ul> </li> </ul>
Rev 4.15.4.1100/ 4.5.2.1100	January, 2016	<ul style="list-style-type: none"> <li>Added the following sections:               <ul style="list-style-type: none"> <li>Section 3.1.1, “Wake-on-LAN (WoL)”, on page 13</li> <li>Section 3.2.1, “Single Root IO Virtualization (SR-IOV)”, on page 16</li> <li>Section 3.2.2, “VXLAN Hardware Offload”, on page 19</li> </ul> </li> <li>Updated the following section:               <ul style="list-style-type: none"> <li>Section 1.3.1.1, “nmlx5_core Parameters”, on page 9</li> <li>Section 2.2, “Installing Mellanox NATIVE ESXi Driver for VMware vSphere”, on page 11 (Step 4)</li> <li>Section 2.4, “Firmware Programming”, on page 12</li> </ul> </li> </ul>
Rev 4.15.2.0	September, 2015	Initial release of the Initial release of this MLNX-NATIVE-ESX-ConnectX-4 version

## About this Manual

This preface provides general information concerning the scope and organization of this User's Manual.

## Intended Audience

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

## Common Abbreviations and Acronyms

**Table 2 - Abbreviations and Acronyms**

Abbreviation / Acronym	Whole Word / Description
B	(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
HCA	Host Channel Adapter
HW	Hardware
LSB	Least significant <i>byte</i>
lsb	Least significant <i>bit</i>
MSB	Most significant <i>byte</i>
msb	Most significant <i>bit</i>
NIC	Network Interface Card
SW	Software
VPI	Virtual Protocol Interconnect
PR	Path Record
RDS	Reliable Datagram Sockets
SDP	Sockets Direct Protocol
SL	Service Level
MPI	Message Passing Interface
QoS	Quality of Service
ULP	Upper Level Protocol
vHBA	Virtual SCSI Host Bus adapter
uDAPL	User Direct Access Programming Library

## Related Documentation

**Table 3 - Reference Documents**

Document Name	Description
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. For further information please refer to the Mellanox website. <a href="http://www.mellanox.com">www.mellanox.com</a> -> Support -> Firmware Download
MFT User Manual	Mellanox Firmware Tools User's Manual. For further information please refer to the Mellanox website. <a href="http://www.mellanox.com">www.mellanox.com</a> -> Products -> Ethernet Drivers -> Firmware Tools
MFT Release Notes	Release Notes for the Mellanox Firmware Tools. For further information please refer to the Mellanox website. <a href="http://www.mellanox.com">www.mellanox.com</a> -> Products -> Ethernet Drivers -> Firmware Tools
VMware vSphere 6.0 Documentation Center	VMware website

# 1 Introduction

Mellanox ConnectX®-4/ConnectX-5 NATIVE ESXi is a software stack which operates across all Mellanox network adapter solutions supporting up to 50Gb/s Ethernet (ETH) and 2.5 or 5.0 GT/s PCI Express 2.0 and 3.0 uplinks to servers.

The following sub-sections briefly describe the various components of the Mellanox ConnectX-4/ConnectX-5 NATIVE ESXi stack.

## 1.1 nmlx5 Driver

`nm1x5` is the low level driver implementation for the ConnectX-4/ConnectX-5 adapter cards designed by Mellanox Technologies. ConnectX-4/ConnectX-5 adapter cards can operate as an InfiniBand adapter, or as an Ethernet NIC. The ConnectX-4/ConnectX-5 NATIVE ESXi driver supports Ethernet NIC configurations. To accommodate the supported configurations, the driver consist of `mlnx5_core` module.

### `nm1x5_core`

A 10/25/40/50GigE driver that handles Ethernet specific functions and plugs into the ESXi uplink layer

## 1.2 Mellanox NATIVE ESXi Package for ConnectX-4/ConnectX-5

### 1.2.1 Software Components

MLNX-NATIVE-ESX-ConnectX-4/ConnectX-5 contains the following software components:

- Mellanox Host Channel Adapter Drivers
  - `nm1x5_core` (Ethernet)

## 1.3 Module Parameters

### 1.3.1 `nm1x5` Module Parameters

To set `nm1x5` parameters:

```
esxcli system module parameters set -m nm1x5_core -p <parameter>=<value>
```

To show all parameters which were set until now:

```
esxcli system module parameters list -m <module name>
```

Parameters which are not set by the user, remain on default value.



### 1.3.1.1 nmlx5\_core Parameters

**Table 1 - nmlx5\_core Module Parameters**

Name	Description	Values
DRSS	<p>Number of hardware queues for Default Queue (DEFQ) RSS.</p> <p><b>Note:</b> This parameter replaces the previously used "drss" parameter which is now obsolete.</p>	<ul style="list-style-type: none"> <li>• 2-16</li> <li>• 0 - disabled</li> </ul> <p>When this value is != 0, DEFQ RSS is enabled with 1 RSS Uplink queue that manages the 'drss' hardware queues.</p> <p><b>Notes:</b></p> <ul style="list-style-type: none"> <li>• The value must be a power of 2.</li> <li>• The value must not exceed num. of CPU cores.</li> <li>• Setting the DRSS value to 16, sets the Steering Mode to device RSS</li> </ul>
enable_nmlx_debug	Enables debug prints for the core module.	<ul style="list-style-type: none"> <li>• 1 - enabled</li> <li>• 0 - disabled (Default)</li> </ul>
max_vfs	Number of PCI VFs to initialize. The number of max_vf is set per port. For example, in case of a dual-port adapter when the max_vf is set to 8, the total number of opened VFs would be 16.	<ul style="list-style-type: none"> <li>• 0 - disabled (Default)</li> </ul> <p>N number of VF to allocate over each port</p>
mst_recovery	Enables recovery mode (only NMST module is loaded).	<ul style="list-style-type: none"> <li>• 1 - enabled</li> <li>• 0 - disabled (Default)</li> </ul>
pfcrx	Priority based Flow Control policy on RX.	<ul style="list-style-type: none"> <li>• 0-255</li> <li>• 0 - default</li> </ul> <p>It is an 8 bits bit mask, where each bit indicates a priority [0-7].</p> <p>Bit values:</p> <ul style="list-style-type: none"> <li>• 1 - respect incoming PFC pause frames for the specified priority.</li> <li>• 0 - ignore incoming pause frames on the specified priority.</li> </ul> <p><b>Note:</b> The pfcrx and pfcpx values must be identical.</p>

**Table 1 - nmlx5\_core Module Parameters**

Name	Description	Values
pfctx	Priority based Flow Control policy on TX.	<ul style="list-style-type: none"> <li>• 0-255</li> <li>• 0 - default</li> </ul> It is an 8 bits bit mask, where each bit indicates a priority [0-7]. Bit values: <ul style="list-style-type: none"> <li>• 1 - generate pause frames according to the RX buffer threshold on the specified priority.</li> <li>• 0 - never generate pause frames on the specified priority.</li> </ul> <b>Note:</b> The pfcrx and pfctx values must be identical.
RSS	Number of hardware queues for NetQ RSS.  <b>Note:</b> This parameter replaces the previously used "rss" parameter which is now obsolete.	<ul style="list-style-type: none"> <li>• 2-8</li> <li>• 0 - disabled</li> </ul> When this value is != 0, NetQ RSS is enabled with 1 RSS uplink queue that manages the 'rss' hardware queues. Notes: <ul style="list-style-type: none"> <li>• The value must be a power of 2</li> <li>• The value must not exceed number of CPU cores</li> </ul>
supported_num_ports	Sets the maximum supported ports.	2-8 Default 4

## 2 Installation

This chapter describes how to install and test the Mellanox ConnectX-4/ConnectX-5 NATIVE ESXi package on a single host machine with Mellanox Ethernet adapter hardware installed.

### 2.1 Hardware and Software Requirements

*Table 2 - Software and Hardware Requirements*

Requirements	Description
Platforms	A server platform with an adapter card based on one of the following Mellanox Technologies' HCA devices: <ul style="list-style-type: none"> <li>• ConnectX®-4 (EN) (firmware: fw-ConnectX4)</li> <li>• ConnectX®-4 Lx (EN) (firmware: fw-ConnectX4-Lx)</li> <li>• ConnectX®-5 (VPI) (firmware: fw-ConnectX5)</li> <li>• ConnectX®-5 Ex (VPI) (firmware: fw-ConnectX5)</li> </ul>
Device ID	For the latest list of device IDs, please visit Mellanox website.
Operating System	ESXi 5.5 U1: 4.5.8.8 ESXi 6.0/6.0 U1/6.0 U2: 4.5.8.8
Installer Privileges	The installation requires administrator privileges on the target machine.

### 2.2 Installing Mellanox NATIVE ESXi Driver for VMware vSphere



Please uninstall any previous Mellanox driver packages prior to installing the new version. See [Section 2.3, “Removing the Previous Mellanox Driver”](#), on page 12 for further information.

➤ **To install the driver:**

1. Log into the ESXi server with root permissions.
2. Install the driver.

```
#> esxcli software vib install -d <path>/<bundle_file>
```

Example:

```
#> esxcli software vib install -d /tmp/MLNX-NATIVE-ESX-ConnectX-4-5_4.5.8.8-10EM-600.0.0.2768847.zip
```

3. Reboot the machine.
4. Verify the driver was installed successfully.

```
# esxcli software vib list | grep mlx
ESX 5.5:
nmlx5-core          4.5.8.8-10EM.550.0.0.1391871    MEL      PartnerSupported  2017-01-31
ESX 6.0:
nmlx5-core          4.5.8.8-10EM.600.0.0.2768847    MEL      PartnerSupported  2017-01-31
```



After the installation process, all kernel modules are loaded automatically upon boot.

## 2.3 Removing the Previous Mellanox Driver



Please unload the driver before removing it.

### ➤ *To remove all the drivers:*

1. Log into the ESXi server with root permissions.
2. List all the existing NATIVE ESXi driver modules. (see [Step 4 in Section 2.2, on page 11](#))
3. Remove each module.

```
#> esxcli software vib remove -n nmlx5-core
```



To remove the modules, the command must be run in the same order as shown in the example above.

4. Reboot the server.

## 2.4 Firmware Programming

1. Download the VMware bootable binary images v4.6.0 from the [Mellanox Firmware Tools \(MFT\)](#) site.
  - **ESXi 5.5 File:** mft-4.6.0.48-10EM-550.0.0.1391871.x86\_64.vib  
**MD5SUM:** e9534bd77824682a94d845d30b7b6ee4
  - **ESXi 6.0 File:** mft-4.6.0.48-10EM-600.0.0.2768847.x86\_64.vib  
**MD5SUM:** e4c6092a0e61e24630127aa506975f48
2. Install the image according to the steps described in the [MFT User Manual](#).



The following procedure requires custom boot image downloading, mounting and booting from a USB device.

## 3 Features Overview and Configuration

### 3.1 Ethernet Network

#### 3.1.1 Wake-on-LAN (WoL)



Please note that Wake-on-LAN (WOL) is applicable only to adapter cards that support this feature.

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:

```
esxcli network nic set -n <nic name> -w g
```

or

```
set /net/pNics/<nic name>/wol g
```

- To disable WoL:

```
vsish -e set /net/pNics/<nic name>/wol d
```

- To verify configuration:

```
esxcli network nic get -n vmnic5
  Advertised Auto Negotiation: true
  Advertised Link Modes: 10000baseT/Full, 40000baseT/Full, 100000baseT/Full, 100baseT/
Full, 1000baseT/Full, 25000baseT/Full, 50000baseT/Full
  Auto Negotiation: false
  Cable Type: DA
  Current Message Level: -1
  Driver Info:
    Bus Info: 0000:82:00:1
    Driver: nmlx5_core
    Firmware Version: 12.18.0356
    Version: 4.5.16.xx
  Link Detected: true
  Link Status: Up
  Name: vmnic5
  PHYAddress: 0
  Pause Autonegotiate: false
  Pause RX: false
  Pause TX: false
  Supported Ports:
  Supports Auto Negotiation: true
  Supports Pause: false
  Supports Wakeon: false
  Transceiver:
  Wakeon: MagicPacket (tm)
```

### 3.1.2 Set Link Speed

The driver is set to auto-negotiate by default. However, the link speed can be forced to a specific link speed supported by ESXi using the following command:

```
esxcli network nic set -n <vmnic> -S <speed> -D <full, half>
```

Example:

```
esxcli network nic set -n vmnic4 -S 10000 -D full
```

where:

- <speed> in ESXi 6.0 can be 10/100/1000Mb/s.
- <vmnic> is the vmnic for the Mellanox card as provided by ESXi
- <full, half> The duplex to set this NIC to. Acceptable values are: [full, half]

The driver can be reset to auto-negotiate using the following command:

```
esxcli network nic set -n <vmnic> -a
```

Example:

```
esxcli network nic set -n vmnic4 -a
```

### 3.1.3 Priority Flow Control (PFC)

Priority Flow Control (PFC) IEEE 802.1Qbb applies pause functionality to specific classes of traffic on the Ethernet link. PFC can provide different levels of service to specific classes of Ethernet traffic (using IEEE 802.1p traffic classes).



When PFC is enabled, Global Pause will be operationally disabled, regardless of what is configured for the Global Pause Flow Control.

#### ➤ *To configure PFC:*

**Step 1.** Enable PFC for specific priorities.

```
esxcfg-module nmlx5_core -s "pfctx=0x08 pfcrx=0x08"
```

The parameters, "pfctx" (PFC TX) and "pfcrx" (PFC RX), are specified per host. If you have more than a single card on the server, all ports will be enabled with PFC (Global Pause will be disabled even if configured).

The value is a bitmap of 8 bits = 8 priorities. We recommend that you enable only lossless applications on a specific priority.

To run more than one flow type on the server, turn on only one priority (e.g. priority 3), which should be configured with the parameters "0x08" = 00001000b (binary). Only the 4th bit is on (starts with priority 0,1,2 and 3 -> 4th bit).

**Note:** The values of "pfctx" and "pfcrx" must be equal.

**Step 2.** Restart the driver.

```
/etc/init.d/sfcbdd-watchdog stop
esxcfg-module -u nmlx5_core
esxcfg-module nmlx5_core
/etc/init.d/sfcbdd-watchdog start
kill -POLL $(cat /var/run/vmware/vmkdevmgr.pid)
```

### 3.1.4 Receive Side Scaling (RSS)

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

#### 3.1.4.1 Default Queue Receive Side Scaling (DRSS)

Default Queue RSS (DRSS) allows the user to configure multiple hardware queues backing up the default RX queue. DRSS improves performance for large scale multicast traffic between hypervisors and Virtual Machines interfaces.

To configure DRSS, use the 'DRSS' module parameter which replaces the previously advertised 'device\_rss' module parameter ('device\_rss' is now obsolete). The 'drss' module parameter and 'device\_rss' are mutually exclusive

If the 'device\_rss' module parameter is enabled, the following functionality will be configured:

- The new Default Queue RSS mode will be triggered and all hardware RX rings will be utilized, similar to the previous 'device\_rss' functionality
- Module parameters 'DRSS' and 'RSS' will be ignored, thus the NetQ RSS, or the standard NetQ will be active

To query the 'DRSS' module parameter default, its minimal or maximal values, and restrictions, run a standard esxcli command.

For example:

```
#esxcli system module parameters list -m nmlx5_core
```

#### 3.1.4.2 NetQ RSS

NetQ RSS is a new module parameter for ConnectX-4 adapter cards providing identical functionality as the ConnectX-3 module parameter 'num\_rings\_per\_rss\_queue'. The new module parameter allows the user to configure multiple hardware queues backing up the single RX queue. NetQ RSS improves vMotion performance and multiple streams of IPv4/IPv6 TCP/UDP/IPSEC bandwidth over single interface between the Virtual Machines.

To configure NetQ RSS, use the 'RSS' module parameter. To query the 'RSS' module parameter default, its minimal or maximal values, and restrictions, run a standard esxcli command.

For example:

```
#esxcli system module parameters list -m nmlx5_core
```



Using NetQ RSS is preferred over the Default Queue RSS. Therefore, if both module parameters are set but the system lacks resources to support both, NetQ RSS will be used instead of DRSS.

### 3.1.4.3 Important Notes

If the 'DRSS' and 'RSS' module parameters set by the user cannot be enforced by the system due to lack of resources, the following actions are taken in a sequential order:

1. The system will attempt to provide the module parameters default values instead of the ones set by the user
2. The system will attempt to provide 'RSS' (NetQ RSS mode) default value. The Default Queue RSS will be disabled
3. The system will load with only standard NetQ queues
4. 'DRSS' and 'RSS' parameters are disabled by default, and the system loads with standard NetQ mode

## 3.2 Virtualization

### 3.2.1 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 in ConnectX-4/ConnectX-5 adapter cards up to 16 virtual instances called Virtual Functions (VFs). These virtual functions can then be provisioned separately. Each VF can be seen as an addition device connected to the Physical Function. It shares the same resources with 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 ESXi environment using Mellanox ConnectX® adapter cards family.

#### 3.2.1.1 System Requirements

To set up an SR-IOV environment, the following is required:

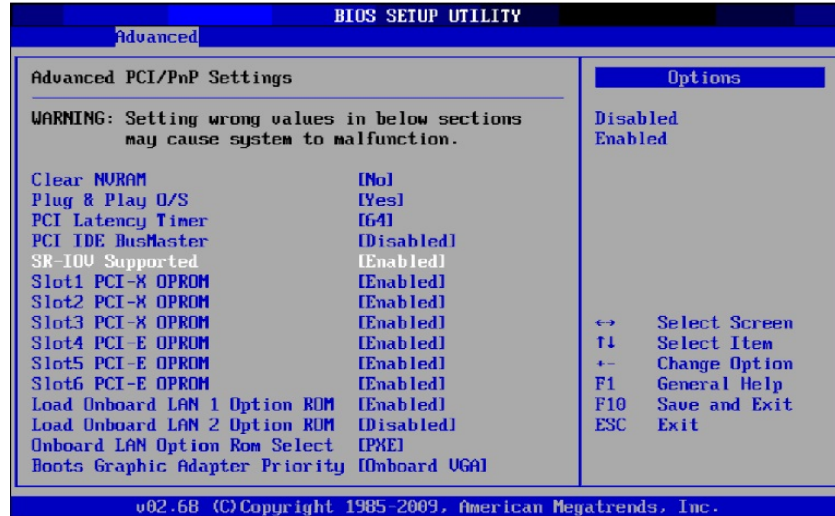
- nmlx5\_core Driver
- A server/blade with an SR-IOV-capable motherboard BIOS
- Mellanox ConnectX® Adapter Card family with SR-IOV capability
- Hypervisor that supports SR-IOV such as: ESXi 5.5 and 6.0



### 3.2.1.2 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:

**Step 1.** Enable "SR-IOV" in the system BIOS.



**Step 2.** Enable "Intel Virtualization Technology".



**Step 3.** Install ESXi 5.5 or ESXi 6.0 that support SR-IOV.

### 3.2.1.2.1 Configuring SR-IOV for ConnectX-4

- Step 1.** Install the MLNX-NATIVE-ESX-ConnectX-4 driver for ESXi that supports SR-IOV.
- Step 2.** Check if SR-IOV is enabled in the firmware.

```
/opt/mellanox/bin/mlxconfig -d /dev/mst/mt4115_pciconf0 q

Device #1:
-----

Device type:    ConnectX4
PCI device:    /dev/mst/mt4115_pciconf0
Configurations:    Current
  SRIOV_EN      1
  NUM_OF_VFS    8
  FPP_EN        1
```

If not, use `mlxconfig` to enable it.

```
mlxconfig -d /dev/mst/mt4115_pciconf0 set SRIOV_EN=1 NUM_OF_VFS=16
```

- Step 3.** Power cycle the server.
- Step 4.** Set the number of Virtual Functions you need to create for the PF using the `max_vfs` module parameter.

```
esxcli system module parameters set -m nmlx5_core -p "max_vfs=8"
```

**Note:** The number of `max_vf` is set per port. See [Table 1, “nmlx5\\_core Module Parameters,” on page 9](#) for further information.

- Step 5.** Reboot the server and verify the SR-IOV is supported once the server is up.

```
lspci | grep Mellanox
08:00.0 Network controller: Mellanox Technologies MT27700 Family [ConnectX-4]
08:00.1 Network controller: Mellanox Technologies MT27700 Family [ConnectX-4]
08:00.2 Network controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual
Function]
08:00.3 Network controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual
Function]
08:00.4 Network controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual
Function]
08:00.5 Network controller: Mellanox Technologies MT27700 Family [ConnectX-4 Virtual
Function]
```

### 3.2.1.3 Assigning a Virtual Function to a Virtual Machine in the vSphere Web Client

After you enable the Virtual Functions on the host, each of them becomes available as a PCI device.

➤ **To assign Virtual Function to a Virtual Machine in the vSphere Web Client:**

- Step 1.** Locate the Virtual Machine in the vSphere Web Client.
  - a. Select a data center, folder, cluster, resource pool, or host and click the Related Objects tab.
  - b. Click Virtual Machines and select the virtual machine from the list.
- Step 2.** Power off the Virtual Machine.

- Step 3.** On the **Manage** tab of the Virtual Machine, select **Settings > VM Hardware**.
- Step 4.** Click **Edit** and choose the **Virtual Hardware** tab.
- Step 5.** From the **New Device** drop-down menu, select **Network** and click **Add**.
- Step 6.** Expand the **New Network** section and connect the Virtual Machine to a port group.  
The virtual NIC does not use this port group for data traffic. The port group is used to extract the networking properties, for example VLAN tagging, to apply on the data traffic.
- Step 7.** From the **Adapter Type** drop-down menu, select **SR-IOV passthrough**.
- Step 8.** From the **Physical Function** drop-down menu, select the **Physical Adapter** to back the passthrough Virtual Machine adapter.
- Step 9.** **[Optional]** From the **MAC Address** drop-down menu, select **Manual** and type the static MAC address.
- Step 10.** Use the **Guest OS MTU Change** drop-down menu to allow changes in the MTU of packets from the guest operating system.  
**Note:** This step is applicable only if this feature is supported by the driver.
- Step 11.** Expand the **Memory** section, select **Reserve all guest memory (All locked)** and click **OK**.  
I/O memory management unit (IOMMU) must reach all Virtual Machine memory so that the passthrough device can access the memory by using direct memory access (DMA).
- Step 12.** Power on the Virtual Machine.

## 3.2.2 VXLAN Hardware Offload

VXLAN hardware offload enables the traditional offloads to be performed on the encapsulated traffic. With ConnectX® family adapter cards, data center operators can decouple the overlay network layer from the physical NIC performance, thus achieving native performance in the new network architecture.

### 3.2.2.1 Configuring VXLAN Hardware Offload

VXLAN hardware offload includes:

- TX: Calculates the Inner L3/L4 and the Outer L3 checksum
- RX:
  - Checks the Inner L3/L4 and the Outer L3 checksum
  - Maps the VXLAN traffic to an RX queue according to:
    - Inner destination MAC address
    - Outer destination MAC address
    - VXLAN ID

VXLAN hardware offload is enabled by default and its status cannot be changed.

VXLAN configuration is done in the ESXi environment via VMware NSX manager. For additional NSX information, please refer to VMware documentation, see:

<http://pubs.vmware.com/NSX-62/index.jsp#com.vmware.nsx.install.doc/GUID-D8578F6E-A40C-493A-9B43-877C2B75ED52.html>.



## 4 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.

### 4.1 General Related Issues

**Table 3 - General Related Issues**

Issue	Cause	Solution
The system panics when it is booted with a failed adapter installed.	Malfunction hardware component	<ol style="list-style-type: none"> <li>1. Remove the failed adapter.</li> <li>2. Reboot the system.</li> </ol>
Mellanox adapter is not identified as a PCI device.	PCI slot or adapter PCI connector dysfunctionality	<ol style="list-style-type: none"> <li>1. Run <code>lspci</code>.</li> <li>2. 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.</li> </ol>
Mellanox adapters are not installed in the system.	Misidentification of the Mellanox adapter installed	<p>Run the command below to identify the Mellanox adapter installed.</p> <pre>lspci   grep Mellanox'</pre>

### 4.2 Ethernet Related Issues

**Table 4 - Ethernet Related Issues**

Issue	Cause	Solution
No link.	Mis-configuration of the switch port or using a cable not supporting link rate.	<ul style="list-style-type: none"> <li>• Ensure the switch port is not down</li> <li>• Ensure the switch port rate is configured to the same rate as the adapter's port</li> </ul>
No link with break-out cable.	Misuse of the break-out cable or misconfiguration of the switch's split ports	<ul style="list-style-type: none"> <li>• Use supported ports on the switch with proper configuration. For further information, please refer to the MLNX_OS User Manual.</li> <li>• Make sure the QSFP break-out cable side is connected to the SwitchX.</li> </ul>
Physical link fails to negotiate to maximum supported rate.	The adapter is running an outdated firmware.	Install the latest firmware on the adapter.

**Table 4 - Ethernet Related Issues**

Issue	Cause	Solution
Physical link fails to come up.	The cable is not connected to the port or the port on the other end of the cable is disabled.	Ensure that the cable is connected on both ends or use a known working cable

## 4.3 Installation Related Issues

**Table 5 - Installation Related Issues**

Issue	Cause	Solution
Driver installation fails.	<p>The install script may fail for the following reasons:</p> <ul style="list-style-type: none"> <li>Failed to uninstall the previous installation due to dependencies being used</li> <li>The operating system is not supported</li> </ul>	<ul style="list-style-type: none"> <li>Uninstall the previous driver before installing the new one</li> <li>Use a supported operating system and kernel</li> </ul>