



Connect. Accelerate. Outperform.™

# **Mellanox HPC-X™ Software Toolkit User Manual**

Rev 1.6

**Powered by CORE-Direct® Technology**

[www.mellanox.com](http://www.mellanox.com)

## NOTE:

THIS HARDWARE, SOFTWARE OR TEST SUITE PRODUCT (“PRODUCT(S)”) AND ITS RELATED DOCUMENTATION ARE PROVIDED BY MELLANOX TECHNOLOGIES “AS-IS” WITH ALL FAULTS OF ANY KIND AND SOLELY FOR THE PURPOSE OF AIDING THE CUSTOMER IN TESTING APPLICATIONS THAT USE THE PRODUCTS IN DESIGNATED SOLUTIONS. THE CUSTOMER'S MANUFACTURING TEST ENVIRONMENT HAS NOT MET THE STANDARDS SET BY MELLANOX TECHNOLOGIES TO FULLY QUALIFY THE 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

Mellanox Technologies, Ltd.  
 Beit Mellanox  
 PO Box 586 Yokneam 20692  
 Israel  
[www.mellanox.com](http://www.mellanox.com)  
 Tel: +972 (0)74 723 7200  
 Fax: +972 (0)4 959 3245

© Copyright 2014. Mellanox Technologies. All Rights Reserved.

Mellanox®, Mellanox logo, BridgeX®, ConnectX®, Connect-IB®, CORE-Direct®, InfiniBridge®, InfiniHost®, InfiniScale®, MetroX®, MLNX-OS®, PhyX®, ScalableHPC®, SwitchX®, UFM®, Virtual Protocol Interconnect® and Voltaire® are registered trademarks of Mellanox Technologies, Ltd.

ExtendX™, FabricIT™, Mellanox Open Ethernet™, Mellanox Virtual Modular Switch™, MetroDX™, TestX™, Unbreakable-Link™ are trademarks of Mellanox Technologies, Ltd.

All other trademarks are property of their respective owners.

# Table of Contents

<b>Table of Contents</b> .....	<b>3</b>
<b>List of Tables</b> .....	<b>5</b>
<b>List of Figures</b> .....	<b>6</b>
<b>Document Revision History</b> .....	<b>7</b>
<b>About This Manual</b> .....	<b>9</b>
Scope .....	9
Intended Audience .....	9
Syntax Conventions .....	9
<b>Chapter 1 HPC-X™ Software Toolkit Overview</b> .....	<b>10</b>
1.1 HPC-X Package Contents .....	10
1.2 HPC-X™ Requirements .....	10
<b>Chapter 2 Installing and Loading HPC-X™</b> .....	<b>11</b>
2.1 Installing HPC-X .....	11
2.2 Loading HPC-X Environment from bash .....	11
2.3 Loading HPC-X Environment from Modules .....	11
<b>Chapter 3 Running, Configuring and Rebuilding HPC-X™</b> .....	<b>12</b>
3.1 Starting FCA Manager from HPC-X .....	12
3.2 Profiling IB verbs API .....	12
3.3 Profiling MPI API .....	12
3.4 Rebuilding Open MPI from HPC-X™ Sources .....	13
3.5 Generating MXM Statistics for Open MPI/OpenSHMEM .....	14
3.6 Generating MXM Environment Parameters .....	15
3.7 Loading KNEM Module .....	26
3.8 Running MPI with FCA v2.5 .....	26
3.9 Running OpenSHMEM with FCA v2.5 .....	26
3.10 Running MPI with FCA v3.5 (hcoll) .....	27
3.11 IB-Router .....	27
<b>Chapter 4 Mellanox Fabric Collective Accelerator (FCA)</b> .....	<b>28</b>
4.1 Overview .....	28
4.2 FCA Installation Package Content .....	30
4.3 Differences Between FCA v2.5 and FCA v3.x .....	31
4.4 Configuring FCA .....	31
4.4.1 Compiling Open MPI with FCA 3.x .....	31
4.4.2 Enabling FCA in Open MPI .....	31
4.4.3 Tuning FCA 3.X Setting .....	32
4.4.4 Selecting Ports and Devices .....	32
4.4.5 Enabling Multicast Accelerated Collectives .....	32
4.4.5.1 Configuring IPoIB .....	32
4.4.6 Enabling HCOLL Topology Awareness .....	33
4.4.7 Enabling SHArP Accelerated Collectives .....	34

<b>Chapter 5 MellanoX Messaging Library</b>	<b>35</b>
5.1 Overview	35
5.2 Compiling Open MPI with MXM	35
5.3 Enabling MXM in Open MPI	36
5.4 Tuning MXM Settings	37
5.5 Configuring Multi-Rail Support	37
5.6 Configuring MXM over the Ethernet Fabric	38
5.7 Running MXM with RoCE	38
5.7.1 Running MXM with RoCE v1	38
5.7.2 Running MXM with RoCE v2	39
5.7.2.1 Using RoCE v2 in Open MPI	39
5.8 Configuring MXM over Different Transports	39
5.9 Configuring Service Level Support	39
5.10 Running Open MPI with pml “yalla”	40
5.11 Adaptive Routing for UD Transport	41
5.12 Support for a Non-Base LID	41
5.13 MXM Performance Tuning	41
5.14 MXM Utilities	43
5.14.1 mxm_dump_config	43
5.14.2 mxm_perftest	43
<b>Chapter 6 PGAS Shared Memory Access Overview</b>	<b>45</b>
6.1 HPC-X Open MPI/OpenSHMEM	45
6.2 Running HPC-X OpenSHMEM	46
6.2.1 Running HPC-X OpenSHMEM with MXM	46
6.2.1.1 Enabling MXM for HPC-X OpenSHMEM Jobs	46
6.2.1.2 Working with Multiple HCAs	46
6.2.2 Running HPC-X™ OpenSHMEM with FCA v2.x	47
6.2.3 Developing Application using HPC-X OpenSHMEM together with MPI	47
6.2.4 HPC-X™ OpenSHMEM Tunable Parameters	48
6.2.4.1 OpenSHMEM MCA Parameters for Symmetric Heap Allocation	49
6.2.4.2 Parameters Used to Force Connection Creation	49
6.3 Tuning OpenSHMEM Atomics Performance	49
6.4 Tuning MTU Size to the Recommended Value	50
6.4.1 HPC Applications on Intel Sandy Bridge Machines	51
<b>Chapter 7 Unified Parallel C Overview</b>	<b>52</b>
7.1 Compiling and Running the UPC Application	52
7.1.1 Compiling the UPC Application	52
7.1.2 Running the UPC Application	52
7.1.2.1 Basic upcrun Options	53
7.1.2.2 Environment Variables	53
7.2 FCA Runtime Parameters	53
7.2.1 Enabling FCA Operations through Environment Variables in HPC-X UPC	54
7.2.2 Controlling FCA Offload in HPC-X UPC using Environment Variables	54
7.3 Various Executable Examples	54

## List of Tables

Table 1:	Syntax Conventions . . . . .	9
Table 2:	Generating MXM Environment Parameters . . . . .	15
Table 3:	MLNX_OFED and MXM Versions . . . . .	36
Table 4:	Runtime Parameters . . . . .	53

## List of Figures

Figure 1: FCA Architecture .....	29
Figure 2: FCA Components .....	30

## Document Revision History

Revision	Date	Description
1.6	June 2016	<ul style="list-style-type: none"> <li>• Added the following sections:               <ul style="list-style-type: none"> <li>• <a href="#">Section 3.11, “IB-Router”, on page 27</a></li> <li>• <a href="#">Section 4.4.6, “Enabling HCOLL Topology Awareness”, on page 33</a></li> <li>• <a href="#">Section 4.4.7, “Enabling SHArP Accelerated Collectives”, on page 34</a></li> </ul> </li> <li>• Updated the following sections:               <ul style="list-style-type: none"> <li>• <a href="#">Section 1.1, “HPC-X Package Contents”, on page 10</a> <ul style="list-style-type: none"> <li>• Updated the FCA version to 3.5</li> <li>• Updated the MXM version to 3.5</li> </ul> </li> <li>• <a href="#">Section 3.4, “Rebuilding Open MPI from HPC-X™ Sources”, on page 13</a> <ul style="list-style-type: none"> <li>• Updated the MXM version to 3.5</li> </ul> </li> </ul> </li> </ul>
1.5	January 2016	<ul style="list-style-type: none"> <li>• Added the following sections:               <ul style="list-style-type: none"> <li>• <a href="#">Section 5.7, “Running MXM with RoCE”, on page 38</a></li> <li>• <a href="#">Section 5.7.1, “Running MXM with RoCE v1”, on page 38</a></li> <li>• <a href="#">Section 5.7.2, “Running MXM with RoCE v2”, on page 39</a></li> <li>• <a href="#">Section 5.7.2.1, “Using RoCE v2 in Open MPI”, on page 39</a></li> </ul> </li> <li>• Updated the following sections:               <ul style="list-style-type: none"> <li>• <a href="#">Section 1.1, “HPC-X Package Contents”, on page 10</a> <ul style="list-style-type: none"> <li>• Updated the Open MPI version to 1.10</li> <li>• Updated the MXM version to 3.4</li> </ul> </li> <li>• <a href="#">Section 3.6, “Generating MXM Environment Parameters”, on page 15</a> <ul style="list-style-type: none"> <li>• Added the following parameter: <code>MXM_IB_GID_INDEX</code></li> </ul> </li> <li>• <a href="#">Section 5.8, “Configuring MXM over Different Transports”, on page 39</a> <ul style="list-style-type: none"> <li>• Updated the MXM version to 3.4</li> </ul> </li> </ul> </li> <li>• Removed section Running MPI with MXM</li> </ul>
1.4	September 2015	<ul style="list-style-type: none"> <li>• Updated the following sections:               <ul style="list-style-type: none"> <li>• <a href="#">Section 4.1, “Overview”, on page 28</a></li> <li>• <a href="#">Section 4.4.2, “Enabling FCA in Open MPI”, on page 31</a></li> <li>• <a href="#">Section 5.8, “Configuring MXM over Different Transports”, on page 39</a></li> </ul> </li> </ul>

Revision	Date	Description
1.3	June 2015	<ul style="list-style-type: none"> <li>• Added the following section:               <ul style="list-style-type: none"> <li>• <a href="#">Section 4.4.5, “Enabling Multicast Accelerated Collectives”</a>, on page 32</li> <li>• <a href="#">Section 4.4.5.1, “Configuring IPoIB”</a>, on page 32</li> </ul> </li> <li>• Updated the following sections:               <ul style="list-style-type: none"> <li>• <a href="#">Section 3.10, “Running MPI with FCA v3.5 (hcoll)”</a>, on page 27</li> <li>• <a href="#">Section 4.4.2, “Enabling FCA in Open MPI”</a>, on page 31</li> <li>• <a href="#">Section 5.10, “Running Open MPI with pml “yalla””</a>, on page 40</li> <li>• <a href="#">Section 5.12, “Support for a Non-Base LID”</a>, on page 41</li> <li>• <a href="#">Section 6.2.1.1, “Enabling MXM for HPC-X OpenSH-MEM Jobs”</a>, on page 46</li> <li>• <a href="#">Section 7.3, “Various Executable Examples”</a>, on page 54</li> </ul> </li> <li>• Removed the following sections:               <ul style="list-style-type: none"> <li>• Runtime Configuration of FCA and its subsections</li> <li>• FCA 3.1 Integration</li> <li>• Configuring Hugepages</li> </ul> </li> </ul>
1.2	August 2014	Initial release



# About This Manual

## Scope

This document describes Mellanox HPC-X™ Software Toolkit acceleration packages. It includes information on installation, configuration and rebuilding of HPC-X packages.

## Intended Audience

This manual is intended for system administrators responsible for the installation, configuration, management and maintenance of the software and hardware.

It is also for users who would like to use the latest Mellanox software accelerators to achieve the best possible application performance.

## Syntax Conventions

**Table 1 - Syntax Conventions**

Prompt	Shell
machine-name%	C shell on UNIX, Linux, or AIX
machine-name#	C shell superuser on UNIX, Linux, or AIX
\$	Bourne shell and Korn shell on UNIX, Linux, or AIX
#	Bourne shell and Korn shell superuser on UNIX, Linux, or AIX
C:\>	Windows command line

# 1 HPC-X™ Software Toolkit Overview

Mellanox HPC-X™ is a comprehensive software package that includes MPI, SHMEM and UPC communications libraries. HPC-X also includes various acceleration packages to improve both the performance and scalability of applications running on top of these libraries, including MXM (Mellanox Messaging) which accelerates the underlying send/receive (or put/get) messages, and FCA (Fabric Collectives Accelerations) which accelerates the underlying collective operations used by the MPI/PGAS languages. This full-featured, tested and packaged version of HPC software enables MPI, SHMEM and PGAS programming languages to scale to extremely large clusters by improving on memory and latency related efficiencies, and to assure that the communication libraries are fully optimized of the Mellanox interconnect solutions.

Mellanox HPC-X™ allow OEM's and System Integrators to meet the needs of their end-users by deploying the latest available software that takes advantage of the features and capabilities available in the most recent hardware and firmware changes.

## 1.1 HPC-X Package Contents

Mellanox HPC-X package contains the following pre-compiled HPC packages:

Components	Description
MPI	<ul style="list-style-type: none"> <li>Open MPI and OpenSHMEM v1.10 (MPI-3 complaint, OpenSHMEM v1.0 compliant). Open MPI and OpenSHMEM are available at: <a href="http://www.open-mpi.org/software/ompi/v1.10/">http://www.open-mpi.org/software/ompi/v1.10/</a></li> <li>MPI profiler (IPM - open source tool from <a href="http://ipm-hpc.org/">http://ipm-hpc.org/</a>)</li> <li>MPI tests (OSU, IMB, random ring, etc.)</li> </ul>
HPC Acceleration Package	<ul style="list-style-type: none"> <li>MXM 3.5 (default)</li> <li>FCA v2.5 (default)</li> <li>FCA v3.5 (code name: "hcoll")</li> <li>knem (High-Performance Intra-Node MPI Communication module from: <a href="http://runtime.bordeaux.inria.fr/knem/">http://runtime.bordeaux.inria.fr/knem/</a>)</li> </ul>
Extra packages	<ul style="list-style-type: none"> <li>Berkeley UPC v2.22.0</li> <li>libibprof (IB verbs profiler)</li> </ul>

## 1.2 HPC-X™ Requirements

The platform and requirements for HPC-X are detailed in the following table:

Platform	Drivers and HCAs
OFED / MLNX-OFED	<ul style="list-style-type: none"> <li>OFED 1.5.3 and later</li> <li>MLNX_OFED 1.5.3-x.x.x, 2.0-x.x.x and later</li> </ul>
HCAs	<ul style="list-style-type: none"> <li>Mellanox ConnectX®-2 / ConnectX®-3 / ConnectX®-3 Pro</li> <li>Mellanox Connect-IB®</li> </ul>

## 2 Installing and Loading HPC-X™

### 2.1 Installing HPC-X

➤ *To install HPC-X:*

**Step 1.** Extract hpcx.tar into current working directory.

```
$ tar zxvf hpcx.tar
```

**Step 2.** Update shell variable of the location of hpc-x installation.

```
$ cd hpcx
$ export HPCX_HOME=$PWD
```

### 2.2 Loading HPC-X Environment from bash

HPC-X includes Open MPI v1.10. Each Open MPI version has its own module file which can be used to load desired version.

The symbolic links `hpcx-init.sh` and `modulefiles/hpcx` point to the default version (Open MPI v1.10).

➤ *To load Open MPI/OpenSHMEM v1.10 based package:*

```
$ source $HPCX_HOME/hpcx-init.sh
$ hpcx_load
$ env | grep HPCX
$ mpirun -np 2 $HPCX_MPI_TESTS_DIR/examples/hello_usempi
$ oshrun -np 2 $HPCX_MPI_TESTS_DIR/examples/hello_oshmem
$ hpcx_unload
```

### 2.3 Loading HPC-X Environment from Modules

➤ *To load Open MPI/OpenSHMEM v1.10 based package:*

```
$ module use $HPCX_HOME/modulefiles
$ module load hpcx
$ mpirun -np 2 $HPCX_MPI_TESTS_DIR/examples/hello_c
$ oshrun -np 2 $HPCX_MPI_TESTS_DIR/examples/hello_oshmem
$ module unload hpcx
```

The instrumentation can be applied on a per-collective basis, and is controlled by the following environment variables:

```
$ export IPM_ADD_BARRIER_TO_REDUCE=1
$ export IPM_ADD_BARRIER_TO_ALLREDUCE=1
$ export IPM_ADD_BARRIER_TO_GATHER=1
$ export IPM_ADD_BARRIER_TO_ALL_GATHER=1
$ export IPM_ADD_BARRIER_TO_ALLTOALL=1
$ export IPM_ADD_BARRIER_TO_ALLTOALLV=1
$ export IPM_ADD_BARRIER_TO_BROADCAST=1
$ export IPM_ADD_BARRIER_TO_SCATTER=1
$ export IPM_ADD_BARRIER_TO_SCATTERV=1
$ export IPM_ADD_BARRIER_TO_GATHERV=1
$ export IPM_ADD_BARRIER_TO_ALLGATHERV=1
$ export IPM_ADD_BARRIER_TO_REDUCE_SCATTER=1
```

By default, all values are set to '0'.

## 3 Running, Configuring and Rebuilding HPC-X™

The sources for BUPC, SHMEM and OMPI can be found at `$HPCX_HOME/sources/`. Please refer to `$HPCX_HOME/sources/` for more information on building details.

### 3.1 Starting FCA Manager from HPC-X

Prior to using FCA, the following command should be executed as root once on all cluster nodes in order to mimic post-install procedure.

```
# $HPCX_HOME/fca/scripts/udev-update.sh
```

The FCA manager should be run on only one machine, and not on one of the compute nodes.

```
$ $HPCX_HOME/fca/scripts/fca_managerd start
```

FCA 2.5 is the default FCA version embedded in the HPC-X package. To install the FCA Manager please refer to FCA 2.5 User Manual section Installing the FCA Manager on a Dedicated Node

### 3.2 Profiling IB verbs API

➤ *To profile IB verbs API:*

```
$ export IBPROF_DUMP_FILE=ibprof_%J_%H_%T.txt
$ export LD_PRELOAD=$HPCX_IBPROF_DIR/lib/libibprof.so:$HPCX_MXM_DIR/lib/lib-
mxm.so:$HPCX_HCOLL_DIR/lib/libhcoll.so
$ mpirun -x LD_PRELOAD <...>
```

For further details on profiling IB verbs API, please refer to libibprof README file.

README file location is:

```
$HPCX_HOME/libibprof/README
```

### 3.3 Profiling MPI API

➤ *To profile MPI API:*

```
$ export IPM_KEYFILE=$HPCX_IPM_DIR/etc/ipm_key_mpi
$ export IPM_LOG=FULL
$ export LD_PRELOAD=$HPCX_IPM_DIR/lib/libipm.so
$ mpirun -x LD_PRELOAD <...>
$ $HPCX_IPM_DIR/bin/ipm_parse -html outfile.xml
```

For further details on profiling MPI API, please refer to:

<http://ipm-hpc.org/>

The Mellanox-supplied version of IMP contains an additional feature (Barrier before Collective), not found in the standard package, that allows end users to easily determine the extent of application imbalance in applications which use collectives. This feature instruments each collective so that it calls `MPI_Barrier()` before calling the collective operation itself. Time spent in this `MPI_Barrier()` is not counted as communication time, so by running an application with and without the Barrier before Collective feature, the extent to which application imbalance is a factor in performance, can be assessed.

The instrumentation can be applied on a per-collective basis, and is controlled by the following environment variables:

```
$ export IPM_ADD_BARRIER_TO_REDUCE=1
$ export IPM_ADD_BARRIER_TO_ALLREDUCE=1
$ export IPM_ADD_BARRIER_TO_GATHER=1
$ export IPM_ADD_BARRIER_TO_ALL_GATHER=1
$ export IPM_ADD_BARRIER_TO_ALLTOALL=1
$ export IPM_ADD_BARRIER_TO_ALLTOALLV=1
$ export IPM_ADD_BARRIER_TO_BROADCAST=1
$ export IPM_ADD_BARRIER_TO_SCATTER=1
$ export IPM_ADD_BARRIER_TO_SCATTERV=1
$ export IPM_ADD_BARRIER_TO_GATHERV=1
$ export IPM_ADD_BARRIER_TO_ALLGATHERV=1
$ export IPM_ADD_BARRIER_TO_REDUCE_SCATTER=1
```

By default, all values are set to '0'.

### 3.4 Rebuilding Open MPI from HPC-X™ Sources

HPC-X package contains Open MPI sources which can be found at `$HPCX_HOME/sources/` folder.

➤ **To build Open MPI from sources:**

```
$ HPCX_HOME=/path/to/extracted/hpcx
$ ./configure --prefix=${HPCX_HOME}/hpcx-ompi --with-knem=${HPCX_HOME}/knem \
  --with-fca=${HPCX_HOME}/fca --with-mxm=${HPCX_HOME}/mxm \
  --with-hcoll=${HPCX_HOME}/hcoll \
  --with-platform=contrib/platform/mellanox/optimized \
  --with-slurm --with-pmi
$ make -j9 all && make -j9 install
```

Open MPI and OpenSHMEM are pre-compiled with MXM v3.5 and use it by default.

### 3.5 Generating MXM Statistics for Open MPI/OpenSHMEM

In order to generate statistics, the statistics destination and trigger should be set.

- Destination is set by `MXM_STATS_DEST` environment variable whose values can be one of the following:

Value	Description
empty string	statistics are not reported
stdout	Print to standard output
stderr	Print to standard error
file:<filename>	Write to a file. Following substitutions are made: %h: host, %p: pid, %c: cpu, %t: time, %e: exe
file:<filename>:bin	Same as previous, but a binary format is used when saving. The file will be smaller, but not human-readable. The <code>mxm_stats_parser</code> tool can be used to parse binary statistics files

#### Examples:

```
$ export MXM_STATS_DEST="file:mxm_%h_%e_%p.stats"
$ export MXM_STATS_DEST="file:mxm_%h_%c.stats:bin"
$ export MXM_STATS_DEST="stdout"
```

- Trigger is set by `MXM_STATS_TRIGGER` environment variables. It can be one of the following:

Environment Variable	Description
exit	Dump statistics just before exiting the program
timer:<interval>	Dump statistics periodically, interval is given in seconds

#### Example:

```
$ export MXM_STATS_TRIGGER=exit
$ export MXM_STATS_TRIGGER=timer:3.5
```



The statistics feature is only enabled for the 'debug' version of MXM which is included in HPC-X. To use the statistics, run the command below from the command line:

```
$ mpirun -x LD_PRELOAD=$HPCX_DIR/mxm/debug/lib/libmxm.so ...
```

## 3.6 Generating MXM Environment Parameters

**Table 2 - Generating MXM Environment Parameters**

Variable	Valid Values and Default Values	Description
MXM_LOG_LEVEL	<ul style="list-style-type: none"> <li>FATAL</li> <li>ERROR</li> <li>INFO</li> <li>DEBUG</li> <li>TRACE</li> <li>REQ</li> <li>DATA</li> <li>ASYNC</li> <li>FUNC</li> <li>POLL</li> <li>WARN (default)</li> </ul>	MXM logging level. Messages with a level higher or equal to the selected will be printed.
MXM_LOG_FILE	String	<p>If not empty, MXM will print log messages to the specified file instead of std-out.</p> <p>The following substitutions are performed on this string:</p> <ul style="list-style-type: none"> <li>%p - Replaced with process ID</li> <li>%h - Replaced with host name</li> </ul> <p>Value: String.</p>
MXM_LOG_BUFFER	1024	Buffer size for a single log message. Value: memory units: <number>[b kb mb gb]
MXM_LOG_DATA_SIZE	0	How much of the packet payload to print, at most, in data mode. Value: unsigned long.
MXM_HANDLE_ERRORS	<ul style="list-style-type: none"> <li>None: No error handling</li> <li>Freeze: Freeze and wait for a debugger</li> <li>Debug: attach debugger</li> <li>bt: print back-trace</li> </ul>	Error handling mode.
MXM_ERROR_SIGNALS	<ul style="list-style-type: none"> <li>ILL</li> <li>SEGV</li> <li>BUS</li> <li>FPE</li> <li>PIPE</li> </ul>	Signals which are considered an error indication and trigger error handling. Value: comma-separated list of: system signal (number or SIGxxx)
MXM_GDB_COMMAND	gdb	If non-empty, attaches a gdb to the process in case of error, using the provided command. Value: string

Variable	Valid Values and Default Values	Description
MXM_DEBUG_SIGNO	HUP	Signal number which causes MXM to enter debug mode. Set to 0 to disable. Value: system signal (number or SIGxxx)
MXM_ASYNC_INTERVAL	50000.00us	Interval of asynchronous progress. Lower values may make the network more responsive, at the cost of higher CPU load. Value: time value: <number>[s us ms ns]
MXM_ASYNC_SIGNO	ALRM	Signal number used for async signaling. Value: system signal (number or SIGxxx)
MXM_STATS_DEST	<ul style="list-style-type: none"> <li>• udp:&lt;host&gt;[:&lt;port&gt;] - send over UDP to the given host:port.</li> <li>• stdout: print to standard output.</li> <li>• stderr: print to standard error.</li> <li>• file:&lt;file-name&gt;[:bin] - save to a file (%h: host, %p: pid, %c: cpu, %t: time, %e: exe)</li> </ul>	Destination to send statistics to. If the value is empty, statistics are not reported.
MXM_STATS_TRIGGER	<ul style="list-style-type: none"> <li>• timer:&lt;interval&gt;: dump in specified intervals.</li> <li>• exit: dump just before program exits (default)</li> </ul>	Trigger to dump statistics Value: string
MXM_MEMTRACK_DEST	<ul style="list-style-type: none"> <li>• file:&lt;filename&gt;: save to a file (%h: host, %p: pid, %c: cpu, %t: time, %e: exe)</li> <li>• stdout: print to standard output</li> <li>• stderr: print to standard error</li> </ul>	Memory tracking report output destination. If the value is empty, results are not reported.



Variable	Valid Values and Default Values	Description
MXM_INSTRUMENT	<ul style="list-style-type: none"> <li>• %h: host</li> <li>• %p: pid</li> <li>• %c: cpu</li> <li>• %t: time</li> <li>• %e: exe.</li> </ul>	File name to dump instrumentation records to. Value: string
MXM_INSTRUMENT_SIZE	1048576	Maximal size of instrumentation data. New records will replace old records. Value: memory units: <number>[b kb mb gb]
MXM_PERF_STALL_LOOPS	0	Number of performance stall loops to be performed. Can be used to normalize profile measurements to packet rate Value: unsigned long
MXM_ASYNC_MODE	<ul style="list-style-type: none"> <li>• Signal</li> <li>• none</li> <li>• Thread (default)</li> </ul>	Asynchronous progress method Value: [none signal thread]
MXM_MEM_ALLOC	<ul style="list-style-type: none"> <li>• cpages: Contiguous pages, provided by Mellanox-OFED.</li> <li>• hugetlb - Use System V shared memory API for getting pre-allocated huge pages.</li> <li>• mmap: Use private anonymous mmap() to get free pages.</li> <li>• libc: Use libc's memory allocation primitives.</li> <li>• sysv: Use system V's memory allocation.</li> </ul>	Memory allocators priority. Value: comma-separated list of: [libc hugetlb cpages mmap sysv]
MXM_MEM_ON_DEMAND_MAP	<ul style="list-style-type: none"> <li>• n: disable</li> <li>• y: enable</li> </ul>	Enable on-demand memory mapping. USE WITH CARE! It requires calling mxm_mem_unmap() when any buffer used for communication is unmapped, otherwise data corruption could occur. Value: <y n>

Variable	Valid Values and Default Values	Description
MXM_INIT_HOOK_SCRIPT	-	Path to the script to be executed at the very beginning of MXM initialization Value: string
MXM_SINGLE_THREAD	<ul style="list-style-type: none"> <li>y - single thread</li> <li>n - not single thread</li> </ul>	Mode of the thread usage. Value: <y n>
MXM_SHM_KCOPY_MODE	<ul style="list-style-type: none"> <li>off: Don't use any kernel copy mode.</li> <li>knem: Try to use knem. If it fails, default to 'off'.</li> <li>autodetect: If knem is available, first try to use knem. If it fails, default to 'off' (default)</li> </ul>	Modes for using to kernel copy for large messages.
MXM_IB_PORTS	*:*	Specifies which Infiniband ports to use. Value: comma-separated list of: IB port: <device>:<port_num>
MXM_EP_NAME	%h:%p	Endpoint options. Endpoint name used in log messages. Value: string
MXM_TLS	<ul style="list-style-type: none"> <li>self</li> <li>shm</li> <li>ud</li> </ul>	Comma-separated list of transports to use. The order is not significant. Value: comma-separated list of: [self shm rc dc ud oob]
MXM_ZCOPY_THRESH	2040	Threshold for using zero copy. Value: memory units: <number>[b kb mb gb]
MXM_IB_CQ_MODERATION	64	Number of send WREs for which a CQE is generated. Value: unsigned
MXM_IB_CQ_WATERMARK	127	Consider ep congested if poll cq returns more than n wqes. Value: unsigned
MXM_IB_DRAIN_CQ	<ul style="list-style-type: none"> <li>n</li> <li>y</li> </ul>	Poll CQ till it is completely drained of completed work requests. Enabling this feature may cause starvation of other endpoints Value: <y n>
MXM_IB_RESIZE_CQ	<ul style="list-style-type: none"> <li>n</li> <li>y</li> </ul>	Allow using resize_cq(). Value: <y n>

Variable	Valid Values and Default Values	Description
MXM_IB_TX_BATCH	16	Number of send WREs to batch in one post-send list. Larger values reduce the CPU usage, but increase the latency because we might need to process lots of send completions at once. Value: unsigned
MXM_IB_RX_BATCH	64	Number of post-receives to be performed in a single batch. Value: unsigned
MXM_IB_MAP_MODE	<ul style="list-style-type: none"> <li>• first: Map the first suitable HCA port to all processes (default).</li> <li>• affinity: Distribute evenly among processes based on CPU affinity.</li> <li>• nearest: Try finding nearest HCA port based on CPU affinity.</li> <li>• round-robin</li> </ul>	HCA ports to processes mapping method. Ports not supporting process requirements (e.g. DC support) will be skipped. Selecting a specific device will override this setting.
MXM_IB_NUM_SLS	1: (default)	Number of InfiniBand Service Levels to use. Every InfiniBand endpoint will generate a random SL within the given range <code>FIRST_SL..(FIRST_SL+NUM_SLS-1)</code> , and use it for outbound communication. applicable values are 1 through 16. Value: unsigned
MXM_IB_WC_MODE	<ul style="list-style-type: none"> <li>• wqe: Use write combining to post full WQEs (default).</li> <li>• db: Use write combining to post doorbell.</li> <li>• flush: Force flushing CPU write combining buffers.wqe</li> <li>• flush (default)</li> </ul>	Write combining mode flags for InfiniBand devices. Using write combining for 'wqe' improves latency performance due to one less wqe fetch. Avoiding 'flush' relaxes CPU store ordering, and reduces overhead. Write combining for 'db' is meaningful only when used without 'flush'. Value: comma-separated list of: [wqe db flush]

Variable	Valid Values and Default Values	Description
MXM_IB_LID_PATH_BITS	-( default) 0 <= value < 2^(LMC) - 1	List of InfiniBand Path bits separated by comma (a,b,c) which will be the low portion of the LID, according to the LMC in the fabric. If no value is given, MXM will use all the available offsets under the value of MXM_IB_MAX_PATH_BITS To disable the usage of LMC, please set this value to zero. Value: comma-separated list of: unsigned
MXM_IB_MAX_PATH_BITS	18: (default)	Number of offsets to use as lid_path_bits in case 2^LMC is larger than this value. This value derives from the number of ports on the switch. Value: unsigned
MXM_IB_FIRST_SL	0-15	The first Infiniband Service Level number to use. Value: unsigned
MXM_IB_CQ_STALL	100	CQ stall loops for SandyBridge far socket. Value: unsigned
MXM_UD_ACK_TIMEOUT	300000.00us	Timeout for getting an acknowledgment for sent packet. Value: time value: <number>[s us ms ns]
MXM_UD_FAST_ACK_TIMEOUT	1024.00us	Timeout for getting an acknowledgment for sent packet. Value: time value: <number>[s us ms ns]
MXM_FAST_TIMER_RESOLUTION	64.00us	Resolution of ud fast timer. The value is treated as a recommendation only. Real resolution may differ as mxm rounds up to power of two. Value: time value: <number>[s us ms ns]
MXM_UD_INT_MODE	rx	Traffic types to enable interrupt for. Value: comma-separated list of: [rx tx]
MXM_UD_INT_THRESH	20000.00us	The maximum amount of time that may pass following an mxm call, after which interrupts will be enabled. Value: time value: <number>[s us ms ns]

Variable	Valid Values and Default Values	Description
MXM_UD_WINDOW_SIZE	1024	The maximum number of unacknowledged packets that may be in transit. Value: unsigned
MXM_UD_MTU	65536	Maximal UD packet size. The actual MTU is the minimum of this value and the fabric MTU. Value: memory units: <number>[b kb mb gb]
MXM_UD_CA_ALGO	<ul style="list-style-type: none"> <li>• none: no congestion avoidance</li> <li>• bic: binary increase (default)</li> </ul>	Use congestion avoidance algorithm to dynamically adjust send window size. Value: [none bic]
MXM_UD_CA_LOW_WIN	0	Use additive increase multiplicative decrease congestion avoidance when current window is below this threshold. Value: unsigned
MXM_UD_RX_QUEUE_LEN	4096	Length of receive queue for UD QPs. Value: unsigned
MXM_UD_RX_MAX_BUFS	-1	Maximal number of receive buffers for one endpoint. -1 is infinite. Value: integer
MXM_UD_RX_MAX_INLINE	0	Maximal size of data to receive as inline. Value: memory units: <number>[b kb mb gb]
MXM_UD_RX_DROP_RATE	0	If nonzero, network packet loss will be simulated by randomly ignoring one of every X received UD packets. Value: unsigned
MXM_UD_RX_OOO	<ul style="list-style-type: none"> <li>• n</li> <li>• y</li> </ul>	If enabled, keep packets received out of order instead of discarding them. Must be enabled if network allows out of order packet delivery, for example, if Adaptive Routing is enabled Value: <y n>
MXM_UD_TX_QUEUE_LEN	128	Length of send queue for UD QPs. Value: unsigned
MXM_UD_TX_MAX_BUFS	32768	Maximal number of send buffers for one endpoint. -1 is infinite. Value: integer

Variable	Valid Values and Default Values	Description
MXM_UD_TX_MAX_INLINE	128	Bytes to reserve in TX WQE for inline data. Messages which are small enough will be sent inline. Value: memory units: <number>[b kb mb gb]
MXM_CIB_PATH_MTU	default	Path MTU for CIB QPs. Possible values are: default, 512, 1024, 2048, 4096. Setting “default” will select the best MTU for the device. Value: [default 512 1024 2048 4096]
MXM_CIB_HARD_ZCOPY_THRESH	16384	Threshold for using zero copy. Value: memory units: <number>[b kb mb gb]
MXM_CIB_MIN_RNR_TIMER	25	InfiniBand minimum receiver not ready timer, in seconds (must be $\geq 0$ and $\leq 31$ )
MXM_CIB_TIMEOUT	20	InfiniBand transmit timeout, plugged into formula: $4.096 \text{ microseconds} * (2^{\text{timeout}})$ (must be $\geq 0$ and $\leq 31$ ) Value: unsigned
MXM_CIB_MAX_RDMA_DST_OPS	4	InfiniBand maximum pending RDMA destination operations (must be $\geq 0$ ) Value: unsigned
MXM_CIB_RNR_RETRY	7	InfiniBand “receiver not ready” retry count, applies <b>ONLY</b> for SRQ/XRC queues. (must be $\geq 0$ and $\leq 7$ : 7 = “infinite”) Value: unsigned
MXM_UD_RNDV_THRESH	262144	UD threshold for using rendezvous protocol. Smaller value may harm performance but excessively large value can cause a deadlock in the application. Value: memory units: <number>[b kb mb gb]
MXM_UD_TX_NUM_SGE	3	Number of SG entries in the UD send QP. Value: unsigned
MXM_UD_HARD_ZCOPY_THRESH	65536	Threshold for using zero copy. Value: memory units: <number>[b kb mb gb]
MXM_CIB_RETRY_COUNT	7	InfiniBand transmit retry count (must be $\geq 0$ and $\leq 7$ ) Value: unsigned

Variable	Valid Values and Default Values	Description
MXM_CIB_MSS	4224	Size of the send buffer. Value: memory units: <number>[b kb mb gb]
MXM_CIB_TX_QUEUE_LEN	256	Length of send queue for RC QPs. Value: unsigned
MXM_CIB_TX_MAX_BUFS	-1	Maximal number of send buffers for one endpoint. -1 is infinite. <b>Warning:</b> Setting this param with value != -1 is a dangerous thing in RC and could cause deadlock or performance degradation Value: integer
MXM_CIB_TX_CQ_SIZE	16384	Send CQ length. Value: unsigned
MXM_CIB_TX_MAX_INLINE	128	Bytes to reserver in TX WQE for inline data. Messages which are small enough will be sent inline. Value: memory units: <number>[b kb mb gb]
MXM_CIB_TX_NUM_SGE	3	Number of SG entries in the RC QP. Value: unsigned
MXM_CIB_USE_EAGER_RDMA	y	Use RDMA WRITE for small messages. Value: <y n>
MXM_CIB_EAGER_RDMA_THRESHOLD	16	Use RDMA for short messages after this number of messages are received from a given peer, must be >= 1 Value: unsigned long
MXM_CIB_MAX_RDMA_CHANNELS	8	Maximum number of peers allowed to use RDMA for short messages, must be >= 0 Value: unsigned
MXM_CIB_EAGER_RDMA_BUFFS_NUM	32	Number of RDMA buffers to allocate per rdma channel, must be >= 1 Value: unsigned
MXM_CIB_EAGER_RDMA_BUFF_LEN	4224	Maximum size (in bytes) of eager RDMA messages, must be >= 1 Value: memory units: <number>[b kb mb gb]
MXM_SHM_RX_MAX_BUFFERS	-1	Maximal number of receive buffers for endpoint. -1 is infinite. Value: integer

Variable	Valid Values and Default Values	Description
MXM_SHM_RX_MAX_MEDIUM_BUFFERS	-1	Maximal number of medium sized receive buffers for one endpoint. -1 is infinite. Value: integer
MXM_SHM_FIFO_SIZE	64	Size of the shmем tl's fifo. Value: unsigned
MXM_SHM_WRITE_RETRY_COUNT	64	Number of retries in case where cannot write to the remote process. Value: unsigned
MXM_SHM_READ_RETRY_COUNT	64	Number of retries in case where cannot read from the shmем FIFO (for multi-thread support). Value: unsigned
MXM_SHM_HUGETLB_MODE	<ul style="list-style-type: none"> <li>• y: Allocate memory using huge pages only.</li> <li>• n: Allocate memory using regular pages only.</li> <li>• try: Try to allocate memory using huge pages and if it fails, allocate regular pages (default).</li> </ul>	Enable using huge pages for internal shared memory buffers Values: <yes no try>
MXM_SHM_RX_BUFF_LEN	8192	Maximum size (in bytes) of medium sized messages, must be >= 1 Value: memory units: <number>[b kb mb gb]
MXM_SHM_HARD_ZCOPY_THRESH	2048	Threshold for using zero copy. Value: memory units: <number>[b kb mb gb]
MXM_SHM_RNDV_THRESH	65536	SHM threshold for using rendezvous protocol. Smaller value may harm performance but too large value can cause a deadlock in the application. Value: memory units: <number>[b kb mb gb]
MXM_SHM_KNEM_MAX_SIMULTANEOUS	0	Maximum number of simultaneous ongoing knem operations to support in shmем tl. Value: unsigned



Variable	Valid Values and Default Values	Description
MXM_SHM_KNEM_DMA_CHUNK_SIZE	67108864	Size of a single chunk to be transferred in one dma operation. Larger values may not work since they are not supported by the dma engine Value: memory units: <number>[b kb mb gb]
MXM_SHM_RELEASE_FIFO_FACTOR	0.500	Frequency of resource releasing on the receiver's side in shmем tl. This value refers to the percentage of the fifo size. (must be $\geq 0$ and $< 1$ ) Value: floating point number
MXM_TM_UPDATE_THRESHOLD_MASK	8	Update bit-mask length for connections traffic counters.(must be $\geq 0$ and $\leq 32$ : 0 - always update, 32 - never update.)# Value: bit count
MXM_TM_PROMOTE_THRESHOLD	5	Relative threshold (percentage) of traffic for promoting connections, Must be $\geq 0$ . Value: unsigned
MXM_TM_PROMOTE_BACKOFF	1	Exponential backoff degree (bits) for all counters upon a promotion, Must be $\geq 0$ and $\leq 32$ . Value: unsigned
MXM_RC_QP_LIMIT	64	Maximal amount of RC QPs allowed (Negative for unlimited). Value: integer
MXM_DC_QP_LIMIT	64	Maximal amount of DC QPs allowed (Negative for unlimited). Value: integer
MXM_DC_RECV_INLINE	<ul style="list-style-type: none"> <li>• 128</li> <li>• 512</li> <li>• 1024</li> <li>• 2048</li> <li>• 4096</li> </ul>	Bytes to reserve in CQE for inline data. In order to allow for inline data, -x MLX5_CQE_SIZE=128 must also be specified. Value: memory units: <number>[b kb mb gb]
MXM_CIB_RNDV_THRESH	16384	CIB threshold for using rendezvous protocol. Smaller value may harm performance but excessively large value can cause a deadlock in the application. Value: memory units: <number>[b kb mb gb]

Variable	Valid Values and Default Values	Description
MXM_SHM_USE_KNEM_DMA	<ul style="list-style-type: none"> <li>n</li> <li>y</li> </ul>	Whether or not to offload to the DMA engine when using KNEM Value: <y n>
MXM_IB_GID_INDEX	<ul style="list-style-type: none"> <li>0</li> </ul>	GID index to use as local Ethernet port address.

### 3.7 Loading KNEM Module

MXM's intra-node communication uses the KNEM module which improves the performance significantly.

➤ **To use the KNEM module:**

- Load the KNEM module.

Please run the following commands on all cluster nodes to enable KNEM intra-node device.

```
# insmod $HPCX_HOME/knem/lib/modules/$(uname -r)/knem.ko
# chmod 666 /dev/knem
```



On RHEL systems, to enable the KNEM module on machine boot, add these commands into the `/etc/rc.modules` script.

Making `/dev/knem` public accessible poses no security threat, as only the memory buffer that was explicitly made readable and/or writable can be accessed read and/or write through the 64bit cookie. Moreover, recent KNEM releases enforce by default that the attacker and the target process have the same UID which prevent any security issues.

### 3.8 Running MPI with FCA v2.5

Make sure FCA manager is running in the fabric.

```
$ mpirun -mca coll_fca_enable 1 <...>
```

For further details on starting FCA manager, please refer to [3.1 “Starting FCA Manager from HPC-X,” on page 12](#)

### 3.9 Running OpenSHMEM with FCA v2.5

Make sure FCA manager is running in the fabric.

```
$ oshrun -mca scoll_fca_enable 1 <...>
```

For further details on starting FCA manager, please refer to [3.1 “Starting FCA Manager from HPC-X,” on page 12](#)

### 3.10 Running MPI with FCA v3.5 (hcoll)



FCA v3.5 is enabled by default in HPC-X.

- Running with default FCA configuration parameters:

```
$ mpirun -mca coll_hcoll_enable 1 -x HCOLL_MAIN_IB=mlx4_0:1 <...>
```

- Running OSHMEM with FCA v3.5:

```
% oshrun -mca scoll_mpi_enable 1 -mca scoll basic,mpi -mca coll_hcoll_enable 1 <...>
```

### 3.11 IB-Router

HPC-X 1.6 supports `ib-router` to allow hosts that are located on different IB subnets to communicate with each other. This support is currently available when using the `'openib btl'` in Open MPI.

To use `ib-router`, make sure `MLNX_OFED v3.3-1.0.0.0` is installed. and then recompile your Open MPI with `'--enable-openib-rdmacm-ibaddr'` (for further information of how to compile Open MPI, refer to [Section 3.4, “Rebuilding Open MPI from HPC-X™ Sources”, on page 13](#))

➤ **To enable routing over IB, please follow these steps:**

1. Configure Open MPI with `--enable-openib-rdmacm-ibaddr`.
2. Use `rdmacm` with `openib btl` from the command line.
3. Set the `btl_openib_allow_different_subnets` parameter to 1.  
It is 0 by default.
4. Set the `btl_openib_gid_index` parameter to 1.

For example - to run the IMB benchmark on `host1` and `host2` which are on separate subnets, i.e. have different `subnet_prefix`, use the following command line:

```
shell$ mpirun -np 2 --display-map --map-by node -H host1,host2 -mca pml ob1 -mca btl self,sm,openib --mca btl_openib_cpc_include rdmacm -mca btl_openib_if_include mlx5_0:1 -mca btl_openib_gid_index 1 -mca btl_openib_allow_different_subnets 1 ./IMB/src/IMB-MPI1 pingpong
```

More information about how to enable and use `ib-router` is here - <https://www.open-mpi.org/faq/?category=openfabrics#ib-router>



When using `"openib btl"`, RoCE and IB router are mutually exclusive. The Open MPI inside HPC-X 1.6 is not compiled with `ib-router` support, therefore it supports RoCE out-of-the-box.

## 4 Mellanox Fabric Collective Accelerator (FCA)

### 4.1 Overview

To meet the needs of scientific research and engineering simulations, supercomputers are growing at an unrelenting rate. As supercomputers increase in size from mere thousands to hundreds-of-thousands of processor cores, new performance and scalability challenges have emerged. In the past, performance tuning of parallel applications could be accomplished fairly easily by separately optimizing their algorithms, communication, and computational aspects. However, as systems continue to scale to larger machines, these issues become co-mingled and must be addressed comprehensively.

Collective communications execute global communication operations to couple all processes/nodes in the system and therefore must be executed as quickly and as efficiently as possible. Indeed, the scalability of most scientific and engineering applications is bound by the scalability and performance of the collective routines employed. Most current implementations of collective operations will suffer from the effects of systems noise at extreme-scale (system noise increases the latency of collective operations by amplifying the effect of small, randomly occurring OS interrupts during collective progression.) Furthermore, collective operations will consume a significant fraction of CPU cycles, cycles that could be better spent doing meaningful computation.

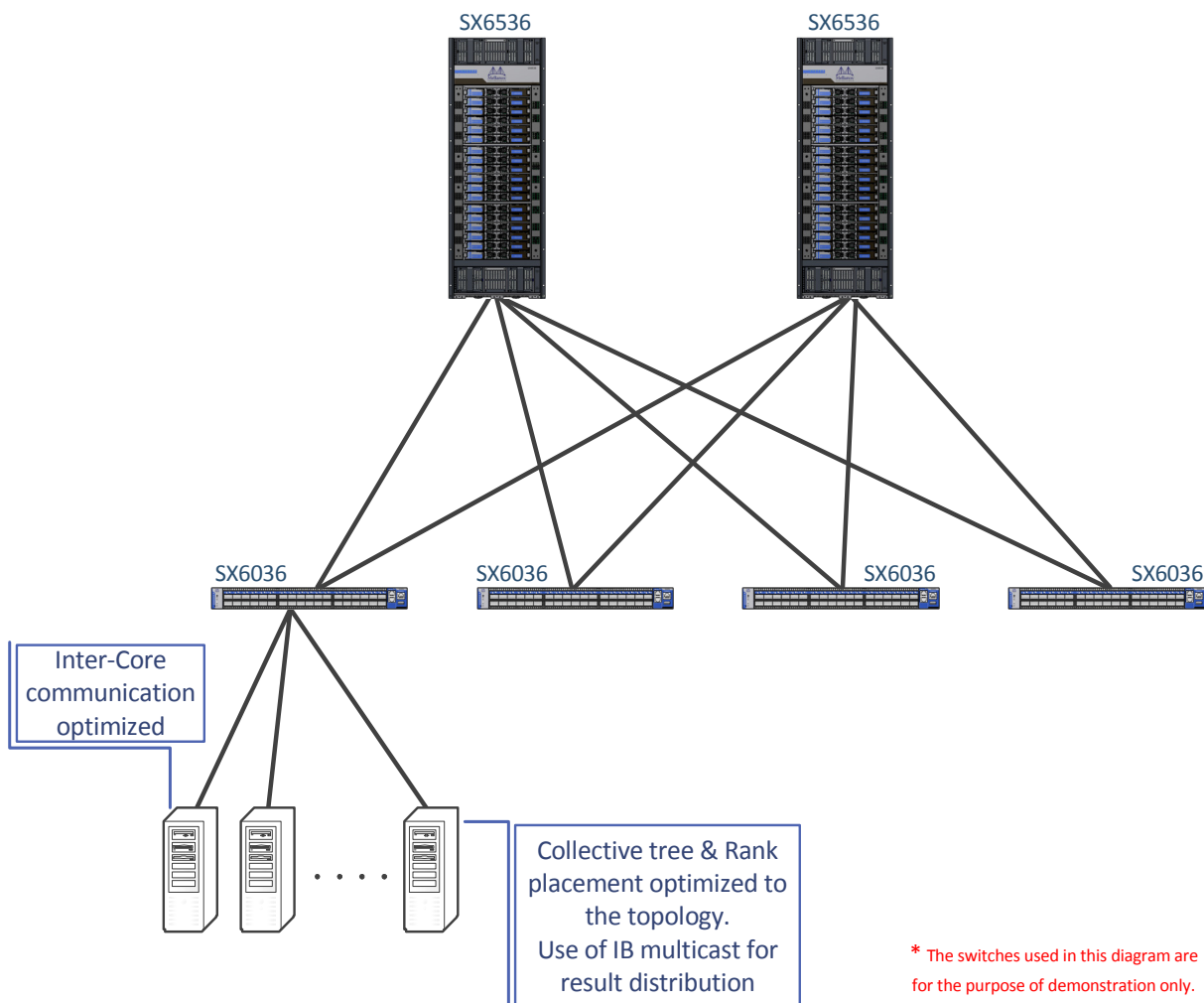
Mellanox Technologies has addressed these two issues, lost CPU cycles and performance lost to the effects of system noise, by offloading the communications to the host channel adapters (HCAs) and switches. The technology, named CORE-Direct® (Collectives Offload Resource Engine), provides the most advanced solution available for handling collective operations thereby ensuring maximal scalability, minimal CPU overhead, and providing the capability to overlap communication operations with computation allowing applications to maximize asynchronous communication.

Additionally, FCA v3.5 also contains support for building runtime configurable hierarchical collectives. We currently support socket and UMA level discovery with network topology available in beta form. As with FCA 2.X, FCA v3.5 leverages hardware multicast capabilities to accelerate collective operations. In FCA v3.5 we also expose the performance and scalability of Mellanox's advanced point-to-point library, MXM, in the form of the "mlnx\_p2p" BCOL. This allows users to take full advantage of new hardware features with minimal effort.

FCA v3.5 and above is a standalone library that can be integrated into any MPI or PGAS runtime. Support for FCA is currently integrated into Open MPI versions 1.7.4 and higher. FCA v3.5 release currently supports blocking and non-blocking variants of "Allgather", "Allreduce", "Barrier", and "Bcast".

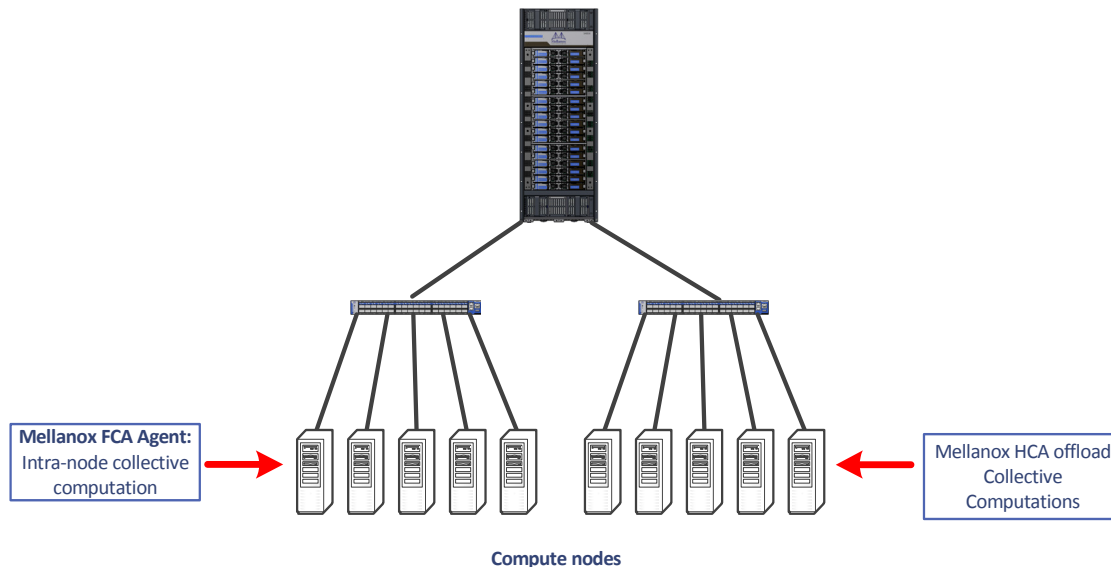
The following diagram summarizes the FCA architecture:

**Figure 1: FCA Architecture**



The following diagram shows the FCA components and the role that each plays in the acceleration process:

**Figure 2: FCA Components**



## 4.2 FCA Installation Package Content



hcoll is part of the HPC-X software toolkit and does not require special installation.

The FCA installation package includes the following items:

- FCA- Mellanox Fabric Collector Accelerator Installation files
  - hcoll-<version>.x86\_64.<OS>.rpm
  - hcoll-<version>.x86\_64.<OS>.tar.gz
- where:
  - <version>: The version of this release
  - <OS>: One of the supported Linux distributions listed in Prerequisites (on page 11).
- Mellanox Fabric Collective Accelerator (FCA) Software: End-User License Agreement
- FCA MPI runtime libraries
- Mellanox Fabric Collective Accelerator (FCA) Release Notes

## 4.3 Differences Between FCA v2.5 and FCA v3.x

FCA v3.x is new software which continues to expose the power of CORE-Direct® to offload collective operations to the HCA. It adds additional scalable algorithms for collectives and supports both blocking and non-blocking APIs (MPI-3 SPEC compliant). Additionally, FCA v3.x (hcoll) does not require FCA manager daemon.

## 4.4 Configuring FCA

### 4.4.1 Compiling Open MPI with FCA 3.x

➤ *To compile Open MPI with FCA 3.x*

**Step 1.** Install FCA 3.x from:

- an RPM.

```
# rpm -ihv hcoll-x.y.z-1.x86_64.rpm
```

- a tarball.

```
% tar jxf hcoll-x.y.z.tbz
```

FCA 3.X will be installed automatically in the /opt/mellanox/hcoll folder.

**Step 2.** Enter the Open MPI source directory and run the following command:

```
% cd $OMPI_HOME
% ./configure --with-hcoll=/opt/mellanox/hcoll --with-mxm=/opt/mellanox/mxm < ... other
configure parameters>
% make -j 9 && make install -j 9
```



libhcoll requires MXM v2.1 or higher.

➤ *To check the version of FCA installed on your host:*

```
% rpm -qi hcoll
```

➤ *To upgrade to a newer version of FCA 3.X:*

**Step 1.** Remove the existing FCA 3.X.

```
% rpm -e hcoll
```

**Step 2.** Remove the precompiled Open MPI.

```
% rpm -e mlnx-openmpi_gcc
```

**Step 3.** Install the new FCA 3.X and compile the Open MPI with it.

### 4.4.2 Enabling FCA in Open MPI

To enable FCA 3.X HCOLL collectives in Open MPI, explicitly ask for them by setting the following MCA parameter:

```
%mpirun -np 32 -mca coll_hcoll_enable 1 -x coll_hcoll_np=0 -x HCOLL_
MAIN_IB=<device_name>:<port_num> -x MXM_RDMA_PORTS=<device_name>:<port_num> -x
HCOLL_ENABLE_MCAST_ALL=1 ./a.out
```





Or you can use `/sbin/ip` for the same purpose

```
hpchead ~ >ip addr show ib0
4: ib0: <BROADCAST,MULTICAST> mtu 2044 qdisc mq state DOWN qlen 1024
    link/infiniband a0:04:02:20:fe:80:00:00:00:00:00:00:00:02:c9:03:00:21:f9:31 brd
    00:ff:ff:ff:ff:12:40:1b:ff:ff:00:00:00:00:00:00:00:ff:ff:ff:ff
    inet 192.168.1.1/24 brd 192.168.1.255 scope global ib0
```

In the example above, the IP is defined (192.168.1.1). If it is not defined, then you can define an IP address now.

#### 4.4.6 Enabling HCOLL Topology Awareness

FCA v3.5 (HCOLL) is able to leverage fabric topology information. To allow HCOLL to take advantage of this capability, the steps below need to be executed during fabric bring-up in order to stage the proper information for HCOLL to retrieve.

1. Generate a `fabric.cache` file.

```
ibnetdiscover --cache fabric.cache
```

In case of several cards connected as independent fabrics, choose the correct card/port to start from:

```
ibnetdiscover -C mlx5_0 -P 1 --cache fabric.cache
```

2. Stage the `fabric.cache` file on NFS.

The `fabric.cache` file can be placed on a shared file system in one of these locations:

- In a user defined location, `/path/to/fabric.cache`. The location should be passed to HCOLL via the command line parameter value

```
-x HCOLL_TOPOLOGY_DATAFILE=/path/to/fabric.cache
```

- The standard location where HCOLL expects to find the `file:$HPCX_HCOLL_DIR/etc/fabric.cache`

- a. Node local placement of `fabric.cache`.

Super user needs to run as root:

```
export cluster=[partition_name]
ibnetdiscover --cache /tmp/${cluster}_fabrictopo.cache > /tmp/${cluster}_fabrictopo.data
sudo pdsh -P $cluster ibstat > /tmp/${cluster}_fabrictopo.map
```



The `fabric.cache` file must be rebuilt if there were changes in the cluster topology (e.g. connections.)

All `fabric.cache` files across the fabric must be identical. For this purpose each rank calculates the `md5sum` of its `fabric.cache` file and exchanges this data with all other ranks.

- **To disable this feature:**

```
-x HCOLL_ENABLE_TOPOLOGY=0
```

#### 4.4.7 Enabling SHArP Accelerated Collectives

HPC-X Rev 1.6 supports SHArP Accelerated Collectives. These collectives are enabled by default if FCA 3.5 (HCOLL) detects that it is running in a supporting environment.

For instructions on how to deploy SHArP in Infiniband fabric, see the SHArP Deployment Guide.

Once SHArP is deployed, you need to only specify the HCA device (`device_name`) and port number (`port_num`) that is connected to the SHArP tree in the following way:

```
-x HCOLL_MAIN_IB=<device_name>:<port_num>
```

➤ **To disable SHArP acceleration:**

```
-x HCOLL_ENABLE_SHARP=0
```

## 5 MellanoX Messaging Library

### 5.1 Overview

MellanoX Messaging (MXM) library provides enhancements to parallel communication libraries by fully utilizing the underlying networking infrastructure provided by Mellanox HCA/switch hardware. This includes a variety of enhancements that take advantage of Mellanox networking hardware including:

- Multiple transport support including RC, DC and UD
- Proper management of HCA resources and memory structures
- Efficient memory registration
- One-sided communication semantics
- Connection management
- Receive side tag matching
- Intra-node shared memory communication

These enhancements significantly increase the scalability and performance of message communications in the network, alleviating bottlenecks within the parallel communication libraries.

The latest MXM software can be downloaded from the [Mellanox website](#).

### 5.2 Compiling Open MPI with MXM



MXM has been integrated into the HPC-X Toolkit package. The steps described below are only required if you have downloaded the `mxm.rpm` from the Mellanox site

**Step 1.** Install MXM from:

- an RPM

```
% rpm -ihv mxm-x.y.z-1.x86_64.rpm
```

- a tarball

```
% tar jxf mxm-x.y.z.tar.bz
```

MXM will be installed automatically in the `/opt/mellanox/mxm` folder.

**Step 2.** Enter Open MPI source directory and run:

```
% cd $OMPI_HOME
% ./configure --with-mxm=/opt/mellanox/mxm <... other configure parameters...>
% make all && make install
```

Older versions of MLNX\_OFED come with pre-installed older MXM and Open MPI versions. Uninstall any old MXM version prior to installing the latest MXM version in order to use it with older MLNX\_OFED versions.

**Table 3 - MLNX\_OFED and MXM Versions**

MLNX_OFED Version	MXM Version
v1.5.3-3.1.0 and v2.0-3.0.0	MXM v1.x and Open MPI compiled with MXM v1.x
v2.0-3.0.0 and higher	MXM v2.x/3.x and Open MPI compiled with MXM v2.x/3.x

To check the version of MXM installed on your host, run:

```
% rpm -qi mxm
```

➤ **To upgrade MLNX\_OFED v1.5.3-3.1.0 or later with a newer MXM:**

**Step 1.** Remove MXM.

```
# rpm -e mxm
```

**Step 2.** Remove the pre-compiled Open MPI.

```
# rpm -e mlnx-openmpi_gcc
```

**Step 3.** Install the new MXM and compile the Open MPI with it.



To run Open MPI without MXM, run:

```
% mpirun -mca mtl ^mxm --mca pml ^yalla <...>
```



When upgrading to MXM v3.5, Open MPI compiled with the previous versions of the MXM should be recompiled with MXM v3.5.

## 5.3 Enabling MXM in Open MPI

MXM is selected automatically starting from any Number of Processes (NP) when using Open MPI v1.8.x and above.

For older Open MPI versions, use the command below.

➤ **To activate MXM for any NP, run:**

```
% mpirun -mca mtl_mxm_np 0 <...other mpirun parameters ...>
```

As of Open MPI v1.8, MXM is selected when the number of processes is higher or equal to 0. i.e. by default.

## 5.4 Tuning MXM Settings

The default MXM settings are already optimized. To check the available MXM parameters and their default values, run the `$HPCX_MXM_DIR/bin/mxm_dump_config -f` utility which is part of the MXM RPM.

MXM parameters can be modified in one of the following methods:

- Modifying the default MXM parameters value as part of the `mpirun`:

```
% mpirun -x MXM_UD_RX_MAX_BUFFS=128000 <...>
```

- Modifying the default MXM parameters value from SHELL:

```
% export MXM_UD_RX_MAX_BUFFS=128000
% mpirun <...>
```

## 5.5 Configuring Multi-Rail Support

Multi-Rail support enables the user to use more than one of the active ports on the card, making better use of system resources, allowing increased throughput.

Multi-Rail support as of MXM v3.5 allows different processes on the same host to use different active ports. Every process can only use one port (as opposed to MXM v1.5).

### ➤ *To configure dual rail support:*

- Specify the list of ports you would like to use to enable multi rail support.

```
-x MXM_RDMA_PORTS=cardName:portNum
```

or

```
-x MXM_IB_PORTS=cardName:portNum
```

For example:

```
-x MXM_IB_PORTS=mlx5_0:1
```

It is also possible to use several HCAs and ports during the run (separated by a comma):

```
-x MXM_IB_PORTS=mlx5_0:1,mlx5_1:1
```

MXM will bind a process to one of the HCA ports from the given ports list according to the `MXM_IB_MAP_MODE` parameter (for load balancing).

Possible values for `MXM_IB_MAP_MODE` are:

- `first` - [Default] Maps the first suitable HCA port to all processes
- `affinity` - Distributes the HCA ports evenly among processes based on CPU affinity
- `nearest` - Tries to find the nearest HCA port based on CPU affinity

You may also use an asterisk (\*) and a question mark (?) to choose the HCA and the port you would like to use.

- \* - use all active cards/ports that are available
- ? - use the first active card/port that is available

For example:

```
-x MXM_IB_PORTS=*:?
```

will take all the active HCAs and the first active port on each of them.

## 5.6 Configuring MXM over the Ethernet Fabric

➤ *To configure MXM over the Ethernet fabric:*

**Step 1.** Make sure the Ethernet port is active.

```
% ibv_devinfo
```



`ibv_devinfo` displays the list of cards and ports in the system. Make sure (in the `ibv_devinfo` output) that the desired port has Ethernet at the `link_layer` field and that its state is `PORT_ACTIVE`.

**Step 2.** Specify the ports you would like to use, if there is a non Ethernet active port in the card.

```
-x MXM_RDMA_PORTS=mlx4_0:1
```

or

```
-x MXM_IB_PORTS=mlx4_0:1
```



Please note that the `MXM_RDMA_PORTS` and `MXM_IB_PORTS` parameters are alias that mean the same used in both Section 5.5 (page 37) and Section 5.6 (page 38).

`MXM_IB_PORTS` and `MXM_RDMA_PORTS` are aliases which mean the same, but one is used in Section 5.5 (page 37) and the second is Section 5.6 (page 38).

## 5.7 Running MXM with RoCE

MXM supports 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.

### 5.7.1 Running MXM with RoCE v1

To use RoCE with MXM, the desired Ethernet port must be specified (with the `MXM_IB_PORTS` parameter).

```
mpirun --mca pml yalla -x MXM_RDMA_PORTS=mlx4_0:2 ...
```

For further information on the `pml yalla`, please refer to Section 5.10, “Running Open MPI with `pml “yalla”`”, on page 40.

If using older versions of MXM, the following additional environment parameter might be required to be set:

```
mpirun --mca pml yalla -x MXM_TLS=self,shm,ud -x MXM_RDMA_PORTS=mlx4_0:2 ...
```

## 5.7.2 Running MXM with RoCE v2

To enable RoCE v2, add the following parameter to the command line.

```
mpirun --mca pml yalla -x MXM_RDMA_PORTS=mlx4_0:2 -x MXM_IB_GID_INDEX=1 ...
```

### 5.7.2.1 Using RoCE v2 in Open MPI

RoCE v2 is supported in Open MPI as of v1.8.8.

To use it, enable RoCE v2 in the command line:

```
$ mpirun --mca pml ob1 --mca btl openib,self,sm --mca btl_openib_cpc_include rdmacm --mca btl_openib_rroce_enable 1 ...
```

## 5.8 Configuring MXM over Different Transports

MXM v3.5 supports the following transports.

- Intra node communication via Shared Memory with KNEM support
- Unreliable Datagram (UD)
- Reliable Connected (RC)
- SELF transport - a single process communicates with itself
- Dynamically Connected Transport (DC)

**Note:** DC is supported on Connect-IB® HCAs with MLNX\_OFED v2.1-1.0.0 and higher.

To use DC set the following:

- in the command line:

```
% mpirun -x MXM_TLS=self,shm,dc
```

- from the SHELL:

```
% export MXM_TLS=self,shm,dc
```

By default the transports (TLS) used are: `MXM_TLS=self,shm,ud`

## 5.9 Configuring Service Level Support

Service Level enables Quality of Service (QoS). If set, every InfiniBand endpoint in MXM will generate a random Service Level (SL) within the given range, and use it for outbound communication.

Setting the value is done via the following environment parameter:

```
export MXM_IB_NUM_SLS=1
```

Available Service Level values are 1-16 where the default is 1.

You can also set a specific service level to use. To do so, use the `MXM_IB_FIRST_SL` parameter together with `MXM_IB_NUM_SLS=1` (which is the default).

For example:

```
% mpirun -x MXM_IB_NUM_SLS=1 -x MXM_IB_FIRST_SL=3 ...
```

where a single SL is being used and the first SL is 3.

## 5.10 Running Open MPI with pml “yalla”

A new pml layer (pml “yalla”) was added to Mellanox's Open MPI v1.8. It is used to reduce overhead by cutting through layers and using MXM directly. Consequently, for messages < 4K in size, it yields a latency improvement of up to 5%, message rate of up to 50% and bandwidth of up to 45%.

Open MPI's default behavior is to run with pml 'yalla'. To directly use MXM with it, perform the steps below.

### ➤ To run MXM with the pml “yalla”:

**Step 1.** Download the HPC-X package from the Mellanox site (Mellanox's Open MPI v1.10 has been integrated into HPC-X package).

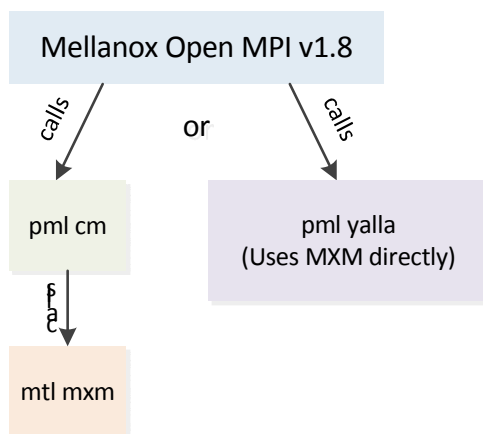
[http://www.mellanox.com/page/products\\_dyn?product\\_family=189&mtag=hpc-x](http://www.mellanox.com/page/products_dyn?product_family=189&mtag=hpc-x)

**Step 2.** Use the new yalla pml.

```
% mpirun -mca pml yalla
```



pml “yalla” uses MXM directly by default. If the pml “yalla” is not used, Open MPI will use MXM through “pml cm” and “mtl mxm” (see figure below).





## 5.11 Adaptive Routing for UD Transport

Adaptive Routing (AR) enables the switch to select the output port based on the port's load. Adaptive Routing for the UD transport layer enable out of order packet arrival.

- When the `MXM_UD_RX_OOO=n`, parameter is set to "n", the out of order packet indicates a packet loss and triggers MXM UD flow control/congestion avoidance.
- When the parameter `MXM_UD_RX_OOO=y`, is set to "y" MXM will queue out of order packets until it is possible to process them in order instead of assuming a packet loss.

To configure adaptive routing one must use OpenSM with adaptive route manager plugin and a switch with Mellanox OS. AR support in UD is set to ON by default.

➤ *To disable Adaptive Routing for UD transport:*

```
% mpirun -x MXM_UD_RX_OOO=y ...
```

## 5.12 Support for a Non-Base LID

MXM enables the user to use multiple LIDs per port when the LID Mask Control (LMC) configured in the fabric is higher than zero. By default, MXM will use all the possible LIDs that are enabled except for the LMC (from zero to  $2^{LMC}$ ) unless their number is higher than the `MXM_IB_MAX_PATH_BITS` parameter in which case the latter value will be used (this number is derived from the number of ports on a switch). MXM also enables the user to specify a list of LIDs to use via the `MXM_IB_LID_PATH_BITS` parameter.

The usage of multiple LIDs enables MXM to distribute the traffic in the fabric over multiple LIDs and therefore achieve a better usage of the bandwidth. It prevents the overloading of single link and results in an improved performance and more balanced traffic.

## 5.13 MXM Performance Tuning

MXM uses the following features to improve performance:

- **Bulk Connections**

The `mxm_ep_wireup` and `mxm_ep_powerdown` functions were added to the MXM API to allow pre-connection establishment for MXM. This will enable MXM to create and connect all the required connections for the future communication during the initialization stage rather than creating the connections between the peers in an on-demand manner.



Please note that the usage of these parameters is only possible when using MXM with "`-mca pml cm --mca mtl mxm`", not when running with the `yalla pml`.

Bulk connection is the default for establishing connection for MXM in the Open MPI, 'mtl mxm' layer.

When running an application that has a sparse communication pattern, it is recommended to disable Bulk Connections.

➤ *To enable Bulk Connections:*

```
% mpirun -mca mtl_mxm_bulk_connect 1 -mca mtl_mxm_bulk_disconnect 1 ...
```

➤ ***To disable Bulk Connections:***

```
% mpirun -mca mtl_mxm_bulk_connect 0 -mca mtl_mxm_bulk_disconnect 0 ...
```

- **Solicited event interrupt for the rendezvous protocol**

The solicited event interrupt for the rendezvous protocol improves performance for applications which have large messages communication overlapping with computation. This feature is disabled by default.

➤ *To enable Solicited event interrupt for the rendezvous protocol:*

```
% mpirun -x MXM_RNDV_WAKEUP_THRESH=512k ...
```

\*<thresh>: Minimal message size which will trigger an interrupt on the remote side, to switch the remote process from computation phase and force it to handle MXM communication.

## 5.14 MXM Utilities



When MXM is used as part of HPC-X software toolkit, the MXM utilities can be found at the `$HPCX_HOME/mxm/bin` directory.

When MXM is used as part of MLNX\_OFED driver, the MXM utilities can be found at the `/opt/mellanox/mxm/bin` directory.

### 5.14.1 mxm\_dump\_config

Enables viewing of all the environment parameters that MXM uses.

To see all the parameters, run: `$HPCX_HOME/mxm/bin/mxm_dump_config -f`.

For further information, run: `$HPCX_HOME/mxm/bin/mxm_dump_config -help`



Environment parameters can be set by using the “export” command.

For example, to set the `MXM_TLS` environment parameter, run:

```
% export MXM_TLS=<...>
```

### 5.14.2 mxm\_perftest

A client-server based application which is designed to test MXM's performance and sanity checks on MXM.

To run it, two terminals are required to be opened, one on the server side and one on the client side.

The working flow is as follow:

1. The server listens to the request coming from the client.
2. Once a connection is established, MXM sends and receives messages between the two sides according to what the client requested.
3. The results of the communications are displayed.

For further information, run: `$HPCX_HOME/mxm/bin/mxm_perftest -help`.

**Example:**

- From the server side run: `$HPCX_HOME/mxm/bin/mxm_perftest`
- From the client side run:  
`$HPCX_HOME/mxm/bin/mxm_perftest <server_host_name> -t send_lat`

Among other parameters, you can specify the test you would like to run, the message size and the number of iterations.

## 6 PGAS Shared Memory Access Overview

The Shared Memory Access (SHMEM) routines provide low-latency, high-bandwidth communication for use in highly parallel scalable programs. The routines in the SHMEM Application Programming Interface (API) provide a programming model for exchanging data between cooperating parallel processes. The SHMEM API can be used either alone or in combination with MPI routines in the same parallel program.

The SHMEM parallel programming library is an easy-to-use programming model which uses highly efficient one-sided communication APIs to provide an intuitive global-view interface to shared or distributed memory systems. SHMEM's capabilities provide an excellent low level interface for PGAS applications.

A SHMEM program is of a single program, multiple data (SPMD) style. All the SHMEM processes, referred as processing elements (PEs), start simultaneously and run the same program. Commonly, the PEs perform computation on their own sub-domains of the larger problem, and periodically communicate with other PEs to exchange information on which the next communication phase depends.

The SHMEM routines minimize the overhead associated with data transfer requests, maximize bandwidth, and minimize data latency (the period of time that starts when a PE initiates a transfer of data and ends when a PE can use the data).

SHMEM routines support remote data transfer through:

- “put” operations - data transfer to a different PE
- “get” operations - data transfer from a different PE, and remote pointers, allowing direct references to data objects owned by another PE

Additional supported operations are collective broadcast and reduction, barrier synchronization, and atomic memory operations. An atomic memory operation is an atomic read-and-update operation, such as a fetch-and-increment, on a remote or local data object.

SHMEM libraries implement active messaging. The sending of data involves only one CPU where the source processor puts the data into the memory of the destination processor. Likewise, a processor can read data from another processor's memory without interrupting the remote CPU. The remote processor is unaware that its memory has been read or written unless the programmer implements a mechanism to accomplish this.

### 6.1 HPC-X Open MPI/OpenSHMEM

HPC-X Open MPI/OpenSHMEM programming library is a one-side communications library that supports a unique set of parallel programming features including point-to-point and collective routines, synchronizations, atomic operations, and a shared memory paradigm used between the processes of a parallel programming application.

HPC-X OpenSHMEM is based on the API defined by the OpenSHMEM.org consortium. The library works with the OpenFabrics RDMA for Linux stack (OFED), and also has the ability to utilize Mellanox Messaging libraries (MXM) as well as Mellanox Fabric Collective Accelerations (FCA), providing an unprecedented level of scalability for SHMEM programs running over InfiniBand.

## 6.2 Running HPC-X OpenSHMEM

### 6.2.1 Running HPC-X OpenSHMEM with MXM

MellanoX Messaging (MXM) library provides enhancements to parallel communication libraries by fully utilizing the underlying networking infrastructure provided by Mellanox HCA/switch hardware. This includes a variety of enhancements that take advantage of Mellanox networking hardware including:

- Multiple transport support including RC, DC and UD
- Proper management of HCA resources and memory structures
- Efficient memory registration
- One-sided communication semantics
- Connection management
- Receive side tag matching
- Intra-node shared memory communication

These enhancements significantly increase the scalability and performance of message communications in the network, alleviating bottlenecks within the parallel communication libraries

#### 6.2.1.1 Enabling MXM for HPC-X OpenSHMEM Jobs

##### ➤ *To enable MXM for SHMEM jobs for any PE:*

- Add the following MCA parameter to the `oshrun` command line.

```
-mca spml_ikrit_np <number>
```

##### ➤ *To force MXM usage:*

- Add the following MCA parameter `oshrun` command line.

```
-mca spml_ikrit_np 0
```

For additional MXM tuning information, please refer to the MellanoX Messaging Library README file found in the [Mellanox website](#).

#### 6.2.1.2 Working with Multiple HCAs

If there several HCAs in the system, MXM will choose the first HCA with the active port to work with. The HCA/port to be used can be specified by setting the `MXM_RDMA_PORTS` environment variable. The variable format is as follow: `MXM_RDMA_PORTS=hca_name:port,...`

For example, the following will cause MXM to use port one on two installed HCAs:

```
MXM_RDMA_PORTS=mlx4_0:1,mlx4_1:1
```



The environment variables must be run via the `oshrun` command line:

```
% oshrun -x MXM_RDMA_PORTS=mlx4_0:1 ...
```

## 6.2.2 Running HPC-X™ OpenSHMEM with FCA v2.x

The Mellanox Fabric Collective Accelerator (FCA) is a unique solution for offloading collective operations from the Message Passing Interface (MPI) or HPC-X OpenSHMEM process onto Mellanox InfiniBand managed switch CPUs. As a system-wide solution, FCA utilizes intelligence on Mellanox InfiniBand switches, Unified Fabric Manager and MPI nodes without requiring additional hardware. The FCA manager creates a topology based collective tree, and orchestrates an efficient collective operation using the switch-based CPUs on the MPI/HPC-X OpenSHMEM nodes.

FCA accelerates MPI/HPC-X OpenSHMEM collective operation performance by up to 100 times providing a reduction in the overall job runtime. Implementation is simple and transparent during the job runtime.



FCA v2.x is disabled by default and *must* be configured prior to using it from the HPC-X OpenSHMEM.

➤ **To enable FCA by default in the HPC-X OpenSHMEM:**

1. Edit the `$HPCX_OSHMEM_DIR/etc/openmpi-mca-params.conf` file.
2. Set the `scoll_fca_enable` parameter to 1.

```
scoll_fca_enable=1
```

3. Set the `scoll_fca_np` parameter to 0.

```
scoll_fca_np=0
```

➤ **To enable FCA in the `oshrun` command line, add the following:**

```
-mca scoll_fca_enable=1
-mca scoll_fca_enable_np 0
```

➤ **To disable FCA:**

```
-mca scoll_fca_enable 0 -mca coll_fca_enable 0
```

For more details on FCA installation and configuration, please refer to the FCA User Manual found in the [Mellanox website](#).

## 6.2.3 Developing Application using HPC-X OpenSHMEM together with MPI

The SHMEM programming model can provide a means to improve the performance of latency-sensitive sections of an application. Commonly, this requires replacing MPI send/recv calls with `shmem_put/ shmem_get` and `shmem_barrier` calls. The SHMEM programming model can deliver significantly lower latencies for short messages than traditional MPI calls. An alternative to `shmem_get /shmem_put` calls can also be considered the MPI-2 `MPI_Put/ MPI_Get` functions.

An example of MPI-SHMEM mixed code.

```
/* example.c */

#include <stdlib.h>
#include <stdio.h>
#include "shmem.h"
#include "mpi.h"

int main(int argc, char *argv[])
{
```

```

MPI_Init(&argc, &argv);
    start_pes(0);

    {
        int version = 0;
        int subversion = 0;
        int num_proc = 0;
        int my_proc = 0;
        int comm_size = 0;
        int comm_rank = 0;

        MPI_Get_version(&version, &subversion);
        fprintf(stdout, "MPI version: %d.%d\n", version, subversion);

        num_proc = _num_pes();
        my_proc = _my_pe();

        fprintf(stdout, "PE#%d of %d\n", my_proc, num_proc);

        MPI_Comm_size(MPI_COMM_WORLD, &comm_size);
        MPI_Comm_rank(MPI_COMM_WORLD, &comm_rank);

        fprintf(stdout, "Comm rank#%d of %d\n", comm_rank, comm_size);
    }

    return 0;
}

```

## 6.2.4 HPC-X™ OpenSHMEM Tunable Parameters

HPC-X™ OpenSHMEM uses Modular Component Architecture (MCA) parameters to provide a way to tune your runtime environment. Each parameter corresponds to a specific function. The following are parameters that you can change their values to change the application's the function:

- memheap - controls memory allocation policy and thresholds
- scoll - controls HPC-X OpenSHMEM collective API threshold and algorithms
- spml - controls HPC-X OpenSHMEM point-to-point transport logic and thresholds
- atomic - controls HPC-X OpenSHMEM atomic operations logic and thresholds
- shmem - controls general HPC-X OpenSHMEM API behavior

➤ **To display HPC-X OpenSHMEM parameters:**

1. Print all available parameters. Run:

```
% oshmem_info -a
```

2. Print HPC-X OpenSHMEM specific parameters. Run:

```

% oshmem_info --param shmem all
% oshmem_info --param memheap all
% oshmem_info --param scoll all
% oshmem_info --param spml all
% oshmem_info --param atomic all

```



### 6.2.4.1 OpenSHMEM MCA Parameters for Symmetric Heap Allocation

SHMEM memheap size can be modified by adding the `SHMEM_SYMMETRIC_HEAP_SIZE` parameter to the `oshrun` file. The default heap size is 256M.

➤ **To run SHMEM with memheap size of 64M. Run:**

```
% oshrun -x SHMEM_SYMMETRIC_HEAP_SIZE=64M -np 512 -mca mpi_paffinity_alone 1 -bynode -
display-map -hostfile myhostfile example.exe
```

Memheap can be allocated with the following methods:

- `sysv` - system V shared memory API. Allocation with hugepages is currently not supported
- `verbs` - IB verbs allocator is used
- `mmap` - `mmap()` is used to allocate memory

By default HPC-X OpenSHMEM will try to find the best possible allocator. The priority is `verbs`, `sysv` and `mmap`. It is possible to choose a specific memheap allocation method by running `-mca sshmem <name>`

### 6.2.4.2 Parameters Used to Force Connection Creation

Commonly SHMEM creates connection between PE lazily. That is at the sign of the first traffic.

➤ **To force connection creating during startup:**

- Set the following MCA parameter.

```
-mca shmem_preconnect_all 1
```

Memory registration (ex: infiniband rkeys) information is exchanged between ranks during startup.

➤ **To enable on demand memory key exchange:**

- Set the following MCA parameter.

```
-mca shmalloc_use_modex 0
```

## 6.3 Tuning OpenSHMEM Atomics Performance

HPC-X OpenSHMEM uses separate communication channel to perform atomic operations. By default this channel is disabled and uses RC transport.



When running on Connect-IB® adapter cards, it is recommended to use DC transport instead of RC.

Atomic tunable parameters:

- `-mca spml_ikrit_hw_rdma_channle 0|1` - default is 1 (enabled)
- `MXM_OSHMEM_HW_RDMA_TLS=rc|dc` - Decides what transport is used for atomic operations. Default is rc

## 6.4 Tuning MTU Size to the Recommended Value



The procedures described below apply to user using MLNX\_OFED 1.5.3.-3.0.0 only.

When using MLNX\_OFED 1.5.3-3.0.0, it is recommended to change the MTU to 4k. Whereas in MLNX\_OFED 3.1-x.x.x and above, the MTU is already set by default to 4k.

➤ **To check the current MTU support of an InfiniBand port, use the *smpquery* tool:**

```
# smpquery -D PortInfo 0 1 | grep -i mtu
```

If the MtuCap value is lower than 4K, enable it to 4K.

Assuming the firmware is configured to support 4K MTU, the actual MTU capability is further limited by the mlx4 driver parameter.

➤ **To further tune it:**

1. Set the `set_4k_mtu` mlx4 driver parameter to 1 on all the cluster machines. For instance:

```
# echo "options mlx4_core set_4k_mtu=1" >> /etc/modprobe.d/mofed.conf
```

2. Restart `openibd`.

```
# service openibd restart
```

➤ **To check whether the parameter was accepted, run:**

```
# cat /sys/module/mlx4_core/parameters/set_4k_mtu
```

To check whether the port was brought up with 4K MTU this time, use the *smpquery* tool again.

## 6.4.1 HPC Applications on Intel Sandy Bridge Machines

Intel Sandy Bridge machines have NUMA hardware related limitation which affects performance of HPC jobs utilizing all node sockets. When installing MLNX\_OFED 3.1-x.x.x, an automatic workaround is activated upon Sandy Bridge machine detection, and the following message is printed in the job's standard output device: "mlx4: Sandy Bridge CPU was detected"

➤ **To disable MLNX\_OFED 3.1-x.x.x Sandy Bridge NUMA related workaround:**

- Set the SHELL environment variable before launching HPC application. Run:

```
% export MLX4_STALL_CQ_POLL=0
% oshrun <...>
```

or

```
oshrun -x MLX4_STALL_CQ_POLL=0 <other params>
```

## 7 Unified Parallel C Overview

Unified Parallel C (UPC) is an extension of the C programming language designed for high performance computing on large-scale parallel machines. The language provides a uniform programming model for both shared and distributed memory hardware. The programmer is presented with a single shared, partitioned address space, where variables may be directly read and written by any processor, but each variable is physically associated with a single processor. UPC uses a Single Program Multiple Data (SPMD) model of computation in which the amount of parallelism is fixed at program startup time, typically with a single thread of execution per processor.

In order to express parallelism, UPC extends ISO C 99 with the following constructs:

- An explicitly parallel execution model
- A shared address space
- Synchronization primitives and a memory consistency model
- Memory management primitives

The UPC language evolved from experiences with three other earlier languages that proposed parallel extensions to ISO C 99: AC, Split-C, and Parallel C Preprocessor (PCP). UPC is not a superset of these three languages, but rather an attempt to distill the best characteristics of each. UPC combines the programmability advantages of the shared memory programming paradigm and the control over data layout and performance of the message passing programming paradigm.

HPC-X UPC is based on Berkely UPC package (see <http://upc.lbl.gov/>) and contains the following enhancements:

- GasNet library used within UPC integrated with Mellanox FCA which off-loads from UPC collective operations.
- GasNet library contains MXM conduit which offloads from UPC all P2P operations as well as some synchronization routines.

### 7.1 Compiling and Running the UPC Application

#### 7.1.1 Compiling the UPC Application

➤ *To build the UPC application:*

- Use the `upcc` compiler.

```
$ upcc -o app app.c
```

➤ *To build the application with a debug info in UPC and GASNet:*

```
$ upcc -g -o app app.c
```

For further information on additional build options, please refer to the `upcc` man page.

#### 7.1.2 Running the UPC Application

The UPC application can be run using the UPC execution command.

```
$ upcrun -shared-heap=<size_per_process> \
-n <total_number_of_processes> \
-N <machines_number_to_be_used> \
-bind-threads <executable>
```

### 7.1.2.1 Basic upcrun Options

Each UPC process uses shared heap. The exact size of the required shared heap is application-specific. The amount of shared memory per UPC thread is controlled by the `-shared-heap` parameter or by the `UPC_SHARED_HEAP_SIZE` environment variable.

A hard limit (per UNIX process) of the shared memory heap can be set using the `-shared-heap-max` parameter or the `UPC_SHARED_HEAP_SIZE` environment variable. This constitutes an upper limit on the `-shared-heap` parameter. Setting this value too high can lead to long application startup times or memory exhaustion on some systems.

For optimal performance, you should bind (a.k.a pin) the UPC threads to the processors using the `-bind-threads` option or using the `UPC_BIND_THREADS` environment variable. These parameters are silently ignored on unsupported platforms.

The following is an example of running a binary with PPN=8 on 8 hosts.

```
$ upcrun -shared-heap=256M -n 64 -N 8 -bind-threads <my_application>
```

For further information on additional usage options, please refer to the `upcrun` man pages or go to <http://upc.lbl.gov/docs/user/upcrun.html>

### 7.1.2.2 Environment Variables

Any command line argument has its environment variable equivalent which is required as the BUPC supports direct execution using the schedulers.

Any environment variable that begins with `UPC_` or `GASNET_` is automatically propagated to all the nodes.

However, there is no Open MPI's equivalent of passing any environment variable to all the nodes, therefore, in order to pass an arbitrary environment variable to the BUPC or to your application, this variable has to be present in the environment where the UPC application is executed.

For example, add your environment variables to `~/.bashrc` or `~/.bash_profile`

## 7.2 FCA Runtime Parameters

The following parameters can be passed to “`upcrun`” in order to change FCA support behavior:

**Table 4 - Runtime Parameters**

Parameter	Description
<code>-fca_enable &lt;0 1&gt;</code>	Disables/Enables FCA support at runtime (default: disable).
<code>-fca_np &lt;value&gt;</code>	Enables FCA support for collective operations if the number of processes in the job is greater than the <code>fca_np</code> value (default: 64).
<code>-fca_verbose &lt;level&gt;</code>	Sets verbosity level for the FCA modules

**Table 4 - Runtime Parameters**

Parameter	Description
<code>-fca_ops &lt;+/-&gt;[op_list]</code>	<p><code>op_list</code> - comma separated list of collective operations.</p> <ul style="list-style-type: none"> <li><code>-fca_ops &lt;+/-&gt;[op_list]</code> - Enables/disables only the specified operations</li> <li><code>-fca_ops &lt;+/-&gt;</code> - Enables/disables all operations</li> </ul> <p>By default all operations are enabled. Allowed operation names are: barrier (br), bcast (bt), reduce (rc), allgather (ag). Each operation can be also enabled/disabled via environment variable:</p> <ul style="list-style-type: none"> <li><code>GASNET_FCA_ENABLE_BARRIER</code></li> <li><code>GASNET_FCA_ENABLE_BCAST</code>,</li> <li><code>GASNET_FCA_ENABLE_REDUCE</code>,</li> </ul> <p><b>Note:</b> All the operations are enabled by default.</p>

### 7.2.1 Enabling FCA Operations through Environment Variables in HPC-X UPC

This method can be used to control UPC FCA offload from environment using job scheduler `srn` utility. The valid values are: 1 - enable, 0 - disable.

➤ *To enable a specific operation with shell environment variables in HPC-X UPC:*

```
% export GASNET_FCA_ENABLE_BARRIER=1
% export GASNET_FCA_ENABLE_BCAST=1
% export GASNET_FCA_ENABLE_REDUCE=1
```

### 7.2.2 Controlling FCA Offload in HPC-X UPC using Environment Variables

➤ *To enable FCA module under HPC-X UPC:*

```
% export GASNET_FCA_ENABLE_CMD_LINE=1
```

➤ *To set FCA verbose level:*

```
% export GASNET_FCA_VERBOSE_CMD_LINE=10
```

➤ *To set the minimal number of processes threshold to activate FCA:*

```
% export GASNET_FCA_NP_CMD_LINE=1
```



HPC-X UPC contains modules configuration file (<http://modules.sf.net>) which can be found at `/opt/mellanox/bupc/2.2/etc/bupc_modulefile`.

## 7.3 Various Executable Examples

The following are various executable examples.

➤ *To run a HPC-X UPC application without FCA support:*

```
% upcrun -np 128 -fca_enable 0 <executable filename>
```

➤ *To run HPC-X UPC applications with FCA enabled for any number of processes:*

```
% export GASNET_FCA_ENABLE_CMD_LINE=1 GASNET_FCA_NP_CMD_LINE=0
% upcrun -np 64 <executable filename>
```

- **To run HPC-X UPC application on 128 processes, verbose mode:**

```
% upcrun -np 128 -fca_enable 1 -fca_np 10 -fca_verbose 5 <executable filename>
```

- **To run HPC-X UPC application, offload to FCA Barrier and Broadcast only:**

```
% upcrun -np 128 -fca_ops +barrier,bt <executable filename>
```



BUPC provided with HPC-X is not re-locatable. In order to use BUPC from a desired location, run post install procedure from HPCX BUPC source tarball:  
bupc/contrib/relocate\_bupc.