

OFA Interoperability Working Group

OFA-IWG Interoperability Test Plan Release 1.45



October 09, 2012
DRAFT

Copyright © 2012 by OpenFabrics - All rights reserved.

This document contains information proprietary to OpenFabrics. Use or disclosure without written permission from an officer of the OpenFabrics is prohibited.

OpenFabrics.org

Revision History

Revision	Release Date	
0.50	Apr 4, 2006	• First FrameMaker Draft of the Interop Test Plan which was used in the March 2006 IBTA-OpenFabrics Plugfest.
0.51	Apr 25, 2006	• Added DAPL and updated MPI.
0.511	June 1, 2006	• Arkady Added iWARP.
0.52	May 30, 2006	• Added Intel MPI.
0.53	June 6, 2006	• Updated uDAPL section provided by Arkady.
0.54	June 13, 2006	• Updated entire Test Spec based on changes made by Arkady to incorporate iWARP into the Test Spec.
0.80	June 14, 2006	• Updated for the OFA conference in Paris and for BoD meeting. Added OFA logo and URL.
1.0	June 21, 2006	• Released after review and approval at the OFA conference in Paris.
1.01	Aug 17, 2006	• Updated the iWARP Equipment requirements in the General System Setup section.
1.02	Oct 31, 2006	<ul style="list-style-type: none"> • Updated Table 4 for iSER, Table 5 for SRP, Table 10 for uDAPL and corresponding info in Tables 17,18 and 22 as per request by Arkady. • Added new test section from Bob Jaworski for Fibre Channel Gateway.
1.03	Dec 10, 2006	<ul style="list-style-type: none"> • Updated test procedures based on the October 2006 OFA Interop Event. • Updated Fibre Channel Gateway test based on changes submitted by Karun Sharma (QLogic). • Added Ethernet Gateway test written by Karun Sharma (QLogic).
1.04	Mar 6, 2007	• Updated test procedures in preparation for the April 2007 OFA Interop Event
1.05	Mar 7, 2007	• Updated iWARP test procedures based on review by Mikkel Hagen of UNH-IOL. Added missing results tables.
1.06	April 3, 2007	• Updated for April 2007 Interop Event based on review from OFA IWG Meeting on 3/27/07.
1.07	April 3, 2007	• Updated for April 2007 Interop Event based on review from OFA IWG Meeting on 4/3/07
1.08	April 4, 2007	• Added list of Mandatory Tests for April 2007 Interop Event.
1.09	April 9, 2007	• Updated Intel MPI based on review by Arlin Davis.
1.10	April 10, 2007	• Updated after final review by Arlin Davis and after the OFA IWG meeting on 4/10/2007

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

Revision	Release Date	
1.11	Sep 7, 2007	<ul style="list-style-type: none">Updated with the latest scripts developed by UNH IOL and based on the results from the April 2007 Interop Event
1.12	Sep 12, 2007	<ul style="list-style-type: none">Updated the documents to embed the test scripts in the document.
1.13	Jan 22, 2008	<ul style="list-style-type: none">Updated the documents for the March 2008 OFA Interop event. IPoIB updated along with Cover Page and the Test Requirements section.
1.14	Feb 11, 2008	<ul style="list-style-type: none">Added the following tests:<ul style="list-style-type: none">1. Ethernet Switch Tests2. IPoIB Connected Mode3. RDMA Interop4. RDS
1.15	Feb 18, 2008	<ul style="list-style-type: none">Updates to the following tests:<ul style="list-style-type: none">1. Ethernet Switch Tests2. IPoIB Connected Mode3. RDMA Interop
1.16	Feb 25, 2008	<ul style="list-style-type: none">Removed all reference to Low Latency Ethernet Switches. This is the version for the March 2008 Interop Event
1.17	March 3, 2008	<ul style="list-style-type: none">Added HP-MPI
1.18	July 22, 2008	<ul style="list-style-type: none">Updated HP-MPI based on results from the March 2008 Interop Event
1.19	July 28, 2008	<ul style="list-style-type: none">Updated HP-MPI URL for the tests.Added section for Open MPIUpdated MPI based on feedback from UNH IOL
1.20	July 30, 2008	<ul style="list-style-type: none">Updated section for Open MPI and added tablesUpdated IB SM Failover as per Nick Wood
1.21	Aug 1, 2008	<ul style="list-style-type: none">Updated SRP call <code>srp_daemon -o -e -n</code>Updated IB SM Failover as Bob JaworskiUpdated HP-MPIUpdated Intel MPIUpdated Open MPI
1.22	Aug 29, 2008	<ul style="list-style-type: none">Added a section for MVAPICH 1 under OSU MPI
1.23	Feb 16, 2009	<ul style="list-style-type: none">Updated Link Init, Fabric Init, SRP, SDP, IPoIB CM, IPoIB DM based on updates received from UNH-IOL

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

Revision	Release Date		
1.24	Feb 23, 2009	<ul style="list-style-type: none"> Updated Intel MPI and Open MPI to reflect the fact that they are not intended to work in a heterogeneous environment. Updated the RDS test procedure Updated the Test Glossary Updated the Mandatory test table for April 2009 	1 2 3 4 5 6 7 8
1.25	Feb 24, 2009	<ul style="list-style-type: none"> Updated the RDS Test after review by the OFA IWG group. 	9 10
1.26	Mar 13, 2009	<ul style="list-style-type: none"> Restructured entire document to accommodate WinOF and OFED Added NFS over RDMA to the test plan. Added WinOF tests Updated HP-MPI Add List of Contributors 	11 12 13 14 15 16
1.27	Mar 17, 2009	<ul style="list-style-type: none"> Updates based on the review from the OFA IWG 	17
1.28	Mar 27, 2009	<ul style="list-style-type: none"> Added links in Chapter 10 to the InfiniBand Test Scripts Added links to HP-MPI installation Packages 	18 19
1.29	Aug 25, 2009	<ul style="list-style-type: none"> Editorial & Technical updates based on April 2009 Interop Event. Updated Mandatory tests for October 2009. Added Topology Check Added new Firmware Policy 	20 21 22 23 24
1.30	Sep 4, 2009	<ul style="list-style-type: none"> Updated Mandatory iWARP tests and several comments based on the review from Harry Cropper Added changes suggested by Jess Robel from QLogic to IPoIB DM and CM and Fabric Init. 	25 26 27 28
1.31	April 6, 2010	<ul style="list-style-type: none"> Added definition of homogenous to Test Glossary Added updates from the November 2009 Interop Event 	29 30
1.32	April 20, 2010	<ul style="list-style-type: none"> Updated after the OFA IWG meeting on 4/6/2010 Updated MPI and MVAPICH based on changes requested by Jeff Laird and Intel 	31 32 33
1.33	April 23, 2010	<ul style="list-style-type: none"> Major changes to Section 8 which describes the Software and Firmware polices 	34 35
1.34	July 20, 2010	<ul style="list-style-type: none"> Changed uDAPL for iWARP to Beta for Aug 2010 GA Event Removed HP MPI which is no longer supported Added -mca mpi_leave_pinned 0 for OpenMPI Add new parameters for MVAPICH2 for iWARP devices. 	36 37 38 39 40

Revision	Release Date	
1.35	July 27, 2010	<ul style="list-style-type: none">Added new parameters for MVAPICH2 for iWARP devices. The parameter is: MV2_USE_RDMA_CM=1
1.36	Feb 22, 2011	<ul style="list-style-type: none">Added Link Init section as per changes provided by Chris Hutchins and approved by OFA IWG.Updated Test Plan Status for April 2011 and October 2011Nick Wood from UNH-IOL updated NFSoRDMAMarty requested that we update SRP Results Table 6 and remove the disconnect commands.
1.37	Oct 4, 2011	<ul style="list-style-type: none">Updated Test Plan Status for November 2011Added new Test Table for OS and OFED versionsNick Wood updated Link Init for IBChris Hutchins updated RDMA Interop and RDMA StressRemoved XANSation testing
1.38	Oct 11, 2011	<ul style="list-style-type: none">Changed Link Init Section from Recommendation to MOIUpdated Section 8 for Firmware, Software and Hardware Policies to bring in line with Logo Program DocumentUpdated InfiniBand Test Table 24
1.39	Oct 24, 2011	<ul style="list-style-type: none">Updated Open MPI as per changes submitted by Nick WoodUpdated RDMA Interop small test: drop iterations from 100000 to 25000Updated RDMA Interop large test, increase iterations from 100 to 300Updated IPoIB Part A:, drop iterations (number of pings) from 100 to 10.
1.40	Oct 25, 2011	<ul style="list-style-type: none">Modified the following sections12.6.9 iwarp client 100000 -> 2500012.6.13olarge read client 65536 -> 1000000olarge write client 65536 -> 1000000Added large send command (section c)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

Revision	Release Date		
1.41	Mar 20, 2012	<ul style="list-style-type: none"> • General Instructions: Added note that the OpenSM will be used to run all mandatory tests in the test plan and the Vendor SM testing will include testing IPoIB, RDMA Interop and Open MPI testing. • General Instructions: The OFILG decided as of April 2012 that the various ULPs contained in this test plan will only be tested if it is supported by the Operating System. • Logo Program Requirements: updated IB and iWARP. Made NFSoRDMA Mandatory and MVAPICH Optional. • IPoIB: Modified the way IPoIB is set to connected or datagram mode • IPoIB: Changed the ping interval in IPoIB tests from 0.01 to 0.2 • IPoIB: Reduced number of frame sizes tested in the Ping Test. • MVAPICH: Made testing of MVAPICH 1 & 2 Optional • NFSoRDMA: Eliminate the need to specify nfs-utils in the NFSoRDMA installation section • NFSoRDMA: Changed the way the servers are mounted in NFSoRDMA • SDP: Eliminated the need for vsftpd in SDP • SDP: Eliminated the environment variables section in SDP • SDP: Changed the way the netperf server is started in SDP • SDP: Made SDP mandatory only for those Operating Systems that support it. • SRP: Mandated that Targets only advertise two volumes in order to reduce the amount of time required to run the tests 	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27
1.42	Apr 3, 2012	<ul style="list-style-type: none"> • Updated Ethernet Test requirements to move NFSoRDMA to Beta for April 2012 • Changed the status of Intel MPI and OSU MVAPICH to deprecated meaning the tests are no longer being run or supported. • Updated SRP notes as per Marty Schlining 	28 29 30 31 32
1.43	Aug 14, 2012	<ul style="list-style-type: none"> • Updated the definition for \$NP in MVAPICH section 12.10.2, 2, ii • Updated Mandatory test tables for iWARP and IB • Cleared all change bars for October 2012 Interop event 	33 34 35 36

Revision	Release Date	
1.44	Sep 18, 2012	<ul style="list-style-type: none">Removed Intel MPI because it is not Open SourceRemoved SDP because no longer supported in OFEDRemoved Ethernet Fabric Initialize, Failover and reconvergence. No longer applicable given DCB etc.Removed TI RDS for iWARP because RDS does not support iWARPRemove iWARP Connectivity - replaced by RDMA Interop test sectionAdded section 8 for OS Installation and OS Policy
1.45	Oct 9, 2012	<ul style="list-style-type: none">Add second test of SRPAdd RoCE test sections

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

List of Contributors

Editor: Rupert Dance

Name	Company
Mark Alan	HP
Harry Cropper	Intel
Rupert Dance	Software Forge
Sujal Das	Mellanox
Arlin Davis	Intel
Johann George	QLogic
Mike Hagen	UNH-IOL
Mitko Haralanov	QLogic
Allen Hubbe	UNH-IOL
Christopher Hutchins	UNH-IOL
Bob Jaworski	QLogic
Arkady Kanevsky	NetApp
Llolsten Kaonga	Software Forge
Amit Krig	Mellanox
Jeff Laird	UNH-IOL
Jon Mason	Open Grid Computing
Edward Mossman	UNH-IOL
Bob Noseworthy	UNH-IOL
Yaroslav Pekelis	Mellanox
Jess Robel	Qlogic
Hal Rosenstock	HNR Consulting
Martin Schlining	DataDirect Networks
Karun Sharma	QLogic
Stan Smith	Intel
Dave Sommers	Intel (NetEffect)
Jeff Squyres	Cisco
Dennis Tolstenko	Lamprey Networks
Steve Wise	Open Grid Computing
Robert Woodruff	Intel

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

Name	Company
Nick Wood	UNH-IOL

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

LEGAL DISCLAIMER

"This version of a proposed OpenFabrics Interop Test Plan is provided "AS IS" and without any warranty of any kind, including, without limitation, any express or implied warranty of non-infringement, merchantability or fitness for a particular purpose.

In no event shall OpenFabrics, IBTA or any member of these groups be liable for any direct, indirect, special, exemplary, punitive, or consequential damages, including, without limitation, lost profits, even if advised of the possibility of such damages."

Conditional text tag *Explanation* is shown in green.

~~Conditional text tag Deleted is shown in red with strike through.~~

Conditional text tag *Proposal* is shown in turquoise (r0_g128_b128).

Conditional text tag *Author* is shown as is.

Conditional text tag Comment is shown in red with underline

1 INTRODUCTION

Server OEM customers have expressed the need for RDMA hardware and software to interoperate.

Specifically, InfiniBand HCA, OpenFabrics host software to interoperate with InfiniBand Switches, gateways, and bridges with management software provided by OEMs, and IB integrated server OEM vendors. And, iWARP RNIC and OpenFabrics host software to interoperate with Ethernet Switches and management software and hardware provided by Ethernet Switch OEMs and iWARP integrated server OEM vendors.

It is necessary that the interoperability test effort be an industry-wide effort where interoperability testing is conducted under the auspices of the appropriate networking organizations. For InfiniBand it is the IBTA, specifically within the charter of the CIWG and for iWARP it is the IETF.

1.1 PURPOSE

This document is intended to describe the production tests step by step explaining each test and its references. The purpose of this test plan is three fold:

- 1) Define the scope, equipment and software needs, and test procedures for verifying full interoperability of RDMA HW and SW. For Infiniband HW it is InfiniBand HCAs using the latest OpenFabrics OFED software with currently available OEM Switches and their management software. The target OEM IB Switch vendors are Intel and Mellanox. For iWARP HW it is iWARP RNICs using the latest OpenFabrics OFED software with currently available OEM Ethernet Switches, Bridges, Gateways, Edge Devices and so on with their management software.
- 2) Serve as a basis for evaluating customer acceptance criteria for OFA host software interoperability and OFA Logo.
- 3) Serve as a basis for extensions to InfiniBand IBTA CIWG test procedures related to interoperability and use of these test procedures in upcoming PlugFest events organized by IBTA.
Serve as a basis for extensions to iWARP test procedures for OpenFabrics software related to interoperability and use of these test procedures in upcoming PlugFest events organized by the UNH IOL OFILG testing service.

1.2 INTENDED AUDIENCE

The following are the intended audience for this document:

- 1) Project managers in OEM Switch, Router, Gateway, Bridge Vendor companies to understand the scope of testing and participate in the extension of this test plan and procedures as necessary to meet their requirements.
- 2) IBTA and CIWG, and iWARP and UNH IOL iWARP testing personnel and companies to evaluate the scope of testing and participate in the extension of this test plan and procedures as necessary to meet their requirements.
- 3) Test engineering and project leads and managers who will conduct the testing based on this document.

- 4) Customers and users of OFA host software who rely on OFA Logo for interoperability. 1
- 5) Integrators and OEM of RDMA products. 2

1.3 TEST PLAN STRUCTURE 3

This test plan is divided into two main sections. 4

- 1) Interoperability testing using OFED for Linux. 5
- a) See Sections 10-12 6
- 2) Interoperability testing using WinOFED for Windows Platforms. 7
- a) See Section 13 8

Sections 1.4 through 1.10 provide an overview of the tests which are described in detail in sections 10 through 13. 9

1.4 INFINIBAND ONLY - TEST OVERVIEW

The tables below list all of the specific test procedures for InfiniBand Devices. See the Transport Independent section for tests that apply to all transports.

Table 1 - IB Link Initialize

Test #	Test	Description
1	Phy link up all ports	Check that all relevant LEDs are on for all HCAs and switches.

Table 2 - IB Fabric Initialization

Test #	Test	Description
1	Fabric Initialization	Run SM from each node in cluster and see that all ports are in Armed or Active state.

Table 3 - IB IPoIB - Connect Mode (CM)

Test #	Test	Description
1	Ping all to all	Run SM from one of the nodes and check all nodes responding. Repeat with all SMs.
2	Connect disconnect host	Run SM from one of the nodes and check all nodes responding.
3	FTP Procedure	Using a 4MB test file, put the file, then get the file and finally compare the file.

Table 4 - IB IPoIB - Datagram Mode (DM)

Test #	Test	Description
1	Ping all to all	Run SM from one of the nodes and check all nodes responding. Repeat with all SMs.
2	Connect disconnect host	Run SM from one of the nodes and check all nodes responding.
3	FTP Procedure	Using a 4MB test file, put the file, then get the file and finally compare the file.

Table 5 - IB SM Tests

Test #	Test	Description
1	Basic sweep test	verify that all SMs are NOT ACTIVE (after receiving the SMSet of SMInfo to DISABLE) and that the selected SM (SM1) is the master (

Table 5 - IB SM Tests

Test #	Test	Description
2	SM Priority test	Verify Subnet and SMs behavior according to the SMs priority.
3	Failover - Disable SM1	Disable the master SM and verify that standby SM becomes master and configures the cluster.
4	Failover - Disable SM2	Disable the master SM and verify that standby SM becomes master and configures the cluster.

Table 6 - IB SRP Tests

Test #	Test	Description
1	Basic dd application	Run basic dd application from SRP host connected to target.
2	IB SM kill	Kill the IB master SM while test is running and check that it completes properly.
3	Disconnect Host	Unload SRP Host and check SRP connection properly disconnected.
4	Disconnect Target	Unload SRP Target and check SRP connection properly disconnected.

Table 7 - IB Ethernet Gateway

Test #	Test	Description
1	Basic Setup	Connect the HCA of the IB host and Ethernet Gateway to the IB fabric. Connect the Ethernet gateway to the Ethernet network or Ethernet device. Start the SM to be used in this test.
2	Start ULP	Determine which ULP your ethernet gateway uses and be sure that ULP is running on the host.
3	Discover Gateway	Restart the ULP or using the tool provided by the ULP, make sure that the host "discovers" the Ethernet Gateway.
4	SM Failover	While the ping is running, kill the master SM. Verify that the ping data transfer is unaffected.
5	Ethernet gateway reboot	Reboot the Ethernet Gateway. After the Ethernet Gateway comes up, verify that the host can discover the Ethernet Gateway as it did before and we are able to configure the interfaces.
6	ULP restart	Restart the ULP used by Ethernet Gateway and verify that after the ULP comes up, the host can discover the Ethernet Gateway and we are able to configure the interfaces.
7	Unload/load ULP	Unload the ULP used by Ethernet Gateway and check that the Ethernet Gateway shows it disconnected. Load the ULP and verify that the Ethernet gateway shows the connection.

Table 8 - IB Fibre Channel Gateway

Test #	Test	Description
1	Basic Setup	Connect the HCA of the IB host to the IB fabric. Connect the FC Gateway to the IB Fabric. Connect the FC Gateway to the FC network or FC device. Start the SM to be used in this test.
2	Configure Gateway	Configure the FC Gateway appropriately (how to do this is vendor specific).

Table 8 - IB Fibre Channel Gateway

Test #	Test	Description
3	Add Storage Device	Use ibsrpdm tool in order to have the host "see" the FC storage device. Add the storage device as target.
4	Basic dd application	Run basic dd application from SRP host connected to target.
5	IB SM kill	Kill the IB master SM while test is running and check that it completes properly.
6	Disconnect Host/Target	Unload the SRP host / SRP Target (target first/host first) and check that the SRP connection is properly disconnected.
7	Load Host/Target	Load the SRP host / SRP Target. Using ibsrpdm, add the target.
8	dd after SRP Host and Target reloaded	Run basic dd application from the SRP host to the FC storage device.
9	Reboot Gateway	Reboot the FC Gateway. After FC Gateway comes up, verify using ibsrpdm tool that the host see the FC storage device. Add the storage device as target.
10	dd after FC Gateway reboot	Verify basic dd works after rebooting Gateway.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

1.5 ETHERNET ONLY - TEST OVERVIEW

The tables below list all of the specific test procedures for iWARP and Ethernet Devices. See the Transport Independent section for tests that apply to all transports.

Table 9 - iWARP Link Initialize

Test #	Test	Description
1	Phy link up all ports	Check that all relevant green LEDs are on for all RN ICs and switches.
2	Verify basic IP connectivity	Verify IP and RDMA connectivity can occur by driving minimum size ICMP echo requests and replies across the link or equivalent traffic.

Table 10 - RoCE Link Initialize

Test #	Test	Description
1	Phy link up all ports	Check that all relevant green LEDs are on for all RCAs and switches.
2	Verify basic IP connectivity	Verify IP and RDMA connectivity can occur by driving minimum size ICMP echo requests and replies across the link or equivalent traffic.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

1.6 TRANSPORT INDEPENDENT - TEST OVERVIEW

The tables below list the test procedures that apply to devices regardless of the transport.

Table 11 - TI iSER

Test #	Test	Description
1	Basic dd application	Run basic dd application from iSER host connected to target.
2	IB SM kill	[IB Specific] - Kill the IB master SM while test is running and check that it completes properly.
3	Disconnect Initiator	Unload iSER Host and check iSER connection properly disconnected.
4	Disconnect Target	Unload iSER Target and check iSER connection properly disconnected.
5	Repeat with previous SM Slave	[IB Specific Test] Repeat steps 1-4 now with the previous slave SM (we did not actually stop the target).

Table 12 - TI NFS Over RDMA

Test #	Test	Description
1	File and directory creation	A total of six files and six directories are created
2	File and directory removal	removes the directory tree that was just created by test1
3	Lookups across mount point	changes directory to the test directory and gets the file status of the working directory
4	Setattr, getattr, and lookup	Permissions are changed (chmod) and the file status is retrieved (stat) for each file
5	Read and write	Creates a file (creat), Gets status of file (fstat) , Checks size of file, Writes 1048576 bytes into the file (write) in 8192 byte buffers, Closes file (close), Gets status of file (stat) , Checks the size of the file
6	Readdir	The program creates 200 files (creat). The current directory is opened (opendir), the beginning is found (rewinddir), and the directory is read (readdir) in a loop until the end is found
7	Link and rename	This program creates ten files. For each of these files, the file is renamed (rename) and file statistics are retrieved (stat) for both the new and old names
8	Symlink and readlink	This program makes 10 symlinks (symlink). It reads (readlink), and gets statistics for (lstat) each, and then removes them (unlink).
9	Statfs	This program changes directory to the test directory (chdir and/or mkdir) and gets the file system status on the current directory (statfs).

Table 13 - TI RDS

Test #	Test	Description
1	rds-ping procedure	Run rds-ping and verify that you can reach all hosts in the cluster

Table 13 - TI RDS

Test #	Test	Description
2	rds-stress procedure	Set up passive receiving instance and an active sender and verify data is exchanged without error

Table 14 - TI uDAPL

Test #	Test	Description
1	Point-to-Point Topology	Connection and simple send receive.
2	Point-to-Point Topology	Verification, polling and scatter gather list.
3	Switched Topology	Verification and private data.
4	Switched Topology	Add multiple endpoints, polling, and scatter gather list.
5	Switched Topology	Add RDMA Write.
6	Switched Topology	Add RDMA Read.
7	Multiple Switches	Multiple threads, RDMA Read, and RDMA Write.
8	Multiple Switches	Pipeline test with RDMA Write and scatter gather list.
9	Multiple Switches	Pipeline with RDMA Read.
10	Multiple Switches	Multiple switches.

Table 15 - RDMA Basic Interop

Test #	Test	Description
1	Small RDMA READ	Create an RDMA command sequence to send a READ operation of one byte.
2	Large RDMA READ	Create an RDMA command sequence to send a READ operation of 10,000,000 bytes
3	Small RDMA Write	Create an RDMA command sequence to send a Write operation of one byte
4	Large RDMA Write	Create an RDMA command sequence to send a Write operation of 10,000,000 bytes
5	Small RDMA SEND	Create an RDMA command sequence to send a SEND operation of one byte.
6	Large RDMA SEND	Create an RDMA command sequence to send a SEND operation of one million bytes
7	Small RDMA Verify	Create an RDMA command sequence to send a VERIFY operation of one byte.
8	Large RDMA Verify	Create an RDMA command sequence to send a VERIFY operation of 10,000,000 bytes

Table 16 - RDMA Stress Tests

Test #	Test	Description
1	Switch Load	For one pair of endpoints generate a stream of RDMA READ operation in one direction and RDMA write operations in the opposite direction. For all remaining endpoint pairs configure an RDMA WRITE operation of 1 byte and have it sent 10000 times on both streams of the endpoint pair.
2	Switch Fan In	Connect all possible endpoint pairs such that data exchanges between pairs must traverse the pair of ports interconnecting the switch

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

1.7 OPEN MPI - TEST OVERVIEW

Table 17 - TI - Open MPI Test Suite Description

Test #	Open MPI TESTs	Open MPI TESTs Suite Description
Phase 1: "Short" tests		
1	2	OMPI built with OpenFabrics support
2	3	OMPI basic functionality (hostname)
3	4.1	Simple MPI functionality (hello_c)
4	4.2	Simple MPI functionality (ring_c)
5	5	Point-to-point benchmark (NetPIPE)
6	6.1.1	Point-to-point benchmark (IMB PingPong multi)
7	6.1.2	Point-to-point benchmark (IMB PingPing multi)
Phase 2: "Long" tests		
8	6.2.1	Point-to-point benchmark (IMB PingPong)
9	6.2.2	Point-to-point benchmark (IMB PingPing)
10	6.2.3	Point-to-point benchmark (IMB Sendrecv)
11	6.2.4	Point-to-point benchmark (IMB Exchange)
12	6.2.5	Collective benchmark (IMB Bcast)
13	6.2.6	Collective benchmark (IMB Allgather)
14	6.2.7	Collective benchmark (IMB Allgatherv)
15	6.2.8	Collective benchmark (IMB Alltoall)
16	6.2.9	Collective benchmark (IMB Reduce)
17	6.2.10	Collective benchmark (IMB Reduce_scatter)
18	6.2.11	Collective benchmark (IMB Allreduce)
19	6.2.12	Collective benchmark (IMB Barrier)
20	6.3.1	I/O benchmark (IMB S_Write_Indv)
21	6.3.2	I/O benchmark (IMB S_IWrite_Indv)
22	6.3.3	I/O benchmark (IMB S_Write_Expl)
23	6.3.4	I/O benchmark (IMB S_IWrite_Expl)
24	6.3.5	I/O benchmark (IMB P_Write_Indv)
25	6.3.6	I/O benchmark (IMB P_IWrite_Indv)
26	6.3.7	I/O benchmark (IMB P_Write_Shared)

Table 17 - TI - Open MPI Test Suite Description

Test #	Open MPI TESTs	Open MPI TESTs Suite Description
27	6.3.8	I/O benchmark (IMB P_IWrite_Shared)
28	6.3.9	I/O benchmark (IMB P_Write_Priv)
29	6.3.10	I/O benchmark (IMB P_IWrite_Priv)
30	6.3.11	I/O benchmark (IMB P_Write_Expl)
31	6.3.12	I/O benchmark (IMB P_IWrite_Expl)
32	6.3.13	I/O benchmark (IMB C_Write_Indv)
33	6.3.14	I/O benchmark (IMB C_IWrite_Indv)
34	6.3.15	I/O benchmark (IMB C_Write_Shared)
35	6.3.16	I/O benchmark (IMB C_IWrite_Shared)
36	6.3.17	I/O benchmark (IMB C_Write_Expl)
37	6.3.18	I/O benchmark (IMB C_IWrite_Expl)
38	6.3.19	I/O benchmark (IMB S_Read_Indv)
39	6.3.20	I/O benchmark (IMB S_IRead_Indv)
40	6.3.21	I/O benchmark (IMB S_Read_Expl)
41	6.3.22	I/O benchmark (IMB S_IRead_Expl)
42	6.3.23	I/O benchmark (IMB P_Read_Indv)
43	6.3.24	I/O benchmark (IMB P_IRead_Indv)
44	6.3.25	I/O benchmark (IMB P_Read_Shared)
45	6.3.26	I/O benchmark (IMB P_IRead_Shared)
46	6.3.27	I/O benchmark (IMB P_Read_Priv)
47	6.3.28	I/O benchmark (IMB P_IRead_Priv)
48	6.3.29	I/O benchmark (IMB P_Read_Expl)
49	6.3.30	I/O benchmark (IMB P_IRead_Expl)
50	6.3.31	I/O benchmark (IMB C_Read_Indv)
51	6.3.32	I/O benchmark (IMB C_IRead_Indv)
52	6.3.33	I/O benchmark (IMB C_Read_Shared)
53	6.3.34	I/O benchmark (IMB C_IRead_Shared)
54	6.3.35	I/O benchmark (IMB C_Read_Expl)
55	6.3.36	I/O benchmark (IMB C_IRead_Expl)
56	6.3.37	I/O benchmark (IMB Open_Close)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

1.8 OSU MPI - TEST OVERVIEW

Table 18 - TI - OSU MPI

Test #	Test	Description
1	Test 1: PingPong	
2	Test 1: PingPing point-to-point	
3	Test 2: PingPong	
4	Test 2: PingPing	
5	Test 2: Sendrecv	
6	Test 2: Exchange	
7	Test 2: Bcast	
8	Test 2: Allgather	
9	Test 2: Allgatherv	
10	Test 2: Alltoall	
11	Test 2: Alltoallv	
12	Test 2: Reduce	
13	Test 2: Reduce_scatter	
14	Test 2: Allreduce	
15	Test 2: Barrier	

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

1.9 REQUIREMENTS FOR OFA INTEROPERABILITY LOGO PROGRAM

The following table indicates the mandatory tests that will be used for Interop Validation during the October 2012 Interop Debug Event and the Interop GA Event using OFED 3.5 GA. Deprecated means that the test is no longer being actively run during the OFA Interop Events.

Table 19 - InfiniBand Transport Test Status for October 2012 Interop Event

Test Procedure	Linux	WinOF
IB Link Initialize	Mandatory	Mandatory
IB Fabric Initialization	Mandatory	Mandatory
IB IPoIB Connected Mode	Mandatory	Not Available - 1
IB IPoIB Datagram Mode	Mandatory	Beta
IB SM Failover/Handover - OpenSM	Mandatory	Beta
IB SM Failover/Handover - Vendor SM	Optional	Optional
IB SRP	Mandatory	Beta
IB Ethernet Gateway	Beta	Not Available - 3
IB Fibre Channel Gateway	Beta	Not Available - 3
TI iSER	Mandatory	Beta
TI NFS over RDMA	Mandatory	Not Available - 1
TI RDS	Mandatory	Not Available - 2
TI uDAPL	Mandatory	Beta
TI Basic RDMA Interop	Mandatory	Not Available - 3
TI RDMA Stress	Mandatory	Not Available - 3
TI MPI Open MPI	Mandatory	Not Available - 2
TI MVAPICH - OSU	Deprecated	Not Available - 2

Not Available means one of three things:

- 1) The feature is not currently supported by the WinOFED stack.
- 2) The ULP application has not been ported to the WinOFED Stack.
- 3) The test has not been updated for WinOFED.

Optional means that this test will not be made mandatory because it depends on proprietary vendor capabilities. The test may be run during the OFA Interop Events and reported in the results but it will not affect eligibility for the OFA Logo List.

Table 20 - iWARP Transport Test Status for October 2012 - OFED 3.5

Test Procedure	Linux
iWARP Link Initialize	Mandatory
TI iSER	Beta
TI NFS over RDMA	Beta
TI uDAPL	Mandatory
TI Basic RDMA Interop	Mandatory
TI RDMA Stress	Mandatory
TI MPI Open MPI	Mandatory
TI MVAICH2 - OSU	Deprecated

Table 21 - RoCE Transport Test Status for October 2012 - OFED 3.5

Test Procedure	Linux
RoCE Link Initialize	Beta
RoCE Fabric Init	TBD
RoCE IPoCE	TBD
RoCE InfiniBand Gateway	TBD
RoCE Fibre Channel Gateway	TBD
TI iSER	Beta
TI NFS over RDMA	Beta
TI uDAPL	Beta
TI Basic RDMA Interop	Beta
TI RDMA Stress	Beta
TI MPI Open MPI	Beta

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

1.10 SUBJECTS NOT COVERED

Table 22 - SUBJECTS NOT COVERED

Number	Subject/ Feature	Reason	Executor	Due Date
1	iWARP peer to peer	Future Testing		TBD
2	IPv6 testing	Future Testing		TBD

1.11 TEST GLOSSARY

Table 23 - Test Glossary

Technical Terms	
DCB	Data Center Bridging (used in RoCE)
HCA	IB Host Channel Adapter
IPoIB	IP over InfiniBand
iSER	iSCSI Extensions for RDMA
MPI	Message Passing Interface
RCA	RoCE Channel Adapter
RDF	Readme File
RDS	Reliable Datagram Sockets
RNIC	RDMA NIC (iWARP Network Interface Card)
RoCE	RDMA over Converged Ethernet
SA	IB Subnet Administration
SDN	Software Defined Network
SDP	Sockets Direct Protocol
SM	IB Subnet Manager
SPB	Shortest Path Bridging (used in RoCE)
SRP	SCSI RDMA Protocol
TD	Test Descriptions
TI	Transport Independent (tests)
TRILL	Transparent Interconnect of Lots of Links is a IETF Standard implemented by devices called RBridges (Routing Bridges) or TRILL Switches (used in RoCE)
uDAPL	User Direct Access Programming Library

1.12 HOMOGENOUS VERSUS HETEROGENEOUS

Heterogeneous & homogeneous clusters are the same with one exception: the end points must be from the same vendor in homogeneous clusters. The table below defines the guidelines for building homogeneous and heterogeneous clusters

Description	Homogenous	Heterogeneous
Mixing switches (both models and vendor products)	Encouraged	Encouraged
The use of any InfiniBand subnet manager	Encouraged	Encouraged
All devices of the same model number shall use the same firmware.	Mandatory	Mandatory
Any mix of products from the same vendor is acceptable - e.g. different model HCAs	Encouraged	Encouraged
A mix of end points (HCA/RNIC) from different OFA vendors	Prohibited	Mandatory
Mixing x86-32 (ix86) and x86_64 Operating System - see notes	Not-Tested	Not-Tested
32 bit architecture and 32 bit OS - see notes	Not-Tested	Not-Tested
Mixing x86-32 and x86-64 user-level application	Optional	Optional
Mixed system architecture - e.g. x86 servers mixed with IA-64 (Itanium) servers	Prohibited	Prohibited
Mixing endianness in system OS	Prohibited	Prohibited
Mixing the quantity of server RAM installed on the hosts	Encouraged	Encouraged
Mixing the server clock speeds	Encouraged	Encouraged
Mixing the number of server cores	Encouraged	Encouraged
Mixing PCIe generations	Encouraged	Encouraged
All servers shall run the same OFED version.	Encouraged	Encouraged
Mixing supported Operating Systems	Encouraged	Encouraged

Notes: Intel drivers do not support 32 bit operating systems

2 USE OF OPENFABRICS SOFTWARE FOR PRE-TESTING

Depending on the schedule of testing and bugs or issues encountered, different snapshots of latest OpenFabrics software will be used during pre-testing prior to the Interoperability Event. Any changes that result in the OpenFabrics software from interoperability testing per this test plan will be deposited back into the OpenFabrics repository so that the OpenFabrics development community will have full access to any bug fixes or feature additions that may result out of this testing effort. The frequency of such deposits will be determined based on completion of adequate testing of the said fixes or feature additions.

3 USE OF OPENFABRICS SOFTWARE FOR IBTA/CIWG COMPLIANCE PLUGFESTS

During the pre-testing phase, UNH-IOL will apply all reasonable effort to ensure that the OpenFabrics source and binary repositories are up-to-date with the latest OFED release. This will enable cable interoperability testing at plugfests to be conducted using software directly sourced from the OpenFabrics tree.

Should there be any issues with the OpenFabrics community not accepting certain bug fixes or features with the time frames matching with Compliance Events, UNH-IOL will inform all participants about the same and offer those bug fixes or features in source code and binary formats directly to the participants and InfiniBand solution suppliers.

4 USE OF OPENFABRICS SOFTWARE FOR OFA IWG INTEROPERABILITY EVENTS

During the pre-testing phase, UNH-IOL will apply all reasonable effort to ensure that the OpenFabrics source and binary repositories are up-to-date with the latest OFED releases chosen by the OFA IWG for use in the Interoperability Event.

Should there be any issues with the OpenFabrics community not accepting certain bug fixes or features with the time frames matching with Interoperability Events, UNH-IOL will inform all participants about the same and offer those bug fixes or features in source code and binary formats directly to the participants and InfiniBand solution suppliers.

5 GENERAL SYSTEM SETUP

Configuration

The test environment for the user interface contains:

5.1 IB HW UNITS

Table 24 - IB Equipment

Equipment	Amount	Details	Check
Servers with OS installed	12 or more	The OS should be supported by OpenFabrics Software.	
4X IB Cables	30 or more	Between 1 meter => 10 meters.	
IB Switches	4	The number and types of switches needed from member companies or OEMs is dependent on variations in subnet management and other IBTA defined management software. For example if the software on Switch A is different from the software used in Switch B, both Switches will be needed. Note that it is not dependent on number of ports supported by a switch.	
IB HCAs	12 or more		

5.2 IB SOFTWARE

5.2.1 LINUX/WINDOWS PLATFORMS

5.2.2 OFED - MOST CURRENT TESTED RELEASE

5.2.3 IB HCA FW – VERSION XXX - VENDOR SPECIFIC

5.2.4 IB SWITCH FW CANDIDATE – VERSION XXX - VENDOR SPECIFIC

5.2.5 IB SWITCH SW – VERSION XXX - VENDOR SPECIFIC

5.3 iWARP HW UNITS

Table 25 - iWARP Equipment

Equipment	Amount	Details	Check
Servers with OS installed	5 or more	The OS should be supported by OpenFabrics Software.	
4X CX4 or SFP Cables	10 or more	Between 1 meter => 10 meters.	
10 GbE Switches	1	At least one 10 GbE switch must be made available to support the various RNICs in the Fabric.. There is no need to have multiple switches if there are enough ports on the primary switches to support all the devices in the fabric.	
iWARP RNIC	5 or more	Each vendor must supply 5 or more RNICs in order to support MPI testing.	

5.4 IWARP SOFTWARE

5.4.1 LINUX PLATFORMS

5.4.2 OFED - MOST CURRENT TESTED RELEASE

5.4.3 IWARP RNIC FW – VERSION XXX - VENDOR SPECIFIC

5.4.4 10GbE SWITCH FW CANDIDATE – VERSION XXX - VENDOR SPECIFIC

5.4.5 10GbE SWITCH SW – VERSION XXX - VENDOR SPECIFIC

5.4.6 VENDOR SPECIFIC NOTES

Note: Currently there is no interoperability between cxgb4 and nes if peer2peer is enabled. Both nes and cxgb4 have their own proprietary ways of doing "client must send the first fpdu". The Chelsio parameter file /sys/module/iw_cxgb4/parameters/peer2peer should be modified on all hosts to contain the appropriate value for each test. For example: the value must be set to '1' for the uDAPL test.

Arlin Davis suggests the following given the current situation:

- 1)The daplttest -T P (performance tests) will always send data from server side first. This test will NOT work reliably with iWARP vendors.
- 2)The daplttest -T T (transaction tests) should work fine with both IB and iWARP vendors given that it always sends from client side first.
- 3)I recommend using only daplttest transaction mode (-T T) in your test plan and removing -T P mode tests.

5.5 ROCE HW UNITS

Table 26 - RoCE Equipment

Equipment	Amount	Details	Check
Servers with OS installed	5 or more	The OS should be supported by OpenFabrics Software.	
4X QSFP+ Cables	10 or more	Between 1 meter => 10 meters.	
GbE DCB Switches	1	At least one 10 or 40 GbE DCB switch must be made available to support the various RCAs in the Fabric. There is no need to have multiple switches if there are enough ports on the primary switches to support all the devices in the fabric.	
RoCE RCA	5 or more	Each vendor must supply 5 or more RCAs in order to support MPI testing.	

5.6 ROCE SOFTWARE

5.6.1 LINUX PLATFORMS

5.6.2 OFED - MOST CURRENT TESTED RELEASE

5.6.3 ROCE FW – VERSION XXX - VENDOR SPECIFIC

5.6.4 10/40 GbE DCB SWITCH FW CANDIDATE – VERSION XXX - VENDOR SPECIFIC

5.6.5 10/40 GbE DCB SWITCH SW – VERSION XXX - VENDOR SPECIFIC

5.7 MPI TESTING

- 1)HCA/RCA/RNIC vendors must provide a minimum of five adapters. The adapters need not be all the same model, but they can be.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

6 IB HW DESCRIPTION & CONNECTIVITY

The test contains two major parts. This description is for each of those parts.

6.1 BASIC CONNECTIVITY (P1P1)

6.1.1 HCA 1 SHOULD BE CONNECTED FROM PORT 1 TO LOWEST PORT NUMBER IN SWITCH

6.1.2 HCA 2 SHOULD BE CONNECTED FROM PORT 1 TO HIGHEST PORT NUMBER IN SWITCH

6.1.3 BOTH WITH COMPLIANT INFINIBAND CABLES

6.2 SWITCHES AND SOFTWARE NEEDED

6.2.1 SWITCHES PROVIDED BY OEMS

It is necessary that Switches provided by OEMs cover the full breadth of software versions supported by the Switch OEMs. Port count is not critical for the tests. It is recommended that OEMs provide six switches covering all variations of software supported on the Switches.

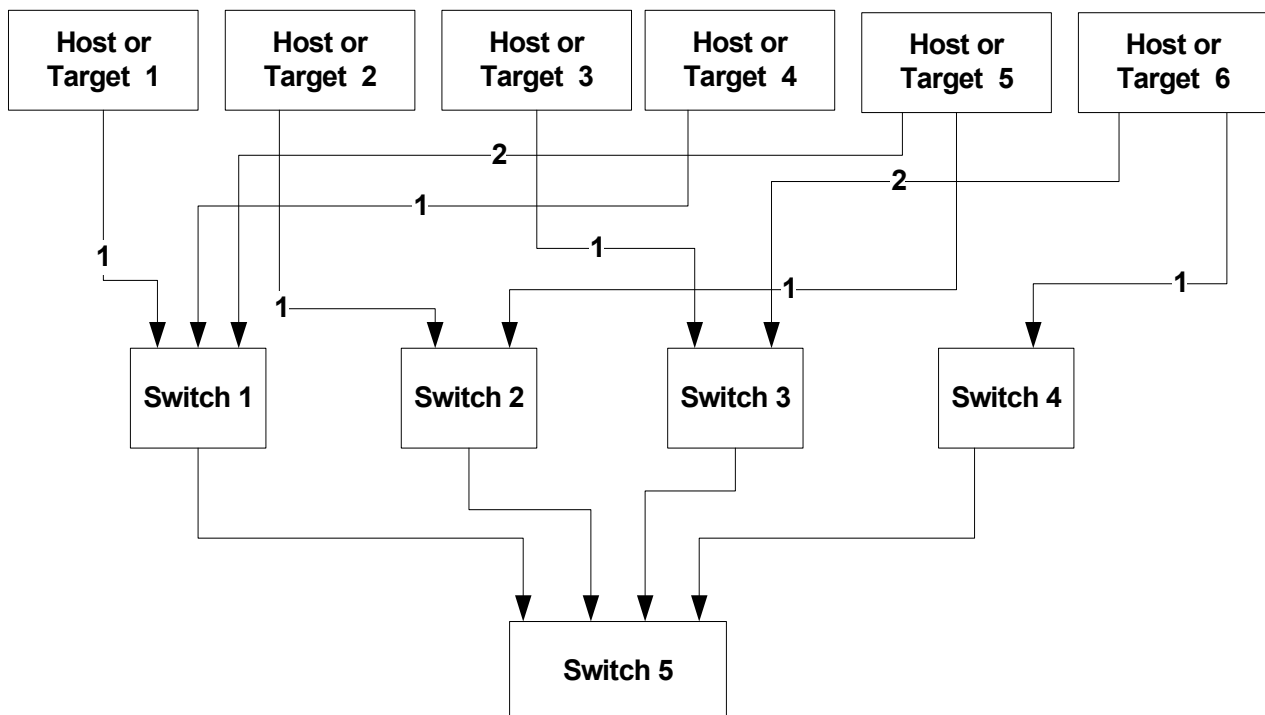
6.2.2 OPENFABRICS SOFTWARE RUNNING ON HOSTS

Where there are dependencies of OEM provided and IBTA defined management software (such as subnet managers and agents, performance managers and agents etc.) with OpenFabrics software running on Hosts, such software should be provided to UNH-IOL for interoperability testing. Any known dependencies should be communicated to UNH-IOL.

6.3 CLUSTER CONNECTIVITY

6.3.1 HOSTS AND TARGETS 1-6 SHOULD BE CONNECTED FROM PORT 1 OR 2 TO PORTS X IN ALL SWITCHES USING COMPLIANT INFINIBAND CABLES.

Figure 1 - Template for IB Interop Setup



7 IWARP HW DESCRIPTION & CONNECTIVITY

7.1 IWARP BASIC CONNECTIVITY (P1P1)

7.1.1 RNIC 1 ON ONE HOST SHOULD BE DIRECTLY CONNECTED TO RNIC 2 ON ANOTHER HOST OR TO A 10GbE SWITCH.

7.1.2 WITH 10GbE CABLES

7.2 SWITCHES AND SOFTWARE NEEDED

7.2.1 SWITCHES PROVIDED BY OEMS

It is necessary that Switches provided by OEMs cover the full breadth of software versions supported by the Switch OEMs. Port count is not critical for the tests. It is recommended that OEMs provide a switch per variations of software supported on the Switch.

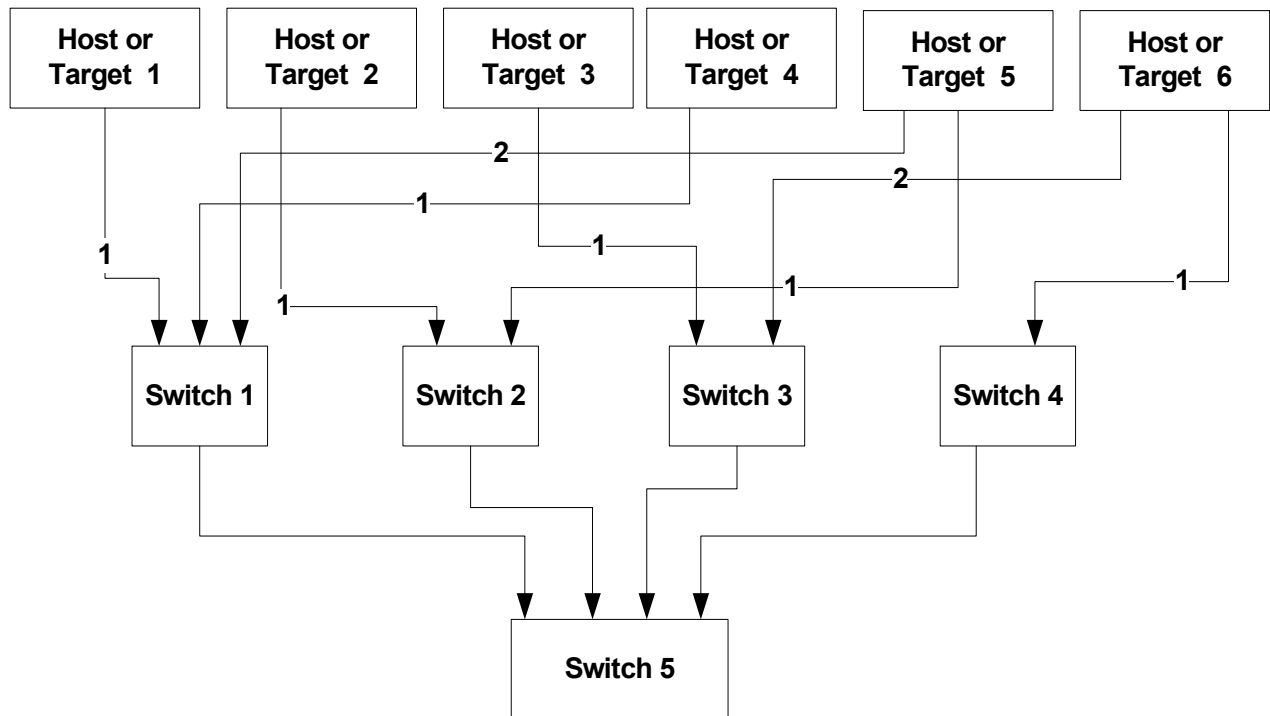
7.2.2 OPENFABRICS SOFTWARE RUNNING ON RNICs

Where there are dependencies of OEM provided with OpenFabrics software running on RNICs, such software should be provided to UNH-IOL for interoperability testing, and any known dependencies should be communicated to UNH-IOL.

7.3 CLUSTER CONNECTIVITY

7.3.1 HOSTS AND TARGETS 1-6 SHOULD BE CONNECTED TO SWITCHES USING 10GbE CABLES.

Figure 2 Template for iWARP Interop Setup



7.4 GATEWAY, BRIDGES, ROUTERS CONNECTIVITY

TBD

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

8 RoCE HW DESCRIPTION & CONNECTIVITY

8.1 RoCE BASIC CONNECTIVITY (P1P1)

8.1.1 RCA 1 ON ONE HOST SHOULD BE DIRECTLY CONNECTED TO RCA 2 ON ANOTHER HOST OR TO A 10/40 GbE SWITCH DCB ENABLED.

8.1.2 CONNECTED WITH 10/40 GbE CABLES

8.2 SWITCHES AND SOFTWARE NEEDED

8.2.1 SWITCHES PROVIDED BY OEMS

RoCE testing is being introduced as of October 2012 and the choice of Ethernet Fabrics such as Fabric Path, QFabric, MLAG, SPB, TRILL and others are initially not being addressed. This allows us to start Beta Testing RoCE with just one 10/40 GbE Ethernet Switch which is DCB enabled. In future Interop events we will consider using multiple switches from vendors such as Brocade, Cisco, Extreme, HP, Mellanox and others which will allow us to test various Ethernet Fabric solutions.

8.2.2 OPENFABRICS SOFTWARE RUNNING ON RCAs

Where there are dependencies of OEM provided with OpenFabrics software running on RCAs, such software should be provided to UNH-IOL for interoperability testing, and any known dependencies should be communicated to UNH-IOL.

8.2.3 RoCE PRIORITY LEVELS

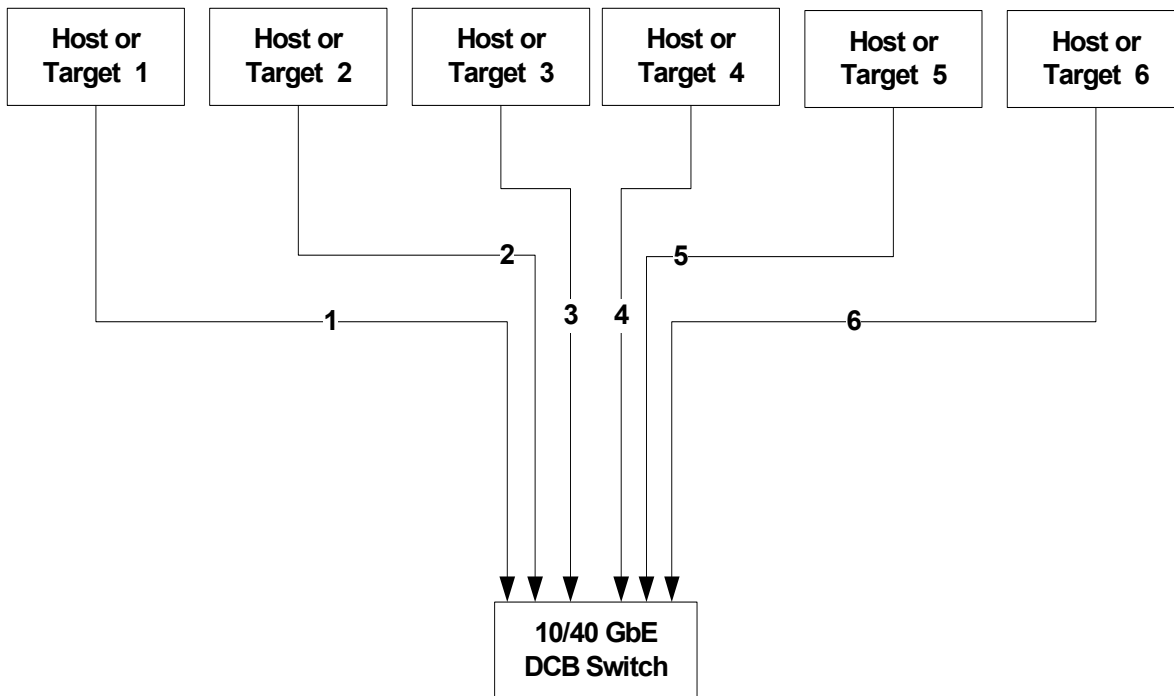
Ethernet provides a construct, called a Priority Level which corresponds conceptually to InfiniBand's SLs. Eight priorities, numbered zero through seven are supported. As in InfiniBand, a verbs consumer accessing a RoCE port specifies its desired service level, which is then mapped to a given Ethernet Priority. The default mapping is as follows:

- SL 0-7 are mapped directly to Priorities 0-7 respectively
- SL 8-15 are reserved.

8.3 FABRIC CONNECTIVITY

8.3.1 HOSTS AND TARGETS 1-6 SHOULD BE CONNECTED TO SWITCHES USING 10/40 GbE CABLES.

Figure 3 Template for RoCE Interop Setup



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

9 FW & SW INSTALLATION

9.1 BURNING THE FW

9.1.1 FIRMWARE POLICY

Firmware Policy during the Interop Debug Event

The firmware used during the Interop Debug Event is at the discretion of the device vendor. Vendors will be allowed to make changes to the firmware during the Interop Debug Event. However changes should be made as early in the event period as possible to reduce the amount of retesting which will result from these changes.

Firmware Policy during the Interop GA Event

The firmware image used during the Interop GA Event must be provided to the UNH-IOL at least one week prior to the event. No firmware changes of any kind are allowed during the Interop GA Event. If the vendor does not provide updated firmware by the deadline, then the UNH-IOL will use the firmware from the Interop Debug Event or from the vendor's website, whichever is more current.

Firmware Policy after the Interop GA Event

The firmware used to obtain the OFA Logo (or a child of this firmware with the same base functionality) must be the default publicly available firmware on the vendor's website and must be the default firmware that is shipped with the product. This must be completed within six months of the Interop GA Event.

9.1.2 PLEASE REFER TO FIRMWARE BURNING TOOLS AND PROCEDURES DOCUMENTATION FROM HCA IB VENDOR

9.2 OPERATING SYSTEM INSTALLATION

9.2.1 OPERATING SYSTEM POLICY

The OS used during an Interop Debug Event will be determined by the OFA IWG and will be none as the primary OS. All available updates will be installed prior to the start of the Interop Debug Event and frozen in place for the duration of the Interop Debug Event.

The OS used during an Interop GA Event will be the same agreed-upon version of the primary OS tested during the Interop Debug Event. The updates applied at the start of the Interop Debug Event will remain frozen in place for the duration of the Interop GA Event.

In addition to the mandatory testing performed using the primary OS, beta testing using the secondary operating systems is performed after completion of mandatory testing. The secondary operating systems are deployed in a similar manner to the primary OS, in that updates are applied at the beginning of the Interop Debug Event and frozen in place for the duration of the Interop GA Event.

9.2.2 OPERATING SYSTEM INSTALLATION

Install the primary OS on all hosts in the cluster. Use a package manager to update all installed packages to their latest versions available as of the start of the Interop Debug Event.

Install the secondary operating systems on all hosts in the cluster. Use a package manager to update all installed packages to their latest versions available as of the start of the Interop Debug Event. Install and test as many secondary operating systems as time permits.

9.3 SW INSTALLATION

9.3.1 SOFTWARE POLICY

Software Policy during an Interop Debug Event

The software used during an Interop Debug Event will be an agreed-upon RC release of the subsequent OFED version. During the Interop Debug Event vendors will be allowed to make changes to the software, provided that the changes are based on the same RC release. Vendors are not allowed to extensively modify the software or completely replace it.

Software Policy during the Interop GA event

The software used during an Interop GA Event will be the GA release of the same OFED version as was used during the Interop Debug Event. No software changes of any kind are allowed during the Interop GA Event. It is the vendor's responsibility to ensure that any changes made during the Interop Debug Event are present in the OFED GA release. Vendors whose products do not use firmware may request that patches be applied to an OFED GA release if that release has known defects that prevent the vendor product from being interoperable. The Arbitration Committee will be responsible for approving the requested patches.

Software Policy after the Interop GA event

All products that are granted the OFA Logo must be distributed by default with the OFED GA version (or a later revision of OFED with the same base functionality).

9.3.2 PLEASE REFER TO SOFTWARE INSTALLATION MANUAL FROM HCA IB VENDOR.

9.3.3 PLEASE REFER TO SOFTWARE INSTALLATION MANUAL FROM RNIC VENDOR.

9.4 SUMMARY

- For the Interop GA Event the vendor cannot update or change any part of the device under test - this includes hardware, firmware and software. The only exception is for an outright hardware failure in which case the hardware may be replaced with an identical piece of hardware with the same SW and FW.
- If an end user requests customized firmware or a modified version of OFED, then the vendor must disclose that this is not an OFA certified configuration.
- The OFA reserves the right to revoke the OFA Logo for products that do not follow these policies.
- These policies will be in effect for the April 2011 Interop Events and all events thereafter.

9.5 HARDWARE POLICY

For MPI testing, HCA/RNIC vendors must provide at least five adapters. The adapters need not be all the same model, but they can be.

9.6 OFED USAGE

- OFED Release Candidates (RC) should be used during the Interop Debug Event. This allows vendors to resolve bugs and issues and commit them to the OFED tree before the OFED General Availability (GA) is released.
- OFED GA versions shall be used for the Interop GA Events.

10 GENERAL INSTRUCTIONS

10.1 FIRST STEP INSTRUCTIONS

- 1) Burn the FW release XXX on all HCAs and RNICs using the above procedure as required by vendor.
- 2) Host and Target Configuration
 - a) Install OFED software on host systems (using a 64 bit OS) configured to run OFED.
 - b) Install WinOF software on host systems (using a 64 bit OS) configured to run WinOF.
 - c) Configure non-OFED systems for use in the cluster as per the vendors instructions.
 - d) Configure iSER/SRP targets for use in the cluster as per the vendors instructions.
- 3) Install the switch or gateway with the candidate SW stack as required by vendor.
- 4) Burn the switch or gateway with the released FW as required by vendor.
- 5) Connect the Hosts and Targets to an appropriate switch following the basic connectivity.

10.2 INFINIBAND SUBNET MANAGERS

- 1) The OpenSM will be used to run all mandatory tests in the test plan
- 2) Vendor SM testing will include testing IPoIB, RDMA Interop and Open MPI testing. In order to reduce the scope of testing, iSER, NFS over RDMA, RDS, SDP, SM Failover and SRP will not be performed using vendor SMs.

10.3 OPERATING SYSTEM CONSIDERATIONS

- 1) The OFILG decided as of April 2012 that the various ULPs contained in this test plan will only be tested if it is supported by the Operating System.
- 2) As a requirement for the OFILG Logo, a vendor's DUT must pass all mandatory testing using an agreed upon primary OS and OpenSM. Additional beta testing is performed using secondary Operating Systems. This beta testing has no bearing on whether the OFILG Logo is granted to a device It is purely informative.

11 INFINIBAND SPECIFIC INTEROP PROCEDURES USING OFED

Note: UNH-IOL has created automated scripts to run many of the OFED based tests. Please contact them at ofalab@iol.unh.edu if you wish to obtain copies of the latest scripts

11.1 IB LINK INITIALIZE USING OFED FOR LINUX

11.1.1 Procedure

- 1) Select a pair of devices to test from the created topology
- 2) Determine the maximum port width and lane speed supported by both devices
- 3) Select a cable to use which has been certified for the link parameters determined by step 2 of section 10.1.1 during an IBTA Plugfest held within the last 6 months
- 4) Disconnect all IB cables from the selected devices
- 5) Shutdown all SMs running on the selected devices
- 6) Connect the selected devices back to back using the cable selected during step 3 of section 10.1.1
- 7) Wait for a physical indication that a link has been established
- 8) Verify that the link created in step 6 of section 10.1.1 has come up with the parameters determined in step 2 of section 10.1.1
- 9) Repeat steps 1-8 with a different device pairing
 - a) All unique device pairs present in the created topology must be tested; except SRP target to SRP target and gateway to SRP target.
 - b) Each device must link at the maximum port width and lane speed supported by both devices in all pairings for said device to pass link initialization testing

11.1.2 Method of Implementation for all Linux OSs

- 1) To perform step 7 of section 10.1.1:
 - a) Look for link LEDs on the ports you are using
- 2) To perform step 8 of section 10.1.1:
 - a) ssh into a device supporting such remote connections and is running the OFED stack; usually a compute node with an HCA
 - b) Run "ibdiagnet -wt <desired-topology-file-name>"
 - c) Check the topology file created by the previous command:
 - i) Match the GUIDs to the devices in the selected pair
 - ii) Verify link width is the highest common denominator of pair capabilities (1x, 4x, 12x)
 - iii) Verify link speed is the highest common denominator of pair capabilities (2.5G, 5G, 10G, 14G)
- 3) To determine switch to SRP target and switch to switch link parameters
 - a) Run the commands outlined by step 2 of section 10.1.2 from a third device

- i) Should be a compute node with an HCA that is linked to a switch that is part of the desired pairing
- ii) Carefully match the GUIDS as you now have more than just two in the topology file

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

11.2 IB FABRIC INITIALIZATION USING OFED

11.2.1 Architect the Network we want to build.

- 1) Develop a cluster diagram based on the devices that have been submitted for Interop Testing and assign IP addresses to the IPoIB interfaces and the ethernet management interfaces.
- 2) See [Figure 4- Sample Network Configuration](#) below.

11.2.2 Procedure

- 1) Connect the HCAs and switches as per the Architected Network and make sure that no SM/SA is running on the Fabric.
- 2) Start an SM on a device and let it initialize (all SM's will need to be tested)
- 3) Visually verify that all devices are in the active state. Verify that the LED is on when the port is active.
- 4) Run "ibdiagnet -wt <file>" to generate a topology file
- 5) Run "ibdiagnet -pc" to clear all port counters
- 6) Wait 17 seconds as per the specifications requirements.
- 7) Run "ibdiagnet -c 1000" to send 1000 node descriptions.
- 8) Run "ibdiagnet" to generate fabric report.
 - a) Use /tmp/ibdiagnet.sm file to determine running sm
 - b) sminfo can also be used to determine the master SM or saquery -s to find all SMs.

Note: "ibdiagnet -r" seg faulted but was fixed in OFED 1.5 according to Bug 1618
- 9) Run "ibchecknet" to build guid list.
- 10) Run "ibdiagnet -t <file>" to compare current topology to the previously generated topology file

11.2.3 Verification Procedures

- 1) Review "PM Counters" section of the fabric report. There should be no illegal PM counters. The Specification says there should be no errors in 17 seconds.
- 2) Review "Subnet Manager " section of the fabric report. Verify that the running SM is the one you started and verify number of nodes and switches in the fabric.
- 3) Review the ibchecknet report and verify that there are no duplicate GUIDs in the fabric
- 4) Verify that step 10 above indicates that the topology before the test and the topology after the test are the same.

Restart all devices in the fabric and follow Sections 10.2.2 and 10.2.3. Run the SM from a different device in the fabric until all SMs present have been used. All SMs on managed switches (including those switches running **opensm**) should be tested and at least one instance of **opensm** on an HCA must be tested. If there are HCAs from more than one vendor, then **opensm** should be run from each vendor's HCA.

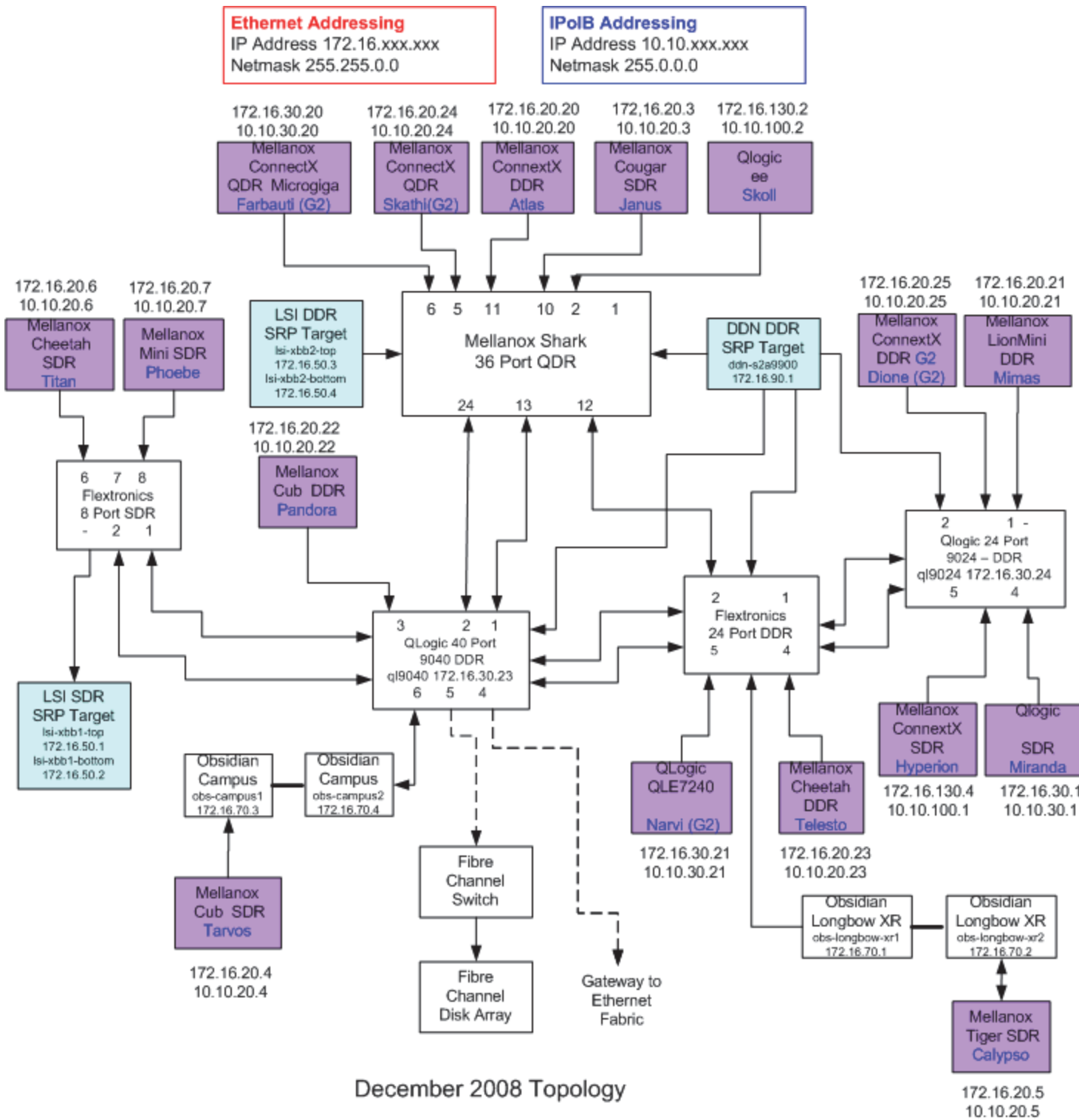
Each device must pass all verification procedures with every SM to pass Fabric Initialization test.

Table 27 - ibdiagnet commands

Commands	Description
ibdiagnet -c 1000	Send 1000 node descriptions
ibdiagnet -h	Help
ibdiagnet -lw 4x - ls 2.5	Specify link width and speed
ibdiagnet - pc	Clear counters
ibdiagnet -t <file>	Compare current topology to saved topology
ibdiagnet -wt	Writes the topology to a file

Note: The topology file is being generated after the SM starts but before any testing has started. The topology comparison is being performed after testing has been completed but before the systems get rebooted. A topology check is performed during every part of every test section that does not specifically state "change the topology". For example Fabric Init only has 1 part so there is only 1 check but RDS has 2 parts so 2 checks are performed. However, IPoIB has 3 parts for each of 2 modes but 1 of those parts specifically says to change the topology so only 4 checks occur.

Figure 4 - Sample Network Configuration



11.3 IB IpoIB CONNECT MODE (CM) USING OFED

11.3.1 SETUP

Connect the HCAs and switches as per the Architected Network and make sure that no SM is running on the Fabric.

This procedure, as the previous ones, will be based on the cluster connectivity. An SM/SA which supports IpoIB (sufficient IB multicast support) will be running on the HCAs, or on a switch with an embedded SM/SA or a third HCA which would only run SM/SA for the partner pair (with a switch in the middle). This procedure has been developed for Linux and may be ported to Windows if there is sufficient vendor support.

Optional: In the procedures below, an IB analyzer can be inserted in the appropriate link to obtain traces and validate the aspects of the procedures specifically detailed below in subsequent sections.

11.3.2 IpoIB INTERFACE CREATION AND IpoIB SUBNET CREATION

- 1) Configure IpoIB address. All addresses must reside on the same subnet.
 - a) Set interfaces to 10.0.0.x/24 (10.0.0.x/netmask 255.255.255.0) using the command `ifconfig ib0 10.0.0.x netmask 255.255.255.0`

11.3.3 .BRINGING THE IpoIB IN CONNECTED MODE

- 1) `echo 'connected' > /sys/class/net/ib0/mode`
- 2) Validate CM mode by checking that `"/sys/class/net/<I/F name>/mode"` equal to **'connected'**
- 3) Repeat steps 1-2 in section 10.3.3 on all nodes being tested.

11.3.4 PING PROCEDURES

Step A

- 1) Stop all SM's and verify that none are running
- 2) Power cycle all switches in the fabric (this insures that the new SM will configure all the links and create the multi-cast join).
- 3) Start an SM (All SM's will need to be tested) and let it initialize
 - a) Visually verify that all devices are in the active state. Verify that the LED is on when the port is active.
 - b) Run `"ibdiagnet -r"` and verify that the SM you started is the one that is running and and that it is the master. You will need to know the GUID of the device since the SM will be reassigned on each reboot.
 - c) Verify that all nodes and switches were discovered.
- 4) Examine the arp table (via `arp -a`) and remove the destination node's ib0 address from the sending node's arp table (via `arp -d`).

- 5) Ping every HCA except localhost with packet sizes of 511, 1025, 2044, 8192, 32768 and 65507.
 - a) ping -i 0.2 -t 3 -c 10 -s <ping size> <destination>
 - i) "-i" - interval 0.2 seconds
 - ii) "-t" - IP Time to Live equals 3 seconds
 - iii) "-c" - count equals 100
 - iv) "-s" - size of the ping
 - v) "destination" - the IP address of the IPoIB interface being pinged.
 - b) Repeat step #4 before issuing each ping command. Every packet size is a new ping command.
 - 6) In order to pass Step A, a reply must be received for every ping sent (without losing a single packet) while using each one of the SMs available in the cluster.
- Step B**
- 1) Bring up all HCAs but one.
 - 2) Start an SM (all SMs will need to be tested).
 - 3) Check for ping response between all node (All to All).
 - a) A response from the disconnected HCA should not be returned.
 - 4) Disconnect one more HCA from the cluster.
 - 5) Ping to the newly disconnected HCA from all nodes (No response should be returned).
 - 6) Connect the first machine (the one that was not connected) and check for ping response from all nodes that are still connected.
 - 7) Connect the disconnected HCA to a different switch on the subnet which will change the topology.
 - 8) Ping again from all nodes (this time we should get a response).
 - 9) Follow Step B, this time bring the interface down and then back up using ifconfig ibX down and ifconfig ibX up commands instead of physically disconnecting the HCAs.
- Note:** Each step must exhibit the expected behavior while using each SM in order for the device to pass Step B overall.

Step C

Follow Step A and B using a different SM until all SM's have been used. Only one instance of each available SM is required. Steps A, B, and C must pass in order for the device to pass 10.3.4 overall.

11.3.5 SFTP PROCEDURE

SFTP procedures require an SFTP server to be configured on each machine in the partner pair. An SFTP client needs to be available on each machine as well. The default RHEL install includes both.

A 4 MB file will be SFTP'd to the partner and then SFTP'd back and binary compared to the original file, this will be done in each direction and then bidirectional using every SM available.

11.3.5.1 SETUP

- 1) Make sure vsftpd is installed on each node for SFTP application.
- 2) A special account for this should be created as follows:
 - b) Username: Interop
 - c) Password: openfabrics

11.3.5.2 PROCEDURE

- 1) Run SFTP server on all nodes.
- 2) Start an SM (all SM's will need to be tested) and let it initialize
 - a) Verify that the running SM is the one you started.
- 3) SFTP:
 - a) Connect an HCA pair via SFTP on IPoIB using the specified user name and password.
 - b) Put the 4MB file to the /tmp dir on the remote host.
 - c) Get the same file to your local dir again.
 - d) Compare the file using the command *cmp tfile tfile.orig*.
 - i) The two must be identical
- 4) Repeat the procedure with a different SM.

Note: Every node must SFTP the 4MB file to all others using all SM's and the files must be identical as determined by the binary compare in order for the device to pass 10.3.5 overall.

Note: Sections 10.3.4 and 10.3.5 must pass using the configuration determined by sections 10.3.1, 10.3.2, and 10.3.3 for the device to pass IPoIB Connected mode overall.

11.4 IB IPoIB DATAGRAM MODE (DM) USING OFED

11.4.1 SETUP

Connect the HCAs and switches as per the Architected Network and make sure that no SM is running on the Fabric.

This procedure, as the previous ones, will be based on the cluster connectivity. An SM/SA which supports IPoIB (sufficient IB multicast support) will be running on the HCAs, or on a switch with an embedded SM/SA or a third HCA which would only run SM/SA for the partner pair (with a switch in the middle). This procedure has been developed for Linux and may be ported to Windows if there is sufficient vendor support.

Optional: In the procedures below, an IB analyzer can be inserted in the appropriate link to obtain traces and validate the aspects of the procedures specifically detailed below in subsequent sections.

11.4.2 IPOIB INTERFACE CREATION AND IPOIB SUBNET CREATION

- 1) Configure IPoIB address. All addresses must reside on the same subnet.
 - a) Set interfaces to 10.0.0.x/24 (10.0.0.x/netmask 255.255.255.0) using the command `ifconfig ib0 10.0.0.x netmask 255.255.255.0`

11.4.3 .BRINGING THE IPOIB IN DATAGRAM MODE

- 1) `echo 'datagram' > /sys/class/net/ib0/mode`
- 2) Validate DM mode by checking that `"/sys/class/net/<I/F name>/mode"` equal to **'datagram'**
- 3) Repeat steps 1-2 in section 10.4.3 on all nodes being tested.

11.4.4 PING PROCEDURES

Step A

- 1) Stop all SM's and verify that none are running
- 2) Power cycle all switches in the fabric (this insures that the new SM will configure all the links and create the multi-cast join).
- 3) Start an SM (All SM's will need to be tested) and let it initialize
 - a) Visually verify that all devices are in the active state. Verify that the LED is on when the port is active.
 - b) Run `"ibdiagnet -r"` and verify that the SM you started is the one that is running and and that it is the master. You will need to know the GUID of the device since the SM will be reassigned on each reboot.
 - c) Verify that all nodes and switches were discovered.

Note: `ibdiagnet` may show more switches than indicated by the physical number of switch platforms present. This is because some switches have multiple switch chips.
- 4) Examine the arp table (via `arp -a`) and remove the destination node's ib0 address from the sending node's arp table (via `arp -d`).
- 5) Issue the command: `sysctl net.ipv4.neigh.ib0.unres_qlen=33`

- a) This sets the qlen variable to 33 which increases the buffer size so that you do not get an initial dropped packet when using ping sizes 8192 and greater. 1
- 6) Ping every HCA except localhost with packet sizes of 511, 1025, 2044, 8192, 32768 and 65507. 2
- a) ping -i 0.2 -t 3 -c 10 -s <ping size> <destination> 3
- i) "-i" - interval 0.2 seconds 4
- ii) "-t" - IP Time to Live equals 3 seconds 5
- iii) "-c" - count equals 100 6
- iv) "-s" - size of the ping 7
- v) "destination" - the IP address of the IPoIB interface being pinged. 8
- b) Repeat step #4 before issuing each ping command. Every packet size is a new ping command. 9
- 7) In order to pass Step A, a reply must be received for every ping sent (without losing a single packet) while using each one of the SMs available in the cluster. 10

Step B

- 1) Bring up all HCAs but one. 11
 - 2) Start an SM (all SMs will need to be tested). 12
 - 3) Check for ping response between all node (All to All). 13
 - a) A response from the disconnected HCA should not be returned. 14
 - 4) Disconnect one more HCA from the cluster. 15
 - 5) Ping to the newly disconnected HCA from all nodes (No response should be returned). 16
 - 6) Connect the first machine (the one that was not connected) and check for ping response from all nodes that are still connected. 17
 - 7) Connect the disconnected HCA to a different switch on the subnet which will change the topology. 18
 - 8) Ping again from all nodes (this time we should get a response). 19
 - 9) Follow Step B, this time bring the interface down and then back up using ifconfig ibX down and ifconfig ibX up commands instead of physically disconnecting the HCAs. 20
- Note:** Each step must exhibit the expected behavior while using each SM in order for the device to pass Step B overall. 21

Step C

- 1) Follow Step A and B using a different SM until all SM's have been used. Only one instance of each available SM is required. Steps A, B, and C must pass in order for the device to pass 10.4.4 overall. 22
- 2) Issue the command: sysctl net.ipv4.neigh.ib0.unres_qlen=3 23
- a) This sets the qlen variable back to the default. 24

11.4.5 SFTP PROCEDURE

SFTP procedures require an SFTP server to be configured on each machine in the partner pair. An SFTP client needs to be available on each machine as well. The default RHEL install includes both. 25

A 4 MB file will be SFTP'd to the partner and then SFTP'd back and binary compared to the original file, this will be done in each direction and then bidirectional using every SM available.

11.4.5.1 SETUP

- 1) Make sure vsftpd is installed on each node for SFTP application.
- 2) A special account for this should be created as follows:
 - b) Username: Interop
 - c) Password: openfabrics

11.4.5.2 PROCEDURE

Run SFTP server on all nodes.

- 1) Start an SM (all SM's will need to be tested) and let it initialize
 - a) Verify that the running SM is the one you started.
- 2) SFTP:
 - a) Connect an HCA pair via SFTP on IPoIB using the specified user name and password.
 - b) Put the 4MB file to the /tmp dir on the remote host.
 - c) Get the same file to your local dir again.
 - d) Compare the file using the command *cmp tfile tfile.orig*.
 - i) The two must be identical
- 3) Repeat the procedure with a different SM.

Note: Every node must SFTP the 4MB file to all others using all SM's and the files must be identical as determined by the binary compare in order for the device to pass 10.4.5 overall.

Note: Sections 10.4.4 and 10.4.5 must pass using the configuration determined by sections 10.4.1, 10.4.2, and 10.4.3 for the device to pass IPoIB Datagram mode overall.

11.5 IB SM FAILOVER AND HANDOVER PROCEDURE USING OFED

11.5.1 SETUP

- 1) Connect HCAs per the selected topology.
- 2) In this test, all active SMs on the fabric which are going to be tested, must be from the same vendor. They will be tested pairwise; two at a time.

11.5.2 PROCEDURE

- 1) Disable all SMs in the cluster then start a SM on either machine in a chosen pair.
- 2) Run "saquery" on a node in the fabric.
 - a) Verify that all nodes in the cluster are present in the output
- 3) Using the ibdiagnet tool with the -r option, verify that the running SM is the master.
- 4) Start a SM on the second machine in the current pair.
- 5) Verify that the SMs behave according to the SM priority rules. Use "ibdiagnet -r" again.
 - a) SM with highest numerical priority value is master and the other is in standby.
 - a) If both SMs have the same priority value then the SM with the smallest guid is master and the other is in standby.
- 6) Run "saquery" on either machine in the current pair.
 - a) Verify that all nodes in the cluster are present in the output.
- 7) Shutdown the master SM.
- 8) Verify the other active SM goes into the master state using "ibdiagnet -r" again.
- 9) Run "saquery" on either machine in the current pair.
 - a) Verify that all nodes in the cluster are present in the output.
- 10) Start the SM you just shutdown.
- 11) Verify that the newly started SM resumes it's position as master while the other goes into standby again.
- 12) Run "saquery" on either machine in the current pair.
 - a) Verify that all nodes in the cluster are present in the output.
- 13) Shutdown the standby SM.
- 14) Verify that the previous master SM is still the master.
- 15) Run "saquery" on either machine in the current pair.
 - a) Verify that all nodes in the cluster are present in the output.
- 16) Repeat steps 1-15 above 2 more times, ensuring that the below criteria is met (total of 3 tests per pair which can be run in any order):
 - a) First SM to be started having highest numerical priority value.
 - b) Second SM to be started having highest numerical priority value.

- c) Both SMs having equal numerical priority values. 1
 - 17) Repeat steps 1-16 until all possible SM pairs from identical vendors in the 2
cluster have been tested. 3
 - 18) All of the "saquery" commands must return the expected list of nodes in 4
order for the SMs in this test to receive a passing grade. 5
- 6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

11.6 IB SRP USING OFED

11.6.1 SETUP

Connect the HCAs and switches as per the Architected Network and make sure that no SM is running on the Fabric.

Note: As of the April 2012 Interop events, one SRP target (i.e.target port) should present 2 or more volumes. All other target ports may be limited to one volume per port. This decision was made in order to reduce the amount of time required to run the tests.

Note: As of October 2012, the SRP Extended Procedure is a Beta test

11.6.2 SRP CORE PROCEDURE - MANDATORY

- 1) Start an SM (all SM's will need to be tested) and let it initialize
 - a) Verify that the running SM is the one that you started
- 2) Choose a node to work with
- 3) Unload the srp module
- 4) Load srp module with cmd_sg_entries=255
 - a) **Example:** modprobe ib_srp cmd_sg_entries=255
 - b) Let it initialize
- 5) Verify that the module loaded correctly
 - a) **Example:** lsmod | grep ib_srp
- 6) Load srp_daemon with -e -o -n options
 - a) **Example:** srp_daemon -e -o -n
 - b) Let it initialize
- 7) Find all volumes from all targets
 - a) Use lsscsi

Note: As of April 2012, the OFILG mandated that the target only include two volumes when doing mandatory testing.
- 8) Perform 6GB read from srp volume to null
 - a) **Example:** dd if=\$drive of=/dev/null count=600 bs=10M
- 9) Perform 6GB write from zero to srp volume
 - a) **Example:** dd if=/dev/zero of=\$drive count=600 bs=10M
- 10) Perform steps #8 and #9 for both volumes found from each target as determined by step #7
- 11) Unload srp module
- 12) Repeat steps 2 through 9 for all HCAs
- 13) Reboot all devices in the fabric and repeat the procedure using a different SM.

Note: An HCA must successfully complete all DD operations to and from all volumes on all targets using all available SM's in order to pass SRP testing.

11.6.3 SRP EXTENDED PROCEDURE - BETA

- 1) Start an SM (all SM's will need to be tested) and let it initialize
 - a) Verify that the running SM is the one that you started
- 2) Choose a node to work with
- 3) Unload the srp module
- 4) Load srp module with `cmd_sg_entries=255 allow_ext_sg=1 indirect_sg_entries=2048`
 - a) **Example:** `modprobe ib_srp cmd_sg_entries=255 allow_ext_sg=1 indirect_sg_entries=2048`
 - b) Let it initialize
- 5) Verify that the module loaded correctly
 - a) **Example:** `lsmod | grep ib_srp`
- 6) Load `srp_daemon` with `-e -o -n` options
 - a) **Example:** `srp_daemon -e -o -n`
 - b) Let it initialize
- 7) Find all volumes from all targets
 - a) Use `ls SCSI`

Note: As of April 2012, the OFILG mandated that the target only include two volumes when doing mandatory testing.
- 8) Perform 6GB read from srp volume to null
 - a) **Example:** `dd if=$drive of=/dev/null count=600 bs=10M`
- 9) Perform 6GB write from zero to srp volume
 - a) **Example:** `dd if=/dev/zero of=$drive count=600 bs=10M`
- 10) Perform steps #8 and #9 for both volumes found from each target as determined by step #7
- 11) Unload srp module
- 12) Repeat steps 2 through 9 for all HCAs
- 13) Reboot all devices in the fabric and repeat the procedure using a different SM.

Note: An HCA must successfully complete all DD operations to and from all volumes on all targets using all available SM's in order to pass SRP testing

11.7 IB ETHERNET GATEWAY USING OFED

11.7.1 PROCEDURE

- 1) Connect the HCA of the IB host to the IB fabric. Connect the Ethernet Gateway to the IB fabric. Connect the Ethernet gateway to the Ethernet network or Ethernet device. Start the SM to be used in this test.
- 2) Determine which ULP your ethernet gateway uses and be sure that ULP is running on the host (VNIC or IPoIB).
- 3) Restart the ULP or using the tool provided by the ULP, make sure that the host "discovers" the Ethernet Gateway. Configure the interfaces and make sure they are up.
- 4) Run ping from the host to the Ethernet device. While the ping is running, kill the master SM. Verify that the ping data transfer is unaffected.
- 5) Reboot the Ethernet Gateway. After the Ethernet Gateway comes up, verify that the host can discover the Ethernet Gateway as it did before and we are able to configure the interfaces.
- 6) Restart the ULP used by Ethernet Gateway and verify that after the ULP comes up, the host can discover the Ethernet Gateway and we are able to configure the interfaces.
- 7) Unload the ULP used by Ethernet Gateway and check that the Ethernet Gateway shows it disconnected. Load the ULP and verify that the Ethernet gateway shows the connection.
- 8) Repeat step 4 by using ssh and scp instead of ping.

11.8 IB FIBRECHANNEL GATEWAY USING OFED

11.8.1 PROCEDURE

- 1) Connect the HCA of the IB host to the IB fabric. Connect the FC Gateway to the IB Fabric (how to do this is determined by the FC Gateway vendor). Connect the FC Gateway to the FC network or FC device. Start the SM to be used in this test.
- 2) Configure the FC Gateway appropriately (how to do this is vendor specific).
- 3) Use ibsrpdm tool in order to have the host "see" the FC storage device. Add the storage device as target.
- 4) Run basic dd application from the SRP host to the FC storage device.
- 5) Run basic dd application from the SRP host to the FC storage device. While the test is running, kill the master SM. Verify that the test completes properly.
- 6) Unload the SRP host / SRP Target (target first/host first) and check that the SRP connection is properly disconnected.
- 7) Load the SRP host / SRP Target. Using ibsrpdm, add the target.
- 8) Run basic dd application from the SRP host to the FC storage device.
- 9) Reboot the FC Gateway. After FC Gateway comes up, verify using ibsrpdm tool that the host see the FC storage device. Add the storage device as target.
- 10) Run basic dd application from the SRP host to the FC storage device.
- 11) Follow steps 1-10 above with each SM to be tested and with each HCA to be tested, until each HCA and each SM has been tested with the FC Gateway.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

12 ETHERNET SPECIFIC INTEROP PROCEDURES USING OFED

12.1 iWARP LINK INITIALIZE USING OFED

12.1.1 PURPOSE

The iWARP Link Initialize test is a validation that all iWARP devices receiving the OFA Logo can link and pass traffic under nominal (unstressed) conditions.

12.1.2 RESOURCE REQUIREMENTS

- 1) Gigabit or 10Gigabit iWARP RNIC,
- 2) Gigabit or 10Gigabit Ethernet Switch
- 3) Compliant Cables

12.1.3 DISCUSSION

The validation of the underlying transport infrastructure is essential to the end-users experience of the operation of the OFED software stack. To this end, this test confirms that iWARP devices receiving the OFA Logo can suitably link and pass traffic in any configuration. Exhaustive compliance testing of BER performance of the channel or electrical signaling of the ports is not performed; however, successful completion of this test provides further evidence of the robustness of the OFA logo bearing device.

12.1.4 PROCEDURE

- 1) Connect the two link partners together utilizing compliant cables.
- 2) Check all relevant LEDs on both ends of the link.
- 3) Verify that basic IP connectivity can occur by driving minimum size ICMP echo requests and replies across the link or equivalent traffic (including RDMA traffic if readily configured, in which case an additional RNIC responder station is required). To verify that an RDMA link has been initialized between Host A and Host B run the following commands:
 - a) Start a server in verbose mode on Host A:
 - i) `rping -sv`
 - b) Start a client on Host B to ping Host A.
 - i) `rping -cv -a Host A RNIC_IP_Address`
 - c) Optional Command for the client
 - i) `rping -cv -a Host A RNIC_IP_Address -C 4 -S 50`

Note: This sends a count of 4 pings and character strings of size 50
- 4) Repeat steps 1-3 for all combinations of 2 RNICs to switches, switch to switch, and RNIC to RNIC link partner combinations. Previously tested combinations resident in the OFILG cluster may be omitted.

12.1.5 OBSERVABLE RESULTS

- 1) Link should be established on both ends of the channel.
- 2) Traffic should pass in both directions. Error rates of 10e-5 or better should be readily confirmed (no lost frames in 10,000).

12.1.6 POSSIBLE PROBLEMS

- 1) Traffic directed to a switches IP management address may not be processed at high speed, in such cases, traffic should be passed across the switch to a remote responder.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

12.2 RoCE LINK INITIALIZE USING OFED

12.2.1 PURPOSE

The RoCE Link Initialize test is a validation that all RoCE devices receiving the OFA Logo can link and pass traffic under nominal (unstressed) conditions.

12.2.2 RESOURCE REQUIREMENTS

- 1) 10 or 40 Gigabit RoCE Channel Adapter (RCA)
- 2) 10 or 40 Gigabit RoCE Switch (DCB Enabled)
- 3) Compliant Cables

12.2.3 DISCUSSION

The validation of the underlying transport infrastructure is essential to the end-users experience of the operation of the OFED software stack. To this end, this test confirms that RoCE devices receiving the OFA Logo can suitably link and pass traffic in any configuration. Exhaustive compliance testing of BER performance of the channel or electrical signaling of the ports is not performed; however, successful completion of this test provides further evidence of the robustness of the OFA logo bearing device.

12.2.4 PROCEDURE

- 1) Connect the two link partners together utilizing compliant cables.
- 2) Check all relevant LEDs on both ends of the link.
- 3) Verify that basic IP connectivity can occur by driving minimum size ICMP echo requests and replies across the link or equivalent traffic (including RDMA traffic if readily configured, in which case an additional RoCE responder station is required). To verify that an RDMA link has been initialized between Host A and Host B run the following commands:
 - a) Start a server in verbose mode on Host A:
 - i) `rping -sv`
 - b) Start a client on Host B to ping Host A.
 - i) `rping -cv -a Host A RCA_IP_Address`
 - c) Optional Command for the client
 - i) `rping -cv -a Host A RCA_IP_Address -C 4 -S 50`
Note: This sends a count of 4 pings and character strings of size 50
- 4) Repeat steps 1-3 for all combinations of 2 RCAs to switches, switch to switch, and RCA to RCA link partner combinations. Previously tested combinations resident in the OFILG cluster may be omitted.

12.2.5 OBSERVABLE RESULTS

- 1) Link should be established on both ends of the channel.
- 2) Traffic should pass in both directions. Error rates of 10e-5 or better should be readily confirmed (no lost frames in 10,000).

12.3 RoCE FABRIC INIT USING OFED

This test will be developed for the April 2013 Interop Debug event

12.4 RoCE IPoCE

The test for IP over Converged Ethernet will be developed for the April 2013 Interop Debug event

12.5 RoCE INFINIBAND GATEWAY

This test will be developed for the April 2013 Interop Debug event

12.6 RoCE FIBRE CHANNEL GATEWAY

This test will be developed for the April 2013 Interop Debug event

1 |
2 |
3 |
4 |
5 |
6 |
7 |
8 |
9 |
10 |
11 |
12 |
13 |
14 |
15 |
16 |
17 |
18 |
19 |
20 |
21 |
22 |
23 |
24 |
25 |
26 |
27 |
28 |
29 |
30 |
31 |
32 |
33 |
34 |
35 |
36 |
37 |
38 |
39 |
40 |
41 |
42 |

13 TRANSPORT INDEPENDENT INTEROP PROCEDURES USING OFED

13.1 TI iSER USING OFED

13.1.1 IB SETUP

Connect initiator/target to switch as well as run one or more SMs (embedded in the switch or host based). If more than one SM, let the SMs split into master and slave.

Optional: In the procedures below, an IB analyzer can be inserted in the appropriate link to obtain traces and validate the aspects of the procedures specifically detailed below in subsequent sections.

13.1.2 IWARP SETUP

Connect iSER host initiator and target RNICs to an 10GbE switch.

13.1.3 RoCE SETUP

Connect iSER host initiator and target RCA to a 10/40 GbE switch which is DCB Enabled.

13.1.4 PROCEDURE

- 1) Load iSER target and iSER initiator to hosts from OpenFabrics tree, check iSER connection.
- 2) Run basic dd application from iSER initiator host connected to target.
- 3) [IB Specific Test] Run basic dd application from iSER initiator host connected to target. Kill the master SM while test is running and check that it completes properly.
- 4) Unload iSER initiator from a Host and check iSER connection properly disconnected on a target host.
- 5) Unload iSER target from a Host and check iSER connection properly disconnected on an initiator host.
- 6) [IB Specific Test] Repeat steps 2-5 now with the previous slave SM (we did not actually stop the target).

13.2 TI NFS OVER RDMA USING OFED

13.2.1 Installation

Note: Steps 2-4 are unneeded if an OFED supported OS is used along with an official OFED release downloaded from <http://www.openfabrics.org>

1) Verify that you are using a Linux kernel with NFS/RDMA on every system used

a) The NFS/RDMA client and server are both included in the mainline Linux kernel version 2.6.25 and later. This and other versions of the 2.6 Linux kernel can be found at: <ftp://ftp.kernel.org/pub/linux/kernel/v2.6/>

Note: OFED supported OS releases of lower kernel revision than mentioned above have been updated by their respected maintainers to allow NFS RDMA to function. Check the `nfs-rdma.release-notes.txt` provided with the OFED release you are using for supported OS releases.

Note: As of OFED 1.5.3 rc2 NFSoRDMA is not installed by default. To do so you must have built OFED from src with `nfsrdma=y` directive contained within the `ofed.conf` file used by the OFED installer.

i) To generate an `ofed.conf` file run the following from within the downloaded OFED src.

1. `$. /install.pl -p`

ii) Add the following directives to the generated `ofed-all.conf` file

1. `nfsrdma=y`

iii) Install OFED

1. `./install.pl -c ofed-all.conf`

2) Configure the RDMA stack on every system used

a) Make sure your kernel configuration has RDMA support enabled. Under Device Drivers => InfiniBand support, update the kernel configuration to enable InfiniBand support.

Note: the option name is misleading. Enabling InfiniBand support is required for all RDMA devices (IB, iWARP, etc.).

b) Enable the appropriate IB HCA support (`mlx4`, `mthca`, `ehca`, `ipath`, `qib`, etc.) or iWARP adapter support (`amso`, `cxgb3`, etc.).

c) If you are using InfiniBand, be sure to enable IP-over-InfiniBand (IPoIB) support.

3) Configure the NFS client

a) Your kernel configuration must also have NFS file system support and/or NFS server support enabled. These and other NFS related configuration options can be found under File Systems => Network File Systems.

4) Build, install, reboot

a) The NFS/RDMA code will be enabled automatically if NFS and RDMA are turned on. The NFS/RDMA client and server are configured via the hidden `SUNRPC_XPRT_RDMA` config option that depends on `SUNRPC` and `INFINIBAND`. The value of `SUNRPC_XPRT_RDMA` will be:

- i) - N if either SUNRPC or INFINIBAND are N, in this case the NFS/RDMA client and server will not be built
 - ii) - M if both SUNRPC and INFINIBAND are on (M or Y) and at least one is M, in this case the NFS/RDMA client and server will be built as modules
 - iii) - Y if both SUNRPC and INFINIBAND are Y, in this case the NFS/RDMA client and server will be built into the kernel
 - b) If you have followed the steps above and turned on NFS and RDMA, the NFS/RDMA client and server will be built.
 - c) Build a new kernel, install it and boot it
- 5) Check RDMA Setup
- a) If you are using InfiniBand, make sure there is a Subnet Manager (SM) running on the network.
 - b) Use IPoIB to ping two hosts.
- 6) Configure NFS exports, start NFS server
- a) Use two machines, one to act as the client and one to act as the server.
 - b) On the server system, configure the /etc/exports file and start the NFS/RDMA server. Export entries with the following formats have been tested:
 - i) /vol0 192.168.0.47(fsid=0,rw,async,insecure,no_root_squash)
 - ii) /vol0 192.168.0.0/255.255.255.0(fsid=0,rw,async,insecure,no_root_squash)
 - c) The IP address(es) is (are) the client's IPoIB address for an InfiniBand HCA or the client's iWARP address(es) for an RNIC.
- Note:** The "insecure" option must be used because the NFS/RDMA client does not use a reserved port. This does not interfere with normal NFS over TCP/IP operations.
- d) The remainder of this section will assume an export of /server
 - e) Start the NFS server
 - i) If the NFS/RDMA server was built as a module (CONFIG_SUNRPC_XPRT_RDMA=m in kernel config), load the RDMA transport module:
 - 1. \$ modprobe svcrdma
 - ii) Regardless of how the server was built (module or built-in), start the server:
 - 1. \$ /etc/init.d/nfs start or service nfs start
 - iii) Instruct the server to listen on the RDMA transport:
 - 1. \$ echo rdma 20049 > /proc/fs/nfsd/portlist
- 7) Check NFS Setup
- a) For the NFS components enabled above (client and/or server), test their functionality over standard Ethernet using TCP/IP or UDP/IP.
 - b) On the client system:

- i) Use this command to mount the NFS server export: 1
 - 1. `$ mount <server name or TCP/IP address>:<export> /<mount path>` 2
- ii) To verify that the mount is using TCP, run "cat /proc/mounts" and check the "proto" field for the given mount. 3
- 8) Check NFS/RDMA Setup 4
- a) For the NFS components enabled above (client and/or server), test their functionality over RDMA. 5
- b) On the client system: 6
 - i) If the NFS/RDMA client was built as a module (CONFIG_SUNRPC_XPRT_RDMA=m in kernel config), load the RDMA client module: 7
 - 1. `$ modprobe xprtrdma` 8
 - ii) Regardless of how the client was built (module or built-in), use this command to mount the NFS server export: 9
 - 1. `$/sbin/mount.nfs <IPoIB server name or address>:<export> /<mount path> -o \ rdma,port=20049` 10
- Note:** OFED will build and install the mount utility needed. The binary is called mount.nfs. Either this binary or the mount binary provided with nfs-utils revision greater than version 1.1 can be used. The remainder of this section will assume mount.nfs built by OFED is used. 11
- iii) To verify that the mount is using RDMA, run "cat /proc/mounts" and check the "proto" field for the given mount. 12
- 9) Connectathon 13
- a) Download the Connectathon test suite from <http://www.connectathon.org/nfstests.html> 14
- b) Install Connectathon on every client to be used 15
 - i) Modify tests.init within the connectathon tarball to suit your environment. 16
 - 1. The MOUNTCMD, UMOUNTCMD and MNTOPTIONS directives are unimportant as we will be calling the runtests connectathon binary directly. 17
 - 2. Be sure to remove the -fwritable-strings option from the CFLAGS variable. Your build will fail if this is used. 18
 - ii) Run make to build the connectathon binaries. 19
- 10) Test the connectathon runtests binary 20
- a) Run `sudo ./runtests -a -t /mnt/` to test the binary against the local file system. 21
- b) All tests should pass but you will see 1 warning. This is ok. 22

13.2.2 NFSoRDMA Test Procedure

- 1) **Note:** IB Only
- a) Start an SM
- 2) Server setup
- a) Add nfs rdma server support to the running kernel if not already present.
 - i) `$ modprobe svcrdma`
- b) Start the server
 - i) `$ /etc/init.d/nfs start`
- c) Tell the server to listen for rdma connection requests on port 20049
 - i) `$ echo rdma 20049 > /proc/fs/nfsd/portlist`
- 3) Client setup
- a) Add nfs rdma client support to the running kernel if not already present.
 - i) `$ modprobe xprtrdma`
- b) Mount the servers export using rdma
 - i) `$ /sbin/mount -t nfs <server IPoIB address>:/server /<mount path> -o \ rdma,port=20049 -i`
Note: <mount path> is assumed to be /mnt/<servername> for the remainder of this section
- c) Verify that the mount is using the rdma protocol
 - i) Verify that the mount is using RDMA, run "cat /proc/mounts" and check the "proto" field for the given mount.
- 4) Run Connectathon's runtests binary
 - a) `$./runtests -a -t /mnt/<servername>/<hostname>`
- 5) Repeat steps 2-4 using a new client-server pair until all nodes have acted as both a server and a client.
- 6) Repeat steps 2-5 using a new SM until all registered SM's have been used.
- 7) All tests run by the connectathon runtests binary must pass on all client nodes rdma mount points from all server nodes using all SM's in order for the device to pass [NFSoRDMA Test Procedure](#) overall.

13.3 TI RELIABLE DATAGRAM SERVICE (RDS) USING OFED

13.3.1 RDS-PING PROCEDURE

Note: RDS does not support iWARP

- 1) Use the command `modprobe rds_rdma` to add RDS support to the kernel
- 2) Verify that the kernel supports RDS by issuing the `rds-info` command.
 - a) The `rds-info` utility presents various sources of information that the RDS kernel module maintains. When run without any optional arguments `rds-info` will output all the information it knows of.

- 3) **[For IB]** Start one of the Subnet Managers in the cluster

Note: RDS is IP based so you need to provide a host address either through an out of band Ethernet connection or through IPoB. RDS also requires the LIDs to be set in an InfiniBand Fabric and therefore an SM must be run.

Note: All SMs in the fabric should be tested.

- 4) Choose a host and use `rds-ping host` to communicate with every other end point in the fabric.

Note: Be sure that you identify the correct host when using the command `rds-ping host`.

- a) `rds-ping` is used to test whether a remote node is reachable over RDS. Its interface is designed to operate in a similar way to the standard `ping(8)` utility, even though the way it works is pretty different.
- b) `rds-ping` opens several RDS sockets and sends packets to port 0 on the indicated host. This is a special port number to which no socket is bound; instead, the kernel processes incoming packets and responds to them.

- 5) Verify that all nodes respond without error.

Note: To avoid losing packets, do not run this while RDS-Stress is running.

13.3.2 RDS-STRESS PROCEDURE

- 1) Choose a host and start a passive receiving session for the RDS Stress test. It only needs to be told what port to listen on.

a) `$ rds-stress -p 4000`

- 2) Chose a second host and start an active sending instance giving it the address and port at which it will find a listening passive receiver. In addition, it is given configuration options which both instances will use.

a) `$ rds-stress -T 5 -s recvhost -p 4000 -t 1 -d 1`

Note: If you repeat the test in less than one minute you may get the error message "Cannot assign requested address" since the port numbers are not immediately reusable. Either wait or change the port number using the `-p` option

Note: The `-t` option is for the number of tasks (child processes), which defaults to 1 so "`-t 1`" is optional. The `-d` option is for the message queue depth, which also defaults to 1 so "`-d 1`" is optional.

- 3) Every second, the parent process will display statistics of the ongoing stress test. If the -T option is given, the test will terminate after the specified time and a summary is printed.
- 4) Verify that the test completes without error.
- 5) Repeat steps 1-4 until all end points in the cluster have been tested.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

13.4 TI uDAPLTEST COMMANDS USING OFED

Server Command: `dapltest -T S -D <ia_name>`

13.4.1 SETUP

- The `/etc/dat.conf` needs to be verified to be sure that the correct interface is used. By default the `dapl` interface for IB is `ib0` and for `iWARP` is `eth2`. If these are not correct for the current cluster then errors will occur.
- It is also important to verify that the desired `dapl` library is being used.
- [For IB] an SM needs to be running.
- [For `iWARP` hosts with Chelsio RNICs] Ensure that `/sys/module/iw_cxgb3/parameters/peer2peer` contains '1' on all hosts.

13.4.2 GROUP 1: POINT-TO-POINT TOPOLOGY

[1.1] 1 connection and simple send/recv:

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -R BE`
- client SR 256 1 server SR 256 1

[1.2] Verification, polling, and scatter gather list:

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 1024 3 -f \
- server SR 1536 2 -f

13.4.3 GROUP 2: SWITCHED TOPOLOGY

InfiniBand Switch: Any InfiniBand switch

iWARP Switch: 10 GbE Switch

RoCE Switch: 10/40 GbE DCB Enabled switch

[2.1] Verification and private data:

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 1024 1 \
- server SR 1024 1

[2.2] Add multiple endpoints, polling, and scatter gather list:

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 10 -V -P -R`
- BE client SR 1024 3 \
- server SR 1536 2

[2.3] Add RDMA Write :

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 256 1 \
- server RW 4096 1 server SR 256 1

[2.4] Add RDMA Read:

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`

- client SR 256 1 \ 1
- server RR 4096 1 server SR 256 1 2

13.4.4 GROUP 3: SWITCHED TOPOLOGY WITH MULTIPLE SWITCHES 3

Note: This test is **not applicable to RoCE** for the October 2012 Events 4

[3.1] Multiple threads, RDMA Read, and RDMA Write: 5 6

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 4 -w 8 -V -P -R BE` 7
- client SR 256 1 \ 8
- server RR 4096 1 server SR 256 1 client SR 256 1 server RR 4096 1 \ 9
- server SR 256 1 10

[3.2] Pipeline test with RDMA Write and scatter gather list: 11

- `dapltest -T P -s <server_name> -D <ia_name> -i 1024 -p 64 -m p RW`
8192 2 12 13

[3.3] Pipeline with RDMA Read: 14

- **InfiniBand:** `dapltest -T P -s <server_name> -D <ia_name> -i 1024 -p 64`
-m p RR 4096 2 15 16
- **iWARP:** `dapltest -T P -s <server_name> -D <ia_name> -i 1024 -p 64 -m`
p RR 4096 1 17 18

[3.4] Multiple switches: 19

- `dapltest -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 10 -V -P -R` 20
- BE client SR 1024 3 \ 21
- server SR 1536 2 22

13.5 TI RDMA BASIC INTEROP

13.5.1 Purpose

To demonstrate the ability of endpoints to exchange core RDMA operations across a simple network path. This test procedure validates the operation of endpoints at the RDMA level, in a simple network configuration.

The Basic RDMA interop test identifies interoperability issues in one of four ways:

- The inability to establish connections between endpoints
- The failure of RDMA operations to complete
- Incorrect data after the completion of RDMA exchanges
- Inconsistent performance levels.

13.5.2 General Setup

The RDMA interop procedure can be carried out using the OFA Verbs API to create RDMA Connections and send RDMA operation.

13.5.3 Topology

The topology of the network that interconnects the switches can be changed to validate operation of the endpoints over different networks paths. It is recommended that this procedure first be executed between endpoints connected by a single switch, and then the process repeated for more complex network configurations.

13.5.4 IB Setup

Connect endpoints to switch and run one or more SMs (embedded in the switch or host based).

13.5.5 iWARP Setup

Connect iWARP RDMA endpoints to an 10GbE switch.

13.5.6 RoCE Setup

Connect RoCE RCAs to a 10/40 GbE switch which is DCB Enabled.

13.5.7 RDMA Connectivity Setup

Each of the tests described below must be run twice with Host A being the server and then Host B being the server. This ensures that the different semantics associated with active and passive sides of the connection are exercised. This way each RDMA interface tested will be sending RDMA data (Requestor) in one test and receiving RDMA data (Target) in the next.

13.5.8 Small RDMA READ Procedure

- 1) Select the two devices that will be tested:
- 2) On the server device issue the following command on command line:
 - a) **[For IB & RoCE]** `ib_read_bw -d <dev_name> -i <port> -m 2048`
 - b) **[For iWARP]** - Not applicable - see 12.6.9
- 3) On the client device issue the following command on command line:

- a) **[For IB & RoCE]** `ib_read_bw -d <dev_name> -i <port> -s 1 -n 25000 -m 2048` 1
- b) **[For iWARP]** - Not applicable - see 12.6.9 2
- 4) Verify that the operation completed without error and the level of performance achieved is reasonable and as expected. 3

13.5.9 Large RDMA READ Procedure 4

- 1) Select the two devices that will be tested: 5
- 2) On the server device issue the following command on command line: 6
- a) **[For IB & RoCE]** `ib_read_bw -d <dev_name> -i <port> -m 2048` 7
- b) **[For iWARP]** - Not applicable - see 12.6.10 8
- 3) On the client device issue the following command on command line: 9
- a) **[For IB & RoCE]** `ib_read_bw -d <dev_name> -i <port> -s 1000000 -n 300 -m 2048` 10
- b) **[For iWARP]** - Not applicable - see 12.6.10 11
- 4) Verify that the operation completed without error and the level of performance achieved is reasonable and as expected. 12

13.5.10 Small RDMA Write Procedure 13

- 1) Select the two devices that will be tested: 14
- 2) On the server device issue the following command on command line: 15
- a) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -m 2048` 16
- b) **[For iWARP]** `rdma_bw -c -s 1 -n 25000` 17
- 3) On the client device issue the following command on command line: 18
- a) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -s 1 -n 25000 -m 2048` 19
- b) **[For iWARP]** `rdma_bw -c -s 1 -n 25000 RNIC_IP_Address` 20
- 4) Verify that the operation completed without error and the level of performance achieved is reasonable and as expected. 21

13.5.11 Large RDMA Write Procedure 22

- 1) Select the two devices that will be tested: 23
- 2) On the server device issue the following command on command line: 24
- a) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -m 2048` 25
- b) **[For iWARP]** `rdma_bw -c -s 1000000 -n 300` 26
- 3) On the client device issue the following command on command line: 27
- a) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -s 1000000 -n 300 -m 2048` 28
- b) **[For iWARP]** `rdma_bw -c -s 1000000 -n 300 RNIC_IP_Address` 29
- 4) Verify that the operation completed without error and the level of performance achieved is reasonable and as expected. 30

13.5.12 Small RDMA SEND Procedure

This procedure may fail due to the inability of a endpoint to repost the consumed buffers.

- 1) Select the two devices that will be tested:
- 2) On the server device issue the following command on command line:
 - a) **[For IB & RoCE]** `ib_send_bw -d <dev_name> -i <port> -m 2048`
 - b) **[For iWARP]** - Not applicable - see 12.6.9
- 3) On the client device issue the following command on command line:
 - a) **[For IB & RoCE]** `ib_writesend_bw -d <dev_name> -i <port> -s 1 -n 25000 -m 2048`
 - b) **[For iWARP]** - Not applicable - see 12.6.9
- 4) Verify that the operation completed without error and the level of performance achieved is reasonable and as expected.

13.5.13 Large RDMA SEND Procedure

This procedure may fail due to the inability of a endpoint to repost the consumed buffers.

- 1) Select the two devices that will be tested:
- 2) On the server device issue the following command on command line:
 - a) **[For IB & RoCE]** `ib_send_bw -d <dev_name> -i <port> -m 2048`
 - b) **[For iWARP]** - Not applicable - see 12.6.10
- 3) On the client device issue the following command on command line:
 - a) **[For IB & RoCE]** `ib_send_bw -d <dev_name> -i <port> -s 1000000 -n 300 -m 2048`
 - b) **[For iWARP]** - Not applicable - see 12.6.10
- 4) Verify that the operation completed without error and the level of performance achieved is reasonable and as expected.

13.5.14 Additional IB Notes

- 1) Alternate read commands available
 - a) Server command: `ib_read_bw -m 2048`
 - b) Client command (small): `ib_read_bw -s 1 -n 25000 IPoIB Address for server -m 2048`
 - c) Client command (large): `ib_read_bw -s 1000000 -n 300 IPoIB Address for server -m 2048`
- 2) Alternate write commands available
 - a) Server command: `ib_write_bw -m 2048`
 - b) Client command (small): `ib_write_bw -s 1 -n 25000 IPoIB Address for server`
 - c) Client command (large): `ib_write_bw -s 1000000 -n 300 IPoIB Address for server -m 2048`

- 3) Alternate send commands available
 - a) Server command: `ib_send_bw -m 2048`
 - b) Client command: `ib_send_bw -s 1 -n 25000 IPoB Address for server -m 2048`
 - c) Client command (large): `ib_send_bw -s 1000000 -n 300 IPoB Address for server -m 2048`
- 4) Explanation of parameters
 - a) "-d" allows you to specify the device name which may be obtained from the command `lane: ibv_devinfo`
 - b) "-i" allows you to specify the port number. This may be useful if you are running the tests consecutively because a port number is not immediately released and this will allow you to specify another port number to run the test.
 - c) "-s" - this is the size of the operation you wish to complete
 - d) "-n" - this is the number of operations you wish to complete.
 - e) "-m" - this specifies the IB PMTU size. AS of 10/3/2011 some devices did not support greater than 2048

13.5.15 Additional iWARP Notes

- 1) The "-c" option specifies to use the `rdma_cm` for connection

IB Example: DevInfo - Server

```
hca_id: mthca0
  fw_ver:          1.2.0
  node_guid:       0002:c902:0020:b4dc
  sys_image_guid:  0002:c902:0020:b4df
  vendor_id:       0x02c9
  vendor_part_id:  25204
  hw_ver:          0xA0
  board_id:        MT_0230000001
  phys_port_cnt:   1
    port: 1
      state:        PORT_ACTIVE (4)
      max_mtu:      2048 (4)
      active_mtu:   2048 (4)
      sm_lid:        1
      port_lid:      2
      port_lmc:      0x00
```

Command Line: `ib_read_bw -d mthca0 -i 1`

DevInfo - Client

```
hca_id: mlx4_0
  fw_ver:          2.2.238
  node_guid:       0002:c903:0000:1894
  sys_image_guid:  0002:c903:0000:1897
```



```
vendor_id:          0x02c9          1
vendor_part_id:    25418            2
hw_ver:            0xA0             3
board_id:          MT_04A0110002    4
phys_port_cnt:    2                 4
  port: 1          5
    state:         PORT_ACTIVE (4)   6
    max_mtu:       2048 (4)          7
    active_mtu:    2048 (4)          8
    sm_lid:        1                 8
    port_lid:      1                 9
    port_lmc:      0x00              10
```

Command Line: ib_send_bw -d mlx4_0 -i 1 10.0.0.1 -s 1 -n 300

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

13.6 TI RDMA STRESS TEST

13.6.1 Purpose

This test is designed to identify problems that arise when RDMA operations are performed over interconnection devices in the fabric. The test is not designed to measure the forwarding rate or switching capacity of a device, but does use performance measures to identify failures.

Test failures are identified by the following events:

- The inability to establish connections between endpoints
- The failure of RDMA operations to complete
- Incorrect data after the completion of RDMA exchanges
- Inconsistent performance levels.

13.6.2 Topology

This test does not define a detailed topology and can be used either on a single switch or across a RDMA fabric that may include gateways to and from other technologies. The test configuration depends on the number of endpoints available to perform the testing.

13.6.3 Switch Load

The switch load test validates proper operation of a switch when processing a large number of small RDMA frames. This test is analogous to normal switch testing.

- 1) Attach a device to each port on the switch.
- 2) Select two ports on the switch to test (This will be your control stream)
- 3) Generate RDMA WRITE Operations of size 1024 bytes 100, 000 times on each device by issuing the following commands
 - a) On the server device issue the following command on command line:
 - i) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -m 2048`
 - ii) **[For iWARP]** `rdma_bw -c -s 1024 -n 25000`
 - b) On the client device issue the following command on command line:
 - i) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -s 1024 -n 25000 -m 2048`
 - ii) **[For iWARP]** `rdma_bw -c -s 1024 -n 25000 RNIC_IP_Address`
- 4) This must be done on both devices at the same time.
- 5) On all other pairs generate RDMA WRITE Operations of size 1 byte continuously until the control stream completes.
- 6) Repeat above steps until all port pairs are tested.

- 7) Repeat the above steps with all endpoint pairs, except the control stream changed such that the size of the RDMA WRITE operation is 1,000,000 bytes (~1 MB)

13.6.4 Switch FAN in

The switch fan in test attempts to validate proper operation of RDMA exchanges in the presence of traffic loads that exceed the forwarding capacity of the switch. The test requires a minimum of two switches that are interconnected by one port pair.

- 1) Connect all possible endpoint pairs such that data exchanges between pairs must traverse the pair of ports interconnecting the switch. The control connections must be across the interconnect network.
- 2) Select two ports such that it has to cross both switches. (This will be your control stream)
- 3) Generate RDMA WRITE Operations of size 1024 bytes 100, 000 times on each device by issuing the following commands
 - a) On the server device issue the following command on command line:
 - i) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -m 2048`
 - ii) **[For iWARP]** `rdma_bw -c -s 1024 -n 25000`
 - b) On the client device issue the following command on command line:
 - i) **[For IB & RoCE]** `ib_write_bw -d <dev_name> -i <port> -s 1024 -n 25000 -m 2048`
 - ii) **[For iWARP]** `rdma_bw -c -s 1024 -n 25000 RNIC_IP_Address`
- 4) This must be done on both devices at the same time.
- 5) On all other pairs generate RDMA WRITE Operations of size 1 byte continuously until the control stream completes.
- 6) Repeat above steps until all port pairs are tested.
- 7) Repeat the above steps with all endpoint pairs, except the control stream changed such that the size of the RDMA WRITE operation is 1,000,000 bytes (~1 MB)

13.7 TI MPI - OPEN MPI USING OFED

The following values are used in examples below:

- \$MPIHOME: The absolute directory location of the Open MPI installation that is common to all systems under test.
- \$NP: The number of MPI processes to use in the test.
- \$HOSTFILE: The absolute filename location of the hostfile
- \$IMBHOME: The absolute directory location of the Intel MPI Benchmark (IMB) tools installation that is common to all systems under test.

13.7.1 CLUSTER SETUP

- 1) Network configuration requirements
 - a) All systems must be reachable by each other over IPoIB.
 - b) All nodes must agree on the IPoIB IP addresses of all systems (e.g., via /etc/hosts, DNS, or some other mechanism).
- 2) The same version of OFED must be installed in the same filesystem location on all systems under test.
- 3) The same version of the Intel MPI Benchmark (IMB) tools must be installed in the same filesystem location on all systems under test.
 - a) IMB can be used from the OFED installation or, if a later version of Open MPI is to be used, IMB can be downloaded from Intel's web site:
<http://software.intel.com/en-us/articles/intel-mpi-benchmarks/?wapkw=intel%20mpi%20benchmarks>
- 4) The same version of Open MPI must be available in the same filesystem location on all systems under test.
 - a) Open MPI can be used from the OFED installation, or, if a later version is required, can be downloaded and installed from the main Open MPI web site:
<http://www.open-mpi.org/>
 - i) If building Open MPI from source, and if the OpenFabrics libraries and headers are installed in a non-default location, be sure to use the --with-openib=<dir> option to configure to specify the OpenFabrics filesystem location.
 - ii) Open MPI can be installed once on a shared network filesystem that is available on all nodes, or can be individually installed on all systems. The main requirement is that Open MPI's filesystem location is the same on all systems under test.
 - iii) If Open MPI is built from source, the --prefix value given to configure should be the filesystem location that is common on all systems under test. For example, if installing to a network filesystem on the filesystem server, be sure to specify the filesystem location under the common mount point, not the "native" disk location that is only valid on the file server.

- iv) **Note** that Open MPI is included in some Linux distributions and other operating systems. Multiple versions of Open MPI can peacefully co-exist on a system as long as they are installed into separate file-system locations (i.e., configured with a different `--prefix` argument). All MPI tests must be built and run with a single installation of Open MPI.
- v) Ensure that the Open MPI installation includes OpenFabrics support:

```
shell$ $MPIHOME/bin/ompi_info | grep openib MCA btl: openib  
(MCA v1.0, API v1.0.1, Component v1.4)
```

The exact version numbers displayed will vary depending on your version of Open MPI. The important part is that a single "btl" line appears showing the openib component.
- b) Basic Open MPI run-time functionality can first be verified by running simple non-MPI applications. This ensures that the test user's rsh and/or ssh settings are correct, etc.

```
shell$ $MPIHOME/bin/mpirun -ssh -np $NP --hostfile $HOSTFILE host-  
name
```

 - i) The output should show the hostname of each host listed in the hostfile; the hostname should appear as many times as there are lines in the hostfile. The list of hostnames may appear in random order; this is normal
 - ii) Note that any serial application can be run; "hostname" is a good, short test that clearly identifies that specific hosts were used
- 5) All systems must be setup with at least one identical user account. This user must be able to SSH or RSH to all systems under test from the system that will launch the Open MPI tests with no additional output to stdout or stderr (e.g., all SSH host keys should already be cached, no password/passphrase prompts should be emitted, etc.).
- 6) The lockable memory limits on each machine should be set to allow unlimited locked memory per process.
- 7) The underlying OpenFabrics network used in the test should be stable and reliable.
- 8) No other fabric interoperability tests should be running during the Open MPI tests.
- 9) MPI tests should be run across at least 5 separate systems to force the use of the OpenFabrics network (vs. using just shared memory for in-system communication).

13.7.2 TEST SETUP

- 1) Create a hostfile (\$HOSTFILE) listing the hostname of each system that will be used in the test. If a system under test can run more than one MPI process (such as multiprocessor or multicore systems), list the hostname as many times as MPI processes are desired. For example, for two systems named node1.example.com and node2.example.com that are each able to run 4 processes:

```
shell$ cat hostfile.txt
```

node1.example.com 1
node1.example.com 2
node1.example.com 3
node1.example.com 4
node2.example.com 5
node2.example.com 6
node2.example.com 7
node2.example.com 8
node2.example.com 9

- 2) Determine the number of Open MPI processes (\$NP) that are to be run determined by the number of host entries in the created hostfile. 11
12
- 3) Open MPI defaults to probing all available networks at run-time to determine which to use. OpenFabrics testing must specifically force Open MPI to *only* use its OpenFabrics stack for testing purposes (e.g., do not fail over to TCP if the OpenFabrics stack is unavailable). To do this add an extra command line parameter; both iWarp and InfiniBand: 13
14
15
16
17
18
- 4) It has been discovered that the following Open MPI command line parameter is required to facilitate multi RDMA adaptor vendor MPI rings; both iWarp and InfiniBand: 19
20
21
22
- 5) It has been discovered that the following Open MPI command line parameter is required to facilitate multi RNIC adaptor vendors MPI rings; iWarp specific: 23
24
25
26
- mca btl openib,self
--mca pml ob1 --mca btl_openib_flags 306
--mca btl_openib_receive_queues P,65536,256,192,128

13.7.3 TEST PROCEDURE

- 1) Create a hostfile listing the MPI ring nodes, process distribution, and total number of processes to use as indicated in steps 1 and 2 of section 12.11.2. The filesystem location of this hostfile is irrelevant. 27
28
29
- 2) Locate the "mpirun" binary that will be used. This determines the version of Open MPI that will be used. 30
31
- 3) Locate the "IMB-MPI1" IMB binary. This must have been built against the version of Open MPI selected above. If using an OFED distribution this build process has already been performed. 32
33
34
- 4) Verify that a subnet manager has configured the fabric. If not, start one. 35
- 5) Verify that all hosts present within the hostfile are online and accessible. 36
- 6) Run the IMB-MPI1 benchmarks 37
- 7) Repeat steps 4-6 using a different subnet manager until all subnet managers under test have been used. 38
39
- 8) All IMB benchmarks must pass successfully using all subnet managers under test in order for the devices under test defined within the hostfile pass. 40
41
42

13.7.4 METHOD OF IMPLEMENTATION FOR ALL LINUX OS'S

- 1) To perform step 4 of section 12.9.3 use "ibdiagnet -r" from a host defined in the mpi hostfile and look for an "SM - Master" entry in the output
- 2) To perform step 5 of section 12.9.3 ping the IPoIB address of all hosts defined in the mpi hostfile from a host defined in said hostfile.
- 3) To perform step 6 of section 12.9.3 use the following command from a host that can access all hosts defined within the hostfile; this host can be part of the hostfile
 - a) For **InfiniBand & RoCE**:
\$MPIHOME/bin/mpirun --mca btl openib,self,sm --mca pml ob1 -mca btl_openib_flags \ 306 -np \$NP -hostfile \$HOSTFILE \$IMBHOME/IMB-MPI1
 - a) For **iWarp**:
\$MPIHOME/bin/mpirun --mca btl openib,self,sm --mca pml ob1 --mca \ btl_openib_flags 306 --mca btl_openib_receive_queues P,65536,256,192,128 -np \ \$NP -hostfile \$HOSTFILE \$IMBHOME/IMB-MPI1

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42

13.8 TI MPI - OHIO STATE UNIVERSITY USING OFED

13.8.1 MVAPICH - SETUP

- 1) Network configuration requirements
 - a) All systems must be reachable by each other a common network that supports TCP (Ethernet, IPoIB, etc.)
 - b) All nodes must agree on the IP addresses for all TCP networks on all systems (e.g., via /etc/hosts, DNS, or some other mechanism).
 - 2) The same version of OFED must be installed in the same filesystem location on all systems under test.
 - 3) MVAPICH is included in OFED distributions. The updated versions of MVAPICH can be obtained from OpenFabrics website.
 - 4) The same version of MVAPICH must be available in the same filesystem location on all systems under test.
 - a) MVAPICH can be installed once on a shared network filesystem that is available on all nodes, or can be individually installed on all systems. The main requirement is that MVAPICH filesystem location is the same on all systems under test.
 - 5) All systems must be setup with at least one identical user account. This user must be able to SSH or RSH to all systems under test from the system that will launch the MVAPICH tests with no additional output to stdout or stderr (e.g., all SSH host keys should already be cached, no password/passphrase prompts should be emitted, etc.).
 - 6) The lockable memory limits on each machine should be set to allow unlimited locked memory per process. This can be achieved by using ulimit command.
 - 7) The underlying IB network(s) used in the test should be stable and reliable. No other fabric interoperability tests should be running during the MVAPICH tests.
 - 8) Multiple versions of MVAPICH can peacefully co-exist on a system as long as they are installed into separate filesystem locations (i.e., configured with a different --prefix argument). All tests must be built and run with a single installation of MVAPICH.
 - 9) MVAPICH tests should be run across at least 5 separate systems to force the use of the IB networks (vs. using just shared memory for in-system communication).
- Note:** MVAPICH is commonly referred to as MVAPICH1 to distinguish it from the new and updated MVAPICH2

13.8.2 MVAPICH - TEST SETUP AND PROCEDURE

- 1) Test Setup
 - a) Create a hostfile listing the hostname of each system that will be used in the test. If a system under test can run more than one MPI process (such as multiprocessor or multicore systems) list the hostname as many times as MPI processes are desired. For example, for two 2 processor systems named host1 and host2

- ```
$ cat hostfile.txt
```
- host1  
host1  
host2  
host2
- b) Download and install Intel® MPI Benchmarks on all nodes from:  
<http://www.intel.com/cd/software/products/asm-na/eng/cluster/mpi/219848.htm>  
Follow the instructions below to install:
- i) untar downloaded archive  
ii) open <natured directory>/src/make\_mpich and fill in the following variables:
- MPI\_HOME=<path to mvapich1 directory> #mine was /usr/mpi/gcc/mvapich-1.0.1
  - CPPFLAGS= -DCHECK
- iii) gmake -f make\_mpich  
This will install the benchmarks inside the MPI\_HOME/tests directory  
**Note:** Intel® MPI Benchmarks are installed with OFED installation by default
- c) Enter all nodes and run the following commands:
- i) echo "PATH=\\$PATH:<path to mvapich1 directory>/bin:<path to mvapich1 directory>/tests/IMB-3.0" >> /<username>/.bashrc # or .cshrc  
ii) echo "ulimit -l unlimited" >> /<username>/.bashrc # or .cshrc  
iii) source /<username>/.bashrc # or .cshrc  
**Note:** these commands may fail or produce unexpected results with a shared \$HOME
- 2) Testing Procedure
- a) The following values are used in the examples below
- i) \$MPIHOME - The absolute directory location of the MVAPICH installation that is common to all systems under test  
ii) \$NP - The number of MPI processes that are to be run determined by the number of host entries in the created hostfile.  
iii) \$HOSTFILE - The absolute location of the hostfile
- b) Run Intel® MPI Benchmarks:
- i) Run the PingPong and PingPing point-to-point tests  
\$MPIHOME/bin/mpirun\_rsh -ssh -np \$NP IMB-MPI1 -multi 0 PingPong PingPing -hostfile \$HOSTFILE  
ii) Run all the tests (PingPong, PingPing, Sendrecv, Exchange, Bcast, Allgather, Allgatherv, Alltoall, Reduce, Reduce\_scatter, Allreduce, Barrier), in non-multi mode.  
\$MPIHOME/bin/mpirun\_rsh -ssh -np \$NP IMB-MPI1 -multi 0 -hostfile \$HOSTFILE

### 13.8.3 MVAPICH2 - SETUP

- 1) Download and install OFED on all nodes from:  
<http://www.openfabrics.org/downloads/OFED>
- 2) Download and install Intel® MPI Benchmarks on all nodes from:  
<http://www.intel.com/cd/software/products/asmo-na/eng/cluster/mpi/219848.htm>  
You will have to accept a license. Follow the instructions below to install.
  - a) untar downloaded archive
  - b) open <untarred directory>/src/make\_mpich and fill in the following variables:
    - i) MPI\_HOME=<path to mvapich2 directory> #mine was /usr/mpi/gcc/mvapich2-1.0.3
    - ii) CPPFLAGS= -DCHECK
  - c) gmake -f make\_mpich  
This will install the benchmarks inside the MPI\_HOME/tests directory
- 3) All nodes should be physically connected.
- 4) Enter all nodes and run the following cmds:
  - a) echo "PATH=\$PATH:<path to mvapich2 directory>/bin:<path to mvapich2 directory>/tests/IMB-3.0" >> /<username>/.bashrc # or .cshrc
  - b) echo "ulimit -l unlimited" >> /<username>/.bashrc;
  - c) source /<username>/.bashrc # or .cshrc
- 5) Create an mpi ring:
  - a) Construct a file called hosts that has the following format. Include as many lines as you have hosts. Be sure to leave a blank line at the end of the file:
    - i) <host>ifhn=<infiniband ip address>
  - b) Run the following commands
    - i) mpdboot -n `cat hosts|wc -l` -f hosts --ifhn=<localhost infiniband ip address>
    - ii) mpdtrace -l #OPTIONAL, shows current ring members.
- 6) MVAPICH tests should be run across at least 5 separate systems to force the use of the IB networks (vs. using just shared memory for in-system communication).

### 13.8.4 MVAPICH2 - TEST PROCEDURE

#### Step A:

[For IB] Run a subnet manager from one node only.

#### Step B

Run Intel® MPI Benchmarks:

- 1) Two sets of tests should be run, with these command lines

[For IB]

a) mpirun\_rsh -ssh -np <number of nodes X number of processors/node>  
IMB-MPI1 -multi 0 PingPong PingPing

b) mpirun\_rsh -ssh -np <number of nodes X number of processors/node>  
IMB-MPI1

**[For iWARP]**

a) mpirun\_rsh -ssh -np <number of nodes X number of processors/node>  
MV2\_USE\_IWARP\_MODE=1 MV2\_USE\_RDMA\_CM=1 IMB-MPI1 -  
multi 0 PingPong PingPing

b) mpirun\_rsh -ssh -np <number of nodes X number of processors/node>  
MV2\_USE\_IWARP\_MODE=1 MV2\_USE\_RDMA\_CM=1 IMB-MPI1

The first command runs just the PingPong and PingPing point-to-point tests,  
but makes all tasks active (pairwise).

The second command runs all the tests (PingPong, PingPing, Sendrecv, Ex-  
change, Bcast, Allgather, Allgather, Alltoall, Reduce, Reduce\_scatter, Allre-  
duce, Barrier), in non-multi mode.

2) **[For IB]** If the test passes shutdown current subnet manager and start an-  
other one on a different node; run both tests again.

3) **[For IB]** Repeat until all nodes have run a subnet manager and passed all  
tests.

## 14 INFINIBAND SPECIFIC INTEROP PROCEDURES USING WINOF

### 14.1 IB LINK INITIALIZE USING WINOF

#### 14.1.1 Setup

**Note:** The WinOF Subnet Manager and diagnostics are still evolving as compared to OFED. Therefore, you must include an OFED Linux node along with the Win

- 1) Disconnect the full topology and select a cable whose length should be a maximum of 15 meters for SDR and 10 meters for DDR when using copper cables. OF node to run diagnostics for this test.
- 2) Verify that no SM is running
- 3) Connect two devices back to back
- 4) ssh to the OFED node.
  - a) Run "ibdiagnet -lw 4x" to verify portwidth
  - b) Run "ibdiagnet -ls 2.5" to check link speed. Interpret output and compare to advertised speed.

**Note:** This command will only produce output if the link speed is anything other than SDR. Keep this in mind during your interpretation of the output.

- 5) Repeat steps 1-3 with a different device pairing.
  - a) All device pairs must be tested except SRP target to SRP target.
    - i) HCA to HCA
    - ii) HCA to Switch
    - iii) HCA to Target
    - iv) Switch to Switch
    - v) Switch to Target

**Note:** HCA to Target and HCA to HCA cannot be tested under WinOF 2.0.2 because there are no utilities available. Switches can be tested by using a Linux Host and the OFED Utilities.
  - b) Each device must link to all other devices in order for the device to pass link init over all.

#### 14.1.2 Recommendations

In order to determine Switch to Target and Switch to Switch link parameters, run commands from an HCA linked to the switch under test. This does require more interpretation of the output to differentiate the reported parameters.

## 14.2 IB FABRIC INITIALIZATION USING WINOF

### 14.2.1 Architect the Network we want to build.

**Note:** The WinOF Subnet Manager and diagnostics are still evolving as compared to OFED. Therefore, you must include an OFED Linux node along with the WinOF node to run diagnostics for this test.

- 1) Design and implement a Cluster Topology.
- 2) End to end IPoIB connectivity is required between all end points. Therefore you must create and assign IP addresses to each IB end point.
- 3) See [Figure 5- Sample Network Configuration](#) below.

### 14.2.2 Procedure

- 1) Connect the HCAs and switches as per the Architected Network and make sure that no SM/SA is running on the Fabric.
- 2) Start an SM on a device and let it initialize (all SMs will need to be tested)
- 3) Visually verify that all devices are in the active state using LEDs (however the vendor decided to implement it).
- 4) The following steps must be done using a Linux OFED end point.
  - a) Run "ibdiagnet -pc" to clear all port counters
  - b) Wait 17 seconds as per the specifications requirements.
  - c) Run "ibdiagnet -c 1000" to send 1000 node descriptions.
  - d) Run "ibdiagnet" to generate fabric report and open report to see results. /tmp/ibdiagnet.sm
  - e) Run "ibchecknet" to build guid list.

### 14.2.3 Verification Procedures

- 1) Review "PM Counters" section of the fabric report. There should be no illegal PM counters. The Specification says there should be no errors in 17 seconds.
- 2) Review "Subnet Manager " section of the fabric report. Verify that the running SM is the one you started and verify number of nodes and switches in the fabric.
- 3) Review the ibchecknet report and verify that there are no duplicate GUIDs in the fabric

**Note:** the reports are located in the /tmp directory

Restart all devices in the fabric and follow Sections 13.2.2 and 13.2.3. Run the SM from a different device in the fabric until all SMs present have been used. All SMs on managed switches and one instance of **opensm** must be used.

Each device must pass all verification procedures with every SM to pass Fabric Initialization test.

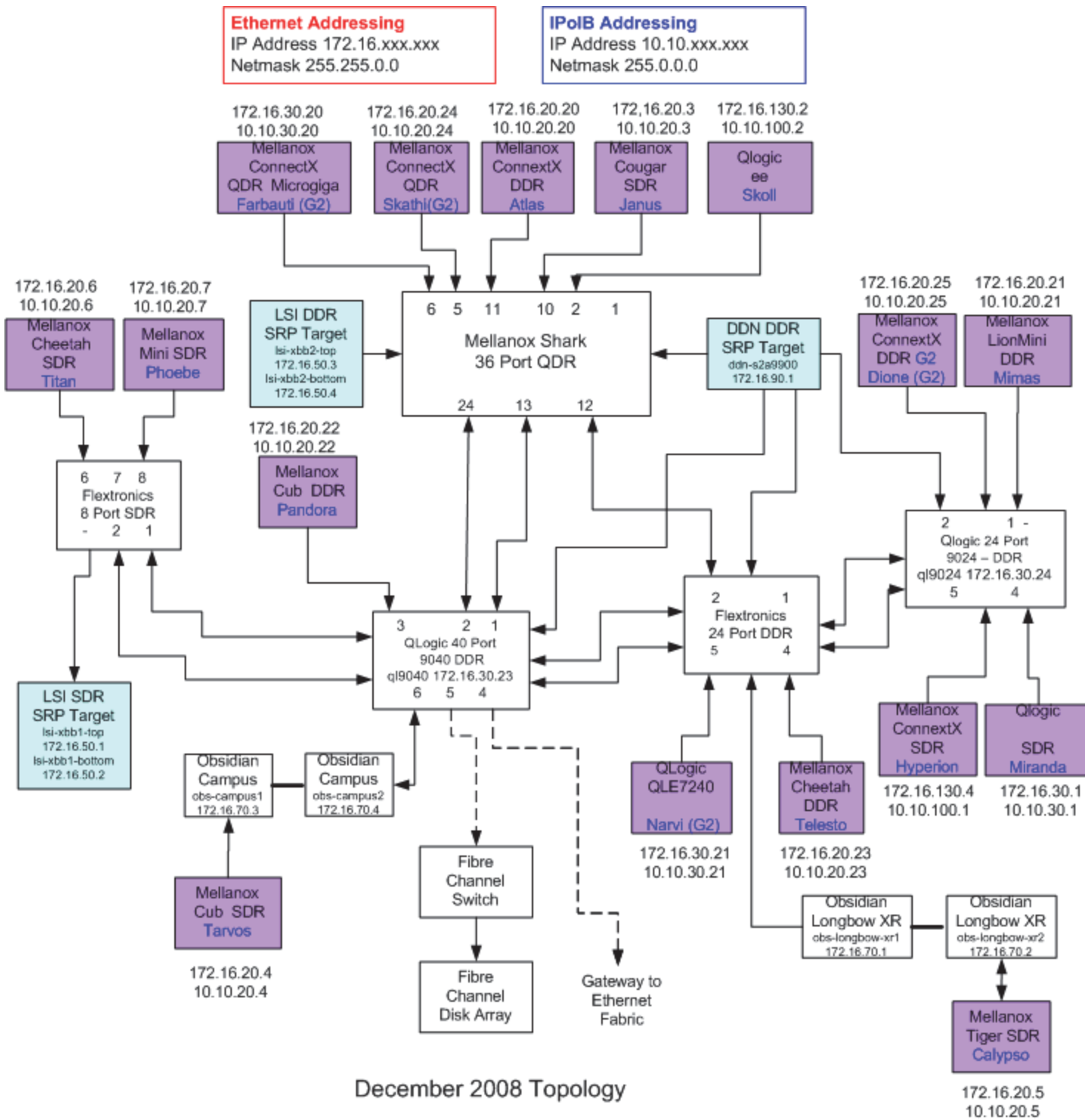
**Table 28 - ibdiagnet commands**

| Commands                  | Description                                |
|---------------------------|--------------------------------------------|
| ibdiagnet -c 1000         | send 1000 Node Descriptions                |
| ibdiagnet -h              | Help                                       |
| ibdiagnet -lw 4x - ls 2.5 | Specify link width and speed               |
| ibdiagnet - pc            | Clear Counter                              |
| ibdiagnet -t <file>       | Compare current topology to saved topology |
| ibdiagnet -wt             | Writes the topology to a file              |

**Note:** The topology file is being generated after the SM starts but before any testing has started. The topology comparison is being performed after testing has been completed but before the systems get rebooted. A topology check is performed during every part of every test section that does not specifically state "change the topology". For example Fabric Init only has 1 part so there is only 1 check but RDS has 2 parts so 2 checks are performed. However, IPoIB has 3 parts for each of 2 modes but 1 of those parts specifically says to change the topology so only 4 checks occur.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

Figure 5 - Sample Network Configuration



## 14.3 IB IPoIB DATAGRAM MODE (DM) USING WINOF

### 14.3.1 SETUP

**Note:** WinOF 2.0.2 only supports IPoIB Datagram Mode. Future WinOF releases will support IPoIB Connected-Mode.

Connect the HCAs and switches as per the Architected Network and make sure that no SM is running on the Fabric.

This procedure, as the previous ones, will be based on the cluster connectivity. An SM/SA which supports IPoIB (sufficient IB multicast support) will be running on the HCAs, or on a switch with an embedded SM/SA or a third HCA which would only run SM/SA for the partner pair (with a switch in the middle). This procedure has been developed for the Windows environment.

**Optional:** In the procedures below, an IB analyzer can be inserted in the appropriate link to obtain traces and validate the aspects of the procedures specifically detailed below in subsequent sections.

### 14.3.2 IPOIB INTERFACE CREATION AND IPOIB SUBNET CREATION

- 1) Configure IPoIB address. All addresses must reside on the same subnet.
- 2) Verify which 'Local Area Connection' the IPoIB interfaces are bound to:
  - a) Start | Server Manager | View Network Connections.
  - b) Find the OpenFabrics IPoIB interfaces (one per HCA port). If your platform has two Ethernet ports, then IPoIB interfaces likely will be assigned '**Local Area Connection 3**' & '**Local Area Connection 4**' as the Ethernet ports are assigned '**Local Area Connection**' and '**Local Area Connection 2**' .
- 3) Set interfaces to 10.0.0.x/24 (10.0.0.x/netmask 255.255.255.0) using the following commands:
  - a) netsh interface ip set address "Local Area Connection 3" static 10.10.4.x 255.255.255.0
  - b) netsh interface ip set address "Local Area Connection 4" static 10.10.4.y 255.255.255.0
- 4) View the IPoIB IP address using the following command
  - a) netsh interface ip show address "Local Area Connection 3"

### 14.3.3 PING PROCEDURES

#### Step A

- 1) Stop all SM's and verify that none are running
- 2) Power cycle all switches in the fabric (this insures that the new SM will configure all the links and create the multi-cast join).
- 3) Start an SM (All SM's will need to be tested) and let it initialize

**Note:** For link testing it is recommended to use an OFED Linux OpenSM as the Windows version of OpenSM does not support all SA queries and functionality of the OFED 1.4 OpenSM.



**Note:** All WinOF installed systems contain a disabled OpenSM windows service. A WinOF installation option/feature is to automatically 'start/enable' the OpenSM service on the local node.

- Start | Server Manager | Configuration | Services | InfiniBand Subnet Manager | Automatic | apply
- Start | Apply will enable the local OpenSM to start and be started upon system boot.
- a) Visually verify that all devices are in the active state. Orange led will be on if the port is active.
- b) From a Linux system, Run "ibdiagnet" and verify that the SM you started is the one that is running and and that it is the master. You will need to know the GUID of the device since the SM will be reassigned on each reboot; the Windows 'vstat' command displays HCA info.
- c) Verify that all nodes and switches were discovered.
- d) WinOF 2.0.2 does not provide a ibdiagnet utility.

**Note:** Ibdagnet may show more switches than indicated by the physical number of switch platforms present. This is because some switches have multiple switch chips.

- 4) Examine the arp table (via arp -a) and remove the destination node's ib0 address from the sending node's arp table (via arp -d).
- 5) Issue the command: sysctl net.ipv4.neigh.ib0.unres\_qlen=18
  - a) This sets the qlen variable to 18 which increases the buffer size so that you do not get an initial dropped packet when using ping sizes 8192 and greater.
- 6) Ping every IPoB interface IPv4 address except localhost with packet sizes of 511, 1025, 2044, 8192, 32768 and 65500. 'ping /?' displays ping help.
  - a) 10 packets of each size will be sent
  - b) Every packet size is a new ping command.

**Note:** Windows does not support 65507 so we used 65500.

**Note:** This is done from the Head Node utility "Run a Command" using the following command:

```
for %i in (64, 511, 2044, 8192, 32768 and 65500) DO %d arp -d %d & ping -i 1 -n 10 -l %i %d & arp -d %d
```
- 7) In order to pass Step A, a reply must be received for every ping sent (without losing a single packet) while using each one of the SMs available in the cluster.

## Step B

- 1) Bring up all HCAs but one.
- 2) Start an SM (all SMs will need to be tested).
- 3) Check for ping response between all node (All to All).
  - a) A response from the disconnected HCA should not be returned.
- 4) Disconnect one more HCA from the cluster.

- 5) Ping to the newly disconnected HCA from all nodes (No response should be returned). 1
  - 6) Connect the first machine (the one that was not connected) and check for ping response from all nodes that are still connected. 2
  - 7) Connect the disconnected HCA to a different switch on the subnet which will change the topology. 3
  - 8) Ping again from all nodes (this time we should get a response). 4
  - 9) Follow Step B, this time bring the interface down and then back up: Start | Server Manager | View Network Connections | IPoB(Local Area connection) disable and enable commands instead of physically disconnecting the HCAs. 5
- Note:** Each step must exhibit the expected behavior while using each SM in order for the device to pass Step B overall. 6

### Step C

- 1) Follow Step A and B using a different SM until all SM's have been used. Only one instance of each available SM is required. Steps A, B, and C must pass in order for the device to pass 13.3.3 overall. 7
- 2) Issue the command: `sysctl net.ipv4.neigh.ib0.unres_qlen=3` 8
- a) This sets the qlen variable back to the default. 9

### 14.3.4 FTP PROCEDURE

FTP procedures requires an FTP server to be configured on each machine in the partner pair. An FTP client needs to be available on each machine as well; an FTP client is a standard Windows component. 10

An FTP server is a component of the IIS '**Internet Information Services**' manger which **not** a part of a standard Windows installation: 11

See Start | Server Manager | Roles | Add IIS. Configure FTP server via IIS manager. 12

#### 14.3.4.1 SETUP

- 1) Make sure ftpd is installed on each node for the FTP application. 13
- 2) A special account for this should be created as follows: 14
- b) Username: Interop 15
- c) Password: openfabrics 16

#### 14.3.4.2 PROCEDURE

Run FTP server on all nodes. 17

- 1) Start an SM (all SMs will need to be tested) and let it initialize (ref MS Network utilities docs) 18
- a) Verify that the running SM is the one you started. 19
- 2) FTP: 20

- a) Connect an HCA pair via FTP on IPoIB using the specified user name and password. 1
  - b) Put the 4MB file to the %windir%\temp folder (generally C:\Windows\Temp) on the remote host. 2
  - c) Get the same file to your local dir again. 3
  - d) Binary compare the file using the Windows command 'fc /B tfile tfile.orig'. 4
  - i) The two must be identical 5
- 3) Repeat the procedure with a different SM. 6
- Note:** Every node must FTP the 4MB file to all others using all SMs and the files must be identical as determined by the binary compare in order for the device to pass 13.3.4 overall. 7
- Note:** Sections 13.3.3 and 13.3.4 must pass using the configuration determined by sections 13.3.1 and 13.3.2 for the device to pass IPoIB Datagram mode overall. 8

## 14.4 IB SM FAILOVER AND HANDOVER PROCEDURE USING WINOF

### 14.4.1 SETUP

- 1) Connect HCAs per the selected topology.
- 2) In this test, all active SMs on the fabric which are going to be tested, must be from the same vendor. They will be tested pairwise: two at a time.

### 14.4.2 PROCEDURE

- 1) Disable all SMs in the cluster.
- 2) Start a SM on either machine in a chosen pair.
  - a) Start | Server Manager | Configuration | Services | InfiniBand Subnet Manager | start | apply
- 3) Run "vstat" on all Windows nodes in the fabric.
  - a) Verify HCA link active in vstat output.
- 4) Verify IPoIB is active on each node
  - a) Verify Local Area Connection assigned to IPoIB interface:
    - i) Start | Control Panel | Network and Sharing Center | Manage Network Connections.
    - b) Show IPv4 address assigned to IPoIB Interface(s):
      - i) netsh interface ip show address "Local Area Connection 3"
      - ii) netsh interface ip show address "Local Area Connection 4"
    - c) Verify the IPoIB devices (one per cabled connected HCA port) are visible & operational from a device driver perspective using Device Manager
      - i) Start | Run | devmgmt.msc
    - d) Ping the IPoIB interface IPv4 address local and remote, verify traffic is actually going in/out over IPoIB 'local area connection x'.
- 5) Start an Open SM on the second machine in the current pair.
- 6) Verify that the SMs behave according to the SM priority rules.
  - a) The Windows OpenSM log file is located at '%windir%\temp\osm.log'.

**Note:** The SM with highest numerical priority value is master and the other is in standby. If both SMs have the same priority value then the SM with the smallest guid is master and the other is in standby.
- 7) Verify that all nodes in the cluster are present - ping all IPoIB interfaces
- 8) Shutdown the master SM.
- 9) Verify the other active SM goes into the master state: see osm.log file.
- 10) Verify that all nodes in the cluster are present - ping all IPoIB interfaces
- 11) Start the SM you just shutdown.
- 12) Verify that the newly started SM resumes it's position as master while the other goes into standby again; see '%windir%\temp\osm.log'.
- 13) Verify that all nodes in the cluster are present - ping all IPoIB interfaces

- 14) Shutdown the standby SM. 1
  - 15) Verify that the previous master SM is still the master; view 2  
'%windir%\temp\osm.log'. 3
  - 16) Verify that all nodes in the cluster are present - ping all IPoIB interfaces 4
  - 17) Repeat proceeding steps [1-16] 2 more times with the same node pair, en- 5  
suring that the below criteria is met (total of 3 tests per pair which can be run 6  
in any order): 7
    - a) First SM to be started having highest numerical priority value. 8
    - b) Second SM to be started having highest numerical priority value. 9
    - c) Both SMs having equal numerical priority values. 10
  - 18) Repeat steps 1-17 until all possible SM pairs from identical vendors in the 11  
cluster have been tested. 12
- 13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

## 14.5 IB SRP USING WINOF

### 14.5.1 SETUP

- 1) Connect the HCAs and switches as per the Architected Network and make sure that no SM is running on the Fabric.
- 2) Configure and Start a Linux OFED SRP target - VDISK BLOCKIO mode; (some assembly required) - <https://wiki.openfabrics.org/tiki-index.php?page=SRPT+Installation>
  - a) assume /dev/sdb1 & /dev/sdc1 are formatted with /sbin/mkfs.msdos
  - b) Setting SRPT\_LOAD=yes in /etc/infiniband/openib.conf is not good enough. It only loads ib\_srpt module and does not load scst and its dev\_handlers.
  - c) modprobe scst
  - d) modprobe scst\_vdisk
  - e) echo "open vdisk0 /dev/sdb BLOCKIO" > /proc/scsi\_tgt/vdisk/vdisk
  - f) echo "open vdisk1 /dev/sdc BLOCKIO" > /proc/scsi\_tgt/vdisk/vdisk
  - g) echo "add vdisk0 0" > /proc/scsi\_tgt/groups/Default/devices
  - h) echo "add vdisk1 1" > /proc/scsi\_tgt/groups/Default/devices

**Note:** For the April 2012 Interop events, the OFILG decided that each target should only advertise two volumes in order to reduce the amount of time required to run the tests

### 14.5.2 WINDOWS PROCEDURE

- 1) Start an SM (all SM's will need to be tested) and let it initialize.
  - a) Verify that the running SM is the one that you started
- 2) Choose a node to work with
- 3) Verify the SRP driver loaded correctly; locate the SRP Miniport.
  - a) Start | Control Panel | Device Manager | Storage Controllers [InfiniBand SRP Miniport]
- 4) Discover + Enable (bring online) the SRP drive(s)
  - a) Start | Server Manager | Storage | Disk Management
- 5) You will find a basic 'unknown' and 'offline' disk; this one of your SRP volume(s).
- 6) Right-click the offline disk and select 'online'.
- 7) Right-click the volume space, assign the drive letter 'T'.
- 8) Right-click the volume space, format the volume.
- 9) Access the SRP drive via assigned drive letter. From a Windows/DOS command prompt window, execute the following commands.
  - a) vol T:
  - b) dir T:\ (should be empty)

- c) mkdir T:\tmp 1
  - d) copy /B WinOF\_wlh\_x64.msi T:\tmp 2
  - e) fc /B WinOF\_wlh\_x64.msi T:\tmp\WinOF\_wlh\_x64.msi 3
  - f) copy /B T:\tmp\WinOF\_wlh\_x64.msi T:\tmp\WOF2.msi 4
  - g) fc /B T:\tmp\WinOF\_wlh\_x64.msi T:\tmp\WOF2.msi 5
  - h) fc /B WinOF\_wlh\_x64.msi T:\tmp\WOF2.msi 6
  - i) copy /B T:\tmp\WOF2.msi WOF3.msi 7
  - j) fc /B WinOF\_wlh\_x64.msi WOF3.msi 8
  - k) del T:\tmp\WOF2.msi 9
  - l) del T:\tmp\WinOF\_wlh\_x64.msi 10
  - m) dir T:\tmp (should be empty) 11
  - n) rmdir T:\tmp 12
  - o) dir T:\ (should be empty) 13
  - p) del WOF3.msi 14
- 10) For each SRP target located in Procedure #4 15
- a) Perform step 9 for each volume found for all targets as determined by 16  
Windows Procedure step #4 - see [Discover + Enable \(bring online\) the 17](#)  
[SRP drive\(s\)](#) 18
- 11) Take SRP drive offline 19
- a) Start | Server Manager | Storage | Disk Management 20
  - b) Right-click the online disk and select 'offline' 21
  - c) dir T:\ (should fail). 22
- 12) Reboot all devices in the fabric and repeat the procedure using a different 23  
SM. 24
- Note:** An HCA must successfully complete all operations to and from all volumes 25  
on all targets using all available SM's in order to pass SRP testing. Two volumes 26  
per target are all that is required. 27

## 14.6 IB uDAPLTEST COMMANDS USING WINOF

Server Command: `dapl2test -T S -D <ia_name>`

### 14.6.1 IB SETUP

- The `%SystemDrive%\DAT\dat.conf` needs to be verified to be sure that the correct interface is used. The DAPL interface for IB is `ibnic0v2`.
- It is also important to verify that the desired `dat/dapl` libraries are available
  - `%windir%\dat2.dll`
  - `%windir%\dapl2.dll`
- To run `dapl2test` on IB, an SM needs to be running.

### 14.6.2 GROUP 1: POINT-TO-POINT TOPOLOGY

[1.3] 1 connection and simple send/recv:

- `dapl2test -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -R BE`
- client SR 256 1 server SR 256 1

[1.4] Verification, polling, and scatter gather list:

- `dapl2test -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 1024 3 -f \
- server SR 1536 2 -f

### 14.6.3 GROUP 2: SWITCHED TOPOLOGY

InfiniBand Switch: Any InfiniBand switch

[2.5] Verification and private data:

- `dapl2test -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 1024 1 \
- server SR 1024 1

[2.6] Add multiple endpoints, polling, and scatter gather list:

- `dapl2test -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 10 -V -P -R BE`
- client SR 1024 3 \
- server SR 1536 2

[2.7] Add RDMA Write :

- `dapl2test -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 256 1 \
- server RW 4096 1 server SR 256 1

[2.8] Add RDMA Read:

- `dapl2test -T T -s <server_name> -D <ia_name> -i 100 -t 1 -w 1 -V -P -R BE`
- client SR 256 1 \



- server RR 4096 1 server SR 256 1

#### 14.6.4 GROUP 3: SWITCHED TOPOLOGY WITH MULTIPLE SWITCHES

[3.5] Multiple threads, RDMA Read, and RDMA Write:

- dapl2test -T T -s <server\_name> -D <ia\_name> -i 100 -t 4 -w 8 -V -P -R BE
- client SR 256 1 \
- server RR 4096 1 server SR 256 1 client SR 256 1 server RR 4096 1 \
- server SR 256 1

[3.6] Pipeline test with RDMA Write and scatter gather list:

- dapl2test -T P -s <server\_name> -D <ia\_name> -i 1024 -p 64 -m p RW 8192 2

[3.7] Pipeline with RDMA Read:

- dapl2test -T P -s <server\_name> -D <ia\_name> -i 1024 -p 64 -m p RR 4096 2

[3.8] Multiple switches:

- dapl2test -T T -s <server\_name> -D <ia\_name> -i 100 -t 1 -w 10 -V -P -R
- BE client SR 1024 3 \
- server SR 1536 2

#### 14.6.5 WINOF DAPL2TEST WRAPPER SCRIPTS

All the specified DAPL tests are conveniently located in the WinOF distributed DAPL test server & client scripts.

- %ProgramFiles(x86)%\WinOF\dt-svr.bat
  - To run the dapl2test Server, to a Windows cmd-prompt window type 'dt-svr'. Only one server is necessary – multiple clients can communicate with a single dapl2test server; multiple servers on different nodes can exist. A single dapl2test client communicates with only one dapl2test server at a time.
  - No further server action is required as the dapl2test server is persistent; looping waiting for dapltest client requests.
- %ProgramFiles(x86)%\WinOF\dt-cli.bat
  - 'dt-cli' no arguments, will display dt-cli command args & options.
  - Dapl2test client invocation: 'dt-cli IPoIB\_IPv4\_server\_address cmd'
  - If the dt-svr command was executed on a system where the IPoIB interface address is 10.10.4.200 then
  - 'dt-cli 10.10.4.200 interop' would run the above dap2tests between the client and server.
  - 'dt-cli 10.10.4.200 conn' is a simple, quick test to verify dapl2test client | server connection is operational.

## 14.7 IB MPI - INTEL MPI USING WINOF

### 14.7.1 Requirements

- 1) Intel MPI is not part of the WinOF installation; acquire Intel MPI installer file from Intel.
- 2) Install same O/S version (Windows Server 2008-HPC) on homogenous x86\_64 systems.
- 3) MPI testing requires a reliable IB fabric without other fabric interop testing occurring.
- 4) Private Ethernet Network configuration
  - a) DNS names must match hostnames in hosts file.
- 5) WinOF Installation requirements
  - a) Install the latest version of WinOF on all systems (double-click WinOF\_wlh\_x64.msi); see
    - i) <http://www.openfabrics.org/downloads/WinOF/README>
    - ii) Select the 'default' set of install features; includes uDAPL.
    - iii) Run OpenSM either on the headnode OR from one of the IB switches.
    - iv) If OpenSM on the headnode, select WinOF install feature 'OpenSM Started'.
  - b) Once WinOF installation on all nodes has completed, configure IPoIB interfaces.
    - i) %windir%\system32\Drivers\etc\hosts should be setup with IB hostnames and static IP addresses.
    - ii) Assign IPv4 address, from hosts file, to each IPoIB interface; Example: Local Area Connection 3 is the 1st IPoIB interface.
      - netsh interface ip set address "Local Area Connection 4" static 10.10.4.y 255.255.255.0  
This allows you to **set** the IPoIB IP address.
      - netsh interface ip show address "Local Area Connection 3"  
This allows you to **view** the IPoIB IP address.
    - iii) Verify by pinging IPoIB interface addresses on all nodes.

### 14.7.2 Setup information for Intel MPI

Install Intel MPI on every cluster node:

- 1) [Intel MPI runtime environment kit](#)
  - a) <http://www.intel.com/cd/software/products/asmo-na/eng/308295.htm>
- 2) [Intel MPI Benchmarks](#) ,
  - a) <http://www.intel.com/cd/software/products/asmo-na/eng/cluster/mpi/219848.htm>
- 3) Add identical user account (%SystemDrive%\users\test) on every node.

- 4) Headnode mount points (%SystemDrive%\test\export) on user accounts.

### 14.7.3 Additional Information

- 1) Go to the individual test directories and follow the steps in the respective README-\*.txt files.
- 2) For Intel MPI Support Services go to:
  - a) <http://software.intel.com/en-us/articles/intel-mpi-library-for-windows/all/1/>
  - b) See [Intel MPI Reference Manual](#) for Additional information

### 14.7.4 Intel MPI (MVAPICH 2) - Test Procedure

- 1) Run a subnet manager from one node only.
- 2) Run Intel® MPI Benchmarks from the HPC head-node:
  - a) Two sets of tests should be run, with these command lines
    - `mpiexec -np <number of nodes X number of processors/node> IMB-MPI1 -multi 0 PingPong PingPing`
    - `mpiexec -np <number of nodes X number of processors/node> IMB-MPI1`

The first command runs just the PingPong and PingPing point-to-point tests, but makes all tasks active (pairwise).

The second command runs all the tests (PingPong, PingPing, Sendrecv, Exchange, Bcast, Allgather, Allgatherv, Alltoall, Reduce, Reduce\_scatter, Allreduce, Barrier), in non-multi mode.
  - b) If the test passes shutdown current subnet manager and start another one on a different node; run both tests again.
- 3) Repeat until all nodes have run a subnet manager and passed all tests.

### 14.7.5 Interpreting the results

- 1) TBA

## 15 BUG REPORTING METHODOLOGY DURING PRE-TESTING

The following bug reporting methodology will be followed during the execution of interoperability pre-testing at UNH-IOL.

- 1) UNH-IOL and the OEMs (e.g. Chelsio, Data Direct, Intel, NetApp, Mellanox) will assign a focal point of contact to enable fast resolution of problems.
- 2) Bug reports will include:
  - a) Detailed fail report with all relevant detail (Test/Application, Topology.).
  - b) **[For IB]** IB trace if needed.
  - c) **[For iWARP]** iWARP, TCP and SCTP traces if needed.
- 3) Bug reports will be sent via email by UNH-IOL to the focal point assigned by the OEM
- 4) Bug reports and suggested fixes will be sent to the OpenFabrics development community - [OFA Bugzilla](#). When such reports are communicated, UNH-IOL will ensure that confidentiality between UNH-IOL and the OEM will be maintained. Bug reports will be generalized and not include any company specific proprietary information such as product name, software name, version etc.
- 5) All bug fixes/issues that are found during testing will be uploaded to the OpenFabrics repository. Documentation related to fixes will not mention any company specific proprietary information.

**Note:** This test plan does not cover how bugs will be reported by IBTA/CIWG or IETF iWARP during or after interoperability testing at plugfests.

## 16 RESULTS SUMMARY

### 16.1 INFINIBAND SPECIFIC TEST RESULTS

Please add a check mark whenever a test case passes and when the system is behaving according to the criteria mentioned below. Otherwise indicate a failure along with a comment explaining the nature of the failure.

**Results Table 1 - IB Link Initialize**

| Test # | Test                  | Pass | Fail | Comment |
|--------|-----------------------|------|------|---------|
| 1      | Phy link up all ports |      |      |         |

**Results Table 2 - IB Fabric Initialization**

| Test # | Test                                               | Pass | Fail | Comment |
|--------|----------------------------------------------------|------|------|---------|
| 1      | Verify that all ports are in Armed or Active state |      |      |         |

**Results Table 3 - IB IPoIB - Connected Mode (CM)**

| Test # | Test                              | Pass | Fail | Comment |
|--------|-----------------------------------|------|------|---------|
| 1      | Ping all to all - Ping using SM 1 |      |      |         |
| 2      | Ping all to all - Ping using SM 2 |      |      |         |
| 3      | Ping all to all - Ping using SM 3 |      |      |         |
| 4      | Ping all to all - Ping using SM 4 |      |      |         |
| 5      | Ping all to all - Ping using SM 5 |      |      |         |
| 6      | Ping all to all - Ping using SM 6 |      |      |         |
| 7      | Ping all to all - Ping using SM x |      |      |         |
| 8      | Connect/Disconnect Host           |      |      |         |
| 9      | FTP Procedure                     |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**Results Table 4 - IB IPoIB - Datagram Mode (DM)**

| Test # | Test                              | Pass | Fail | Comment |
|--------|-----------------------------------|------|------|---------|
| 1      | Ping all to all - Ping using SM 1 |      |      |         |
| 2      | Ping all to all - Ping using SM 2 |      |      |         |
| 3      | Ping all to all - Ping using SM 3 |      |      |         |
| 4      | Ping all to all - Ping using SM 4 |      |      |         |
| 5      | Ping all to all - Ping using SM 5 |      |      |         |
| 6      | Ping all to all - Ping using SM 6 |      |      |         |
| 7      | Ping all to all - Ping using SM x |      |      |         |
| 8      | Connect/Disconnect Host           |      |      |         |
| 9      | FTP Procedure                     |      |      |         |

**Table 5 - IB SM Failover/Handover**

| Test # | Test                        | Pass | Fail | Comment |
|--------|-----------------------------|------|------|---------|
| 1      | Basic sweep test            |      |      |         |
| 2      | SM Priority test            |      |      |         |
| 3      | Failover test - Disable SM1 |      |      |         |
| 4      | Failover test - Disable SM2 |      |      |         |

**Results Table 6 - IB SRP**

| Test # | Test                 | Pass | Fail | Comment |
|--------|----------------------|------|------|---------|
| 1      | Basic dd application |      |      |         |
| 2      | IB SM kill           |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**Results Table 7 - Fibre Channel Gateway - (IB Specific)**

| Test # | Test                                  | Pass | Fail | Comment |
|--------|---------------------------------------|------|------|---------|
| 1      | Basic Setup                           |      |      |         |
| 2      | Configure Gateway                     |      |      |         |
| 3      | Add Storage Device                    |      |      |         |
| 4      | Basic dd application                  |      |      |         |
| 5      | IB SM kill                            |      |      |         |
| 6      | Disconnect Host/Target                |      |      |         |
| 7      | Load Host/Target                      |      |      |         |
| 8      | dd after SRP Host and Target reloaded |      |      |         |
| 9      | Reboot Gateway                        |      |      |         |
| 10     | dd after FC Gateway reboot            |      |      |         |

**Results Table 8 - Ethernet Gateway - (IB Specific)**

| Test # | Test                    | Pass | Fail | Comment |
|--------|-------------------------|------|------|---------|
| 1      | Basic Setup             |      |      |         |
| 2      | Start ULP               |      |      |         |
| 3      | Discover Gateway        |      |      |         |
| 4      | SM Failover             |      |      |         |
| 5      | Ethernet gateway reboot |      |      |         |
| 6      | ULP restart             |      |      |         |
| 7      | Unload/load ULP         |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

## 16.2 ETHERNET SPECIFIC TEST RESULTS

Results Table 9 - iWARP Link Initialize

| Test # | Test                         | Pass | Fail | Comment |
|--------|------------------------------|------|------|---------|
| 1      | Phy link up all ports        |      |      |         |
| 2      | Verify basic IP connectivity |      |      |         |

Table 10 - RoCE Link Initialize

| Test # | Test                         | Pass | Fail | Comment |
|--------|------------------------------|------|------|---------|
| 1      | Phy link up all ports        |      |      |         |
| 2      | Verify basic IP connectivity |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42



### 16.3 TRANSPORT INDEPENDENT TEST RESULTS

**Results Table 11 - TI iSER**

| Test # | Test                          | Pass | Fail | Comment |
|--------|-------------------------------|------|------|---------|
| 1      | Basic dd application          |      |      |         |
| 2      | IB SM kill                    |      |      |         |
| 3      | Disconnect Initiator          |      |      |         |
| 4      | Disconnect Target             |      |      |         |
| 5      | Repeat with previous SM Slave |      |      |         |

**Results Table 12 - TI NFS Over RDMA**

| Test # | Test                         | Pass | Fail | Comment |
|--------|------------------------------|------|------|---------|
| 1      | File and directory creation  |      |      |         |
| 2      | File and directory removal   |      |      |         |
| 3      | Lookups across mount point   |      |      |         |
| 4      | Setattr, getattr, and lookup |      |      |         |
| 5      | Read and write               |      |      |         |
| 6      | Readdir                      |      |      |         |
| 7      | Link and rename              |      |      |         |
| 8      | Symlink and readlink         |      |      |         |
| 9      | Statfs                       |      |      |         |

**Results Table 13 - TI RDS**

| Test # | Test                 | Pass | Fail | Comment |
|--------|----------------------|------|------|---------|
| 1      | rds-ping procedure   |      |      |         |
| 2      | rds-stress procedure |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**Results Table 14 - TI uDAPL**

| Test # | Test                                                                       | Pass | Fail | Comment |
|--------|----------------------------------------------------------------------------|------|------|---------|
| 1      | P2P - Connection & simple send receive                                     |      |      |         |
| 2      | P2P - Verification, polling & scatter gather list                          |      |      |         |
| 3      | Switched Topology - Verification and private data                          |      |      |         |
| 4      | Switched Topology - Add multiple endpoints, polling, & scatter gather list |      |      |         |
| 5      | Switched Topology - Add RDMA Write                                         |      |      |         |
| 6      | Switched Topology - Add RDMA Read                                          |      |      |         |
| 7      | Multiple Switches - Multiple threads, RDMA Read, & RDMA Write              |      |      |         |
| 8      | Multiple Switches - Pipeline test with RDMA Write & scatter gather list    |      |      |         |
| 9      | Multiple Switches - Pipeline with RDMA Read                                |      |      |         |
| 10     | Multiple Switches - Multiple switches                                      |      |      |         |

**Results Table 15 - TI RDMA Basic Interop**

| Test # | Test              | Pass | Fail | Comment |
|--------|-------------------|------|------|---------|
| 1      | Small RDMA READ   |      |      |         |
| 2      | Large RDMA READ   |      |      |         |
| 3      | Small RDMA Write  |      |      |         |
| 4      | Large RDMA Write  |      |      |         |
| 5      | Small RDMA SEND   |      |      |         |
| 6      | Large RDMA SEND   |      |      |         |
| 7      | Small RDMA Verify |      |      |         |
| 8      | Large RDMA Verify |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**Results Table 16 - TI RDMA Stress Tests**

| Test # | Test          | Pass | Fail | Comment |
|--------|---------------|------|------|---------|
| 1      | Switch Load   |      |      |         |
| 2      | Switch Fan In |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

## 16.4 OPEN MPI TEST RESULTS

**Results Table 17 - TI MPI - Open MPI**

| Test #                        | Test Suite                                    | Pass | Fail | Comment |
|-------------------------------|-----------------------------------------------|------|------|---------|
| <b>Phase 1: "Short" tests</b> |                                               |      |      |         |
| 2                             | OMPI built with OpenFabrics support           |      |      |         |
| 3                             | OMPI basic functionality (hostname)           |      |      |         |
| 4.1                           | Simple MPI functionality (hello_c)            |      |      |         |
| 4.2                           | Simple MPI functionality (ring_c)             |      |      |         |
| 5                             | Point-to-point benchmark (NetPIPE)            |      |      |         |
| 6.1.1                         | Point-to-point benchmark (IMB PingPong multi) |      |      |         |
| 6.1.2                         | Point-to-point benchmark (IMB PingPing multi) |      |      |         |
| <b>Phase 2: "Long" tests</b>  |                                               |      |      |         |
| 6.2.1                         | Point-to-point benchmark (IMB PingPong)       |      |      |         |
| 6.2.2                         | Point-to-point benchmark (IMB PingPing)       |      |      |         |
| 6.2.3                         | Point-to-point benchmark (IMB Sendrecv)       |      |      |         |
| 6.2.4                         | Point-to-point benchmark (IMB Exchange)       |      |      |         |
| 6.2.5                         | Collective benchmark (IMB Bcast)              |      |      |         |
| 6.2.6                         | Collective benchmark (IMB Allgather)          |      |      |         |
| 6.2.7                         | Collective benchmark (IMB Allgatherv)         |      |      |         |
| 6.2.8                         | Collective benchmark (IMB Alltoall)           |      |      |         |
| 6.2.9                         | Collective benchmark (IMB Reduce)             |      |      |         |
| 6.2.10                        | Collective benchmark (IMB Reduce_scatter)     |      |      |         |
| 6.2.11                        | Collective benchmark (IMB Allreduce)          |      |      |         |
| 6.2.12                        | Collective benchmark (IMB Barrier)            |      |      |         |
| 6.3.1                         | I/O benchmark (IMB S_Write_Indv)              |      |      |         |
| 6.3.2                         | I/O benchmark (IMB S_IWrite_Indv)             |      |      |         |
| 6.3.3                         | I/O benchmark (IMB S_Write_Expl)              |      |      |         |
| 6.3.4                         | I/O benchmark (IMB S_IWrite_Expl)             |      |      |         |
| 6.3.5                         | I/O benchmark (IMB P_Write_Indv)              |      |      |         |
| 6.3.6                         | I/O benchmark (IMB P_IWrite_Indv)             |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**Results Table 17 - TI MPI - Open MPI**

| Test # | Test Suite                          | Pass | Fail | Comment |
|--------|-------------------------------------|------|------|---------|
| 6.3.7  | I/O benchmark (IMB P_Write_Shared)  |      |      |         |
| 6.3.8  | I/O benchmark (IMB P_IWrite_Shared) |      |      |         |
| 6.3.9  | I/O benchmark (IMB P_Write_Priv)    |      |      |         |
| 6.3.10 | I/O benchmark (IMB P_IWrite_Priv)   |      |      |         |
| 6.3.11 | I/O benchmark (IMB P_Write_Expl)    |      |      |         |
| 6.3.12 | I/O benchmark (IMB P_IWrite_Expl)   |      |      |         |
| 6.3.13 | I/O benchmark (IMB C_Write_Indv)    |      |      |         |
| 6.3.14 | I/O benchmark (IMB C_IWrite_Indv)   |      |      |         |
| 6.3.15 | I/O benchmark (IMB C_Write_Shared)  |      |      |         |
| 6.3.16 | I/O benchmark (IMB C_IWrite_Shared) |      |      |         |
| 6.3.17 | I/O benchmark (IMB C_Write_Expl)    |      |      |         |
| 6.3.18 | I/O benchmark (IMB C_IWrite_Expl)   |      |      |         |
| 6.3.19 | I/O benchmark (IMB S_Read_Indv)     |      |      |         |
| 6.3.20 | I/O benchmark (IMB S_IRead_Indv)    |      |      |         |
| 6.3.21 | I/O benchmark (IMB S_Read_Expl)     |      |      |         |
| 6.3.22 | I/O benchmark (IMB S_IRead_Expl)    |      |      |         |
| 6.3.23 | I/O benchmark (IMB P_Read_Indv)     |      |      |         |
| 6.3.24 | I/O benchmark (IMB P_IRead_Indv)    |      |      |         |
| 6.3.25 | I/O benchmark (IMB P_Read_Shared)   |      |      |         |
| 6.3.26 | I/O benchmark (IMB P_IRead_Shared)  |      |      |         |
| 6.3.27 | I/O benchmark (IMB P_Read_Priv)     |      |      |         |
| 6.3.28 | I/O benchmark (IMB P_IRead_Priv)    |      |      |         |
| 6.3.29 | I/O benchmark (IMB P_Read_Expl)     |      |      |         |
| 6.3.30 | I/O benchmark (IMB P_IRead_Expl)    |      |      |         |
| 6.3.31 | I/O benchmark (IMB C_Read_Indv)     |      |      |         |
| 6.3.32 | I/O benchmark (IMB C_IRead_Indv)    |      |      |         |
| 6.3.33 | I/O benchmark (IMB C_Read_Shared)   |      |      |         |
| 6.3.34 | I/O benchmark (IMB C_IRead_Shared)  |      |      |         |
| 6.3.35 | I/O benchmark (IMB C_Read_Expl)     |      |      |         |
| 6.3.36 | I/O benchmark (IMB C_IRead_Expl)    |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**Results Table 17 - TI MPI - Open MPI**

| Test # | Test Suite                     | Pass | Fail | Comment |
|--------|--------------------------------|------|------|---------|
| 6.3.37 | I/O benchmark (IMB Open_Close) |      |      |         |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

**16.5 OSU MPI TEST RESULTS**

**Results Table 18 - TI MPI - OSU**

| Test # | Test                            | Pass | Fail | Comment |
|--------|---------------------------------|------|------|---------|
| 1      | Test 1: PingPong                |      |      |         |
| 2      | Test 1: PingPing point-to-point |      |      |         |
| 3      | Test 2: PingPong                |      |      |         |
| 4      | Test 2: PingPing                |      |      |         |
| 5      | Test 2: Sendrecv                |      |      |         |
| 6      | Test 2: Exchange                |      |      |         |
| 7      | Test 2: Bcast                   |      |      |         |
| 8      | Test 2: Allgather               |      |      |         |
| 9      | Test 2: Allgatherv              |      |      |         |
| 10     | Test 2: Alltoall                |      |      |         |
| 11     | Test 2: Alltoallv               |      |      |         |
| 12     | Test 2: Reduce                  |      |      |         |
| 13     | Test 2: Reduce_scatter          |      |      |         |
| 14     | Test 2: Allreduce               |      |      |         |
| 15     | Test 2: Barrier                 |      |      |         |

**Results Table 19 Remarks**

|                                                                                                                                             |
|---------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>General Remarks:</b> Comments about the set-up, required updates to the TD, and any other issues that came up during the testing.</p> |
|                                                                                                                                             |
|                                                                                                                                             |
|                                                                                                                                             |
|                                                                                                                                             |
|                                                                                                                                             |
|                                                                                                                                             |

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42