



10GbE MPI BENCHMARKING REPORT

- Chelsio 10GbE RDMA NIC (RNIC)
- Arastra DX7124-S Switch

Executive Summary

IN TESTS conducted at the Chelsio facility, results demonstrate successful interoperability between Chelsio's latest SFP+ 10GbE adapter and Arastra's new 10GbE SFP+ switch. Performance results show that the application can sustain line rate and latency numbers are ~7 μ sec.

Chelsio Communications • www.chelsio.com • sales@chelsio.com • +1-408-962-3600

Test Configuration

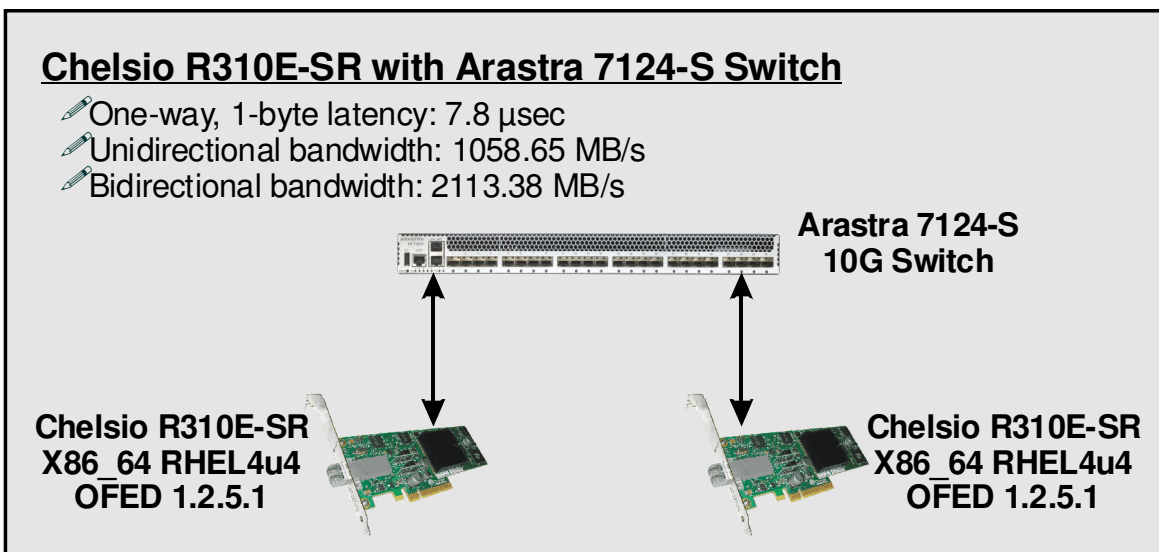
Host Systems

System1:

- CPU: 2 x dual-core Intel Xeon® 5160 @ 3 GHz
- Root Complex: Intel 5000 Chipset
- Memory: 8GB
- OS: Linux 2.6.9-42
- OFED 1.2.5.1
- x86_64

System2:

- CPU: 1 x dual-core Intel Xeon® 3070 @ 2.66 GHz
- Root Complex: Intel 5000 Chipset
- Memory: 8GB
- OS: Linux 2.6.9-42
- OFED 1.2.5.1
- x86_64



Performance Measurement Software

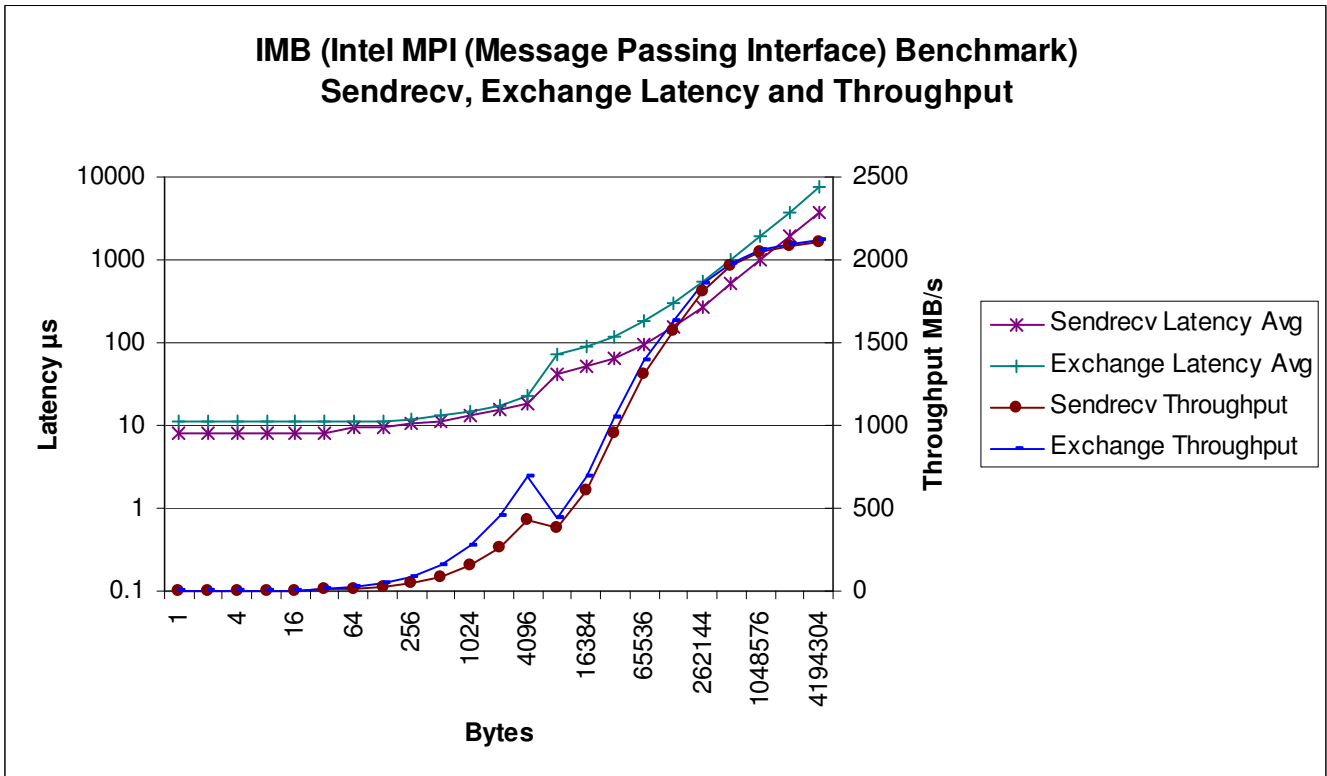
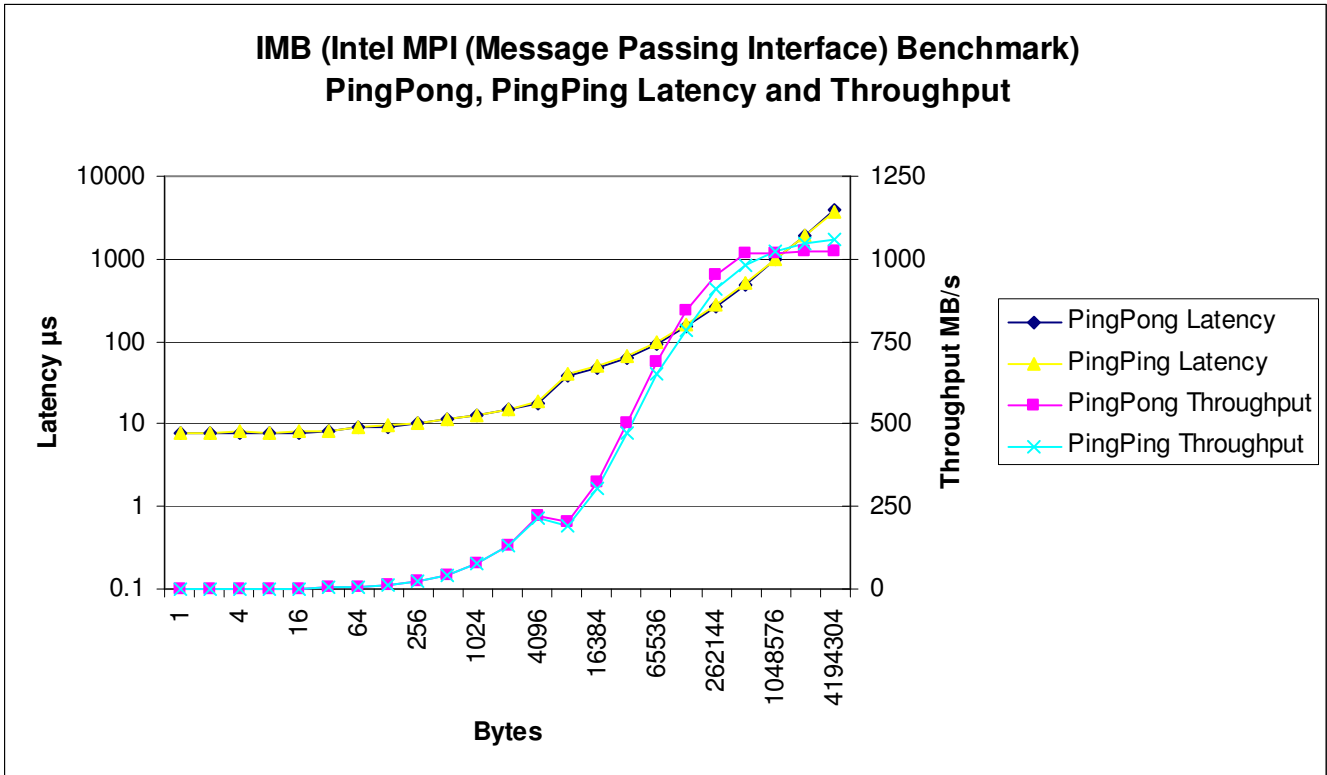
- Test tool: Intel® MPI Benchmark Suite V2.3

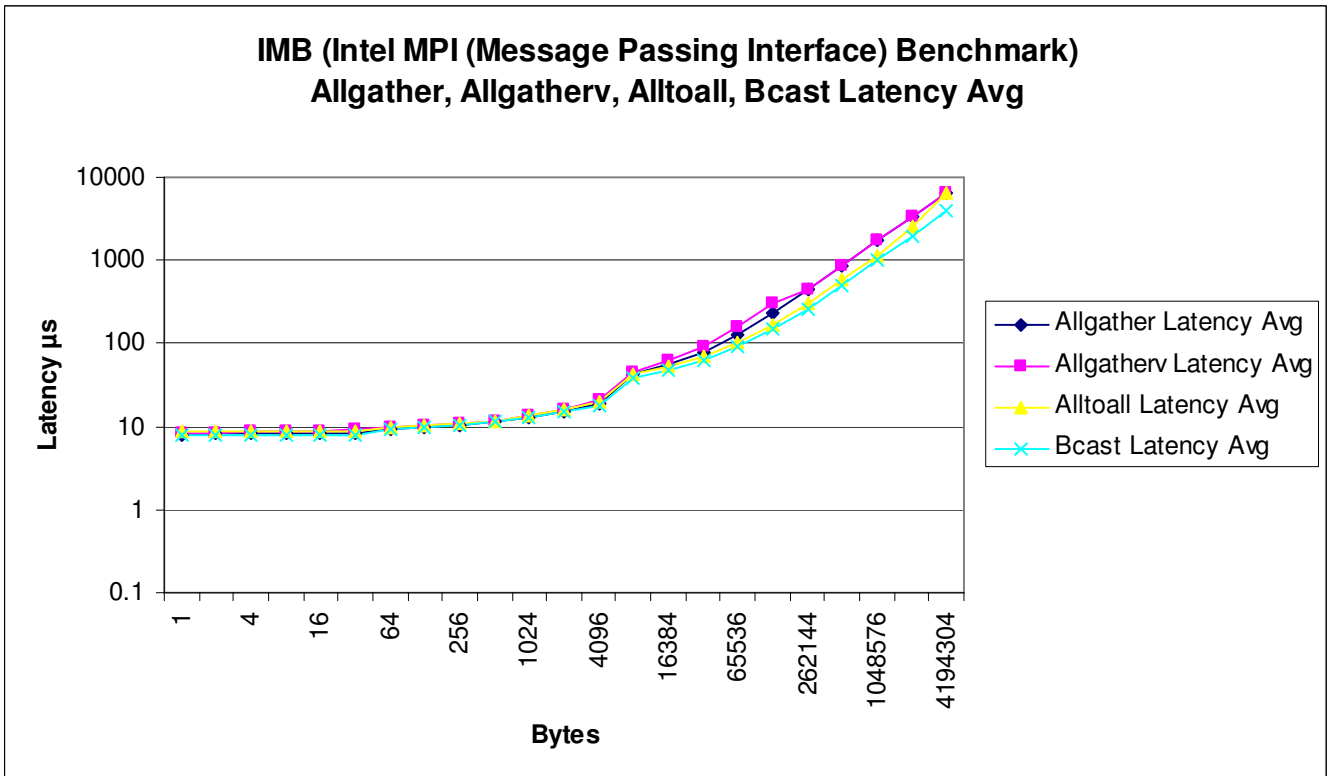
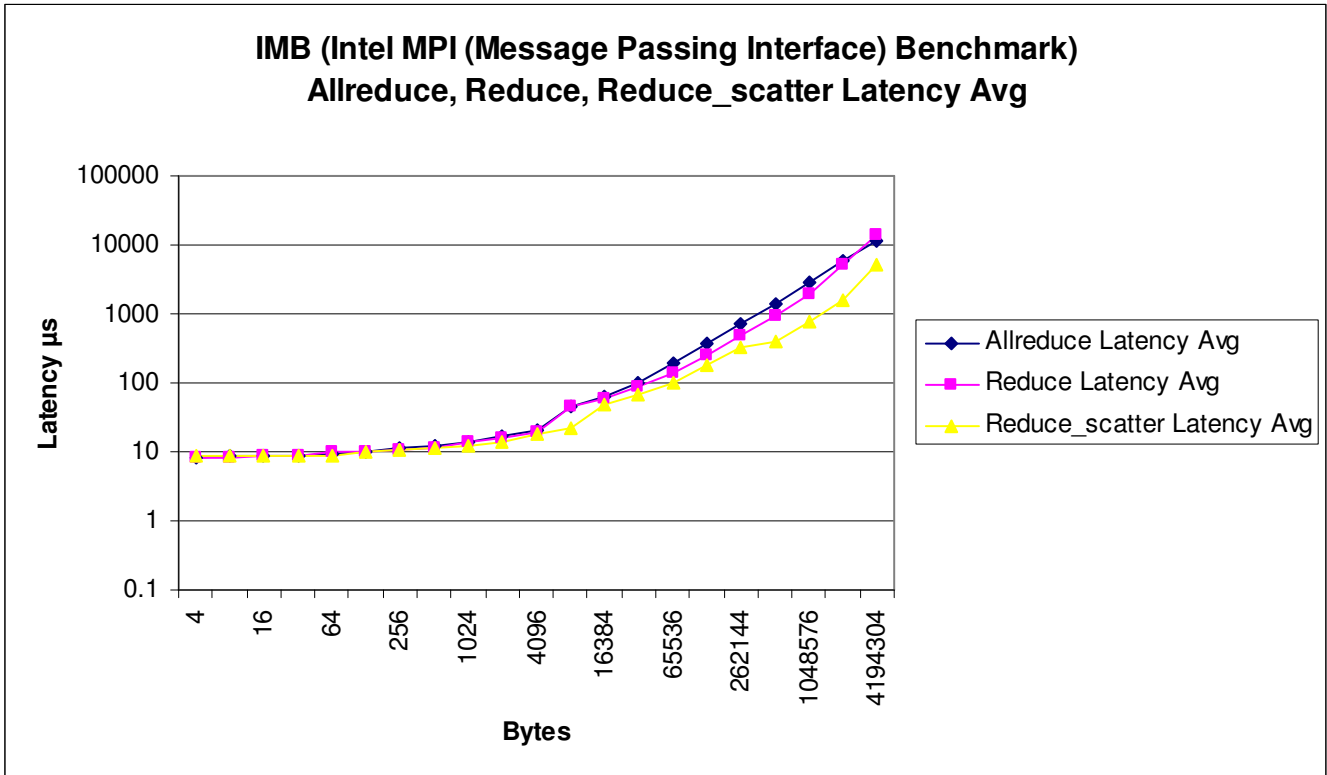
Results

	<i>Latency (μsec)</i>	<i>Throughput (MB/s)</i>
Ping-Pong	7.7	1025.13
Ping-Ping	7.8	1058.65
SendRecv	8.0	2109.43
Exchange	11.1	2113.38
Allreduce	8.4	
Reduce	8.4	
Reduce Scatter	8.7	
Allgather	8.0	
Allgatherv	8.4	
Alltoall	8.5	
Bcast	7.9	

10/18/07

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH CHELSIO PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN CHELSIO'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, CHELSIO ASSUMES NO LIABILITY WHATSOEVER, AND CHELSIO DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF CHELSIO PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT. CHELSIO PRODUCTS ARE NOT INTENDED FOR USE IN MEDICAL, LIFE SAVING, OR LIFE SUSTAINING APPLICATIONS. CHELSIO MAY MAKE CHANGES TO SPECIFICATIONS AND PRODUCT DESCRIPTIONS AT ANY TIME, WITHOUT NOTICE.





Test Notes

- **PingPing:** Measures the startup (latency) and throughput of a single message sent between two processes, with each node issuing the commands concurrently.
- **Sendrecv:** Creates a chain of sending and receiving between upstream and downstream neighbors.
- **Exchange:** Creates a chain of sending and receiving between upstream and downstream neighbors, with each node issuing the send/receive commands concurrently.
- **Allreduce:** Measures the time to perform the allreduce collective operator where each process in a group owns a vector whose elements are combined with the corresponding elements owned by other processes, producing a new vector of the same length.
- **Reduce:** Measures the time to perform the allreduce collective operator where each process in a group owns a vector whose elements are combined with the corresponding elements owned by other processes, producing a new vector of the same length, changing the root of the process cyclically.
- **Reduce_scatter:** Measures the time to perform the allreduce collective operator, split between all processes, where each process in a group owns a vector whose elements are combined with the corresponding elements owned by other processes, producing a new vector of the same length.
- **Allgather:** Measures the time that it sends X bytes and gathers X*processes bytes.
- **Allgatherv:** Measures the time that it sends X bytes and gathers X*(#processes) bytes, and sees if MPI produces more overhead due to the more complicated situation.
- **Alltoall:** Measures the time that it takes for every process to input X*(#processes) bytes (X for each process) and receives X*(#processes) bytes (X for each process).
- **Bcast:** Measures the time that it takes a root process to broadcast X bytes to all while the root of the operation is changed cyclically.