

Infiniband in ALICE HLT

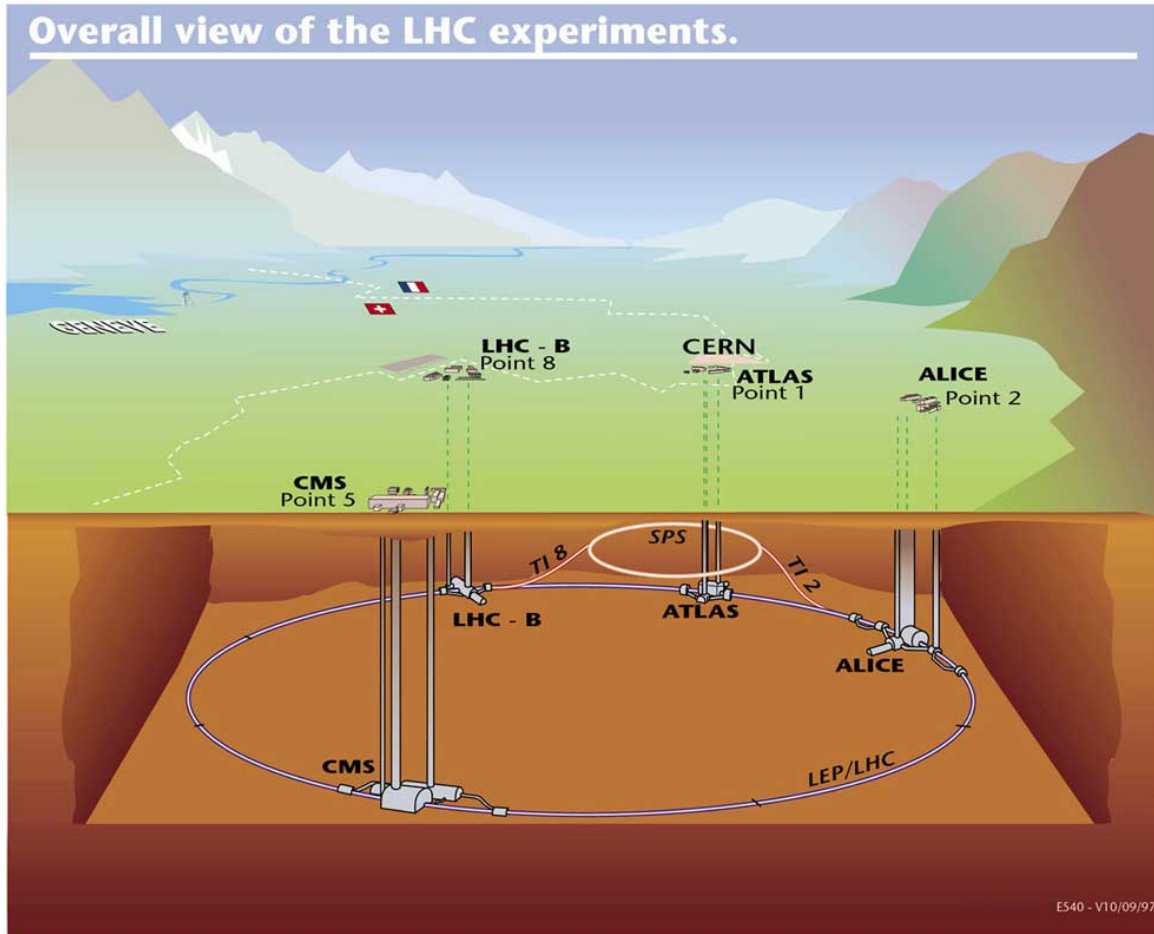
Timo Breitner for the ALICE HLT project

Frankfurt Institute of Advanced Studies

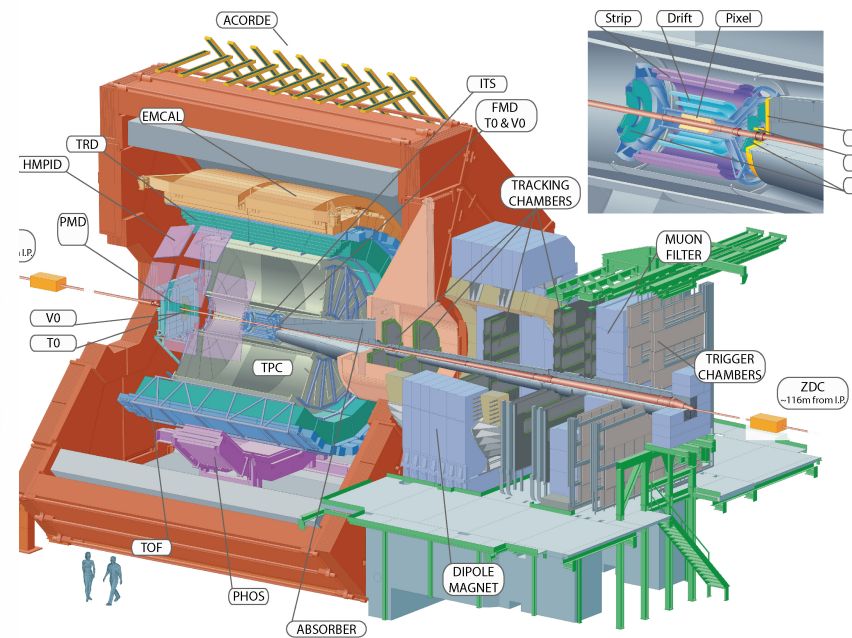
University of Frankfurt, Computer Science Department

tbreitne@cern.ch

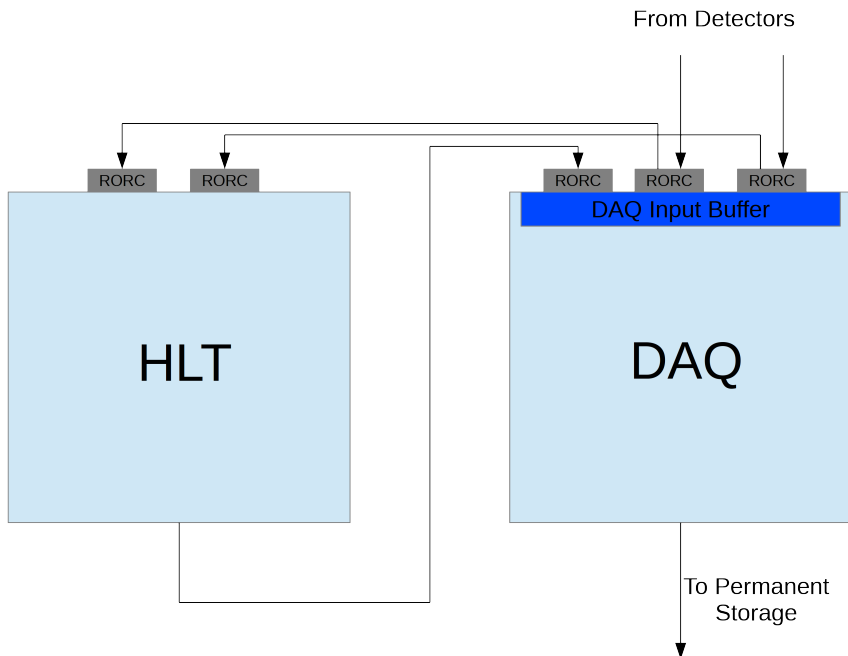
- Infiniband@ALICE HLT installed summer 2010, in production since late 2010
- HLT software stack based on IP
 - Only IPoIB, (currently) no RDMA
- Topics of this talk
 - HLT layout, bandwidth requirements
 - Motivation for using IB
 - Performance numbers
 - Experience



LHC experiment focused on heavy ion physics



Avg. event size of central PbPb collision: $\sim 70\text{MB}$
 TPC largest contributor ($\sim 85\%$)



Purpose: partial/full data reconstruction

Trigger/filter

Data compression

Online QA

Input: Raw detector data, via optical fibers (DDL), direct copy from DAQ

Output: Trigger decision, compressed data forwarded to DAQ via optical fibers

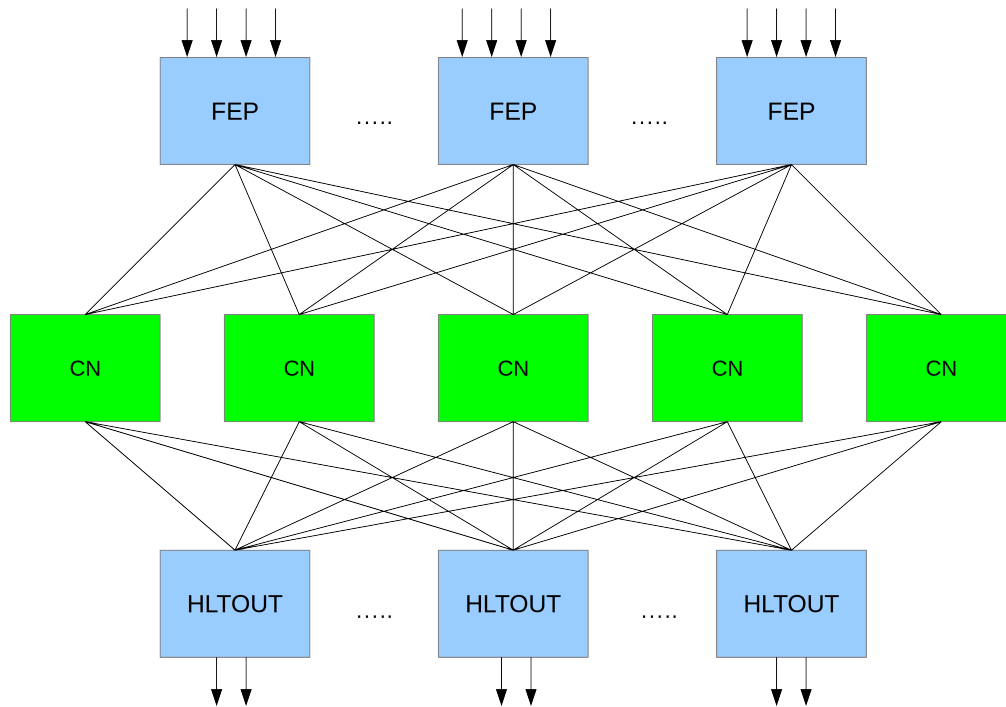
Current std. Operation mode: TPC data reduction

- Local Pattern Recognition (cluster finding)

- Huffman encoding

- Compression factor 4-5, more possible

Input layer: 54 x Front End Processors
 2 x AMD Opteron 2378 (Shanghai, 4 cores)
 4 DDL@160MB/s

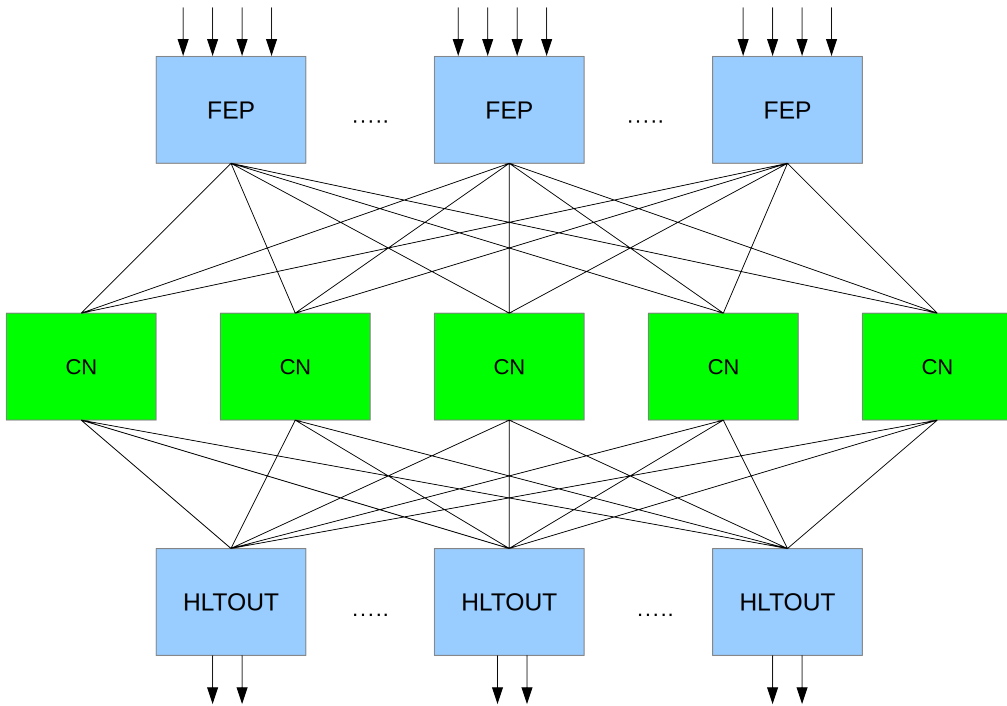


Processing layer: 45 x Compute nodes
 2 x AMD Opteron 6172 (MagnyCours, 12 cores)
 Nvidia GPU GTX480/580
 Track reconstruction, vertexing, data compression

Output layer: 14 x HLTOut nodes
 2 x AMD Opteron 2378 (Shanghai, 4 cores)
 2 DDL@160MB/s

Total nodes in system:
 100 FEP nodes
 80 CN nodes

- HLT data transport framework based on SHM (intra-node) and TCP/IP (inter-node)
- Initial deployment: 1Gbit Ethernet
 - Enough for first years (pp)
 - Still in use for system management
- Towards PbPb runs 2010:
 - Discussion on suitable network infrastructure
 - Particle multiplicity (event size) not clear
 - Event rate not clear
- Worst case scenario assumed!



FEP in: $4 \text{ DDL} * 160 \text{ MB/s} \rightarrow 640 \text{ MB/s}$
 FEP out: 640 MB/s

Network: $54 * 640 \text{ MB/s} \rightarrow 34.56 \text{ GB/s}$

CN in: $34.56 \text{ GB/s} / 45 \text{ nodes} \rightarrow 768 \text{ MB/s}$
 CN out: $768 \text{ MB/s} / 4 \rightarrow 192 \text{ MB/s}$

Network: $45 * 192 \text{ MB/s} \rightarrow 8.64 \text{ GB/s}$

HLTOUT in: $8.64 \text{ GB/s} / 14 \rightarrow 617 \text{ MB/s}$

Per-node Bandwidth: up to 800 MB/s

Network bandwidth: up to 45 GB/s

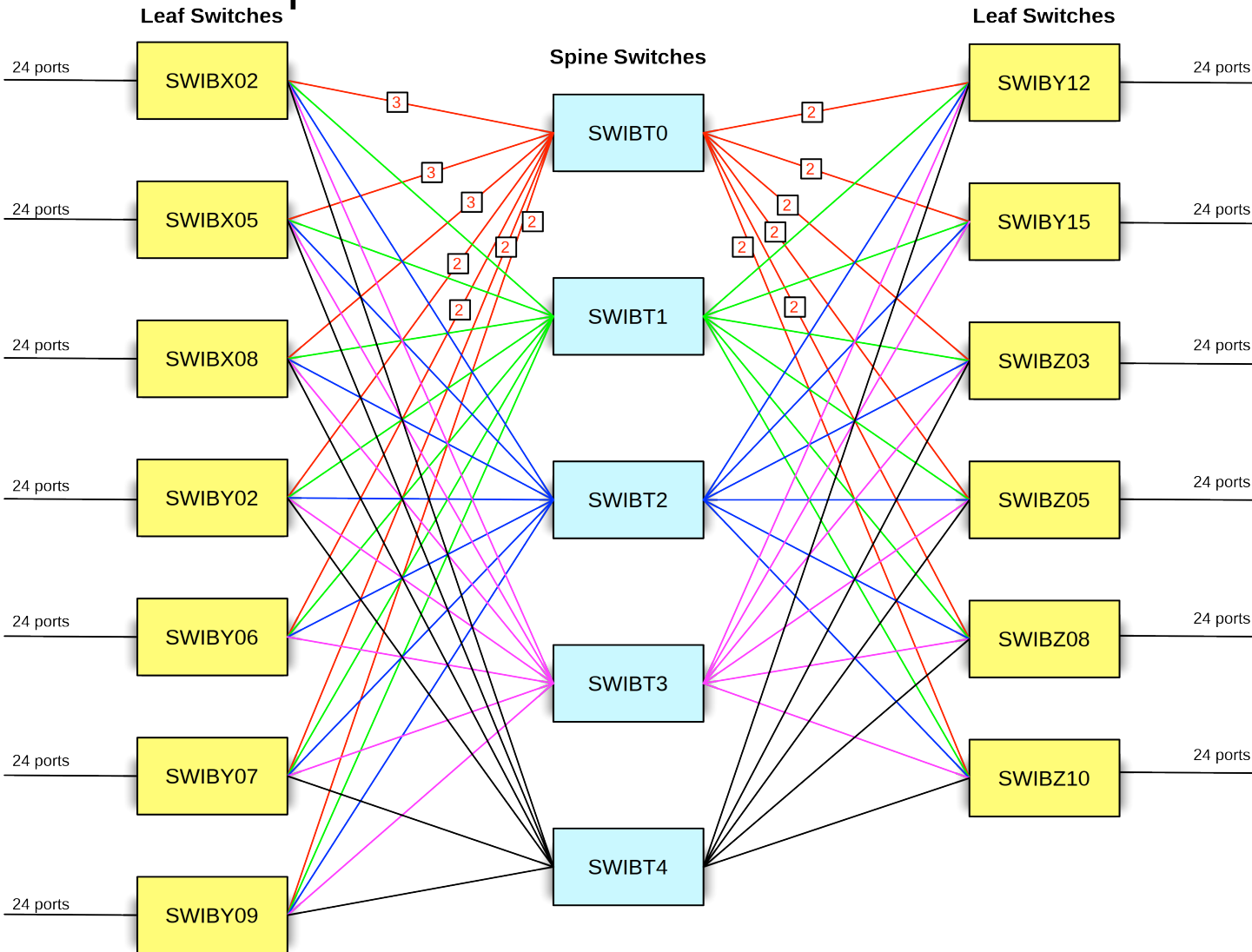
10 Gbit Ethernet OR Infiniband QDR

Note: latency is not an issue

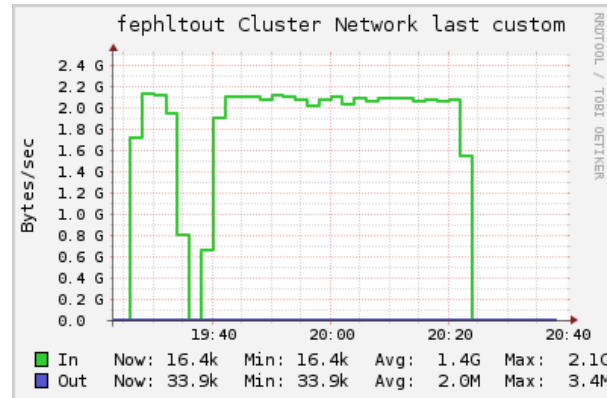
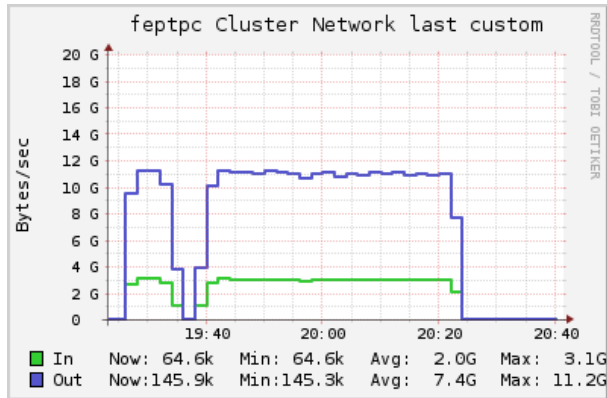
Why Infiniband

- Had already some experience
 - Small test setup: backbone net, proof-of-concept
- Price/performance
 - Nominally 4 x faster than 10GBit ethernet (only with “native” IB mode possible)
 - IPoIB: ~1.5 x faster
 - Cheaper (back then, probably still is)
 - Israeli program to fund acquisition of Israeli Hi-tech products by CERN

5 x 36 port spine switches
 13 x 36 port leaf switches: 24 node connections, 12 uplinks



2/3 uplinks per leaf switch
 Mix of FEP and CN nodes per leaf switch



Typical data rates:

- Input from TPC: 9 GByte/s

- Output to DAQ: 2.2 GByte/s

Maximum data rates:

- Input from TPC: 14.9 GByte/s

- Output to DAQ: 3.5 Gbyte/s

Worst case scenario assumed twice the bandwidth

- 3 types of nodes in use @HLT:
 - Shanghai: AMD Opteron 2378, 2 x 4 cores
 - MagnyCours: AMD Opteron 6172, 2 x 12 cores
 - Nehalem: Intel Xeon E5520, 2 x 4 cores + HT

- All nodes equipped with Mellanox ConnectX PCIe 2.0 QDR (Onboard or add-on card)

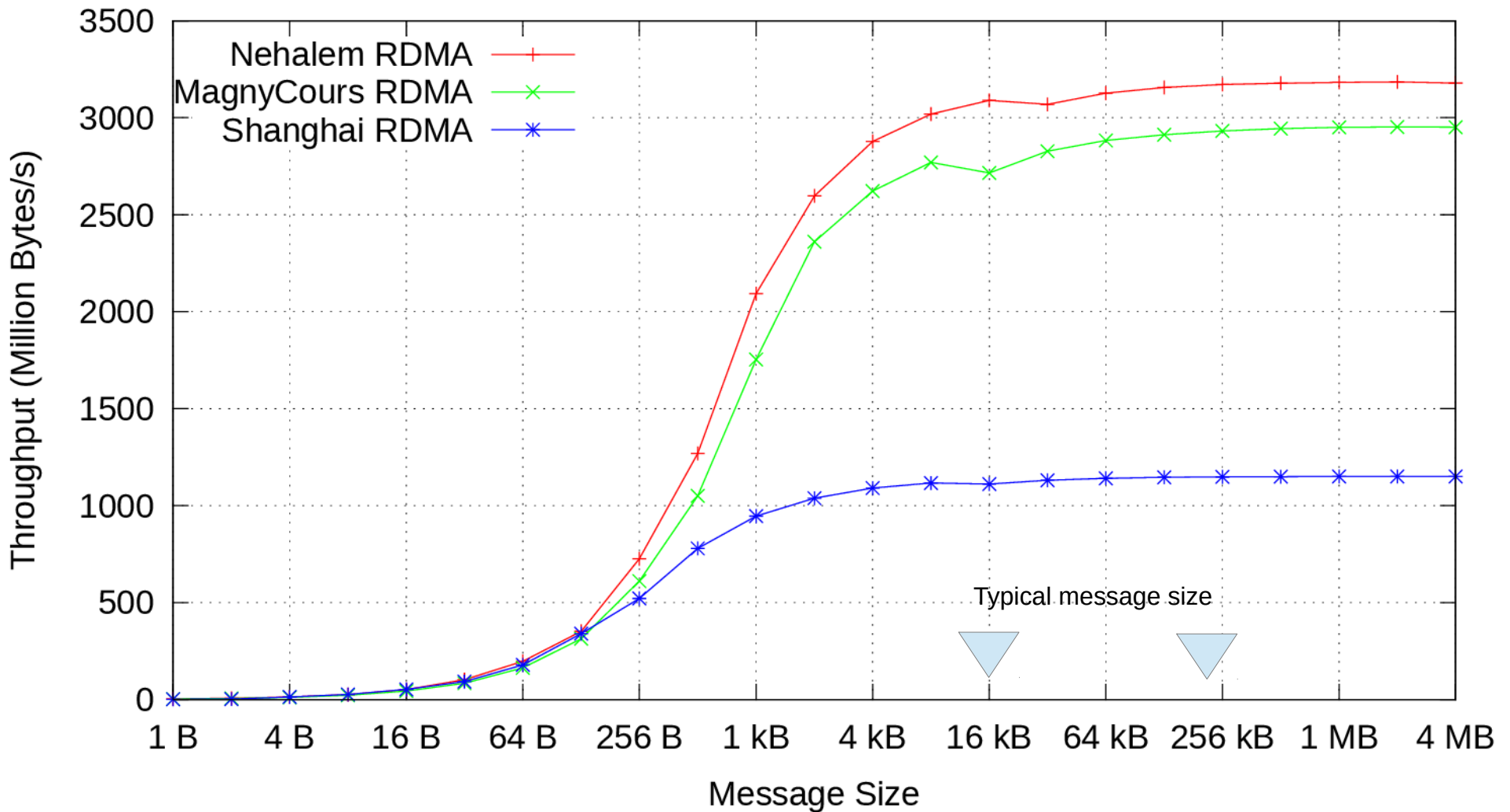
- osu micro-benchmark suite, MVAPICH2

(osu_bw, osu_mbw_mr)

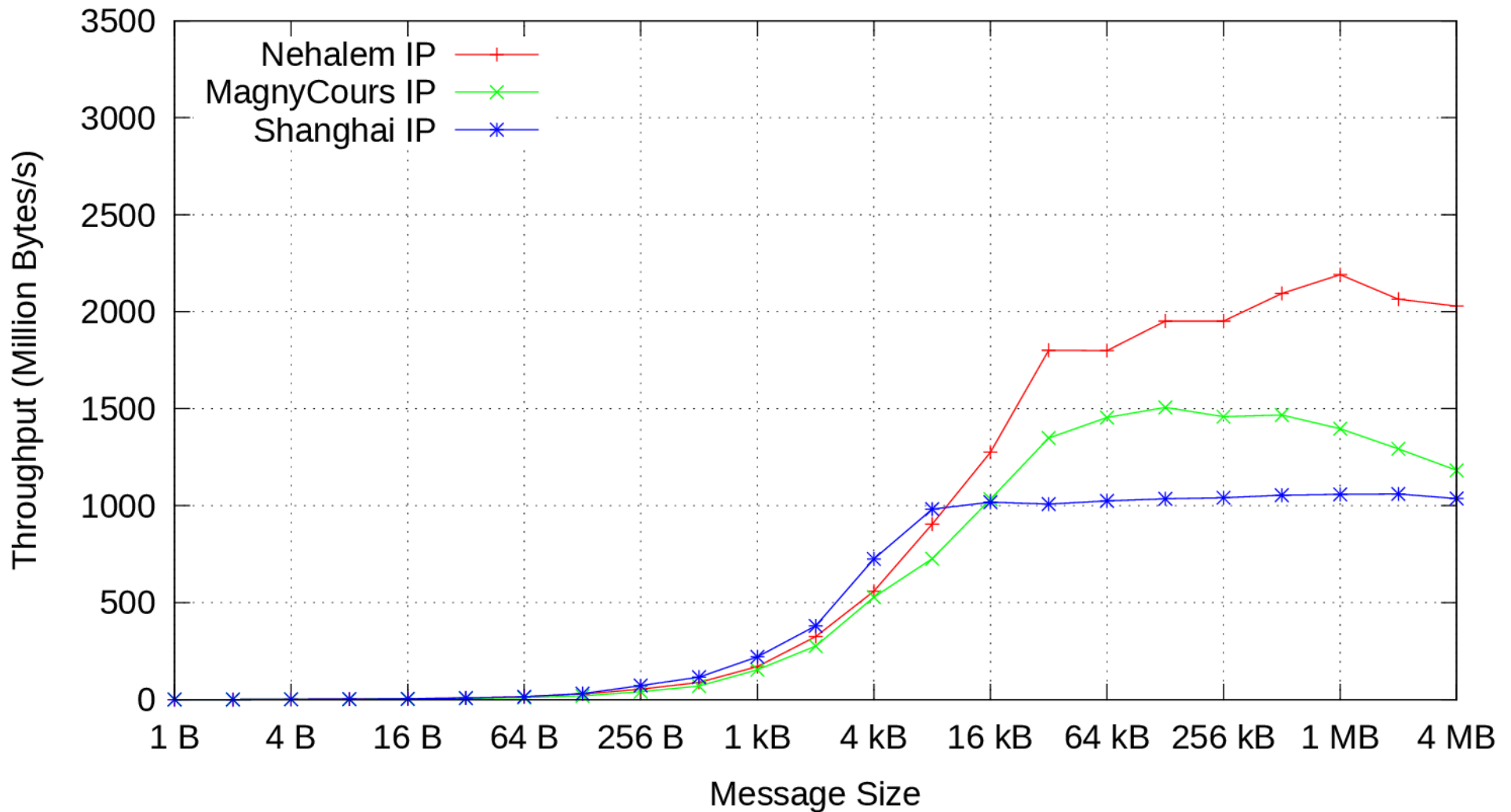
<http://mvapich.cse.ohio-state.edu/benchmarks/>

Point-to-point performance

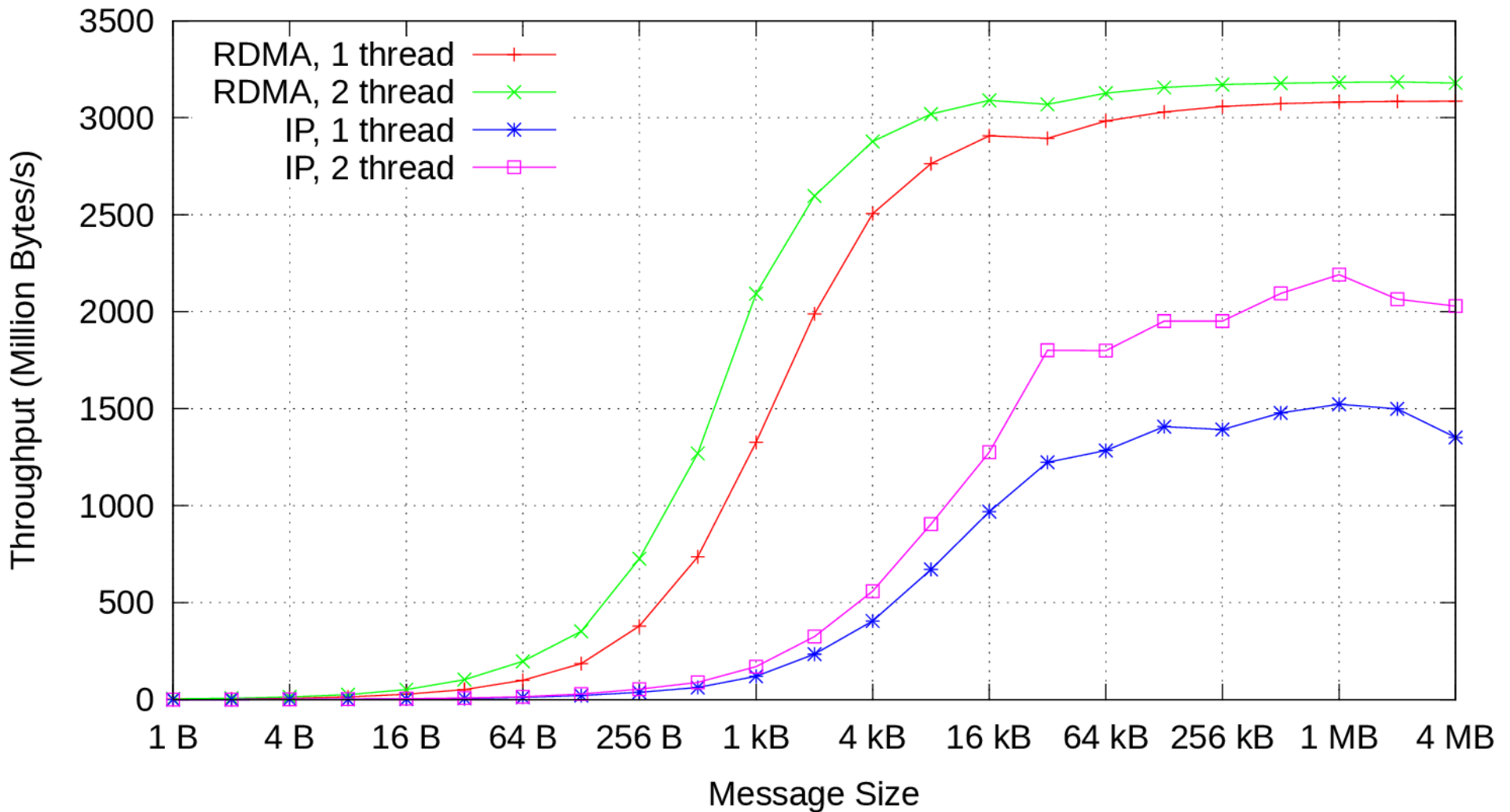
Network Throughput vs. Message Size



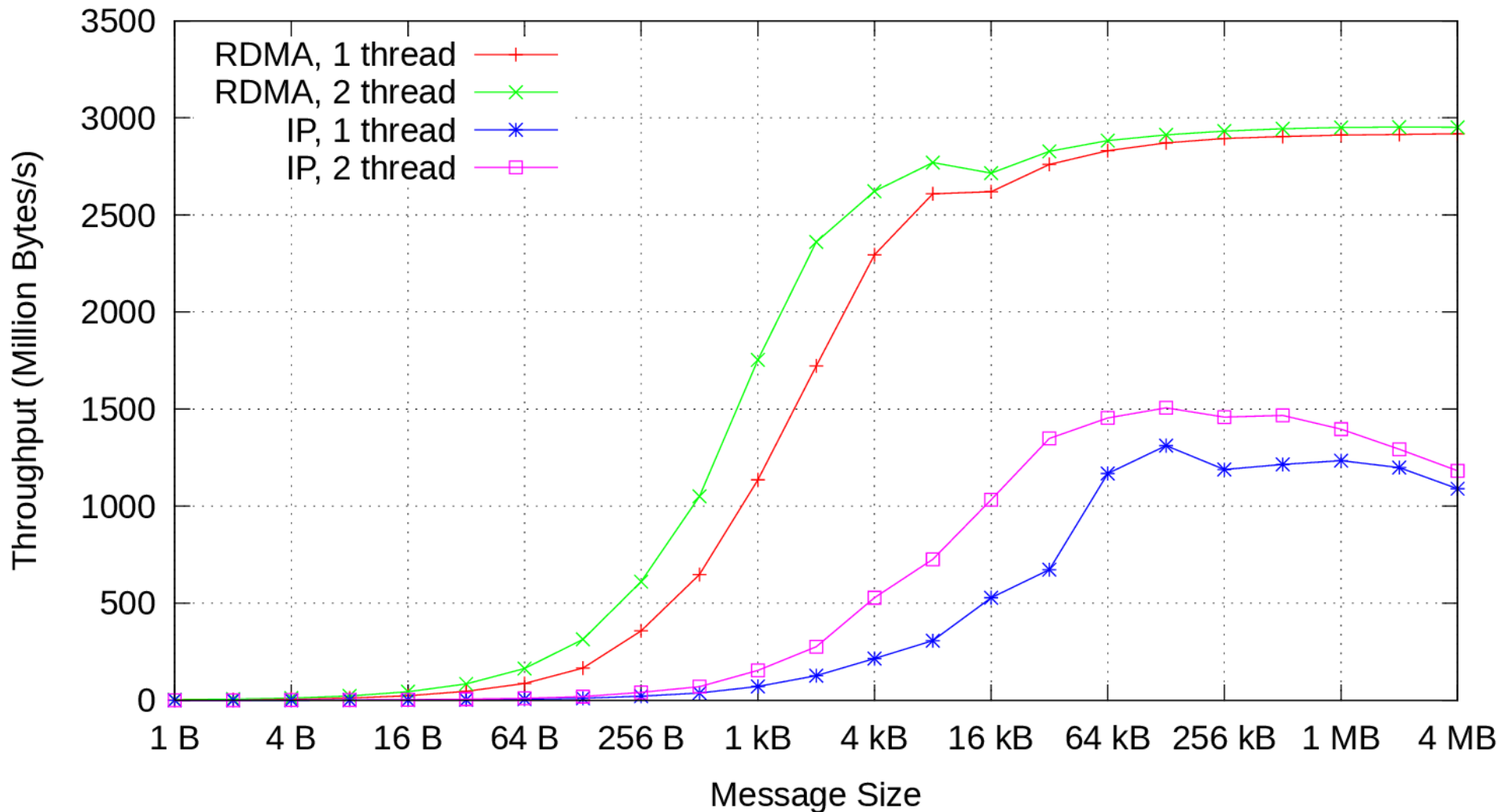
Network Throughput vs. Message Size



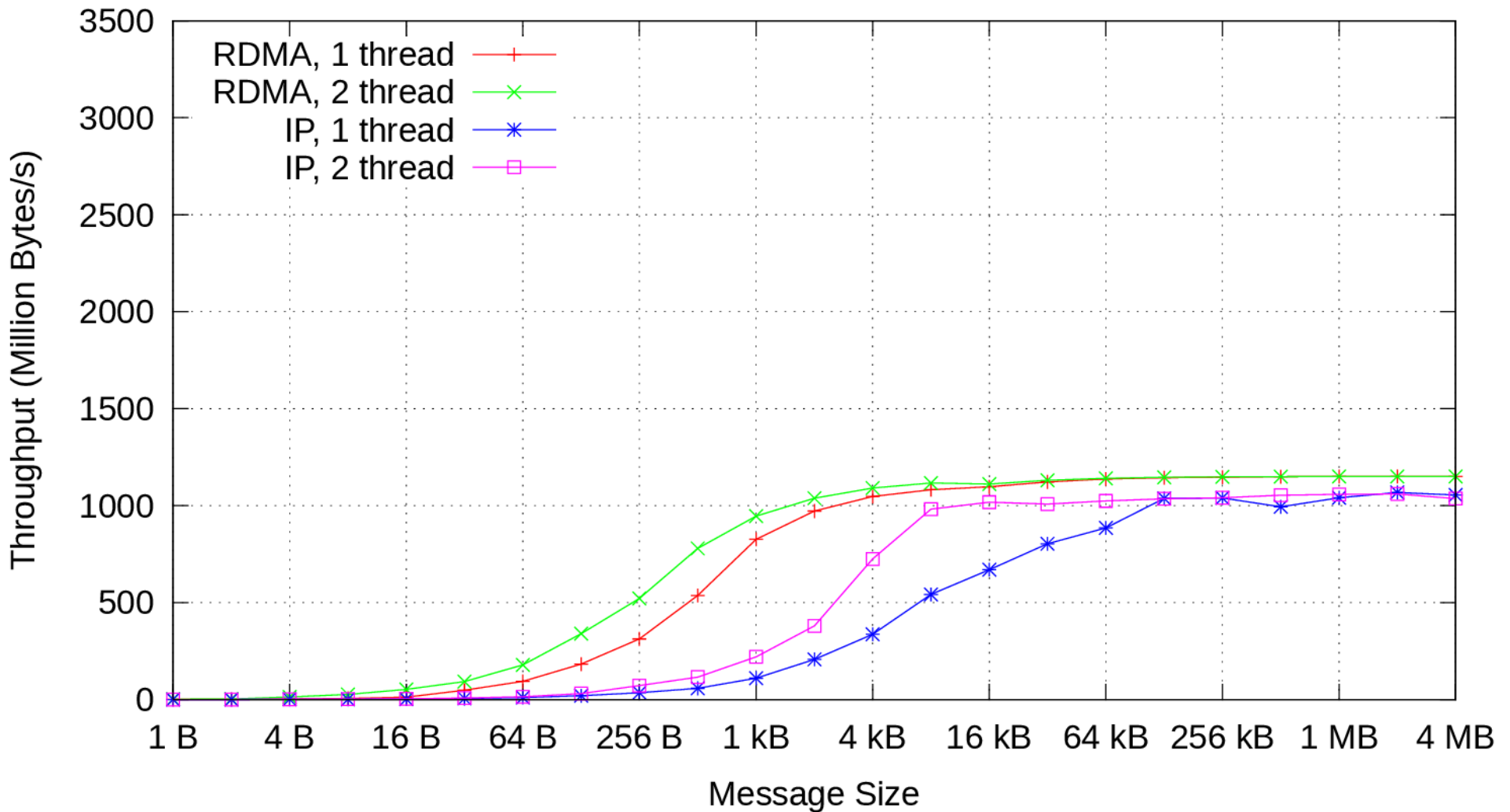
Network Throughput vs. Message Size (Nehalem-2-Nehalem)



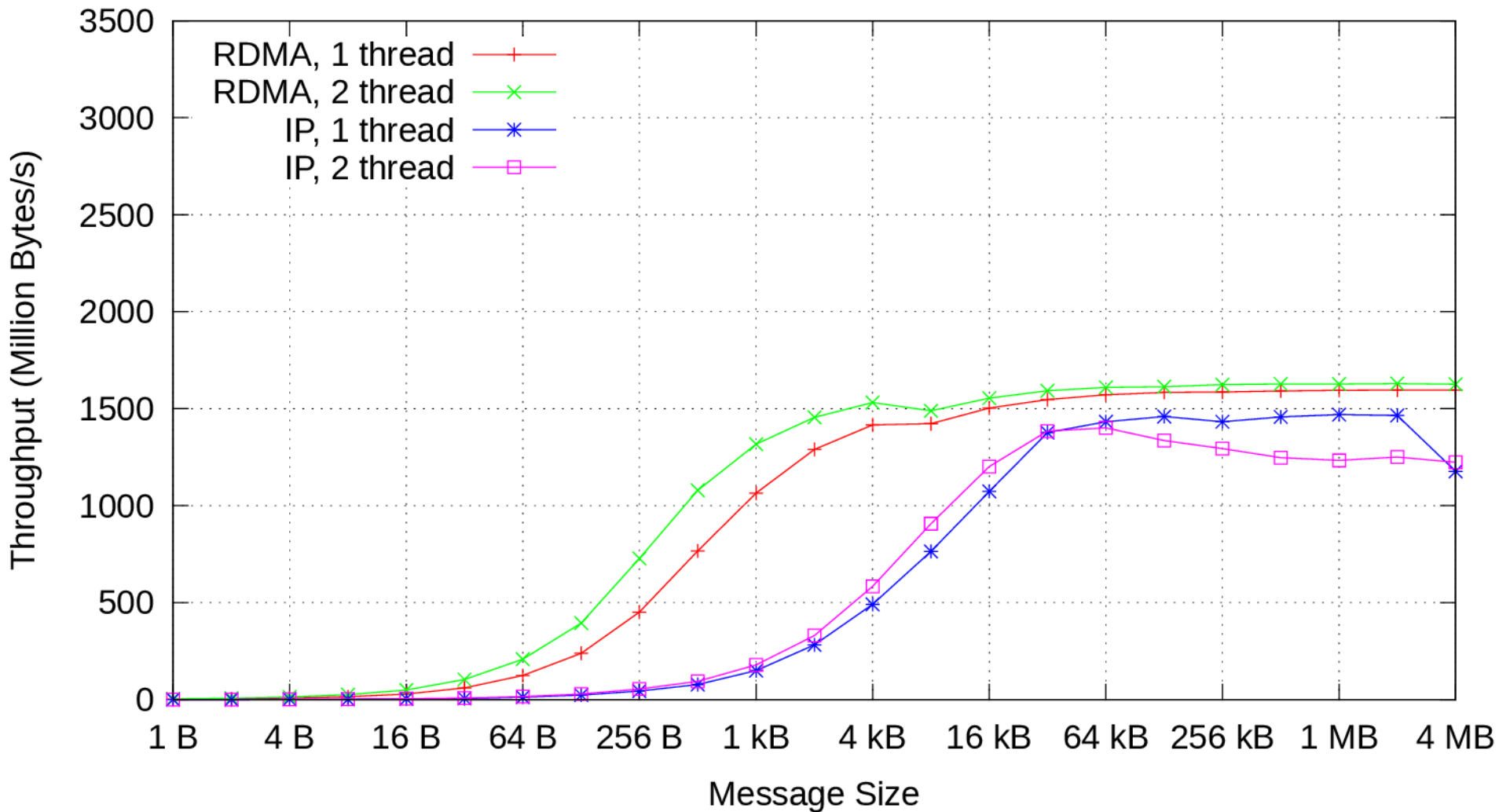
Network Throughput vs. Message Size (MagnyCours-2-MagnyCours)



Network Throughput vs. Message Size (Shanghai-2-Shanghai)

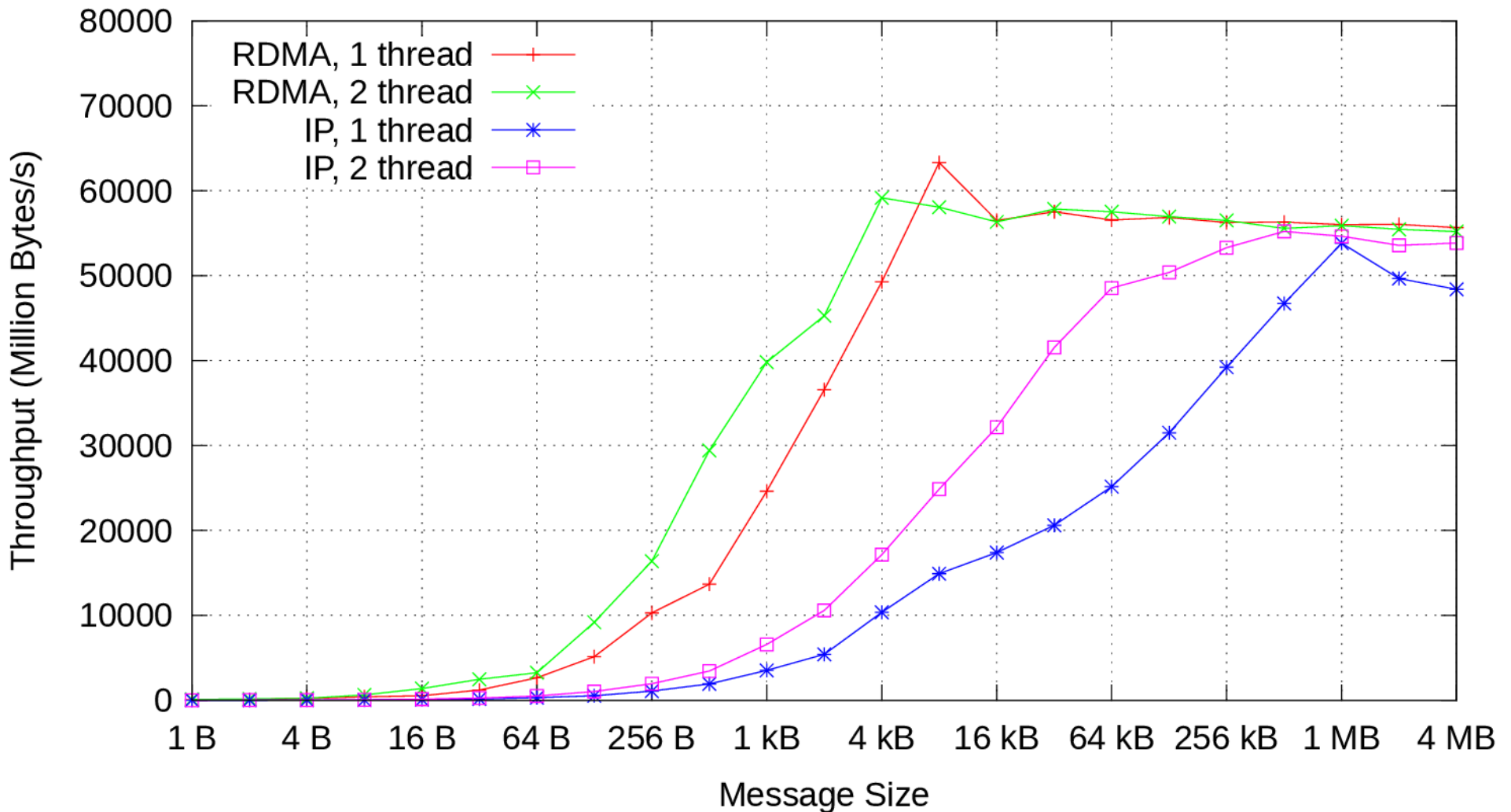


Network Throughput vs. Message Size (Shanghai-2-Nehalem)

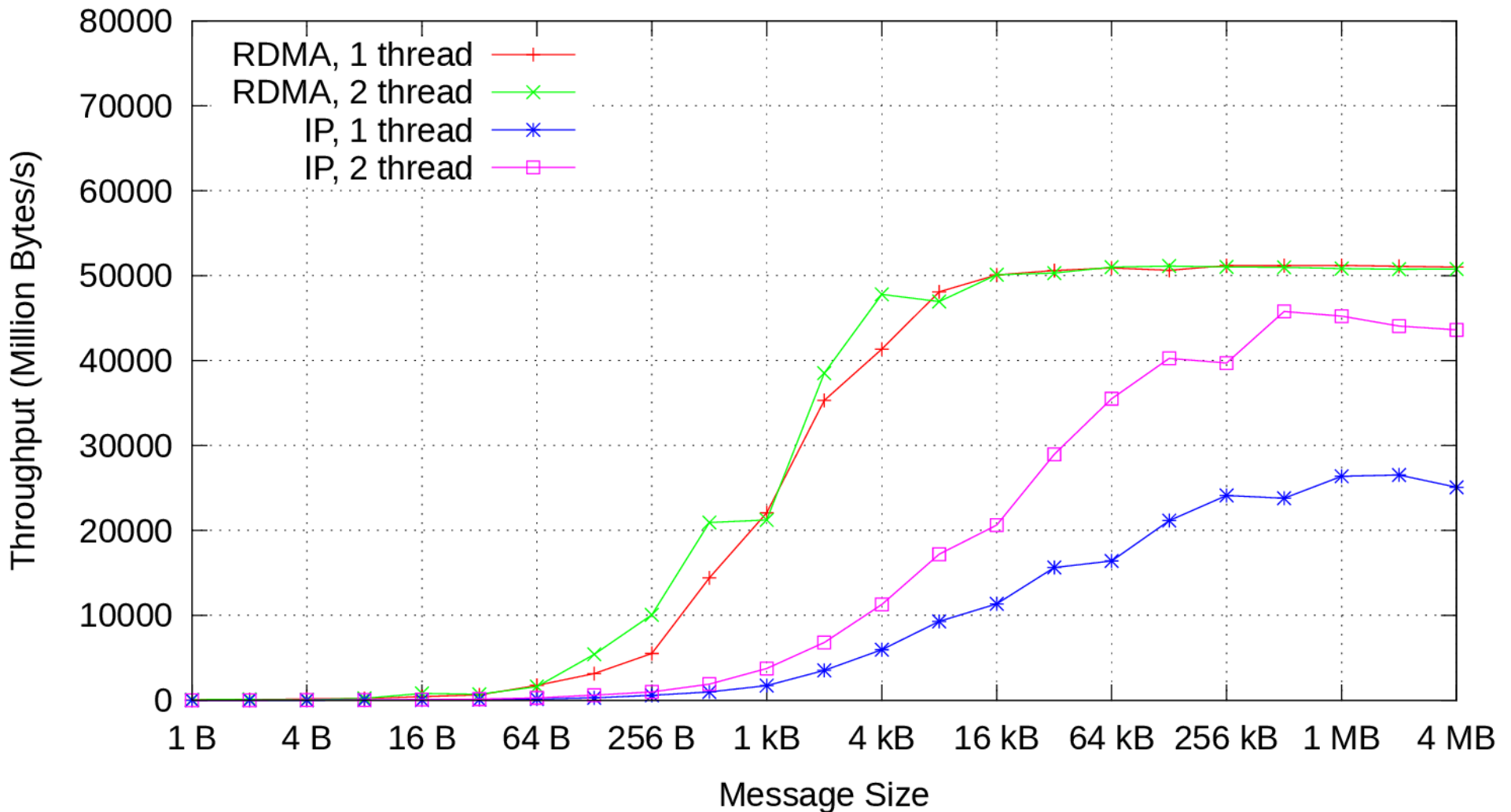


Aggregate bandwidth performance

Network Throughput vs. Message Size (ALL-90pairs)



Network Throughput vs. Message Size (CN-40pairs)



- General experience very good
 - Reliable hardware, no failures/losses
 - Massive network bandwidth, nowhere close to saturation in current situation
- Several issues to be resolved during commissioning:
 - Hardware installation (cables)
 - Software/OS setup
 - Firmware

- Detailed plan of counting room layout and switch location required
 - Cables are rather expensive (esp. longer ones)
 - Copper cables very inflexible (esp. longer ones)
 - Copper cables are quite heavy and can break HCA
- With FDR not so much of a problem
 - Copper cables only up to 3m
 - Fibers for anything beyond (easier to handle, more expensive)

- OFED stack easy to install
- Lots of tunables (OS settings), some are crucial
- OFED packages provide install scripts setting up most of the important things
 - Only RHEL, SLES officially supported
 - Track config files with Chef/Puppet
- Network setup issue during boot:
 - Takes rather long for link state to become “active”
 - Necessary to delay any network access

- Issue experienced during HPC cluster setup at GSI
- Lots of issues with FDR network, large error counts in switches
- Firmware problem: link training was not robust enough/failed
- Had to manually set link speed to QDR
Issue now solved by firmware fix
- Firmware update (switch and HCA) first thing to do
- Mellanox firmware does not fit 3rd party vendor devices

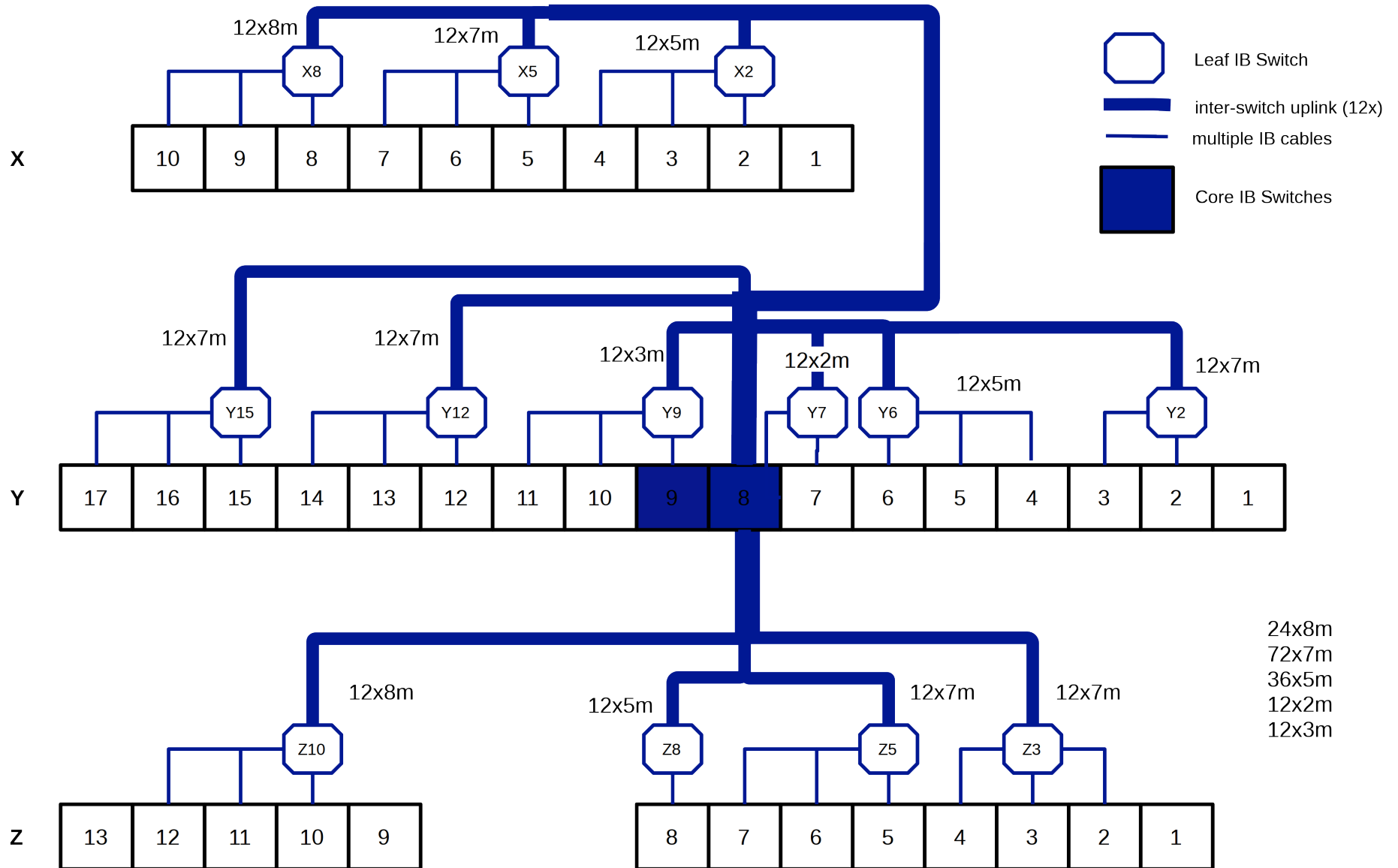
- Reliable and more than enough bandwidth for our current needs
- Some difficulties during setup
- Even more bandwidth possible by using more efficient protocol (RDMA)
- Future ALICE upgrade with combined DAQ/HLT system has much higher requirements (up to 1 TB/s)
- IB also an option there, but has to compete with other technologies in both price and performance



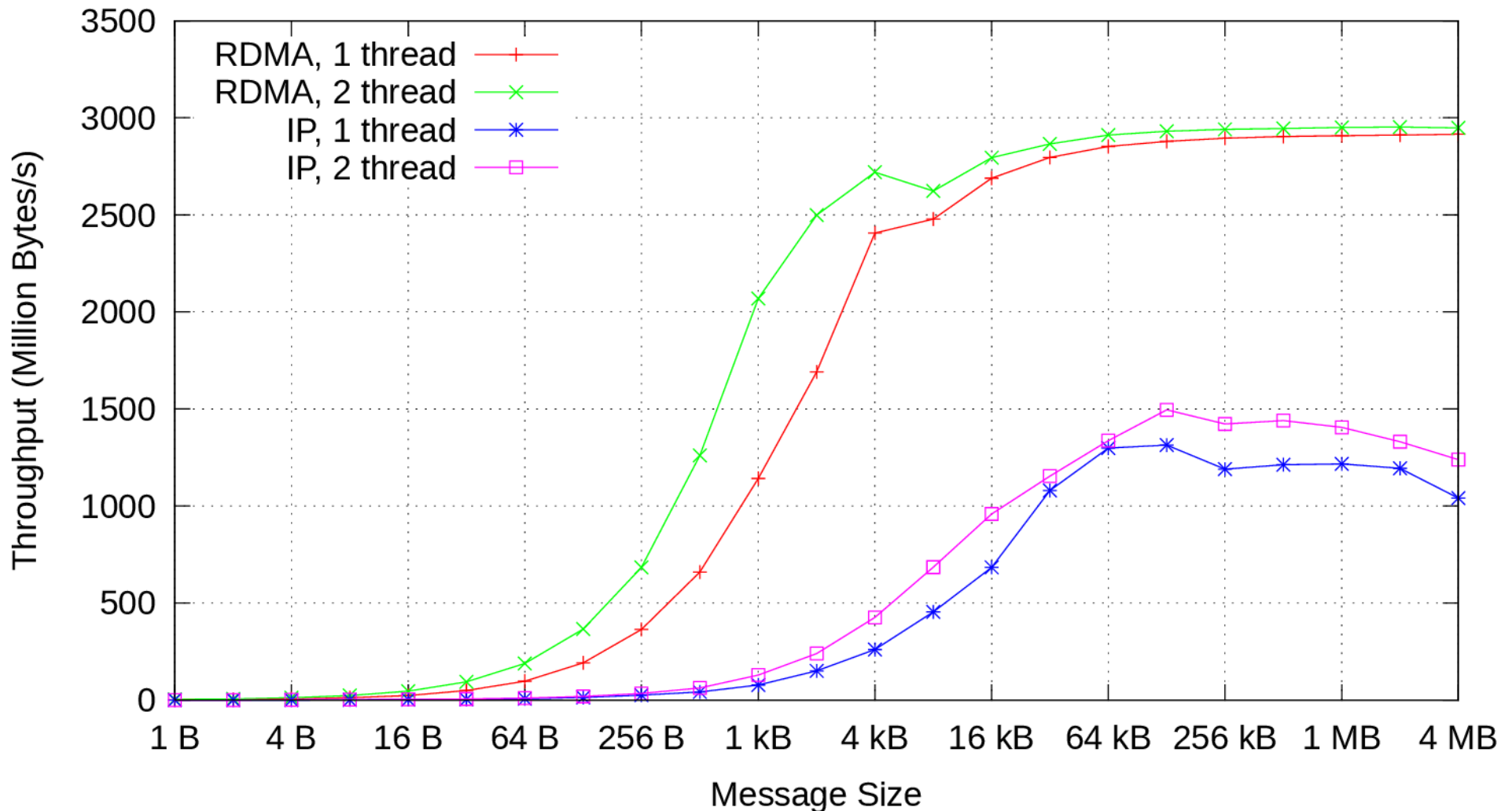
Backup



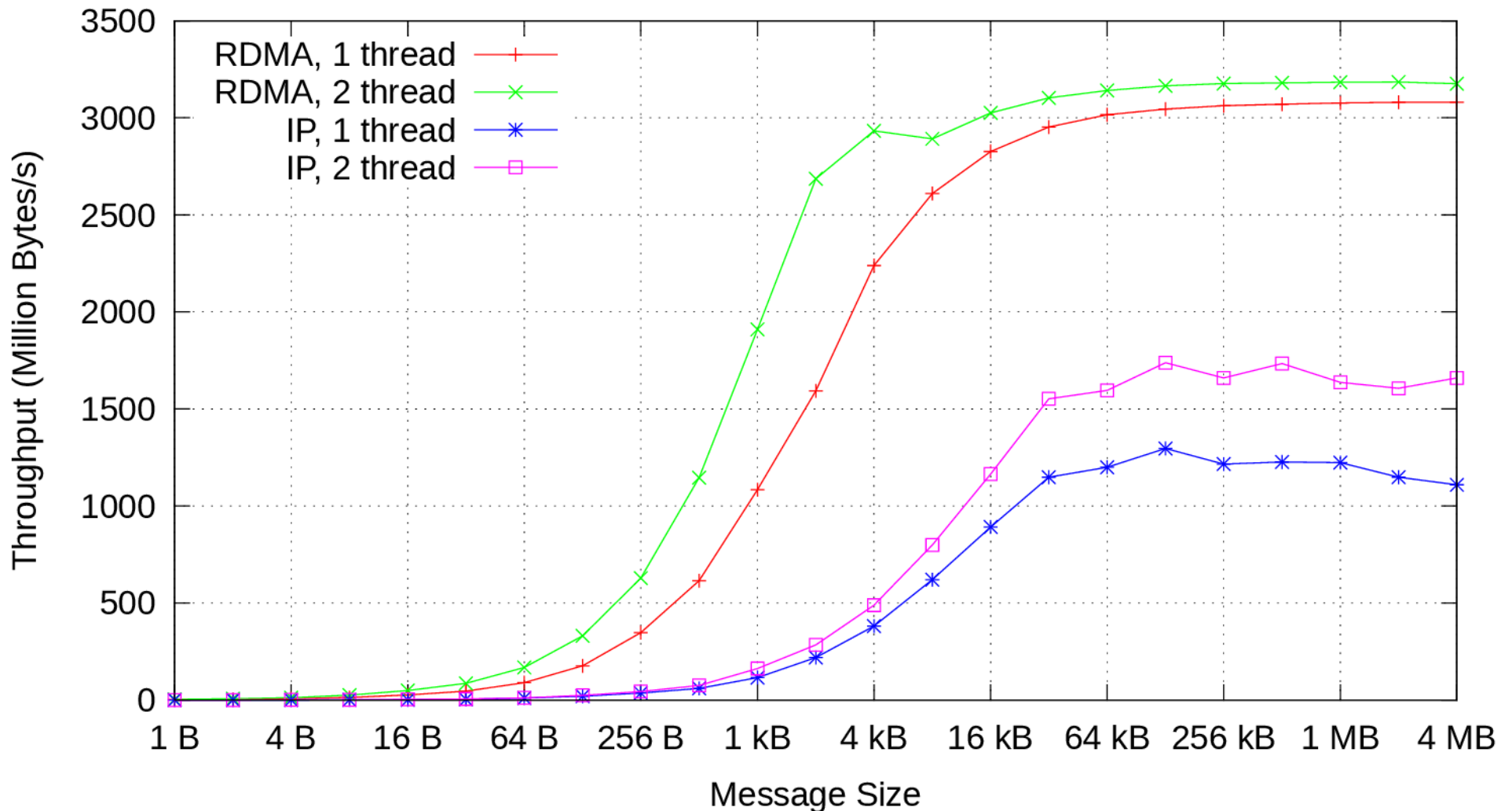
Counting Room Layout



Network Throughput vs. Message Size (MagnyCours-2-Nehalem)



Network Throughput vs. Message Size (Nehalem-2-MagnyCours)



Network Throughput vs. Message Size (Nehalem-2-Shanghai)

