

Scalable HPC Communication Capabilities

Rich Graham HPC 2014, Cetraro, Italy



Mellanox Connect. Accelerate. Outperform."

HPC The Challenges





Exascale-Class Computer Platforms – Communication Challenges

Challenge	Solution focus
Very large functional unit count ~10M	Scalable communication capa point & collectives
Large on-"node" functional unit count ~500	Scalable HCA architecture
Deeper memory hierarchies	Cache aware network access
Smaller amounts of memory per functional unit	Low latency, high b/w capabili
May have functional unit heterogeneity	Support for data heterogeneity
Component failures part of "normal" operation	Resilient and redundant stack
Data movement is expensive	Optimize data movement
Independent remote progress	Independent hardware progre
Power costs	Power aware hardware

© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014





abilities: point-to-





Challenges of Scalable HPC Communication Libraries

Truly asynchronous communication capabilities

- Two-sided point-to-point
- One-sided point-to-point
- Collective Communication
- Truly one-sided remote progress
- System noise mitigation
- Topology awareness



4



Standardize your Interconnect





Standard !

Standardized wire protocol

- Mix and match of vertical and horizontal components
- Ecosystem build together
- Open-source, extensible interfaces
 - Extend and optimize per applications needs









INFINIBAND" TRADE ASSOCIATION







MPI / SHMEM / UPC

Vendor Extensions

HW D

Speed, Speed, Speed





Technology Roadmap – Always One-Generation Ahead



© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014



Mellanox Connect-IB The World's Fastest Adapter

- The 7th generation of Mellanox interconnect adapters
- World's first 100Gb/s interconnect adapter (dual-port FDR 56Gb/s InfiniBand)
- Delivers 137 million messages per second 4X higher than competition
- Support the new innovative InfiniBand scalable transport Dynamically Connected





nfiniBand) tition cally Connected

Connect-IB Provides Highest Interconnect Throughput



Gain Your Performance Leadership With Connect-IB Adapters

© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014











Scalability

Collective Operations

Topology Aware

Hardware Multicast

Separate Virtual Fabric

Offload

Scalable algorithms









The DC Model

- Dynamic Connectivity
- Each DC Initiator can be used to reach any remote DC Target
- No resources' sharing between processes
 - process controls how many (and can adapt to load)
 - process controls usage model (e.g. SQ allocation policy)
 - no inter-process dependencies

Resource footprint

- Function of node and HCA capability
- Independent of system size

Fast Communication Setup Time





cs – concurrency of the sender cr=concurrency of the responder

Dynamically Connected Transport

Key objects

- DC Initiator: Initiates data transfer
- DC Target: Handles incoming data





Targets (destinations)

Dynamically Connected Transport – Exascale Scalability





- 10K nodes
- 100K nodes

Reliable Connection Transport Mode





Dynamically Connected Transport Mode





More...



© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014





Network Offload

6 5 h b 8 5 9 0 Ø



41	U U	Ų	4	7	1	Ţ
66	2 1	1	8	2	3	4
981	06	9	9	9	5	
904	4 5	4		3	4	8
84	1 1	9	9	1	0	0
89	9 0	8			6	1
231	0 0	9	1	8	6	3
421	06	7	2	9	5	6
510	65	5	2	7	4	5
16	1 1	9	7	9	6	4
45	4 6	8	8	5	4	9
66	6 5	5	5	9	6	6
55	5 1	0	9	6	5	8
11	1 3	1	7	2	4	4
33.	4 2	6	8		6	6
5 2	27	5	2	6		
45	5 5	0	1	4	4	4
2 1	1 9	5	6	2	6	9
States of the local division of the						

Scalability of Collective Operations



© 2014 Mellanox Technologies



Scalability of Collective Operations - II

Offloaded Algorithm

Nonblocking Algorithm







CORE-Direct

- Scalable collective communication
- Asynchronous communication
- Manage communication by communication resources
- Avoid system noise

- Task list
- Target QP for task
- Operation
 - Send
 - Wait for completions
 - Enable
 - Calculate





NULL

Example – Four Process Recursive Doubling



© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014



Four Process Barrier Example – Using Managed Queues – Rank 0



© 2014 Mellanox Technologies



25

Nonblocking Alltoall (Overlap-Wait) Benchmark







Optimizing Non Contiguous Memory Transfers

Support combining contiguous registered memory regions into a single memory region. H/W treats them as a single contiguous region (and handles the non-contiguous regions)

- For a given memory region, supports non-contiguous access to memory, using a regular structure representation – base pointer, element length, stride, repeat count.
 - Can combine these from multiple different memory keys

Memory descriptors are created by posting WQE's to fill in the memory key

- Supports local and remote non-contiguous memory access
 - Eliminates the need for some memory copies



Optimizing Non Contiguous Memory Transfers







On Demand Paging

- No memory pinning, no memory registration, no registration caches!
- Advantages
 - Greatly simplified programming
 - Unlimited MR sizes
 - Physical memory optimized to hold current working set



ODP promise: IO virtual address mapping == Process virtual address mapping

© 2014 Mellanox Technologies





Connecting Compute Elements



Xeon

AMD









GPUDirect RDMA



HPC Advisory Council, Europe 2014



Mellanox PeerDirect[™] with NVIDIA GPUDirect RDMA

- HOOMD-blue is a general-purpose Molecular Dynamics simulation code accelerated on GPUs
- GPUDirect RDMA allows direct peer to peer GPU communications over InfiniBand
 - Unlocks performance between GPU and InfiniBand
 - This provides a significant decrease in GPU-GPU communication latency ullet
 - Provides complete CPU offload from all GPU communications across the network
- Demonstrated up to 102% performance improvement with large number of particles

HOOMD-blue Performance (LJ Liquid Benchmark, 512K Particles)











Storage

InfiniBand RDMA Storage – Get the Best Performance!

- Transport protocol implemented in hardware
- Zero copy using RDMA



5-10% the latency under 20x the workload



Data Protection Offloaded

- Used to provide data block integrity check capabilities (CRC) for block storage (SCSI)
- Proposed by the T10 committee
- DIF extends the support to main memory



© 2014 Mellanox Technologies





Power



Motivation For Power Aware Design

- Today networks work at maximum capacity with almost constant dissipation
- Low power silicon devices may not suffice to meet future requirements for low-energy networks
- The electricity consumption of datacenters is a significant contributor to the total cost of operation
 - Cooling costs scale as 1.3X the total energy consumption

Lower power consumption lowers OPEX Annual OPEX for 1KW is ~\$1000



* According to Global e-Sustainability Initiative (GeSI) if green network technologies (GNTs) are not adopted



Increasing Energy Efficiency with Mellanox Switches



Low Energy COnsumption NETworks

HPC Advisory Council, Europe 2014





Prototype Results - Power Scaling @ Link Level

Speed Reduction Power Save [SRPS]



Width Reduction Power Save [WRPS]





Prototype Results - Power Scaling @ System Level

Internal port shutdown (Director switch) - ~1%/port





Fan system analysis for improved power algorithm





number of closed ports

Monitoring and Diagnostics



Unified Fabric Manager

Automatic Discovery





Central Device Mgmt

Fabric Dashboard



Health & Perf Monitoring



Advanced Alerting

th .									
erit.	GBI	Alarm	Log File	Script	SMMP	Threshold	TTL(Sec)	Severity	
dware									
-optimal link with	4	4	4	8	8	1 +	12	😣 Hinor	*
rgested Bandwidth (%) Threshold Reached	4	2	4			10 ÷	300 +	i Hinor	*
t Bandwidth (%) Threshold Reached	4	2	1			195 🗄	300 +	i Hinor	*
dule Temperature Threshold Reached	2	2	12			60 🗄	300 🗄	O Hinor	*
k Error Recovery	12	1	1			1 ‡	300 👶	O Hinor	*
nboi Error	12	2	10			200 🗘	301 🗘	O Warning	*
t Raceive Errora	12	1	10			5 \$	301	O Minor	*
t Downed	×.	1	10			6.0	301	O Critical	•
cal Link Integrity Errors	×.	1	10			5 \$	301	Minor	•
-optimal Link Speed	×.	10	N.			1	14	😝 Minor	-
t Raceive Remote Physical Errors	V.		×.			s -	301	😝 Minor	-
w Control Update Watchdog Timer Expired	V		N.			0	1	😝 Marning	•
pical Model									
work interface Added	V.		V.			0	1	😋 Info	-
work Interface Added	10		10			1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		0.00	

Fabric Health Reports

n. uiti - Fudvic Health Report		08
ada:: 2019-15-01 11:10:50	Eabric Health Report	
Treated By :: admin	rablic fiediti Report	
legert Sammary		0
aline Summary		0
Ren unique and Zero US Values - Completed Reconstruity. 🥥		
Ren unique Rode Descriptions - Completed Sussessfully, 2 Warrings Found	9	0
9 SM Status - Completed Successfully, See details below 🥥		0
🕽 Bed Linka - Completed Seconstally. 🥥		
Las Weth - Completed Seconshilly.		
Lash Speed-Completed Successfully. 🥥		
Formate Versions - Completed Suscensibility. 🔘		
Signal Integrity (only available for QDR) - Completing Successifully. 🥥		
9 Mill Alarma - Ro Open Alarma 🚇		
		-
		Case

© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014

Congestion Analysis

Service Oriented Provisioning

ExaScale

Mellanox InfiniBand Connected Petascale Systems

Connecting Half of the World's Petascale Systems Mellanox Connected Petascale System Examples

© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014

The Only Provider of End-to-End 40/56Gb/s Solutions

From Data Center to Metro and WAN

X86, ARM and Power based Compute and Storage Platforms

The Interconnect Provider For 10Gb/s and Beyond

© 2014 Mellanox Technologies

HPC Advisory Council, Europe 2014

Thank You

Mellanox Connect. Accelerate. Outperform.™