OPEN POSSIBILITIES.

Introducing Cloud-Native Supercomputing: Bare-Metal, Secured Supercomputing Architecture



NETWORKING

Introducing Cloud-Native Supercomputing: Bare-Metal, Secured Supercomputing Architecture

Gilad Shainer, SVP Networking, NVIDIA Dhabaleswar Panda, Professor and University Distinguished Scholar, Ohio State University





Expanding Universe of High Performance Computing







Expanding Universe of High Performance Computing



GLOBAL SUMMIT NOVEMBER 9-10, 2021



In-Network Computing Accelerated Supercomputing Software-Defined, Hardware-Accelerated, InfiniBand Network



Networking Services

	High	Extremely	High
	Throughput	Low Latency	Message Rate
End-to-End	RDMA	GPUDirect RDMA	GPUDirect Storage
	Adaptive	Congestion	Smart
	Routing	Control	Topologies

Data MPI All-to-All Reductions **Tag Matching** Adapter/DPU (SHARP) Switch Data Programmable Self processing Datapath Healing units Accelerator Network (Arm cores) Data security / tenant isolation

OCP GLOBAL SUMMIT NOVEMBER 9-10, 2021

OPEN POSSIBILITIES.

In-Network Computing

Scalable Hierarchical Aggregation and Reduction Protocol (SHARP)

- In-network Tree Based Aggregation Mechanism
- Multiple Simultaneous Outstanding Operations
- Small Message and Large Message Reduction
- Barrier, Reduce, All-Reduce, Broadcast and More
- Sum, Min, Max, Min-loc, max-loc, OR, XOR, AND
- Integer and Floating-Point, 16/32/64 bits





Cloud-Native Supercomputing

- Bare-metal Secured Infrastructure
- Higher Application Performance
- From the Edge to the Main Data Center



PERFORMANCE



MULTI TENANCY



CONNECTX NETWORK INTERFACE PM CORES CCELERATION

BLUEFIELD DPU









SOLUTION PROVIDER®



Cloud-Native Supercomputing Infrastructure

Traditional Supercomputing









Multi-Tenant Isolation - Zero-Trust Architecture

- Secured Network Infrastructure and Configuration
- Storage Virtualization
- Tenant Service Level Agreement (SLA)
- 32K Concurrent Isolated Users on Single Subnet







OPEN POSSIBILITIES.

NETWORKING

Higher Applications Performance



- Parallel applications
- Collective Offloads
- Active Messages
- Smart MPI Progression
- Data Compression
- User-defined Algorithms





Requirements for Next-Generation MPI Libraries

- Message Passing Interface (MPI) libraries are used for HPC and AI applications
- Requirements for a high-performance and scalable MPI library:
 - Low latency communication

- High bandwidth communication
- Minimum contention for host CPU resources to progress non-blocking collectives
- High overlap of computation with communication
- CPU based non-blocking communication progress can lead to subpar performance as the main application has less CPU resources for useful application-level computation





Can MPI Functions be Offloaded?

- The area of network offloading of MPI primitives is still nascent and cannot be used as a universal solution
- State-of-the-art BlueField DPUs bring more compute power into the network
- Can we exploit additional compute capabilities of modern BlueField DPUs into existing MPI middleware to extract
 - Peak pure communication performance
 - Overlap of communication and computation
- For dense non-blocking collective communications?







Overview of BlueField-2 DPU

- ConnectX-6 network adapter with 200Gbps InfiniBand
- System-on-chip containing eight 64bit ARMv8 A72 cores with 2.75 GHz each
- 16 GB of memory for the ARM cores
- How can one re-design an MPI library to take advantage of DPUs and accelerate scientific applications?





NETWORKING



MVAPICH2 Software Architecture with DPU

High Performance Parallel Programming Models				
Message Passing Interface	PGAS	Hybrid MPI + X		
(MPI)	(UPC, OpenSHMEM, CAF, UPC++)	(MPI + PGAS + OpenMP/Cilk)		





Proposed Offload Framework

- Non-blocking collective operations are offloaded to a set of "worker processes"
- BlueField is set to separated host mode
- Worker processes are spawned to the ARM cores of BlueField
- Once the application calls a collective, host processes prepare a set of metadata and provide it to the Worker processes
- Using these metadata, worker processes can access host memory through RDMA
- Worker processes progress the collective on behalf of the host processes
- Once message exchanges are completed, worker processes notify the host processes about the completion of the non-blocking operation







Proposed Non-blocking Alltoall Design

- Worker process performs RDMA Read to receive the data chunk from host main memory
- Once data is available in the ARM memory, worker process performs RDMA Write to the remote host memory





NETWORKING

Proposed Non-blocking Alltoall Design

- Example: Scatter Destination Algorithm
- Focus is on medium and large messages
- Message chunking and pipelining is utilized to reduce the overheads of staging





NETWORKING

Experimental Setup for Performance Evaluation

- HPC Advisory Council High-Performance Computing Center
- Cluster has 32 compute-node with Broadwell series of Xeon dualsocket, 16-core processors operating at 2.60 GHz with 128 GB RAM
- NVIDIA BlueField-2 adapters are equipped with 8 ARM cores operating at 2.0 GHz with 16 GB RAM
- Based on the MVAPICH2-DPU MPI library
- OSU Micro Benchmark for nonblocking Alltoall and P3DFFT Application





NFTWORKING

OSU Micro benchmark ialltoall

- osu_ialltoall benchmark metrics
- Pure communication time
 - Latency t is measured by calling MPI_Ialltoall followed by MPI_Wait
- Total execution time
 - Total T = MPI_Ialltoall + synthetic compute + MPI_Wait
- Overlap
 - Benchmark creates a synthetic computation block that takes t microsecond to finish. Before starting compute, MPI_Ialltoall is called and after that MPI_Wait. Overlap is calculated based on total execution time and compute time.
- Part of the standard OSU Micro-Benchmark

OPEN POSSIBILITI<mark>ES</mark>.





Overlap of Communication and Computation with osu_Ialltoall (32 nodes)





NETWORKING

Total Execution Time with osu_Ialltoall (32 nodes)







32 Nodes, 16 PPN

32 Nodes, 32 PPN

Benefits in Total execution time (Compute + Communication)



P3DFFT Application Execution Time (16 nodes)





NETWORKING

OPEN POSSIBILITI<mark>ES</mark>.

Total Execution Time with osu_Iallgather (16 nodes)

Total Execution Time, BF-2 (osu_iallgather) MVAPICH2 MVAPICH2-DPU 7.00 41% 6.00 5.00 Comm. Time (ms) 00°8 24% 39% 2.00 23% 1.00 0.00 1M 2M 256K 512K Message Size

Total Execution Time, BF-2 (osu_iallgather)



16 Nodes, 1 PPN

)PEN POSSIBILITIES.

GLOBAL SUMMIT NOVEMBER 9-10, 2021



Total Execution Time with osu_Ibcast (16 nodes)



NETWORKING





MVAPICH2 MVAPICH2-DPU



16 Nodes, 16 PPN

PEN POSSIBILITIES.

16 Nodes, 32 PPN



NVIDIA DOCA SDK

- Software Application Framework for BlueField DPUs
- DOCA is for DPUs What CUDA is for GPUs
- Protects Developer Investment for Future DPUs
- Certified Reference Applications, APIs & Partner Solutions
- Rich Partner Ecosystem Across Industries and Workloads
- <u>https://developer.nvidia.com/networking/doca</u>







OPEN POSSIBILITI<mark>ES</mark>.

Call to Action – More Information

- Technical overview <u>https://nvdam.widen.net/s/plhzlwmtrg/tech-overview-infiniband-cloud-native-supercomputing-web</u>
- Community open source software development -<u>https://ucfconsortium.org/projects/opensnapi/</u>
- The MVAPICH2-DPU MPI library <u>https://x-scalesolutions.com/mvapich2-dpu/</u>
- News announcements and more links -<u>https://www.hpcwire.com/2021/04/14/gtc21-dell-building-cloud-native-</u> <u>supercomputers-at-u-cambridge-and-durham/</u>



Thank you!

