



MVAPICH2 at Azure: Enabling High Performance on Cloud

OSU Booth, SC 2022

Jithin Jose, Microsoft
jijos@microsoft.com

Agenda

- Overview of Azure HPC SKUs
 - Azure HBv4, NDv4
- Feature Highlights
- MVAPICH2 on HBv4
- MVAPICH2 GDR on NDv4
- Azure HPC VM Image
- Performance Highlights
- Conclusion

Azure HPC & AI breakthroughs

Demonstrating innovation leadership for cloud HPC

- 2019 ● 20,000 cores MPI job - 1st for cloud
- 2019 ● AMD EPYC Rome InfiniBand HPC clusters - 1st for cloud
- 2019 ● 200 Gb/s HDR InfiniBand with adaptive routing - 1st for cloud
- 2020 ● 80,000 cores MPI job - 1st for cloud, 12x higher than any other cloud
- 2020 ● Top5-scale supercomputer for OpenAI (CPU+GPU) - 1st for cloud
- 2020 ● 1 TB/sec Parallel File System - 1st for cloud
- 2021 ● 1.6 Tb/s InfiniBand for NVIDIA A100 clusters - 1st for cloud
- 2021 ● HBv3 Milan GA in Azure *at* AMD Launch - 1st for cloud *or* on-prem
- 2021 ● 10-year Supercomputing-as-a-Service for Met Office - 1st for cloud
- 2021 ● #10 supercomputer + 4 more in Top40 - 1st for cloud
- 2021 ● #1 cloud on MLPerf benchmark + #2 overall - 1st for cloud
- 2022 ● HBv4 Genoa in Azure *at* AMD Launch - 1st for cloud *or* on-prem
- 2022 ● 400 GB/s NDR InfiniBand proven HPC interconnect - 1st for cloud
- 2022 ● *More to come ... and it will probably be 1st for cloud!*

Azure HPC/AI VM Series with InfiniBand



Azure HPC VMs

Standard HPC Applications

High Compute/Memory + InfiniBand

*HPC SKUs: H, HB, HC, HBv2, HBv3, HBv4, Hxv1



Azure AI/Deep Learning VMs

Deep Learning, AI workloads

Visualization SKUs: NV series

*Deep Learning/AI SKUs (InfiniBand): NC, ND series

- "r" in VM type indicates RDMA support
- InfiniBand exposed to VMs using SR-IOV, offers full host bypass with full feature support
- *InfiniBand/RDMA enabled VMs: One VM per Host

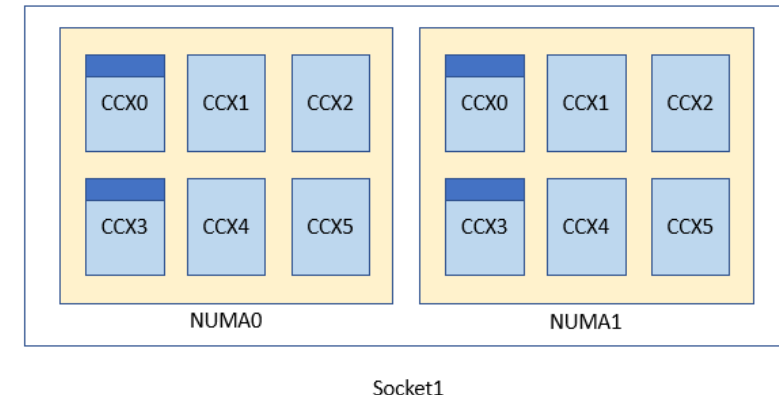
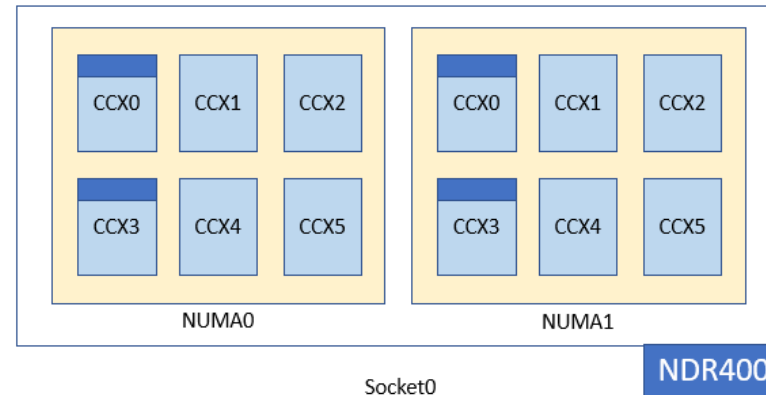
Azure HBv4, HXv1 with Genoa (public preview upcoming)



AMD Genoa



NVIDIA
InfiniBand
NDR 400Gbps



■ Hyper-V Partition (4 cores per NUMA)

• VM Specs:

- AMD Genoa (NPS = 2)
- VM Cores: 176
- Memory: 704 GB, 1408 GB
- Local Disk: 2 x 1.8 TB NVMe SSD
- **Network: 400 Gbps NDR (SR-IOV)**

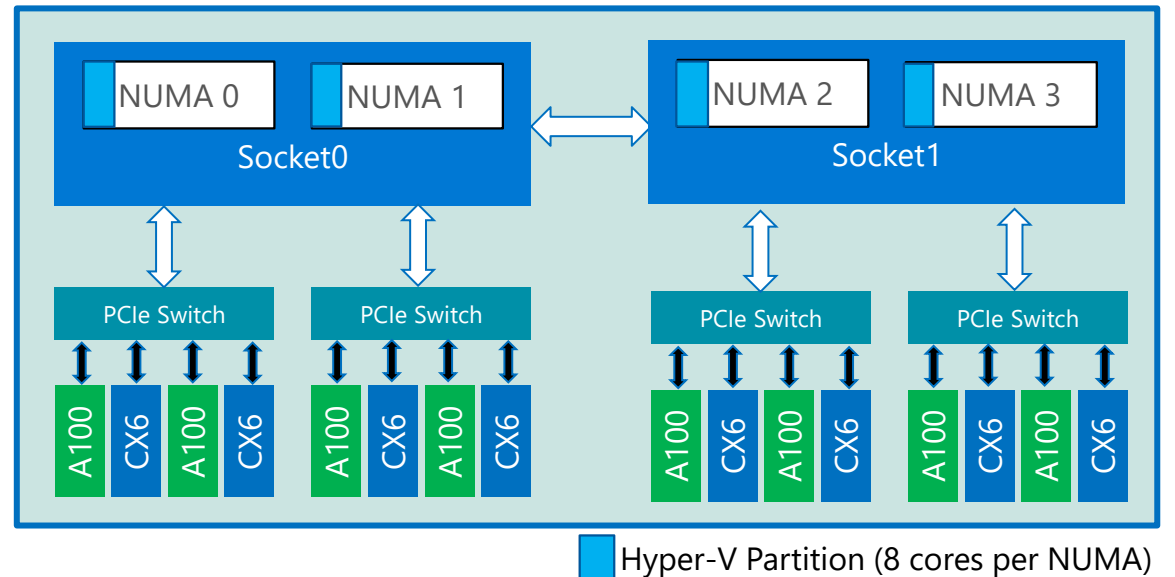
HBv4 VM Sizes (one VM per Host):

- Standard_HB176rs_v4 (all 120 cores)
- Standard_HB176-144rs_v3 (6 cores per CCD)
- Standard_HB176-96rs_v3 (4 cores per CCD)
- Standard_HB176-48rs_v3 (2 cores per CCD)
- Standard_HB176-24rs_v3 (1 cores per CCD)

Azure NDv4

- VM Specs:

- AMD Rome (NPS=2)
- VM Cores: 96 (48 per socket)
- Memory: 900 GB
- 8 x NVIDIA A100 GPUs (NVLink 3.0)
- 8 x HDR 200Gbps InfiniBand
- Local Disk: 6.4 TB local NVMe SSD



Standard_ND96asr_v4 (NDv4)

Ideal for AI/Deep learning workloads

Agenda

- Overview of Azure HPC SKUs
 - Azure HBv4, NDv4
- **Feature Highlights**
- MVAPICH2 on HBv4
- MVAPICH2 GDR on NDv4
- Azure HPC VM Image
- Performance Highlights
- Conclusion

InfiniBand Features in Azure

- **HB, HC, NDv2:**

- EDR 100 Gb/s InfiniBand
- Up to 200 M messages/second

ConnectX-5

- **HBv2, HBv3, NDv4:**

- HDR 200 Gb/s InfiniBand
- Up to 215 M messages/second

ConnectX-6

- **HBv4, HXv1:**

- NDR 400 Gb/s InfiniBand
- Up to 330 M messages/second

ConnectX-7

- **Dynamically Connected Transport (DCT)**

- Reliable and scalable transport
- Lesser Memory footprint

- **Hardware offload**

- Collectives offload framework
- Hardware tag matching

- **UD multicast (MCAST)**

- Unreliable datagram (UD) based multicast
- Create a mcast group and broadcast

- **SHARP**

- Switch based collectives

- **Dynamic Routing**

- Advanced Congestion Control
- Adaptive Routing

- **Better Reliability**

- SHIELD detects link failures and reroutes

GPUDirect RDMA

- Available on Azure NDv4
- Direct data path b/w A100 GPU and HDR200
- Each NIC/GPU pair gets peak b/w simultaneously
- Combined GPUDirect RDMA b/w of **1.6 Tbps**
- Supports **all** GDR capable MPI libraries/middleware (including MVAPICH2-GDR)

```
hpcadmin@compute000000:~$ ./test_ib_gpu.sh compute000000 compute000001 cpu
Pair 0:
8388608 2922 0.00 196.09 0.002922
8388608 2920 0.00 195.96 0.002920
Pair 1:
8388608 2928 0.00 196.49 0.002928
8388608 2930 0.00 196.63 0.002930
Pair 2:
8388608 2894 0.00 194.21 0.002894
8388608 2896 0.00 194.34 0.002896
Pair 3:
8388608 2883 0.00 193.47 0.002883
8388608 2881 0.00 193.34 0.002881
Pair 4:
8388608 2893 0.00 194.14 0.002893
8388608 2895 0.00 194.28 0.002895
Pair 5:
8388608 2883 0.00 193.47 0.002883
8388608 2885 0.00 193.61 0.002885
Pair 6:
8388608 2922 0.00 196.09 0.002922
8388608 2920 0.00 195.96 0.002920
Pair 7:
8388608 2916 0.00 195.48 0.002916
8388608 2915 0.00 195.62 0.002915
```

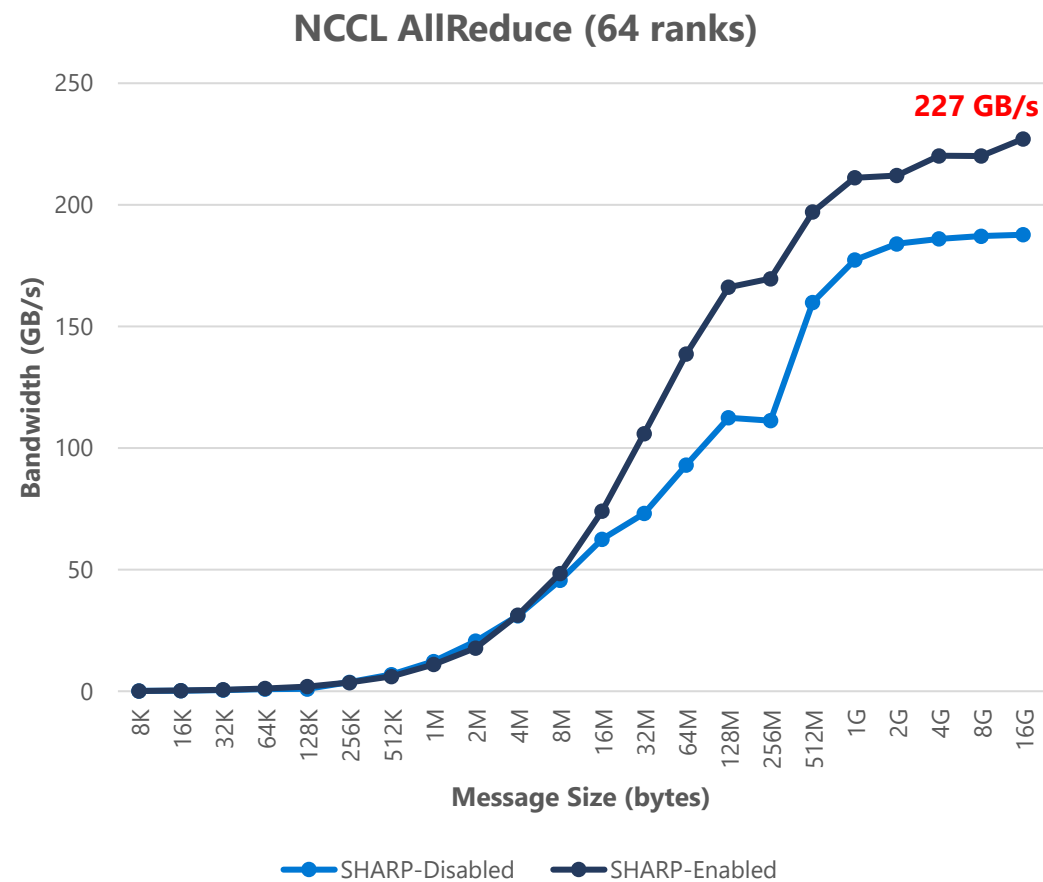
RDMA (Host Memory)

```
hpcadmin@compute000000:~$ ./test_ib_gpu.sh compute000000 compute000001 gpu
Pair 0:
8388608 2913 0.00 195.49 0.002913
8388608 2913 0.00 195.49 0.002913
Pair 1:
8388608 2914 0.00 195.55 0.002914
8388608 2914 0.00 195.55 0.002914
Pair 2:
8388608 2914 0.00 195.55 0.002914
8388608 2914 0.00 195.55 0.002914
Pair 3:
8388608 2915 0.00 195.62 0.002915
8388608 2915 0.00 195.62 0.002915
Pair 4:
8388608 2914 0.00 195.55 0.002914
8388608 2914 0.00 195.55 0.002914
Pair 5:
8388608 2915 0.00 195.62 0.002915
8388608 2915 0.00 195.62 0.002915
Pair 6:
8388608 2914 0.00 195.55 0.002914
8388608 2914 0.00 195.55 0.002914
Pair 7:
8388608 2915 0.00 195.62 0.002915
8388608 2915 0.00 195.62 0.002915
hpcadmin@compute000000:~$
```

GPUDirectRDMA (GPU Memory)

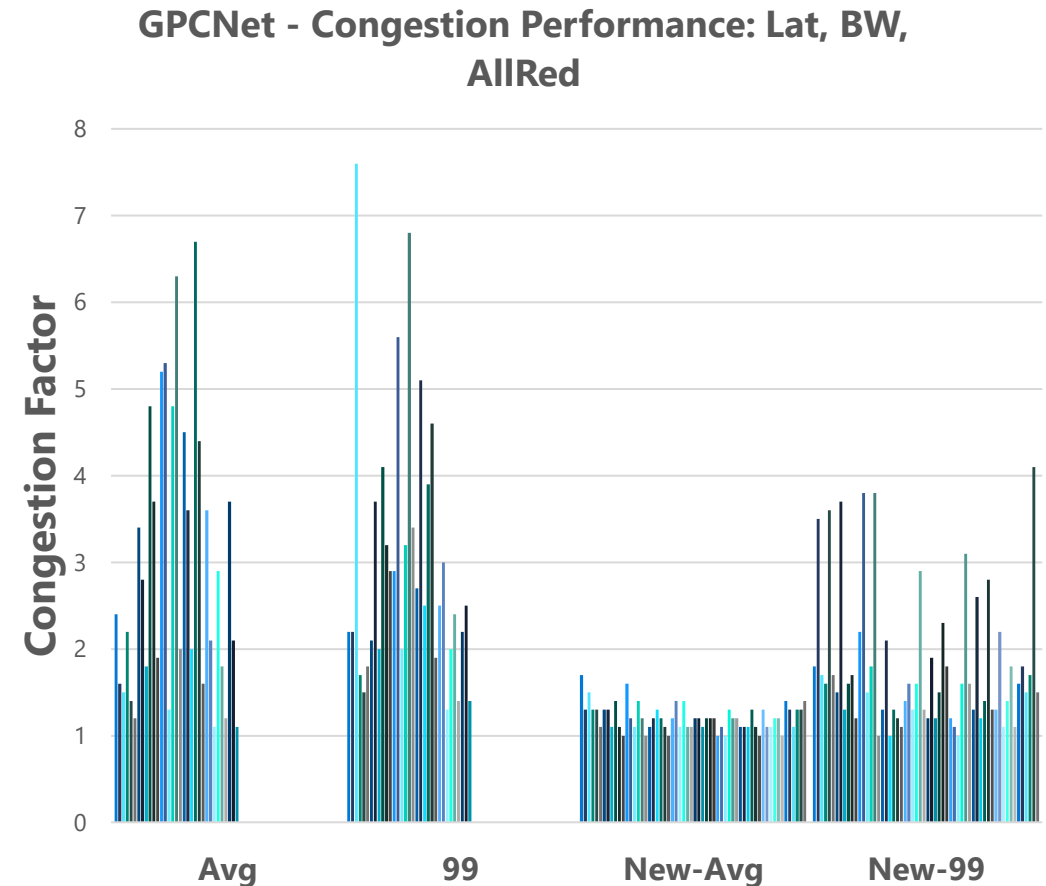
SHARP

- Enabled on dedicated NDv4 clusters
- UCX-based Sharp-AM / SharpD communication
- Optimized SHARP tree initialization
- Connection keepalive
- GRH support



Congestion Control

- Available on all VM Series with HDR200
- Transparent to customer applications
- Improve tail latencies
- Critical in public multi-customer environments

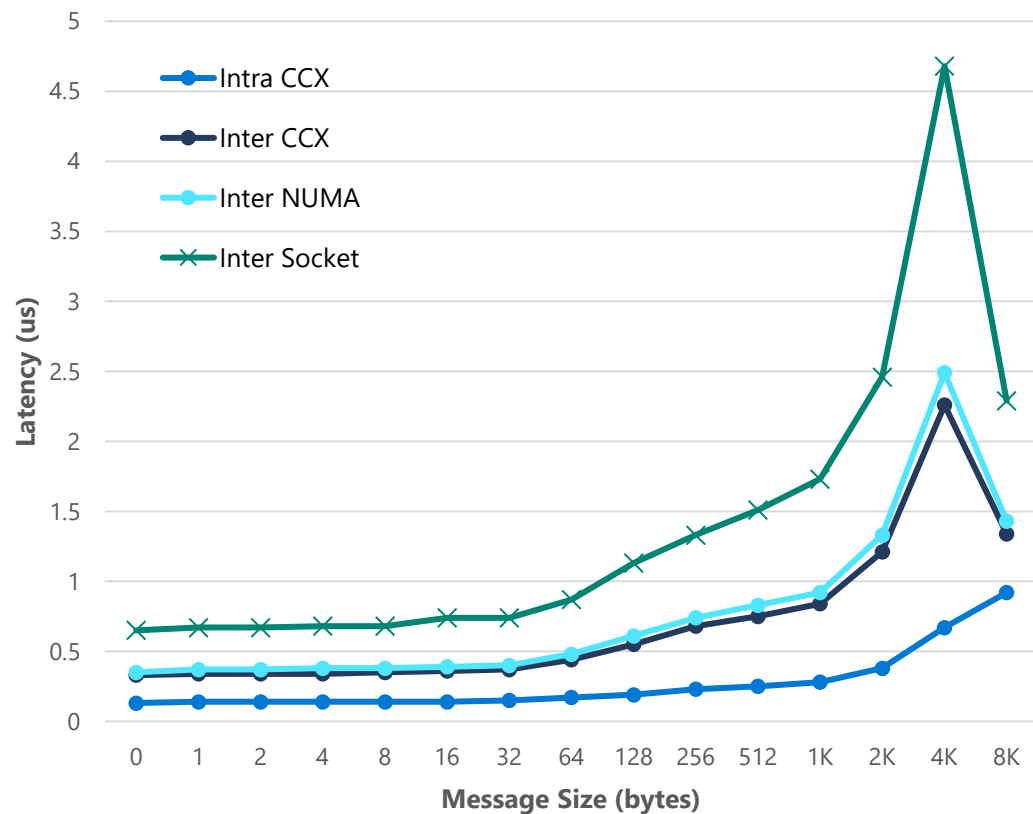


Agenda

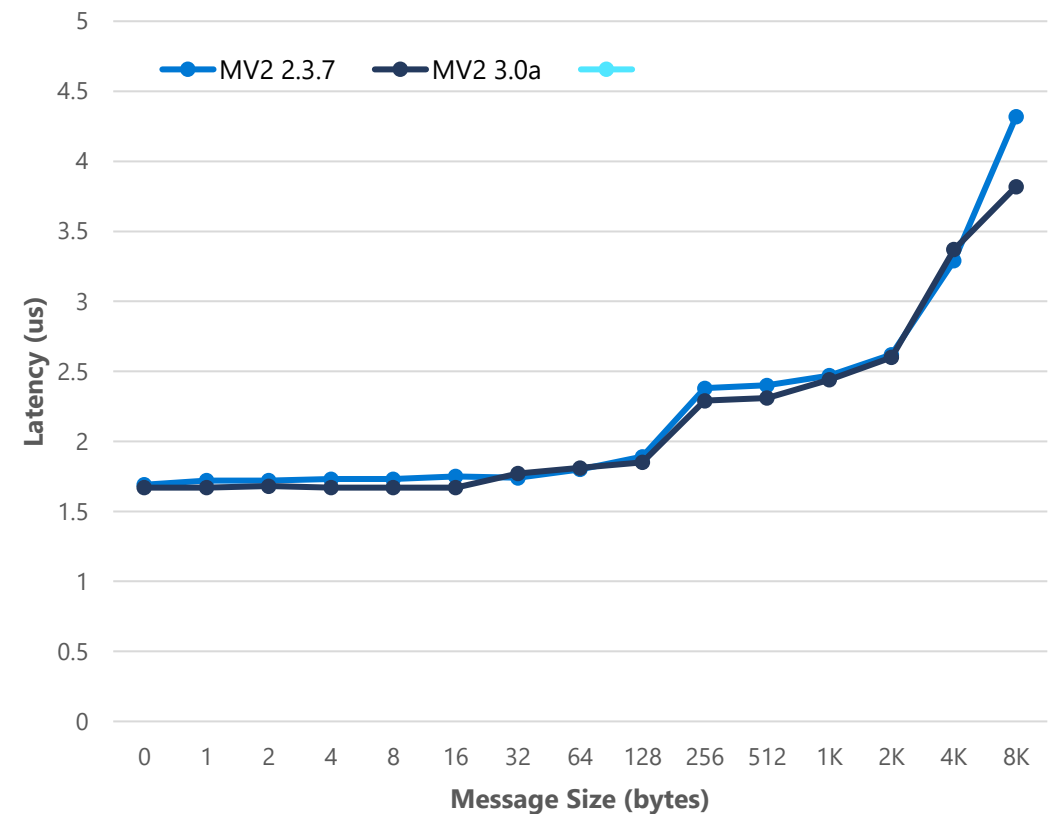
- Overview of Azure HPC SKUs
 - Azure HBv4, NDv4
- Feature Highlights
- **MVAPICH2 on HBv4**
- MVAPICH2 GDR on NDv4
- Azure HPC VM Image
- Performance Highlights
- Conclusion

MVAPICH2 on HBv4: MPI Latency

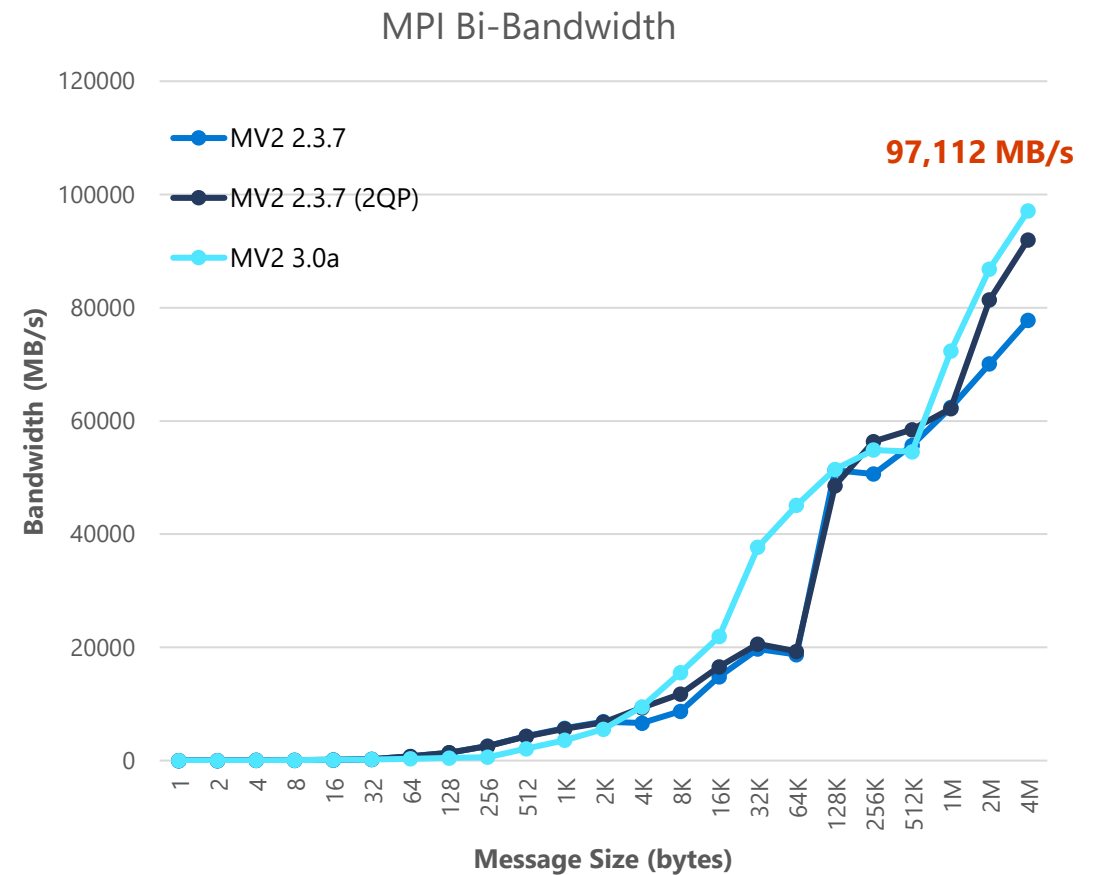
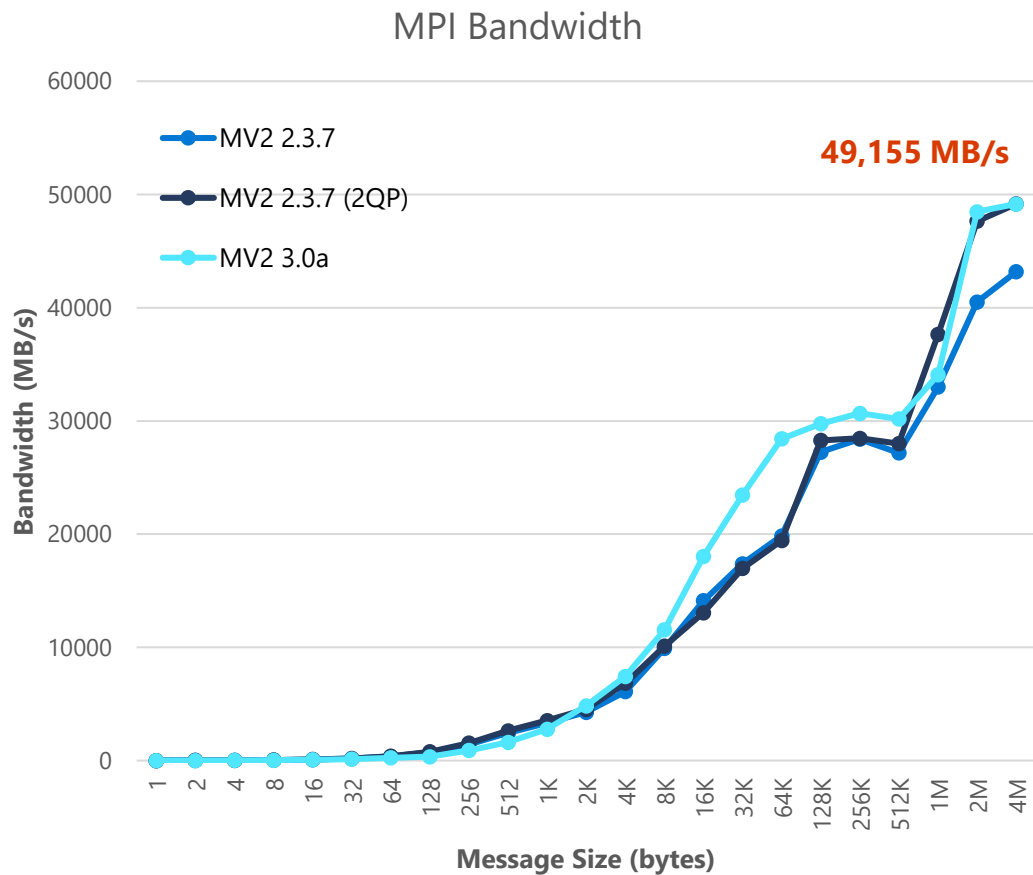
MPI Latency (small messages - intra node)



MPI Latency (small messages - inter node)



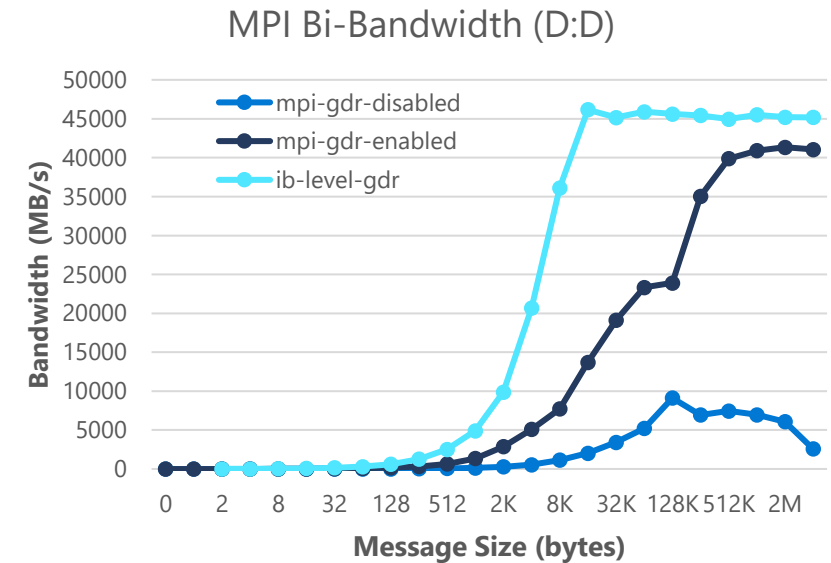
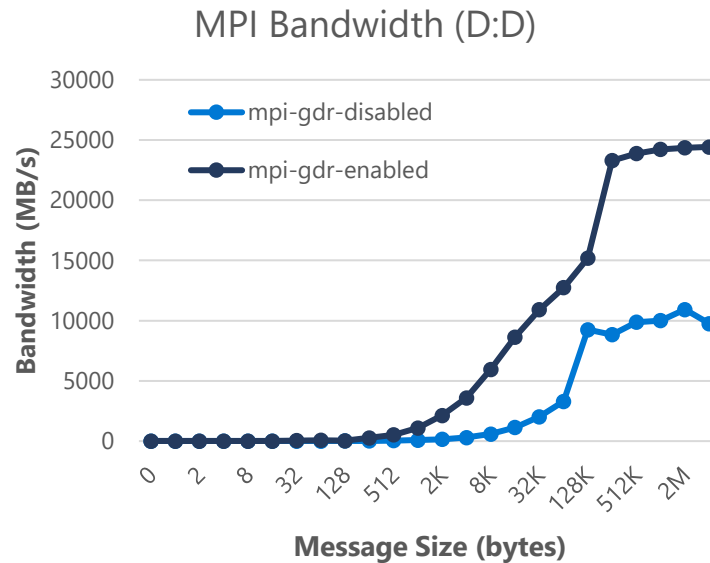
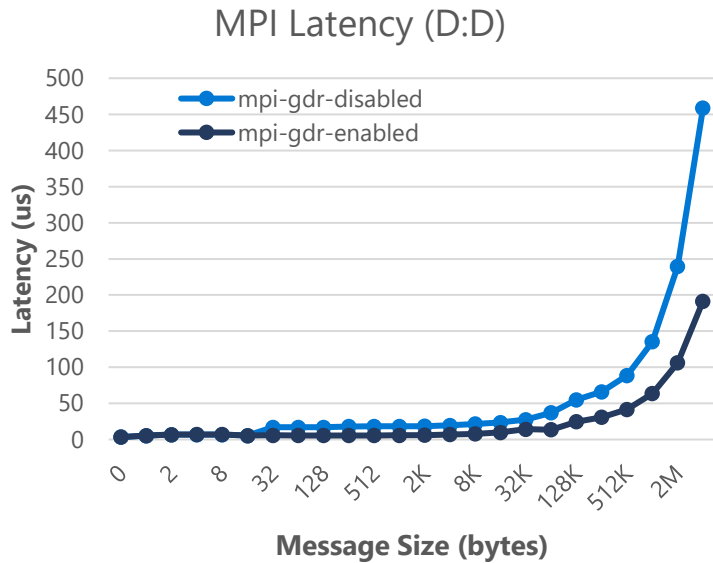
MVAPICH2 on HBv4 (NDR): MPI Bandwidth



Agenda

- Overview of Azure HPC SKUs
 - Azure HBv4, NDv4
- Feature Highlights
- MVAPICH2 on HBv4
- **MVAPICH2 GDR on NDv4**
- Azure HPC VM Image
- Performance Highlights
- Conclusion

MVAPICH2-GDR on NDv4



Software Configuration: MVAPICH2 2.3.7-GDR on Azure [Ubuntu-HPC 18.04 VM Image](#)

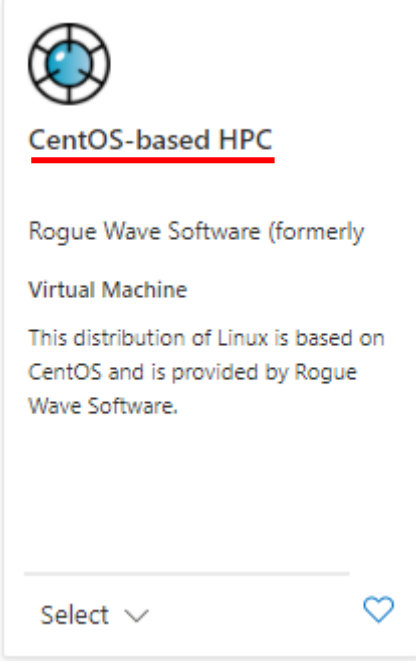
Environment parameters: MV2_NUM_QP_PER_PORT=4 MV2_IBA_EAGER_THRESHOLD=66560 MV2_VBUF_TOTAL_SIZE=66560 MV2_RNDV_PROTOCOL=RPUT
MV2_CUDA_BLOCK_SIZE=131072 MV2_USE_GPUDIRECT_RDMA=1 MLX5_RELAXED_PACKET_ORDERING_ON=all MV2_GPUDIRECT_LIMIT=4194304 MV2_USE_CUDA=1
MV2_IBA_HCA=mlx5_ib0 CUDA_VISIBLE_DEVICES=0



Agenda

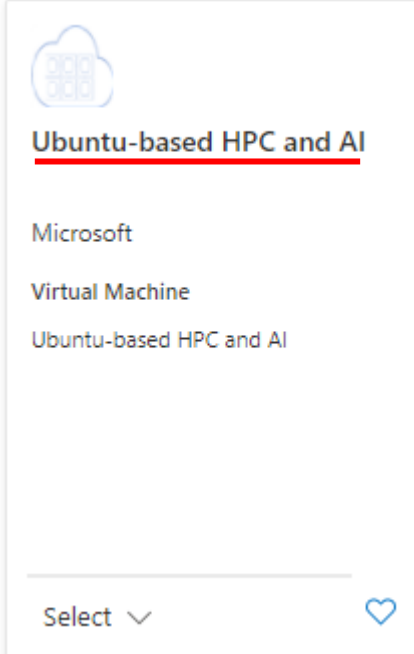
- Overview of Azure HPC SKUs
 - Azure HBv4, NDv4
- Feature Highlights
- MVAPICH2 on HBv4
- MVAPICH2 GDR on NDv4
- [Azure HPC VM Images](#)
- Performance Highlights
- Conclusion



Azure HPC VM Images

- Optimized VM Images for HPC/AI workloads
- Mellanox OFED
- Pre-configured IPoIB InfiniBand based MPI Libraries
 - HPC-X, IntelMPI, **MVAPICH2**, OpenMPI
- Communication Runtimes
 - Libfabric, OpenUCX
- Optimized libraries
 - Blis, FFTW, Flame, MKL
- Recommended Compilers
- GPU Drivers
- NCCL, NCCL RDMA Sharp Plugin, SharpD
- Other optimizations




CentOS-based HPC
Rogue Wave Software (formerly
Virtual Machine
This distribution of Linux is based on
CentOS and is provided by Rogue
Wave Software.
Select 



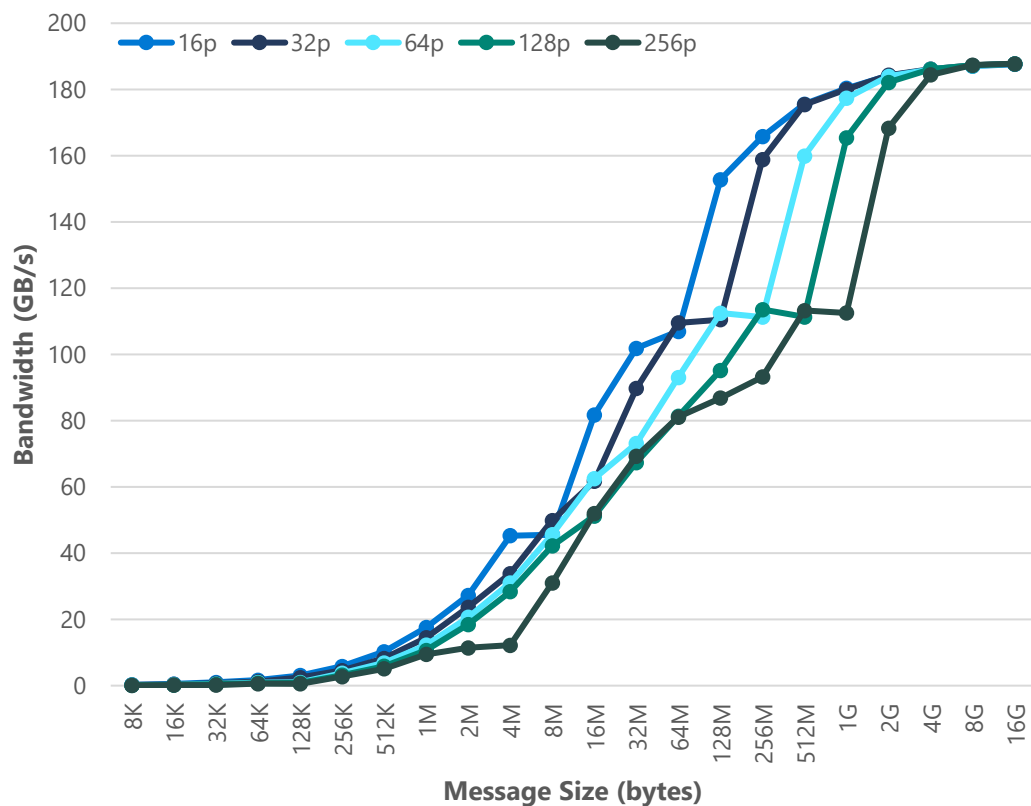

Ubuntu-based HPC and AI
Microsoft
Virtual Machine
Ubuntu-based HPC and AI
Select 

Agenda

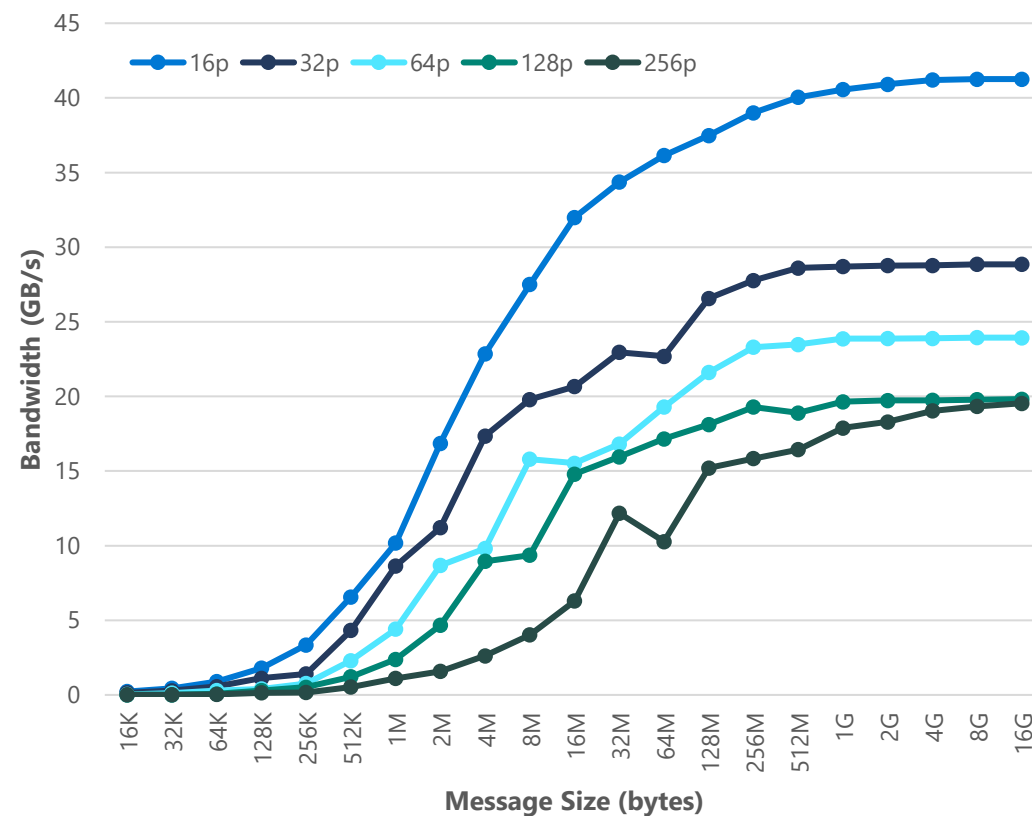
- Overview of Azure HPC SKUs
 - Azure HBv3, NDv4
- Feature Highlights
- MVAPICH2 on HBv3
- MVAPICH2 GDR on NDv4
- Azure HPC VM Images
- Performance Highlights
- Conclusion

NCCL on NDv4

NCCL AllReduce (w/o SHARP)

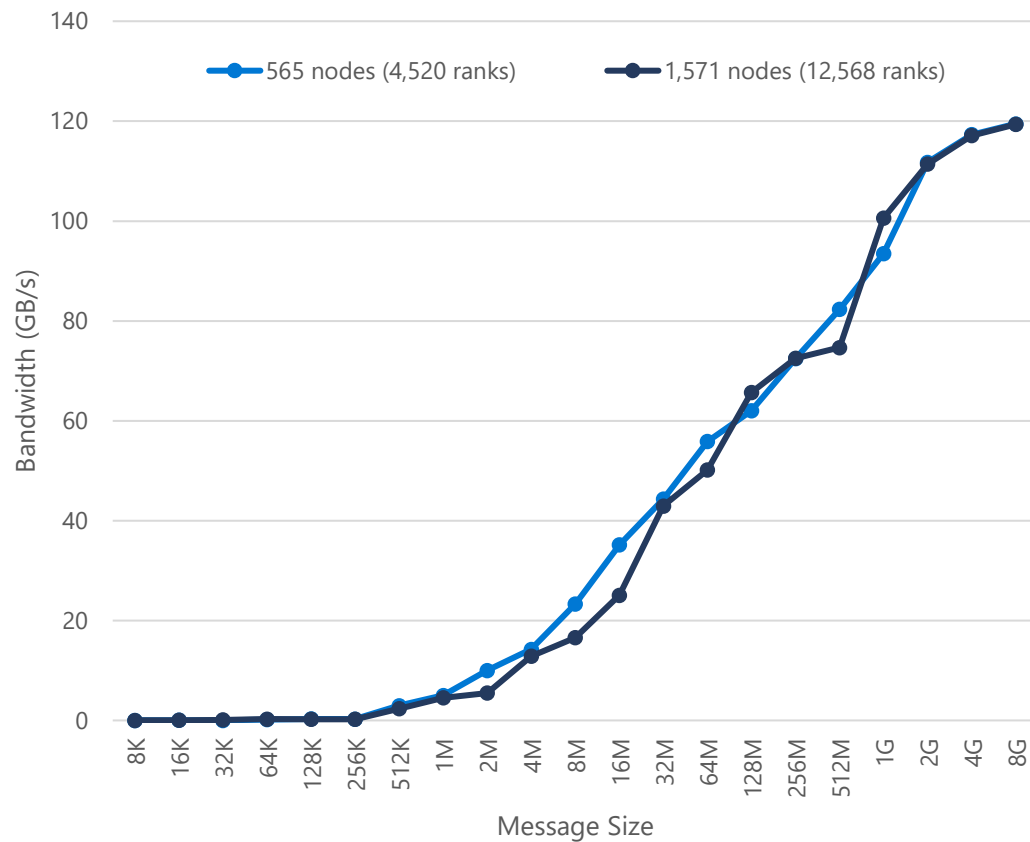


NCCL AlltoAll

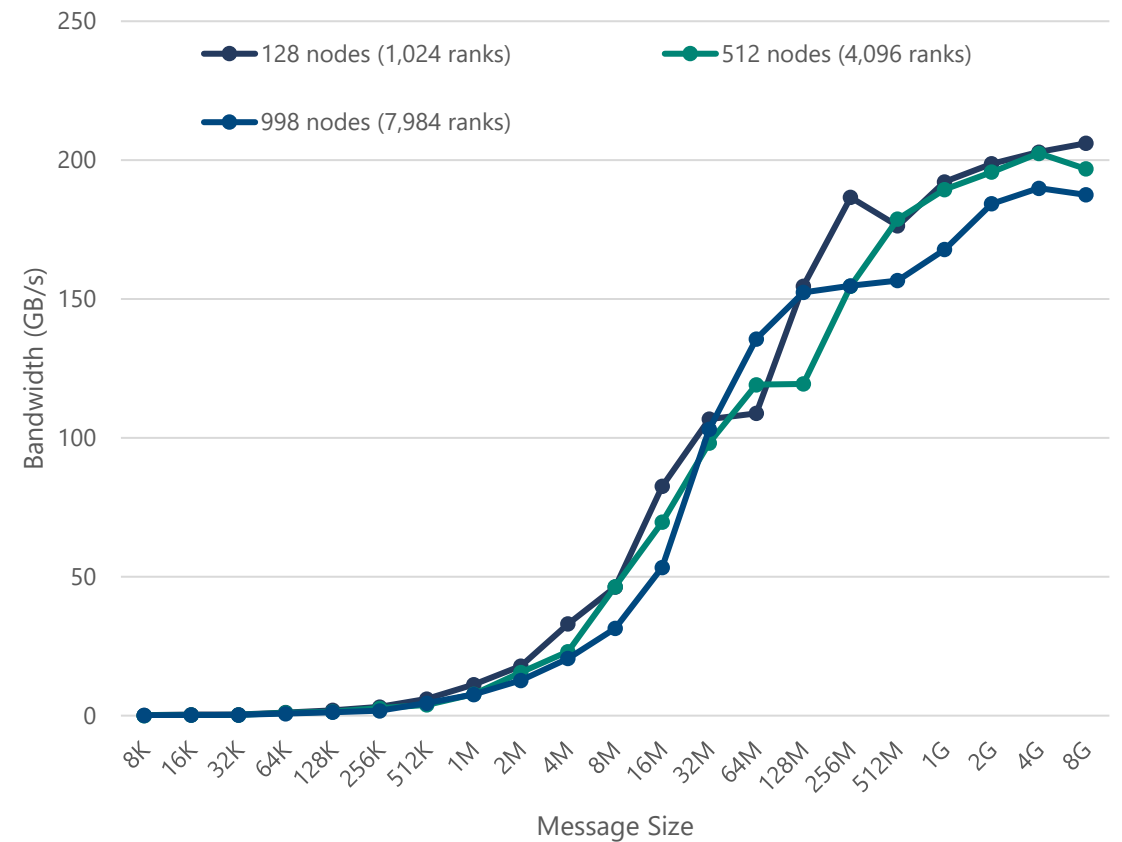


NCCL at Scale on NDv4

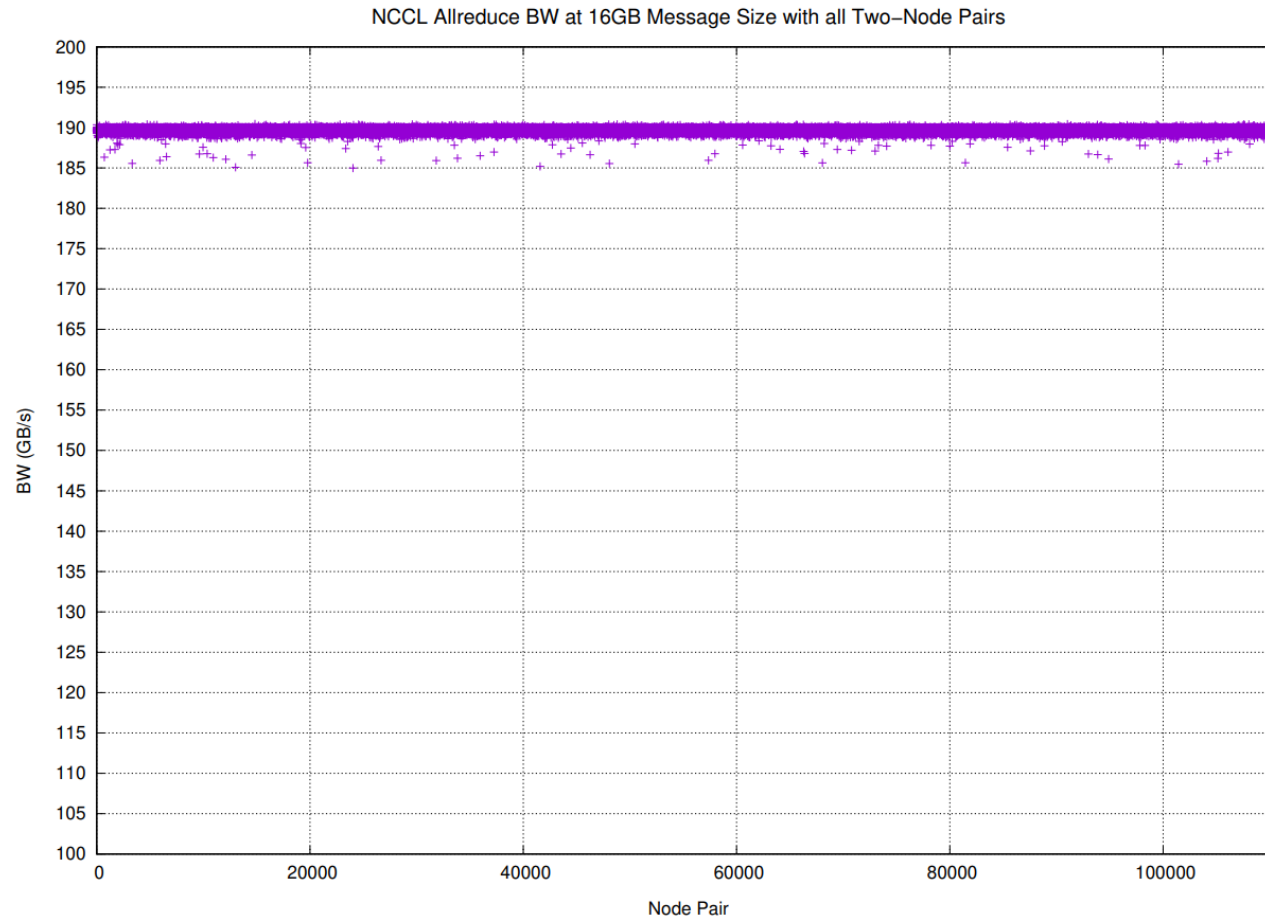
NCCL AllReduce (w/o SHARP)



NCCL AllReduce w/ SHARP



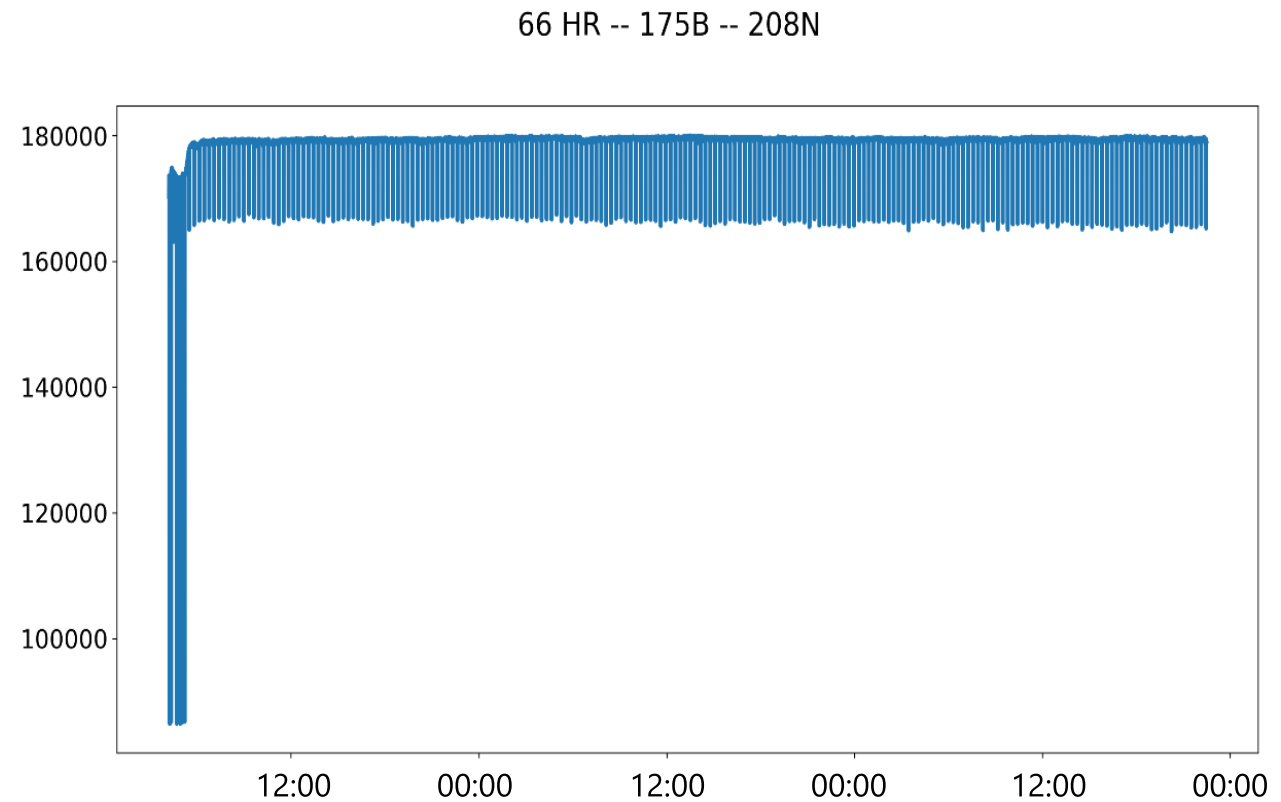
NCCL AllReduce Bandwidth Distribution



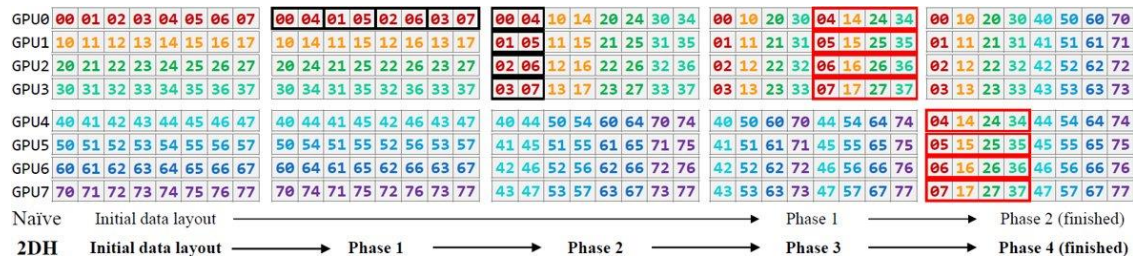
- Azure InfiniBand Clusters deploy Non-blocking (under-subscribed) fat-tree topology
- Evaluation using all-pair NCCL AllReduce benchmark
- Cluster size = ~470 NDv4 (8 x A100, 8 x 200 Gbps HDR) nodes
- Multiple pairs ($N/2$) communicating at the same time
- 100% pairs achieve > 186 GB/s

Sustained Multi-node Training Performance

- MetaSeq Training Workload on NDv4
 - [facebookresearch/metaseq: Repo for external large-scale work \(github.com\)](https://github.com/facebookresearch/metaseq)
- 175B OPT Model
- 208 nodes (1664 A100 GPUs) over InfiniBand
- Delivered sustained training throughput over multiple days

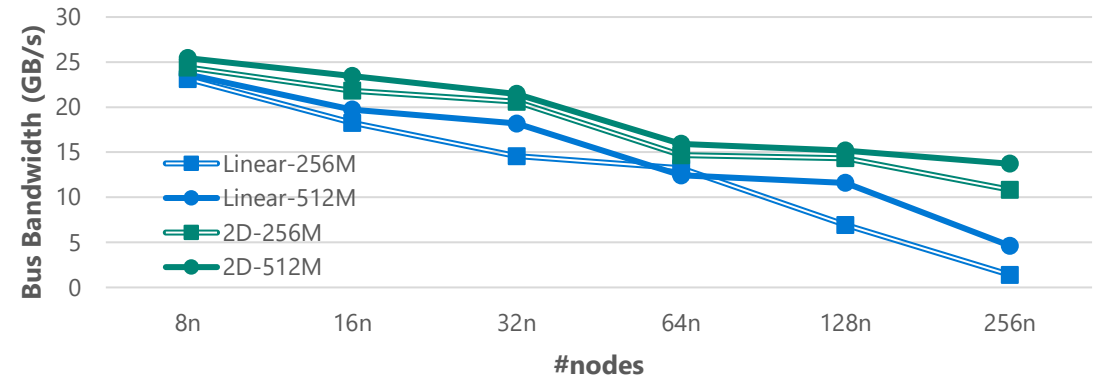


Tutel: Adaptive MoE at Scale

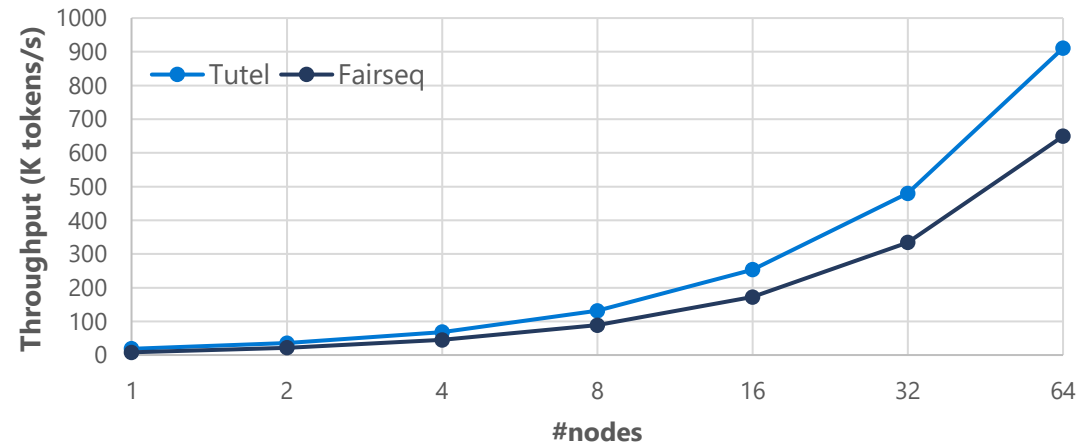


- New AlltoAll algorithm optimized for NDv4/NDmv4 cluster
 - Larger slice through IB => 8x slice size in large scale
 - Only 1-1 IB interconnection required in inter-node aggregation phase
 - Open-source on github.com/microsoft/msccl
 - Achieve **>6.7x** gain on 256MiB and **>1.9x** gain on 512MiB with 256 NDmv4 nodes
- New AlltoAll algorithm + Other framework optimizations: > 40% E2E performance improvement

AlltoAll Bus Bandwidth (Linear vs 2D Hierarchical Algorithm)

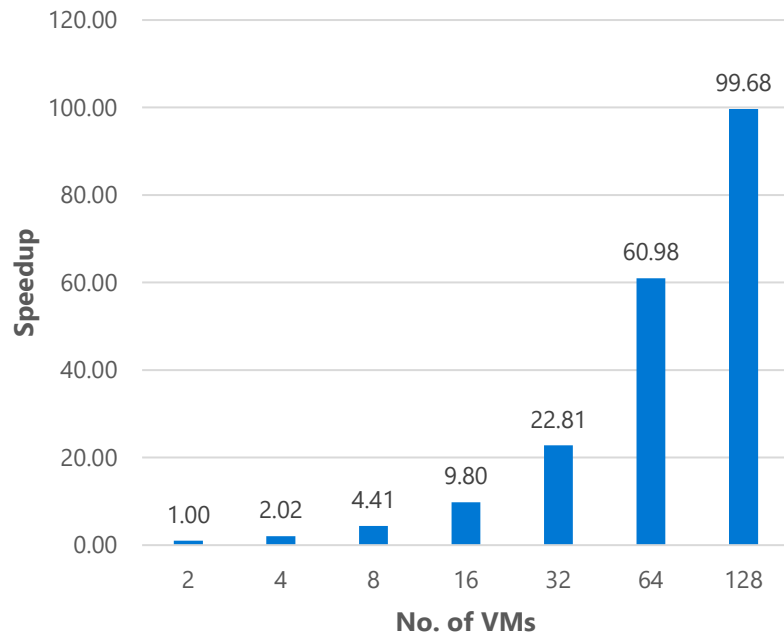


Meta GPT-3 MoE Model - E2E Throughput



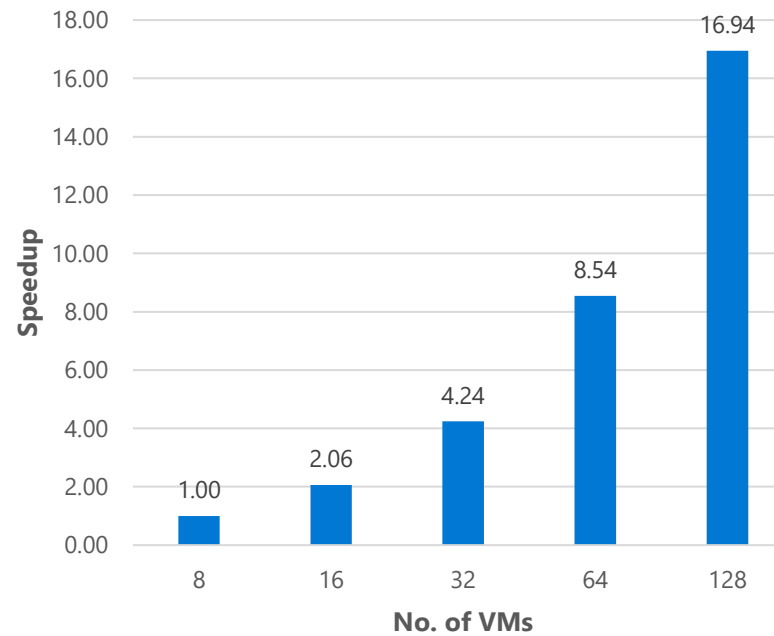
Scaling Efficiency on HBv3 (Milan-X)

Ansys Fluent 2021 R1
f1_racecar_140m



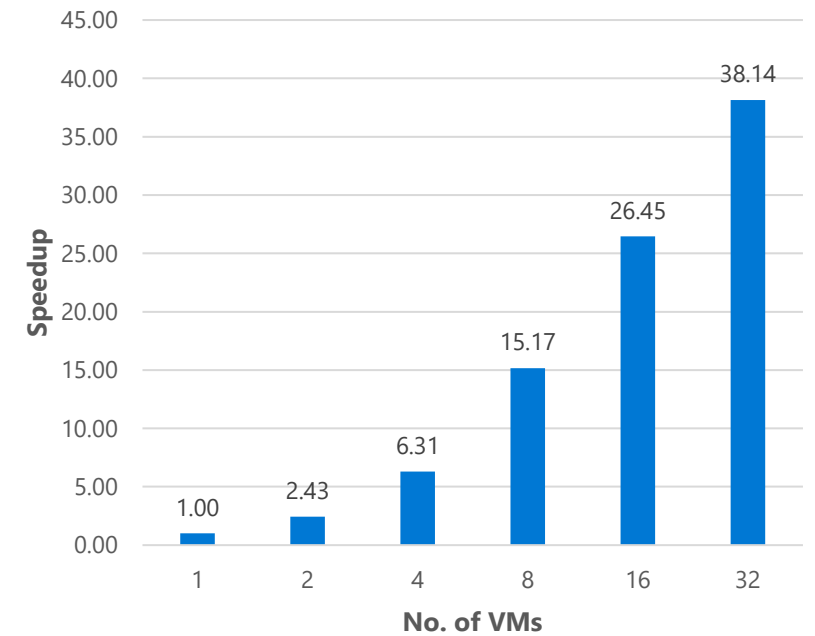
156% scaling efficiency

Ansys Fluent 2021 R1
f1_combustor_830m



106% scaling efficiency

OpenFOAM v. 1912
Motorbike 28m



119% scaling efficiency

<https://aka.ms/MilanXPerf>

Agenda

- Overview of Azure HPC SKUs
 - Azure HBv4, NDv4
- Feature Highlights
- MVAPICH2 on HBv4
- MVAPICH2 GDR on NDv4
- Azure HPC VM Images
- Performance Highlights
- **Conclusion**

Azure: *the* cloud purpose-built for HPC & AI

- ✓ **Genuine HPC approach**
platforms, benchmarks, people,
and end-to-end experience
- ✓ **Purpose-built platforms** for
best performance, and best price-
performance, and differentiated solutions
- ✓ **Leading time-to-market** for key hardware
innovations to accelerate
time-to-solution for customers
- ✓ **Partnering with customers** for the long
term to solve HPC and business needs



Supercomputing
for the most
demanding
applications

**InfiniBand
HPC & AI**
clusters for best
performance on
real workloads

Compute
optimized VMs
with "low"
latency networks

Azure

Azure is the
only public cloud
provider offering
the full range of
HPC and AI
capabilities

Compute
optimized VMs
with "low"
latency networks

Other clouds

Conclusion

- Supercomputer on Cloud is real!
- Azure HPC Cloud made into Top500, Graph500
- High Performance middleware such as MVAPICH2 enables cutting edge technology
 - Deliver High Scalability and Performance

Pointers

Getting Started

- [High Performance Computing \(HPC\) on Azure](#)

HPC VM Series

- [Azure VM sizes - HPC - Azure Virtual Machines](#)

GPU VM Series

- [Azure VM sizes - GPU - Azure Virtual Machines](#)

HPC VM Images

- [Azure HPC VM Images](#)
- [GitHub Repository](#)

HPC VM Deployment

- [Sample HPC VM deployment scripts](#)
- [Azure CycleCloud](#)
- [MUG '20 Tutorial](#)

Azure HPC Blogs

- [Azure Compute - Microsoft Tech Community](#)