



# UCX on Azure HPC and AI Clusters

Jithin Jose, Microsoft  
jjjos@microsoft.com

# Agenda



## Overview of Azure HPC



Azure HBv3, NDv4



Network features



Azure HPC VM Images



Performance Highlights

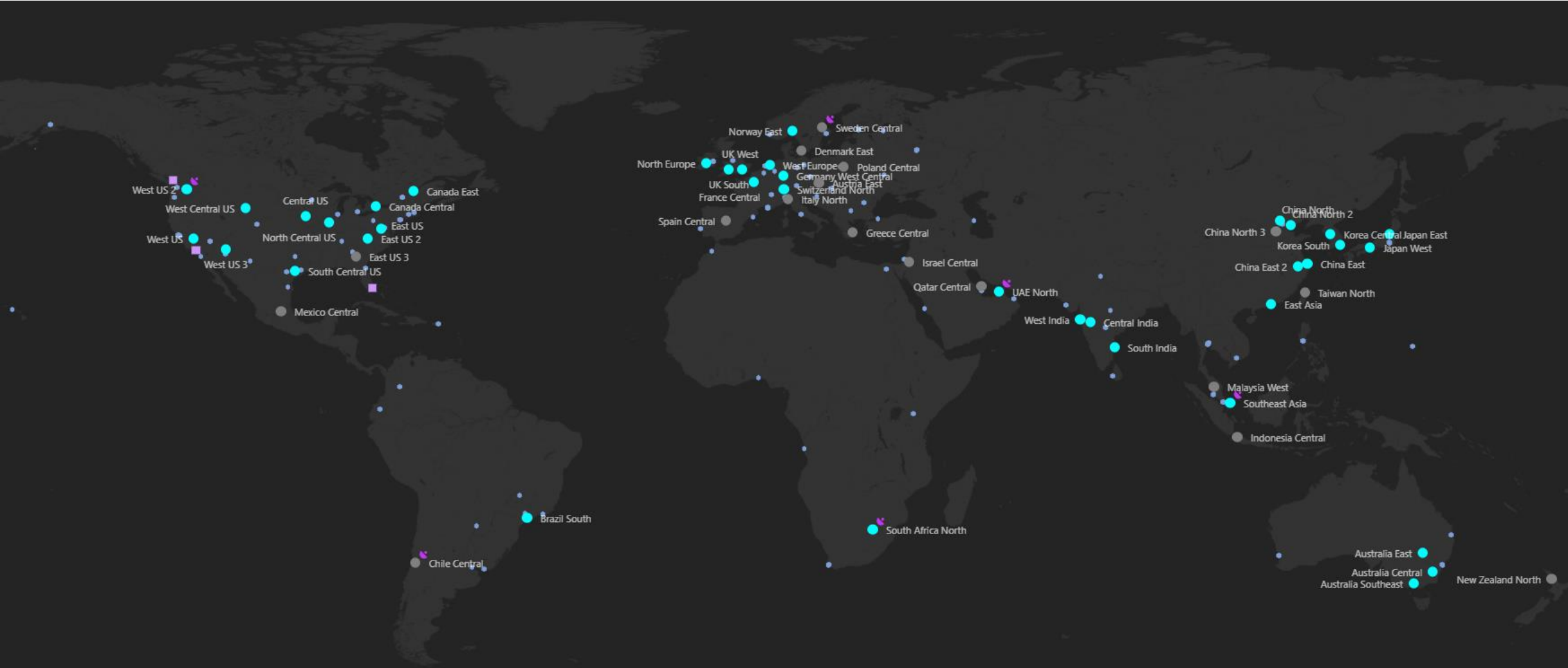
UCX on HBv3

UCX on NDv4



Conclusion

# Global Infrastructure



Compliance and data residency



Service availability



Pricing

# Azure HPC/AI VM Series



## Standard HPC VMs

Standard HPC Applications  
High Compute/Memory + InfiniBand  
HPC SKUs: HB, HC, HBv2, HBv3



## GPU VMs

Deep Learning, AI workloads

Visualization SKUs:  
NV series

Deep Learning/AI SKUs  
NC, ND series

- "r" in VM type indicates RDMA support (InfiniBand)
- InfiniBand/RDMA enabled VMs: One VM per Host
- InfiniBand exposed to VMs using SR-IOV, offers full host bypass with full feature support
- Partition Key (P-key) based isolation

# Agenda



Overview of Azure HPC



**Azure HBv3, NDv4**



Network features



Azure HPC VM Images



Performance Highlights

UCX on HBv3

UCX on NDv4



Conclusion

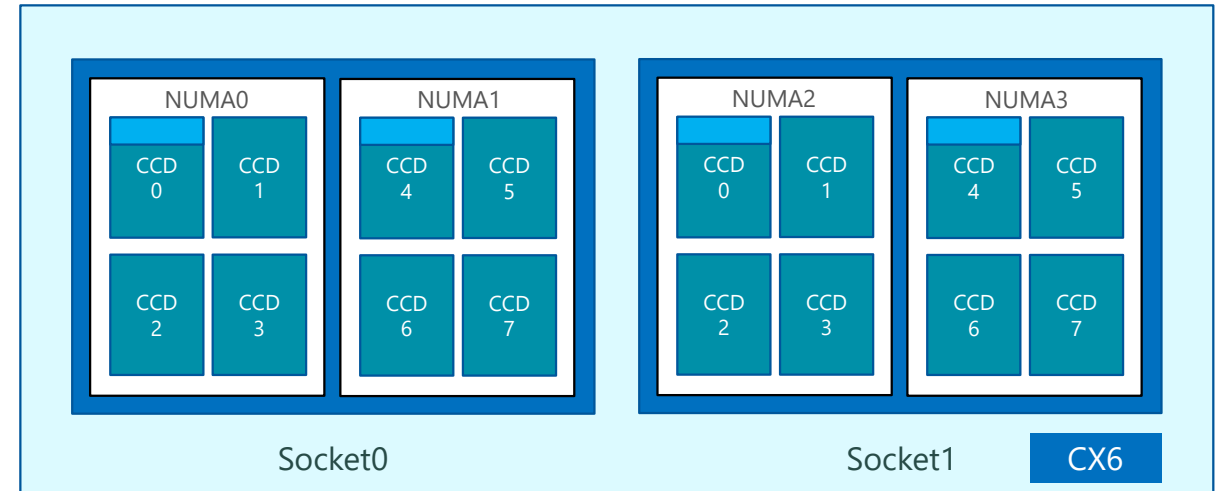
# Azure HBv3 (Milan-X)



AMD EPYC  
Milan-X



NVIDIA®  
InfiniBand HDR  
200Gbps



■ Hyper-V Partition (2 cores per NUMA)

## • VM Specs:

- AMD Milan-X (NPS = 2)
- VM Cores: 120
- L3 Cache: 1.5 GB per VM
- Memory: 448 GB
- Local Disk: 2 x 900 GB NVMe SSD
- Network: 200 Gbps HDR (SR-IOV)

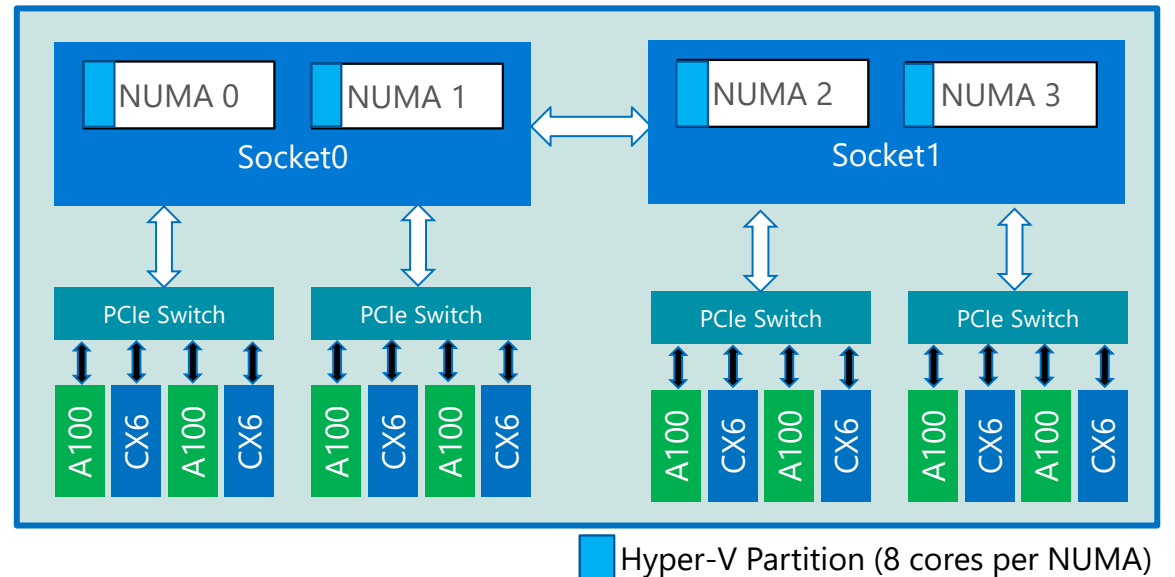
## HBv3 VM Sizes (one VM per Host):

- Standard\_HB120rs\_v3 (all 120 cores)
- Standard\_HB120-96rs\_v3 (6 cores per CCD)
- Standard\_HB120-64rs\_v3 (4 cores per CCD)
- Standard\_HB120-32rs\_v3 (2 cores per CCD)
- Standard\_HB120-16rs\_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
- 8 x HDR 200Gbps InfiniBand
- Local Disk: 6.4 TB local NVMe SSD



Standard\_ND96asr\_v4 (NDv4)

# Agenda



Overview of Azure HPC



Azure HBv3, NDv4



**Network features**



Azure HPC VM Images



Performance Highlights

UCX on HBv3

UCX on NDv4



Conclusion



# InfiniBand Features in Azure

- **HB, HC, NDv2:**



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

- **HBv2, HBv3, NDv4:**



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

- **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

- **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

```
hpcadmin@compute000000:~$ ./test_ib_gpu.sh compute000000 compute000001 gpu /
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 0.002913
8388608 2915 0.00 195.62 0.002915
```

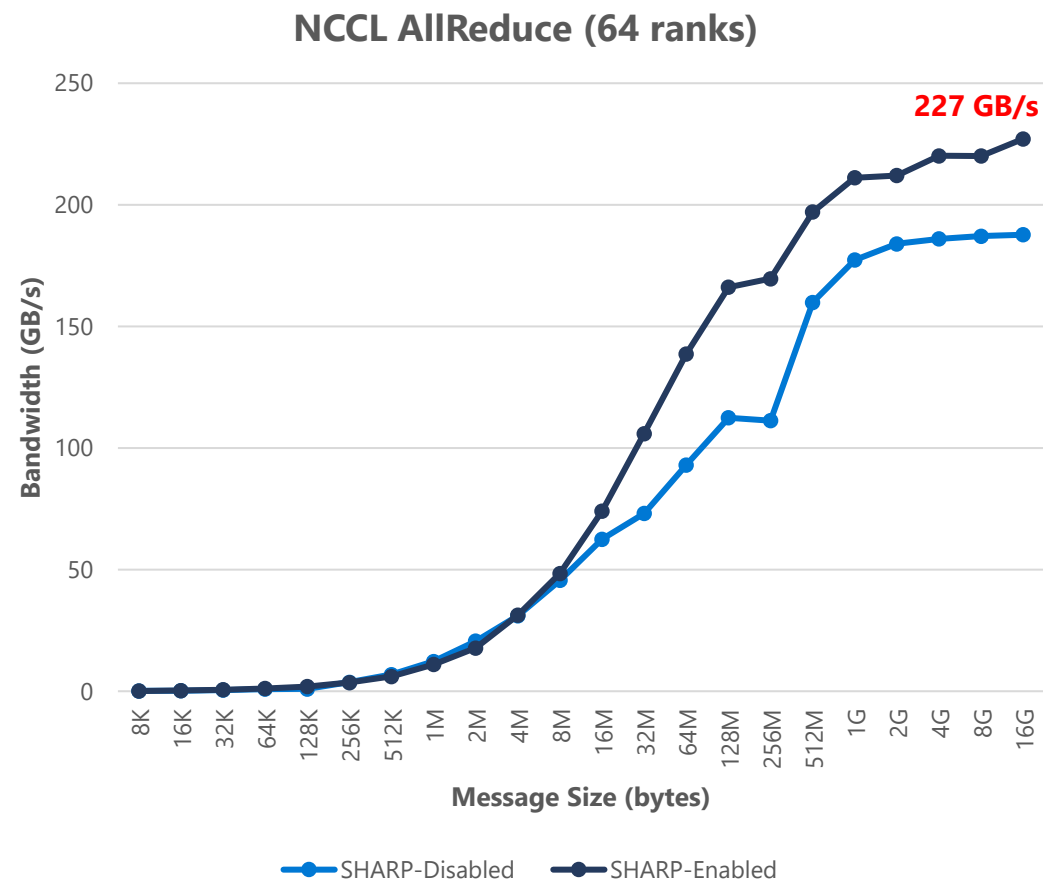
RDMA (Host Memory)

```
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

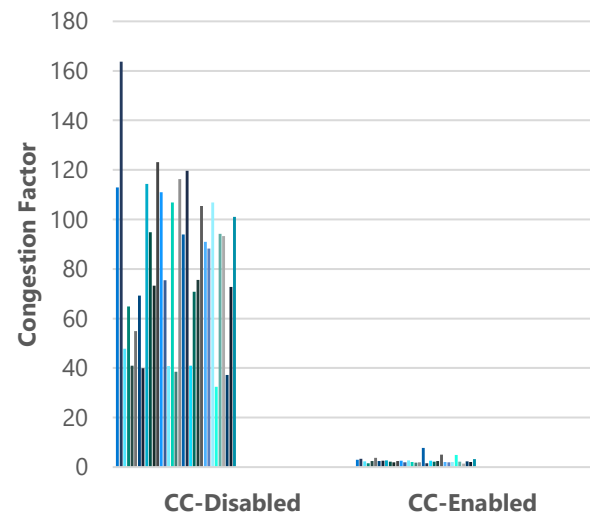
- Enabled on dedicated NDv4 clusters
- UCX-based Sharp-AM / SharpD communication
- Optimized SHARP tree initialization
- Connection keepalive
- GRH support



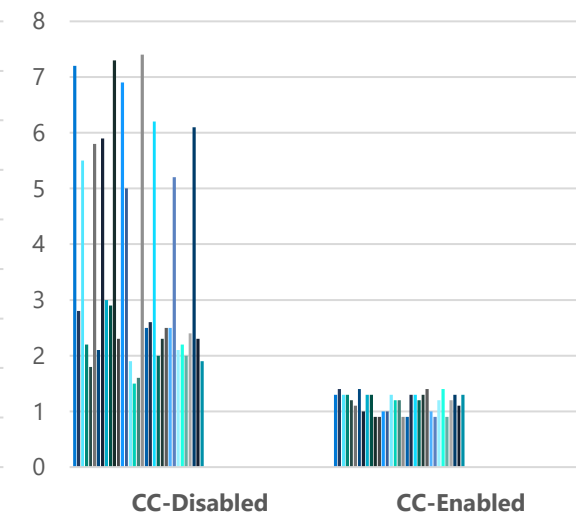
# Adv. Congestion Control

- Available on all VM Series with HDR
- Transparent to customer applications
- Avoids congestion, Improve tail latencies
- Critical in public multi-customer environments

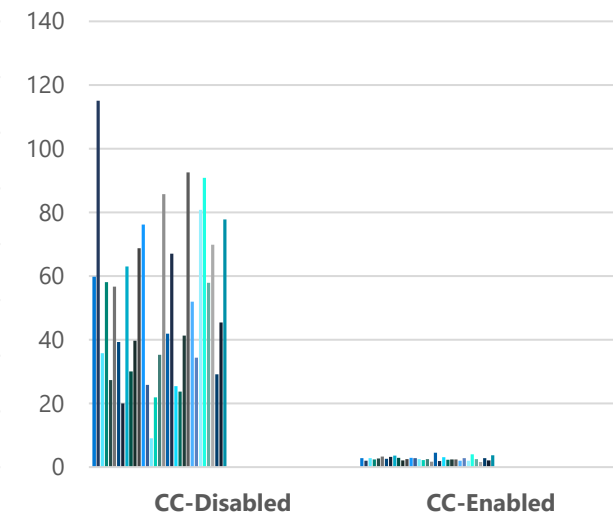
GPCNet: RR Two-sided Latency (Avg)



GPCNet: RR Two-sided BW (Avg)



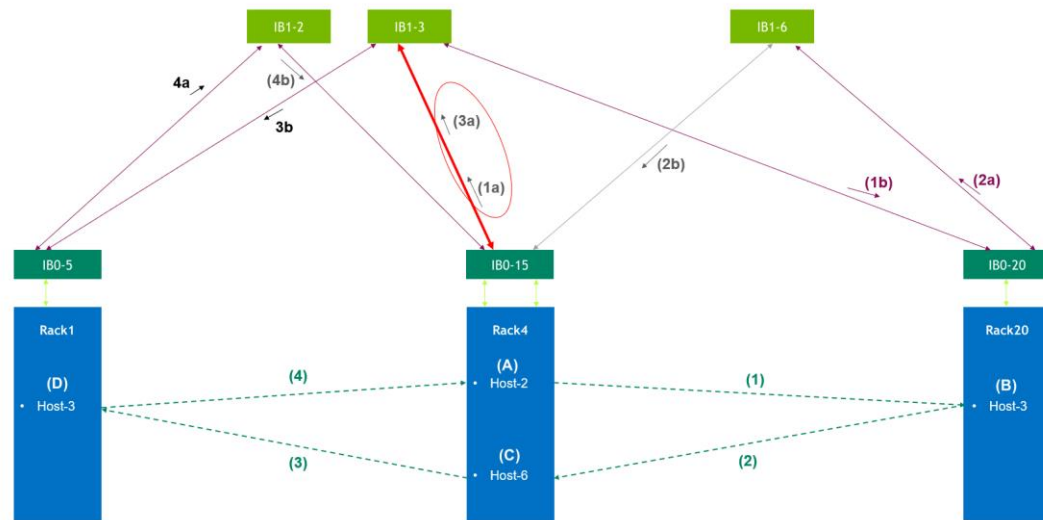
GPCNet: RR 2-sided AllReduce (Avg)



More results on Thursday's (12/02) session (10.30 am PST):

**"Cloud-Native Supercomputing Performance isolation"**

# Adaptive Routing

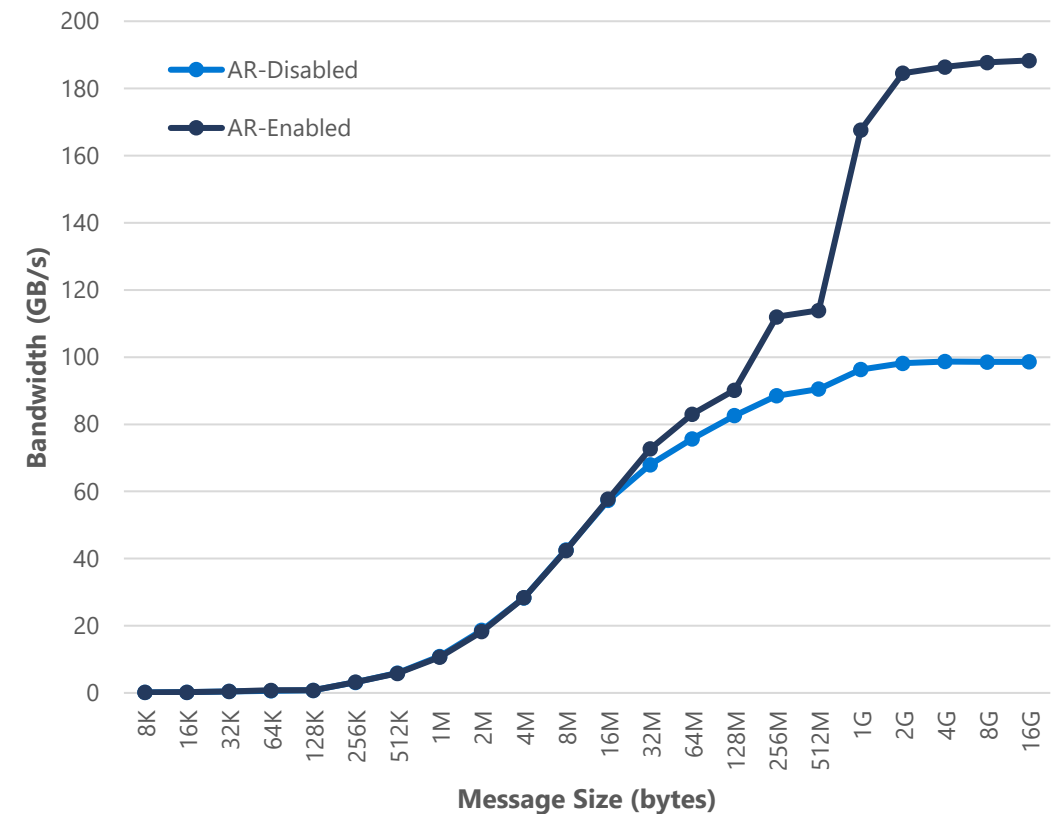


Communication paths during NCCL AllReduce

## Impact of Adaptive Routing

- Congestion can happen with static routing if a single link is being used by two or more communicating pairs
- AR avoids congestion and offers stable performance
- More details: [Adaptive Routing on Azure HPC Clusters](#)

## NCCL AllReduce Bandwidth



# Agenda



Overview of Azure HPC



Azure HBv3, NDv4



Network features



**Azure HPC VM Images**



Performance Highlights

UCX on HBv3

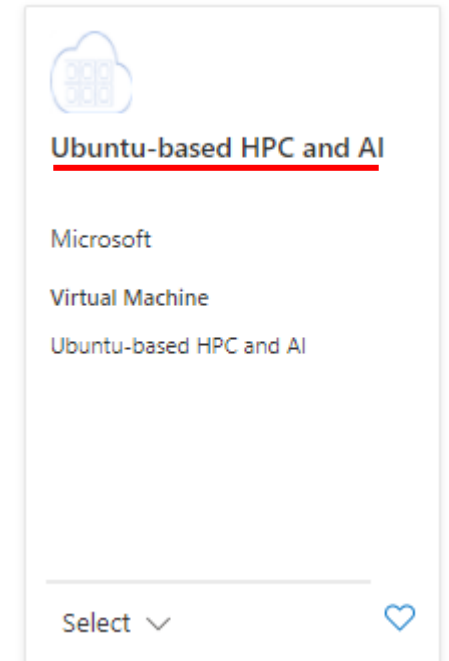
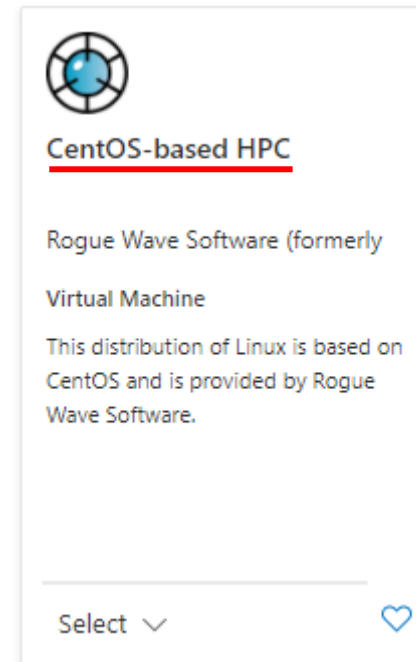
UCX on NDv4



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, **UCX**
- Optimized HPC libraries
  - Blis, FFTW, Flame, MKL
- Recommended Compilers
- GPU Drivers
- NCCL, NCCL RDMA Sharp Plugin, SharpD
- Other platform optimizations



<https://github.com/Azure/azhpc-images>

# Agenda



Overview of Azure HPC



Azure HBv3, NDv4



Network features



Azure HPC VM Images



**Performance Highlights**

UCX on HBv3

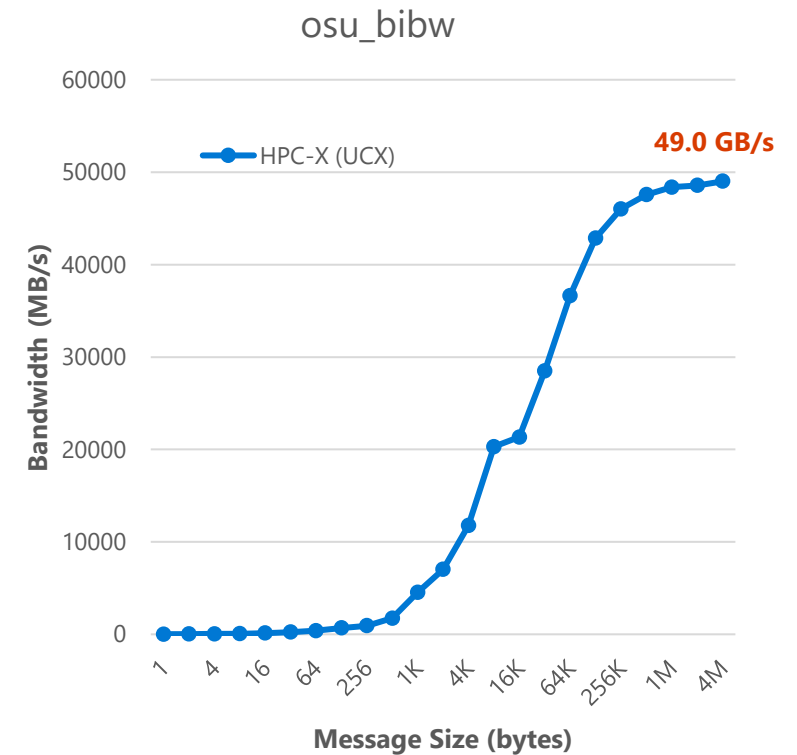
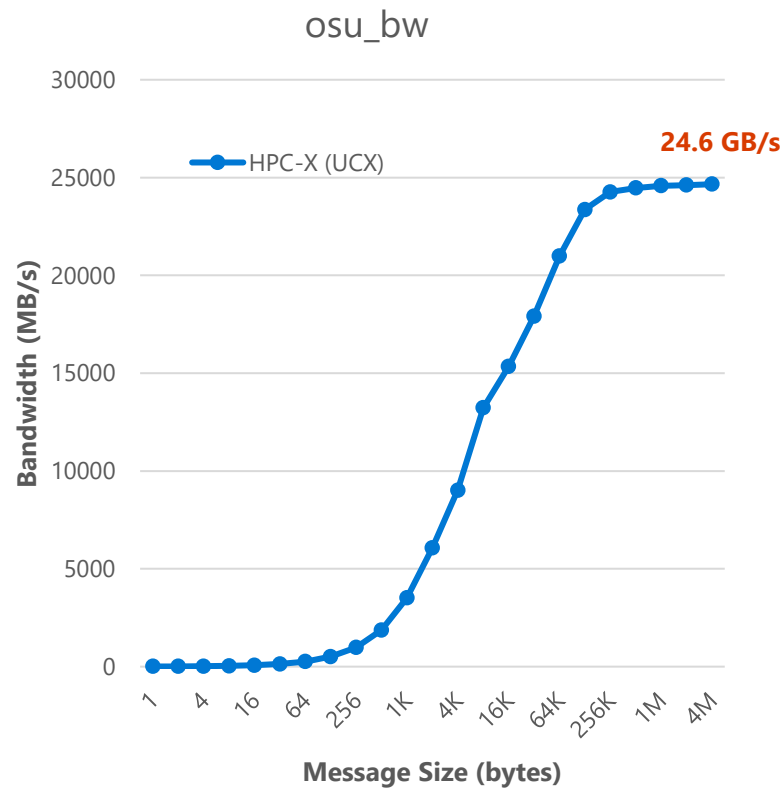
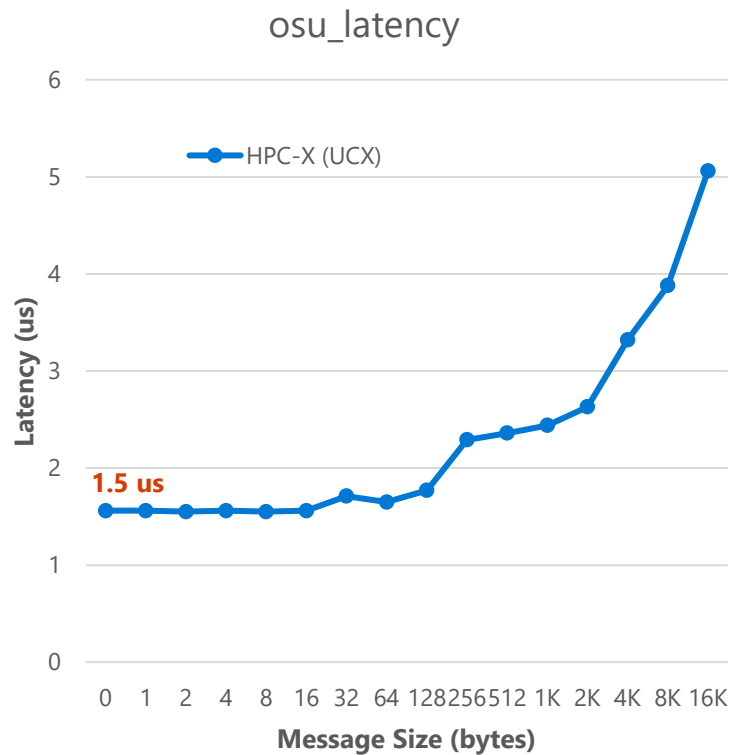
UCX on NDv4



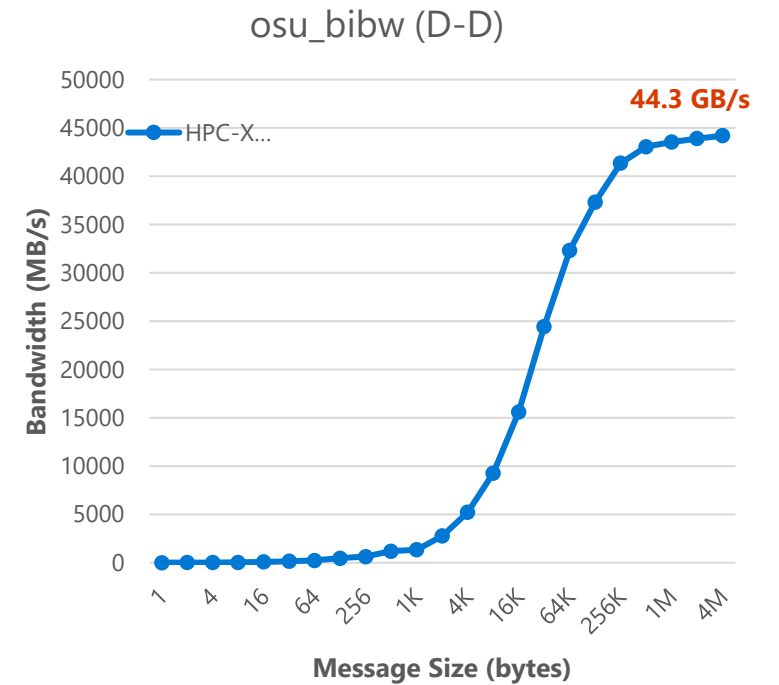
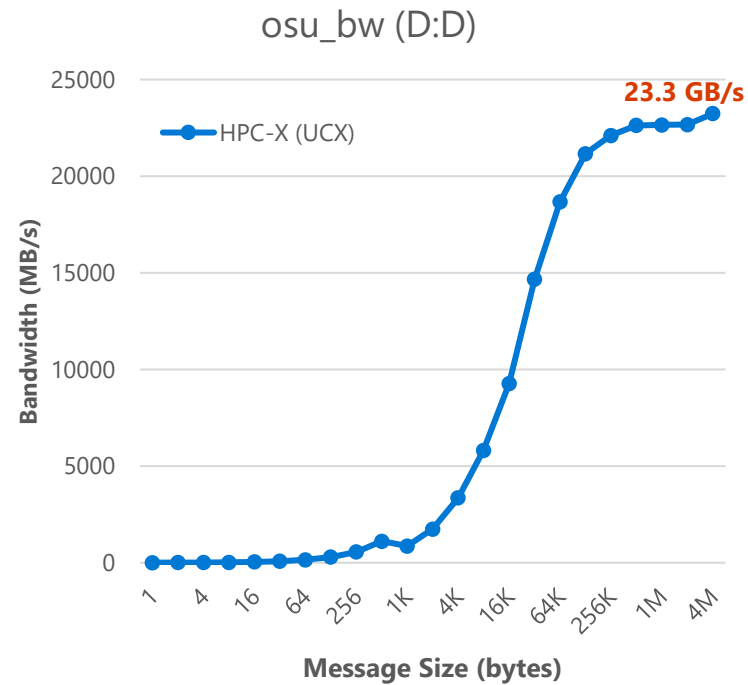
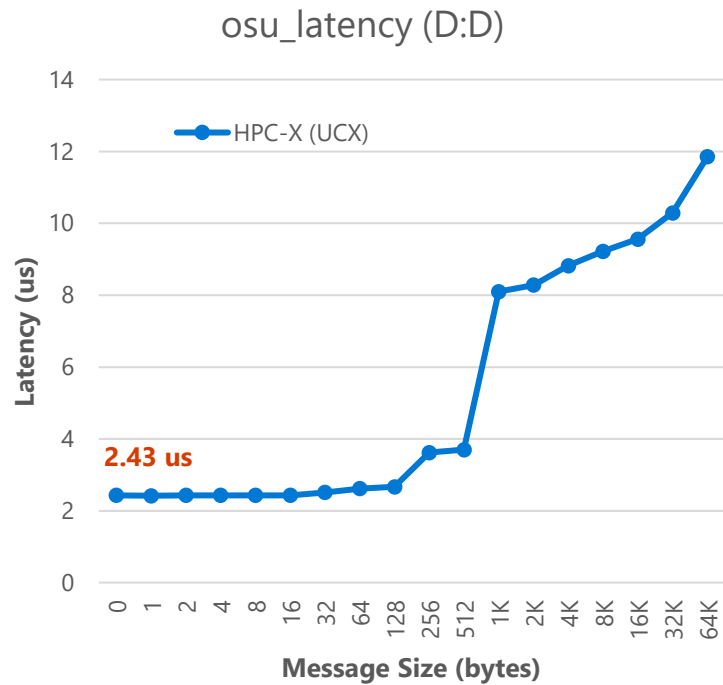
Conclusion



# HBv3 MPI (UCX) Performance Characteristics

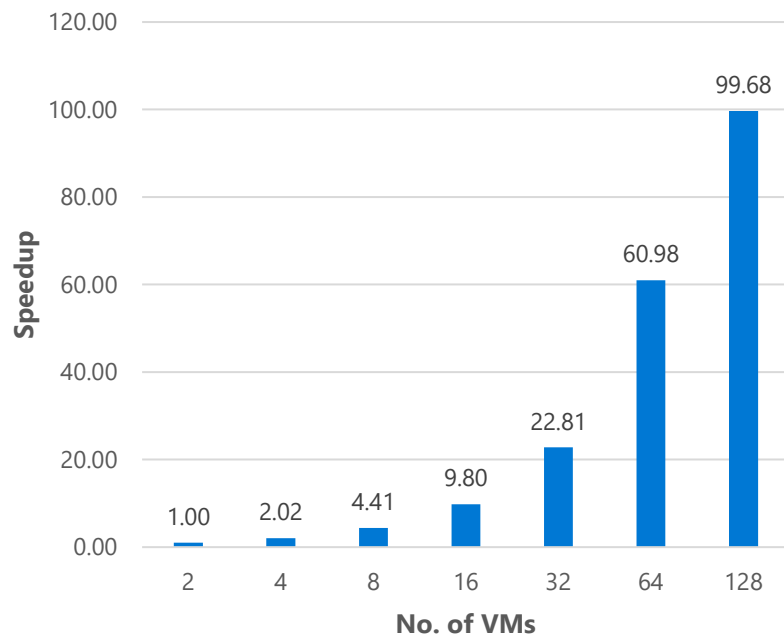


# NDv4 MPI (UCX) Performance Characteristics (D:D)



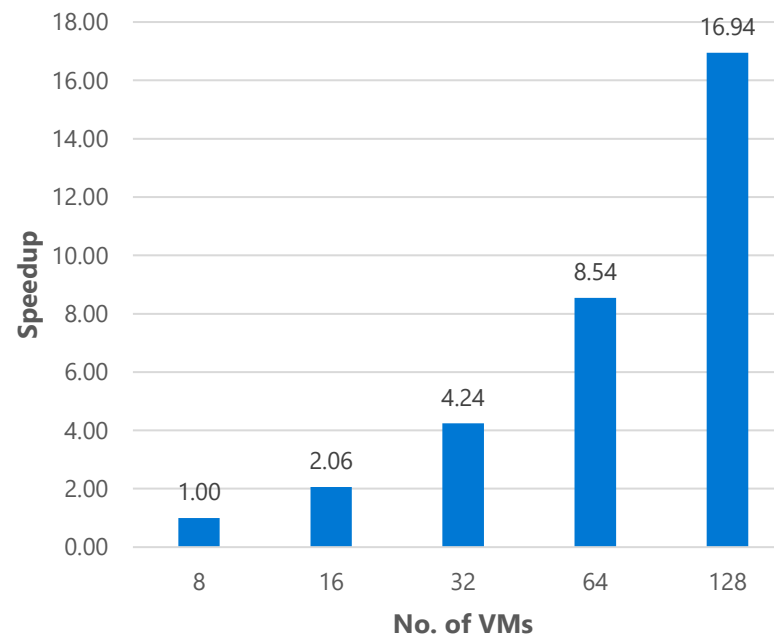
# Scaling Efficiency on HBv3 (Milan-X) using UCX

Ansys Fluent 2021 R1  
f1\_racecar\_140m



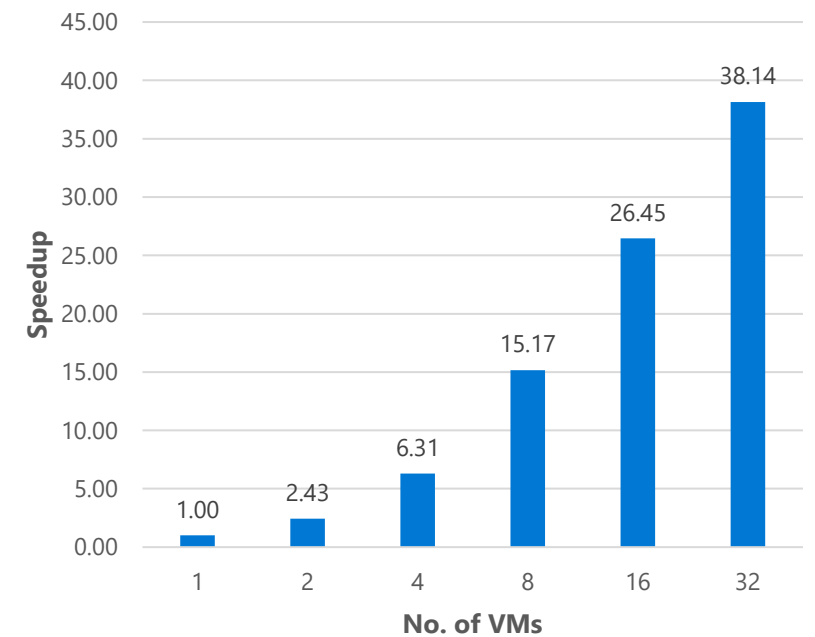
156% scaling efficiency

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



Azure HBv3, NDv4



Network features



Azure HPC VM Images



Performance Highlights

UCX on HBv3

UCX on NDv4



**Conclusion**

# Conclusion

- Supercomputer on Cloud is real!
- Azure HPC Cloud in Top500, Graph500 top rankings
  - Rank 2 overall in MLPerf Dec. 2021
  - Rank 10 in Top500 Nov. 2021
  - Rank 17 in Graph500 Nov. 2020
- Cloud democratize Supercomputer
- High Performance middleware such as UCX enables cutting edge technology
  - Deliver High Scalability and Performance

# Resources

## 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](#)

## Azure HPC Blogs

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



Thank you