



## The InfiniBand Advantage

Joshua S. Ladd, PhD

Staff Engineer, HPC Software R&D Group

OpenSHMEM 2015, Annapolis, Maryland



# Exascale-Class Computer Platforms – Communication Challenges



Challenge	Solution focus
Very large functional unit count ~10M	Scalable communication capabilities: point-to-point & collectives Scalable Network: Adaptive routing
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 capabilities
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 progress
Power costs	Power aware hardware



The Future is Here



# Enter the World of Scalable Performance

# At the Speed of 100Gb/s!

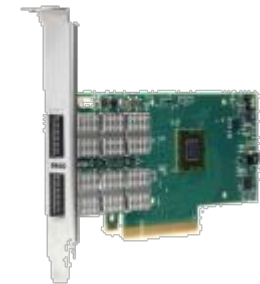


## Entering the Era of 100Gb/s

Adapters

ConnectX<sup>®</sup> 4

100Gb/s Adapter, 0.7us latency  
150 million messages per second  
(10 / 25 / 40 / 50 / 56 / 100Gb/s)



Switch

SwitchIB<sup>™</sup>

36 EDR (100Gb/s) Ports, <90ns Latency  
Throughput of 7.2Tb/s



Interconnect

LinkX<sup>™</sup>



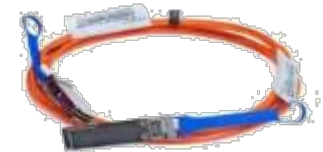
Copper (Passive, Active)



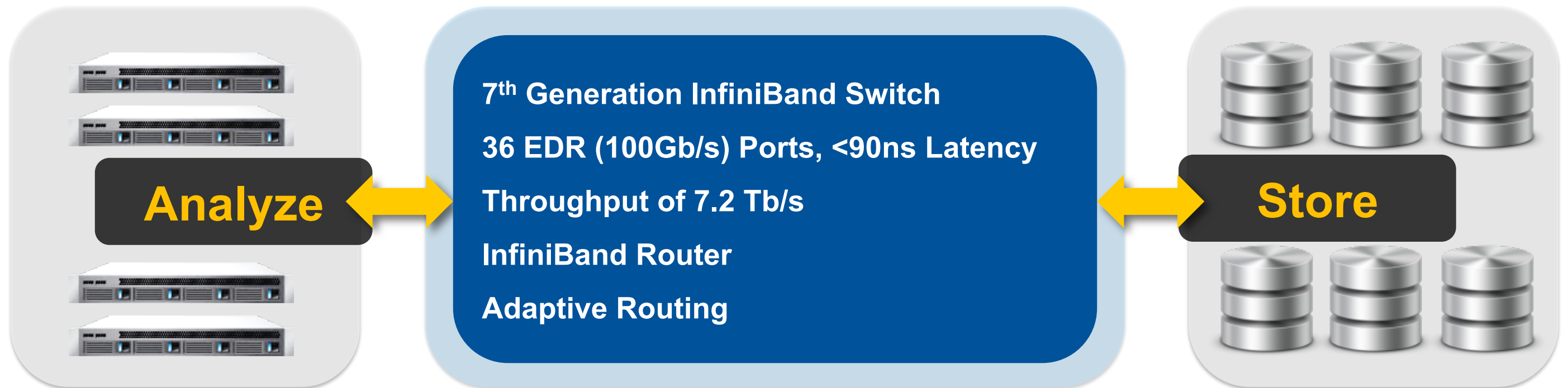
Optical Cables (VCSEL)



Silicon Photonics



## Switch-IB: Highest Performance Switch in the Market



SwitchIB™

## ConnectX-4: Highest Performance Adapter in the Market

InfiniBand: SDR / DDR / QDR / FDR / EDR

Ethernet: 10 / 25 / 40 / 50 / 56 / 100GbE

100Gb/s, <0.7us latency

150 million messages per second

OpenPOWER CAPI technology

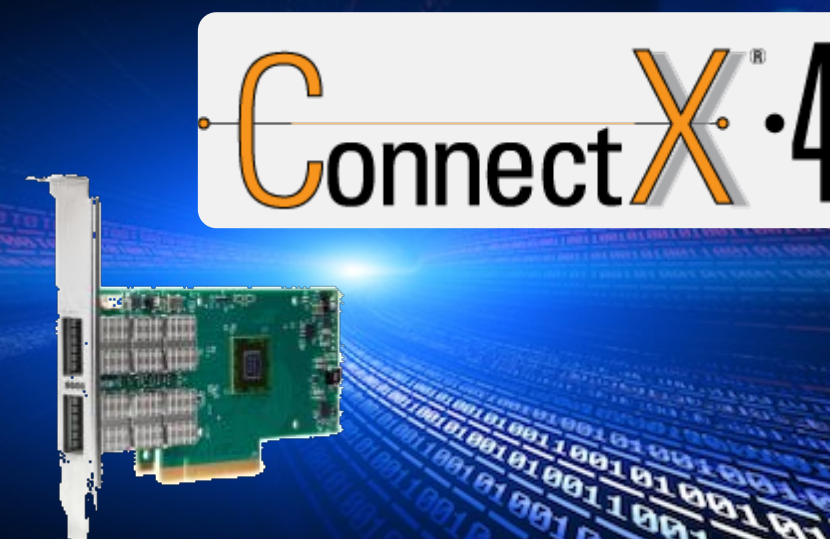
CORE-Direct technology

GPUDirect RDMA

Dynamically Connected Transport (DCT)

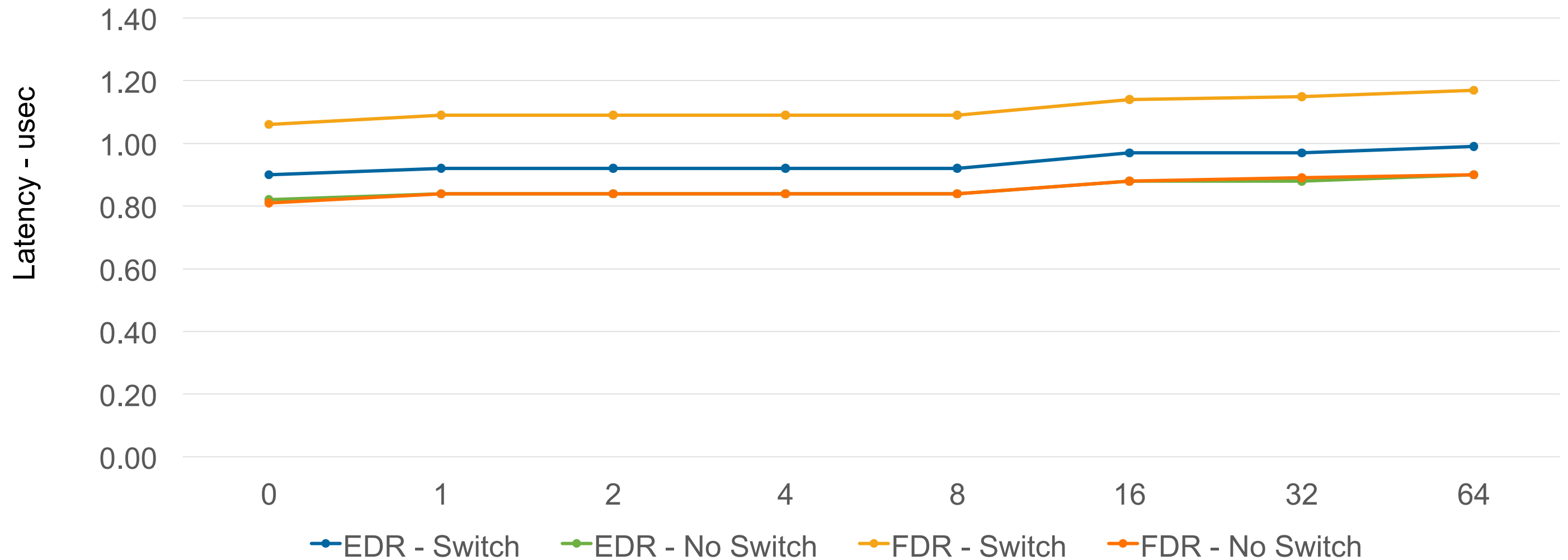
Ethernet offloads (HDS, RSS, TSS, LRO, LSOv2)

Connect. Accelerate. Outperform



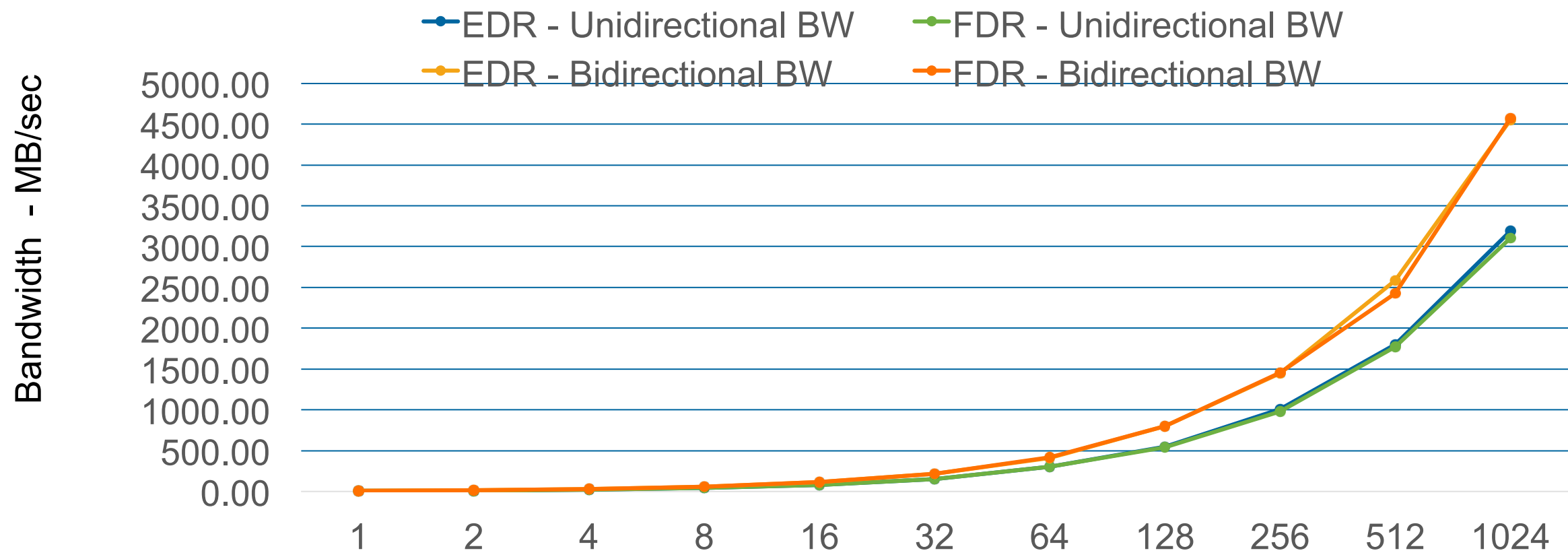
# Point-to-Point Data

# MPI Latency – OSU Latency test

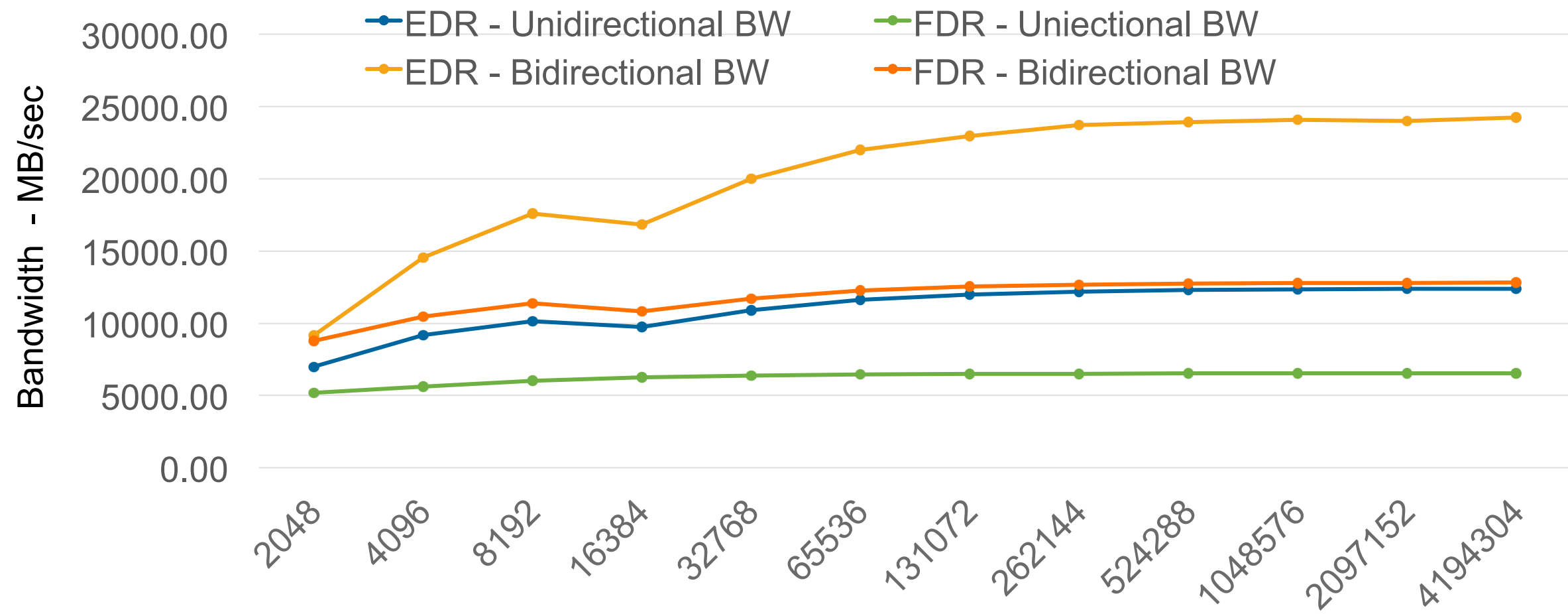




# MPI Bandwidth – OSU Bandwidth test

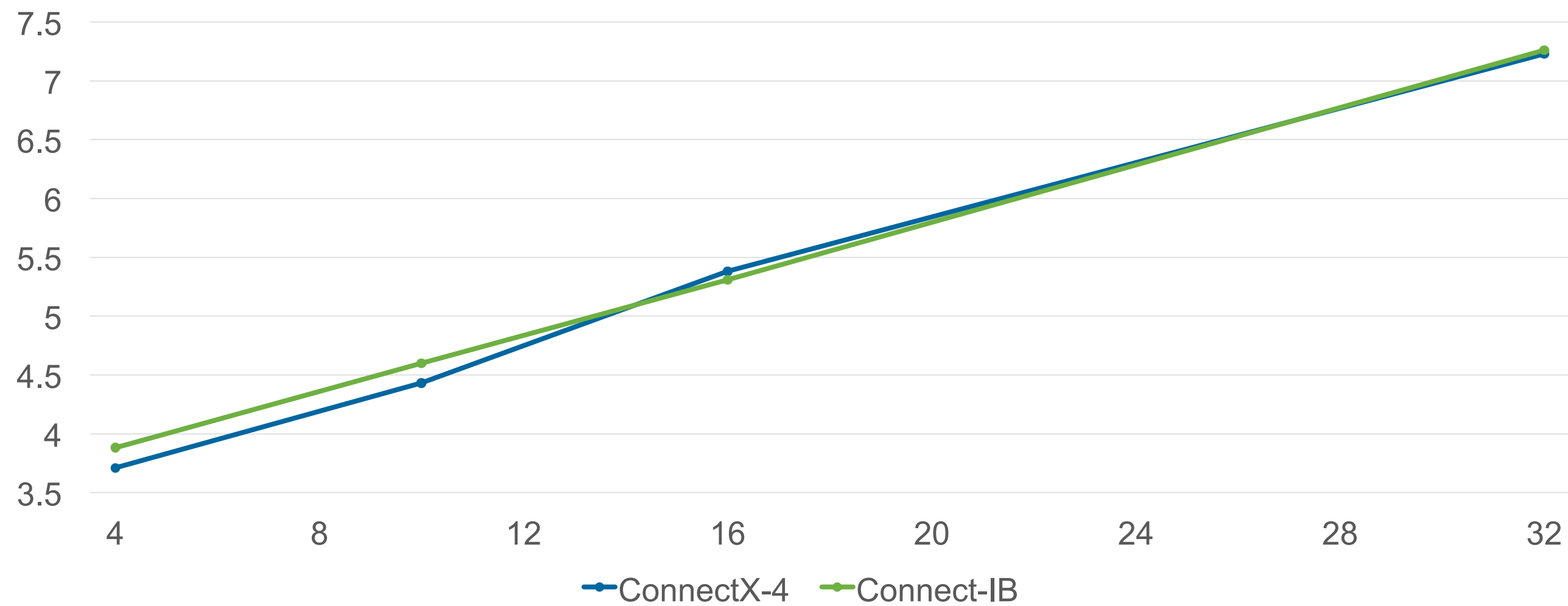


# MPI Bandwidth – OSU Bandwidth test

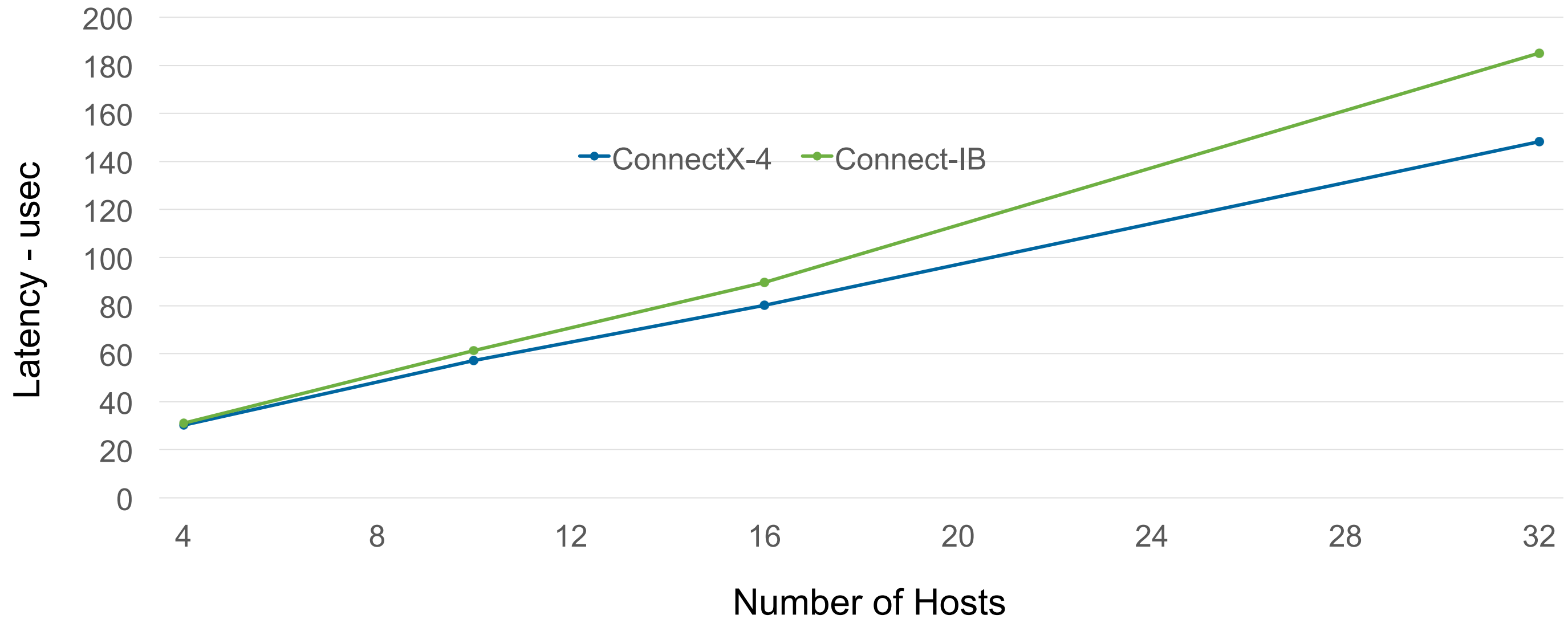


# Collective Communication

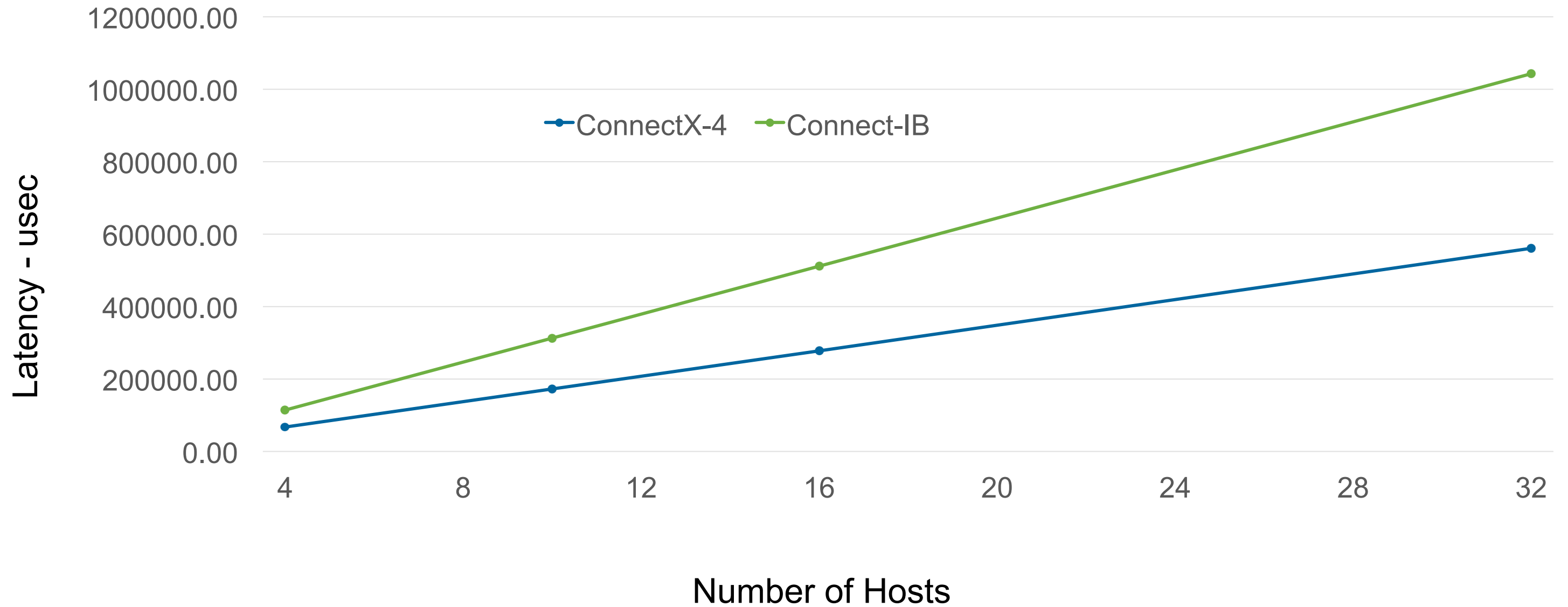




# All-to-All - 8 bytes

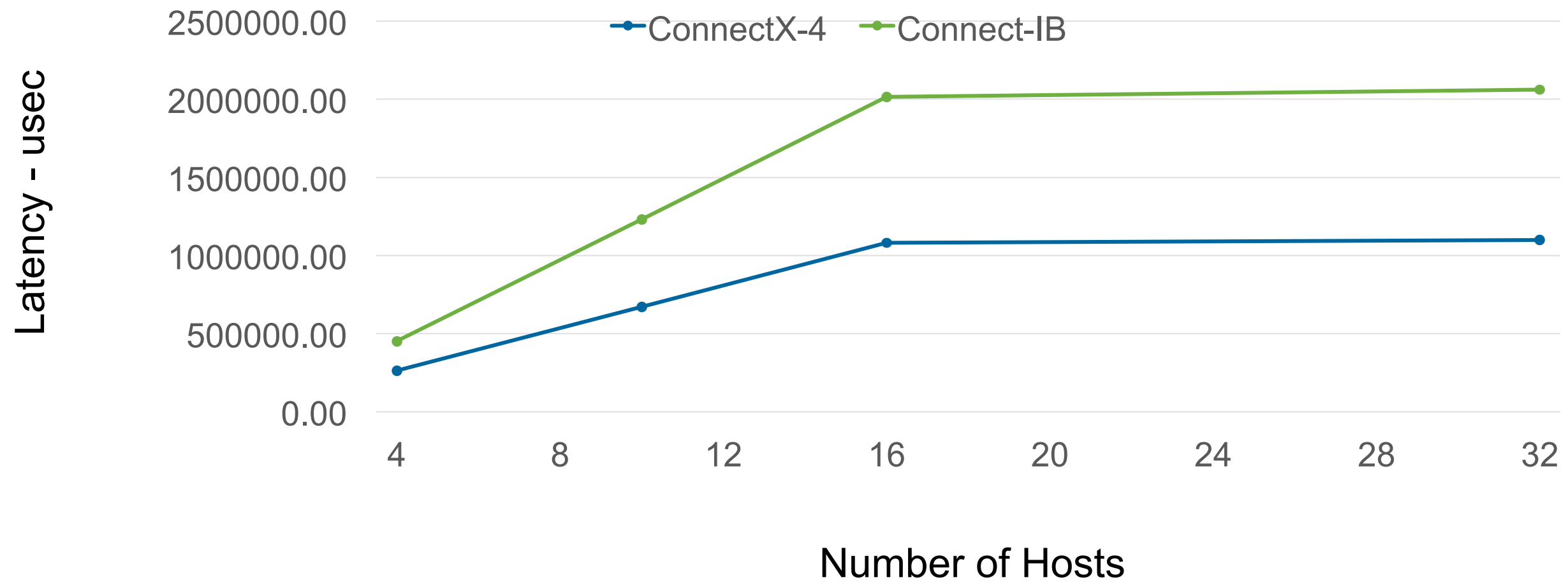


# All-to-All – 256 Kbytes

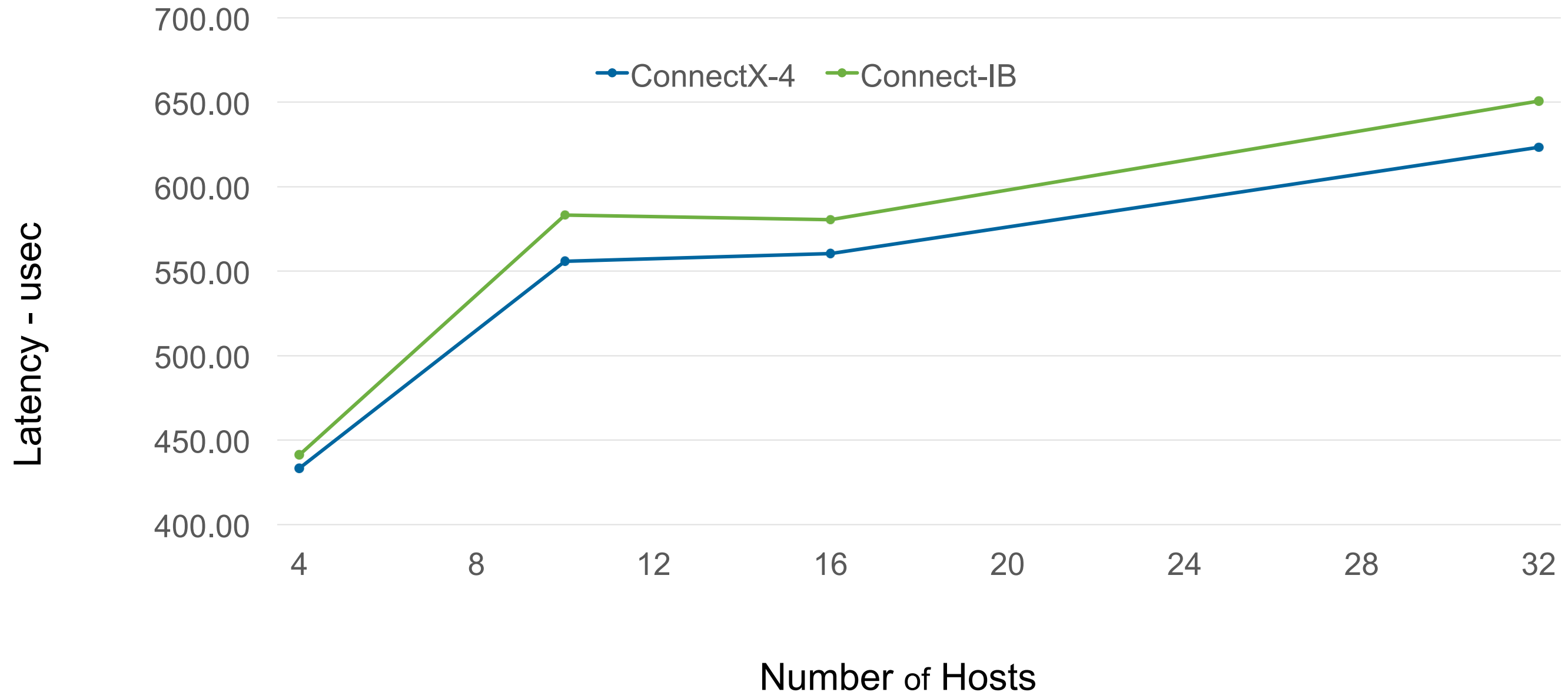




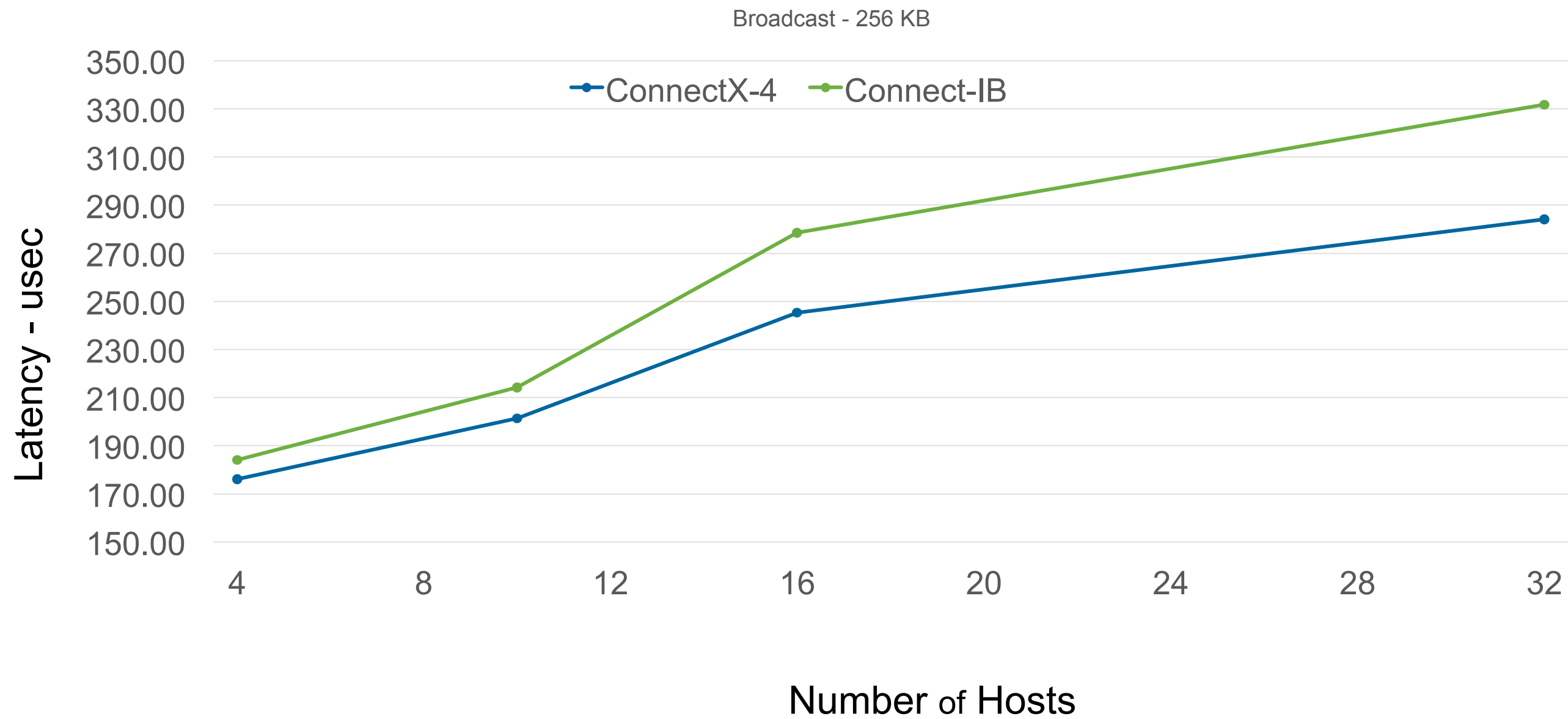
# All-to-All – 1 Mbyte



# Allreduce – 256K Bytes

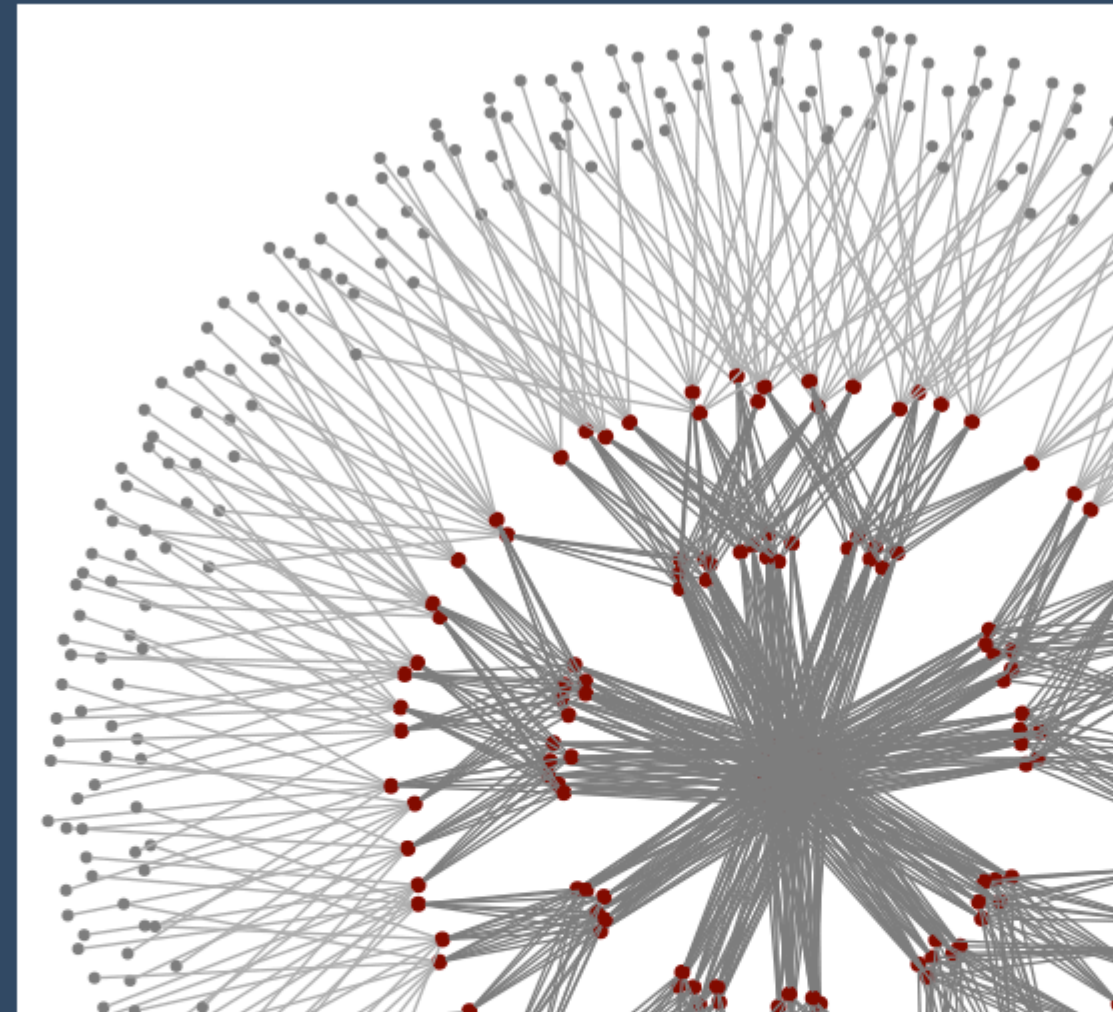


# Broadcast – 256K Bytes



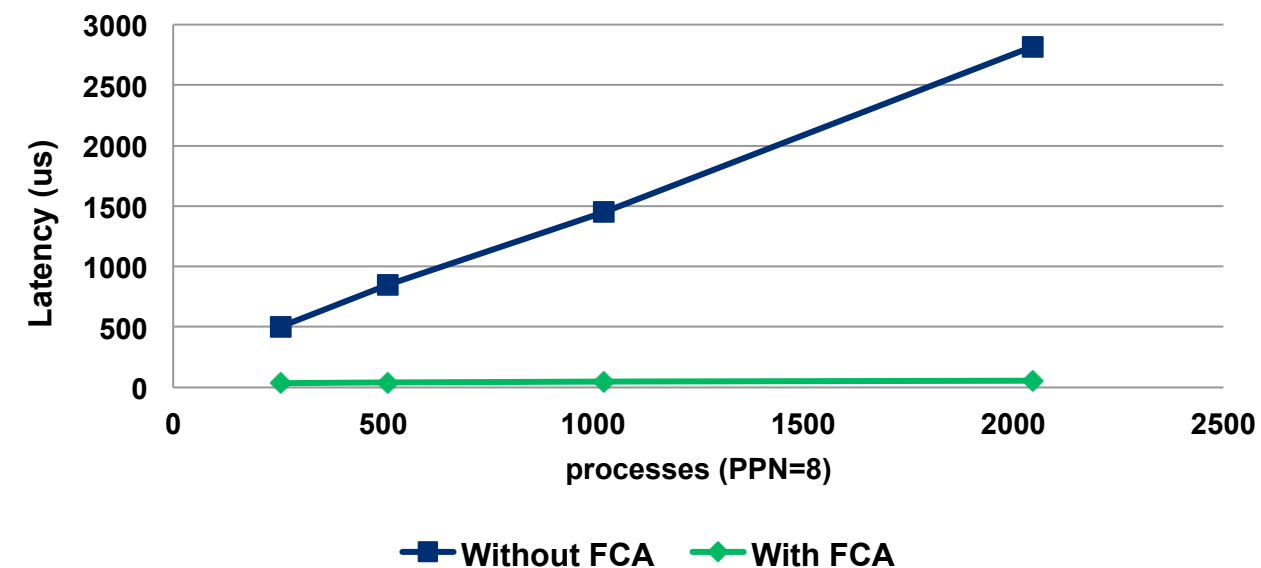


# Scalability

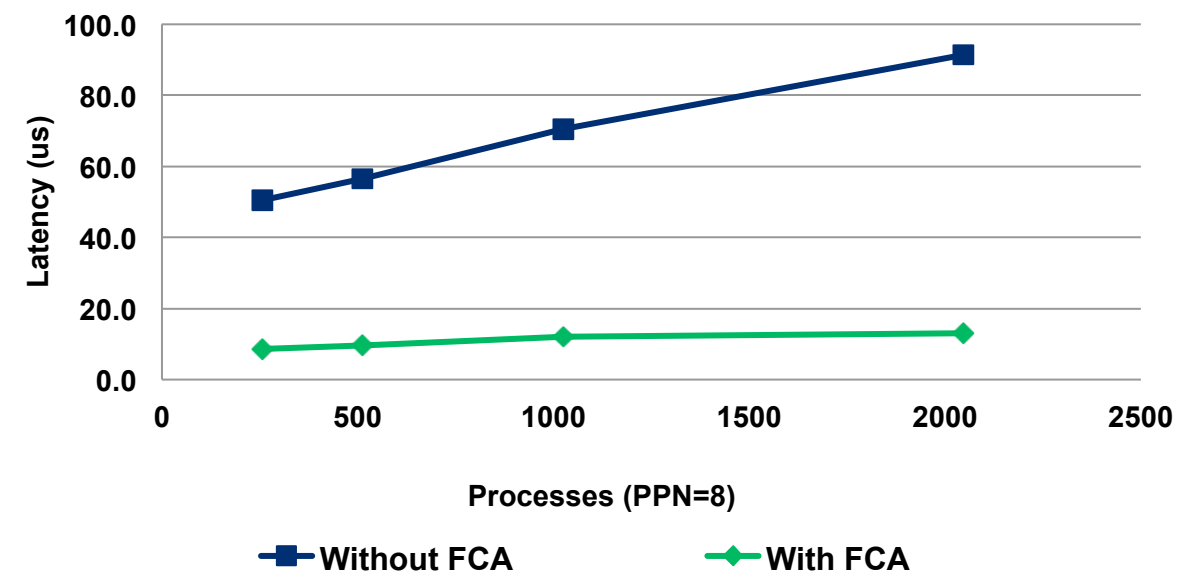


- Topology Aware
- Hardware Multicast
- Offload
- Scalable algorithms

## Reduce Collective

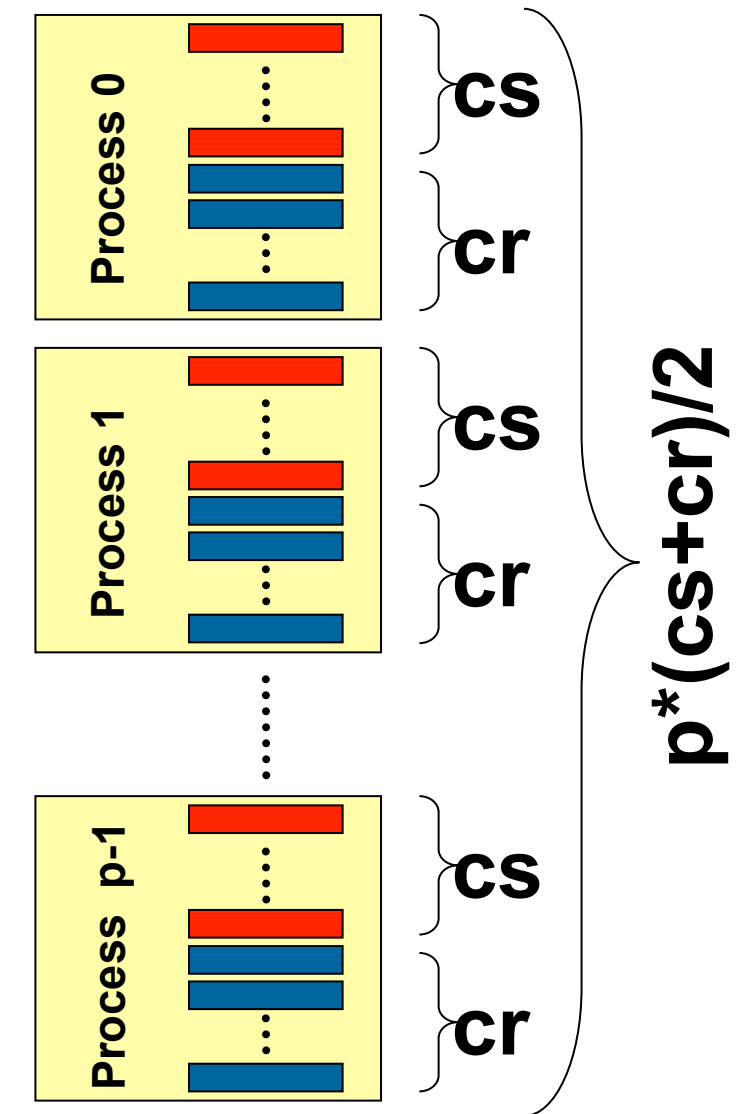


## Barrier Collective



# The Dynamically Connected Transport Model

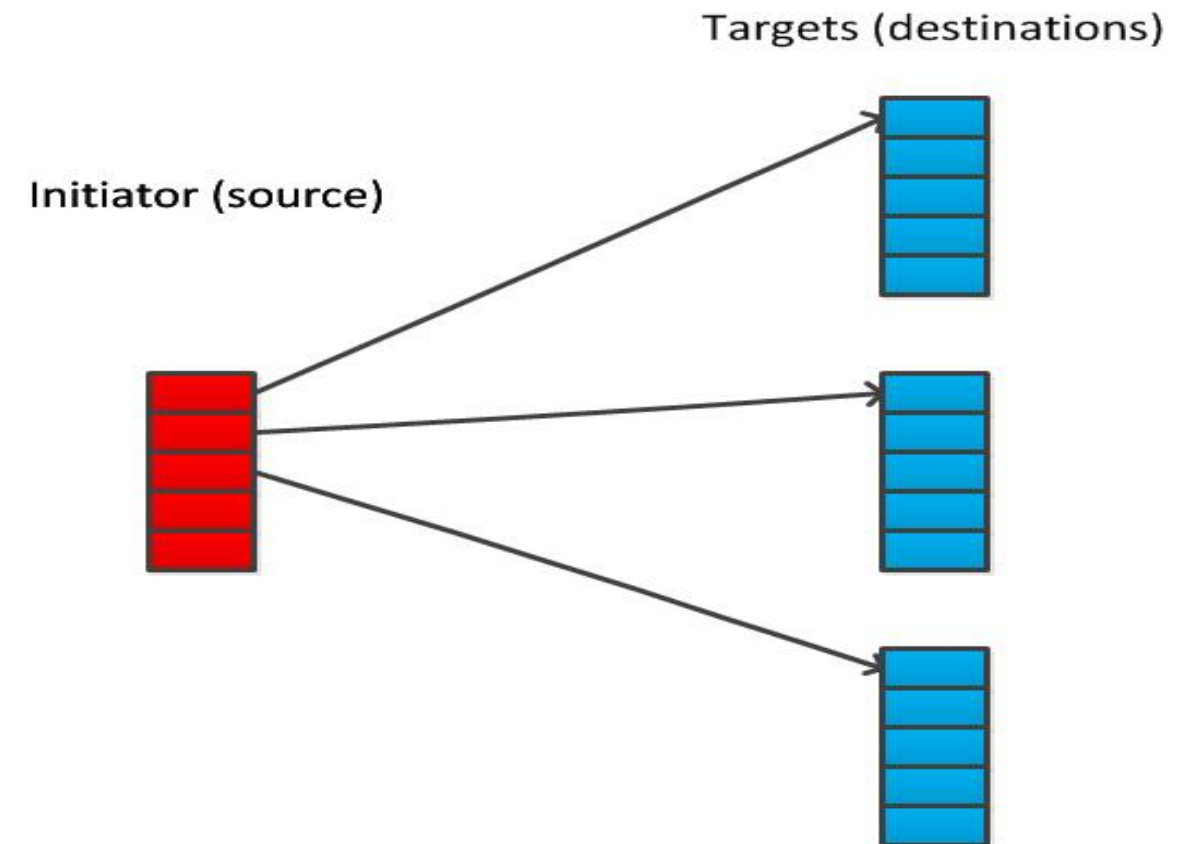
- Dynamic Connectivity
- Each DC Initiator can be used to reach any remote DC Target
- No resources' sharing between processes
  - process controls how many (and can adapt to load)
  - process controls usage model (e.g. SQ allocation policy)
  - no inter-process dependencies
- Resource footprint
  - Function of node and HCA capability
  - Independent of system size
- Fast Communication Setup Time



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

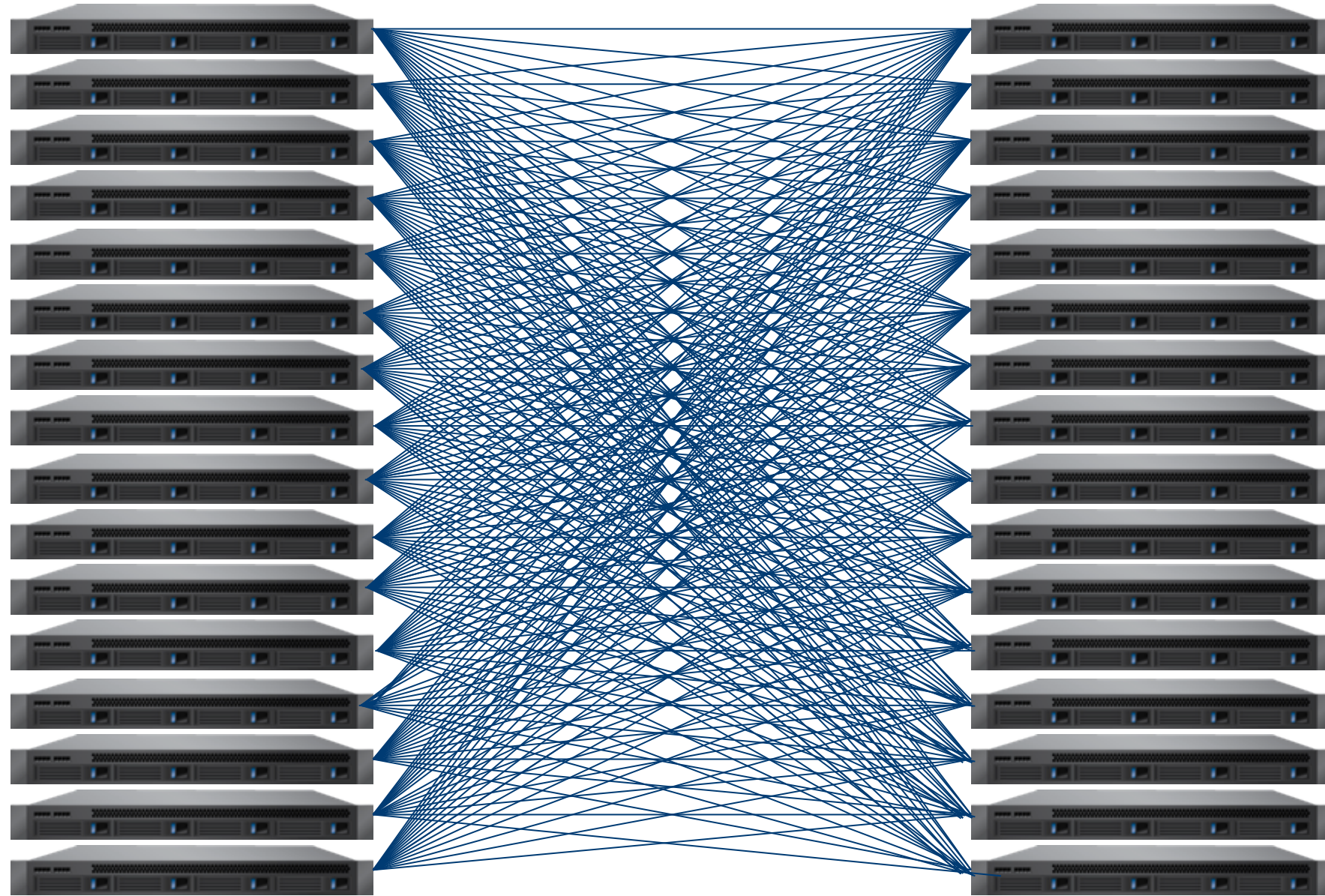
## ■ Key objects

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



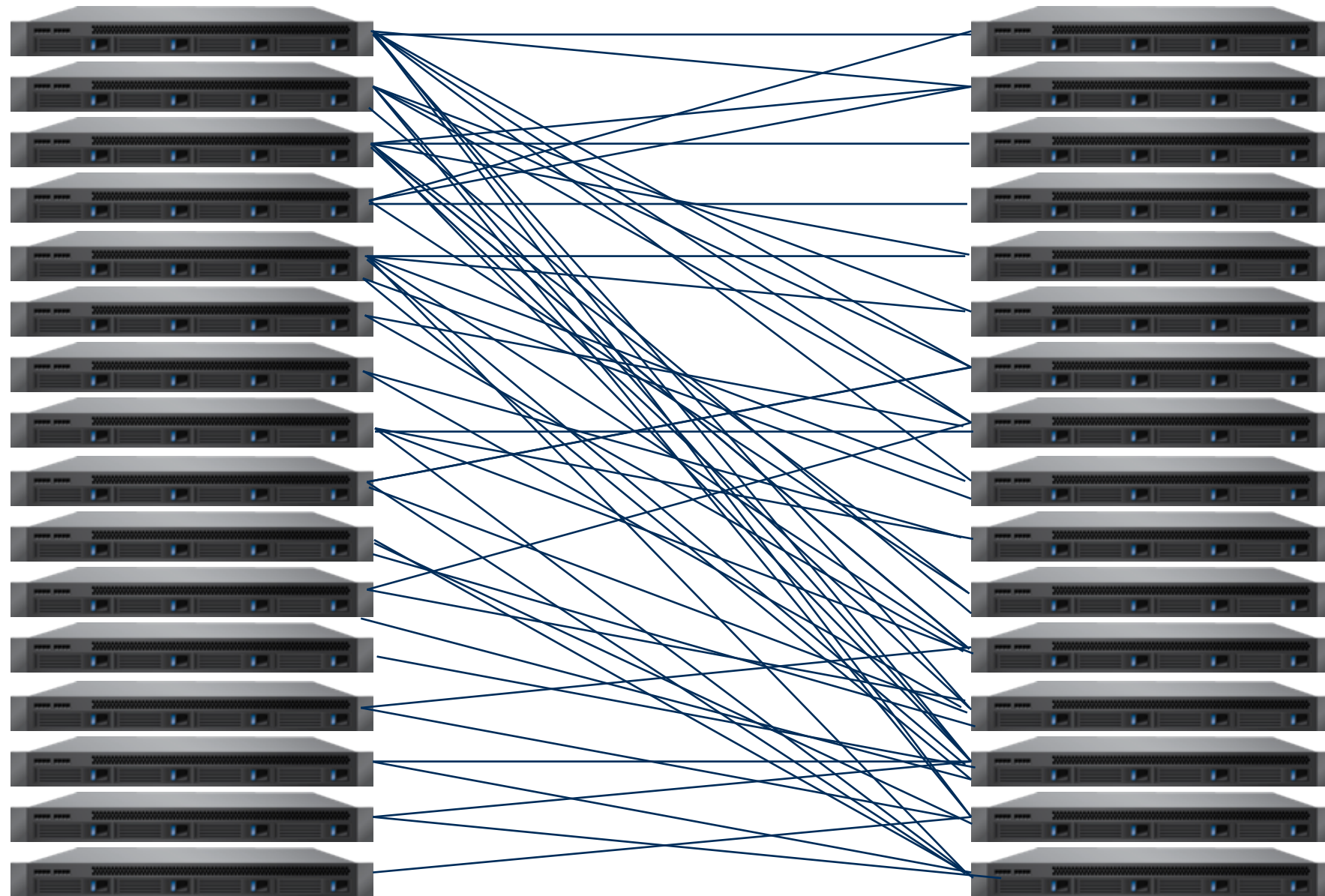


# Reliable Connection Transport Mode

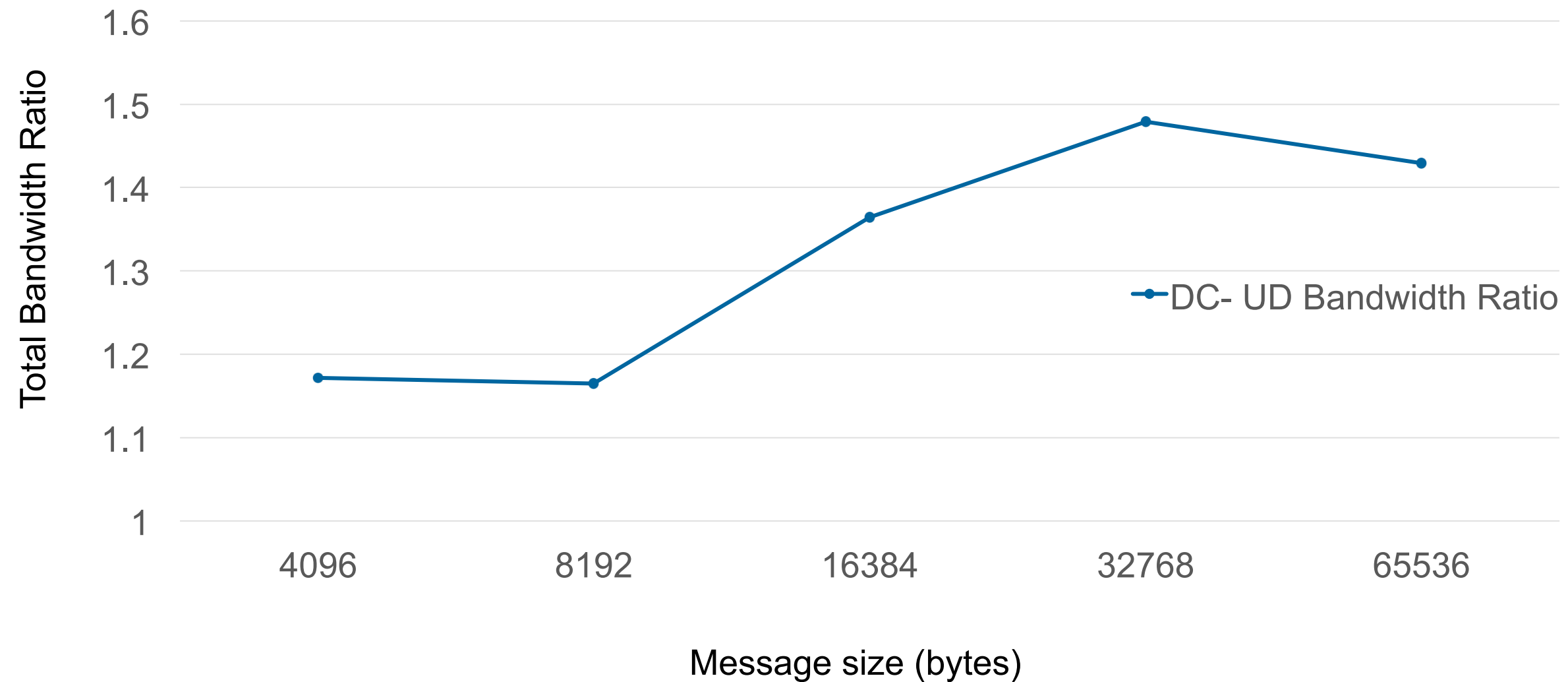




# Dynamically Connected Transport Mode



# All-To-All Performance





# Scalability Under Load



0 4 7 1 7 4 3 7	0 4 7 8 4 1 0 0	1 8 1 4 7 5 5 0	8 7 2 9 4 1 0 0	0 4 7 1 7
1 8 2 3 4 1 0 4	1 8 2 4 6 6 2 1	6 4 0 9 9 4 1 3	7 8 7 8 6 6 2 1	1 8 2 3 4
9 9 9 5 9 9 3 9	9 9 9 1 9 8 0 6	8 1 3 8 8 6 6 2	4 5 4 4 9 8 0 6	9 9 9 5 9
4 4 3 4 8 8 5 8	4 4 3 9 9 0 4 5	4 5 2 6 7 5 5 0	5 4 8 5 9 0 4 5	4 4 3 4 8
9 9 1 0 0 7 4 4	9 9 1 8 8 4 1 1	9 1 0 5 9 4 0 6	8 0 5 4 8 4 1 1	9 9 1 0 0
8 8 4 6 1 9 6 5	8 8 4 4 8 9 9 0	8 3 5 1 6 1 4 4	2 8 7 6 8 9 9 0	8 8 4 6 1
6 0 8 2 3 0 0 2	3 5 3 0 9 4 2 0	6 0 8 2 3 0 0 2	0 2 1 1 2 3 0 0	9 1 8 6 3
1 5 5 3 2 6 1 0	6 4 5 3 8 9 0 1	1 5 5 3 2 6 1 0	1 1 5 1 4 2 0 6	7 2 9 5 6
5 4 9 0 1 5 4 1	5 9 1 6 7 8 1 9	5 4 9 0 1 5 4 1	6 6 6 6 5 1 6 5	5 2 7 4 5
4 4 4 1 8 1 9 4	4 8 3 5 4 7 4 8	4 4 4 1 8 1 9 4	5 5 5 5 1 6 1 1	9 7 9 6 4
4 9 9 4 7 6 6 4	9 0 2 4 9 8 4 7	4 9 9 4 7 6 6 4	1 1 1 3 4 5 4 6	8 8 5 4 9
9 8 8 9 4 4 5 9	6 6 1 9 8 4 9 4	9 8 8 9 4 4 5 9	4 4 2 2 6 6 6 5	5 5 9 6 6
8 4 4 8 9 5 1 8	8 5 6 8 0 0 8 9	8 4 4 8 9 5 1 8	6 5 0 1 5 5 5 1	0 9 6 5 8
7 6 5 7 8 1 6 4	4 4 5 4 4 2 4 8	7 6 5 7 8 1 6 4	5 1 5 5 1 1 1 3	1 7 2 4 4
4 5 4 8 4 3 5 1	6 9 7 6 8 1 1 0	4 5 4 8 4 3 5 1	1 0 1 6 3 3 4 2	6 8 1 6 6
6 1 5 9 6 0 1 6	5 5 4 5 7 6 6 4	6 1 5 9 6 0 1 6	6 2 6 5 5 2 2 7	5 2 6 8 5
2 6 3 0 5 5 4 5	4 1 9 3 9 5 5 6	2 6 3 0 5 5 4 5	8 0 5 1 4 5 5 5	0 1 4 4 4
1 8 4 1 0 6 6 1	9 3 8 1 9 0 1 4	1 8 4 1 0 6 6 1	5 1 4 5 2 1 1 9	5 6 2 6 9

## ■ Purpose

- Improved Network utilization: choose alternate routes on congestion
- Network resilience: Alternative routes on failure

## ■ Supported Hardware

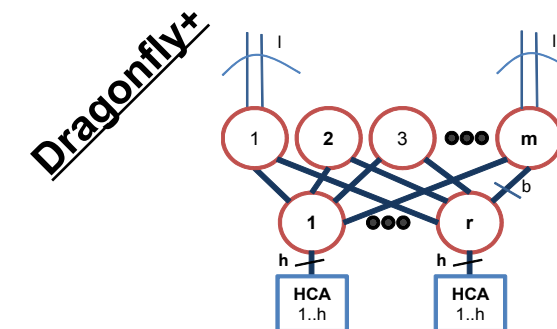
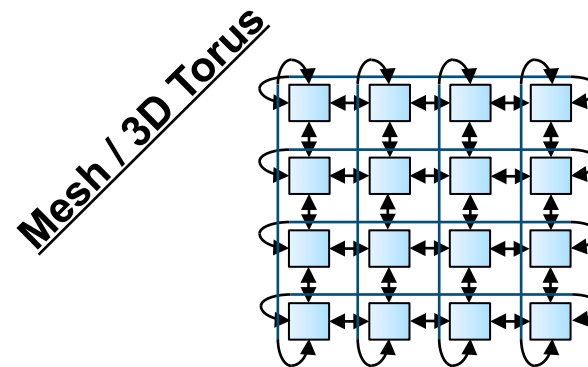
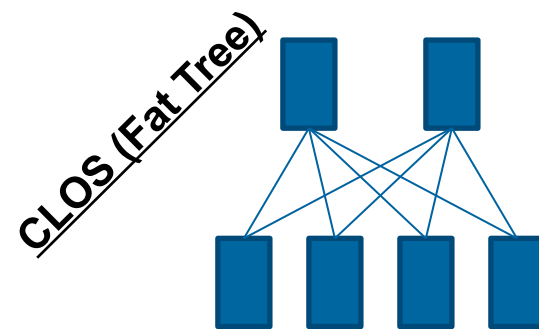
- SwitchX-2
- Switch-IB: Adaptive routing notification added



- Mellanox hardware is NOT topology specific
  - SDN concept separates the configuration plane from the data plane
  - Every feature is software controlled
  - Fat-Tree, Dragonfly and Dragonfly+ are fully supported
  - New hardware features introduced to support Dragonfly and Dragonfly+



## Topologies





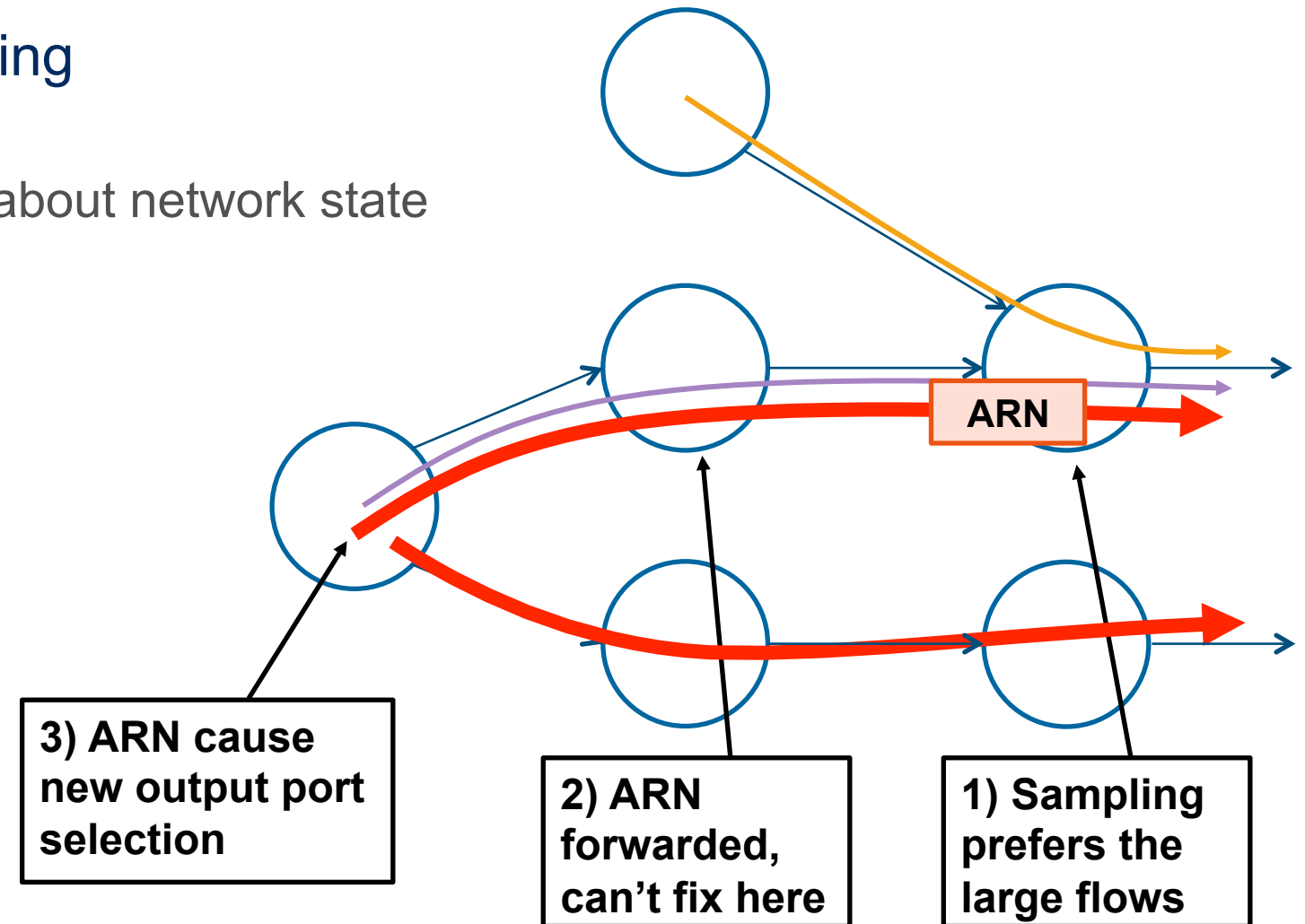
# Is the Packet Allowed to Adapt?



- For every incoming packet the adaptive routing process has two main stages
  - Route Change Decision (to adapt or not to adapt)
  - New output port selection
- AR Modes
  - Static – traffic is always bound to a specific port
  - Time-Bound – traffic is bound to the last port used if not more than  $T_b$  [sec] passed since that event
  - Free – traffic may select a new out port freely
- Packets are classified to be either Legacy, Restricted or Unrestricted
- Destinations are classified to be either Legacy, Restricted, Timely-Restricted or Unrestricted
- A matrix maps possible combinations of packet and destination based classification to AR modes

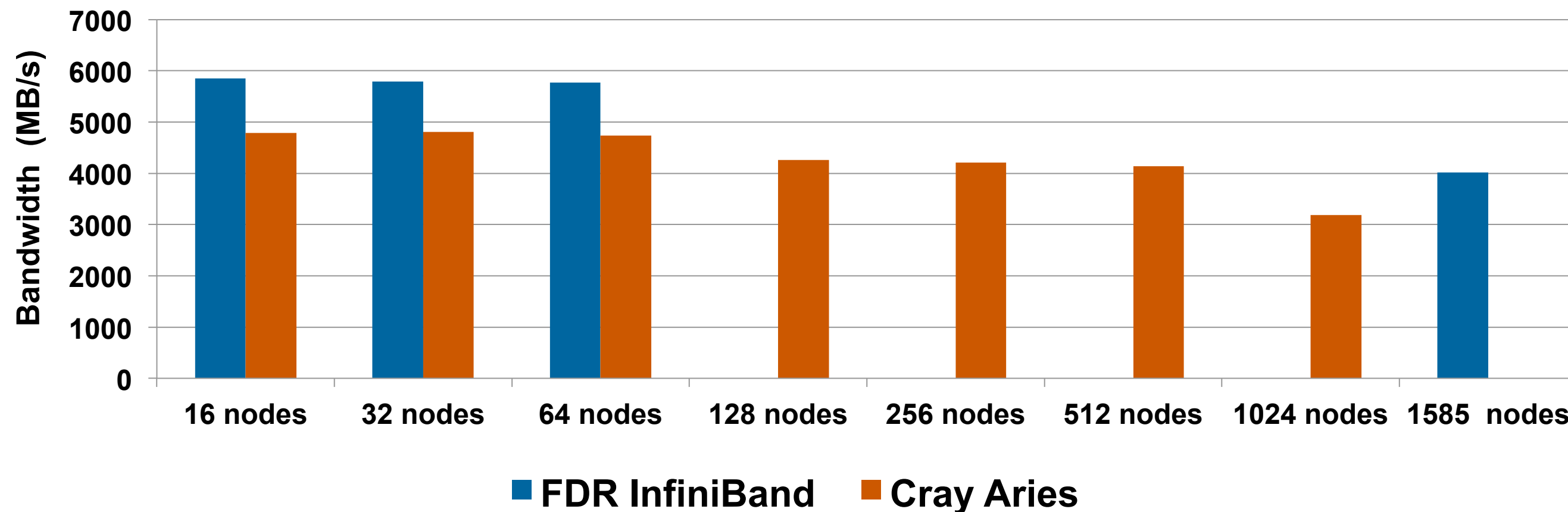
# Mellanox Adaptive Routing Notification (ARN)

- The “reaction” time is critical to Adaptive Routing
  - Traffic modes change fast
  - A “better” AR decision requires some knowledge about network state
- Internal switch to switch communications
- Faster convergence after routing changes
- Fast notification to decision point
- Fully configurable (topology agnostic)



**Faster Routing Modifications, Resilient Network**

## B\_Eff Benchmark



**Higher Performance and Better Network Utilization**



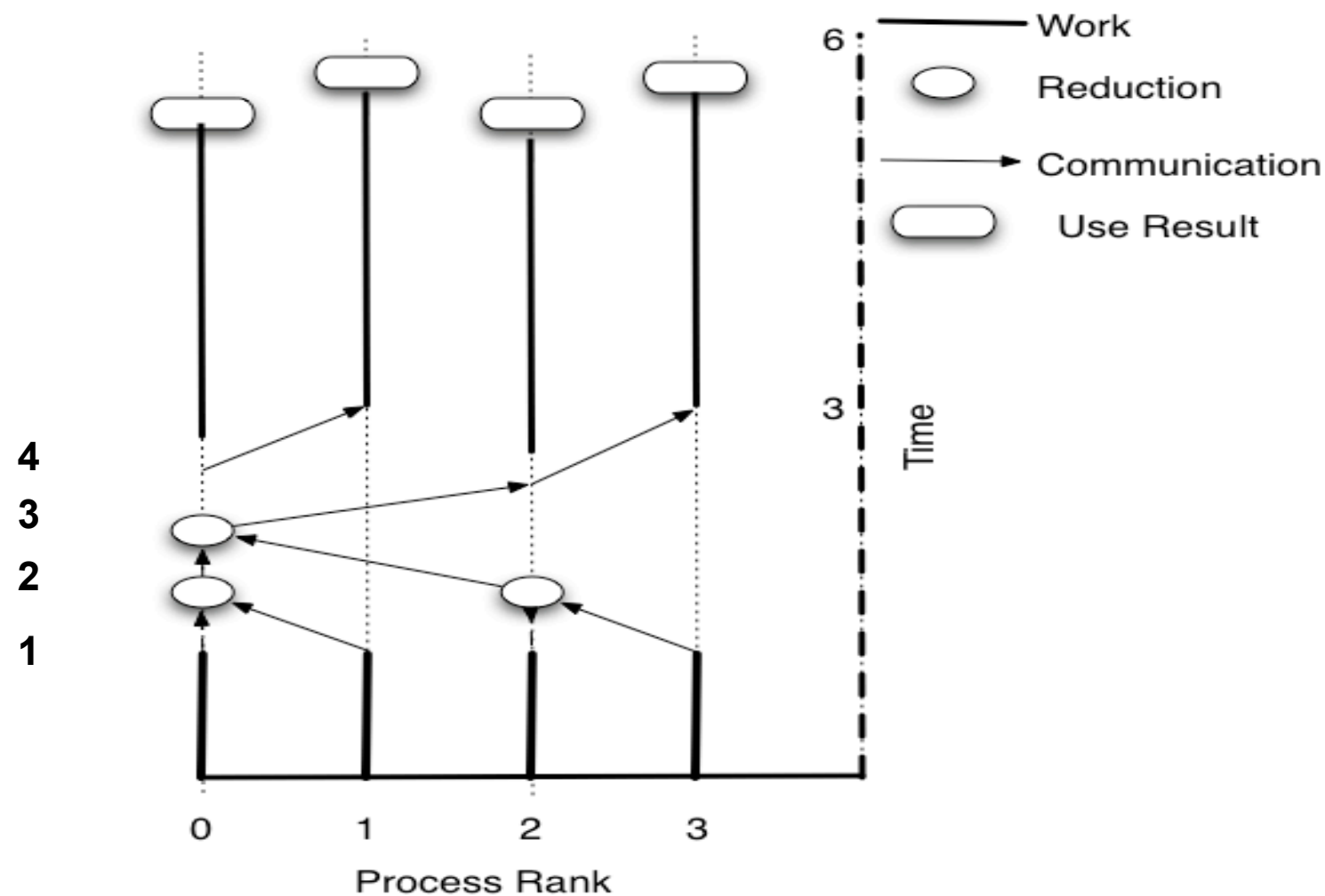
# Network Offload

```
0 4 7 1 7 4 3 7 0 4 7 8 4 1 0 0 1 8 1 4 7 5 5 0 8 7 2 9 4 1 0 0 0 4 7 1 7
1 8 2 3 4 1 0 4 1 8 2 4 6 6 2 1 6 4 0 9 9 4 1 3 7 8 7 8 6 6 2 1 1 8 2 3 4
9 9 9 5 9 9 3 9 9 9 9 1 9 8 0 6 8 1 3 8 8 6 6 2 4 5 4 4 9 8 0 6 9 9 9 5 9
4 4 3 4 8 8 5 8 4 4 3 9 9 0 4 5 4 5 2 6 7 5 5 0 5 4 8 5 9 0 4 5 4 4 3 4 8
9 9 1 0 0 7 4 4 9 9 1 8 8 4 1 1 6 9 1 0 5 9 4 0 6 8 0 5 4 8 4 1 1 6 9 9 1 0 0
8 8 4 6 1 9 6 5 8 8 4 4 8 9 9 0 5 8 3 5 1 6 1 4 4 2 8 7 6 8 9 9 0 5 8 8 4 6 1
6 0 8 2 3 0 0 2 3 5 3 0 9 4 2 0 6 0 8 2 3 0 0 2 0 2 1 1 2 3 0 0 9 1 8 6 3
1 5 5 3 2 6 1 0 6 4 5 3 8 9 0 1 1 5 5 3 2 6 1 0 1 1 5 1 4 2 0 6 7 2 9 5 6
5 4 9 0 1 5 4 1 5 9 1 6 7 8 1 9 5 4 9 0 1 5 4 1 6 6 6 6 5 1 6 5 0 5 2 7 4 5
4 4 4 1 8 1 9 4 4 8 3 5 4 7 4 8 4 4 4 1 8 1 9 4 5 5 5 5 1 6 1 1 9 7 9 6 4
4 9 9 4 7 6 6 4 9 0 2 4 9 8 4 7 9 4 9 9 4 7 6 6 4 1 1 1 3 4 5 4 6 8 8 5 4 9
9 8 8 9 4 4 5 9 6 6 1 9 8 4 9 4 8 9 8 8 9 4 4 5 9 4 4 2 2 6 6 6 5 5 5 9 6 6
8 4 4 8 9 5 1 8 4 8 5 6 8 0 0 8 9 7 8 4 4 8 9 5 1 8 4 6 5 0 1 5 5 5 1 0 9 6 5 8
7 6 5 7 8 1 6 4 9 4 4 5 4 4 2 4 8 4 7 6 5 7 8 1 6 4 9 5 1 5 5 1 1 1 3 1 7 2 4 4
4 5 4 8 4 3 5 1 8 6 9 7 6 8 1 1 0 4 5 4 8 4 3 5 1 8 1 0 1 6 3 3 4 2 6 8 1 6 6
6 1 5 9 6 0 1 6 4 5 5 4 5 7 6 6 4 6 1 5 9 6 0 1 6 4 6 2 6 5 5 2 2 7 5 2 6 8 5
2 6 3 0 5 5 4 5 4 1 9 3 9 5 5 6 2 6 3 0 5 5 4 5 4 8 0 5 1 4 5 5 5 0 1 4 4 4
1 8 4 1 0 6 6 1 9 3 8 1 9 0 1 4 1 8 4 1 0 6 6 1 5 1 4 5 2 1 1 9 5 6 2 6 9
```

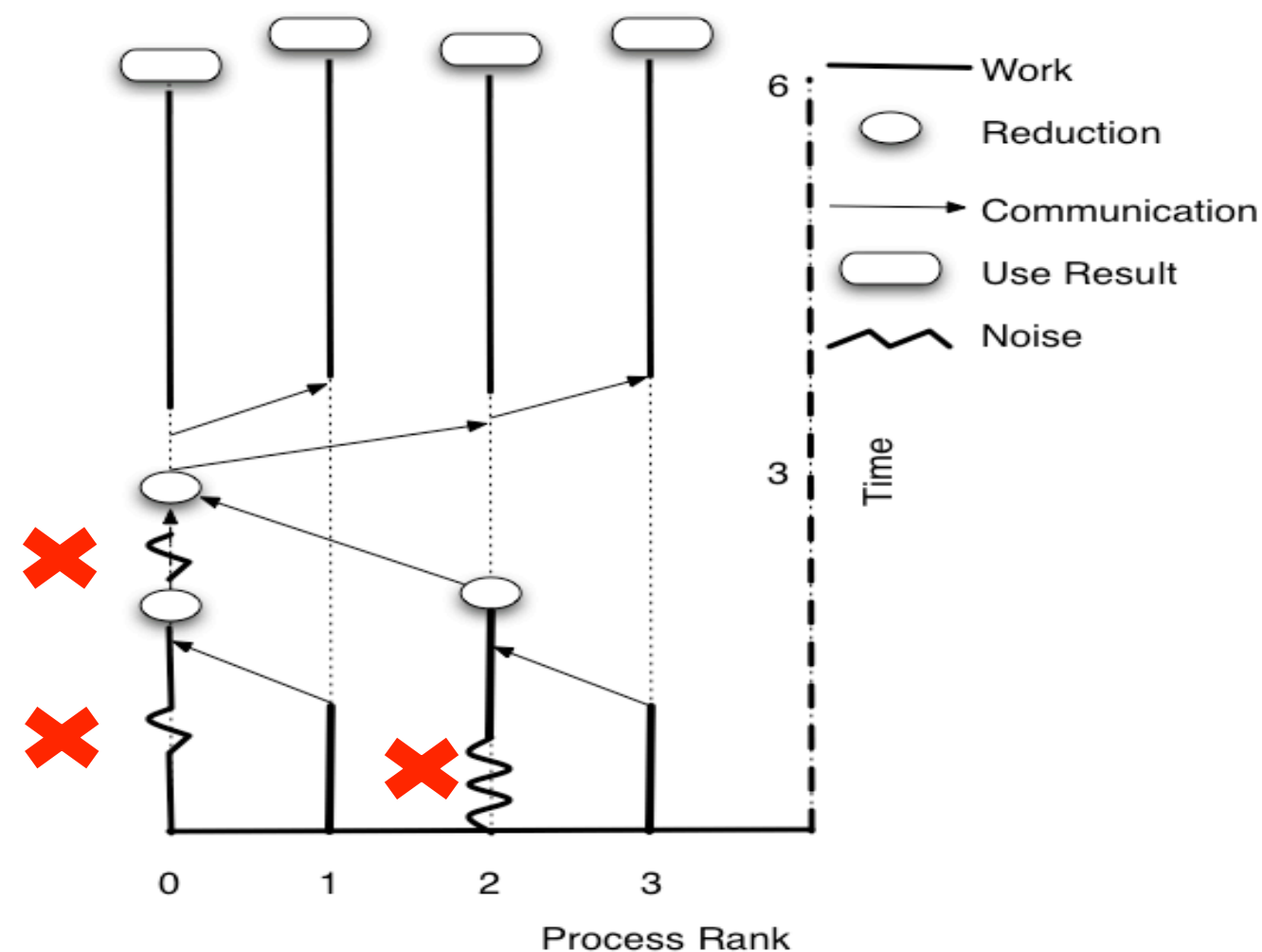
# Cross Channel Synchronization



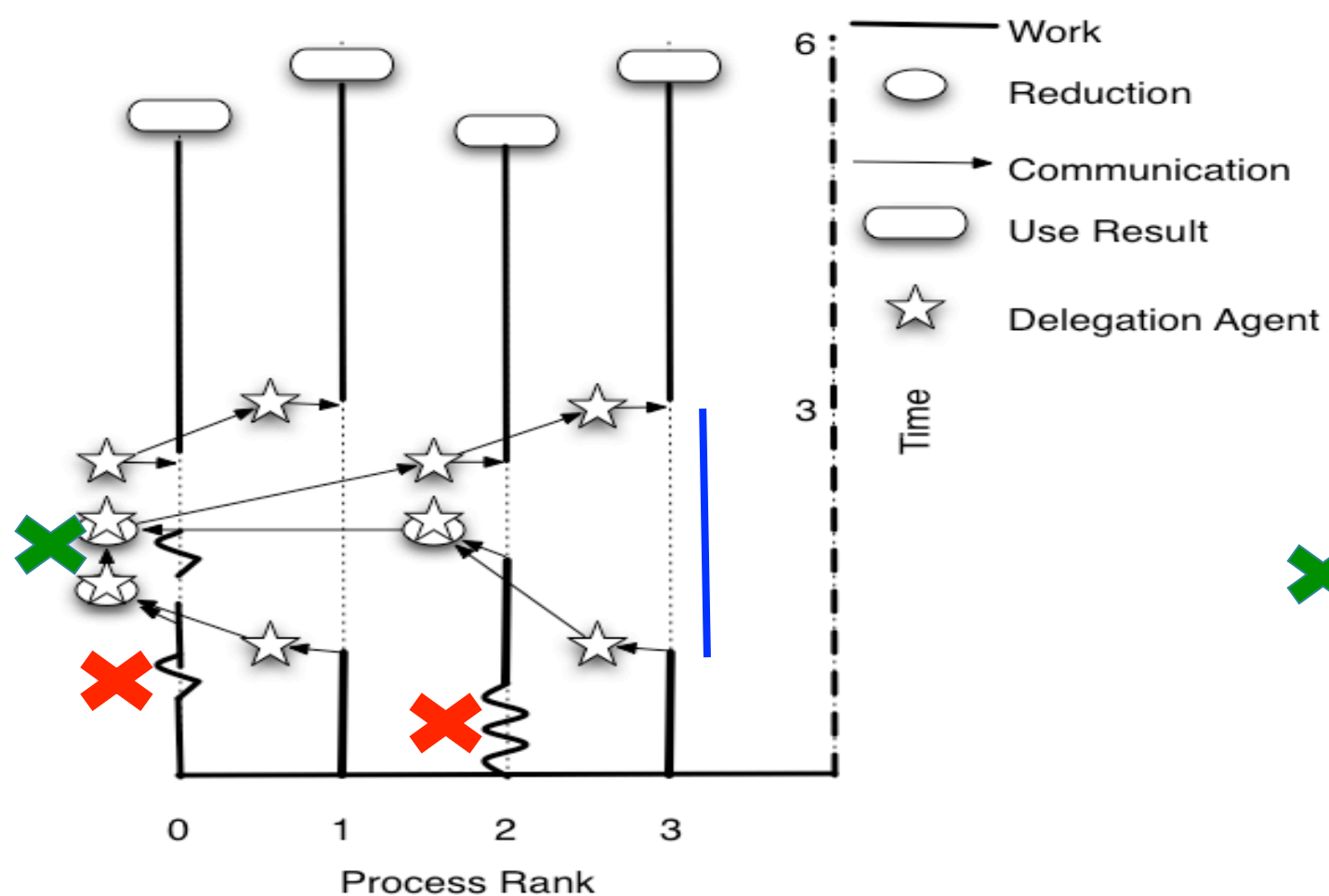
## Ideal Algorithm



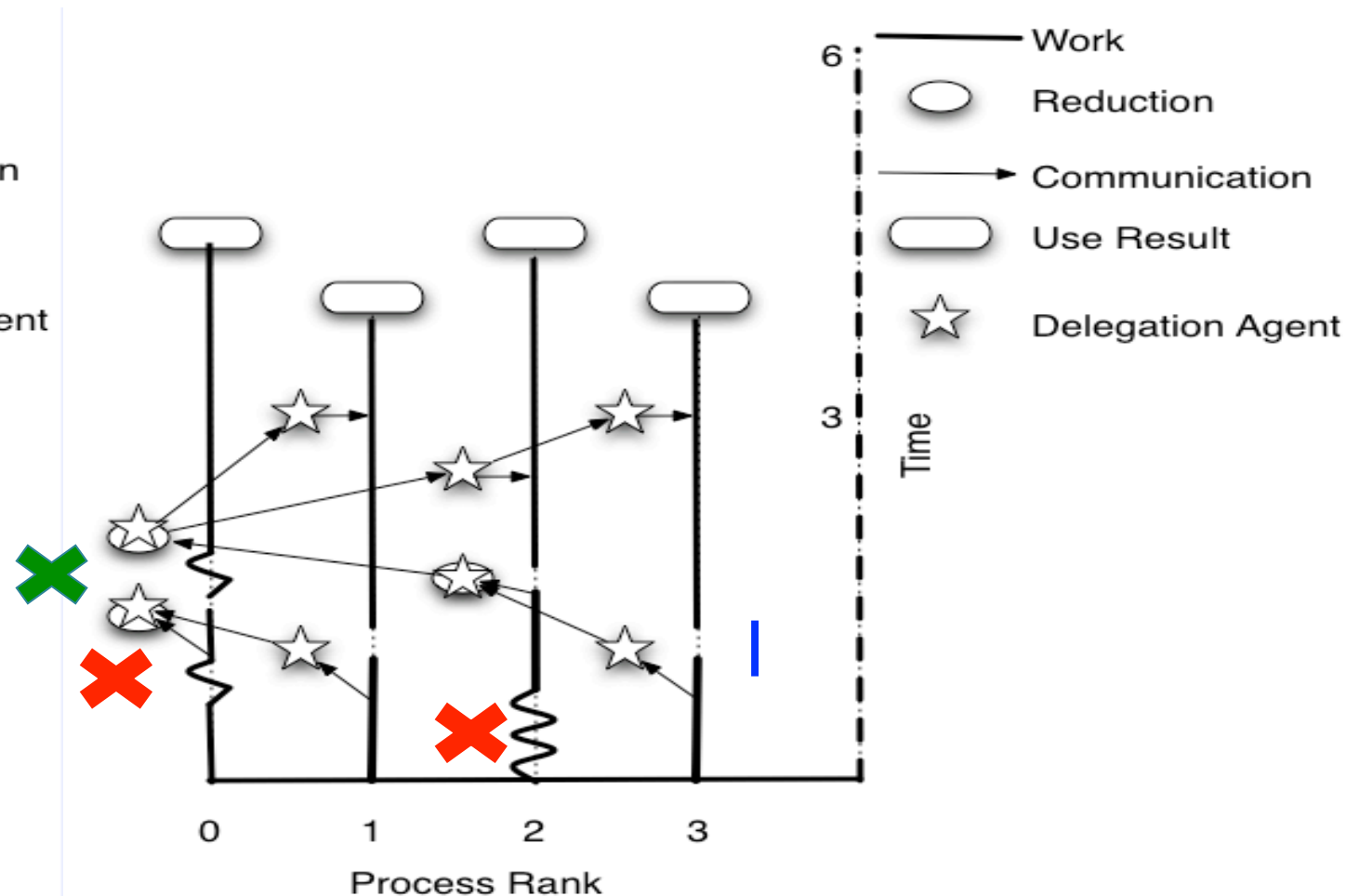
## Impact of System Noise



## Offloaded Algorithm



## Nonblocking Algorithm

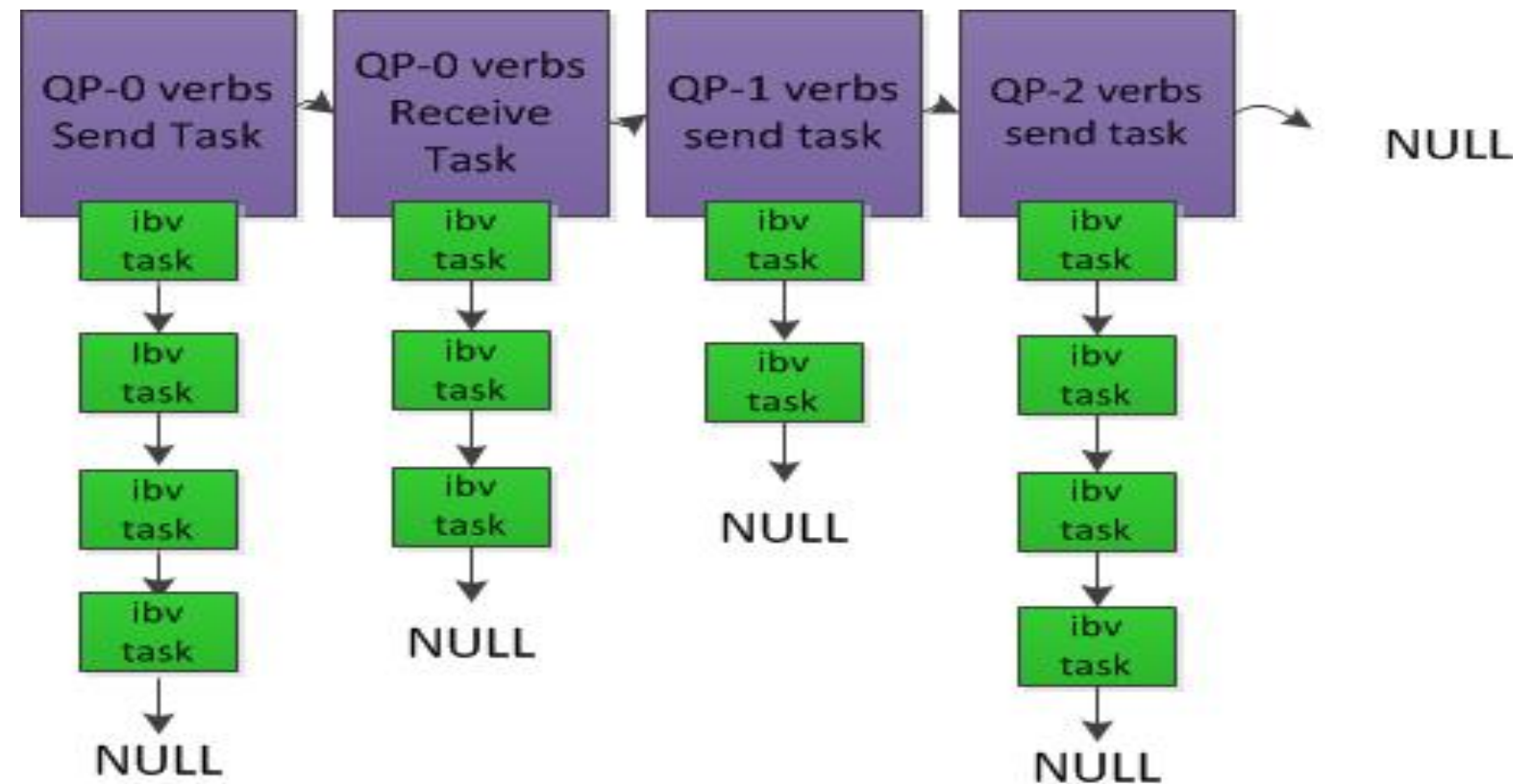


- Communication processing

# Cross Channel Synchronization (aka CORE-Direct)

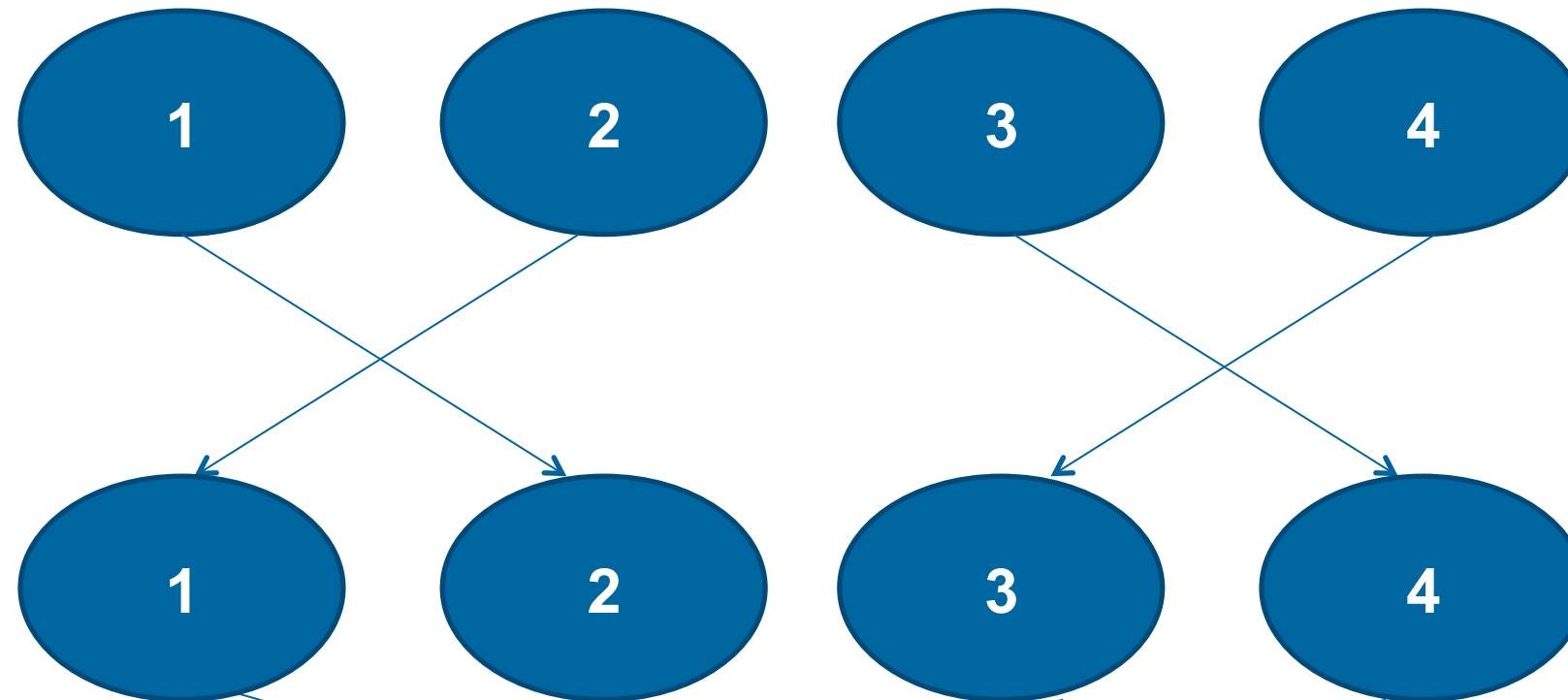
- Scalable collective communication
- Asynchronous communication
- Manage communication by communication resources
- Avoid system noise

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

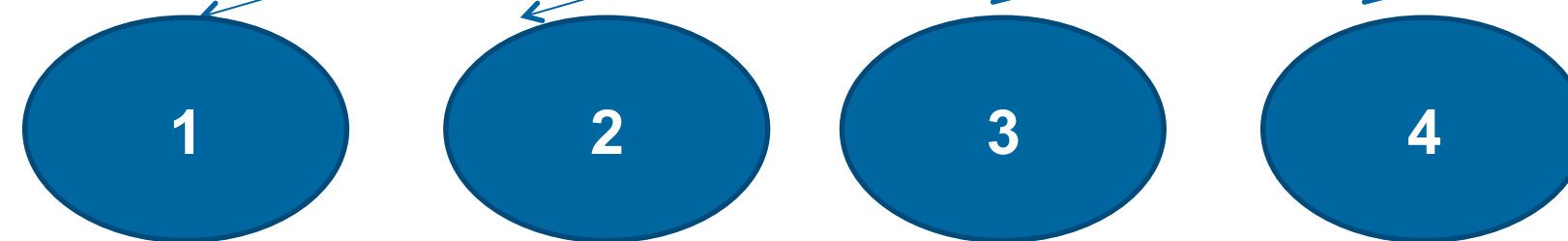


# Example – Four Process Recursive Doubling

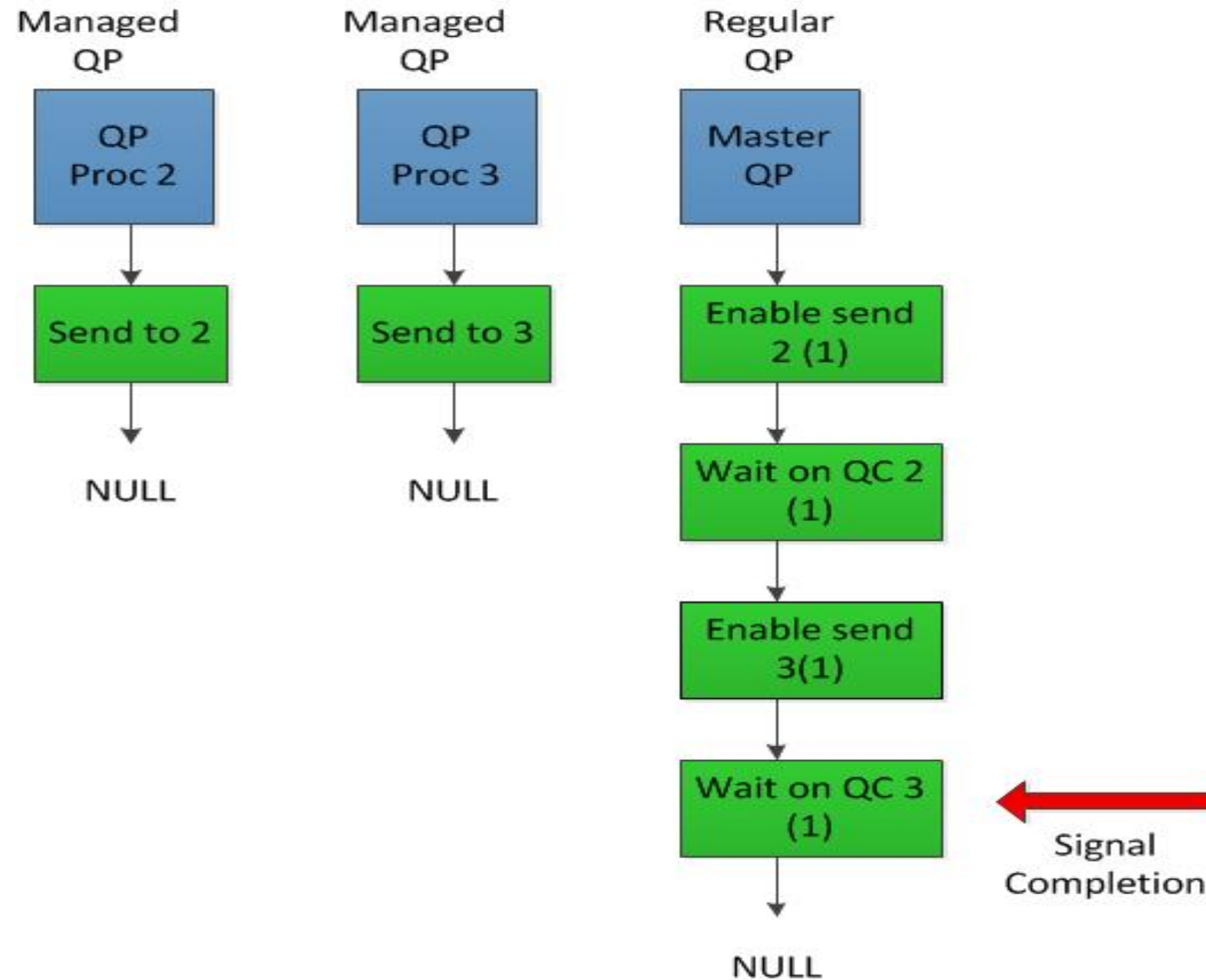
**Step 1**



**Step 2**

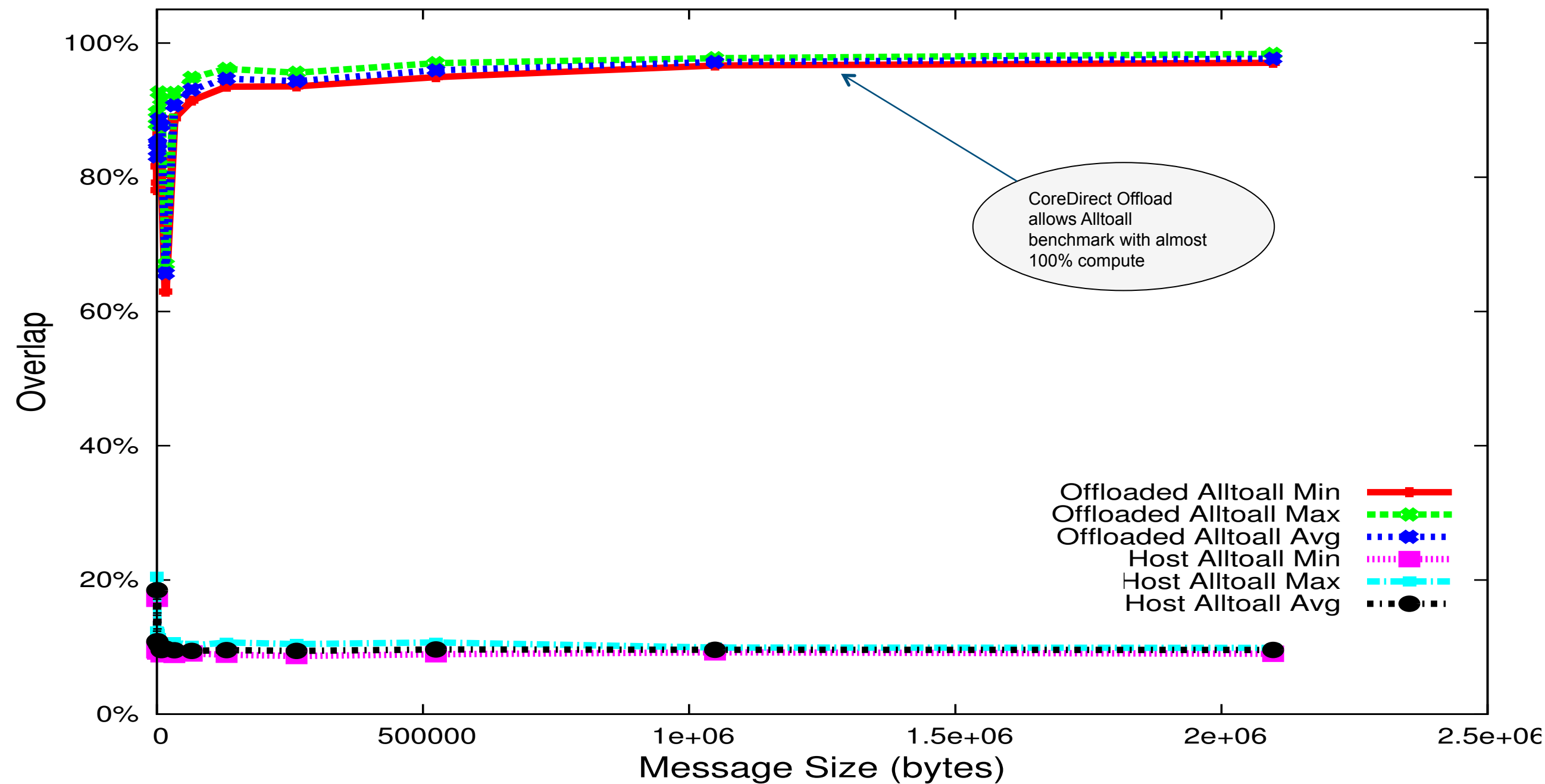


# Four Process Barrier Example – Using Managed Queues – Rank 0





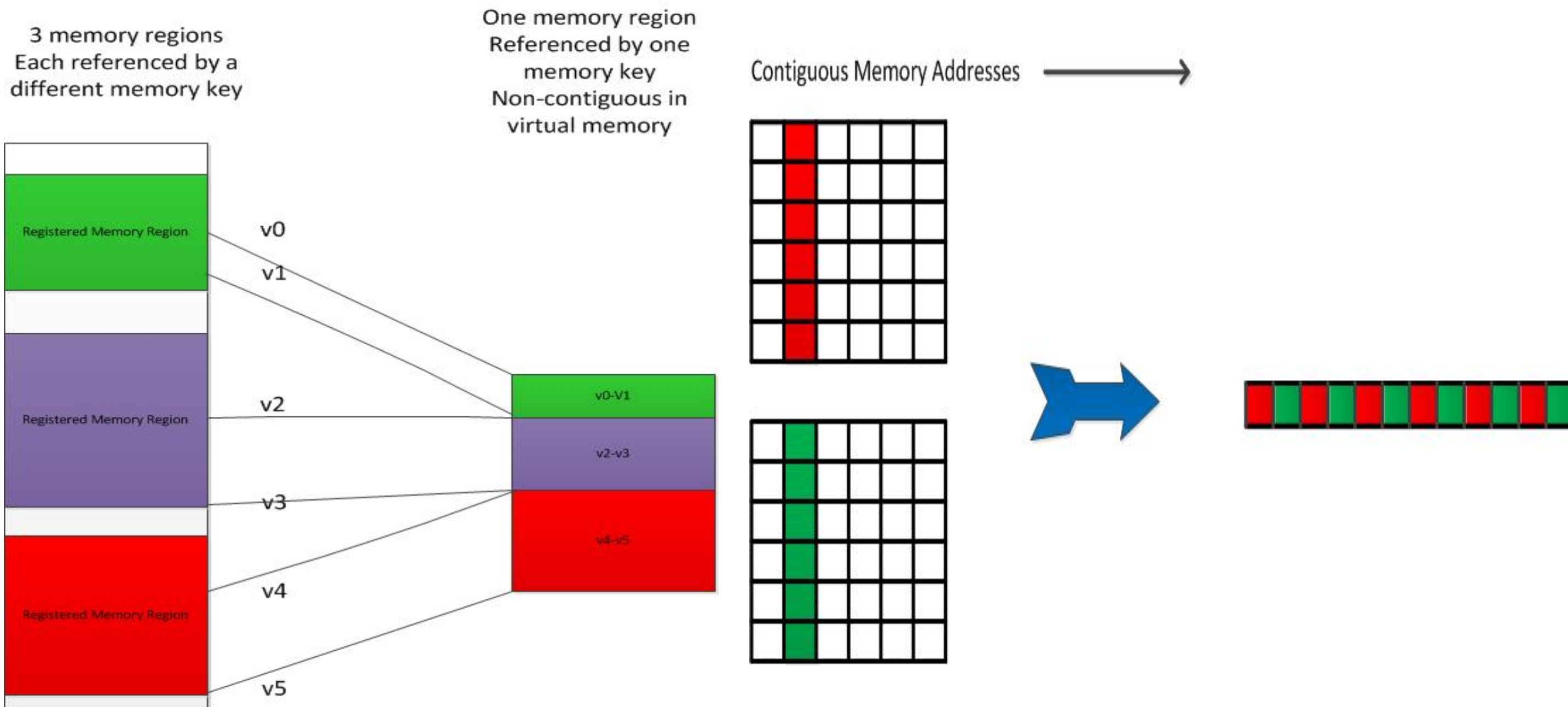
# Nonblocking Alltoall (Overlap-Wait) Benchmark



# Non-Contiguous Data

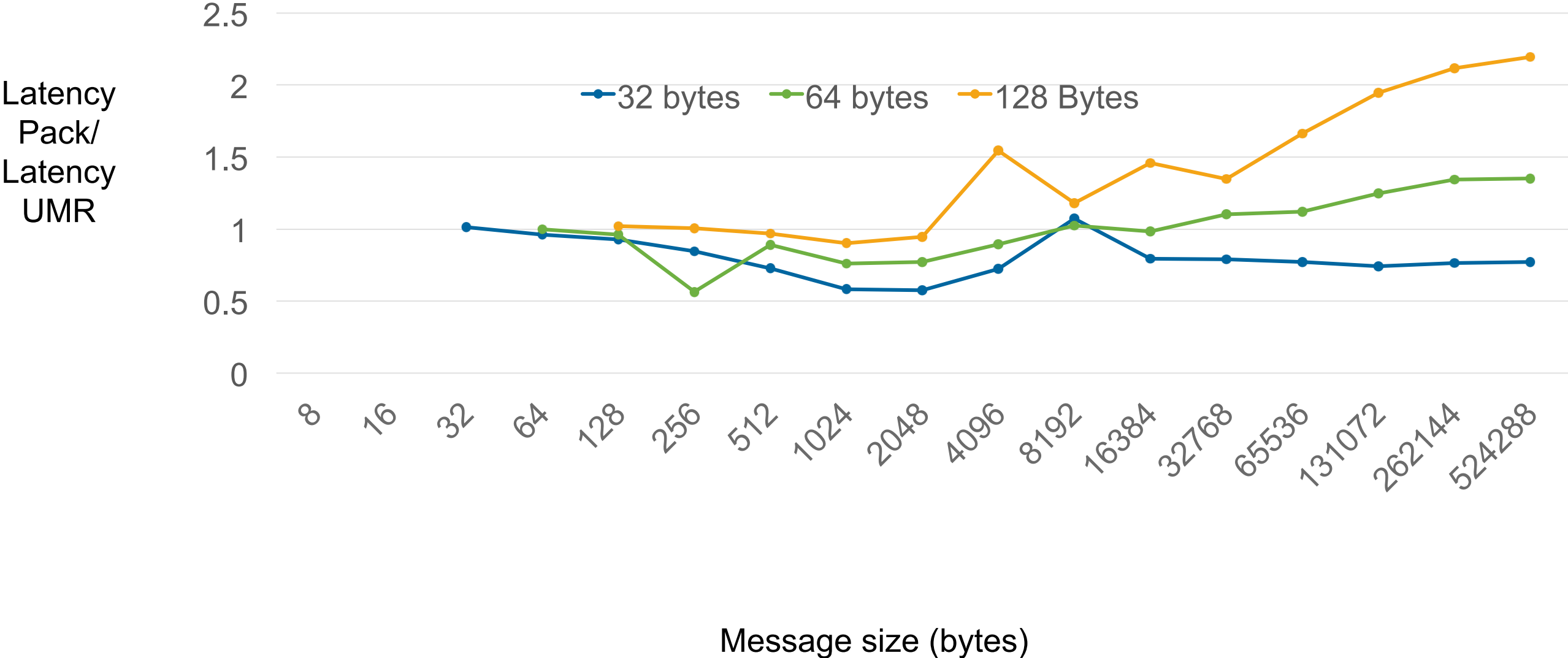
- Support combining contiguous registered memory regions into a single memory region. H/W treats them as a single contiguous region (and handles the non-contiguous regions)
- For a given memory region, supports non-contiguous access to memory, using a regular structure representation – base pointer, 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





# Hardware Gather/Scatter Capabilities – Regular Structure – Ping-Pong latency



# New Effort – Application Optimization



- Starting up effort to work on improving application performance
  - In house application domain experts
  - In house performance optimization experts
  - Looking for interested partners



# Thank You