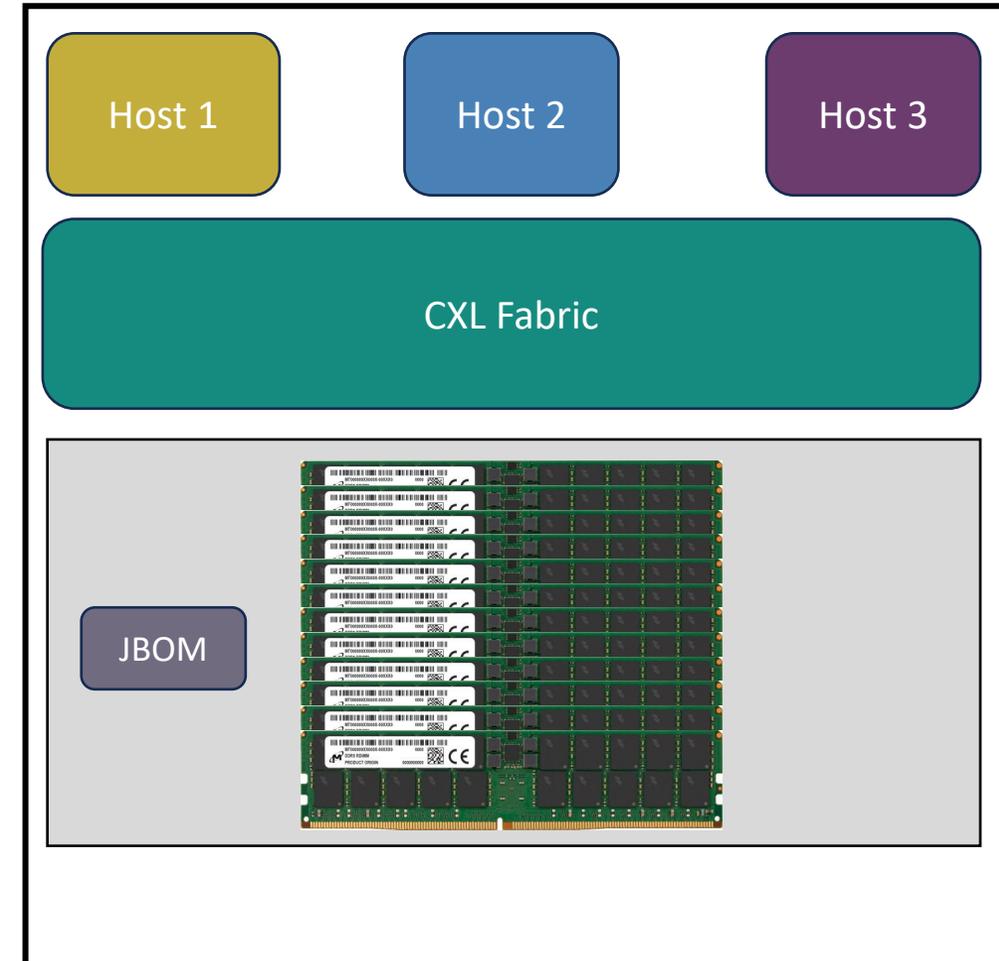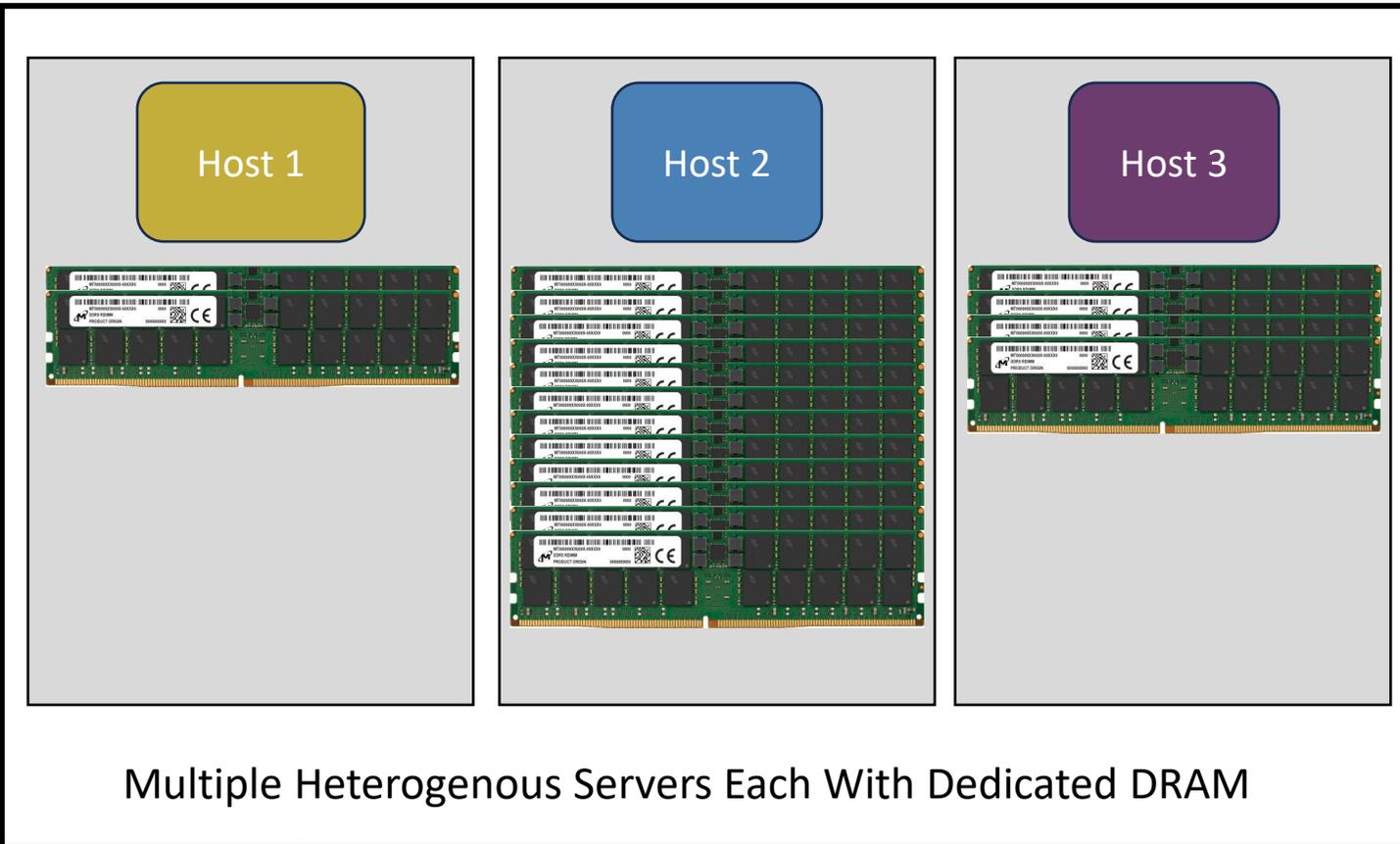# CXL Switch for Scalable & Composable Memory Pooling/Sharing

Presented by: JP Jiang

SVP, XConn Technologies

# CXL Enables DRAM Disaggregation for Usage Optimization

Host 1

Host 2

Host 3

Multiple Heterogenous Servers Each With Dedicated DRAM

Host 1

Host 2

Host 3

CXL Fabric

JBOM

# CXL Switch for Scalable Disaggregation

XConn Tech has developed
CXL2.0 (XC50256) & PCIe
5.0 (XC51256) switch IC

2,048 GB/s total
BW with 256 lanes



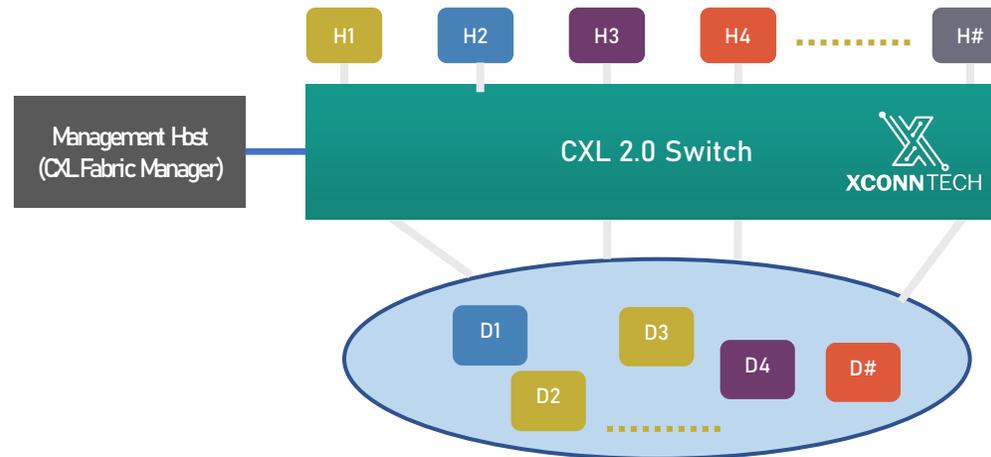Lowest port-to-port latency

Lowest power consumption/port

Reduced PCB area
Lower TCO

- Works with CXL 1.1 server processors, CXL memory devices.
- Works with the upcoming CXL 2.0 processors.
- Works in <u>hybrid</u> mode (CXL/PCIe mixed).

# Scalable Memory Pooling & Sharing Enabled by CXL 2.0 Switch

Memory Pooling/Sharing with
CXL 1.1/2.0 Hosts & Single Logical Devices



- A single CXL Switch connects to 32 combined hosts/devices
- Fully support CXL Fabric Manager
- Support switch cascading for a larger size memory pool
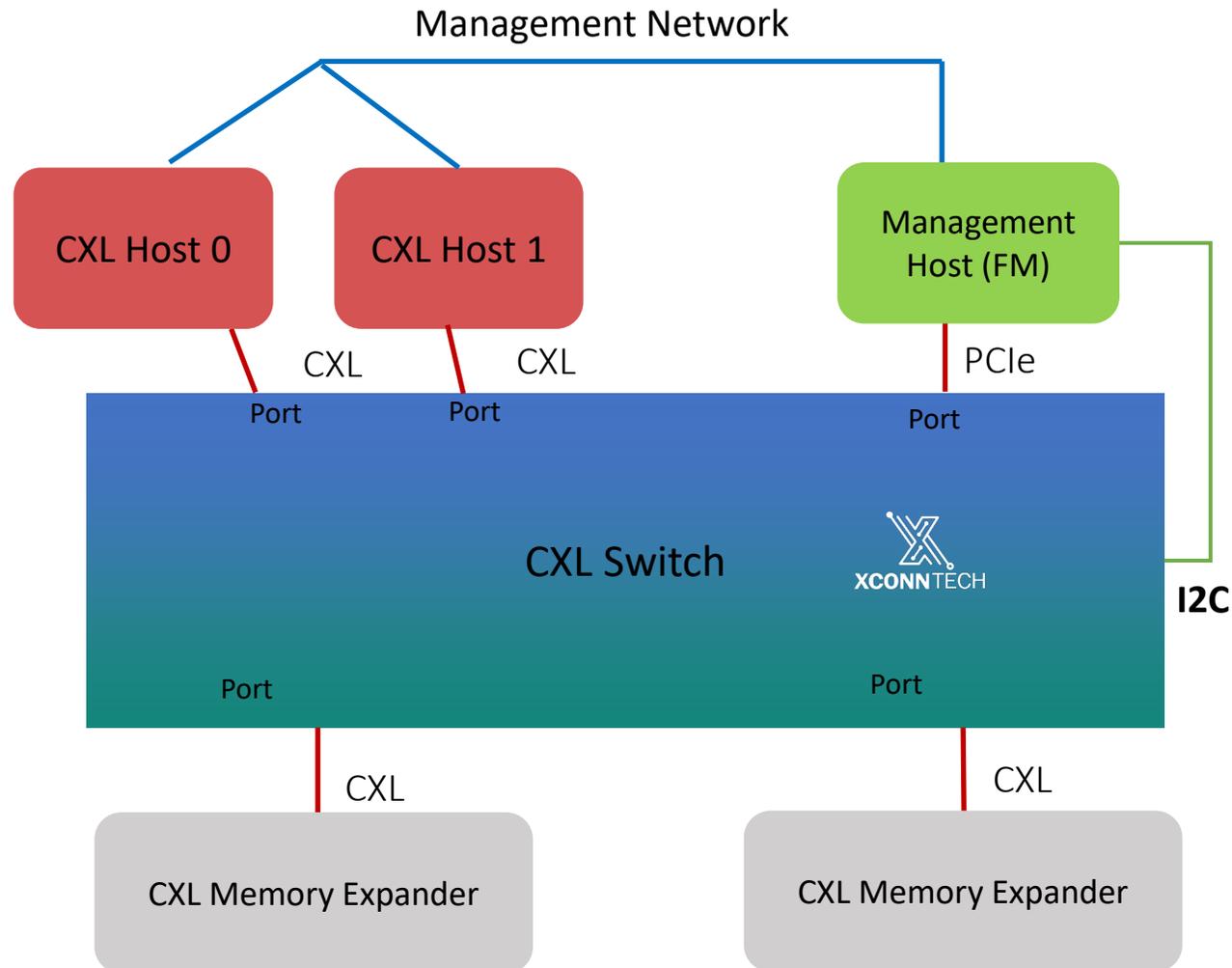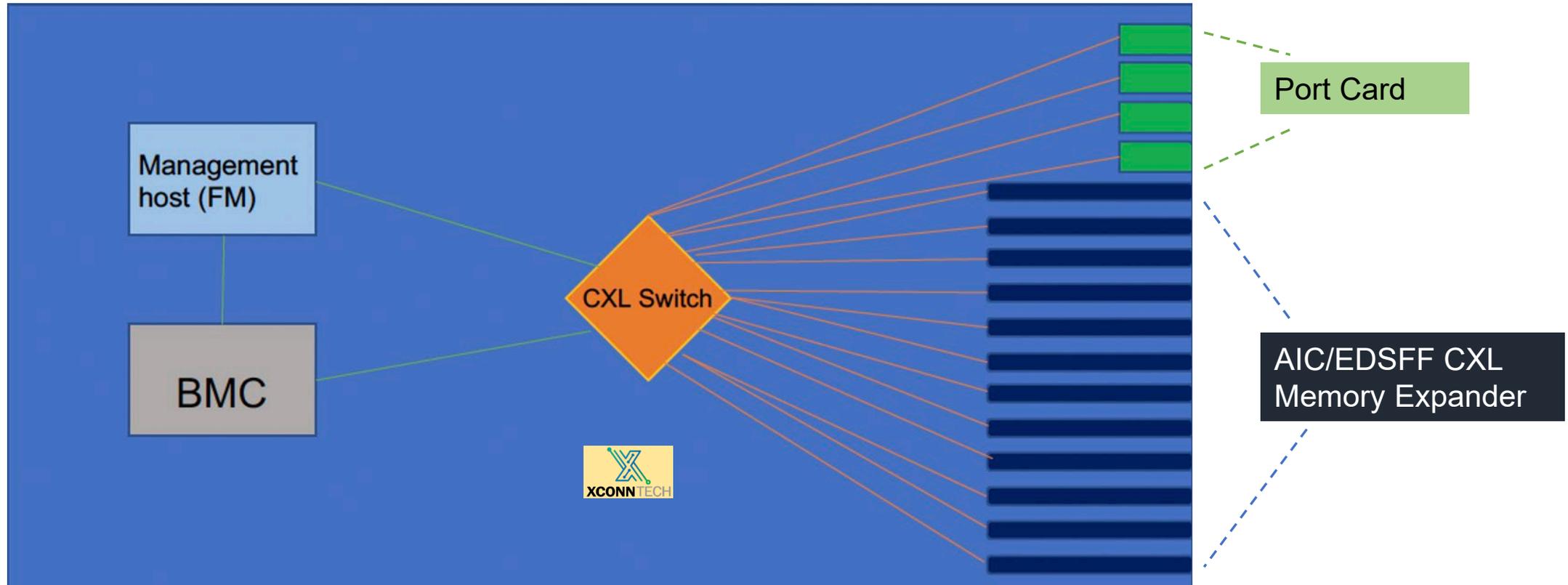
# Applications with Memory Pooling

- In memory database
    - Database requires huge amount of memory (up to 100TB) to maintain performance
    - Sharing memory is ideal for databases running on multiple hosts
    - Load/store offered by CXL is more efficient than RDMA
    - Significant SAP-HANA performance increase with pooled CXL memory(Samsung, MemCon 2024)
- AI Inference
    - Inference requires larger memory capacity in order to keep performance
    - Xconn is collaborating with partners to utilize CXL memory to enhance inference performance
- Solution for "memory wall" and lowering TCO
    - CXL memory expansion/pooling to address "memory wall"
    - Hosts sharing a large memory pool while keeping minimal size of local memory
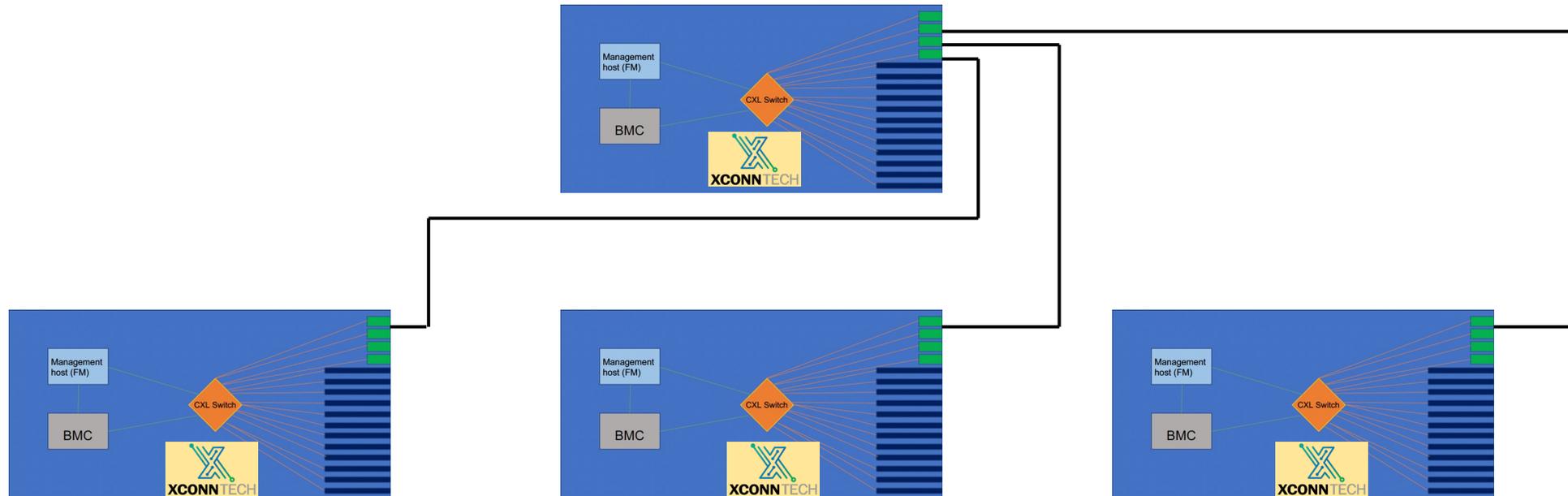    - Reuse DRAMs (such as DDR4) from replaced servers

# Memory Pooling/Sharing PoC Topology

# Memory Pooling/Sharing System

# Scalable Memory Pooling/Sharing System

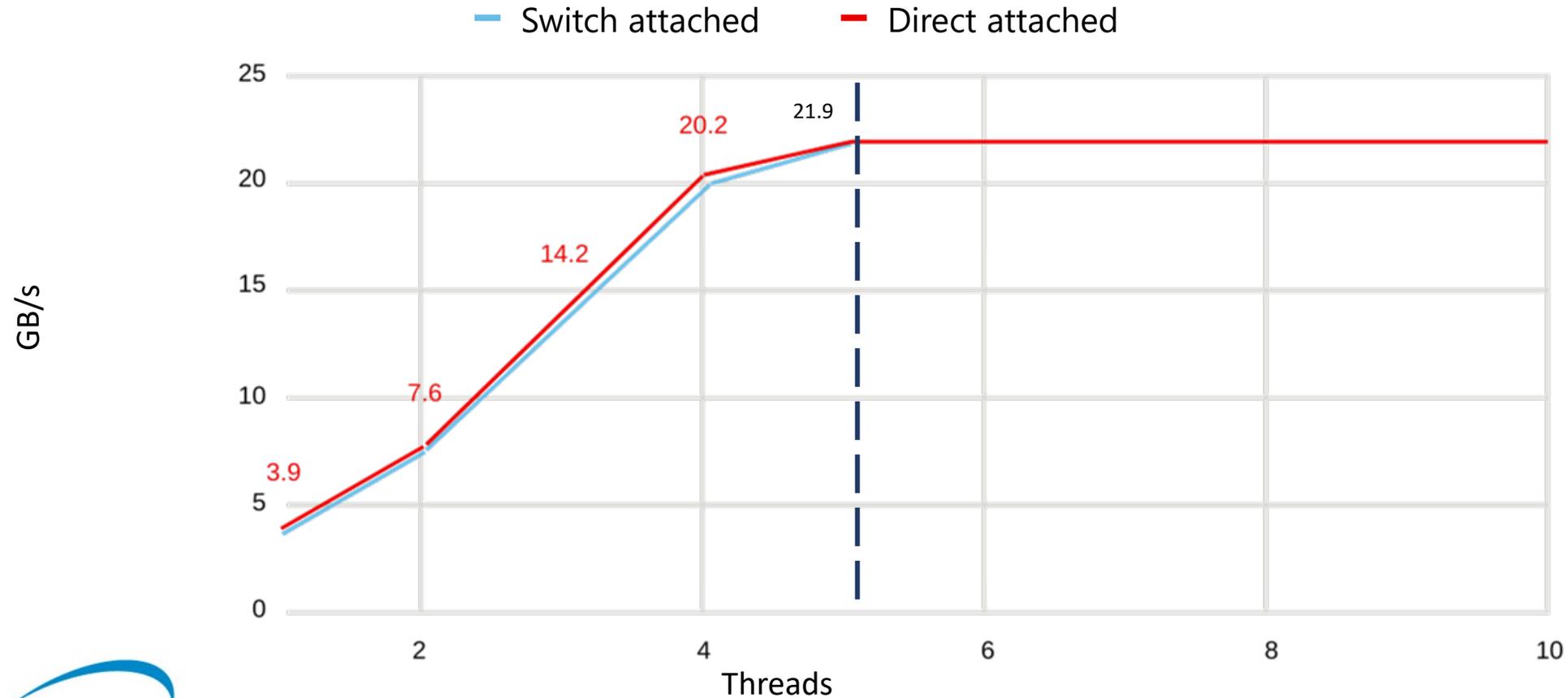# Rack Scale Memory Pooling Appliance



Samsung's CMM-B is a revolutionary system:

- CXL 2.0 spec compliant memory pooling
- Scalable up to 16TB
- Software-managed memory with fabric manager
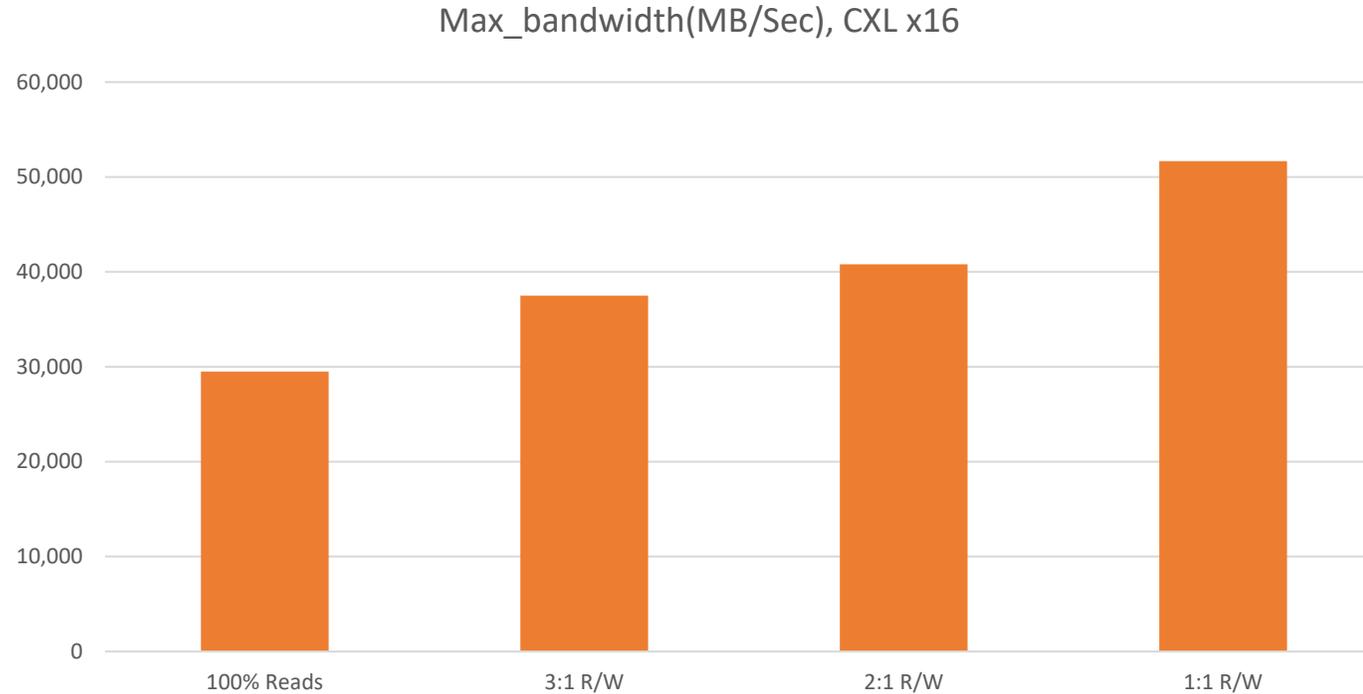- Power efficiency
- Cost effective

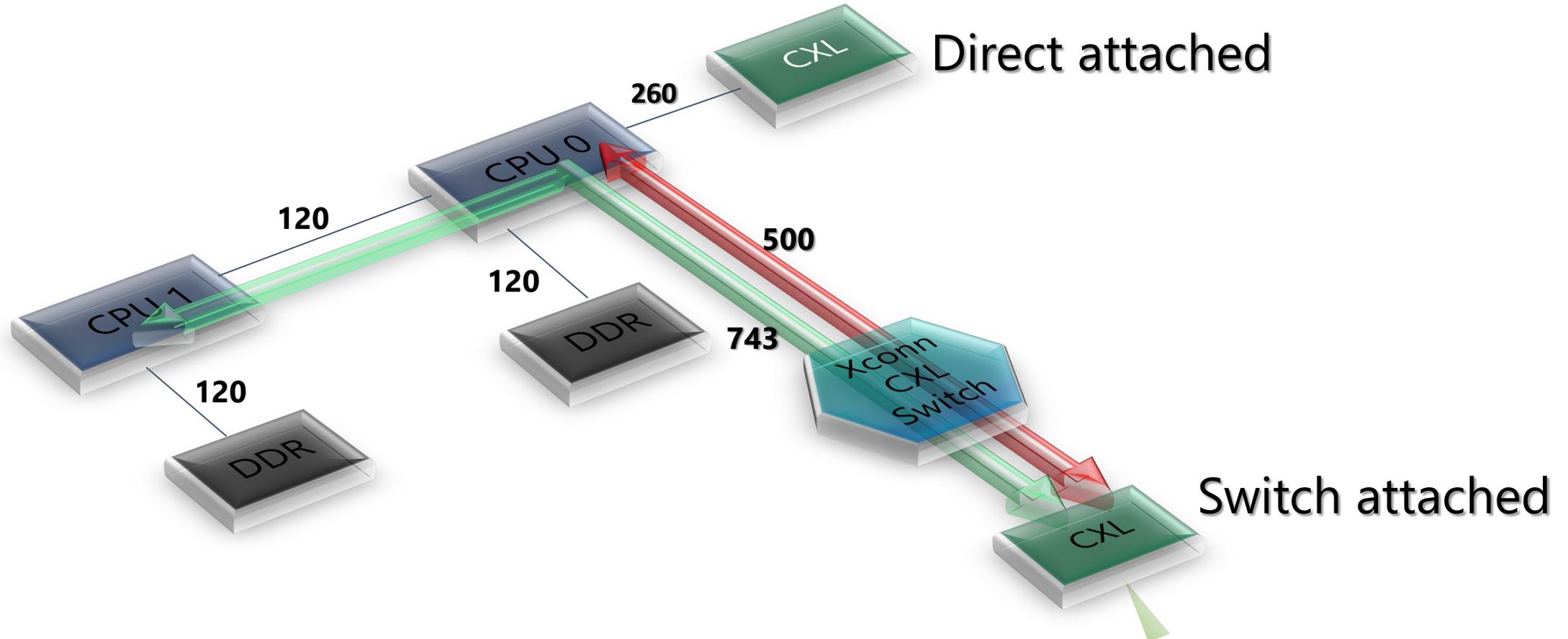# Bandwidth— Switch Attach vs Direct Attach

Bandwidth measured by MLC, CXL x8

# More Bandwidth Results— Switch Attached

Max_bandwidth(MB/Sec), CXL x16
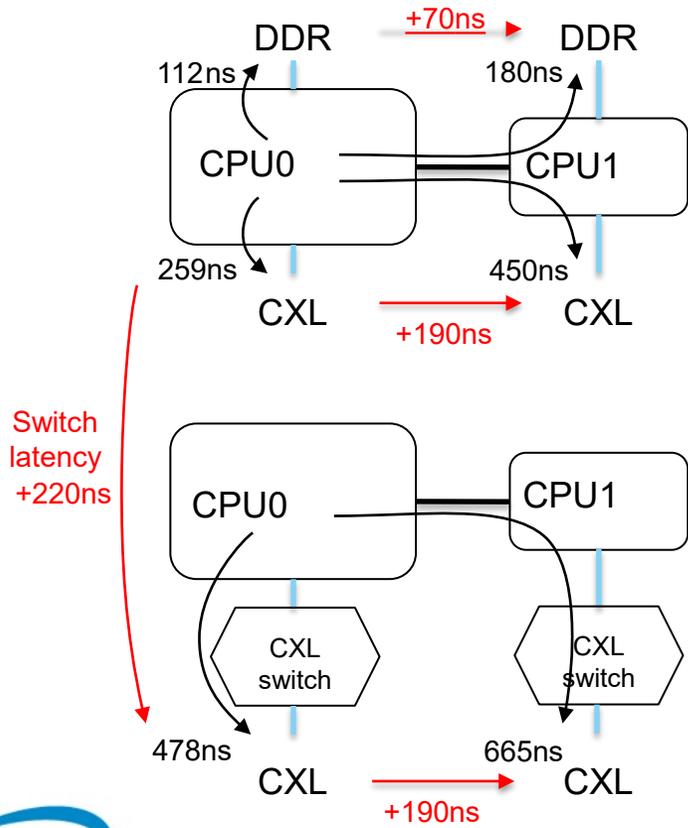
# Switch Idle Latency (ns)



Direct attached

260

120

500

120

743

Switch attached

120

Samsung 128GB CXL CMM-D

# More Latency Testing Results

# In Closing

- CXL switch enabled memory pooling provided a solution for "Memory Wall" for AI and HPC computing

- Software enabled memory sharing finds plethora of applications in database and AI inferencing

- XConn CXL 2.0 switch and fabric manager (FM) provide scalable and composable memory pooling/sharing solution with decent performance

# Thank You!

Https://www.xconn-tech.com
Email: JP.Jiang@xconn-tech.com