

# Octopus: Enhancing CXL Memory Pods via Sparse Topology

Yuhong Zhong 

Shuwei Teng 

Fiodar Kazhamiaka 

Rodrigo Fonseca 

Daniel S. Berger  

Pantea Zardoshti 

Mark D. Hill 

# Executive Summary: CXL as “Pod Area” Network

## Background:

- Memory is expensive (~50% of CapEx) but underutilized due to stranding
- CXL enables a “pod” of servers to reach a common **memory pool**
- CXL also enables intra-pod **communication at memory-like latency**
- However, existing CXL pods are either **expensive** (using CXL switches) or **small** (using multi-port CXL devices, each with  $\leq 4$  ports)

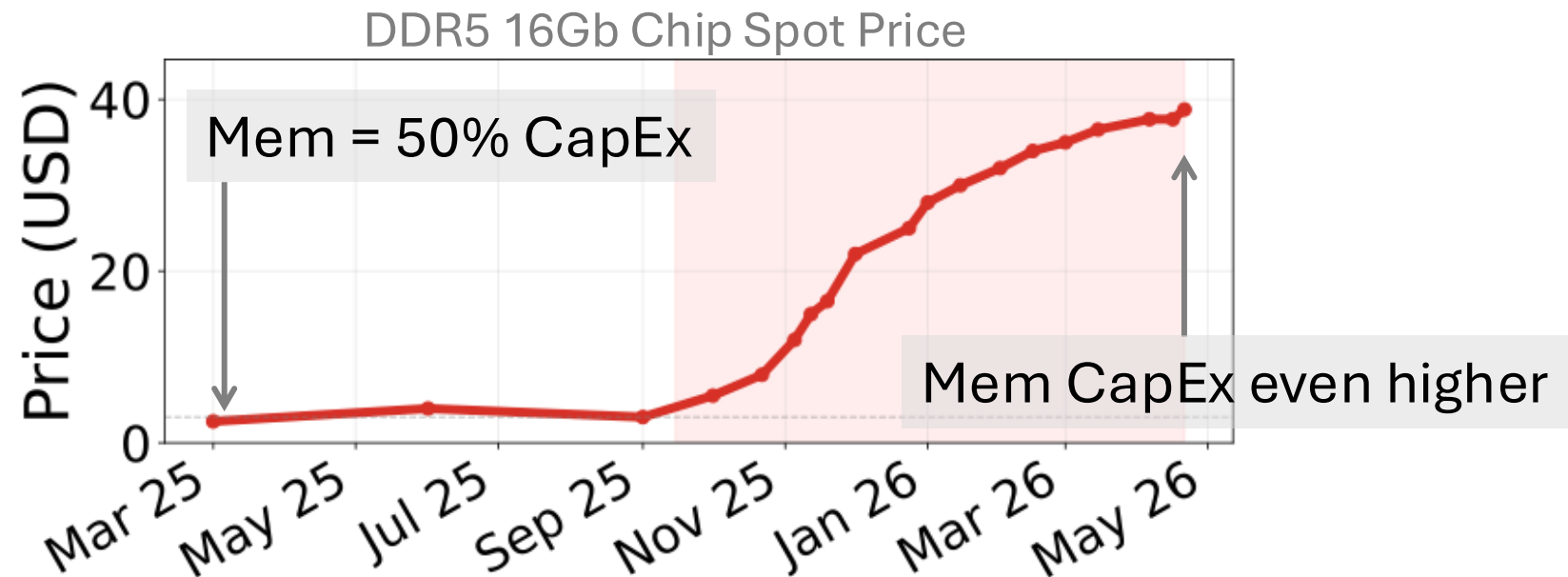
## Contributions:

- Octopus: a low-cost, large-scale CXL pod which **sparsely** connects servers with low-port CXL devices, without using costly switches
  - To balance conflicting demands of memory pooling and intra-pod communication, Octopus partitions a pod into “**islands**”
- Results: 16% memory saving,  $< 1 \mu\text{s}$  intra-pod communication

# Memory Accounts for Up to 50% Server CapEx

In 2023, memory accounts for **50% of server CapEx**

Recent demand surge increases DRAM prices by 10×

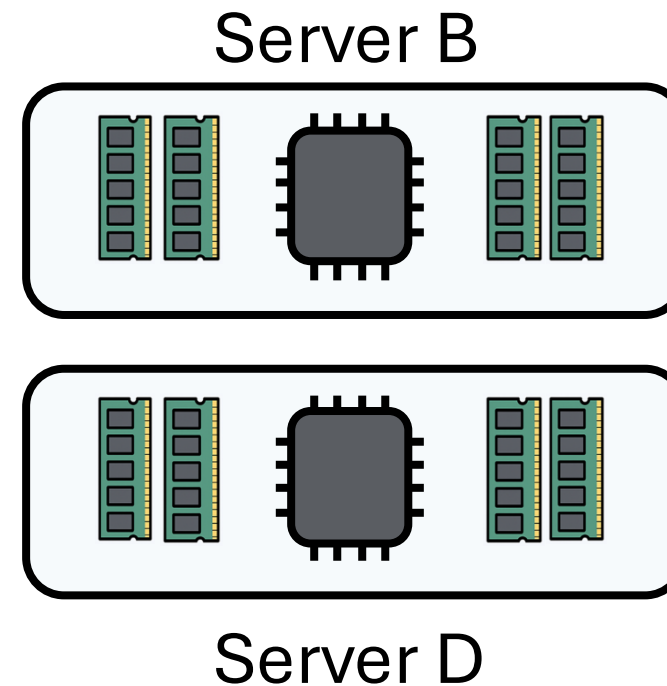
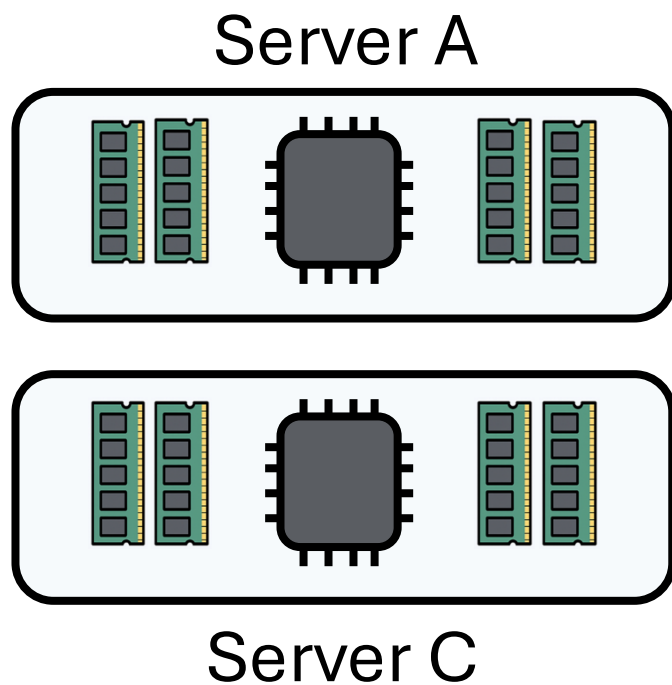


Memory capacity is expensive and limited: **utilization** is critical!

# CXL to Rescue Memory Efficiency via Pooling

Memory capacity is overprovisioned for **peak demand per-server**

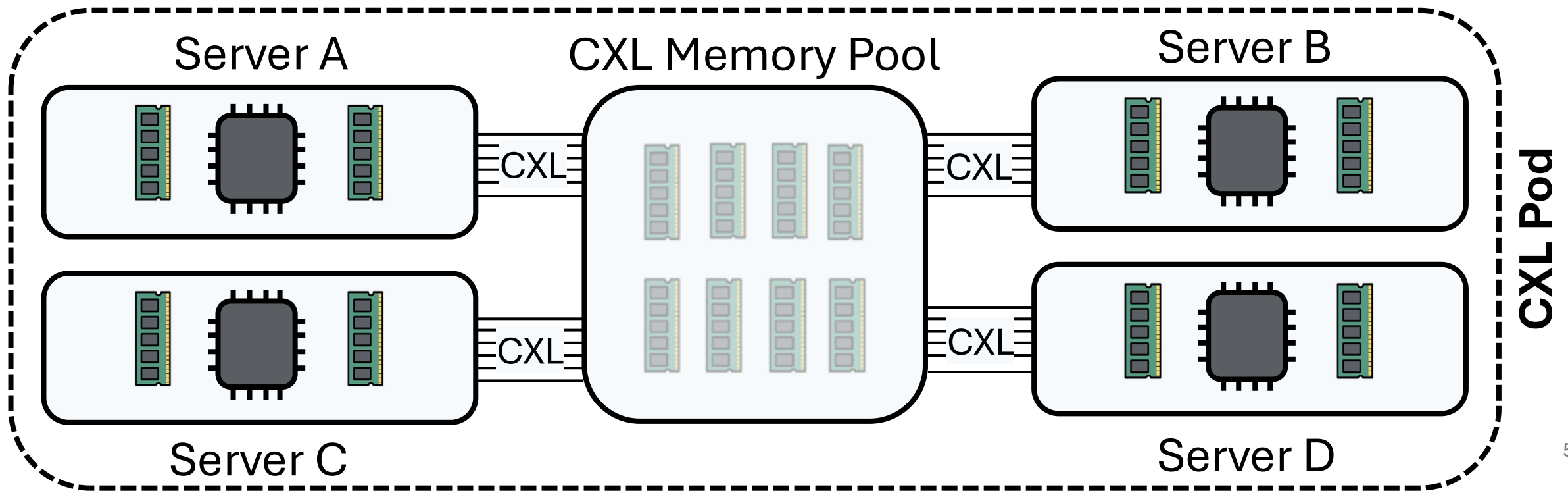
- Idle memory cannot be used by other servers, being “stranded”



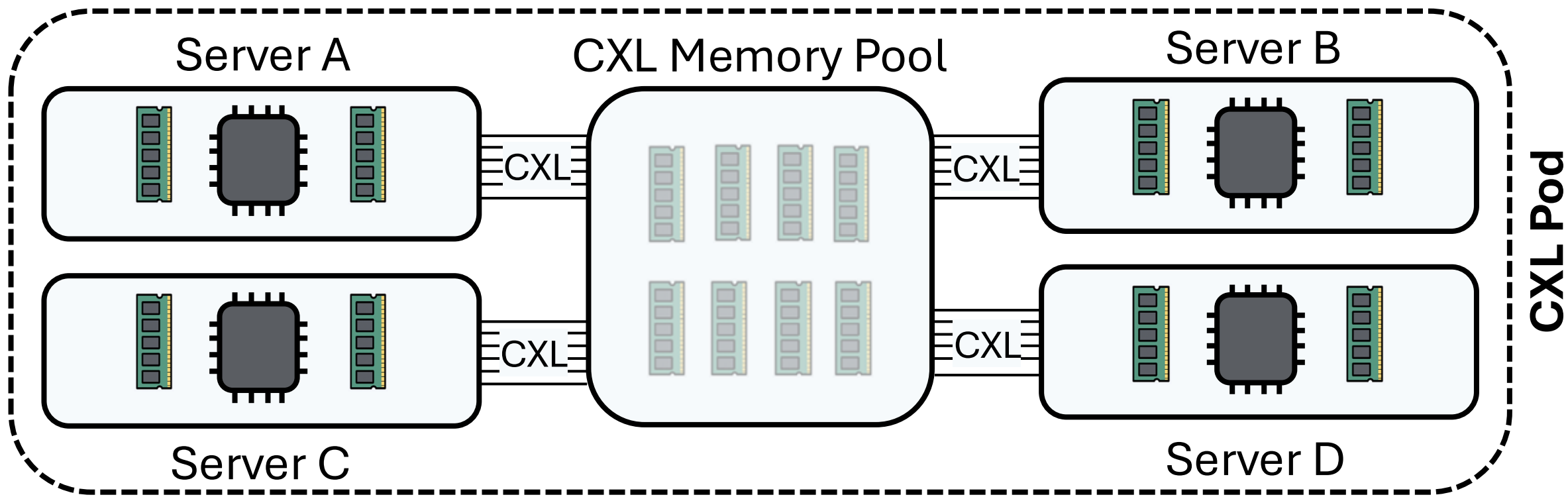
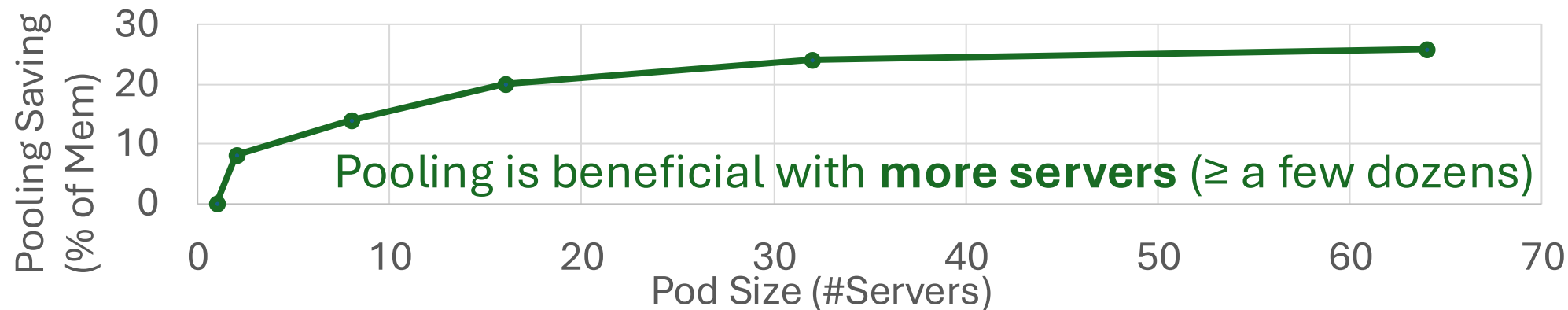
# CXL to Rescue Memory Efficiency via Pooling

CXL memory pools **multiplex** memory demand across servers

- CXL memory allocated *on demand*, less stranded capacity
- Provision capacity for *aggregate peak*, not  $\text{peak} \times N$



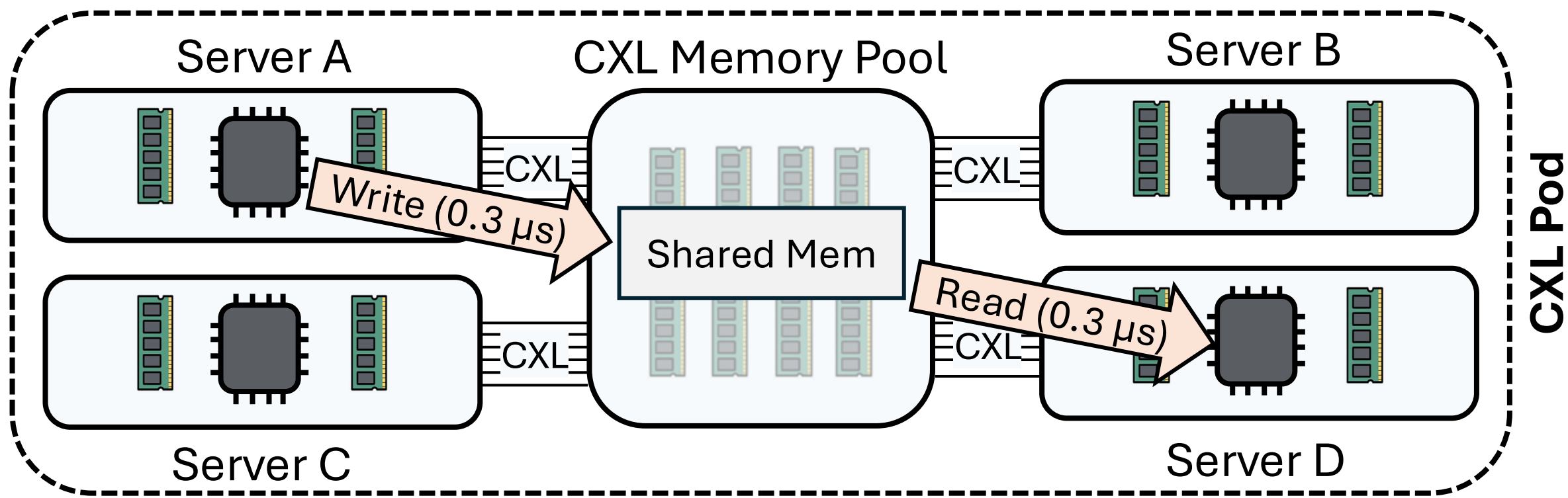
# CXL to Rescue Memory Efficiency via Pooling



# CXL Enables Low-Latency Communication

For communication, CXL ( $0.6 \mu\text{s}$ ) is much faster than RDMA ( $\sim 25 \mu\text{s}$ )

- Especially useful for distributed systems (e.g., consensus protocols, collective communication, and scale-out databases)



# Design CXL Pods to Enable Pooling and Comm.

Key CXL use cases: **Memory pooling** and **low-latency communication**

For memory pooling:

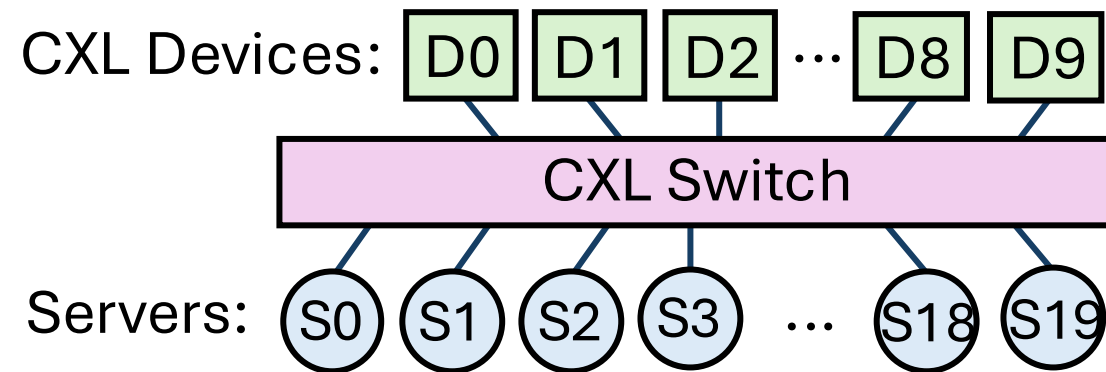
- (#1) **Large CXL pods** to pool across many servers (a few dozens)
- (#2) **Low CXL latency** to pool memory for more workloads
  - Lower latency → more workloads can tolerate CXL → more memory to pool

For communication at memory-like latency:

- (#3) **Common CXL devices** between servers to use shared CXL mem

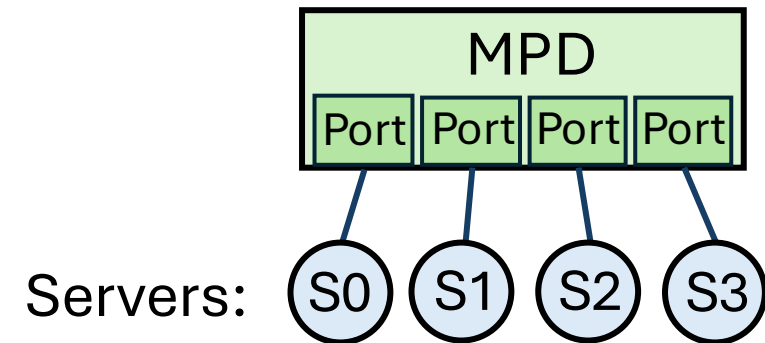
# Challenge: Building CXL Pods With Today's CXL Devices

## Option 1: CXL Switches



- ✗ High latency (490-600 ns)  
Less memory to pool (35% of total mem)
- ✗ Expensive to deploy (\$5000+)
- ✓ High fanout (20 servers)  
Larger CXL pods to pool mem effectively

## Option 2: Multi-Ported Devices (MPD)

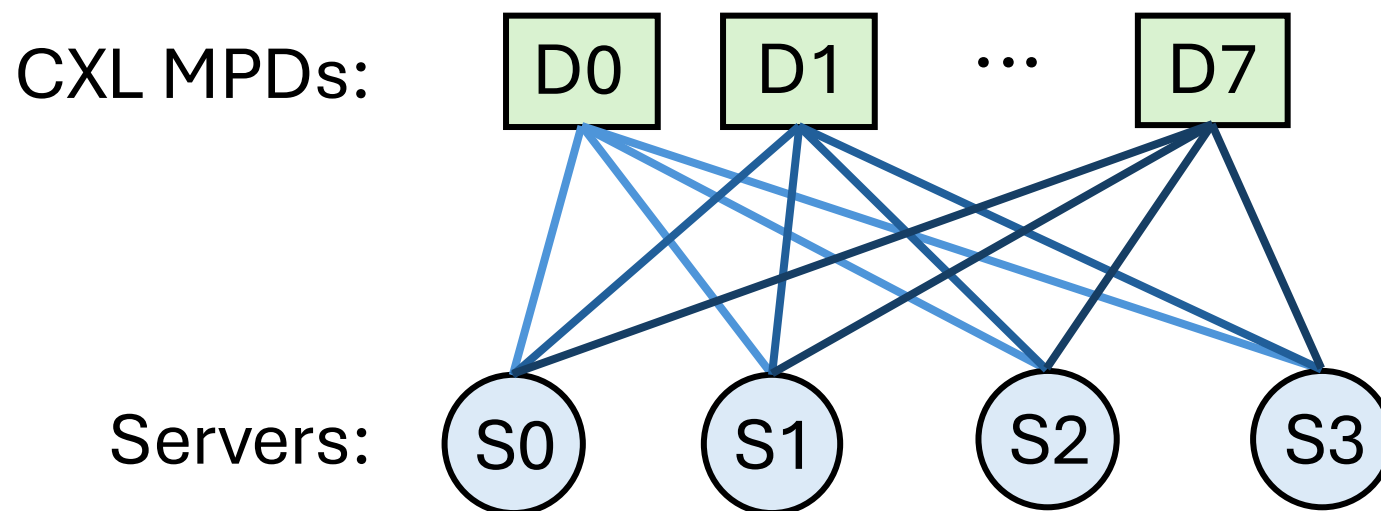


- ✓ Low latency (260-300 ns)  
Pool up to 65% of total memory
- ✓ Cheap to deploy (no switch)
- ⚠ Low fanout (2-4 servers)  
Smaller pods have limited pooling saving

# Prior Work: Fully-Connected Servers and MPDs

Prior work: **fully-connected** CXL pods (every MPD connects to every server)

- CXL pod size = MPD port count (2-4 ports)



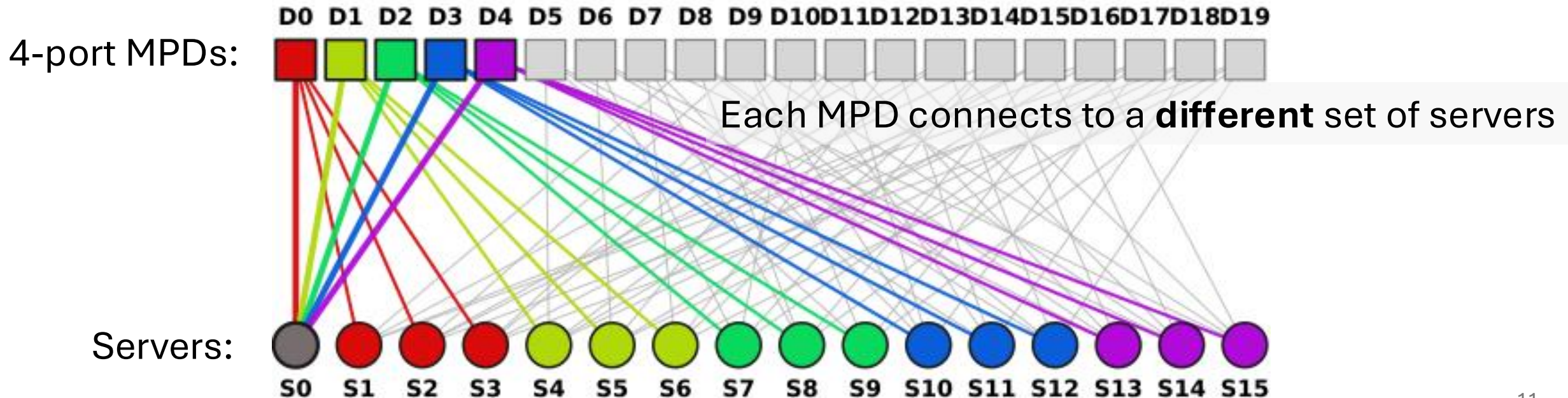
Each MPD connects to the **same** set of servers

# Octopus: *Sparsely* Connect Servers and MPDs

Prior work: **fully-connected** CXL pods (every MPD connects to every server)

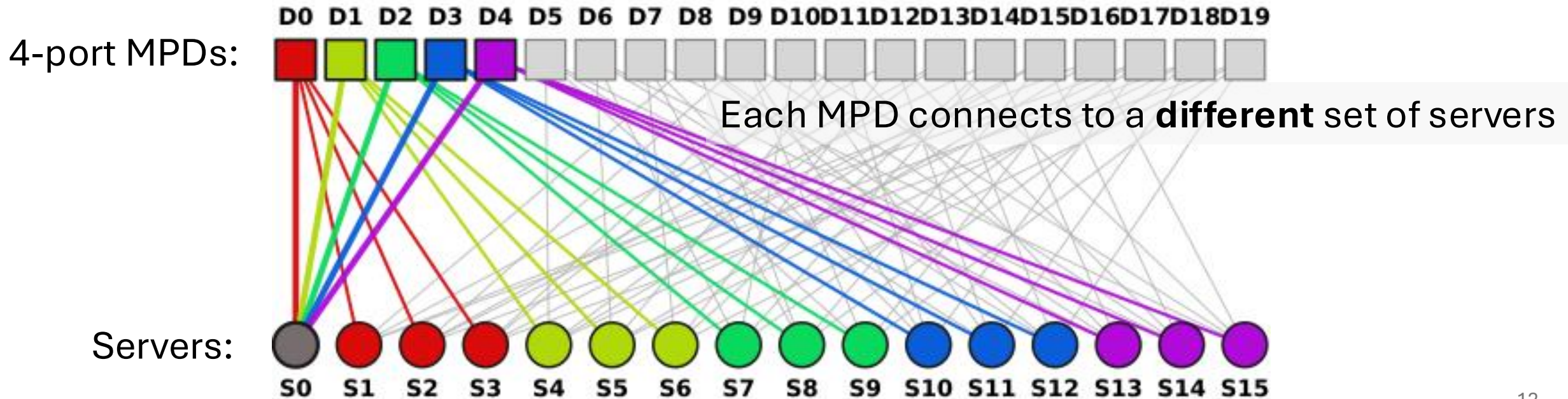
- CXL pod size = MPD port count (2-4 ports)

Key design: **sparsely** connect servers and MPDs to build **larger** CXL pods



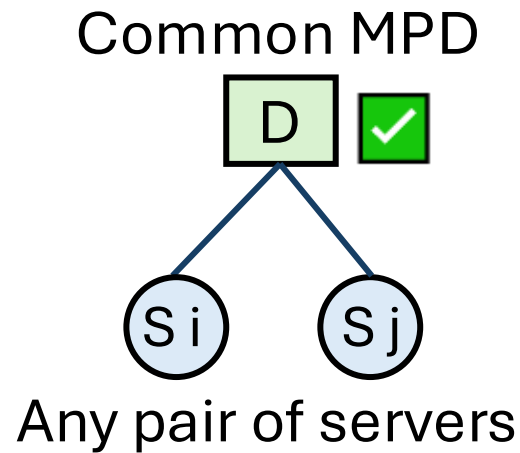
# Octopus: ***Sparse*** Connect Servers and MPDs

Key question: constructing a **sparse topology** between servers and MPDs to satisfy both memory pooling and low-latency communication



# Low-Latency Communication and Memory Pooling Are **Conflicting**

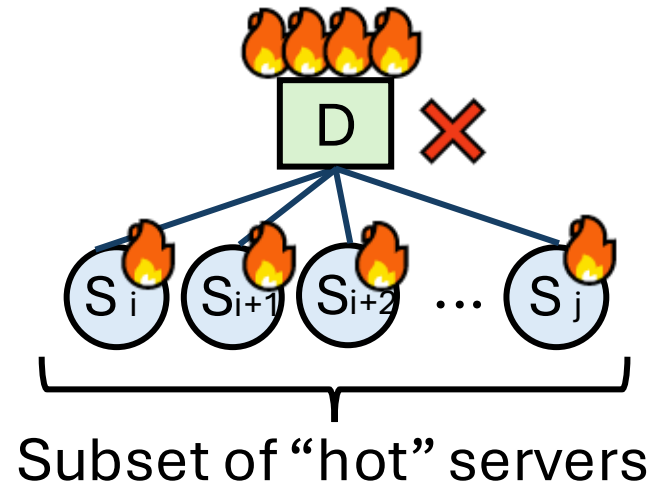
## Low-latency communication:



Property: **Pairwise MPD overlap**

Optimal: Balanced Incomplete Block Design (BIBD)

## Memory pooling:

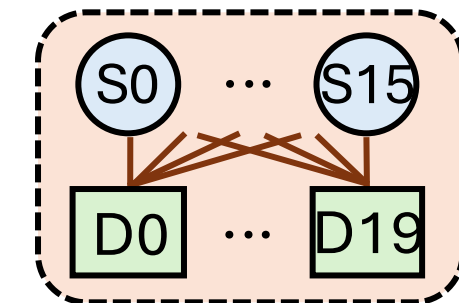


Property: **Avoid MPD overlap** to spread load into as many MPDs as possible

Optimal\*: Expander graphs (e.g., random or Ramanujan graphs)

# Reconcile MPD Overlap With “Islands”

Insight: Low-latency communication is typically needed only at modest scale (e.g.,  $\leq 12$  servers)



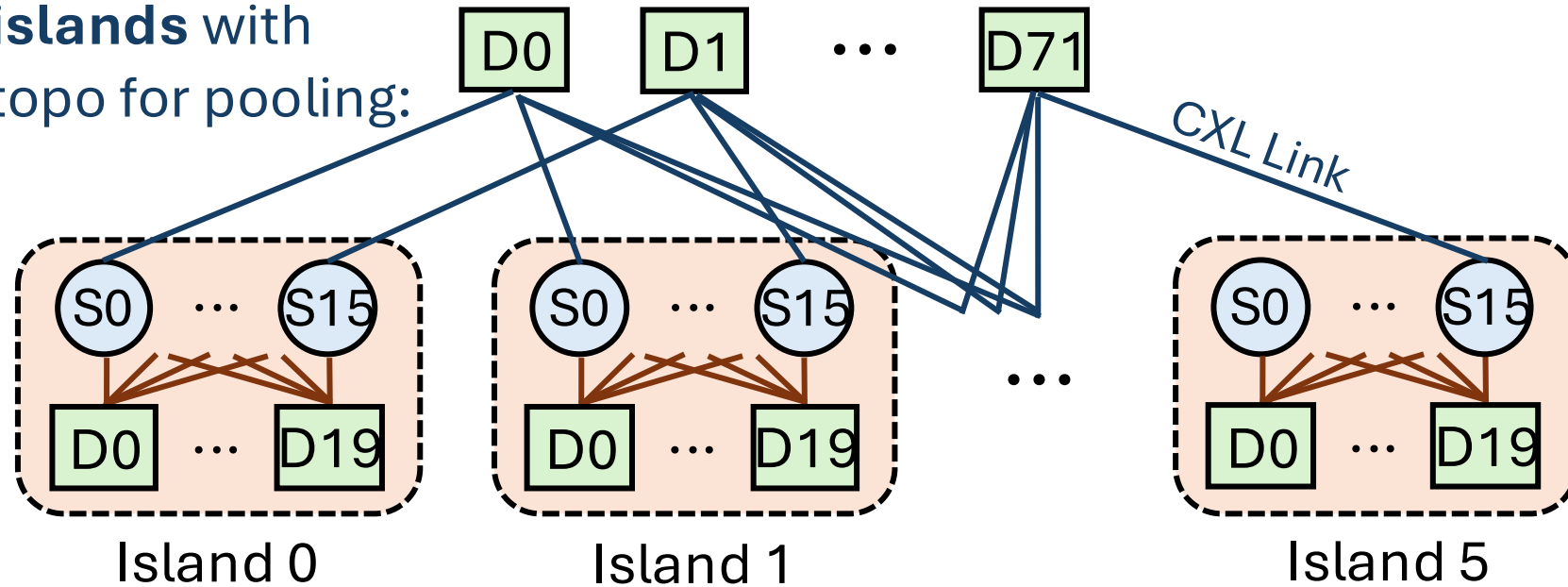
16-Server Island

Ensure pairwise MPD overlap  
**within each island** using BIBD

# Reconcile MPD Overlap With “Islands”









Insight: Low-latency communication is typically needed only at modest scale (e.g.,  $\leq 12$  servers)

**Interconnect islands** with expander-like topo for pooling:

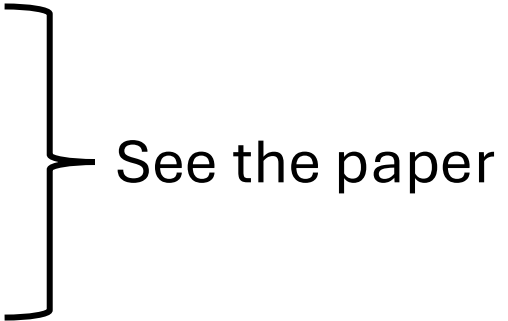


Ensure pairwise MPD overlap  
**within each island** using BIBD

# Octopus Reconciles Pooling and Comm.

	Memory Pooling Saving	Communication Latency
<b>Fully-connected</b> (Pod size: 4 servers)	 <b>Poor (limited pod size)</b>	 <b>Low among 4 servers</b>
<b>BIBD</b> (Pod size: 25 servers)	 <b>Poor (MPD overlap)</b>	 <b>Low among 25 servers</b>
<b>Expander</b> (Pod size: 96 servers)	 <b>Optimal</b>	 <b>High (multi-hop)</b>
<b>Octopus</b> (Pod size: 96 servers)	 <b>Near-Optimal</b>	 <b>Low among 16 servers</b>

# Evaluation

- **Can Octopus achieve significant savings by pooling memory?**
  - Can Octopus provide memory-like communication latency within a pod?
  - How do CXL link failures affect Octopus?
- 
- See the paper

# Octopus Reduces Server CapEx by Pooling

	CXL Switch	Fully-connected (MPD)	<b>Octopus (MPD)</b>
Pool-able Memory	35%	65%	65%
Saving Within Pooled Memory	44%	8%	25%
<b>Total Memory Saving</b>	<b>16%</b>	<b>5%</b>	<b>16%</b>
CXL Device Cost (\$/server)	\$3460	\$1548	\$1548
<b>Server CapEx</b>	<b>+3.3%</b>	<b>+2.5%</b>	<b>-3.0%</b>

# Call for Research on CXL Pod Topologies

CXL expansion already deployed in production, pooling next to come

## **Many research questions in designing CXL pods:**

- Can we borrow ideas from datacenter network topologies?
- Can we design sparse topologies for CXL switches?
- How to design CXL pods for different use cases (e.g., ring topologies for collective communications)?
- CXL 3.0 supports multi-level switching, how to design CXL pods?
- For smart/computational CXL devices, how to design CXL pods?

# Executive Summary: CXL as “Pod Area” Network

## Background:

- Memory is expensive (~50% of CapEx) but underutilized due to stranding
- CXL enables a “pod” of servers to reach a common **memory pool**
- CXL also enables intra-pod **communication at memory-like latency**
- However, existing CXL pods are either **expensive** (using CXL switches) or **small** (using multi-port CXL devices, each with  $\leq 4$  ports)

## Contributions:

- Octopus: a low-cost, large-scale CXL pod which **sparsely** connects servers with low-port CXL devices, without using costly switches
  - To balance conflicting demands of memory pooling and intra-pod communication, Octopus partitions a pod into “**islands**”
- Results: 16% memory saving,  $< 1 \mu\text{s}$  intra-pod communication

Paper

