

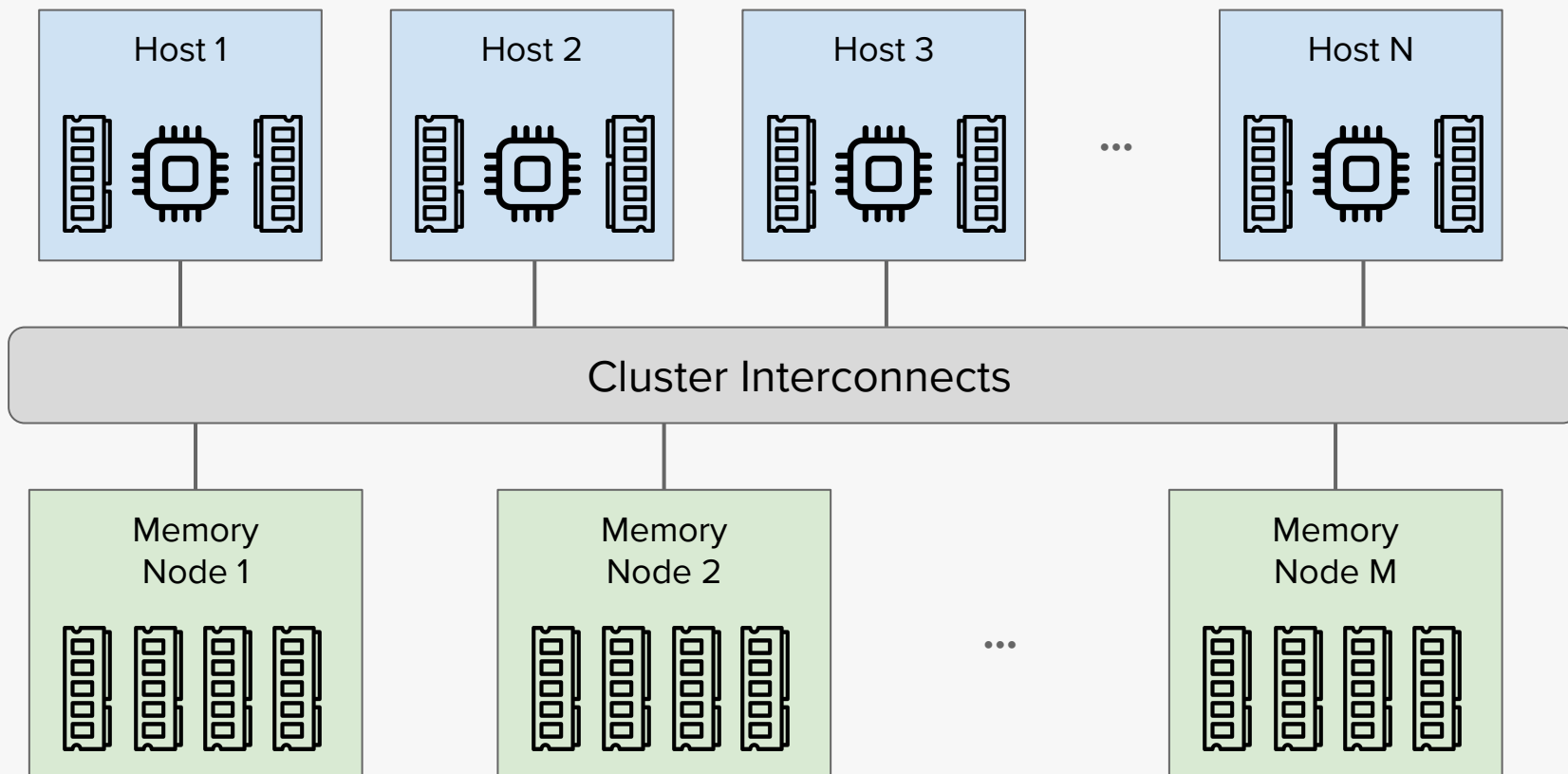
Building A CSEQ-Inspired Transport for Switched CXL Memory Pooling

Zerui Guo, Emily Shriver, Ming Liu

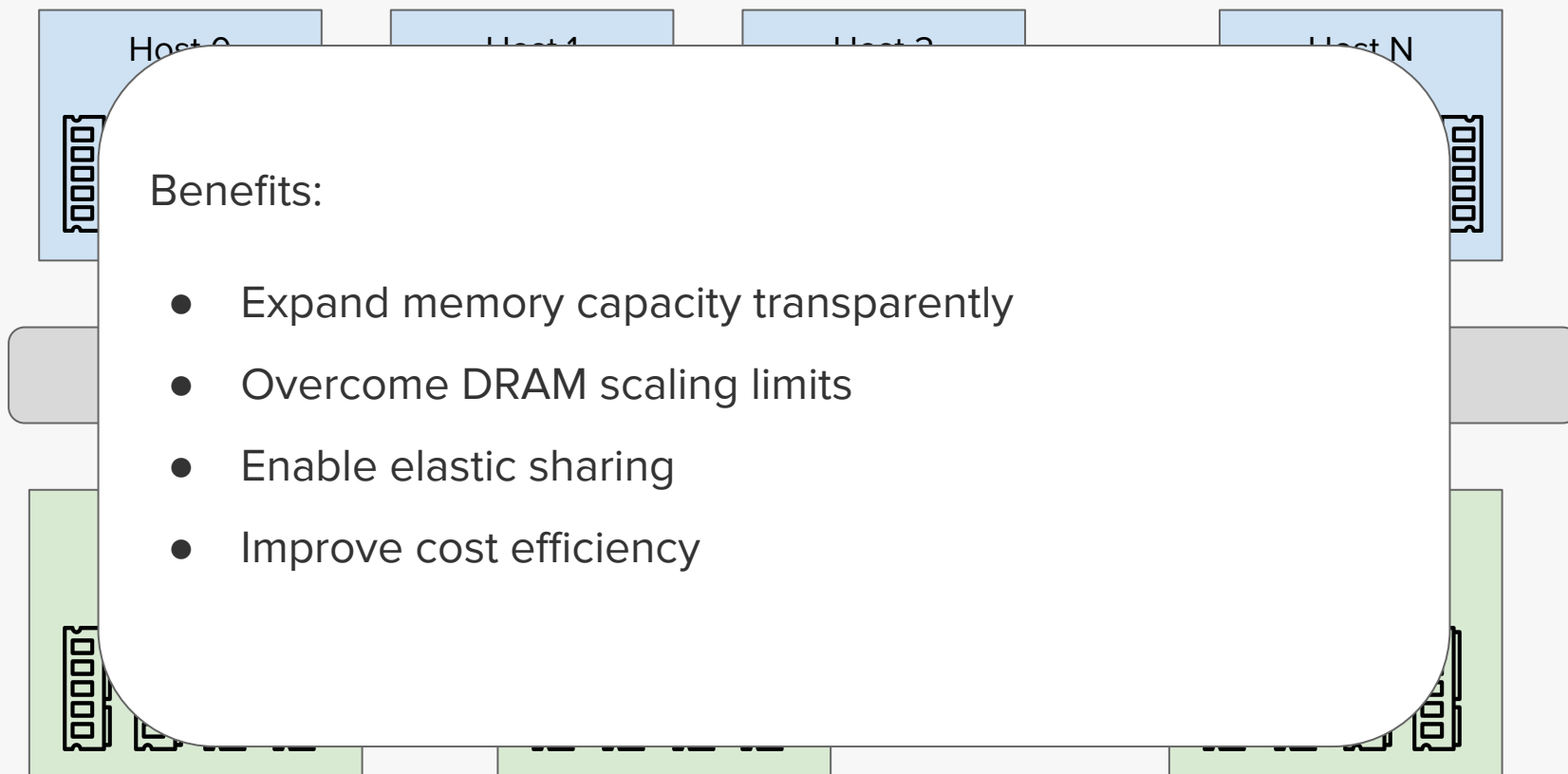
NetLab@UW-Madison



Memory Pooling

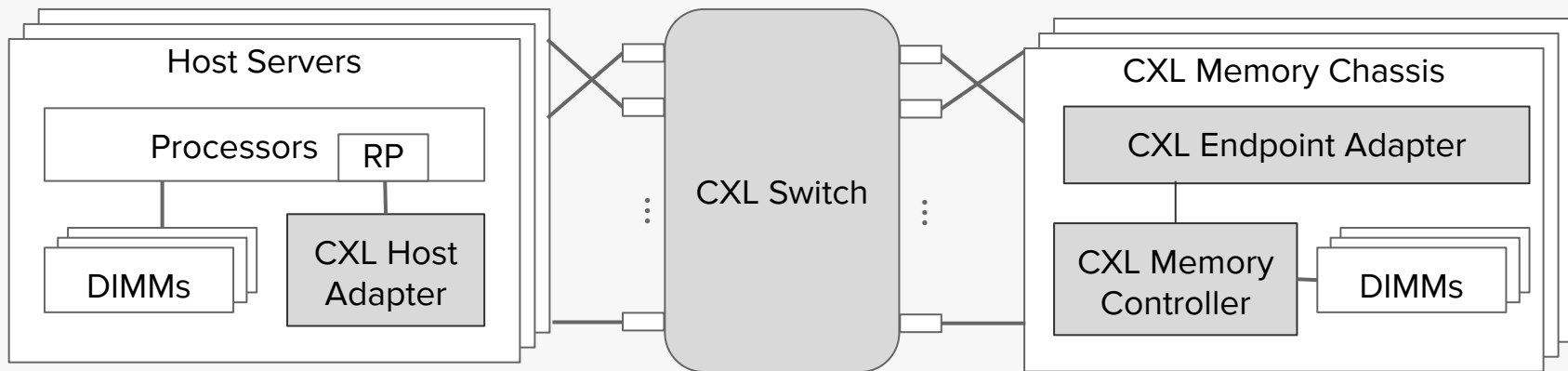


Memory Pooling



Compute Express Link (CXL)

- An emerging cluster interconnect offering load/store access model
- Enable rack-scale composable memory pooling



Compute Express Link (CXL)

- XConn Titan Evaluation Platform
 - An XConn B2 Apollo CXL switch
 - 14 upstream and downstream ports
 - ASIC-based host and endpoint adapters
 - PCIE Gen 5 x16
 - A memory chassis
 - MCIO backplane
 - Up to 12 CXL DIMMs

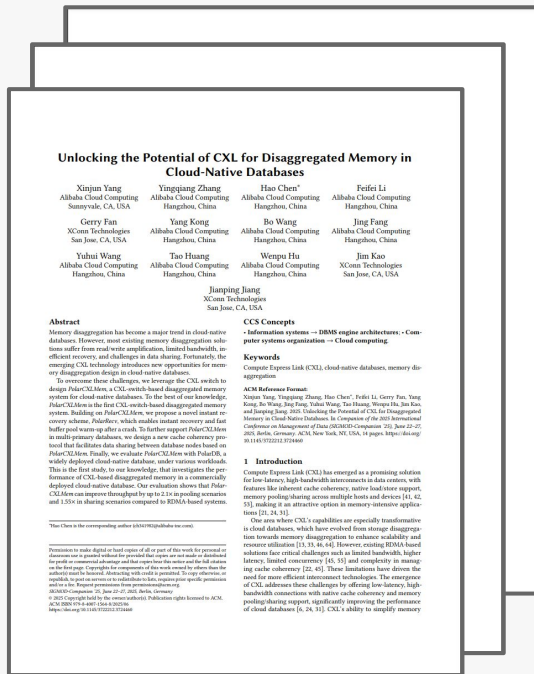


CXL-based Memory Pooling Being Deployed

Cloud Service Providers:



Alibaba Cloud



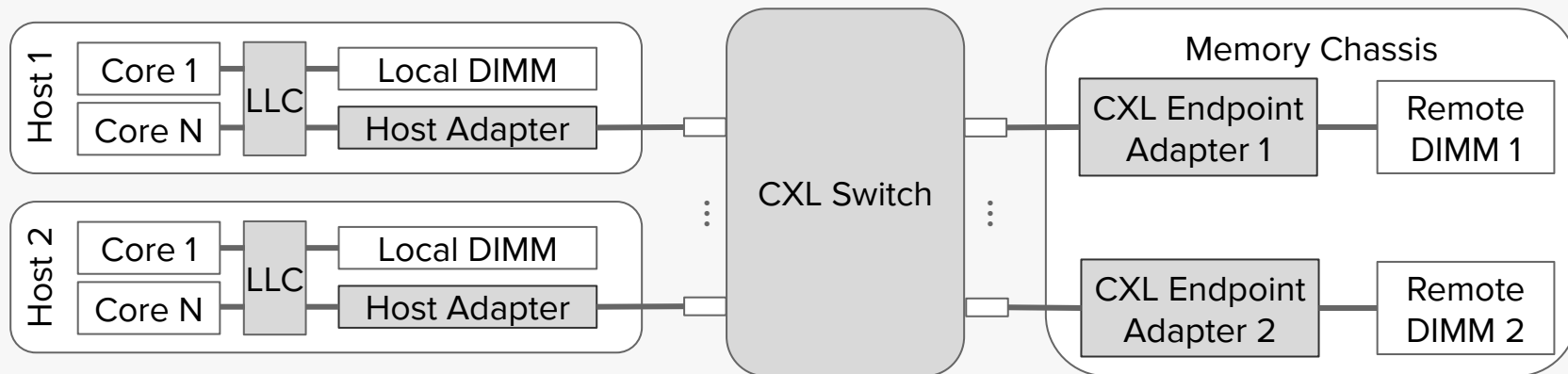
Outline

- **Characterize Switched CXL Memory Pooling**
- MemChannel Design and Implementation
- MemChannel Evaluation
- Conclusion

Experimental Methodology

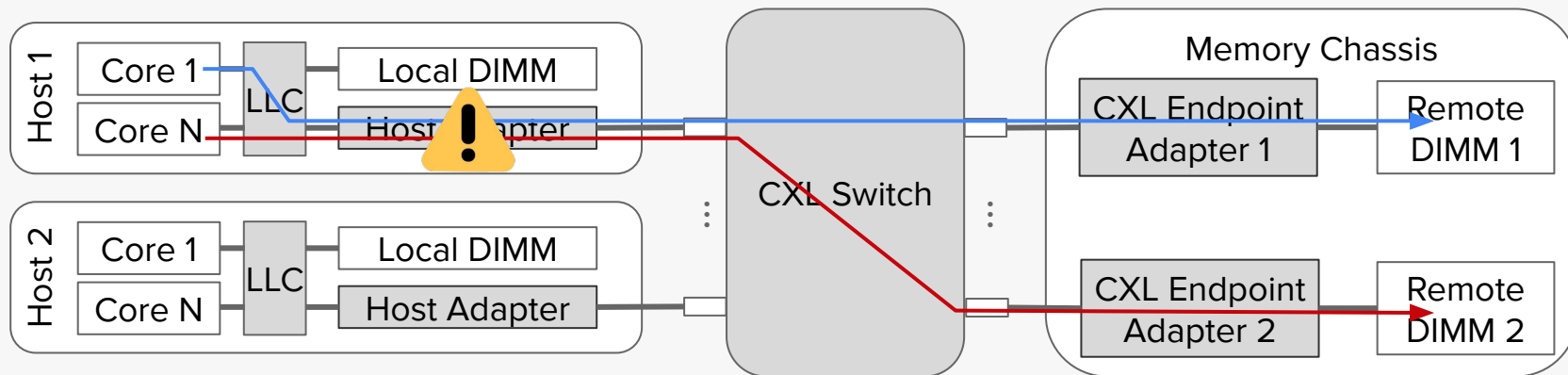
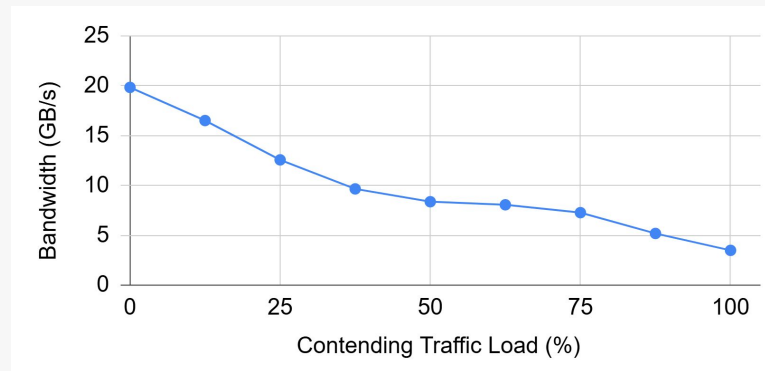
- **Benchmarking framework**

- Launch memory streams to local/CXL DIMMs
- Configurable access patterns (e.g., stride width, temporal/non-temporal store, random/sequential)



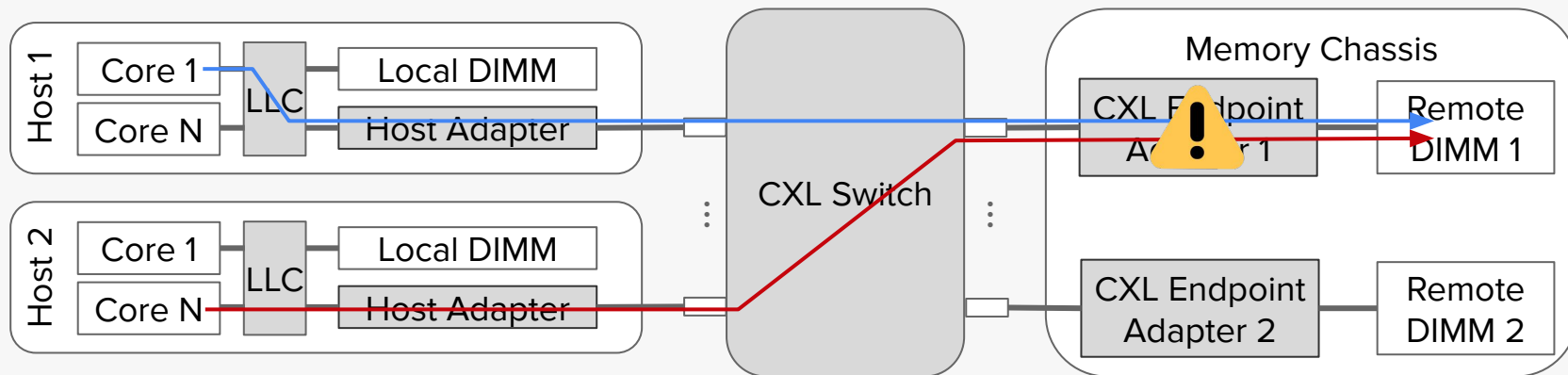
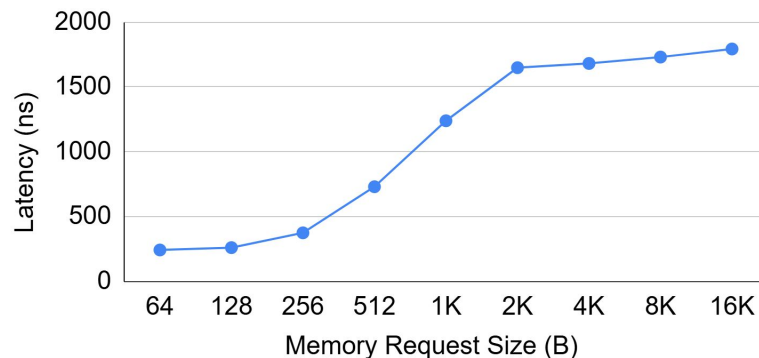
Issue #1: Intra-host Contention

- Host adapter congestion
 - Bandwidth dropped by 82.6%
 - CXL FLITs scheduled in a best-effort fashion



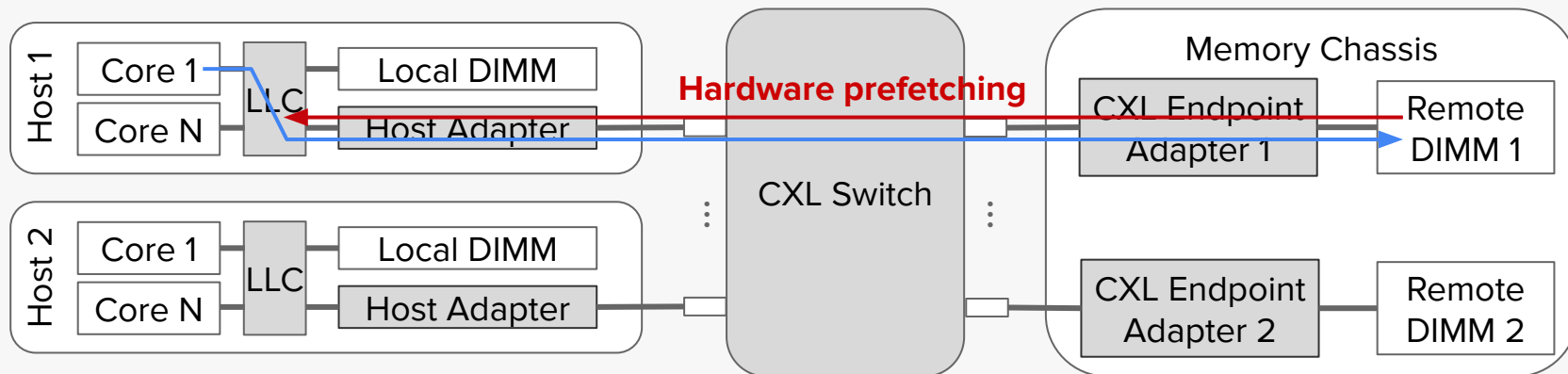
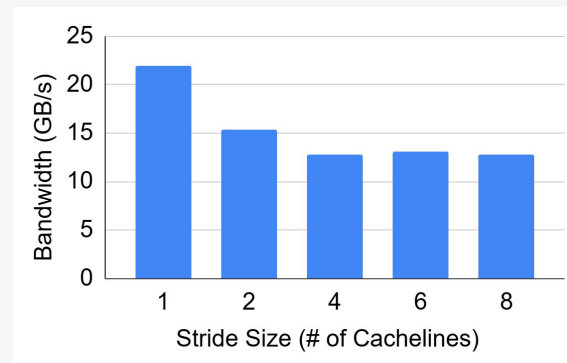
Issue #2: In-fabric Congestion

- Interleaved hetero. memory requests
 - Cacheline-size accesses: 6.9x slowdown
 - Head-of-line (HoL) blocking



Issue #3: Unmanaged Host-rDIMM Interaction

- **Unpredicted bandwidth oversubscription**
 - Throughput dropped by 41.6%
 - Host prefetching generates more traffic than expected.



Root Cause: Unmanaged data streams in switched CXL memory pooling cause severe performance interference.

Outline

- Characterize Switched CXL Memory Pooling
- **MemChannel Design and Implementation**
- MemChannel Evaluation
- Conclusion

MemChannel Overview

- **Overall design:**
 - Construct a data pipe between a host core and a CXL DIMM
 - Mitigate interference from contending memory streams
 - Regulate remote memory access based on the CXL resources and demands
- **Design goals:**
 - High utilization
 - Efficient multi-tenancy
 - Low overheads

Challenges

- **C1: Hardware non-programmability and opaqueness**
 - Impractical to implement in-network transport layer
- **C2: Implicit transmission**
 - Non-trivial to track individual CXL load/store transactions
- **C3: Fine-grained cross-layer traffic monitoring**
 - Difficult to pinpoint the congested component

Proposed Techniques

Challenges

C1: Hardware non-programmability

C2: Implicit transmission

C3: Fine-grained traffic monitoring

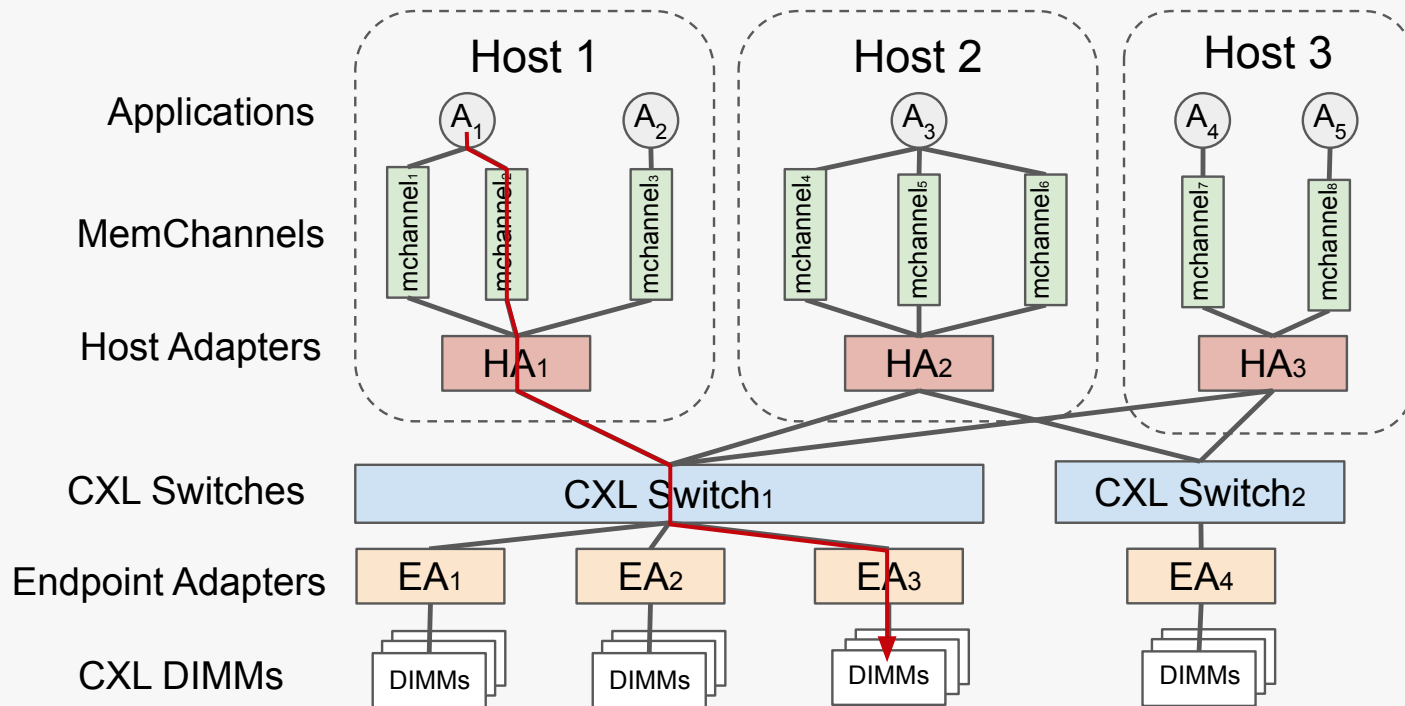
Techniques

T1: Sender-driven admission control

T2: CSFQ-inspired transport stack

T3: Runtime hardware capacity probing

System Model

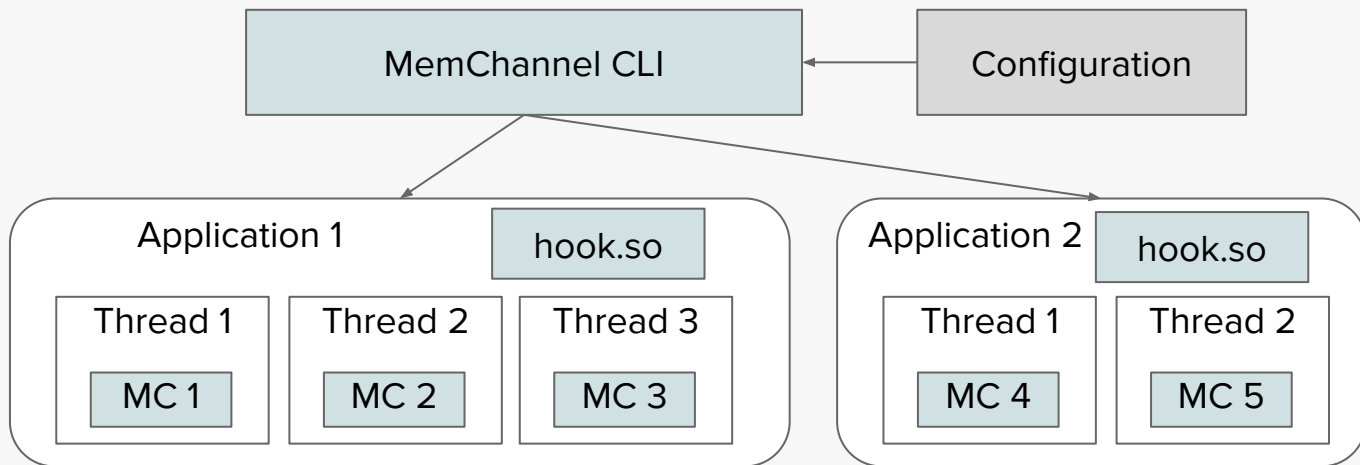


System Model

Interfaces API	Description
mchannel_init	Initialize mchannels from user configs, identify CXL data paths, and allocate shared metadata
mchannel_open	Allocate thread resources, bind CXL memory, and arm timers for mchannel_sched
mchannel_close	Close the thread's mchannel to the remote DIMM, and free allocated resources
mchannel_sched	Enforce bandwidth limits by monitoring CXL requests and controlling thread execution

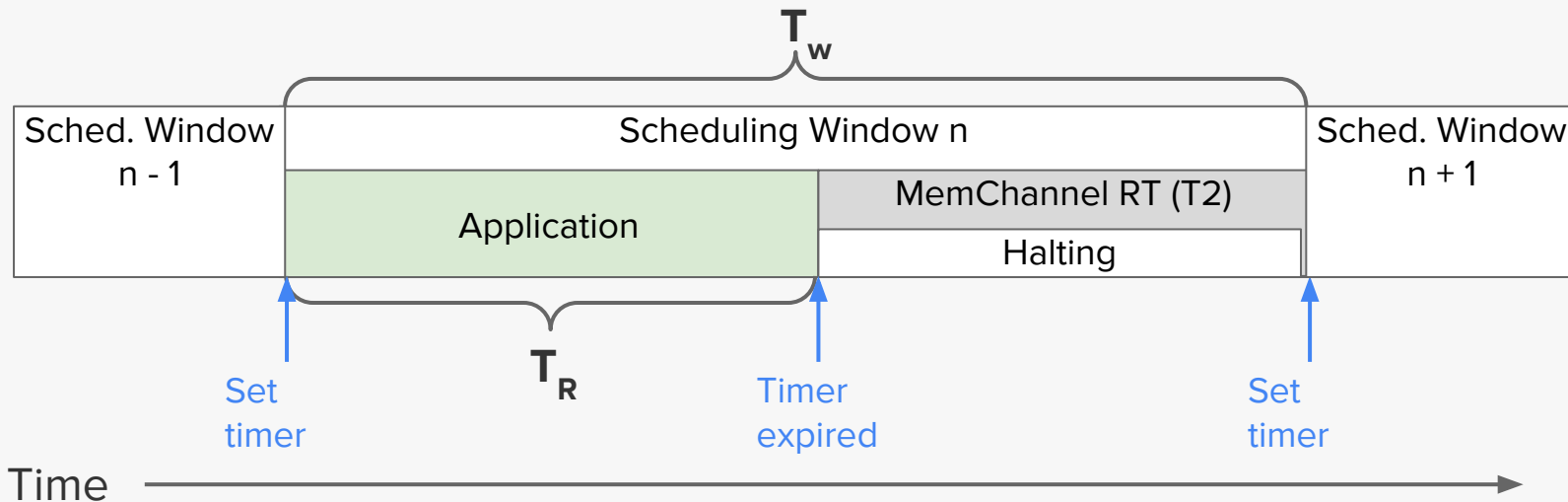
T1: Sender-driven Admission Control

- Support unmodified applications
 - Launch applications as subprocesses based on the configurations
 - Hook into applications (LD_PRELOAD) to create *mchannels*



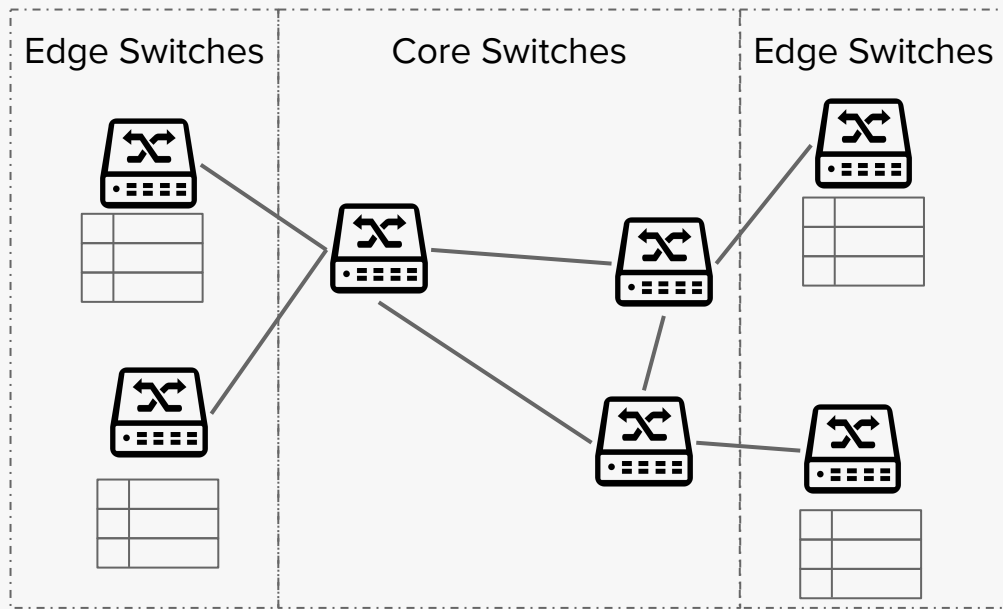
T1: Sender-driven Admission Control

- Sender-driven time-based rate control
 - Adjust application running time based on the transport algorithm
 - Arm a POSIX timer



Core-Stateless Fair Queueing Background

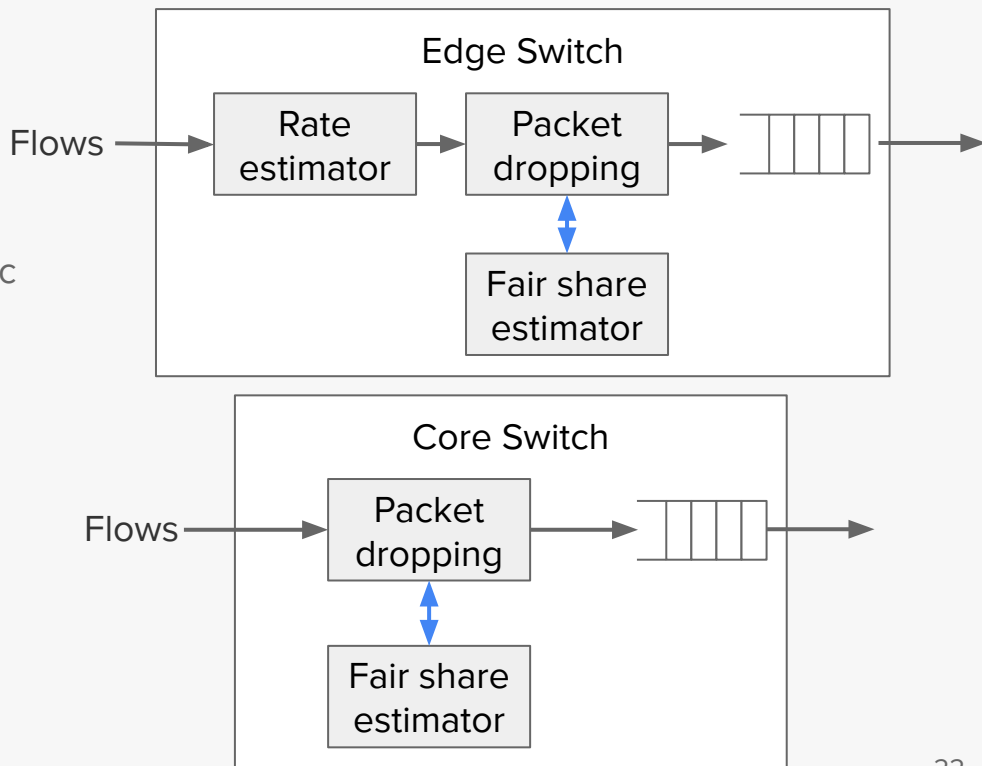
- Edge switches
 - Maintain per-flow state
 - Enforce max-min fairness
- Core switches
 - Operate on aggregated data
 - Enforce max-min fairness



Core Stateless Fair Queueing

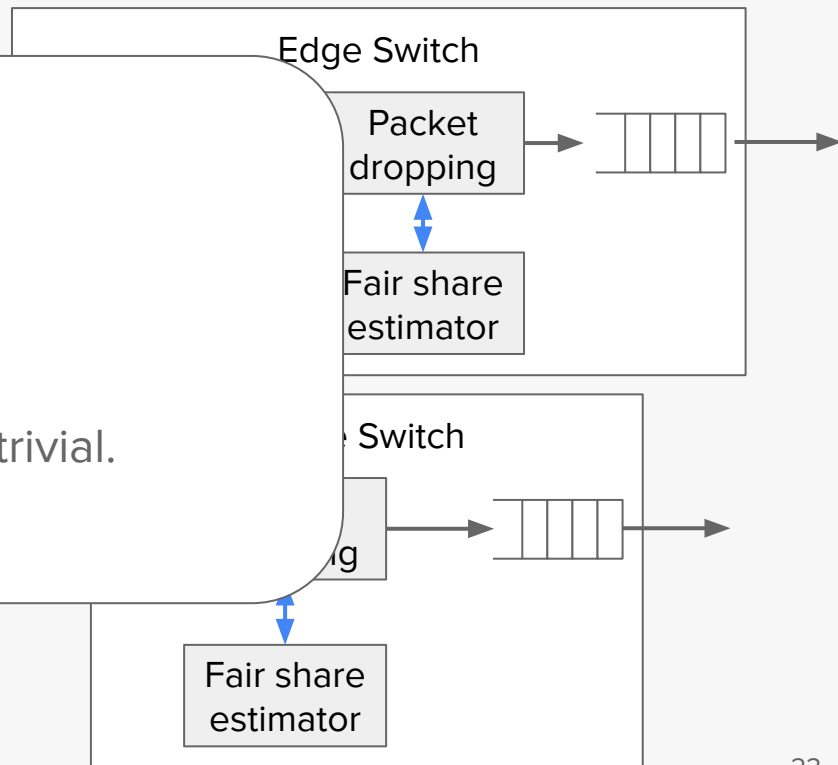
CSEQ Techniques

- Rate estimation
 - Compute individual host flow rates
- Fair share estimation
 - Calculate fair-share based on total traffic vs. capacity.
- Traffic control
 - Regulate flow via probabilistic dropping



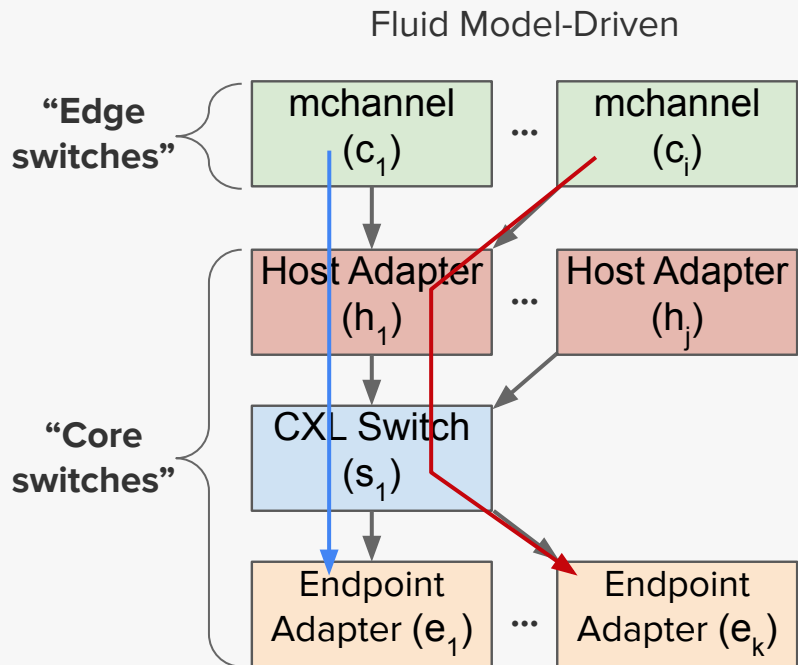
CSFQ Techniques

- T1: rate estimator
 -
- T2
 - Why CSFQ for CXL?
 - Scalable
 - Edge-driven
 - Lightweight
- T3
 -
 - However, applying CSFQ is non-trivial.



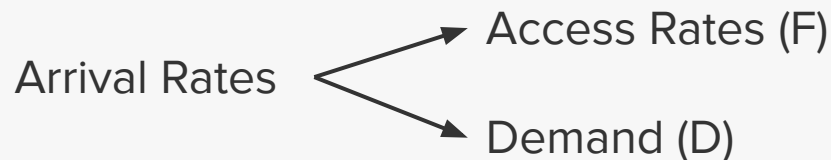
T2: CSEQ-inspired Transport Stack

- Model mchannels as edge switches
- Model CXL devices as core switches
 - Host/Endpoint adapters, CXL switches
- Flows in our context are:
 - Directional memory streams from applications
 - Start from mchannels and reach CXL DIMMs



T2: CSEQ-inspired Transport Stack

- Rate estimation

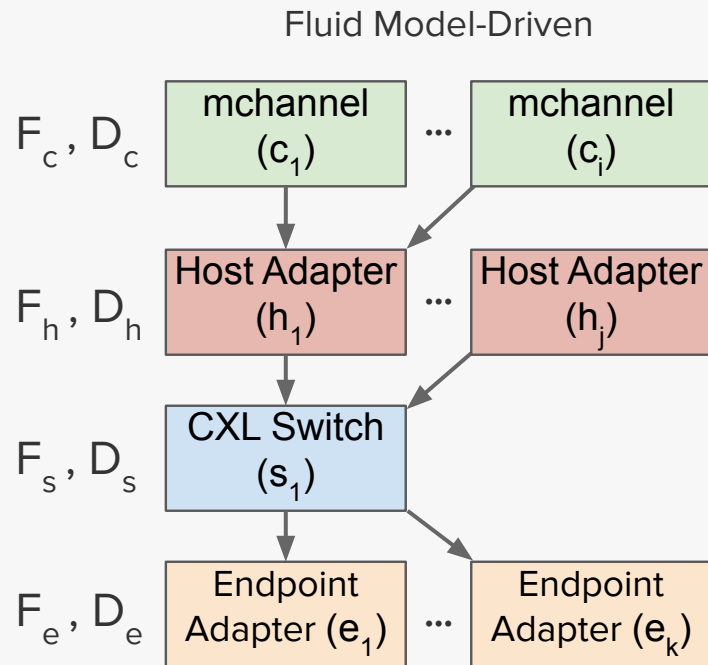


- CXL bandwidth monitoring

- Read off-core response (OCR) event counters
- Divide by window size (T_w) and running time (T_R)

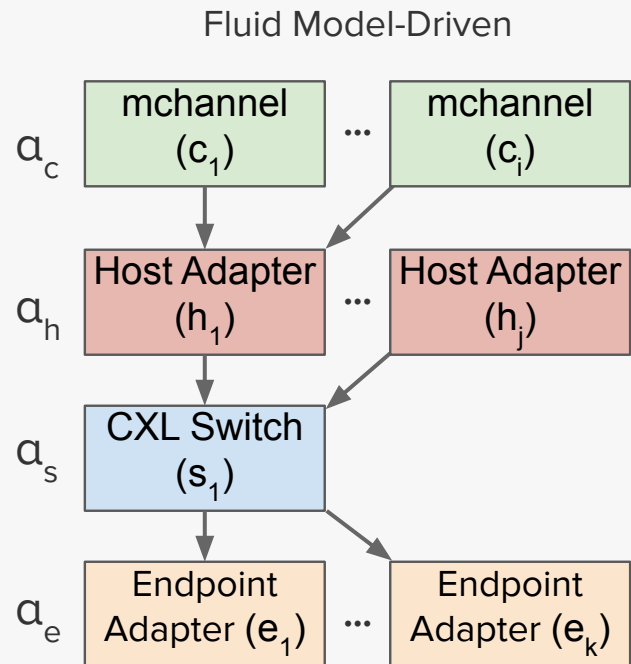
- Calculate aggregated rates directly

- Access rate
- Demand



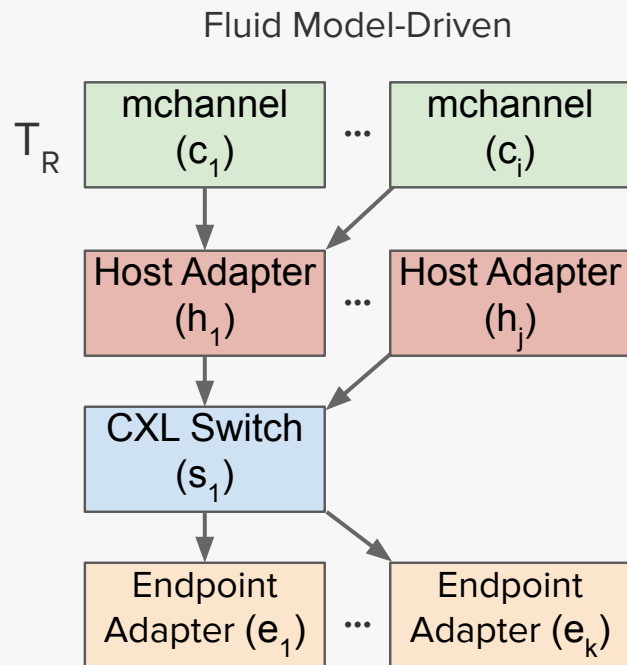
T2: CSEQ-inspired Transport Stack

- Fair share estimation
- Calculate hardware fair share on hosts
- Estimate fair share rate (α) iteratively
 - Aggregated access rate (F)
 - Aggregated demand (D)
 - Hardware capacity (C)
- *mchannel* fair share
 - Minimum fair share along the flow path



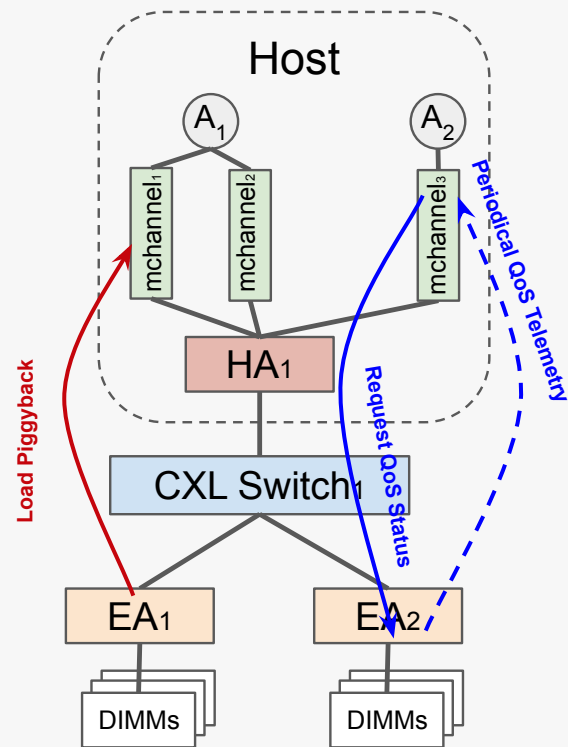
T2: CSEQ-inspired Transport Stack

- Traffic control
- Adjust application running time (T_R)
 - Fair-share-to-demand ratio
 - Allocate running time proportionally



T3: Runtime Hardware Capacity Probing

- Congestion Signals
 - Device internal load
 - Piggyback in memory responses
 - QoS telemetry
 - Report periodically
- Translate congestion signals into loading factors
- Load-based congestion control
 - Congestion free: additive increase
 - Congestion avoidance: multiplicative decrease



Outline

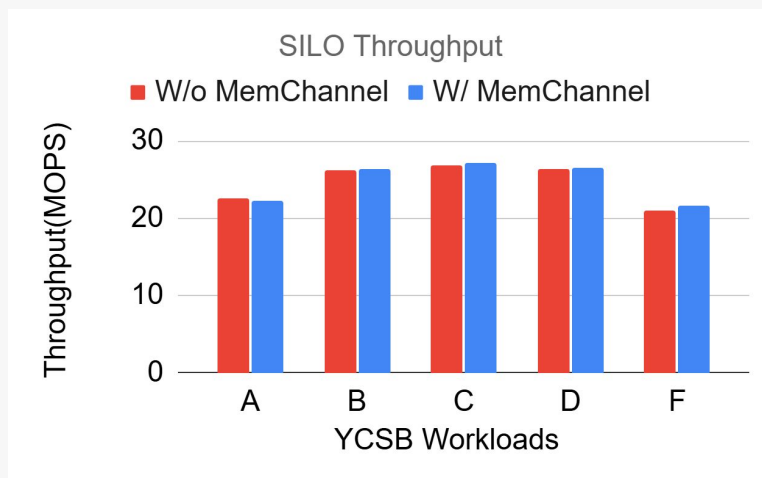
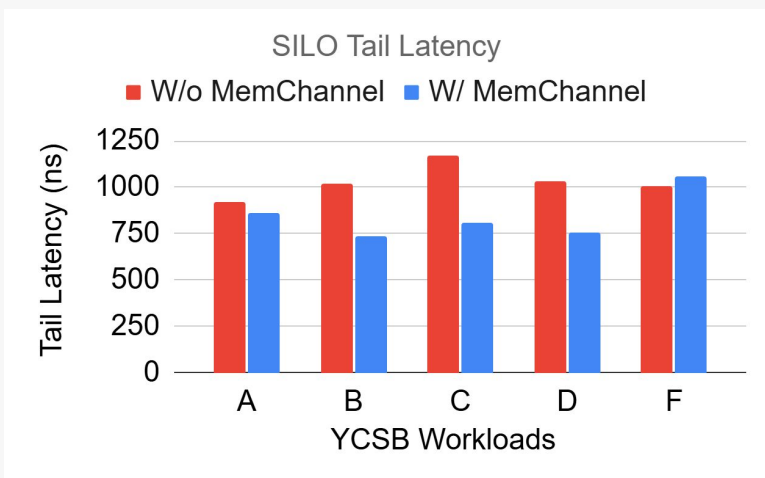
- Characterize Switched CXL Memory Pooling
- MemChannel Design and Implementation
- **MemChannel Evaluation**
- Conclusion

Experimental Setup

- Applications
 - MICA, SILO, GAP benchmark suite, SPEC CPU 2017
- Performance metrics
 - Application throughput
 - Average/tail latency
 - CXL memory bandwidth
- Code
 - <https://github.com/netlab-wisconsin/MemChannel>

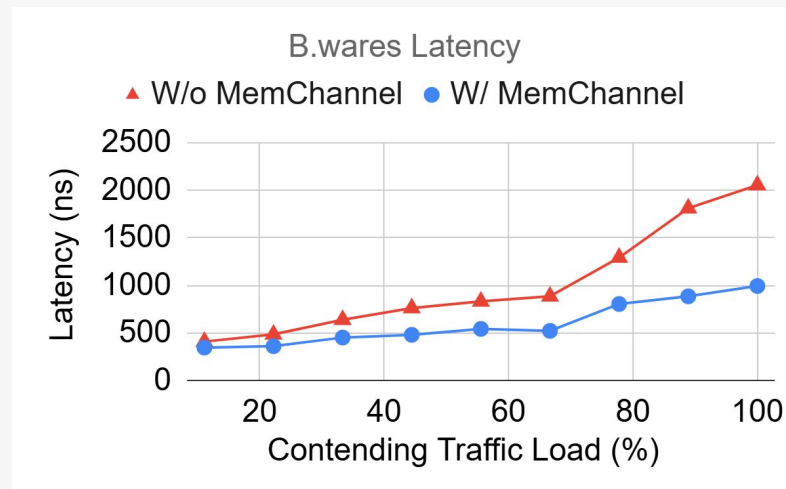
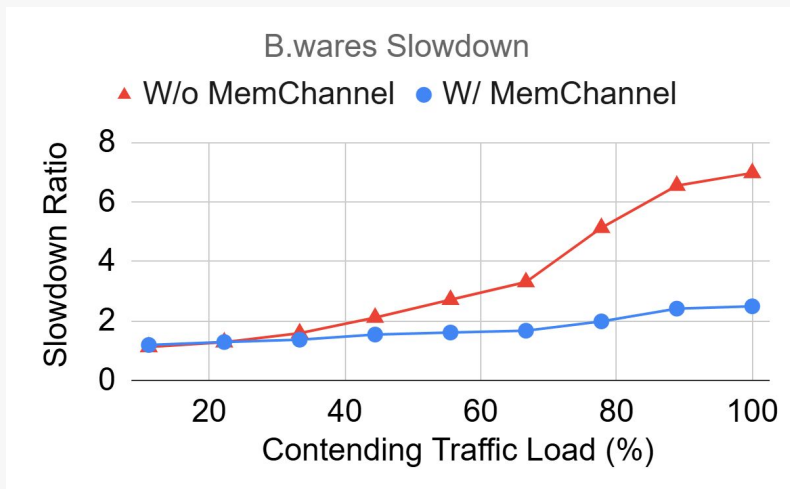
Application Performance

- Evaluate SILO performance under five YCSB workloads
- Reduce tail latency by 17.4% with similar throughput



Performance Isolation

- Evaluate B.wares(CPU SPEC) performance under contending memory flows
- 4.5x throughput increase and 51.6% latency reduction



More Evaluations

- Inter-Host Performance Isolation
- Fairness
- System Overheads
- MemChannel Adaptiveness
-

Conclusion

- Unmanaged data streams causing severe performance interference
- MemChannel: a transport layer for switched CXL memory pooling
- Data pipes between cores and CXL DIMMs with controlled CXL requests
 - Sender-driven admission control
 - CSFQ-inspired transport stack
 - Runtime hardware capacity probing

<https://github.com/netlab-wisconsin/MemChannel>

Thanks!