

CXL Memory Use Cases: Insights into Expansion and Pooling

PRESENTER

Minseong Kim



TABLE OF CONTENTS

Why do we need CXL Memory?

- Memory Capacity Requirement
- Memory Capacity Gap in AI applications & GPUs

CXL Usage Model Tree

- Single Server - Bandwidth Expansion and Capacity Expansion
- Multiple Servers - Pooling & Sharing

CXL Memory Usage Model Survey

- Bandwidth Expansion / Capacity Expansion / Tiering
- Memory Pooling & Sharing

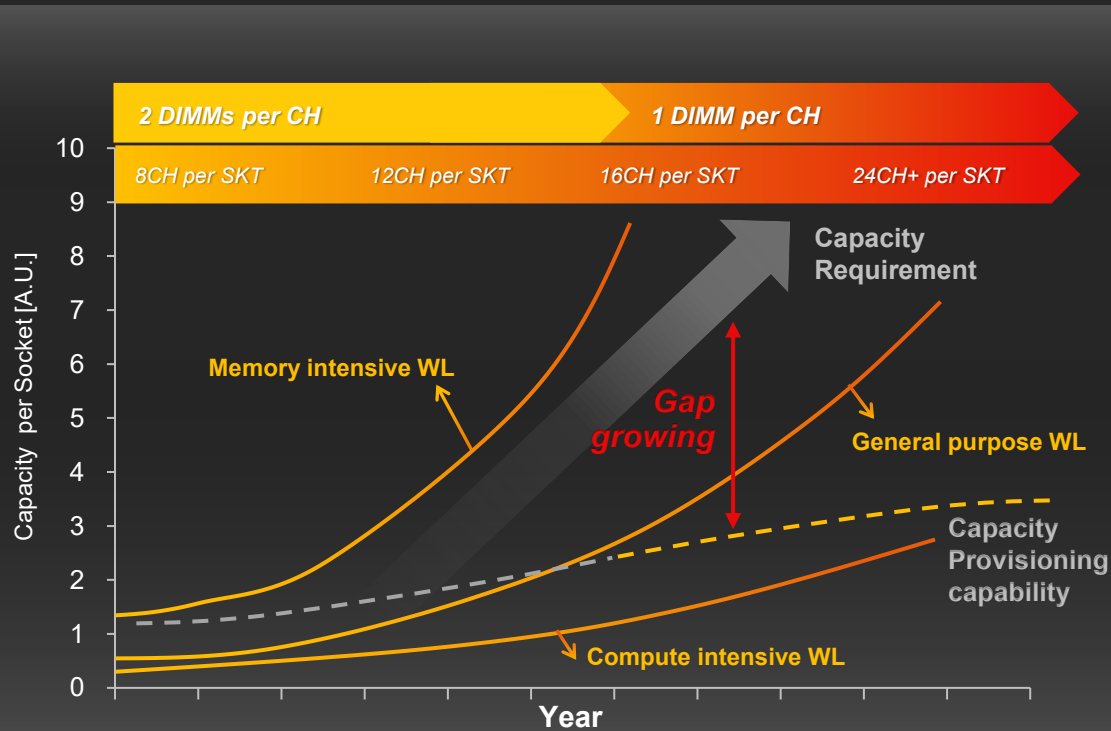
Experimental Results

- Bandwidth Expansion Case – Single Server with LLM Inference
- Bandwidth Expansion Case – Single Server with HPC Workloads
- Capacity Expansion Case – Single Server with Redis IMDB
- Memory Pooling Case – Multiple Servers with Redis IMDB

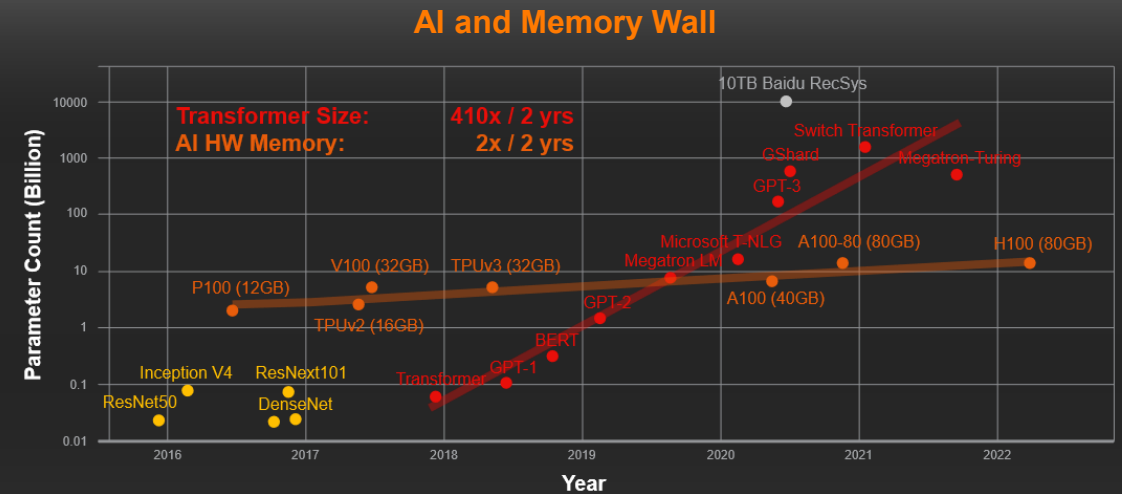
Why do we need CXL Memory?

Memory wall exists relative to CPU Core Count, and memory wall also exists in GPU-based AI Memory System
Need to overcome the difficulty of high-capacity scaling with pooling/switching-based scale-up/scale-out

Memory Capacity Requirement (CPU)



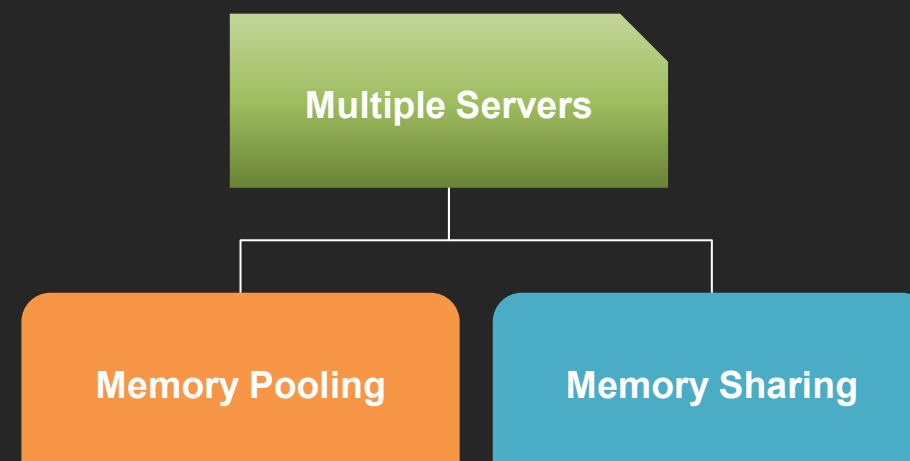
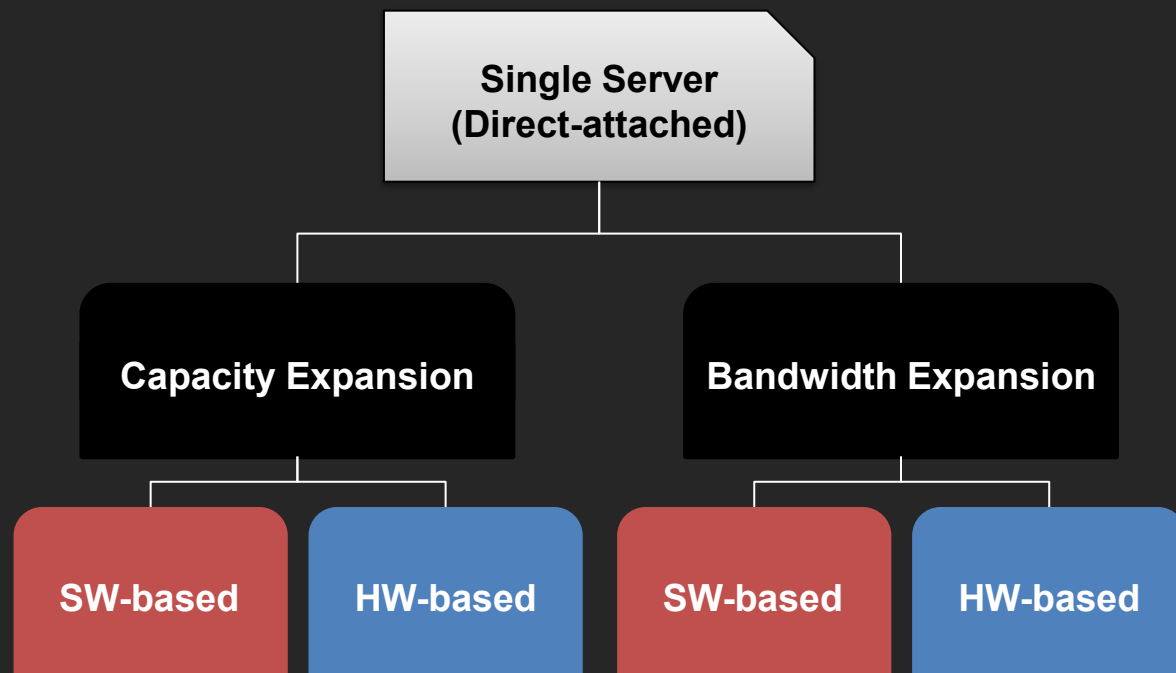
Memory Capacity Gap in AI applications & GPUs



Application fragmentation makes it difficult to know sweet spot capacity, but TB to 10TB+ of high capacity memory is needed within TCO budget








Reference: AI and Memory Wall, A. Gholami et al, IEEE Micro

CXL Memory Usage Model Tree



✓ We explore primary approaches to utilizing the CXL Memory Module.

CXL Memory Usage Model Survey

Category	Capacity Expansion	Bandwidth Expansion	Tiered Memory	Pooling/Sharing	Notes
	In-memory Database	LLM Inference HPC	In-memory Database, LLM Inference, HPC	In-memory DB In-memory Analytics	
	DLRM inference/training ¹⁾	DLRM inference/training ¹⁾		In-memory DB ²⁾ (TPC-DS on SAP HANA)	IMDG Databases & Caches AI/ML Workloads Financial Services
 ³⁾ 	In-memory DB (MSSQL + TPC-H)	HPC (CloverLeaf)	Apache Spark Based Machine Learning (SVM) (Big Data Workload)		
 ⁴⁾	1) PostgreSQL + TPC-H 2) RocksDB + db_bench 3) RAG Pipeline				
			Weaviate Vector DB on gist dataset (ANNS) ⁵⁾		
	Azure Service (PoC) Various Benchmark			Azure Service (PoC) Various Benchmark	Pond (ASPLOS 23) Octopus (2025)

✓ *We focus performance analysis of IMDB and LLM Inference across various scenarios.*

System Configuration

Single Server

CPU	Intel 6 th Generation Xeon Scalable Server Processor (Granite Rapids, GNR)
Board	Intel CRB (Customer Reference Board)
DRAM Only	DDR5 6400Mbps 128GB * 2Ch
DRAM+CXL	DDR5 6400Mbps 128GB * 2Ch CMM-DDR5 6400Mbps 128GB * 1Ch
Bandwidth Expansion Case	Llama.cpp + Llama 3.1 (70B Q8) / HMSDK v2.0 SPEC CPU 2017 / HMSDK v1.1 (numa interleaving)
Capacity Expansion Case	Redis + YCSB / Local Preferred (Linux)

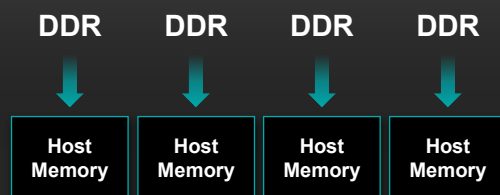
Multiple Servers

CPU	5 th Generation AMD EPYC™ Processors (Turin)
Board	AMD CRB (Customer Reference Board)
DRAM Only	DDR5 6400Mbps 64GB * 1Ch
Memory Pooling Case	CMM-DDR5 6400Mbps 128GB * 22Ch
SW Configuration	Redis + YCSB

Bandwidth Expansion Case – Single Server (1/2)

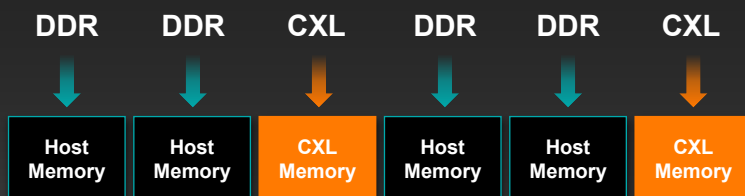
We ask a question to AI model, CMM+DDR shows 30% better performance of system.
Next CMM has 2x better BW compared to AI model result of current product.

<Native DDR Only>



Utilizing Host DRAM 2-channel only

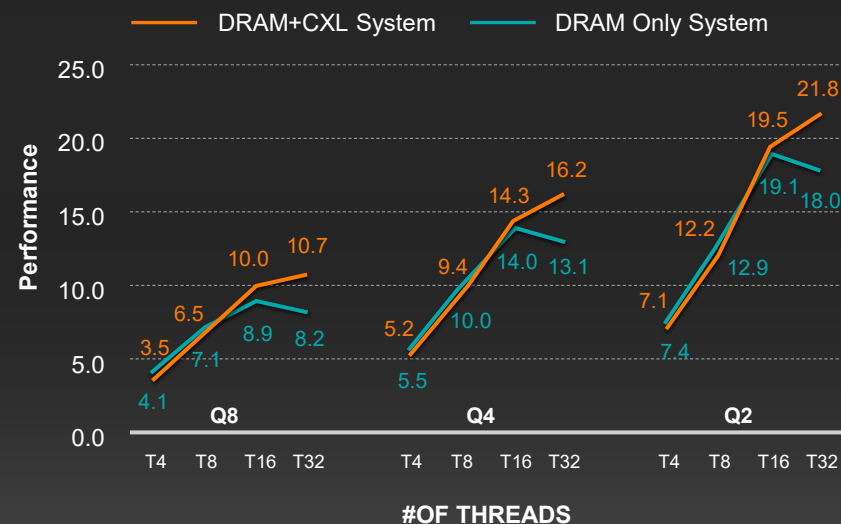
<Native DDR+CXL>



Utilizing Host DRAM 2-channel + CXL 1-channel

Enhanced System Performance

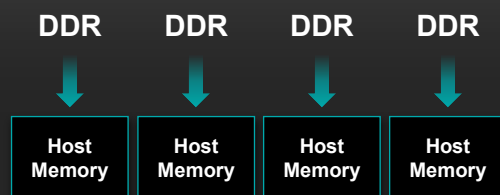
Throughput Comparison of AI Inference Workload



Bandwidth Expansion Case – Single Server (2/2)

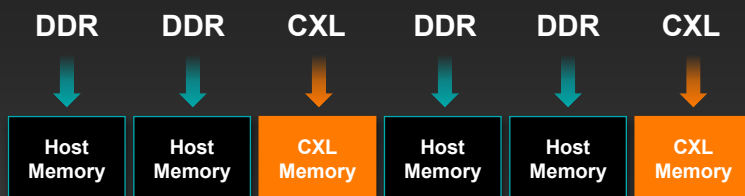
CMM+DDR unlocks up to 33% higher performance on SPEC CPU 2017 memory sensitive items.
This improvement is enabled by the higher memory bandwidth of CMM+DDR system compared to DRAM Only system.

<Native DDR Only>



Utilizing Host DRAM 2-channel only

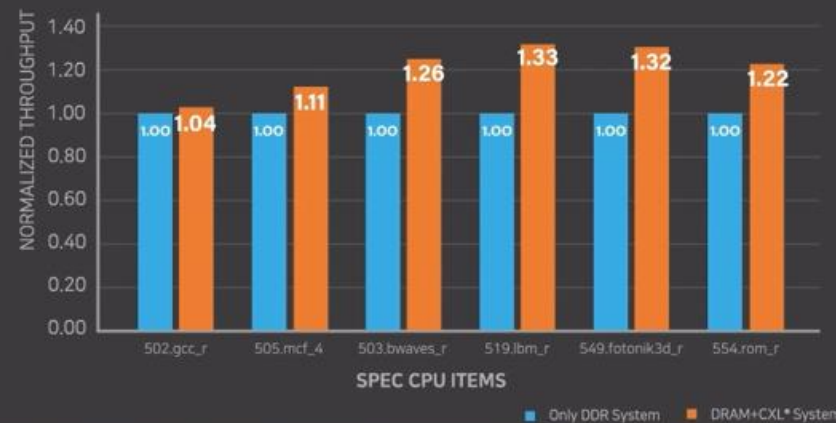
<Native DDR+CXL>



Utilizing Host DRAM 2-channel + CXL 1-channel

Enhanced System Performance

Performance Comparison of SPEC CPU Workload

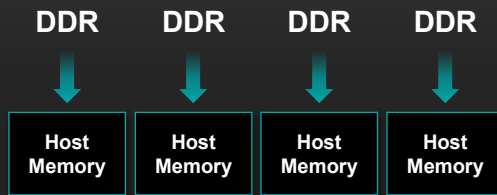


Capacity Expansion Case – Single Server

We use the YCSB benchmark for the Redis in-memory database.

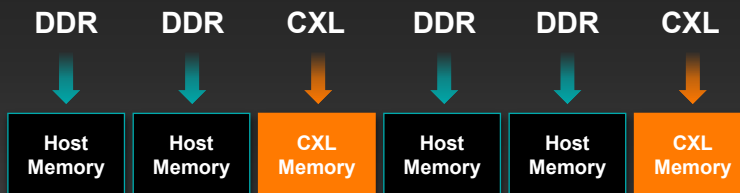
CMM allows us to load a database up to the additional CMM capacity on a single server without throughput drop.

<Native DDR Only>



Utilizing Host DRAM 2-channel only

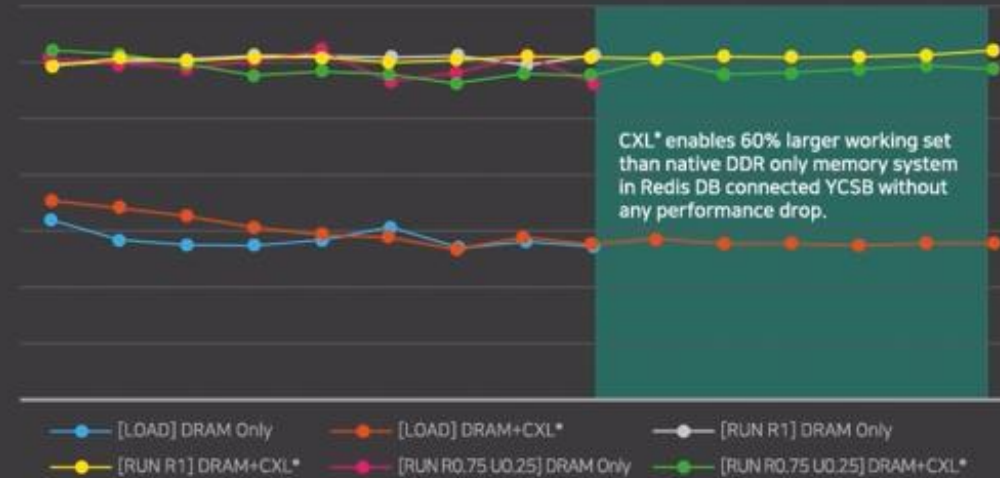
<Native DDR+CXL>



Utilizing Host DRAM 2-channel + CXL 1-channel

No Significant Performance Drop + Load More Database

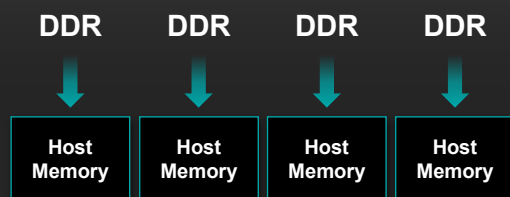
Capacity Expansion Case (Redis + YCSB)



Memory Pooling Case – Multiple Servers

We evaluated the performance of the Redis using the YCSB benchmark with CXL Memory Pooling. We found no significant performance difference. This confirms that IMDB is a suitable use case for CXL memory pooling.

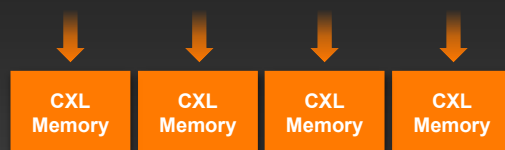
<Native DDR Only>



Utilizing Host DRAM 1-channel only

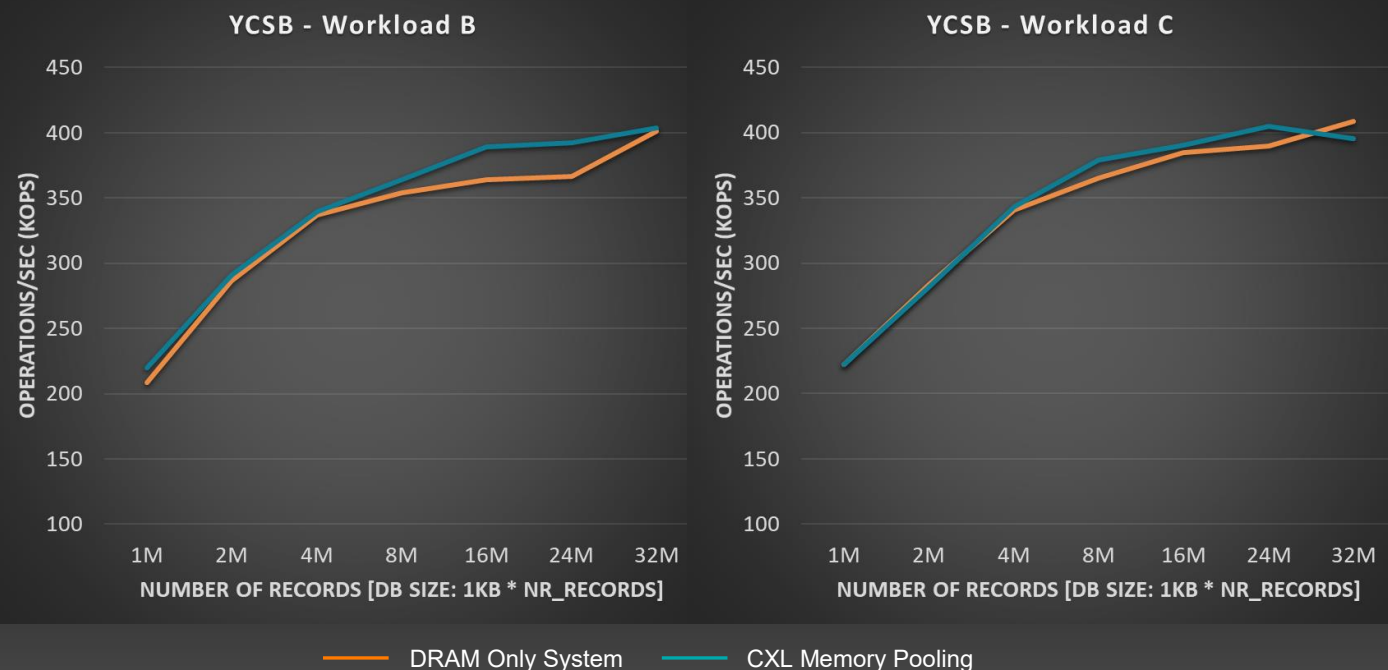
<CXL Memory Pooling>

CXL Pooling



Utilizing CXL Memory Pooling

No Significant Performance Drop in CXL Pooling System



- CXL memory allows systems to scale more effectively by solving memory bandwidth and capacity bottlenecks.
- AI and HPC workloads show clear performance improvement when CXL is used.
- CXL server can handle a much larger database, which reduces the total number of servers.
- Memory pooling with CXL showed good and reasonable results with IMDB (redis case)
- We plan to explore more use cases in the future.

Booth #207

Meet the future of memory.
Just steps from the entrance.

Innovation starts here,
Literally.

SK hynix



Questions?