SNIA Compute, Memory, and Storage

Unlocking CXL's Potential: Revolutionizing Server Memory and Performance

Live Webinar April 2, 2025 10:00 am PT/ 1:00 pm ET

Our Speakers







Arthur Sainio

Co-Chair SNIA Persistent **Memory Special** Interest Group

Jim Handy **General Director** Objective Analysis

Mahesh Natu

Systems and Software Working Group Co-Chair

CXL Consortium



Torry Steed

Senior Product Marketing Manager SMART Modular Technologies



2 | ©SNIA. All Rights Reserved.

Agenda

- CXL is Exciting, but Where Is It Headed?
- CXL Consortium Update
- CXL Benefits, Implementation, and Ecosystem
- Final Thoughts and Q&A



The SNIA Community





4 | ©SNIA. All Rights Reserved.

SNIA Score Compute, Memory, and Storage

What We Do

- Engage technology users
- Educate on compute, memory, and storage technologies
- Accelerate SNIA standards
- Propel technology adoption

How We Do It

- Demonstrate at industry events
- Host Programming Workshops and Hackathons
- Present educational webinars, podcasts, and blogs

Learn more at www.snia.org/groups/cms



SNIA Legal Notice

- The material contained in this presentation is copyrighted by the SNIA unless otherwise noted.
- Member companies and individual members may use this material in presentations and literature under the following conditions:
 - Any slide or slides used must be reproduced in their entirety without modification
 - The SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of the SNIA.
- Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be, construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.
- The information presented herein represents the author's personal opinion and current understanding of the relevant issues involved. The author, the presenter, and the SNIA do not assume any responsibility or liability for damages arising out of any reliance on or use of this information.

NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.





CXL is Exciting, but Where Is It Headed?

Jim Handy



CXL Can Do So Many Things, But What is it Really For?

- Maintaining coherency?
- Eliminating stranded memory?
- Expanding memory size?
- Increasing memory bandwidth?
- Supporting persistent memory?
- Hiding DDR4/DDR5/DDR6 differences?
- Passing messages between xPUs?



CXL Technology Basics

- Memory-speed access over PCIe physical layer
- Supports new architectures:
 - Disaggregated memory
 - Pooled memory
 - Switches for memory fabrics
 - Shared memory
 - Persistent memory



CXL DRAM Is A Lot Like An SSD





CXL DRAM Is A Lot Like An SSD









What Do Potential Users Say?

- **Google:** Stranded memory is not important
- IBM/Georgia Tech: DDR is a poor answer
- Al Providers: We need enormous memories
 Also fast loads of GPU HBM
- Hyperscalers: "Any-to-Any" xPU connections
- PC OEMs: CXL is not immediately useful



CXL Forecast: Slow Steady Growth



13 | ©SNIA. All Rights Reserved.

Long-Term Impact

- Re-thinking system architecture
 - Disaggregated memory
 - Processor arrays with mesh networks
 - Memory agnostic
- Better memory bandwidth & size vs. worse latency
 - Design-arounds will optimize for this
 - CXL Hot-Range Monitoring Unit (CHMU)



CXL Looks for the Perfect Home



New report from Objective Analysis

- Covers all perspectives
 - Where CXL is useful, and where it isn't
 - Demand drivers for CXL DRAM modules
 - Opportunities outside of DRAM
 - Forecast (Revenues, units, ASP)
- Available for immediate download:
 - <u>Objective-Analysis.com/reports</u>





CXL Consortium Update



Mahesh Natu



CXL Specification Release Timeline





CXL 3.x Themes

- Scaling
- Memory Hierarchy
- Security



SNIA SNIA SNIA COMPUTITIV

CXL: Scaling new heights



Compute Express Link[®] and CXL[®] are trademarks of the Compute Express Link Consortium.

SNIA SILA COMPUTITIV

CXL Memory Introduces New Tiers





CXL Hot-Range Monitoring Unit (CHMU) for Tiering



More efficient SW Memory Tiering Better Perf, lower TCO

Challenges faced by the today's SW tiering solutions

- Requires application changes
- Must trade-off accuracy against perf overhead
- CPU vendor specific

CHMU addresses these problems

- Enables generic OS based solutions
- Works for simple and pooling memory devices
- Tracking offloaded to CXL memory device
- Configurable and extensible



CXL Trusted Security Protocol (TSP)



Confidential computing is a pre-requisite to migrating sensitive WLs to Cloud

TSP extensions enable CXL devices to participate in Confidential Computing

• CXL memory device can hold proprietary model parameters



2025 CXL Events

- April
 - CXL DevCon 2025, April 29-30
- August
 - FMS: The Future of Memory and Storage, August 5-7 at the Santa Clara Convention Center
 - Visit CXL kiosk at Booth No. 725
 - Attend CXL presentation on **Wednesday**, **August 6 at 9:45 am PT** to learn about • ready to deploy CXL solutions!
- September
 - Al Infra Summit, September 9-11 at the Santa Clara Convention Center
- November •
 - SC'25 (Supercomputing 2025), November 16-21 in St. Louis, MO
 - Visit the CXL Pavilion at Booth No. 817.
- Other events will be added as opportunities arise Follow the CXL Consortium via LinkedIn (@CXL ٠
 - Consortium) for updates!

ryand AL mem Programming Workshop Instruction Video CMSI AND STORAGE

Persistent Memory and CXL.mem Programming Workshop

https://github.com/pmemhackathon/hackathon

Presented by Igor Chorazewicz

Learn how to program **CXL** Memory Modules in SNIA's **CXL.mem Programming Workshop** Live online and at CXL DevCon! https://www.snia.org/pmsummit/hackathon





CXL Benefits, Implementation, and Ecosystem

Torry Steed



CXL[®] Memory Expansion Advantages



Increased Memory Capacity



Reduced Memory Cost



Increased Performance





DDR5 DIMM Price

	Capacity	Price/G b
	64GB	1x
Г	96GB	1.1x
	128GB	1.4x
L	256GB	2.0x





Applications Benefiting from Large Memory

Real Time Data Analytics Retrieval and Processing



In-Memory Databases

- Keep the entire dataset in memory for real time/fast processing
- Examples: SAP-Hana, Redis
- Real time Feature Stores for Edge Inference applications



Big Data Analytics, AI and Deep Learning

- Newer analytics and machine learning utilizing very large datasets
- Enables rapid query and model training



Climate modeling

- Genomics research
- Fluid dynamics
- Particle Physics



Financial Modeling

- High frequency trading complex risk analysis
- Processing, ingesting vast amounts of market data in real-time



CXL[®] Generational Improvements

CXL 1.1 In-server memory expansion DIMMs DDR CPU DDR CXL via PCIe 5.0 CXL Memory Memory ž Controller **CXL** Memory Module



CXL 3.0 Fully disaggregated memory pool





27 | ©SNIA. All Rights Reserved.

CXL[®] Memory Adoption

In-Server Memory Expansion



In-rack Memory Expansion

CXL 2.0 Switching





CXL[®] and DRAM Latency





CXL[®] and DRAM Bandwidth



Real World Results

In-Memory Database (MS SQL + TPC-H)	Machine Learning (Apache Spark™ SVM) ບໍ່ມີ	High Performance Computing (CloverLeaf)
 Consumes only 50% of memory bandwidth Capacity limited 	 Memory bound workload Sensitive to latency & capacity 	 Bandwidth intensive workload 80% mapped to DRAM 20% mapped to CXL
 44%-88% reduction in SSD paging I/Os 23% performance improvement 	 2.21x performance improvement compared to DRAM only 	 CXL increases memory bandwidth by 33% 17% performance improvement

Source: CXL Memory Expansion: A Closer Look on Actual Platform by Micron and AMD



National Yang Ming Chiao Tung University 1.5B 7B 📕 w/o CXL 📕 w/ CXL w/o CXL w/ CXL **Example Configuration** 40000 10000 30000 7500 256GB Throughput (token/s) Throughput (token/s) 20000 5000 2500 10000 0 Ω 1 x H100 | 8K 2 x H100 | 8K 1 x H100 | 128K 2 x H100 | 128K 1 x H100 | 8K 2 x H100 | 8K 32B 14B w/o CXL w/ CXL 📕 w/o CXL 📕 w/ CXL 2500 5000 2000 4000 en/s) (token/s) 1500 3000 Ę ughput roughput 2000 1000 128GB 128GB 1000 500 0 0 1 x H100 | 8K 2 x H100 | 8K 1 x H100 | 128K 2 x H100 | 128K 1 x H100 | 8K 2 x H100 | 8K

NYCU Taiwan

Study on CXL Affects on LLM Training with CPU Offloading

1 x H100 | 128K 2 x H100 | 128K





Level1Techs Real World Demo



"Benchmarking the Xeon 6 6787P in the Supermicro 222H-TN"





SMART 4-DIMM AIC





Source: Benchmarking the Xeon 6 6787P in the Supermicro 222H-TN

CXL[®] Ecosystem "Connecting the Dots"

Year	AMD EPYC	CXL
2021	Milan	No
2022	Genoa	CXL 1.1
2023	Bergamo	CXL 1.1
2024	Turin	CXL 2.0
2025	Venice	CXL 3.0

BIOS Support

finsyde Full support for CXL

Intel[®] Server M50FCP Family CXL BIOS version v0.1.02.0002

(intel)					
Year	Intel Xeon	CXL			
2021	Ice Lake	No			
2023	Sapphire Rapids	CXL 1.1			
2023	Emerald Rapids	CXL 1.1			
2024	Granite Rapids	CXL 2.0			
2025	Diamond Rapids	CXL 3.0			

AMIBIOS supports CXL



OS Support





Windows Server Beta testing CXL

VMWare supports CXL



Where does that leave CXL[®]?



- CXL memory expansion is available now
- "Off the shelf" systems are here and working
- CXL benefits are real
- CXL rack pooling appliances are coming soon
- CXL rack to rack sharing is on the horizon





Questions?



Next Steps...

- Please rate this webinar and provide us with feedback
- Look for a Blog answering questions from this webinar at <u>www.sniablog.org</u>
- Get more SNIA education!
 - Live
 - ✤ FMS Future of Memory and Storage, August 4-7, 2025 Santa Clara CA fms.com
 - SNIA Developer Conference, September 15-17, 2025 Santa Clara CA sniadeveloper.org
 - SC25, November 16-21, St. Louis MO sc25.supercomputing.org

o Online

- This webinar and many other videos and presentations on today's topics are in the SNIA Educational Library https://snia.org/educational-library
- Join SNIA and the Persistent Memory Special Interest Group
 - www.snia.org/join
 - https://www.snia.org/groups/cms





Thank You

