



# OCP

FUTURE  
TECHNOLOGIES  
SYMPOSIUM

## OCP Global Summit

November 8, 2021 | San Jose, CA

# Innovations in Memory System Architecture: PIM and CXL-Memory

---

**David Wang, Ph.D.**  
**Director Memory Product Planning**

# Disclaimer

This presentation and/or accompanying oral statements by Samsung representatives collectively, the “Presentation”) is intended to provide information concerning the SSD and memory industry and Samsung Electronics Co., Ltd. and certain affiliates (collectively, “Samsung”). While Samsung strives to provide information that is accurate and up-to-date, this Presentation may nonetheless contain inaccuracies or omissions. As a consequence, Samsung does not in any way guarantee the accuracy or completeness of the information provided in this Presentation.

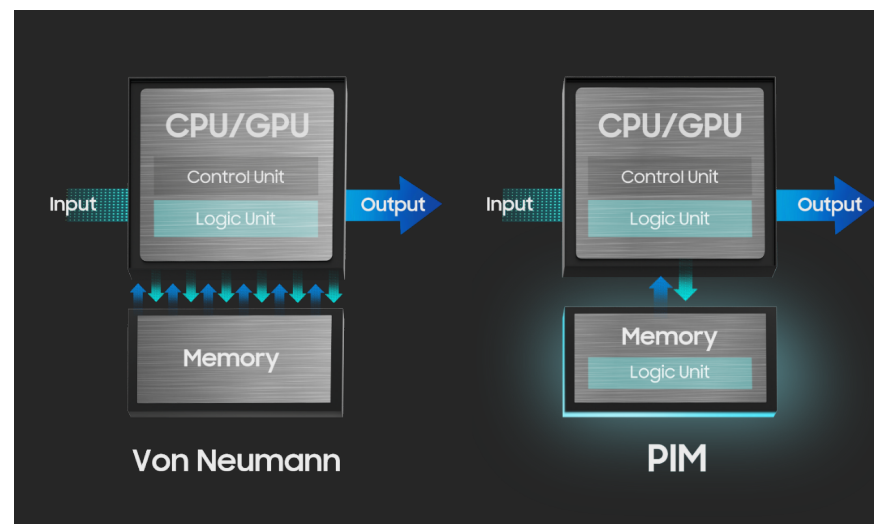
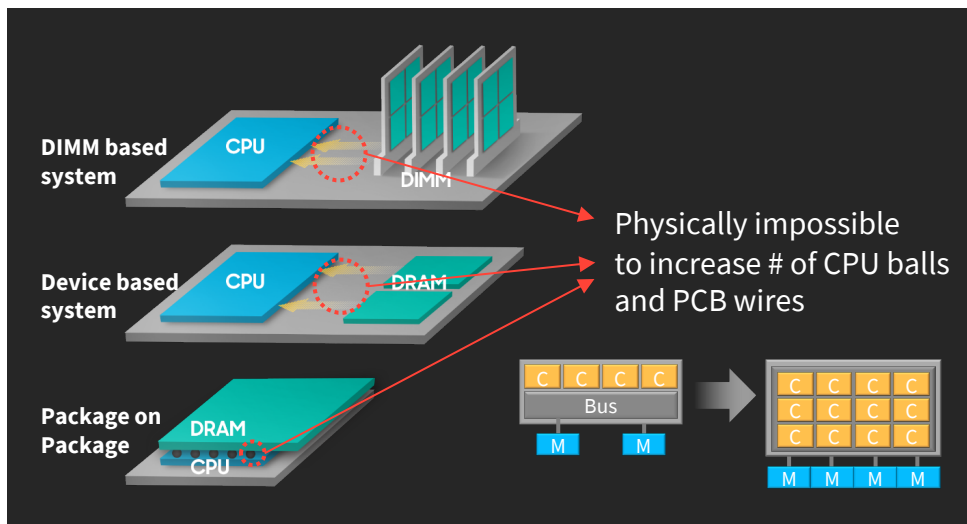
This Presentation may include forward-looking statements, including, but not limited to, statements about any matter that is not a historical fact; statements regarding Samsung’s intentions, beliefs or current expectations concerning, among other things, market prospects, technological developments, growth, strategies, and the industry in which Samsung operates; and statements regarding products or features that are still in development. By their nature, forward-looking statements involve risks and uncertainties, because they relate to events and depend on circumstances that may or may not occur in the future. Samsung cautions you that forward looking statements are not guarantees of future performance and that the actual developments of Samsung, the market, or industry in which Samsung operates may differ materially from those made or suggested by the forward-looking statements in this Presentation. In addition, even if such forward-looking statements are shown to be accurate, those developments may not be indicative of developments in future periods.

# Agenda

- PIM (Processing In Memory)
- AXDIMM (Processing Near Memory)
- CXL Memory
- SMDK (Scalable Memory Development Kit)
- Concluding Remarks

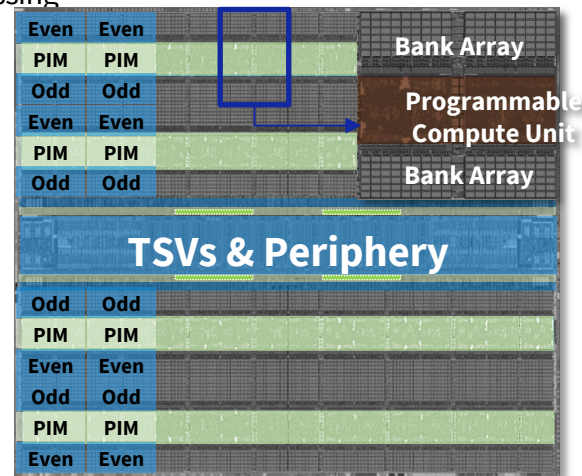
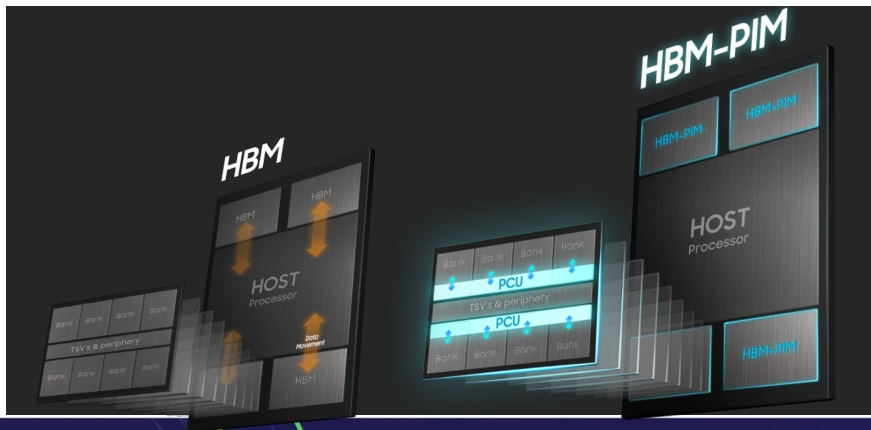
# PIM (Processing-in Memory) Rationale

- System-level performance is constrained by bandwidth scaling
  - Limited by # of PCB wires, # of CPU ball, and thermal constraints
- Proposing to use PIM to improve performance of bandwidth-intensive workloads and improve energy efficiency by reducing computing-memory data movement.



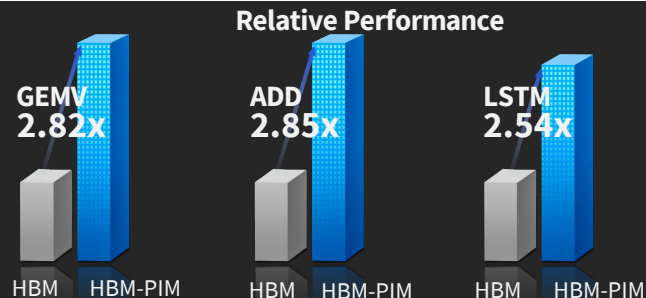
# Aquabolt-XL, 1<sup>st</sup> In-system Demonstrated PIM

- The first demonstrator vehicle of PIM is based on HBM2 Aquabolt, used in leading edge AI and HPC systems.
- HBM-PIM targeted to complement xPUs for optimal system balance and performance per watt for memory-bound AI workloads
  - Accelerate FP32/FP16/INT16 data processing capability in memory
- Programmable Computing Unit (PCU) integrated with memory core within HBM to enable parallel processing and minimize data movement.
- Improves the performance and energy efficiency of the system with in-DRAM processing
  - Performance : Utilize up to 4× higher in-DRAM bandwidth by multi-bank parallel operation
  - Energy Efficiency : Reduce data movement energy by utilizing in-DRAM data processing unit



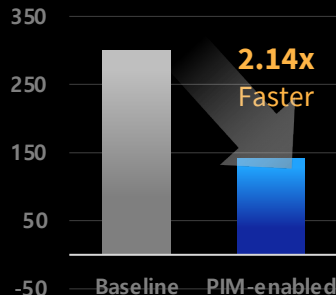
Aquabolt-XL silicon die photo

# Xilinx Alveo U280: HBM-PIM Evaluation

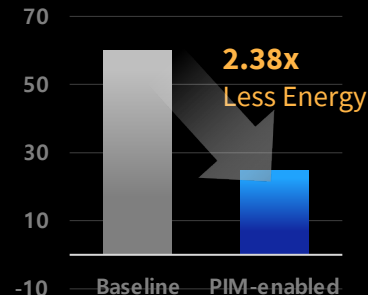


RNN-T Performance and Energy results

RNN-T Latency (ms)

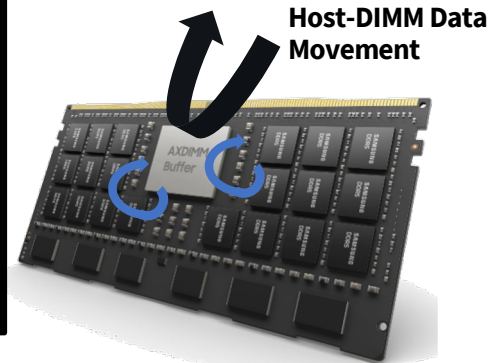
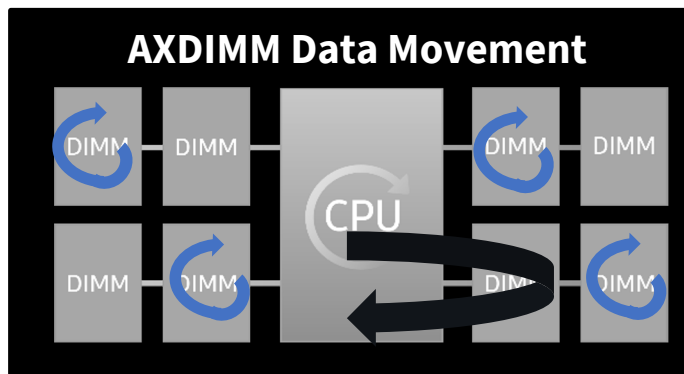
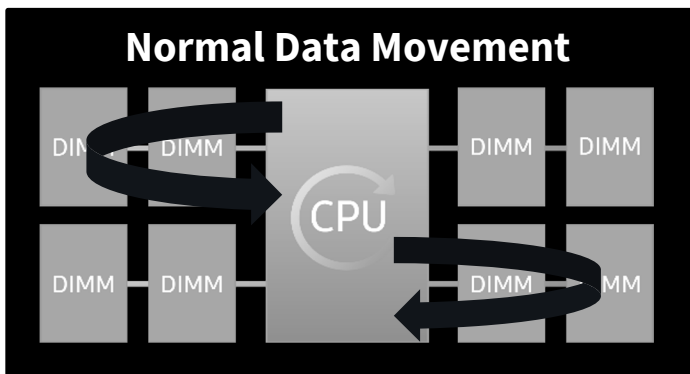
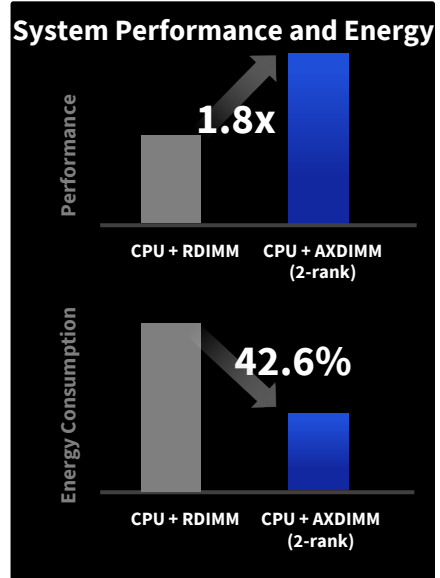


System Energy Consumption (J)



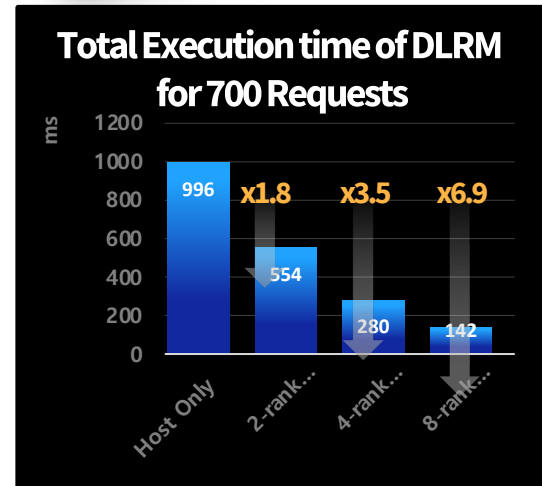
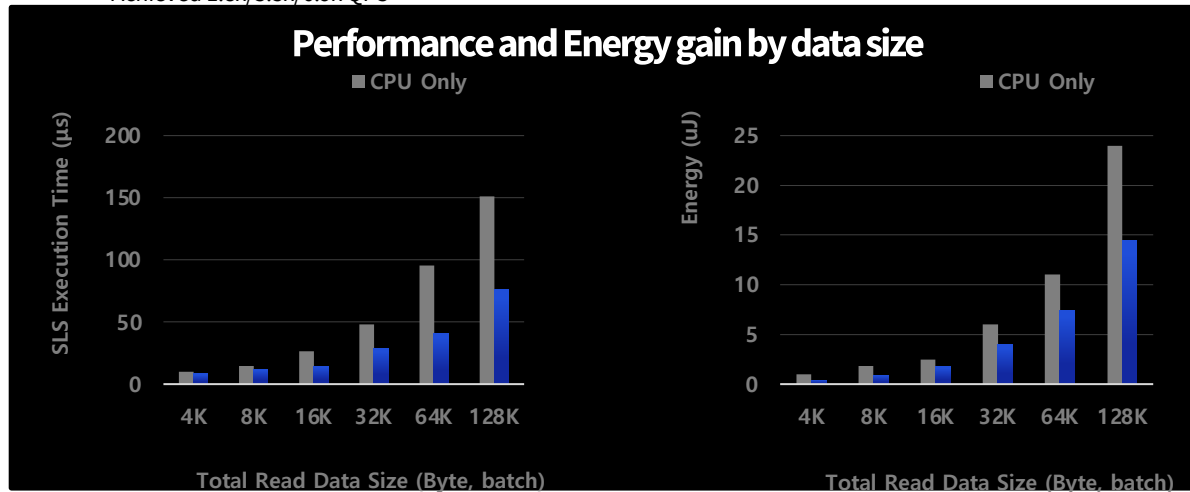
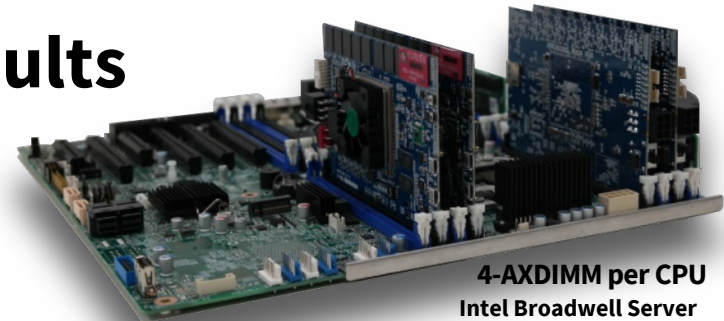
# AXDIMM – Near Memory Compute

- CPU-memory data movement bottlenecks system performance → Use rank-level parallelism
  - Instead of processing in-memory, AXDIMM places compute on buffer with commodity DRAM on memory module
  - Near-Memory Compute instead of In-Memory Compute. Same idea – alleviate von Neumann memory bus bottleneck
  - Samsung to provide AXDIMM SW stack to offload the acceleration functions in AXB(AXDIMM Buffer)
- Improve the performance and energy efficiency of the system with in-DIMM processing
  - Utilize up to higher in-DIMM bandwidth by multi-Rank parallel operation, 1.8x by 2-rank
  - Reduce data movement energy by utilizing in-DIMM data processing unit, -42.6% by 2-rank



# AXDIMM Evaluation System and Results

- x86 based platform with Xilinx Zynq Ultrascale+ FPGA Chip configured as AXDIMM
- Enabled RecNMP\* logic
  - Achieved 1.8x SLS execution speed-up from HW
- Modified DLRM\*\* application utilizing AXDIMM
  - Achieved 1.8x/3.5x/6.9x QPS



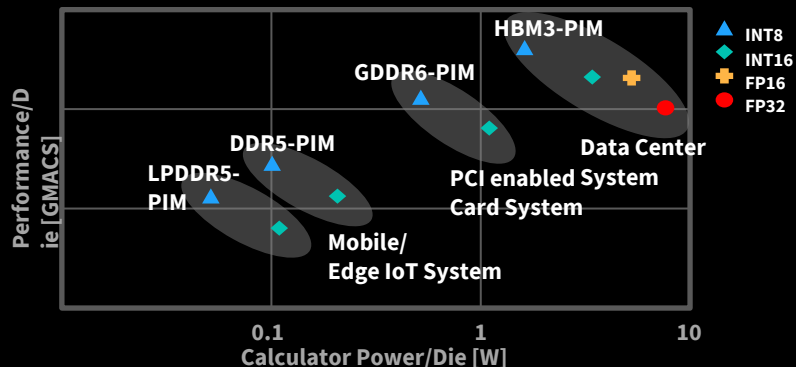
\*: RecNMP - L. Ke *et al.*, "RecNMP: Accelerating Personalized Recommendation with Near-Memory Processing," 2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)

\*\*: DLRM - M. Naumov *et al.*, "Deep learning recommendation model for personalization and recommendation systems," arXiv preprint, 2019

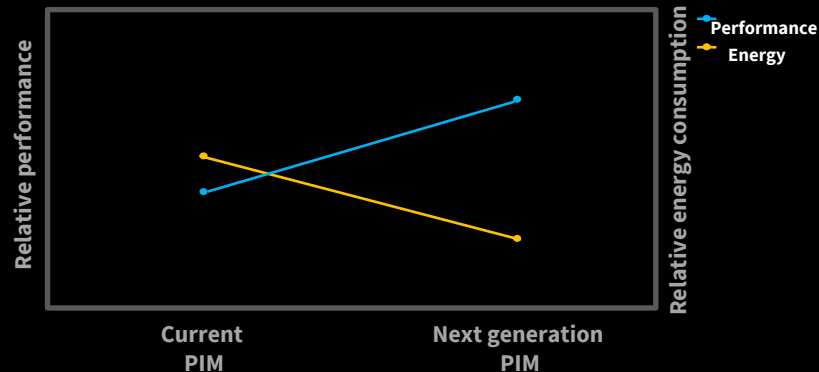
# In/Near Memory Compute Future

- Wider target applications
  - PIM unit supporting multiple functions
- Various DRAM types
  - LPDDR5, DIMM-DDR5, GDDR6, HBM3
- New standards for PIM
  - Command truth table/timing for PIM
- Addendum or addition to current product specs, not new generations
  - Enhanced performance
  - Reduced energy
- Collaborate with industry
  - Supporting custom functions
  - Investigating CXL enablement

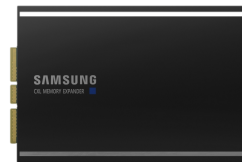
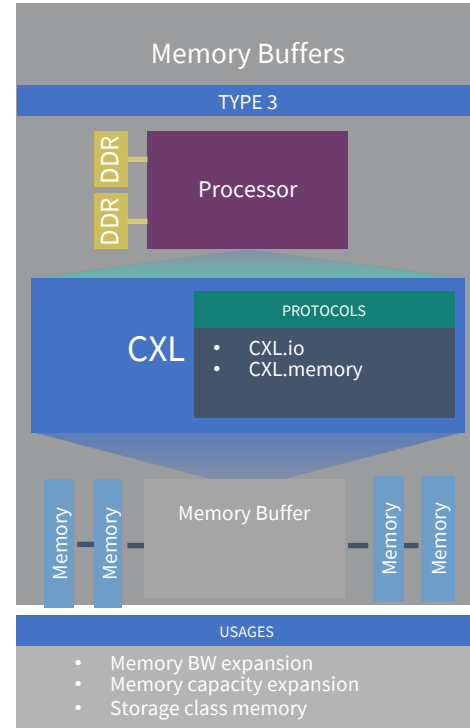
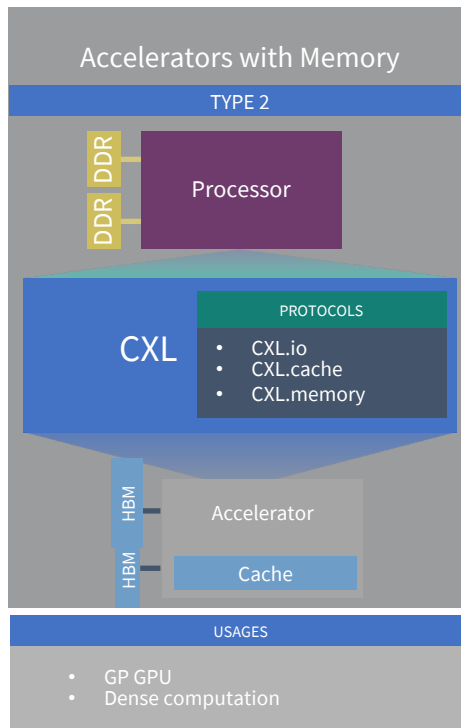
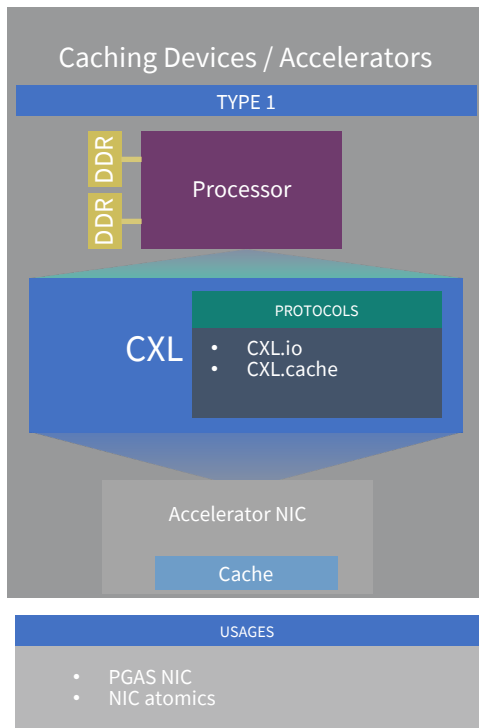
## Prospective PIM-supported data format



## Performance/Energy of next generation PIM



# CXL Device Types



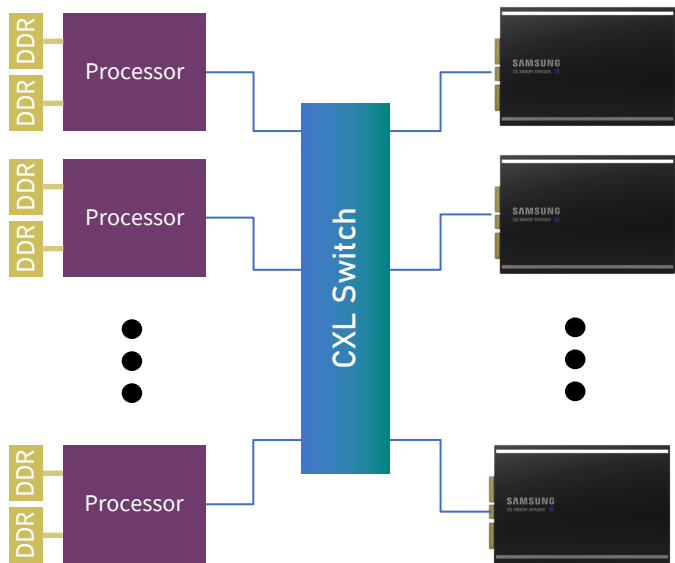
<https://news.samsung.com/global/samsung-unveils-industry-first-memory-module-incorporating-new-cxl-interconnect-standard>

# Samsung CXL Type 3 Memory

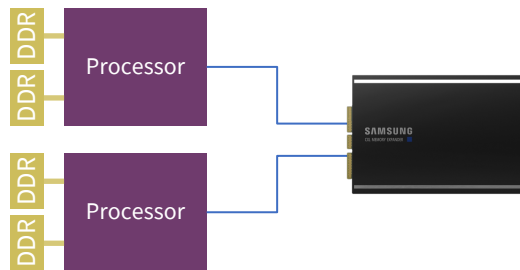
- Form Factor
  - (Under consideration) E3.S, E1.S, E3.L
- CXL Specification
  - 2.0
- Media
  - DDR5
- Capacity
  - 128 GB to TBD
- Other Features (under consideration)
  - Dual Port
- <https://news.samsung.com/global/samsung-unveils-industry-first-memory-module-incorporating-new-cxl-interconnect-standard>



# Memory Pooling with CXL Type 3 Memory

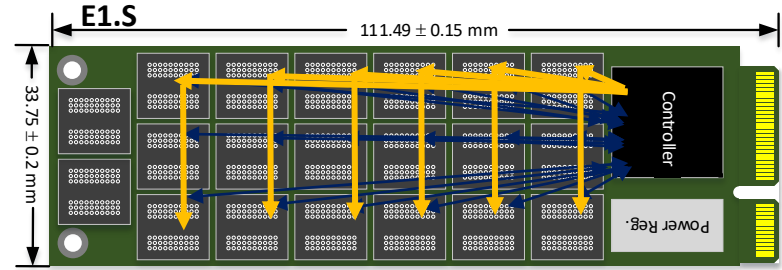
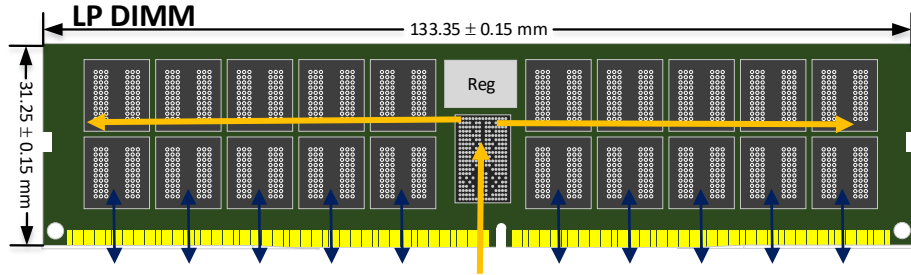


Multiple hosts using pooled memory, attaining memory capacity for exclusive use for certain duration, as needed – through CXL switch



Two hosts using pooled memory, attaining memory capacity for exclusive use, as needed

# DRAM on LP RDIMM vs DRAM on E1.S



## New learnings in

- Component placement
- Routing
- Signal Integrity
- Power Integrity
- Thermal
- Firmware Development

## Form Factors

- HDD: 5.25" → 3.5" → U.2
- SSD: U.2 → EDSFF
- CXL: EDSFF → ??

# CXL Memory Development

- Bandwidth
- Latency
- Cost
- Power
- Capacity
- QoS



- Media controller ownership
  - RAS Features
  - Security Features
  - Thermal Management
  - Persistence Management
  - Media Management
    - DRAM
      - DDRx
      - LPDDRx
    - Future Memory

# SMDK: Scalable Memory Development Kit

## Compatible API

: Supports Memory Expander application without application SW modification

## Optimization API

: Supports high-level optimization by modifying the SW application

## Intelligent Tiering Engine

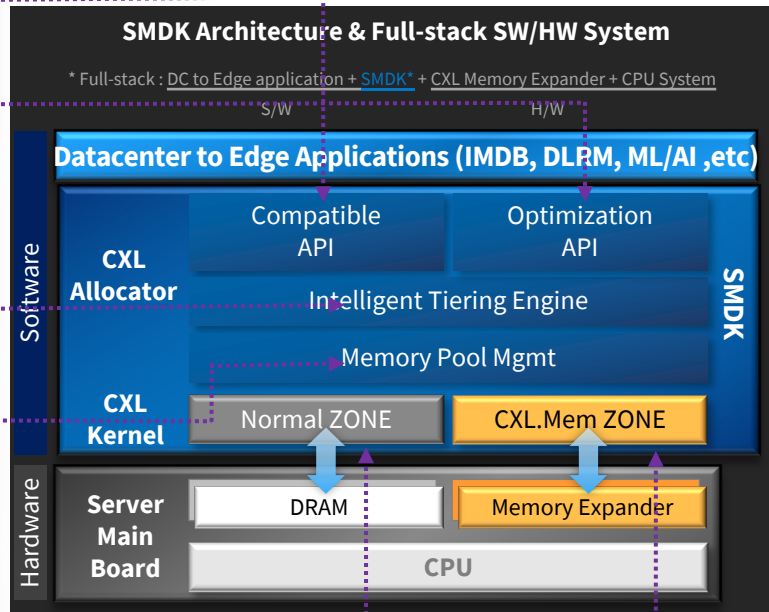
: Supports different use cases by tiering priorities / capacities / bandwidth among memories

## Memory Pool Management

: Separate management of normal memory and CXL.Mem  
: Supports scalability according to CPU memory capacity change

## Memory Zone

: Optimized management preventing mixed usage of different memories



# SMDK Summary

- Open-source software tool to
  - Facilitate CXL memory deployment without the need to modify existing applications
  - Or, allow application programmers to optimize use of memory with different BW/latency characteristics
- Development framework to lays the foundation to managing memory tiering, Intelligent data flow, advanced RAS features
- Now available on a limited basis for initial testing and optimization and will be open-sourced within the first half of next year
- CXL Memory Demo w/SMDK @ OCP
- <https://news.samsung.com/global/samsung-introduces-industrys-first-open-source-software-solution-for-cxl-memory-platform>

# CXL Memory OCP Demo

## 1. Memory Allocation Demo

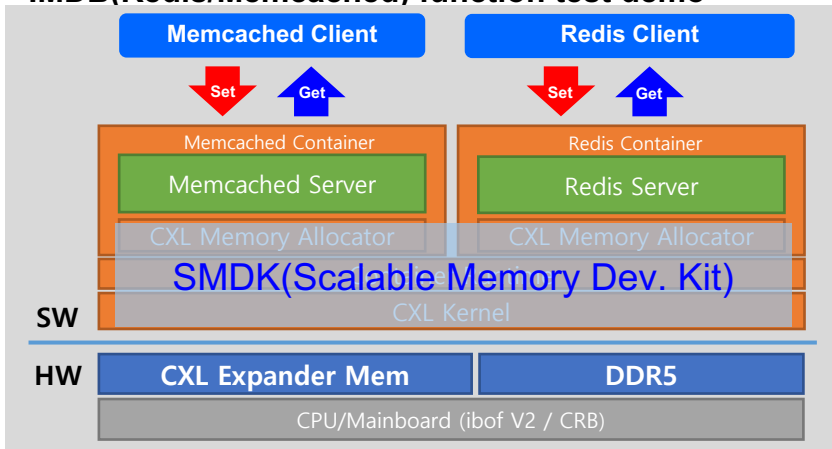
- CXL memory zone allocation with SMDK (CXL zone vs Main memory zone)

## 2. IMDB(Redis) Functional Test Demo

- *Running Redis and Memcached on CXL memory, CXL+DRAM with SMDK(w/o performance test)*
- SMDK is Samsung's s/w development kit which contains cxl memory allocator

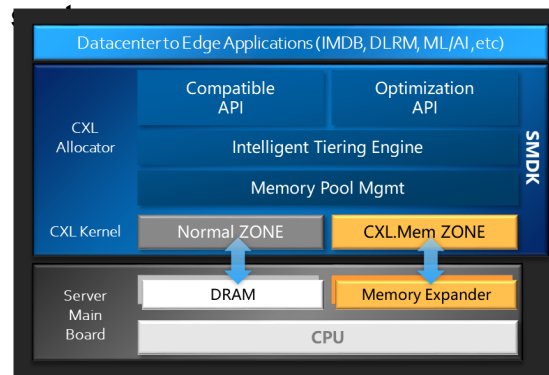
### OCF Demo Scenario

#### IMDB(Redis/Memcached) function test demo



### SMDK Architecture

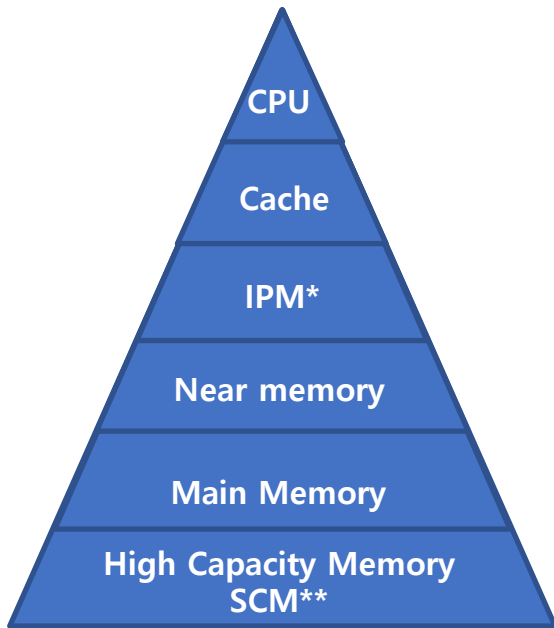
#### Scalable memory development kit allows for easy and optimized deployment of heterogeneous memory



- ① Compatible & Optimization API
- ② Intelligent Memory Tiering (priority/capacity/bandwidth)
- ③ Separate MGMT. of Normal/CXL Memory

# Concluding Remarks

- Must understand application requirements in context of (ever evolving) Memory Hierarchy
- Move data where it is needed, when needed
- Don't move data if it can be avoided
  - Local, heterogeneous computing
- Software (as always) key to manage (or reduce/eliminate) data movement for performance and energy efficiency improvement
- CXL will be key interface for enabling development of new media and SDM system architecture, as well as heterogeneous computing system architecture



\* : In package memory

\*\* : Storage Class Memory



**OCP**  
FUTURE  
TECHNOLOGIES  
SYMPOSIUM



# OCP

## FUTURE TECHNOLOGIES SYMPOSIUM

2021 OCP Global Summit | November 8, 2021, San Jose, CA