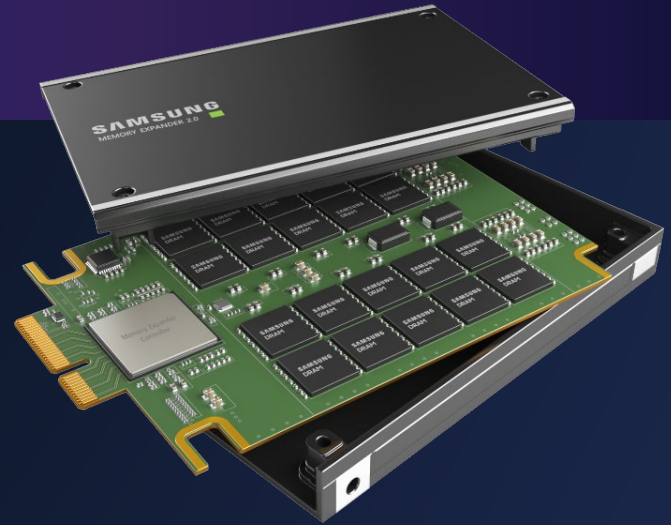


# Scaling of Memory Performance and Capacity with CXL Memory Expander

August, 2022 | Samsung Electronics Co., Ltd.

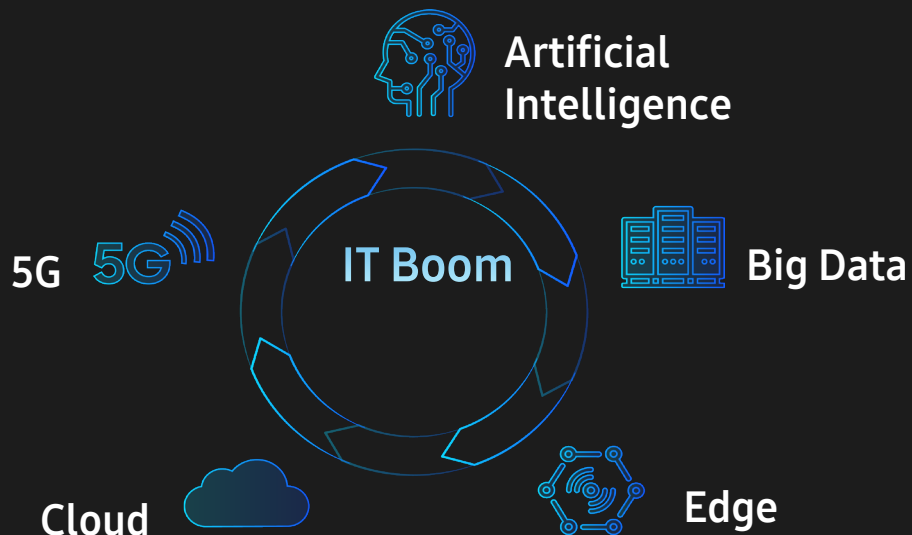
S. J. Park, K.-S. Kim, H. Kim, J. So, J. Ahn, J. Jung, I. Yun, S. Ryu, W.-J. Lee, J.-G. Lee, H.-Y. Ryu, C. Y. Lee, J. Prout, K.-C. Ryoo, S.-J. Han, M.-K. Kook, J.S. Choi, J. Gim, Y. S. Ki, S. Ryu, C. Park, D.-G. Lee, J. Cho, H. Song, and J. Y. Lee



# Agenda

- Industry Trends and Challenges
- Introduction of CXL (Compute Express Link)
- CXL Memory Expander Features
- SMDK: Unified Software Solution for CXL
- Application Benchmark Test Results
- Summary and Future Plan

# Industry Trends and Challenges



Massive demand for data-centric technologies and applications

Memory bandwidth and density not keeping up with increasing CPU core count

Need a next gen interconnect for heterogeneous computing and server disaggregation

# The Fast-Growing Computing Workloads

- Large-scale adoption of AI and ML

Smarter devices

Hyper-connected networks

Super-intelligent services

Digital transformation

Pandemic



First Era

Modern Era

Perceptron

1960

NETtalk

ALVINN

RNN for Speech

TD-Gammon v2.1

LeNet-5

BiLSTM for Speech

Deep Belief Nets & Layer-wise pretraining

AlexNet

ResNets

GPT-3

LaMBDA

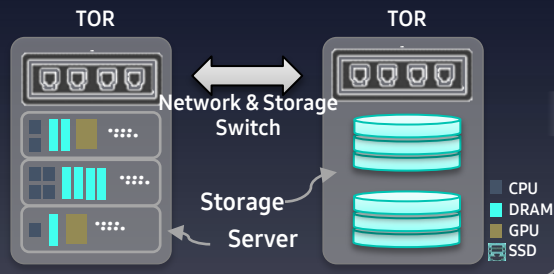
2020

# Evolution of Hyperscale Computing Environment

- From Converged to Composable Architecture

## Converged Architecture

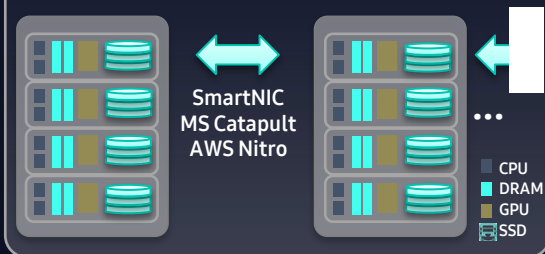
### TOR based Rack Scalable Architecture



Network Challenge

## Hyper-Converged Architecture

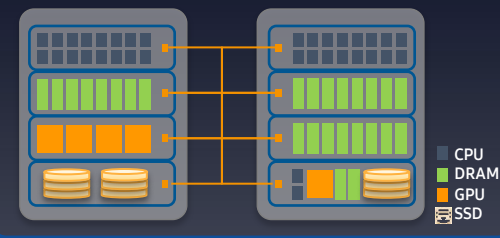
### Server & Storage Combined Architecture



Divergence Challenge

## Disaggregated / Composable Architecture

### Pooled Arch. : Memory, Compute, Storage



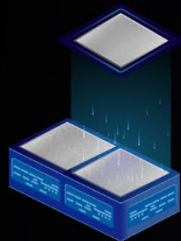
Interconnect Challenge

# The Rising Need for Better Connectivity

- Can be tailored and optimized for various AI applications

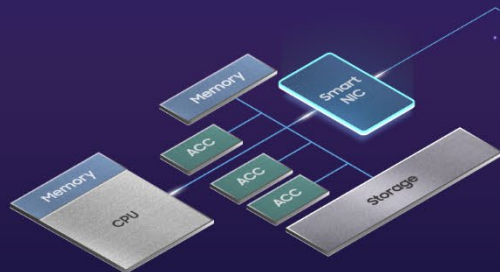
## SoC Interconnect

DIE / PACKAGE



## Processor Interconnect

NODE



**CXL** Compute  
Express  
Link™

A new class of interconnect  
for device connectivity in the era of AI

## Data Center Interconnect

DATA CENTER



## Customer Interconnect

MOBILE / BROADBAND



# CXL: Solution for the Era of HPC

- CXL as the core of composable computing infrastructure

## Key Features of CXL Interface



Cache Coherence



Connectivity



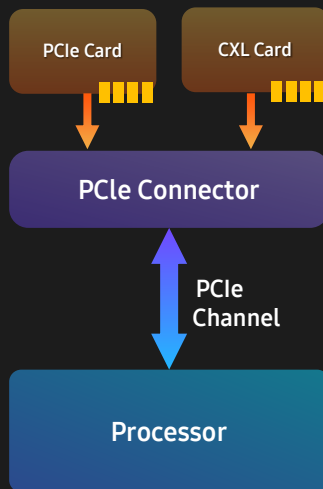
Byte Addressable



Low Latency

# CXL Features

CXL is a high-performance, low-latency protocol that leverages PCIe physical layer



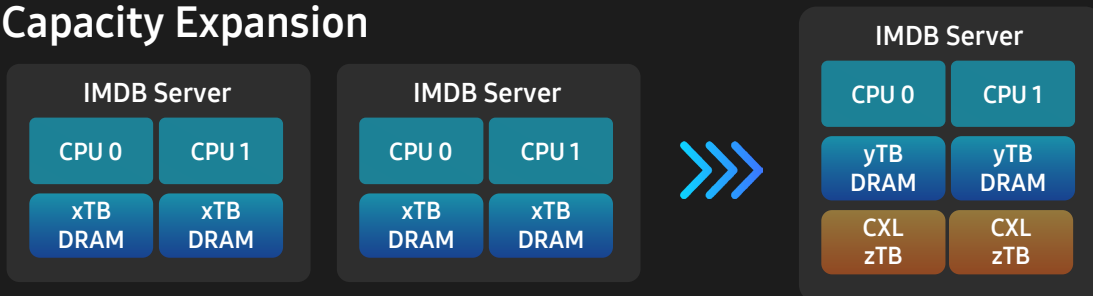
- High-speed and low-latency interconnect
- Leverages PCIe Physical layer (PCIe 5.0, PCIe 6.0)
- Supports various types of memories (volatile, non-volatile)
- CPU and CXL device memory coherency
- Supports switching and memory pooling
- Supports link level integrity and data encryption
- Open standard (non-proprietary)
- Broad industry support in CXL consortium
- Regular specification updates (CXL 1.1, CXL 2.0, CXL 3.0)



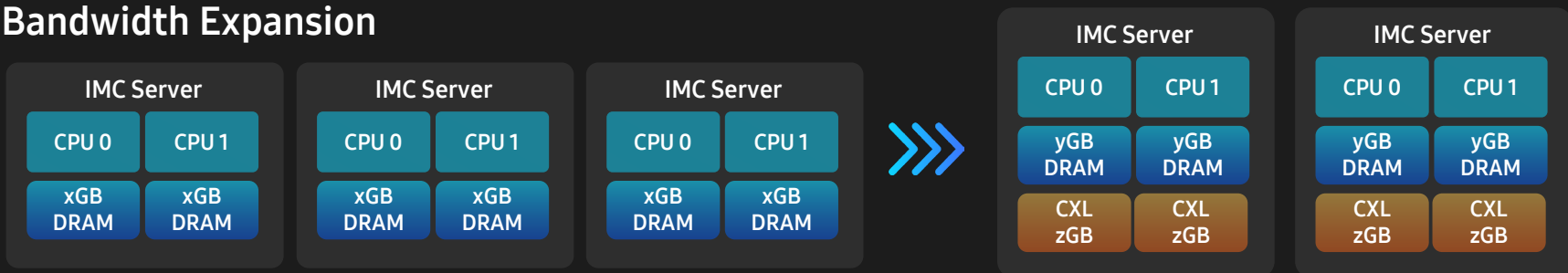
# CXL Use Cases (1/2)

- Capacity and Bandwidth Expansion

## Capacity Expansion



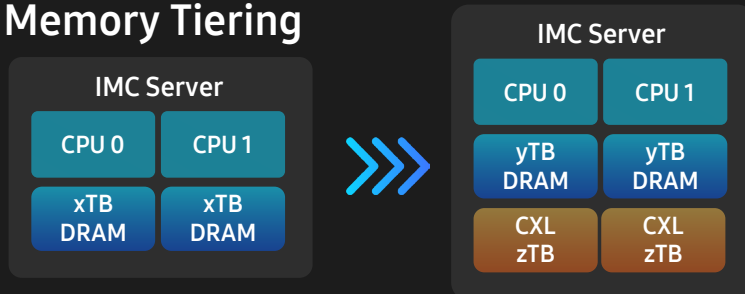
## Bandwidth Expansion



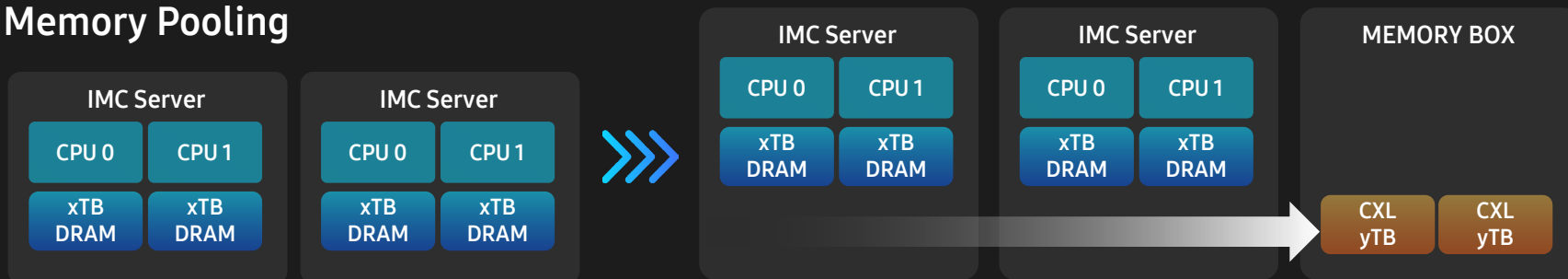
# CXL Use Cases (2/2)

- Tiering and Pooling

## Memory Tiering



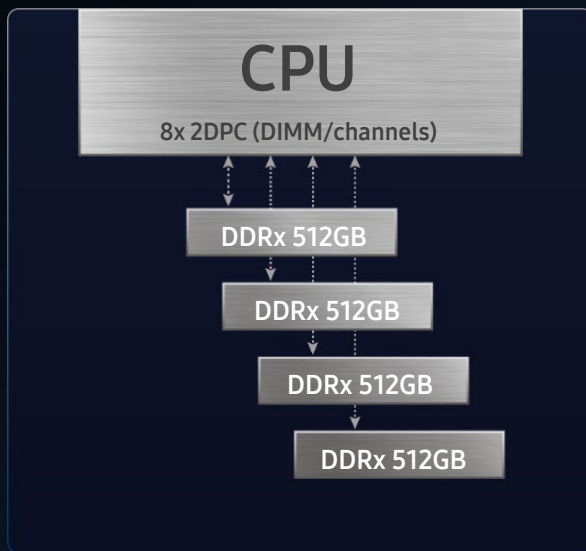
## Memory Pooling



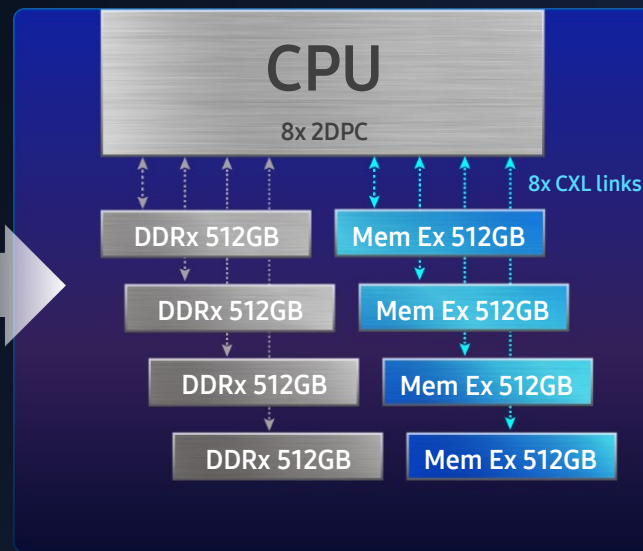
# CXL Memory Expansion Solution

- Doubled Capacity than Conventional Memory

Max. 8TB for 1CPU



Max. 16TB for 1CPU



Note: Max capacity varies with system configurations

# CXL Memory Expander

- New Solution for Memory Dominant Applications

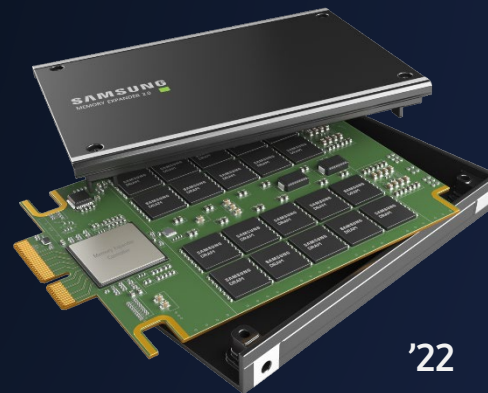


# CXL Memory Expander Line-up

- Built with FPGA and ASIC Controller



'21



'22

FPGA  
PCIe 3.0 (x16)

Host  
(Controller)

ASIC  
PCIe 5.0 (x8)

DDR4  
3200, 128GB

Media  
(DRAM)

DDR5  
4800+, 512GB

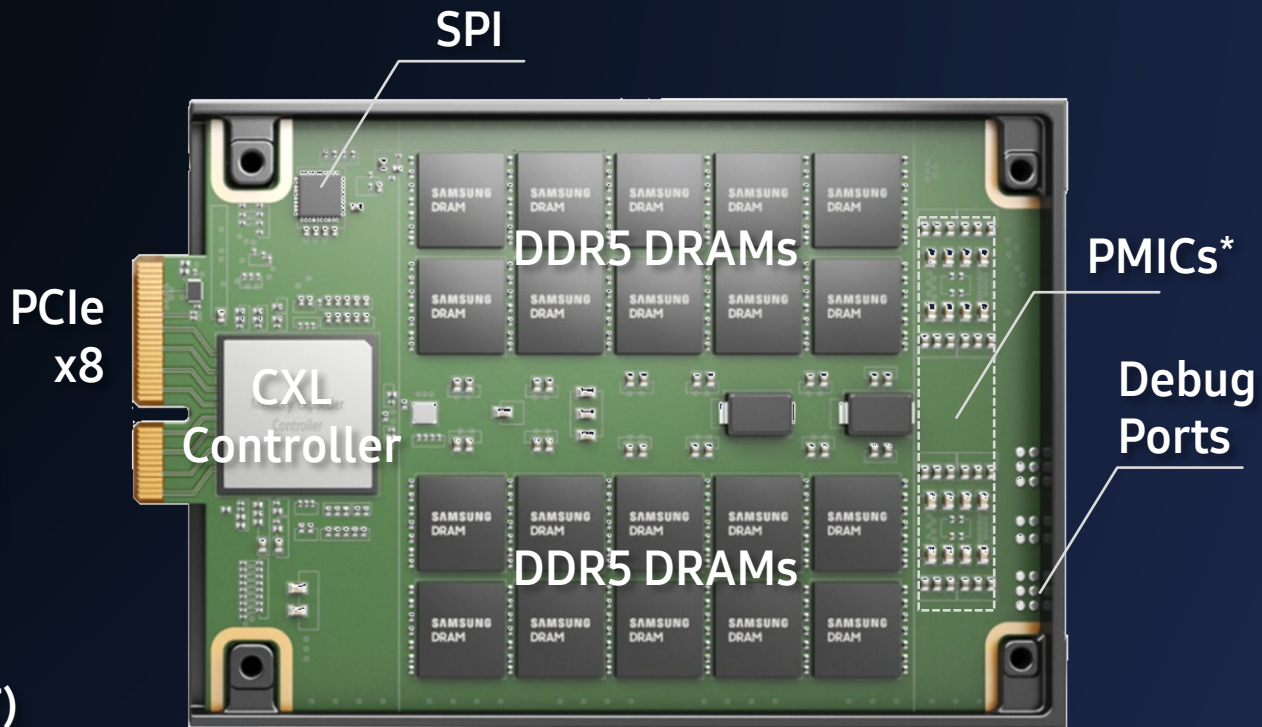
As of August, 2022

# CXL Memory Expander (1/3)

- Solution Overview



Enclosure (2T)



PCIe  
x8

SPI

CXL  
Controller

DDR5 DRAMs

PMICs\*

Debug  
Ports

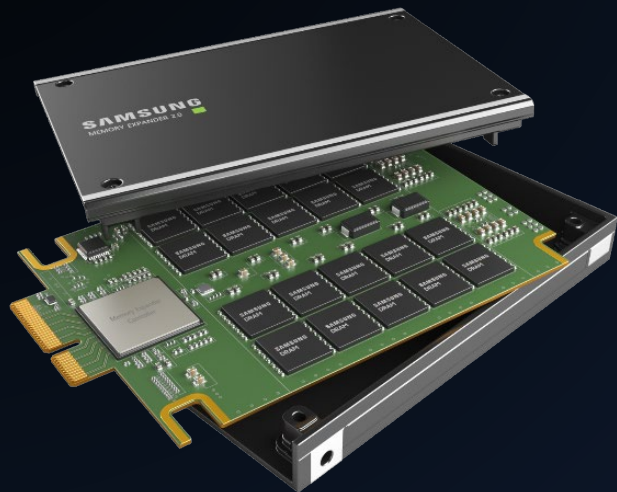
DDR5 DRAMs

E3.S Form Factor

\* Bottom-side

# CXL Memory Expander (2/3)

- Product Features

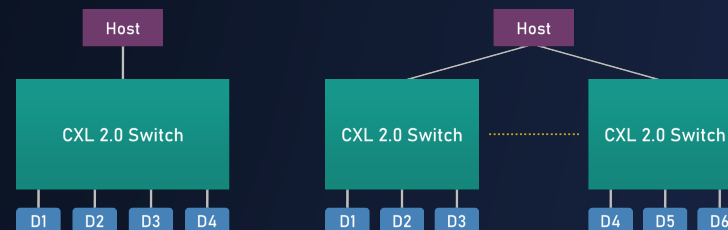


- 
- Form Factor - EDSFF (E3.S)
  - Media - DDR5 4800
  - Module Capacity - Max 512 GB
  - CXL Link Width - x8
  - Maximum CXL Bandwidth - 32GB/s (PCIe 5.0)
  - Other Features - RAS, Interleaving, Diagnostics etc.
  - Availability - Q3'22 for evaluation/testing
-

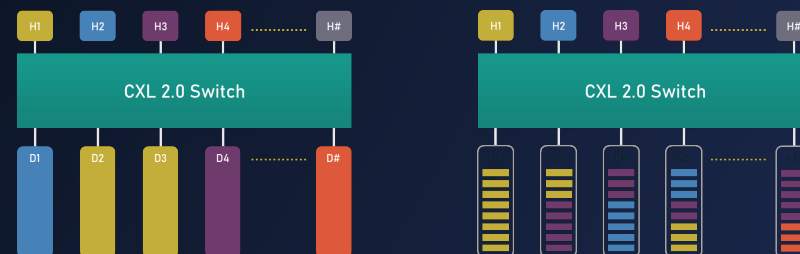
# CXL Memory Expander (3/3)

## Supported Features

- CXL 2.0
- Device Type: Type 3
- Support viral and data poisoning
- Memory error injection
- Multi-symbol ECC
- Media scrubbing
- Post package repairs (hard/soft)



## CXL 2.0 Switching Benefits



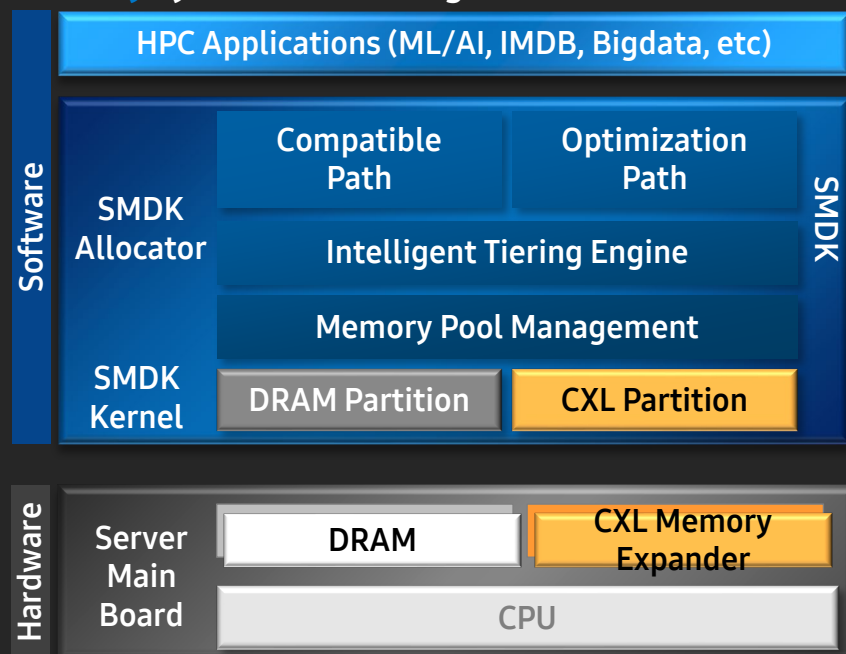
\* Image Source: CXL Consortium



# SMDK\*, Unified Interface for Memory

\* Scalable Memory Development Kit

SW development kit to enable **Software-Define Memory** system on heterogeneous memories



## Plugin

Two selectable paths, Compatible and Optimization Path, without or with modification of application SW

Intelligent Tiering Engine supports memory tiering scenarios with priority, capacity, bandwidth, and so on

Memory Pool Management supports scalability reflecting memory request status and system resource

## Kernel

Memory Partitioning allows logical memory views for heterogeneous physical DRAM and CXL memory

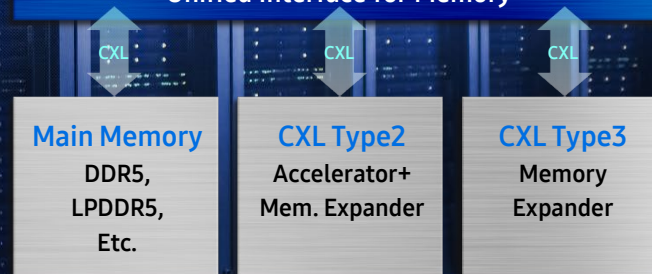
# Benefits of SMDK

- **Unified SW Solution**  
Full-stack SW all about heterogenous memory system
- **Client Experience**  
Transparent as well as Optimized Memory uses
- **CXL Ecosystem**  
OSS for CXL Industry and Research field

Differentiated  
Cloud Performance

**SMDK**

Unified Interface for Memory

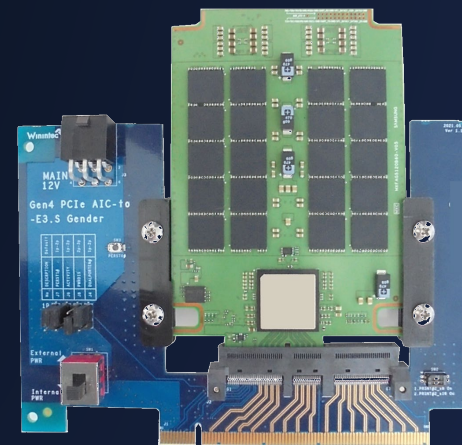
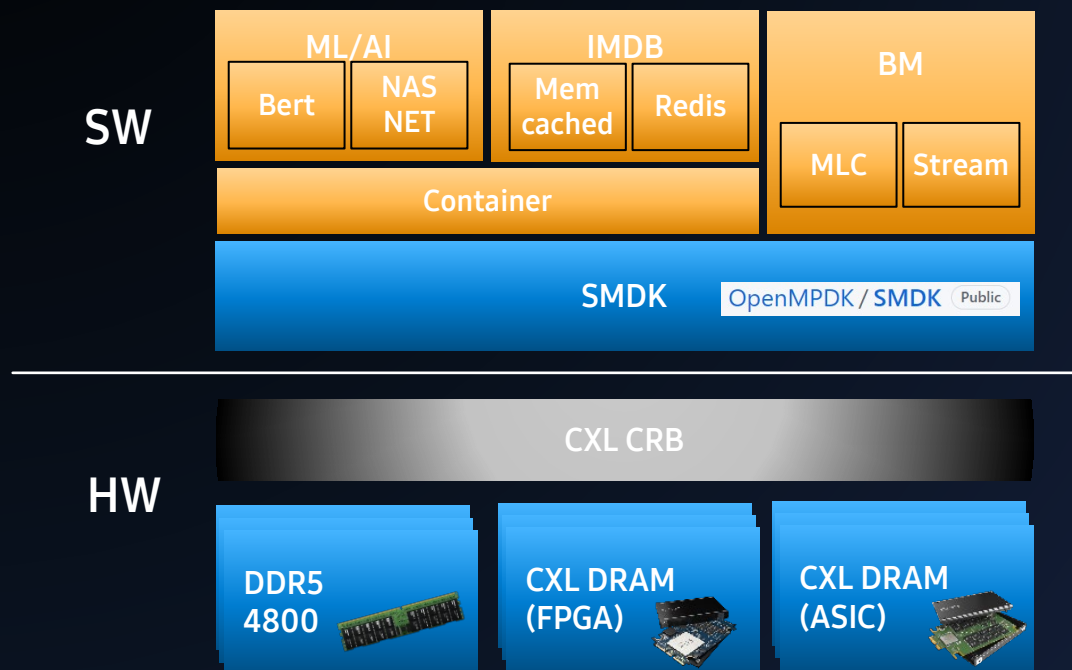


SMDK is available as open source on GitHub

➤ [https://github.com/OpenMPDK/SMDK/releases/tag/smdk\\_v1.1](https://github.com/OpenMPDK/SMDK/releases/tag/smdk_v1.1)

# Experimental Setup

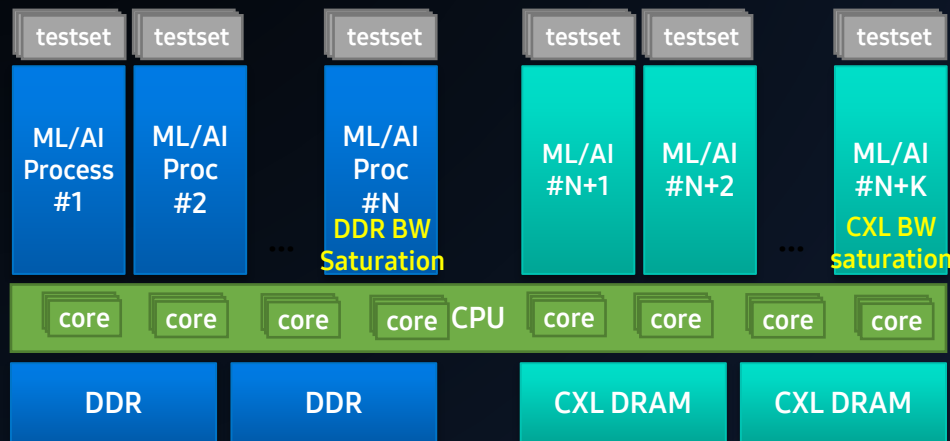
- Configuration of Test Bed



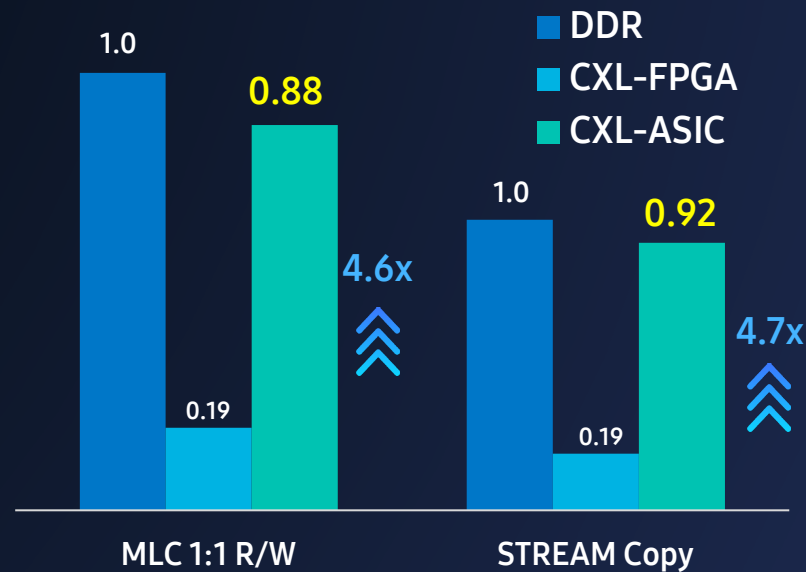
Memory Expander  
w/ EDSFF Riser Card

# Memory Benchmark Test Results

- Comparable Performance with DDR Memory

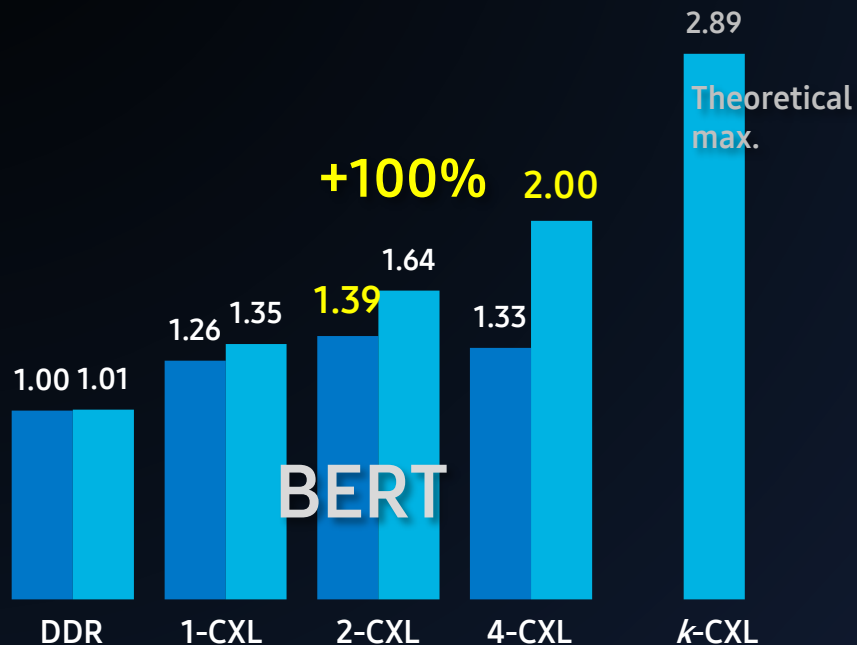


Normalized Bandwidth



# System Test Results (ML/AI)

## ML/AI Applications (BERT\* & NASNet\*\*)



\* Bidirectional Encoder Representations from Transformers

## Inferences per Minutes (Normalized)

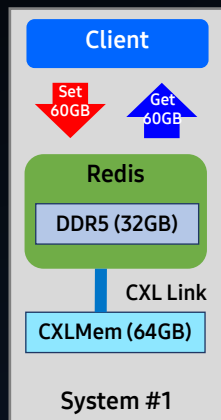


\*\* Neural Architecture Search Network

# System Test Results (IMDB)

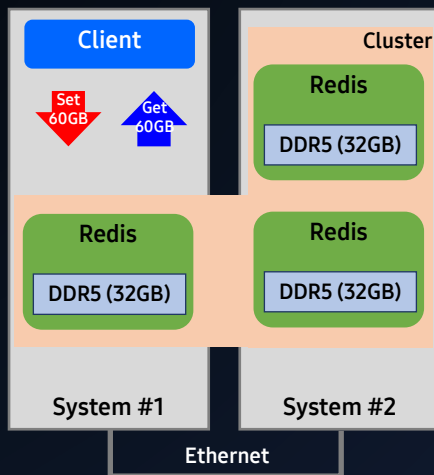
- IMDB Redis\* Memory Usage (Scale-up vs Scale-out)

Single Node  
(DDR+CXL FPGA)



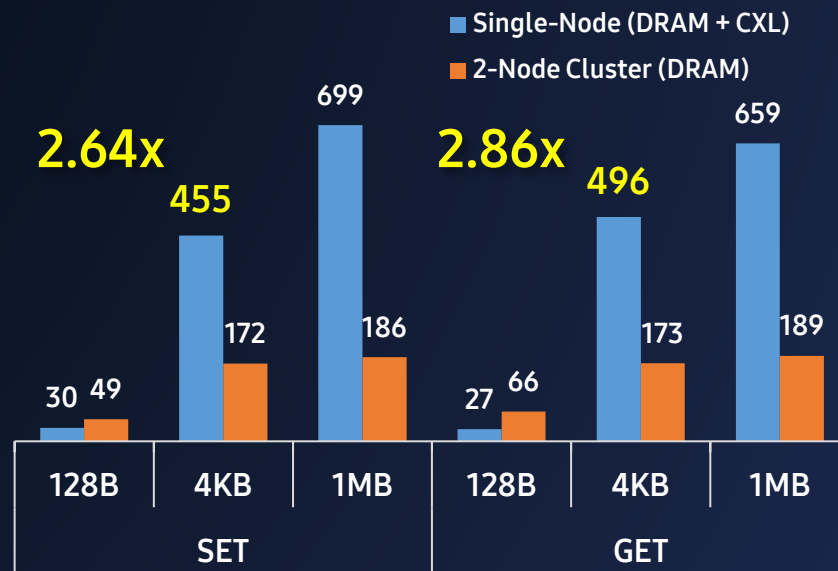
VS

2-Node Cluster  
(DDR x 3)



Performance [MB/s]

Scale-up vs Scale-out



\* Remote dictionary server

(See appendix for detail test condition)

# A Proven Memory Expansion Solution



Increasing System  
Memory Capacity

2X  
Increase



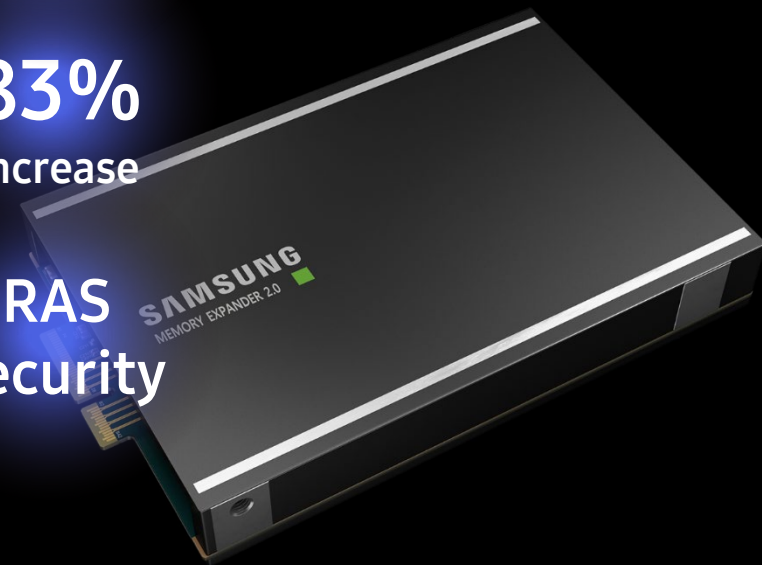
Widening  
Memory Bandwidth

83%  
Increase



Supporting RAS/Security  
based on Memory Controller

RAS  
Security



## Summary and Future Plan

- AI and pandemic drive demand for memory bandwidth and capacity, and new interconnect standard CXL allows expansion of memory
- Samsung developed the industry's first ASIC-based 512GB CXL memory expander, which will be available for early evaluation this quarter
- Memory intensive applications such as IMDB and AI/ML have been tested on CXL memory expander with an open-source software, SMDK
- Samsung to cooperate further on CXL 3.0 and beyond, and provide next-gen memory solutions like memory disaggregation, SDM\*, and more

\* Software-defined memory





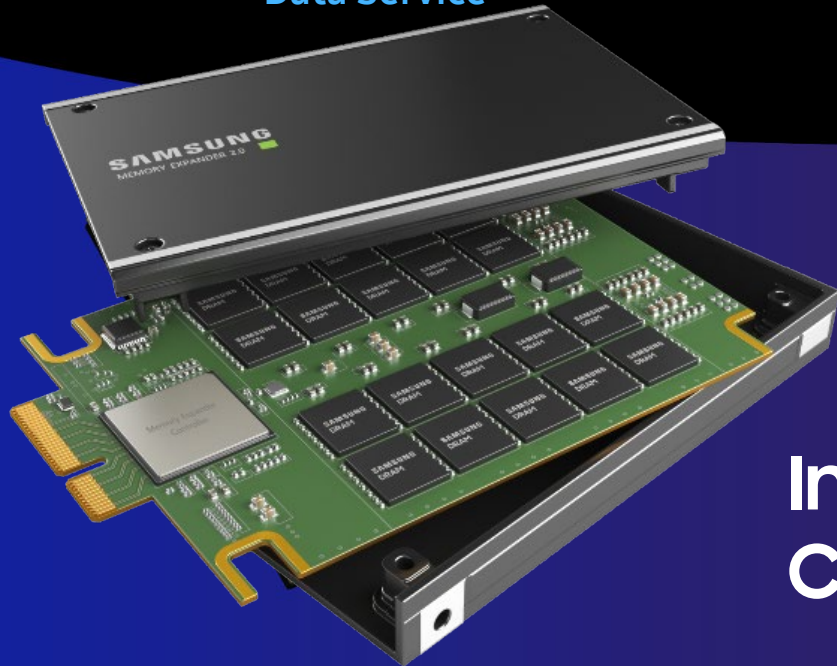
Enhanced  
Data Service



AI/ML  
NLP, Recommendation



Edge Computing



Industry First  
CXL™ Memory Expanders

**SAMSUNG**

# Appendix

## • Test Condition (ML/AI and IMDB)

### ML/AI

#### For BERT and Nasnet

TensorFlow (CPU) >= 1.11.0+, Python ~ 3.7, Numpy < 1.20.0

#### For BERT

Multi-process, 3 cores/process, batch-size:128,  
 max\_seq\_num:256, num-test-data/process: 512  
 dataset=CoLA  
 do\_train=true, do\_eval=true, data\_dir=\$GLUE\_DIR/CoLA  
 vocab\_file=vocab.txt  
 init\_checkpoint=\$BERT\_BASE\_DIR/bert\_model.ckpt  
 max\_seq\_length=128, train\_batch\_size=32  
 learning\_rate=2e-5, num\_train\_epochs=3.0

#### For NASNet

Multi-process, 3 cores/process, batch-size: 100,  
 eval\_image\_size:236, num-test-data/process: 200  
 dataset\_name=imagenet, num\_preprocessing\_threads=4  
 labels\_offset=0, model\_name=inception\_v3  
 preprocessing\_name=inception\_v3  
 moving\_average\_decay=None, quantize=False, use\_grayscale=False

### IMDB

#### For scale-up vs scale-out

```
Redis-server :master
cluster-enabled yes
cluster-node-timeout 300000
save ""
stop-writes-on-bgsave-error yes
rdbcompression yes
rdbchecksum yes
rdb-del-sync-files no
repl-diskless-sync no
rdb-del-sync-files no
replica-serve-stale-data yes
replica-read-only yes
repl-diskless-sync-delay 5
repl-diskless-load disabled
repl-disable-tcp-nodelay no
replica-priority 100
client-output-buffer-limit replica 0 0 0
maxclients 1000000
maxmemory-policy noeviction
maxmemory-samples 10
maxmemory-eviction-tenacity 100
repl-diskless-sync no #master-replica disk-based sync
rdb-del-sync-files no
replica-serve-stale-data yes
replica-read-only yes
replica-priority 100
client-output-buffer-limit replica 0 0 0
io-threads 4
io-threads-do-reads yes
```

```
Redis-server :replica
save ""
port 6380
replicaof 127.0.0.1 6379
replica-read-only yes
stop-writes-on-bgsave-error yes
rdbcompression yes
rdbchecksum yes
rdb-del-sync-files no
repl-diskless-sync no
rdb-del-sync-files no
replica-serve-stale-data yes
replica-read-only yes
repl-diskless-sync-delay 5
repl-diskless-load disabled
repl-disable-tcp-nodelay no
replica-priority 100
client-output-buffer-limit replica 0 0 0
maxclients 1000000
maxmemory-policy noeviction
maxmemory-samples 10
maxmemory-eviction-tenacity 100
replica-lazy-flush no
lazyfree-lazy-user-del no
lazyfree-lazy-user-flush no
oom-score-adj no
oom-score-adj-values 0 200 800
disable-thp yes
```