# Memory Centric 시대

## - SK Hynix Solution 방향 (Storage 중심)

2024.10
SKHY / Solution AT
Junghyun Joh

# AI Pipeline and Storage Workload

## AI Data Pipeline

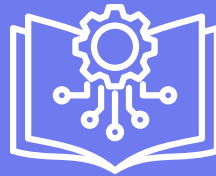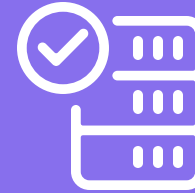| Data Ingestion | Data Preparation | Training | Checkpoint & Restore | Inference |
|---|---|---|---|---|
| Sequential Write | Mostly Sequential | Mostly Random | Mostly Sequential Read/Write | Mostly Sequential Read |

# Challenges to AI Deployment

## AI Datacenter



**Ingest** | **Prep** | **Training** | **Checkpoint** | **Inference**

**STORAGE**
Limited rack power & space, growing dataset

**COMPUTE**
Power optimized SSD as cache device

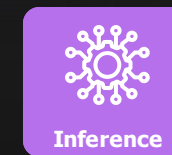## On-device AI



**Inference**

**USER EXPERIENCE**
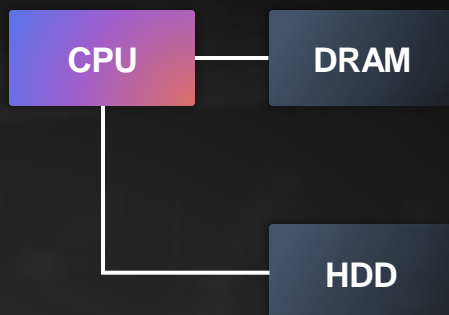AI model loading time and user experience
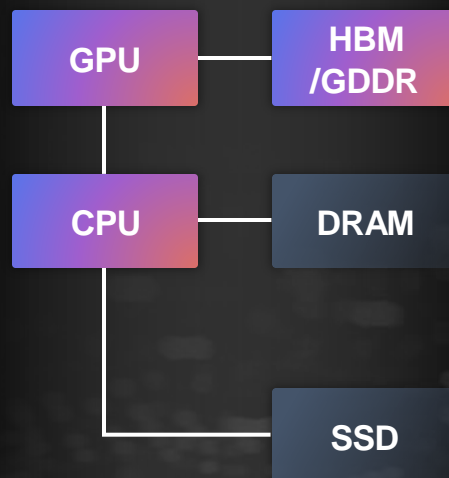
# Memory Centric

## Memory is system competitive

### Past System
Calculator : Data Simple Calculation

```
CPU ─── DRAM
  │
  └──── HDD
```

### Present System
Workable : Data/information Reading

```
GPU ─── HBM/GDDR
  │
CPU ─── DRAM
  │
  └──── SSD
```

### Future System
Thinkable AI : Data/information Creation

```
HBM/GDDR        DRAM        HBM/GDDR
  GPU                          PIM***
  ASIC          CPU
                             NPU**
  HBM/GDDR   Memory        Custom
             Expansion     Memory

  CXL

Memory Pool        Storage Pool
```

* CXL: Compute eXpress Link
** NPU: Neural Processing Unit
***PIM: Processing in Memory

# Contents

**SK** hynix

# Facing Problems

## Risk increasing due to Energy, Cost, Environmental issue and etc.



**Cloud Cover: Cloud Prices Rise as the Era of Generative AI Dawns**

NOVEMBER 29, 2023

By Jordan Galhardo-Burnett, Pankaj Sherawat, Peter Dankert,
Engelhar...
Scognam...

BCG

**Microsoft's and Google's AI plans clouded by concerns of rising costs**

Tech giants tout new tools that will need significant investment as the technology takes hold

FINANCIAL TIMES

**Meta's Costs Rise Rapidly as Zuckerberg Vows to Keep Spending on AI Arms Race**

Shares in the social-media company fell more than 12% after it revealed AI investment plans while reporting record revenue

By Salvador Rodriguez  Follow
Updated April 24, 2024 6:01 pm ET

THE WALL STREET JOURNAL.

**AI Is Pushing The World Toward An Energy Crisis**

Ariel Cohen  Contributor ⓘ
I cover energy, security, Europe, Russia/Eurasia & the Middle East

**Sam Altman Invests in Energy Startup Focused on AI Data Centers**

Investment by OpenAI CEO highlights artificial intelligence's electricity appetite

By Amrith Ramkumar  Follow
April 22, 2024 5:00 am ET

THE WALL STREET JOURNAL.

**AI boom sparks concern over Big Tech's water consumption**

Microsoft, Google and Meta are using more water to cool down data centres that power artificial intelligence products

FINANCIAL TIMES

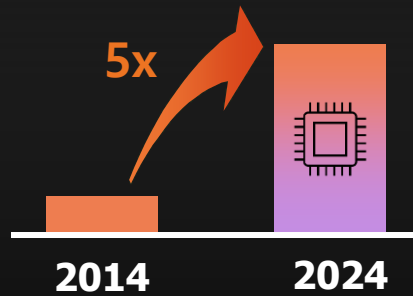**AI Is Accelerating the Loss of Our Scarcest Natural Resource: Water**

Cindy Gordon  Contributor ⓘ
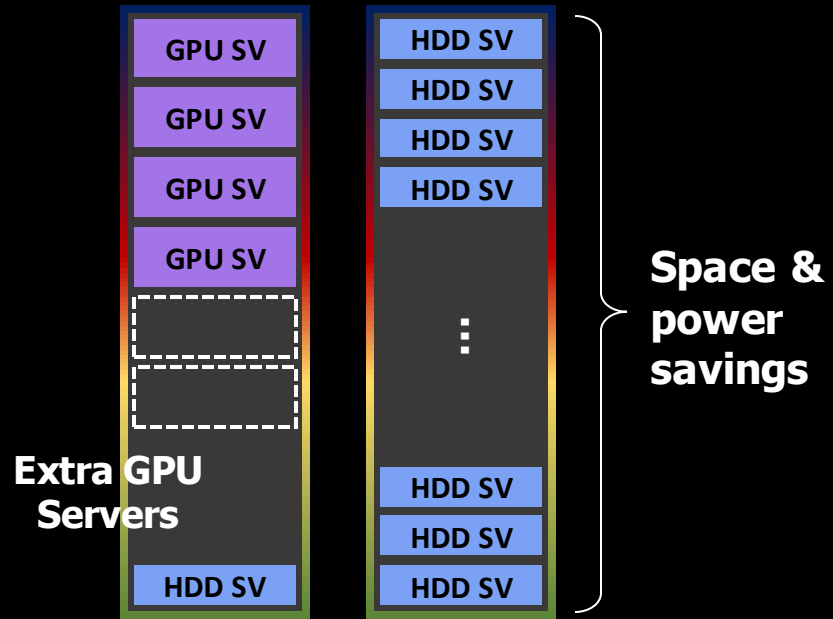CEO, Innovation Leader Passionate about Modernizing via AI
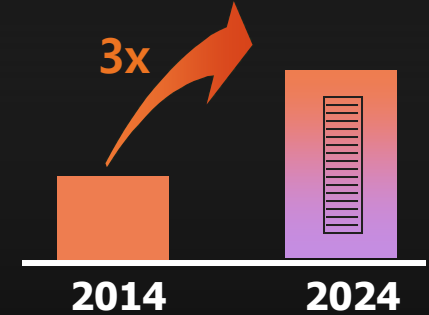
Forbes

# AI Datacenter Challenges: Storage

## In the AI datacenter:

**Avg. processor active power[1]**

**5x**

2014     2024

| GPU SV | | HDD SV |
|--------|--|--------|
| GPU SV | | HDD SV |
| GPU SV | | HDD SV |
| GPU SV | | HDD SV |

**Extra GPU Servers**

HDD SV

⋮

HDD SV

HDD SV

HDD SV

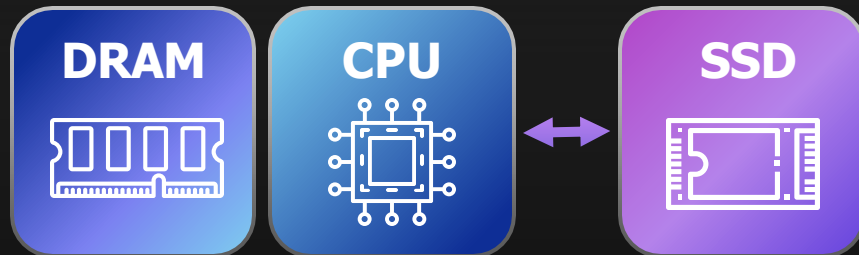**Space & power savings**

**Avg. rack power[1]**

**3x**

2014     2024

## Limited Datacenter Power and Space

1) Source: SDxCentral

# AI Datacenter Challenges: Compute
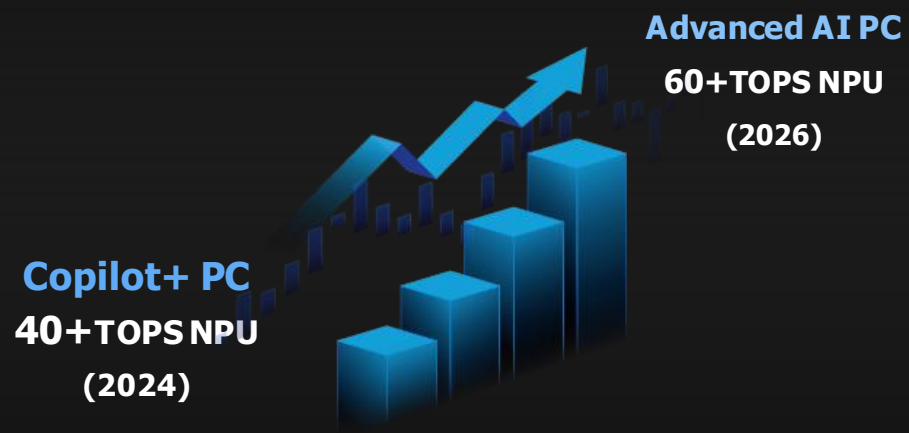


**General-purpose server for Data Prep (ETL)**

DRAM — CPU ↔ SSD

**GPU Server for Training/ Inference**

DRAM — CPU ↔ SSD
HBM — GPU

**Need Something with optimized power & high performance**

# On-device AI Challenges: PC/Mobile

## Computing Power

**Advanced AI PC**

**60+TOPS NPU**

**(2026)**

**Copilot+ PC**

**40+TOPS NPU**

**(2024)**

**High-performance SSD required**

## User Needs

**Minimal AI Model load time (<1s) from Storage to all memory states**

**Dirty**

**50%** **Random I/O**

**Clean**

**1 sec**

**Sequential Performance**

**Sequential + Random Performance**

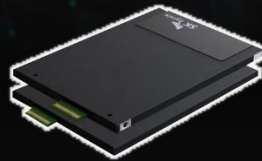(based on SK hynix mainstream Gen4 SSD)

# Contents

SK hynix

# GenAI-ready Solutions

## For AI Datacenter:

**PCIe Gen5 Enterprise SSD**
Best-in-class Performance, IOPS/W

**61TB QLC Enterprise SSD**
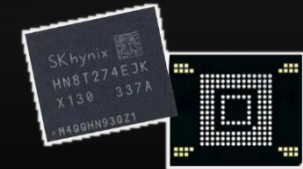World's highest-capacity PCIe SSD

## For On-device AI:

**PCIe Gen5 Client SSD**
World-first for mainstream client

**Zoned UFS**
World-first, vertical optimized mobile storage

# High-density QLC eSSD: vs. HDD

| | HDD | 30TB HDD | | 122TB QLC SSD | QLC SSD |
|---|---|---|---|---|---|
| | | 4 | # GPU servers configurable | 5 (+1 vs. HDD rack) | |
| | | 4U x 13ea | # Storage servers configurable | 2U x 7ea | |
| | | 39.78kW Compute 65% Storage 35% | Power Consumed | 36.7kW Compute 89% Storage 11% | |
| | | 2 Racks | # Racks required | 1 Rack | |

**QLC SSD = +1 GPU server (more compute) per rack; reduced power consumption & rack footprint**

# PCIe Gen5 cSSD: Gen-on-Gen Advanced

## Performance

**2x**

**1.4x**

Seq. Read    Ran. Read

*Gen4 vs. Gen5*

## Power Efficiency

**30% ↑**

*Gen4 vs. Gen5*

## AI Model[1] Launch

**50% faster**

LOADING 100% @ Clean

*Gen4 vs. Gen5*

1) AI Model: Llama2 7BQ8

# PCIe Gen6 & UFS 5.0 Solutions

## PCIe Gen6 SSD

**Performance[1]**

2x (Read)
1.5x (Write)

Read    Write

**Power Efficiency**

2x

Gen5    Gen6

## UFS 5.0[2]

**Random Performance**

5x

UFS4.0    UFS5.0

**Power Efficiency**

+40%

UFS4.0    UFS5.0

**Using SK hynix-proprietary, H-TPU™ architecture**

**1)** Performance bound by max. power: 25W on Gen5 and Gen6 (estimate)    2) Estimates

# Appendix] High-density QLC eSSD: Leading Capacities

## Specification

- PCIe 4.0, NVMe 1.4c, OCP 2.0

## Capacity

- 7.68TB, ~ 61.44TB
- 122TB upcoming, early '25

## Performance

- 7GB/s, 1M IOPs

**"D5-P5336" QLC eSSD**

U.2          E3.S

E1.L

# Appendix] PCIe Gen5 eSSD: BIC Performance

## Specification

- PCIe 5.0, NVMe 2.0, OCP 2.0

## Capacity

- 1.92TB ~ 15.36TB (Read-Intensive)
- 1.6TB ~ 12.8TB (Mixed-Use)

## Performance (Max.)

- Seq. R 14.5GB/s , Seq. W 9.3GB/s
- Ran. R 3,200K IOPS , Ran. W 400K IOPS
- Up to 12% better sequential write of vSAN HCI bench vs. competition

**"PS10x0" Compute eSSD**



**U.2/U.3**          **E3.S**

# Appendix] PCIe Gen5 cSSD: World's First Mainstream

## Specification

- PCIe 5.0, NVMe 2.0c
- HYPERWRITE™ Cache Technology

## Capacity

- 512GB, 1TB, 2TB

## Performance

- Seq. R 14GB/s (@10W)
- Seq. W 12GB/s (@10W)



"PCB01" mainstream cSSD

# Appendix] Zoned UFS: World's First for Advanced Mobile

## Specification

- UFS 4.0

## Capacity

- 512GB, 1TB

## ZUFS Spec

- Zone type SWR[1]

- # Max Open Zones: 6

## Vertical Optimized Mobile Storage



**Zoned UFS 4.0**

1) SWR : Sequential write required zone type

SK hynix

# Appendix] Zoned UFS: Optimized for Advanced Mobile

## F2FS File system



Conventional
UFS 4.0

Zoned
UFS 4.0

**Application Launch Time[1]**
0.45 x

**Memory Allocation Time[2]**
0.67 x

**Product Lifetime**
**40% longer**

1) Condition: long-hours use

2) Avg. of memory & storage backup/restore time

# Contents

SK hynix

# Complete AI Memory Solutions Provider

**Total Memory Solution Provider**

Cache

Main Memory

CXL

SSD / UFS

**MEMORY, THE POWER OF AI**

**Revolutionary Path**

**LPDDR AiM**

**CMS (Computational Memory Solution)**

**Evolutionary Path**

**PCIe Gen6 SSD**
**UFS 5.0**

**CSD (Computational Storage)**

# Memory Centric

Computing in All memory layer

CPU

C

C' Main Memory

C' CXL Memory

C' SSD

C' : computing unit

job distribution

data
Cache
data
DRAM
data
SCM
data
SSD

Data placement by data & workload characteristics

Processor design

● : computing element (CPU/GPU/FPGA/HWA…)

▬ : Data Storage (HDD/SSD/SCM/DRAM/Cache)

# Memory & Computing Fusion

## Various Memory + Computing Solution from AiM to CSD

**Workload density**  **Workload**

**Memory Die level (AiM)**

Bank
C → C →
C → C →

- Computing in memory die
- High density workload

Transformer

**Memory Card level (CMS)**

DDR Devices
C

- Computing in memory card
- Low density workload

Embedding

Data Analytics

**Storage level (CSD)**

NAND Devices
C

- Computing in SSD
- Low density & high capacity workload

Data Analytics

# Appendix] Processing in Memory

- Even with multi batch processing, due to the attention processing for long context, there is huge portion of memory intensive function

## LLM block architecture



* FC : Fully connected layer

■ : FC weight GEMV
- In Multi batch, weight data can be reused
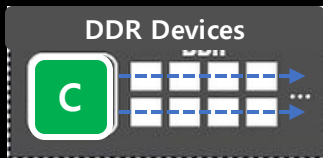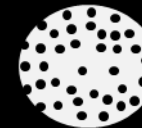→ computing intensive & high GPU utilization

■ : Attention token GEMV
- Even in Multi batch, token data cannot be reused
→ still memory intensive & low GPU utilization

■ : Misc. functions
- Regardless batch
→ computing intensive & high GPU utilization

## Processing time portion

1 Batch, 2K token, MHA

| 20% | 50% | 30% |

GPT3-13B, Last token processing time

32 Batch, 100K token, GQA(8)

| ~1% | 99% |

- - - : computing intensive

- - - : memory intensive

**LLM service trend : longer context**

# Appendix] AiMX system from DC to Edge device LLM

- To provide different level of customer experience and save your operating cost

| Datacenter LLM | Edge device LLM |
|---|---|

Larger LLM size
Multi-batch
Longer context

**99%**

**75%**

Smaller LLM size
Single-batch
Long context

**Both system has memory intensive portion as a majority**

Boosting your Datacenter LLM
by adding AiMX as an attention accelerator

Boosting your Edge device LLM
by replacing memory with AiM



+

x 4
**Performance**

33%
**Energy consumption**

AiM AiM AiM AiM
AiM CTRL  **Mobile AP**  AiM CTRL
AiM AiM AiM AiM

**Power Consumption**

Baseline    Maintain Power    Max Performance
■ Memory  ■ AP  ■ Others

**Execution time**

Baseline    Maintain Power    Max Performance
■ memory intensive  ■ compute intensive

# Appendix] Computational Memory Solution

## CMS



**2** CXL

Core **1** Mem Mem Mem Mem / Mem Mem Mem Mem

| 1 | **Performance Improvement** |
|---|---|
| | By utilizing the higher Bandwidth inside the CXL module |

| 2 | **Energy Efficiency Improvement** |
|---|---|
| | By minimizing data movement between host and CXL module |

## Use case

### Recommendation Model



Compute-Intensive Layer — Host : x86 Server
Top MLP
Feature Interaction
Bottom MLP — Embedding • • • Embedding
PCIe
CMS Card  CMS Card  CMS Card  CMS Card
Memory-Intensive Layer (GB ~ TB) — CMS card : FPGA prototype

### Lightning DB for LLM



Search Keywords/Sentences — Vision-Language Model
GPU
2. Query inference — 1. Image inference
User
Vision data
Lightning DB
CMS
4. Sorted similar images — 3. Search similar images
SSD  SSD
Datalake

# Appendix] CXL™ Computational Memory Solution (CMS)

## CMS 2.0



- **CXL Memory Expansion**
- **Near-memory Processing**
- **Reduce data movement between CPU & CXL**
- **Free up CPUs to do other useful work**

## Data Analytics Performance

### Perf. & CPU Utilization

CPU Utilization
91.4%

15.9%

80%

CPU Only

2-CMS

### Energy Efficiency

22%

CPU Only

2-CMS

*Benchmark : TPC-DS Q28 (Dataset : 38.67GB)*

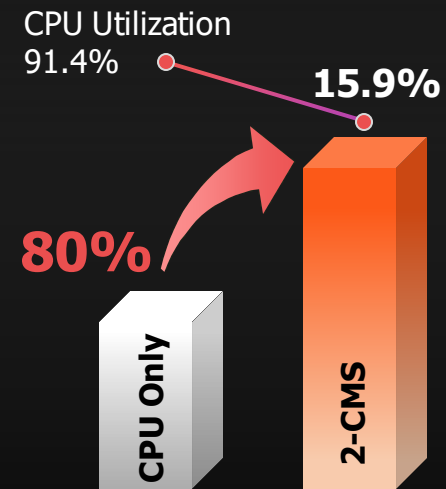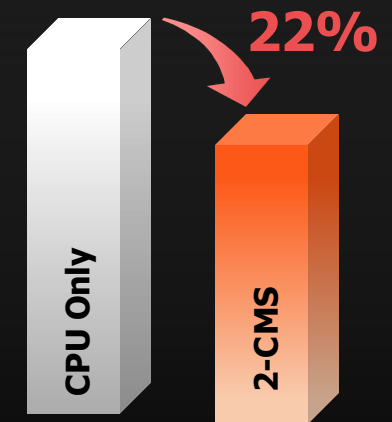# Appendix] Computational Storage Device

## CSD



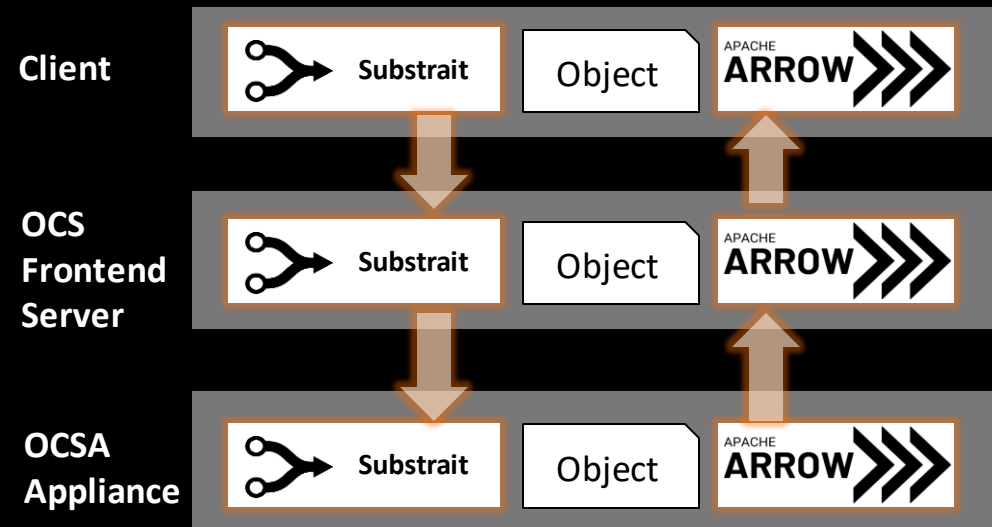| | |
|---|---|
| **1** | **Performance Improvement**<br>By avoiding interface bottleneck |
| **2** | **Energy Efficiency Improvement**<br>By minimizing data movement |
| **3** | **Reduce CPU overhead**<br>By offloading tasks from Host CPU |

## Use case

**Object based Computational Storage for Big Data Analysis**

**SK** hynix

# Appendix] Object-based Computational Storage (OCS)

## Speed and efficiency from data-awareness

**Compute Node**

| Schedule | Join | Filtering |

Object ID(Key) ... Object ID(Key)

**Offloading**

**Filtering**

Object ... Object

Object ID

**Metadata**

Data

**Object based Computational Storage**

---

**Enhanced Data Analytics**

Base Analytics(Legacy)

**3.8x Faster**

OCS

**Execution Time [sec]**

---

**Saving Host Resource**

**99.9% Reducing**

**Data Movement [MB]**

# Memory Centric

Artificial Intelligence

Think like a humans

Machine Learning

Learn like a humans

Deep Learning

Decide like a humans

Transformer

chatGPT

LLM

3rd Boom

Keep Boom up?

3rd winter?

??

AI Service Cost
Energy Problem
Environment Issue
...

Energy Efficiency
(SK hynix Solution)

2nd Boom

1st Boom

1st winter

2nd winter

1950's          1970's          1990's          2010's                    2030's

# Total Memory & Storage Portfolio



HBM3E

AiM

LPDDR5T

NAND

Mainstream cSSD

Boot eSSD

CMM DDR5 CMM MDS

DDR5 RDIMM

PCIe Gen5 16ch eSSD

Value TLC cSSD

Value PCIe eSSD

CXL CMS

AiMx

MCRDIMM

UFS / ZUFS / uMCP

Value QLC cSSD

High Capacity QLC eSSD

SATA eSSD

SK hynix

# End of Document