# Exploring Innovations in Storage &

# Memory Compression at Hyperscale

Presenter:
John Kim, Sr. Director, AI Memory Solutions, SK Hynix
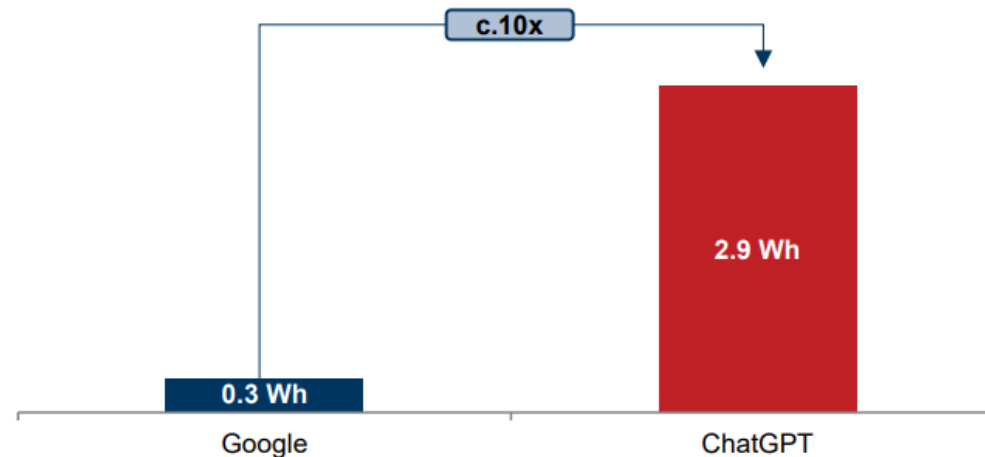Nilesh Shah, VP Business Development, ZeroPoint Technologies

**FMS**

*the Future of Memory and Storage*

# Problem: The AI Energy Efficiency Challenge

| Schneider Electric estimate | 2023 | 2028 |
|---|---|---|
| Total data center power consumption | 57 GW | 93 GW |
| AI power consumption | 4.5 GW | 14.0-18.7 GW |
| AI power consumption (% of total) | 8% | 15-20% |
| AI workload (Training vs Inference) | 20% Training, 80% Inference | 15% Training, 85% Inference |

*Source: Schneider Electric White Paper 110, The AI Disruption: Challenges and Guidance for Data Center Design*

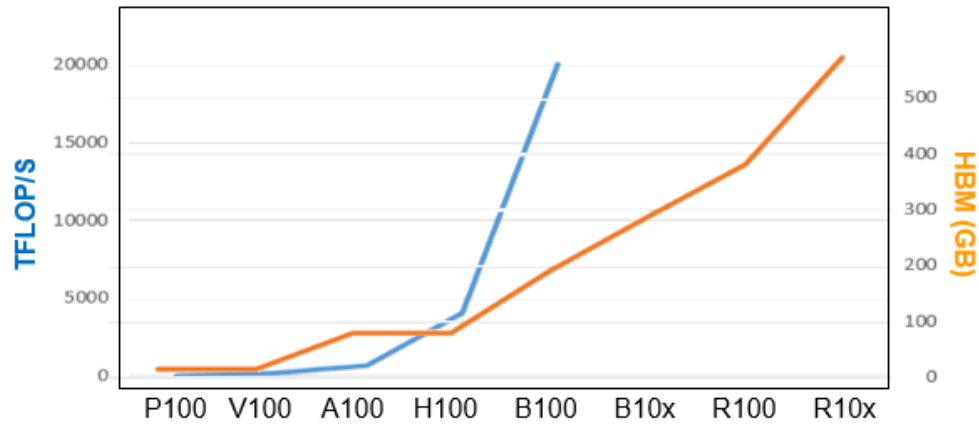ChatGPT queries are 10x as power intensive as Google searches
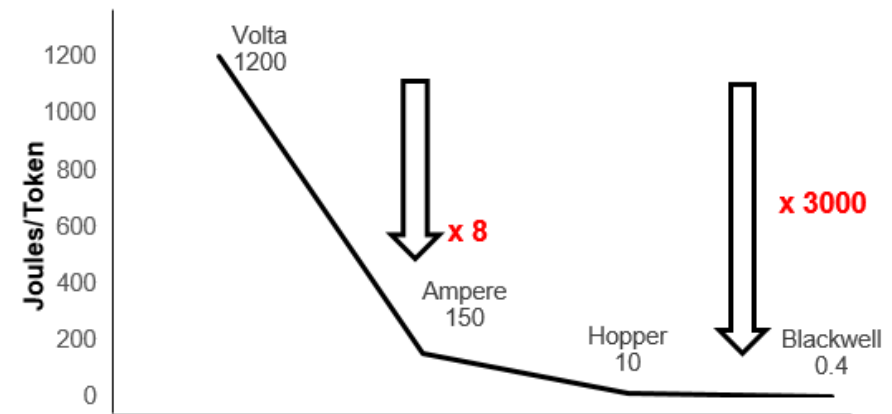Power consumption per query/search, Watt-hour (Wh)

c.10x

2.9 Wh

0.3 Wh

Google

ChatGPT

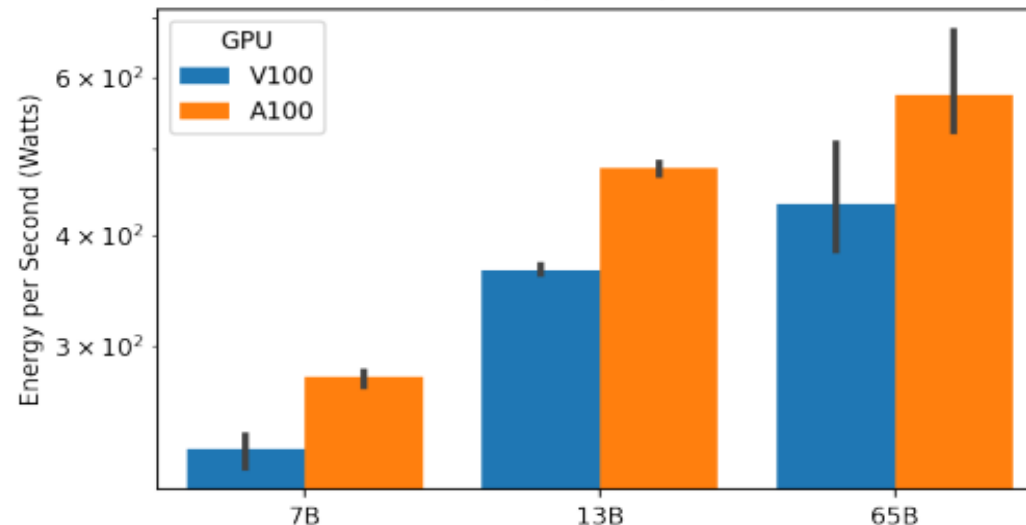*Source: Google, SemiAnalysis, compiled by Goldman Sachs GIR.*

# Can GPU solve the energy program?



Nvidia GPU Generation



Nvidia GPU Inferencing Energy (GPT-1.8T)



Llama 7B/13B/65B inferencing Power for fixed # of Tokens

TDP
V100: 300W
A100: 400W

Transformer Inference
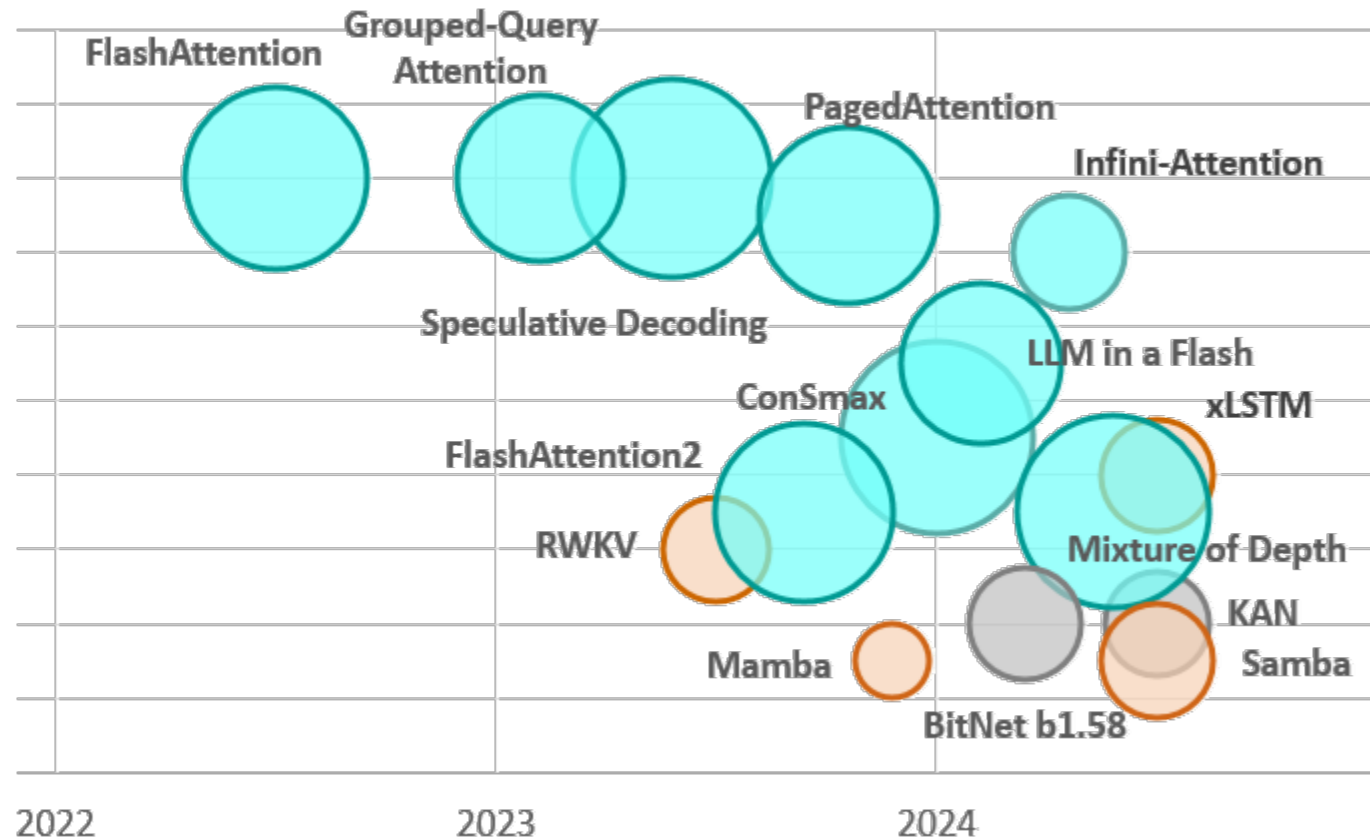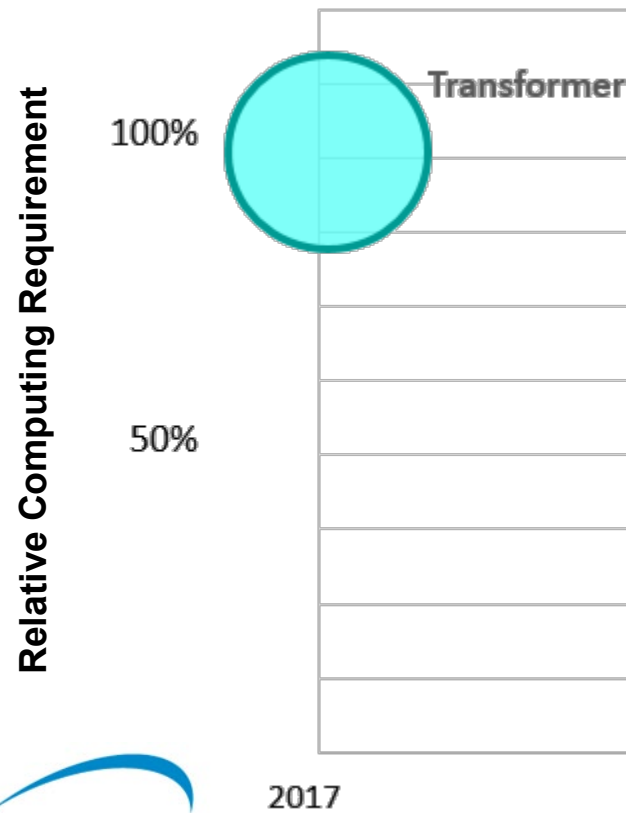Arithmetic Intensity is very low
→ Memory bound

*Source: Samsi et al (MIT). IEEE HPEC`23*

1

# AI Algorithm Evolution



*The circle size is relative memory requirement

# Energy Efficient Memory Solutions

- Compute-In-Memory

- Process-In-Memory

- Processing-Near-Memory (Computational Memory System)

# Compute-In-Memory

Computation happens in the memory cells

```
┌──────────────────────────┐    ┌──────────────────────────────┐
│   Compute-In-Memory      │────│  Compute-In New Memory (NVM)   │
│        (CIM)             │    │   (ReRAM, MRAM, FeRAM…)        │
└──────────────────────────┘    └──────────────────────────────┘
```

High cell density & capacity
Bit-serial computation (Limited OP/s)
Massive parallel operation

```
┌──────────────────┐         ┌──────────────────┐
│  Compute In DRAM │         │  Compute In SRAM │
└──────────────────┘         └──────────────────┘
```

Low cell density (reached its limit)
Limited capacity
Great TOPS/W (~ 10X than GPU)

```
┌──────────┐ ┌──────────┐   ┌──────────┐ ┌──────────┐
│ Digital  │ │ Analog   │   │ Digital  │ │ Analog   │
└──────────┘ └──────────┘   └──────────┘ └──────────┘
```

**CIM**

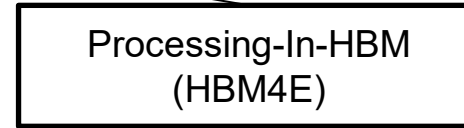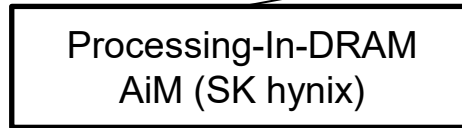| *Computing Cell Array* | *Computing Cell Array* |
| *Computing Cell Array* | *Computing Cell Array* |

PHY

PHY

*Memory Controller*

1

# Process-In-Memory (PIM)

Computation happens in the memory package

```
Processing-In-Memory
(PIM)
```

Monolithic Die
Limited PU area & lower cell density

```
Processing-In-DRAM          Processing-In-HBM
AiM (SK hynix)              (HBM4E)
```

3D Packaging (Hybrid Bonding)
Thermal issues
Connection issues



DRAM Die

Logic Die (Processing Unit)

**PIM**

Memory Cell Array — Compute Unit
Memory Cell Array — Compute Unit
Memory Cell Array — Compute Unit
Memory Cell Array — Compute Unit

PHY

PHY

Memory Controller

# Computational Memory System



**CXL Memory/Storage**

Memory Module — PHY — PHY — Memory Module

Memory Module — PHY — PHY — Memory Module

PHY

**Computing Unit**

CXL I/F

CXL I/F

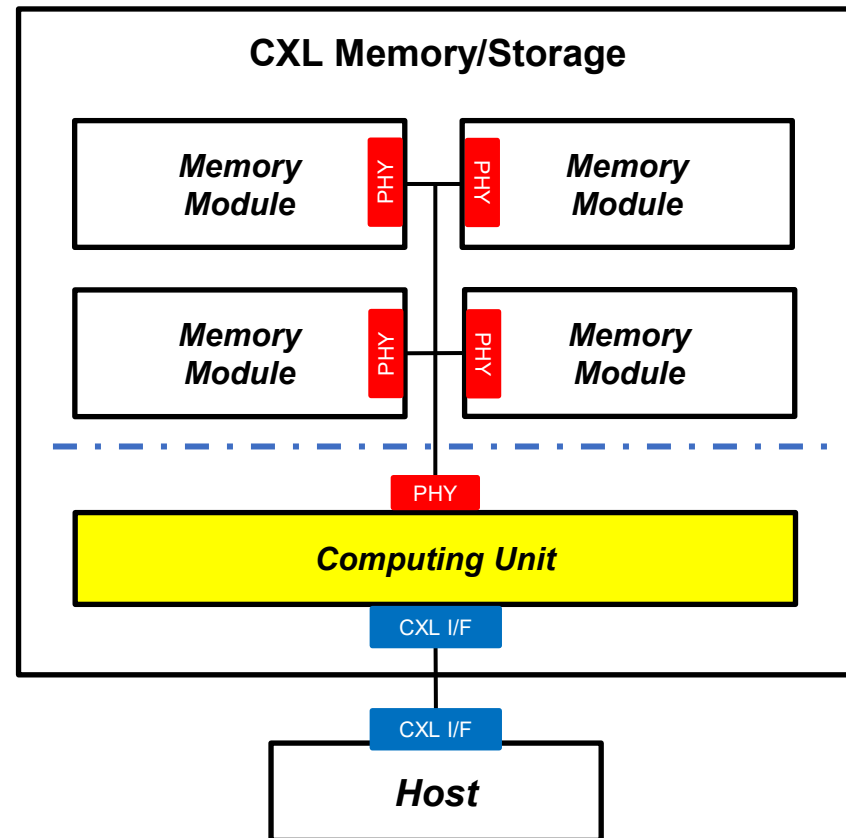**Host**

Bandwidth Improvement
- (De)compression
- Model training computation

# Hyperscale Data Challenge/ compression

Hyperscalers spending significant $$ on software based compression

Hyperscaler requirement: **hardware accelerated compression** is a **MUST-have**

**CPU cycles used for compression**

**4.6%**

**3%**



Hyperscale CXL Tiered Memory Expander Specification

Revision 1

Version 1.0
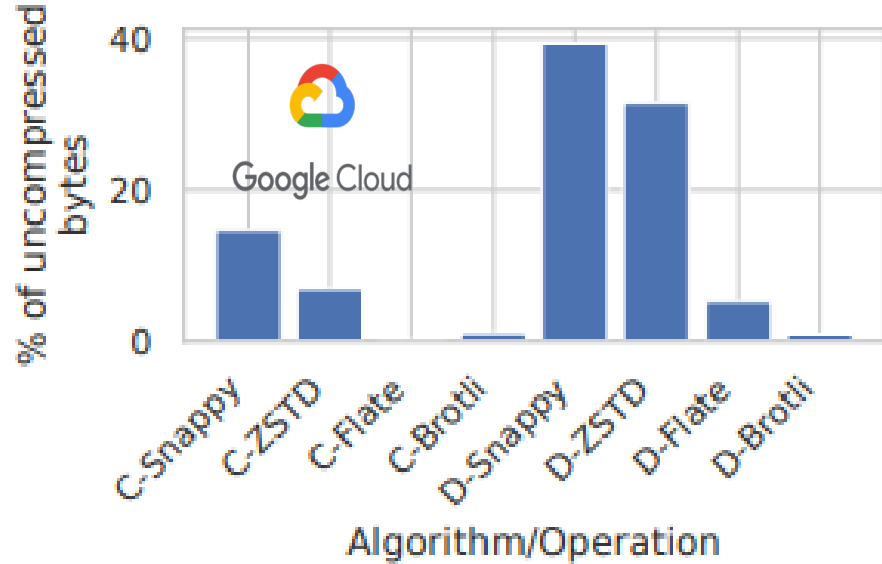**Base Specification Template v1.2**
Effective October 27, 2023
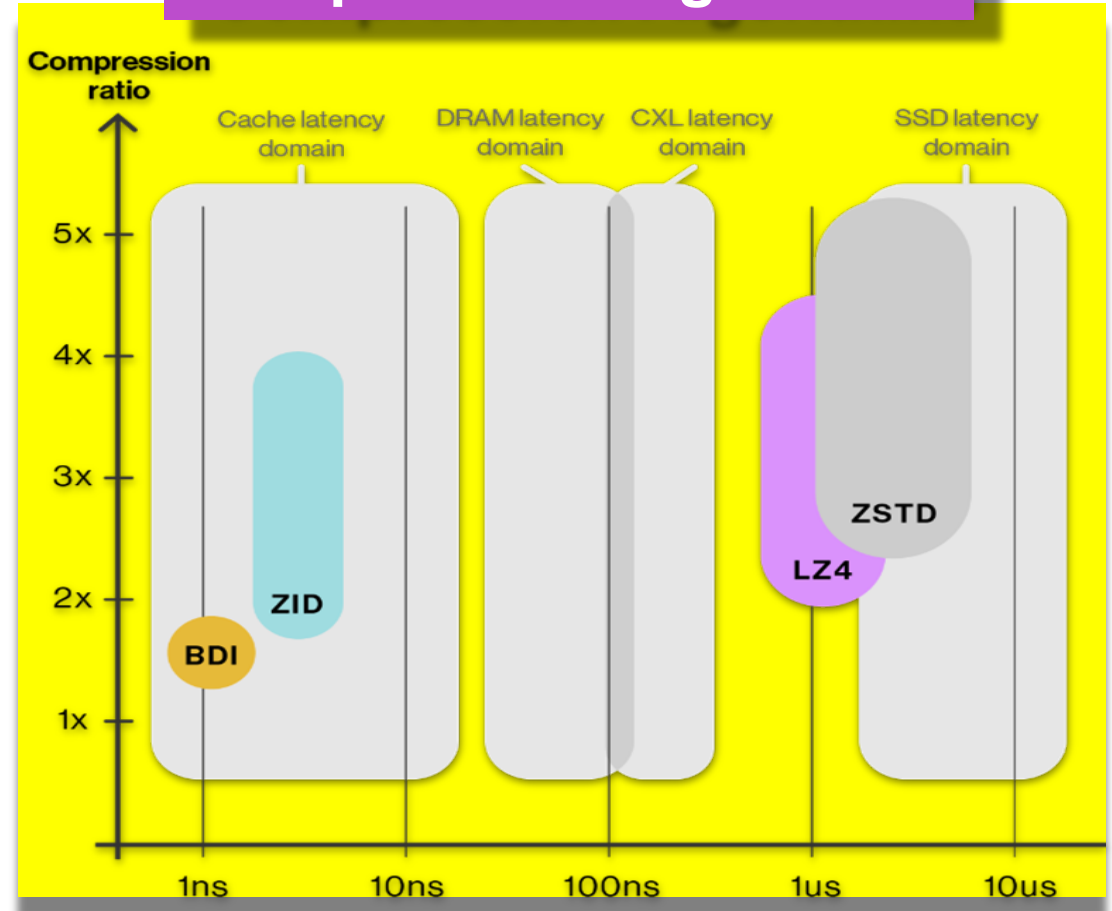
**OCP Hyperscale CXL Tiered Memory Expander Spec**

# Compression Algorithms: Memory, storage

Google Fleetwide De(compression)



(a) Fleet-wide uncompressed bytes handled by (de)compression, broken down by algo.

**Compression Algorithms**



LZ4, ZSTD IP

**Room for Innovation**

**Legacy**

https://www.zeropoint-tech.com/products/nvme-expansion-flash-mx

1

# Case Study: Implementations
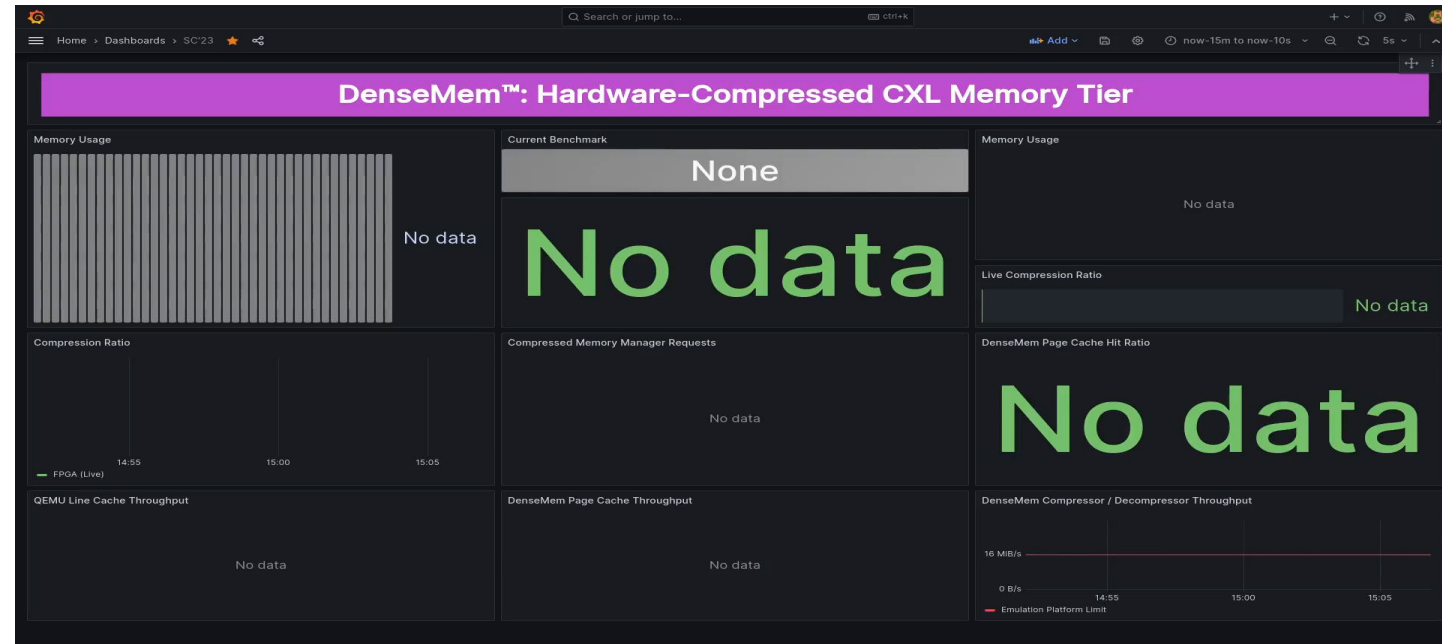


Intel QAT: CPU integrated (De)compression engine



Nvidia Blackwell : 800GB/s embedded Decompression engine

# Solutions, Challenges and considerations

Compaction, Software Transparent
   Address Translation – Dynamically
   Adjusting Capacity

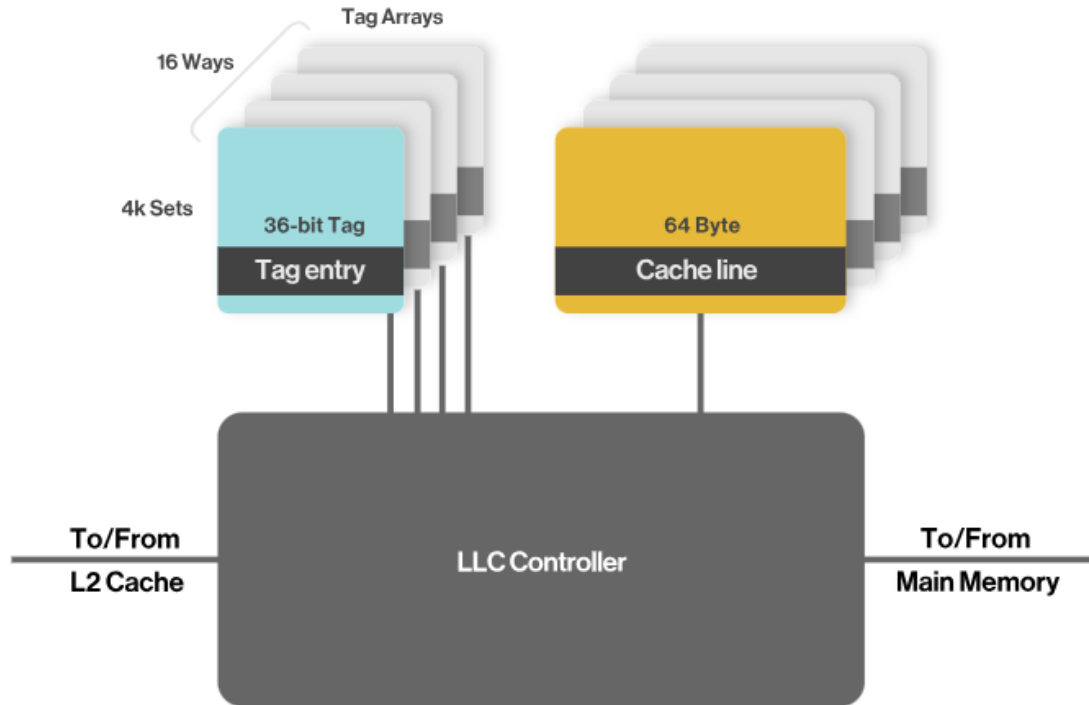Computational Programming
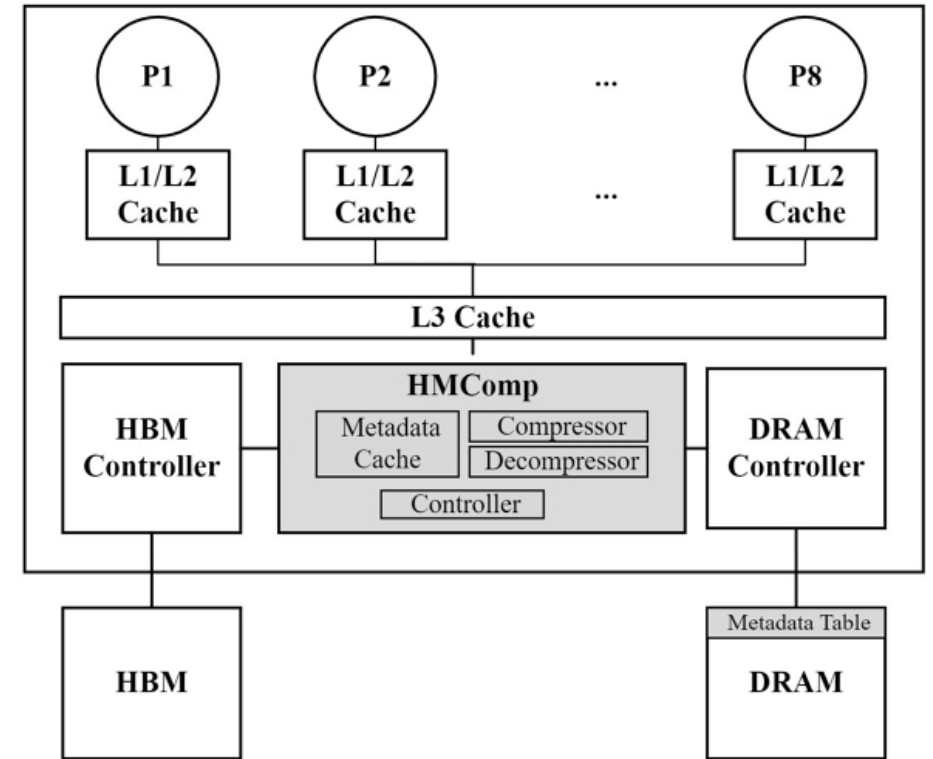
Standardized Data Movers



SNIA: Data Mover Spec

OCP CMS Computational Programming Work Group

# Future Trends / opportunities



CacheMX Cache (De)compressor IP



HMComp: Extending Near-Memory Capacity using Compression in Hybrid Memory

**On chip SRAM cache compression**

**HBM Compression**

# Summary Call to Action

- Joint Collaboration to accelerate Energy Efficient Solution deployment

- Start with Low hanging fruit : (De)Compression

- Prove out other use cases, influence/ shape standards