# CXL® Attached Flash Memory Economics

**Mahinder Saluja**
**Director of Strategy, SSD BU**
**KIOXIA America, Inc.**

CXL and Compute Express Link are registered trademarks of the Compute Express Link Consortium, Inc.
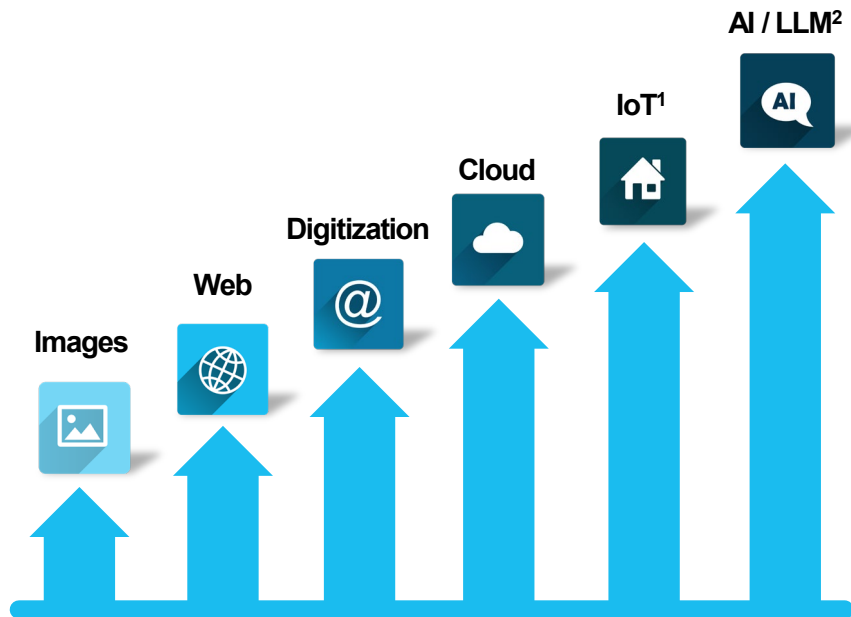
**KIOXIA**

# Agenda

- **Why Memory Expansion?**

- **CXL® Attached Flash Memory**

- **Performance and Cost**

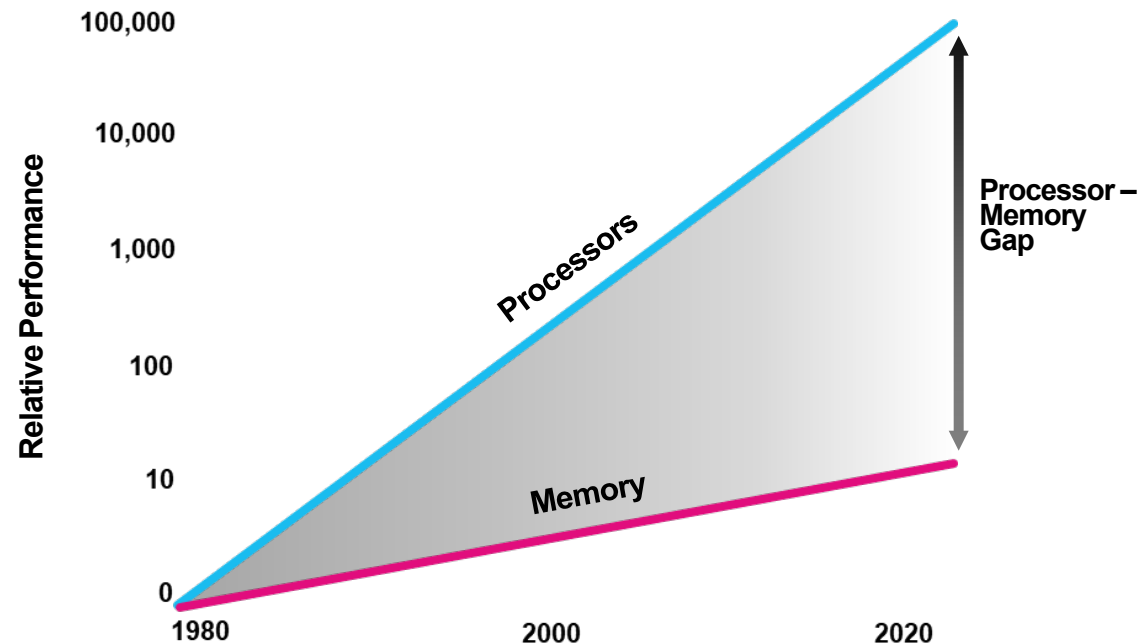- **Challenges and Opportunities**

# Why Memory Expansion?

**Today's data center workloads require increased memory capacity and bandwidth.**

**Central processing unit (CPU) / graphics processing unit (GPU) core idles due to gap in density growth compared to memory capacity and bandwidth**

- Memory ~40%-50%* of total server cost

AI / LLM[2]

IoT[1]

Cloud

Digitization

Web

Images



Relative Performance

100,000

10,000

1,000

100

10

0

Processors

Memory

Processor – Memory Gap

1980    2000    2020

[1] that Internet of Things (IoT)  [2] Artificial Intelligence (AI) / Large Language Models (LLM)

**Image source**: created by KIOXIA

KIOXIA

# Memory Expansion with Flash

- **CXL® Benefits**
  - Cost effective memory capacity and bandwidth expansion
  - Enables memory pooling and sharing with DRAM
  - Abstracts memory media interface

- **CXL® Technology Creates the Perfect Opportunity to…**
  - Explore alternative to costly DRAM
  - Flash media to jump over the semantic wall

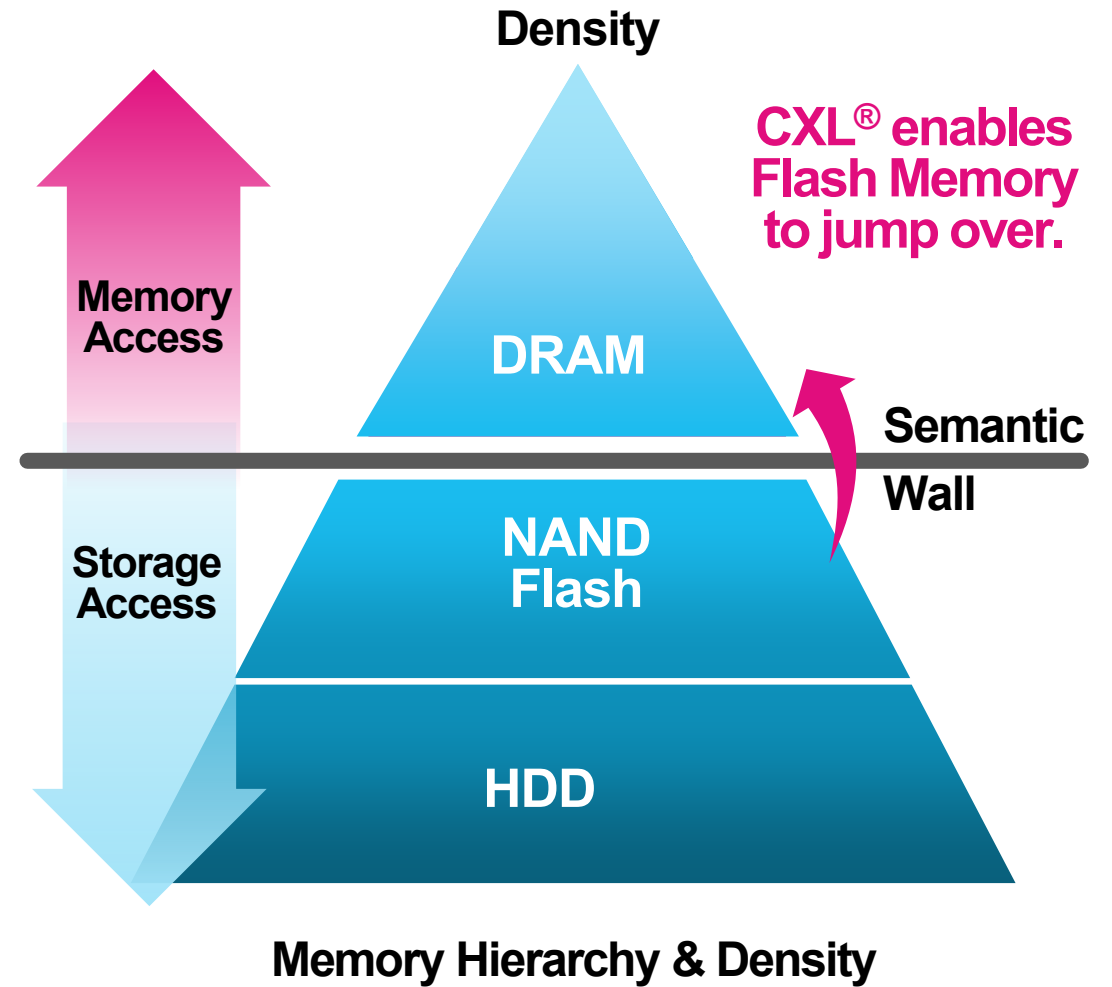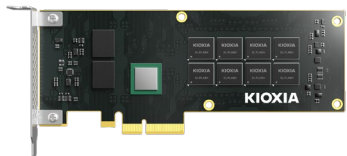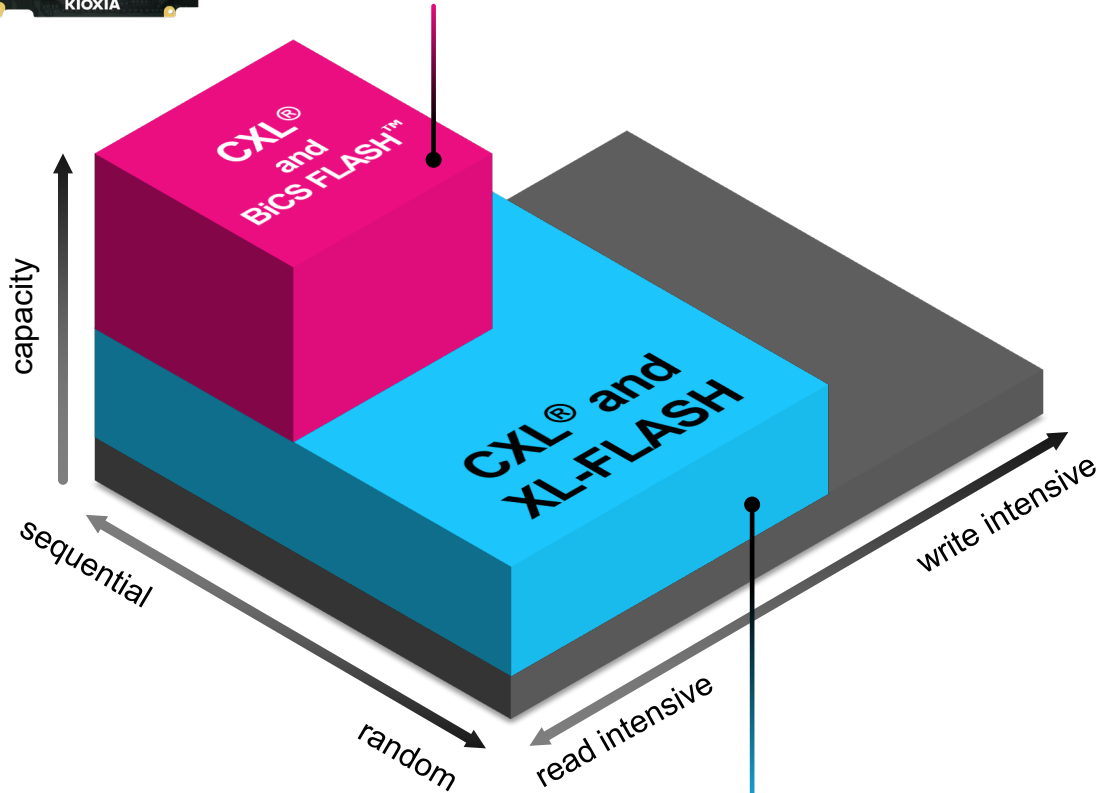**But can flash media jump over the sematic wall and reduce cost?**

**Density**

**Memory Access**

**Storage Access**

**DRAM**

**NAND Flash**

**HDD**

**CXL® enables Flash Memory to jump over.**

**Semantic Wall**

**Memory Hierarchy & Density**

Image source: created by KIOXIA

# KIOXIA CXL® Flash Memory Exploration

**Combination of CXL® and BiCS FLASH™ technologies**
Read intensive high capacity and high bandwidth memory

capacity

CXL®
and
BiCS FLASH™

CXL® and XL-FLASH

sequential

random

read intensive

write intensive

**Combination of CXL® and XL-FLASH technologies**
Random read/write access memory expansion

|  | CXL® and XL-FLASH Technologies | CXL® and BiCS FLASH™ Technologies |
|---|---|---|
| **Media** | BiCS FLASH™ (XL-FLASH) | BiCS FLASH™ |
| **Value Pillar** | Low latency (<5us (single-level cell), <10us (multi-level cell);  DRAM cache tier) | High bandwidth and High capacity |
| **CXL Access** | CXL.mem, CXL.io | CXL.mem,CXL.io |
| **Capacity** | >256 gigabytes (GB) | > 1 terabytes (TB) |
| **Suitable Applications** | Artificial intelligence (AI) / machine learning (ML) inference, In-memory data bases (DB), graph processing, cache, tiering | AI/ML training & inference, big data analytics |

**Image and table source**: *created by KIOXIA*

# CXL® and BiCS FLASH™ Application

## Generative Artificial Intelligence (AI) Inference & Training

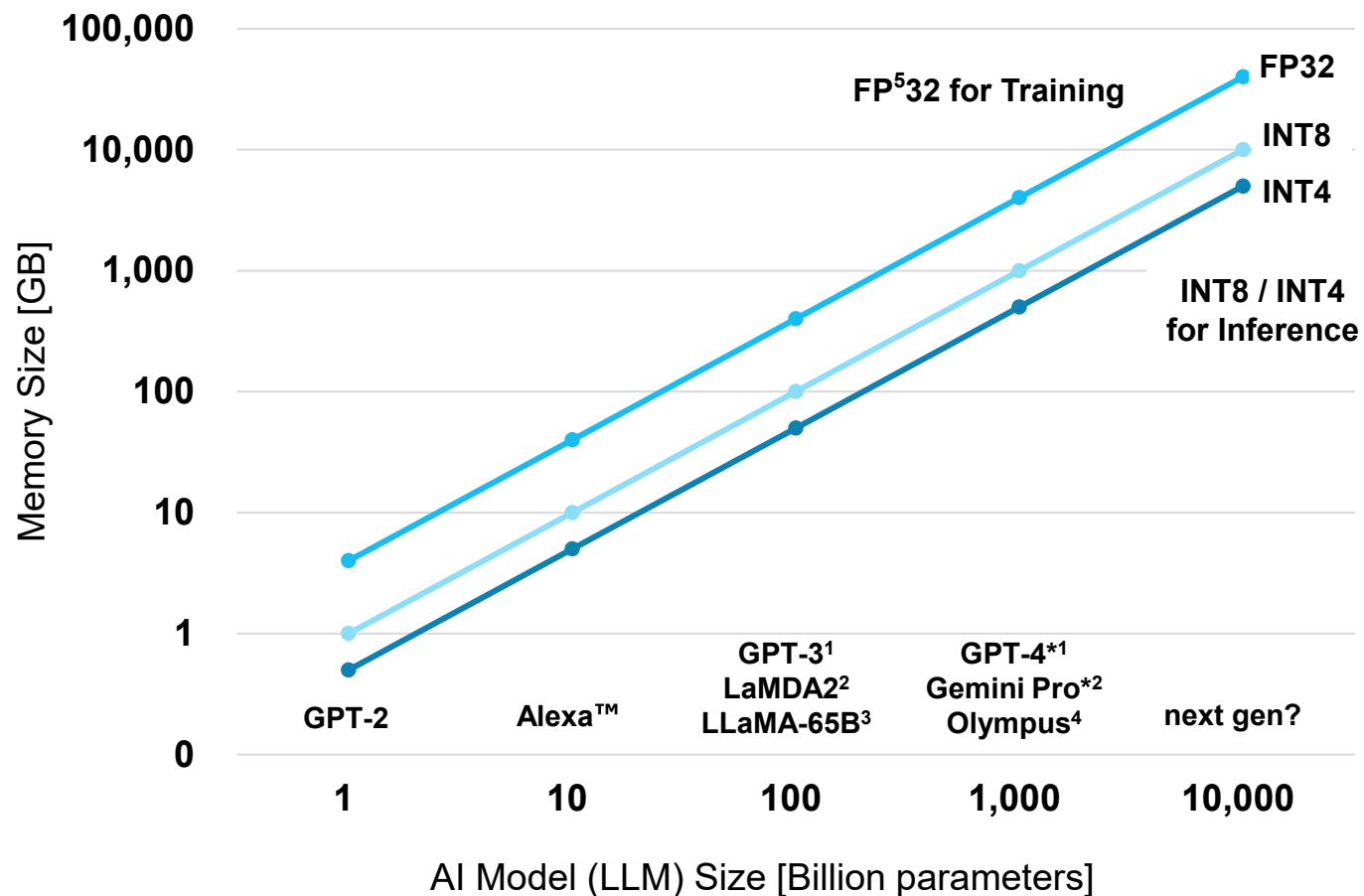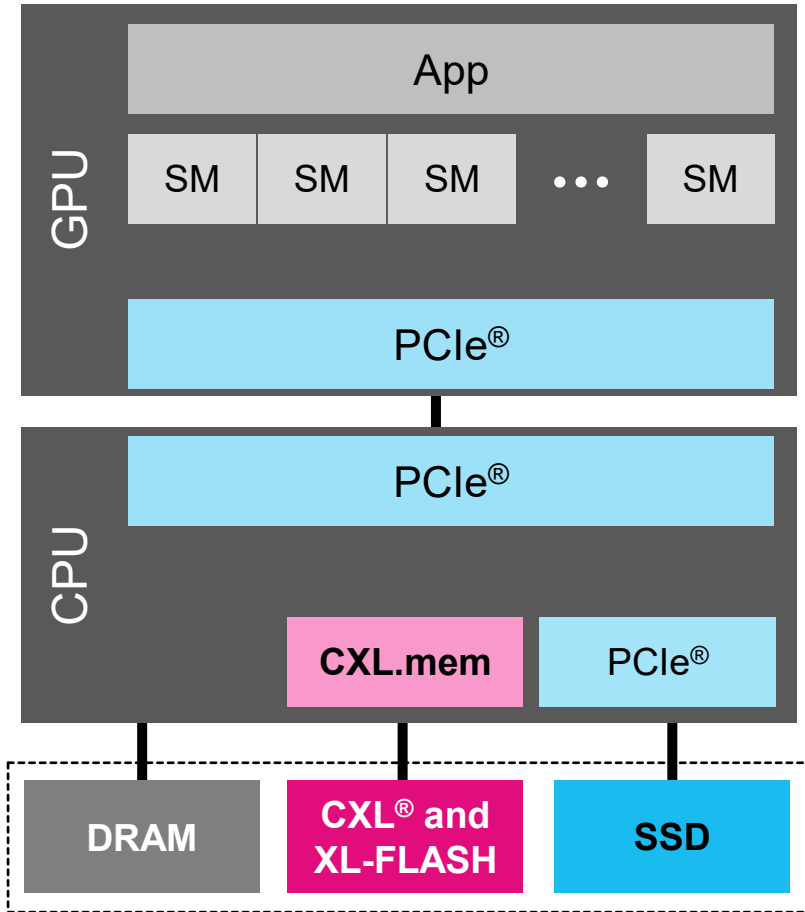Large Language Models (LLMs) require large memory to refill HBM DRAM on GPU

| GPU/xPU with HBM DRAM |
| --- |

| DRAM |
| --- |

**Offload AI Models & Training Data**

| NAND Flash with CXL® Interface |
| --- |

### Memory Requirements by AI Model Size



FP[5]32 for Training

FP32

INT8

INT4

INT8 / INT4 for Inference

Memory Size [GB]

100,000

10,000

1,000

100

10

1

0

GPT-2

Alexa™

GPT-3[1]
LaMDA2[2]
LLaMA-65B[3]

GPT-4*[1]
Gemini Pro*[2]
Olympus[4]

next gen?

1    10    100    1,000    10,000

AI Model (LLM) Size [Billion parameters]

Image and graph source: created by KIOXIA

* Some model sizes are estimated by KIOXIA

KIOXIA

# GPU Graph Processing with Low Latency Flash

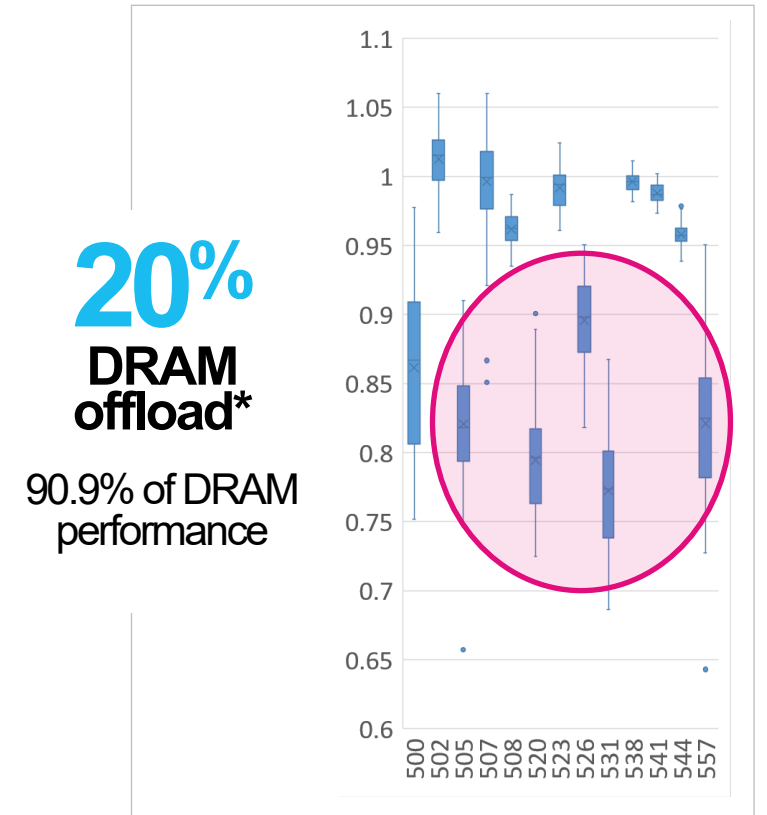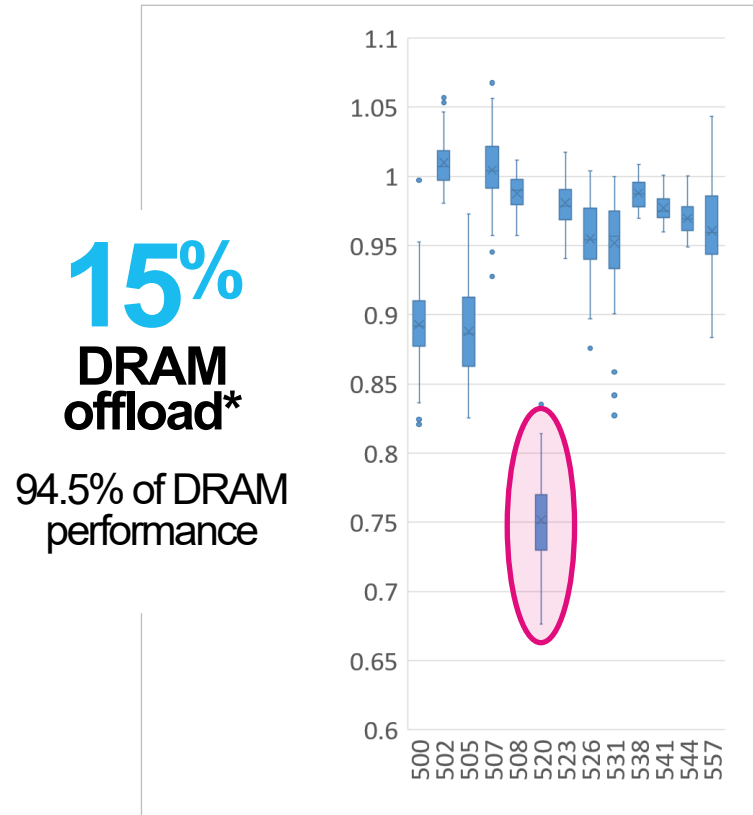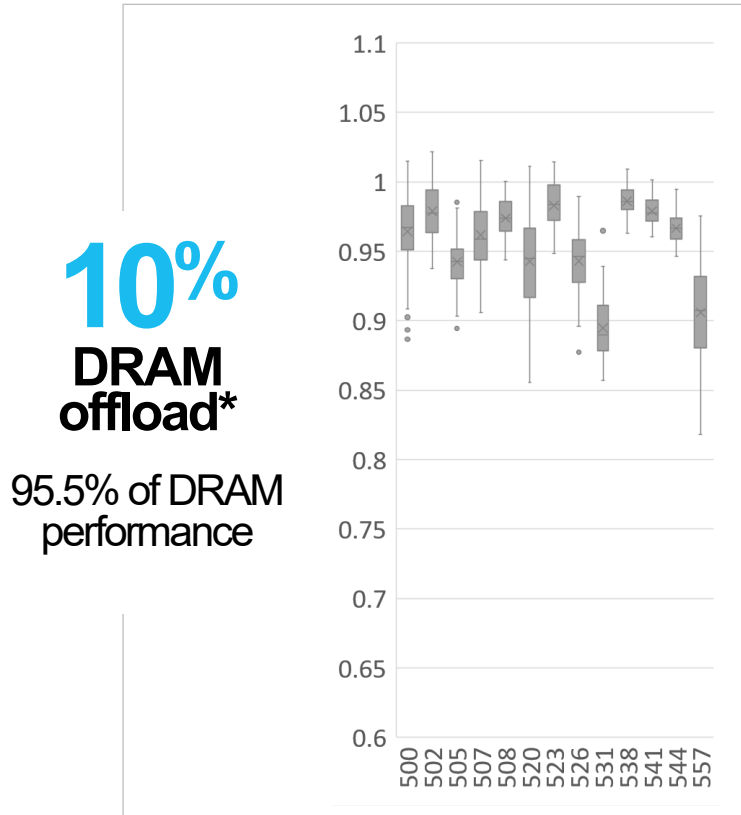

External Memory for GPU

## BFS (Breadth First Search)



GPU Graph Processing on CXL®-Based Microsecond-Latency External Memory (SC23 MSTA)

**Low latency XL-FLASH with cache can deliver DRAM-like application performance.**

Image and graph source: created by KIOXIA

# SPEC CPU® Benchmark with Low Latency Flash

**SPEC CPU:** SPEC CPU is a benchmark suite designed to measure and compare the performance of CPUs, memory subsystems, and compilers through a series of compute-intensive tests.

**Test Setup**: 32 copies x 4 hours, 10%, 15% & 20% offload

**10%**
**DRAM offload***

95.5% of DRAM performance

**15%**
**DRAM offload***

94.5% of DRAM performance

**20%**
**DRAM offload***

90.9% of DRAM performance

*FPGA emu. for memory offload

Image and graph source: created by KIOXIA

**Low Latency XL-FLASH can offload memory with nominal performance degradation.**
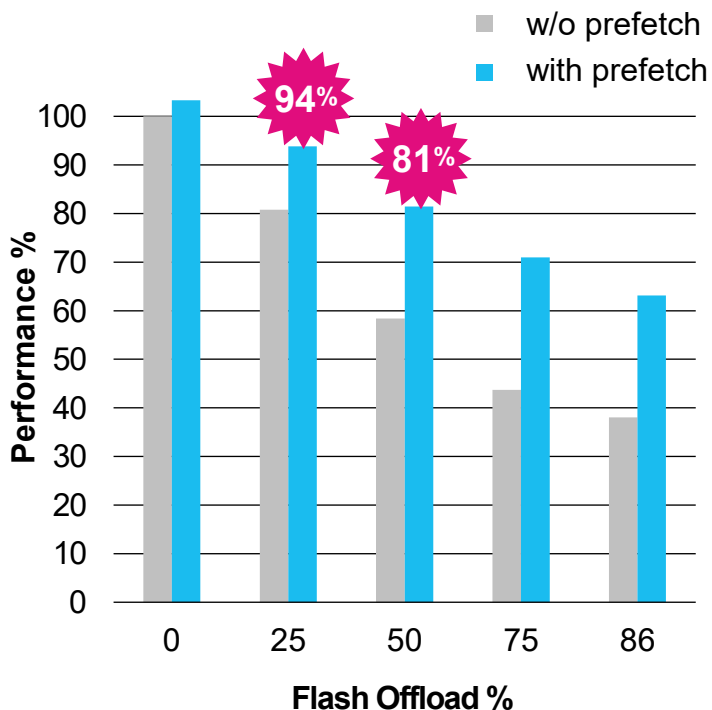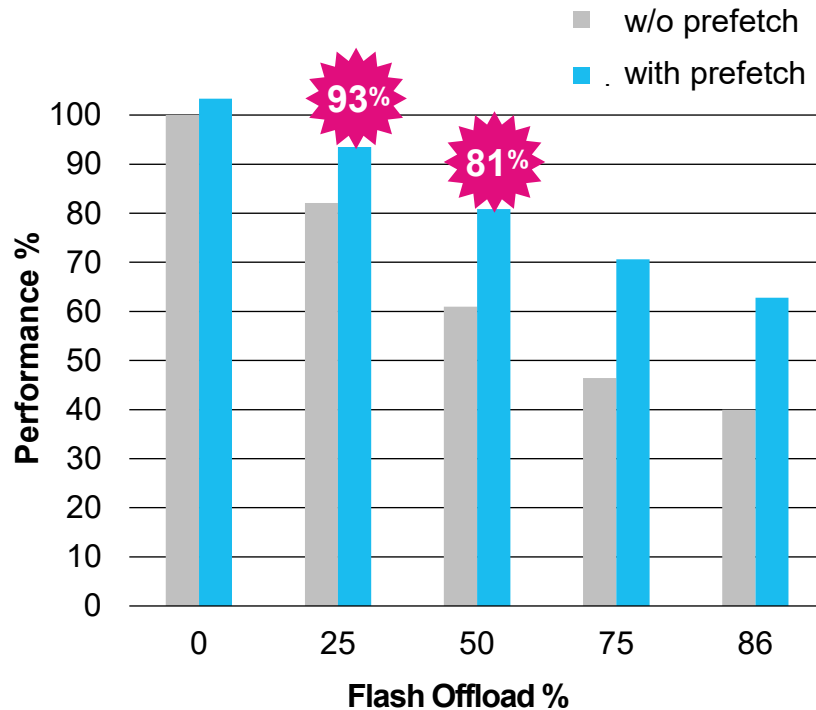
# Redis™ In-Memory Database with Low Latency Flash

Tested with Yahoo!™ Cloud Serving Benchmark (**YCSB**) **tool**
**Setup: 10M records(14 GB)**, 32 client threads

Data Type: **100B*10 fields/record**
**Offload with Linux® TPP (**Transparent Page Placement**)**
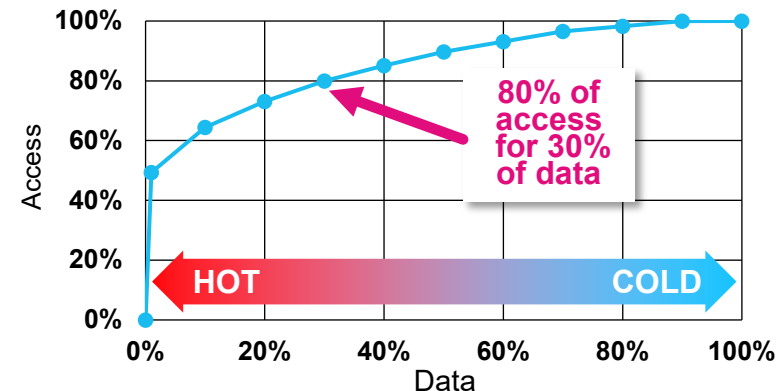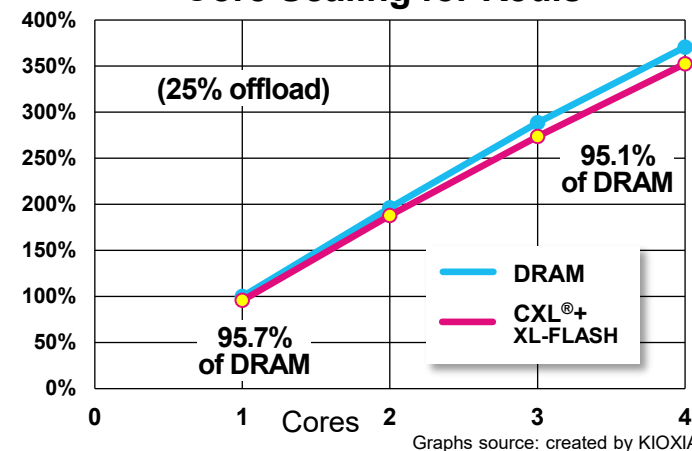
## Test C: Get 100%



Performance % vs Flash Offload %
- w/o prefetch
- with prefetch
- 94%
- 81%

## Test A: Get 50%, Put 50%



Performance % vs Flash Offload %
- w/o prefetch
- with prefetch
- 93%
- 81%

Source: KIOXIA

## Zipf Distribution Workload A,C



Access vs Data

**80% of access for 30% of data**

HOT — COLD

## Core Scaling for Redis



(25% offload)

95.1% of DRAM

95.7% of DRAM

- DRAM
- CXL®+ XL-FLASH

Cores

Graphs source: created by KIOXIA

**YCSB demonstrates CXL® and XL-FLASH technologies can offload 25% of memory with ~5% of performance degradation.**
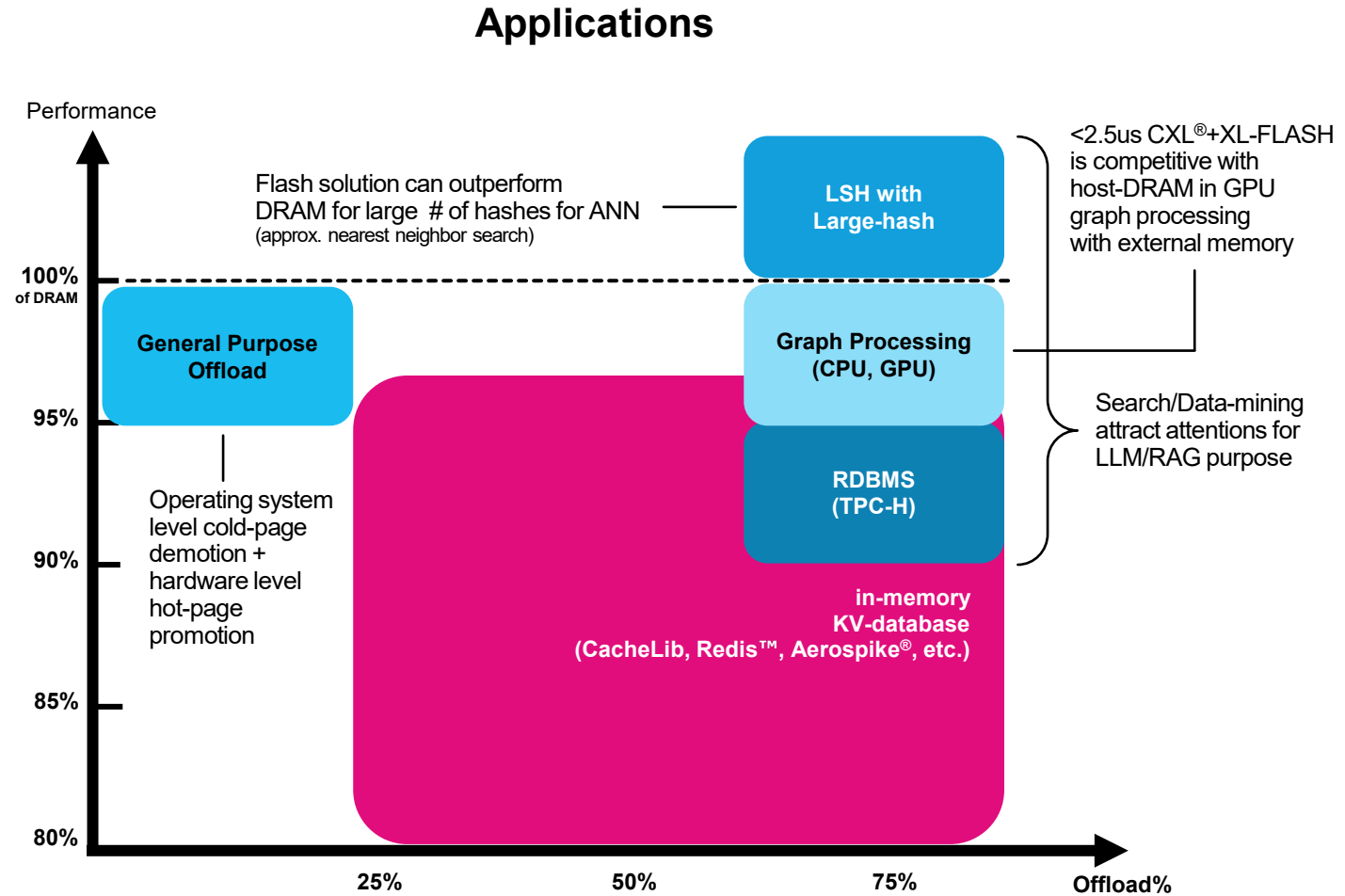
# Challenges and Opportunities with CXL® Attached Flash Memory

## All Applications Are Not The Same

- It is not suitable for latency/bandwidth sensitive applications

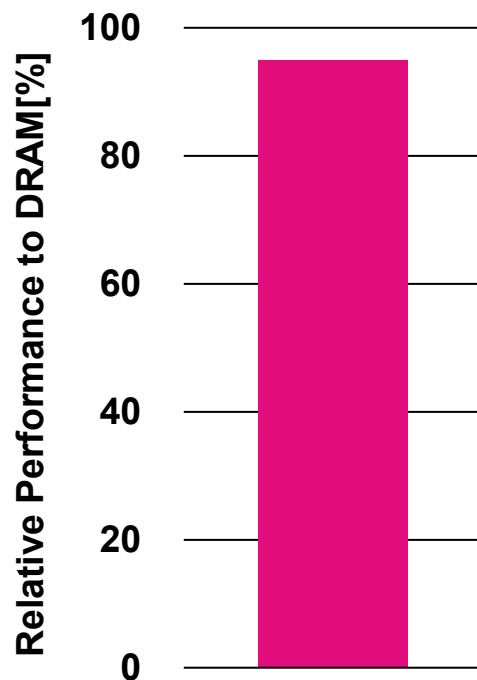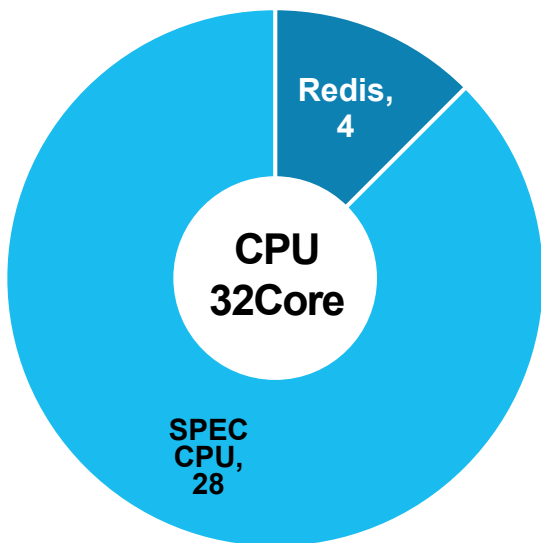- Application not tuned for leveraging memory hierarchy optimally

## Leverage Industry Efforts

- Transparent Page Placement - Automatically manages large memory pages

- Transparent memory tiering solutions optimizes data placement across different memory types

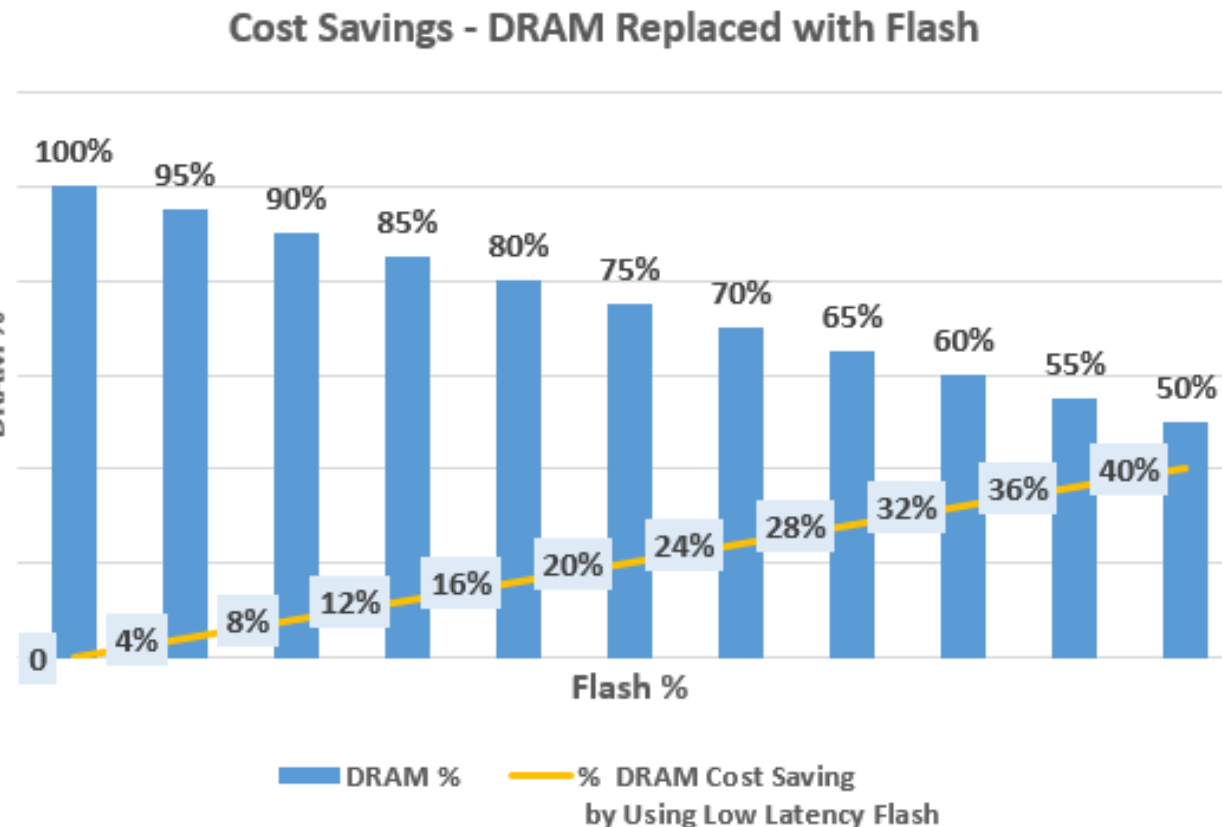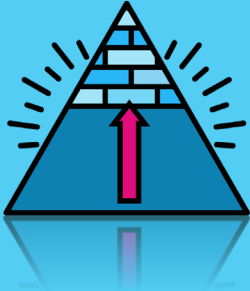- Application specific libraries can further increase the efficiency and reduce cost

**Applications**

Performance

Flash solution can outperform DRAM for large # of hashes for ANN (approx. nearest neighbor search)

**LSH with Large-hash**

<2.5us CXL®+XL-FLASH is competitive with host-DRAM in GPU graph processing with external memory

100% of DRAM

**General Purpose Offload**

**Graph Processing (CPU, GPU)**

Search/Data-mining attract attentions for LLM/RAG purpose

95%

Operating system level cold-page demotion + hardware level hot-page promotion

**RDBMS (TPC-H)**

90%

**in-memory KV-database (CacheLib, Redis™, Aerospike®, etc.)**

85%

80%

25%     50%     75%     Offload%

Source: KIOXIA

# Application and TCO



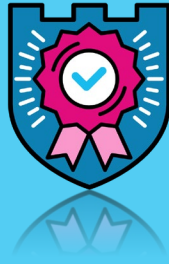| Case B | Offload | Performance |
|--------|---------|-------------|
| Redis™ x 4 cores | 25% | 94% |
| SPEC CPU® x 28 cores | 15% | 95% |
| **total** | **23%** | **95%** |

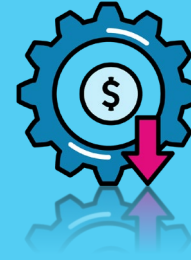KIOXIA assumption: Price of Low Latency Flash is 1/5 DDR Price.

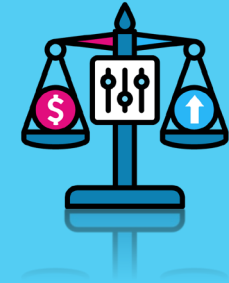Image and graph source: created by KIOXIA

Flash memory can jump the semantic wall.

Flash memory is proven and reliable media.

Flash memory lowers the system TCO.

Flash memory can further perform and reduce cost with software.

If you are working on large memory intensive applications like **Data Mining, Artificial Intelligence (AI), Machine Learning (ML), Analytics, High Performance Computing (HPC), Graph Processing Applications,** Please visit **KIOXIA Booth #307** for collaboration opportunities.

Images and icons source: licensed to KIOXIA from Shutterstock

The product images shown are a representation of the design models and not an accurate product depiction.

KIOXIA