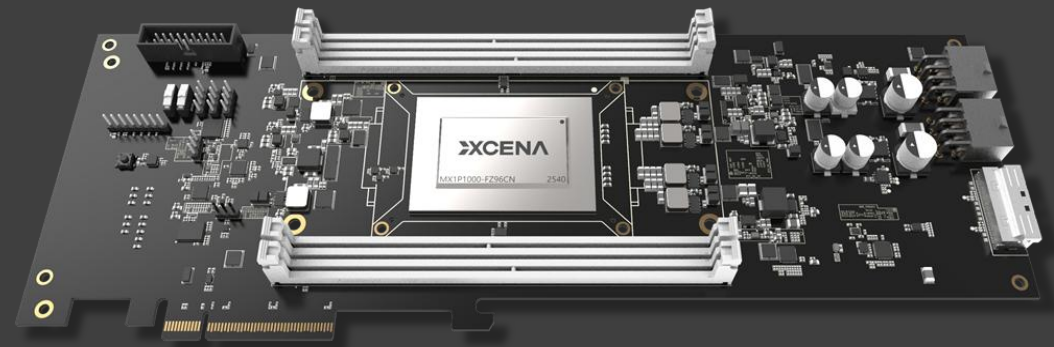


AI+DATA, CXL Memory and Memory Centric Computing

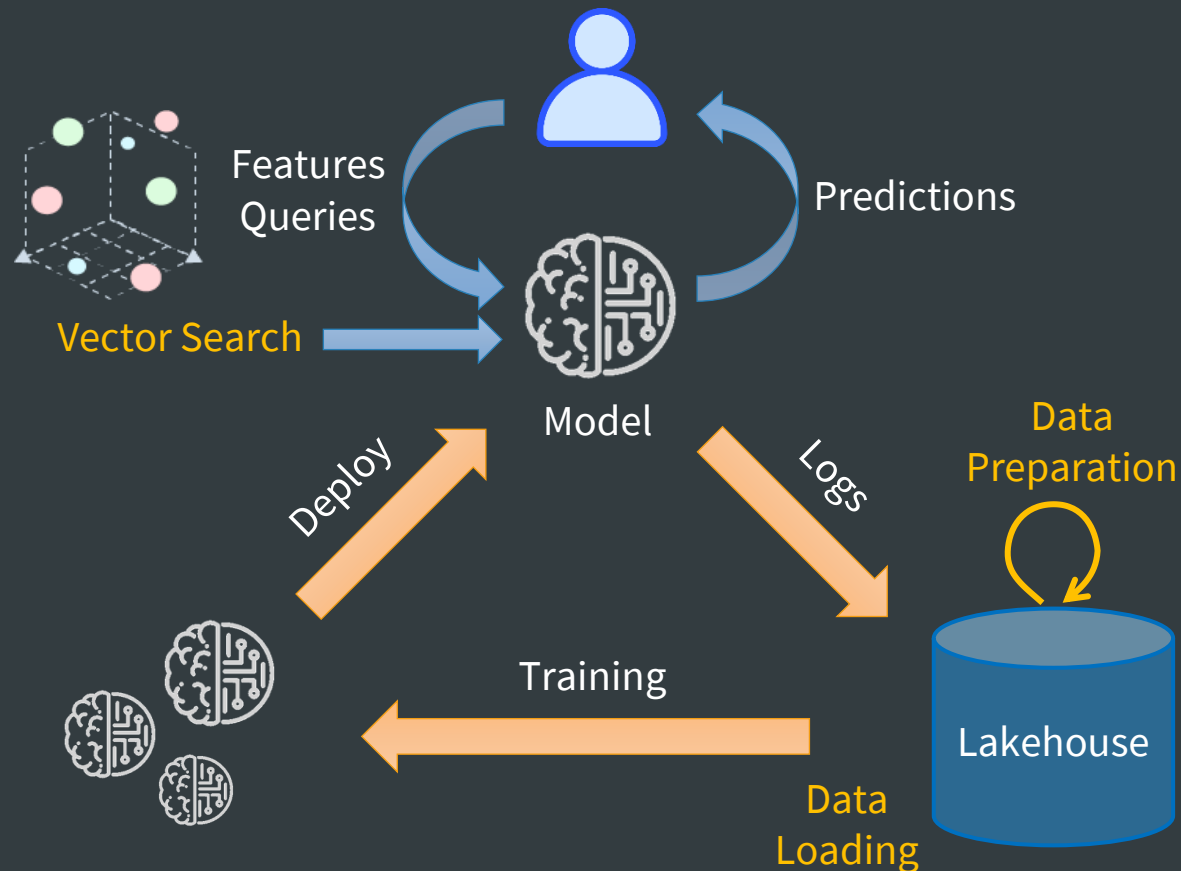
Harry Juhyun Kim

CPO and Co-founder



AI+DATA WORKFLOWS

AI/ML IS THE BIGGEST CONSUMER OF DATA TODAY
THE LARGEST TABLES AND LAKEHOUSES ARE BUILT FOR AI



Data Preparation

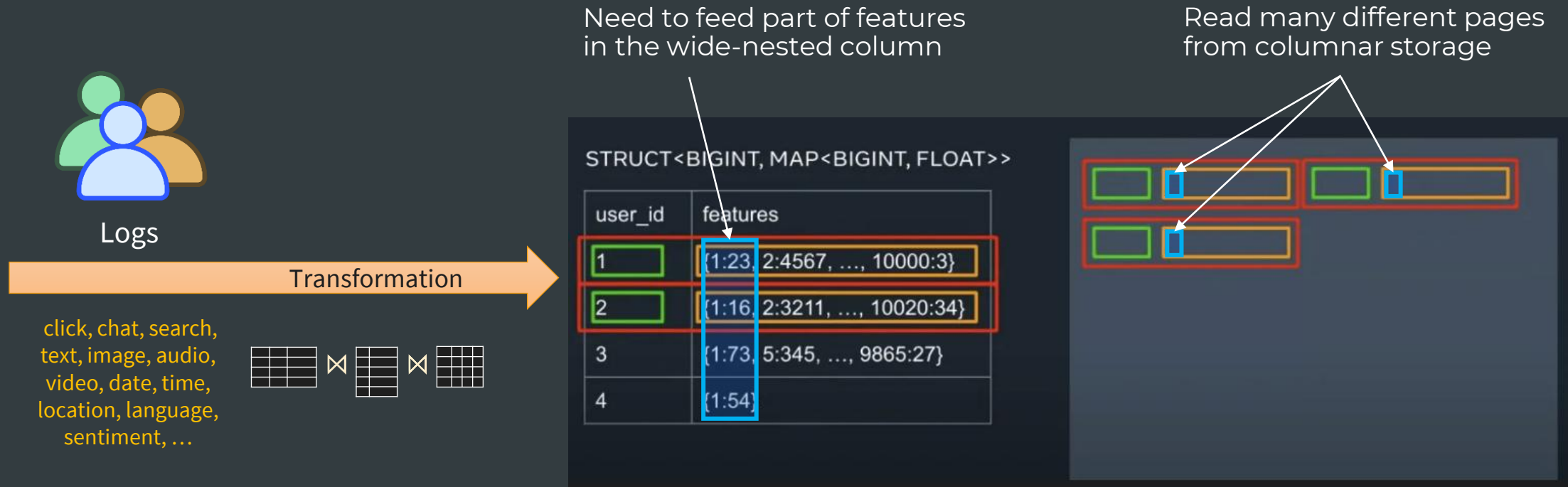
- Preparing tables
- Embeddings, long sequences of features
- Batch ETL pipelines similar to traditional data transformation
- But extremely inefficient
- Need to understand characteristic of “feature engineering”

Data Loading

- Feeding GPUs
- Latency sensitive
- Order sensitive (Training)
→ random/chronological/...
- New challenges for data processing

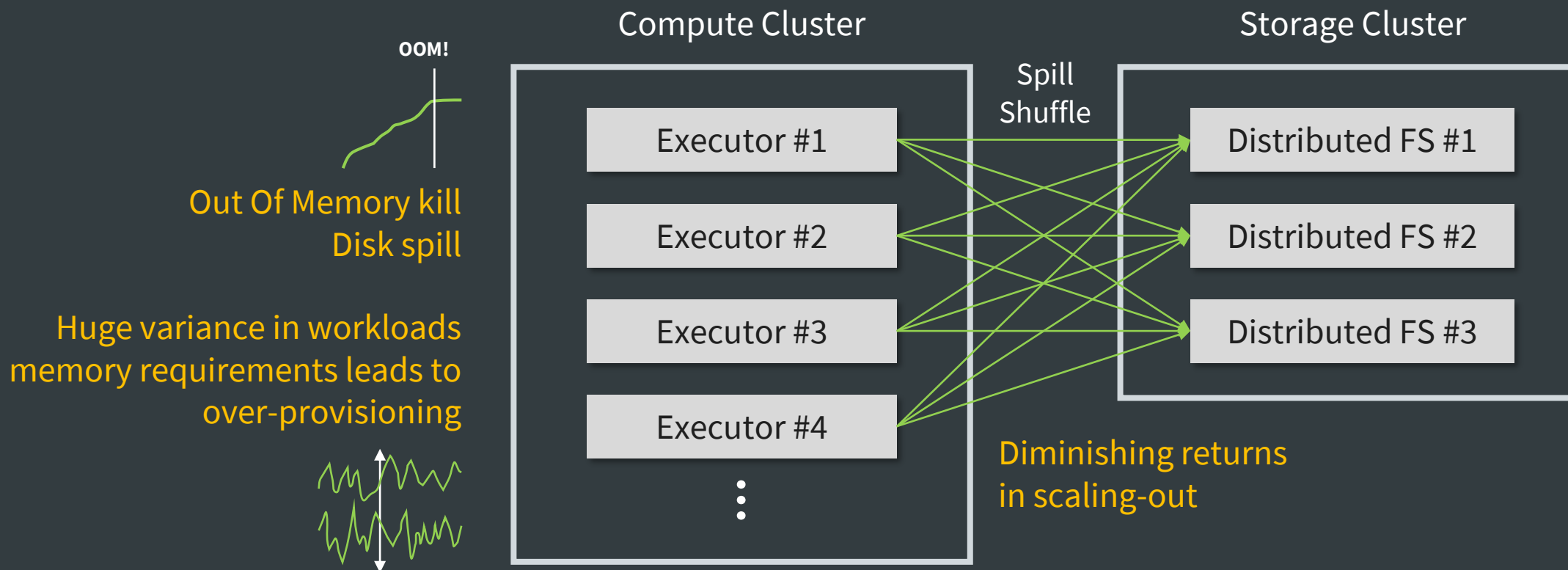
WHAT AI DATA LOOKS LIKE

Features are frequently added and removed in a wide column with nested schema. Transformation, feature engineering, storage and retrieval are all challenging.



CHALLENGES IN DATA ANALYTICS

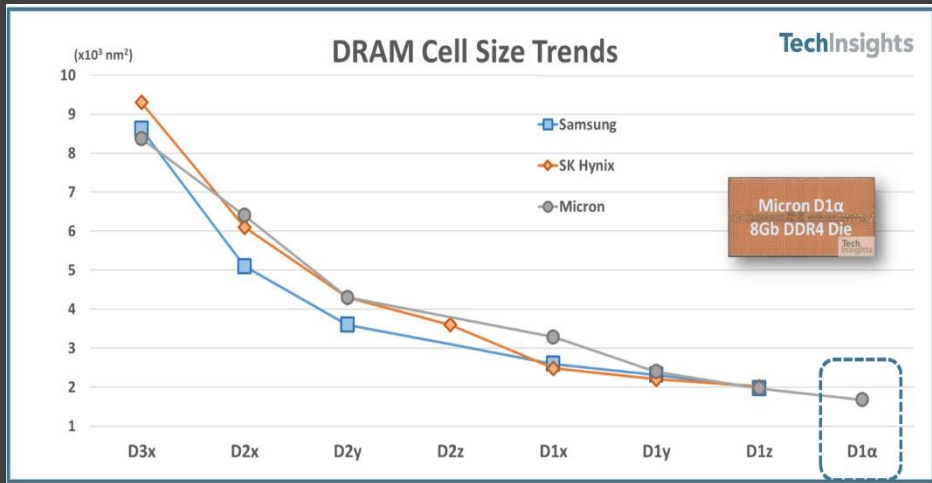
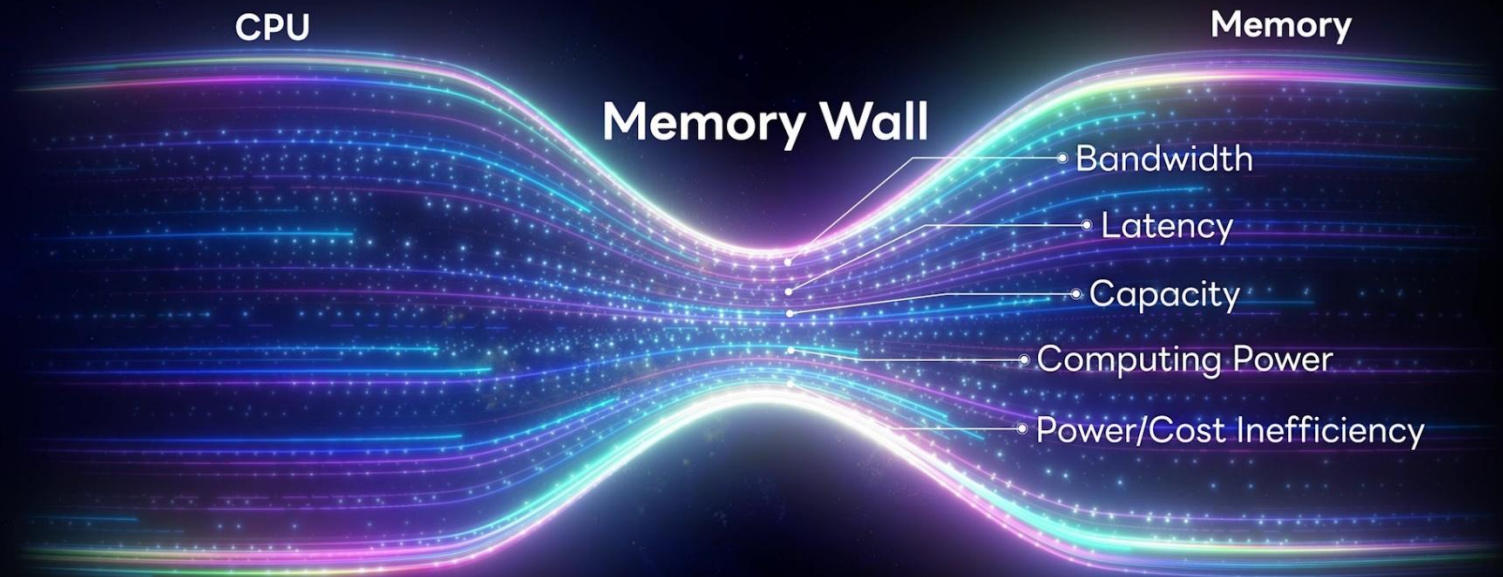
MEMORY INEFFICIENCIES AND DIMINISHING RETURNS IN SCALING
CALL FOR A NEW MEMORY SOLUTION



PROBLEM

DATA IS RAPIDLY
INCREASING,
WHILE
MEMORY TECHNOLOGY
STRUGGLES TO KEEP UP

XCENA AIMS TO BREAK
THE MEMORY WALL PROBLEM



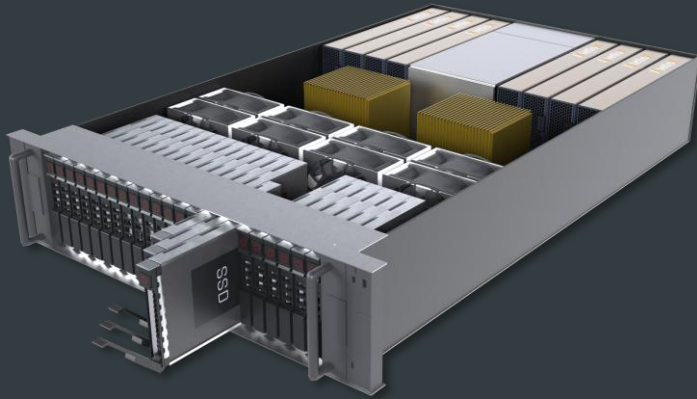
MEMORY WALL: ONE OF THE BIGGEST HEADACHES OF DATA CENTERS

- SIZE: CHALLENGING TO KEEP UP WITH THE RAPID DATA INCREASE
- PERFORMANCE: DATA MOVEMENT IS THE BIGGEST BOTTLENECK
- COST: MEMORY IS REALLY EXPENSIVE (HALF OF THE SERVER COSTS)
- UTILIZATION: MEMORY UTILIZATION IS VERY LOW (OVER-PROVISIONED)

SOLUTION?

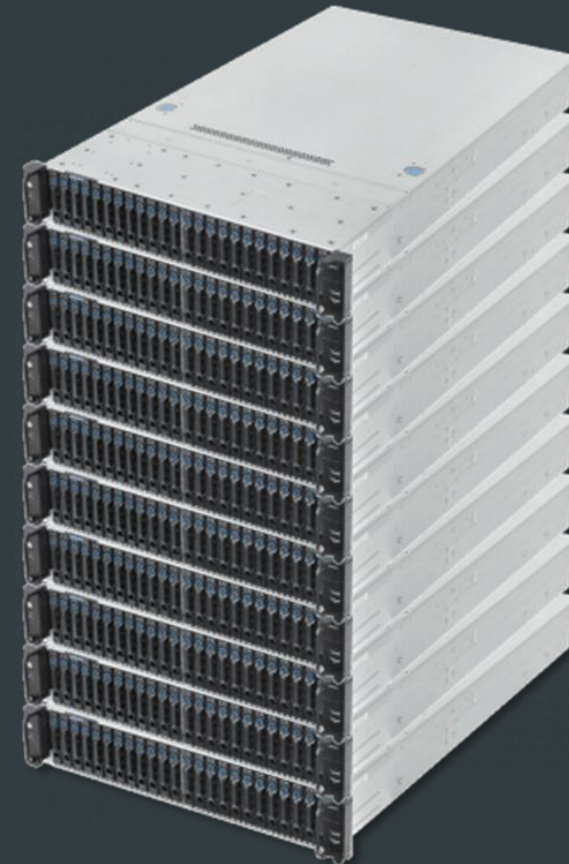
WHAT IF A SINGLE DENSE MEMORY AND COMPUTE NODE
COULD REPLACE A CLUSTER OF 10?

Dense Memory and Compute Node



=

Cluster of Nodes



DATA PROCESSING QUADRANT

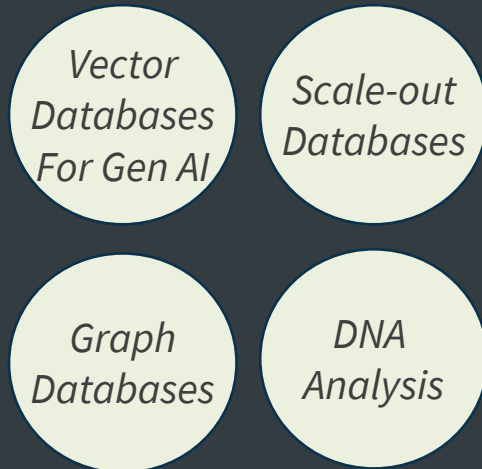
EMBARRESSINGLY PARALLEL, LOW COMPUTE INTENSITY, MEMORY BOUNDED

Disaggregation, Computational Memory

DDR
CXL MEMORY

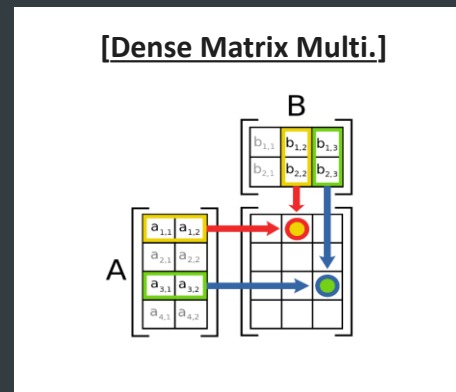
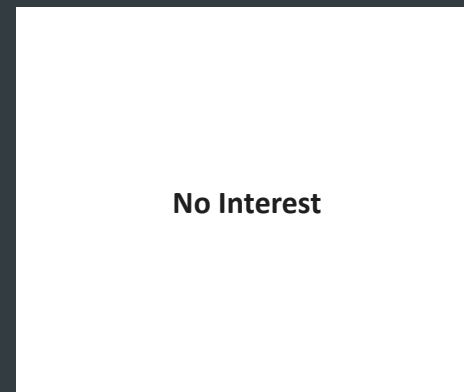
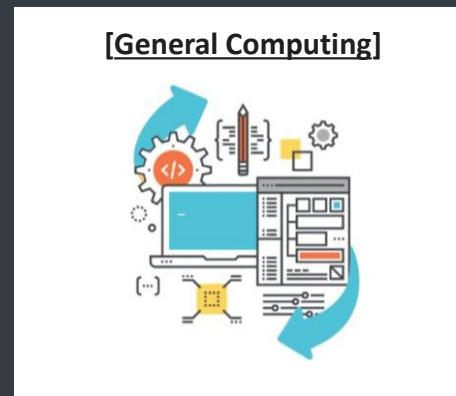
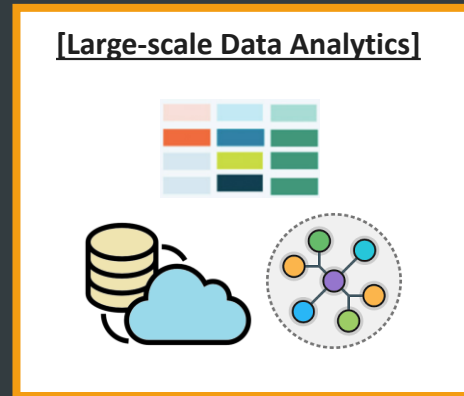


Big Data Processing



Computational Intensity Per Memory Access

Low



High

Bandwidth/Capacity Expansion

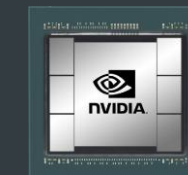


DDR
LPDDR
CXL MEMORY

General Computing

Operation Diversity
Per Memory Access

HBM, Explicit Memory



HBM
GDDR

Artificial Intelligence

COMPUTE EXPRESS LINK

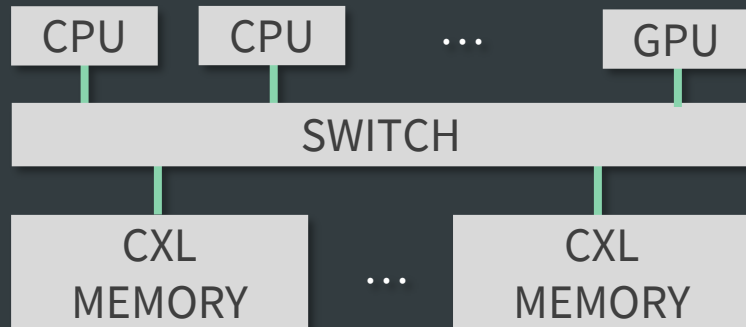


THE NEXT GENERATION MEMORY
INTERCONNECT PROTOCOL

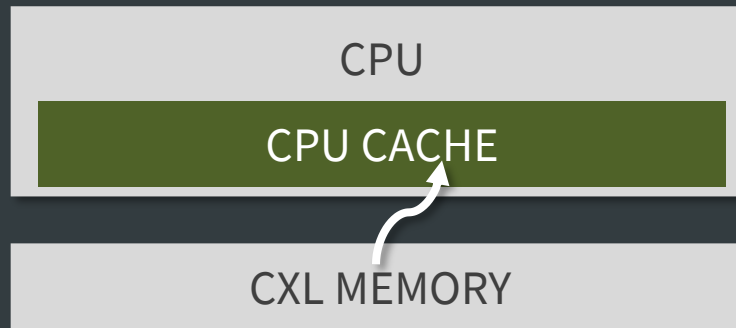
Built on PCIe that lets you plug in memory as the same way you plug in a GPU card

CXL memory looks like regular DRAM with additional latency due to protocol overhead

Memory pooling and sharing through Switch



HIGHLY SCALABLE
MEMORY EXPANSION



CACHE COHERENCE

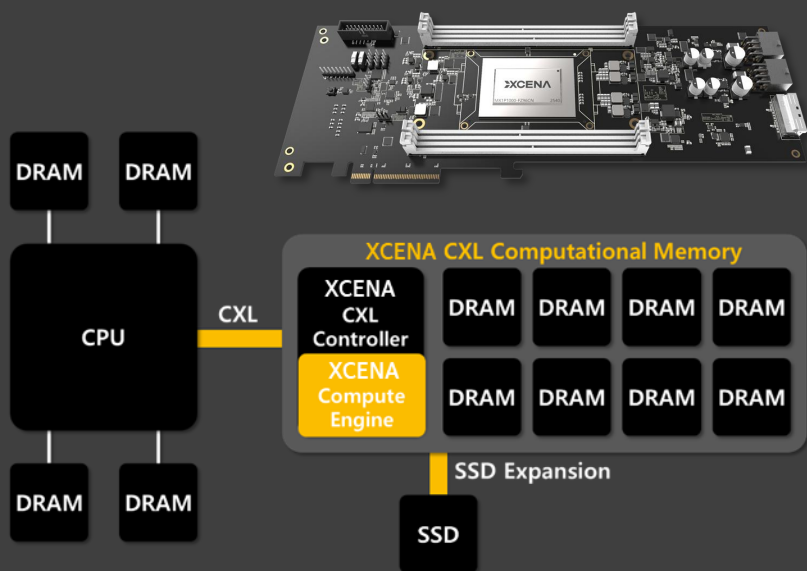


HETEROGENEOUS COMPUTING

XCENA

XCENA

CXL COMPUTATIONAL MEMORY



BEYOND JUST ANOTHER CXL MEMORY EXPANDER

DATA PROCESSING LIBRARY
DATA ANALYTICS ACCELERATION

TCO SAVINGS
SEAMLESS INTEGRATION

COMPUTING HW-SW SOLUTION
NEAR-MEMORY PROCESSING

LESS DATA MOVEMENT
LOWER CPU UTIL
LOWER POWER
BETTER PERFORMANCE

CXL MEMORY

**HIGH PERFORMANCE
CXL DRAM EXPANSION**

GENERAL CXL
PLAYERS

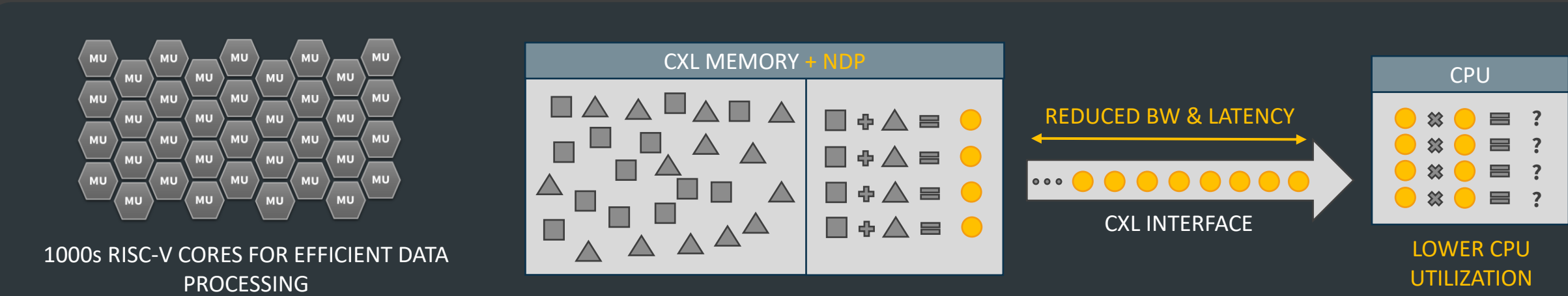
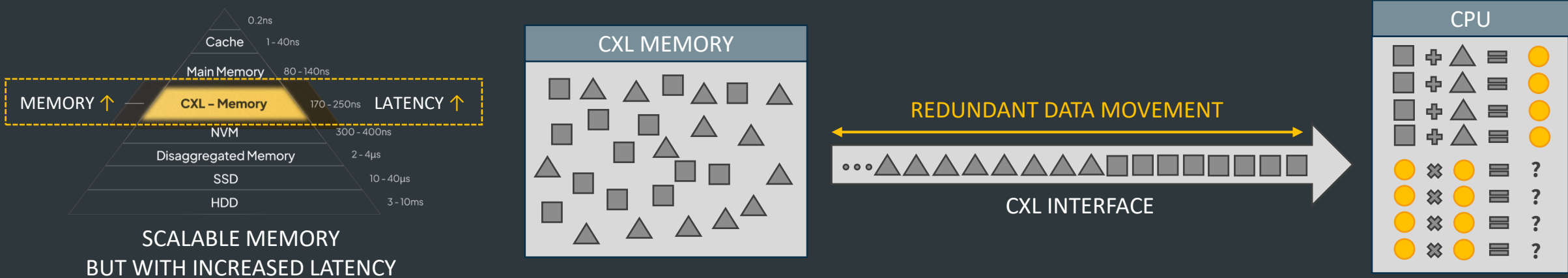
INFINITE MEMORY

**PB-SCALE MEMORY
WITH SSD EXPANSION**

MORE MEMORY
CHEAPER MEMORY

NEAR-MEMORY PROCESSING

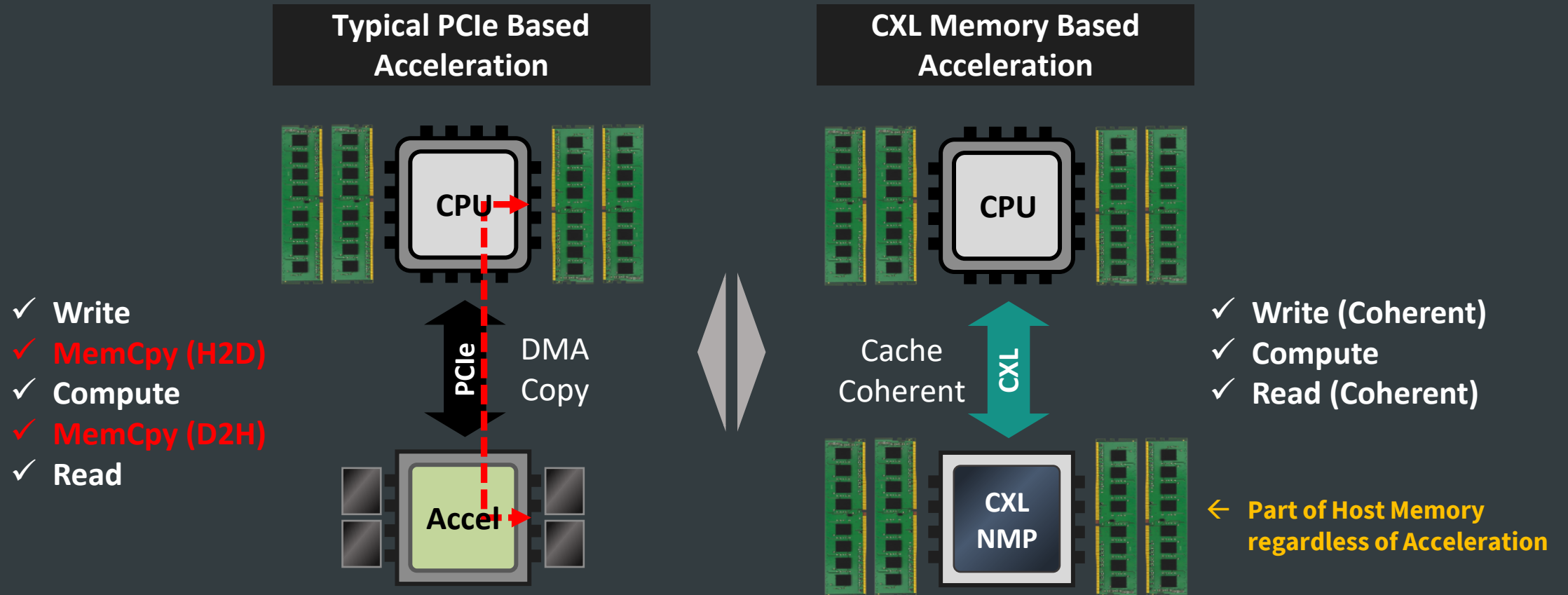
More Capacity, Extra Bandwidth, +Latency
Reducing data movement by Near-Memory Processing



KEY DIFFERENCES

PCIe vs. CXL

CXL MEMORY IS A PART OF HOST SYSTEM MEMORY,
WHEREAS PCIe DEVICE MEMORY IS NOT VISIBLE TO THE HOST.



CXL MEMORY + DATA PROCESSING

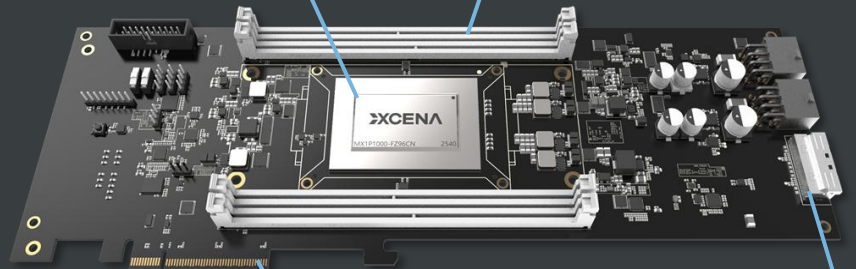


CXL COMPUTATIONAL MEMORY FOR LARGE-SCALE DATA AVAILABLE SOON

Novel CXL Hardware

1000s of Custom
RISC-V Cores
+
TFLOPS
Vector Engine

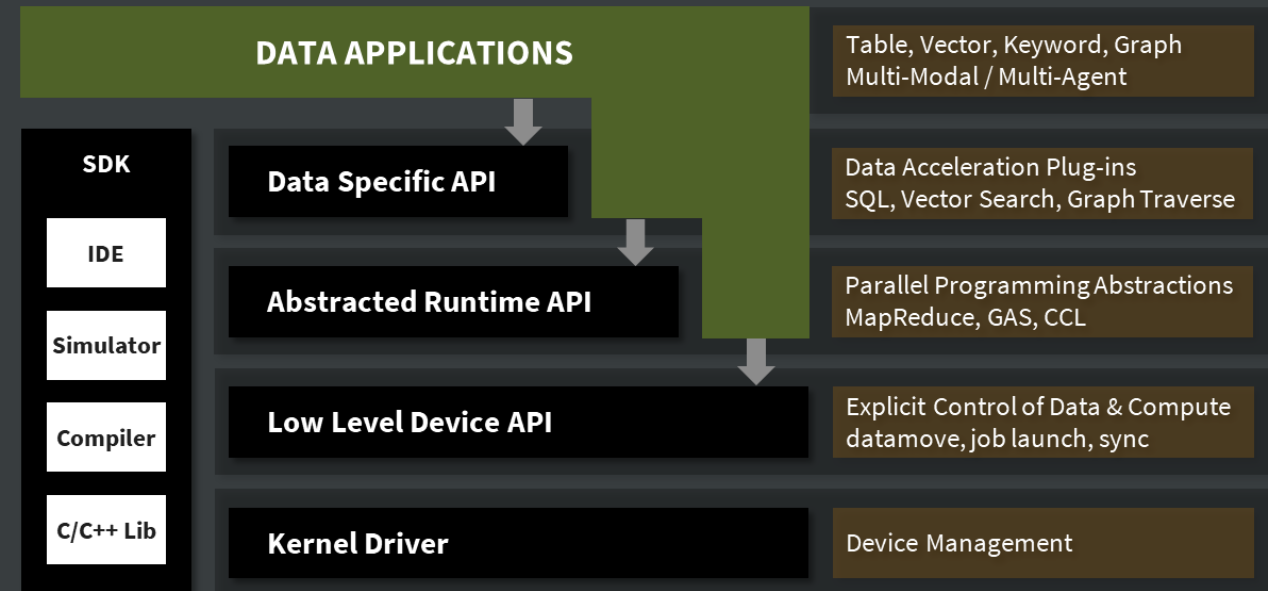
DDR5 x 4Ch
~1TB



CXL 3.0 HDM-DB
with Back-invalidation
Cache Coherence

SSD-backed
CXL Expansion

Rich Software Framework



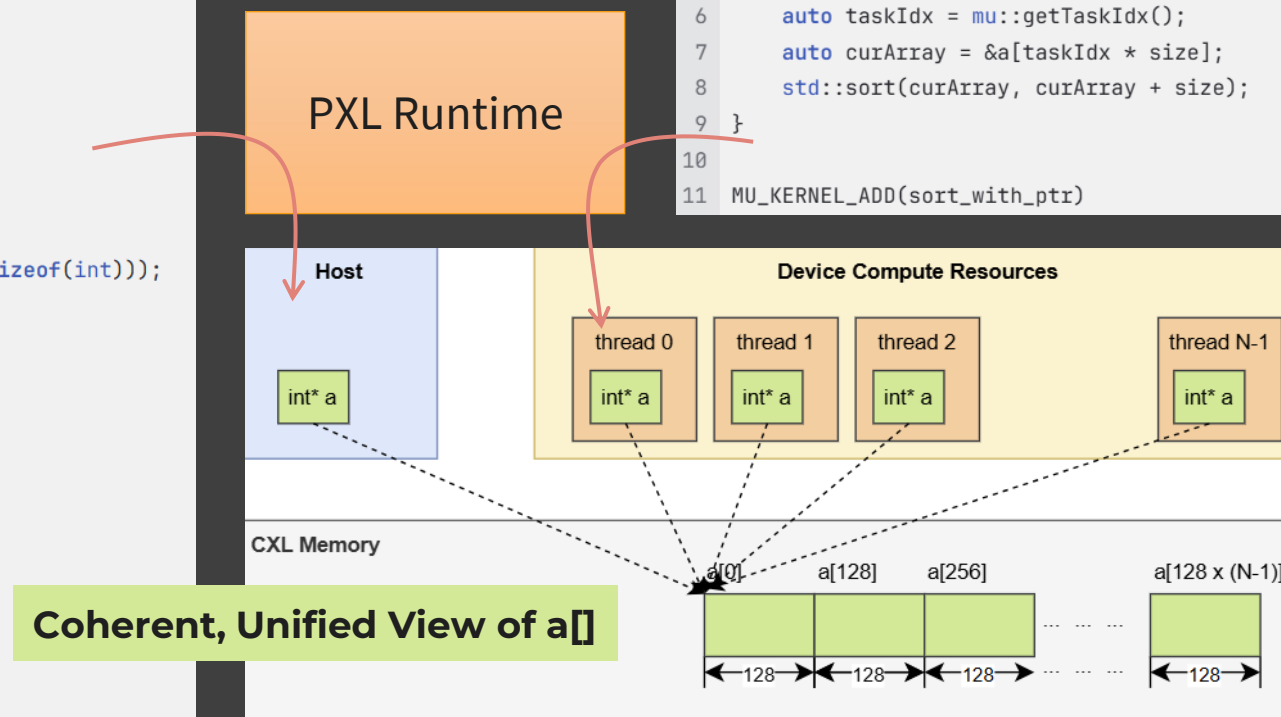
PARALLEL PROGRAMMING

Parallel Xceleration Library (PXL)

Map API encapsulates all the complexities of device management and provides a FRAMEWORK of thought for developers.

```
1 #include "pxl/pxl.hpp"
2
3 // setup the device
4 const char* filename = "mu_kernel/mu_kernel.mubin";
5 const char* muFuncName = "sort_with_ptr";
6 auto context = pxl::runtime::createContext(0);
7 auto module = pxl::createModule(filename);
8 auto job = context->createJob();
9 job->load(module);
10 auto muFunction = module->createFunction(muFuncName);
11
12 // allocate cxl memory
13 int* a = reinterpret_cast<int*>(context->memAlloc(N * 128 * sizeof(int)));
14 // setup data
15 for (size_t i = 0; i < N; i++)
16 {
17     for (size_t j = 0; j < 128; j++)
18     {
19         a[i * 128 + j] = 128 - j;
20     }
21 }
22
23 auto map = job->buildMap(muFunction, N);
24 auto ret = map->execute(a, 128);
25 map->synchronize();
```

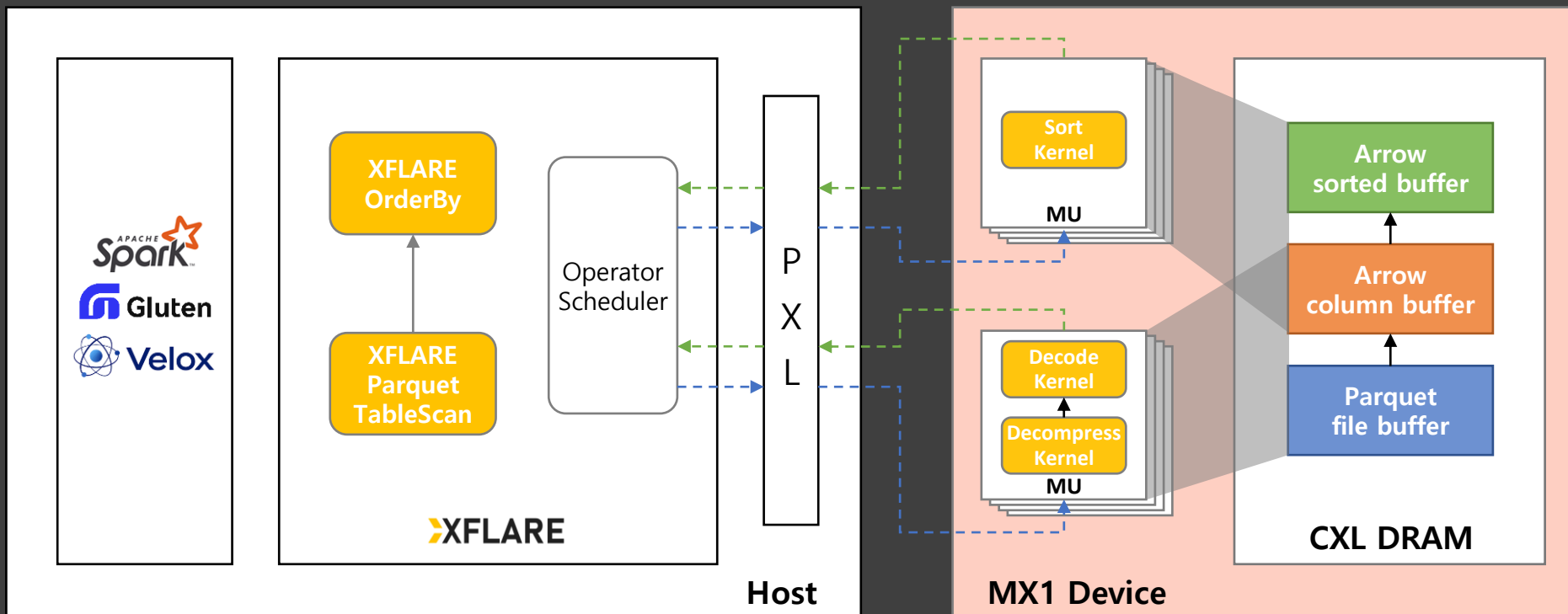
```
1 #include <algorithm>
2 #include "mu/mu.hpp"
3
4 void sort_with_ptr(int* a, int size)
5 {
6     auto taskIdx = mu::getTaskIdx();
7     auto curArray = &a[taskIdx * size];
8     std::sort(curArray, curArray + size);
9 }
10
11 MU_KERNEL_ADD(sort_with_ptr)
```



DATA SPECIFIC LIBRARY

Spark SQL query is compiled and planned by Spark, translated into Velox plans by Gluten, and executed by XFLARE on MX1.

```
1 spark.read.parquet("persons.parquet").createOrReplaceTempView("persons")
2 spark.sql("""
3     SELECT  *
4     FROM    persons
5     ORDER BY age DESC
6 """).show()
```



DON'T MOVE DATA, MOVE COMPUTATIONS!

CXL LOAD/STORE

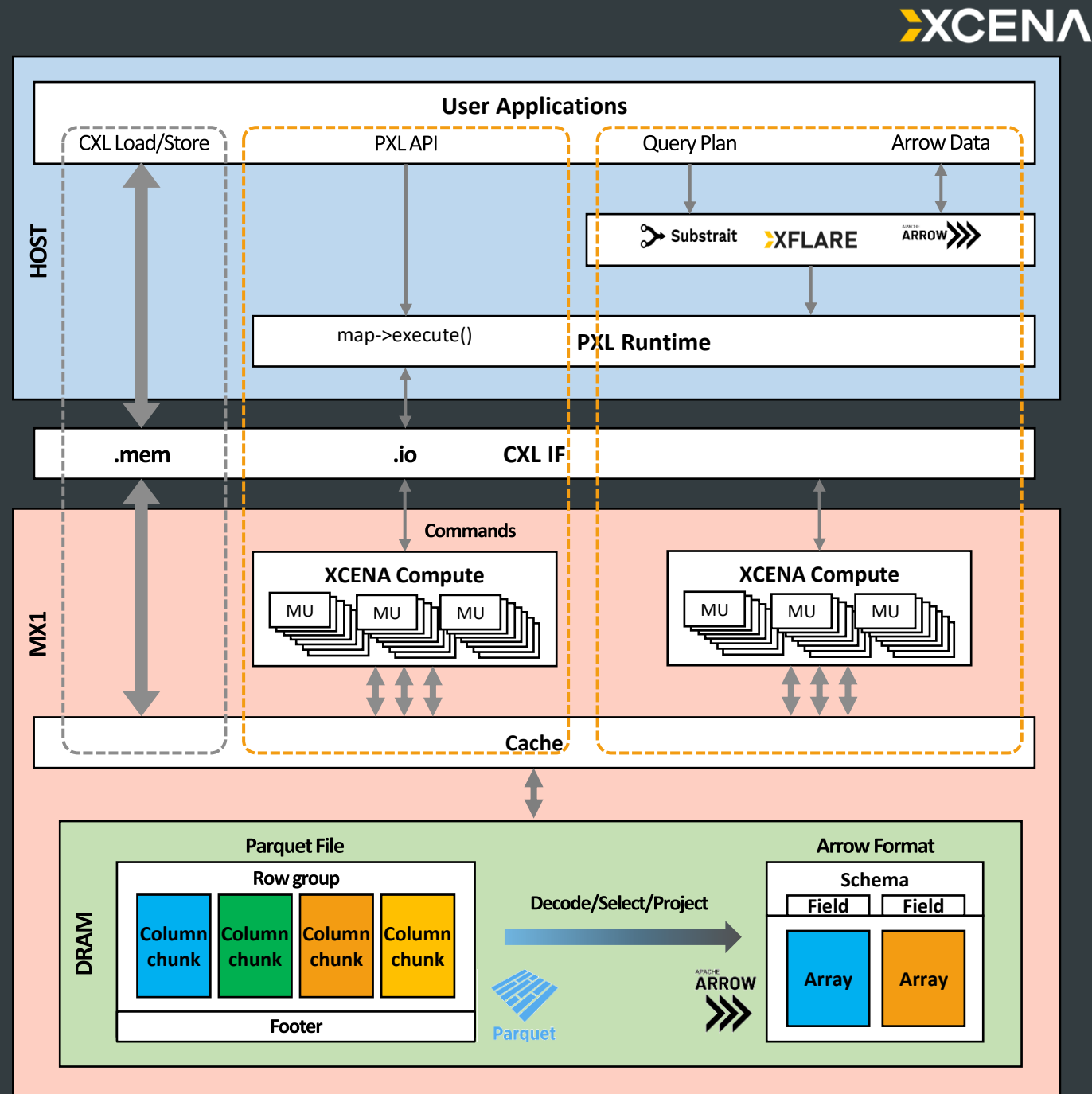
Use CXL memory as just another part of host memory.

PXL API

PXL allows you to manage resources and execute user kernels in parallel.

XFLARE

Common formats (Substrait, Arrow) and high-level APIs simplify query offloading.



[Introduction](#)[Install](#)[Workflow](#)[CXL Key Concepts](#)[Get Started](#)[Tutorials](#)[Documentation](#)[Release Notes](#)[Overview](#)

XCENA Software Development Kit (SDK) is a comprehensive software framework designed to facilitate easy and efficient utilization of XCENA hardware for data-intensive applications such as AI, machine learning, and data-intensive computing. By providing multi-level APIs and a suite of development tools, the XCENA SDK allows for seamless integration with CXL computational memory, enabling developers to accelerate their applications with minimal effort.

[Features](#)

- Seamless integration with computational CXL memory, specialized for data-intensive applications such as AI, machine learning, and data-intensive computing.
- Multi-level APIs for flexible, efficient development
- Support for emulation and simulation with general tools like QEMU and our proprietary simulators, ensuring smooth development and testing environments

[Software Ecosystem](#)

- **Applications:** Main applications running on the host system that leverage XCENA hardware for

XCENA.COM/SDK

QEMU CXL 3.0 Host + Computational Memory

XCENA SDK & Runtime Library

Example Codes and Tutorials

→ ALL IN A DOCKER IMAGE

THANK YOU

HARRY KIM, CPO

harry.kim@xcena.com

<http://xcena.com>

<https://www.linkedin.com/company/xcena/>