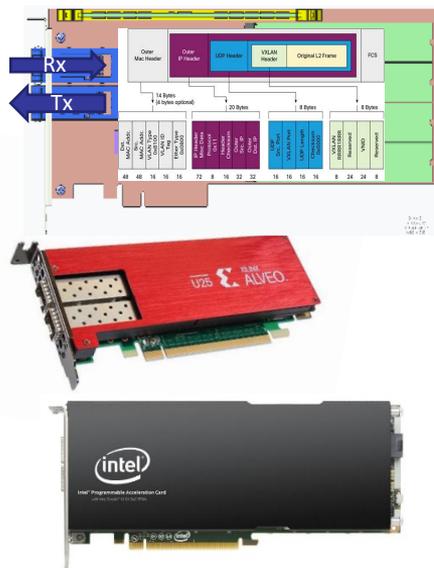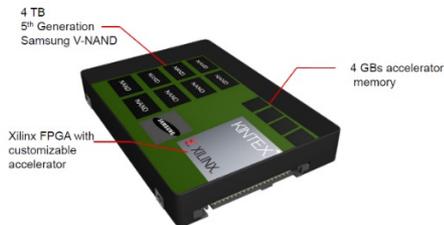# Types of Hardware Accelerators
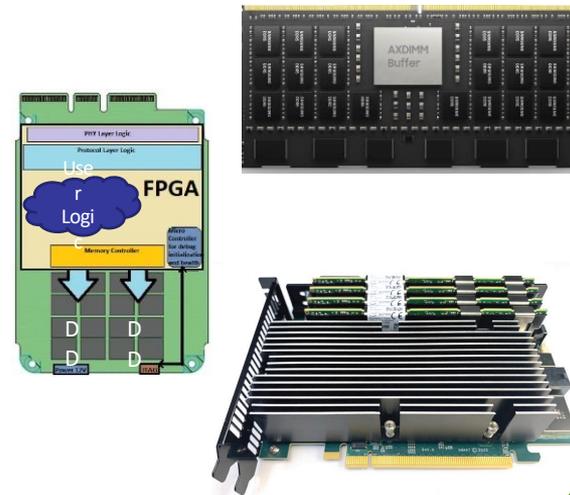


## Network Accelerators (Smart NIC)

## Storage Accelerators Computational Storage (CS)

SmartSSD® CSD

## Memory Accelerators Computational Memory (CM) Processing in Memory (PIM)
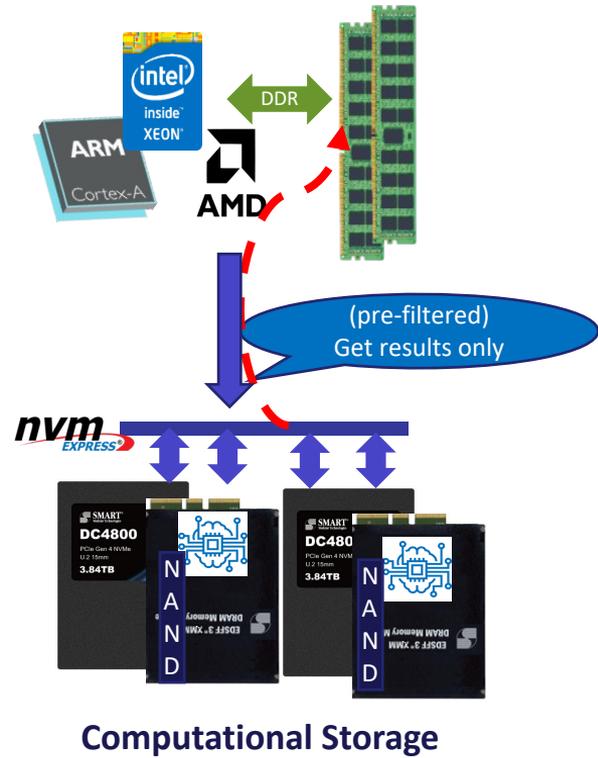
Near Data Processing (NDP)

[1a] https://blocksandfiles.com/2020/10/28/amd-xilinx-smartnic-data-centre/
[1b] https://www.servethehome.com/intel-fpga-pac-d5005-high-end-drop-in-accelerator-launched
[2a] https://www.servethehome.com/xilinx-samsung-smartssd-computational-storage-drive-launched/
[2b] https://www.servethehome.com/intel-fpga-pac-d5005-high-end-drop-in-accelerator-launched/
[3a] https://tekdeeps.com/samsung-also-sees-the-future-in-memories-that-also-perform-calculations/

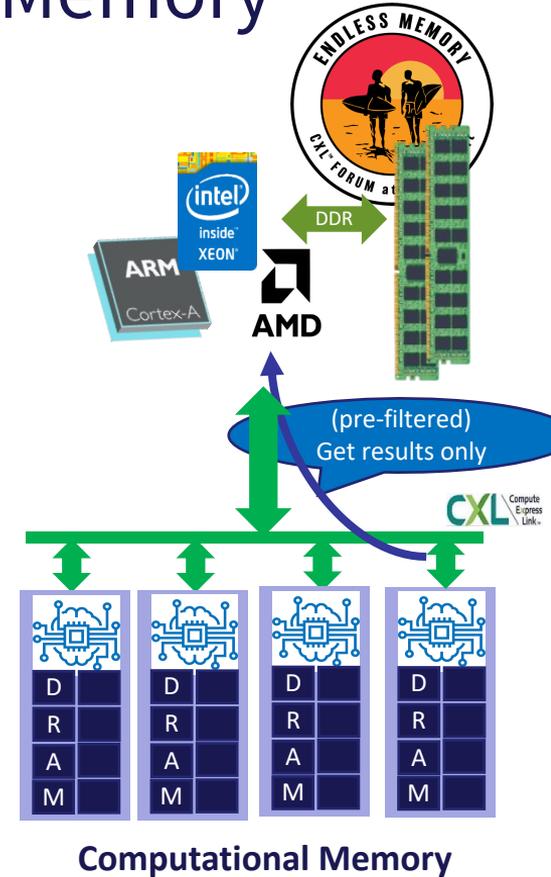# Computation Storage vs Computational Memory



**Protocol specific differences**

| IO Mapped | Direct addressable |
| (high latency) | (low latency) |

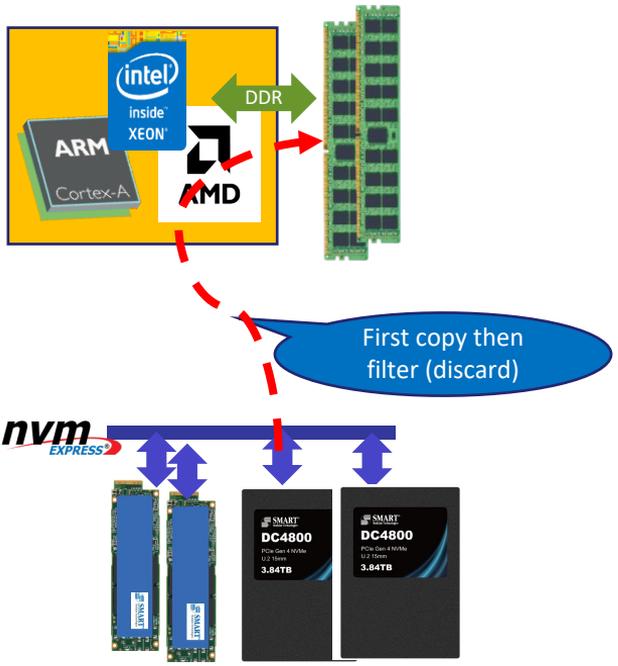| Block access | Memory semantics |
| (read/write) | (load/store) |

**(Media specific difference)**

Limited Endurance | Infinite Endurance
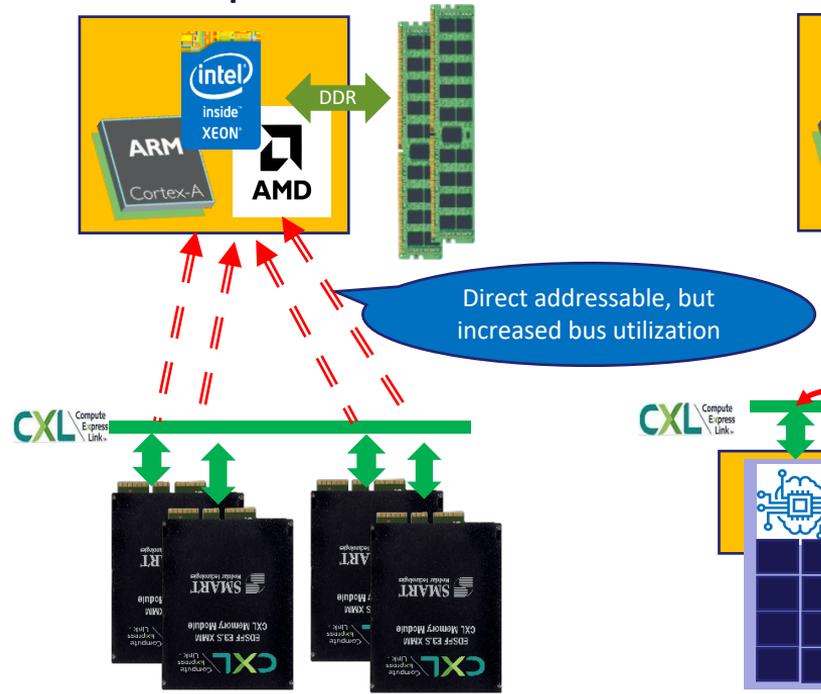Non-Volatile | Volatile/ Persistent

(pre-filtered)
Get results only

**Computational Storage**

**Computational Memory**

# Benefits of Moving Compute near the Data



**Conventional Systems**

First copy then filter (discard)

- Performance limited by I/O latency
- Inefficient bus utilization

**In development**

Direct addressable, but increased bus utilization

- Lower latency but more accesses
- Inefficient bus utilization

**Planning for Future**

Returns pre processed Data

- Distributed Compute
- Efficient bus utilization
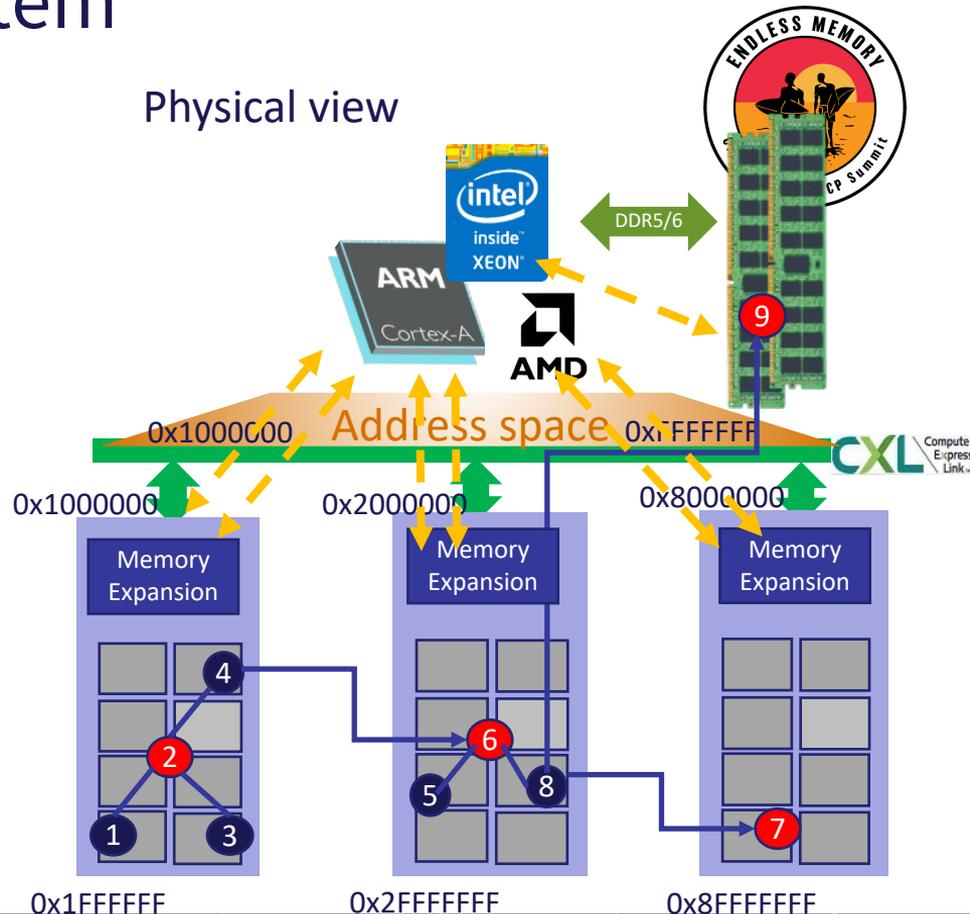
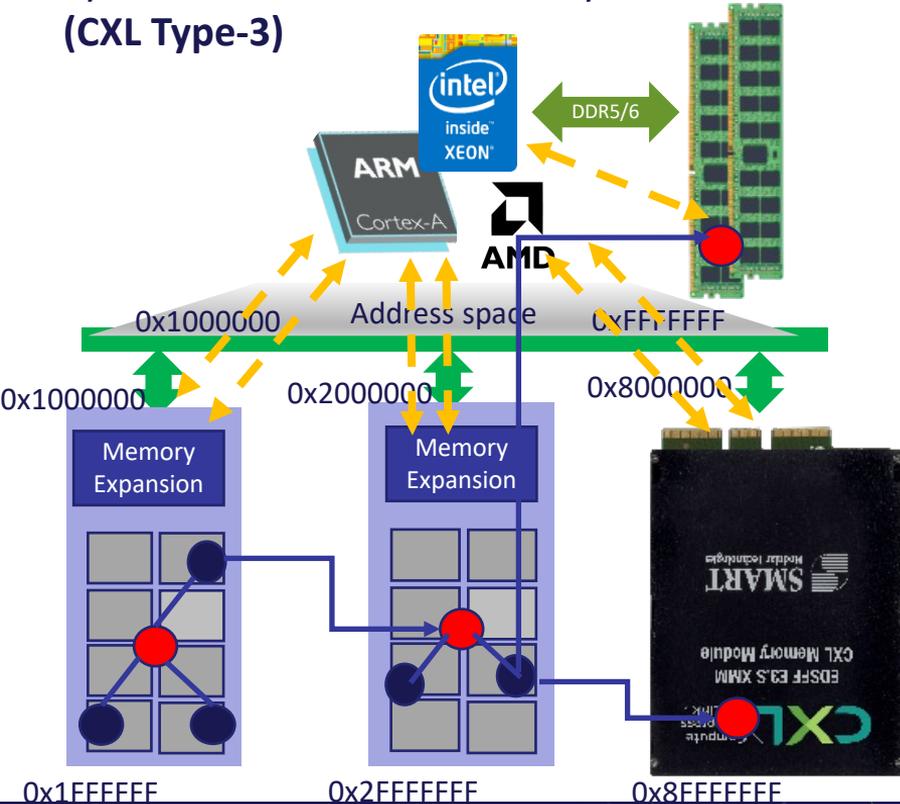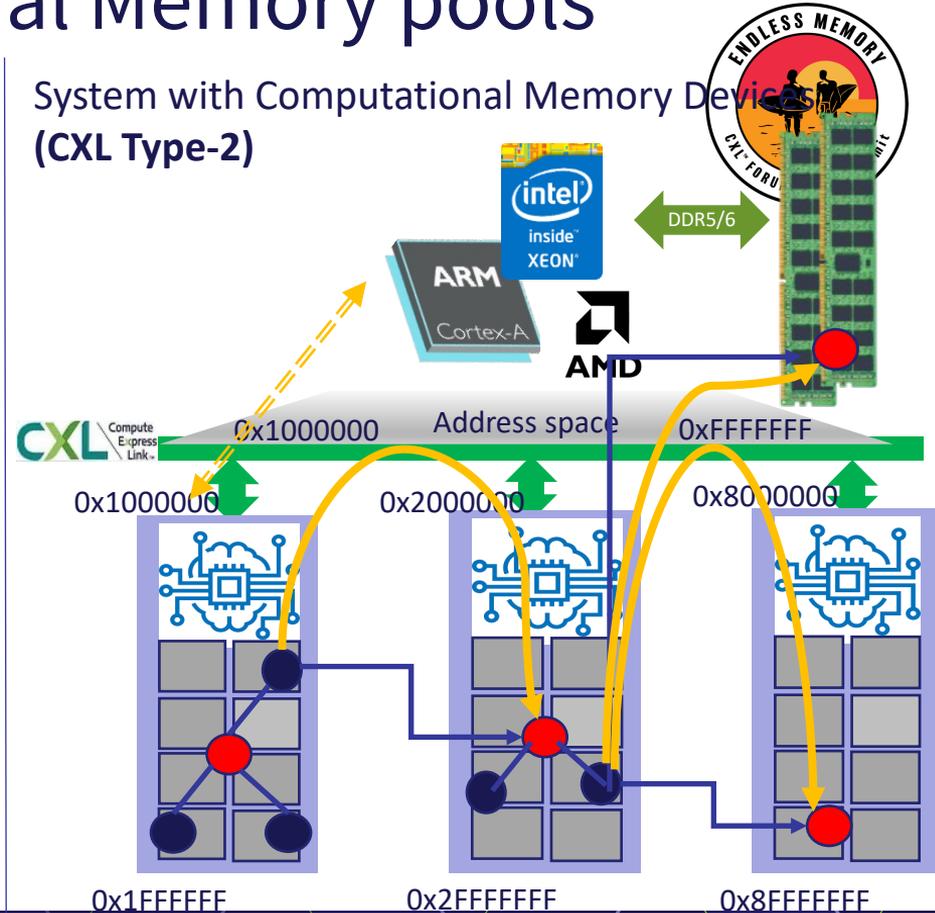# Data in Pooled Memory System

## Logical view



## Physical view

# Standard vs Computational Memory pools



System with Standard Memory Devices
**(CXL Type-3)**

System with Computational Memory Devices
**(CXL Type-2)**

# Combining Best of Both the Worlds

System with Computational Memory Devices (**CXL Type-2**) and Standard Memory Devices (**CXL Type-3**)

✓ Free up Host CPU.
   Computational Memory manages all prefetching and Data consolidation.

✓ Provides deterministic latency.
   Data is pre-filtered and only required response is returned to Host.

✓ Reduces cost by sharing of Memory resources.
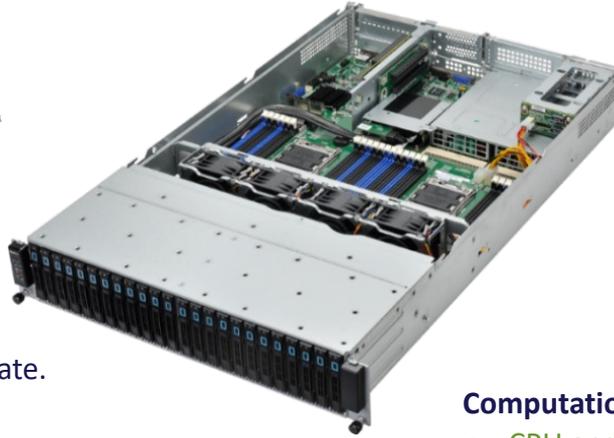
# Selecting the Right Form Factor



**Computational Memory in DIMM form-factor**
- **Parallel attached, Lowest latency**
- CPU and Platform dependent
- Acceleration offload limited by Thermal budged
- Choice of controller silicon limited by PCB real-estate.
- *May degrade speed of entire DDR channel*

**Computational Memory in PCIe attached form-factor**
- CPU and Platform Agnostic
- Flexibility of vendor and capacity selection.
- Cannot Hot Plug. Limited serviceability.

**Computational Memory in EDSFF (like E1.S or E3.S)**
- **Serial attached**. CPU agnostic. Media agnostic
- **Hot pluggable. Improves serviceability**
- Scalable. Pluggable in same slots as SSD.
- Serialization of data-bus adds latency
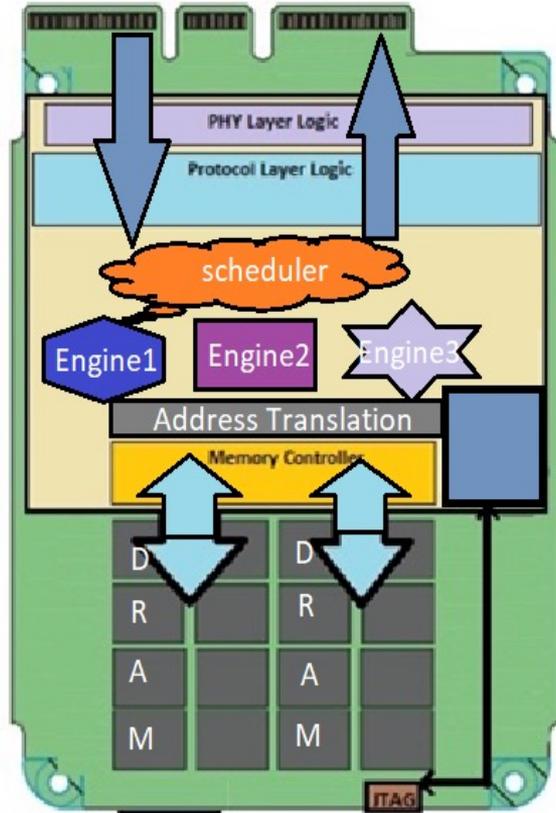- Capacity and acceleration limited by PCB and Thermal profile.

# Anatomy of a Computational Memory Device



**Protocol Command Parser**
- Standard CXL compliant logic for Type-3/2 device
- Extracts information from standard CXL commands
- Device initialization, enumeration and housekeeping
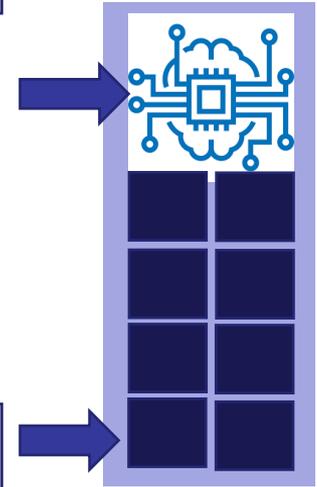
**Programmable Compute Engine(s)**
- Breaks complex commands in atomic executable tasks
- Execute individual tasks and collate results.

**CXL.Mem decoder**
- Converts Host Addresses into Device local Addresses
- Memory Media manager

**Memory Media**
DDR4/5/6, Persistent Memory and even NAND Flash

# Key Take-Aways

1. Computational Memory is different from Computational Storage.

2. CXL is enabling new way to distribute and parallelize Compute.

3. Choosing the right Mechanical Form-factor is important for scalability.

Call for Action

- Standardize API and framework for "Computational Memory"

- Join OCP's Software Defined Memory (SDM) initiative

    - OCP Software_Defined_Memory

# Thank you!

## EMPOWERING OPEN.

OCP GLOBAL SUMMIT

OCTOBER 18-20, 2022
SAN JOSE, CA