

Software Defined Memory (SDM) with Faster Interconnects and Tiered Memory

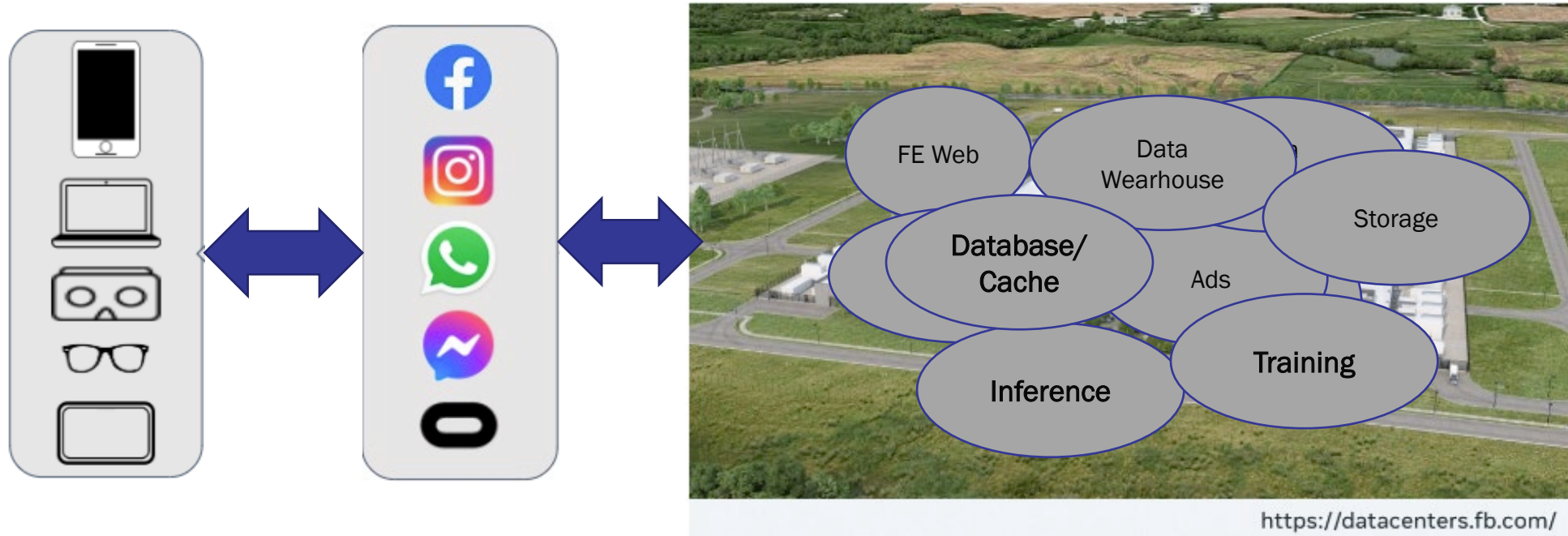
Manoj Wadekar, Hardware Systems Technologist, Meta

Anjaneya “Reddy” Chagam, Sr. Principal Engineer, Intel

Agenda

- Memory Challenges in Hyperscale Infrastructure
- Need for Software Defined Memory (SDM)
- SDM with Compute Express Link (CXL)
- SDM use-cases
- Summary

Hyperscale Infrastructure



- Application performance and growth depends on
 - DC, System, Component performance and growth
 - Compute, Memory, Storage, Network..
- Focusing on Memory discussion

Memory Challenges



Bandwidth and Capacity

- The Gap between bandwidth and capacity is widening
- Applications ready to trade between bandwidth and capacity



Power

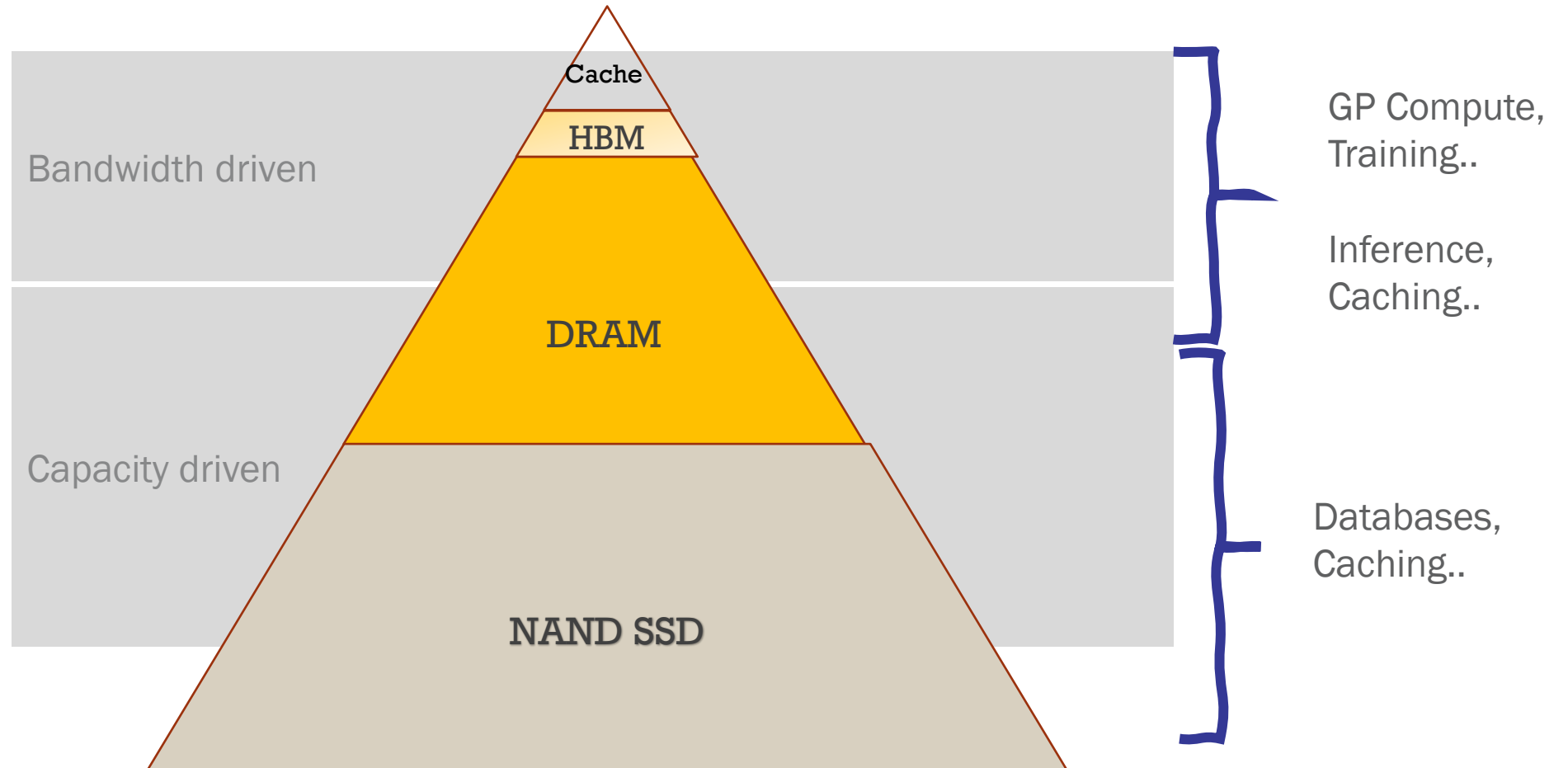
- DIMMs consume significant share of rack power
 - DDR5 exacerbates this
- Applications co-design to achieve higher capacity at optimized power



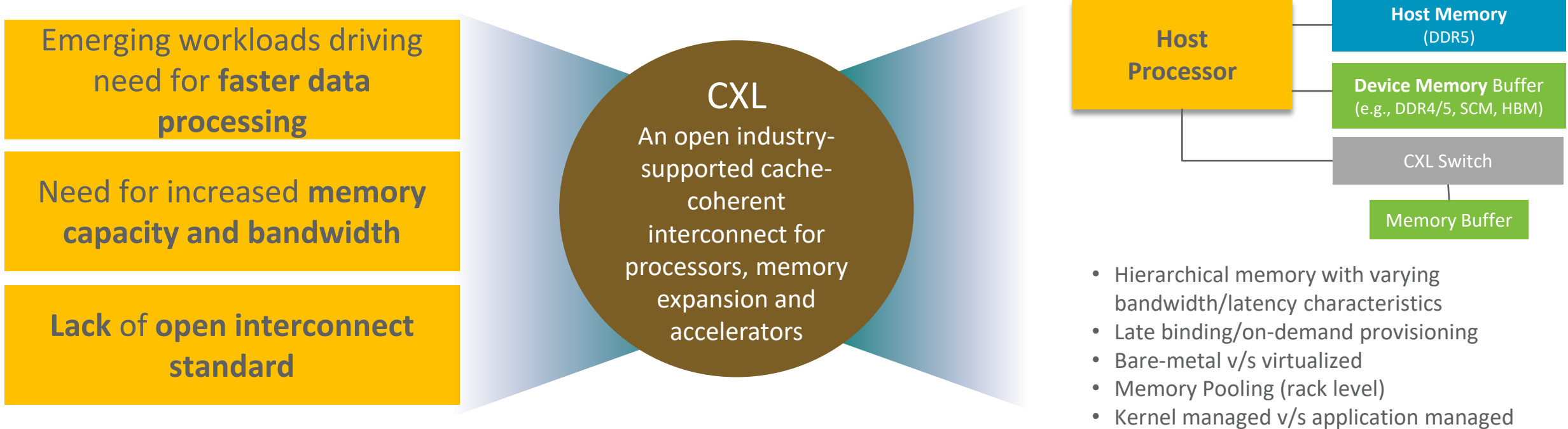
TCO

- Cost impact of min capacity increase and Die/ECC overheads
- Applications can trade performance/capacity to achieve optimal TCO

"Memory" Pyramid today



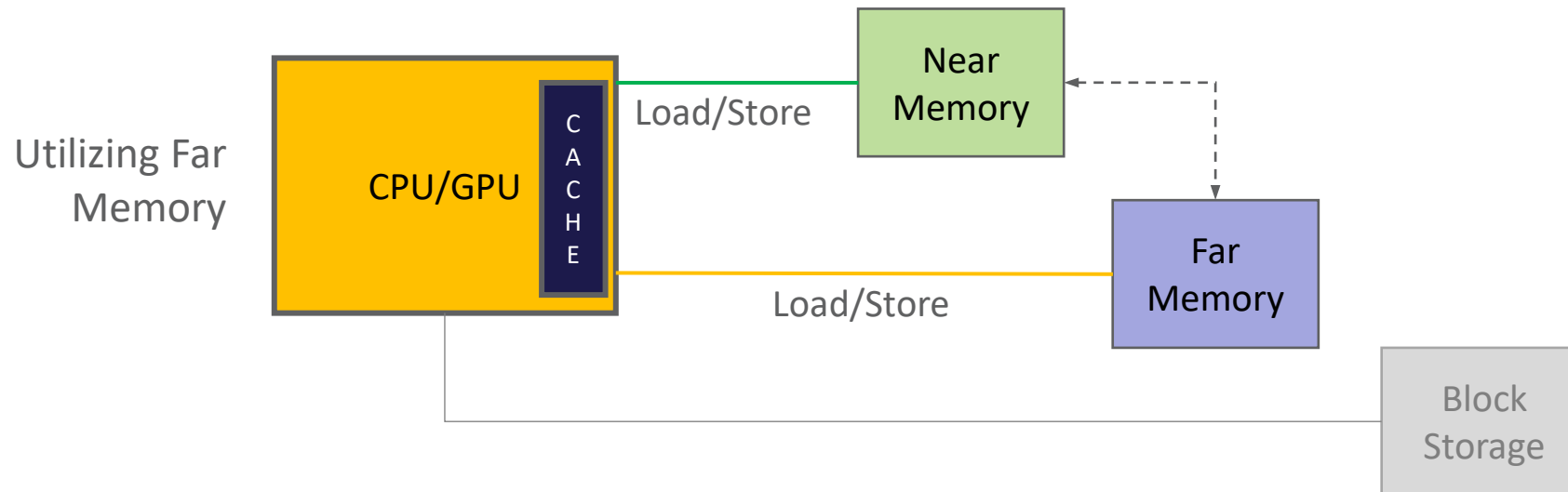
Hierarchical Memory: Need for SDM



Software-Defined Memory (SDM) is needed for optimizing resource utilization to deliver lower TCO and power-efficiency

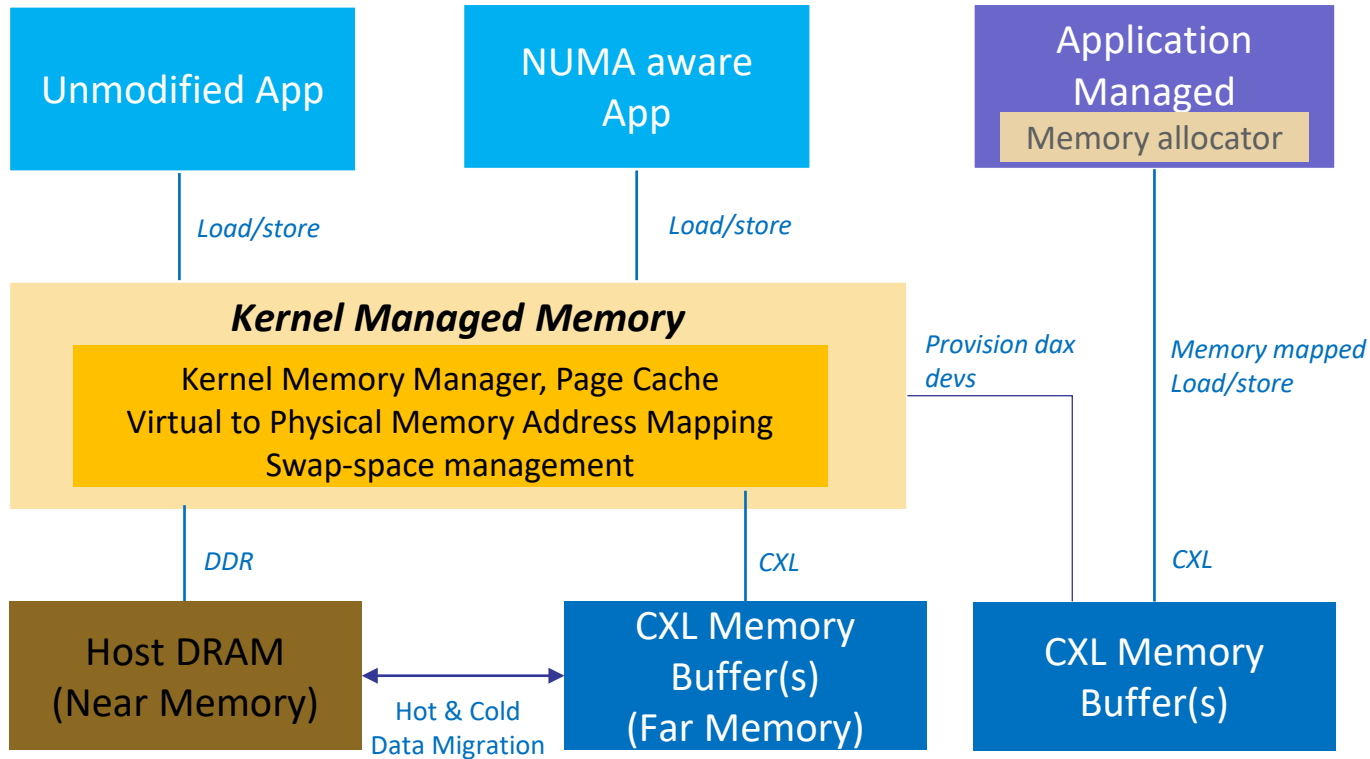
SDM - Introducing Far Memory

- Memory needs are growing faster than underlying memory technology
- Tiered memory can provide additional capacity at appropriate performance (with load/store) to sustain application needs
- Initial SDM opportunity is around CXL 1.1 Memory expansion for workload acceleration and lower costs with CXL connected DRAM



CXL Memory tier can enable performance/capacity trade off to achieve TCO gains

SDM - Kernel & Application Managed Memory



Kernel Managed

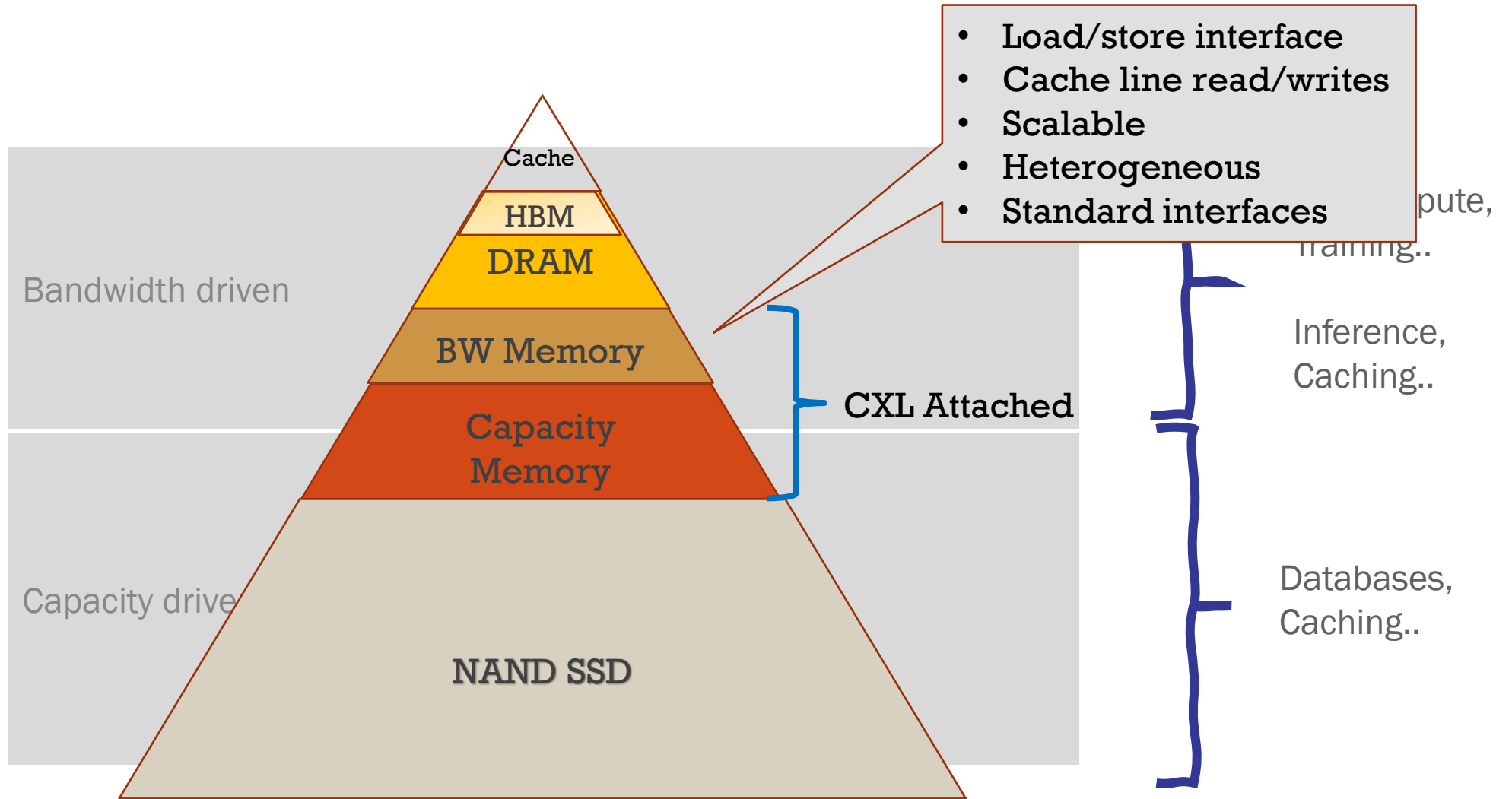
- OS can **map far memory** into application's **virtual address space**
- Applications **can execute from pages in far memory** (albeit more slowly)
- Kernel memory manager can implement **varying policies for migrating hot & cold pages** between tiers

Application Managed

- **Allows apps to access CXL memory as memory-mapped files**
- **Built on top of DAX** (Direct Access) file system

Linux Kernel tiering in early development stages

"Tiered Memory" Pyramid with CXL



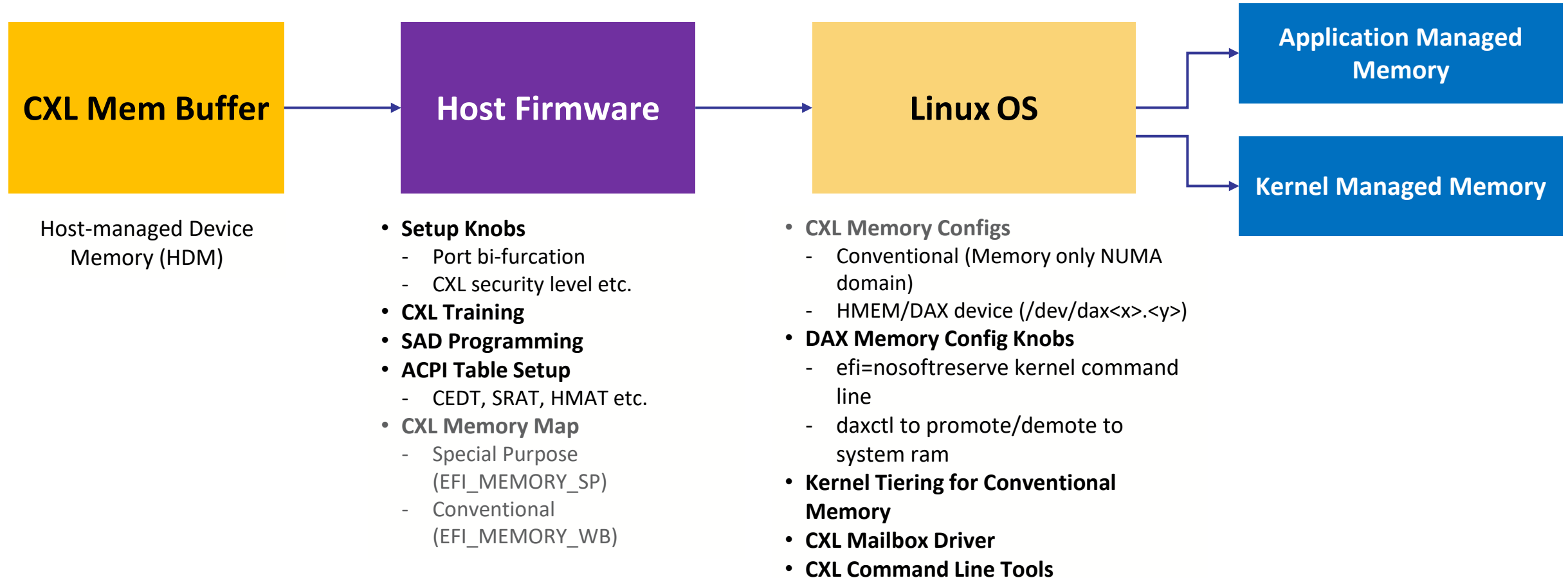
Memory Technologies

	Compute	Storage	Training	Inference
DDR4	✓	✓	✓	✓
DDR5	✓	✓	✓	✓
HBM			✓	✓
CXL+DDR	✓	✓		✓
SCM (PCIe/CXL) [Exploration Phase]		✓	✓	✓

Use Case Examples

- **Caching (e.g. Memcache/Memtier (Cachelib), Redis etc.)**
 - Need to achieve higher QPS while satisfying “retention time”
 - Higher memory capacity needed
 - Current solutions include “tiered memory” with DRAM+NAND, but need load/store
- **Databases (E.g. RocksDB/MongoDB etc.)**
 - Need to achieve efficient storage capacity per node and deliver QPS SLA
 - Higher amount of memory enables more storage per node
- **Inference (E.g. DLRM)**
 - Petaflops and Number of parameters are increasing rapidly
 - AI Models are scaling faster than the underlying memory technology
 - Current solutions include “tiered memory” with DRAM+NAND, but need load/store

Intel 4th Gen Xeon Processor (Sapphire Rapids): CXL Provisioning (*Caching Use Case*)



Intel 4th Gen Xeon Processor (Sapphire Rapids): Caching Use Case (CacheBench)

CacheBench

Caching Application (e.g. CacheBench)

CacheLib API

find,
allocate

ItemHandle

Insert Or
Replace

load/store

block IO

Memory Cache

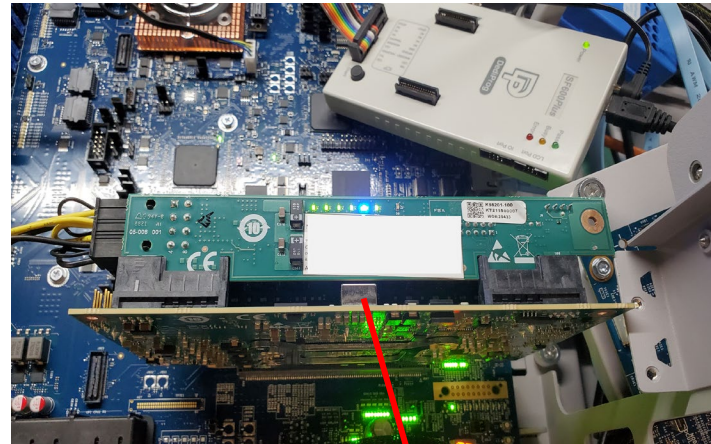
Block Cache

Policy
based
migration

Github: <https://github.com/facebook/CacheLib>

Open-source benchmarking tool for
evaluating caching performance

Intel Sapphire Rapids Pre-production Platform



Intel Pre-production CXL
FPGA Memory Buffer

CXL Memory in Linux OS

```
root@ac2i:~# root@ac2i:~# lspci -v -s 0000:49:00.0
-bash: root@ac2i:~# command not found
root@ac2i:~#
root@ac2i:~# cat y.out
root@ac2i:~# numactl -H
available: 3 nodes (0-2)
node 0 cpus:
node 0 size: 15725 MB
node 0 free: 13261 MB
node 1 cpus:
node 1 size: 16104 MB
node 1 free: 14219 MB
node 2 cpus:
node 2 size: 64509 MB
node 2 free: 63651 MB
node distances:
node  0  1  2
0:  10  21  14
1:  21  10  24
2:  14  24  10

root@ac2i:~# lspci -v -s 0000:49:00.0
49:00.0 Memory controller [0502]: Intel Corporation Device 0d93 (rev 01) (prog-if 10)
Flags: fast devsel, IRQ 255, IOMMU group 442
Memory at b3200000 (32-bit, non-prefetchable) [disabled] [size=2M]
Memory at 204eff080000 (64-bit, prefetchable) [disabled] [size=64K]
Memory at 204efc000000 (64-bit, non-prefetchable) [disabled] [size=16M]
Capabilities: [40] Express Root Complex Integrated Endpoint, MSI 00
Capabilities: [80] MSI: Enable- Count=1/4 Maskable+ 64bit+
Capabilities: [a0] Power Management version 3
Capabilities: [100] Advanced Error Reporting
Capabilities: [200] Multi-Function Virtual Channel <?>
Capabilities: [300] Virtual Channel
Capabilities: [550] Multicast
Capabilities: [588] Latency Tolerance Reporting
Capabilities: [5b0] Transaction Processing Hints
Capabilities: [6e0] Address Translation Service (ATS)
Capabilities: [700] Physical Resizable BAR
Capabilities: [b20] Page Request Interface (PRI)
Capabilities: [b40] Process Address Space ID (PASID)
Capabilities: [b50] Precision Time Measurement
Capabilities: [b80] Single Root I/O Virtualization (SR-IOV)
Capabilities: [c00] Device Serial Number 00-00-00-00-00-00-00-00
Capabilities: [d00] Vendor Specific Information: ID=0040 Rev=1 Len=06c <?>
Capabilities: [e00] Designated Vendor-Specific: Vendor=1e98 ID=0000 Rev=0 Len=56: CXL
Capabilities: [e40] Designated Vendor-Specific: Vendor=1e98 ID=0008 Rev=0 Len=36 <?>
Capabilities: [ed0] Designated Vendor-Specific: Vendor=8086 ID=0050 Rev=0 Len=12 <?>
Capabilities: [fee0] Vendor Specific Information: ID=0043 Rev=0 Len=010 <?>
```

OCP SDM activity and progress

- SDM's focus: Apply emerging memory technologies in the development of use cases
- The OCP SDM group has three real-world memory focus areas:
 - Databases/Caching
 - AI/ML & HPC
 - Virtualized Servers
- SDM Team Members: AMD, ARM, Intel, Meta, Micron, Microsoft, Omdia, Samsung, VMWare
- Vendors are demonstrating CXL Capable CPUs and devices
- Meta and others are investigating solutions to real world memory problems

Call to Action: Join OCP SDM workstream!

Notices & Disclaimers

Intel technologies may require enabled hardware, software or service activation. Your costs and results may vary.

No product or component can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. For more complete information about performance and benchmark results, visit <http://www.intel.com/benchmarks>.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit <http://www.intel.com/benchmarks>.

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings. Circumstances will vary. Intel does not guarantee any costs or cost reduction.

Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.