

STORAGE DEVELOPER CONFERENCE



Fremont, CA  
September 12-15, 2022

*BY Developers FOR Developers*

A **SNIA** Event

# Software Defined Memory with CXL and Tiered Memory to Enable Hyperscale Use Cases

Manoj Wadekar, Hardware Systems Technologist, Meta  
Anjaneya "Reddy" Chagam, Cloud Architect, Intel

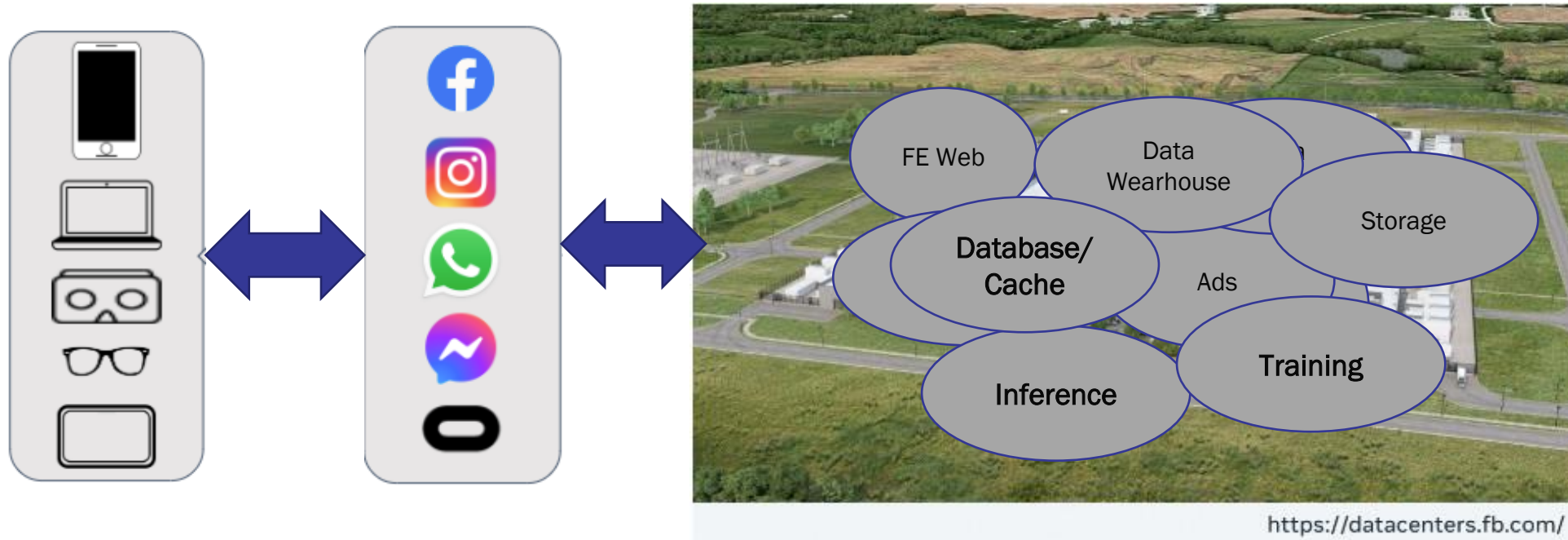
# Agenda

- Memory Challenges in Hyperscale Infrastructure
- Need for Software Defined Memory (SDM)
- SDM Use-cases
- SDM on Intel 4th Gen Xeon Processor
- Summary & Call to Action

# Agenda

- Memory Challenges in Hyperscale Infrastructure
- Need for Software Defined Memory (SDM)
- SDM Use-cases
- SDM on Intel 4th Gen Xeon Processor
- Summary & Call to Action

# 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

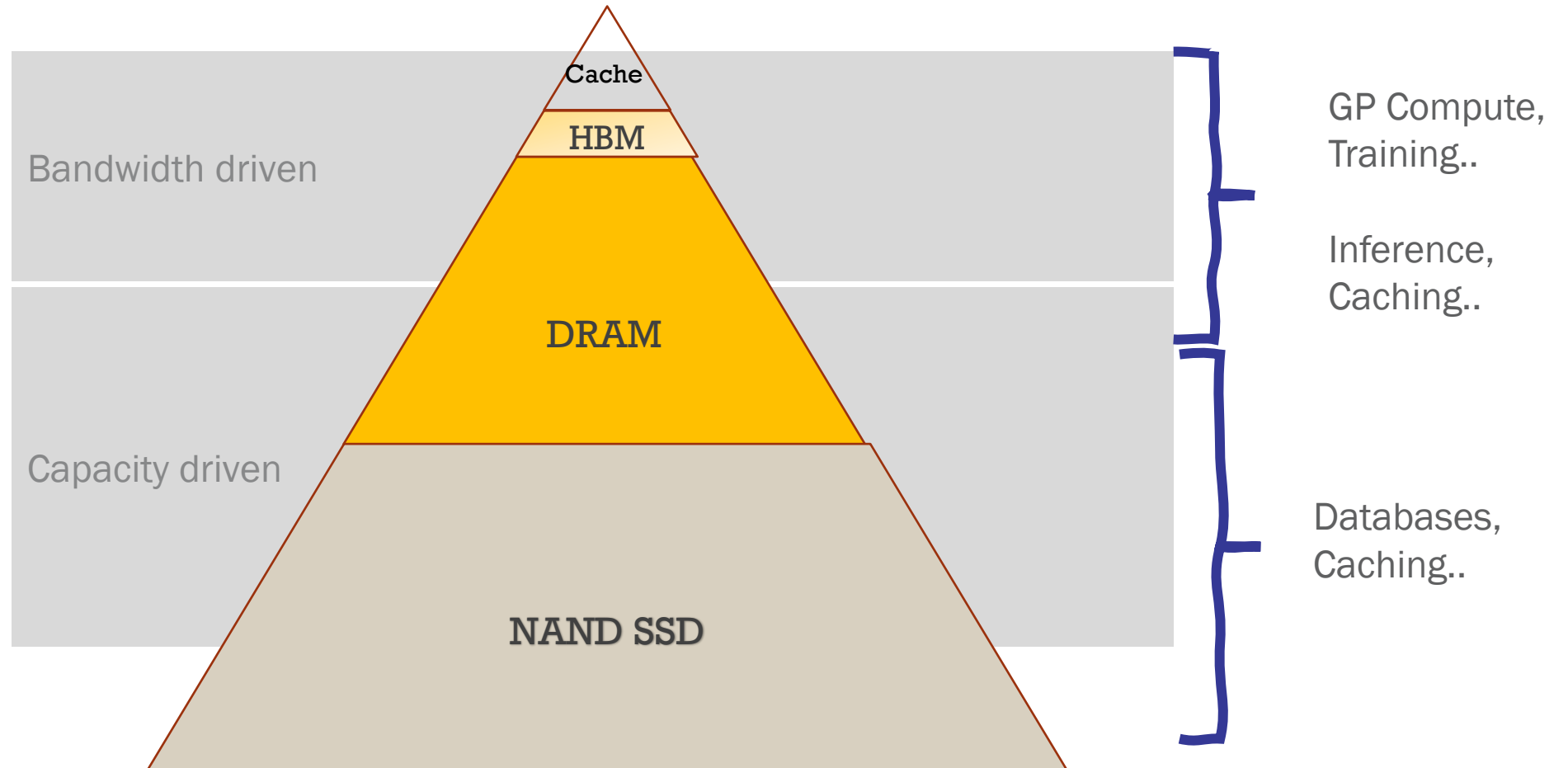
## Power

- 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



# Agenda

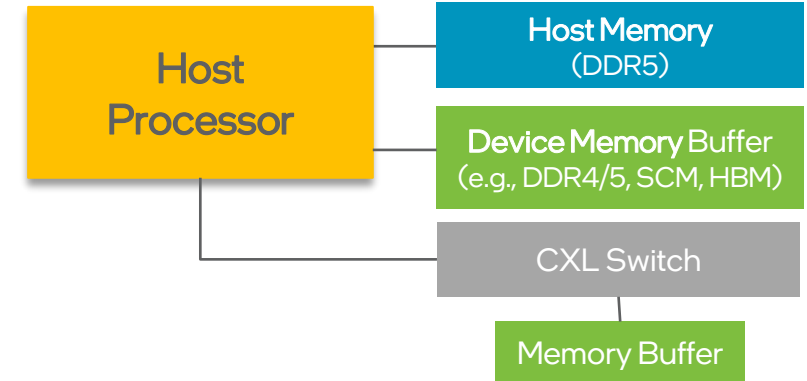
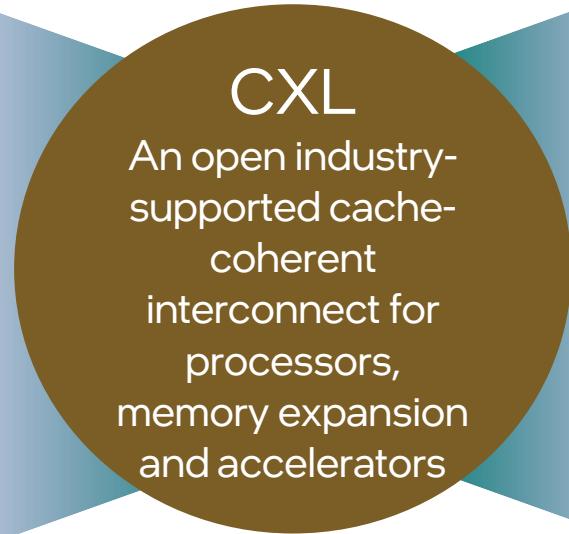
- Memory Challenges in Hyperscale Infrastructure
- **Need for Software Defined Memory (SDM)**
- SDM Use-cases
- SDM on Intel 4th Gen Xeon Processor
- Summary & Call to Action

# Hierarchical Memory: Need for SDM

Emerging workloads driving need for **faster data processing**

Need for increased **memory capacity and bandwidth**

**Lack of open interconnect standard**



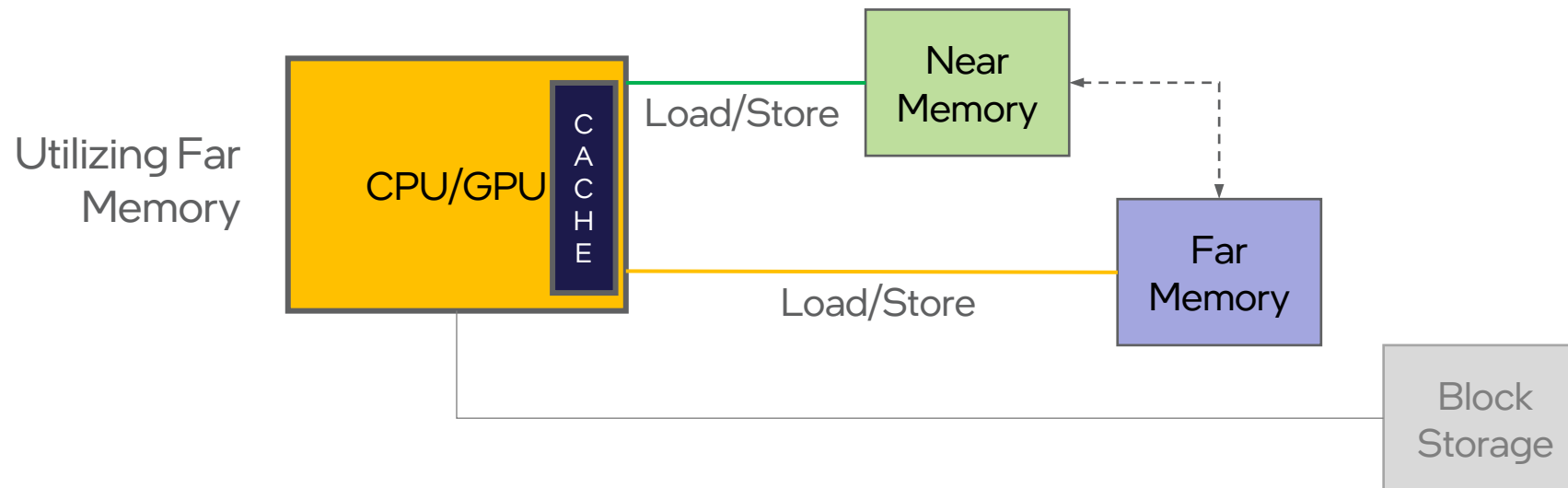
- Hierarchical memory with varying bandwidth/latency characteristics
- Late binding/on-demand provisioning
- Bare-metal v/s virtualized
- Memory Pooling (rack level)
- Kernel managed v/s application managed

**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

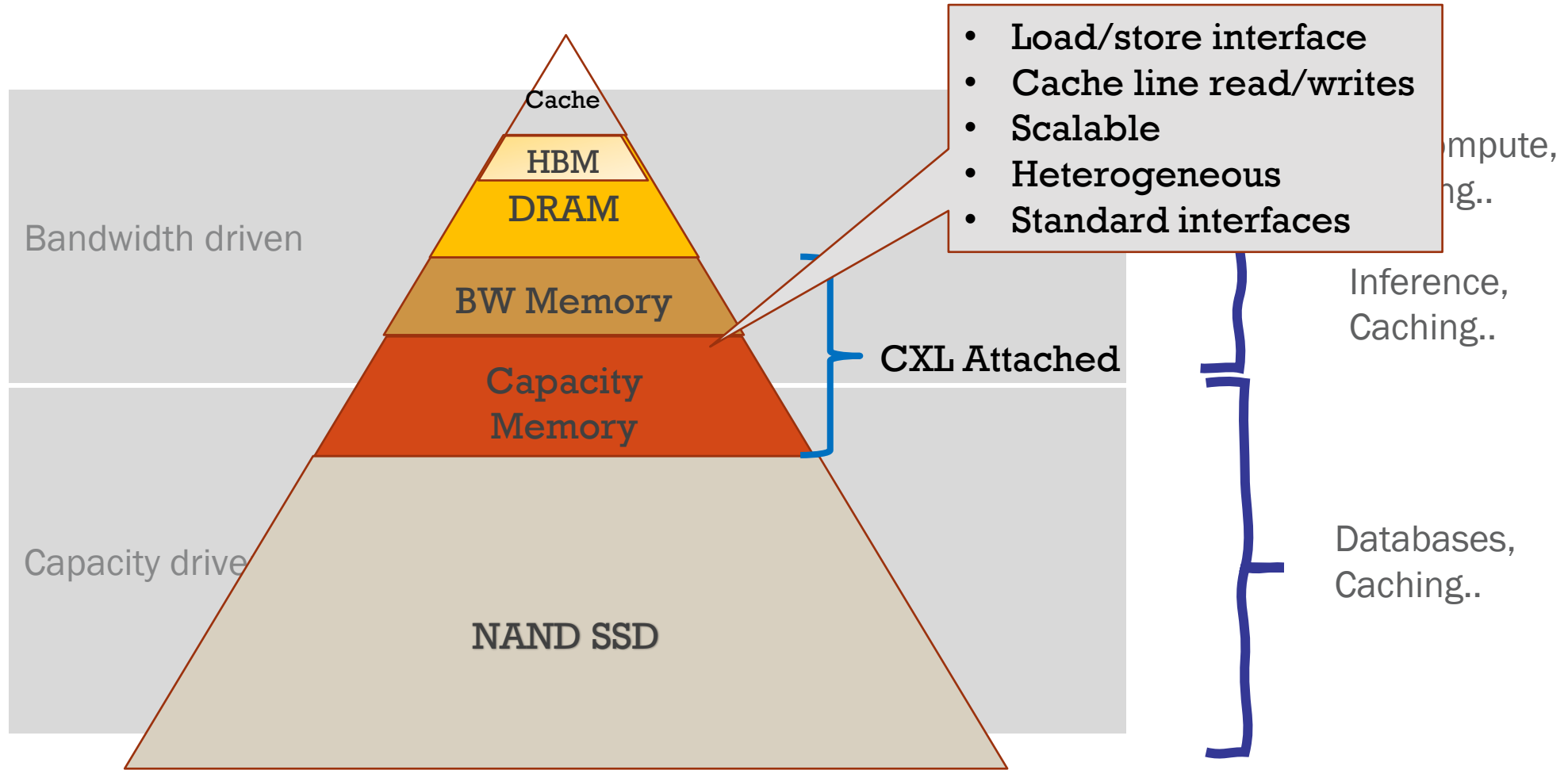


*CXL Memory tier can enable performance/capacity trade off to achieve TCO gains*

# Agenda

- Memory Challenges in Hyperscale Infrastructure
- Need for Software Defined Memory (SDM)
- **SDM Use-cases**
- SDM on Intel 4th Gen Xeon Processor
- Summary & Call to Action

# "Tiered Memory" Pyramid with CXL



# Memory Technologies

	Compute	Storage	Training	Inference
DDR4	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DDR5	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
HBM			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
CXL+DDR	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>		<input checked="" type="checkbox"/>
SCM (PCIe/CXL) [Exploration Phase]		<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

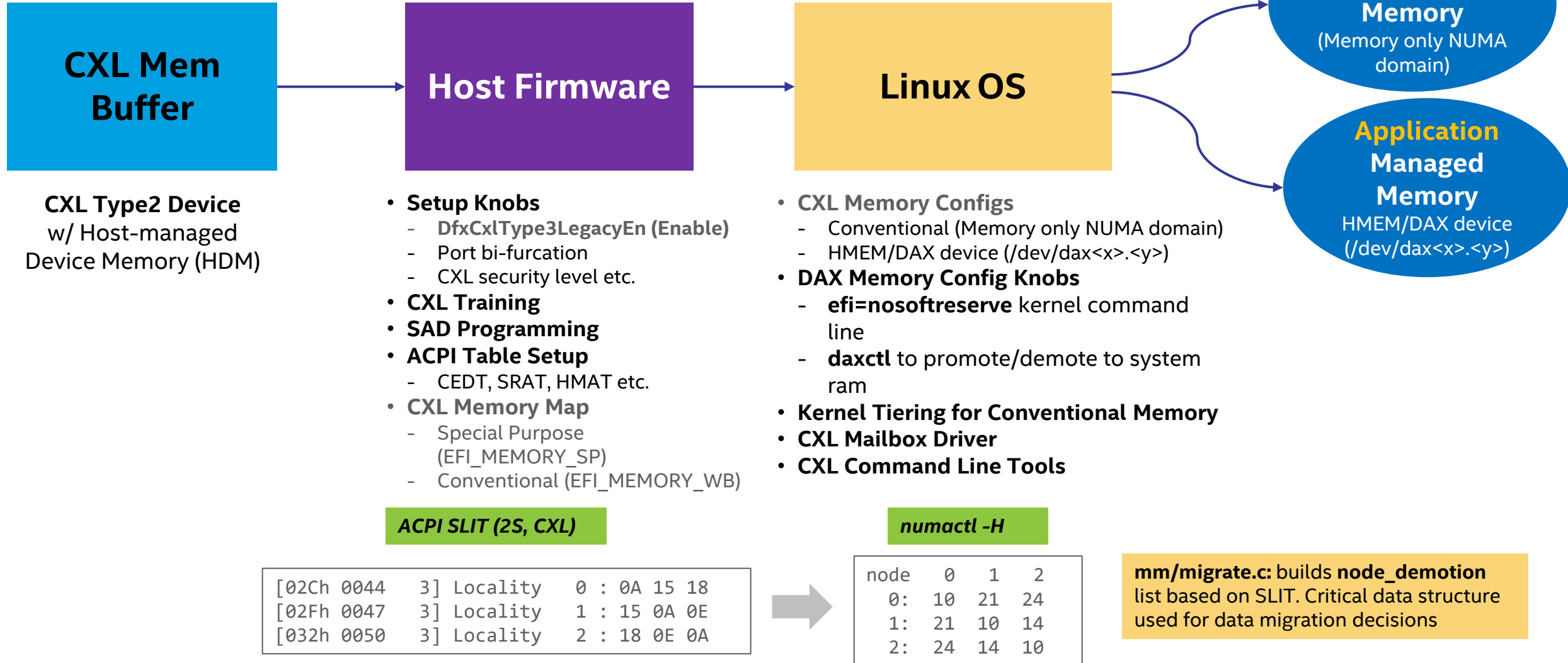
# Use Case Examples

- **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

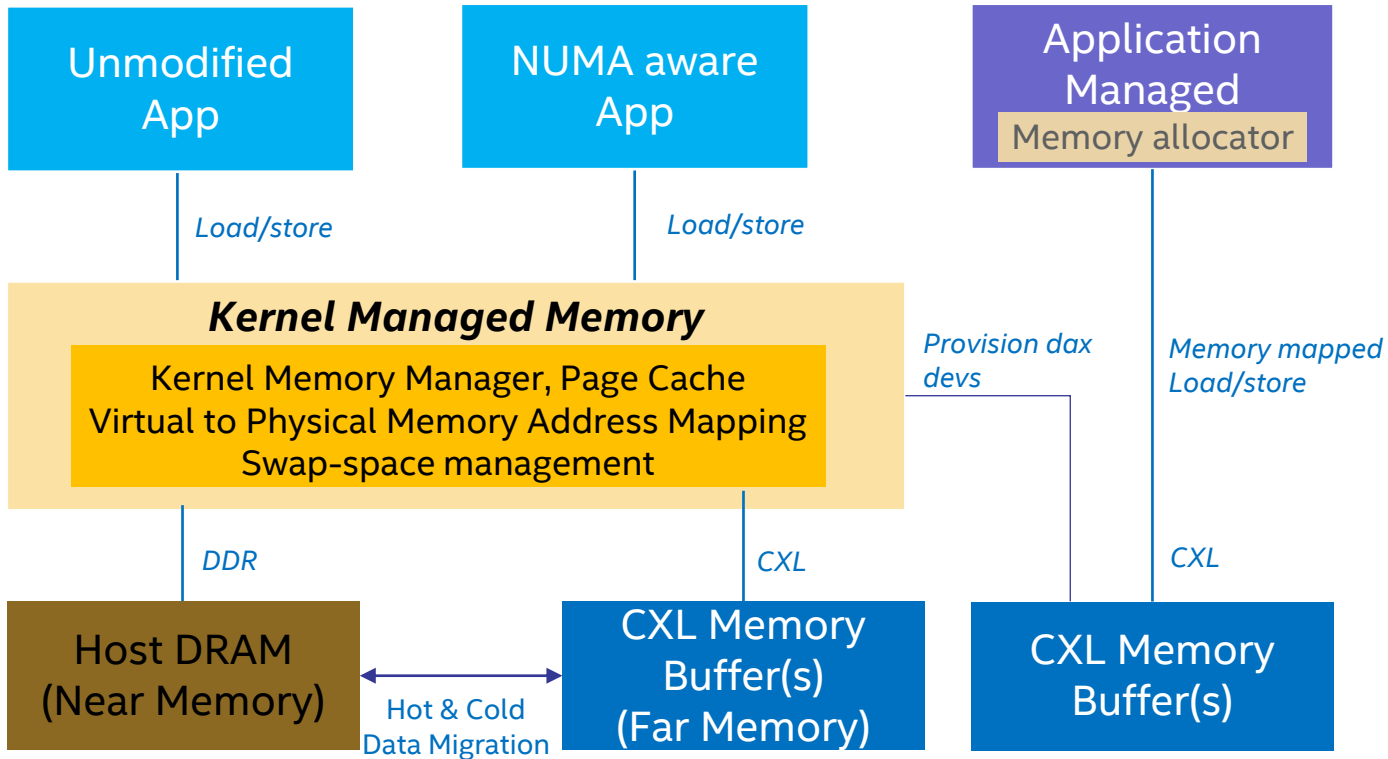
# Agenda

- Memory Challenges in Hyperscale Infrastructure
- Need for Software Defined Memory (SDM)
- SDM Use-cases
- **SDM on Intel 4th Gen Xeon Processor**
- Summary & Call to Action

# Intel 4th Gen Xeon Processor CXL Memory Buffer Provisioning



# 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

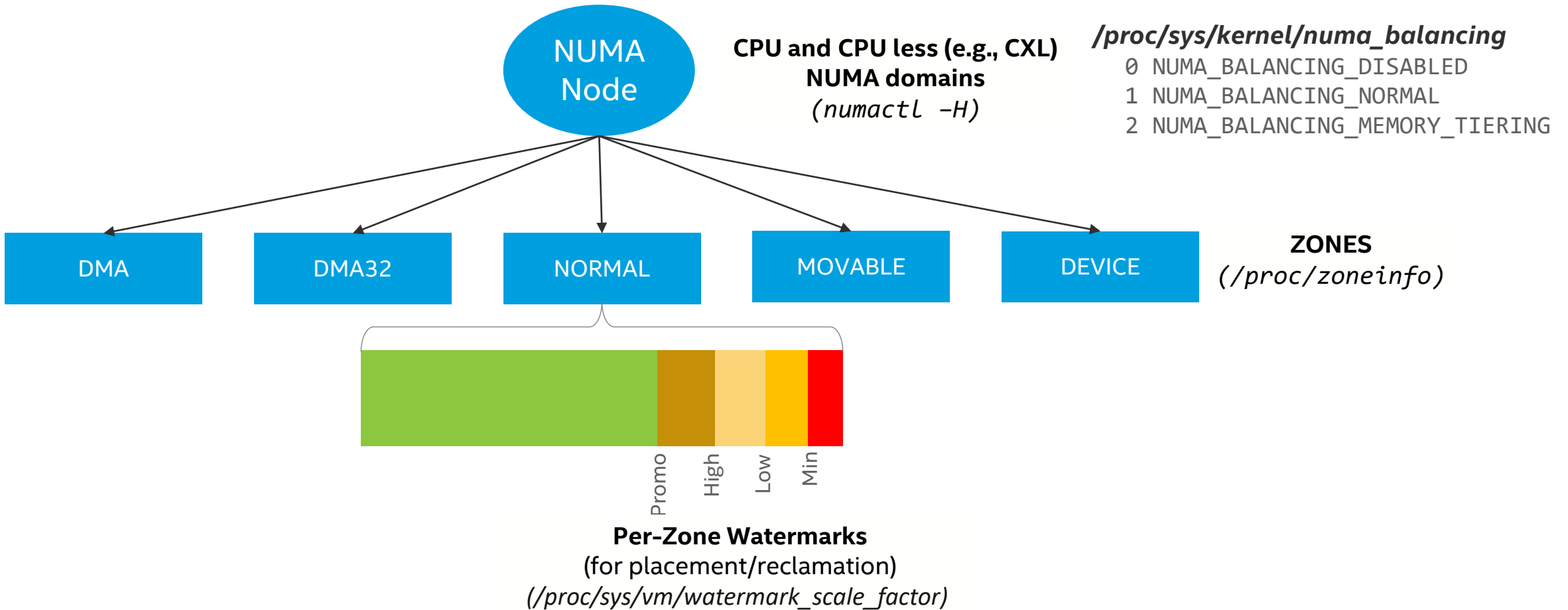
## Application Managed

- 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*

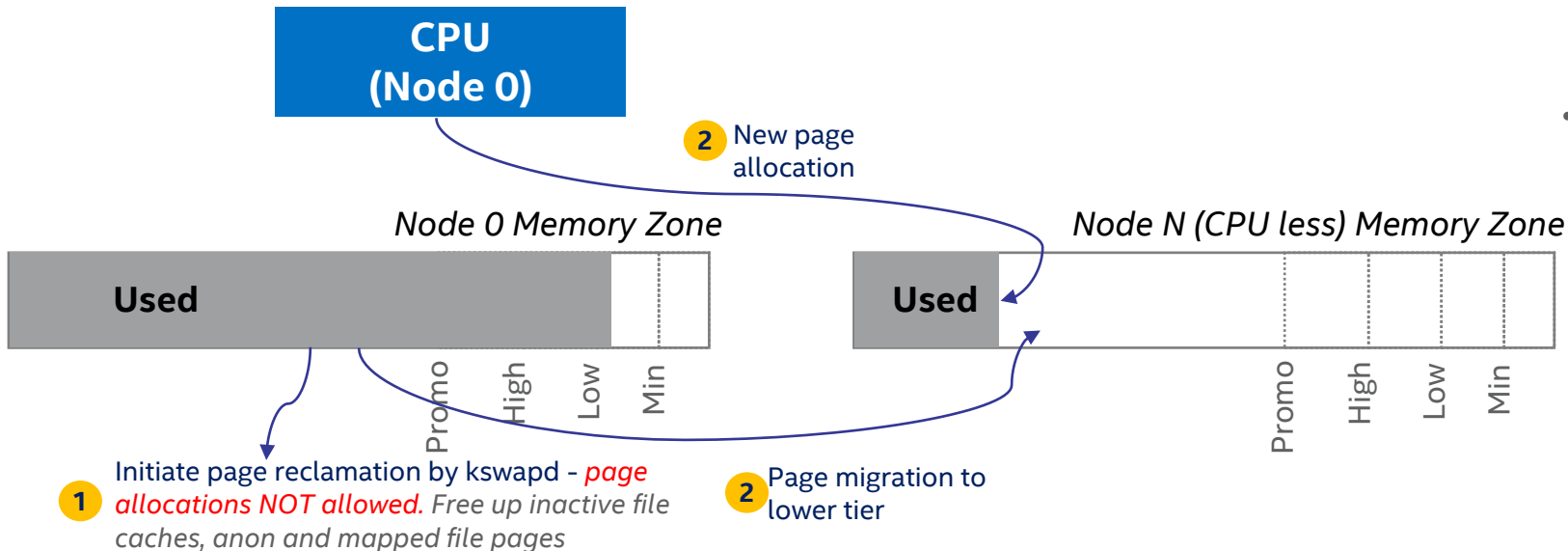
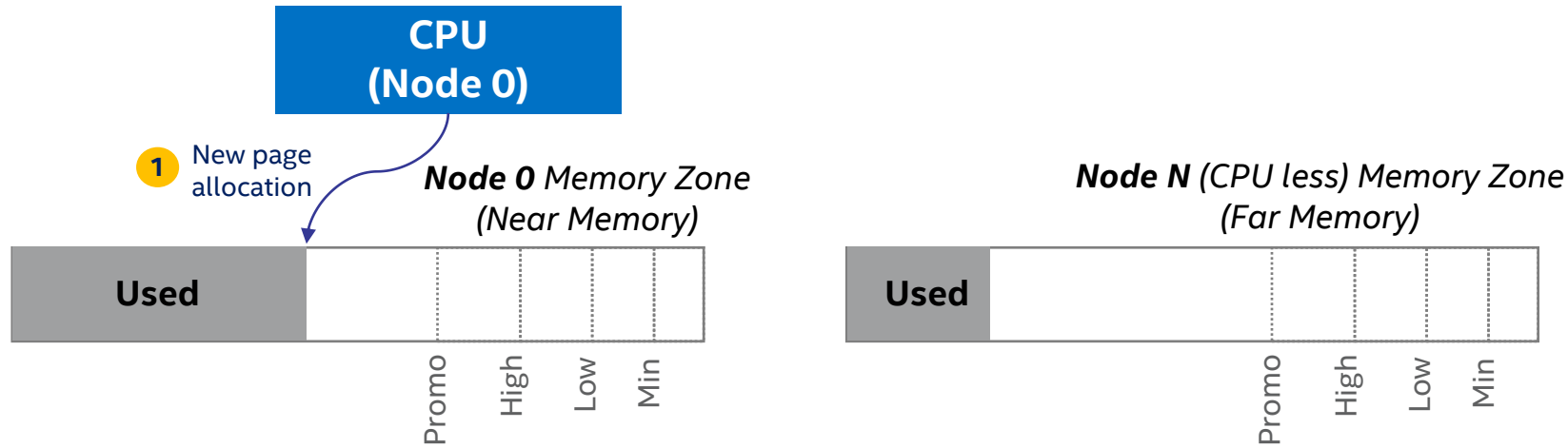


# Linux Kernel – Memory Layout



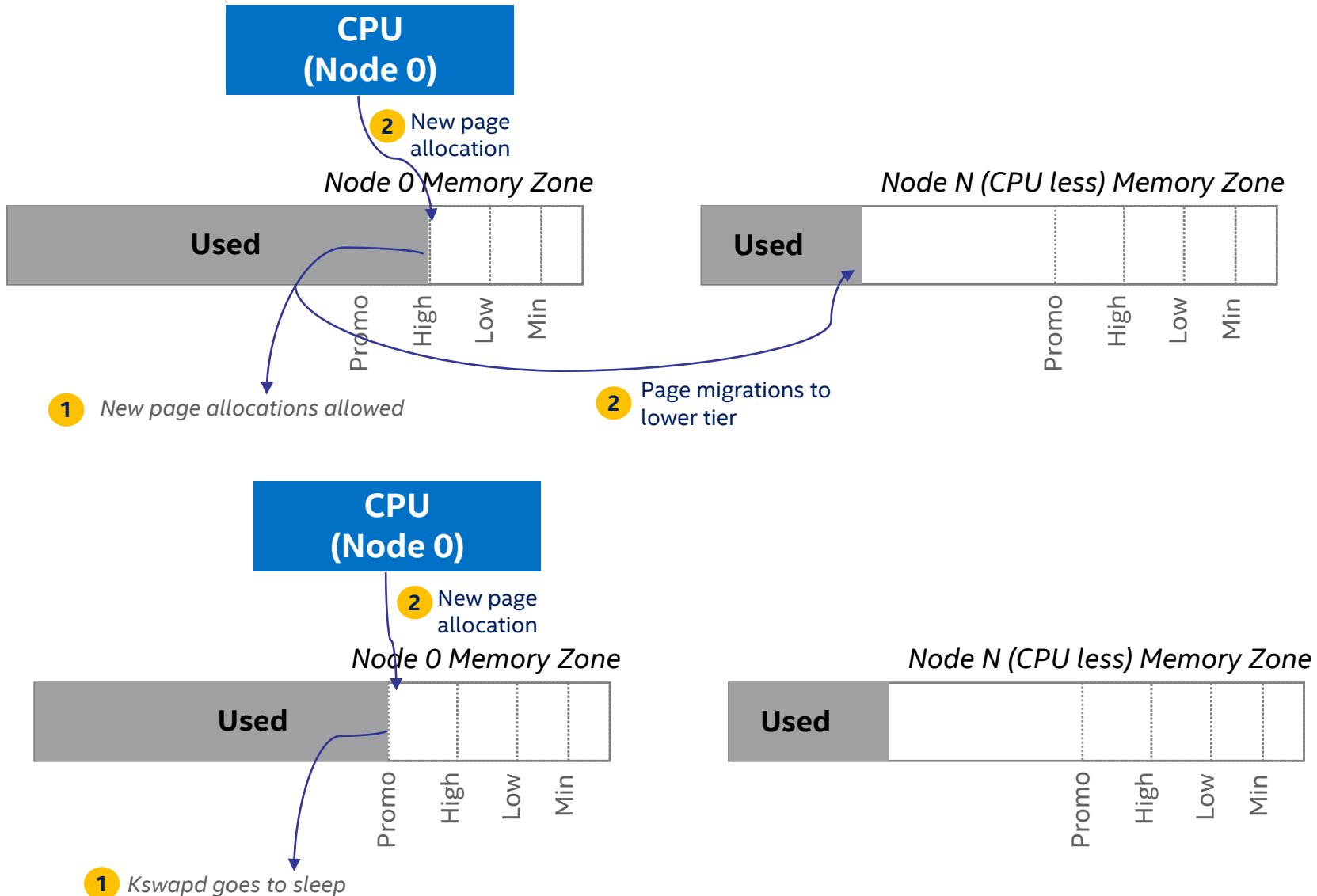
- **Low** - when #of free pages reaches **low** watermark, kswapd is woken up to start freeing pages
- **Min** - when #of free pages reaches **min** watermark, page allocator will do the kswapd work in a synchronous fashion (i.e. direct-reclaim path)
- **High** - zone is not “balanced” until #of free pages reaches high watermark and considered for placement
- **Promo** - reclamation watermark to keep enough buffer for placement and promotion from lower tiers. Once the watermark has been reached, kswapd will go back to sleep. *NOTE: Enabled only when `/proc/sys/kernel/numa_balancing` is `NUMA_BALANCING_MEMORY_TIERING (2)`*

# Linux Kernel – Placement/Reclamation



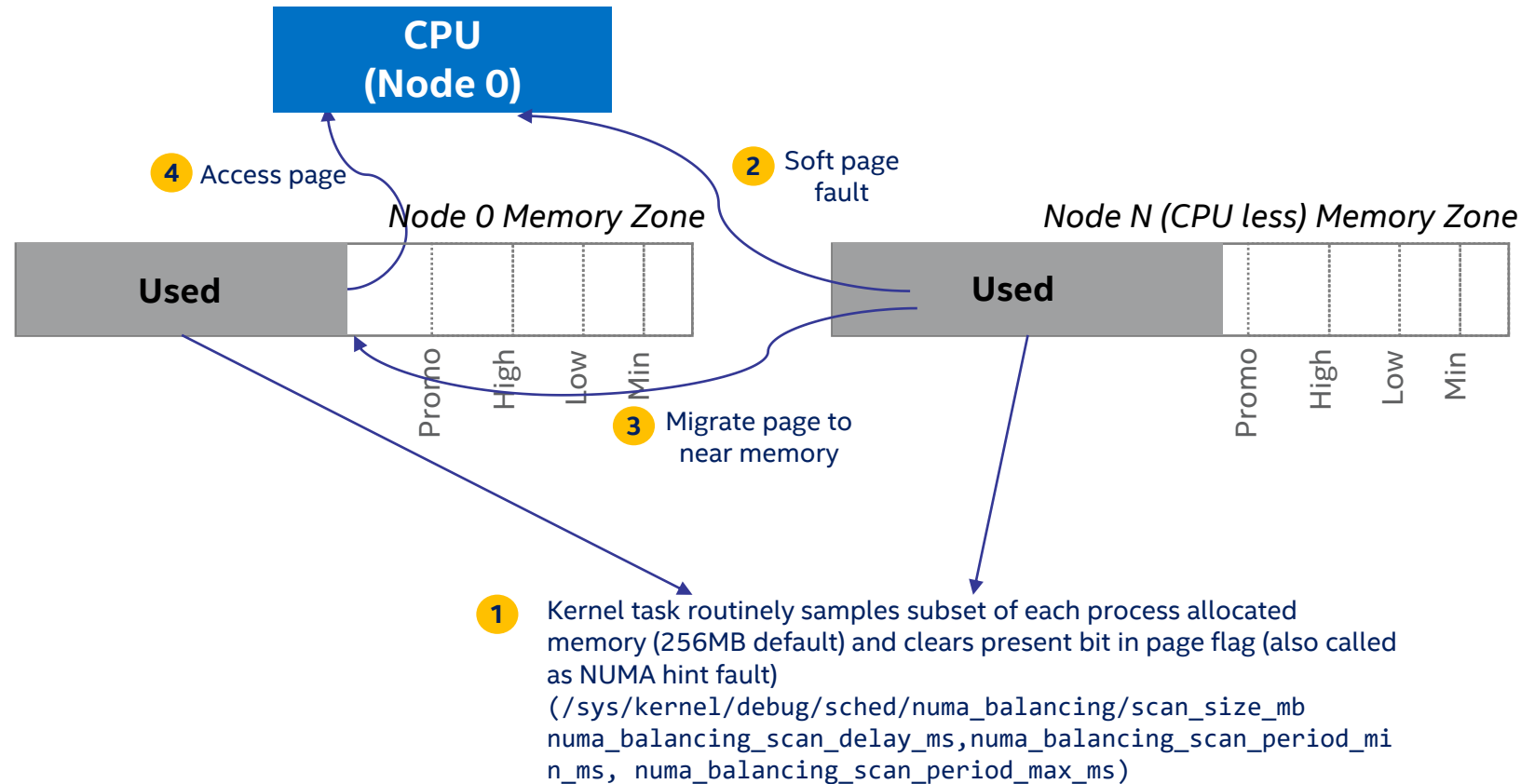
- `/sys/kernel/mm/numa/demotion_enabled` needs to be enabled (default is false) to demote pages to lower tier during reclamation
- **numastat** to show per-NUMA-node memory statistics for processes and the operating system

# Linux Kernel – Placement/Reclamation (Contd)



- **stats: /proc/vmstat**
  - **numa\_pages\_migrated** - # of pages migrated among tiers which includes inter-socket NUMA balancing as well
  - **numa\_pte\_updates**: # of pages marked as inaccessible which are cleared later by a NUMA hinting faults
  - **numa\_huge\_pte\_updates** - # of huge pages marked for NUMA hinting faults
  - **numa\_hint\_faults** - #of NUMA hint faults)
  - **numa\_hint\_faults\_local** - #of NUMA hint faults to local nodes
  - **pgdemote\_kswapd/pgdemote\_direct** - #of pages demoted to lower node

# Linux Kernel – Promotion

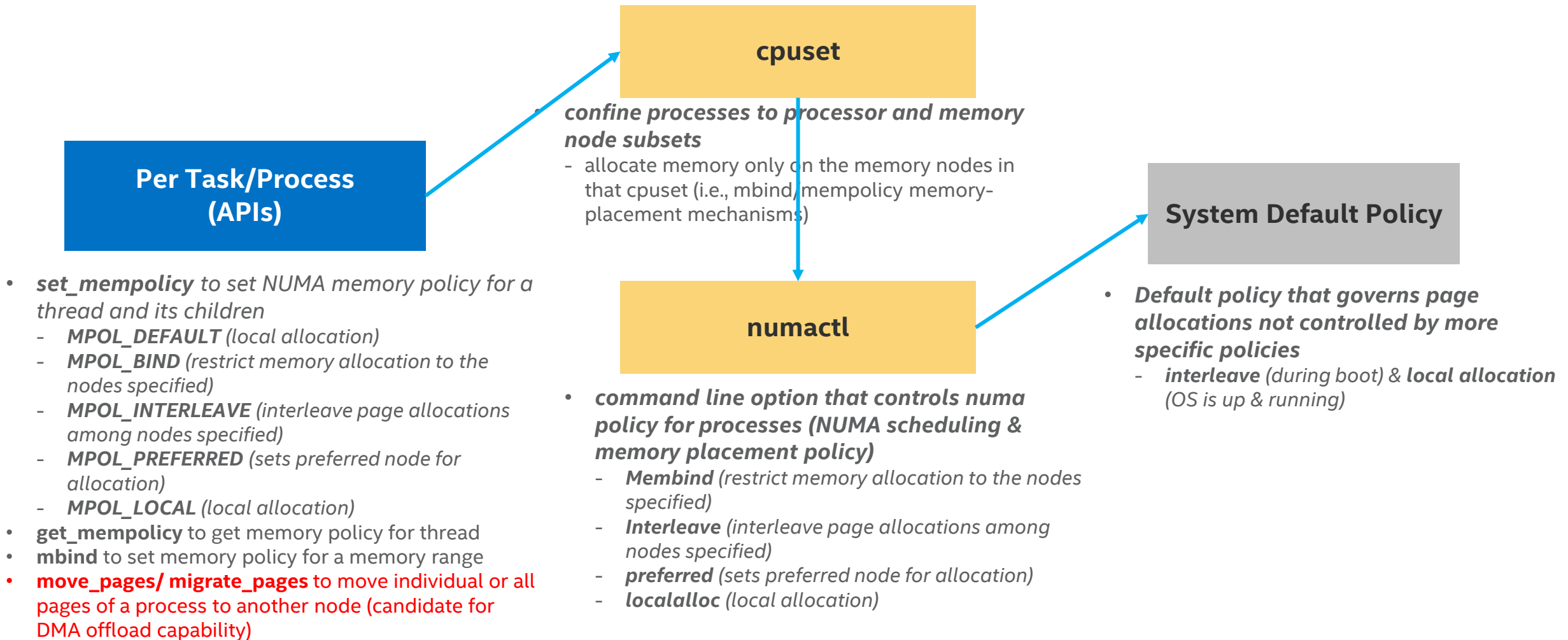


- **stats: /proc/vmstat**
  - **pgpromote\_success** - #of promoted pages to top tier node

## Meta TPP patches (pending):

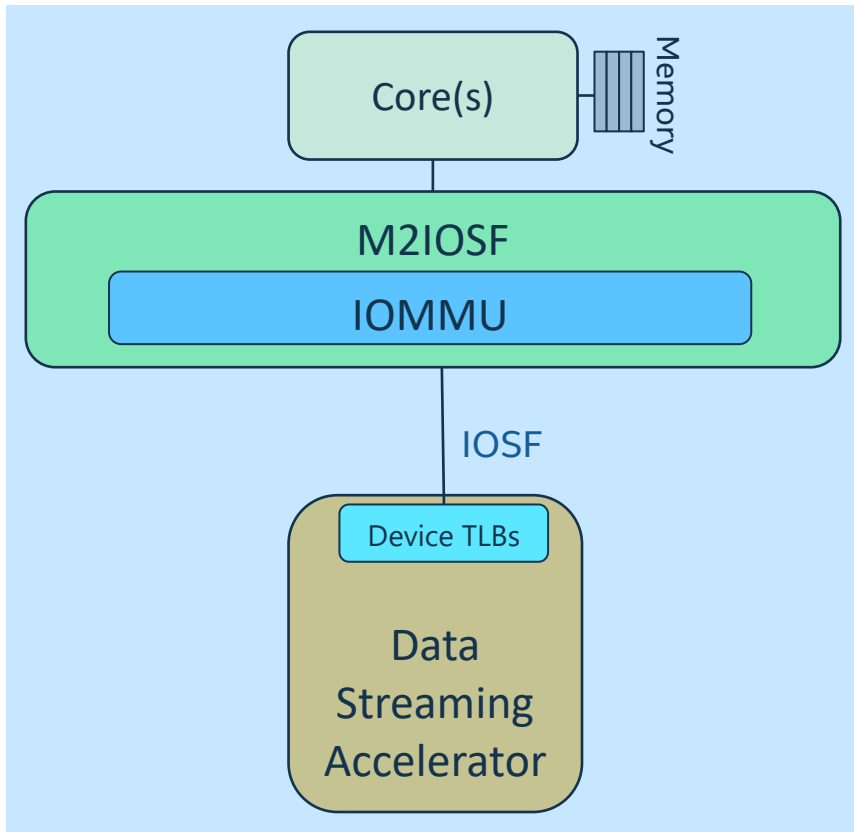
- **LRU-filtered promotion** to avoid ping-pong. If the faulted page is in inactive LRU, we do not consider the page as the promotion candidate instantly as it might be an infrequently accessed page. We consider the faulted page as a promotion candidate only if it is found in the active LRUs (either of anon or file active LRU)  
<https://github.com/hnaz/linux/commit/8040b0b808b69ba680f53979ca9c5db800b2c8a6>
- **Promotion and demotion related statistics:**  
<https://lore.kernel.org/lkml/e3db05f4bd7dd9363c2a895875505088d3273bb8.1637778851.git.hasanamaruf@fb.com/#r>

# Linux Kernel – NUMA Memory Policy



**Meta TPP patch (pending) - mm: mempolicy: N:M interleave policy for tiered memory nodes**  
(<https://github.com/hnaz/linux/commit/efb59df4f19ec0da7a89131cc23b55658437e72d>)

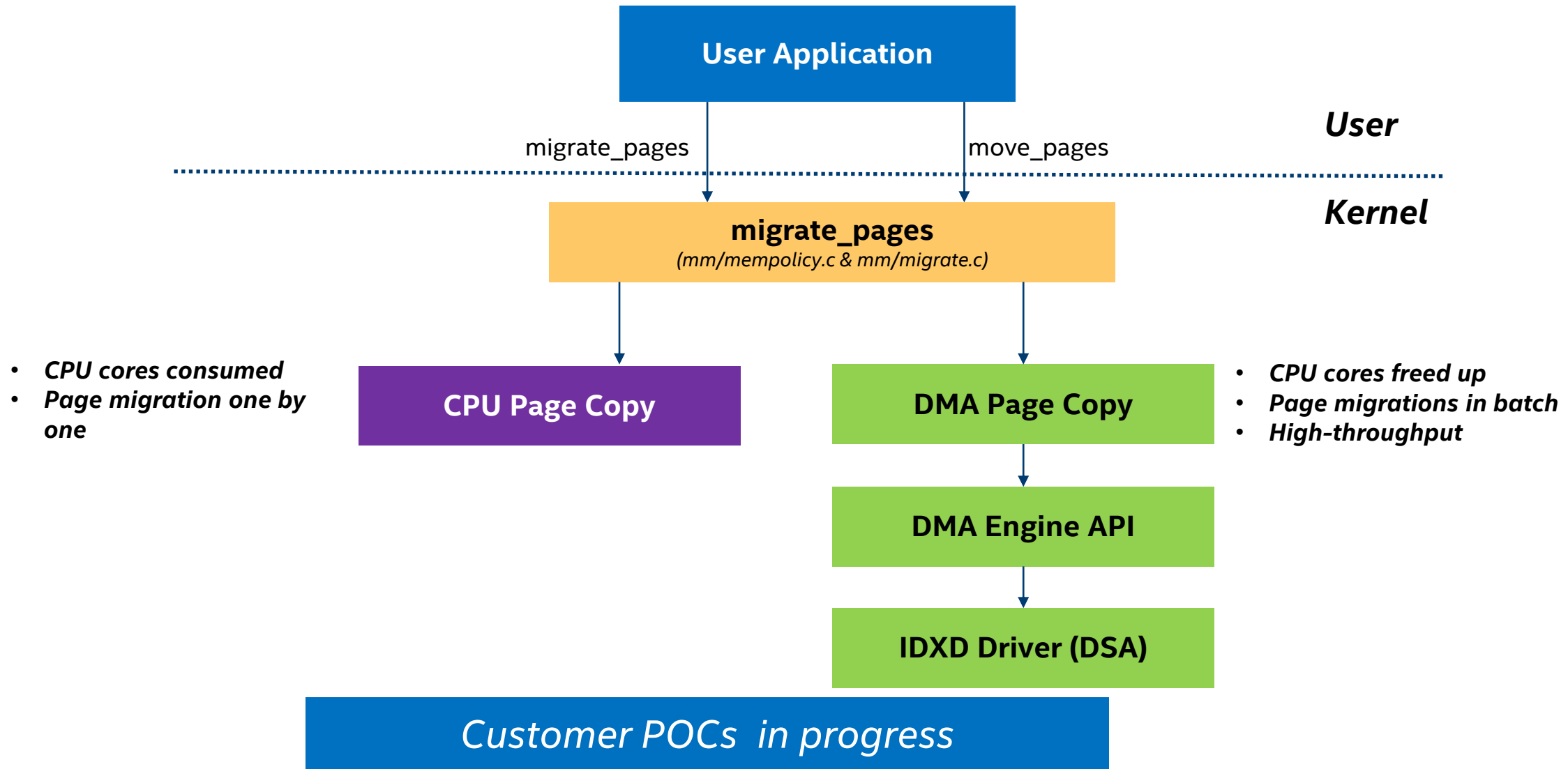
# Intel 4th Gen Xeon Processor - Data Streaming Accelerator (DSA) Overview



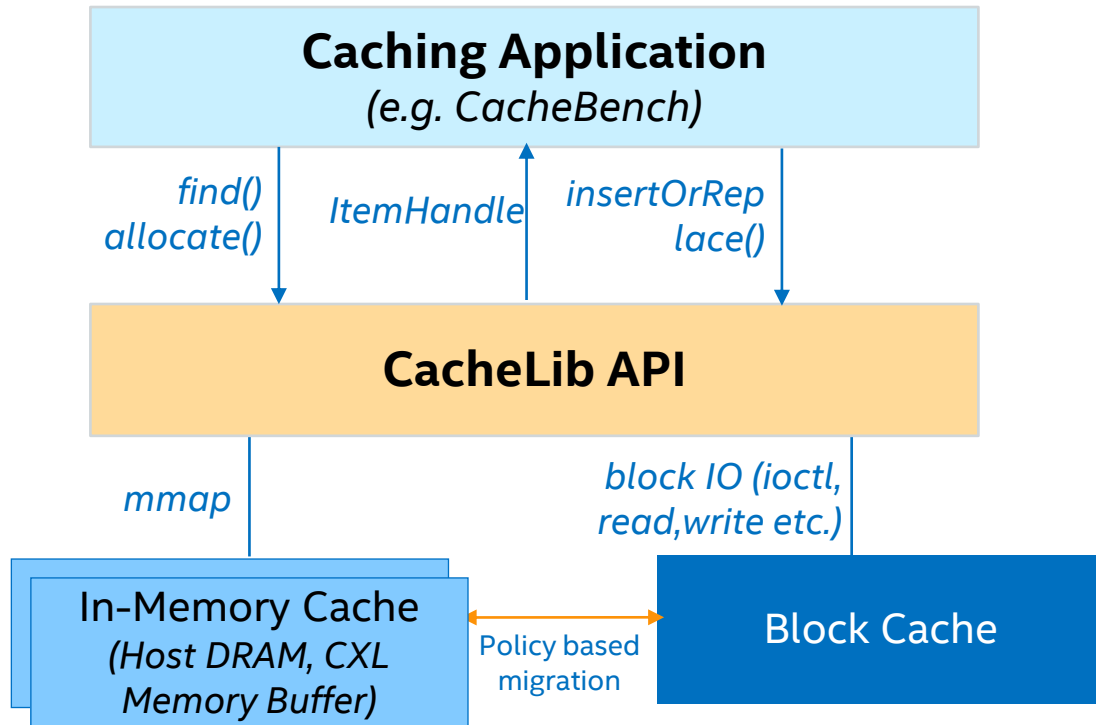
**Logical Representation**

Capability	DSA 1.0 (SPR)
Max Instances per socket	4
Max Throughput (each direction)	32 GB/s
Shared Virtual Memory support	Yes
Low latency dispatch (AIA ISA)	Yes
Recoverable page faults	Yes
DDIO caching controls	Yes
Persistent memory support	Yes
Cross-address space support	No
Num. WQs & Engines per device	8 WQs; 4 Engines
Virtualization support with SIOV	Yes
VM Migration support	Yes
Confidential compute (TDX-IO)	No
Quality of Service	WQ-Engine Groups, Read buffer controls, VCs for ATS/DRAM/PMEM
Data targets supported	DDR, PCI-E, NTB, CXL

# Linux Kernel – Page Copy Offload using DSA



# CacheLib - Overview



Github: <https://github.com/facebook/CacheLib>  
Intel development fork: <https://github.com/intel/CacheLib>

## CacheLib

- **Pluggable in-process open-source caching engine** to build and scale high-performance services
- C++ Library
- Thread-safe API
- Manages DRAM and Block Caching transparently
- Decoupled from underlying medium
- Policy based

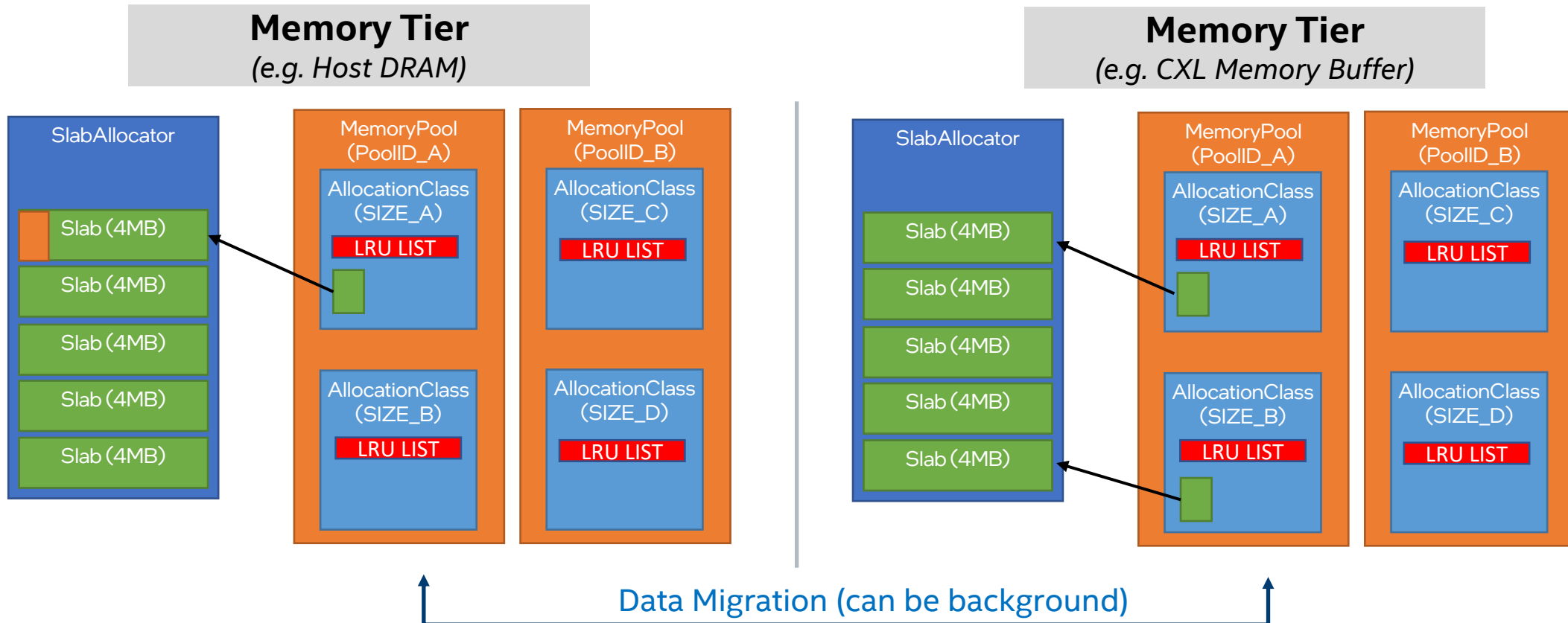
## CacheBench

- **Benchmarking tool** for evaluating caching performance
- **NUMA aware cachelib tiering or Kernel tiering** for **CXL memory buffer benchmarking**

*Memory tiering and DMA (DSA) offload development in progress*



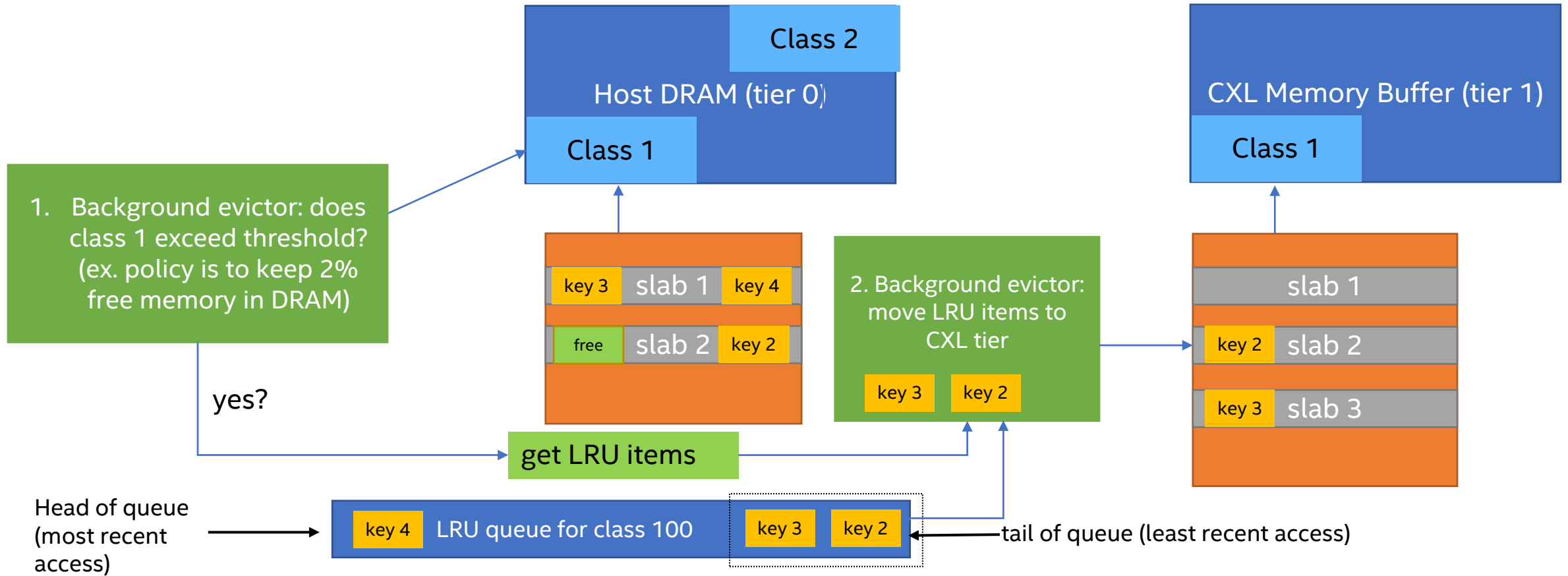
# CacheLib – Memory Tiering (Data Structures)



- Items are initially allocated in tier 1 (DRAM tier) and then migrated to tier 2 (CXL) via eviction (which can be during new item allocation or happen in background)
- Slabs are contiguous chunks of virtual address space
- Slabs are assigned to an allocation class to store items of a given size range
- Pools are independent allocators – an abstraction to support multi-tenancy
- For details refer to Cachelib architecture guide:

[https://cachelib.org/docs/Cache\\_Library\\_Architecture\\_Guide/overview\\_a\\_random\\_walk](https://cachelib.org/docs/Cache_Library_Architecture_Guide/overview_a_random_walk)

# CacheLib – Memory Tiering (Demotion Example)



NOTE: Promotion works the same way but in reverse direction

CPU and DSA offload options will be available

# Intel 4th Gen Xeon Processor: Caching Use Case

## CacheBench

### Caching Application (e.g. CacheBench)

find()  
allocate()

ItemHandle

insertOrReplace()

### CacheLib API

mmap

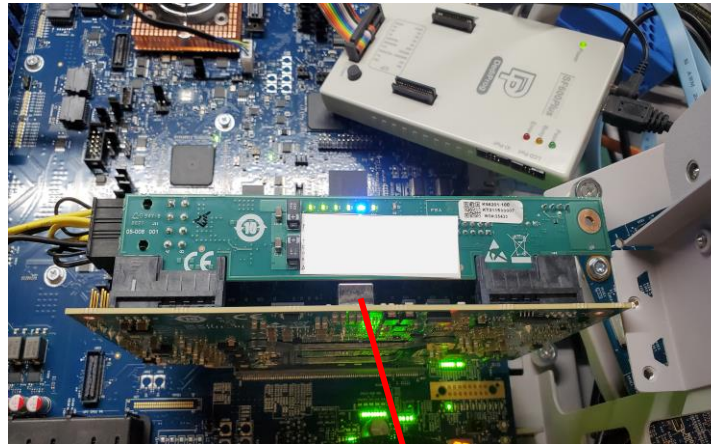
block IO (ioctl,  
read,write etc.)

In-Memory Cache  
(Host DRAM, CXL  
Memory Buffer)

Policy based  
migration

Block Cache

## Intel 4<sup>th</sup> Gen Xeon Processor Pre-production Platform



Intel Pre-production CXL  
FPGA Memory Buffer

## CXL Memory in Linux OS

```
root@ac2:~# root@ac2:~# lspci -v -s 0000:49:00.0
-bash: root@ac2:~#: command not found
root@ac2:~#
root@ac2:~#
root@ac2:~# cat y.out
root@ac2:~# 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@ac2:~# 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 204ef080000 (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: [598] 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: [e00] Designated Vendor-Specific: Vendor=8086 ID=0050 Rev=0 Len=12 <?>
Capabilities: [ee0] Vendor Specific Information: ID=0043 Rev=0 Len=010 <?>
```

# Agenda

- Memory Challenges in Hyperscale Infrastructure
- Need for Software Defined Memory (SDM)
- SDM Use-cases
- SDM on Intel 4th Gen Xeon Processor
- **Summary & Call to Action**

# Summary and Call to Action

- OCP Foundation SDM initiative is focused on applying emerging memory technologies for cloud use cases
- Newer memory technologies (e.g., HBM) and industry standard interconnects (e.g., CXL) are key components of SDM
- Kernel and application tiering provide basic abstraction to underlying memory and storage resources
- 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.



**Please take a moment to rate this session.**

Your feedback is important to us.