

# CXL and why networking people should care

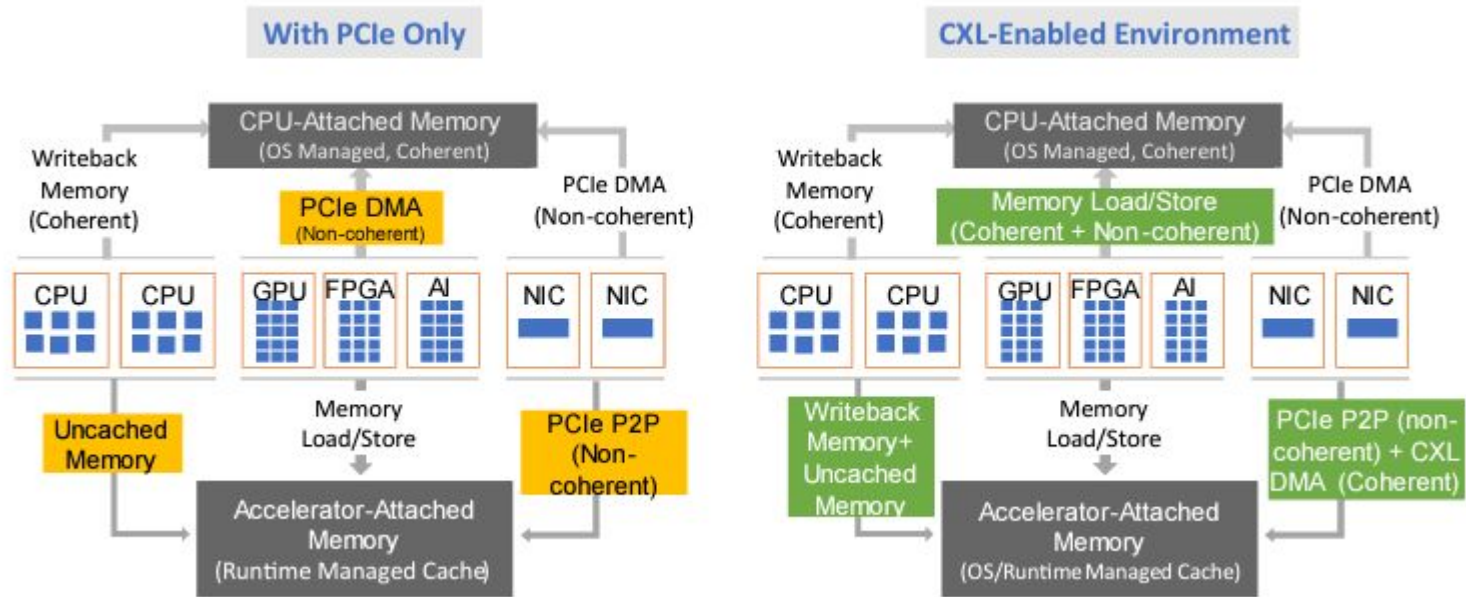


Shrijeet Mukherjee, CTO Enfabrica.

[shrijeet@enfabrica.net](mailto:shrijeet@enfabrica.net)

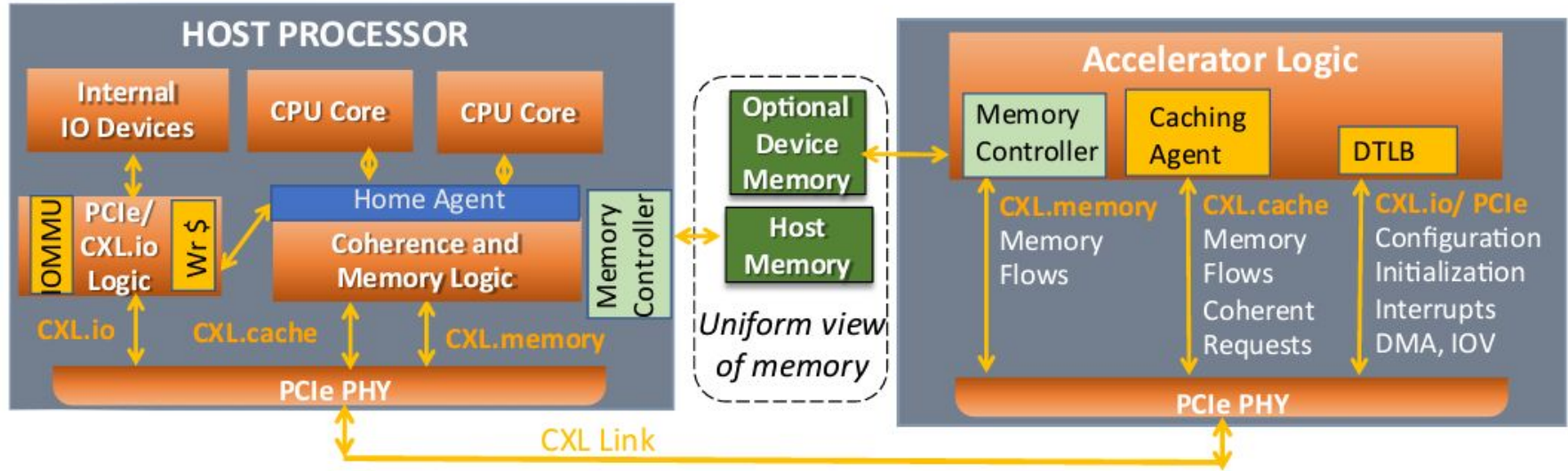
Netdev 0x17, 2023  
Vancouver

# :: CXL : A new beginning

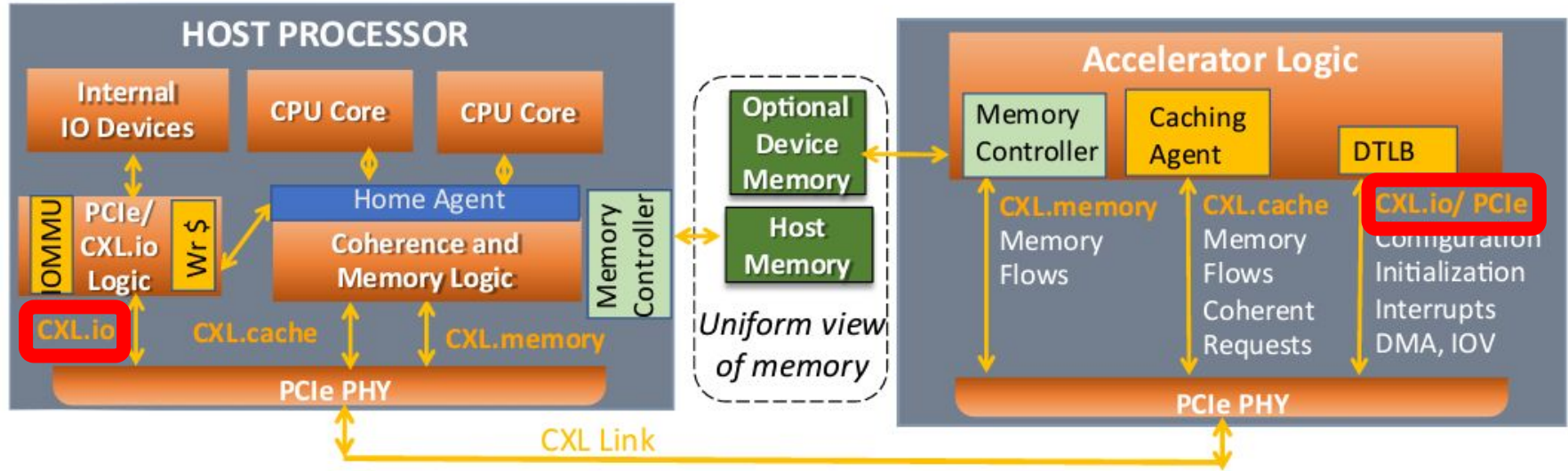


Ref : <https://arxiv.org/pdf/2306.11227.pdf>

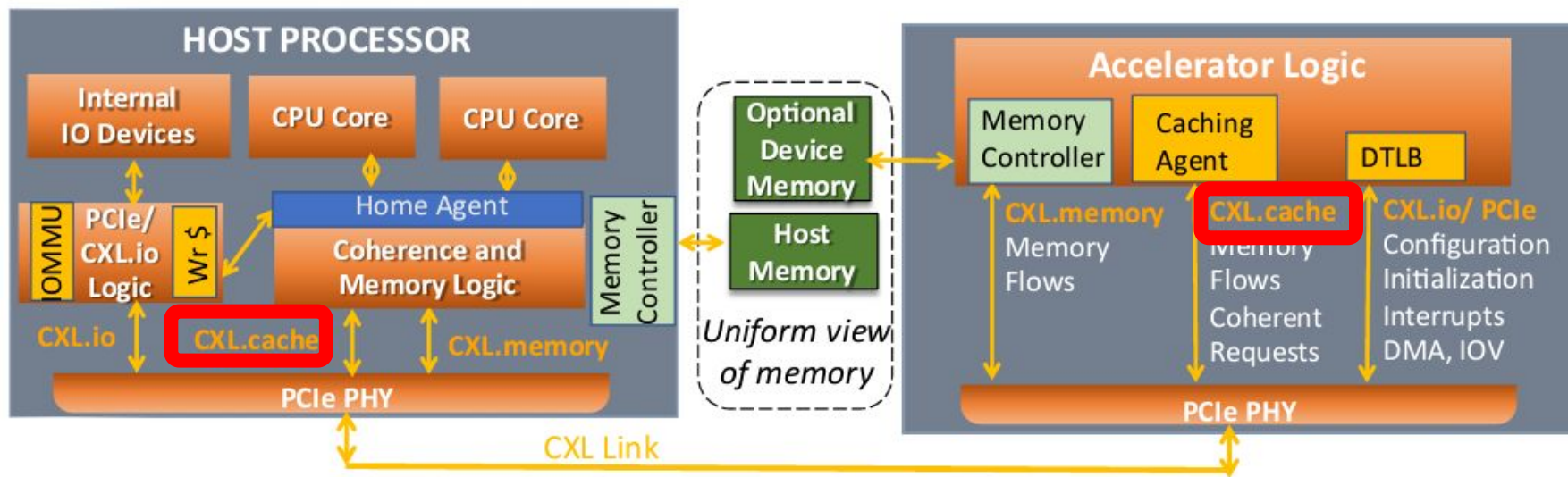
# :: the 3 horsemen of CXL



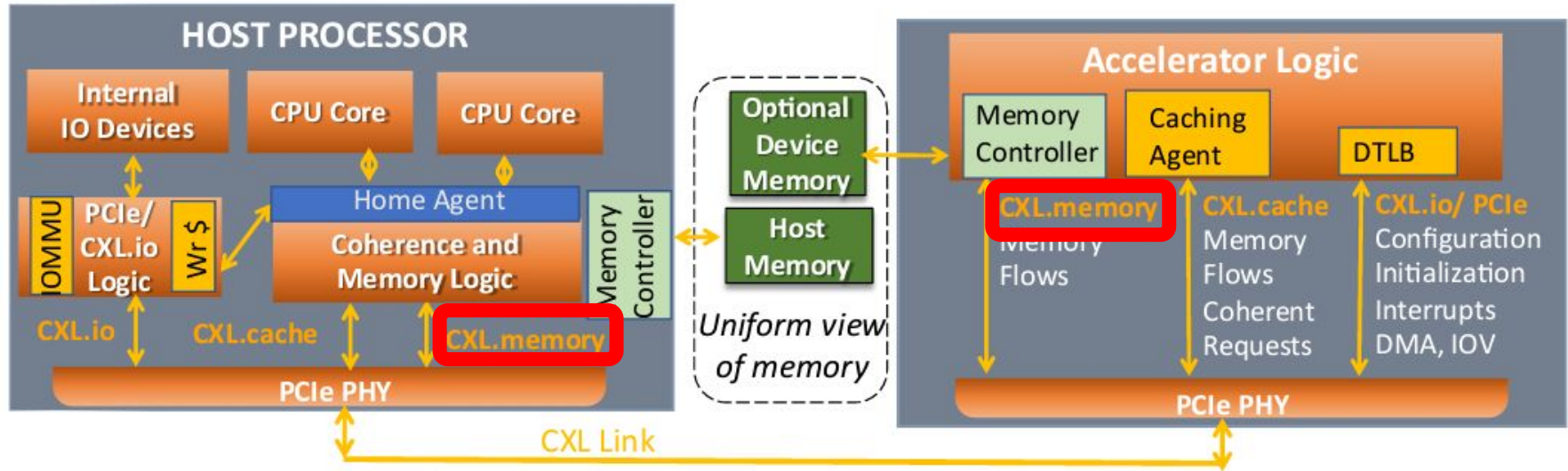
:: the 3 horsemen of CXL : io - PCIe with relaxed ordering @64B



:: the 3 horsemen of CXL : cache - distributed control plane/state



:: the 3 horsemen of CXL : memory - expansion of memory/buffers



# :: cxl target usage models : type 1

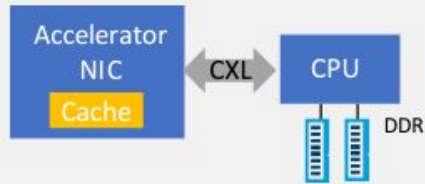
## Type 1: Accelerator w/o Memory

### Usages:

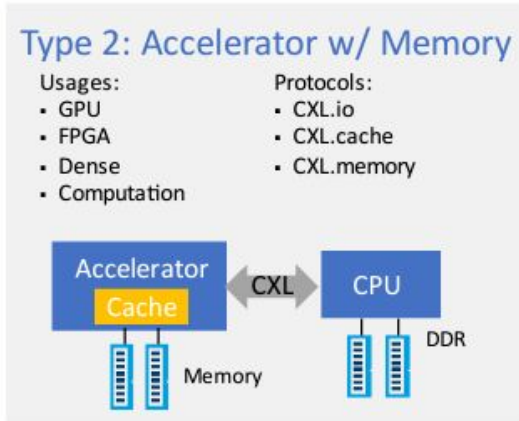
- PGAS NIC
- NIC atomics

### Protocols:

- CXL.io
- CXL.cache

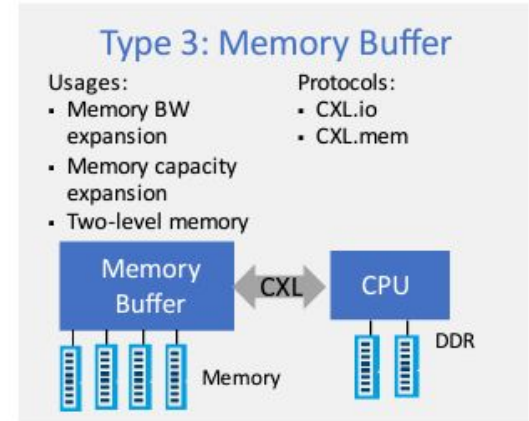


## :: cxl target usage models : type 2

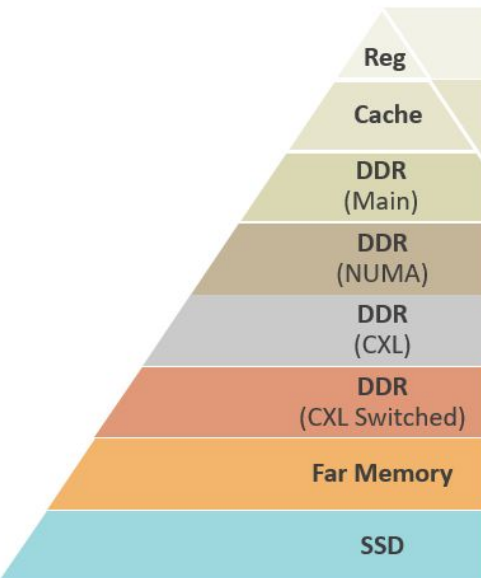




# :: cxl target usage models : type 3



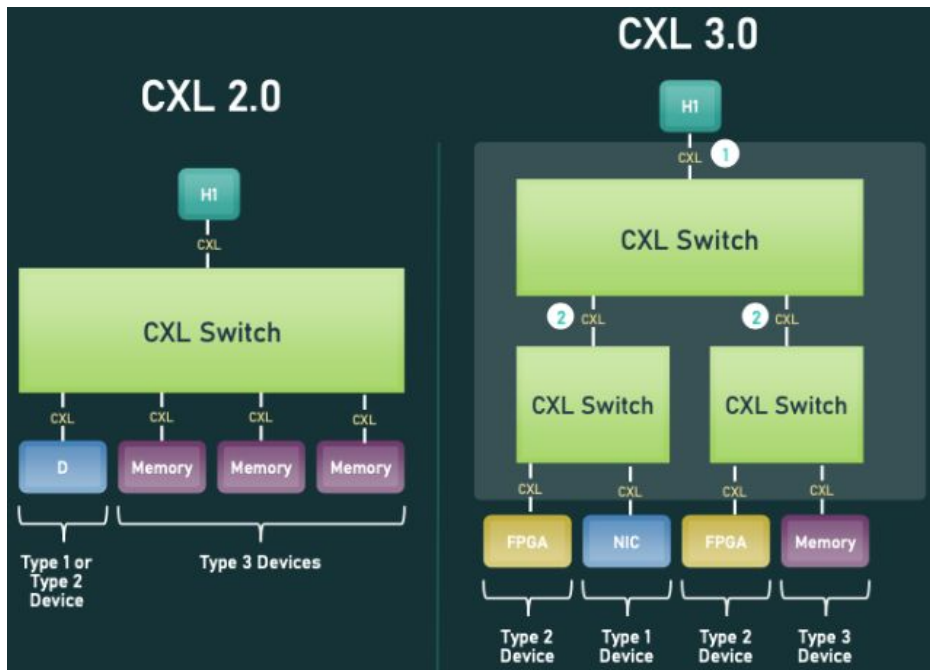
:: so where does it fit in the memory solar system



A pyramid diagram on the left side of the table, with each row of the table corresponding to a level of the pyramid. The pyramid is divided into sections: Reg (top, light green), Cache (light green), DDR (Main) (olive green), DDR (NUMA) (brown), DDR (CXL) (grey), DDR (CXL Switched) (orange), Far Memory (orange), and SSD (bottom, light blue).

	Latency	Bandwidth / Channel	Max Capacity*	Significance	Programmers View
Reg	0.2ns		KB	In CPU	L1 - dereference pointer
Cache	40ns		KB		L2 - dereference pointer
DDR (Main)	80-140ns	32-51.2 GB/s (DDR5)	Up to 4TB		L2 - dereference pointer high perf <u>memcpy</u>
DDR (NUMA)	170-250ns	32-51.2 GB/s (DDR5)	Up to 8TB	CPU independent but local	L3 - dereference pointer high perf <u>memcpy</u> , swap
DDR (CXL)	170-250ns	32-51.2 GB/s (DDR5)	2-4 TB		L4 - <u>memcpy</u> , swap
DDR (CXL Switched)	300-400ns	32-51.2 GB/s (DDR5)	64TB	Network attached	L5 - <u>memcpy</u> , swap
Far Memory	2-4us	100 GB/s (800g ethernet)	infinite		
SSD	50-100us				

# :: cxl - upwards, onwards and every direction possible



## 2.0 :

- Resource migration, memory scaling, and device fanout.
- Support for memory pooling. Adds SLD and MLD support
- Link-level Integrity and Data Encryption (CXL IDE)

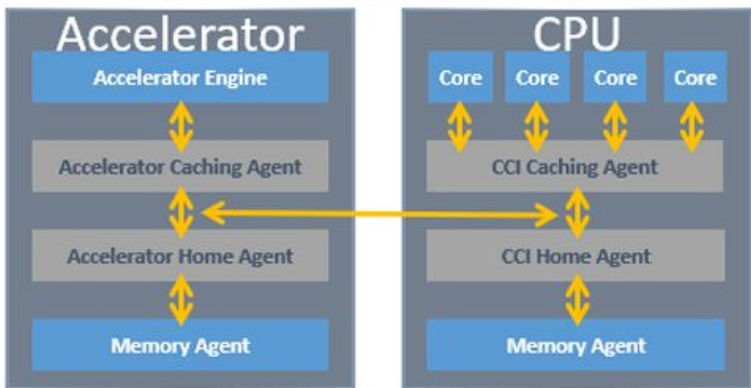
## 3.0 :

- Fabric capabilities
  - Multi-head and Fabric Attached Devices
  - Enhanced fabric management
  - Composable disaggregated infrastructure
  - Enhanced memory pooling
  - Multi-level switching
  - New enhanced coherency capabilities

Ref : [https://community.cadence.com/cadence\\_blogs\\_8/b/breakfast-bytes/posts/hot-chips-cxl-tutorial](https://community.cadence.com/cadence_blogs_8/b/breakfast-bytes/posts/hot-chips-cxl-tutorial)

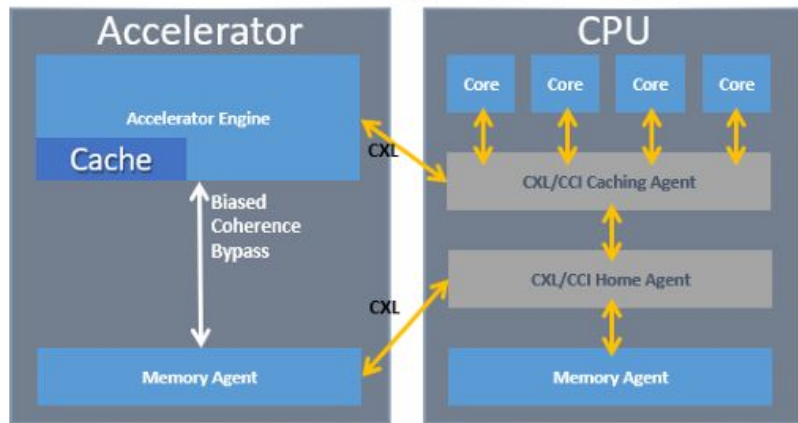
:: cxl is asymmetric and that is a good thing

### CCI\* Model – Symmetric CCI Protocol



\*Cache Coherent Interface

### CXL Model – Asymmetric Protocol

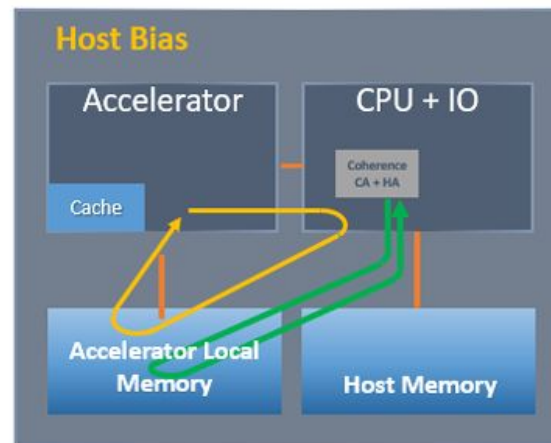
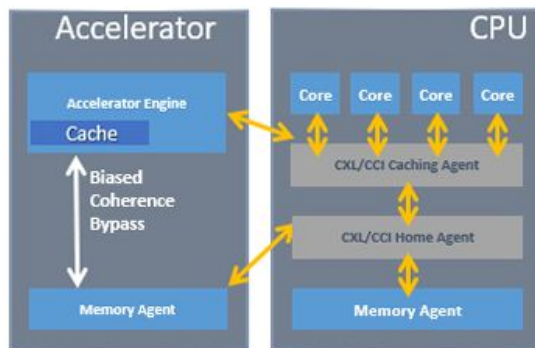
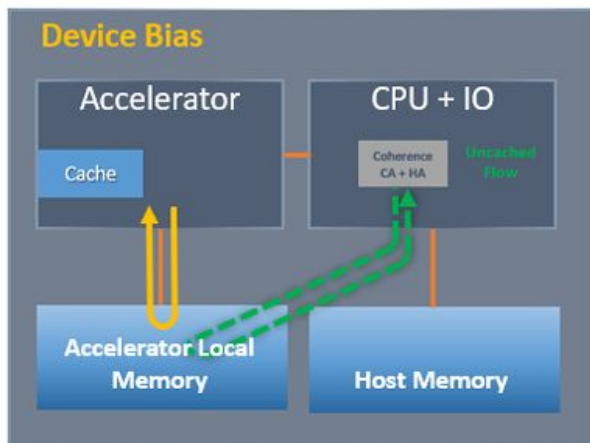


#### **CXL Key Advantages:**

- Avoid protocol interoperability hurdles/roadblocks
- Enable devices across multiple segments (e.g. client / server)
- Enable Memory buffer with no coherency burden
- **Simpler, processor independent device development**

Ref : [https://www.dmtf.org/sites/default/files/CXL\\_Overview\\_Virtual\\_DMTF\\_APTS\\_July\\_2020.pdf](https://www.dmtf.org/sites/default/files/CXL_Overview_Virtual_DMTF_APTS_July_2020.pdf)

# :: asymmetry needs bias



Critical access class for accelerators is "device engine to device memory"

"Coherence Bias" allows a device engine to access its memory coherently without visiting the processor

Two driver managed modes or "Biases"

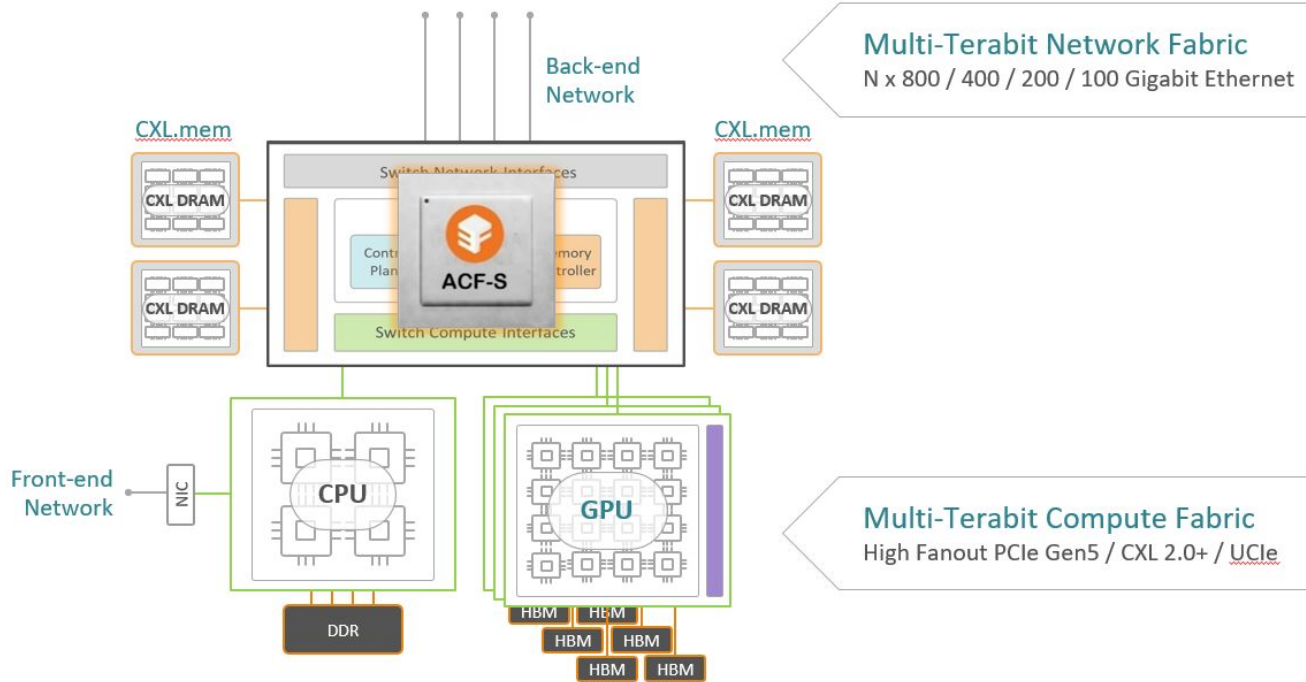
**HOST BIAS:** pages being used by the host or shared between host and device  
**DEVICE BIAS:** pages being used exclusively by the device

Both biases guaranteed correct/coherent

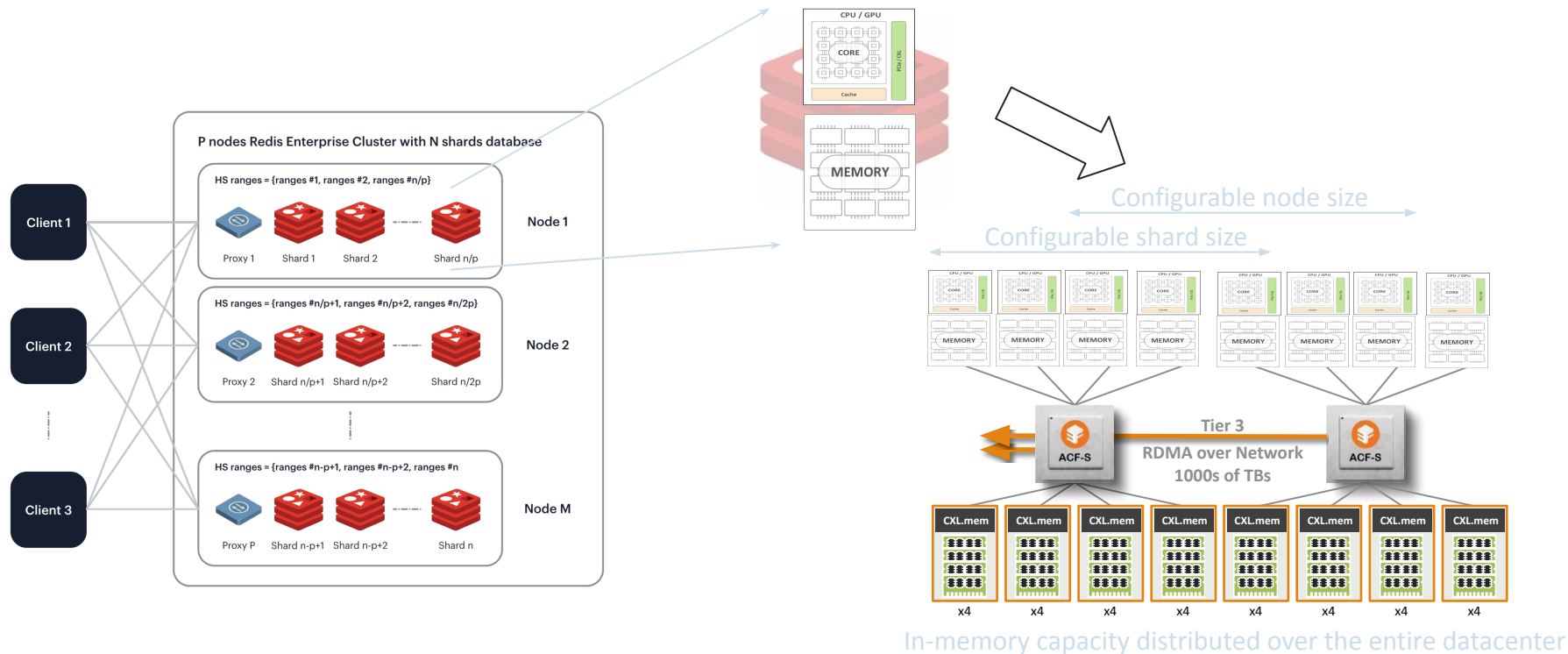
Guarantee applies even when software bugs or speculative accesses unexpectedly access device memory in the "Device Bias" state.

Ref : [https://www.dmtf.org/sites/default/files/CXL\\_Overview\\_Virtual\\_DMTF\\_APTS\\_July\\_2020.pdf](https://www.dmtf.org/sites/default/files/CXL_Overview_Virtual_DMTF_APTS_July_2020.pdf)

# :: cxl as imagined by enfabrica

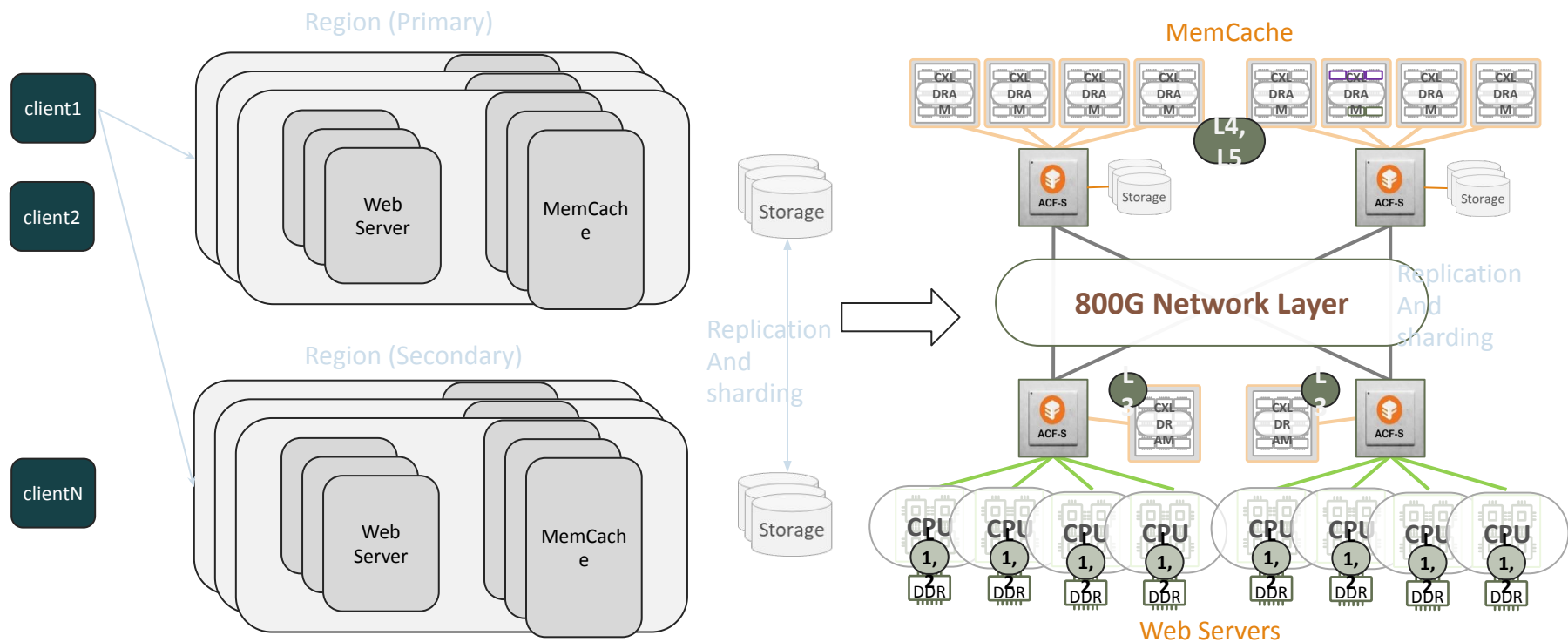


# :: in-memory databases



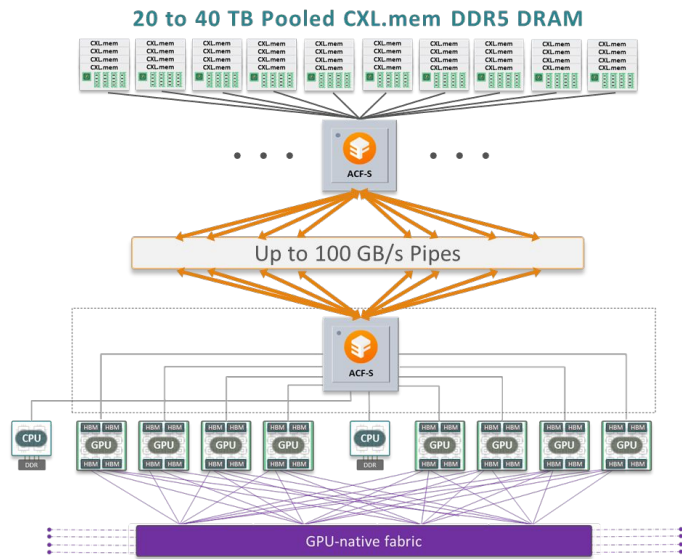
<https://redis.com/redis-enterprise/technology/linear-scaling-redis-enterprise/>

# :: memcached and horizontal scaling



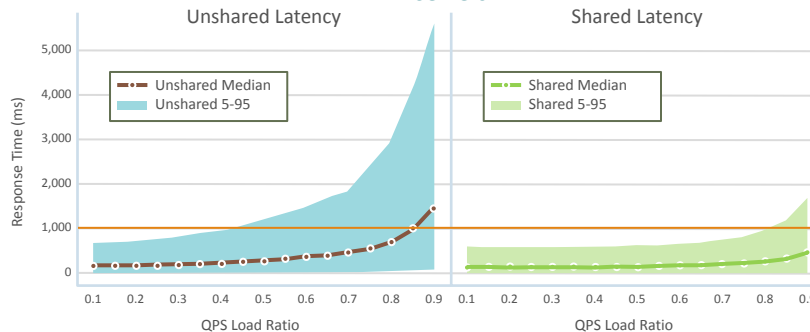


:: reducing inference fleet tco : 50% fewer gpus



Example LLM	Required QPS	# Required GPUs	# Required CPUs	User Context Capacity
User Context in CPU DRAM	1K	128	16	80K
User Context in ACF DRAM	1K	64	8	80K

Response Time Distribution, 200ms Avg Response Time with 4 DGX-style Servers



Scaling memory headlessly with ACF-S reduces total xPU spend for LLM inference