

Disaggregated & Heterogeneous Platform for Data Management

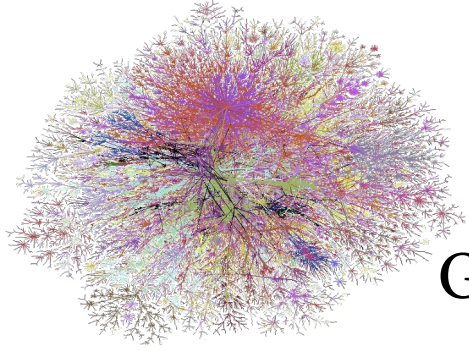
M. Tamer Özsu
Cheriton School of Computer Science
tamer.ozsu@uwaterloo.ca



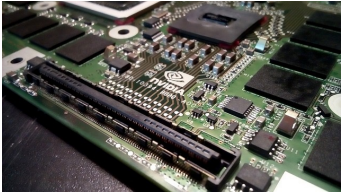
UNIVERSITY OF
WATERLOO



My research

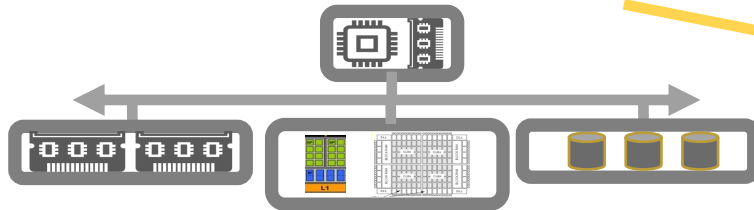


Graph Processing



Graph Processing
On GPUs^{1,2}

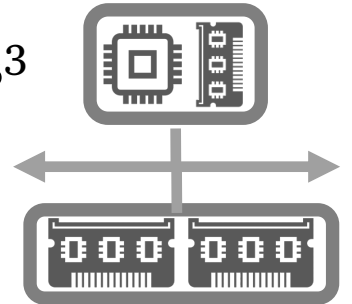
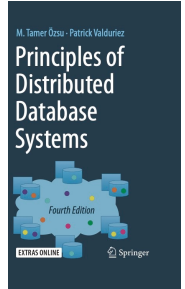
Graph Processing on
Heterogeneous Platform



Graph Processing on Disaggregated
Heterogeneous Platform

Distributed/Parallel
Data Management

Disaggregated Systems³



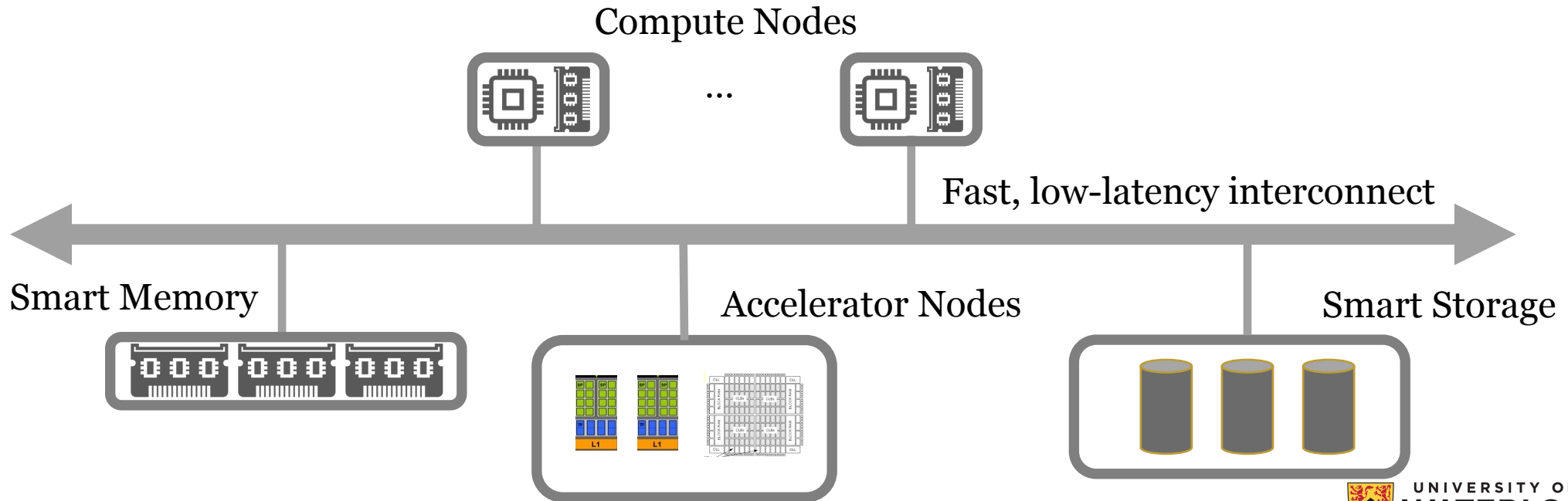
¹L. Hu, L. Zou and M. T. Özsu. GAMMA: A Graph Pattern Mining Framework for Large Graphs on GPU, *Proc. ICDE*, 2023.

²L. Zeng, et al. GSI:GPU-friendly Subgraph Isomorphism, *Proc. ICDE*, 2020.

³R.Wang et al. The Case for Shared-Memory Databases with RDMA-Enabled Memory Disaggregation, *Proc. VLDB*, 2022.

Disaggregated Heterogeneous Platform

- Disaggregated
 - Separate components with a fast interconnect
- Heterogeneous
 - “Executing programs on a computing platform with computing nodes of different characteristics.”¹

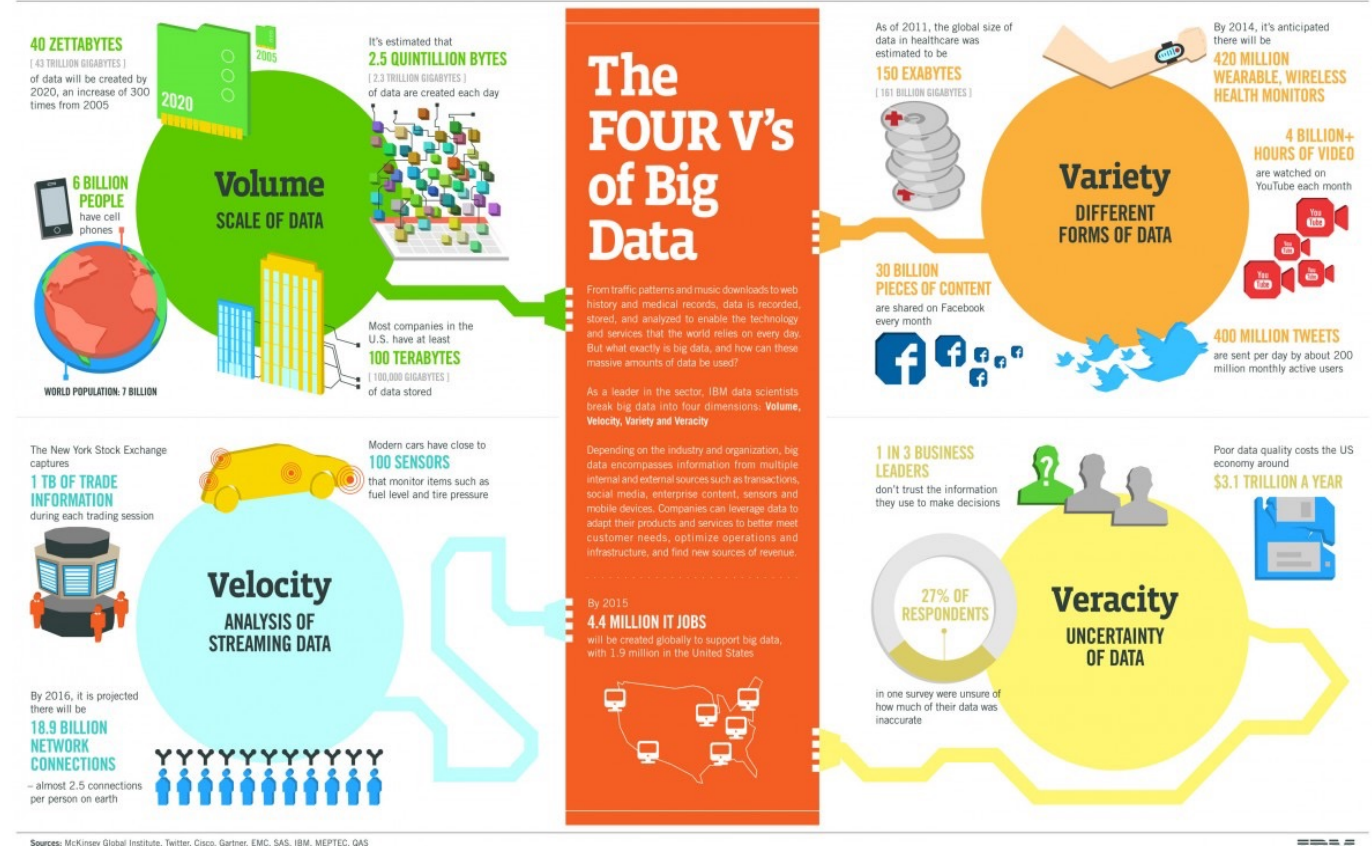


Why?

New Application Demands
Technological Changes

New Application Demands

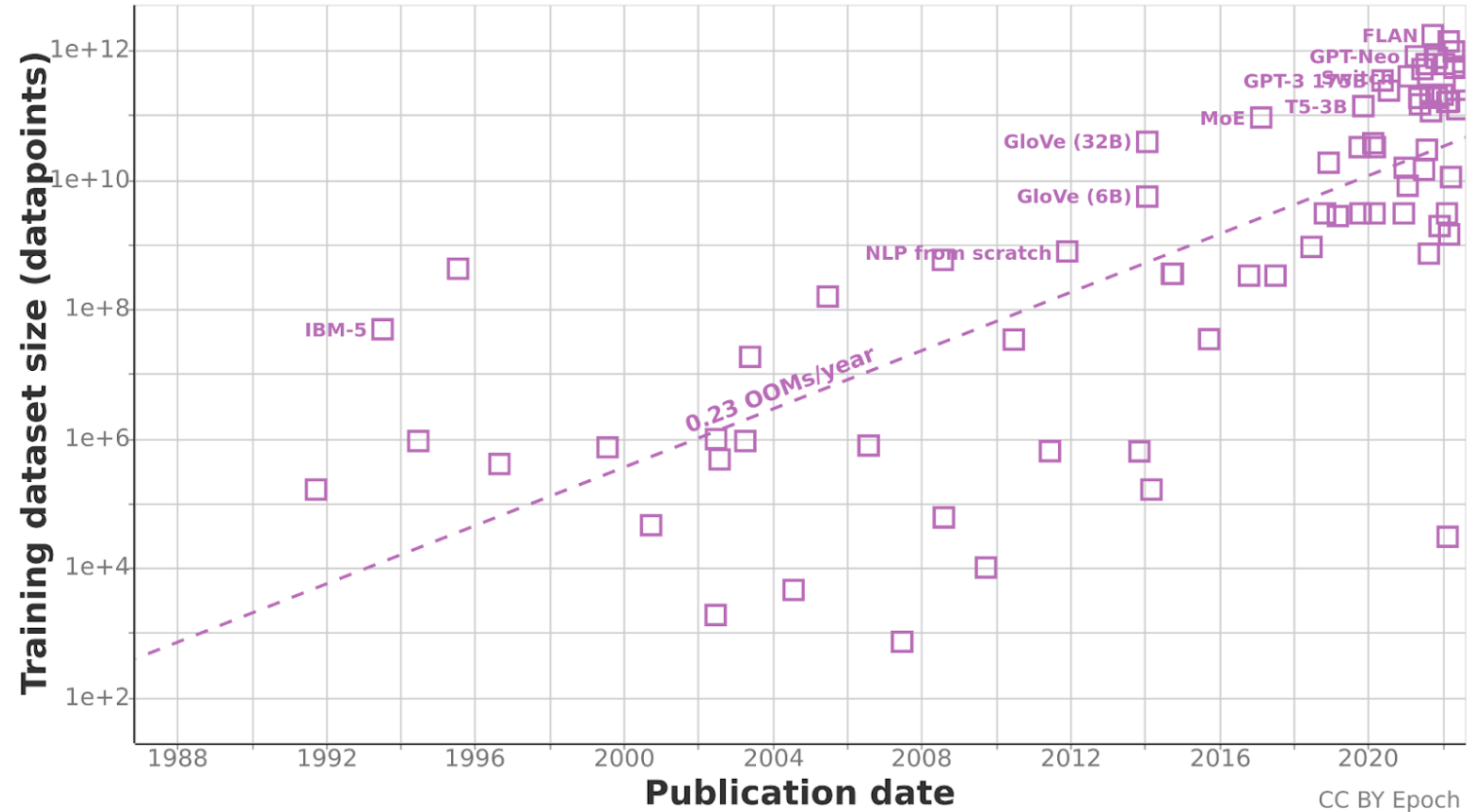
Big data applications



New Application Demands

Big data applications

- Very large storage
 - 10x rule in ML workloads



<https://epochai.org/blog/trends-in-training-dataset-sizes>

New Application Demands

Big data applications

- Very large storage
 - 10x rule in ML workloads
- Repetitive computation

```
class Vertex:
    def __init__(self, name):
        self.name = name
        self.children = []
        self.parents = []
        self.auth = 1.0
        self.hub = 1.0
        self.pagerank = 1.0

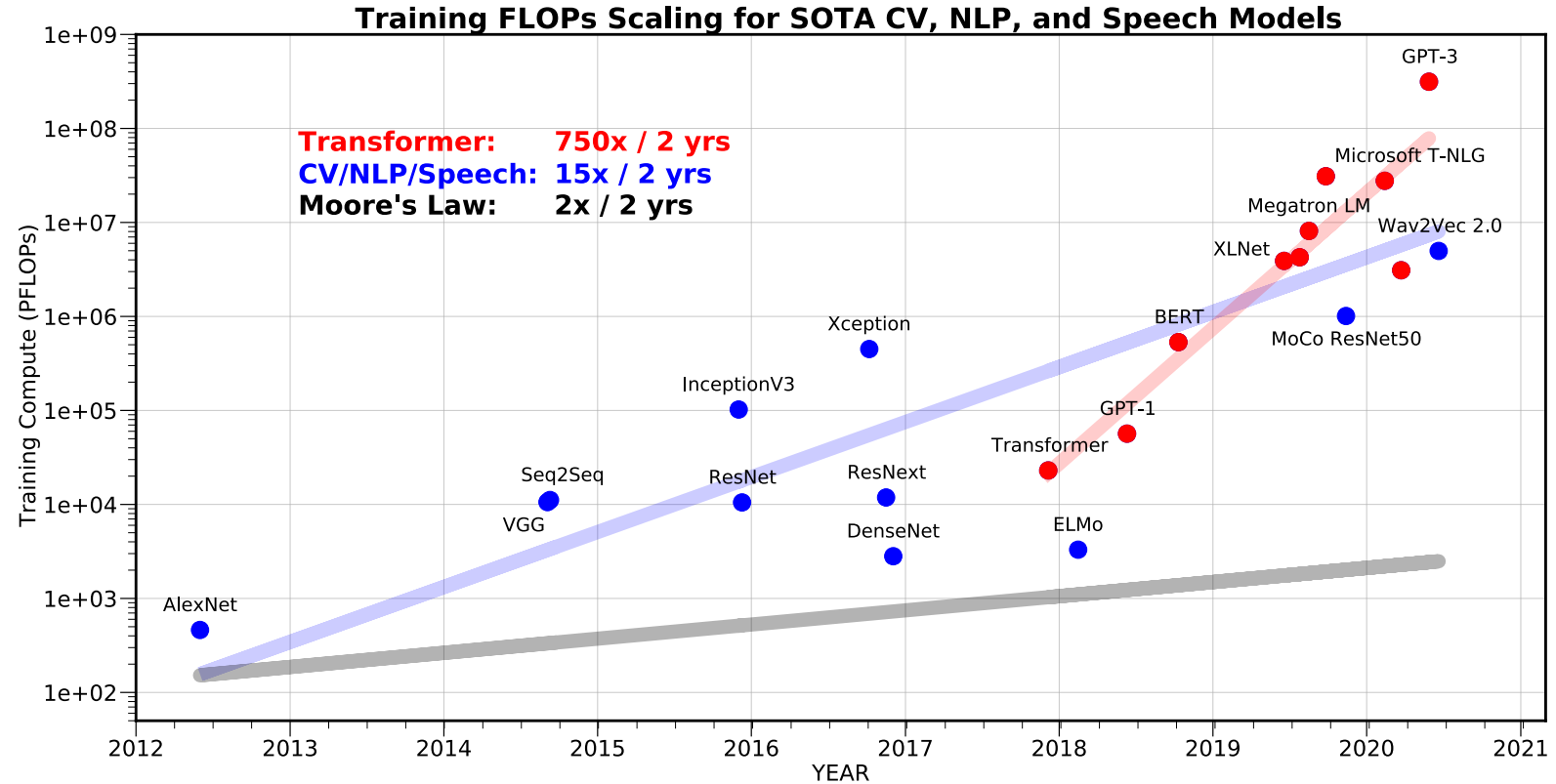
#One iteration computation
def PageRank_one_iter(graph, d):
    vertex_list = graph.vertices
    for vertex in vertex_list:
        vertex.update_pagerank(d, len(graph.vertices))

#Go over all vertices until convergence
def PageRank(graph, d):
    finished = False
    while not finished:
        PageRank_one_iter(graph, d)
        finished = converge()
```

New Application Demands

Big data applications

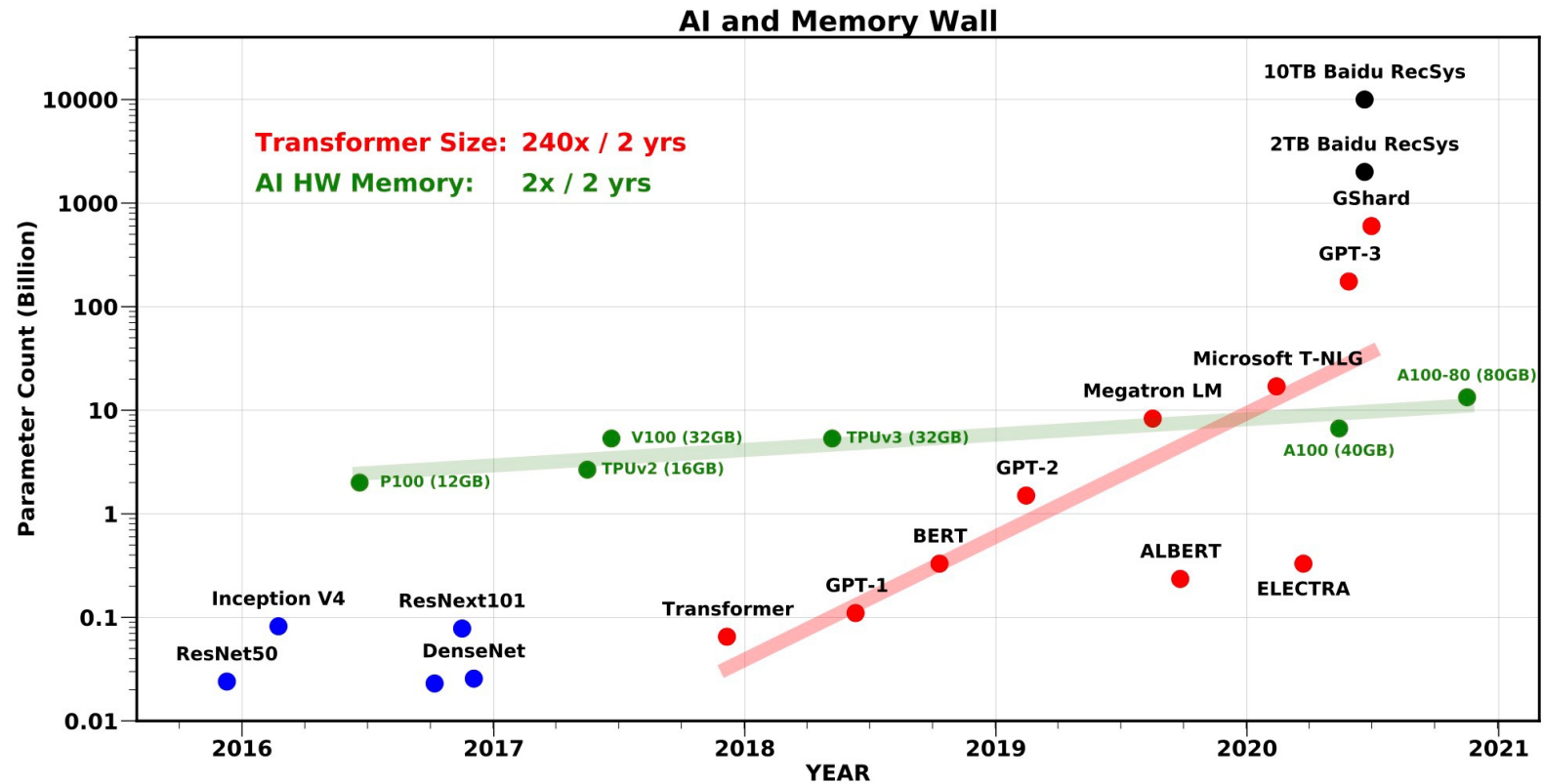
- Very large storage
 - 10x rule in ML workloads
- Repetitive computation
 - Excessive computing cycles



New Application Demands

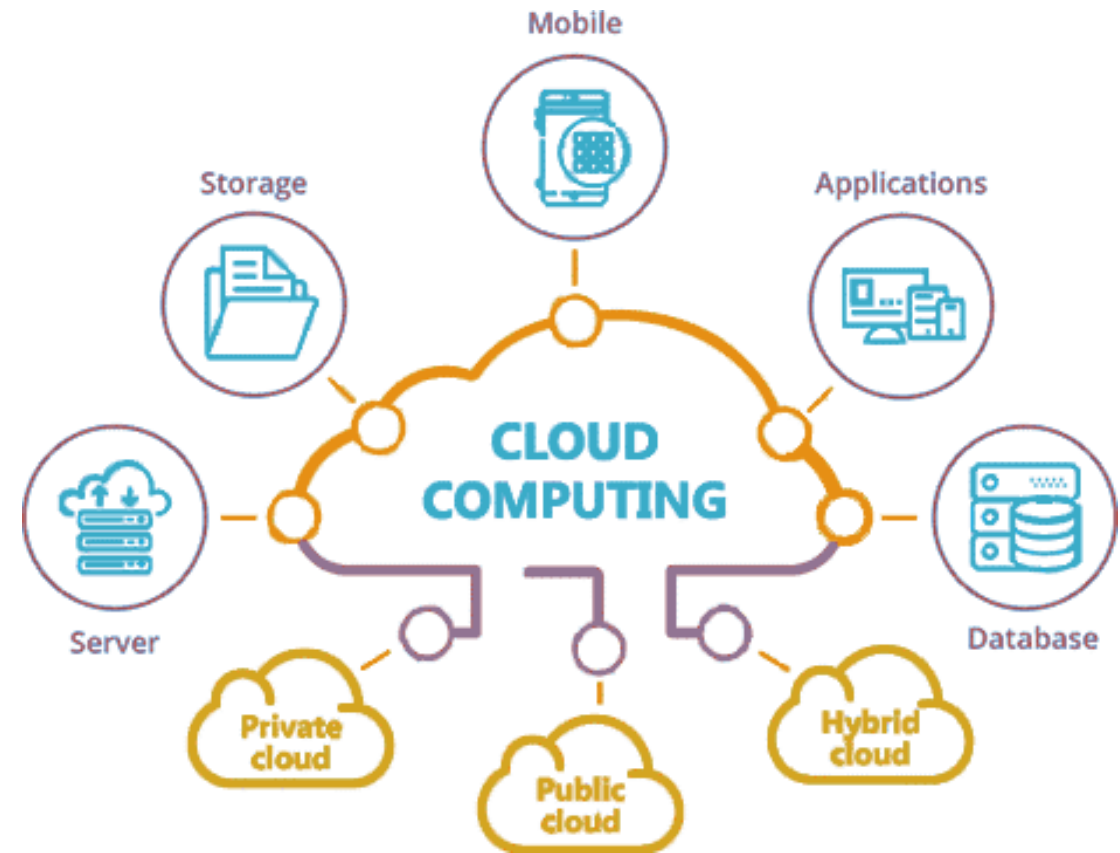
Big data applications

- Very large storage
 - 10x rule in ML workloads
- Repetitive computation
 - Excessive computing cycles
- Large memory



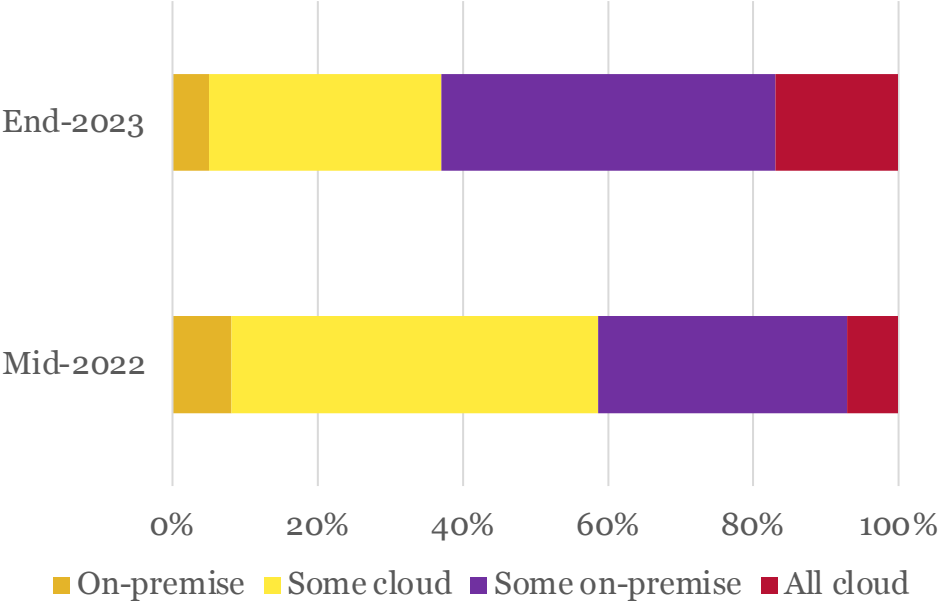
TECHNOLOGY CHANGES

Move to the Cloud

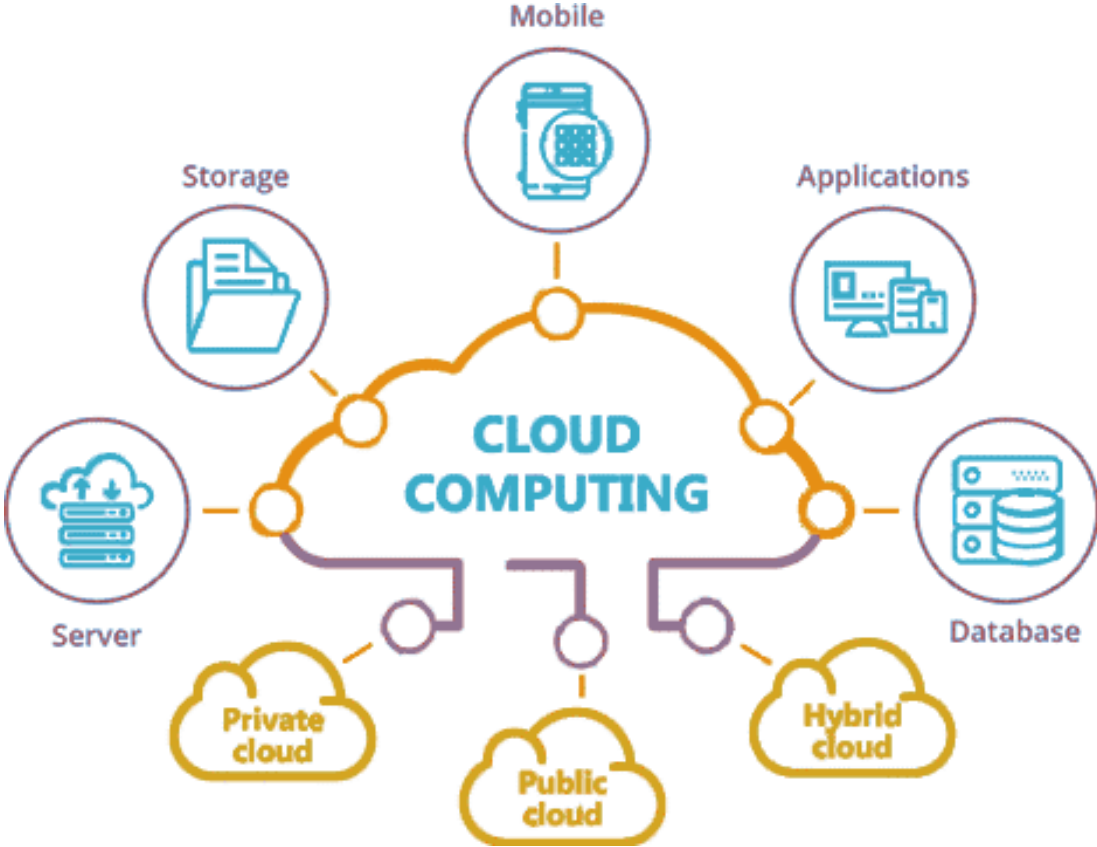


Move to the Cloud

- Move is expanding

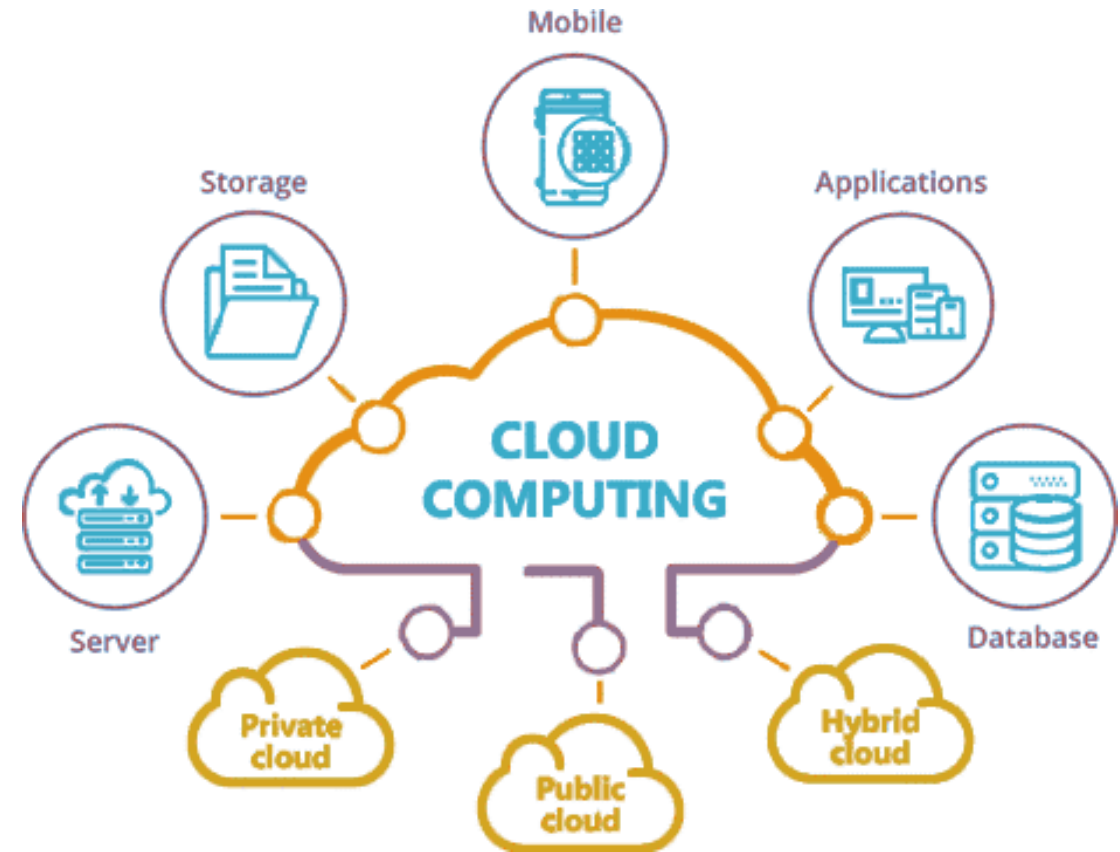


Foundry Cloud Computing Study, 2022



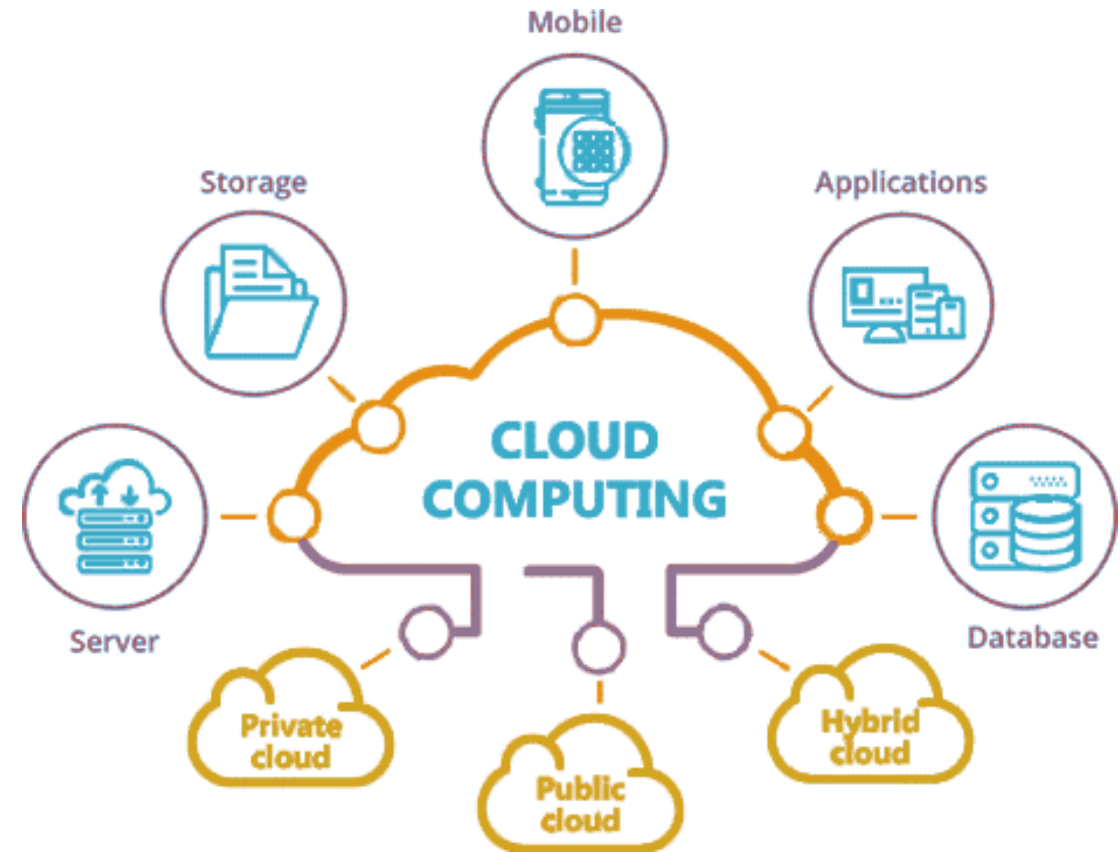
Move to the Cloud

- Move is expanding
- Reasons
 - Elasticity
 - Availability
 - Cost savings
 - ...



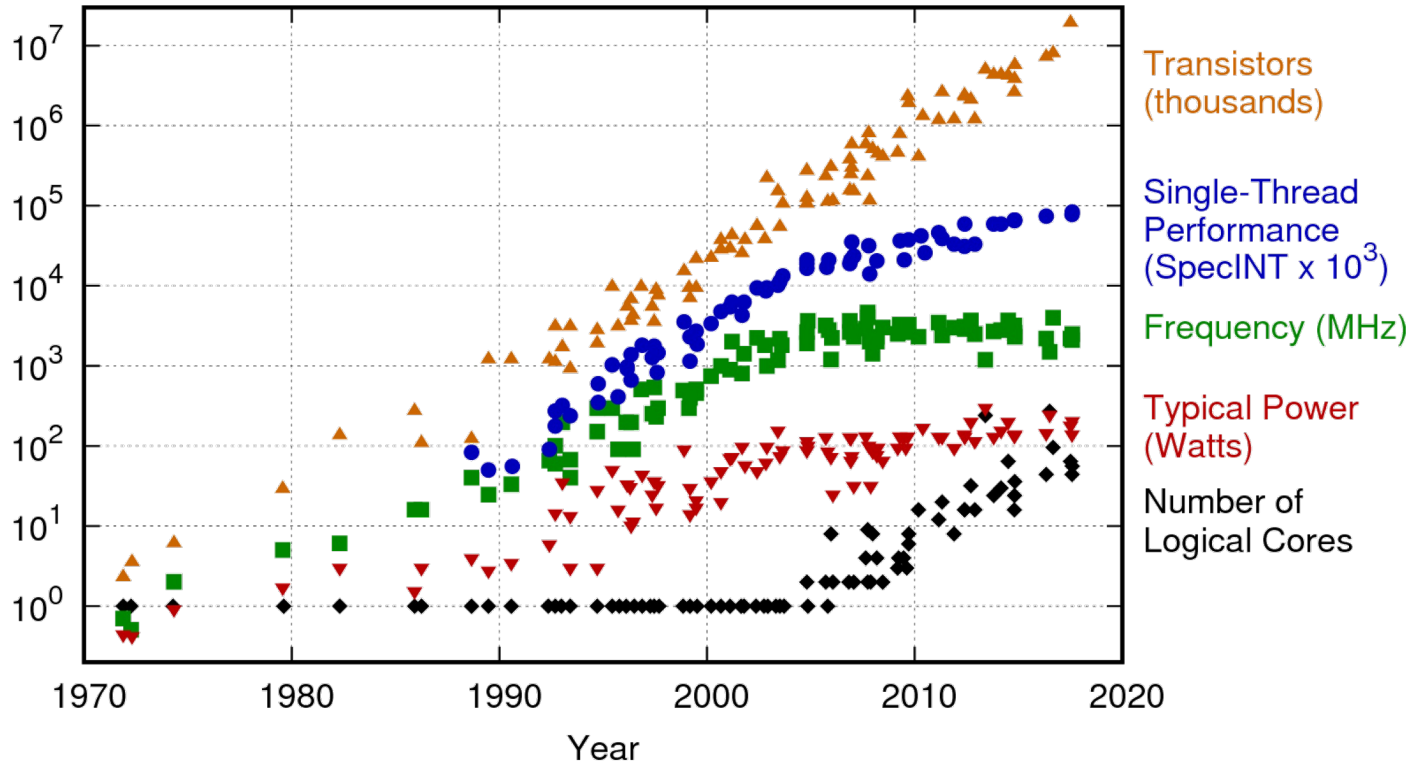
Move to the Cloud

- Move is expanding
- Reasons
 - Elasticity
 - Availability
 - Cost savings
 - ...
- Stress points
 - Elasticity in demand
 - Configuration to meet SLAs
 - ...



Is Traditional Processor Keeping Up?

42 Years of Microprocessor Trend Data



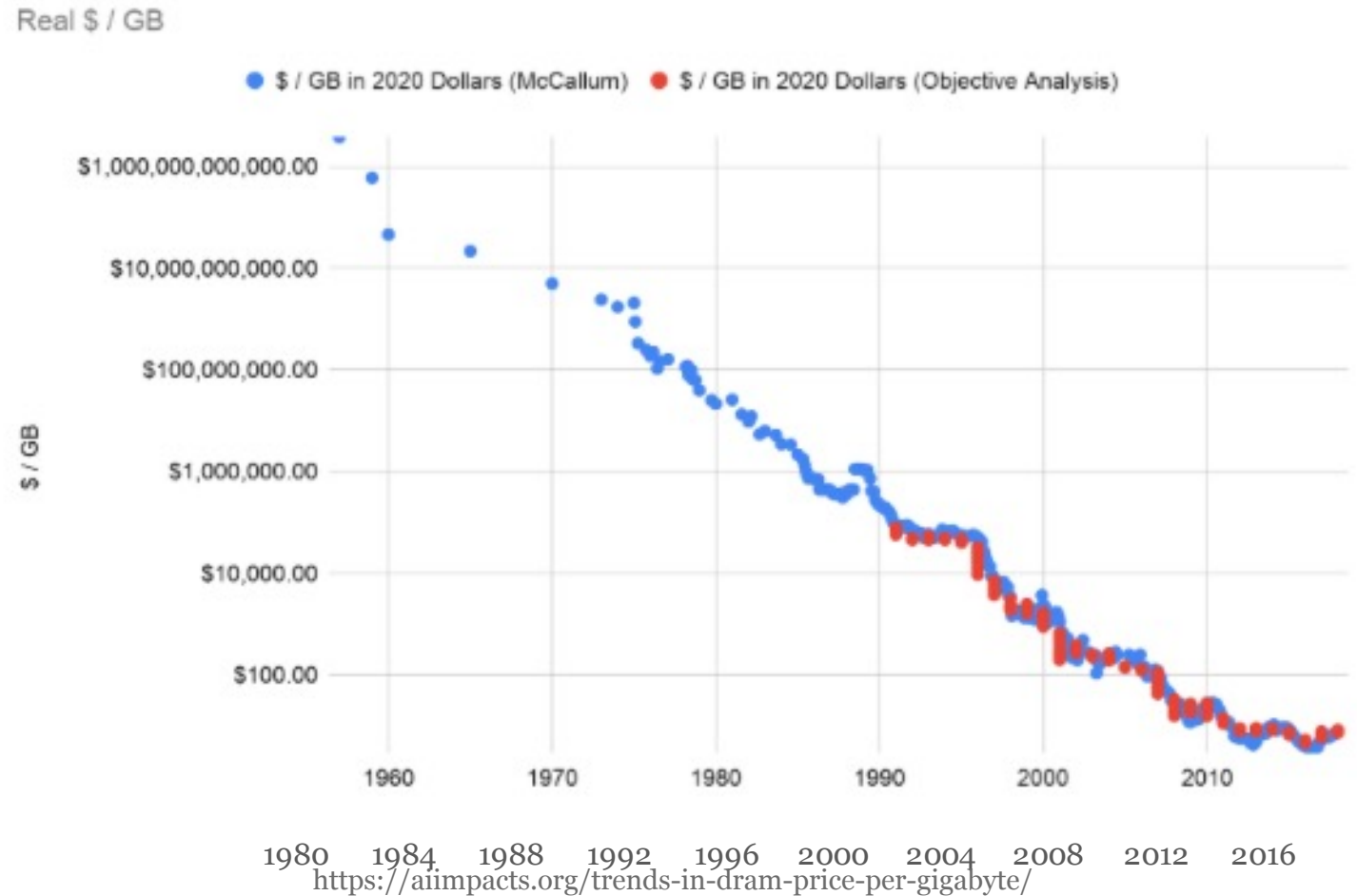
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

<https://www.karlrupp.net/2015/06/40-years-of-microprocessor-trend-data/>

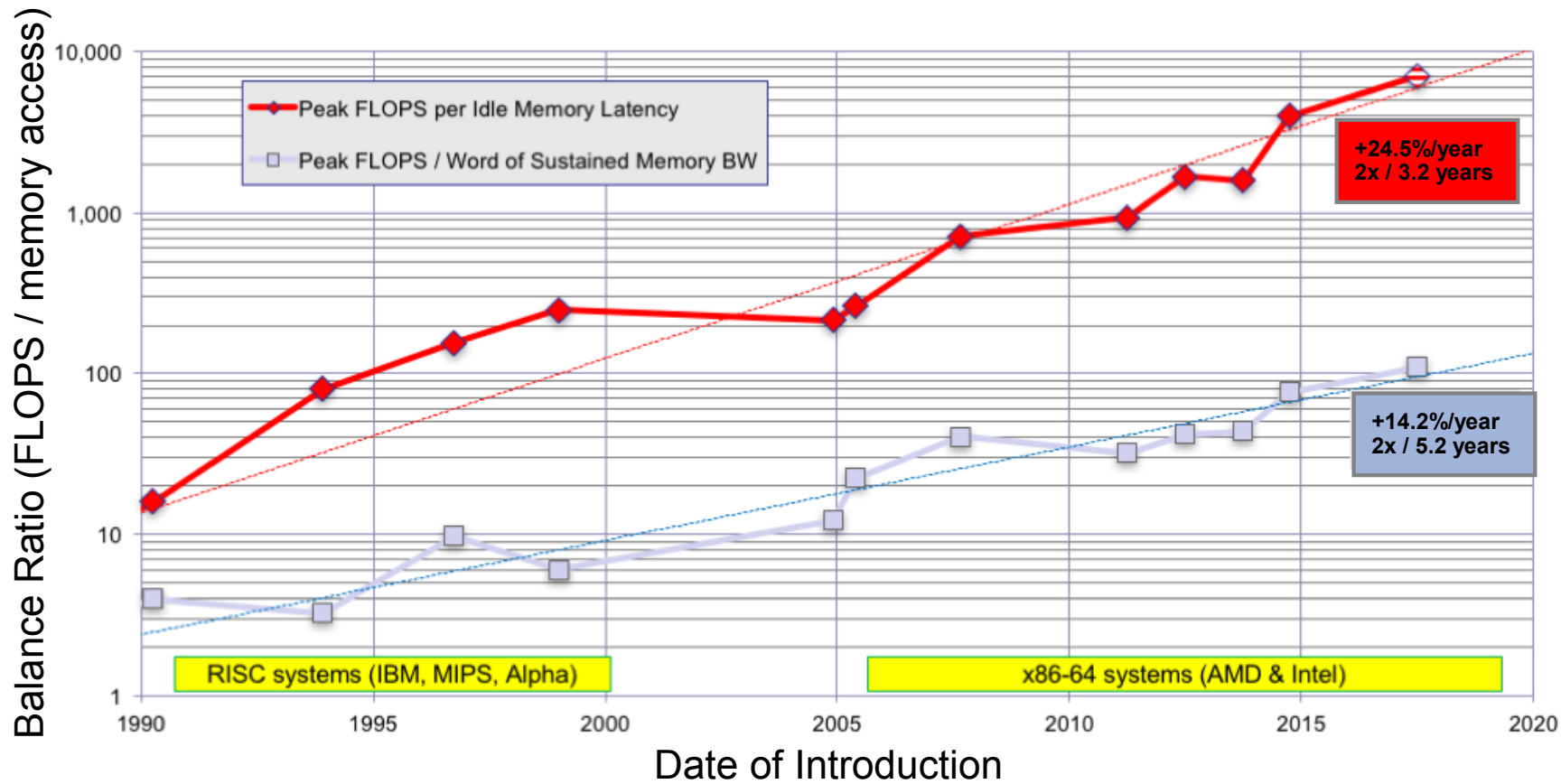
- Moore's Law application questionable
 - *The number of transistors on a microchip doubles every two years*
- Pollack's Rule
 - *Microprocessor performance increase due to microarchitecture advances is roughly proportional to the square root of the increase in complexity.*

Changes in RAM

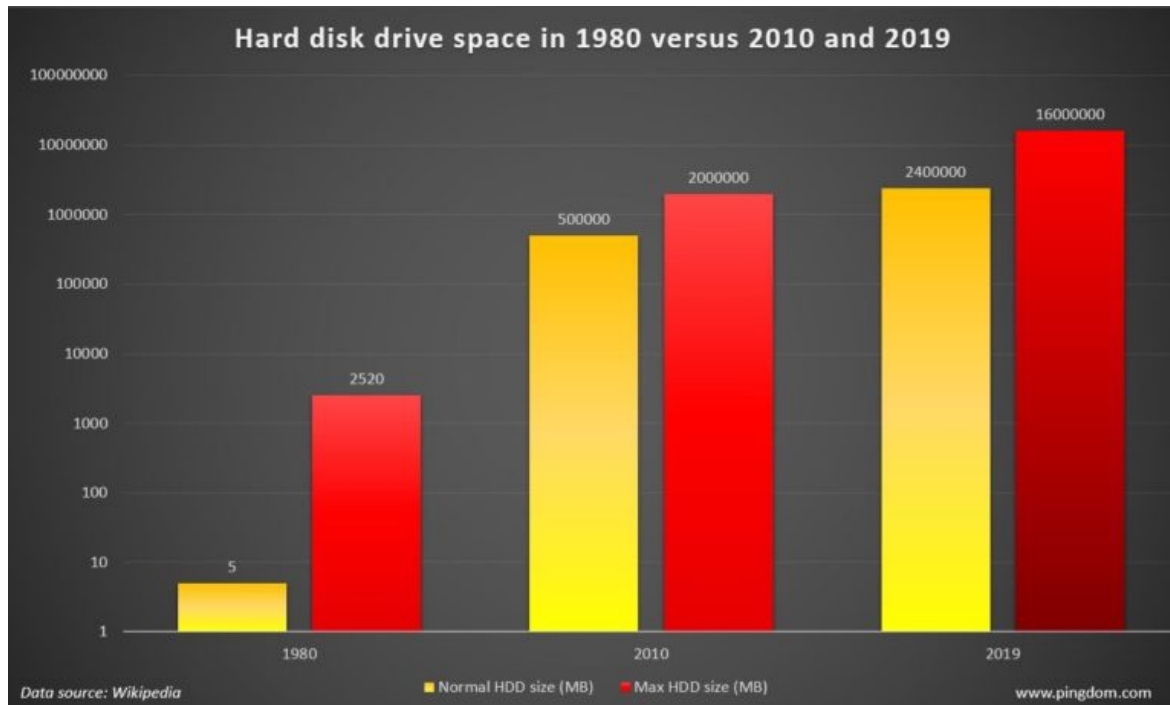
- RAM capacities increasing 10x every four years (?)
 - Not single memory chip
- RAM prices are going down
- Memory bandwidth increasing ~23% per year¹
- Memory latency **increasing** ~4% per year¹



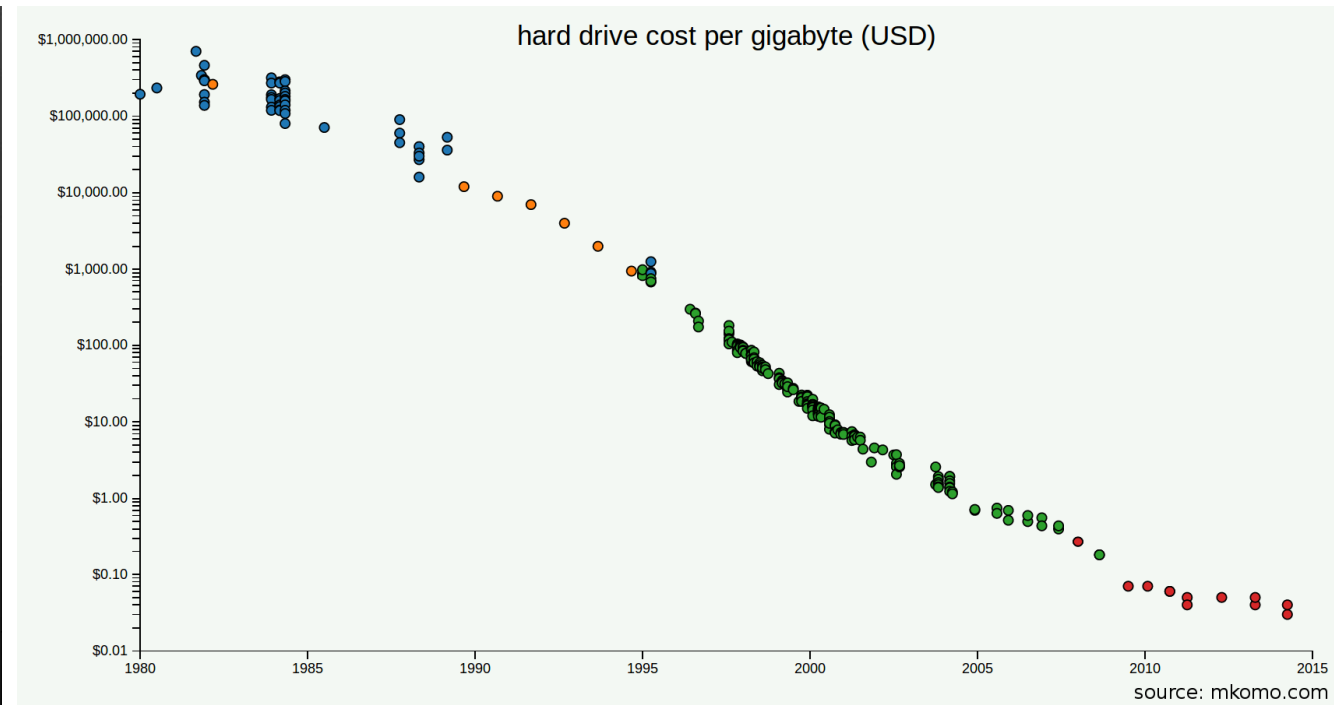
Is The Memory Keeping Up?



Storage Capacity and Price Over Time



<https://www.pingdom.com/>



<https://aiimpacts.com/>

Storage Devices

- 10TB disks are about US\$600
 - *Approximately 4,000 DVD movies, 1.8M digital photos or 2.5M mp3 music files*



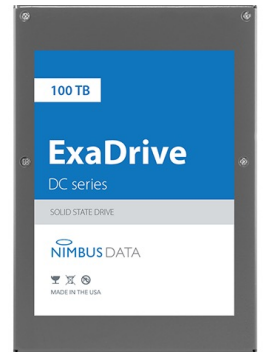
Storage Devices

- 10TB disks are about US\$600
 - *Approximately 4,000 DVD movies, 1.8M digital photos or 2.5M mp3 music files*
- 10TB single drives are now available



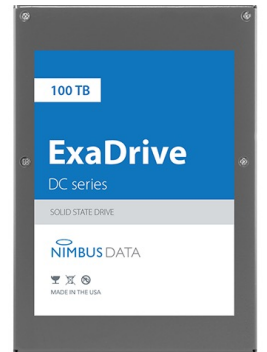
Storage Devices

- 10TB disks are about US\$600
 - *Approximately 4,000 DVD movies, 1.8M digital photos or 2.5M mp3 music files*
- 10TB single drives are now available
- Nimbus ExaDrive DC100: 100TB SSD



Storage Devices

- 10TB disks are about US\$600
 - *Approximately 4,000 DVD movies, 1.8M digital photos or 2.5M mp3 music files*
- 10TB single drives are now available
- Nimbus ExaDrive DC100: 100TB SSD
- Flash is now mainstream
 - 1TB flash is common
 - Different storage hierarchy

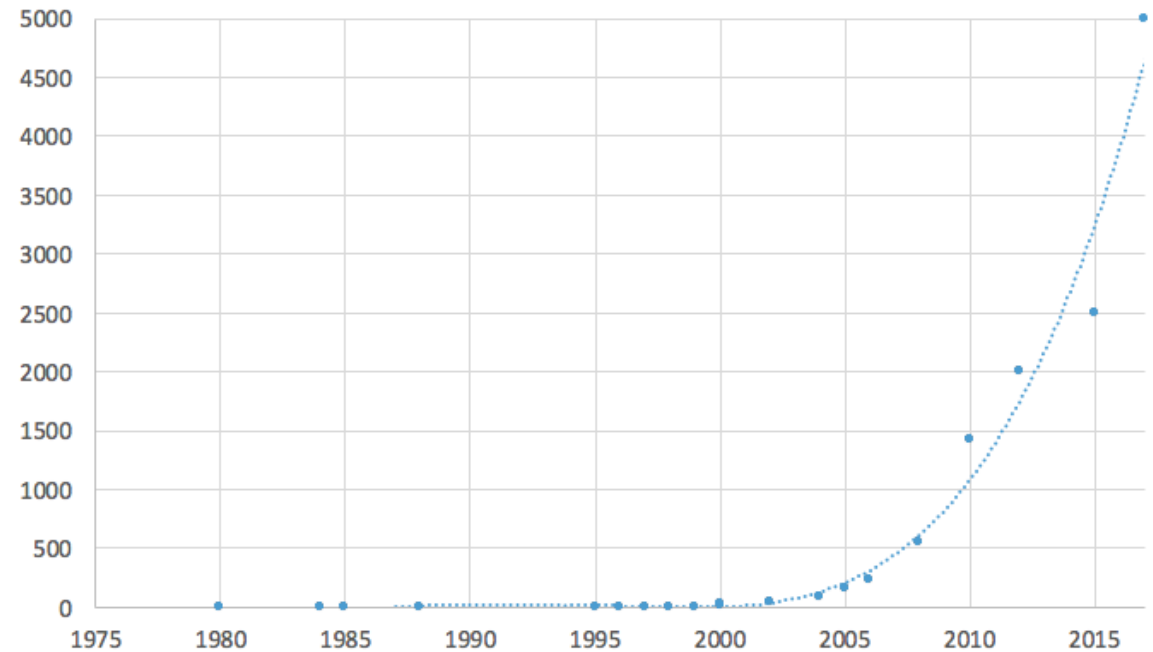


Storage Devices

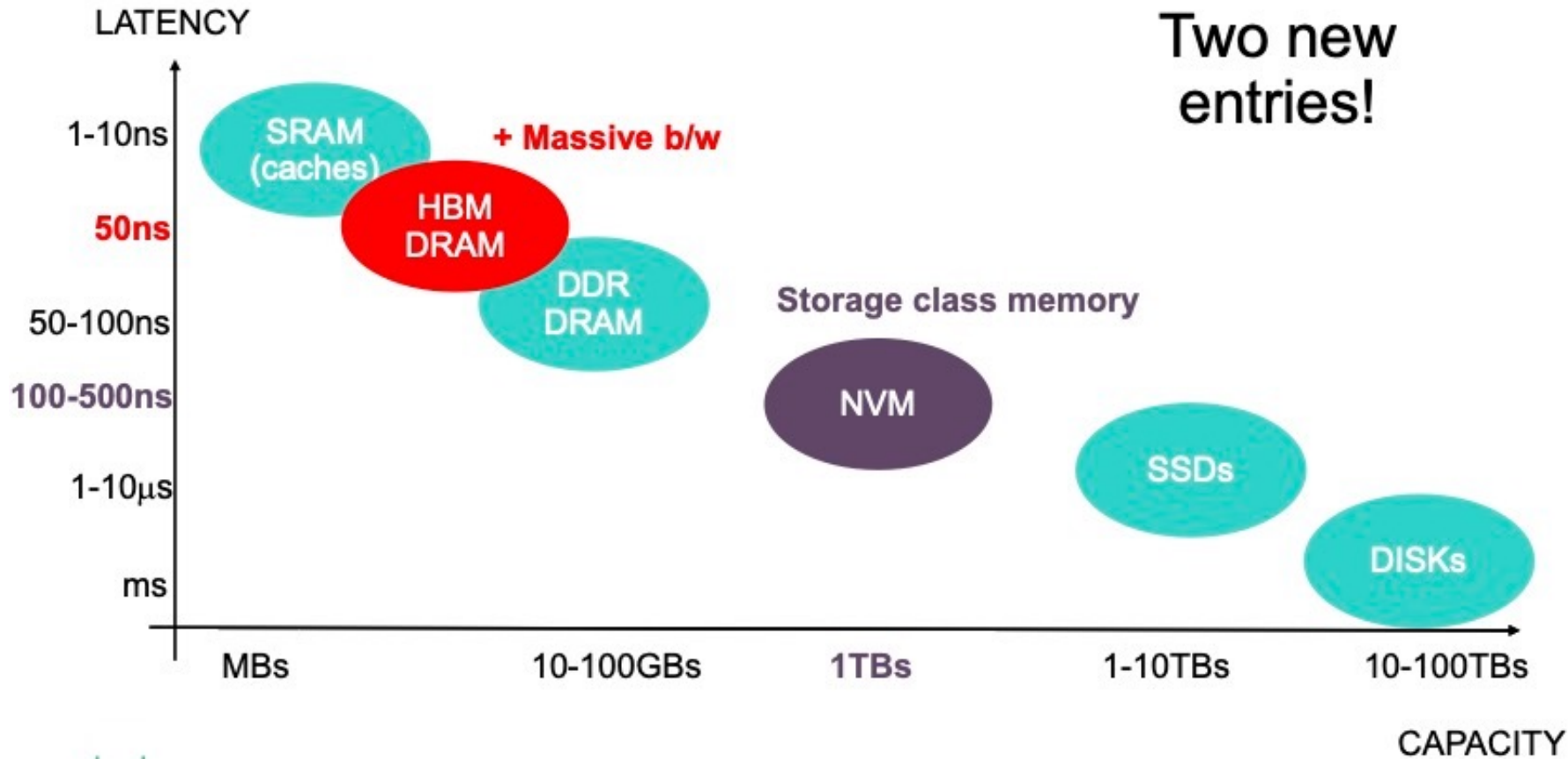
- 10TB disks are about US\$600
 - *Approximately 4,000 DVD movies, 1.8M digital photos or 2.5M mp3 music files*
- 10TB single drives are now available
- Nimbus ExaDrive DC100: 100TB SSD
- Flash is now mainstream
 - 1TB flash is common
 - Different storage hierarchy
- Money goes a long way



Gigabytes of storage you can buy with \$100

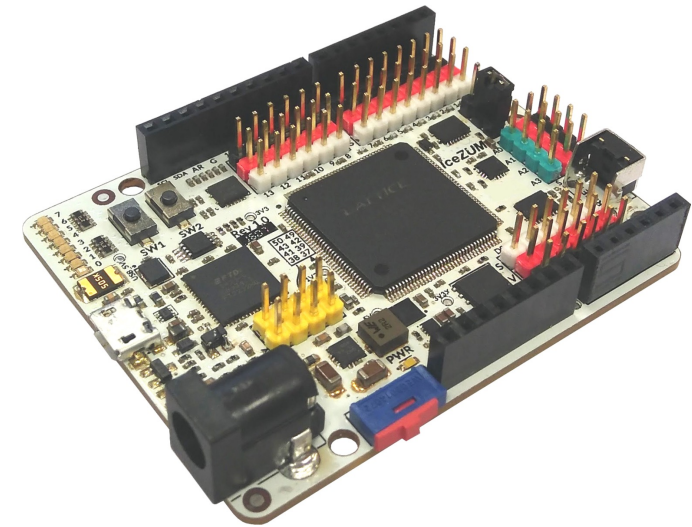
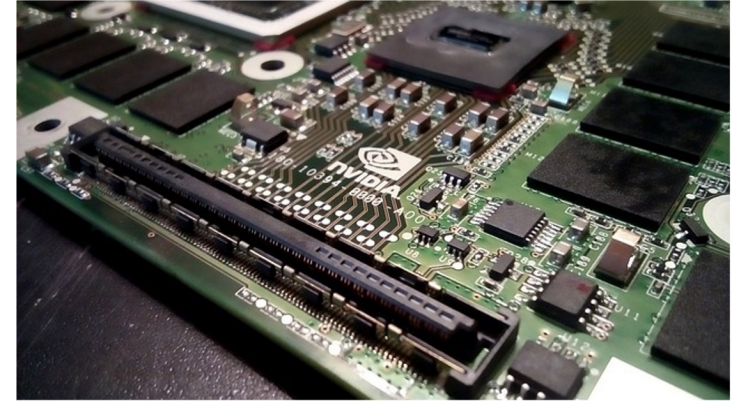


Changing Memory+Storage Hierarchy



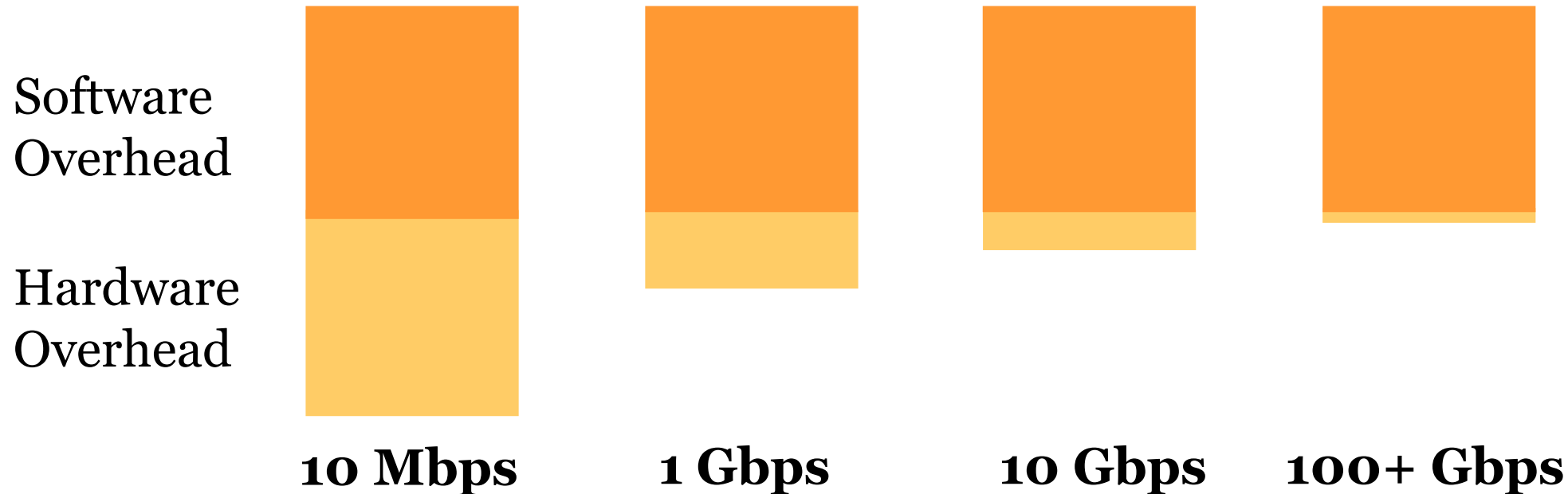
Accelerators

- GPU
 - Large thread parallelism
 - SIMD computation
 - Limited on-chip memory
 - Large global memory w/ increasing bandwidth
- FPGA
 - A set of programmable logic blocks
 - Logic blocks can be configured to perform complex functions
 - On-chip memory configurable
- ASIC
 - Once you know what you are doing



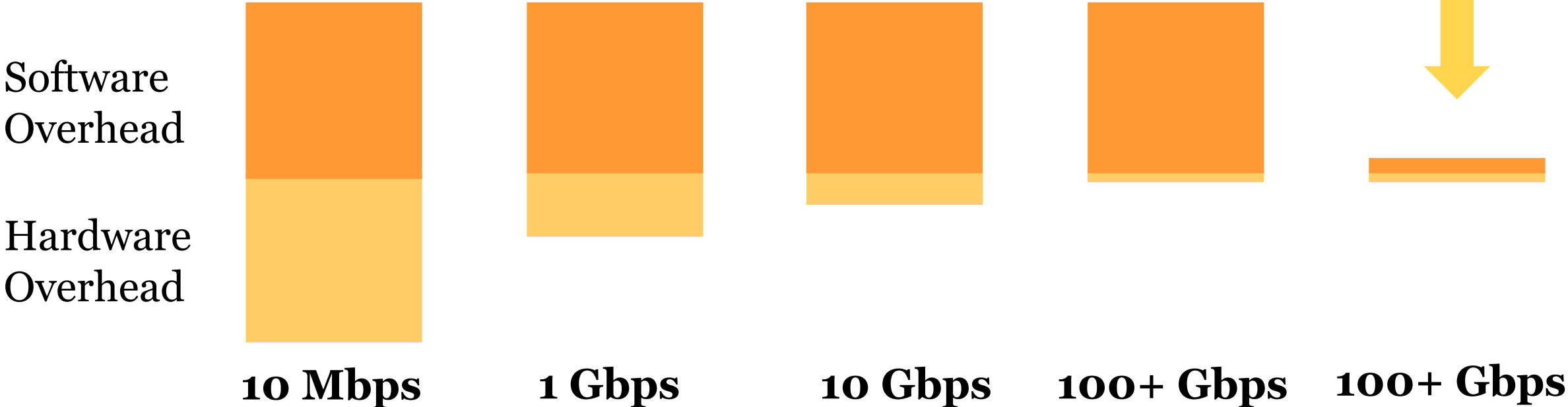
Networking Capabilities

- Bandwidth is increasing considerably
- Transmission speeds have improved (optical networking, etc)
- Messaging overhead is *still* an issue



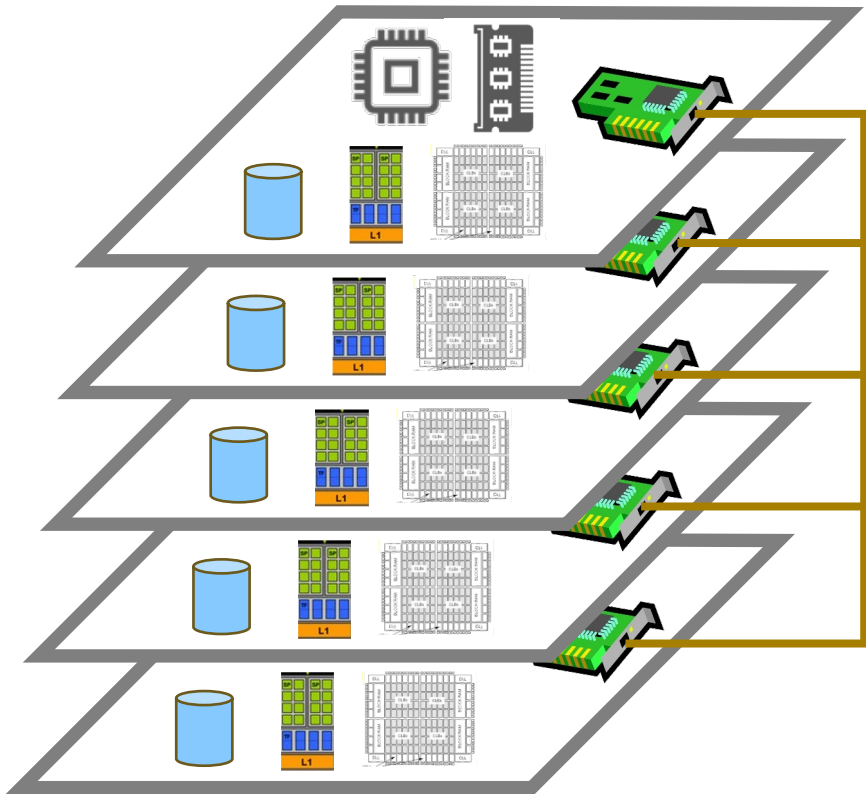
Networking Capabilities

- Bandwidth is increasing considerably
- Transmission speeds have improved (optical networking, etc)
- Messaging overhead is *still* an issue

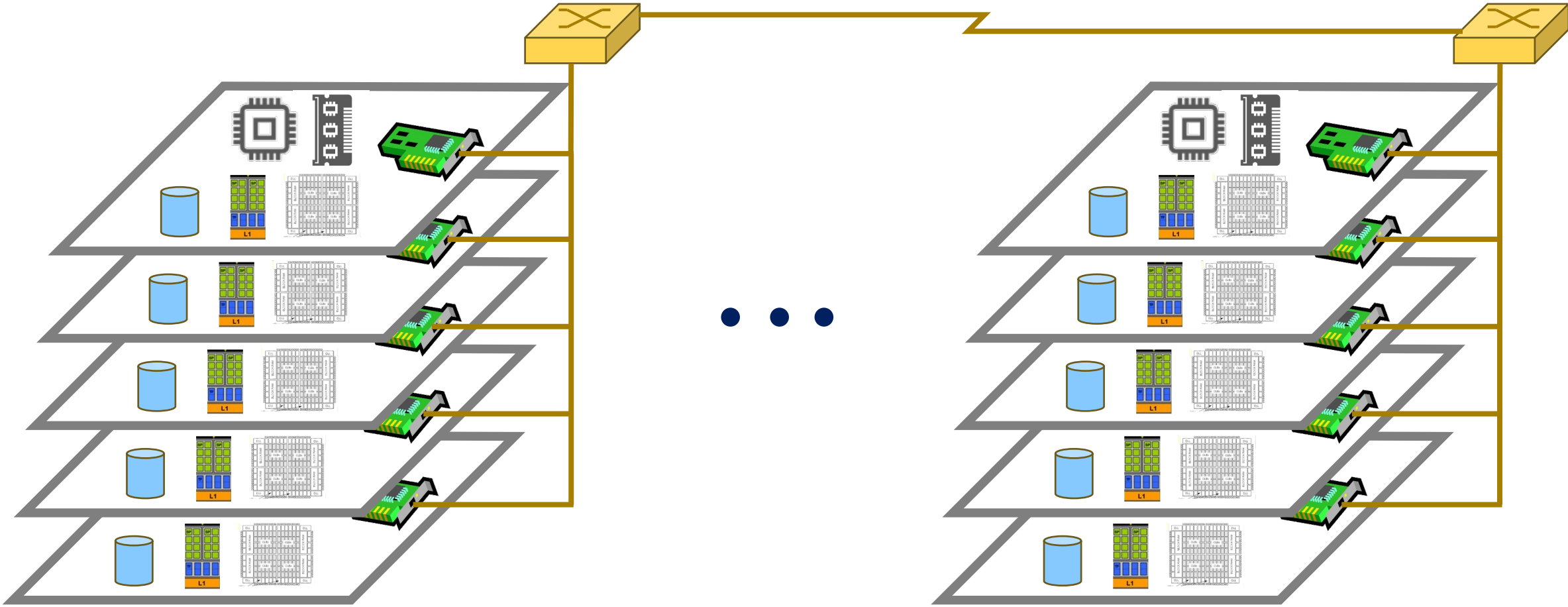


CURRENT ENVIRONMENT

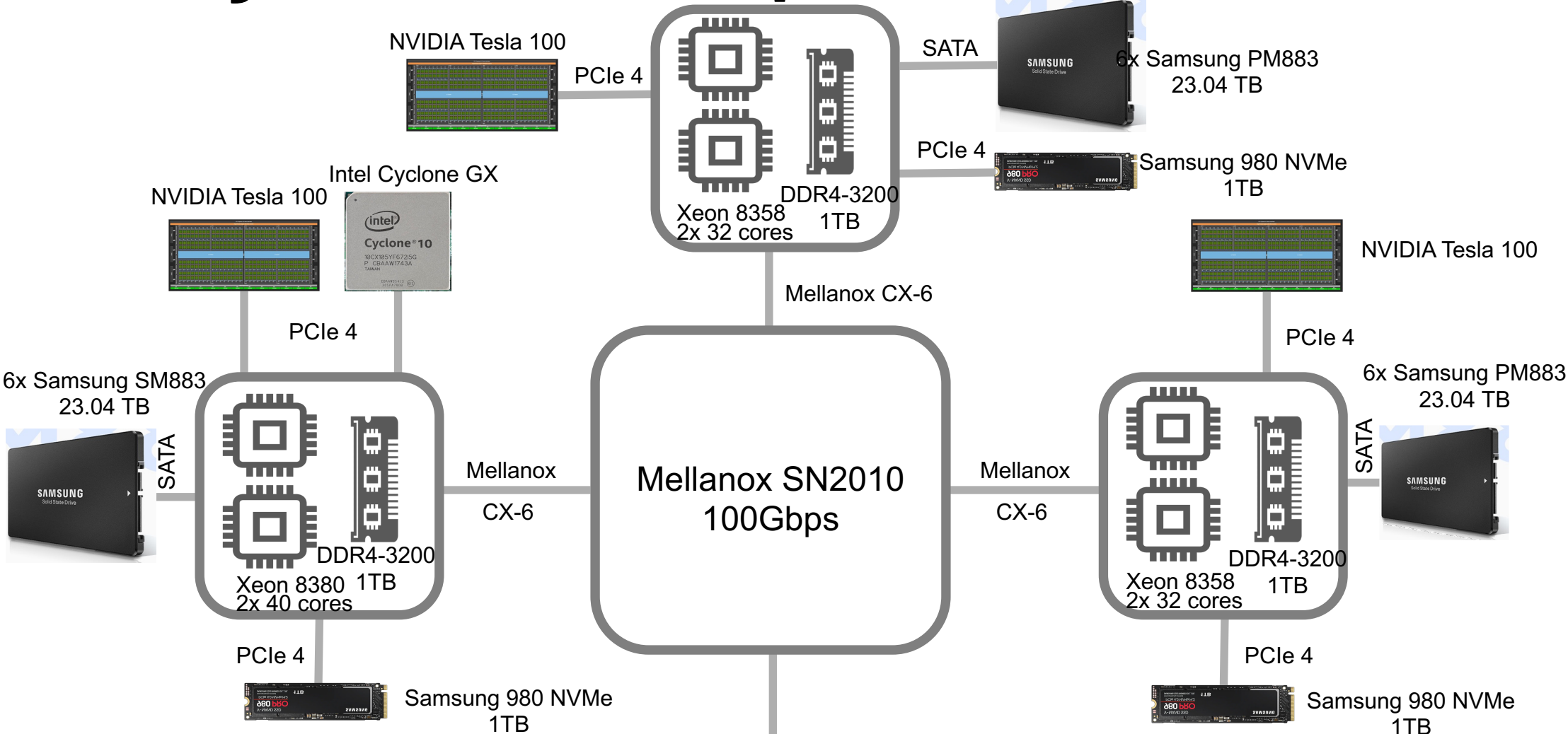
Current Data Centre Rack and Blade Design



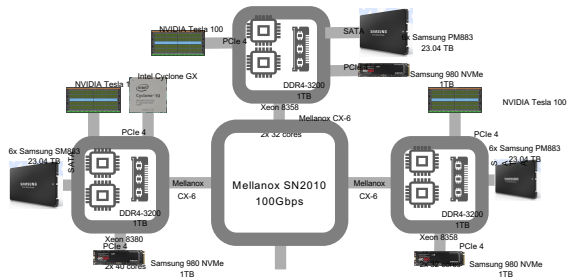
Current Data Centre Rack and Blade Design



Our Existing Platform as Example



Our Existing Platform as Example

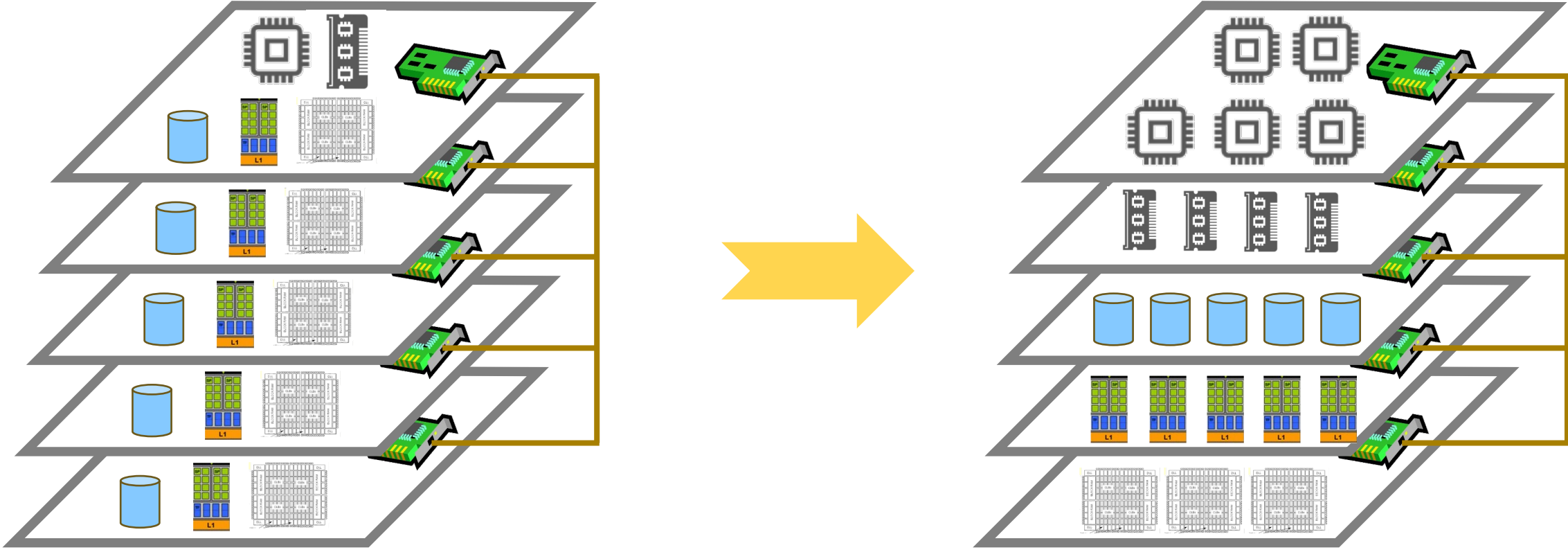


DISAGGREGATED HETEROGENEOUS PLATFORM

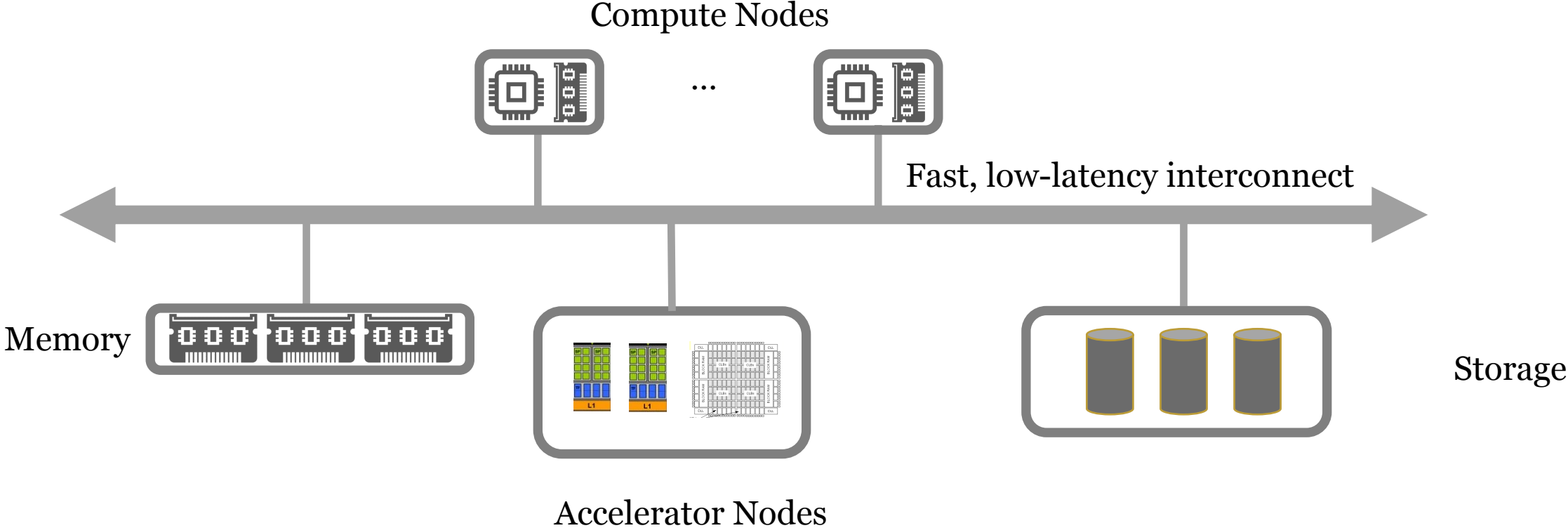
DISAGGREGATED HETEROGENEOUS PLATFORM

Disaggregation

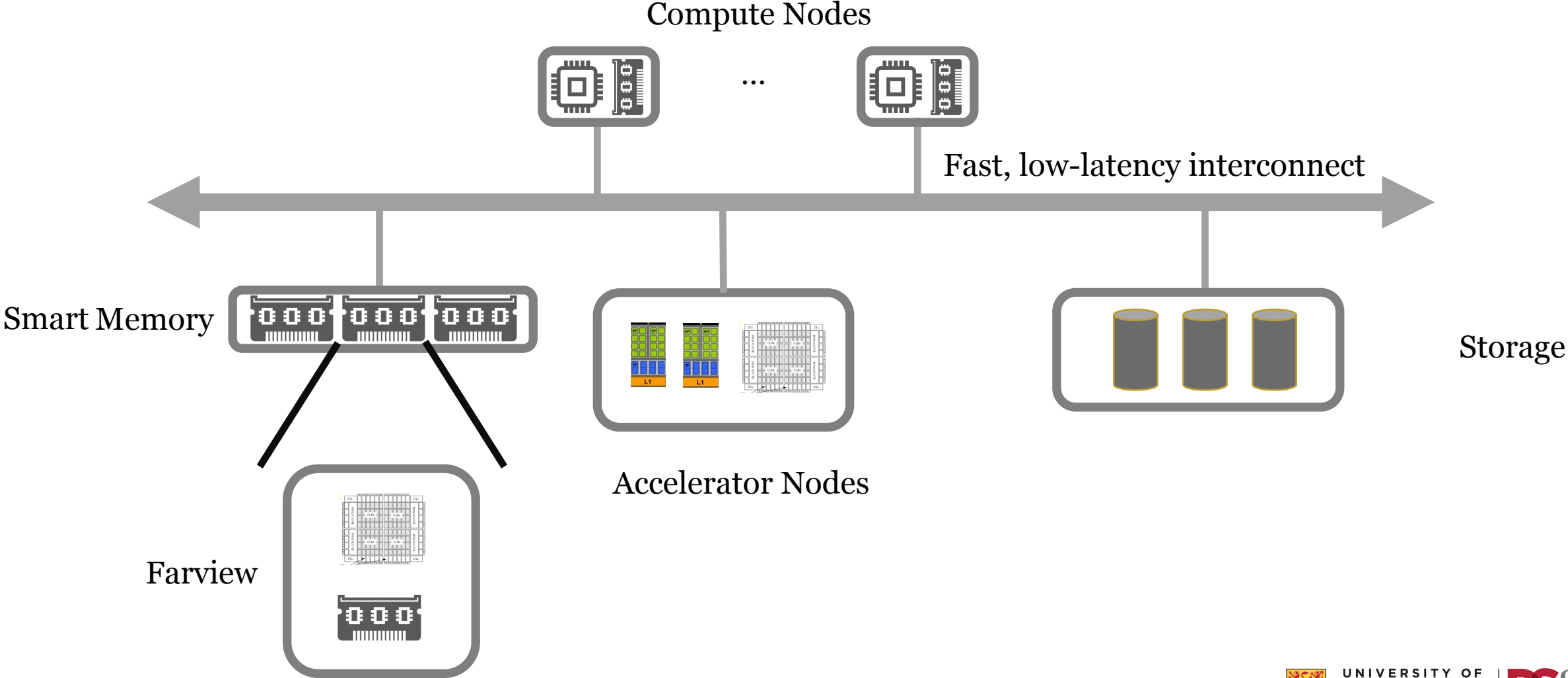
Disaggregated Rack and Blade Design



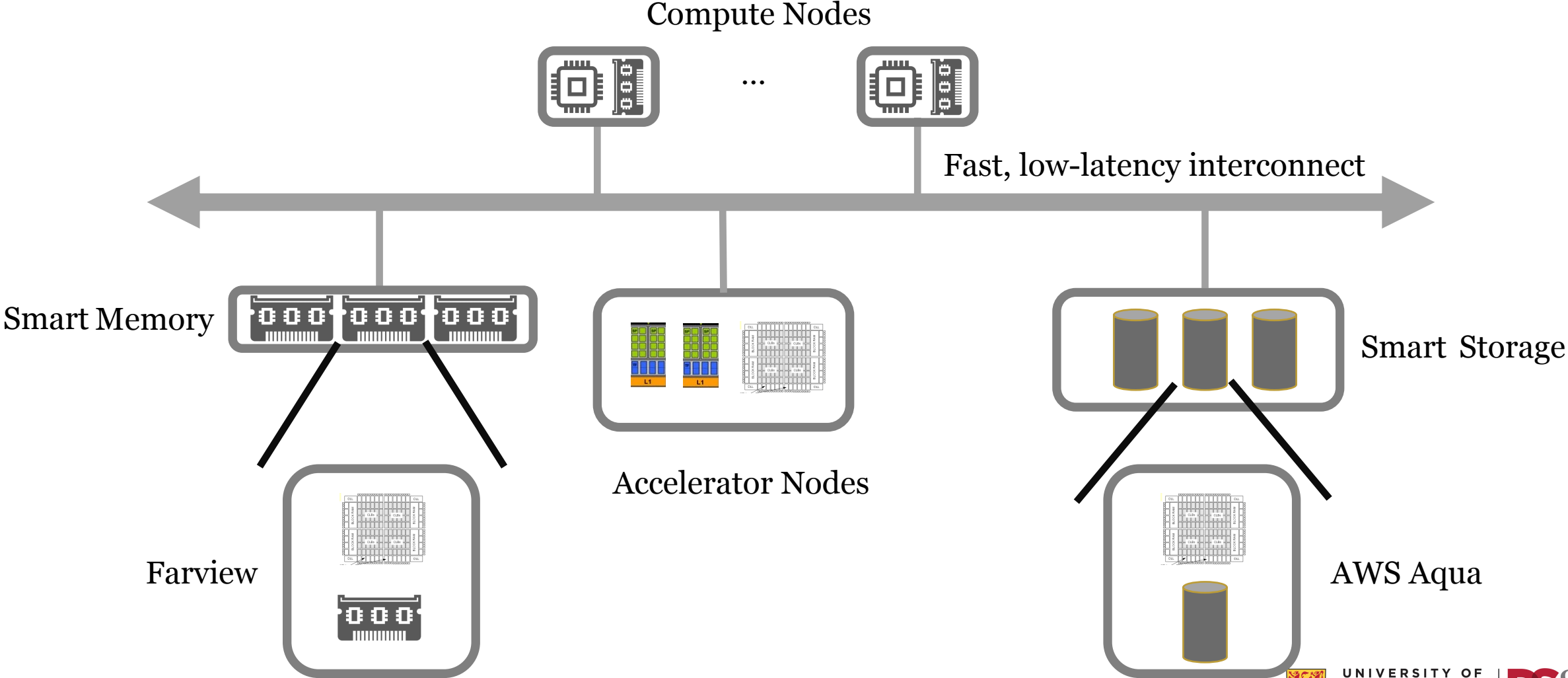
Disaggregated Heterogeneous Platform



Disaggregated Heterogeneous Platform



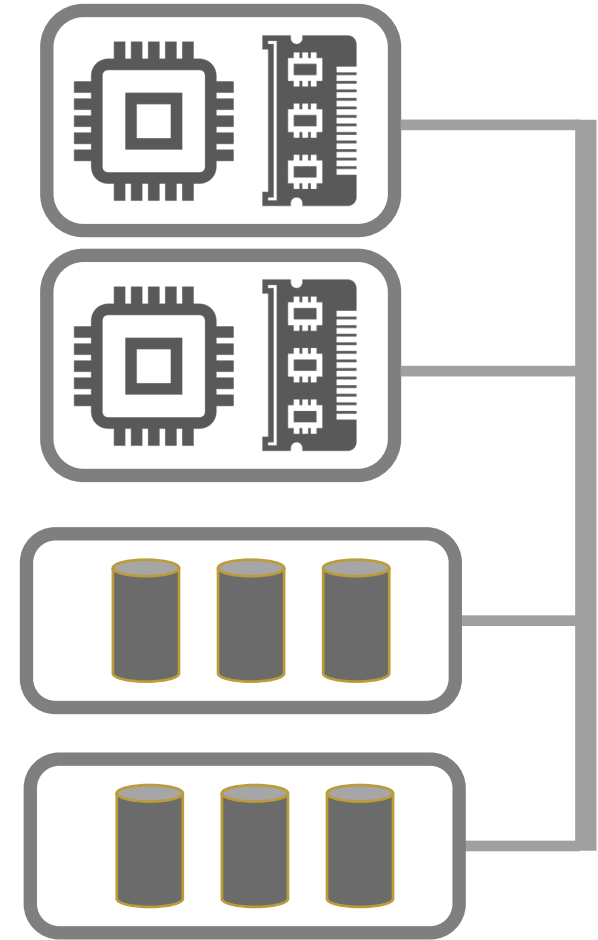
Disaggregated Heterogeneous Platform



What Are We Disaggregating?

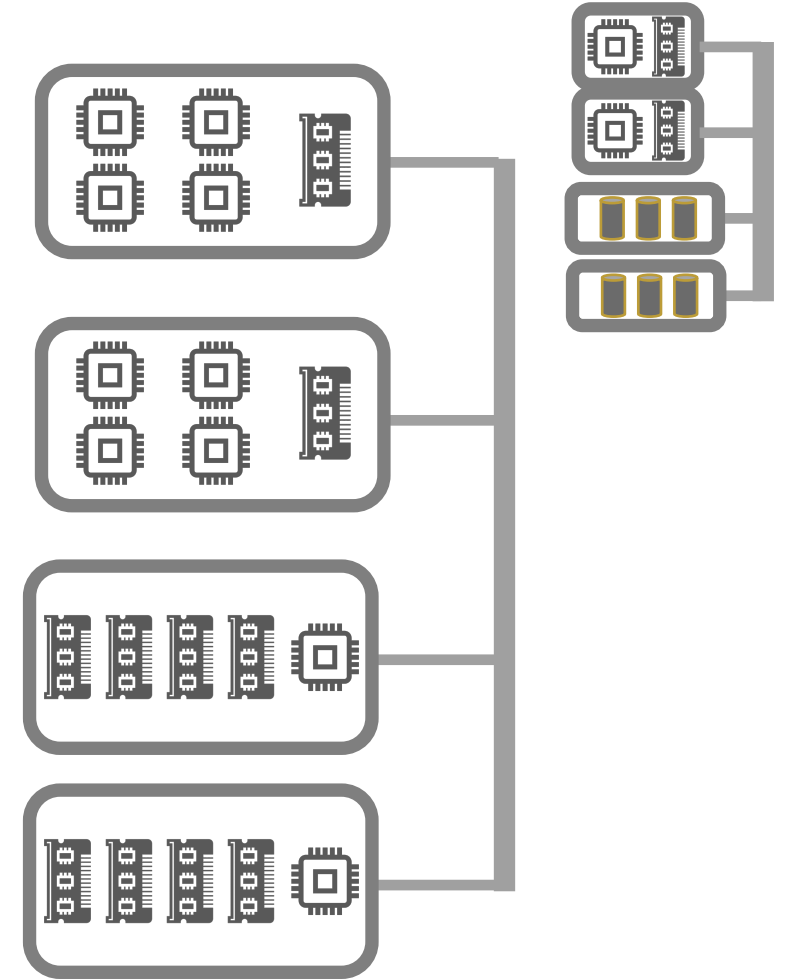
What Are We Disaggregating?

- Storage



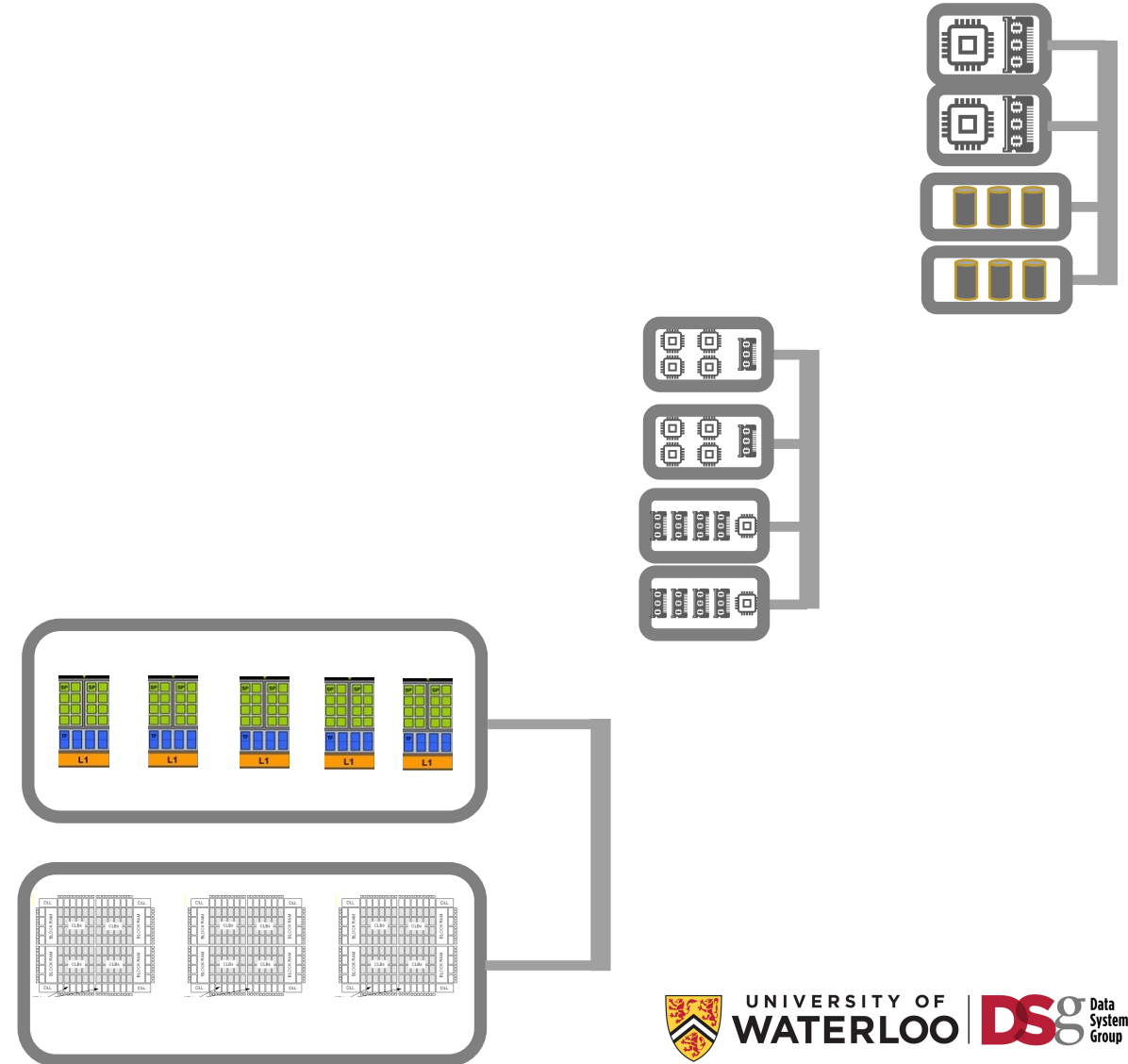
What Are We Disaggregating?

- Storage
- Memory



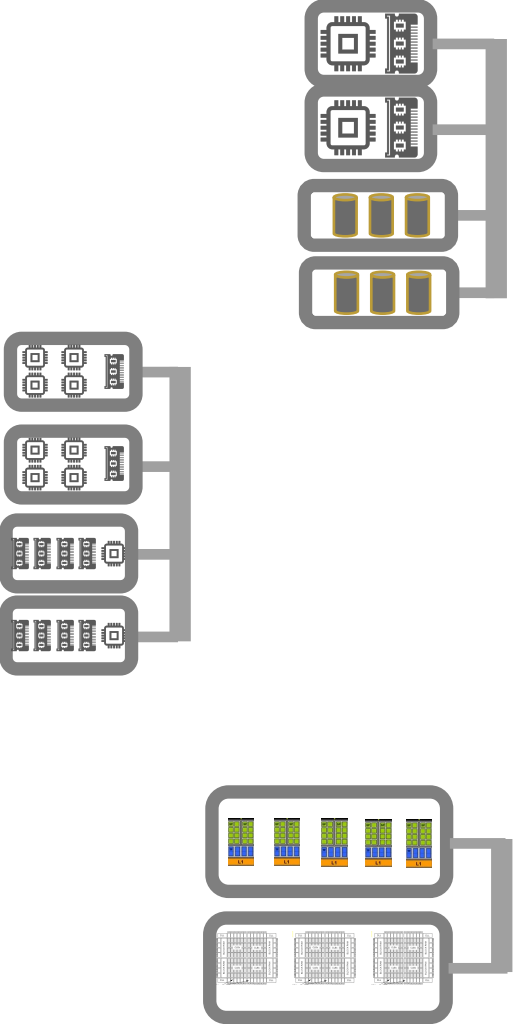
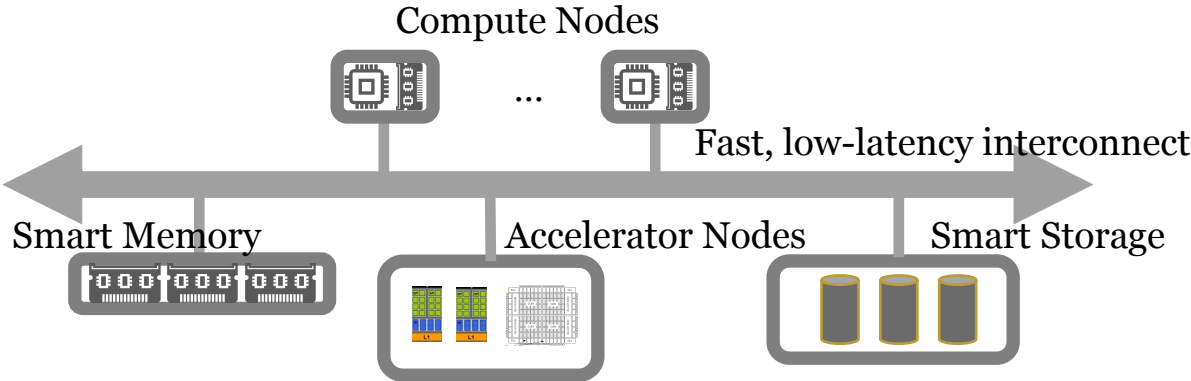
What Are We Disaggregating?

- Storage
- Memory
- Accelerators



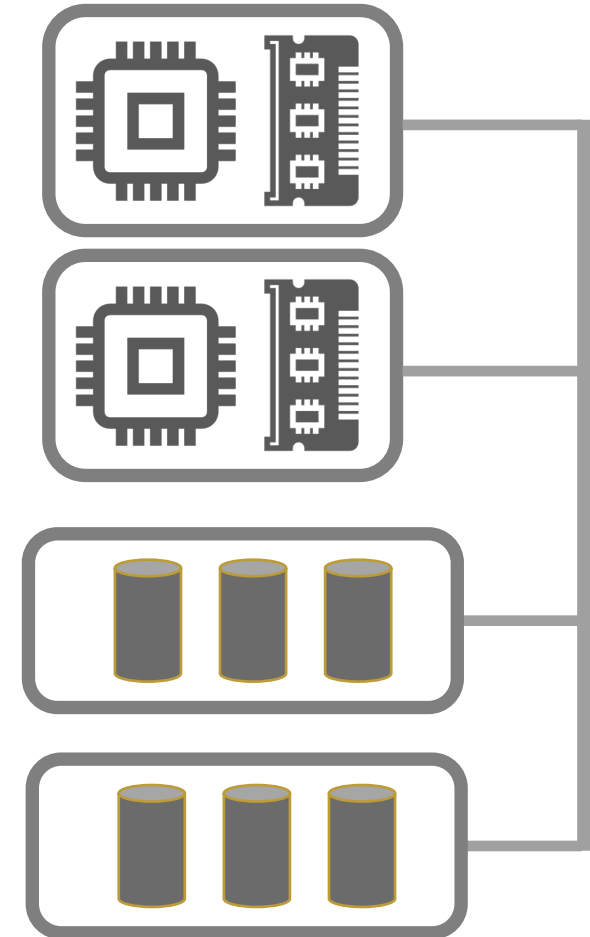
What Are We Disaggregating?

- Storage
- Memory
- Accelerators



Storage Disaggregation

- Cloud deployments prefer shared disk¹
 - Easier for crash recovery; no need for migration
- Easier for serverless computing
 - Faster startup, easier shutdown
- Independent and easy scaling of storage
- TCP/IP for connection may be sufficient



Storage Disaggregation is Well Understood

OLTP Systems

- Amazon Aurora¹



- Microsoft Socrates²



- Google AlloyDB³



- Alibaba PolarDB⁴



OLAP Systems

- Snowflake⁵



- Amazon Redshift⁶



¹ A. Verbitsky et al. Amazon Aurora: Design Considerations for High Throughput Cloud-Native Relational Databases, *Proc. SIGMOD*, 2017.

SiftDB'23/45 ² P. Antonopoulos et al. Socrates: The New SQL Server in the Cloud, *Proc. SIGMOD*, 2019.

³ <https://cloud.google.com/blog/products/databases/alloydb-for-postgresql-intelligent-scalable-storage>

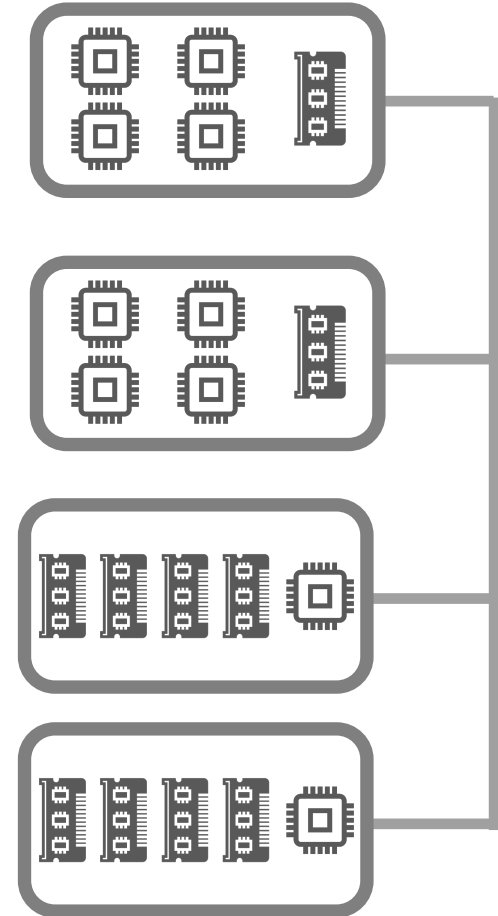
⁴ F. Li. Cloud-Native Database Systems at Alibaba: Opportunities and Challenges, *Proc. VLDB Endow.*, 2019.

⁵ B. Dageville et al. The Snowflake Elastic Data Warehouse, *Proc. SIGMOD*, 2016.

⁶ N. Armenatzoglou et al. Amazon Redshift Re-invented, *Proc. SIGMOD*, 2022.

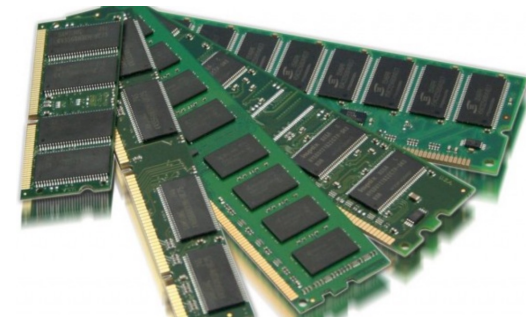
Memory Disaggregation

- Newer direction



Memory Disaggregation

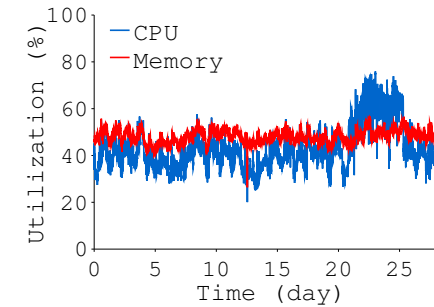
- Newer direction
- The case for memory disaggregation
 - DRAM is an **expensive** resource in the cloud – 50% of server cost on Azure¹



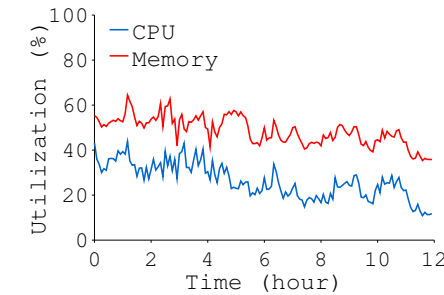
¹ H. Li et al. Pond: CXL-Based Memory Pooling Systems for Cloud Platforms, arXiv 2203.00241, 2022.

Memory Disaggregation

- Newer direction
- The case for memory disaggregation
 - DRAM is an **expensive** resource in the cloud – 50% of server cost on Azure¹
 - Memory utilization is low in the (current) cloud²



Google Cluster



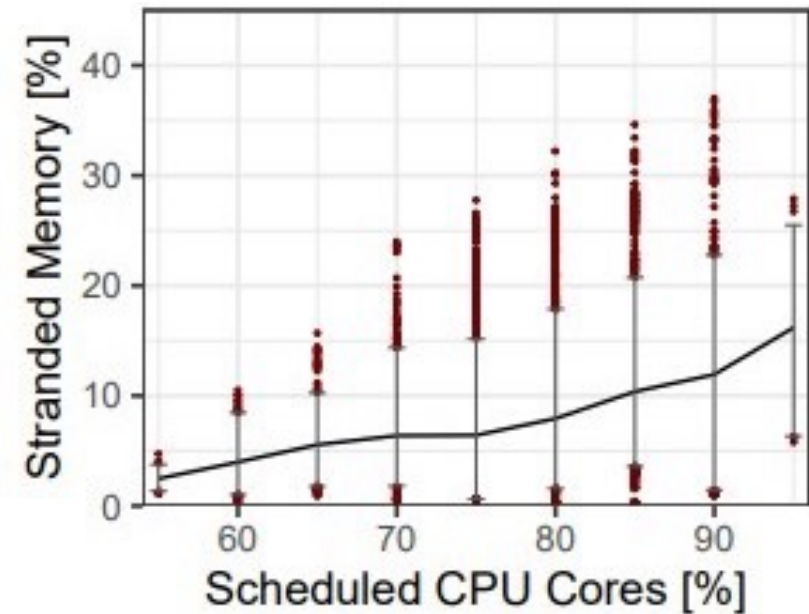
Alibaba Cluster

¹ H. Li et al. Pond: CXL-Based Memory Pooling Systems for Cloud Platforms, arXiv 2203.00241, 2022.

² Y. Shen et al. Pond: LegoOS: A Disseminated, Distributed OS for Hardware Resource Disaggregation, *Proc. OSDI*, 2018.

Memory Disaggregation

- Newer direction
- The case for memory disaggregation
 - DRAM is an **expensive** resource in the cloud – 50% of server cost on Azure¹
 - Memory utilization is low in the (current) cloud²
 - Stranded memory (shows Azure³)



¹ H. Li et al. Pond: CXL-Based Memory Pooling Systems for Cloud Platforms, arXiv 2203.00241, 2022.

² Y. Shen et al. Pond: LegoOS: A Disseminated, Distributed OS for Hardware Resource Disaggregation, *Proc. OSDI*, 2018.

³ <https://www.nextplatform.com/2022/07/11/microsoft-azure-blazes-the-disaggregated-memory-trail-with-znuma/>

Memory Disaggregation

- Newer direction
- The case for memory disaggregation
 - DRAM is an **expensive** resource in the cloud – 50% of server cost on Azure¹
 - Memory utilization is low in the (current) cloud²
 - Stranded memory (shows Azure³)
 - Independent and elastic scaling of compute and memory⁴

AWS EC2 Instance	vCPU	DRAM (GB)
r6g.medium	1	8
r6g.large	2	16
r6g.xlarge	4	32
r6g.2xlarge	8	64
r6g.4xlarge	16	128
r6g.8xlarge	32	256
r6g.16xlarge	64	512

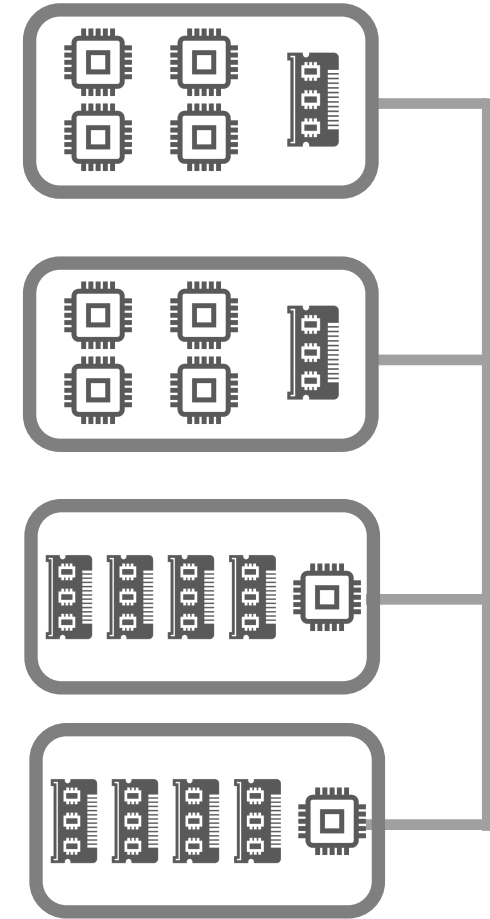
¹ H. Li et al. Pond: CXL-Based Memory Pooling Systems for Cloud Platforms, arXiv 2203.00241, 2022.

² Y. Shen et al. Pond: LegoOS: A Disseminated, Distributed OS for Hardware Resource Disaggregation, *Proc. OSDI*, 2018.

³ <https://www.nextplatform.com/2022/07/11/microsoft-azure-blazes-the-disaggregated-memory-trail-with-znuma/>

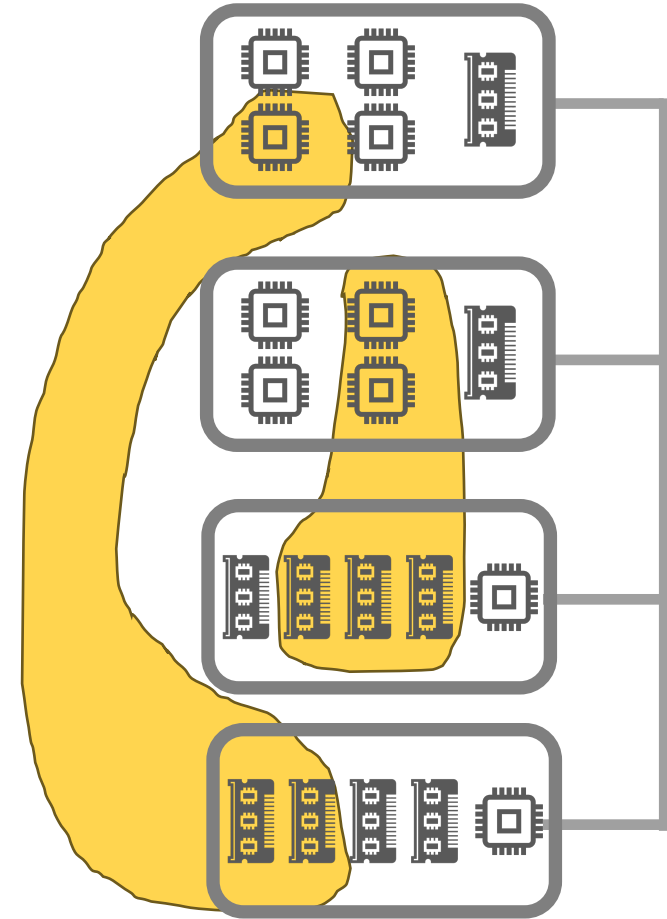
⁴ <https://aws.amazon.com/ec2/instance-types/>

The Case for Memory Disaggregation



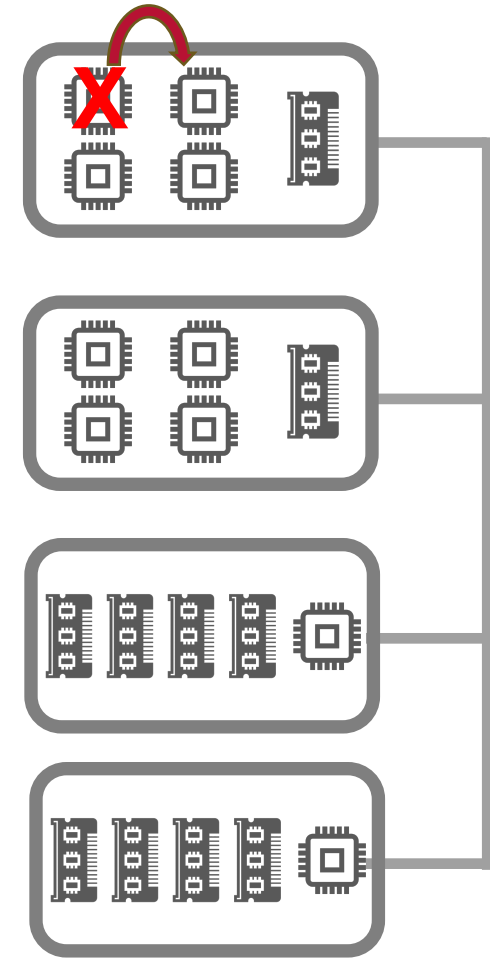
The Case for Memory Disaggregation

- Flexibility in VM configuration
 - Independent allocation



The Case for Memory Disaggregation

- Flexibility in VM configuration
 - Independent allocation
- Fault tolerance
 - Independent failure

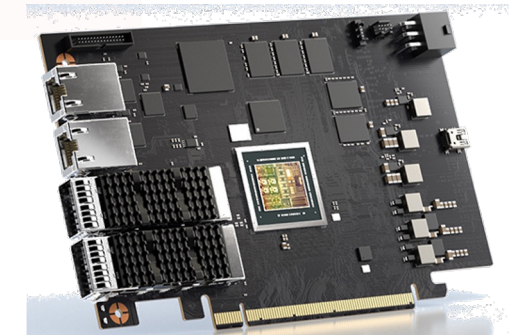


Critical Component in Disaggregation (Beyond Storage)

Low latency, high bandwidth
interconnect

Networking Hardware Improvements

- Basic/Foundational NIC
 - Simple network connection
 - Usually 1Gbps – 25Gbps
 - Relies on CPU for protocol processing – $\geq 30\%$ server CPU for higher speeds
- Smart NIC
 - Offload network protocol processing
 - Have their own processor, memory & OS
 - $\geq 50\text{Gbps}$
- DPU
 - Smart NIC + security + storage + ...
 - Custom chips and/or FPGAs



Networking Software Improvements



100+ Gbps

Networking Software Improvements

- High overhead of TCP



100+ Gbps

Networking Software Improvements

- High overhead of TCP
 - “Data center tax”¹



100+ Gbps

Networking Software Improvements

- High overhead of TCP
 - “Data center tax”¹
 - “Modern networking hardware enables roundtrip times of a few microseconds for short messages. The transport protocol must not add significantly to this latency”²

It's Time to Replace TCP in the Datacenter

John Ousterhout
Stanford University

arXiv 2210.00714

January 18, 2023



100+ Gbps

¹ L. Barroso, et al. Attack of the Killer Microseconds. *Commun. ACM*, 60(4):48–54, March 2017.

² J. Ousterhout. It's Time to Replace TCP in the Datacenter, arXiv 2210.00714, 2023.

Networking Software Improvements

- High overhead of TCP
 - “Data center tax”¹
 - “Modern networking hardware enables roundtrip times of a few microseconds for short messages. The transport protocol must not add significantly to this latency”²
- Low-overhead protocol
 - RDMA (Infiniband, RoCE)
 - CXL

It’s Time to Replace TCP in the Datacenter

John Ousterhout
Stanford University
January 18, 2023

arXiv 2210.00714



100+ Gbps

¹ L. Barroso, et al. Attack of the Killer Microseconds. *Commun. ACM*, 60(4):48–54, March 2017.

² J. Ousterhout. It’s Time to Replace TCP in the Datacenter, arXiv 2210.00714, 2023.

Networking Software Improvements

- High overhead of TCP
 - “Data center tax”¹
 - “Modern networking hardware enables roundtrip times of a few microseconds for short messages. The transport protocol must not add significantly to this latency”²
- Low-overhead protocol
 - RDMA (Infiniband, RoCE)
 - CXL
- Efficient RPC protocol – higher level

It's Time to Replace TCP in the Datacenter

John Ousterhout
Stanford University
January 18, 2023

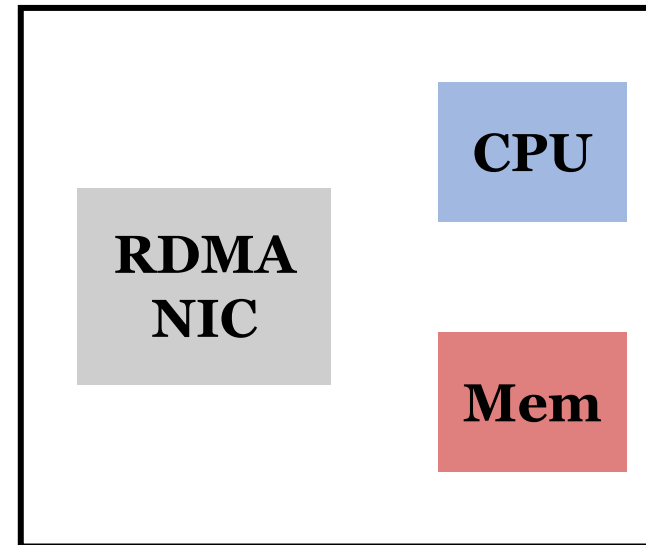
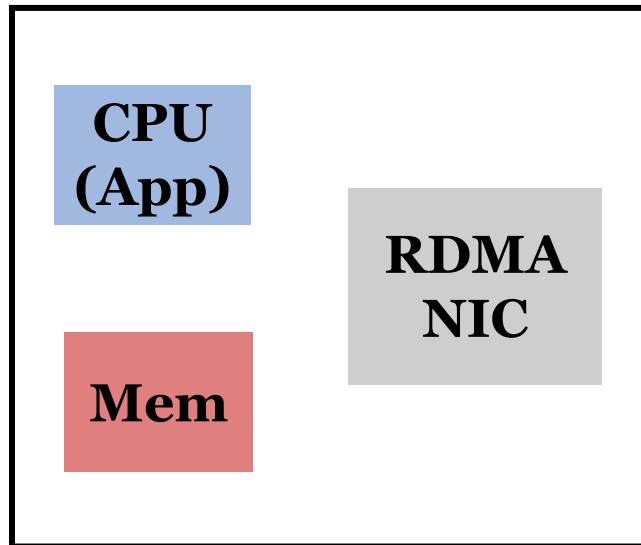
arXiv 2210.00714



100+ Gbps

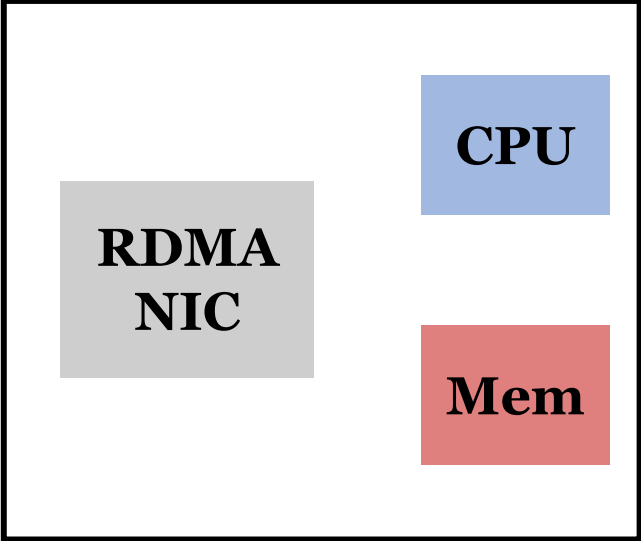
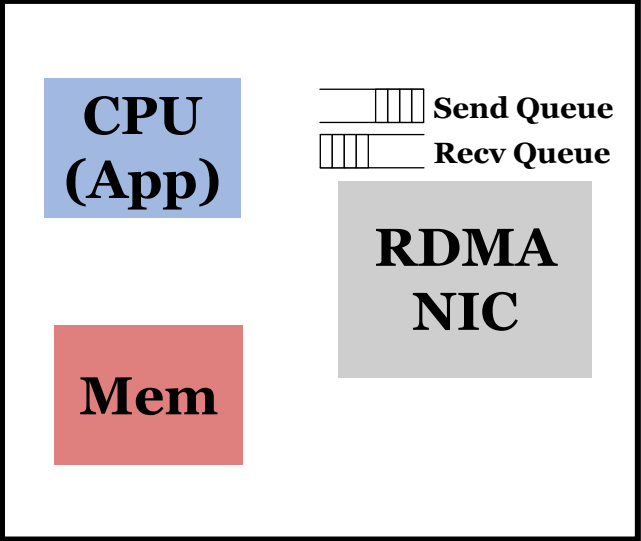
RDMA - Remote Direct Memory Access

Physical View



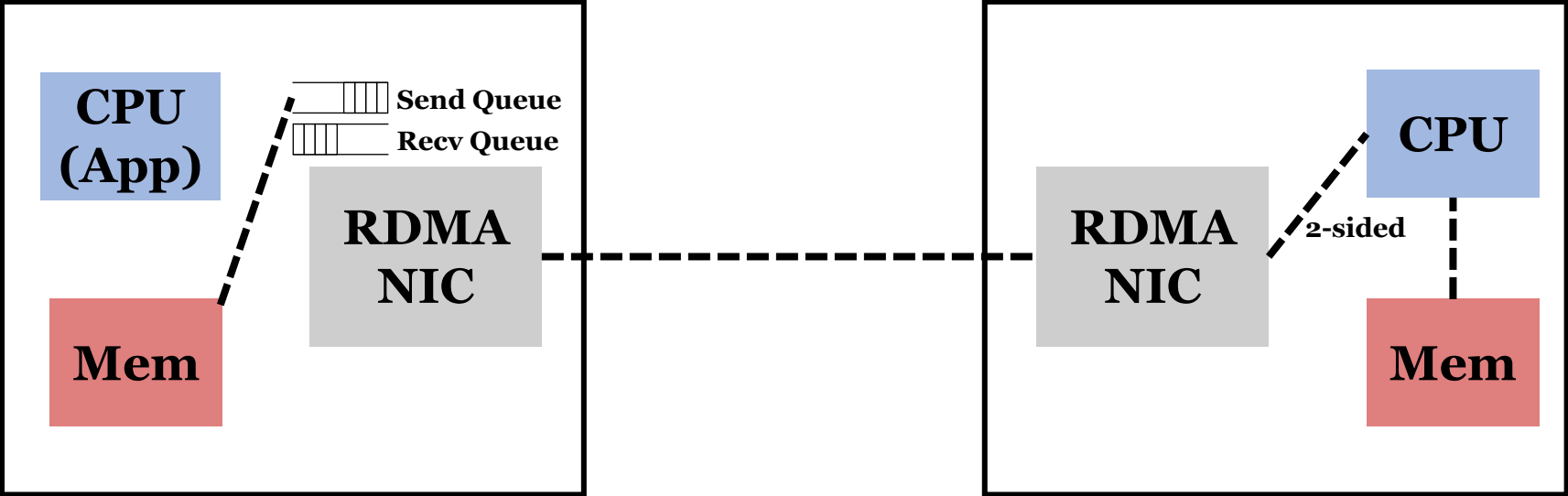
RDMA - Remote Direct Memory Access

Physical View



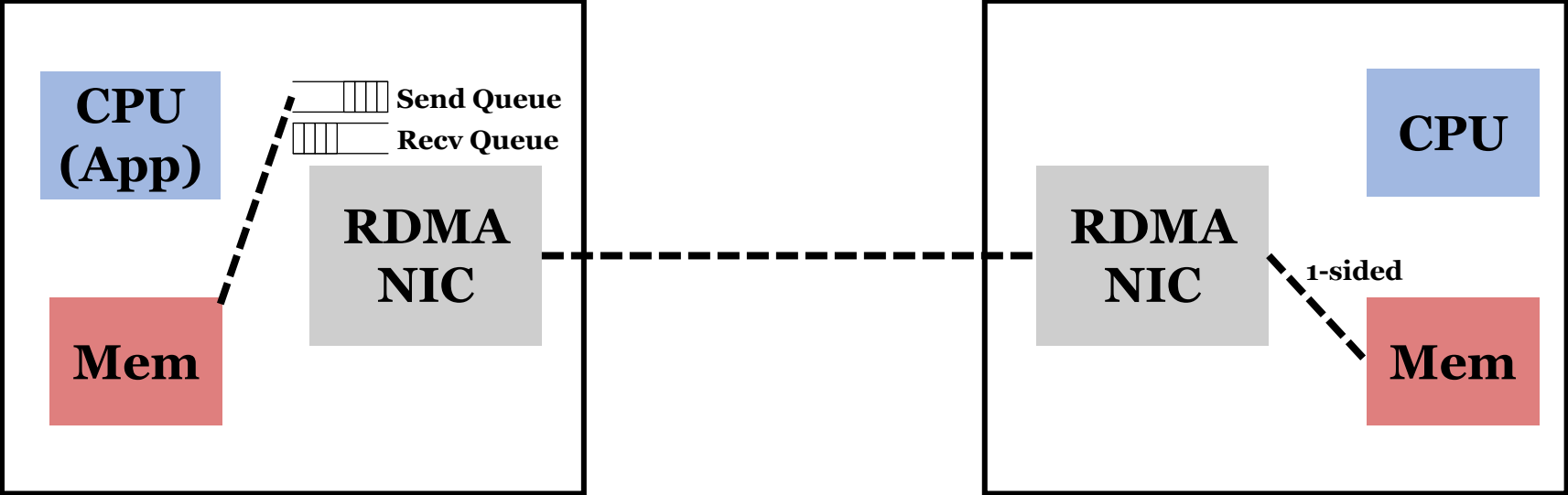
RDMA - Remote Direct Memory Access

Physical View



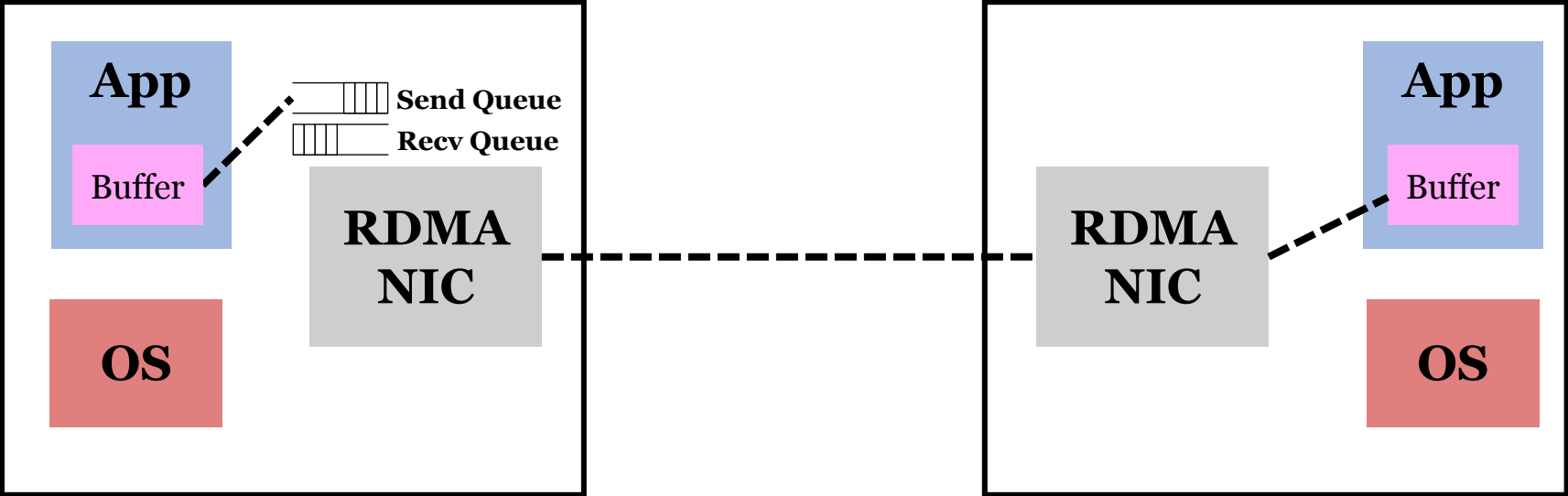
RDMA - Remote Direct Memory Access

Physical View



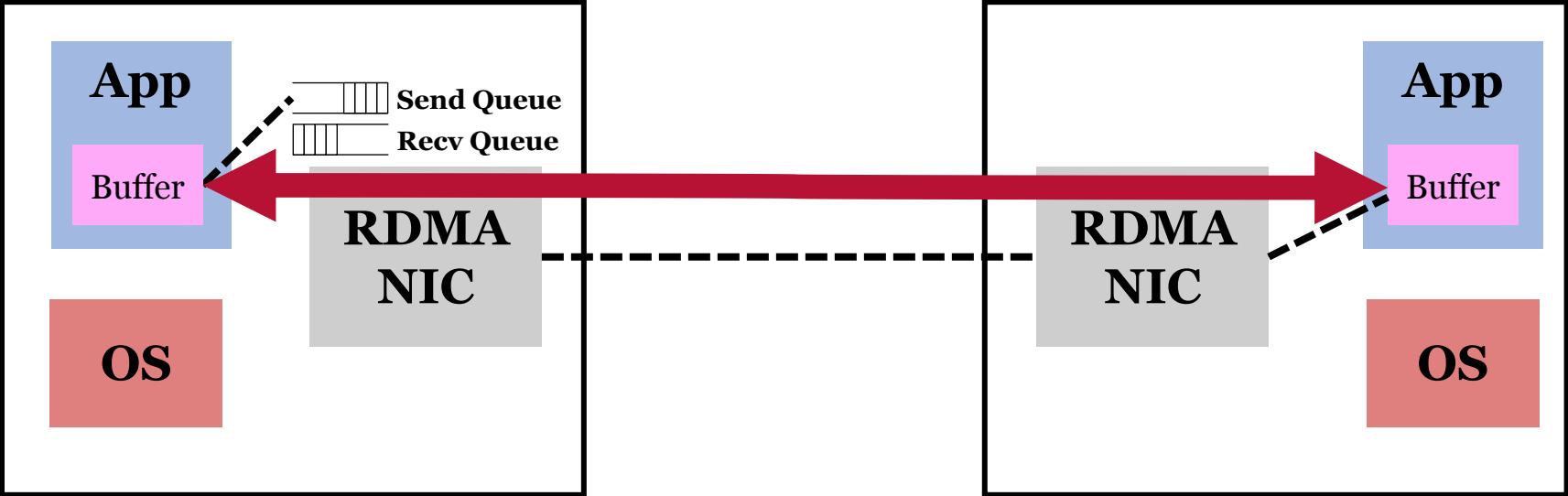
RDMA - Remote Direct Memory Access

Process View



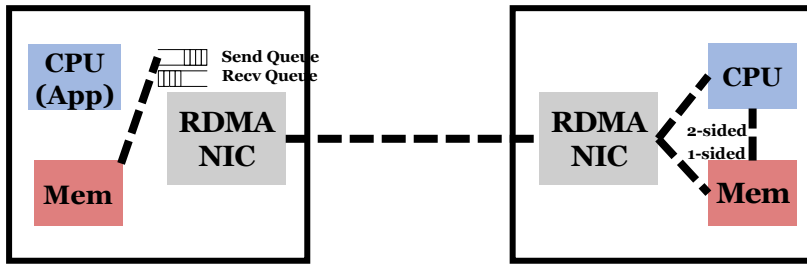
RDMA - Remote Direct Memory Access

Process View



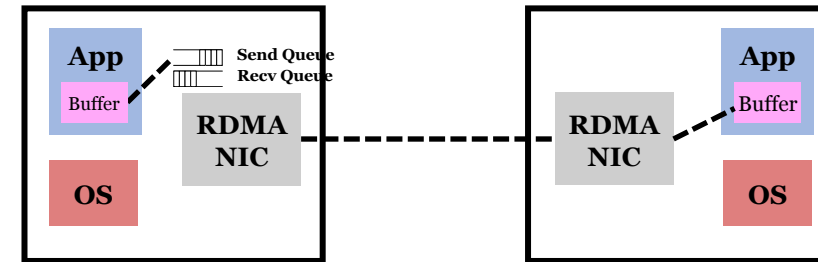
RDMA - Remote Direct Memory Access

Physical View



- Zero-copy networking
- OS bypass

Process View



- CPU bypass
- Message-based communication

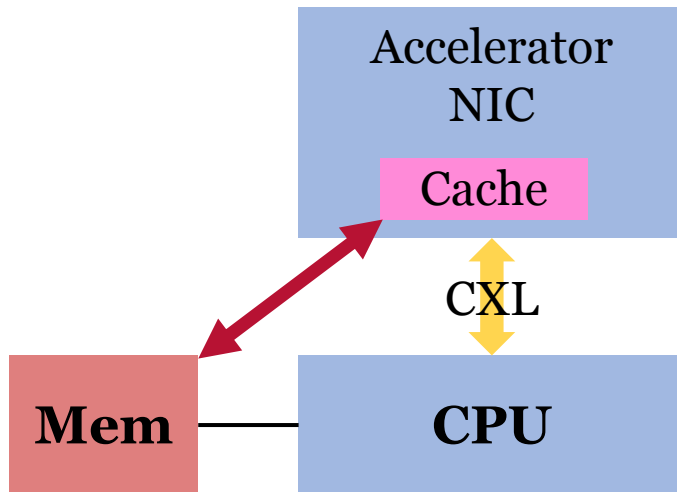
CXL - Compute Express Link

- High bandwidth, low latency, cache coherent interconnect
- Avoids memory copy between application memory and NIC buffers
- Does not require a CPU (or controller) on memory nodes

CXL - Compute Express Link

- High bandwidth, low latency, cache coherent interconnect
- Avoids memory copy between application memory and NIC buffers
- Does not require a CPU (or controller) on memory nodes

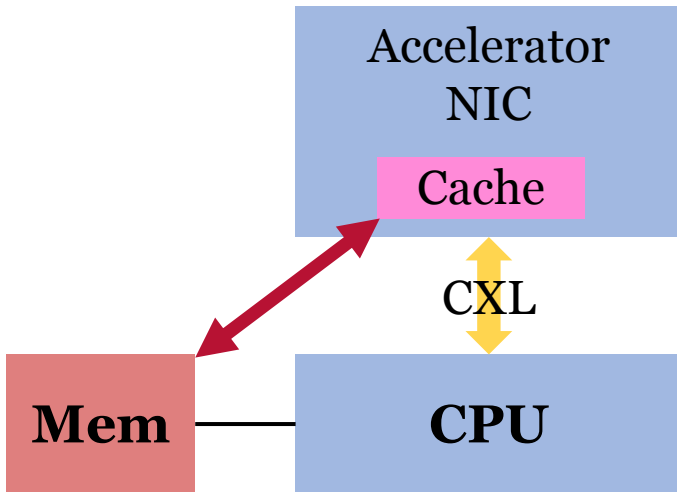
Type 1



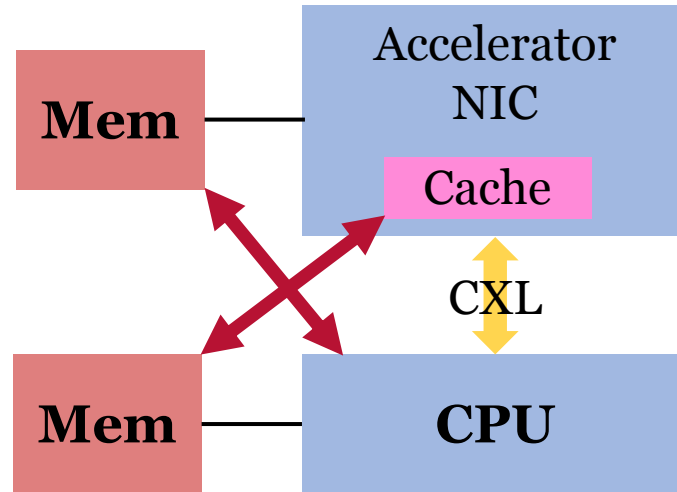
CXL - Compute Express Link

- High bandwidth, low latency, cache coherent interconnect
- Avoids memory copy between application memory and NIC buffers
- Does not require a CPU (or controller) on memory nodes

Type 1

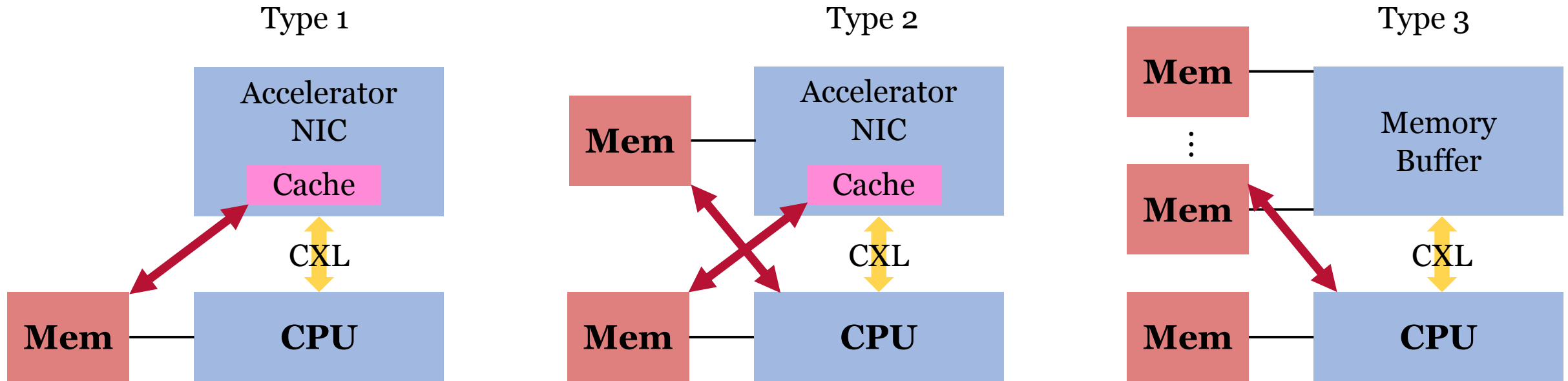


Type 2



CXL - Compute Express Link

- High bandwidth, low latency, cache coherent interconnect
- Avoids memory copy between application memory and NIC buffers
- Does not require a CPU (or controller) on memory nodes



CXL vs RDMA

RDMA

- Copies data from application memory to NIC buffers
- Copies data across the network
- Requires handling of cache coherence (in some configurations)
- Requires a CPU/controller at memory nodes

CXL

- Does not copy data into NIC buffers
- Does not copy data across the network, accesses remote memory
- Provides a hardware supported coherent cache
- No need for CPU/controller at memory nodes

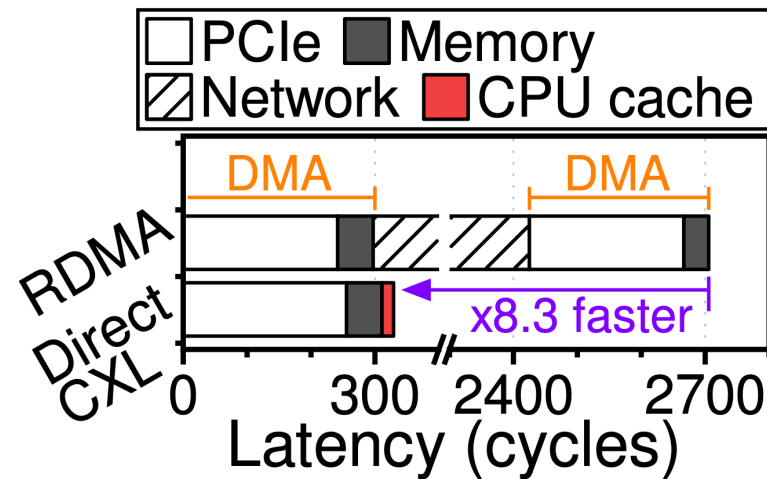
CXL vs RDMA

RDMA

- Copies data from application memory to NIC buffers
- Copies data across the network
- Requires handling of cache coherence (in some configurations)
- Requires a CPU/controller at memory nodes

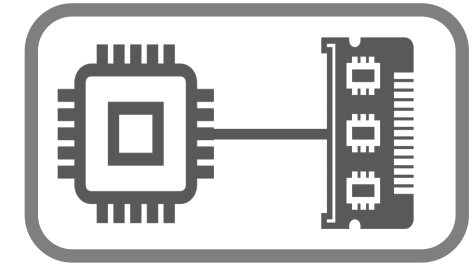
CXL

- Does not copy data into NIC buffers
- Does not copy data across the network, accesses remote memory
- Provides a hardware supported coherent cache
- No need for CPU/controller at memory nodes

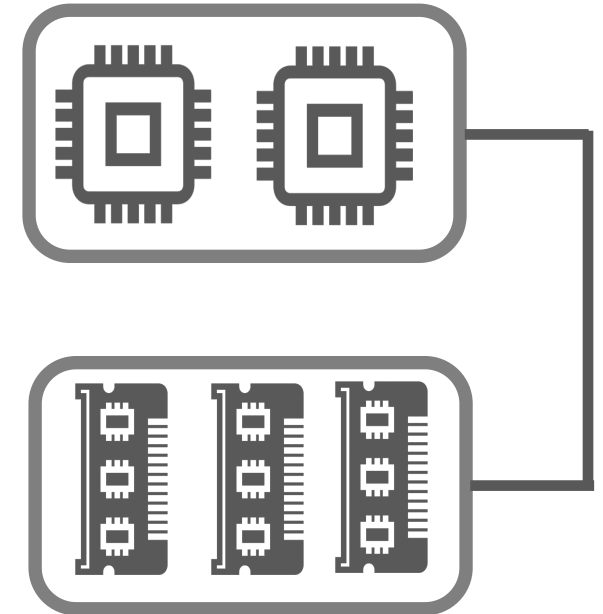


Main Issue with Memory Disaggregation

Access Latency



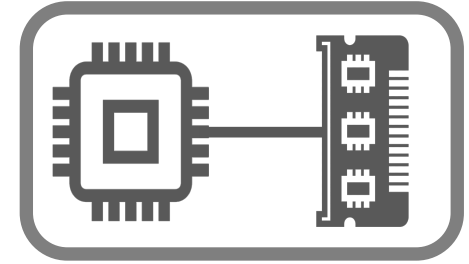
VS



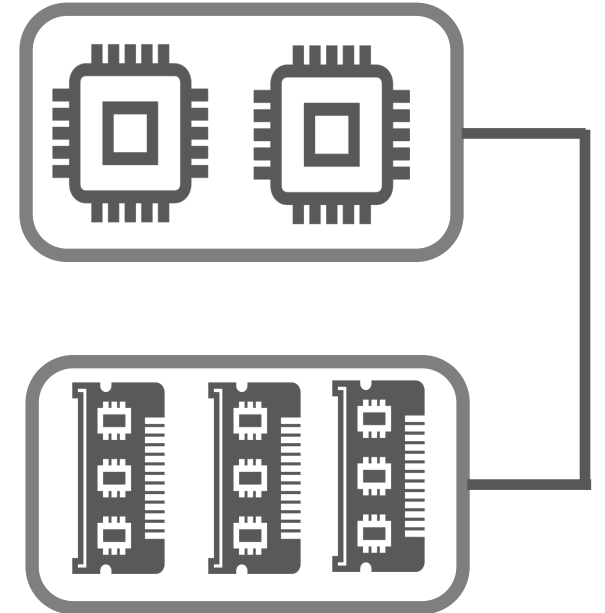
Main Issue with Memory Disaggregation

Access Latency

Near-data Processing
+
Optimized Design

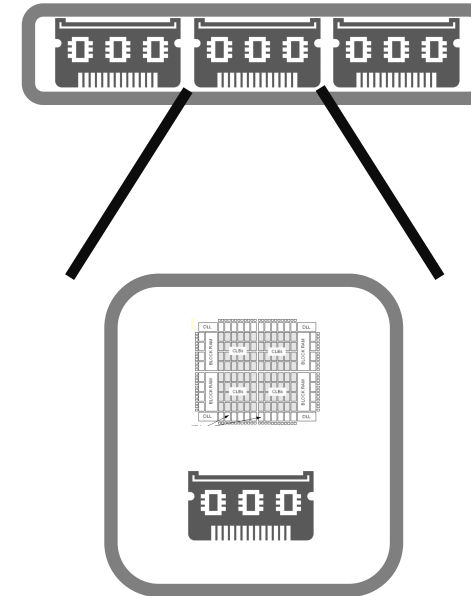


VS



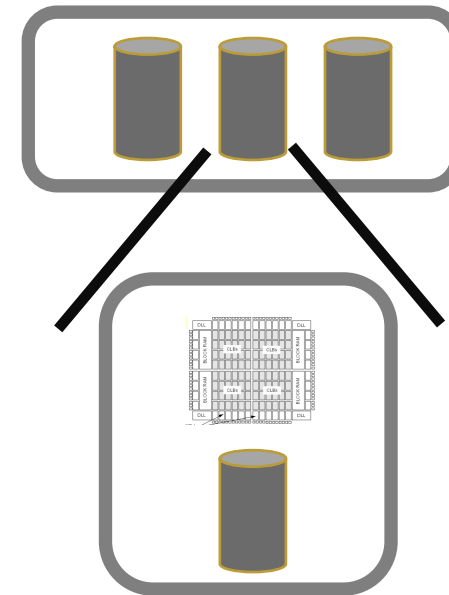
Near-Data Processing: The Case for Smart Memory

- Minimize data movement from memory to CPU
- Also
 - DRAM growing but not at the same speed as the demand
 - The independent provisioning issue again
- Farview: buffer pool on disaggregated memory
 - Push down operators to memory → reduce data movement
 - Centralize buffer cache on disaggregated memory → reduce memory requirement on compute nodes



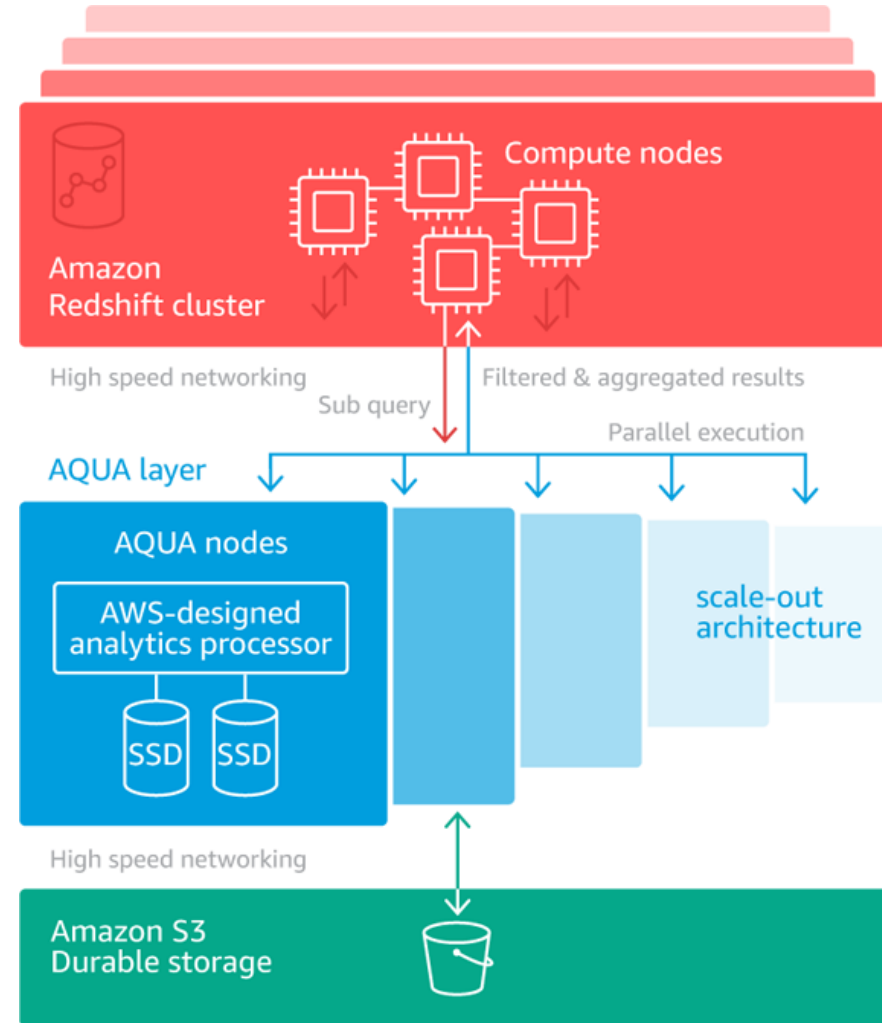
Near-Data Processing: The Case for Smart Storage

- Further restrict data movement
 - Pushing near-data computation to the storage
- Similar to old idea of database machines



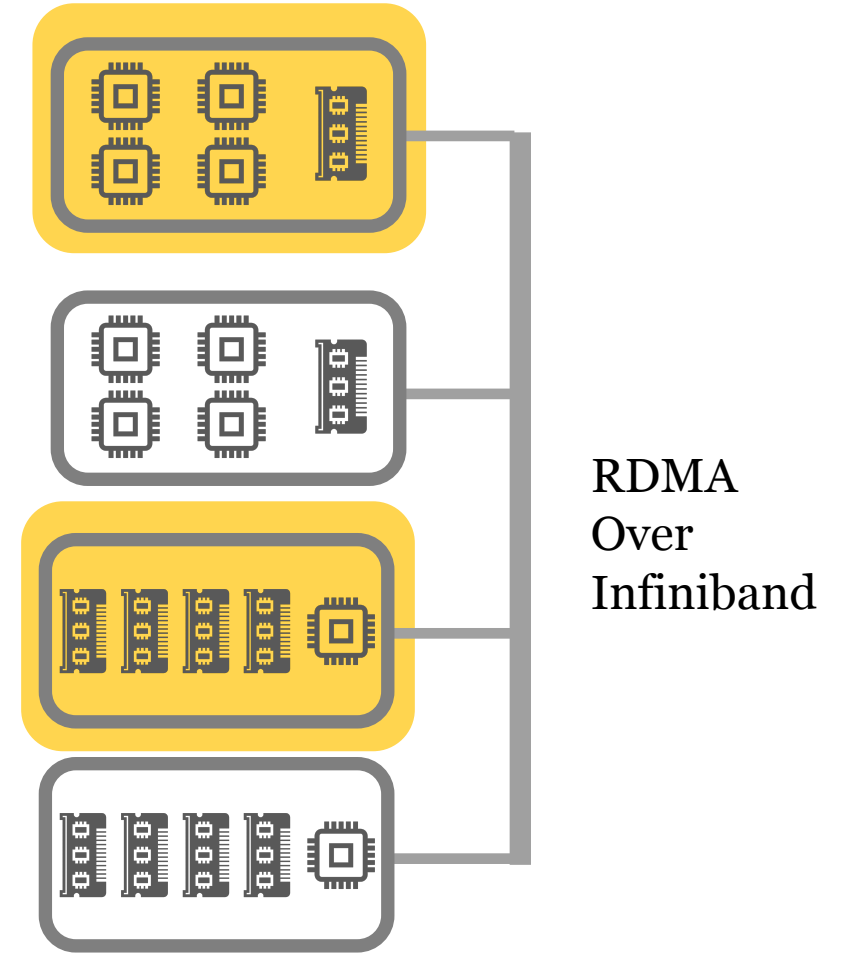
Near-Data Processing: The Case for Smart Storage

- Further restrict data movement
 - Pushing near-data computation to the storage
- Similar to old idea of database machines
- AWS Aqua nodes

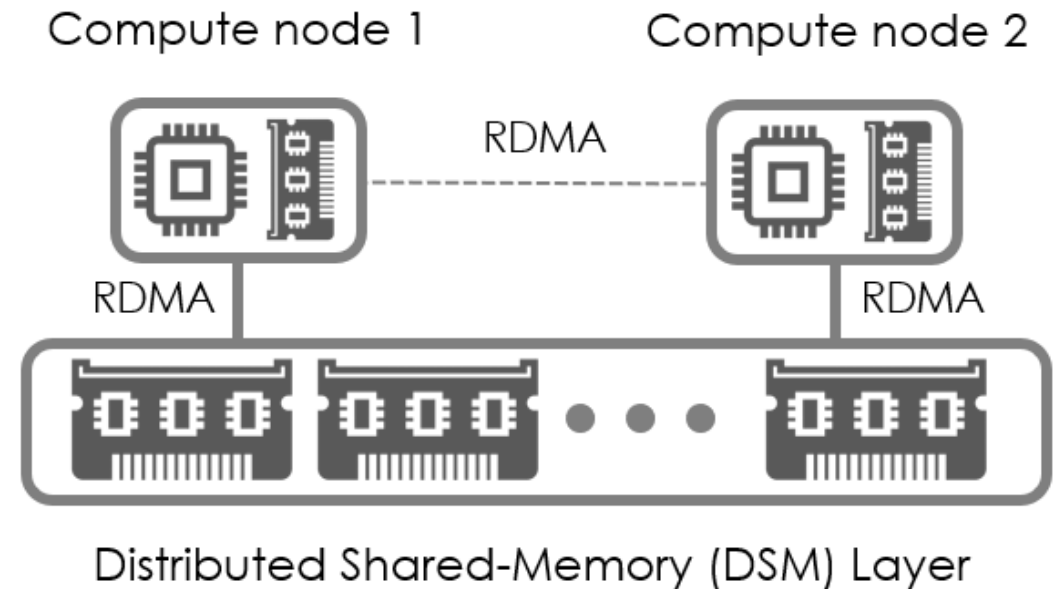


Our Work

- Scope: in-memory DBMS
- Compute nodes
 - Strong computing, but limited local memory
 - SQL parser, optimization, transaction, buffer
 - Each compute can read/write
- Memory nodes
 - Form a DSM layer accessed by all compute nodes

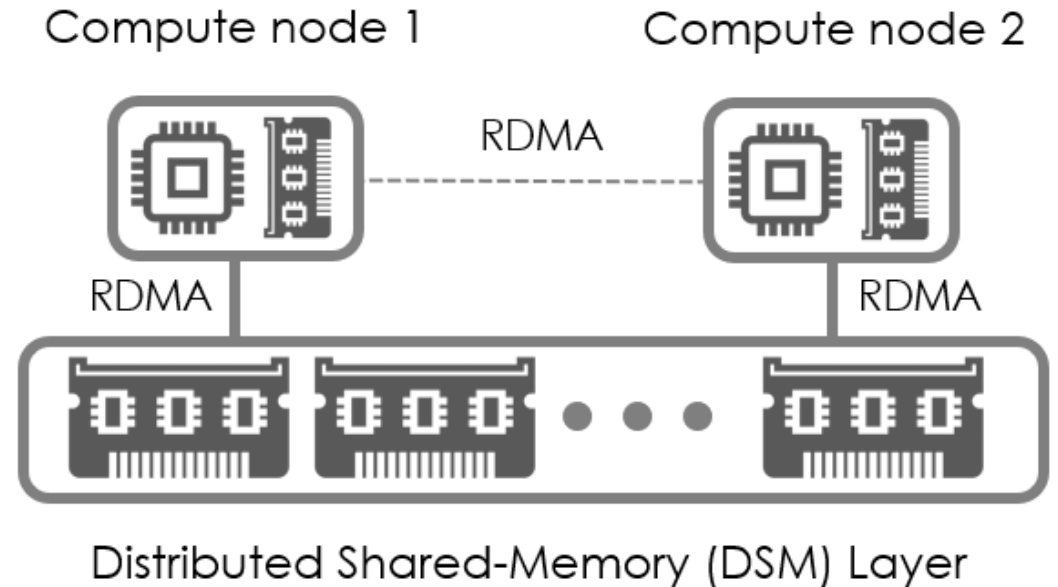


Transaction Execution Design



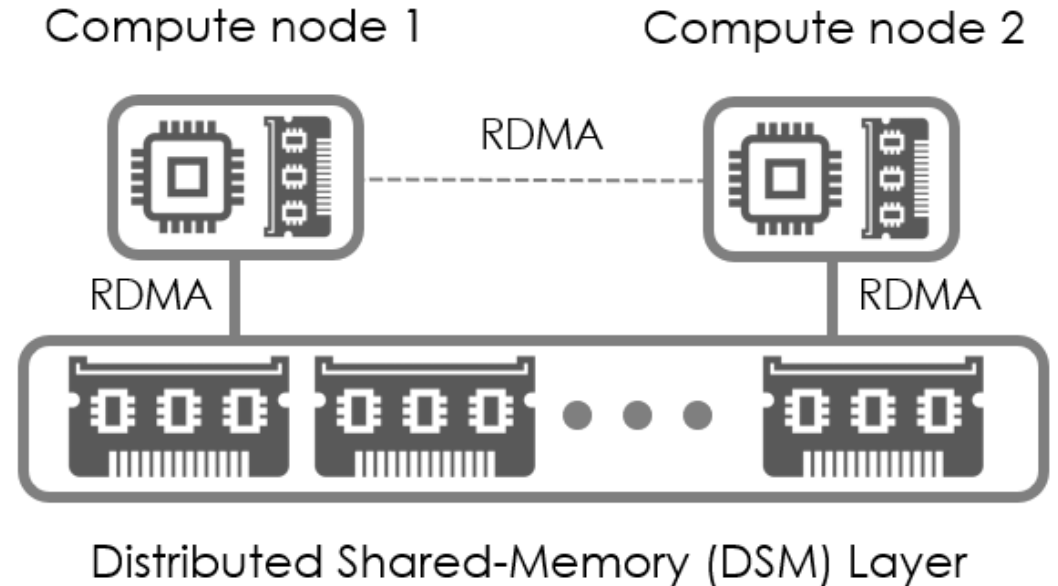
Transaction Execution Design

- Challenge: Cache coherence
 - This is not an issue if CXL is used
 - RDMA: No hardware-supported cache coherence between compute nodes



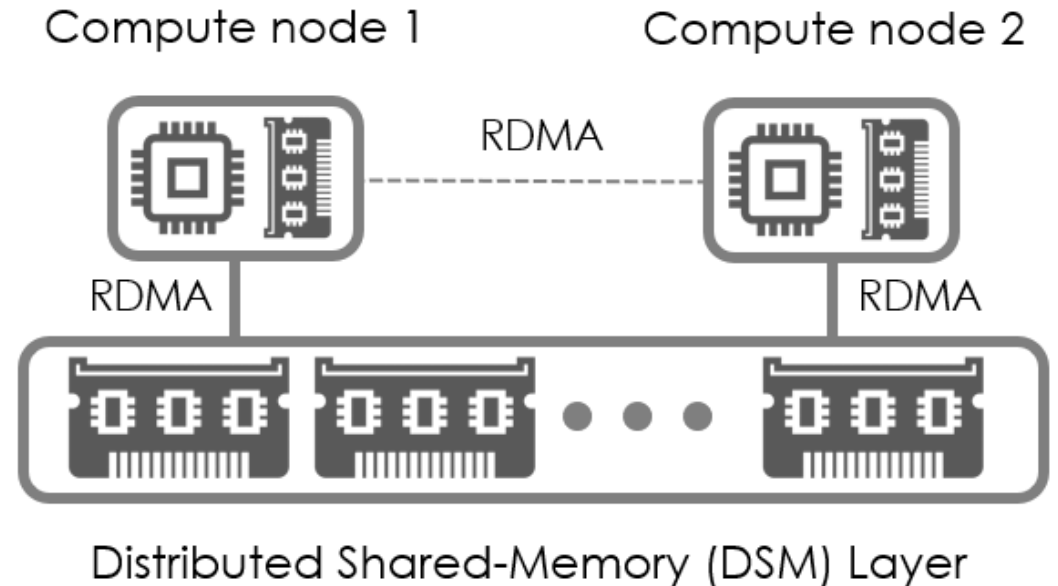
Transaction Execution Design

- Challenge: Cache coherence
 - This is not an issue if CXL is used
 - RDMA: No hardware-supported cache coherence between compute nodes
- Challenge: Distributed commit
 - Depends on the model



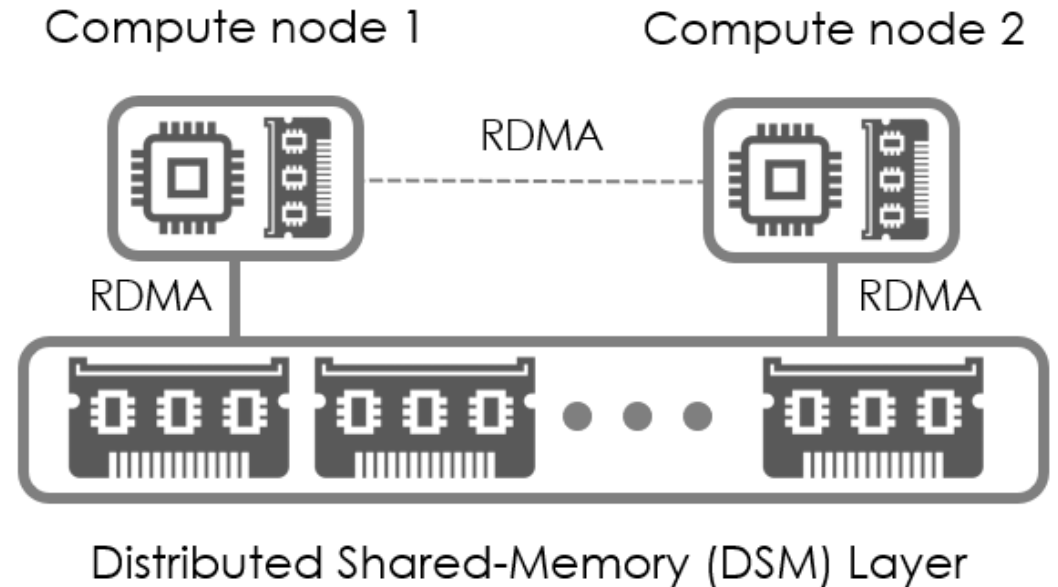
Transaction Execution Design

- Challenge: Cache coherence
 - This is not an issue if CXL is used
 - RDMA: No hardware-supported cache coherence between compute nodes
- Challenge: Distributed commit
 - Depends on the model
- Challenge: Concurrency Control Algs
 - Cannot directly use existing CC protocols, e.g., 2PL, MVCC, OCC



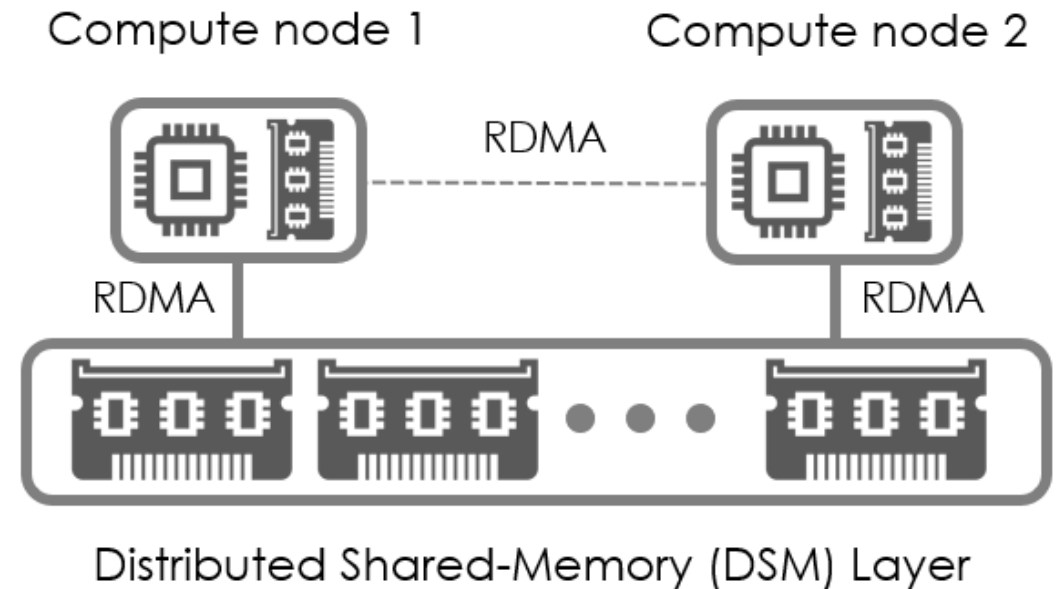
Transaction Execution Design

- Challenge: Cache coherence
 - This is not an issue if CXL is used
 - RDMA: No hardware-supported cache coherence between compute nodes
- Challenge: Distributed commit
 - Depends on the model
- Challenge: Concurrency Control Algs
 - Cannot directly use existing CC protocols, e.g., 2PL, MVCC, OCC
- Challenge: Massive concurrency
 - Main difference with multi-core single-node architecture



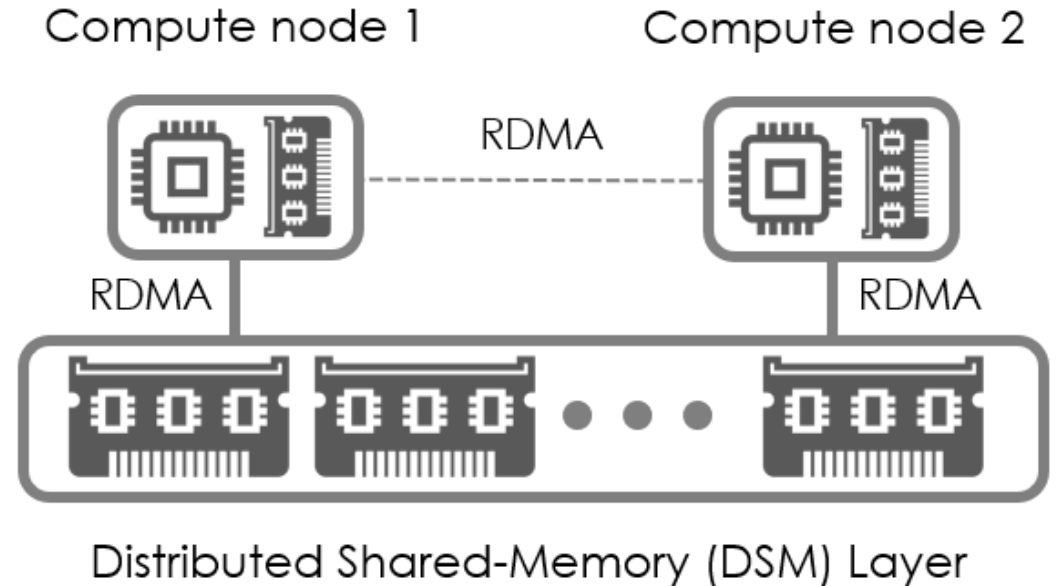
Cache Coherence

- Do we use the *local buffer* in each compute node?



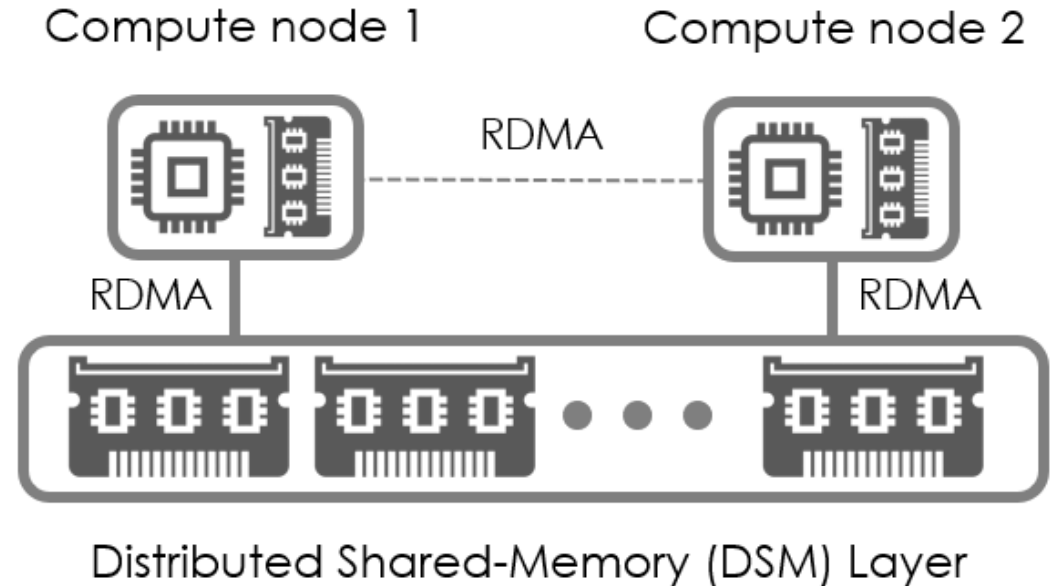
Cache Coherence

- Do we use the *local buffer* in each compute node?
 - If yes, need to address cache-coherence problem via software-level → overhead
 - If not, performance hit
 - many remote accesses



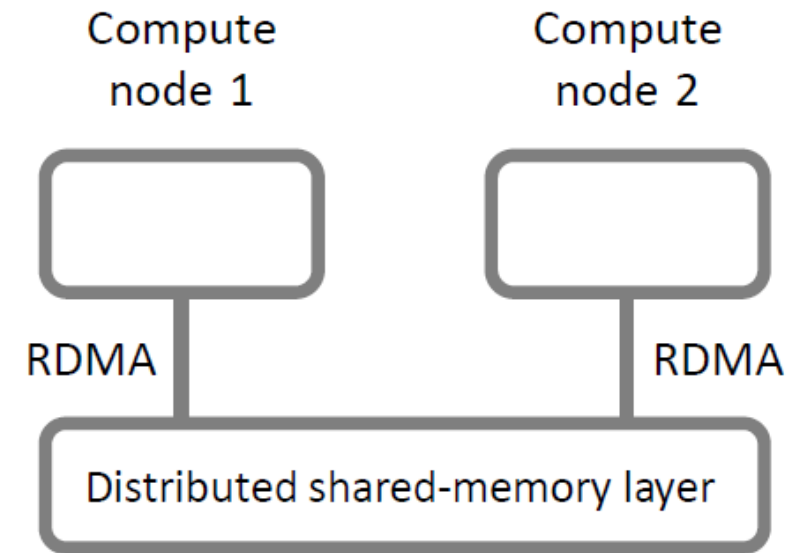
Cache Coherence

- Do we use the *local buffer* in each compute node?
 - If yes, need to address cache-coherence problem via software-level → overhead
 - If not, performance hit
 - many remote accesses
- Do we use *sharding* between compute nodes?



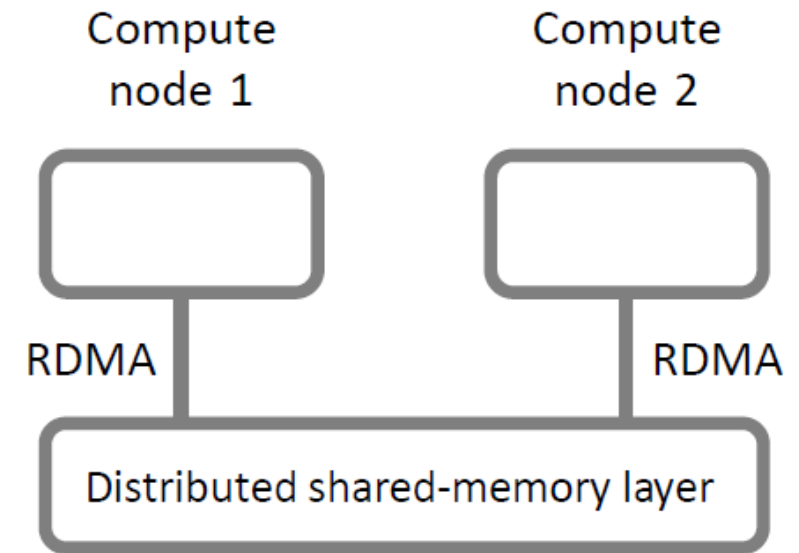
Alternative #1: No Sharding, No Cache

- Each compute node reads/writes all data
- But doesn't store any data in local buffer → No cache coherence issue
- Data pages are stored in DSM



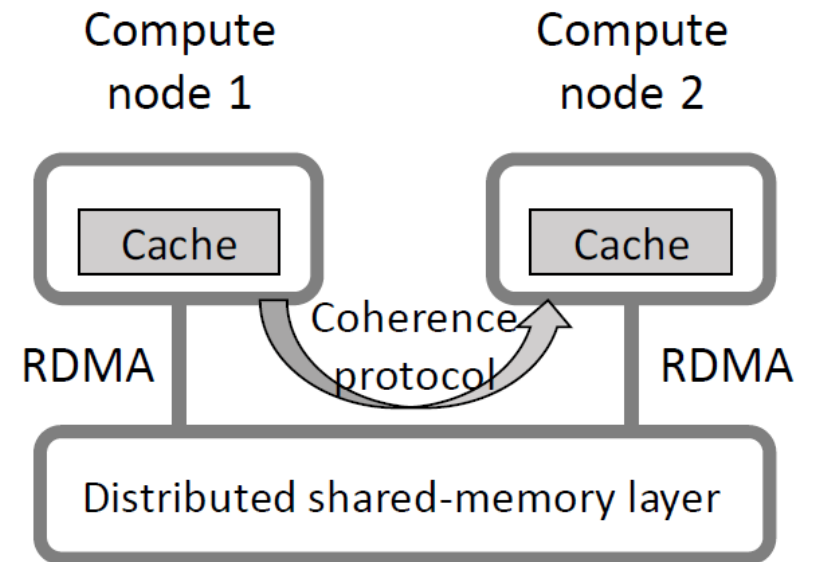
Alternative #1: No Sharding, No Cache

- Each compute node reads/writes all data
- But doesn't store any data in local buffer → No cache coherence issue
- Data pages are stored in DSM
- All compute nodes use RDMA CAS (compare & swap) to acquire a lock first
- Not realistic → many remote accesses



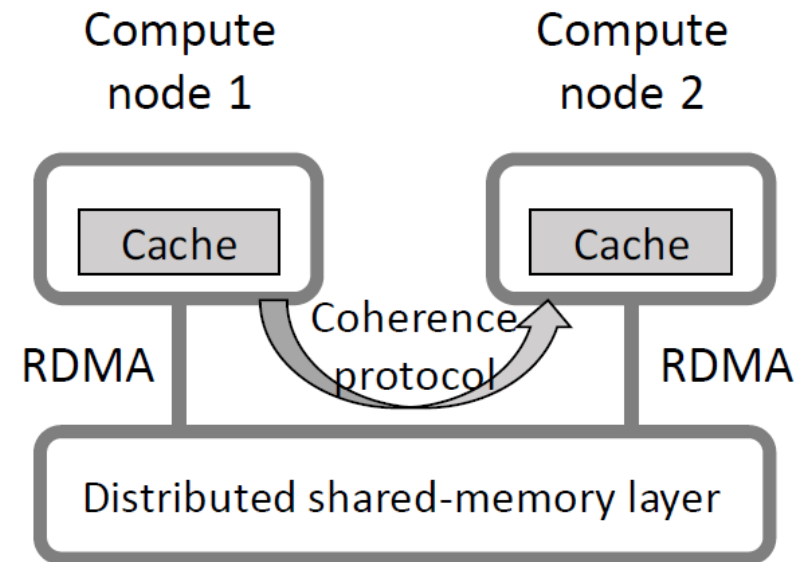
Alternative #2: No Sharding, With Cache

- Each compute node reads/writes all data
- Each compute node caches local data
- How to resolve conflicts?



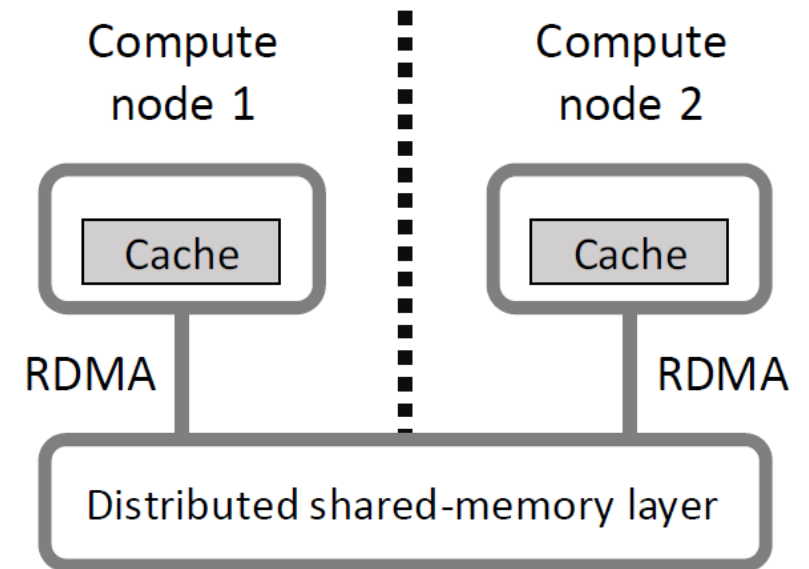
Alternative #2: No Sharding, With Cache

- Each compute node reads/writes all data
- Each compute node caches local data
- How to resolve conflicts?
- Develop a software-level cache coherence protocol
 - Update-based
 - Invalidation-based
- Overhead has to be measured and considered



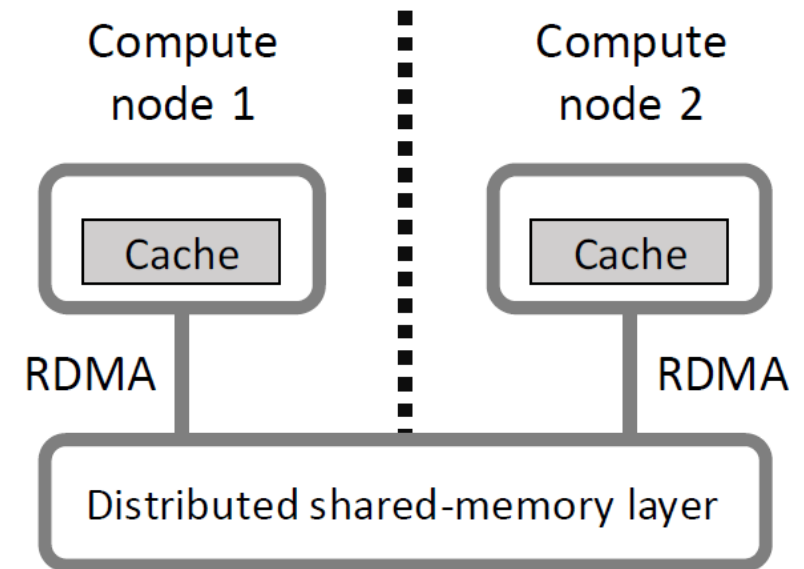
Alternative #3: With Sharding, With Cache

- Logical sharding among compute nodes
- Each compute node accesses a partition
- Use local buffer to cache data
- Bypass cache coherence issue
- But depending on workloads



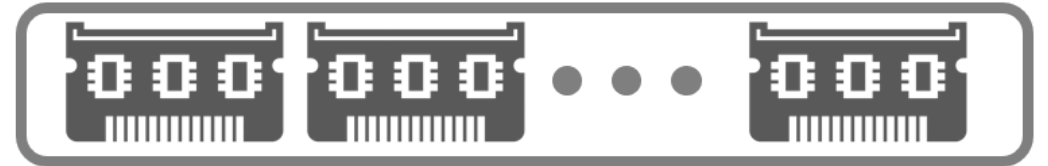
Alternative #3: With Sharding, With Cache

- Logical sharding among compute nodes
- Each compute node accesses a partition
- Use local buffer to cache data
- Bypass cache coherence issue
- But depending on workloads
- Similar to distributed shared-nothing
 - We do logical sharding (not physical sharding)
 - Also, the use of DSM layer can address distributed transaction and data skewness
 - Support elasticity very well



Distributed Shared-Memory (DSM) Design

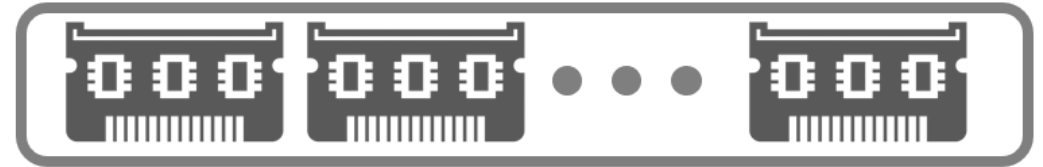
- Goal: manage a cluster of memory nodes, and expose necessary APIs to compute nodes
 - Looks like a single unified memory space



Distributed Shared-Memory (DSM) Layer

Distributed Shared-Memory (DSM) Design

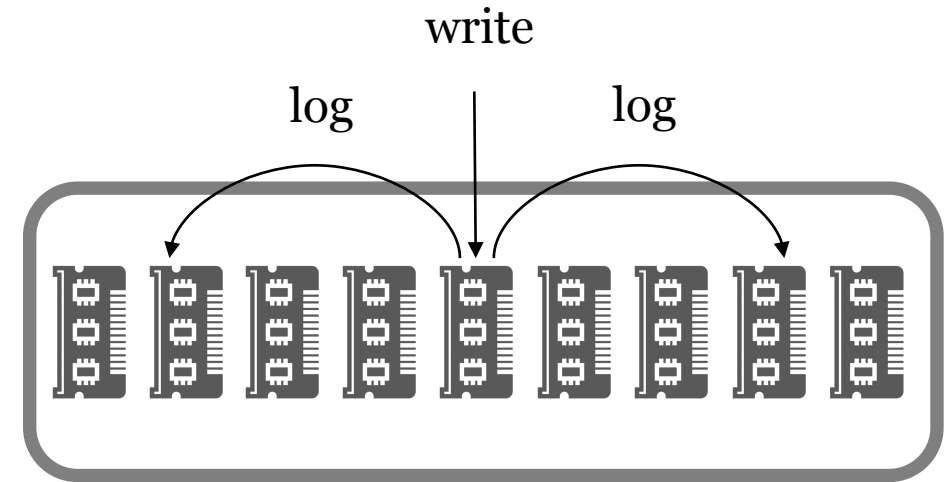
- Goal: manage a cluster of memory nodes, and expose necessary APIs to compute nodes
 - Looks like a single unified memory space
- Challenge: Durability & Availability
- Challenge: Abstract APIs
 - Memory access APIs (also memory space representation)
 - Data transmission APIs (for RDMA)
 - Function offloading APIs (for DBMS ops)



Distributed Shared-Memory (DSM) Layer

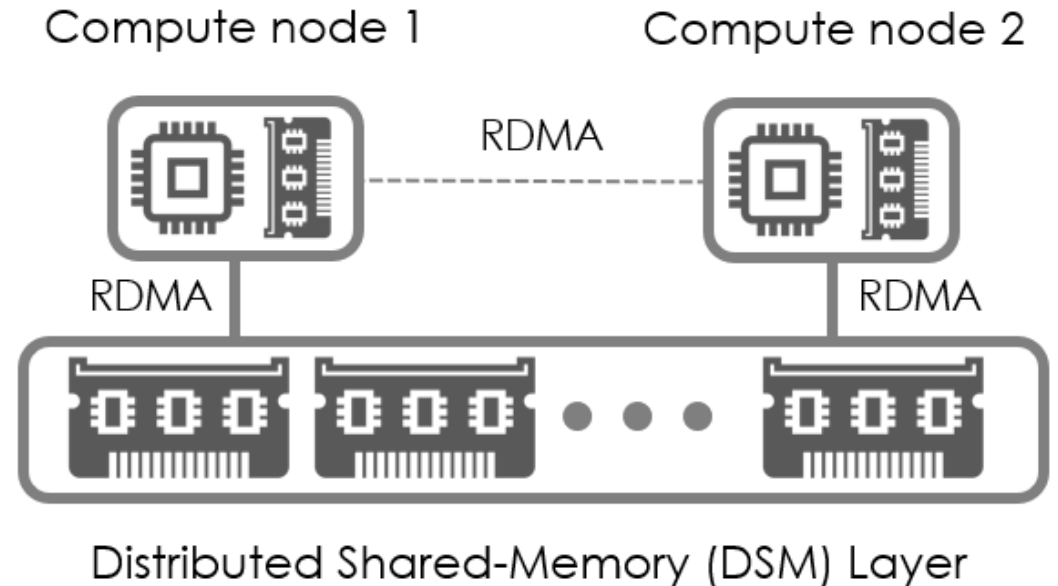
Durability & Availability

- Memory-based system
- Replication of data
 - Additional memory space + write performance
- Log-based solution (RAMCloud¹)
 - Separate logs from data (log store + data store)
 - Only logs are replicated in three memory nodes
 - A transaction is committed once all the logs are replicated and acknowledged
 - Data pages are asynchronously materialized
 - Data pages are stored only once in memory
 - Periodically checkpoint data pages to shared-storage to improve availability



Index Design Challenges

- Challenge: Leveraging RDMA characteristics (e.g., RDMA primitives)
- Challenge: How to use the buffer memory in compute nodes?
- Challenge: Leveraging RDMA byte-addressability
- Challenge: Leveraging near-data computing
- Challenge: How to support multi-node concurrent accesses?



Rethinking Existing Indexes

Disk-based indexes

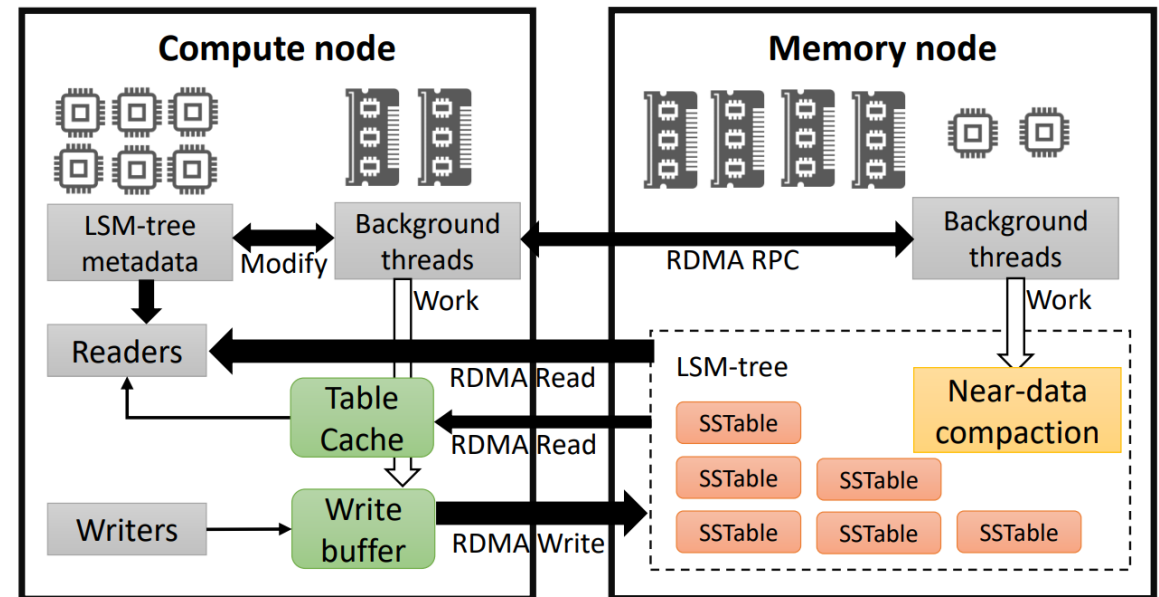
- E.g., B-tree, LSM-tree
- Not designed for RDMA (e.g., which primitive)
- Not optimized for byte-addressability
- Not leveraging near-data computing
- Multi-node concurrent accesses
 - E.g., use a single lock for both read and write or use read-write locks?

Memory-based indexes

- E.g., Bw-tree, Mass-tree, ART
- Not designed for RDMA (e.g., which primitive)
- Not leveraging local buffer
- Not leveraging near-data computing
- Multi-node concurrent accesses
 - E.g., lock-free vs. lock-based

dLSM

- Compute node: MemTable, immutable, index blocks, bloom filters
- Memory node: SSTables (memory-optimized)
- One-sided RDMA read/write
- Reduce software overhead
 - Synchronization
 - Flushing
- Near-data compaction
- Optimize for byte-addressability
 - Can access a single KV pair

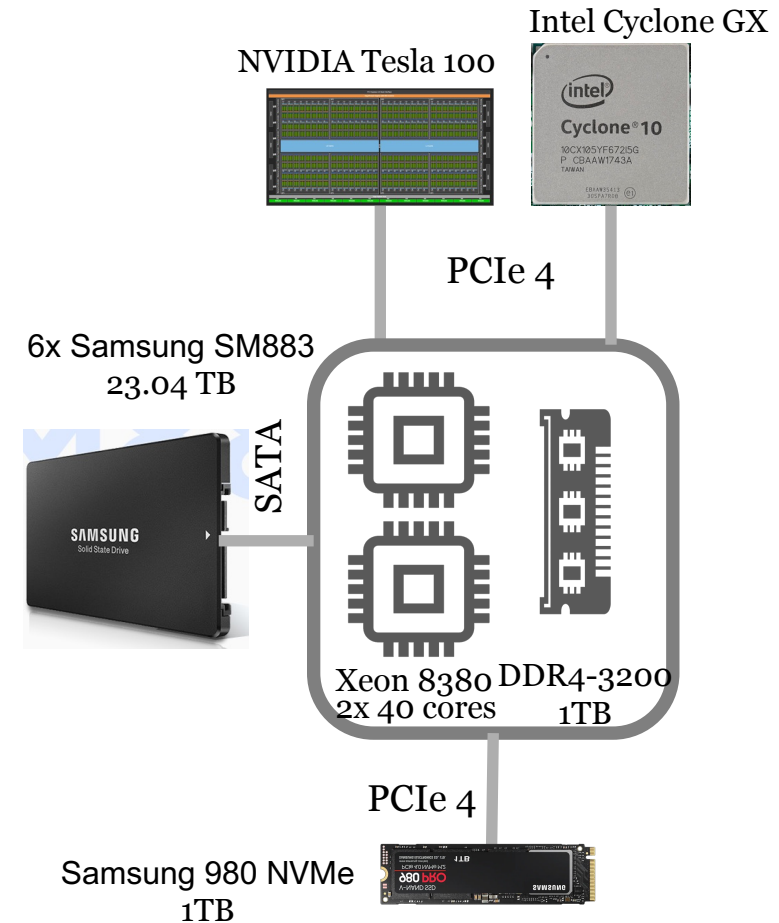


DISAGGREGATED HETEROGENEOUS PLATFORM

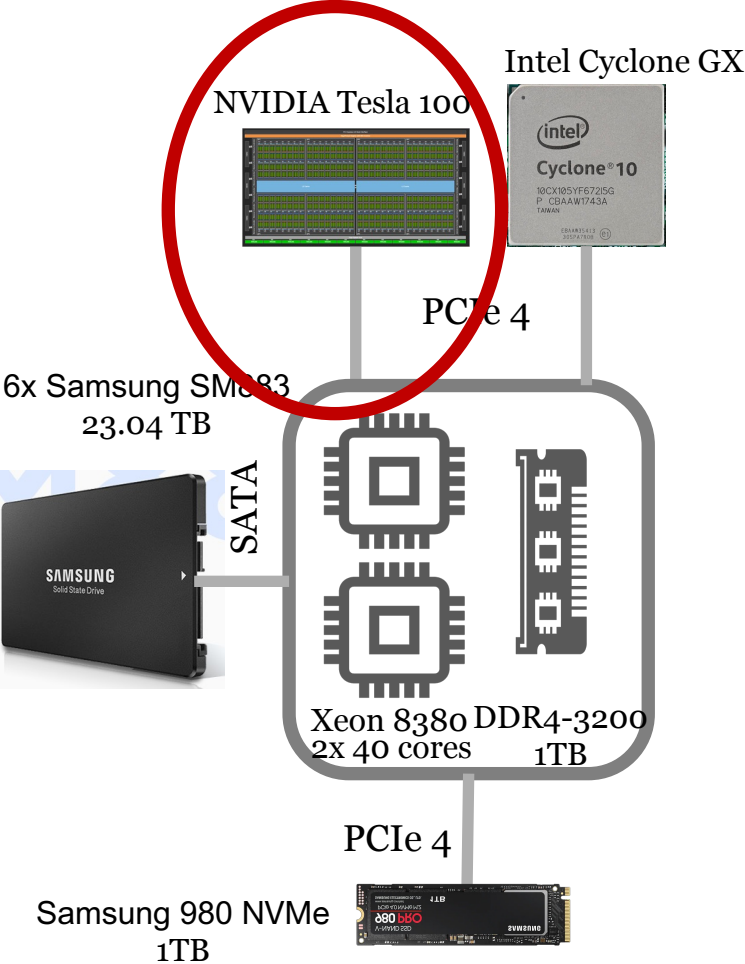
Heterogeneity

Case for Heterogeneous Computing

- Computing platforms are already heterogeneous
- Remember that CPU is not keeping up
- Data parallel computations of new workloads perform well on GPUs
- FPGAs are more flexible and can serve to experiment with solutions
- High volume demand for same task might be best met by ASICs

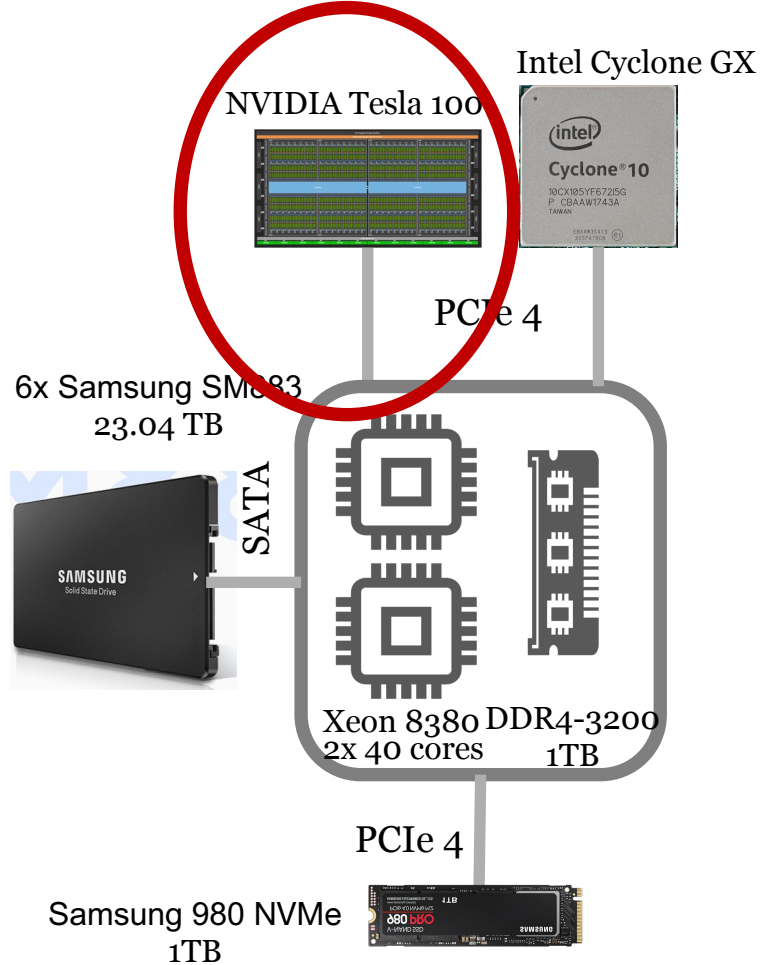


Data Management over GPU/CPU-GPU

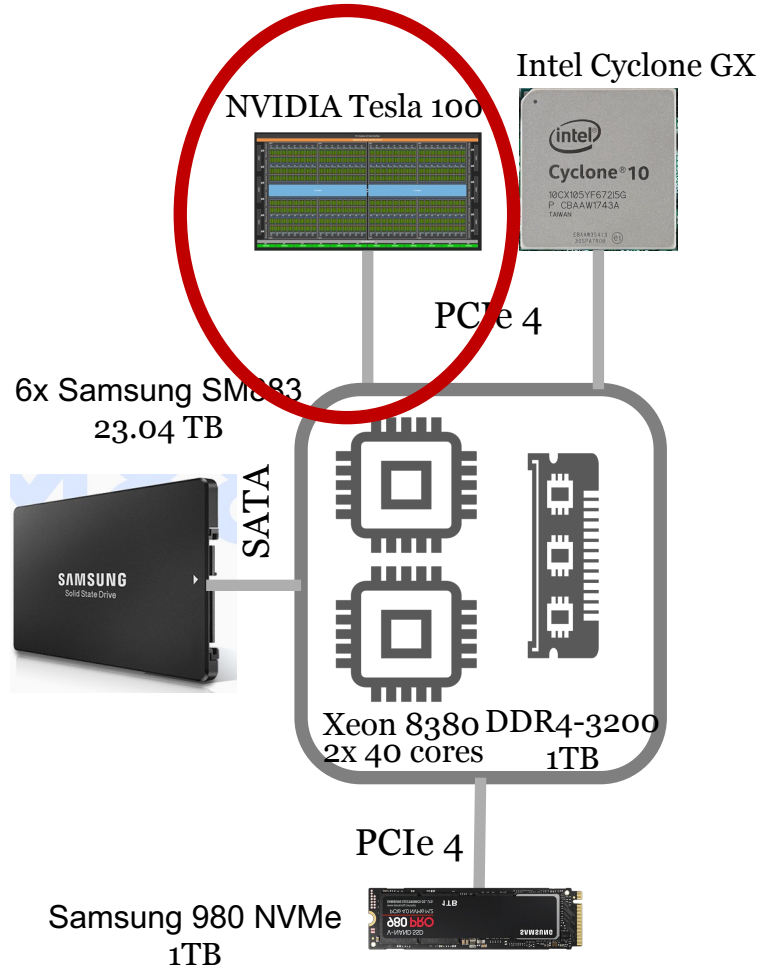


Data Management over GPU/CPU-GPU

Two approaches



Data Management over GPU/CPU-GPU



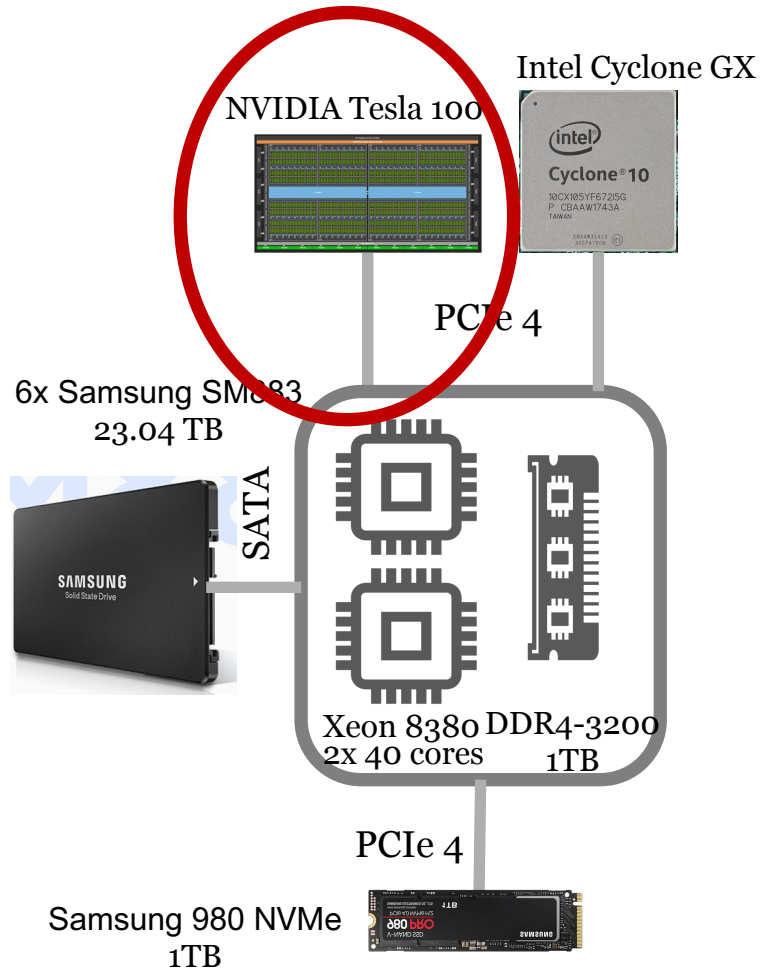
Two approaches

- Transfer data to GPU on-demand¹
 - PCIe bandwidth is an issue
 - PCIe bottleneck sometimes causes performance < optimized CPU implementations
 - Store working set directly on GPU²

¹ J. Li et al. HippogriffDB: Balancing I/O and GPU Bandwidth in Big Data Analytics, *Proc. VLDB*, 2016.

² A. Shanbhag, S. Madden and X. Yu. A Study of the Fundamental Performance Characteristics of GPUs and CPUs for Database Analytics, *Proc. SIGMOD*, 2020.

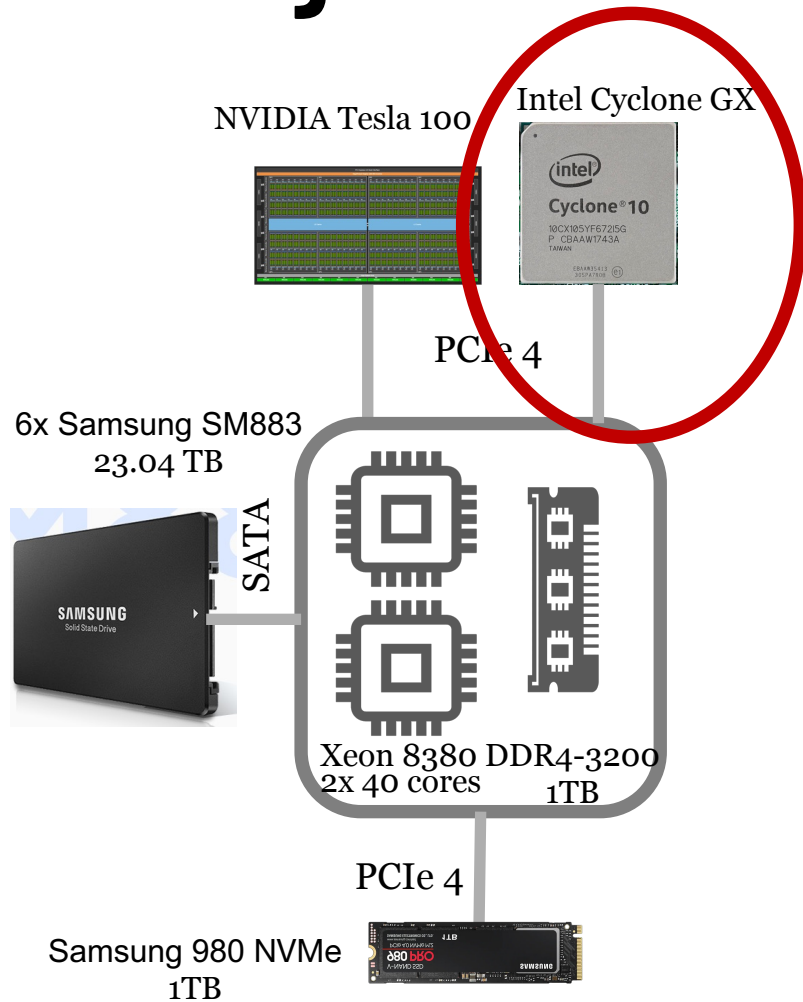
Data Management over GPU/CPU-GPU



Two approaches

- Transfer data to GPU on-demand¹
 - PCIe bandwidth is an issue
 - PCIe bottleneck sometimes causes performance < optimized CPU implementations
 - Store working set directly on GPU²
- CPU-GPU heterogeneous processing^{3,4}
 - Design more complicated
 - Data placement – treat GPU memory as cache⁴
 - Finer granularity cache → more complex query execution
 - Cost model to decide where to execute what

Data Management over FPGA/CPU-FPGA



“FPGAs ... configure different parts of the design space offering advantages that other current options do not have:

- line rate processing;
- enabling processing streams of data out of the network, disks, or memory without performance loss;
- architectural flexibility in that they can be inserted in places where the other type of processor cannot be used such as in NICs, in storage, in memory, etc. ...; and
- the customizable nature of reconfigurable computing with the FPGA serving equally well to accelerate network function virtualization, data reorganization in a database, or to accelerate joins.”

FINAL WORDS

Environment is Getting Interesting

- Other h/w developments that can help (but I don't know much about):
 - Tensor Processing Unit (TPU)
 - Data Processing Unit (DPU)
 - Non-Volatile Memory (NVM)
- Design space is varied and large
 - But start somewhere



UNIVERSITY OF
WATERLOO



Data
Systems
Group

Thank you