

Software in the sixth Epoch of Distributed Computing

Tim Mattson

Human Learning Group



Five Epochs of Distributed Computing*: the first three

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	100 ms Low bandwidth High latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server	10 ms 10 mbps	Computer to computer
Third 1990	Networking wall	MPP HPC, three-tier datacenter networks	1 ms 100 mbs → 1 Gbs	Services to services People to static data

*The five Epochs of distributed computing, Amin Vahdat of Google: SIGCOMM Lifetime achievement award keynote, 2020.
<https://www.youtube.com/watch?v=27zuReojDVw>

The Eight Fallacies of Distributed Computing

(Peter Deutsch of Sun Microsystems, 1994 ... item 8 added in 1997 by James Gosling)

Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause *big* trouble and *painful* learning experiences.

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous

The Eight Fallacies of Distributed Computing

(Peter Deutsch of Sun Microsystems, 1994 ... item 8 added in 1997 by James Gosling)

Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause *big* trouble and *painful* learning experiences.

1. The network is reliable
2. Latency is **low and fixed**
3. Bandwidth is **high and fixed**
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is **negligible**
8. The network is homogeneous

The Eight Fallacies of Distributed Computing

(Peter Deutsch of Sun Microsystems, 1994 ... item 8 added in 1997 by James Gosling)

Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause *big* trouble and *painful* learning experiences.

Data Center

- ~~X~~ ~~The network is reliable~~
- ~~X~~ ~~Latency is low and fixed~~
- ~~X~~ ~~Bandwidth is high and fixed~~
- ~~X~~ ~~The network is secure~~
- ~~X~~ ~~Topology doesn't change~~
- ~~X~~ ~~There is one administrator~~
- ~~X~~ ~~Transport cost is negligible~~
- ~~X~~ ~~The network is homogeneous~~

HPC Cluster

- ✓1. The network is reliable
- ✓2. Latency is low and fixed
- ✓3. Bandwidth is high and fixed
- ✓4. The network is secure
- ✓5. Topology doesn't change
- ✓6. There is one administrator
- ~~X~~ ~~Transport cost is negligible~~
- ✓8. The network is homogeneous

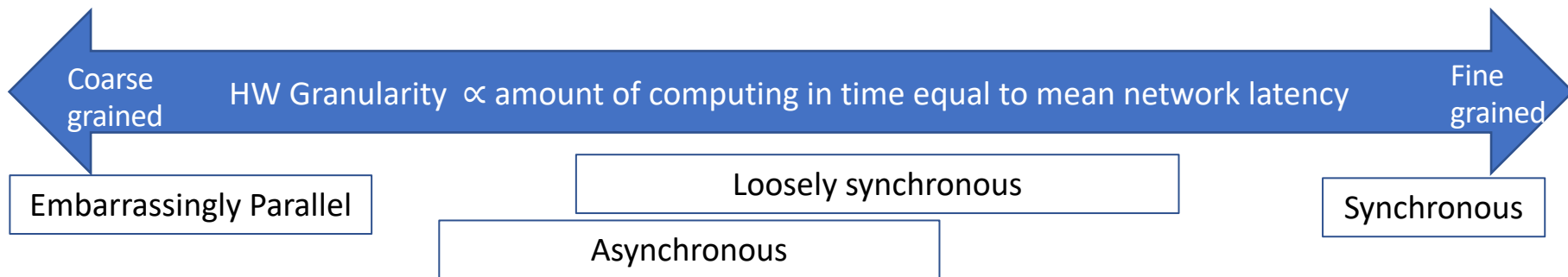
Where can we run HPC applications?

Data Center

- ~~✗~~ ~~The network is reliable~~
- ~~✗~~ ~~Latency is low and fixed~~
- ~~✗~~ ~~Bandwidth is high and fixed~~
- ~~✗~~ ~~The network is secure~~
- ~~✗~~ ~~Topology doesn't change~~
- ~~✗~~ ~~There is one administrator~~
- ~~✗~~ ~~Transport cost is negligible~~
- ~~✗~~ ~~The network is homogeneous~~

HPC Cluster

- ✓1. The network is reliable
- ✓2. Latency is low and fixed
- ✓3. Bandwidth is high and fixed
- ✓4. The network is secure
- ✓5. Topology doesn't change
- ✓6. There is one administrator
- ~~✗~~ ~~Transport cost is negligible~~
- ✓8. The network is homogeneous



Five Epochs of Distributed Computing*

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	100 ms Low bandwidth High latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server	10 ms 10 mbps	Computer to computer
Third 1990	Networking wall	MPP HPC, three-tier datacenter networks	1 ms 100 mbs → 1 Gbs	Services to services People to static data
Fourth 2000	Dennard scaling wall ... per core plateau	Web search, planet-scale services	100 μ s 10 Gbps flash	People to people People to interactive results

*The five Epochs of distributed computing, Amin Vahdat of Google: SIGCOMM Lifetime achievement award keynote, 2020.

FTP: File Transfer Protocol,

MPP: Massively parallel processor,

RPC: Remote Procedure Call

The cloud takes over

- On premises data centers give way to distributed resources “in the cloud”.
- Starts with virtual machines (infrastructure as a service) and grows into a whole new architecture for software ... microservices.

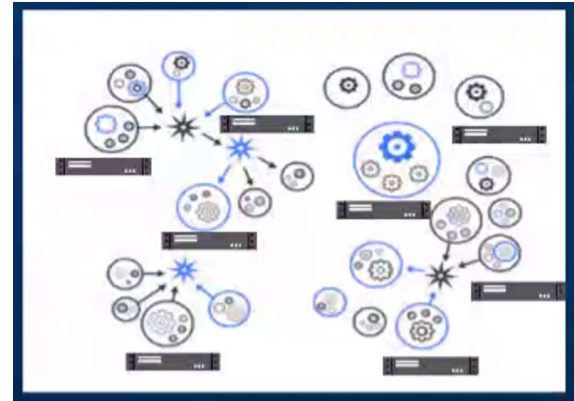
The old way: Monolithic applications



Single program composed of many functions interacting through memory/messages.

Latencies $O(\text{microseconds})$ or less

The new, cool way: microservices

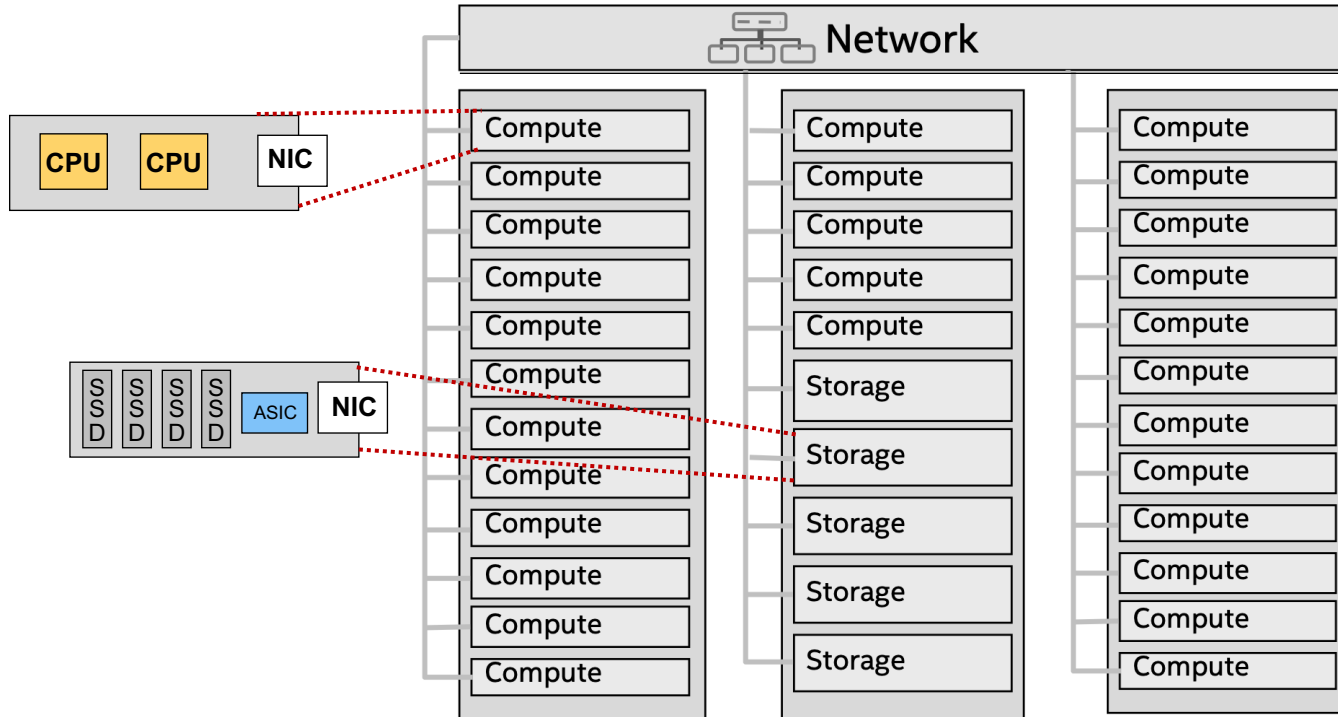


Many small independent programs interacting through remote procedure calls (RPC).

Latencies $O(\text{milliseconds})$ to $O(\text{seconds})$...
the tail latency problem

Hardware in a cloud data center

- The unit of hardware replication is the server node with memory, I/O, and CPUs.

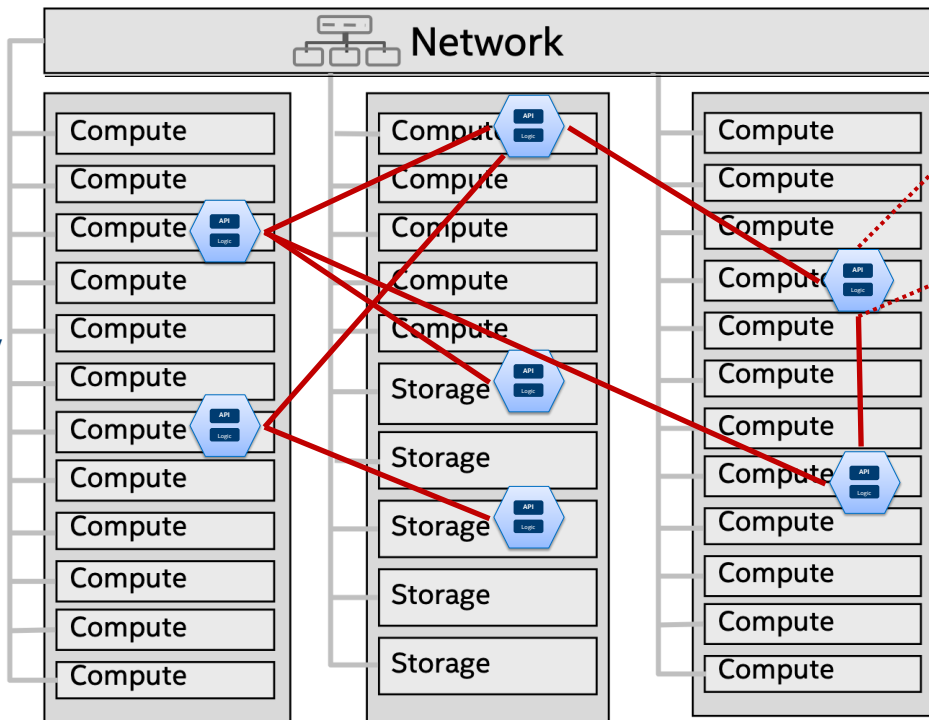


Hardware in a cloud data center

- Applications composed of many ($O(100)$) small independent programs ... ship functions as needed. They interact through remote procedure calls (RPCs).

Apps@scale as collections of microservices

- Discrete units of functionality
- Continuous integration
- Continuous upgrades
- Resilient workflows



A Program embedded inside a microservice

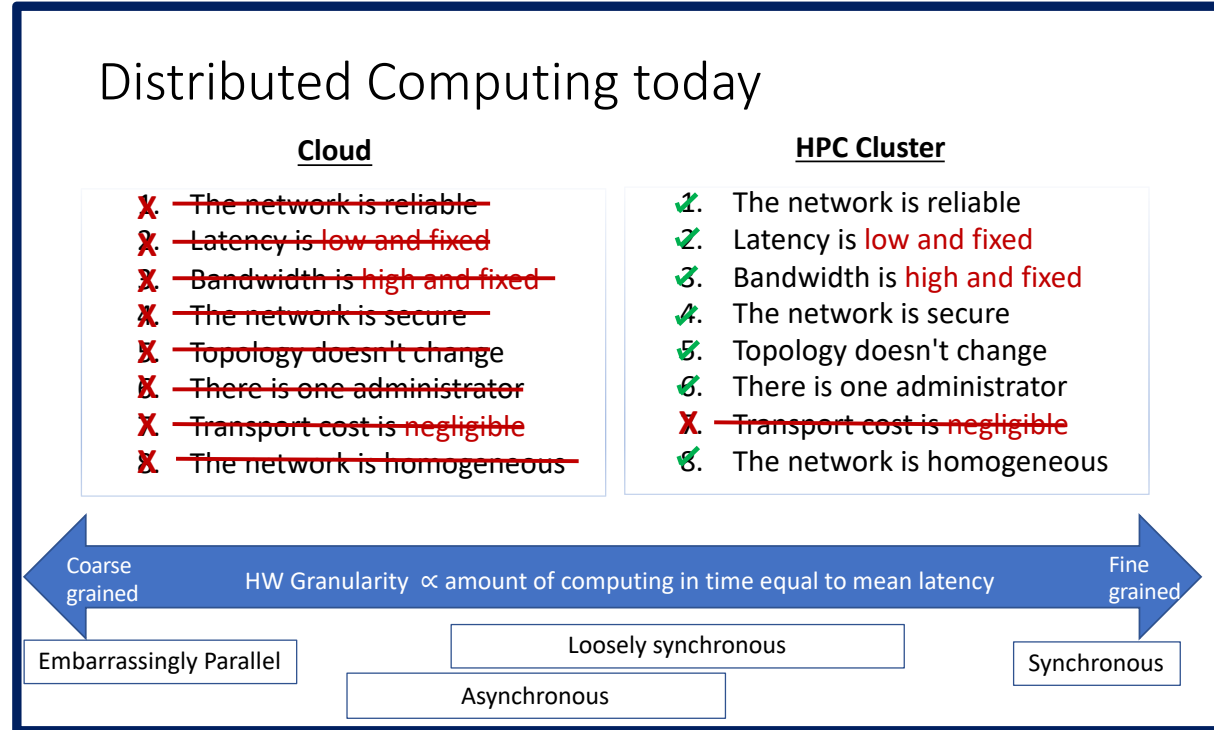
- Tightly coupled parallelism.
- Interacts directly with the Hardware

Programming Distributed computers

There is a clean split between applications that run in the cloud and those that need a dedicated HPC cluster.

This is reflected in the programming models used:

- Cloud: Remote Procedure Call (RPC), distributed object store distinct from tasks, execution flows as task graphs for Function as a service. Heavy use of microservices.
- HPC Cluster: SPMD design pattern with MPI ... also PGAS with SHMEM.



The three domains of parallel programming

Platform*	Laptop or server	HPC Cluster	Cloud
Execution Agent	Threads	Processes	Microservices
Memory	Single Address Space	Distributed memory, local memory owned by individual processes	Distributed object store (in memory) backed by a persistent storage system
Typical Execution Pattern	Fork-join	SPMD	Event driven tasks, FaaS, and Actors

Laptop/server and cluster models work well together.

An impenetrable wall separates them from the cloud-native world

Five Epochs of Distributed Computing*

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	100 ms Low bandwidth High latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server	10 ms 10 mbps	Computer to computer
Third 1990	Networking wall	MPP HPC, three-tier datacenter networks	1 ms 100 mbs → 1 Gbs	Services to services People to static data
Fourth 2000	Dennard scaling wall ... per core plateau	Web search, planet-scale services	100 μ s 10 Gbps flash	People to people People to interactive results
Fifth 2015	Per socket wall ... accelerators take off	Machine Learning, data centric computing	10 μ s 200 Gbps → 1 Tbps	People to insights

*The five Epochs of distributed computing, Amin Vahdat of Google: SIGCOMM Lifetime achievement award keynote, 2020.

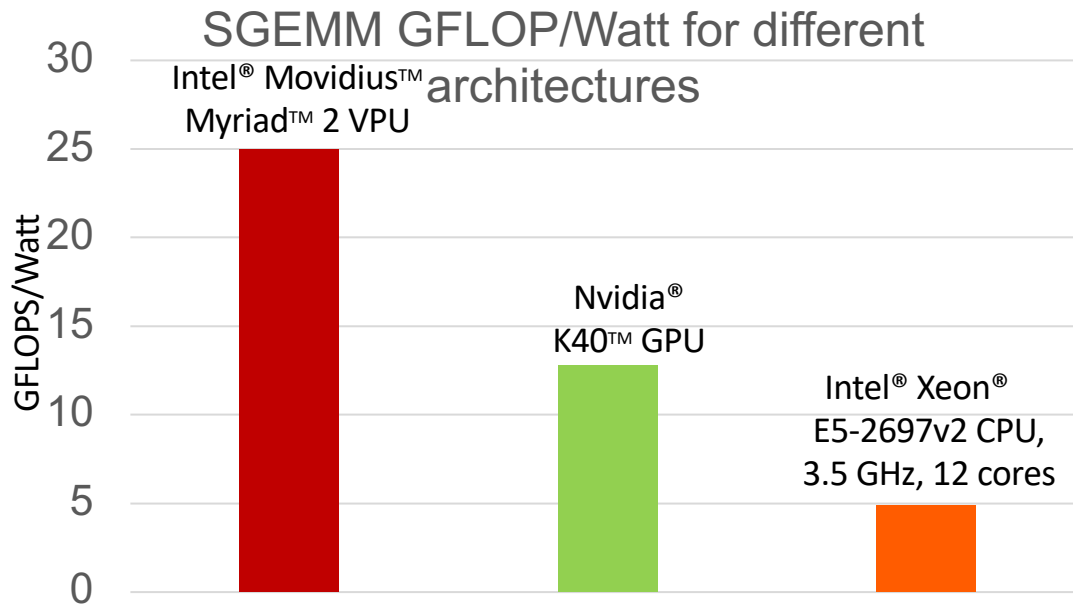
FTP: File Transfer Protocol,

MPP: Massively parallel processor,

RPC: Remote Procedure Call

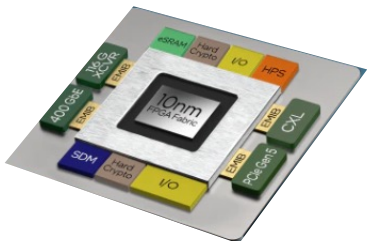
If you care about power, the world is heterogeneous?

Specialized processors doing operations suited to their architecture are more efficient than general purpose processors.



Hence, future systems will be increasingly heterogeneous ... GPUs, CPUs, FPGAs, and a wide range of accelerators

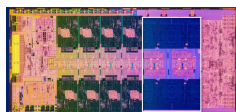
A New Golden Age for Computer Architecture



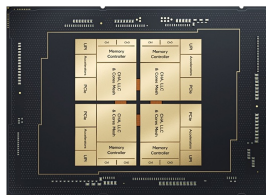
Intel® Agilex™ FGAs



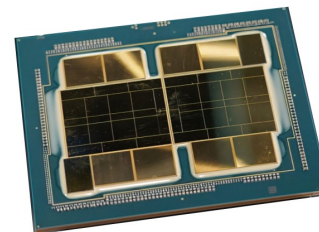
Habana® Gaudi® 2
deep learning accelerator



13th gen Intel® Core™ CPU
Hybrid Architecture with 16 efficiency
cores and 8 performance cores +
integrated GPU



4th Gen Intel® Xeon® CPU with 56
cores and Novel On-Die
Accelerators

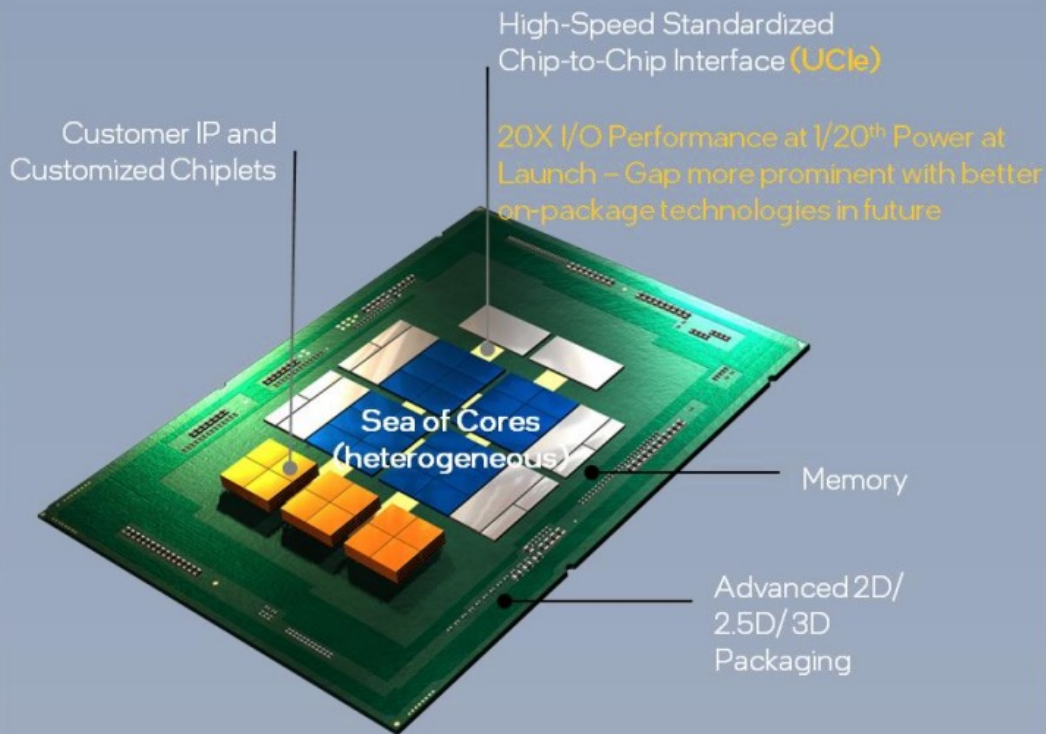


Intel® Xe^e HPC GPU

A New Golden Age for Computer Architecture

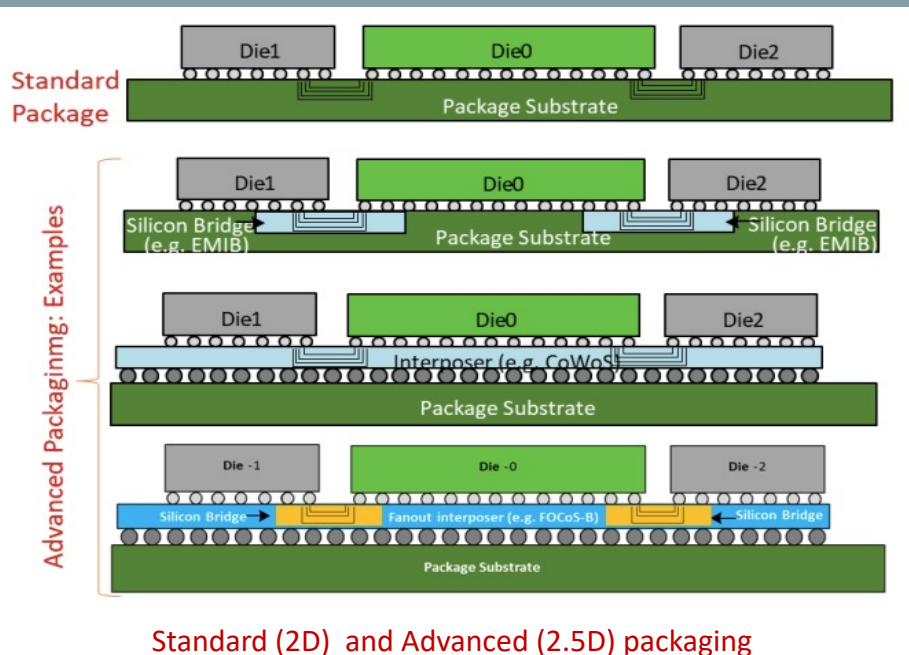


Open Chiplet: Platform on a Package



- Chiplet-based architectures ... building a package placed in a socket composed of distinct little chips (the “chiplet”).
- Connected by high speed in package interconnects ... lets chiplets from multiple fabs fit into one package.
- The Universal Chiplet Interconnect Express effort defines a standard for how to connect chiplets.
- The result ... multi-chiplet packages in a socket with heterogeneous devices from multiple vendors.

A New Golden Age for Computer Architecture



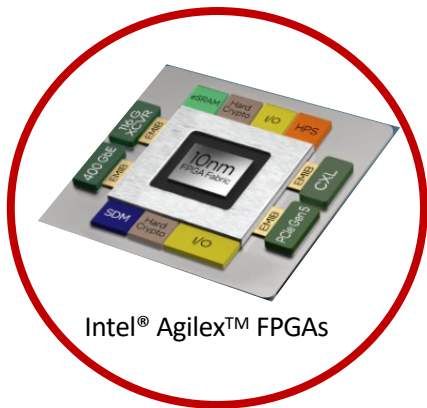
- Chiplet-based architectures ... building a package placed in a socket composed of distinct little chips (the “chiplet”).
- Connected by high speed in package interconnects ... lets chiplets from multiple fabs fit into one package.
- The Universal Chiplet Interconnect Express effort defines a standard for how to connect chiplets.
- The result ... multi-chiplet packages in a socket with heterogeneous devices from multiple vendors.

Membership includes Intel, Nvidia, Samsung, Qualcomm, Alibaba, Google, Microsoft, Meta, Arm ...

A New Golden Age for Computer Architecture



Multi-chiplet packages with components from multiple vendors

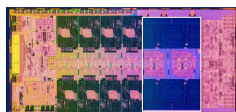


Intel® Agilex™ FPGAs

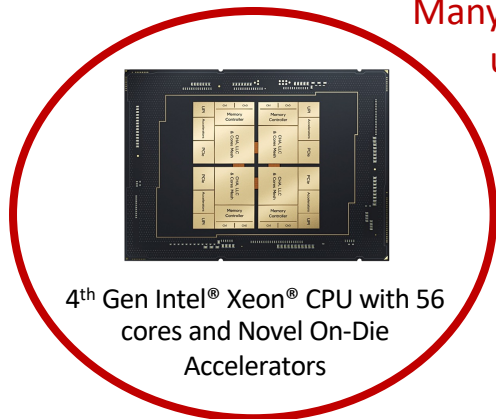


Habana® Gaudi® 2 deep learning accelerator

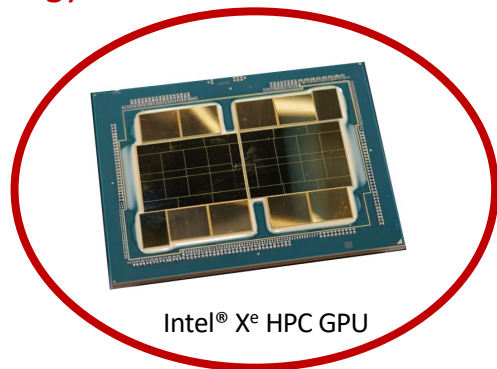
Many of Intel's products already use chiplet technology



13th gen Intel® Core™ CPU Hybrid Architecture with 16 efficiency cores and 8 performance cores + integrated GPU



4th Gen Intel® Xeon® CPU with 56 cores and Novel On-Die Accelerators

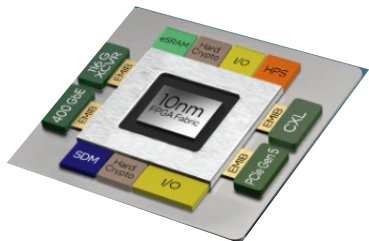


Intel® Xe® HPC GPU

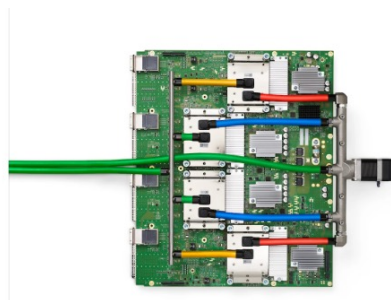
A New Golden Age for Computer Architecture



Multi-chiplet packages with components from multiple vendors



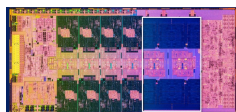
Intel® Agilex™ FPGAs



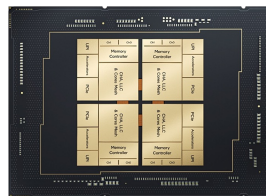
Google® Tensor Processing Unit



Habana® Gaudi® 2 deep learning accelerator



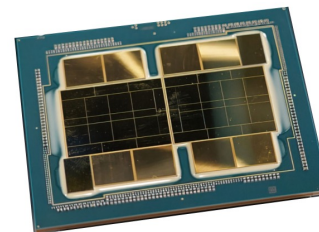
13th gen Intel® Core™ CPU Hybrid Architecture with 16 efficiency cores and 8 performance cores + integrated GPU



4th Gen Intel® Xeon® CPU with 56 cores and Novel On-Die Accelerators



- CPUs
- Discrete GPU
- CPU + integrated GPU

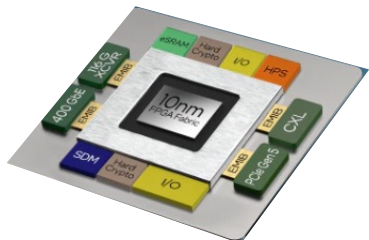


Intel® Xe HPC GPU

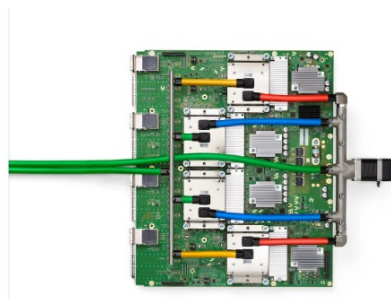
A New **Dark** Age for Computer **Programmers**



Multi-chiplet packages with components from multiple vendors



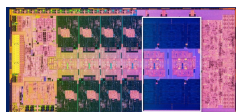
Intel® Agilex™ FPGAs



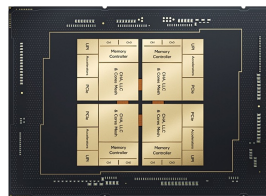
Google® Tensor Processing Unit



Habana® Gaudi® 2 deep learning accelerator



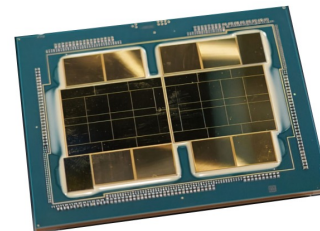
13th gen Intel® Core™ CPU Hybrid Architecture with 16 efficiency cores and 8 performance cores + integrated GPU



4th Gen Intel® Xeon® CPU with 56 cores and Novel On-Die Accelerators



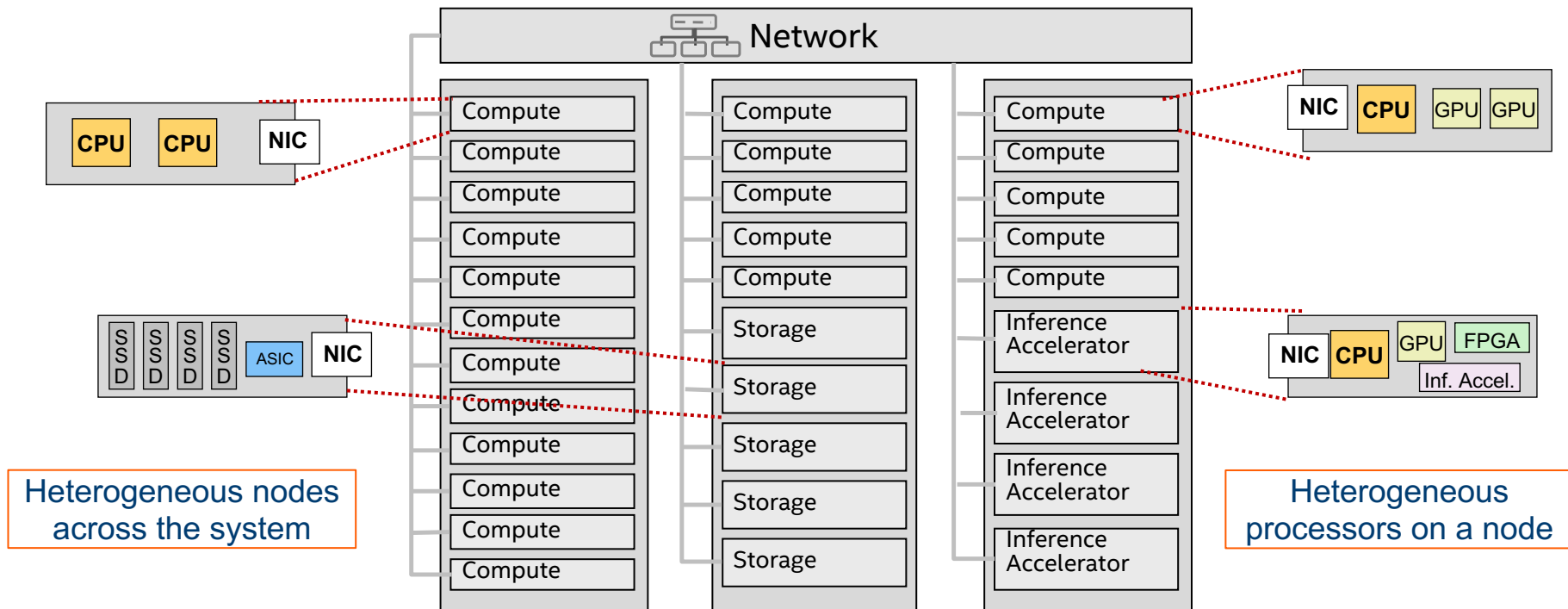
- CPUs
- Discrete GPU
- CPU + integrated GPU



Intel® Xe HPC GPU

Hardware in a cloud data center

- The unit of hardware replication is the server node with memory, I/O, processors, and accelerators.

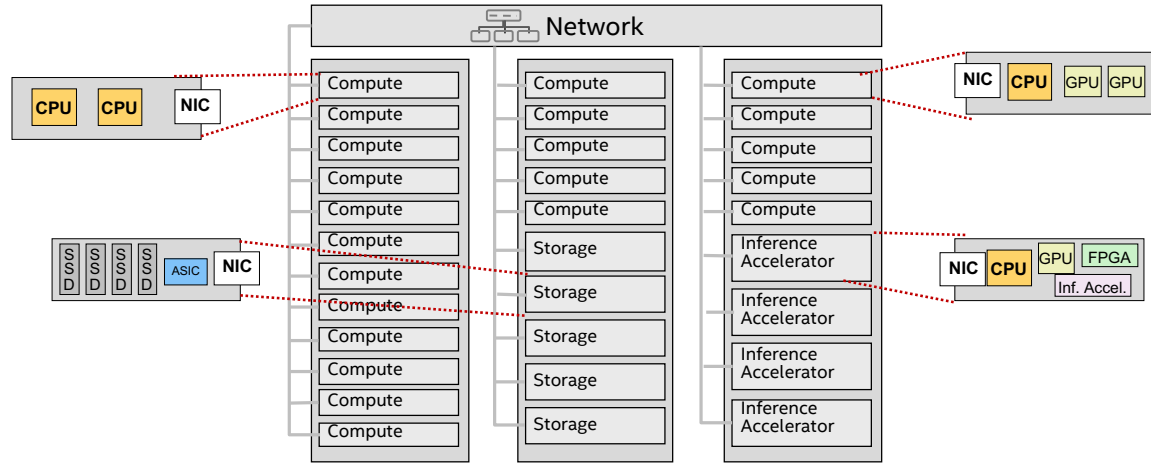


**... But there are problems in our
cloudy paradise**

Wasted Resources in the cloud

How resources are made available is associated with hardware on a node.

For example, if all the cores are allocated to customer VMs, memory on the node not used by the VMs is unavailable for other jobs. It is **stranded**.



Analysis of production traces from Azure, shows that stranding is the dominant source of memory waste. They found that up to 25% of DRAM at any time is stranded (i.e. wasted).*

A similar analysis[§] of traces from Google's data centers found the average DRAM utilization was only 40%.

* "Pond: CXL based memory pooling system for cloud platforms" Huaicheng Lie, et. al. <https://arxiv.org/abs/2203.00241>

§ "Borg: the Next Generation", Muhammad Tirmazi, et. al., <https://storage.googleapis.com/pub-tools-public-publication-data/pdf/5bf4ebfb98ead7f6ee1552860fab88e75a5ed7e.pdf>

Cost break down for a typical server node

For Azure, DRAM can be 50% of server cost*

Other analyses finds it at 40% of the cost

Whether its 40% or 50% of total cost, it is still an expensive resource to waste

* "Pond: CXL based memory pooling system for cloud platforms"
Huaicheng Lie, et. al. <https://arxiv.org/abs/2203.00241>

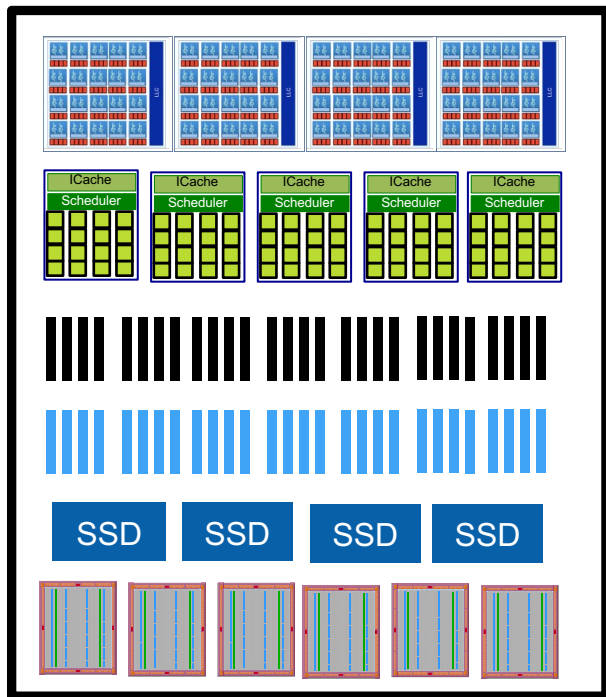
2x Intel Sapphire Rapids Server	
Component	AI Server
CPU	\$ 1,850
8 GPU + 4 NVSwitch Baseboard	\$ -
Memory	\$ 3,930
Storage	\$ 1,536
SmartNIC	\$ 654
Chassis (Case, backplanes, cabling)	\$ 395
Motherboard	\$ 350
Cooling (Heatsinks+fans)	\$ 275
Power Supply	\$ 300
Assembly and Test	\$ 495
Markup	689
Total Cost	\$ 10,474
DRAM BOM %	37.5%
NAND BOM %	14.7%
Memory BOM %	52.2%

*Dylan Patel and Gerald Wong, May 29, 2023.
<https://www.semianalysis.com/p/ai-server-cost-analysis-memory-is>

Disaggregated Computing for SW Defined Servers (SDS)

Consider a Rack composed of multiple pools

Dynamically compose across pools to match a software defined server to the workload



CPU pool

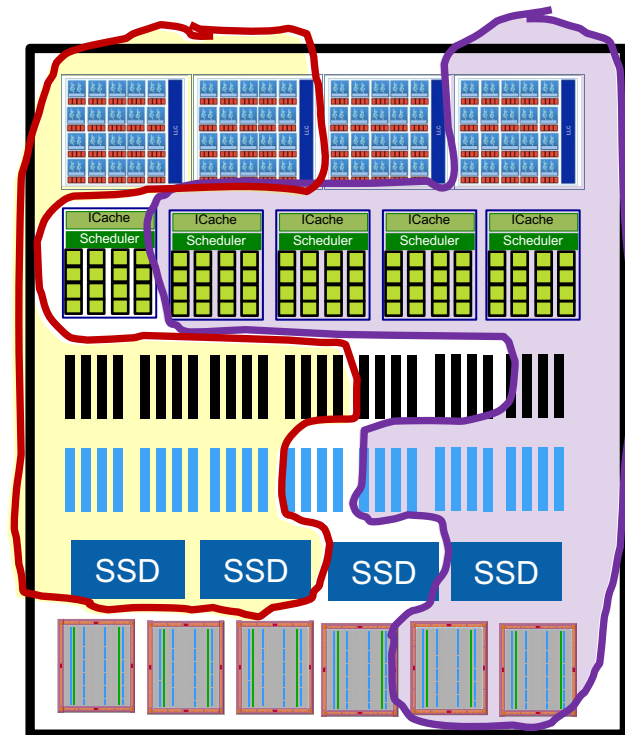
GPU pool

DRAM pool

NVRAM pool

SSD pool

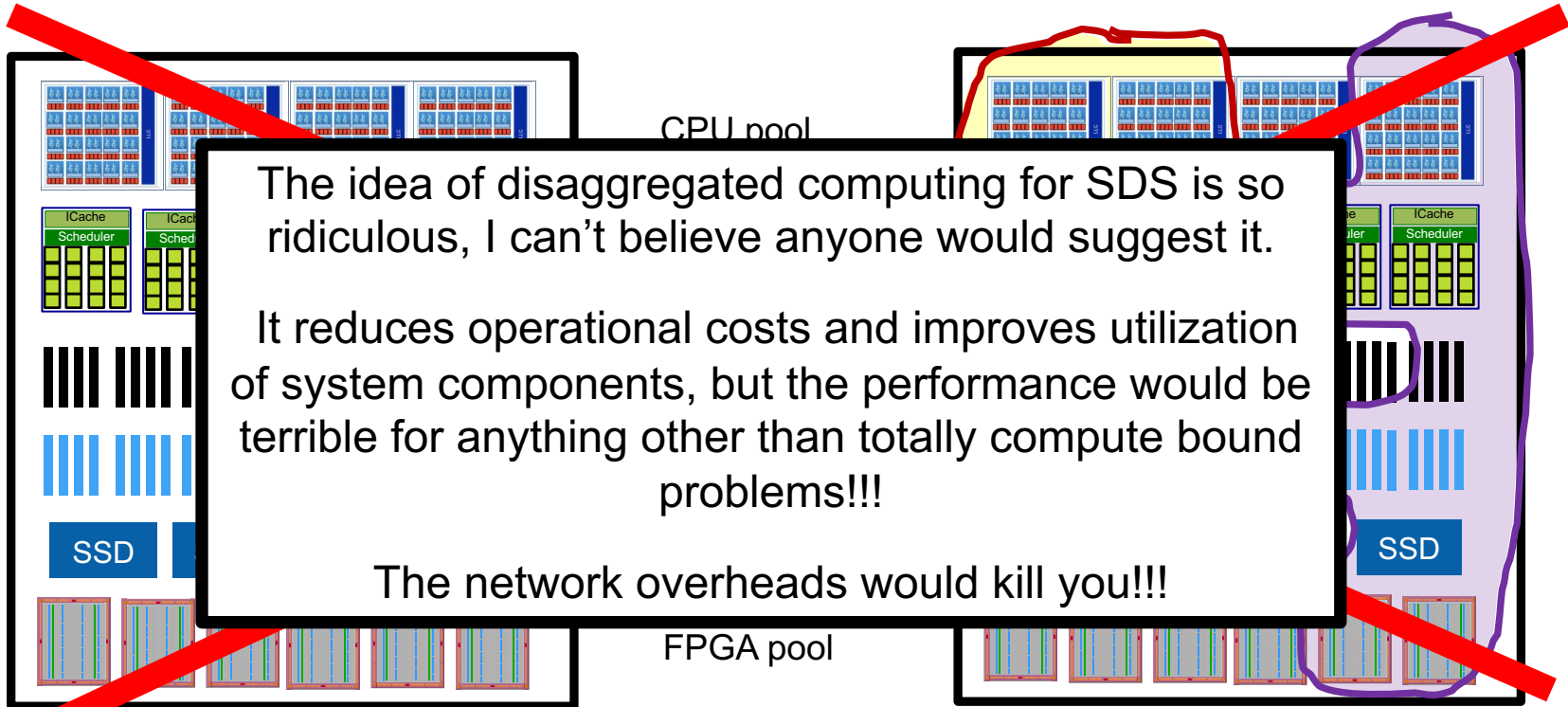
FPGA pool



Disaggregated Computing for SW Defined Servers (SDS)

Consider a Rack composed of multiple pools

Dynamically compose across pools to match a software defined server to the workload

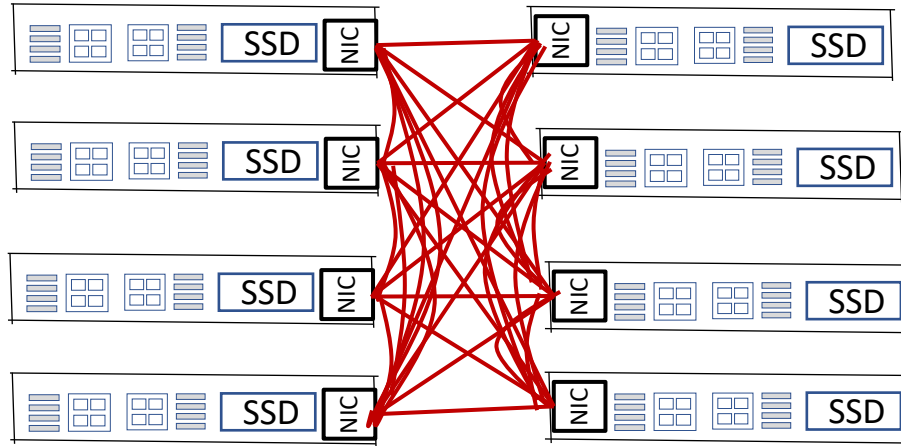


The sixth Epoch of Distributed Computing

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	100 ms Low bandwidth High latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server	10 ms 10 mbps	Computer to computer
Third 1990	Networking wall	MPP HPC, three-tier datacenter networks	1 ms 100 mbs → 1 Gbs	Services to services People to static data
Fourth 2000	Dennard scaling wall ... per core plateau	Web search, planet-scale services	100 μ s 10 Gbps flash	People to people People to interactive results
Fifth 2015	Per socket wall ... accelerators take off	Machine Learning, data centric computing	10 μ s 200 Gbps → 1 Tbps	People to insights
Sixth 2025	Speed of light	Dynamic, real-time AI from data-center to the edge with SDE*	100 ns 10 Tbs	People to experiences

* SDE: Software defined Everything, i.e. software defined networking, software defined infrastructure, software defined servers ... All at the same time ... to dynamically construct systems to meet the needs of workloads.

Networking technology... replace generic data center network with a cluster of cliques



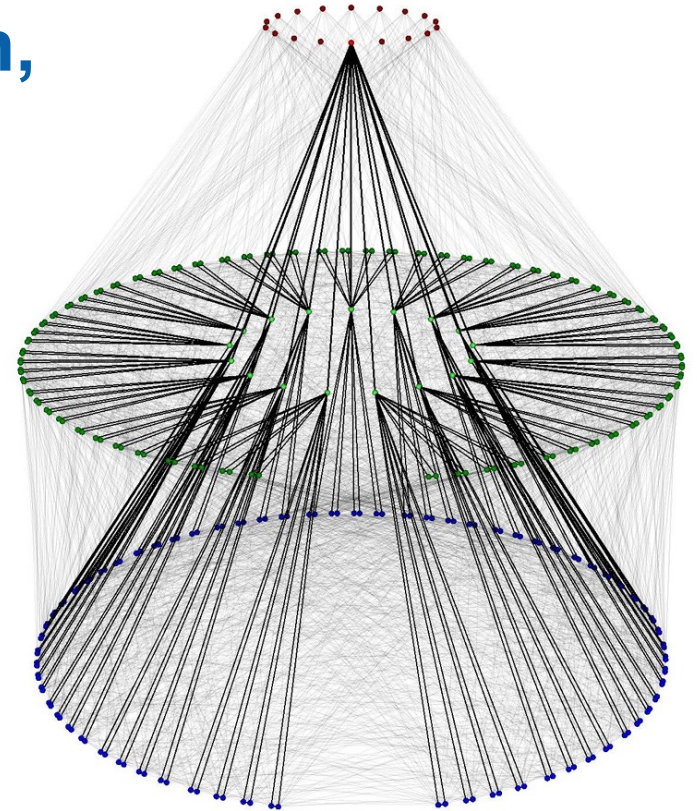
A clique: A graph where every vertex is connected to every other vertex

A Clique: a network of diameter one with $O(\frac{1}{4}N^2)$ bisection bandwidth

Combine with next generation optical networks to hit latencies close to DRAM latencies (100 ns)

Future: scale networks across the data center with next generation, point-to-point optical networks

PolarFly network: Over 150 nodes connect by a network of diameter 2 ... that is, any pair of nodes can be reached by two hops on the network



Latencies every engineer should know ...

L1 cache reference 1.5 ns

L2 cache reference 5 ns

Branch misprediction 6 ns

Uncontended mutex lock/unlock 20 ns

L3 cache reference 25 ns

Main memory reference 100 ns

“Far memory”/Fast NVM reference 1,000 ns (1us)

Read 1 MB sequentially from memory 12,000 ns (12 us)

SSD Random Read 100,000 ns (100 us)

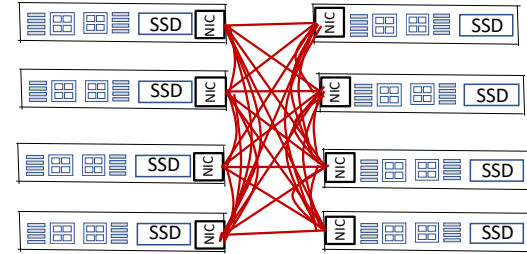
Read 1 MB bytes sequentially from SSD 500,000 ns (500 us)

Read 1 MB sequentially from 10Gbps network 1,000,000 ns (1 ms)

Read 1 MB sequentially from disk 10,000,000 ns (10 ms)

Disk seek 10,000,000 ns (10 ms)

Send packet California→Netherlands→California (150 ms)



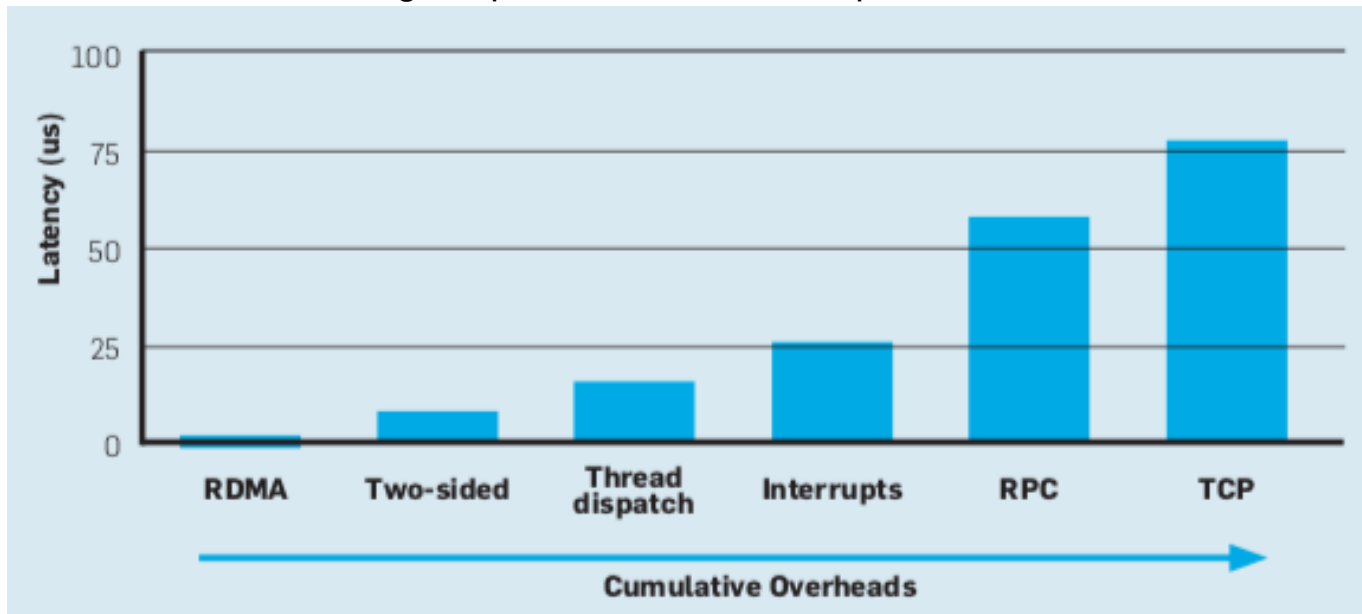
A cluster of nodes with a Clique network topology and low latency optical network...

Yields one hop network latencies on par with DRAM access latencies.

Source: **The Datacenter as a Computer: Designing Warehouse-Scale Machines**, Luiz Andre Barroso, Urs Holzle, Parthasarathy Ranganathan, 3rd edition, Morgan & Claypool, 2019.

Take out the big stuff & you're left with lots of μ s overheads

All those SW overheads add up ... like bricks that combine to build a networking-wall ... turning a 2 μ s network into a 100 μ s network...



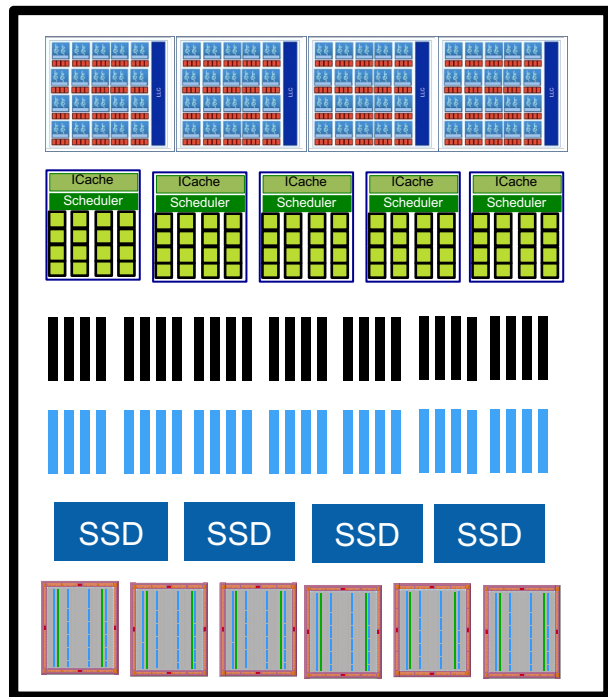
Computer Scientists need to rethink system SW stacks to minimize latencies ... fast RDMA, reduce sync contention, low latency interrupt handlers, and more All to hit $O(\mu$ s) latencies.

Disaggregated Computing for SW Defined Servers (SDS):

This idea works in the sixth Epoch

Consider a Rack composed of multiple pools

Dynamically compose across pools to match a software defined server to the workload



CPU pool

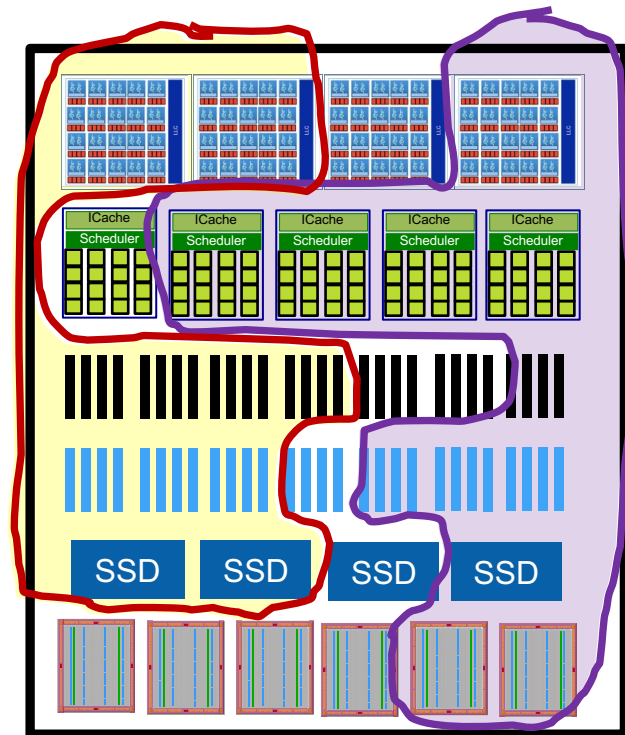
GPU pool

DRAM pool

NVRAM pool

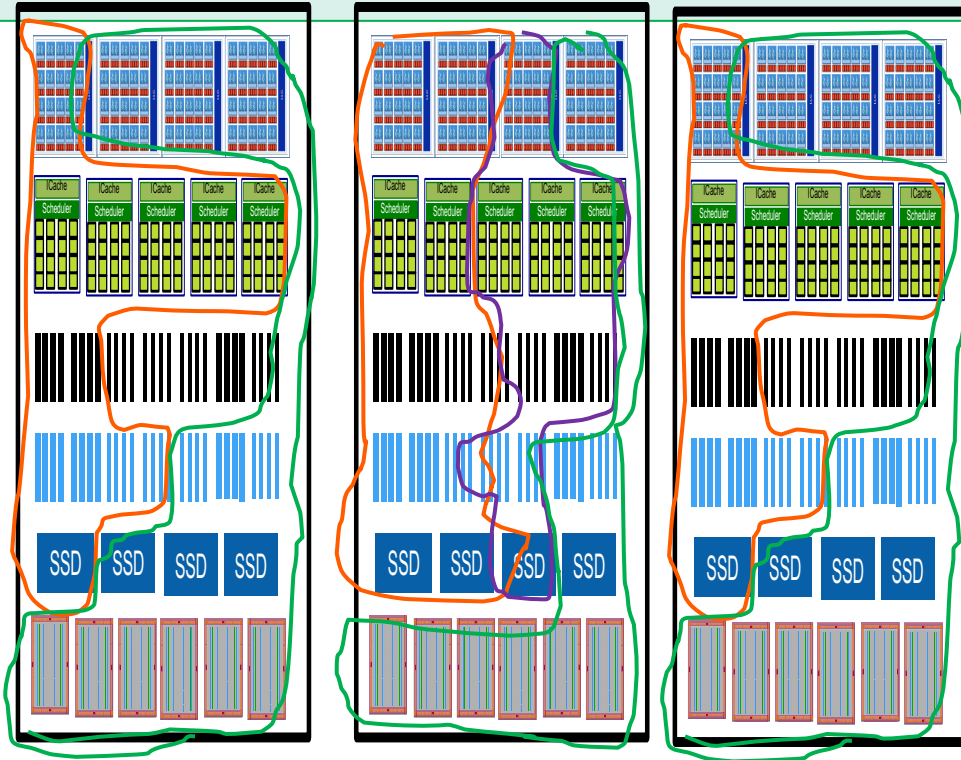
SSD pool

FPGA pool



SW Defined clusters of SW defined Servers

Low latency, high bandwidth network between cliques

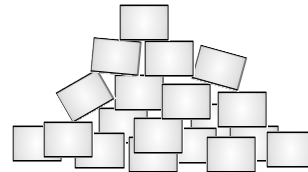


- Dynamic ... changing from one job to the next.
- SW defined servers composed of heterogeneous components
- Dynamically composed into a cluster
- Integrated over a 5G network to devices (and people) at the edge

The three domains of parallel programming

Platform*	Laptop or server	HPC Cluster	Cloud
Execution Agent	Threads	Processes	Microservices
Memory	Single Address Space	Distributed memory, local memory owned by individual processes	Distributed object store (in memory) backed by a persistent storage system
Typical Execution Pattern	Fork-join	SPMD	Event driven tasks, FaaS, and Actors

Advances in networking technology plus low-overhead software stacks optimized to reduce tail-latency will shatter the wall between cloud computing and the other domains



Six Epochs of Distributed Computing ... SW implications

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	High latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server		Computer to computer
Third 1990		MPP HPC, three-tier datacenter networks		
Fourth 2000	Dennard Scaling Wall ... per core plateau	Web search, planet-scale services	100 μ s	People to
Fifth 2015	Per socket wall ... accelerators take off	Machine Learning, data centric computing	Flash	
Sixth 2025	Speed of light	Dynamic AI, integrated from data-center to the edge.	100	

Concurrency: none

Concurrency: event driven

Concurrency: Multithreaded for latency hiding

Concurrency: Multithreaded for performance

Concurrency: Multi-server, multi-threaded

Concurrency: real-time, hierarchical & heterogenous

Writing Parallel Distributed Applications

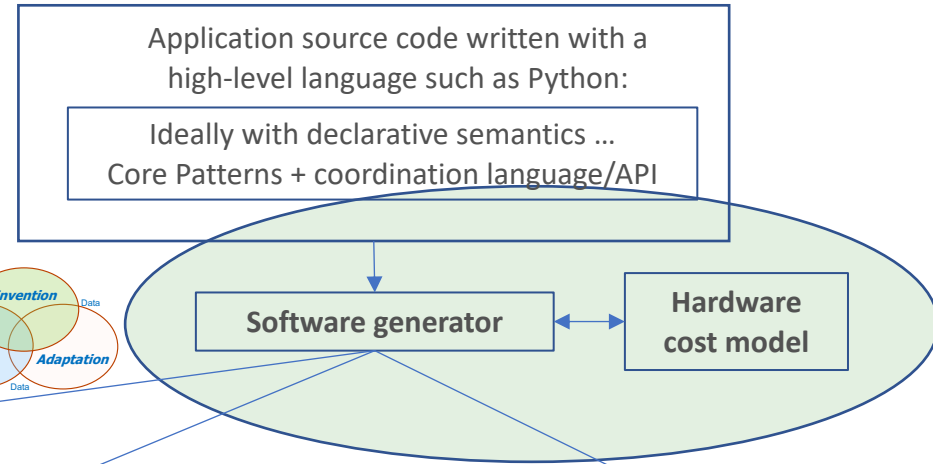
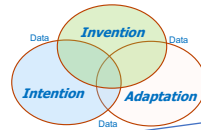
Network technology evolution:

- Lower and more predictable latencies
- Erase distinctions between HPC clusters & the cloud

In response ... we must support:

- One code base → multiple execution models

We call this **machine programming**



- Application task-groups → microservices
- Data structures → distributed object store
- Durable store: Persistent cloud store (e.g. S3)

Cloud Native HPC

- Application task-groups → processes
- Data structures → in process memory
- Durable Store: Cluster file system

HPC Cluster

- Applications task-groups → threads
- Data structures → process heap
- Durable store: local file system

Laptop/server

Conclusion

- We are on the threshold of the sixth Epoch of distributed computing.
- This will shatter the wall separating traditional HPC from the cloud.
- As optical networks become the norm, much of what we do in HPC will be driven to the cloud.
- The issue for the HPC community ... will our software be ready for this brave new world?