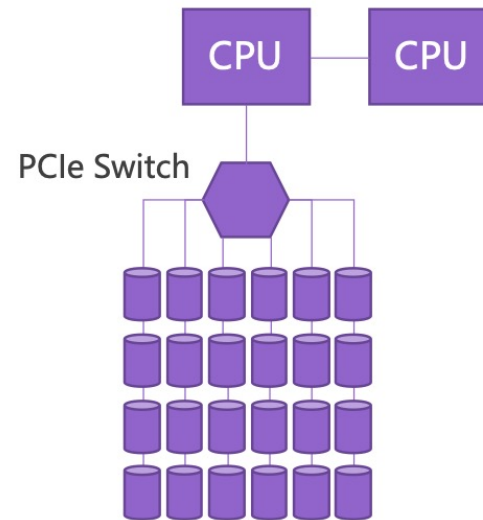


PCIe Net enable Object Storage at Edge

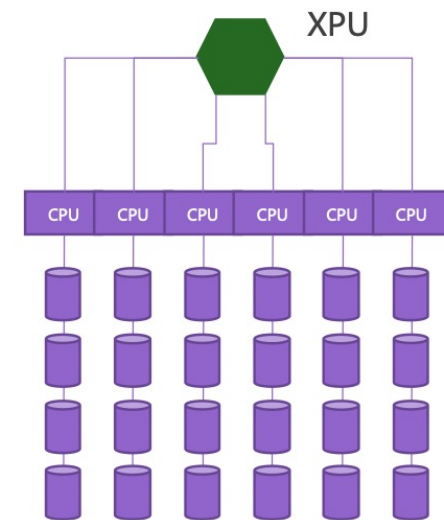


Traditional Storage
Elephant



- High cost
- High power
- Not cloud native

Micro-server Architecture
Ants



- Low cost
- Low power
- Cloud native

A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion - fuelled by internet of things and the use of connected devices - are hard to comprehend, particularly when looked at in the context of one day

500m

tweets are sent every day

Twitter



4PB

of data created by Facebook, including

350m photos

100m hours of video watch time

Facebook Research

320bn

emails to be sent each day by 2021

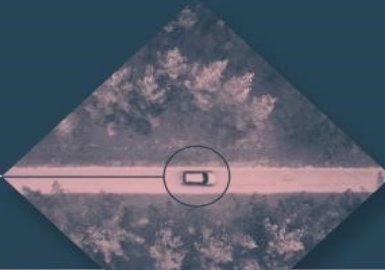
306bn

emails to be sent each day by 2020

294bn

billion emails are sent

Radicati Group



4TB

of data produced by a connected car

Intel

65bn

messages sent over WhatsApp and two billion minutes of voice and video calls made

Facebook



463EB

of data will be created every day by 2025

IDC

95m

photos and videos are shared on Instagram

Instagram Business



28PB

to be generated from wearable devices by 2020

Statista



DEMYSTIFYING DATA UNITS

From the more familiar 'bit' or 'megabyte', larger units of measurement are more frequently being used to explain the masses of data

Unit	Value	Size
b bit	0 or 1	1/8 of a byte
B byte	8 bits	1 byte
KB kilobyte	1,000 bytes	1,000 bytes
MB megabyte	1,000 ² bytes	1,000,000 bytes
GB gigabyte	1,000 ³ bytes	1,000,000,000 bytes
TB terabyte	1,000 ⁴ bytes	1,000,000,000,000 bytes
PB petabyte	1,000 ⁵ bytes	1,000,000,000,000,000 bytes
EB exabyte	1,000 ⁶ bytes	1,000,000,000,000,000,000 bytes
ZB zettabyte	1,000 ⁷ bytes	1,000,000,000,000,000,000,000 bytes
YB yottabyte	1,000 ⁸ bytes	1,000,000,000,000,000,000,000,000 bytes

*A lowercase "b" is used as an abbreviation for bits, while an uppercase "B" represents bytes.



Smart Insights

ACCUMULATED DIGITAL UNIVERSE OF DATA

4.4ZB

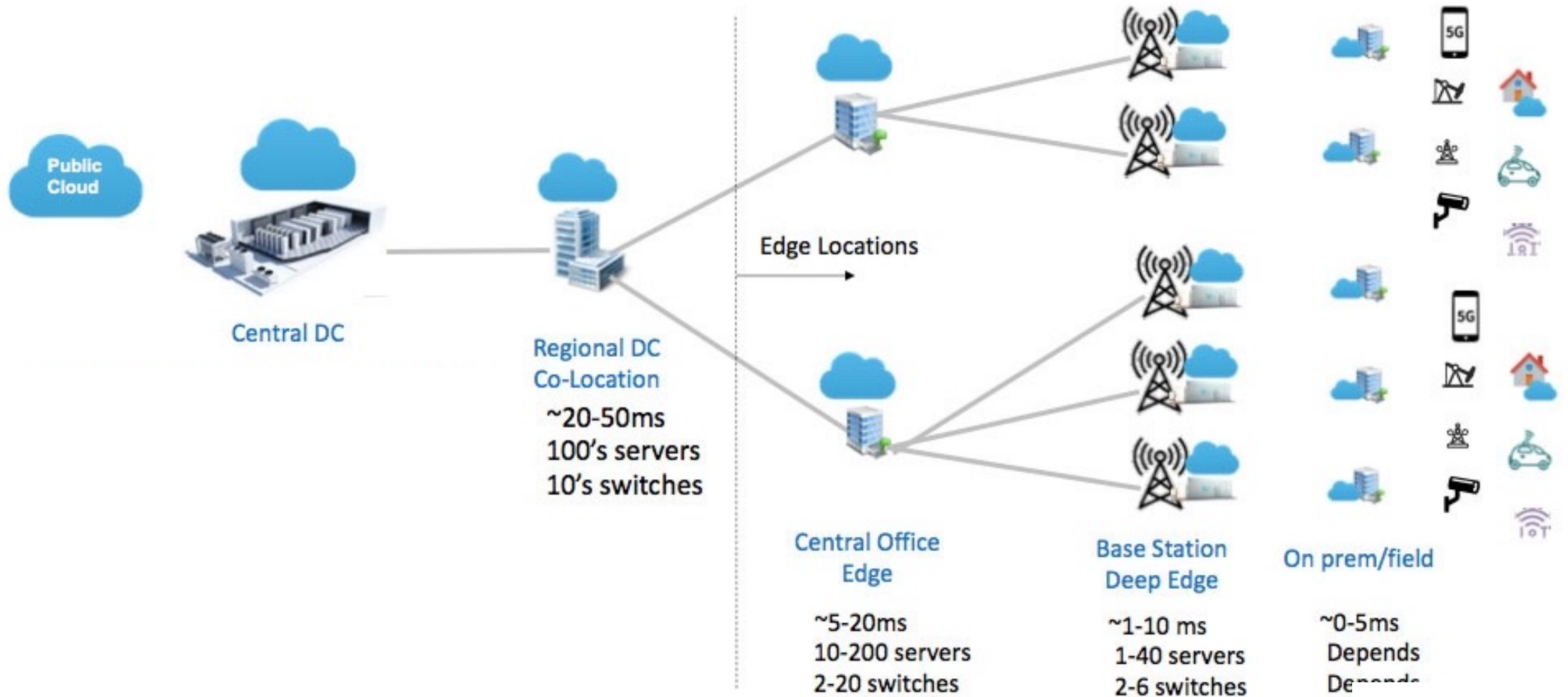
44ZB

PwC

2013

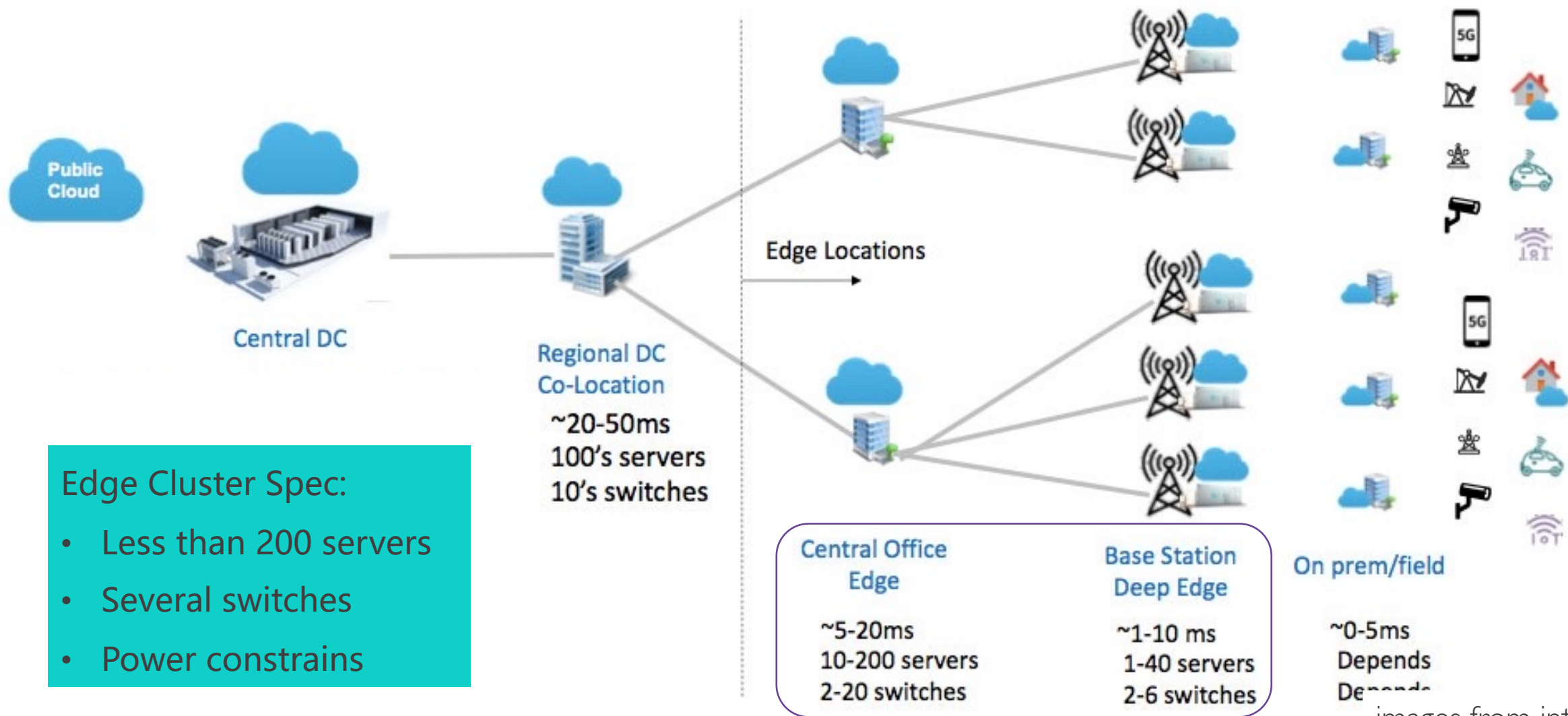
2020

Data Flow from Cloud to Deep Edge



images from internet

Edge Cluster Scale-in



Edge Cluster Spec:

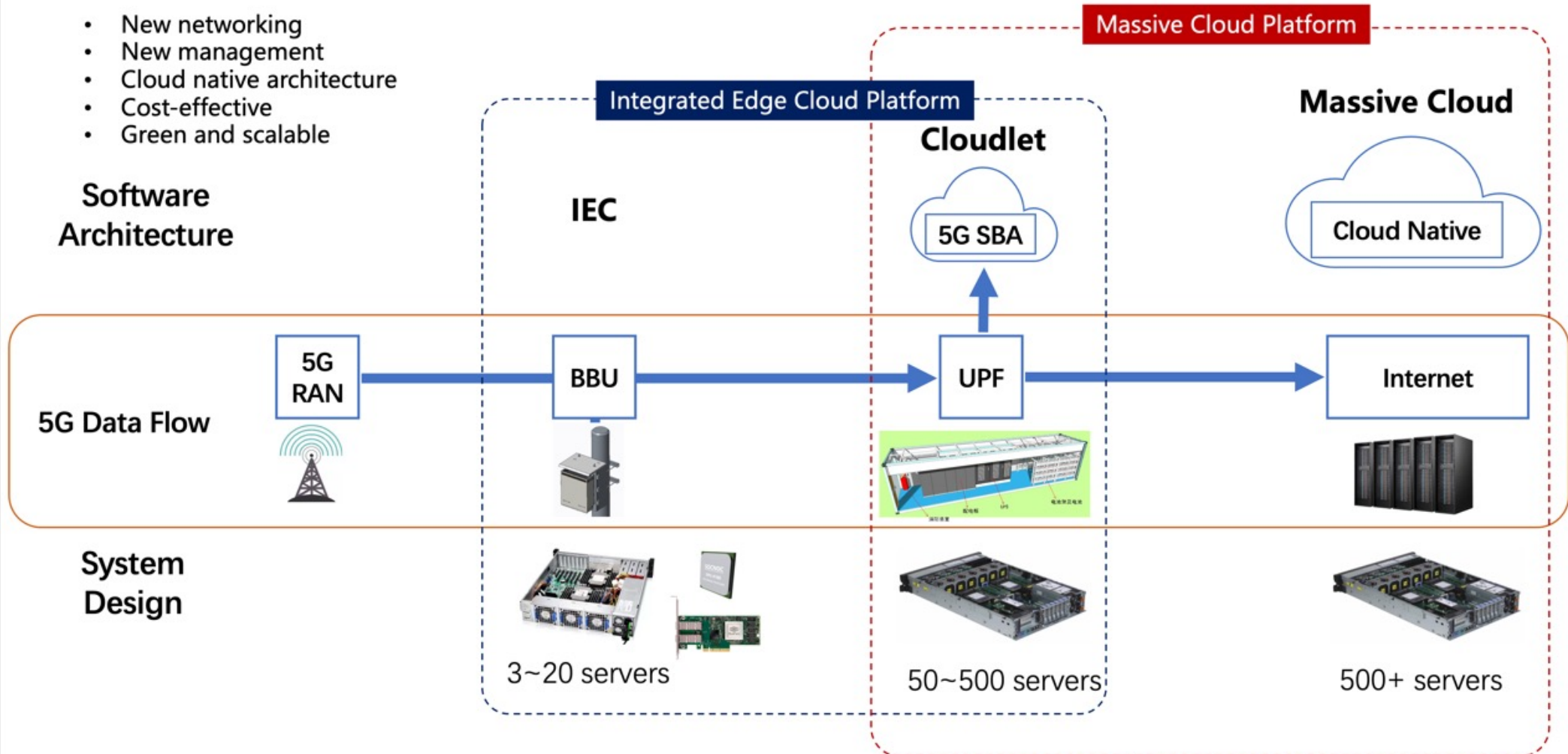
- Less than 200 servers
- Several switches
- Power constrains

images from internet

Integrated Edge Cloud (IEC) - Akraino Type 5

Integrated Edge Cloud for Small Size Networked Cluster

- New networking
- New management
- Cloud native architecture
- Cost-effective
- Green and scalable



Scale-in Cluster at Edge

OLFEDGE

“Due to the dynamic nature of edge deployments, it is critical that networks are able to adapt to current operational context in order to optimize performance and uptime.”

and Industrial worlds, many systems communicate over legacy local area networks such as 4-20mA current loops, serial and CAN Bus, as well as modern low-power wireless technologies such as Bluetooth and LoRa. IoT gateways serve the function of converting these transports into IP traffic.

In terms of application-level protocols, there is a wide array of choices to contend with when developing edge solutions. While there are tens that matter in the IT world (e.g. REST, MQTT), there are literally hundreds if not thousands in the OT/Industrial world, with examples including Modbus, BACnet, PROFINET and EtherCAT. Edge solutions often have to comprehend a blend of these application-level protocols and LF Edge projects like EdgeX Foundry and Fledge are focused on simplifying data flow in heterogeneous environments.

Due to the dynamic nature of edge deployments, it is critical that networks are able to adapt to current operational context in order to optimize performance and uptime. This can include dynamically switching between available connections.

In contrast to centralized data centers, the networking in the compact integrated edge cloud needs to be reconsidered, as the networking challenge shifts from solving for the *scale-out* to a massive number of connected servers to enabling the *scale-in* for connecting a small number of servers in an edge location. The time-honored methods of increasing the *port density* in a switch or utilizing the high throughput NIC won't work in a small scale cluster with less than 32 servers. Hence a novel way to rebuild the networking architecture for the integrated edge cloud is needed not only in terms of cost but also, and even more importantly, due to energy constraints, given the expected large number of deployments of the integrated edge cloud sites.

Project Contributions for Edge Connectivity

The LF Edge projects are addressing connectivity needs both at the application level for protocol normalization to facilitate IoT interoperability and transport level in terms of network virtualization and optimization. The following are examples of each project's focus in the area of connectivity.



Akraino blueprints can provide an end to end EdgeStack to support Virtualized Network Elements (NFVI) per Open-RAN (O-RAN) requirements. The Akraino project advocated collaborating with O-RAN Alliance's specification workgroup 6 (six) responsible for cloud specifications, to align with and publish multiple blueprints to support various RAN use cases for Radio Edge cloud, including ORAN-Software Community's Near-RT RIC software, Network Cloud with RS-IOV or OVS-DPDK, Integrated Cloud Native, Kubernetes Native Infrastructure provider Access edge and more.

There are a number of blueprints that provide interesting examples for edge connectivity that can be applied in different parts of the edge ecosystem. As an example, the Network Cloud with Tungsten Fabric (TF) blueprint provides a fully distributed networking stack based on a microservices architecture, implementing a distributed networking framework for Edge computing. TF SDN Controller provides seamless and full integration between different types of workloads VNFs, CNFs and PNFs using a common networking stack integrated with different orchestration platforms like OpenStack and Kubernetes. The TF SDN Controller works as single entity running at the core, distributed core or edge sites, or public cloud (AWS, Azure, GCP or Equinix Metal) and fully integrated with OpenStack Neutron Plugin, Kubernetes CNI, for all types of Edge computing workloads. The solution provides the Tungsten Fabric Kernel vRouter, DPDK vRouter, and support for SR-IOV and SmartNIC.

“The Network Cloud with Tungsten Fabric (TF) blueprint provides a fully distributed networking stack based on a microservices architecture, implementing a distributed networking framework for Edge computing.”

Another example of an innovative Akraino blueprint is the Integrated Edge Cloud (IEC) for compact edge with PCIe networking and Cloud-on-Board (CoB). In contrast to centralized data centers, the networking in the compact integrated edge cloud needs to be reconsidered, as the networking challenge shifts from solving for the scale-out to a massive number of connected servers to enabling the scale-in for connecting a small number of servers in an edge location. The time-honored methods of increasing the port density in a switch or utilizing the high throughput NIC won't work in a small scale cluster with less than 32 servers. Hence a novel way to rebuild the networking architecture for the integrated

Compact Networking for Small Cluster

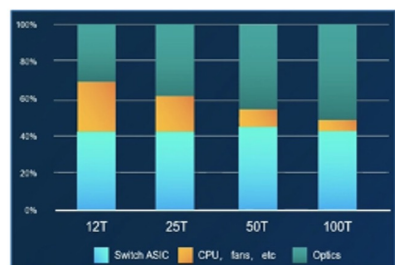
XPU Based Cloud Native Server:

Architecture, Implementation & Applications

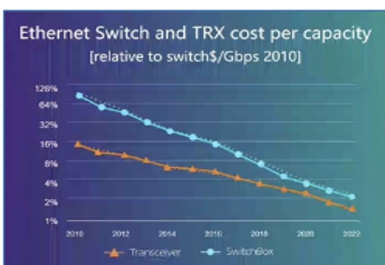
Dr. Fu Li (LEO)

Real Problem Computing Cluster Faced!

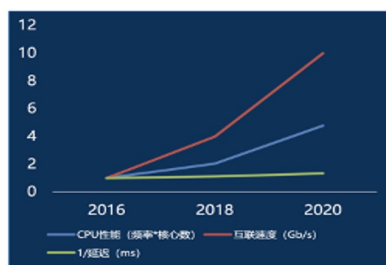
› Optics are too expensive both in **power** and **cost**



光模块功耗占比已经超过50%



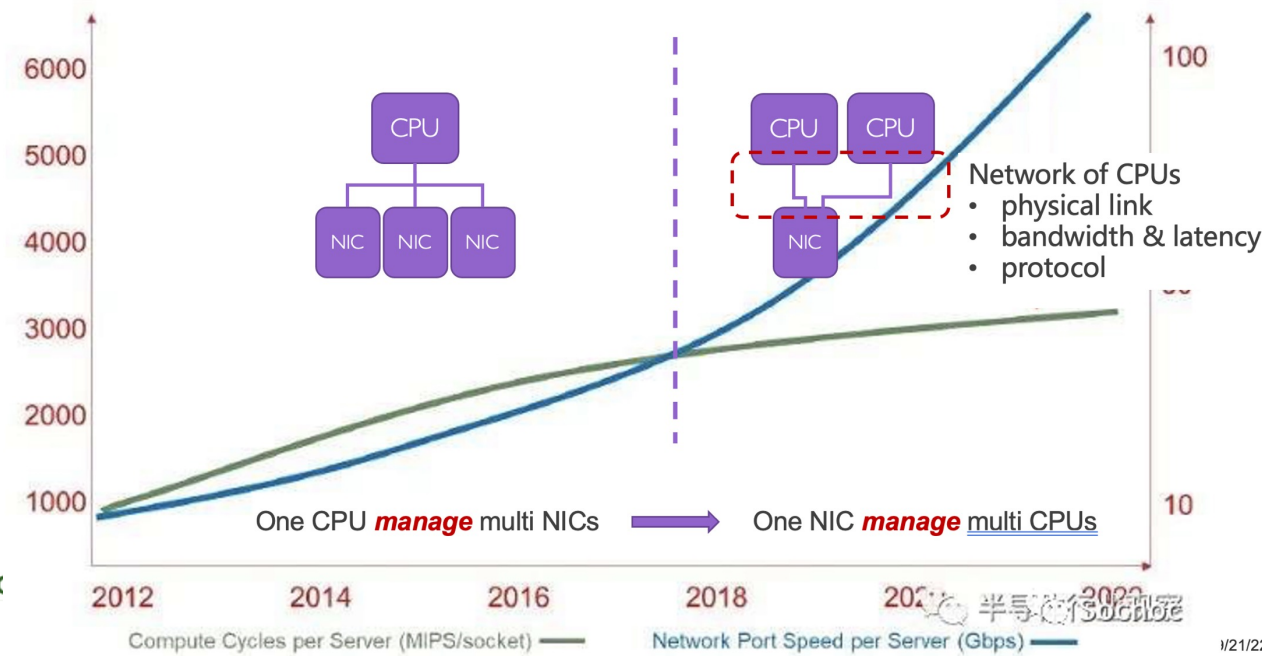
光模块成本已经超过通道成本



延迟降低进展缓慢



Problem Revisit: Clustered CPU on Board



Data Fabric Landscape

TCP/IP



PCIe/CXL
[cxl.io/cxl.cache](#)



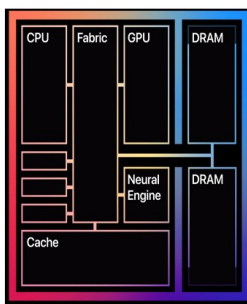
UCIe

● Chiplet



I/O module

● System-in-Package



Interposer Fabric

In-Package Links

● System-on-Board

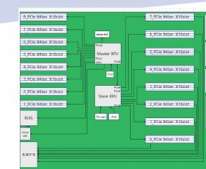


PCIe Bus based link

Off-Package Links

NEW

● Cloud-on-Board

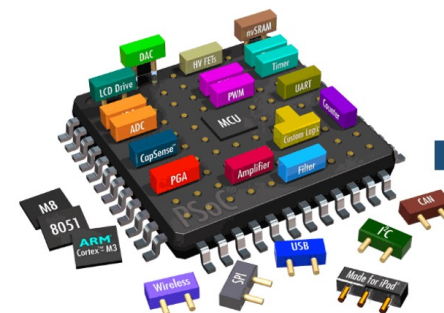


PCIe Net based Fabric

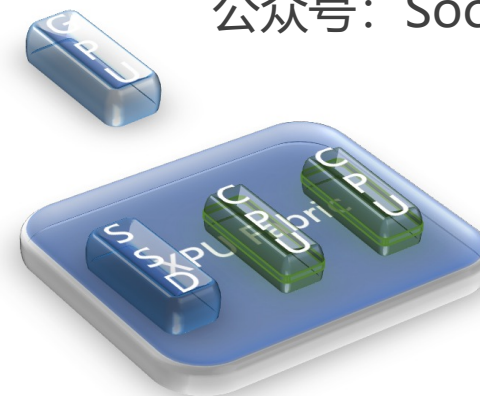


公众号: Socrac

XPU: Center of Data Fabric

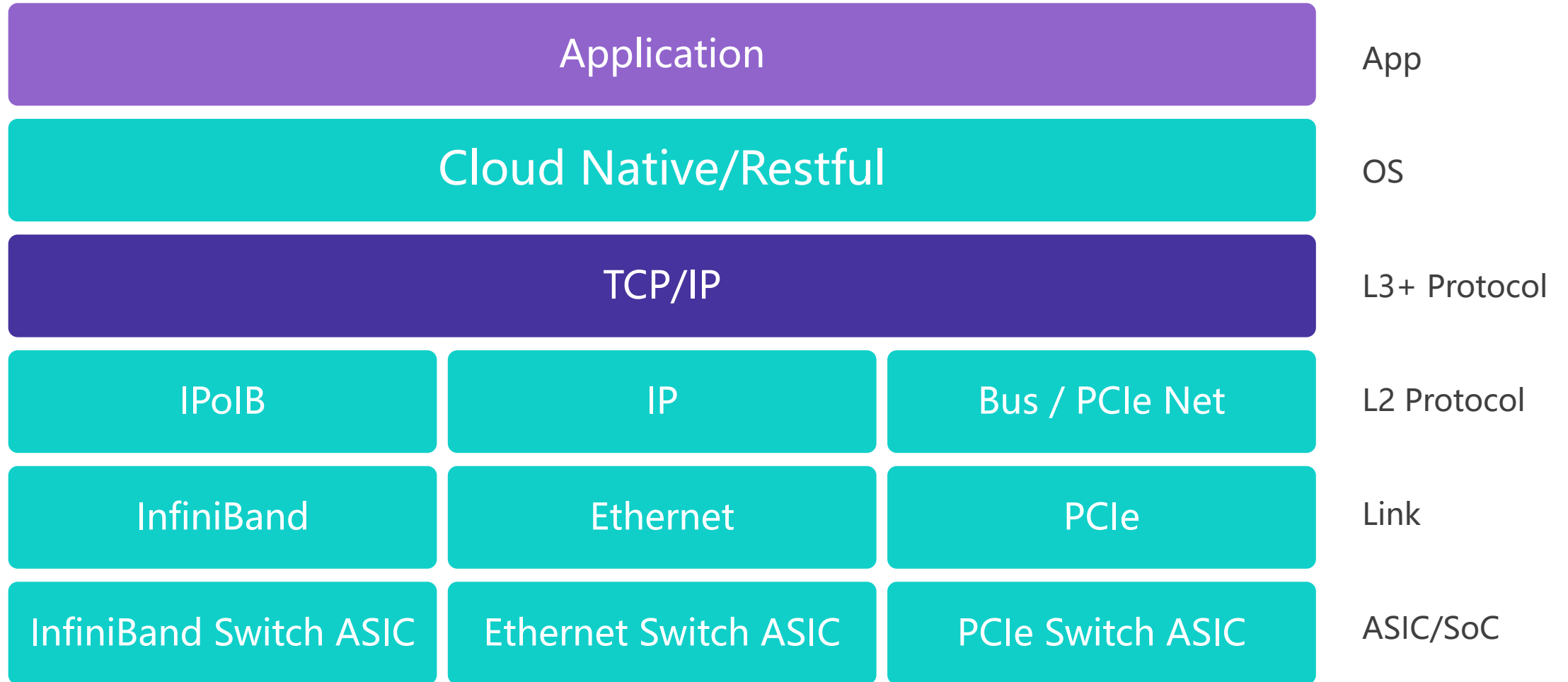


In Package Technology



Off Package Technology

Networking Technologies for Clusters

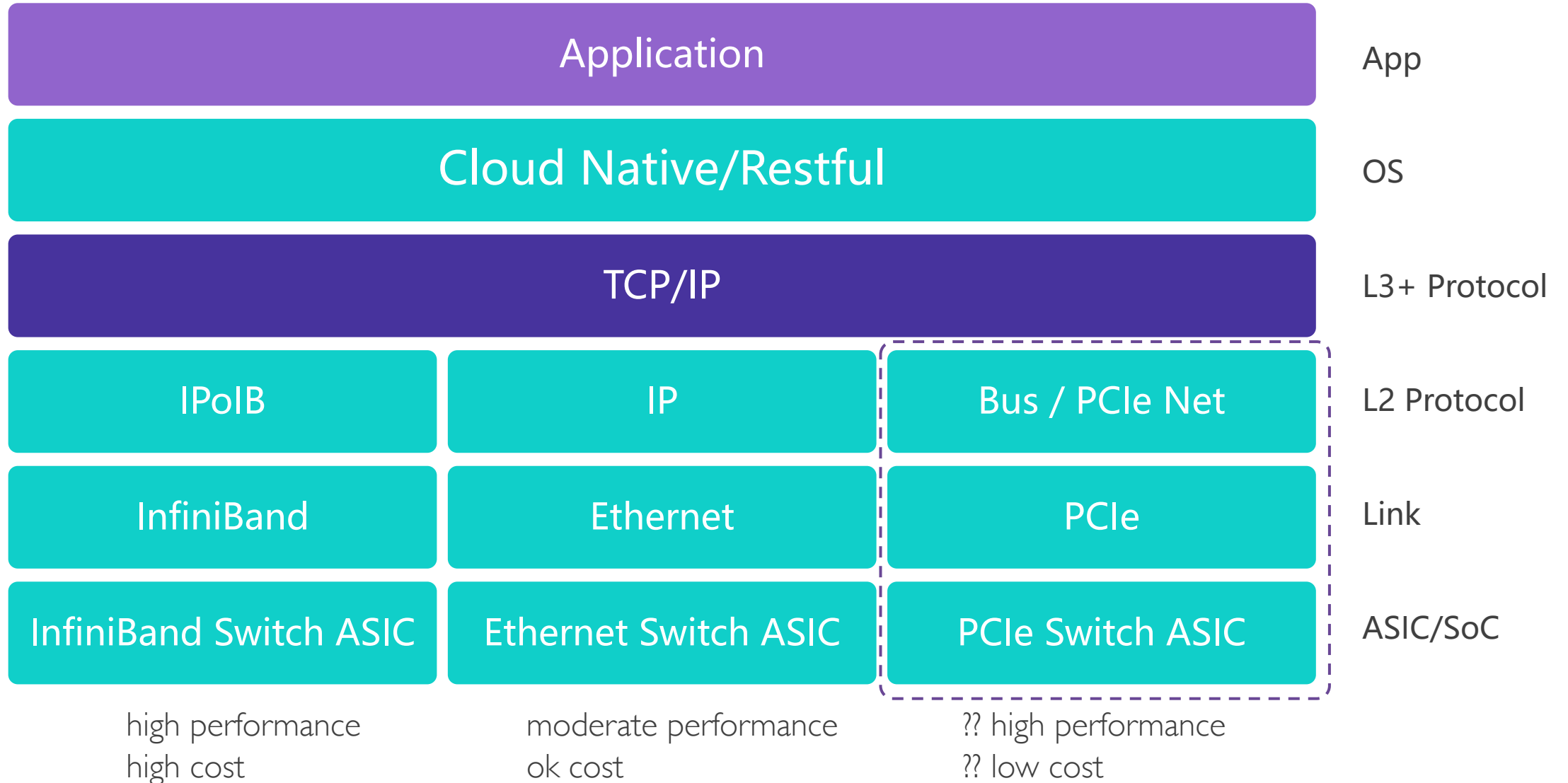


high performance
high cost

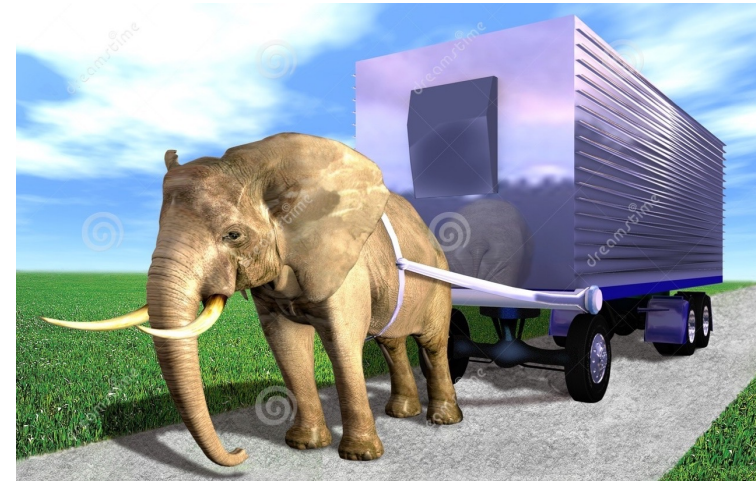
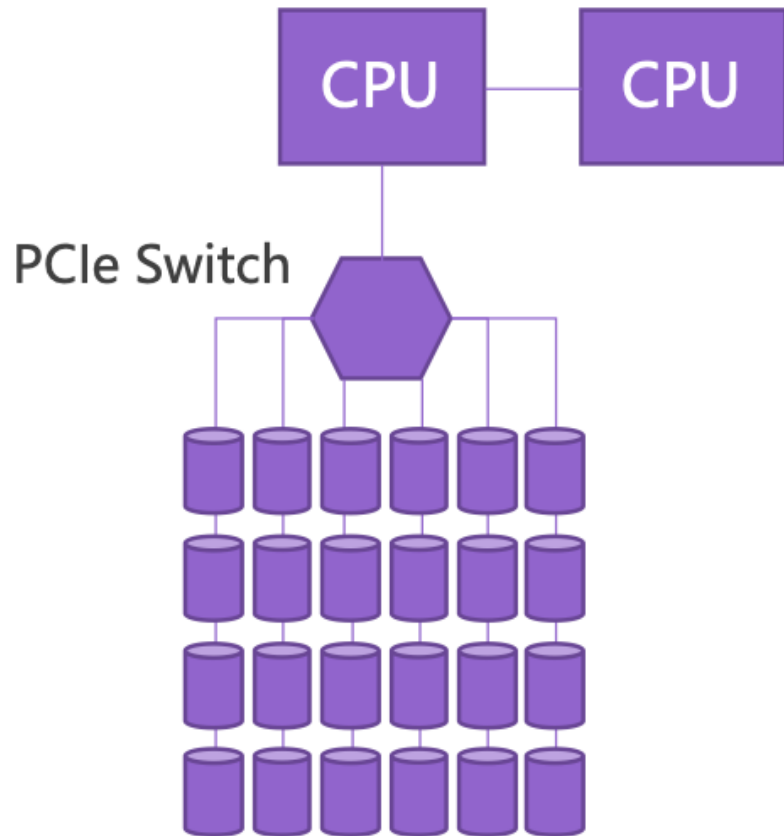
moderate performance
ok cost

?? high performance
?? low cost

PCIe Net and System Bus

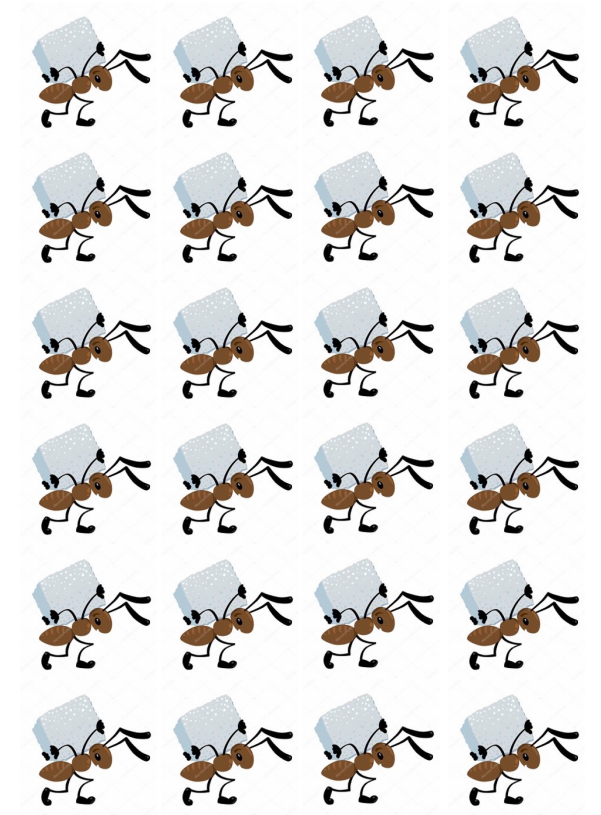
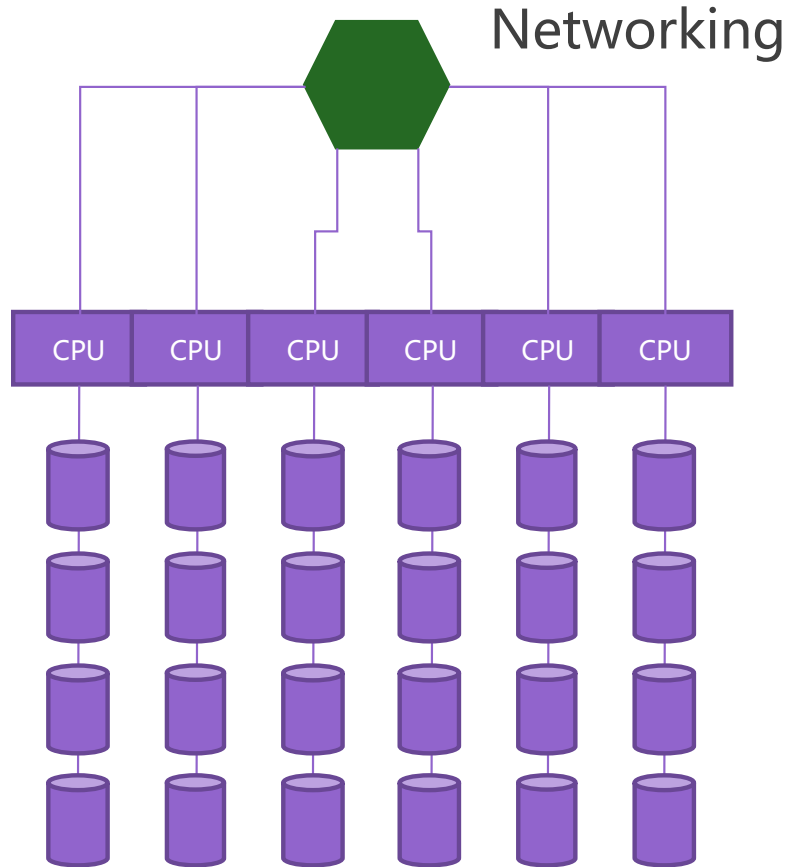


PCIe Fanout Storage System (Elephant)



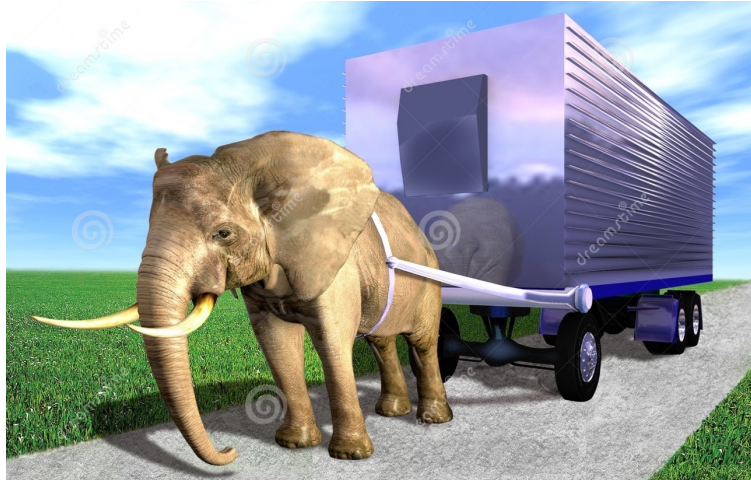
images from internet

Micro server Storage System (Networked Ants)

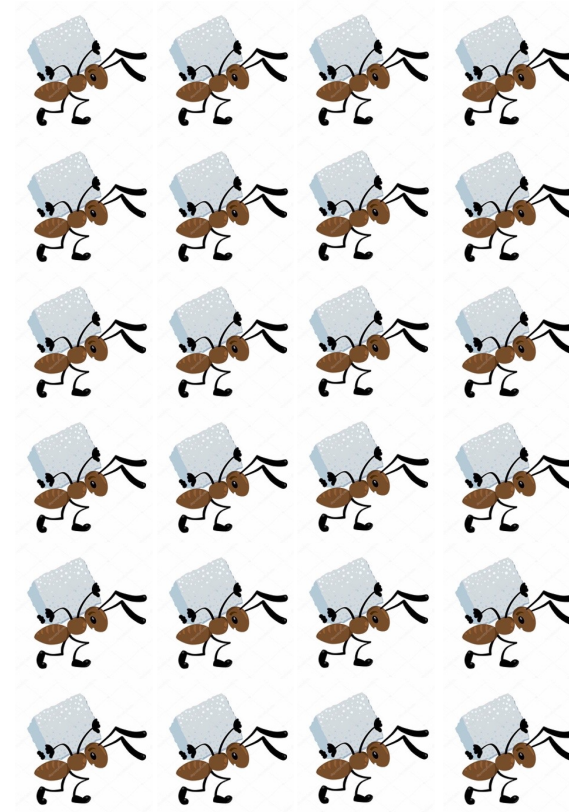


from partner and internet

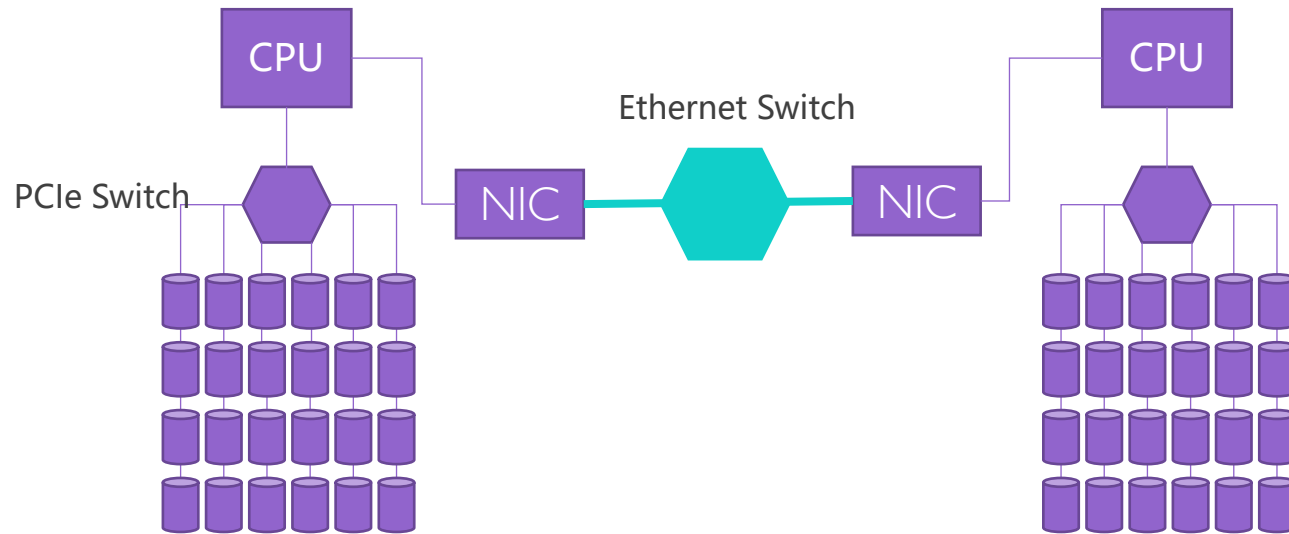
Centralized vs Distributed System



VS

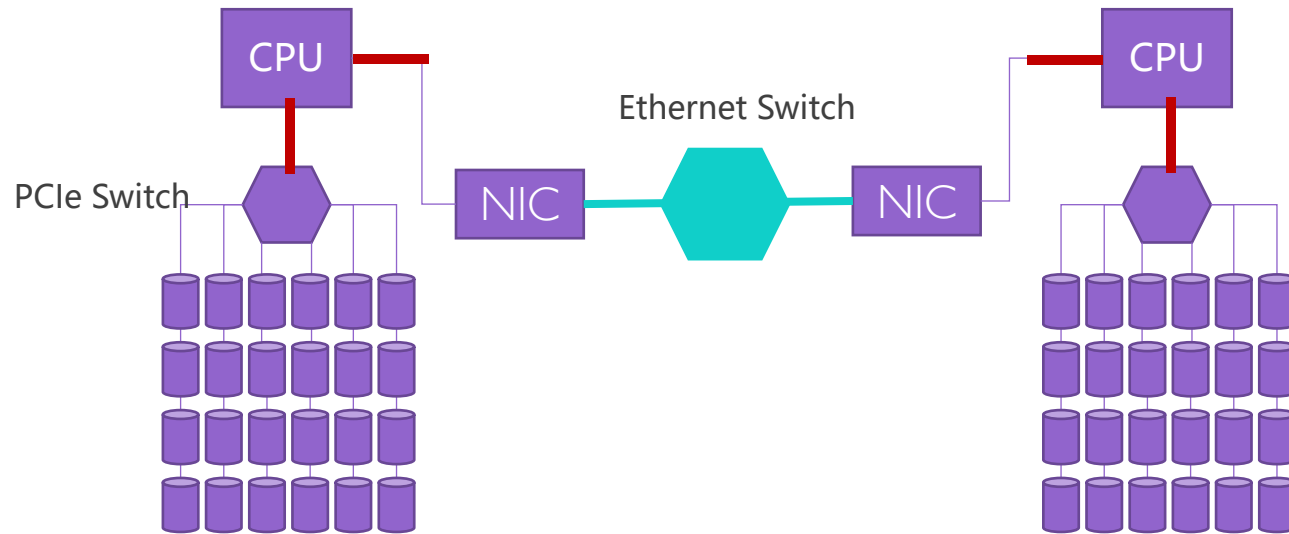


Deep Thinking on System Architecture



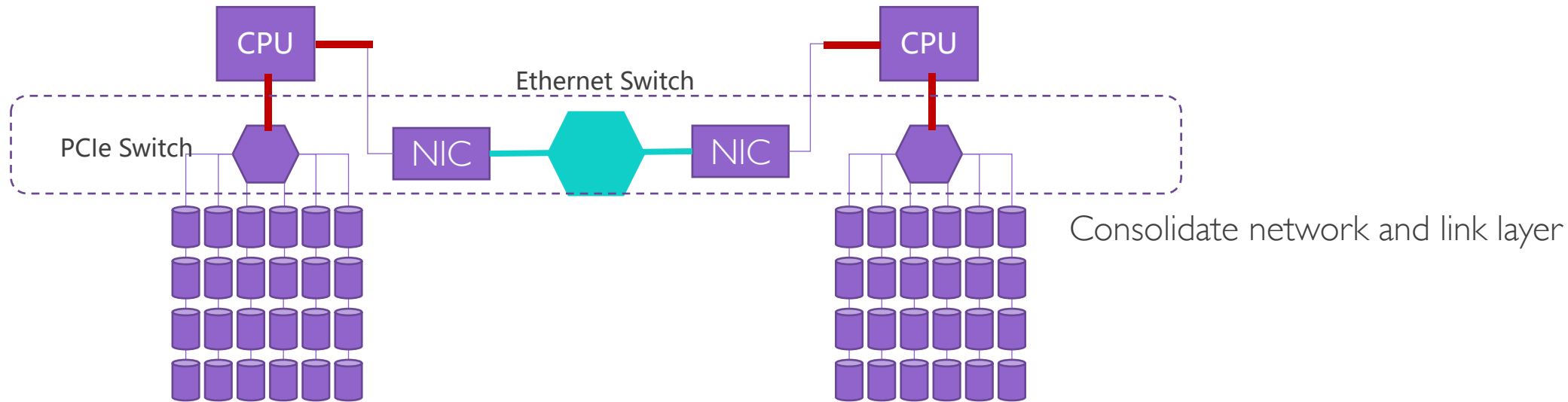
PCIe Switch for Devices, Ethernet Switch for CPUs/Servers

Architecture Bottleneck



	CPU/SoC #	PCIe Ports	Ethernet Switch Ports	Bottleneck
Elephant (HPC server)	2~4	24~36	2~4	Ethernet / CPU
Ants (microserver)	8~24	2-4	8~24	Ethernet Latency

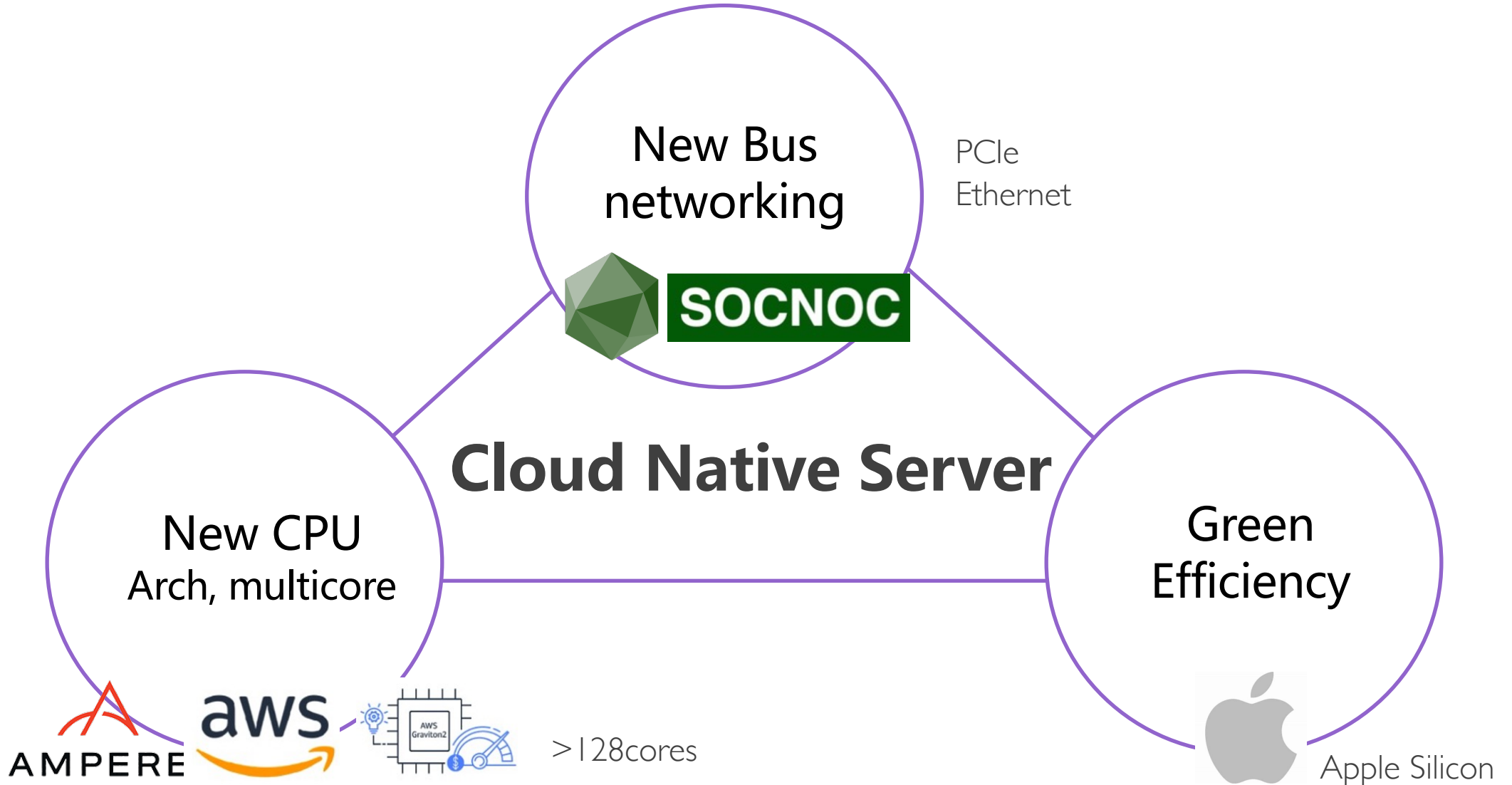
Architecture Innovation?



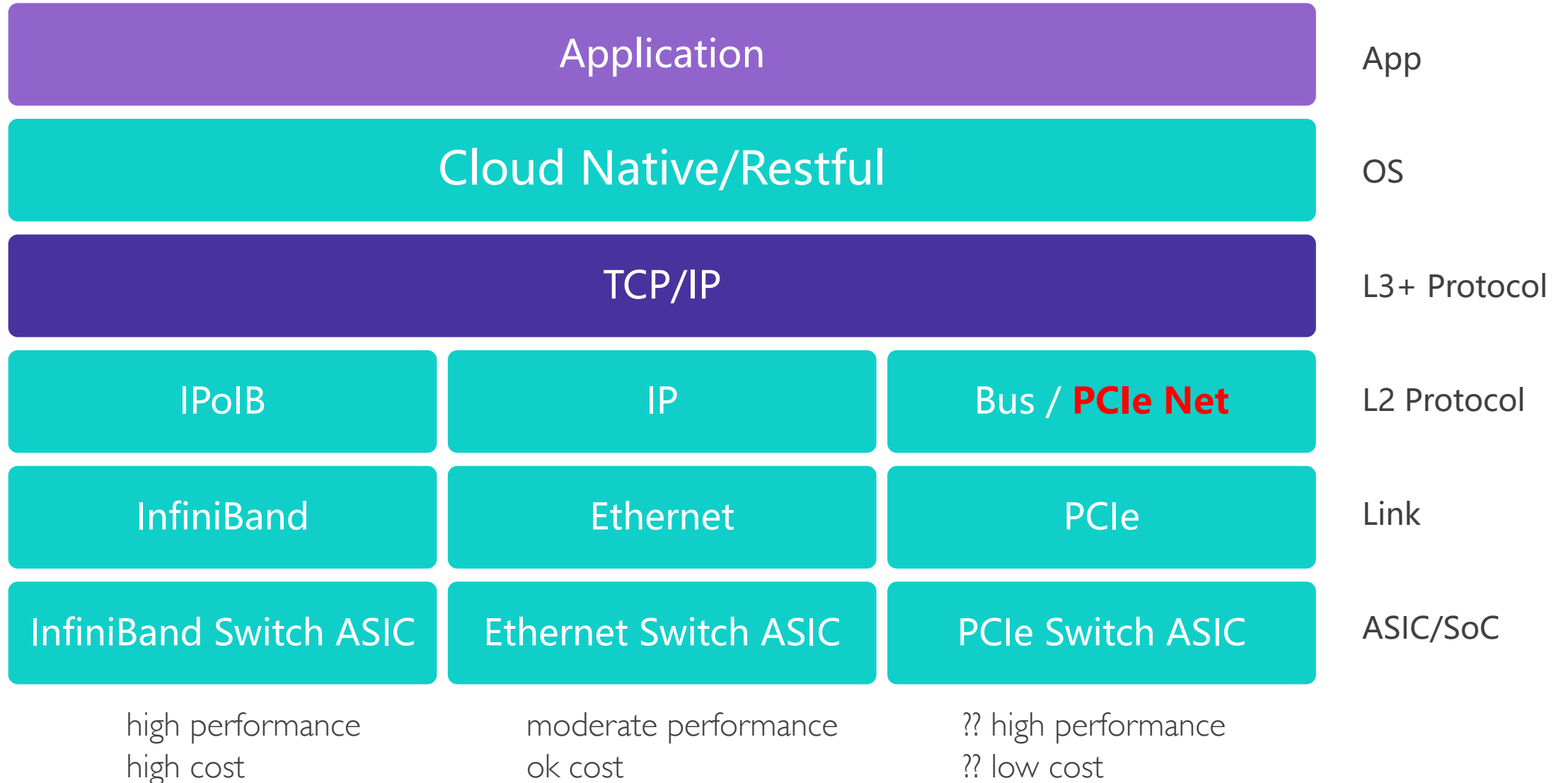
	CPU/SoC #	PCIe Ports	Ethernet Switch Ports	Bottleneck
Elephant (HPC server)	2~4	24~36	2~4	Ethernet / CPU
Ants (microserver)	8~24	2-4	8~24	Ethernet Latency

Solution: Merge PCIe Bus and Ethernet Net PCIe Net

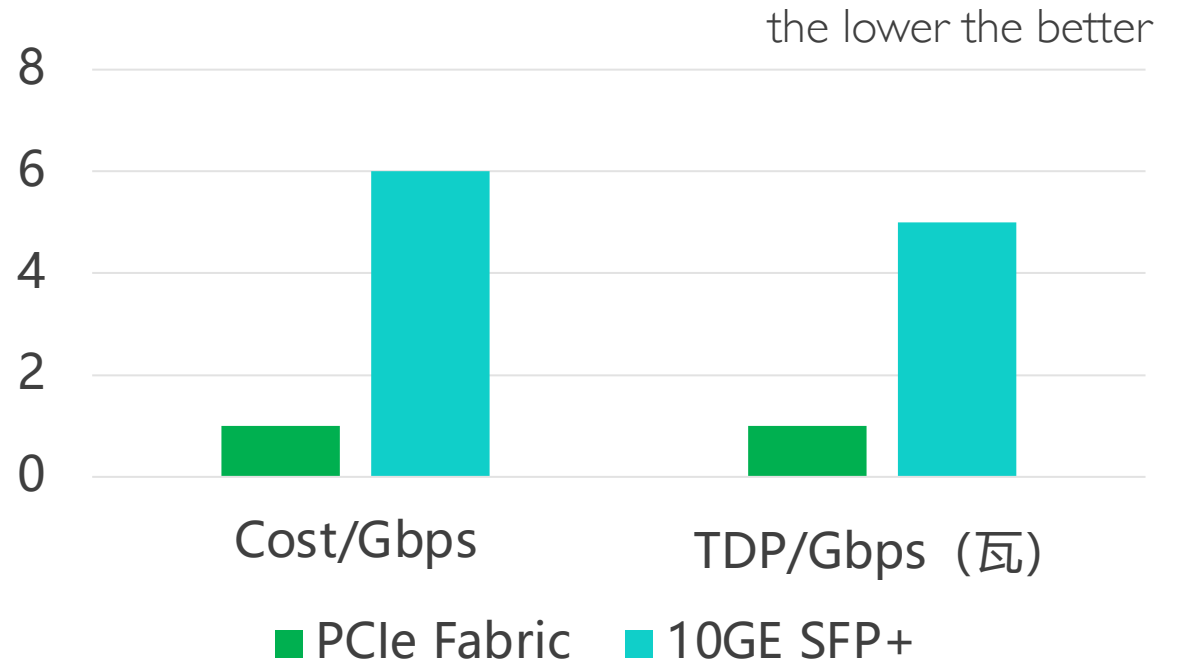
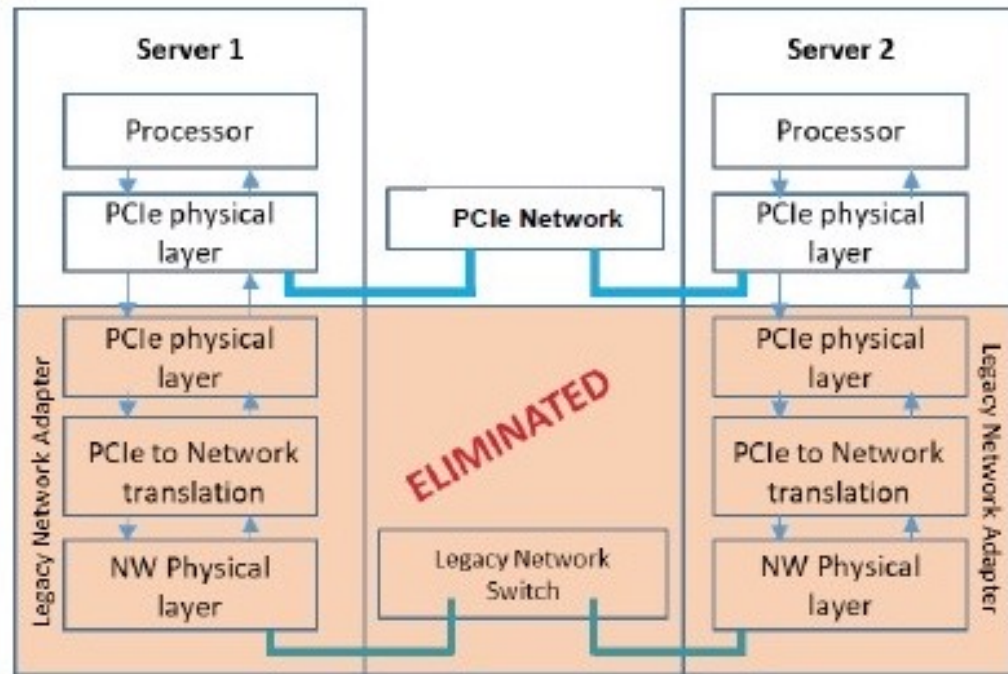
Next-generation System Design



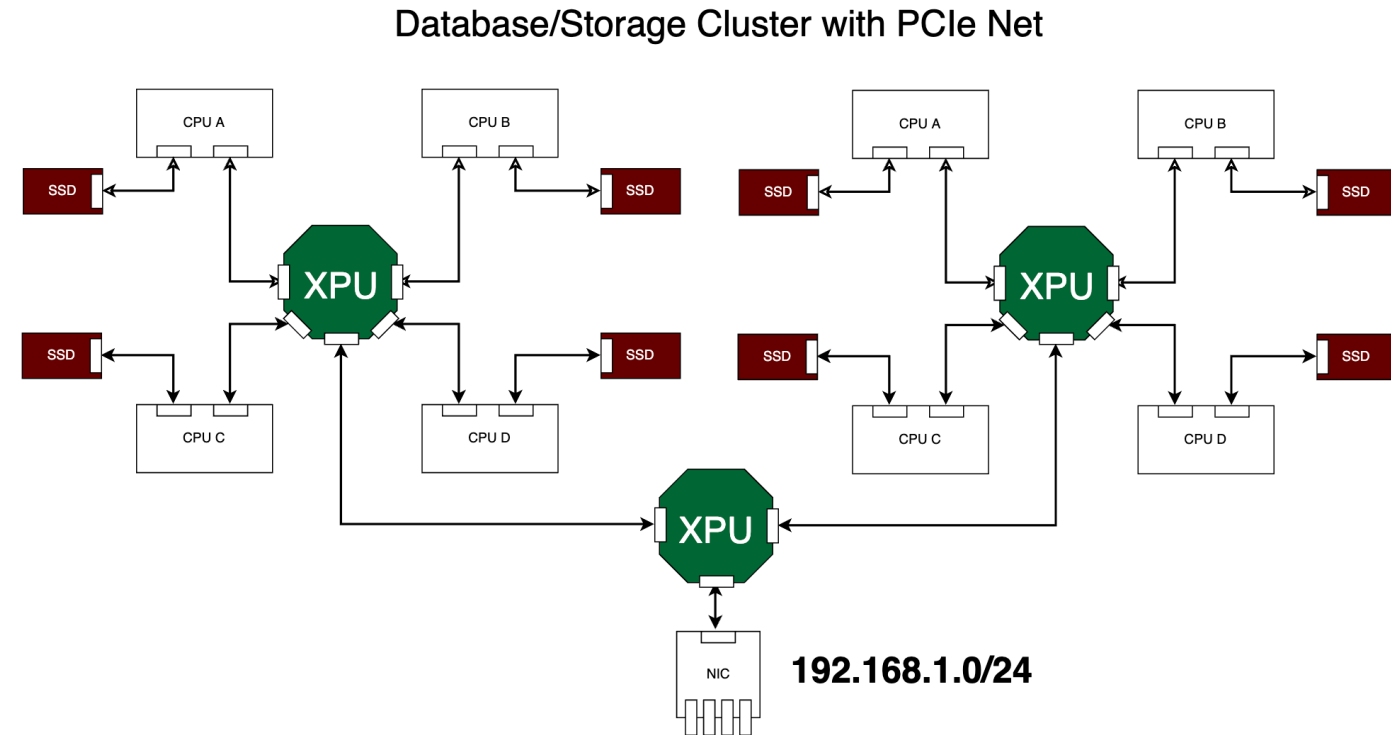
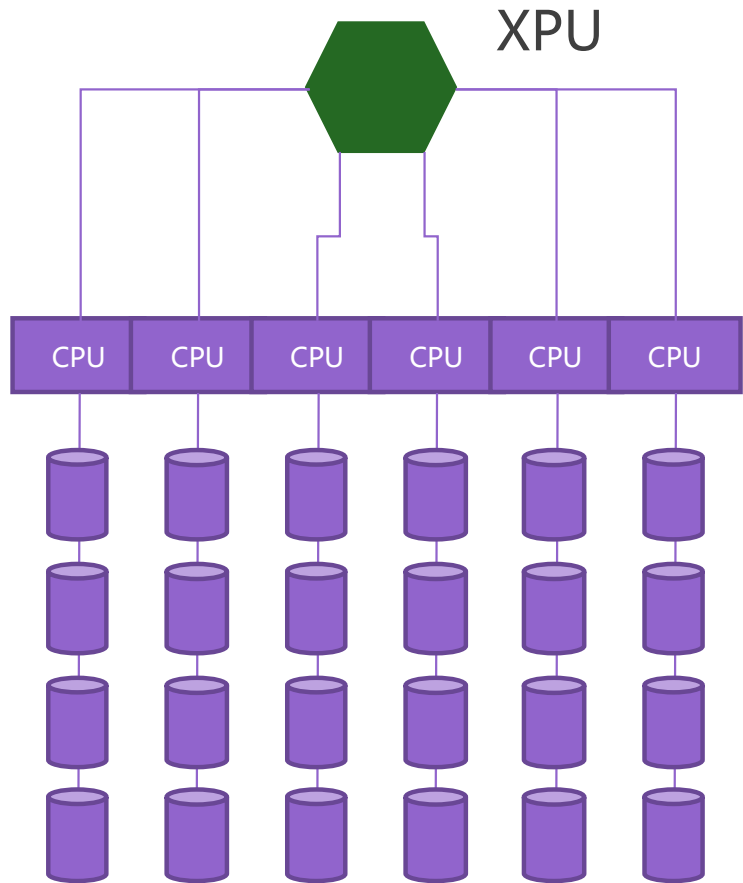
Networking Technologies for Clusters



Extending PCIe Transport with Virtual NIC

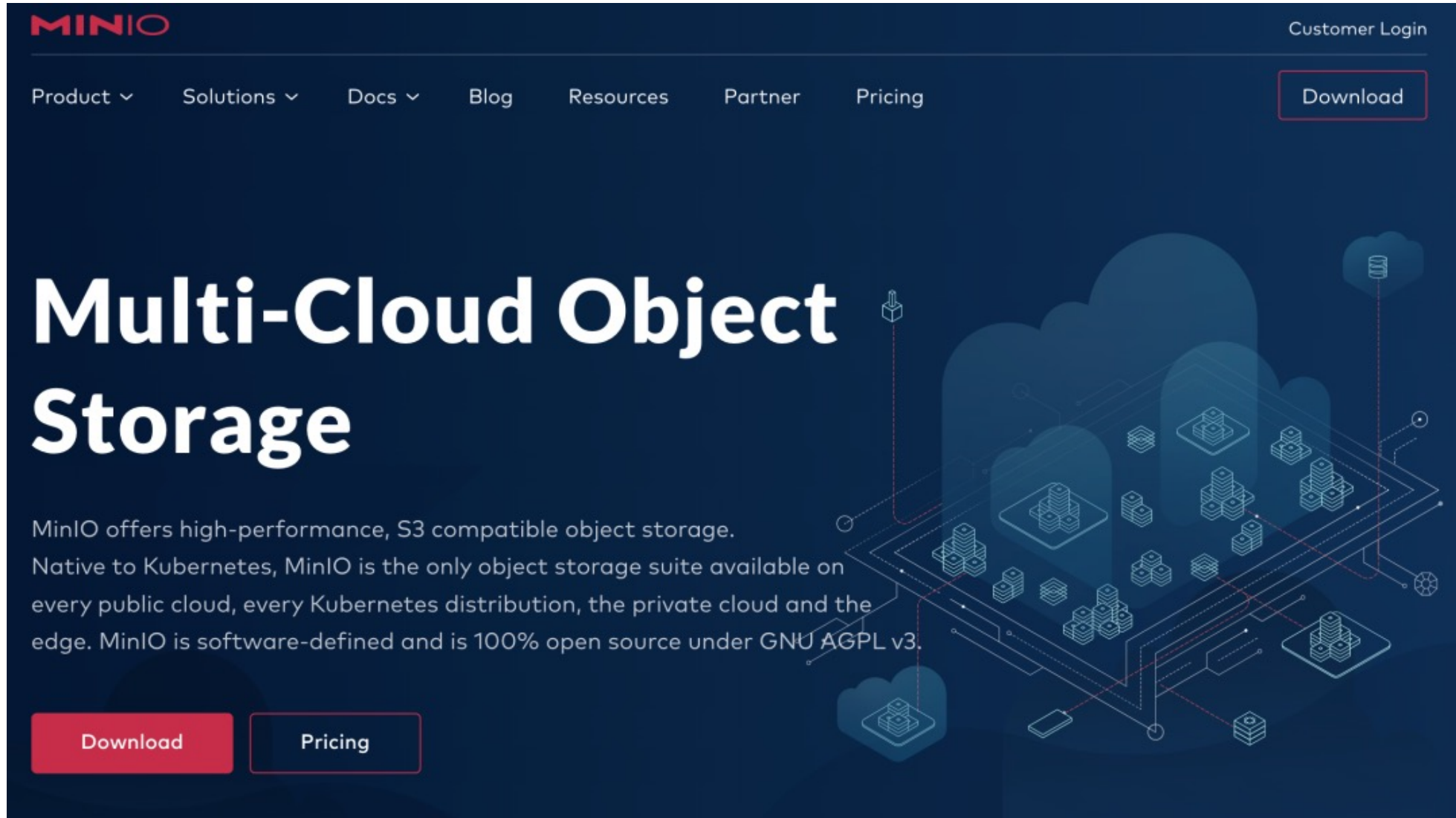


PCIe Net Cluster: Continuous Memory Space



each CPU "share" the same memory space, however communicating with TCP/IP

Demo: Object Storage MinIO over PCIe Net

The image is a screenshot of the MinIO website homepage. The background is a dark blue with a technical illustration of a multi-cloud storage architecture. The illustration shows several cloud icons connected by lines to server racks and storage units, representing a distributed storage system across different environments. The MinIO logo is in the top left corner. The navigation menu includes 'Product', 'Solutions', 'Docs', 'Blog', 'Resources', 'Partner', and 'Pricing'. A 'Customer Login' link is in the top right. A 'Download' button is highlighted in the navigation menu. The main heading is 'Multi-Cloud Object Storage'. Below it, a paragraph describes MinIO's features: high-performance, S3 compatible, native to Kubernetes, and 100% open source. At the bottom, there are two buttons: 'Download' and 'Pricing'.

MINIO

Customer Login

Product ▾

Solutions ▾

Docs ▾

Blog

Resources

Partner

Pricing

Download

Multi-Cloud Object Storage

MinIO offers high-performance, S3 compatible object storage.

Native to Kubernetes, MinIO is the only object storage suite available on every public cloud, every Kubernetes distribution, the private cloud and the edge. MinIO is software-defined and is 100% open source under GNU AGPL v3.

Download

Pricing

Testbed: Four nodes (16 SSD)

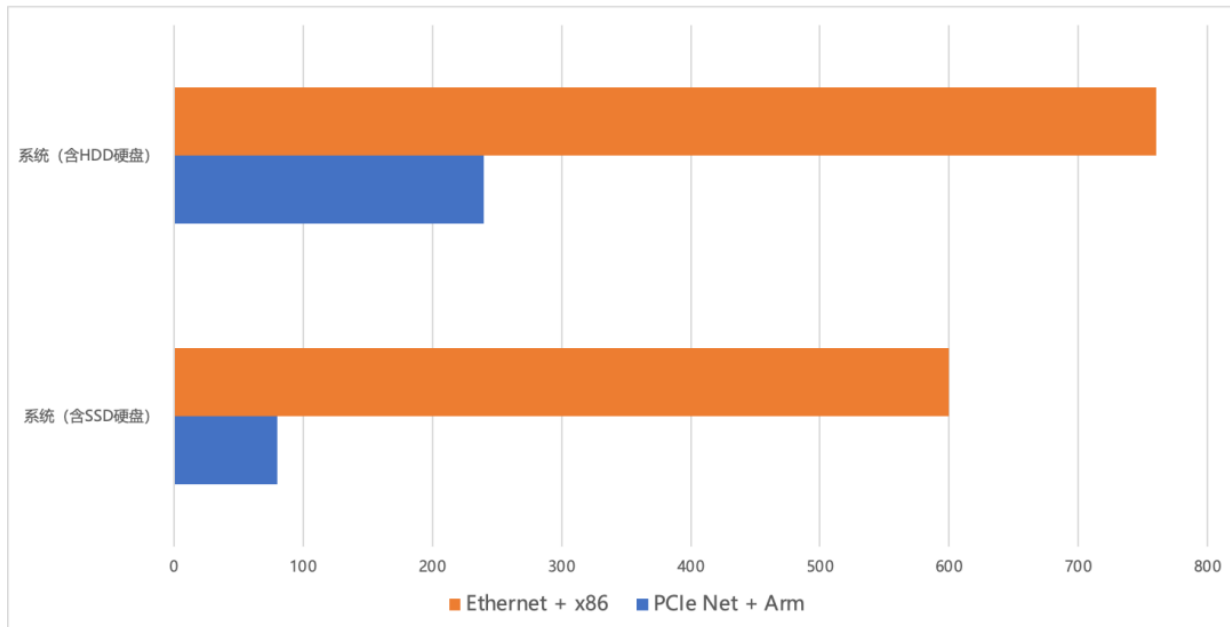
	#node	CPU	MEM	Storage	网络
VPC (x86)	4	24x2.5GHz	48GB	4x250GB	10GbE
IEC (arm)	4	24x1.0GHz	32GB	4x250GB	PCIe Net

	VPC (x86)	IEC (arm)
os	ubuntu 20.04	ubuntu 20.04 (arm)
minIO	RELEASE.2022-07-30	RELEASE.2022-07-30

```
minio server http://min{1...4}.pcie.net:9000/mnt/disk{1...4}/minio
```

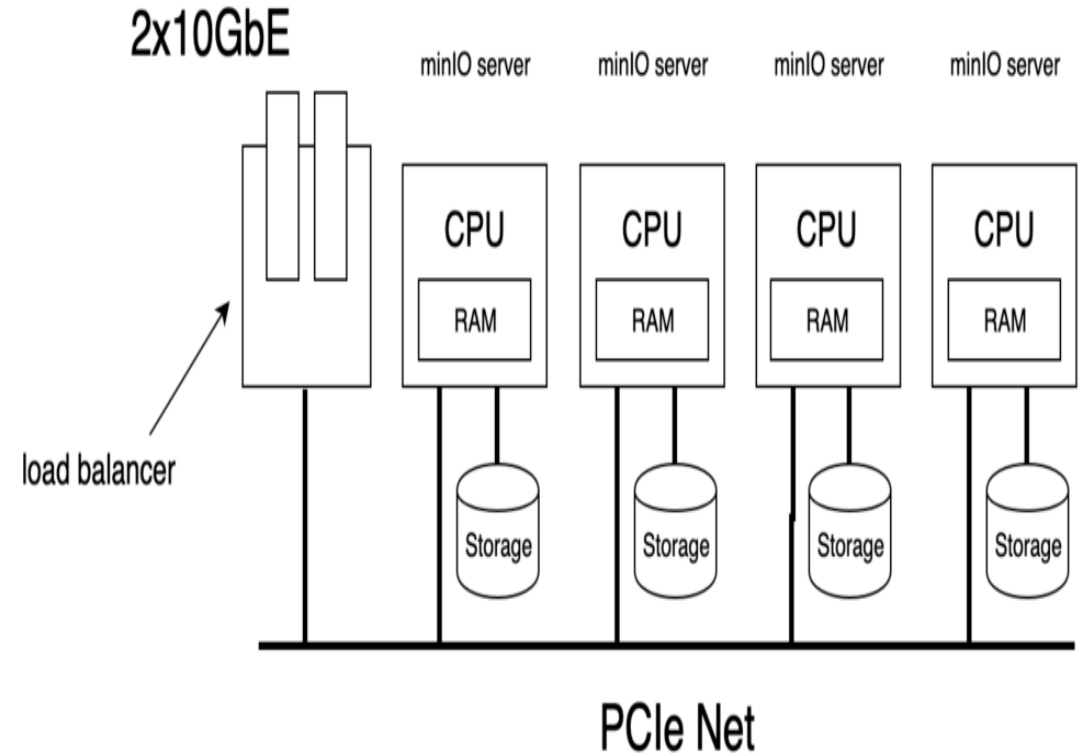
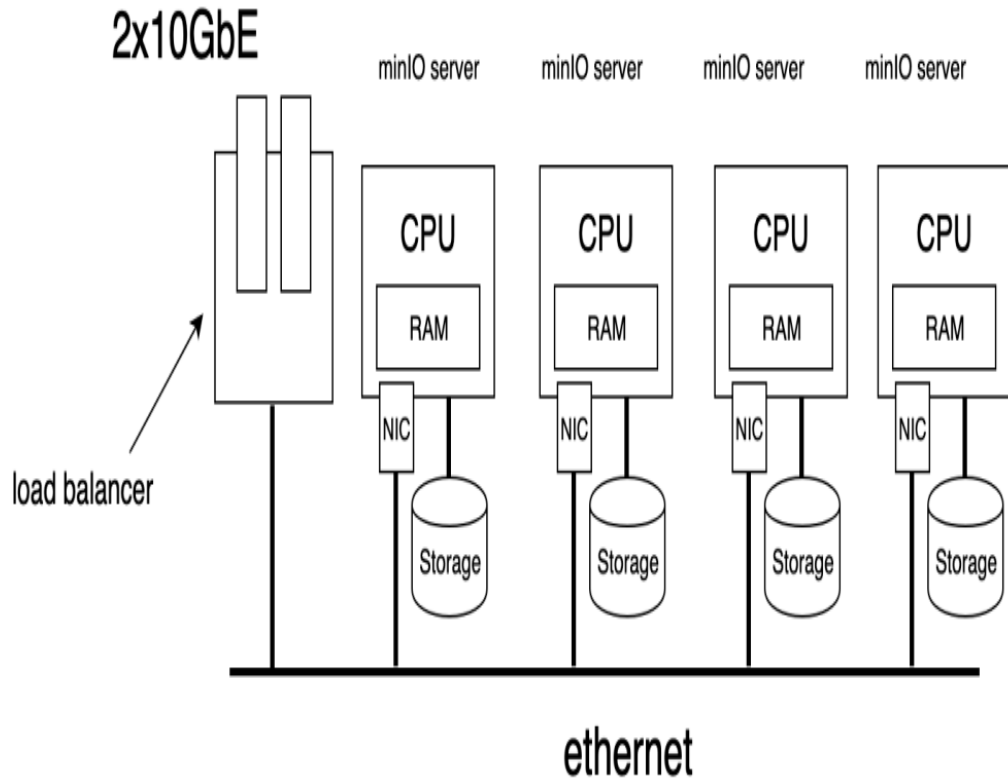
Power

功耗 (瓦)	系统 (含SSD硬盘)	系统 (含HDD硬盘)
PCIe Net + Arm	80	240
Ethernet + x86	600	760



PCIe Net + arm CPU
save lots of energy

Network



Speed Test

Speed Test results: Object Size 8 MB Duration 100 s Retest

GET ↓ **2** GIB/s
293 Obj/s/S

PUT ↑ **1010** MIB/s
132 Obj/s/S

Detailed Results: Download Expand

Speed Test results: Object Size 64 MB Duration 20 s Retest

GET ↓ **2** GIB/s
39 Obj/s/S

PUT ↑ **1** GIB/s
18 Obj/s/S

Detailed Results: Download Expand

Nodes: 4 Drives: 16 Concurrent: 32 MinIO VERSION 2022-07-30T05:21:40Z

Servers	GET	PUT
http://min1.pcie.net:9000	640 MIB/s.	311 MIB/s.
http://min2.pcie.net:9000	610 MIB/s.	228 MIB/s.
http://min3.pcie.net:9000	558 MIB/s.	265 MIB/s.
http://min4.pcie.net:9000	579 MIB/s.	311 MIB/s.

PCIe Net + arm

Speed Test results: Object Size 8 MB Duration 40 s Retest

GET ↓ **3** GIB/s
522 Obj/s/S

PUT ↑ **1** GIB/s
251 Obj/s/S

Detailed Results: Download Expand

Nodes: 4 Drives: 16 Concurrent: 32 MinIO VERSION 2022-08-13T21:54:44Z

Speed Test results: Object Size 64 MB Duration 40 s Retest

GET ↓ **3** GIB/s
61 Obj/s/S

PUT ↑ **1** GIB/s
30 Obj/s/S

Detailed Results: Download Expand

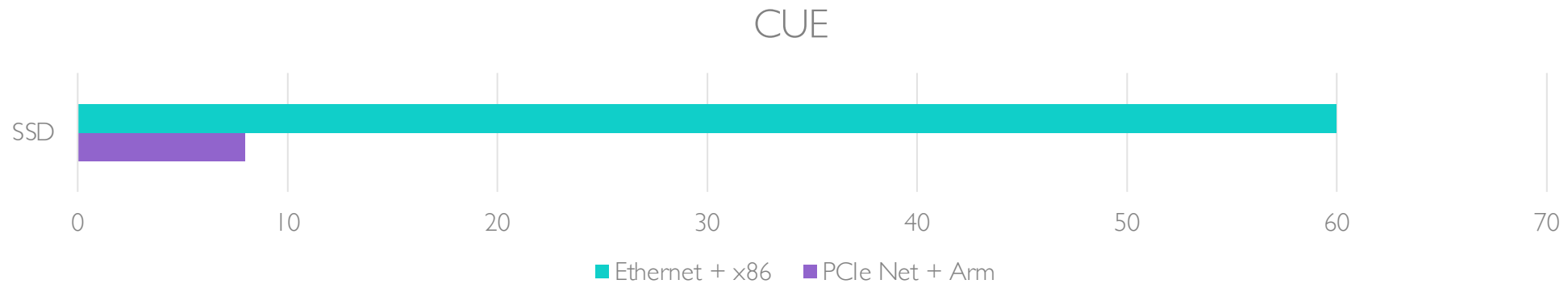
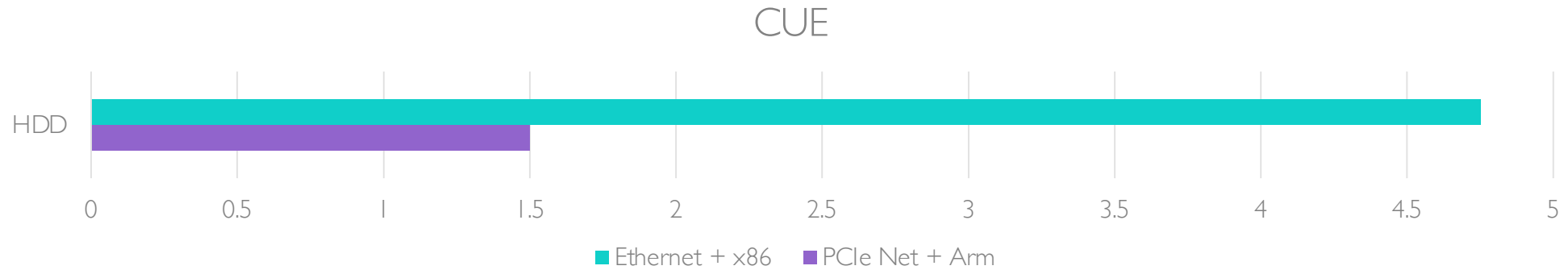
Nodes: 4 Drives: 16 Concurrent: 32 MinIO VERSION 2022-08-13T21:54:44Z

Servers	GET	PUT
http://10.9.8.10:9000	952 MIB/s.	503 MIB/s.
http://10.9.8.3:9000	915 MIB/s.	453 MIB/s.
http://10.9.8.5:9000	941 MIB/s.	442 MIB/s.
http://10.9.8.6:9000	965 MIB/s.	436 MIB/s.

Ethernet + x86

Carbon Usage Effectiveness (CUE)

$$CUE = \frac{\text{Hard Drive TDP} + \text{Bare System TDP} + \text{Networking TDP}}{\text{Hard Drive TDP}}$$

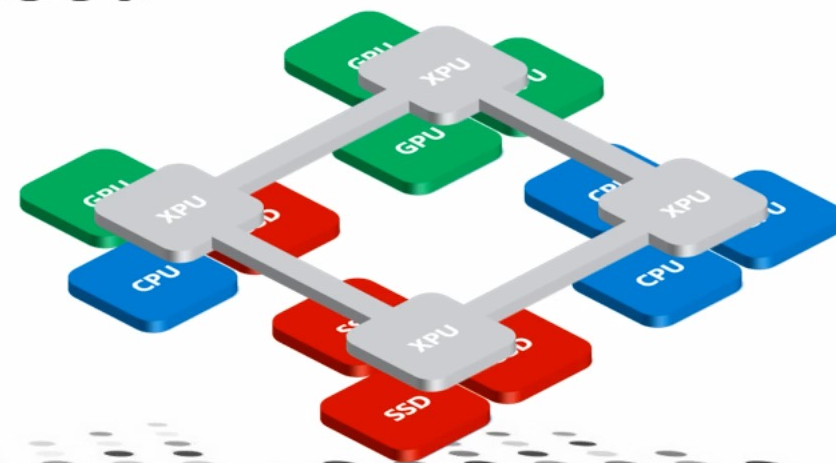


Distributed storage solution based on PCIe Net

All in PCIe: Just Connecting

SOCNOC Co., LTD

<https://www.socnoc.ai>



Thank you !