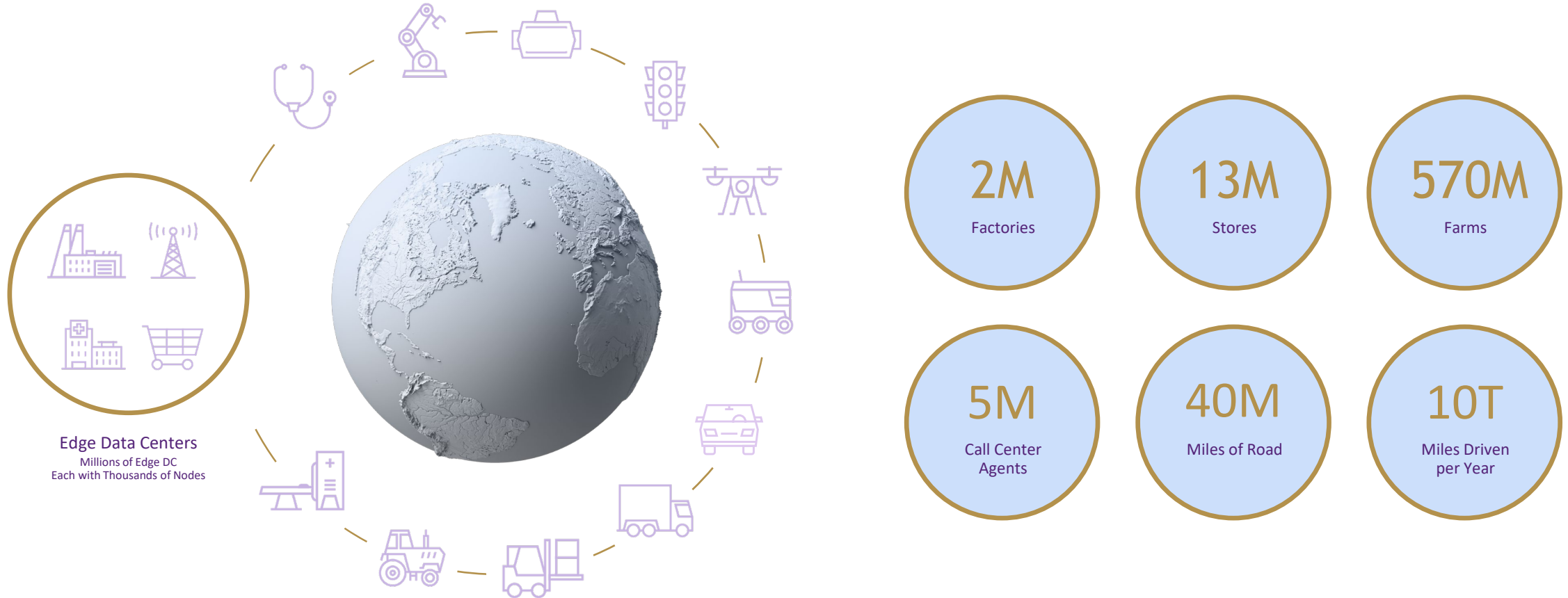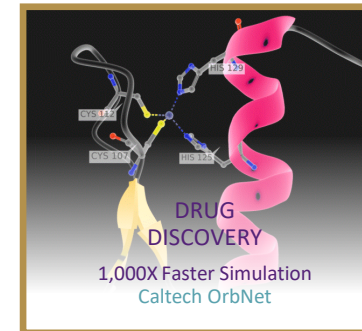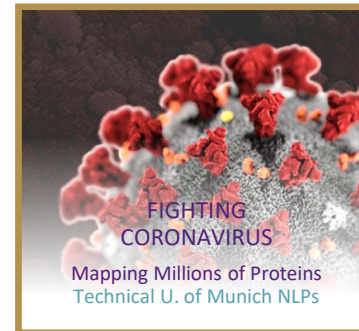# Distributed AI and Computational Storage

Michael Kagan

CTO, NVIDIA

# The New Age of AI

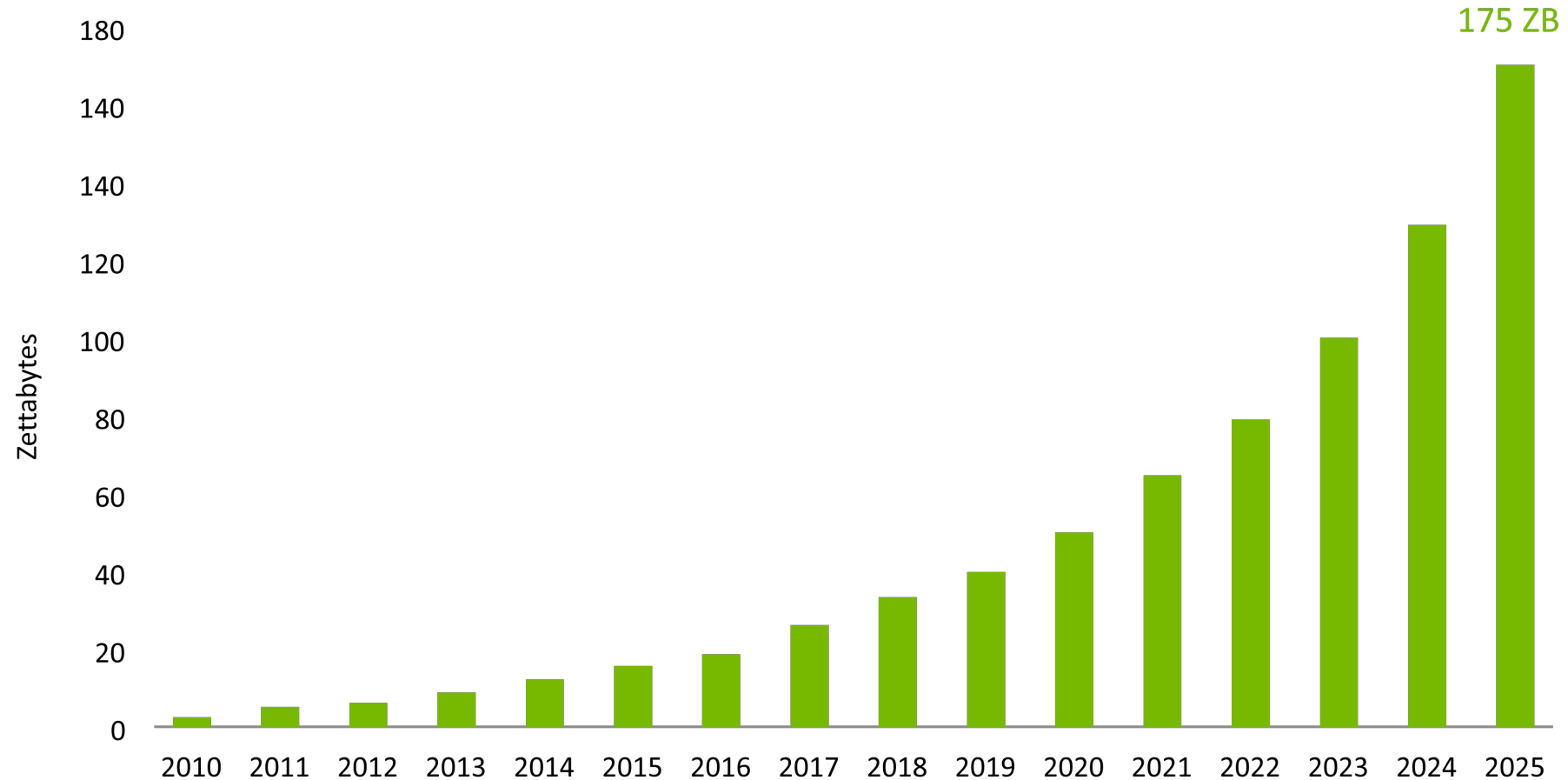Everything is Connected. Data Generated and Consumed Everywhere

**Edge Data Centers**
Millions of Edge DC
Each with Thousands of Nodes

**2M** Factories

**13M** Stores

**570M** Farms

**5M** Call Center Agents

**40M** Miles of Road

**10T** Miles Driven per Year

# Exponential Growth in AI Model Complexity

## EXPLODING MODEL COMPLEXITY
### 30,000X in 5 Years | Now Doubling Every 2 Months

**Petaflop/s - Days** (y-axis: 1.E-03, 1.E-02, 1.E-01, 1.E+00, 1.E+01, 1.E+02, 1.E+03, 1.E+04)

x-axis: 2012, 2014, 2017, 2020

Data points: AlexNet, ResNet, BERT, GPT..., Megatron-GPT2, Megatron-BERT, Turing NLG, GPT-3

**HUMAN-LEVEL READING COMPREHENSION**
24% Better than Avg Human
NVIDIA Megatron BERT

**HUMAN-LIKE CHATBOTS**
Half of Users Preferred AI
Facebook BlenderBot

**DRONE DELIVERIES**
6X Safer Landings
Caltech Neural Lander

**FIGHTING CORONAVIRUS**
Mapping Millions of Proteins
Technical U. of Munich NLPs

**DRUG DISCOVERY**
1,000X Faster Simulation
Caltech OrbNet

**ALL-TERRAIN WALKING ROBOT**
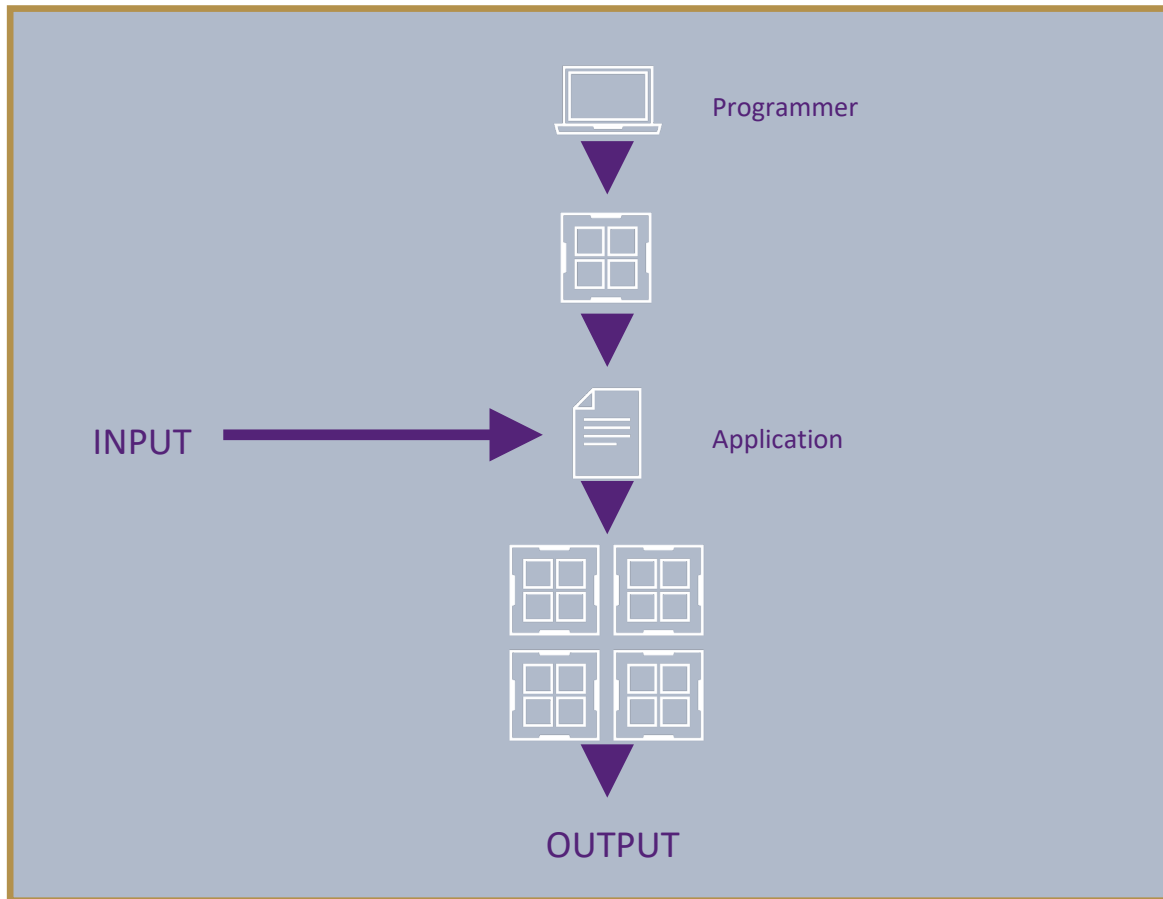On the Fly Surface Adaptation
NVIDIA RL Controller

# Data Grows Exponentially

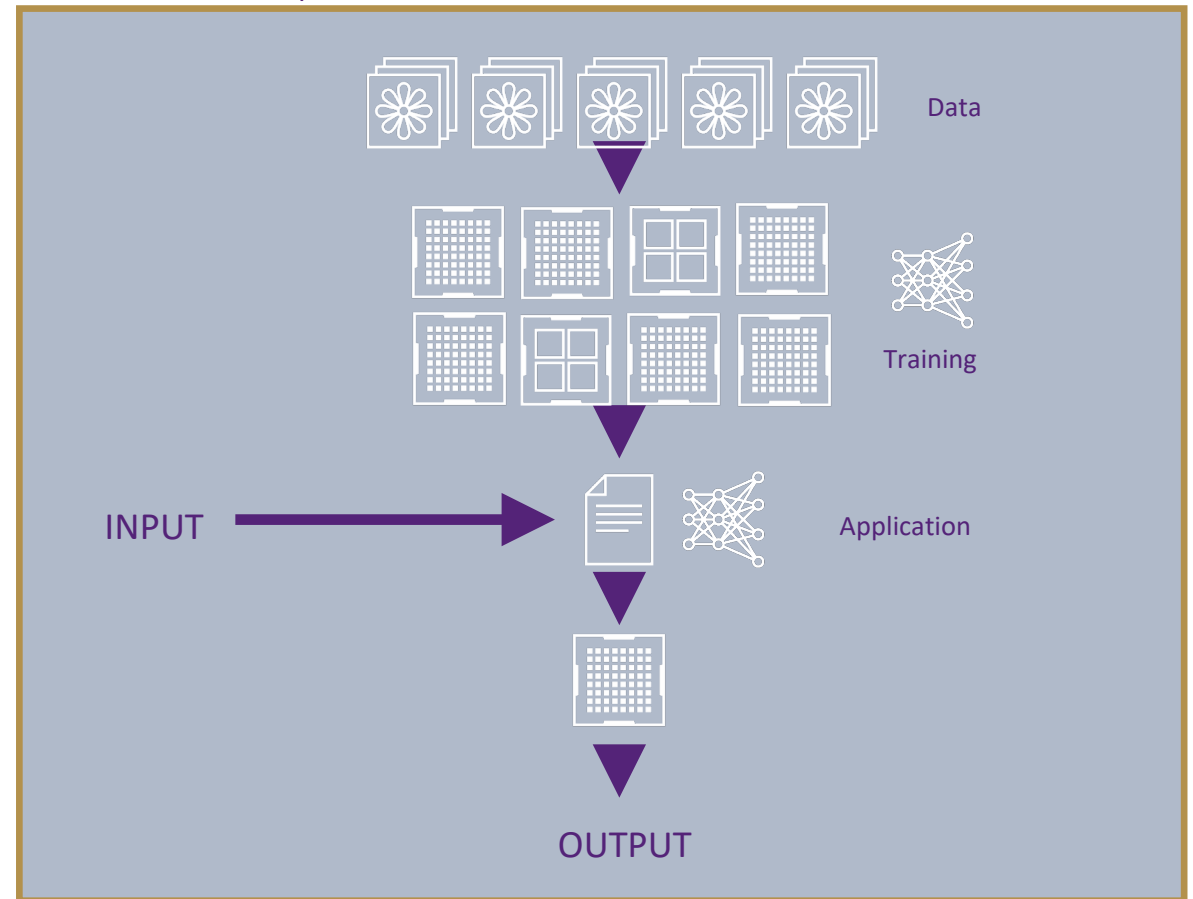More Data → New Models → New Applications → Even More Data



175 ZB

# Software Writes Software

Traditional Model: Programmer uses small computer to create an application that runs on a large centralized computer
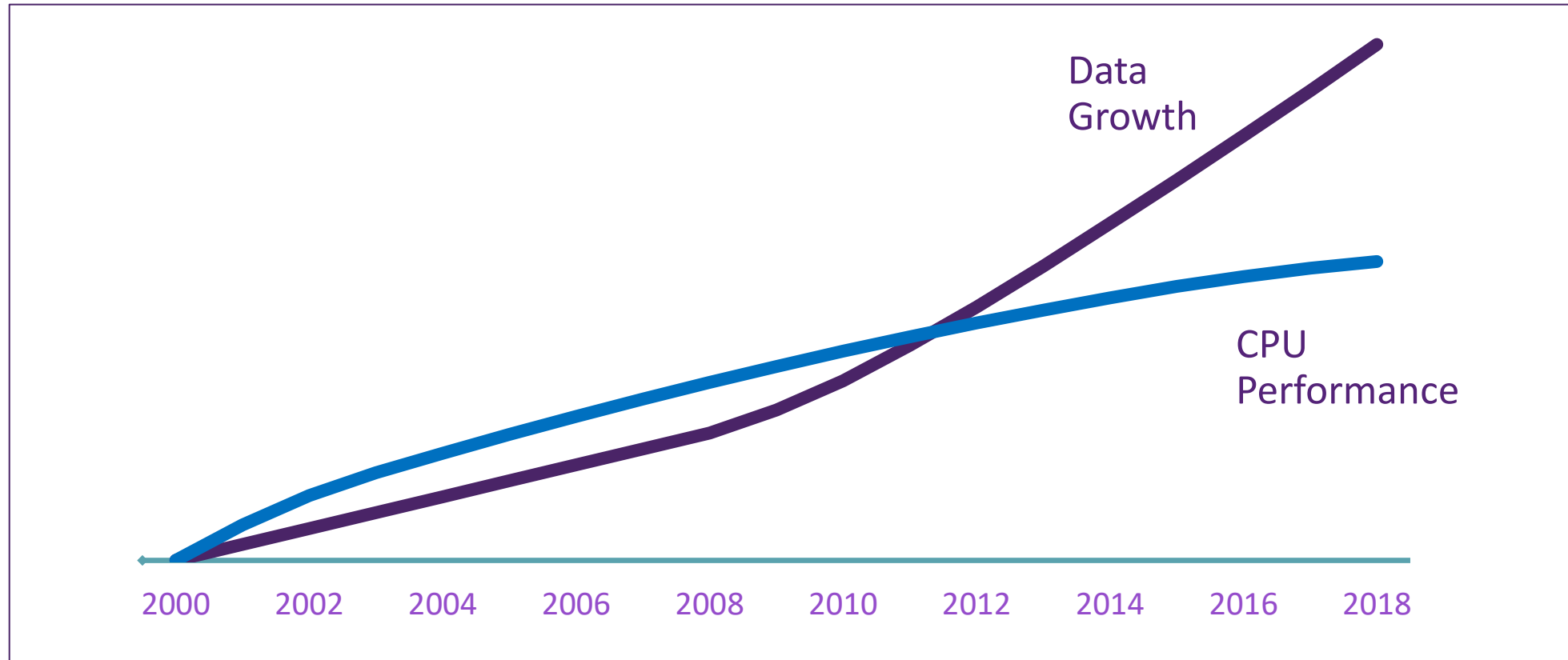
New AI Model: A large computer uses data to train applications that run on distributed small computers
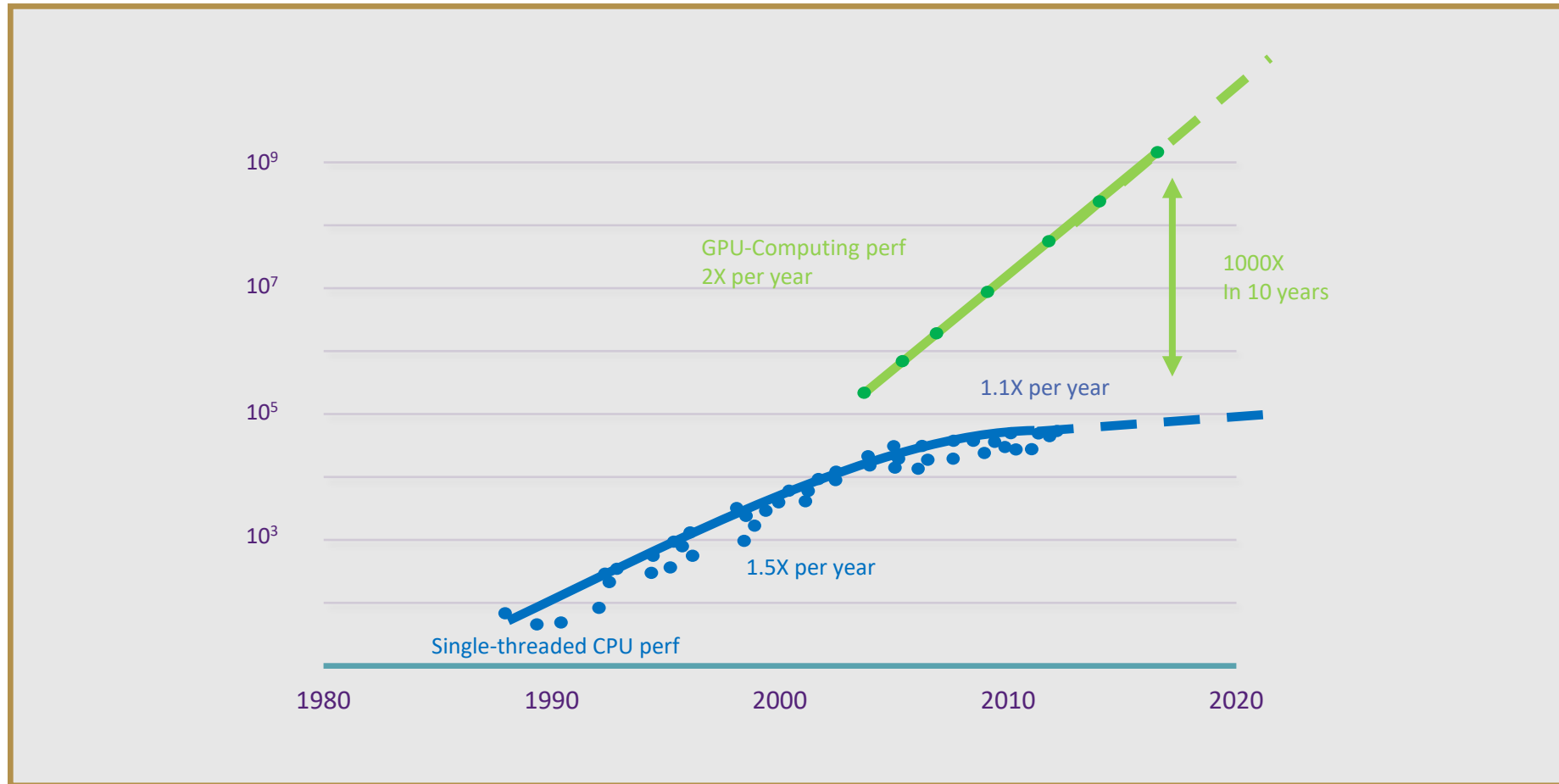
# Now Moore's Law Is Slowing Down

Fails to Keep up with Data Growth, Model Complexity, Memory Speeds, etc.
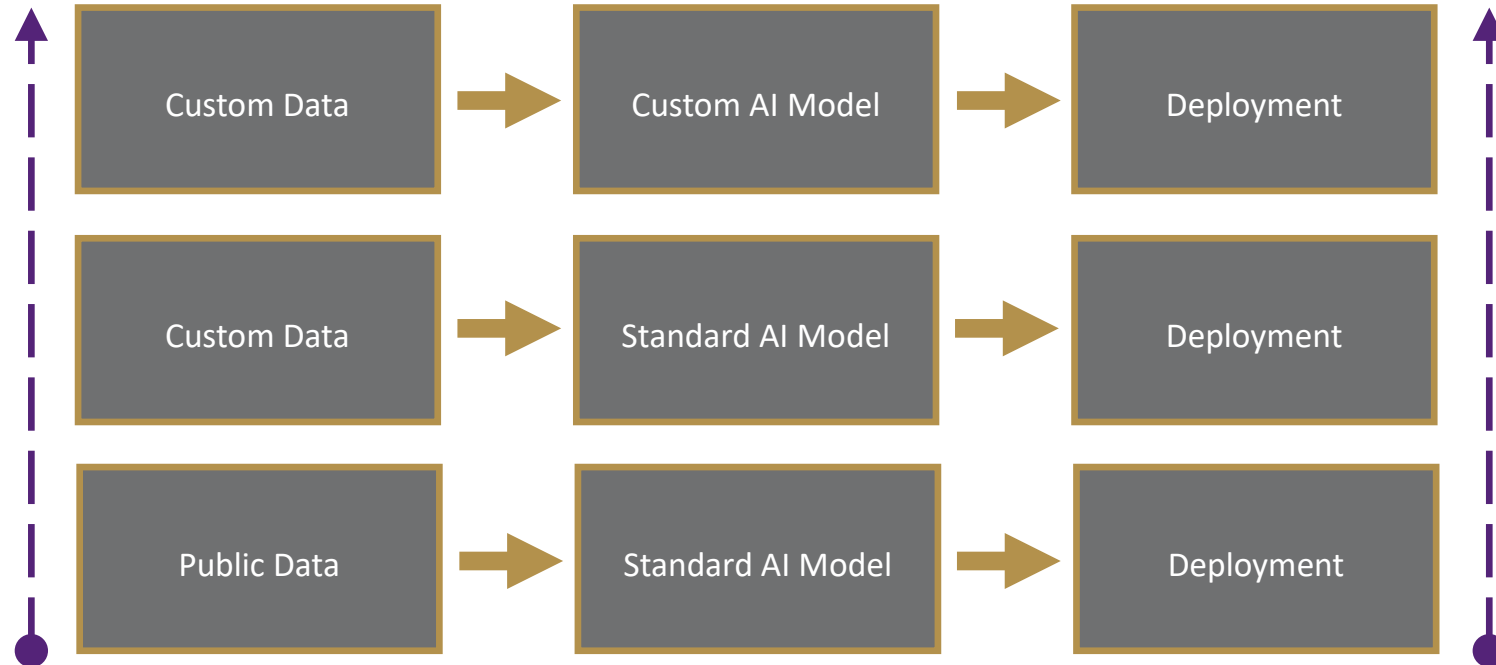


Data Growth

CPU Performance

2000  2002  2004  2006  2008  2010  2012  2014  2016  2018

# New Compute Growth Engines

GPU Becomes the General Processing Unit

GPU-Computing perf
2X per year

1000X
In 10 years

1.1X per year

1.5X per year

Single-threaded CPU perf

$10^9$

$10^7$

$10^5$

$10^3$

1980    1990    2000    2010    2020

# Democratization of AI



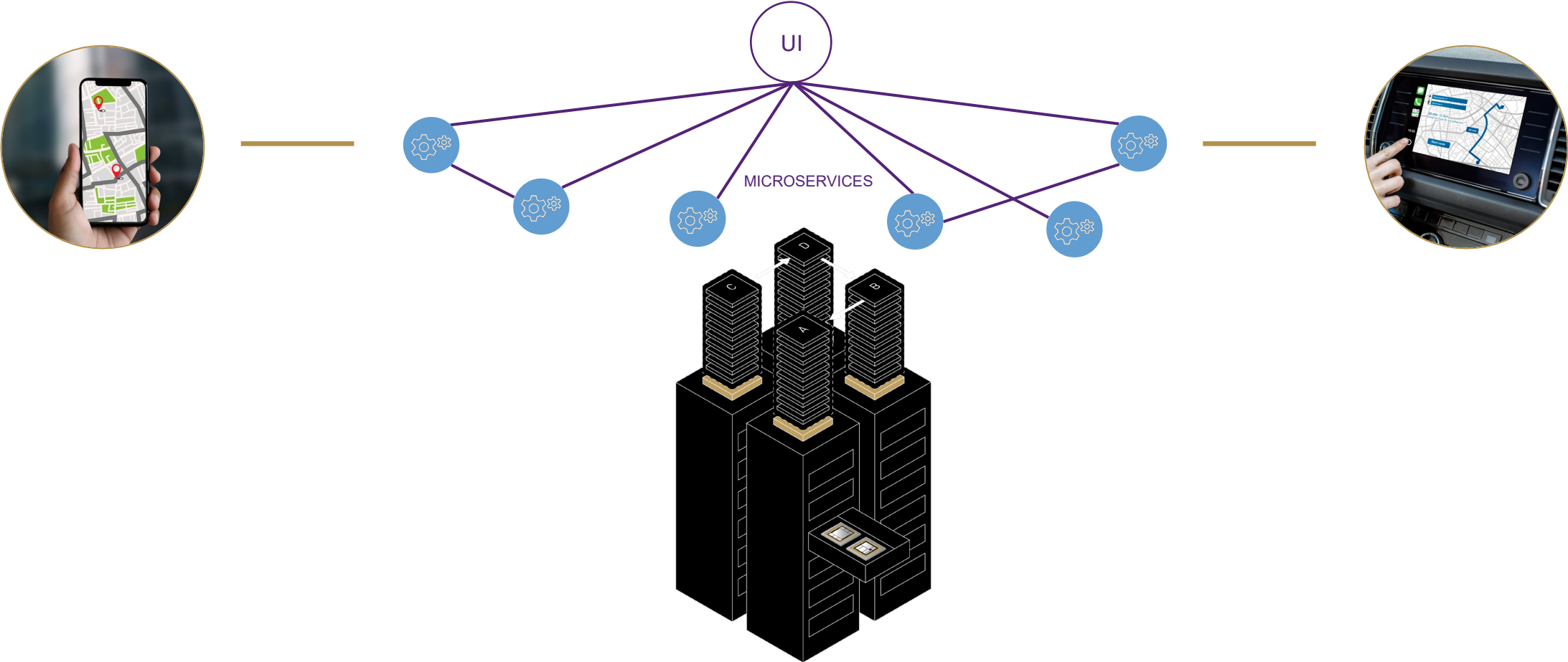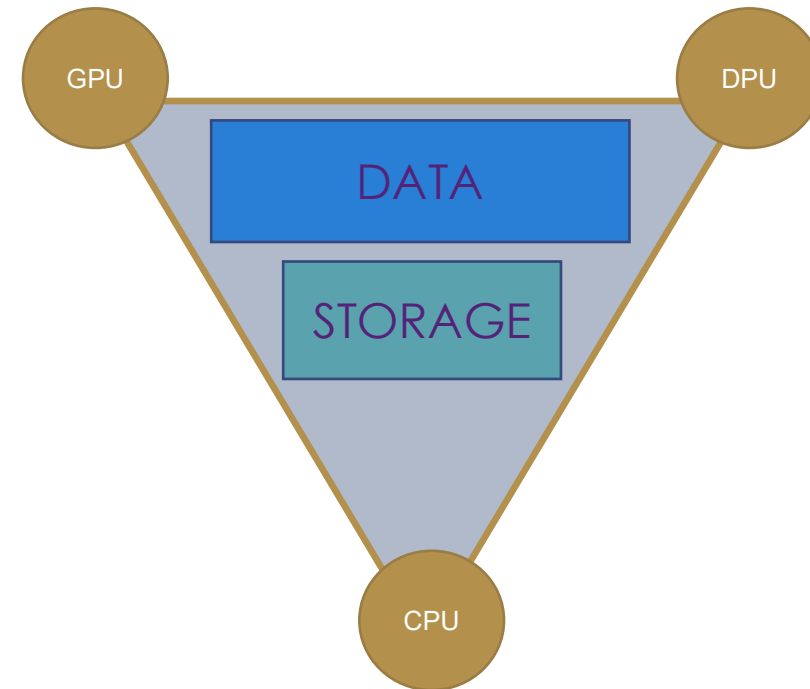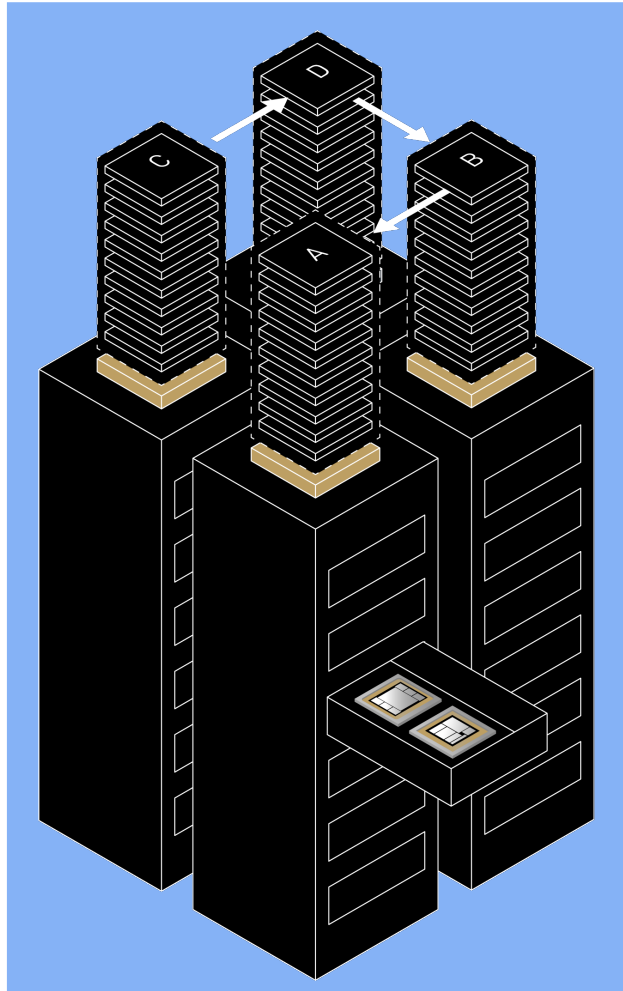| | | |
|---|---|---|
| Custom Data | Custom AI Model | Deployment |
| Custom Data | Standard AI Model | Deployment |
| Public Data | Standard AI Model | Deployment |

# Data Center Is the New Unit of Computing

Distributed Compute with Scale-Out Microservices
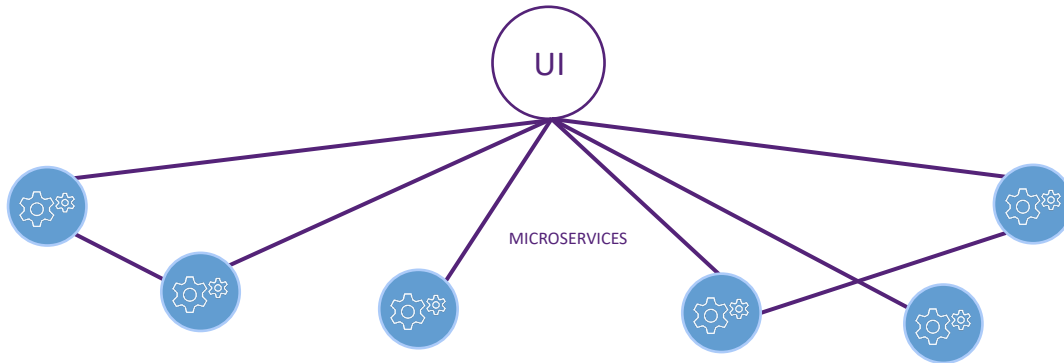


MICROSERVICES

# Powering the New Unit of Computing
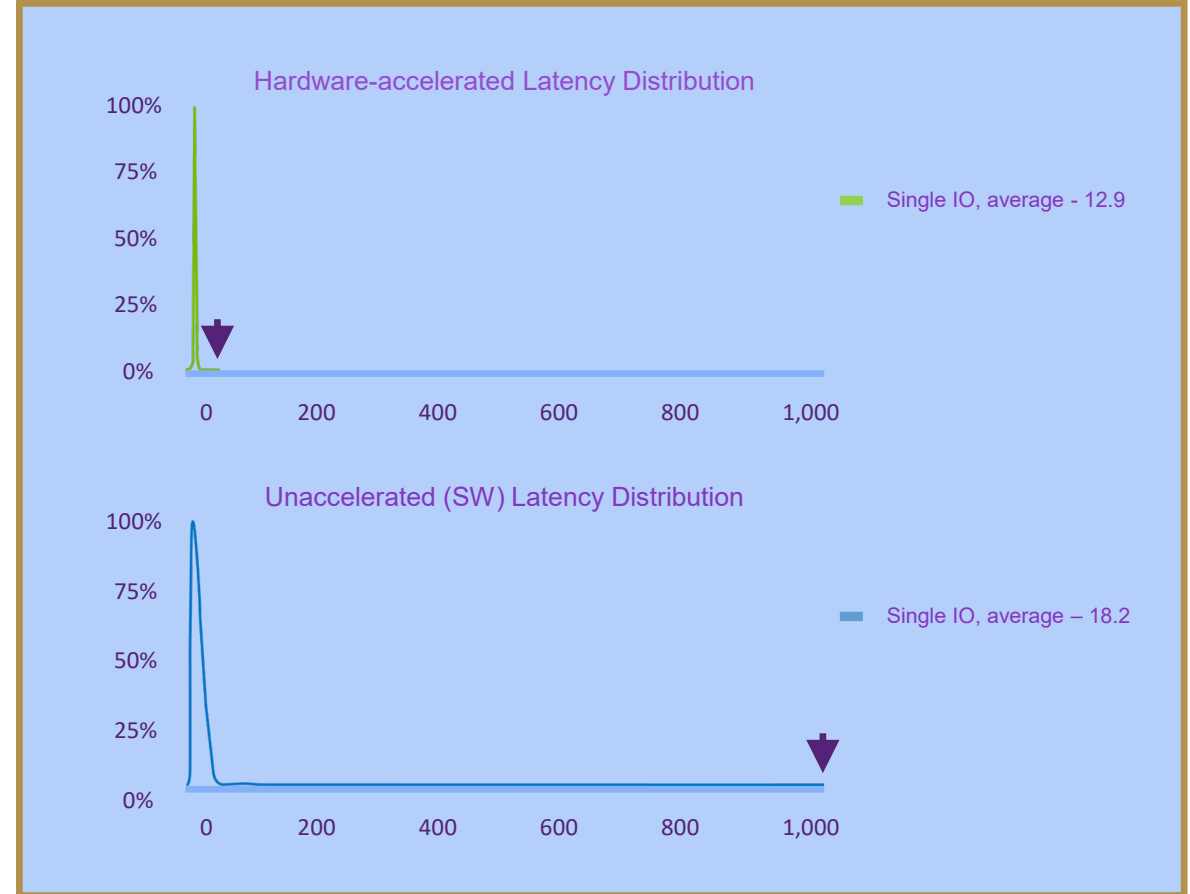
GPU, CPU and DPU are the Engines; Data is the Fuel

# Latency Challenge for Distributed Computing

Average Latencies are Misleading



UI

MICROSERVICES
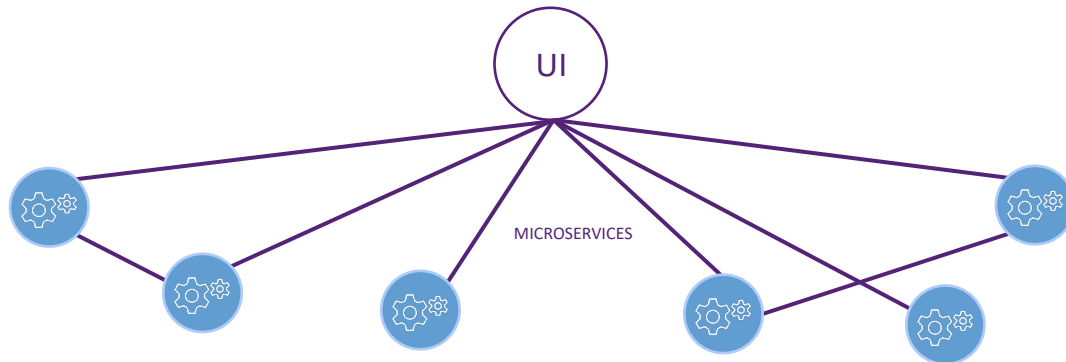
Storage Tail Latency Can be Much Higher than Average Latency When Microservices and Infrastructure Tasks Run Entirely in Software

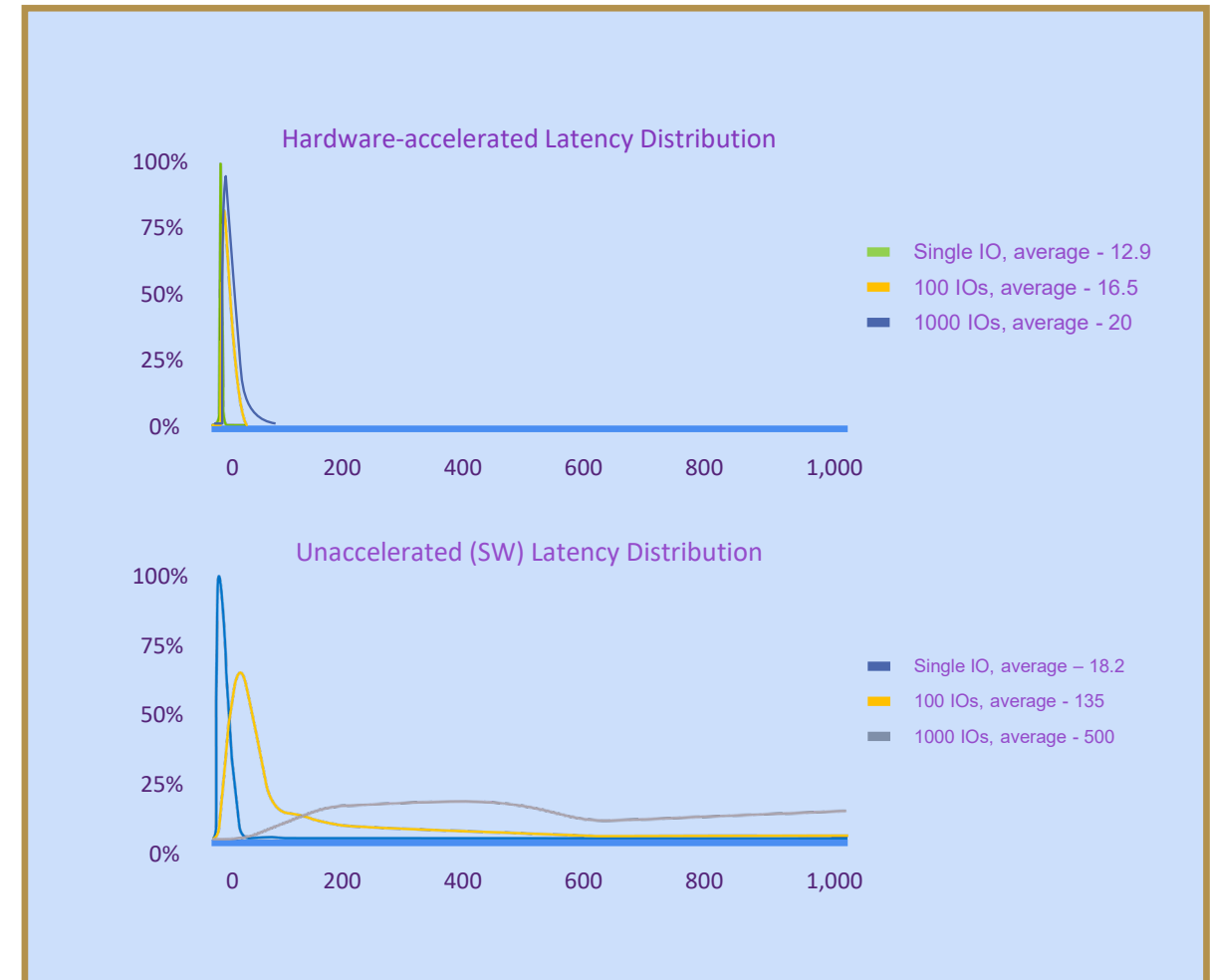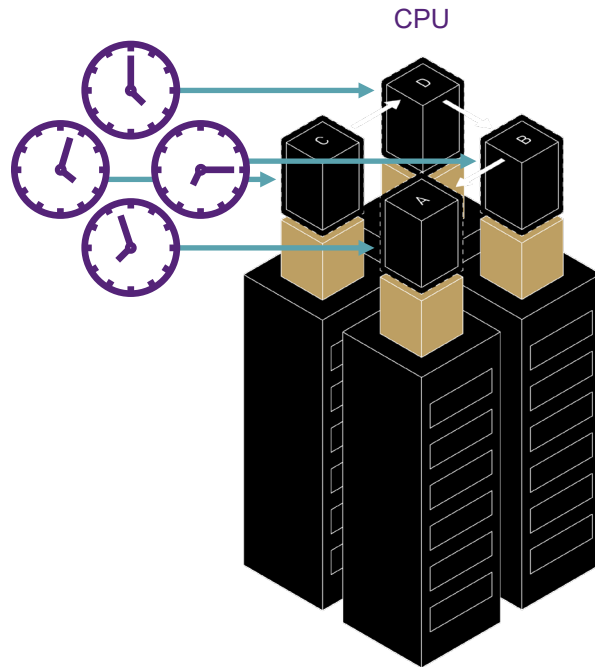### Hardware-accelerated Latency Distribution

Single IO, average - 12.9

### Unaccelerated (SW) Latency Distribution

Single IO, average – 18.2

# The Tail at Scale

Distributed Application Response Time Limited by Tail Latency of the Slowest Microservice

UI

MICROSERVICES

Hardware Accelerations Cuts Storage Tail Latency so Distributed Applications Can Meet SLAs

### Hardware-accelerated Latency Distribution

- Single IO, average - 12.9
- 100 IOs, average - 16.5
- 1000 IOs, average - 20

### Unaccelerated (SW) Latency Distribution

- Single IO, average – 18.2
- 100 IOs, average - 135
- 1000 IOs, average - 500
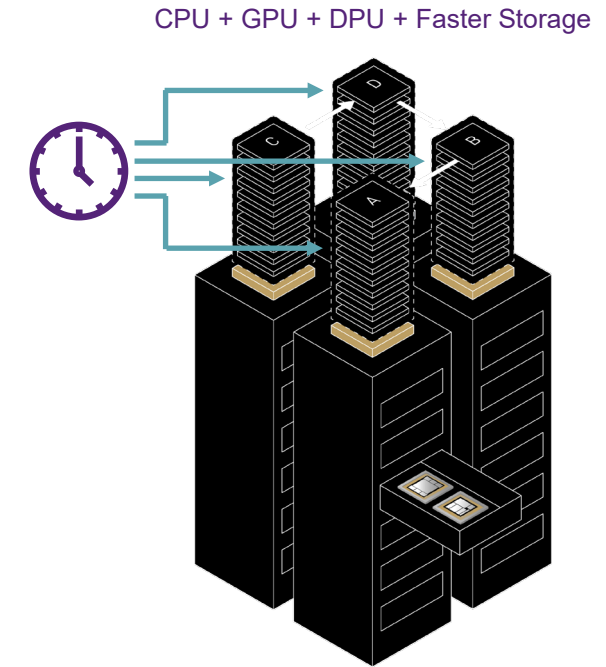
# Accelerating Compute, Infrastructure, Storage

CPU

CPU + GPU

CPU + GPU + DPU + Faster Storage

Disaggregated,
Microservices, Scaled Out

GPU-Accelerated
Computing

GPU- and DPU-Accelerated
Data Center Also Needs Accelerated
Storage Access

# Must Bring Compute and Data Closer Together

## AI Challenges

- More data, more complex models
- AI growth, End of Moore's Law
- Hungry CPUs and GPUs
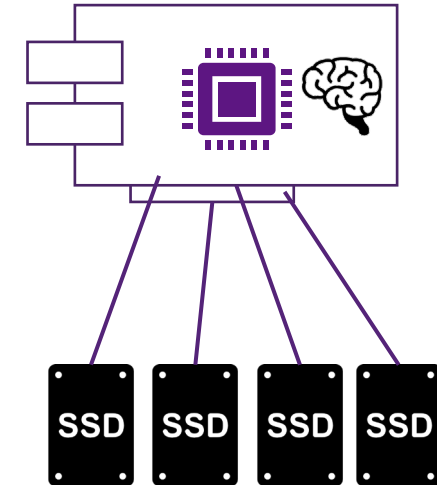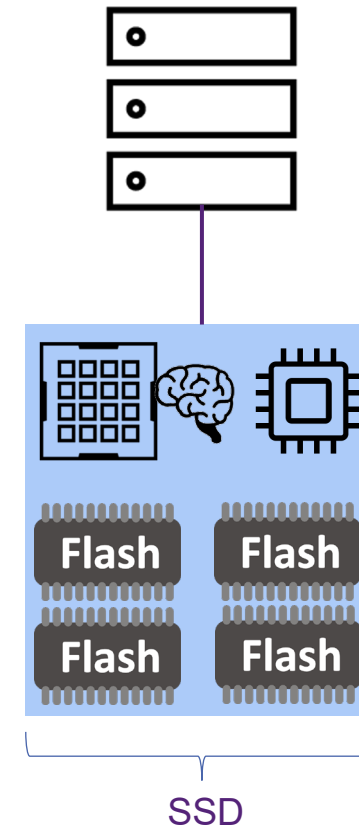- Tail latency, Distributed compute

## Data Solutions

- Computational storage devices
- Smarter storage controllers
- Faster interconnects
- In-network computing

# Computational Storage Devices

- Storage device with embedded compute
- Move results instead of moving data
- Details in earlier sessions

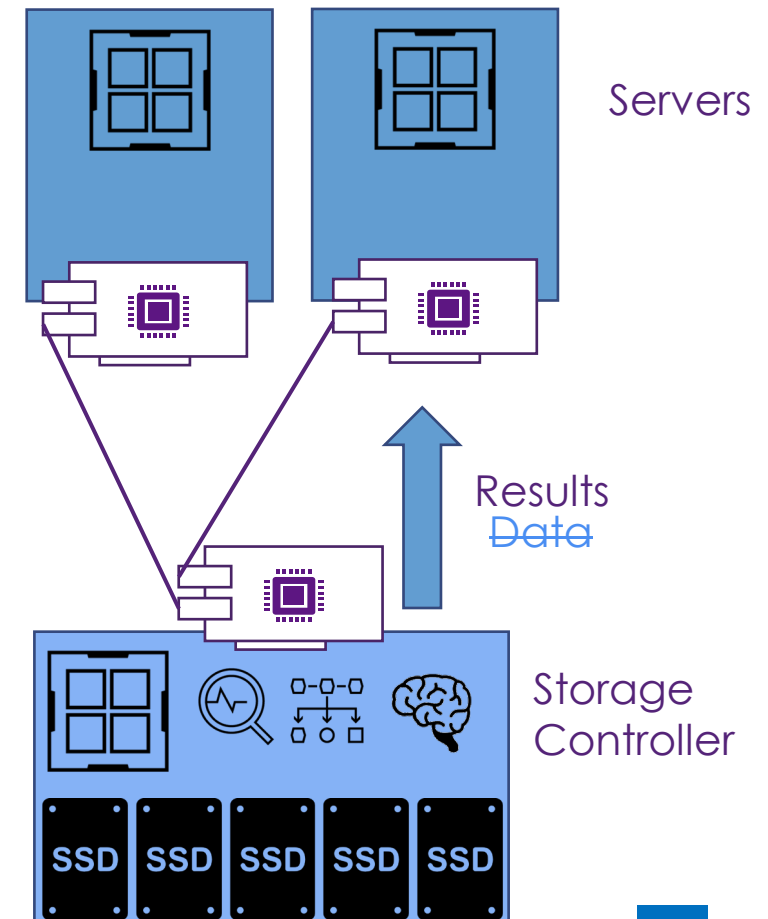Energy-efficient Arm CPUs make it easier to distribute compute

DPU can control SSDs directly

SSD

# Adding Intelligence to Storage Controllers

- Storage controllers perform application compute
  - Not just storage task compute
  - Return results/analysis, not just data

- It moves compute closer to storage
  - Distributing computing to the storage
  - Can debate if it's truly computational storage
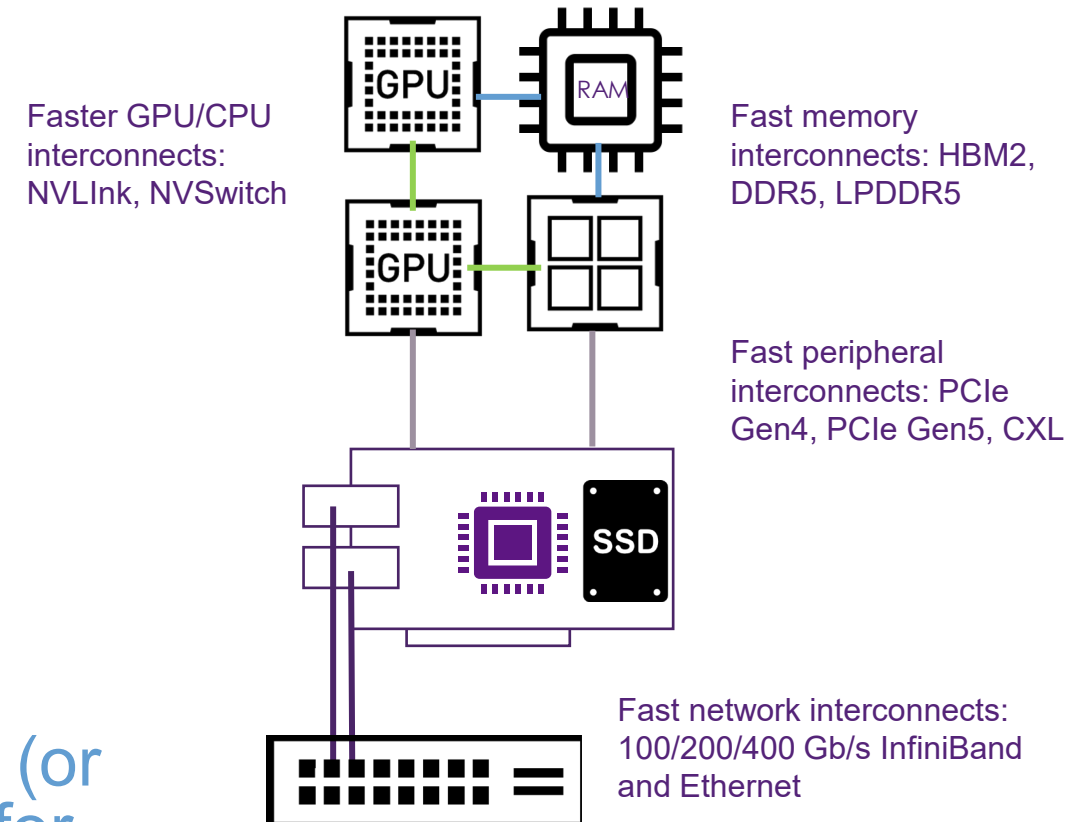  - Solves part of the distributed AI problem

Data Processing Units (DPUs) add compute to storage controllers

Servers

Results
~~Data~~

Storage
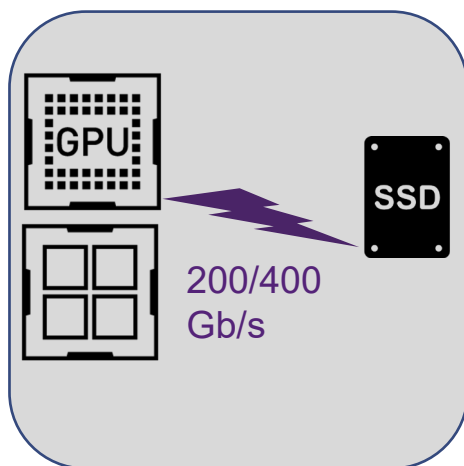Controller

SSD  SSD  SSD  SSD  SSD

# Faster Interconnects Remove Data Bottlenecks

- CPU-to-GPU, GPU-to-GPU
- Memory
- Peripherals
- Networking
- Converged hardware
  - CPU + GPU in one package
  - DPU + Flash in one card

Faster connections from compute to storage (or persistent memory) create distributed proxy for computational storage
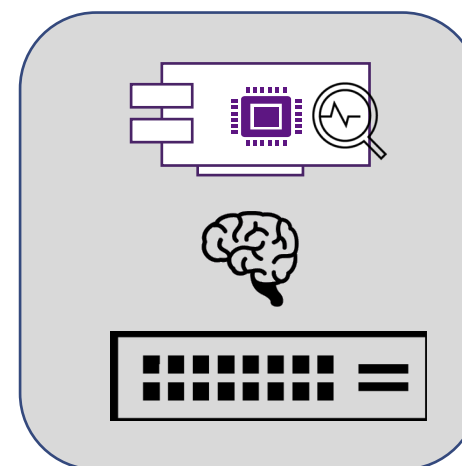
Faster GPU/CPU interconnects: NVLInk, NVSwitch

Fast memory interconnects: HBM2, DDR5, LPDDR5

Fast peripheral interconnects: PCIe Gen4, PCIe Gen5, CXL

Fast network interconnects: 100/200/400 Gb/s InfiniBand and Ethernet

# In-Network Computing

**GPU** **SSD**

200/400 Gb/s

**Network moves data faster**

Faster data access for CPUs, GPUs and storage



Networking | Security
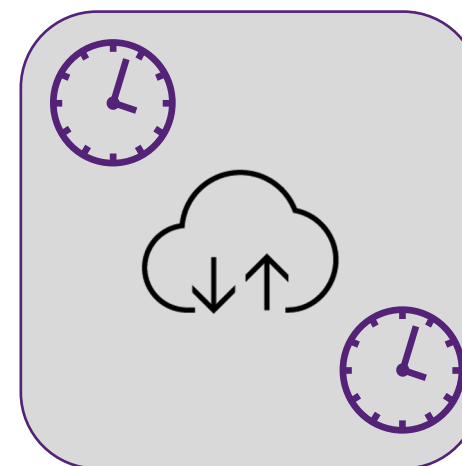Storage | Management

**DPU accelerates infrastructure**

Improved efficiency and security with reduced tail latency



**Network Performs Compute**

Faster AI / ML / HPC problem solving; reduced tail latency



**Time-Synchronized Datacenter**

Faster distributed compute; improved telemetry & security

- AI is everywhere and requires more data

- Data growing faster than CPU power

- New data center distributes compute

- Reducing tail latency becomes critical

Distributed application performance

- Move compute and storage closer together
  - Add compute to storage devices or controllers
  - Accelerate storage access
  - Compute in the storage network

Tail latency

# Thank You

Please visit www.snia.org/pm-summit for presentations

# Title and Abstract

- Date:  Thursday April 22, 11:15-11:45am

- Title:   Why Distributed AI Needs Computational Storage

- Abstract: Artificial Intelligence is increasingly being used in every type of business and industry vertical including finance, telco, healthcare, manufacturing, automotive, and retail. The nature of AI is becoming distributed across multiple nodes in the data center but also across the cloud and edge. Traditional local and networked storage solutions often struggle to meet the needs of AI running on many different types of devices in various locations. Computational storage can solve the challenge of data locality for distributed AI. These solutions include smart storage devices, adding data processing to storage arrays, and deploying new types of compute processors in the storage, next to the storage, or even in the storage network.