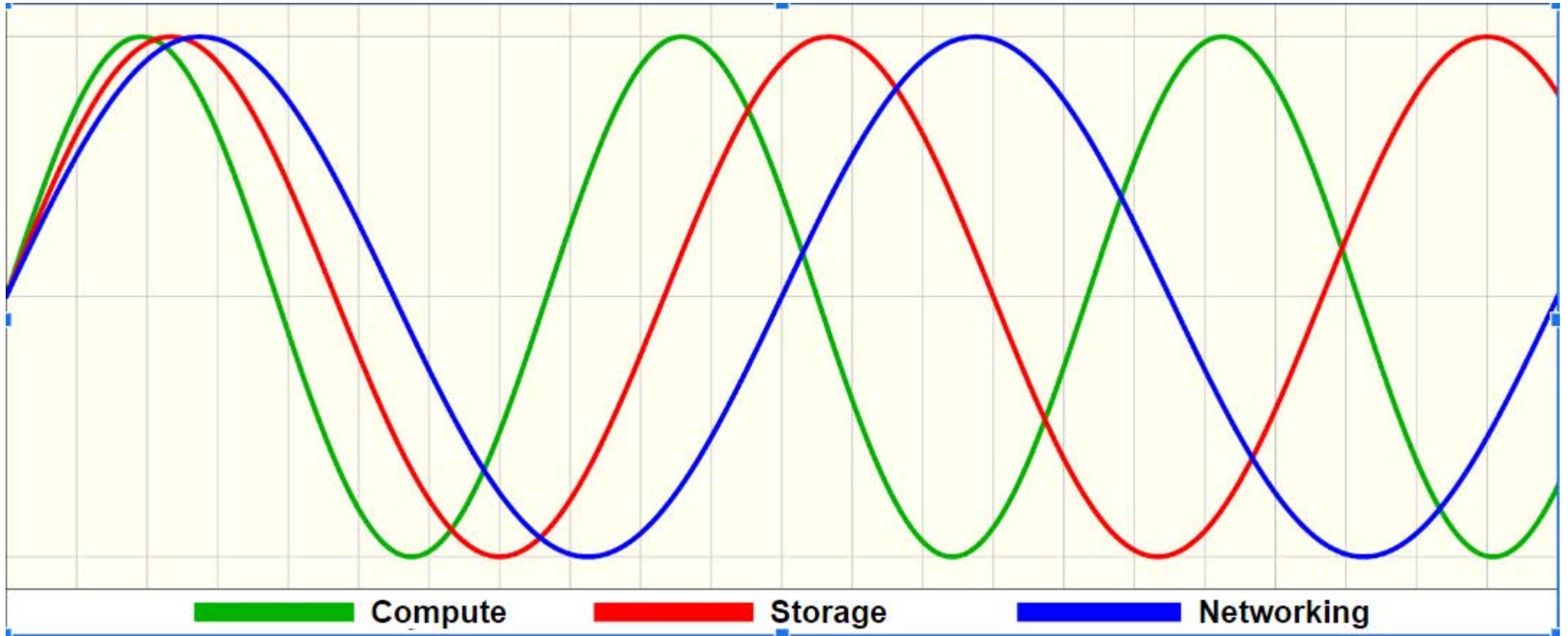# Accelerate Everything.

**How NVM Express and Computational Storage can make your AI Applications Shine!**

Stephen Bates, Chief Technology Officer, MSST 2019

*It's all about the software. Until you reach the limits of the hardware. Then it's all about the hardware [1].*
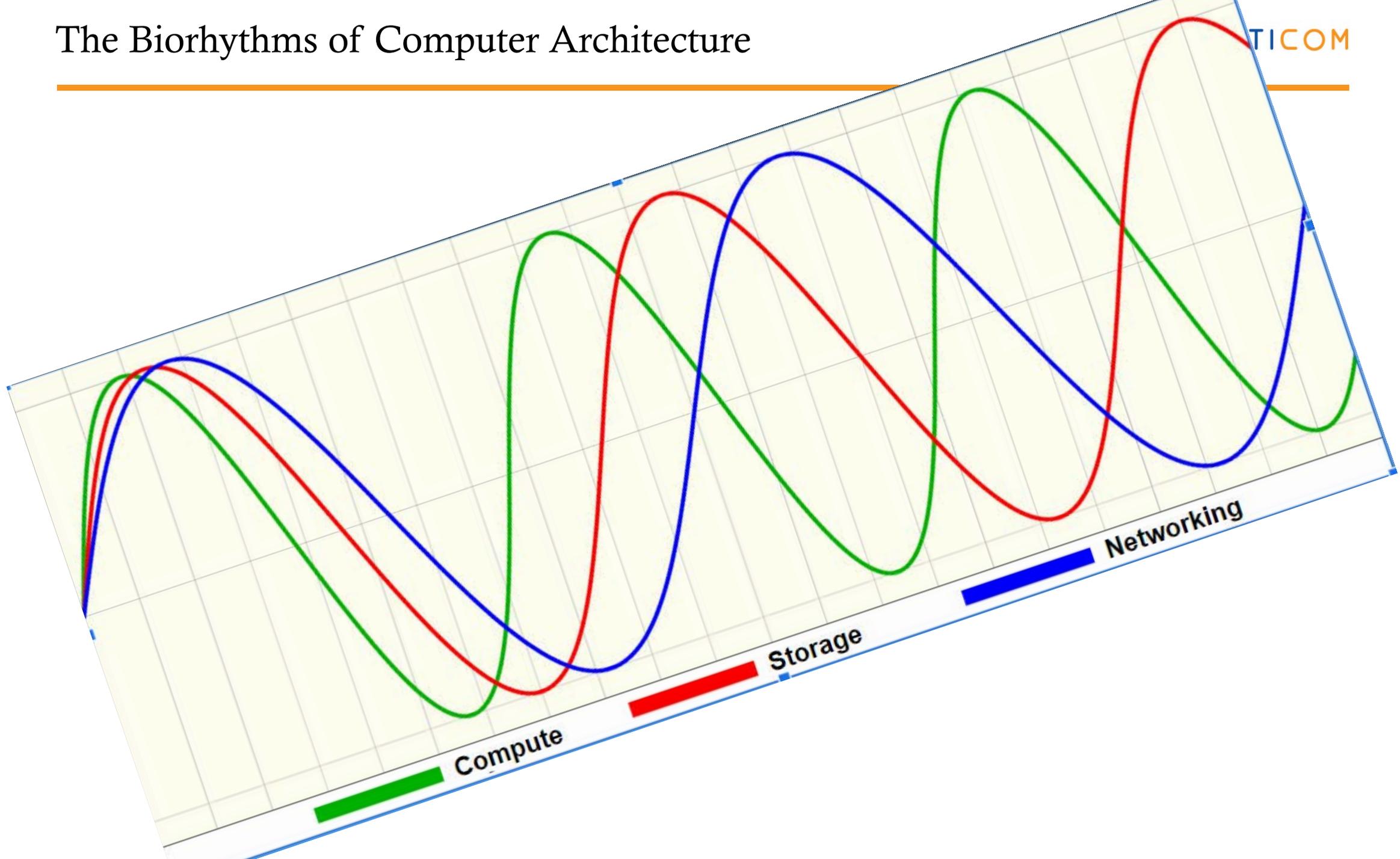
[1] some geek, 2017.

**1955: 5MB, 1 million USD**
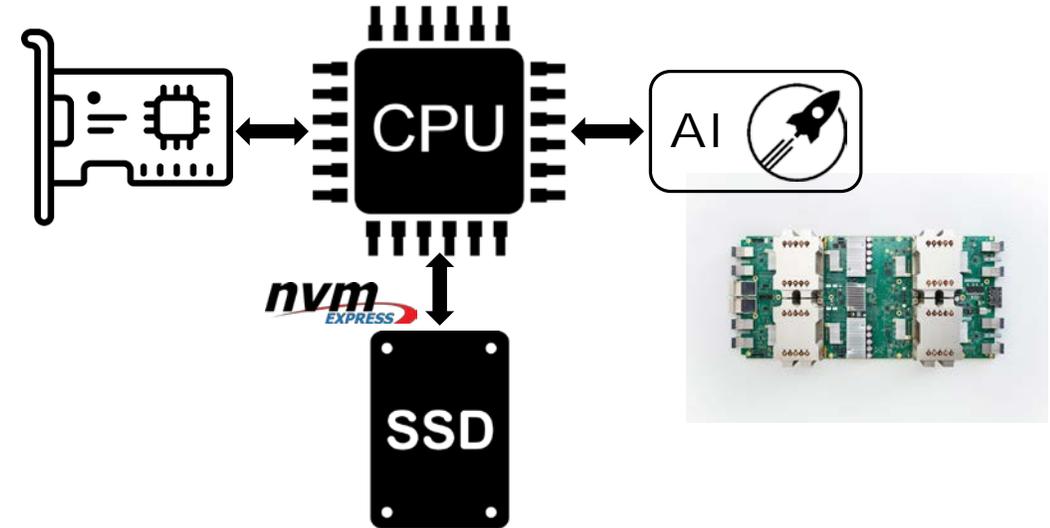


**2018: 1TB (1000000MB), 500 USD**
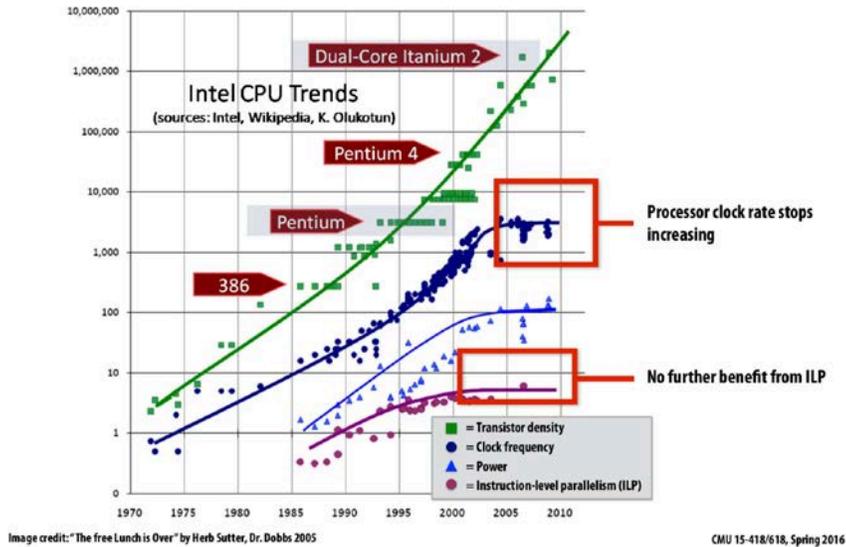
**2019: 8TB, 240MB/s, 5ms**

**2010: 6TB, 2GB/s, 200us**

**HDDs are cheap (0.02/GB) but slow (throughput and latency). SSDs are expensive (0.20/GB) but fast (throughput and latency). For hot-data (which includes AI) SSDs looking more and more attractive.**

# AI has broken some things….



**AI algorithms suck on CPUs**

**AI algorithms shine on Accelerators
(but data movement sucks)**

EIDETICOM

## NVMe for AI: A powerful pairing

NVMe storage capabilities provide the bandwidth and low latency that demanding AI and machine learning applications need to access and manage the massive amounts of data they use.

**John Edwards**

AI and machine learning systems have long relied on traditional compute architectures and storage technologies to meet their

Home › Blog › AI Needs an NVMe-Optimized File System

## AI Needs an NVMe-Optimized File System

Posted on July 16, 2018 by George Crump

Analytics is evolving from big data, machine learning to artificial intelligence. Machine learning is the analysis of data at rest, artificial intelligence (AI) is the analysis of data in real-time. Machine

A next-generation NVMe-native parallel filesystem for accelerating AI workloads (sponsored by WekaIO)

*Liran Zvibel* (WekaIO)
11:55am-12:35pm Thursday, September 6, 2018

Home > GPUs

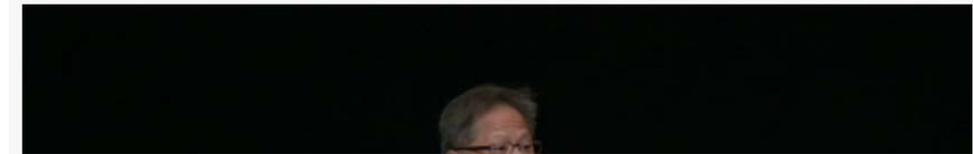## NVIDIA's DGX-2: Sixteen Tesla V100s, 30 TB of NVMe, only $400K

by Ian Cutress on March 27, 2018 2:00 PM EST

**28 Comments**

+ Add A Comment

Posted in GPUs | Systems | Enterprise | NVIDIA | Volta | Servers | HBM2 | GV100 | GTC 2018 | V100 | DGX-2

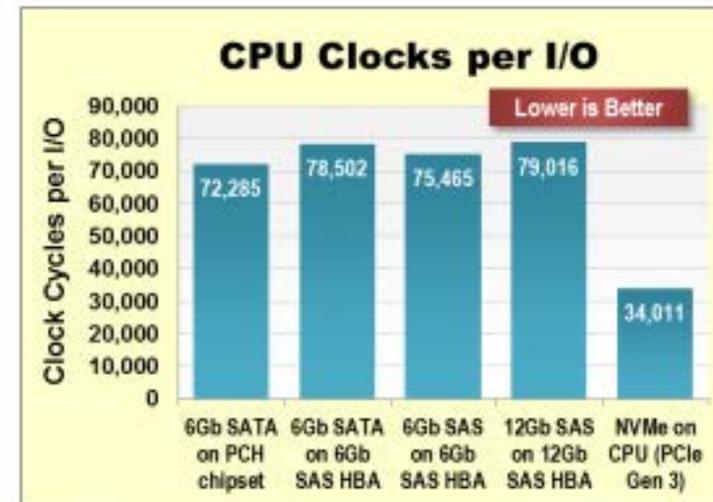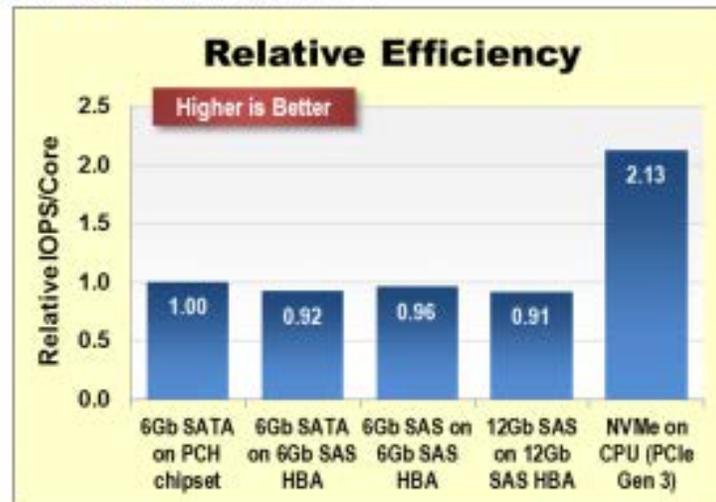EIDETICOM

- **15000 image inferences per second per card (using ResNet50) at 1.3ms latency.**
- **15000 images = 3GB/s (@ 200kB/image)!!**
- **8 cards in a 2U server = 24GB/s (200Gb/s)!**

# EIDETICOM

- Initially a protocol that sat on PCIe and used to talk to Non-Volatile Memory (NVM)
  - Avoided vendor specific solutions for PCIe attached flash.
  - Inherently parallel, which is good for NAND and for multi-core CPUs
    - Out of order execution of commands on a queue.
    - Many queues allowed.
  - Lightweight (less CPU instructions per IO compared to SCSI)
  - Leveraged PCIe which comes directly out of every CPU worth caring about.

**Relative Efficiency**

Higher is Better

| | | | | |
|---|---|---|---|---|
| 1.00 | 0.92 | 0.96 | 0.91 | 2.13 |
| 6Gb SATA on PCH chipset | 6Gb SATA on 6Gb SAS HBA | 6Gb SAS on 6Gb SAS HBA | 12Gb SAS on 12Gb SAS HBA | NVMe on CPU (PCIe Gen 3) |

Relative IOPS/Core

**CPU Clocks per I/O**

Lower is Better

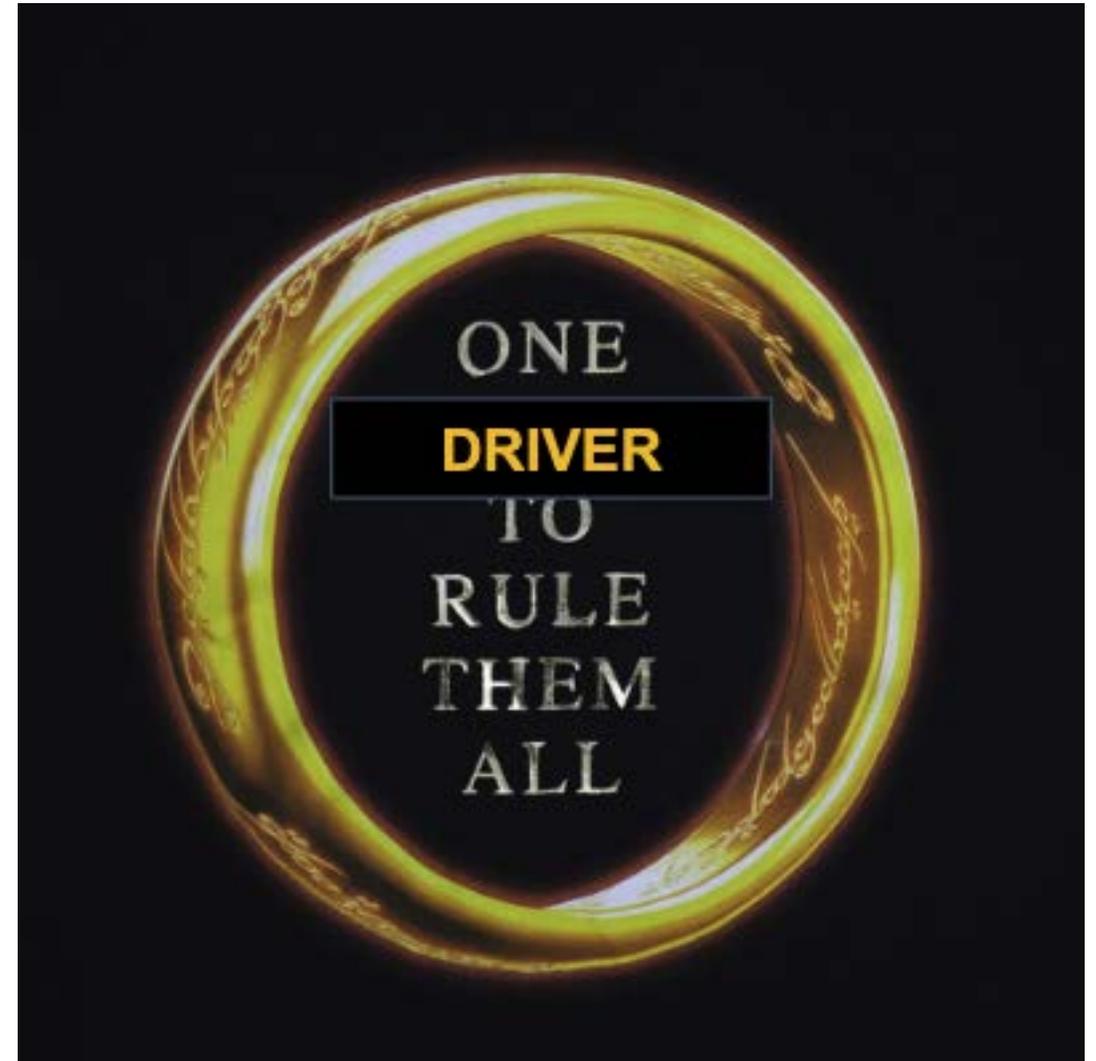| | | | | |
|---|---|---|---|---|
| 72,285 | 78,502 | 75,465 | 79,016 | 34,011 |
| 6Gb SATA on PCH chipset | 6Gb SATA on 6Gb SAS HBA | 6Gb SAS on 6Gb SAS HBA | 12Gb SAS on 12Gb SAS HBA | NVMe on CPU (PCIe Gen 3) |

Clock Cycles per I/O

# NVMe™ is a transport

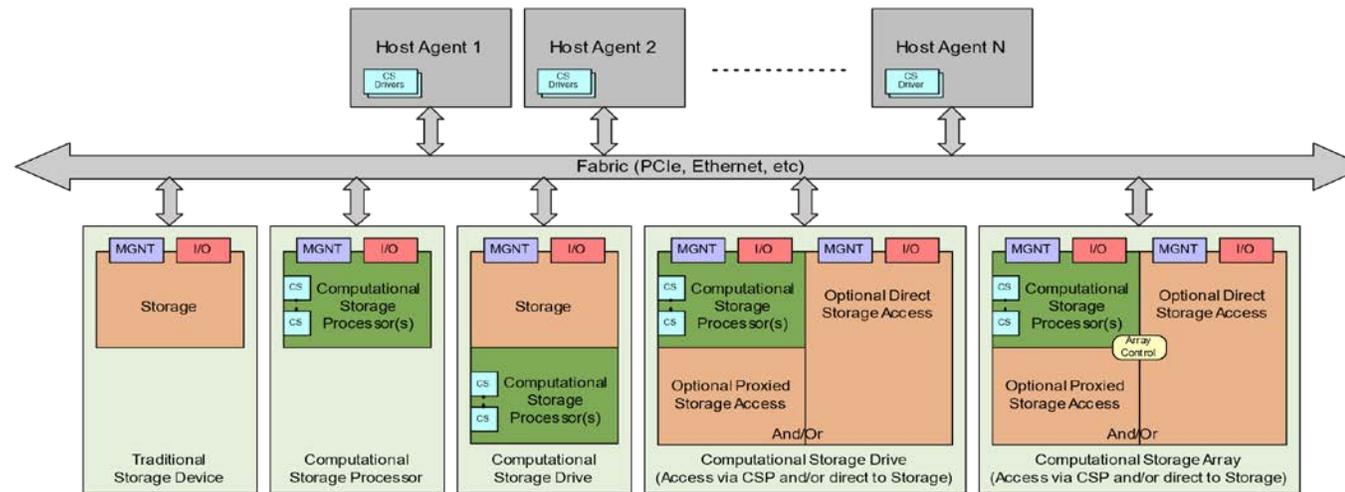Michael Corwell, GM Storage, Microsoft Azure, Dec 5th 2018

# One Driver to Rule Them All?!

- NVMe™ has been incredibly successful as a storage protocol.
- Also being used for networking (NVMe-oF™ and things like AWS Nitro and Mellanox's Sexy NVMe Accelerator Platform (SNAP)).
- Why not extend NVMe to compute and make it *the one driver to rule them all?*

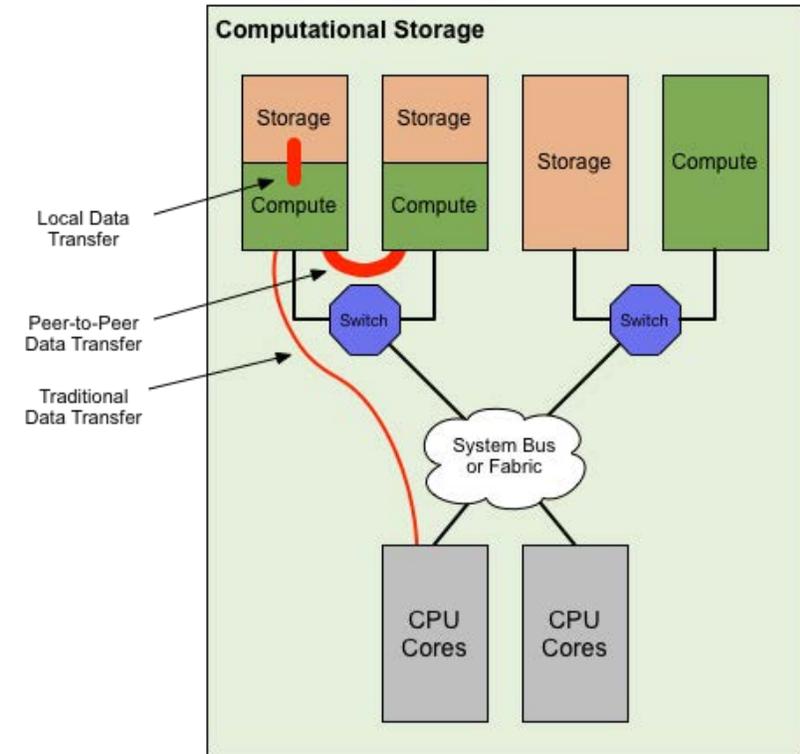# A New Product Category

◆ **Computational Storage Device (CSx)**



◆ **Computational Storage Drive (CSD)**

◆ **Computational Storage Processor (CSP)**

◆ **Computational Storage Array (CSA)**

# What is Computational Storage?
# SNIA has Defined the Following

**Computational Storage Drive (CSD)**: A component that provides persistent data storage and computational services

**Computational Storage Processor (CSP)**: A component that provides computational services to a storage system without providing persistent storage

**Computational Storage Array (CSA)**: A collection of computational storage drives, computational storage processors and/or storage devices, combined with a body of control software

# 40+ Participating Companies
# 128+ Individual Members

SNIA™ | COMPUTATIONAL STORAGE

EIDETICOM

NGD systems
*Bringing Intelligence to Storage*

SAMSUNG

ScaleFlux™

arm

CALYPSO Systems

GIGAIO

inspur

KALRAY

Lenovo

MARVELL

Micron

NetApp

NETINT

NYRIAD

ORACLE

Red Hat

SK hynix

SiliconMotion

TOSHIBA

Western Digital

XILINX

BROADCOM

DELL EMC

FADU

HITACHI

HUAWEI

IBM

intel

Mellanox TECHNOLOGIES

PURESTORAGE

Microsemi
a MICROCHIP company

NEC

SUPERMICRO

SEAGATE

vmware

# Why NVMe™?

- Accelerators require:
  - Low latency
  - High throughput
  - Low CPU overhead
  - Multicore awareness
  - Management at scale
  - QoS awareness

- NVMe provides:
  - Low latency
  - High throughput
  - Low CPU overhead
  - Multicore awareness
  - Management at scale
  - QoS awareness

## Real question is "Why not NVMe?"

**EIDETICOM**

**NoLoad® CSP U.2**

- Standard U.2 SSD form-factor: Utilizing SFF-8639 connector.

**NoLoad® CSP Alveo**

- Standard GPU form-factor: x16 PCIe
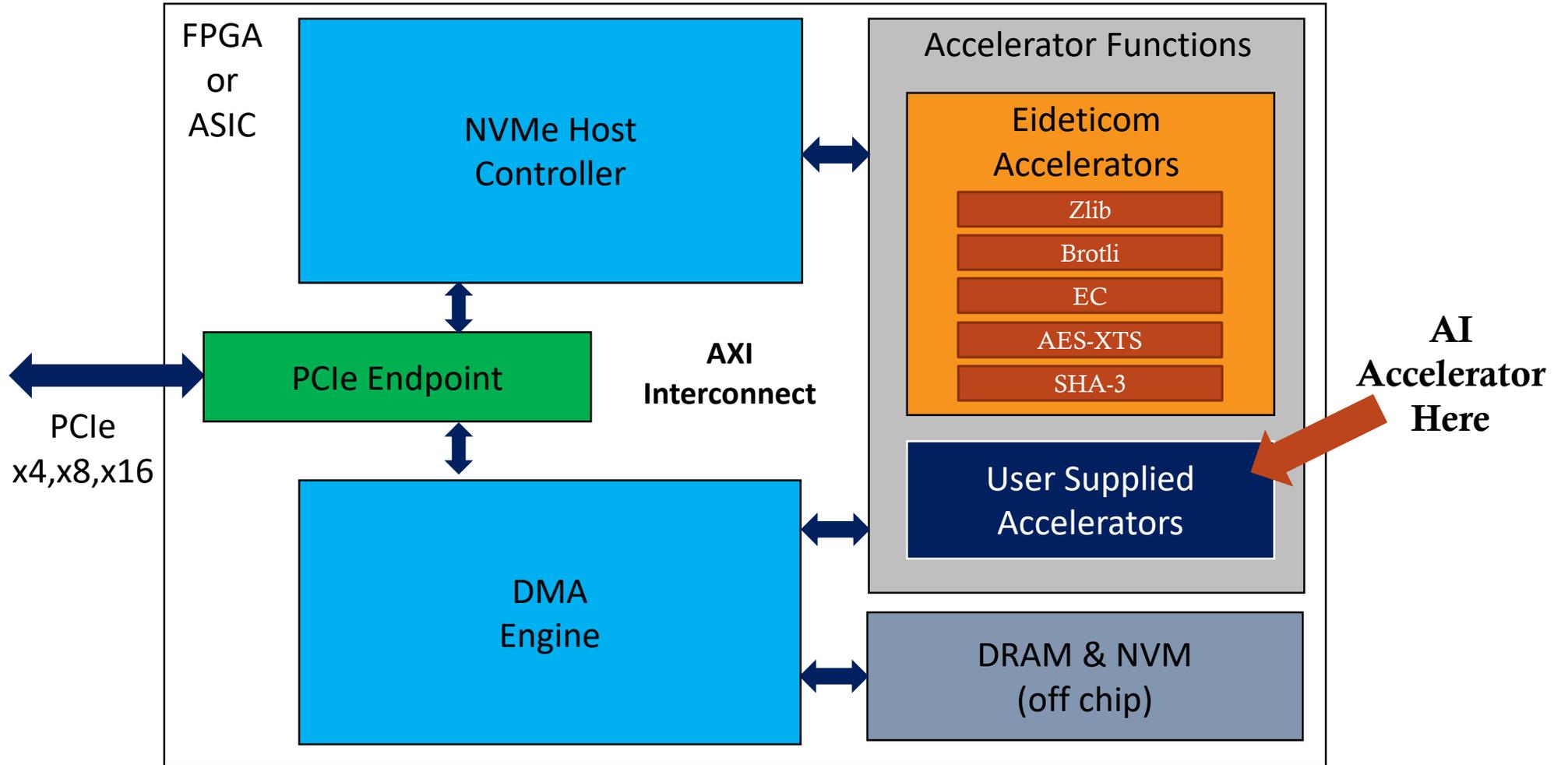- Deployed on Xilinx Alveo U200, 250 or U280

**PCIe Gen4**

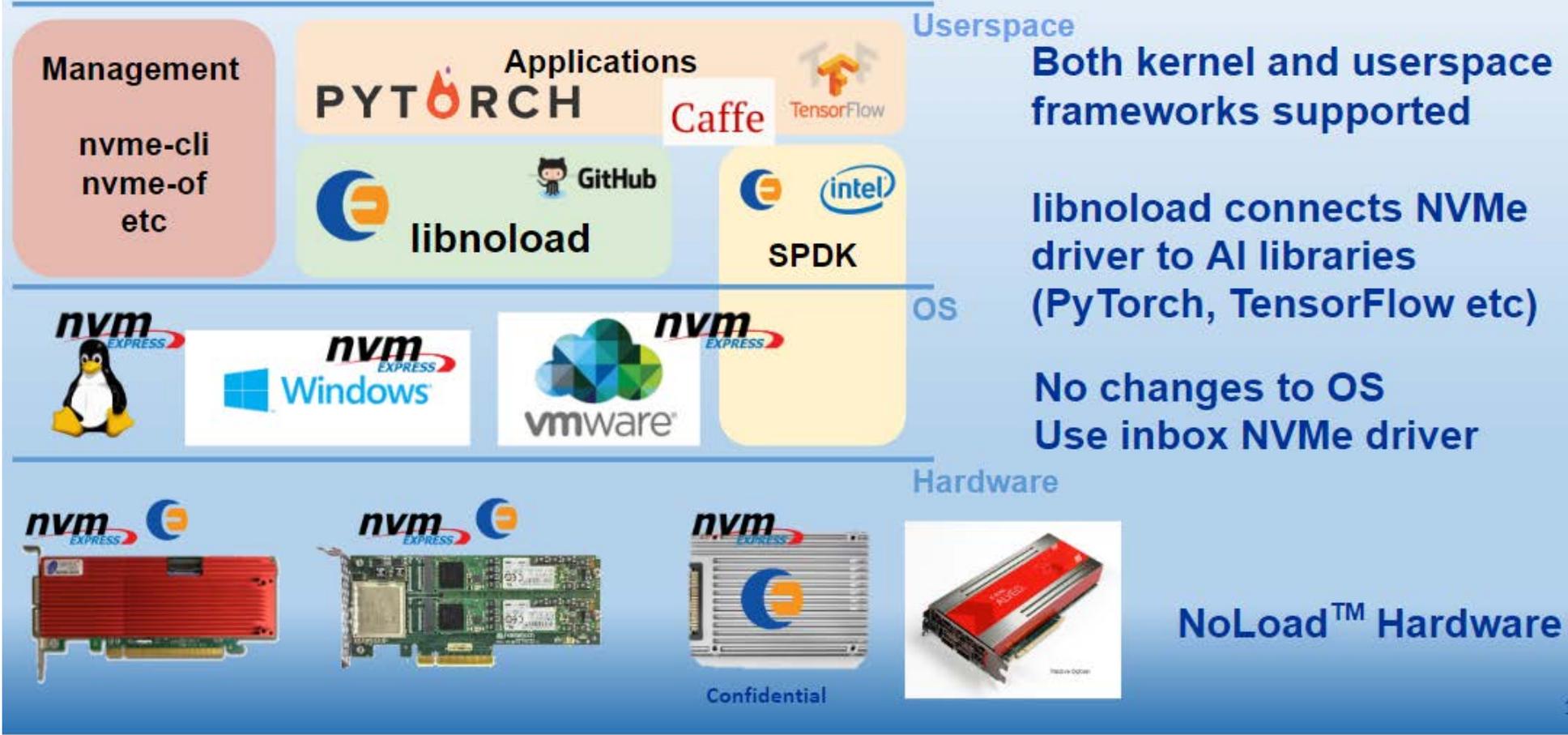- 16GB/s of data ingestion/egestion.

**Eideticom NoLoad IP:**

- NVM Express end-point
- Storage and Compute Accelerators
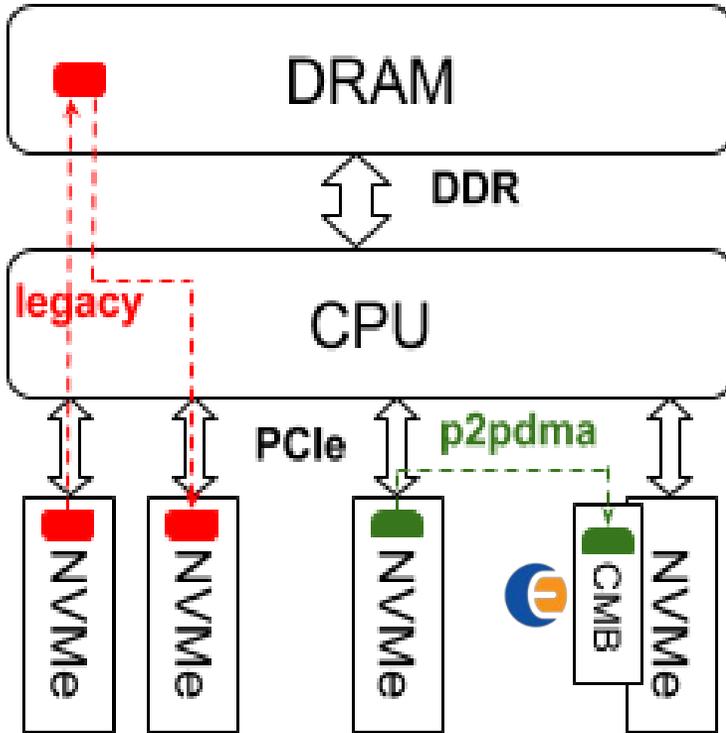- NVMe SGL support
- CMB and P2P support

**Available Now**

NoLoad™ Software for AI

Management

nvme-cli
nvme-of
etc

Applications
PYTORCH    Caffe    TensorFlow

GitHub
libnoload

intel
SPDK

Linux    Windows    vmware

Userspace

OS

Hardware

NoLoad™ Hardware

Both kernel and userspace
frameworks supported

libnoload connects NVMe
driver to AI libraries
(PyTorch, TensorFlow etc)

No changes to OS
Use inbox NVMe driver

Confidential

12

# NoLoad® CSP – Peer to Peer (P2P)
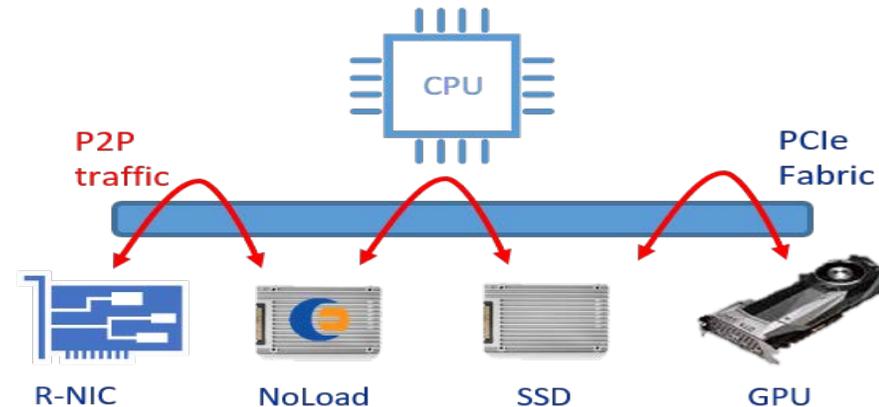
**DRAM** — DDR — **CPU** — legacy — PCIe — p2pdma — NVMe / NVMe / NVMe / NVMe CMB
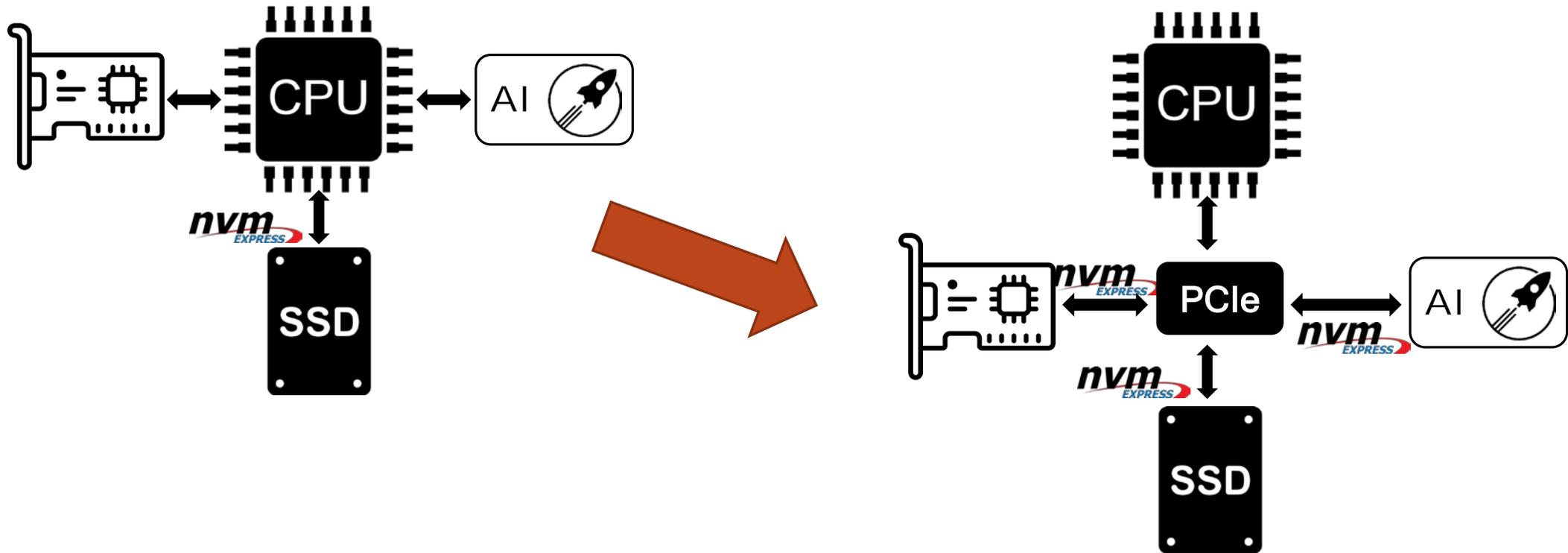
- PCIe End-Points (EPs) are getting faster e.g. NVMe SSDs, RDMA NICs & GPGPUs

- Bounce buffering all IO data through system memory is a waste of system resources

- NoLoad P2P allows PCIe EPs to DMA to each other whilst under host CPU control

- CPU/OS still responsible for security, error handling etc.

- 99.99% of DMA traffic now goes direct between EPs

- P2P avoids CPU's memory subsystem
- NoLoad implements a high performance NVMe CMB which can be used for P2P DMAs.



P2P traffic — CPU — PCIe Fabric — R-NIC — NoLoad — SSD — GPU

- Reduced data movement
- CPU offload of processing
- CPU offload of DMA traffic
- Standards based drivers and software
- Vendor neutral
- Management

**Advantages of NVMe Computational Storage for AI**
- Offloads inference to CSx which improves cost, power and space efficiency.
- p2pdma decouples data-plane from CPU further improving cost and power savings.
- CSx is a NVMe controller and can be **managed at scale using NVMe-MI**
- NoLoad Inference engine is a NVMe namespace and **can be shared over Ethernet** using NVMe-oF. Inference disaggregation reduces cost, saves space and saves power.
- **No proprietary drivers**. In time NVMe Computational Storage standard will allow for vendor-neutral approach.
- SNIA developing a **vendor-neutral interface** so hardware from different vendors is supported by the same SW stack.

**RocksDB Acceleration using NoLoad® CSP**

Eideticom's NoLoad® CSP deployed on Alveo

- 6x more transactions per sec
- 2.5x more efficient
- 4x reduced NAND costs
- Improved QoS

**Similar results achievable for AI inference.**

## NVM Express for AI

- Efficient and parallel PCIe protocol for talking to PCIe-attached NVM.
- Performance metrics aligned with needs of AI (GB/s, sub ms latency).
- Need filesystems to catch-up!

## Computational Storage for AI

- Offload key AI tasks to accelerators in a standards based way.
- Move the computation closer to the data to reduce data movement and improve efficiency.
- Standardize accelerator interfaces and management via NVMe.

Eideticom HQ          Eideticom (Bay Area)
3553 31$^{st}$ NW,        168 South Park,
Calgary, AB,           San Francisco, CA 94107
Canada T2L 2K7        USA

www.eideticom.com

Contact: sales@eideticom.com