# NVMe SSD Computational Storage

**Seamless Programming, Compute Acceleration**

Driving Compute and Storage Throughout the Datasphere!!

# The Market Evolution and Need for Local Compute

**Our Friends at Gartner Say it best…**
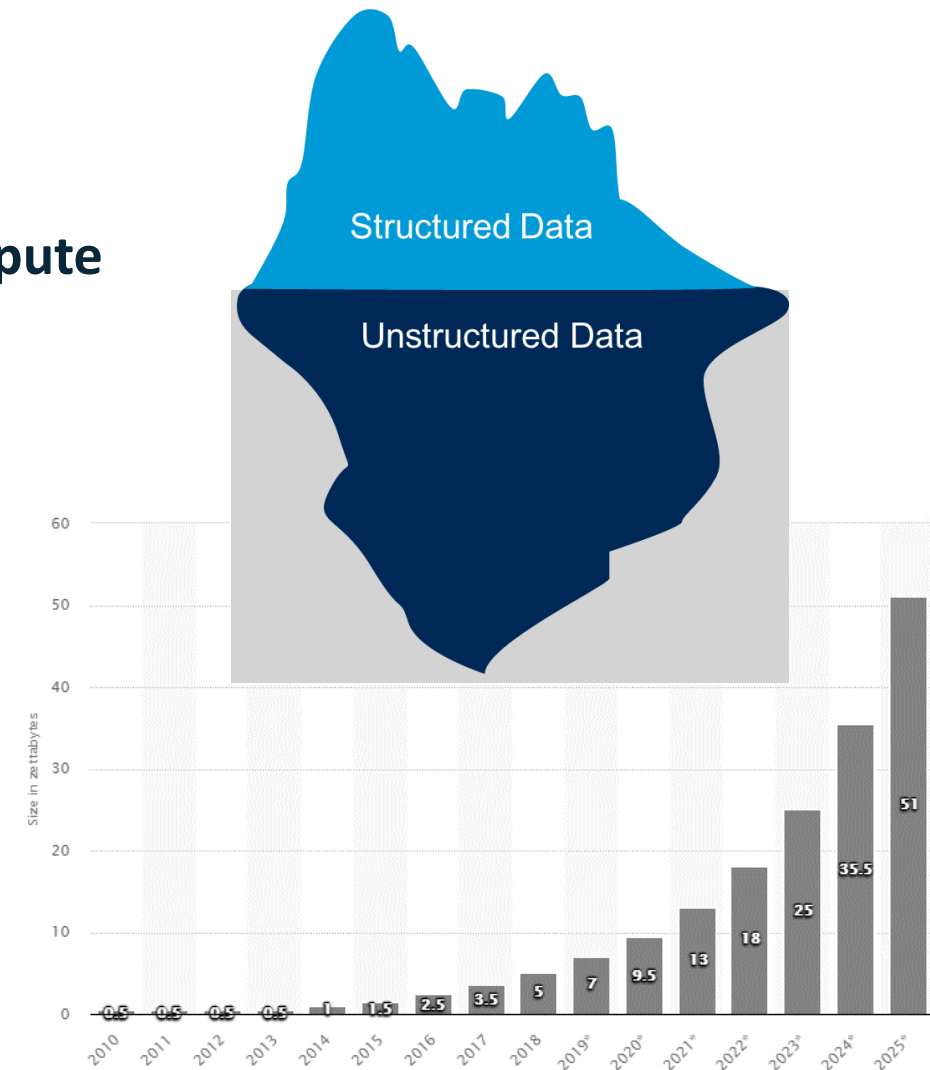
<u>Structured Data</u> is great for current infrastructure
  Allows for ease of data movement, location, access, compute
  Only a small subset of the real data Iceberg

<u>Unstructured Data</u> is the greatest threat to results
  As more and more data is generated, it is more random
  Needs to manage this data locally are key
  Edge Computing is not able to scale at data growth pace
  A new way to compute on random, local data is needed

<u>The Global DataSphere (Statista.com) shows how the
data growth is overshadowing the compute growth</u>

# The Market Needs a New Way to Look at Storage.

## Pain Points

Physical Space
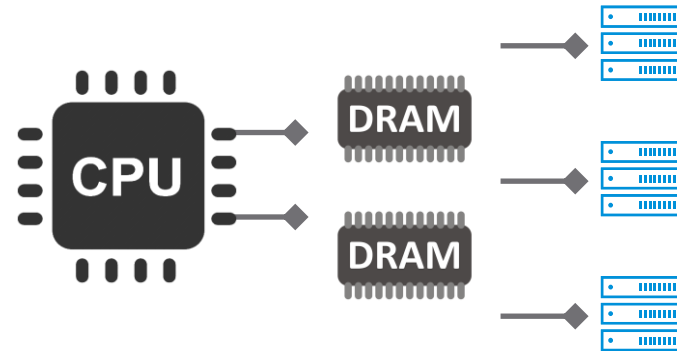
Available Power

Scaling  Mismatch

Bottleneck Shuffle

Traditional storage architectures are in **trouble**.

Scaling requirements are **not met** with existing solutions

One CPU to many storage devices **creates bottlenecks**

These bottlenecks exist, we currently just shift where they reside

Technologies that '**compose**' these elements just move the bottleneck

A way to augment and support **without wholesale change** is needed

# The Path to Compute Solutions is Paved with Smart Intentions

**Finding paths to compute is easy… But one thing is very lacking in these 'Smart Things'**

> **Compute is Needed, DATA is Mandatory!**
>
> **CPU – The Brain of the operations, starved for data, overwhelmed with requests**
>
> **GPU – The Parallel processing Master, Nothing Persistent about it**
>
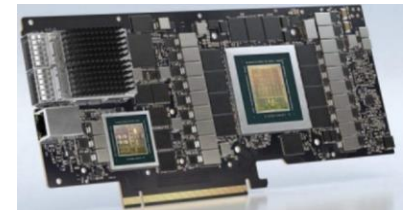> **NIC – The great Mover, not so great at processing**
>
> **Smart NIC – The intelligent mover, but still doesn't know what it is moving**
>
> **DPU – The Processor closer to data, but still not persistent, still Volatile!**

**All these pieces are needed parts of the new ecosystem.**

**But NONE of them address the Real Issue…**

# The Data, where it is, where it comes from, and how to Store it!
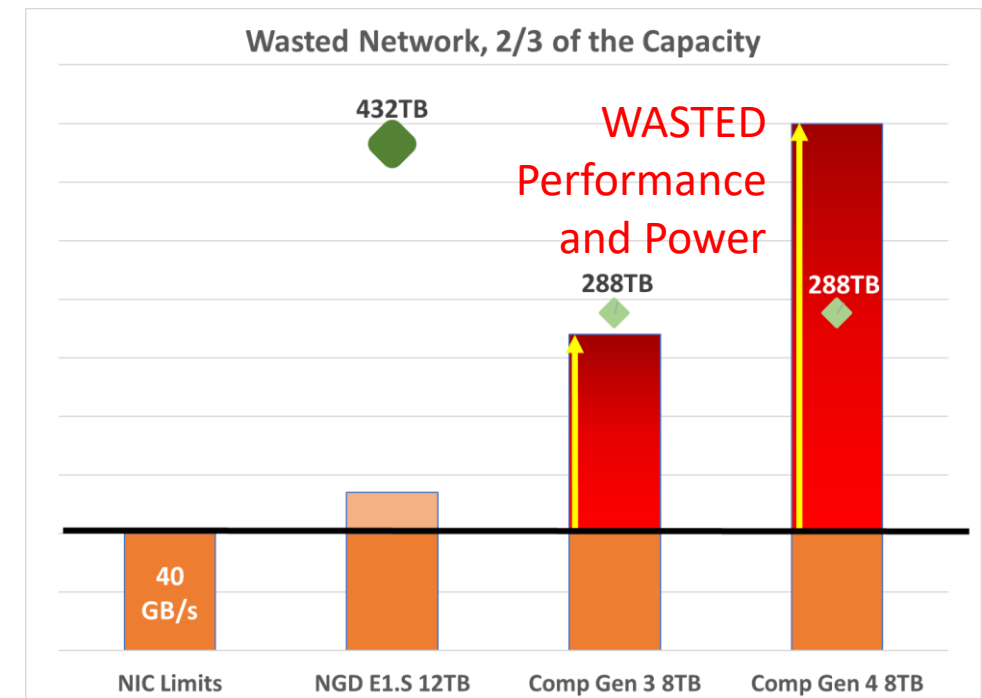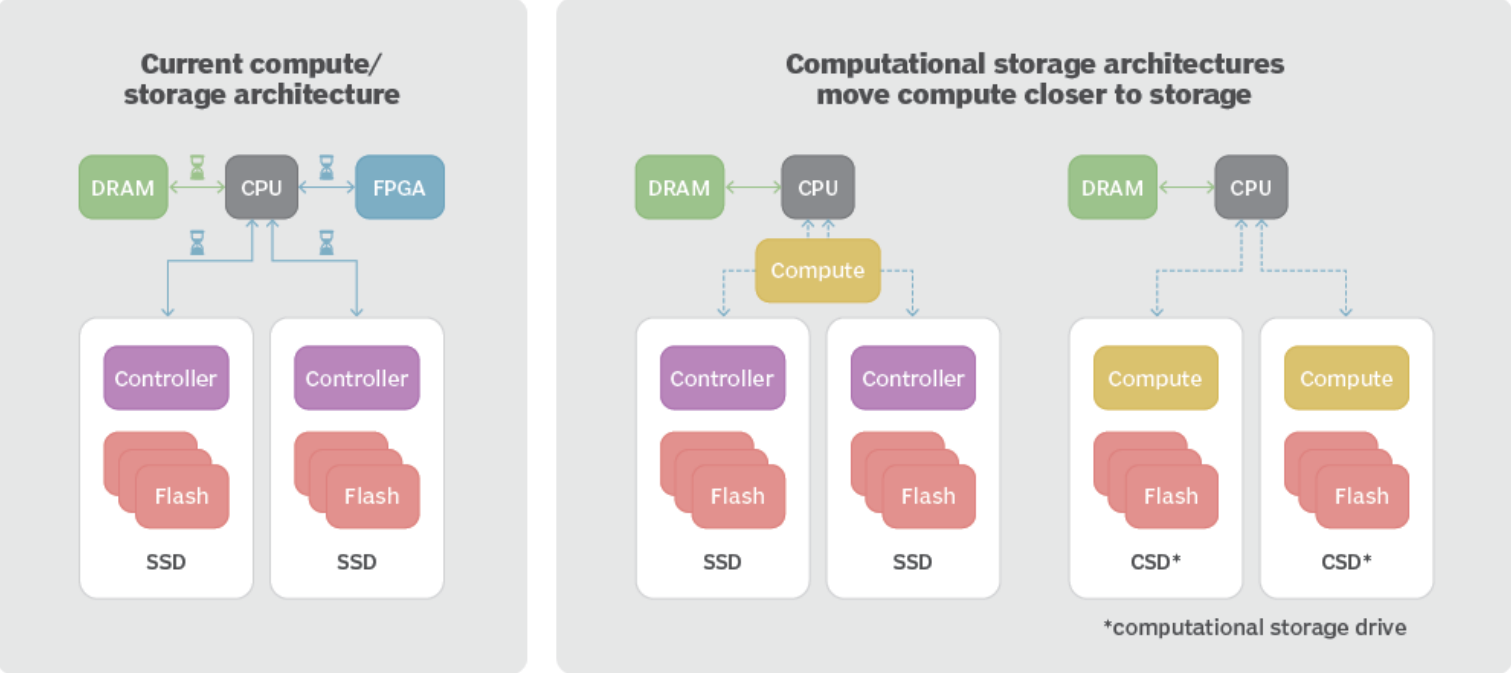
# PCIe is Simply Not Enough.

## More Lanes, More Traffic, No Solutions

- **Storing** the Raw Data is easy!!

- **Working** on the Data will be Difficult

- **MOVING** the Data is Impossible

- **Solve** this with Computational Storage
  – Standardized, Open, Flexible

IDC predicts we will churn out <u>175 zettabytes</u> of data in 2025



Wasted Network, 2/3 of the Capacity

432TB

WASTED Performance and Power

288TB    288TB

40 GB/s

NIC Limits    NGD E1.S 12TB    Comp Gen 3 8TB    Comp Gen 4 8TB

# The Value of Computational Storage – Core Count!



Current compute/storage architecture

Computational storage architectures move compute closer to storage

*computational storage drive

## Computational Arm Cores
### 5376 Cores per 42U rack
### 16PB of Storage per rack



### Value of Computational Cores

✓ Distributed Processing

✓ Near-Data Processing

✓ Smaller System Footprint

32 x E1.S hot swap NVMe SSDs
4 cores each
**128 Additional Cores per 1U Server**



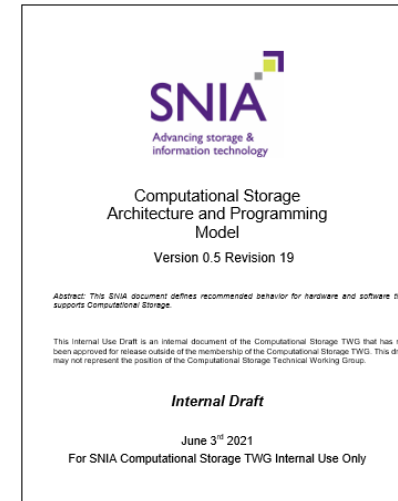## NGD Systems High-Capacity, NVMe Computational Storage

# What the Market is Doing to Drive Computational Storage

**SNIA is driving for an Architectural Solutions**

**NVM Express is working on an Initial Instructions**

**Prototyping and Deploying Now**

SNIA
Advancing storage & information technology

Computational Storage
Architecture and Programming
Model

Version 0.5 Revision 19

Abstract: This SNIA document defines recommended behavior for hardware and software that supports Computational Storage.

This Internal Use Draft is an internal document of the Computational Storage TWG that has not been approved for release outside of the membership of the Computational Storage TWG. This draft may not represent the position of the Computational Storage Technical Working Group.

*Internal Draft*

June 3rd 2021
For SNIA Computational Storage TWG Internal Use Only

**NVMe Computational Storage Task Group**

The charter of Computational Storage Task Group is to develop features associated with the concept of **Computational Storage on NVM Express devices**.

The target audience consists of the vendors and customers of **NVMe Storage Devices** that support computational features.

# Market Readiness for Computational Storage

- ## Industry Investigations Grows

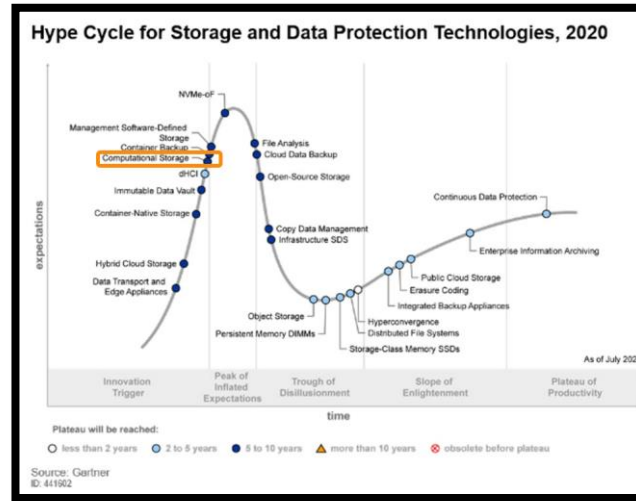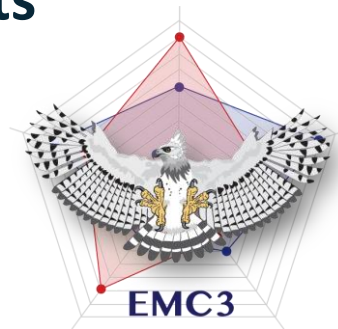- ## Industry Analysts
  - Gartner , IDC, Others



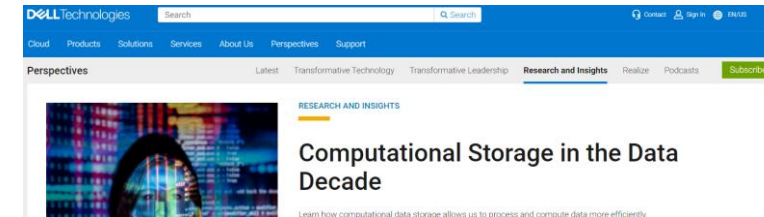Hype Cycle for Storage and Data Protection Technologies, 2020



CW Developer Network

## Computational storage: A Computer Weekly analysis series

Computational storage series: Evaluator Group - Speculations, expectations & extrapolations

CW Developer Network

Computational storage: NGD Systems / SNIA - Icebergs at the Edge

## Computational Storage in the Data Decade

- ## Customer Sponsored Efforts

**Los Alamos National Laboratory Welcomes NGD Systems to the Efficient Mission Centric Computing Consortium**

Monday, August 2 2021



*The collaborative effort will explore high capacity NVMe computational storage drive, and scalable computational offloads for HPC and scalable computing uses*
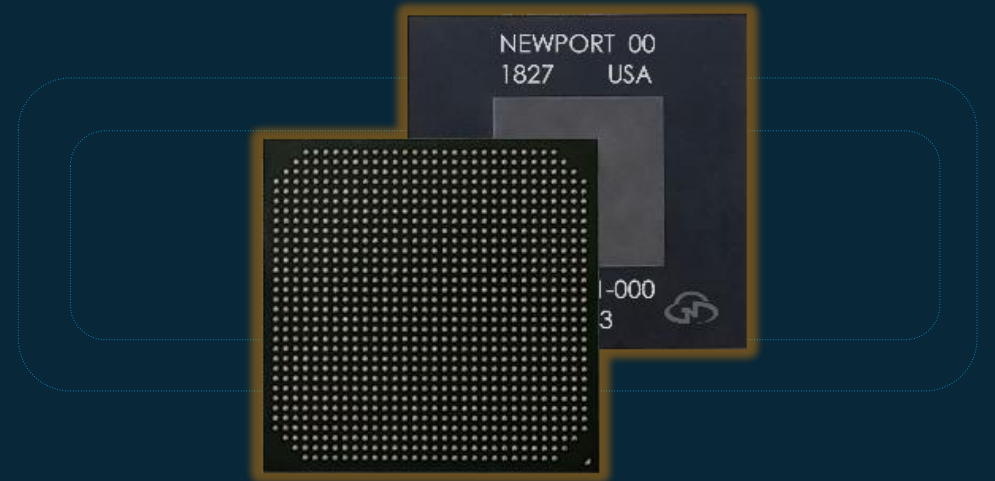
### Key Solution Elements

**Computational Storage**

Embed compute with storage, offloading main server, improving performance on smaller systems by reducing data transfer to main system and enabling on-chip intelligence

**Parallel Database with integrated Analytics**

Query across NVMe devices in parallel, making effective use of computational storage. Embedded analytics allowing analytics free of resources on the main system. Seamless replication of data to backup host.

**vSphere & Bitfusion**

Ability to offer Edge resiliency with vSAN, HA, FT. GPU acceleration for computational storage w/ Bitfusion. Effective use of limited host resources.
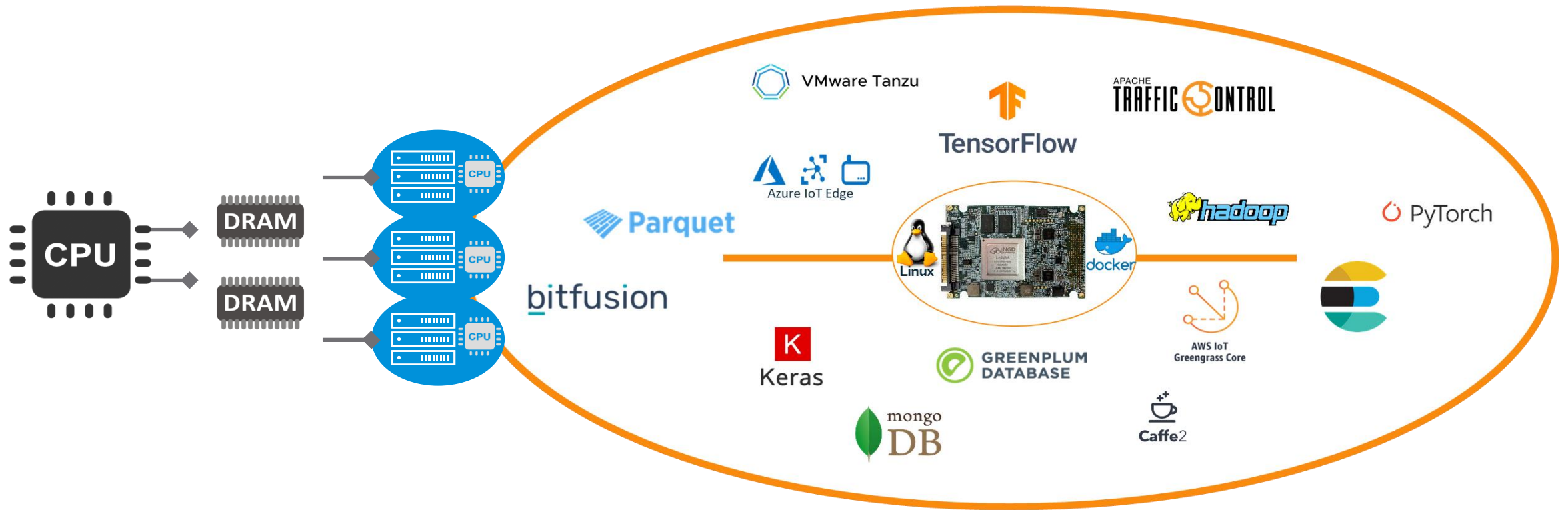
# How To Do It

## Keep It Simple and Seamless

ASIC-based, Single-chip, All in one Solution

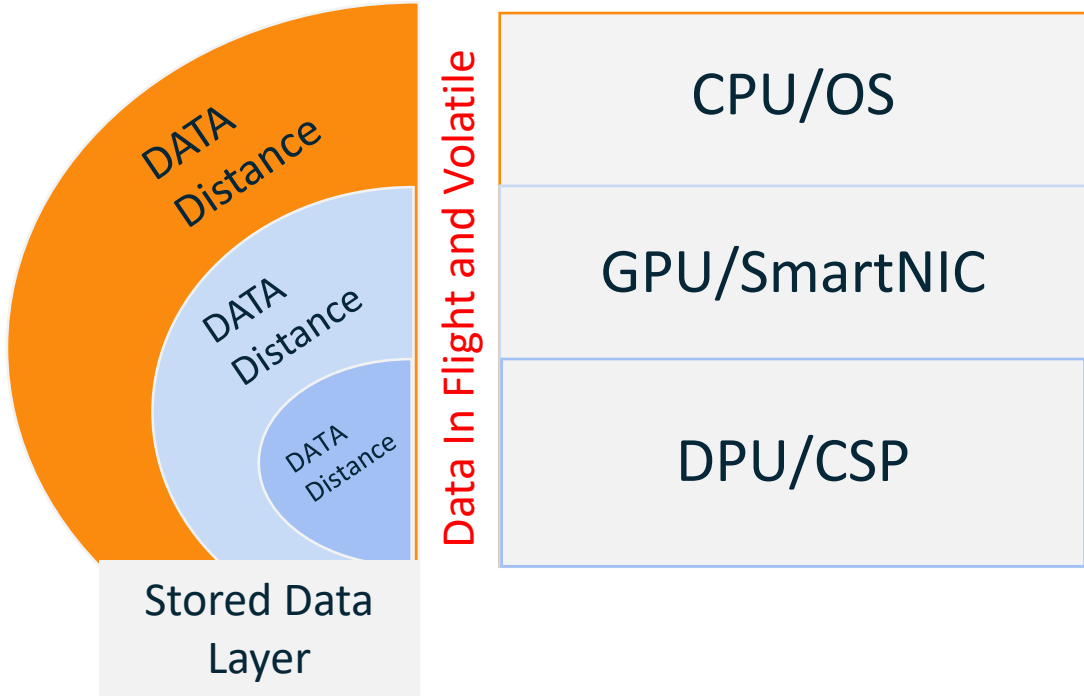# Some Examples of Linux Deployments with K.I.S.S.

- **Keep It Simple & Seamless**
  - The best way to move technology forward is to leverage architectures already in use
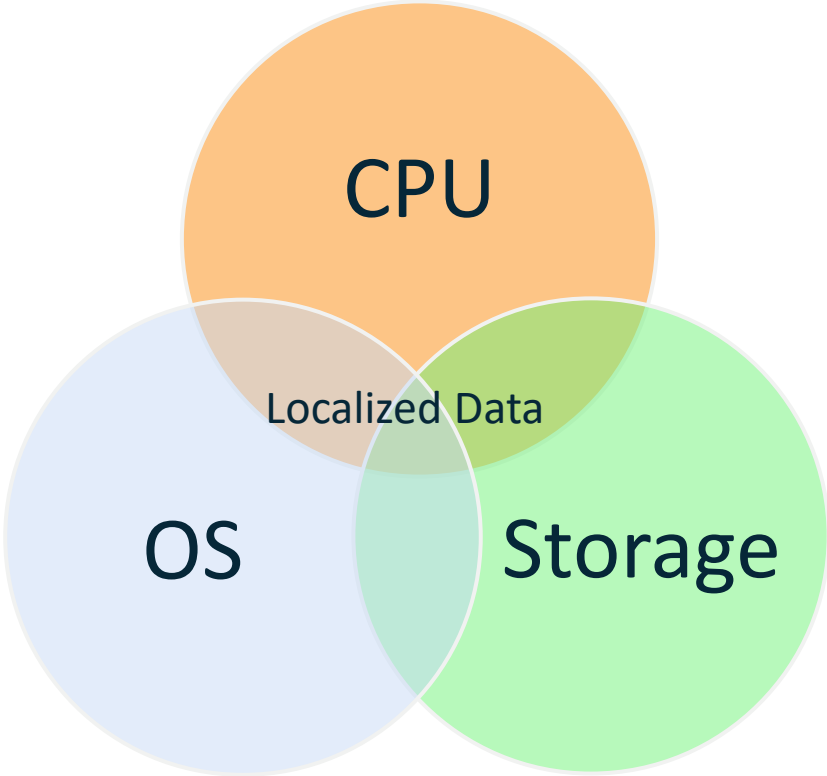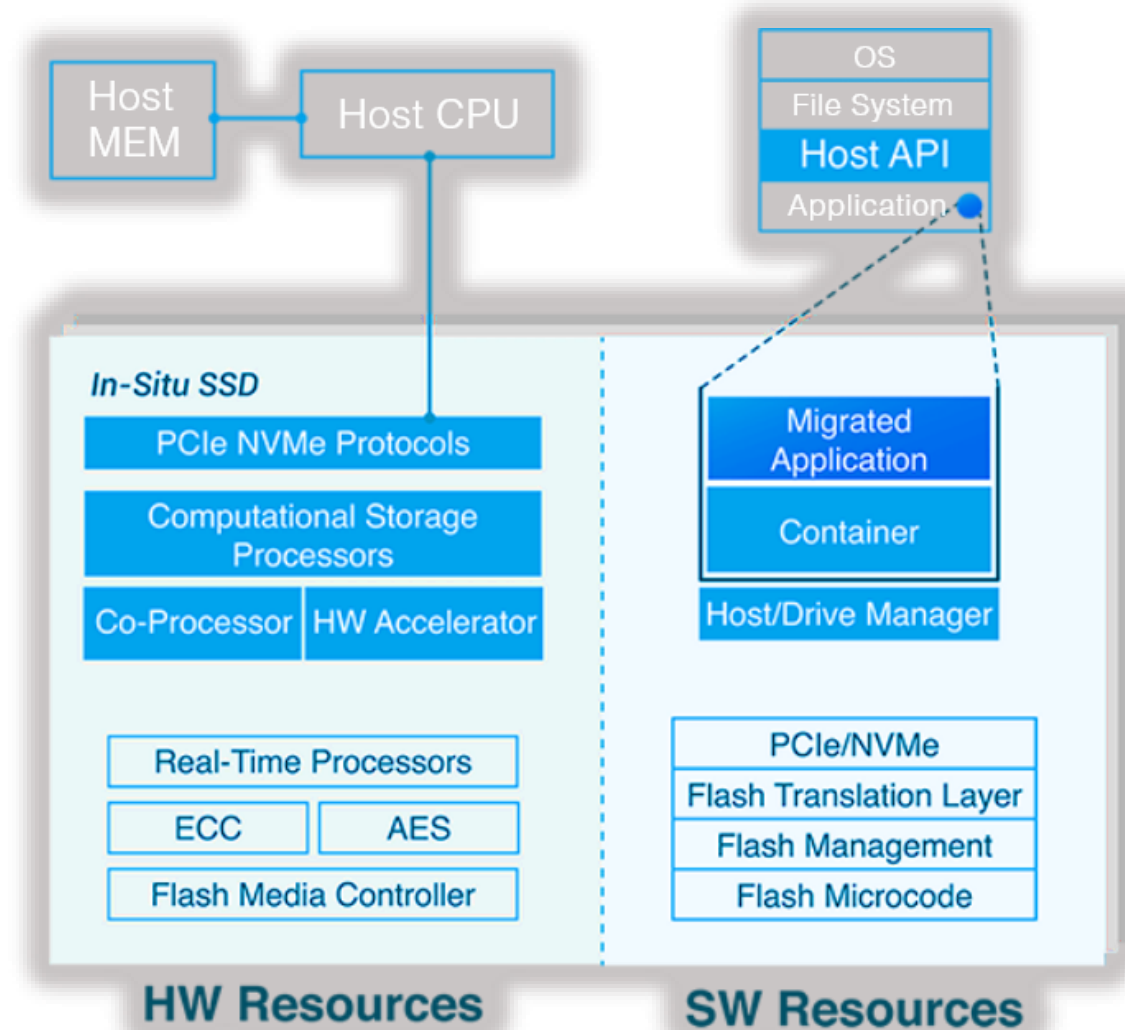
# A Comparison of Compute Infrastructure – Why CSDs?

Today's Standard Infrastructure – **Data Distant**

Linux-Based Computational Storage Drive – **Data Locality**

# A Look at the Hardware and Software of a Linux-Based CSD

The Data Lives on Storage. Why Not Work on it There?

One Host Many Drives

NGD Systems

Server Host Processor Complex

X86  GPU

host OS

host agent

application

NVMe

Computational Storage Drive

media controller

Arm quad-core

shared DRAM

drive OS
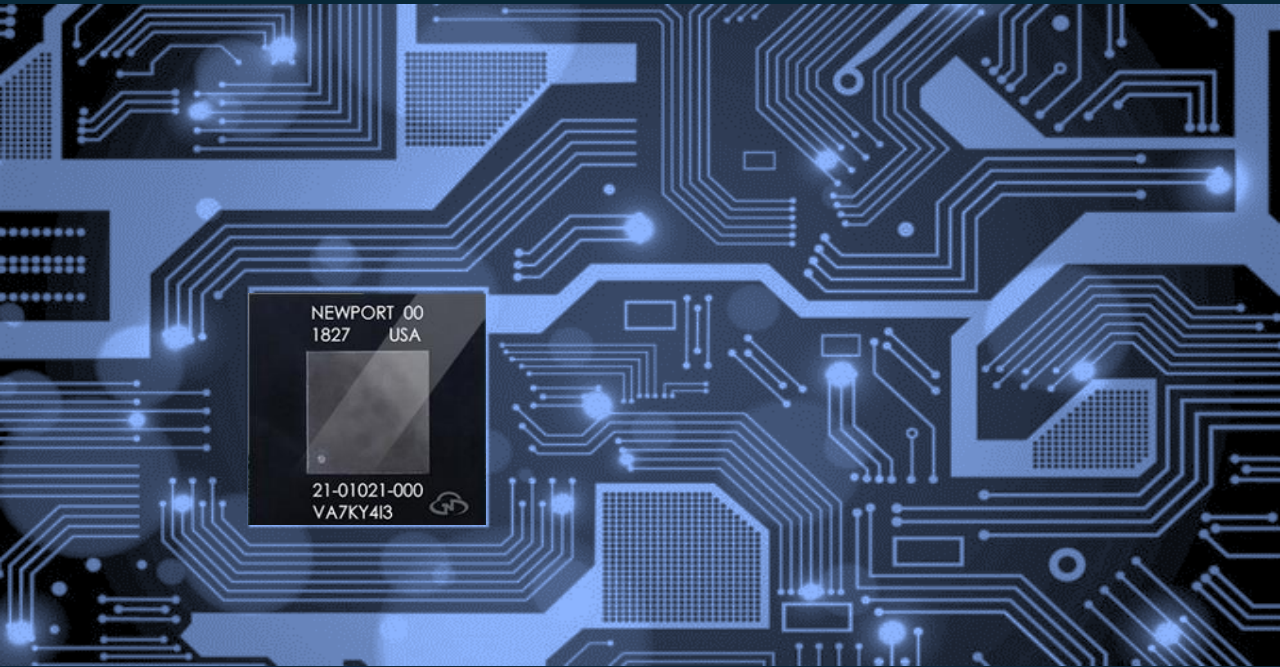
app

NAND media

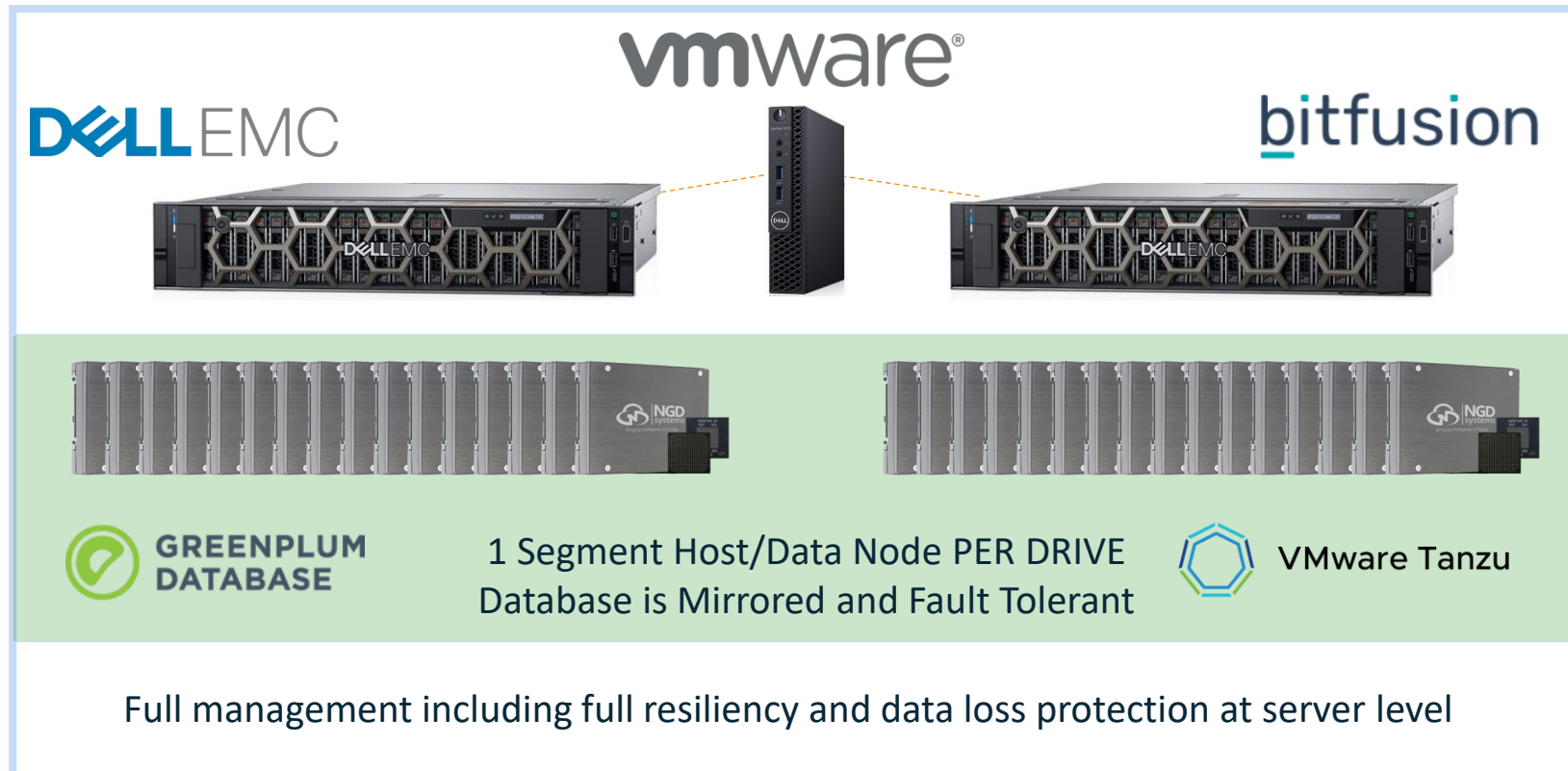# Computational Storage, Some Real World Results.

- **Compression Acceleration – GZIP**

- **CDN-in-a-Box – Real Customer Lab Results**

- **AI Inference in the Datacenter – FAISS**

- **Inferencing at the Edge - WiSARD**

- **Distributed Machine Learning - Stannis**

- **Data Search – Elasticsearch**

- **Distribute Processing – Hadoop**

# Edge Analytics – Live Demo with VMware – xLab 52

**Computational Storage allows it to be drive level.**

**Reducing footprint, server cost, while still offering full fault tolerance**



1 Segment Host/Data Node PER DRIVE
Database is Mirrored and Fault Tolerant

Full management including full resiliency and data loss protection at server level

**Showcased at Vmworld 2020  - Session ID [OCTO478] –**
Computational Storage, Tanzu Greenplum, vSphere Bitfusion

# Finding the Needles in Haystacks with AI and CSDs
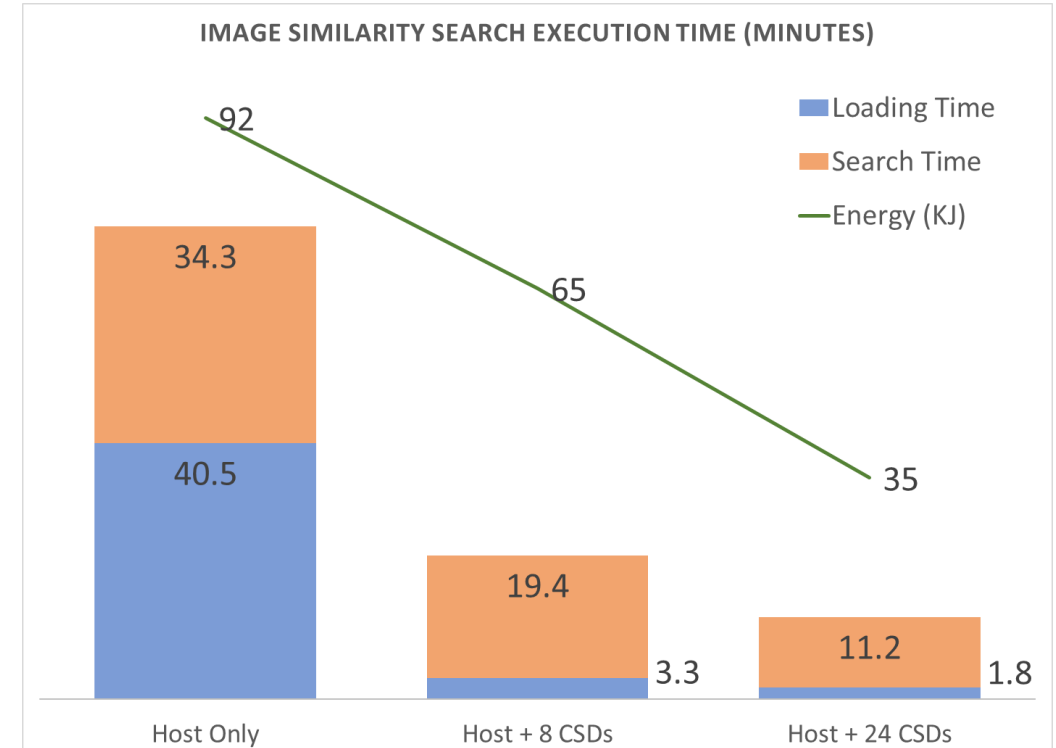
## Problem Statement

- **Databases growing at exponential rates**

| 10 M | 1 Billion | 1 Trillion |
|------|-----------|------------|
| 2007 | 2017 | 2021 |

- **Load and Search time key blocks in getting results**

## Computational Storage Solution

- **Determine best way to increase performance**
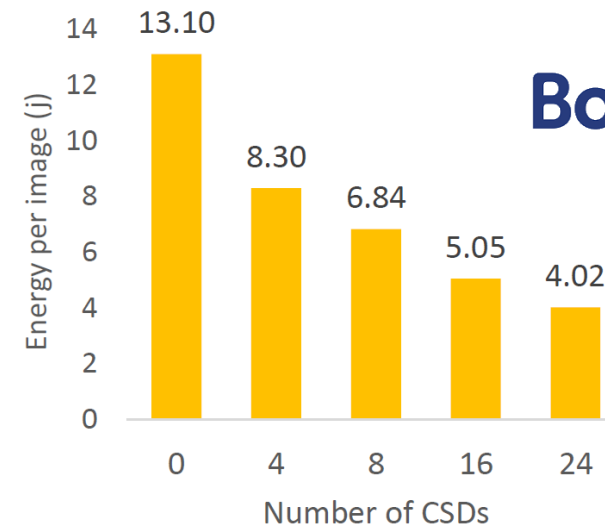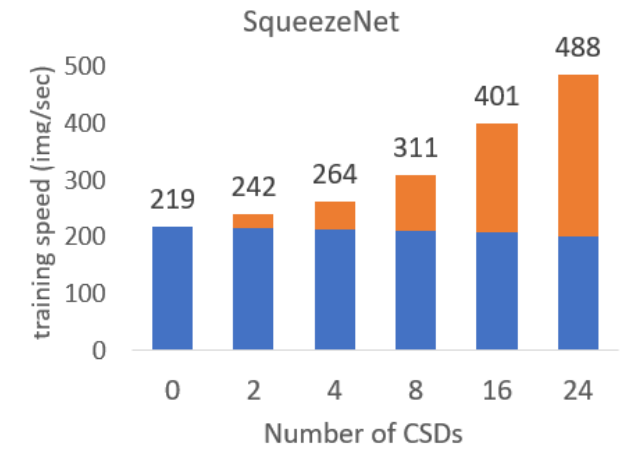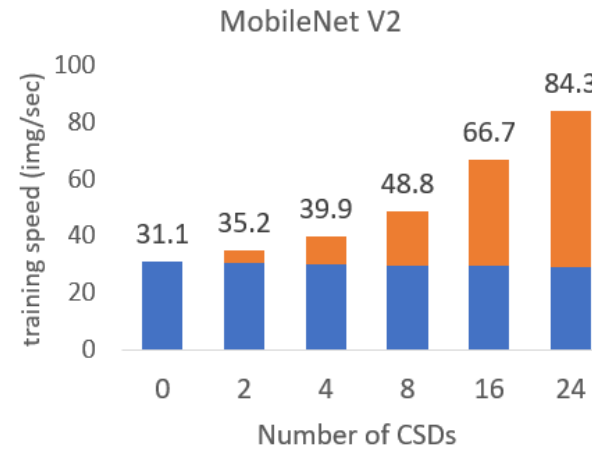- **Load Time Reductions due to CSD Offload of AI code**

## Results are Proven:

- **Load Time Reduced**       **> 95%**
- **Search Time Reduced**      **> 60%**
- **Power Savings of**         **> 60%**



IMAGE SIMILARITY SEARCH EXECUTION TIME (MINUTES)

- Loading Time
- Search Time
- Energy (KJ)

Host Only: Loading Time 40.5, Search Time 34.3, Energy 92
Host + 8 CSDs: Loading Time 3.3, Search Time 19.4, Energy 65
Host + 24 CSDs: Loading Time 1.8, Search Time 11.2, Energy 35

Association for Computing Machinery

Microsoft

Technical paper published in the ACM journal on Computational Storage

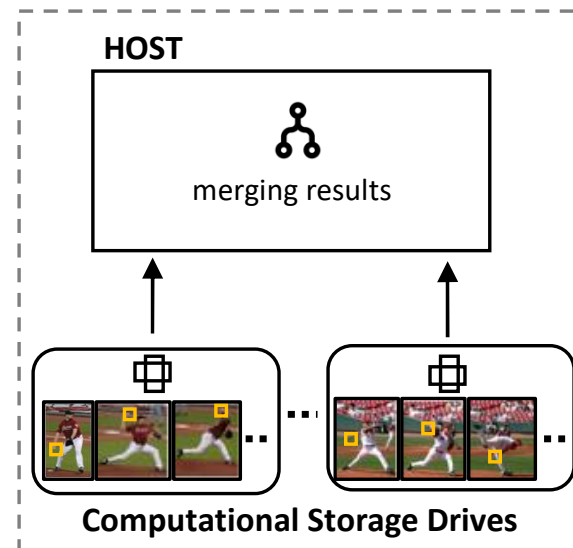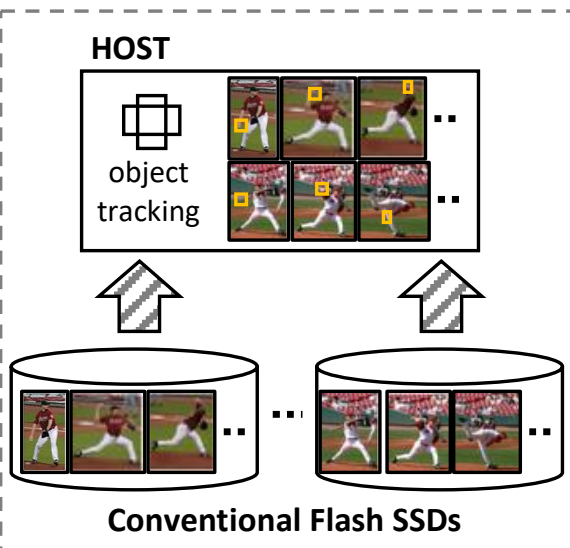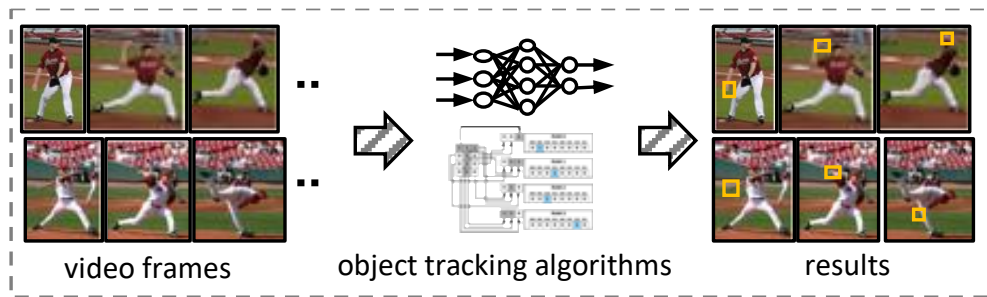# Machine Learning At Scale – Not Just One Way

- Four neural networks Evaluated
  - **MobilenetV2**
  - NASNet
  - **SqueezeNet**
  - InceptionV3Quad-core

- Tested with 24 CSDs

- Training data stored on CSDs

- Using an AIC 2U-FB201-LX server
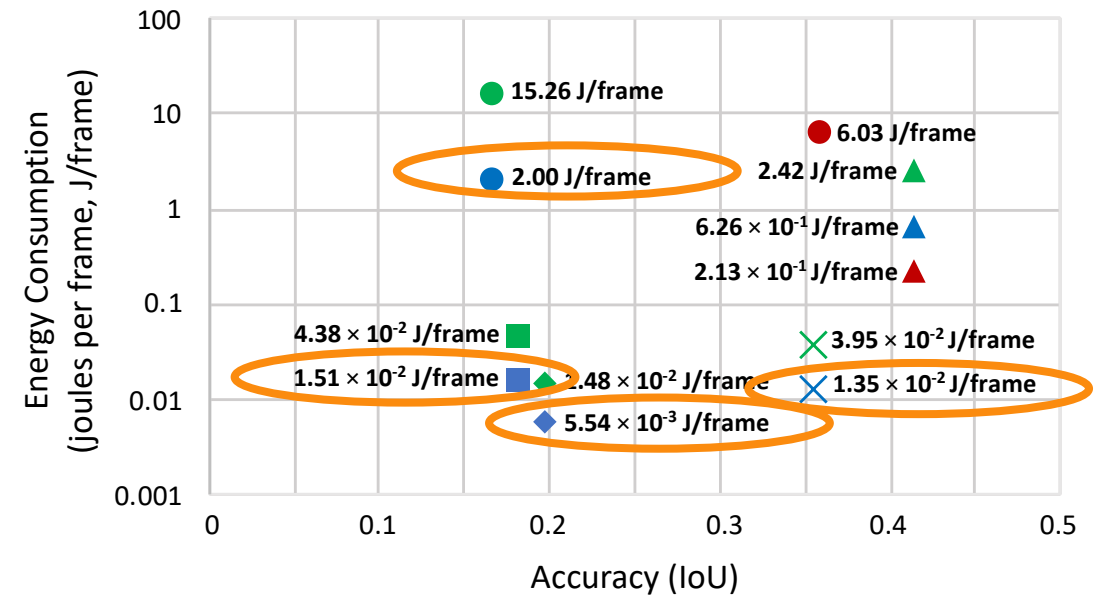  - Intel® Xeon® Silver 4108 CPU
  - 32GB DRAM

# Computer Vision – Lower Power, Same Results.


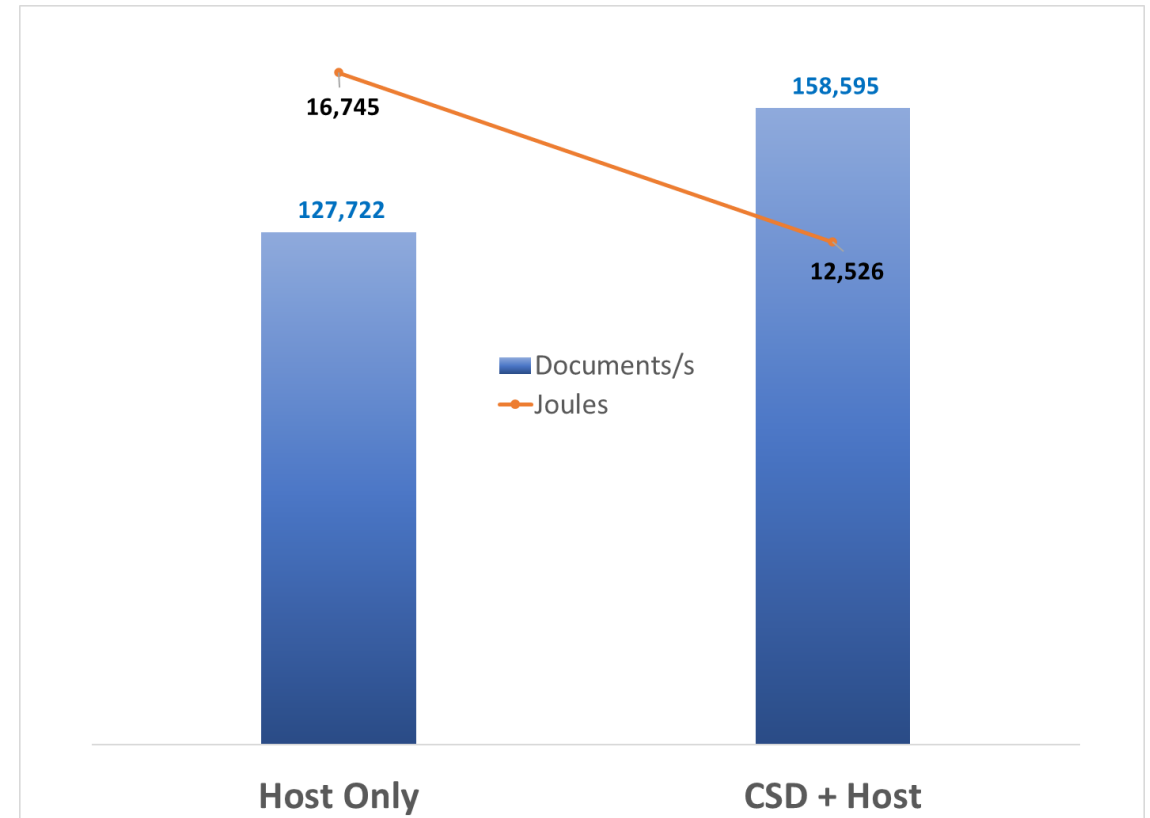
- **Scalable solution**
- **Energy efficiency gains (~10x)**

video frames → object tracking algorithms → results

**HOST**
object tracking
**Conventional Flash SSDs**

**HOST**
merging results
**Computational Storage Drives**

Legend:
- ● YOLO GPU
- ▲ GOTURN GPU
- ● YOLO CPU
- ▲ GOTURN CPU
- ■ KCF CPU
- ◆ MOSSE CPU
- ✕ WiSARD CPU
- ● YOLO CSD
- ▲ GOTURN CSD
- ■ KCF CSD
- ◆ MOSSE CSD
- ✕ WiSARD CSD

Plot — Energy Consumption (joules per frame, J/frame) vs Accuracy (IoU):
- ● 15.26 J/frame
- ● 6.03 J/frame
- ● 2.00 J/frame
- ▲ 2.42 J/frame
- ▲ $6.26 \times 10^{-1}$ J/frame
- ▲ $2.13 \times 10^{-1}$ J/frame
- ■ $4.38 \times 10^{-2}$ J/frame
- ✕ $3.95 \times 10^{-2}$ J/frame
- ■ $1.51 \times 10^{-2}$ J/frame
- $2.48 \times 10^{-2}$ J/frame
- ✕ $1.35 \times 10^{-2}$ J/frame
- ◆ $5.54 \times 10^{-3}$ J/frame

**CSD Results Are equal in accuracy with Less Power**

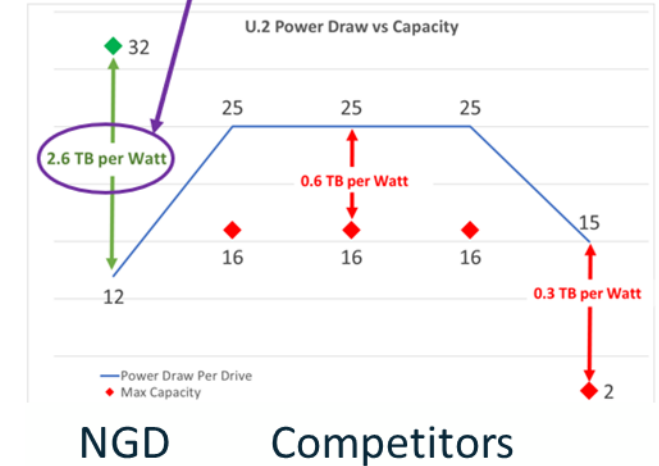# Hybrid Configuration Performance Results

- **Total Performance Improves**
  - **20% Better** Results

- **Reduced Power Consumption**
  - **30% LESS** Power

- **DRAM Usage Reduced by >50%**
  - Host Only used    **25GB**
  - Hybrid used    **12GB**

- **CPU Usage Utilization Reduced by >50%**
  - Host Only used    **24%**
  - Hybrid used    **10%**

# NGD NVMe CSD Products at a Glance.

- Large breadth of **SSD** solutions and capacity options
- Leading **TB/W** Energy Efficiency
- Only **16-Channel** 14nm SSD SoC & 100% Made in the USA
- Industry's Largest capacity NVMe SSDs
- Quad-Core **Computational Storage CPUs**

| Form Factor | Availability | Raw Capacity TLC (TB) | MAX Power (W) |
|---|---|---|---|
| M.2 2280 | CQ3'20 | up to 4 | 8 |
| M.2 22110 | NOW | up to 8 | 8 |
| U.2 15mm | NOW | up to 32 | 12 |
| EDSFF E1.S | NOW | up to 12 | 12 |
| EDSFF E3 | Planned | up to 64 | 12 |



**NGD is the ONLY Provider of Capacity over Power**
Another Paradigm Shift in the Market

U.2 Power Draw vs Capacity

NGD — Competitors

U.2

M.2

E1.S

# How Can We Help?