



# Computational Storage Real-World Deployments

Stephen Bates  
CTO, Eideticom



Flash Memory Summit



# Computational Storage: State of the Nation



- State of the Nation
- Real-World Deployments
  - NGD Systems
  - Samsung
  - ScaleFlux
  - Eideticom

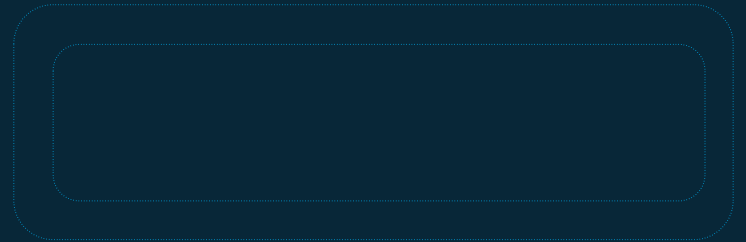


# Computational Storage: Real World Deployments

- Computational Storage standardization is happening!
  - SNIA Computational Storage: 45+ member companies, 224 members
    - ❖ High-level architecture
    - ❖ User-space library
  - NVMe Computational Storage: 25+ member companies, 78 members
    - ❖ NVMe commands for Computational Storage



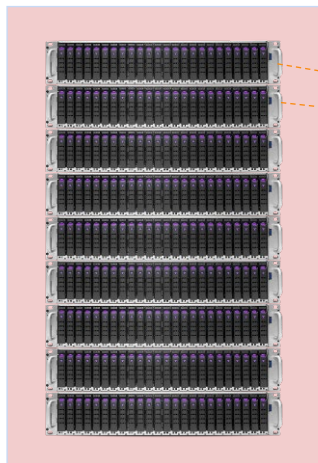
# VMware Demo



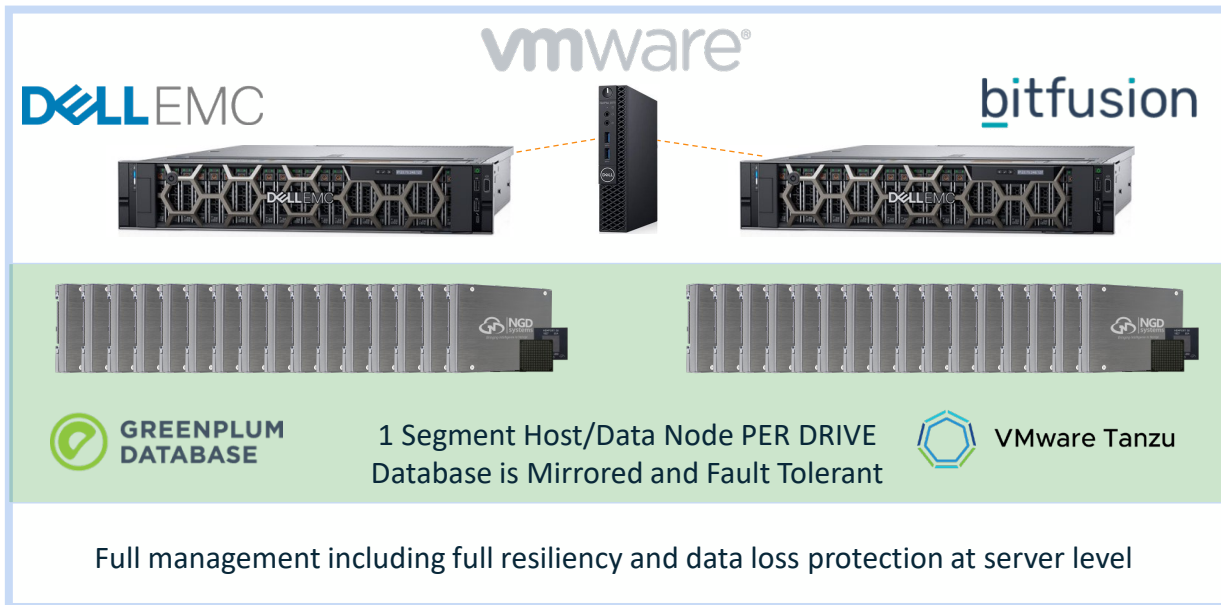
# Edge Analytics –Demo with VMware – xLab Platform

Computational Storage allows it to be drive level.

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

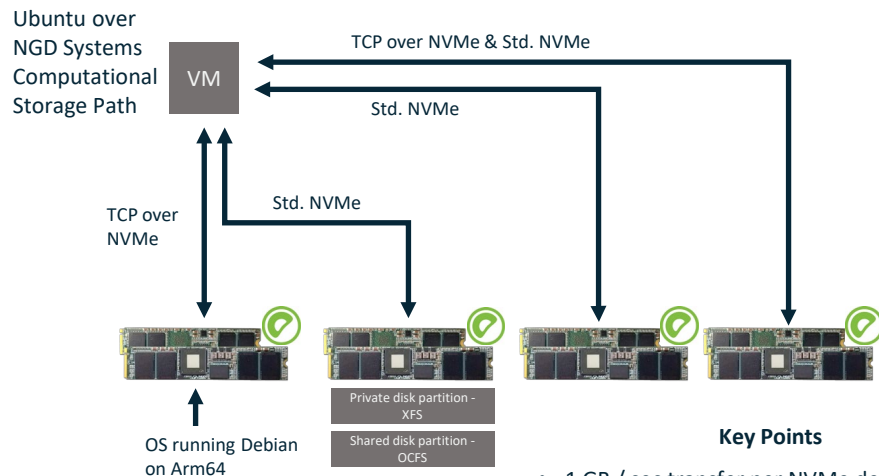


Traditionally Segment Host and Data Nodes were at the server level.



# Integration of NGD Systems Devices to vSphere

[LINK TO ONLINE DEMO – VMWorld 2020](#)



## Key Points

- 1 GB / sec transfer per NVMe device
- 16TB capacity per device
- Simultaneous addressing as storage device & as remote compute node
- PCI passthrough allows native use by VMs
- Greenplum running on each node

1. VM Directpath IO for NVMe devices
  1. up to 15 into a VM
2. TCP connected jump box allows addressing of devices from network.
3. Two partitions
  1. one shared w/OCFS
  2. one dedicated to Greenplum

## 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.



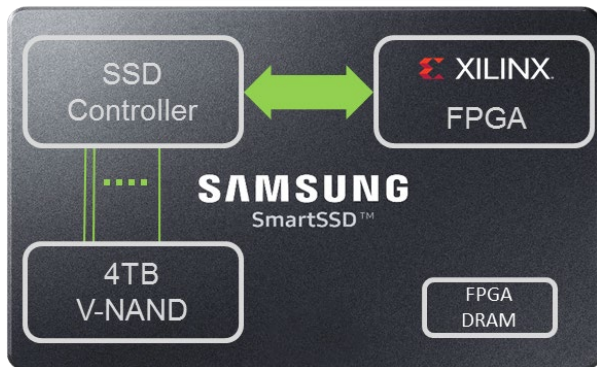
# Samsung SmartSSD



# SmartSSD® CSD Scales to Accelerate Data-Rich Workloads

Flash Memory Summit

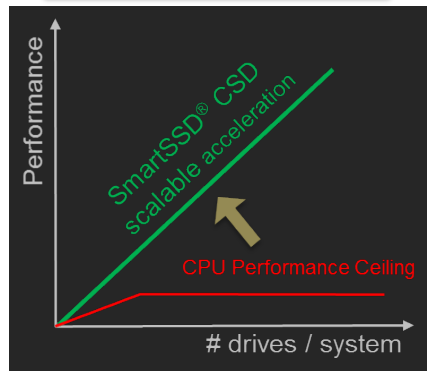
## SmartSSD U.2 Platform



### Computational Storage

- ✓ **3 & 6 GBps internal BW per device:**  
Minimize external data movement
- ✓ **FPGA:** Each device has 3x~10x core equivalents for offload/acceleration
- ✓ **4TB storage, 4 GB FPGA DRAM:**  
For Inline and Data@Rest processing

## Acceleration Concept



### Scalable Performance

- ✓ **Near Data Processing:** Data format conversion, Filtering, Metadata management, DB Analytics, Video processing
- ✓ **New Services:** Secure content, Edge acceleration

## Eideticom + SmartSSD CSD



### NVMe Computational Storage

- ✓ **Standards Based:** Uses NVMe to access both the storage and the computation
- ✓ **Consumable:** Leverages the inbox NVMe drivers and software available in all modern OSes

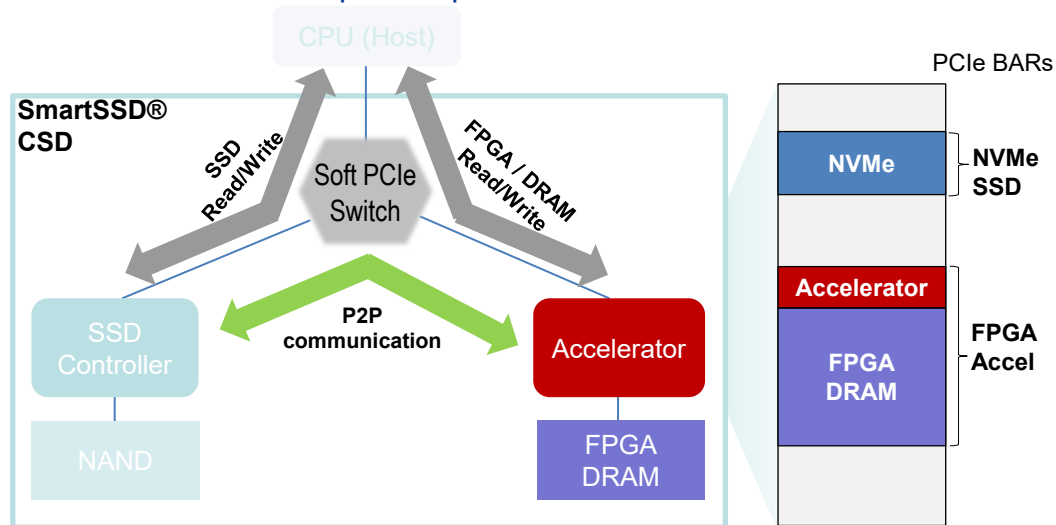
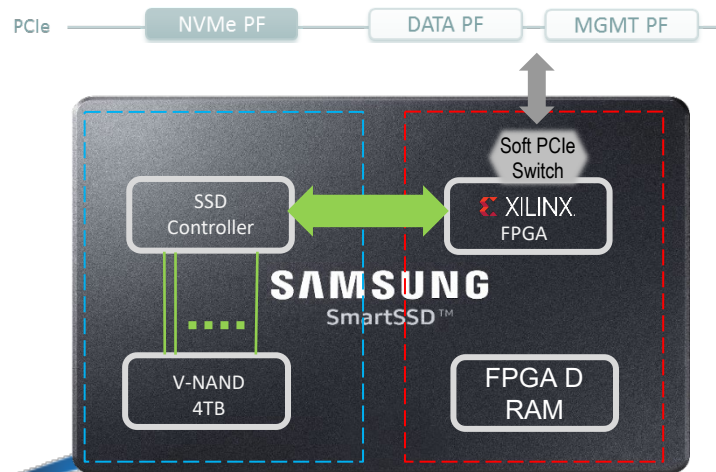




# SmartSSD® CSD HW Architecture

## Flash Memory Summit

- Peer-to-peer (P2P) communication enables unlimited concurrency
  - SSD:Accelerator data transfers use internal data path
    - Save precious L2:DRAM Bandwidth (Compute Nodes) • Scale without costly x86 frontend (Storage Nodes)
    - Avoid the unnecessary funneling and data movement of standalone accelerators
  - FPGA DRAM is exposed to Host PCIe address space
    - NVMe commands can securely stream data from SSD to FPGA peer-to-peer





# ScaleFlux

# Data Filtering in the Real World

## Alibaba Cloud & ScaleFlux

POLARDB Meets Computational Storage: Efficiently Support Analytical Workloads in Cloud-Native Relational Database

Wu Cao<sup>1</sup>, Yang Liu<sup>1</sup>, Zhaohu Cheng<sup>2</sup>, Ning Zheng<sup>1</sup>, Wei Li<sup>1</sup>, Weiqin Wu<sup>1</sup>, Lingyang Duang<sup>1</sup>, Peng Wang<sup>1</sup>, Jing Wang<sup>1</sup>, Ren Gao<sup>1</sup>, Zhenyu Lu<sup>1</sup>, Peng Shi<sup>1</sup>, Yang Zhang<sup>1</sup>

<sup>1</sup> Alibaba Group, Hang Zhou, Zhejiang, China

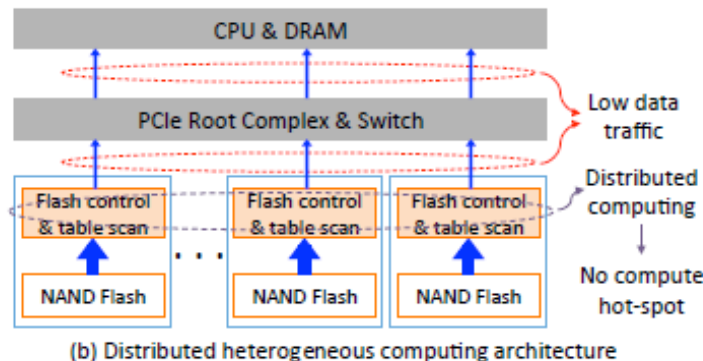
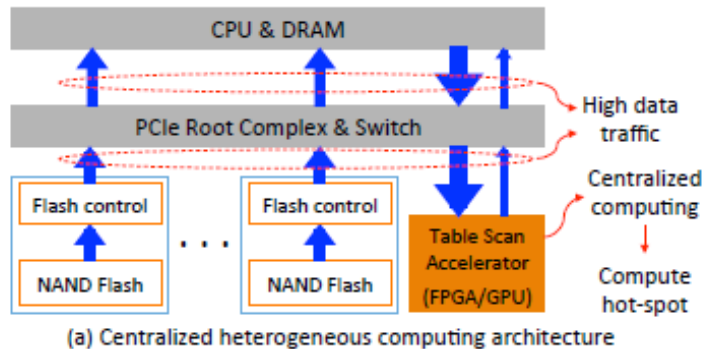
<sup>2</sup> ScaleFlux Inc., San Jose, CA, USA



**Abstract**  
This paper reports the deployment of computational storage in Alibaba Cloud to enable cloud-native database cost-effectively support analytical workloads in compute-storage decoupled architectures by leveraging cloud-native hardware data storage table scans from front-end database nodes to back-end storage nodes. The emerging storage nodes can leverage the in-storage computing capability to perform table scans, reducing the network overhead of data transfer. This paper presents a holistic implementation for a cloud-native relational database (POLARDB). To the best of our knowledge, this is the first real-world deployment of cloud-native database with computational storage nodes ever reported in the open literature.

**1 Introduction**  
Relational database is an essential building block in modern information technology infrastructure. Traditionally, all the database operations are performed on the database nodes (CPU and DRAM). In this paper, we propose a new architecture for cloud-native database with computational storage nodes.

This work applies heterogeneous computing in POLARDB design to efficiently support analytical workloads. The key idea is simple: push table scans from front-end database nodes to back-end storage nodes. This design becomes a computational storage after [1] that can carry out table scans in the I/O path. Compared with old storage table scans to a dedicated storage device (e.g., a FPGA-based storage device), distributing table scans across all the storage devices can minimize the data traffic across the network, thereby and reduce data processing hot-spot. This simple concept is not new and has been discussed



### Challenges with baseline architecture:

1. High Data Traffic
2. Data processing hot-spots

### Solution:

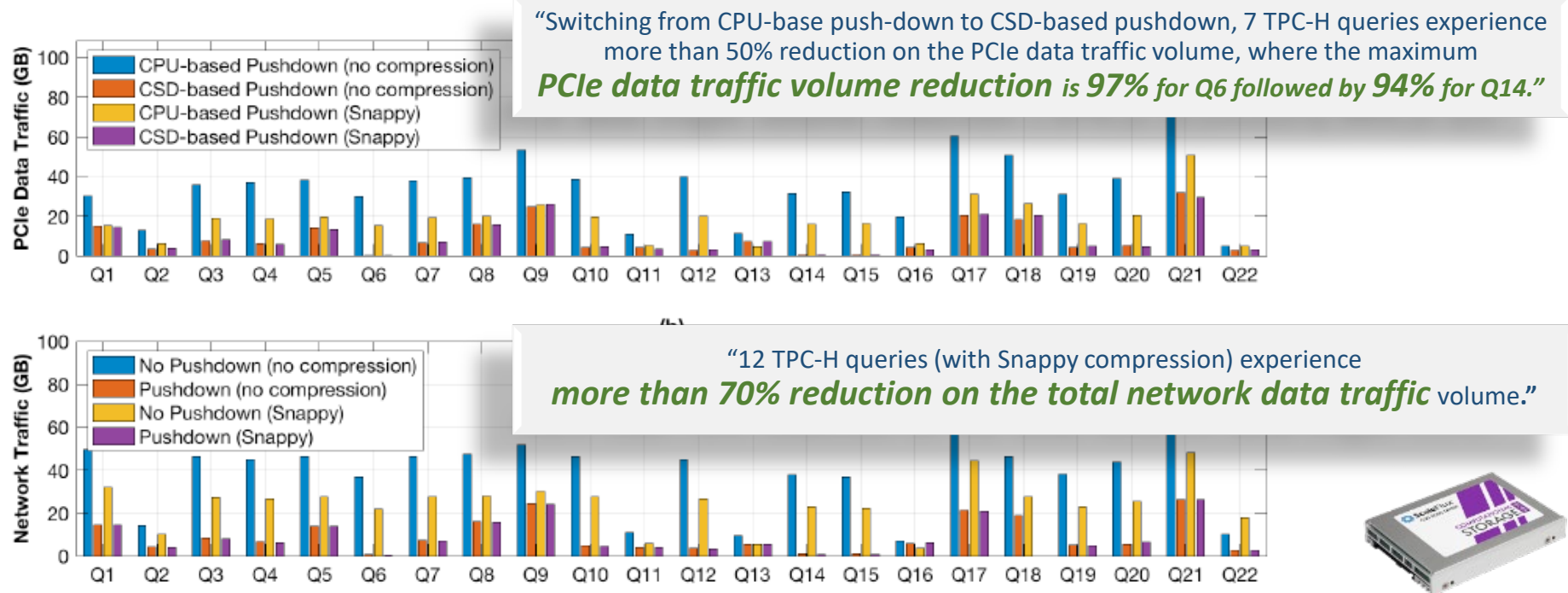
- Distributed heterogeneous compute architecture
- Push table scans to the CSD



POLARDB

# Data Filtering in the Real World

## Alibaba Cloud & ScaleFlux



POLARDB

Figure 11: (a) PCIe data traffic inside storage nodes and (b) network data traffic in the POLARDB cluster.



# Eideticom

## Eideticom's NoLoad<sup>®</sup> CSP

Purpose built for acceleration of storage and compute-intensive workloads

### 1) NoLoad Software Stack

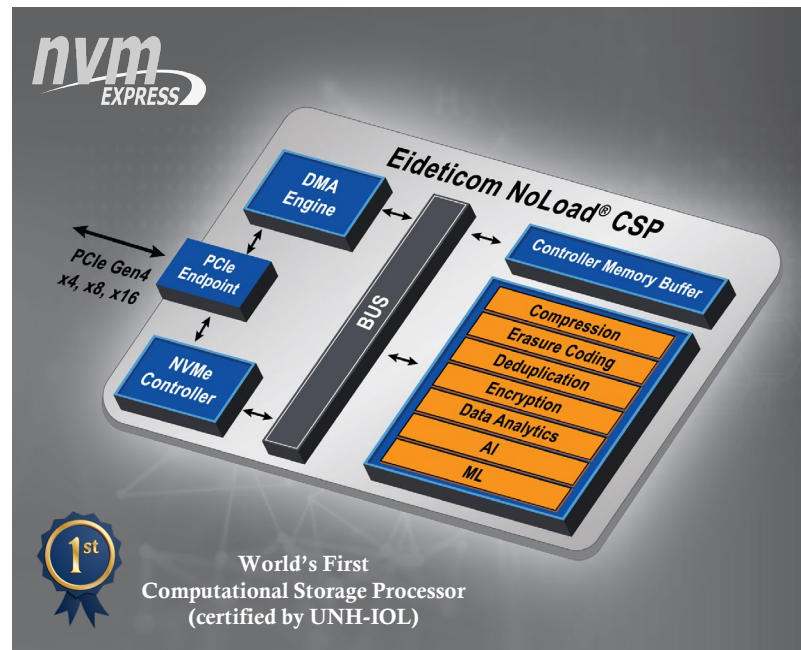
- End-to-end computational storage solution providing transparent computational offload
- Complete Software and IP core stack

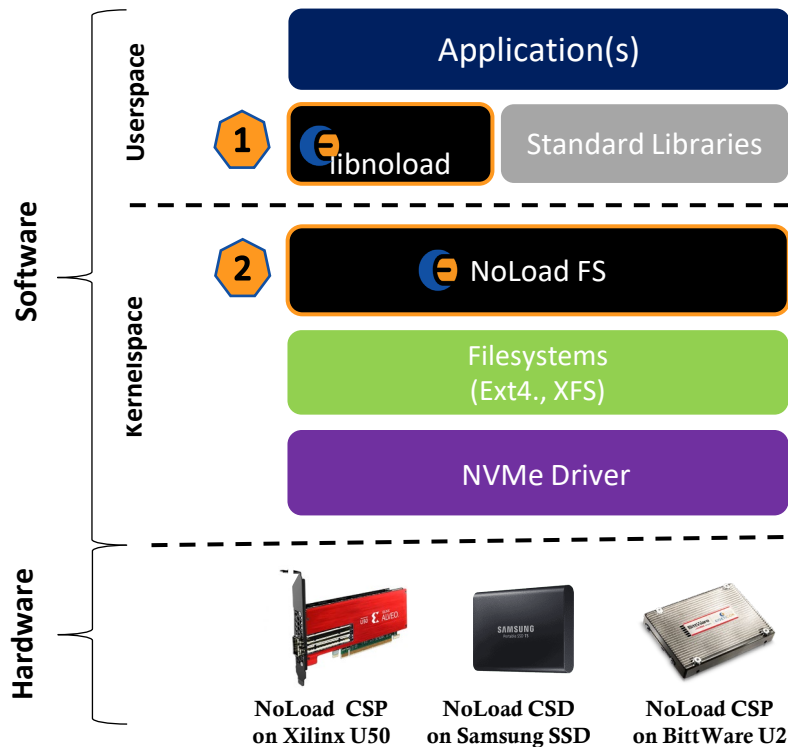
### 2) NoLoad NVMe Front End

- NVMe compliant, standards-based interface
- High performance interface tuned for computation

### 3) NoLoad Computational Accelerators

- **Storage Accelerators:** Compression, Encryption, Erasure Coding, Deduplication
- **Compute Accelerators:** Data Analytics, Video Codec, AI and ML

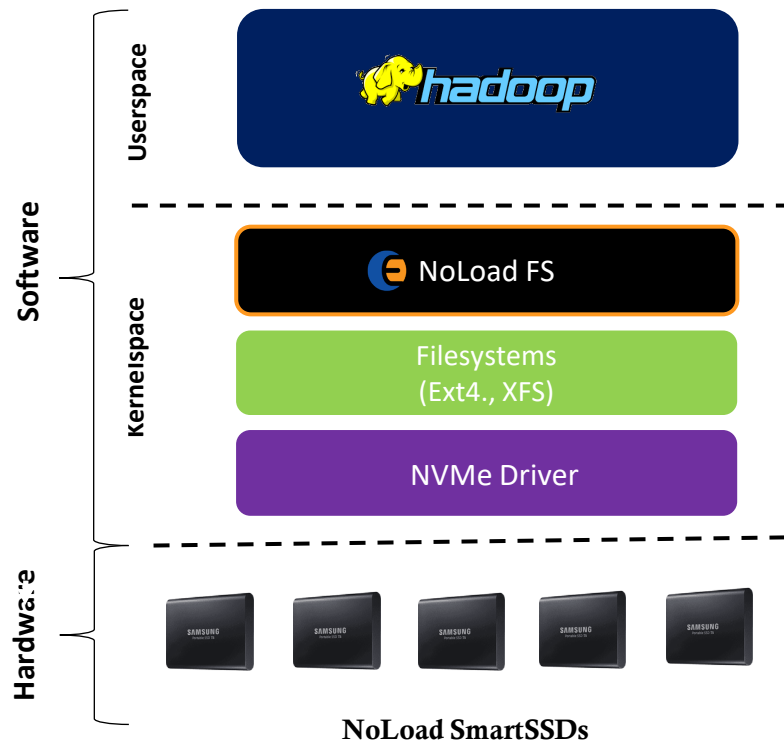
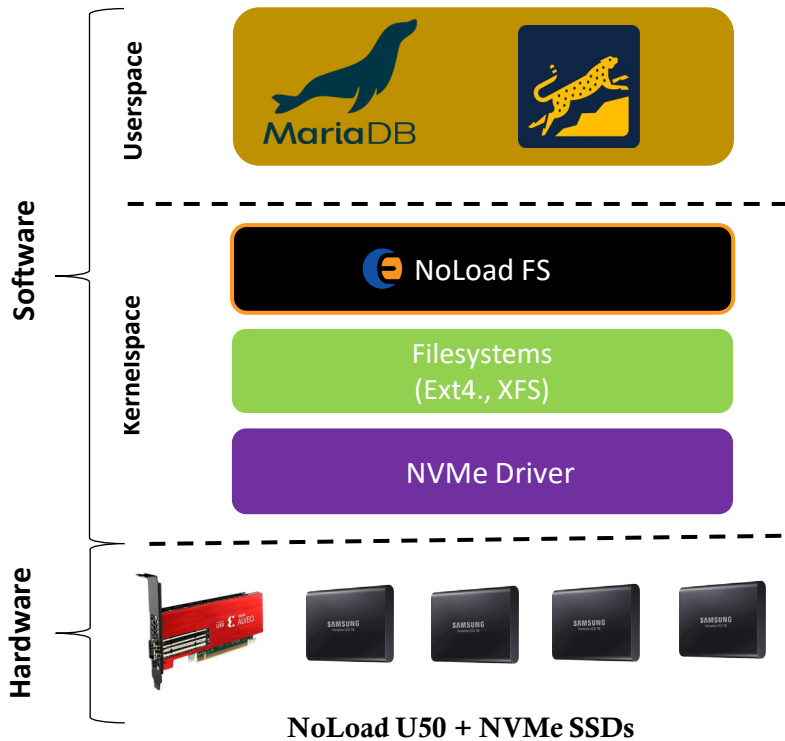




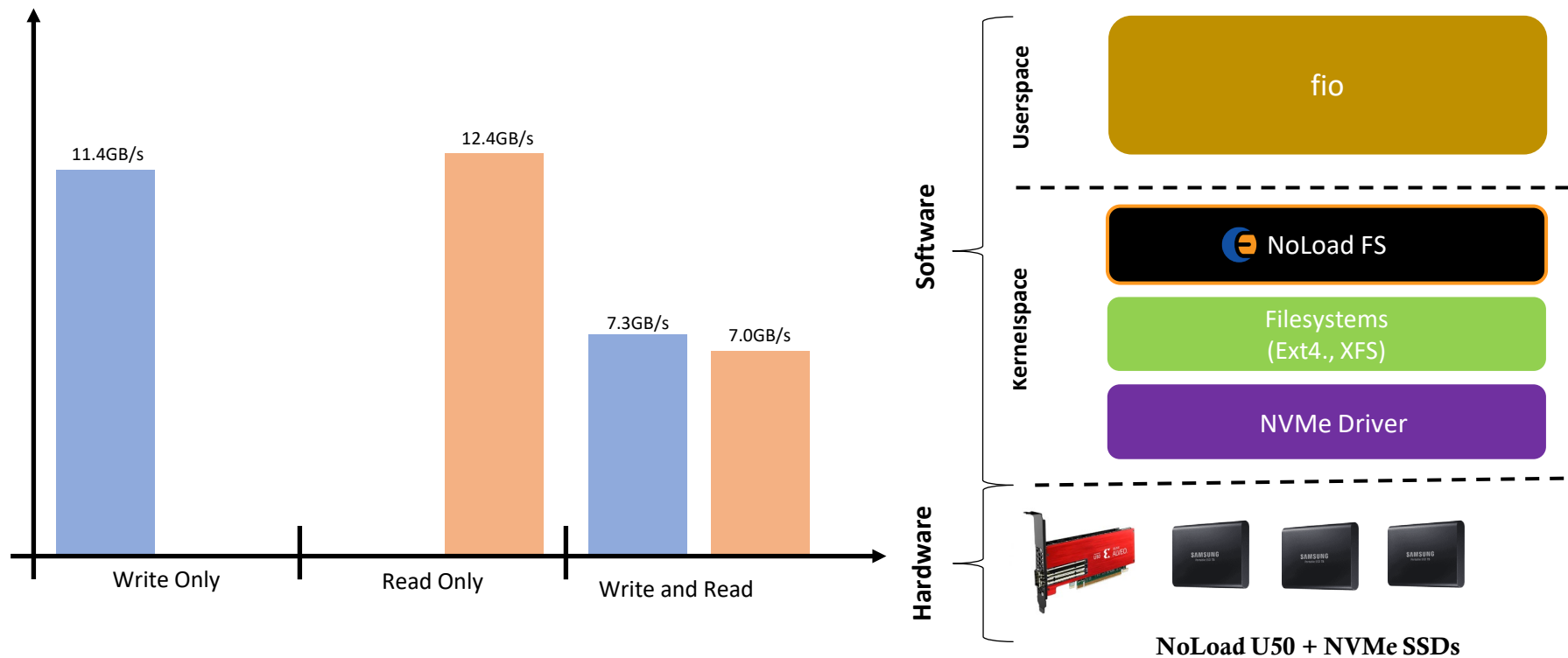
## Eideticom Software Components

- 1 libnoload: User-space library**
  - Acceleration w/o Operating System changes
- 2 NoLoad FS: Stacked-filesystem**
  - Acceleration w/o Application changes
  - Customer chooses their preferred filesystem











# Computational Storage: Real World Deployments

- Computational Storage product deployment is happening!
- Standardization will increase adoption
  - Vendor-neutral interfaces
  - Open-source software ecosystem
  - NVM Express
- NVM Express market expected to grow at >40% CAGR between 2020 and 2025. Computational storage will be a huge part of that.

01100101 00100000 01100011 01101111 01101101 01100101 01110011 00100000 01110100 01101111 00100000 01110001 01110101 01100001 01101110 01110100 01110101 01101101 00100000 01100011 01101111 01101101 01110000 01110101 01110100

# Thank You!



## Flash Memory Summit

Everything You Need To Know  
For Success