# EIDETICOM

Accelerate Everything

---

**Computational Storage Solutions Over Fabrics for ZFS**

Kelly Ursenbach

# Acknowledgements

This work is all a part of a successful partnership between:

- Aeon Computing
- Eideticom
- Nvidia
- Los Alamos National Laboratory (LANL)
- SK hynix

Much of the content provided in this talk can be attributed to:

- Brad Settlemyer – Nvidia
- Roger Bertschmann, Sean Gibb, Andrew Maier, Martin Oliveira – Eideticom
- Jeff Johnson, Doug Johnson - Aeon Computing
- Dominic Manno, Gary Grider, Jason Lee, Brian Atkinson - LANL

- Motivation

- A flexible solution
    - Accelerated Box of Flash (ABOF)

- Performance Analysis

- Outlook

# Motivation

Get maximum milage from flash storage

- Capacity

- Bandwidth

Memory bandwidth limitations observed

Get the benefits of compression without losing performance

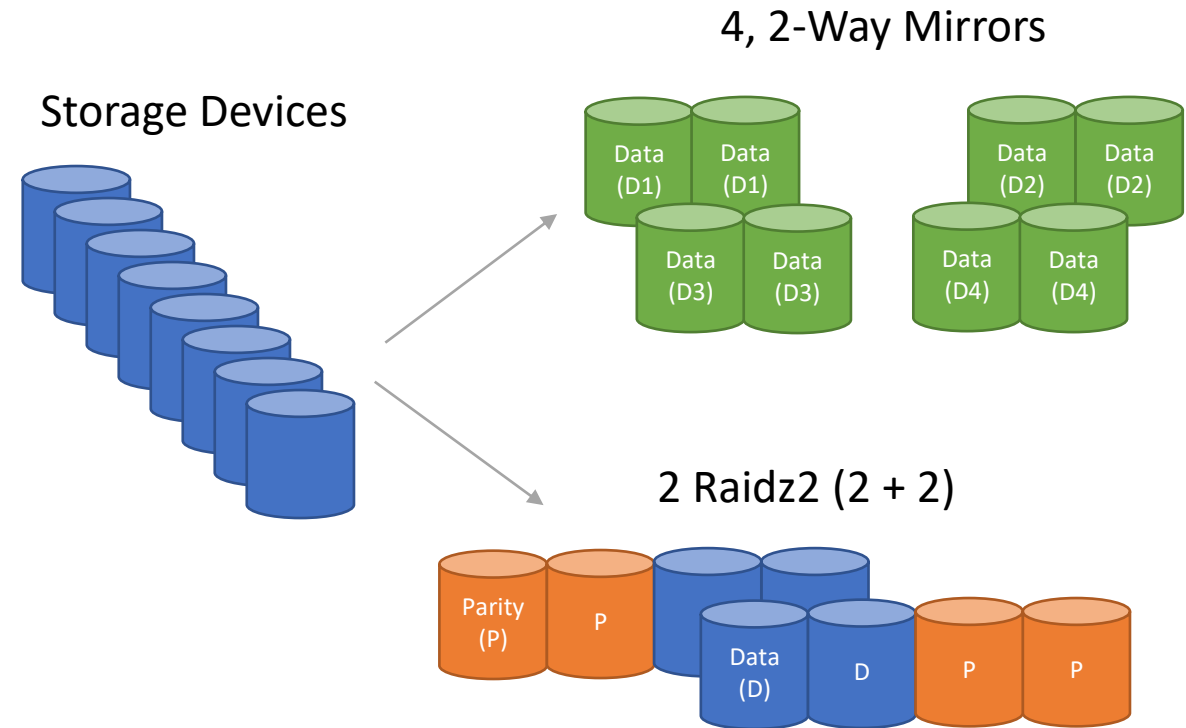# Why ZFS?

Feature Rich

- Compression
- Deduplication
- Encryption

High Integrity

- Erasure Coding (EC)
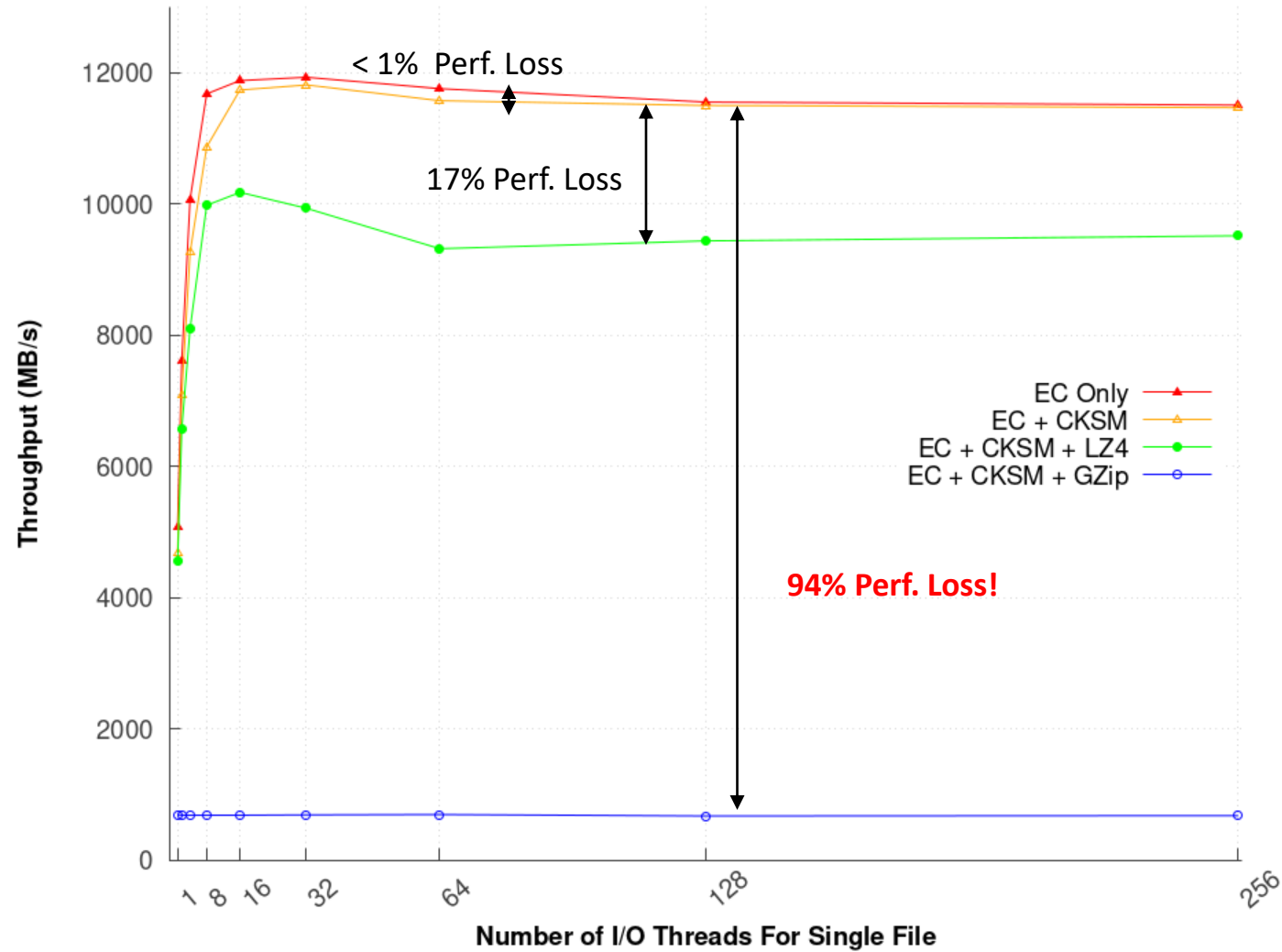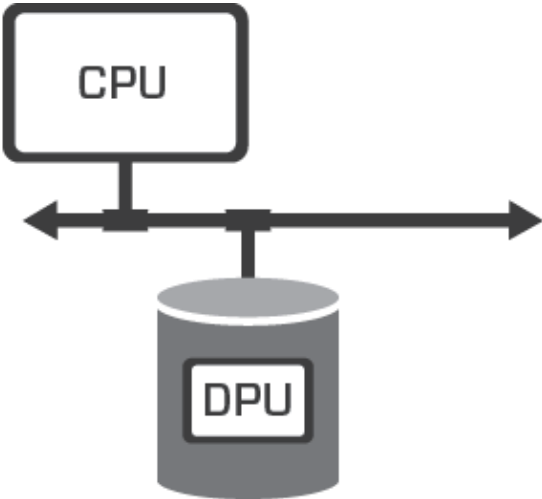- Mirrors
- Snapshots
- Checksums

Lustre over ZFS



OpenZFS

Storage Devices

4, 2-Way Mirrors

Data (D1)  Data (D1)
Data (D3)  Data (D3)

Data (D2)  Data (D2)
Data (D4)  Data (D4)

2 Raidz2 (2 + 2)

Parity (P)  P  Data (D)  D  P  P

# GZIP is Expensive on CPU



**Throughputs of 1MB Writes For Single File Using ZFS Raidz2 (10+2) Using NVMe-oF from Host to Target**

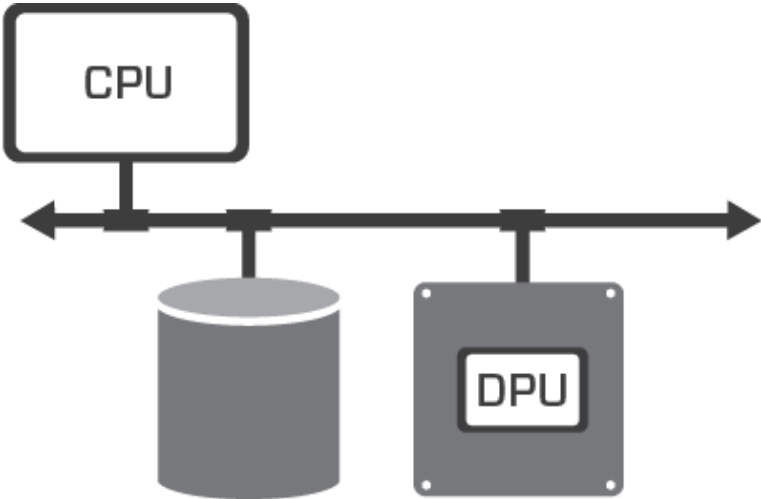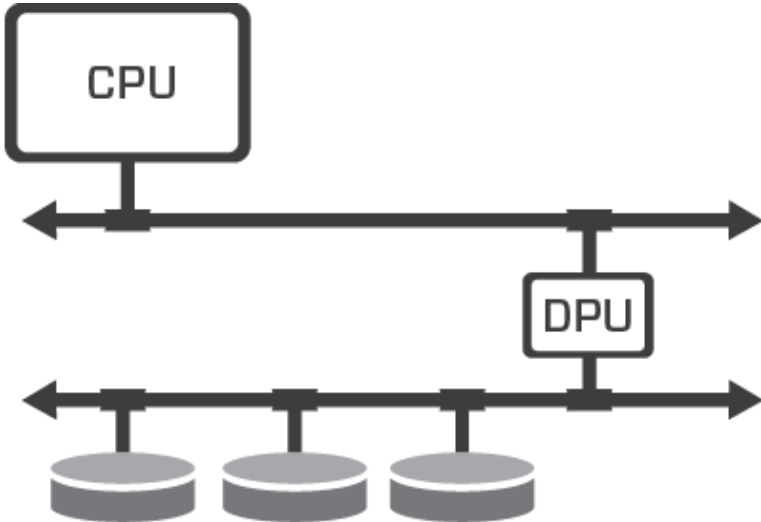< 1% Perf. Loss

17% Perf. Loss

**94% Perf. Loss!**

Legend:
- EC Only
- EC + CKSM
- EC + CKSM + LZ4
- EC + CKSM + GZip

Y-axis: Throughput (MB/s)

X-axis: Number of I/O Threads For Single File

Computational Storage
Device (CSD)

Computational Storage
Processor (CSP)

Computational Storage
Array (CSA)

# Benefits of Compression Offload

Improved storage bandwidth
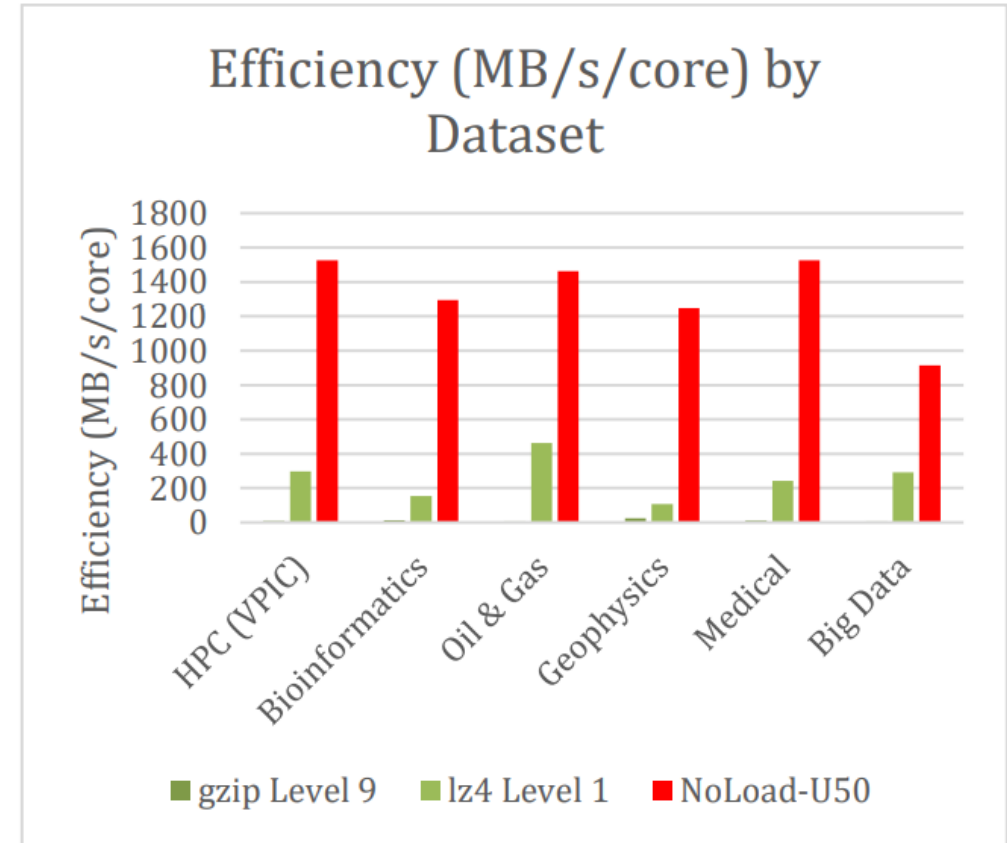
- Dedicated hardware performs near PCIe line rate

Reduced Storage Cost

- Lower power

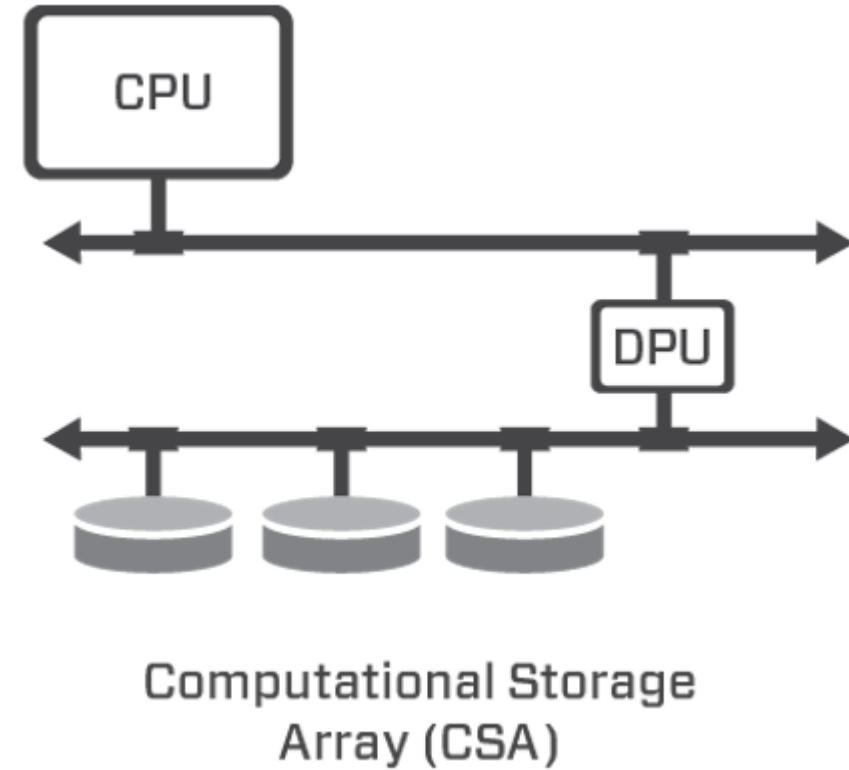- Increase effective storage via compression

Scalability

- Disaggregating compute and storage into independently scalable resources
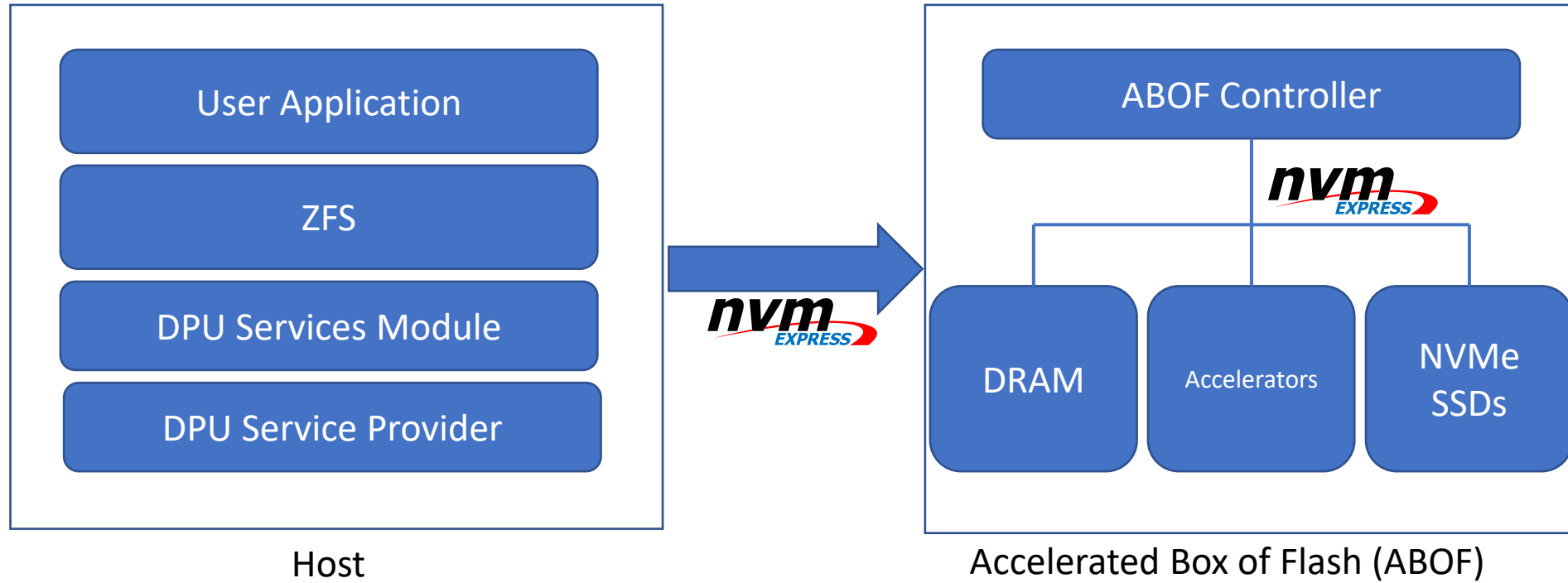
Save CPU cycles on compute nodes



Efficiency (MB/s/core) by Dataset

gzip Level 9 ■ lz4 Level 1 ■ NoLoad-U50

# Encapsulate Storage and Storage Specific Compute (CSA)

- Contain storage software stack

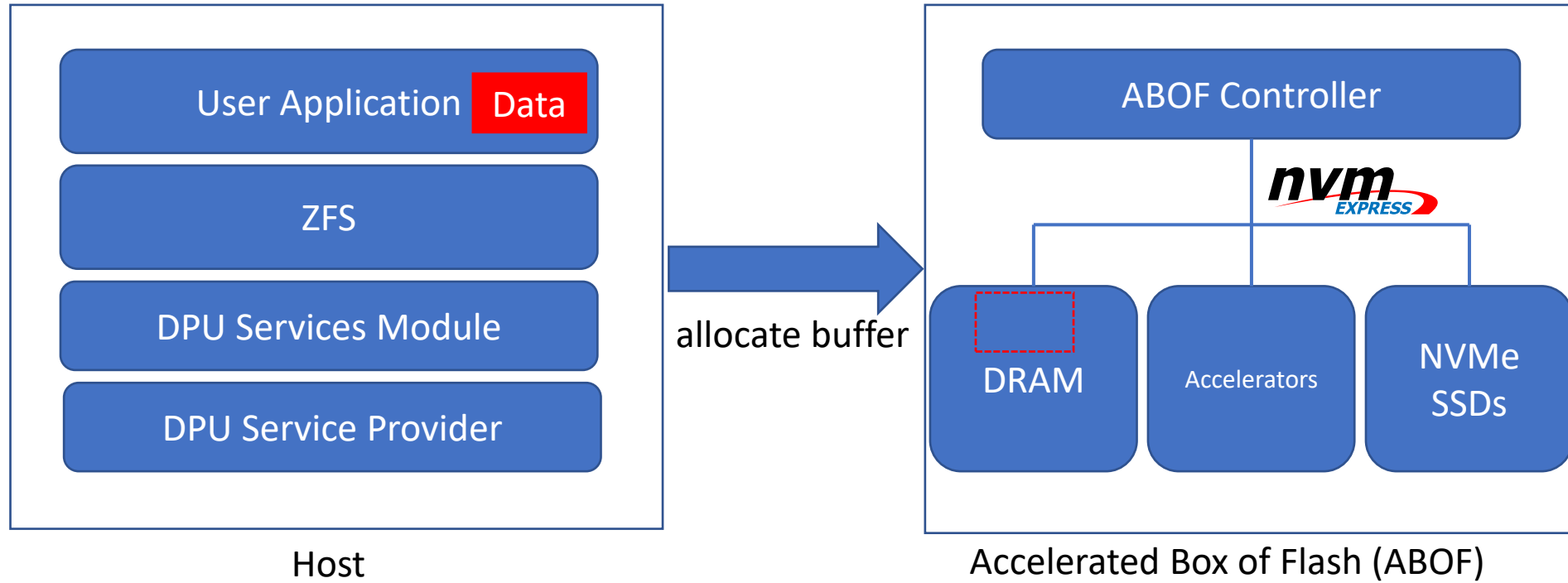- Eliminate stranded resources

- Free memory bandwidth on compute nodes



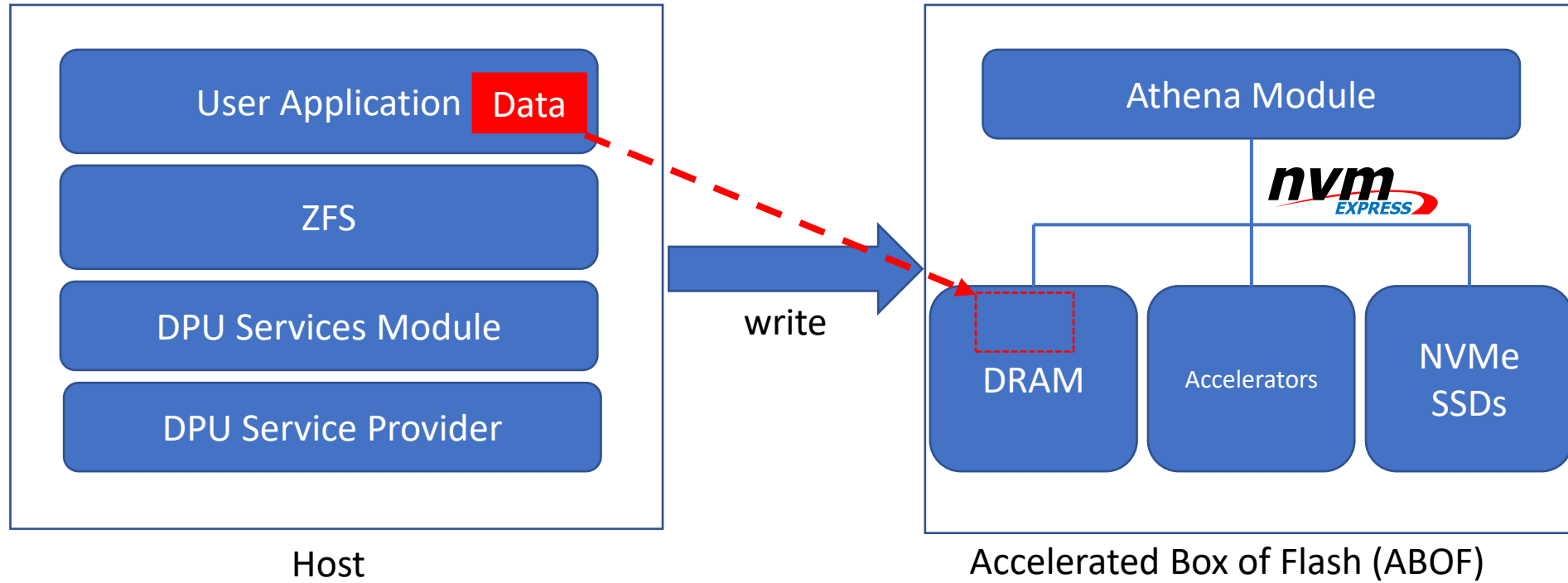Computational Storage
Array (CSA)

Host

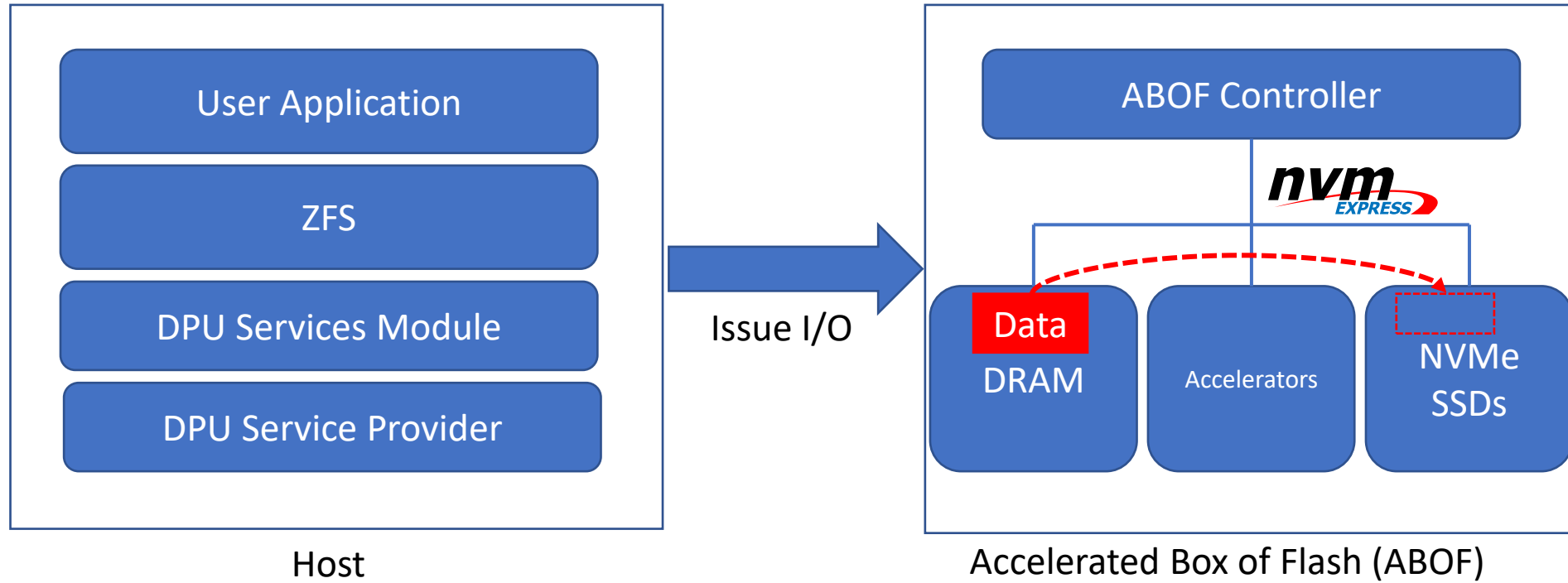Accelerated Box of Flash (ABOF)

# ABOF Operation Overview (write)



Host

Accelerated Box of Flash (ABOF)

# ABOF Operation Overview (compress)

EIDETICOM

**Host**

- User Application
- ZFS
- DPU Services Module
- DPU Service Provider

GZIP compress →

**Accelerated Box of Flash (ABOF)**

- ABOF Controller
- nvm EXPRESS
- Data — DRAM
- Accelerators
- NVMe SSDs

# ABOF Operation Overview (issue I/O)



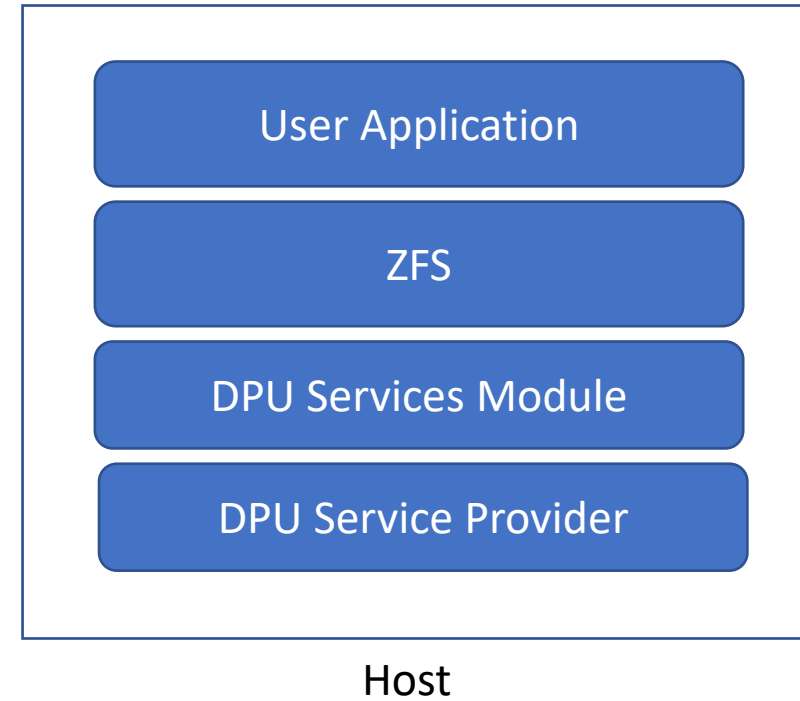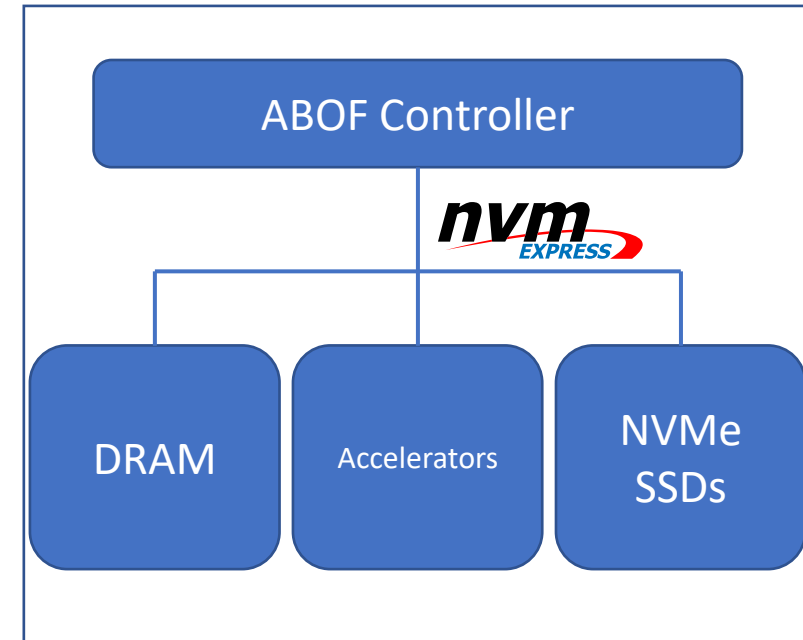Host

Accelerated Box of Flash (ABOF)

- DPUSM is an API bridge to access computational storage features via DPUSM providers

- DPUSM providers links requested operations to available architecture

- ABOF provider offloads Accelerator requests via NVMe-oF

User Application

ZFS

DPU Services Module

DPU Service Provider

Host

Use a set of vendor specific op codes via NVMe-oF to:

- create/free buffer
- load/Store buffer from disk
- read/write buffer
- Perform operation on buffer
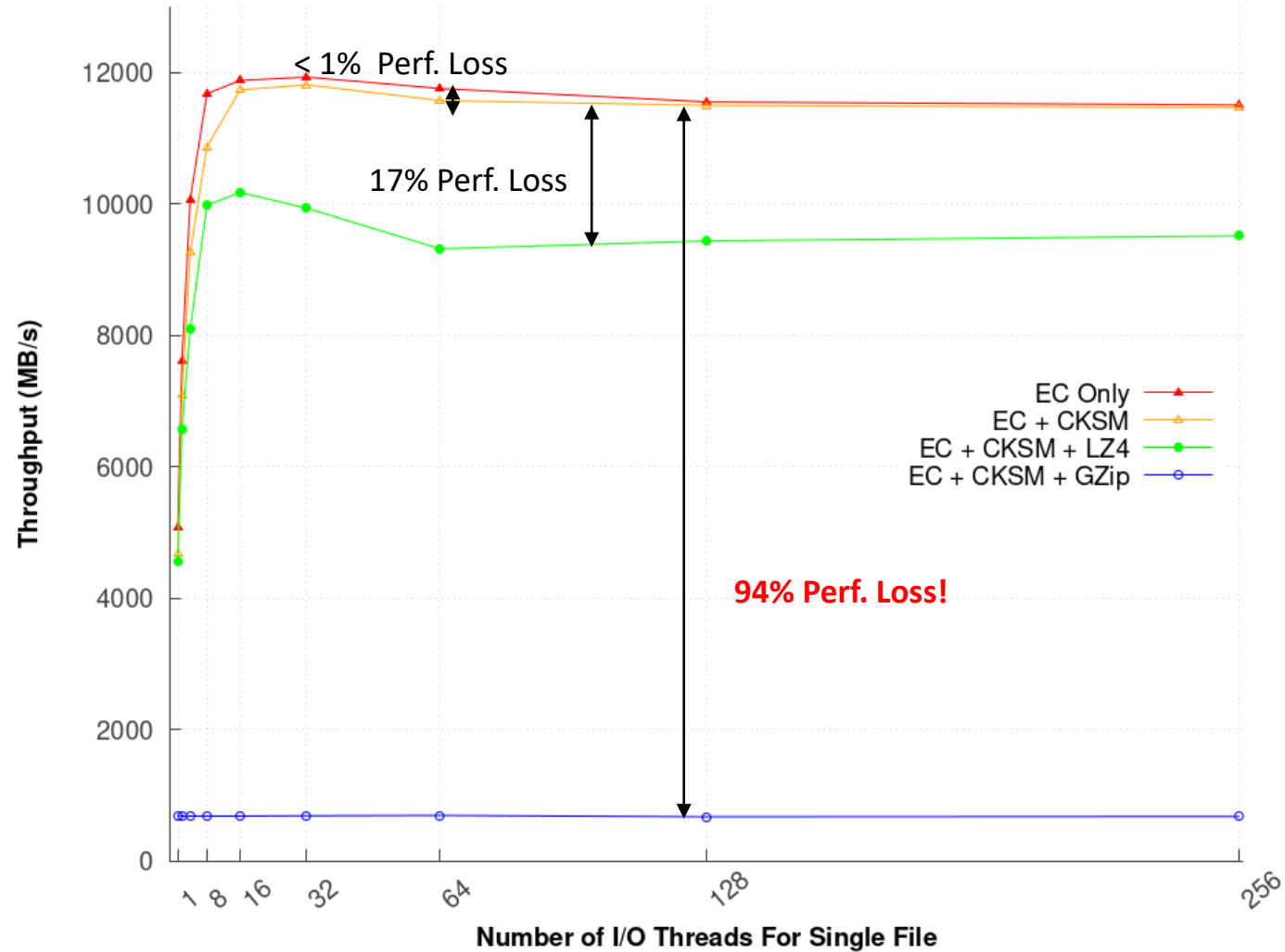  - Compress/Decompress
  - Checksum
  - EC

Kernel Module or SPDK (userspace) implementation
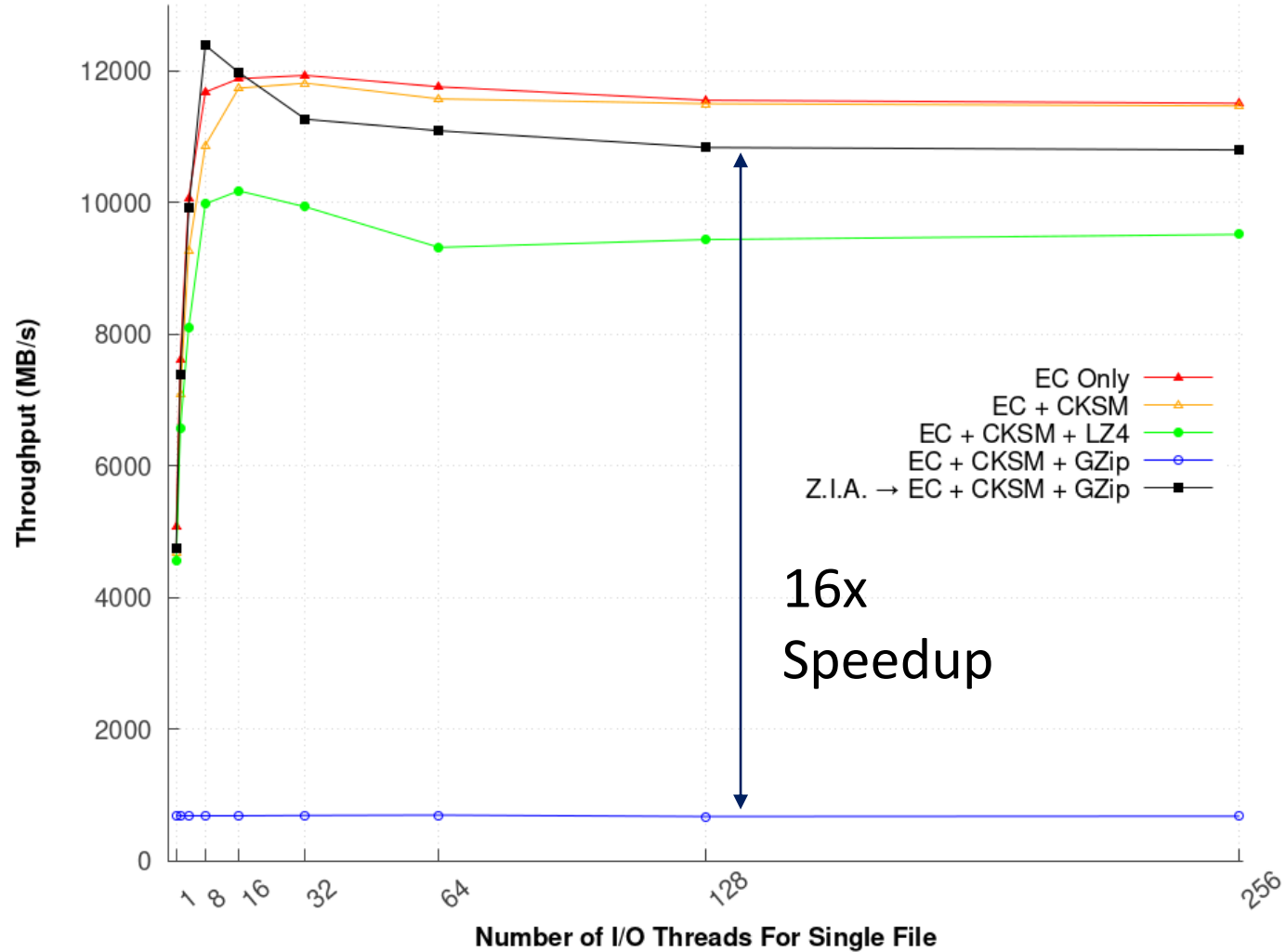


Accelerated Box of Flash (ABOF)

**Throughputs of 1MB Writes For Single File Using ZFS Raidz2 (10+2) Using NVMe-oF from Host to Target**

< 1% Perf. Loss

17% Perf. Loss

94% Perf. Loss!

Legend:
EC Only
EC + CKSM
EC + CKSM + LZ4
EC + CKSM + GZip

Y-axis: Throughput (MB/s)
X-axis: Number of I/O Threads For Single File

Throughputs of 1MB Writes For Single File Using ZFS Raidz2 (10+2) with Z.I.A. Using NVMe-oF from Host to Target

# Conclusions and Future

- Standardization
  - Reduce the cost of integration

- Upcoming ZFS direct IO feature

- Beyond ZFS
  - Additional computational offloads (Analytics, AI, etc.)

- Faster, Higher, Stronger
  - PCIe Gen5
  - 400G+ Networking