# NoLoad Filesystem: A Stacked Filesystem for NVMe Computational Storage

**David Sloan, Logan Gunthorpe,**
**Stephen Bates**

**EIDETICOM**

# Topics

- Introduction: The case for compression

- NoLoad NVMe-Based CSx

- Transparent Compression using an NVMe CSx

- Stacked Filesystems

- NoLoad Filesystem

    - Compression Method

    - Storage Architecture

    - Interface

    - Performance

EIDETICOM

# Introduction

# Introduction: The Case For Compression

- Data volumes are exploding

- NAND is getting cheaper but not as cheap as HDDs.

- NAND provides x1000 the performance of HDDs wrt IOPS, throughput and latency

- Compression can bridge the cost gap

- But it has to be performant, efficient and easy to consume!

EIDETICOM

# NoLoad
## An NVMe-Based Computational Storage Device

# NoLoad® Computational Storage Device (CSx)

## Eideticom's NoLoad® CSx

Purpose built for acceleration of storage and compute-intensive workloads
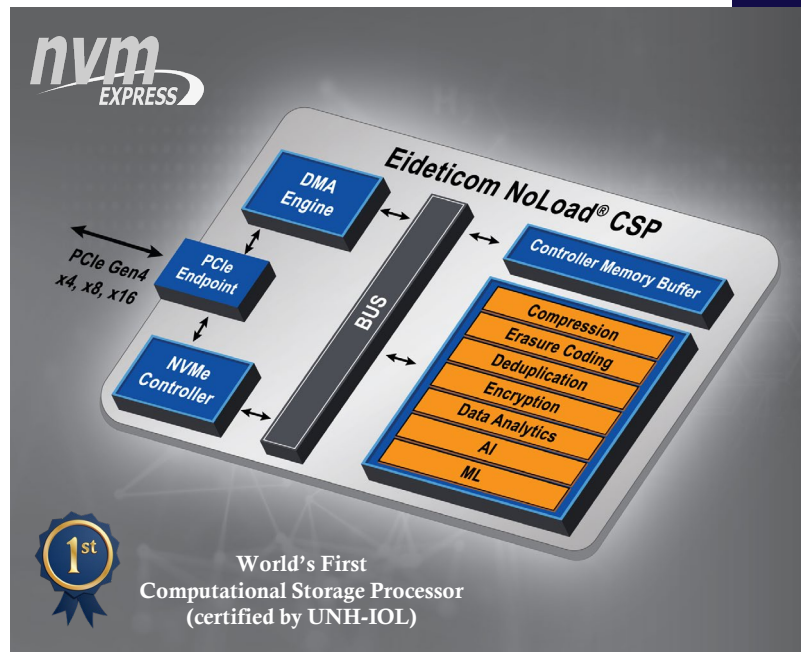
- **NoLoad Software Stack**
  - End-to-end computational storage solution providing transparent computational offload
  - Complete Software and IP core stack
- **NoLoad NVMe Front End**
  - NVMe compliant, standards-based interface
  - High performance interface tuned for computation
- **NoLoad Computational Accelerators**
  - Storage Accelerators: Compression, Encryption, Erasure Coding, Deduplication
  - Compute Accelerators: Data Analytics, Video Codec, AI and ML



*World's First Computational Storage Processor (certified by UNH-IOL)*

2020 Storage Developer Conference. © Eidetic Communications Inc. All Rights Reserved.

# SNIA Computational Storage Terminology

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



NoLoad on Xilinx Alveo U50

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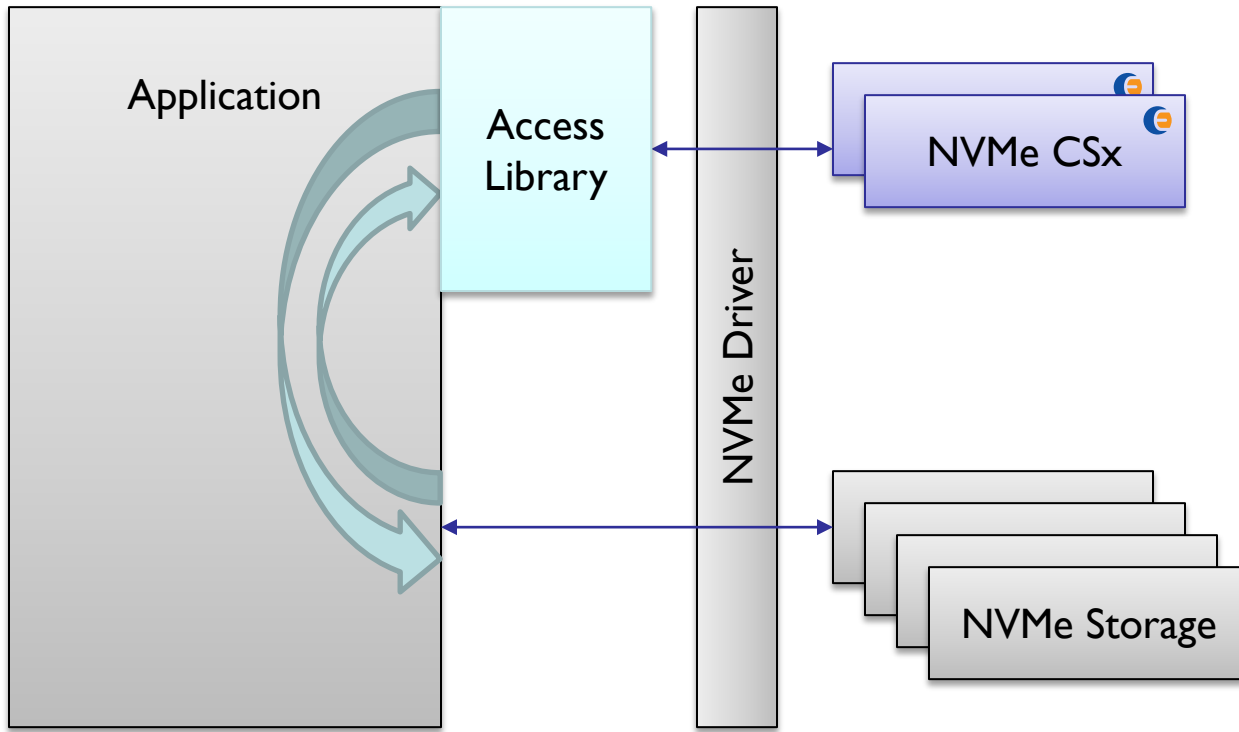A component that provides persistent data storage and computational services
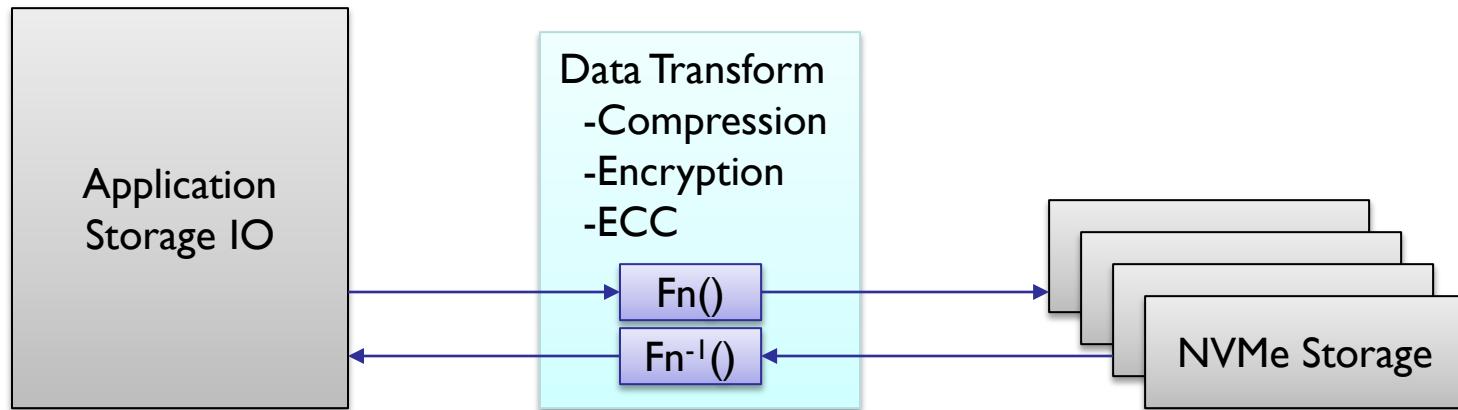


NoLoad on Samsung SmartSSD

# Transparent Compression

- Transparent compression is compression that **the *application* is unaware of**.

- This can happen in one of four places:
  - On the device (multiple vendors)
  - In the block layer (e.g. VDO)
  - In the filesystem (e.g. ZFS)
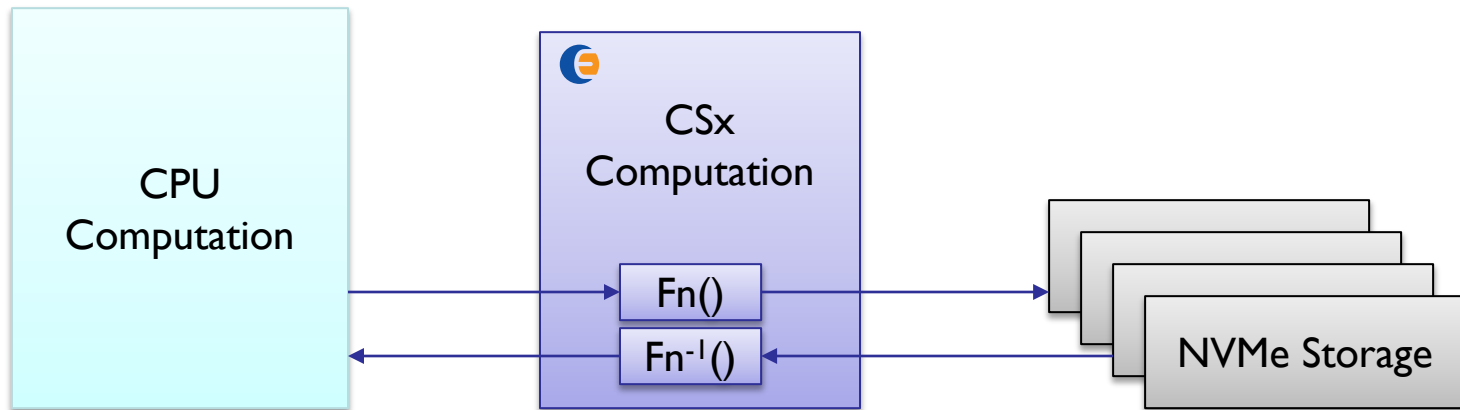  - In a stacked filesystem (e.g. this work)
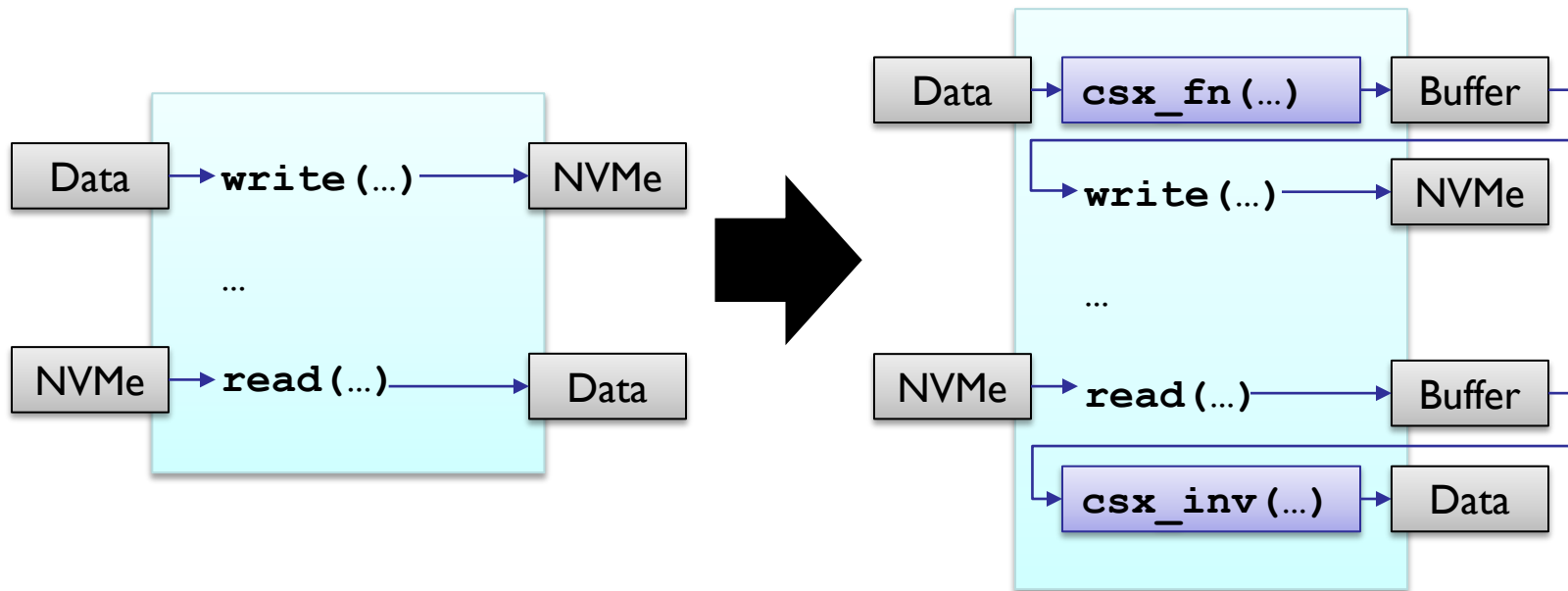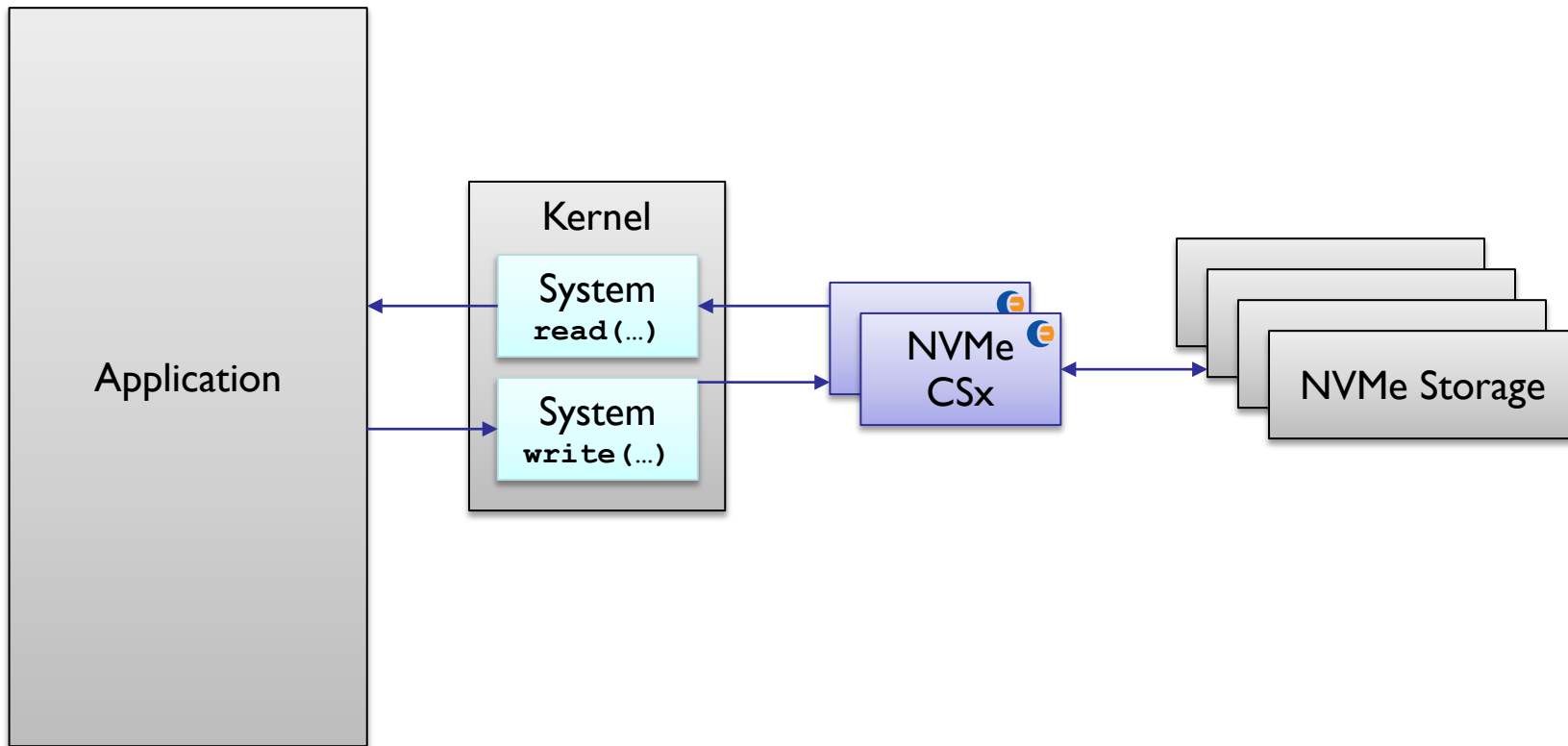
# CSx as an Accelerator



2020 Storage Developer Conference. © Eidetic Communications Inc. All Rights Reserved.

# Transparent Computation

# Transparent Compression

CPU
Computation

CSx
Computation

Fn()

Fn$^{-1}$()

NVMe Storage

EIDETICOM

# Non-Transparent Compression



2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.
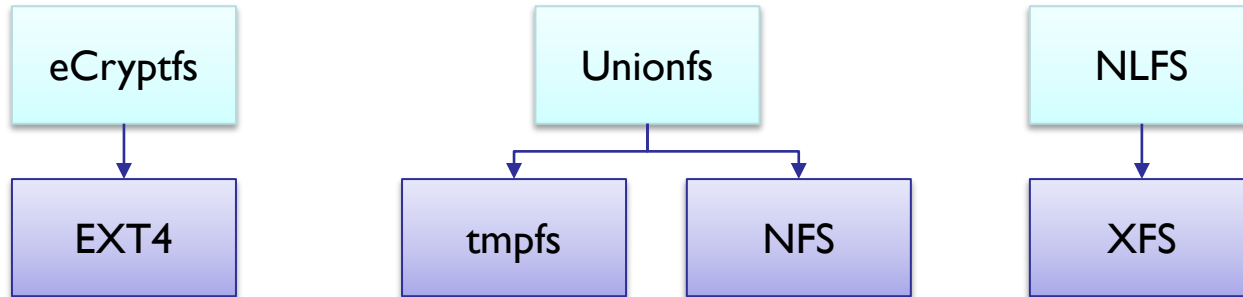
# Transparent Compression

# Stacked Filesystems
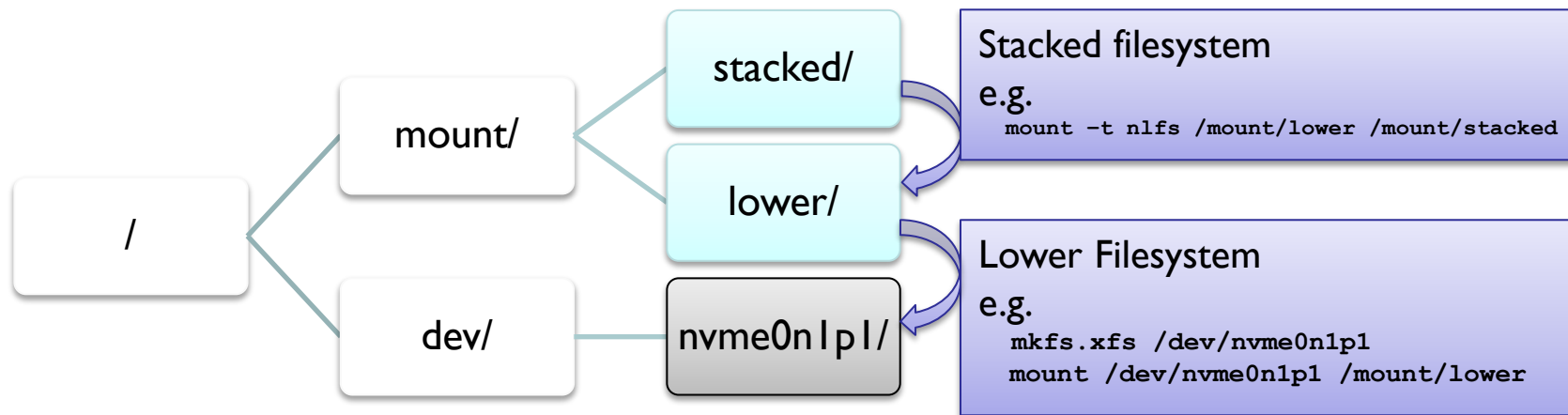
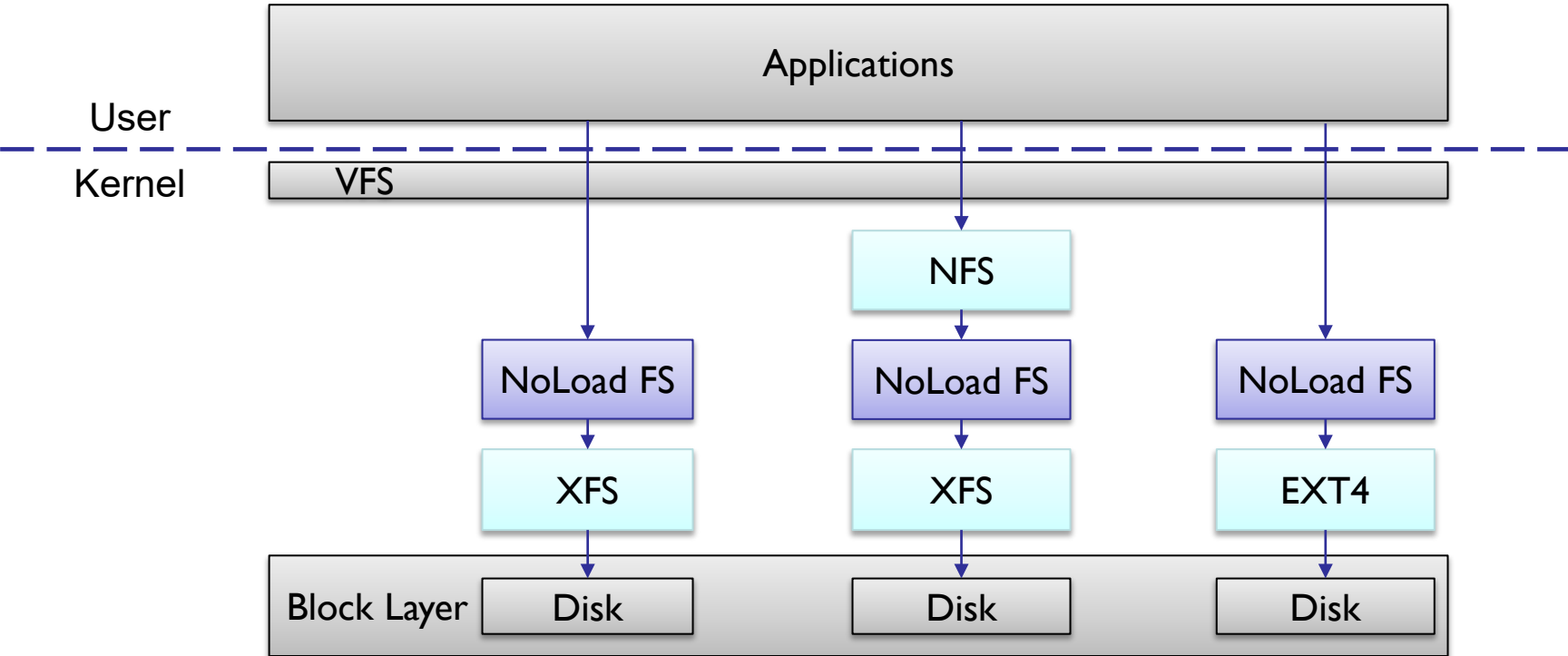EIDETICOM

# Stacked Filesystems

- Provide a code injection point between applications and the filesystem
- Adds functionality to existing filesystems
- Easy to integrate into existing SW stacks!

# Stacked Filesystem



2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.

# Stacked Filesystems



2020 Storage Developer Conference. © Eidetic Communications Inc. All Rights Reserved.

# NoLoad Filesystem

Application

Filesystem

NoLoad Filesystem

XFS, EXT4, etc.

NVMe Driver

NoLoad CSD

Accelerator

NVMe Storage

NoLoad CSP

NVMe Storage

User | Kernel

EIDETICOM

# NoLoad Filesystem



2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.

# NoLoad Filesystem

**NoLoad Filesystem**

**Thin IO Check**

File ops

Has nlfs_xattr?

No

Pass Through

**NoLoad FS Core**

xattr
getxattr(…)
setxattr(…)
…

IO ops
read(…)
write(…)
…

Yes

Compute State

NVMe CSx

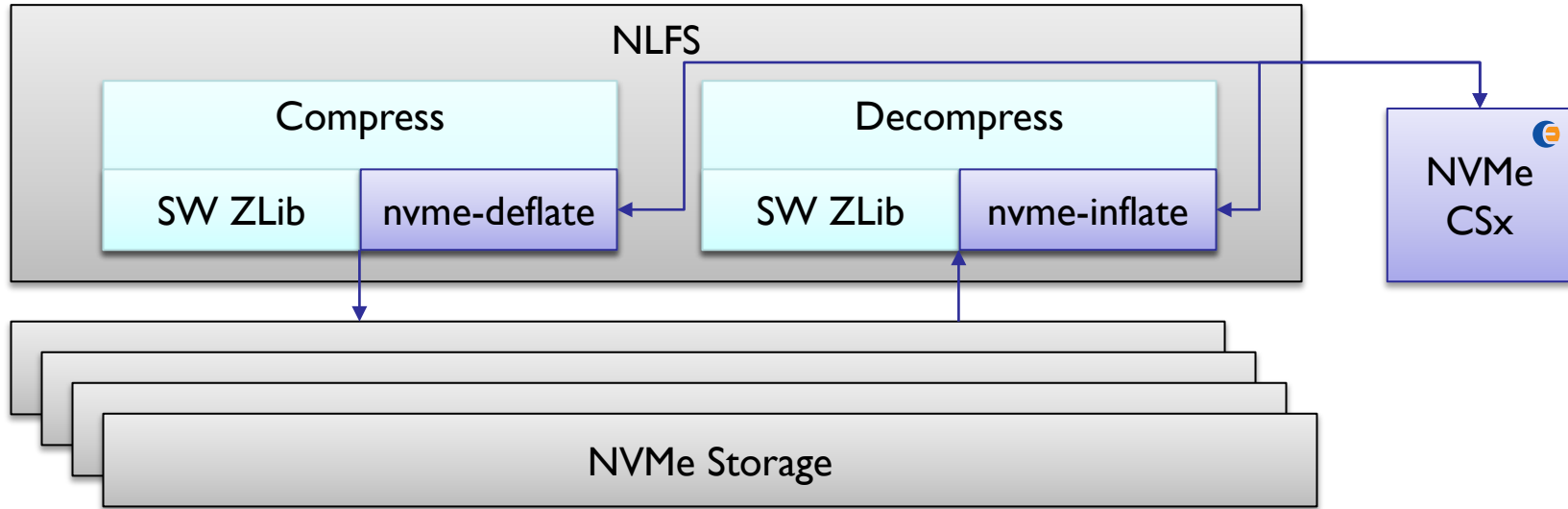**Lower Filesystem**

# NoLoad FS: Compression Method
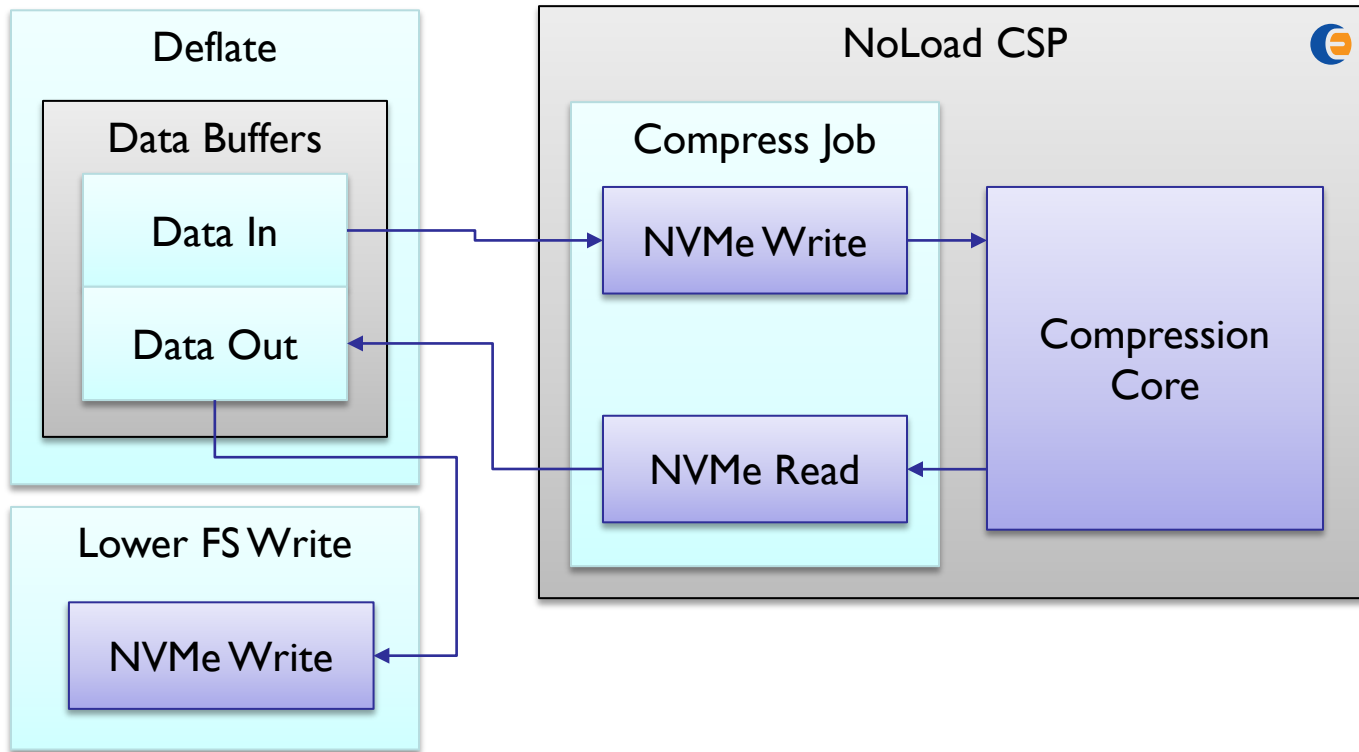
# NLFS: Compressed Filesystem

- Use existing infrastructure, Don't reinvent the wheel
- SW Compatible
  - Z-Lib encoding
- Record-based

# ZLib-Encoding



- SW Recovery
- Asymmetric topologies
- Use Existing Kernel Infrastructure

# NoLoad CSP

# NoLoad CSD



2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.

# NoLoad FS: Storage Architecture

# Storing Compressed Data



- Hole punching informs lower fs where data can be omitted

# Compression Records



- 1 MiB -> 7 kiB rounded up to 8 kiB

2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.

# Compression Records
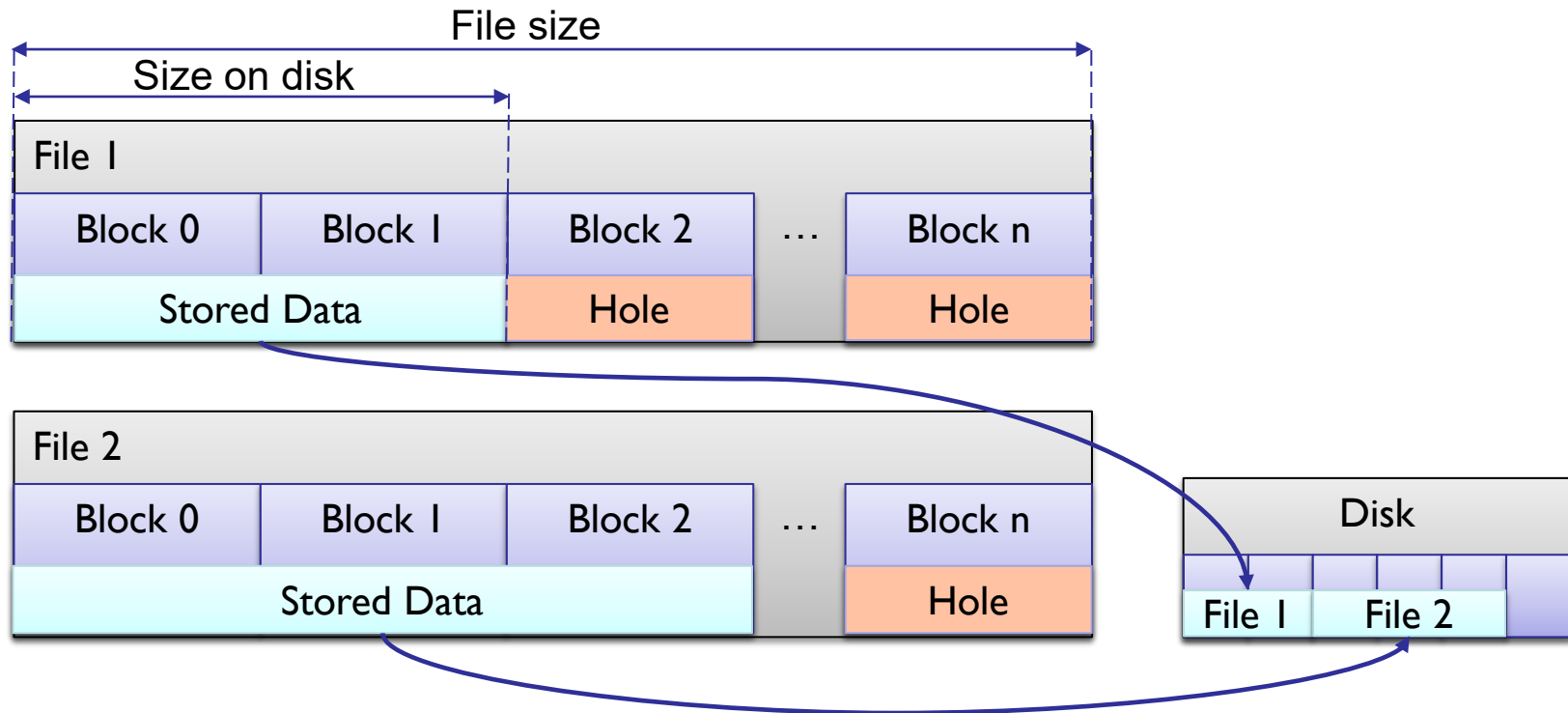
Z-Lib Compression

Data Record

- Low compressibility data is stored raw to improve read performance

Compressed Data

Lower FS

| Block 0 | Block 1 | Block 2 | ... | Block n |
|---------|---------|---------|-----|---------|

Raw Data

Raw Data

EIDETICOM

# Compression Records

# System Stats

```
$ls –lhs /nlfs/
total 128
128K –rw-r--r-- 1 dummy dummy 1.0M Aug 20 10:58 rand_ac.dat
```

Size on Disk                                    File Size

- Report compression info using built in tools
  - ls's -s option displays stored size, which can be less than the file size

# Implementation Specific Stats

```
$xattr -l rand_ac.dat
user.nlfs_recordsize: 65536
user.nlfs_recordstat:
0000 ... 8129.4340.65536.
...
00f0 ... 8129.4340.65536.
```

Stored Size    Record Size

- Allows the reporting and modification of implementation specific settings on a per-file basis

# NLFS Debug Info

NLFS Mount ID

```
$ls /sys/kernel/debug/nlfs/0
compress_records_actions_zlib      compress_ratio
compress_records_actions_nvme      compress_records_write
compress_records_actions_read      compress_records_read
decompress_records_actions_zlib    lower_path   raw_records_read
                                   mnt_path     raw_records_write
```

- compress_ratio

  - Compression ratio of mount as a whole (file size/size on disk)

- [lower|mnt]_path

  - Paths of lower fs and nlfs mount respectively

EIDETICOM

# NLFS Debug Info

```
$ls /sys/kernel/debug/nlfs/0
compress_records_actions_nvme      compress_ratio
decompress_records_actions_nvme    compress_records_write
compress_records_actions_zlib      compress_records_read
decompress_records_actions_zlib    lower_path    raw_records_read
                                   mnt_path      raw_records_write
```

- [de]compress_records_actions_x
  - Number of actions taken for record type
    - Records processed using CSx
    - Records processed using SW

# NLFS Debug Info

```
$ls /sys/kernel/debug/nlfs/0
compress_records_actions_zlib        compress_ratio
compress_records_actions_nvme        compress_records_write
compress_records_actions_read        compress_records_read
decompress_records_actions_zlib   lower_path   raw_records_read
                                  mnt_path     raw_records_write
```

- compress_records_x

  - Number of compressed records read/written to disk

- raw_records_x

  - Number of raw records read/written to disk

# Usage Example

```
$ls /xfs/
raw.dat
$mount –t nlfs /xfs /nlfs –o recordsize=1M
$touch /nlfs/sequential.dat
$touch /nlfs/rand_ac.dat

$xattr –w user.nlfs_recordsize 65536 /nlfs/rand_ac.dat
$cp /nlfs/raw.dat /nlfs/sequential.dat
$cp /nlfs/raw.dat /nlfs/rand_ac.dat
$sync

$ls –lhs /nlfs/
total 1.2M
128K –rw-r--r-- 1 dummy dummy 1.0M Aug 20 11:23 rand_ac.dat
8.0K –rw-r--r-- 1 dummy dummy 1.0M Aug 20 11:23 sequential.dat
1.0M –rw-r--r-- 1 dummy dummy 1.0M Aug 20 11:22 raw.dat
```

- Compressible file on lower fs
- Optimize file for random access
- Create new compressed files
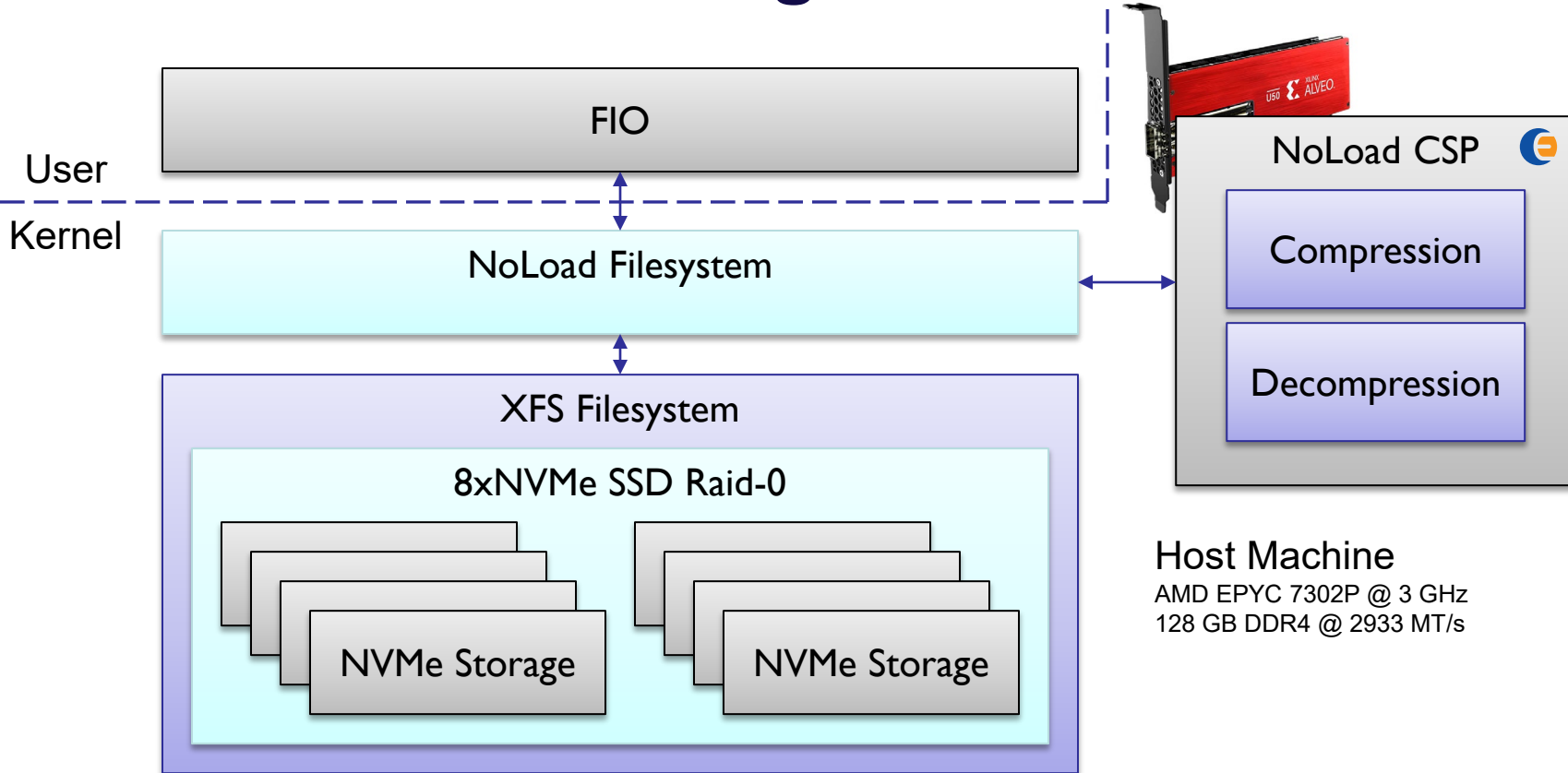- Flush any cached data to disk
- Small record file
- Large record file
- Pre-existing raw file

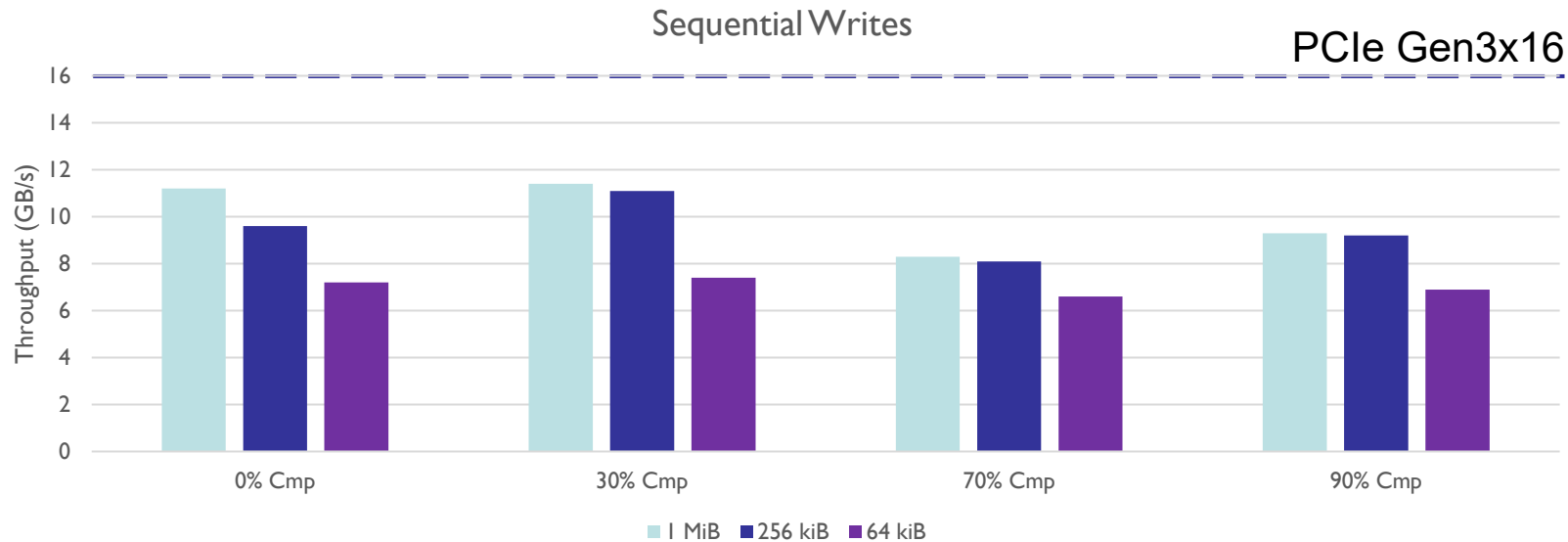# NoLoad FS: Performance

# Configuration



FIO

User

Kernel

NoLoad Filesystem

XFS Filesystem

**8xNVMe SSD Raid-0**

NVMe Storage

NVMe Storage

NoLoad CSP

Compression

Decompression

## Host Machine
AMD EPYC 7302P @ 3 GHz
128 GB DDR4 @ 2933 MT/s

# Write Throughput vs Record Size



Sequential Writes

PCIe Gen3x16

Legend: 1 MiB, 256 kiB, 64 kiB

X-axis categories: 0% Cmp, 30% Cmp, 70% Cmp, 90% Cmp

Y-axis: Throughput (GB/s)

# Read Throughput vs Record Size



Sequential Reads

2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.

# Write Latency vs Record Size



Random Writes

2020 Storage Developer Conference. © Eidetic Communications Inc.  All Rights Reserved.

# Read Latency vs Record Size



Random Reads

2020 Storage Developer Conference. © Eidetic Communications Inc. All Rights Reserved.

# Conclusions

- Data volumes are exploding, and NAND is getting cheaper. Hardware accelerated compression can help bridge the gap

- NVMe CSx devices allow high speed injection of HW accelerated data transforms into applications

- Stacked filesystems allow for greater configuration flexibility and leverage existing filesystem optimizations

- A Stacked filesystem can be inserted into existing infrastructure with minimal effort

Please take a moment
to rate this session.

Your feedback matters to us.