

Turbo-charging the Cloud with Zoned Storage

Dennis Maisenbacher
Principal, Emerging System Architectures
Western Digital



Overview

- Cloud-Native Storage Challenges & Zoned Storage
- Work on Enabling Zoned Storage as a 1st Class Cloud Citizen
 - Longhorn
 - OpenEBS - Mayastor
 - spdk-csi - CSAL
- Showcase

Cloud-Native Storage Challenges & Zoned Storage

Cloud-Native Storage at Scale

- Cloud Service Providers (**CSPs**) are constantly **challenged** with **large volumes of data** and increasing **customer demand** for **cost-effective storage and high performance**.
 - Measurements: IOPS/TB, GB/\$, TB/Rack
- To solve these challenges, CSPs provide **software-defined storage solutions** by exposing the **data as a service** that implements **multiple storage tiers** to provide the appropriate benefits.

CSP's Challenges – Storage Hardware

Storage tiers have inherent **physical characteristics**, which, if to be used in the cloud, need to be **effectively managed**

Flash-based SSDs (Performance Tier):

- Media type (TLC/QLC), WAF, DRAM, OP → **Cost (GB/\$)**
- Relatively expensive compared to HDDs. The characteristics of conventional SSDs are such that they exhibit:
 - **Write amplification** which impacts **lifetime, performance, and overall cost**.
 - Typical **overheads** include **2x cost** (CacheLib paper) and **3-4x reduced lifetime** (rw workload *), reduced lifetime (contradicts the requirement of deploying capacity storage (i.e., QLC) for 5-7 years)
 - Bandwidth **penalties** in single device - **multi tenant** setups.

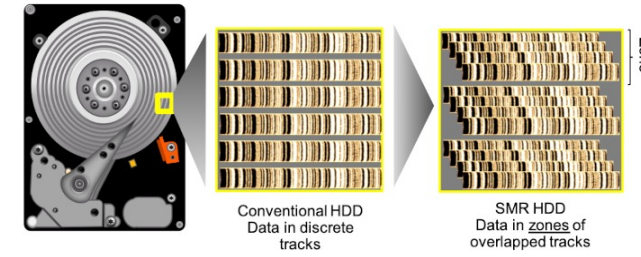
→ **ZNS SSDs** solve these issues while also maintaining **high performance**.

HDDs (Capacity Tier):

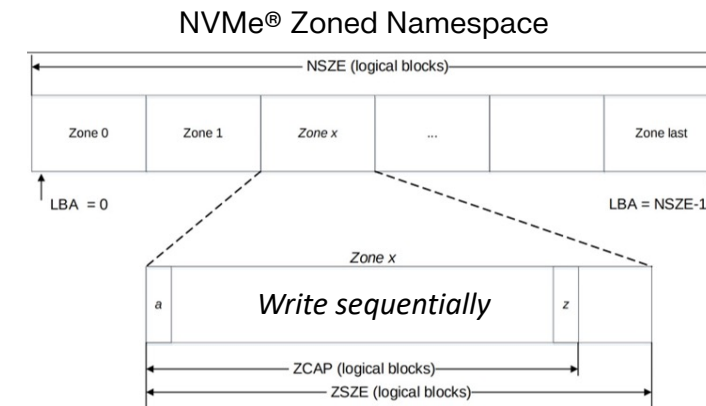
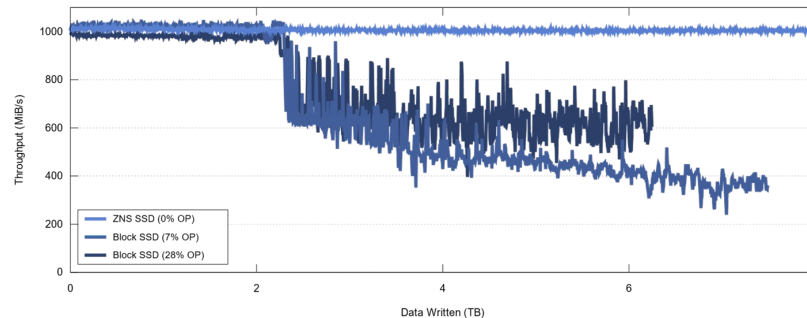
- Rotation speed, actuator and platter count → **Cost and Throughput per GB (IOPS/GB & GB/\$)**
- **Large-scale HDD deployments** (power and space constrained) are **improved** using SMR technology (22TB vs 26TB)

The Benefits of Zoned Storage

SMR HDDs offer 18% additional capacity with the newest 26TB SMR HDD over conventional CMR HDD.



SSDs with Zoned Namespace support offer 3-4x higher performance, 7-28% higher capacity, and usage of QLC in write-heavy workloads.



- Deployed heavily at hyperscalers and in large-scale storage clusters (SMR HDDs and ZNS SSDs)
- To take advantage of SSDs with **Zoned Namespace** support and **SMR** HDDs, the software requires support

Todays Linux® Eco-System

Official Linux® support since 2016

- Zoned API from kernel version 4.10 (Feb 2017)
- ZNS support added in kernel version 5.9 (Oct 2020)
- UFS support to be available when standardized (~2023)

5+ Linux® Distributions with Zoned Storage Support

- RHEL 9+, CentOS 7+, Fedora 33+, Debian 11+, and Ubuntu 21.04+

Two File-systems supports Zoned Storage

- f2fs (client - UFS) and btrfs (enterprise - ZNS/SMR)

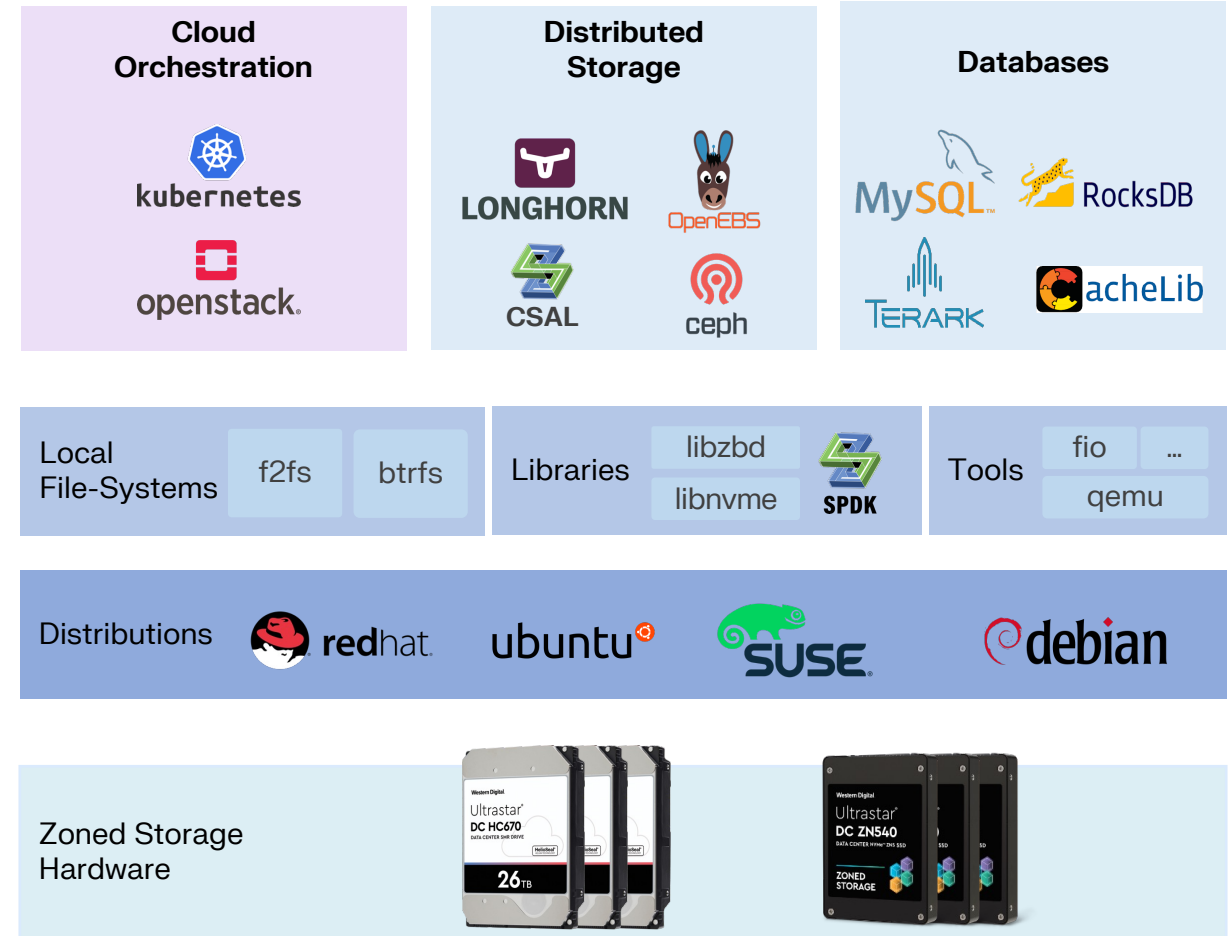
Library/Tools support

- libzbd, libnvme, SPDK, fio, qemu, blkzone, blktests, ...

Broader Enablement

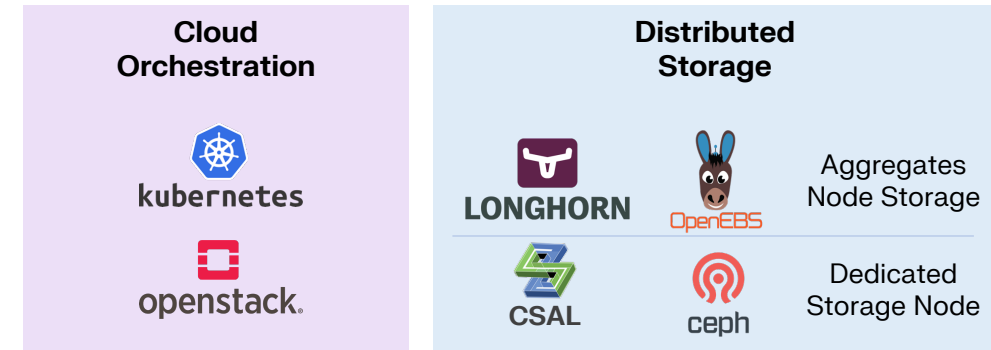
- Cloud Orchestration, Databases, Distributed Storage, Databases, Caching, Generic storage

Mature, robust, and adopted by some of the biggest consumers of storage



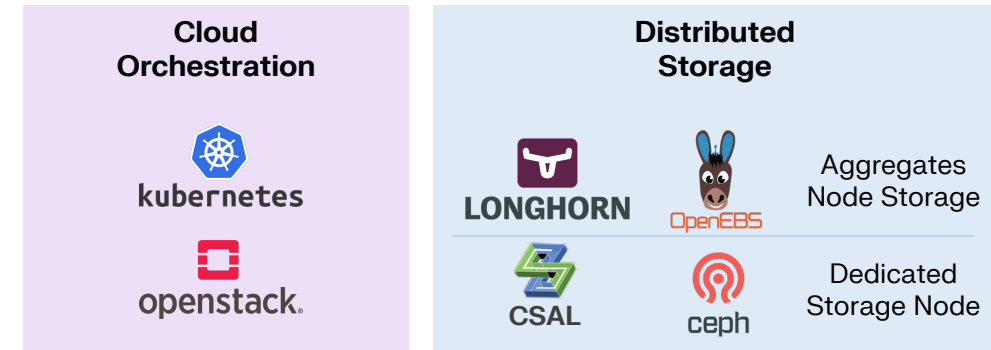
Previously Missing Link for Cloud Applications

- Tight integration with the Cloud Orchestration Platform required.
- The Storage as a Service is provided by the specific CSP (Azure, AWS, GCP) which implements internal solutions to gain the advantages.
- But how to get the benefits with on-premise clouds and/or hybrid clouds?



Previously Missing Link for Cloud Applications

- Tight integration with the Cloud Orchestration Platform required.
- The Storage as a Service is provided by the specific CSP (Azure, AWS, GCP) which implements internal solutions to gain the advantages.
- But how to get the benefits with on-premise clouds and/or hybrid clouds?



Enablements

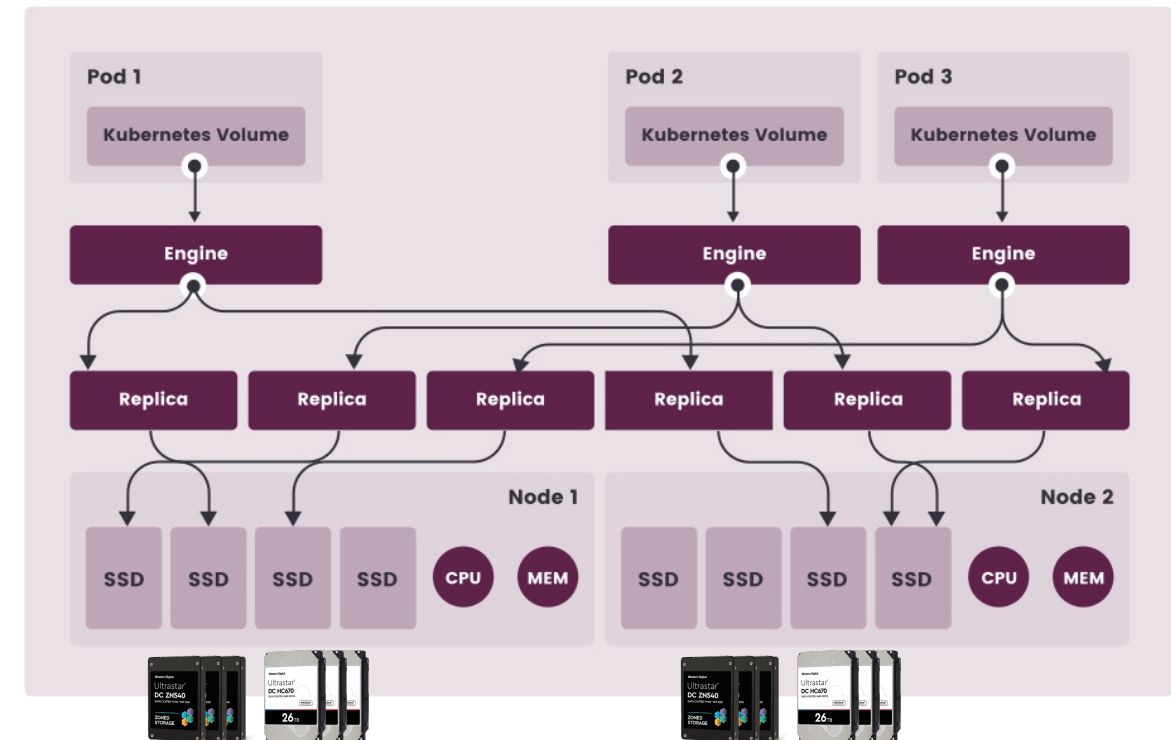
- **Longhorn** and **CSAL** provide a **generic block device/file-system** abstraction for containers using zoned storage
- With **Mayastor**, we made it such that **zoned storage** devices can be **explicitly exposed** to a container!

→ Ceph - FMS Talk tomorrow by *Aravind Remesh*
“ZNS in cloud with Ceph Crimson-OSD”


Zoned Cloud Native Storage

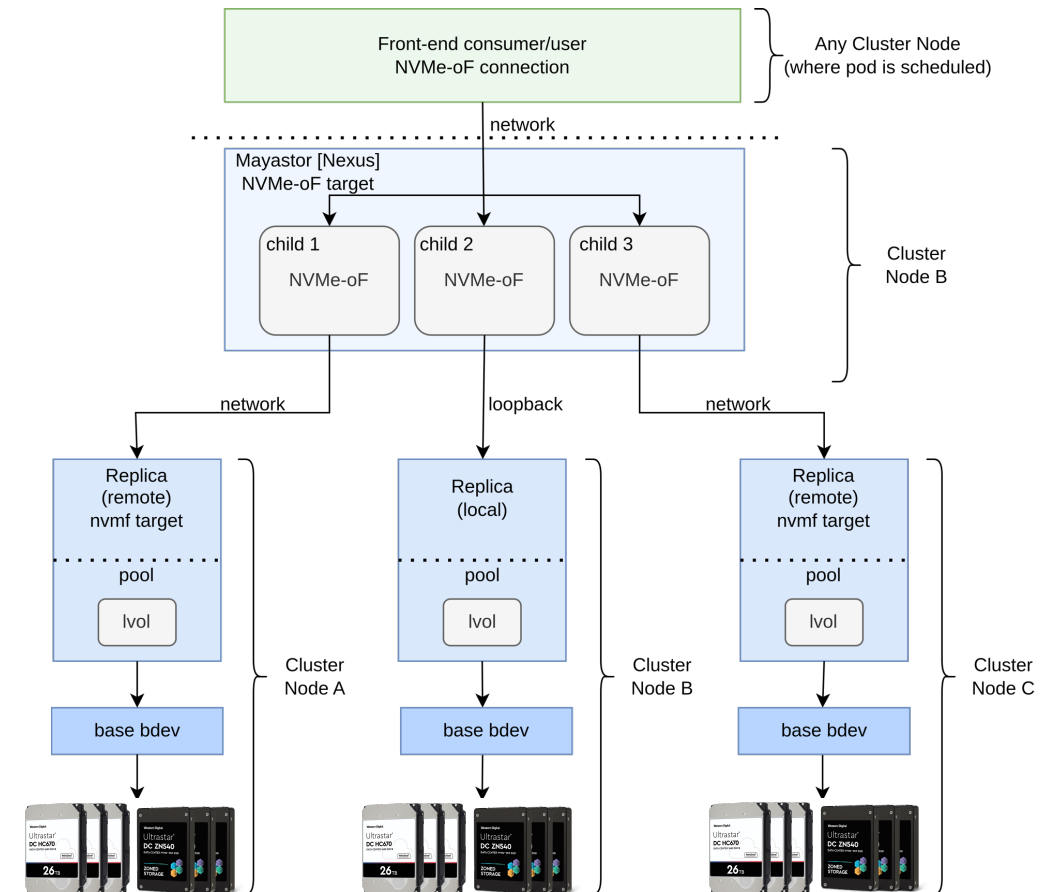
1) Longhorn – Zoned Cloud Native Storage

- Longhorn is a persistent distributed storage system with no single point of failure
 - A replica is accessing the SSD through the file system abstraction
 - Containers accesses a POSIX compatible filesystem as the k8s volume or optionally a conventional block device
- Changes necessary
 - Make use of the native support for zoned storage in BTRFS
 - Multi-year effort on enablement now makes it easy to integrate
 - Works, with minimal configuration, out of the box



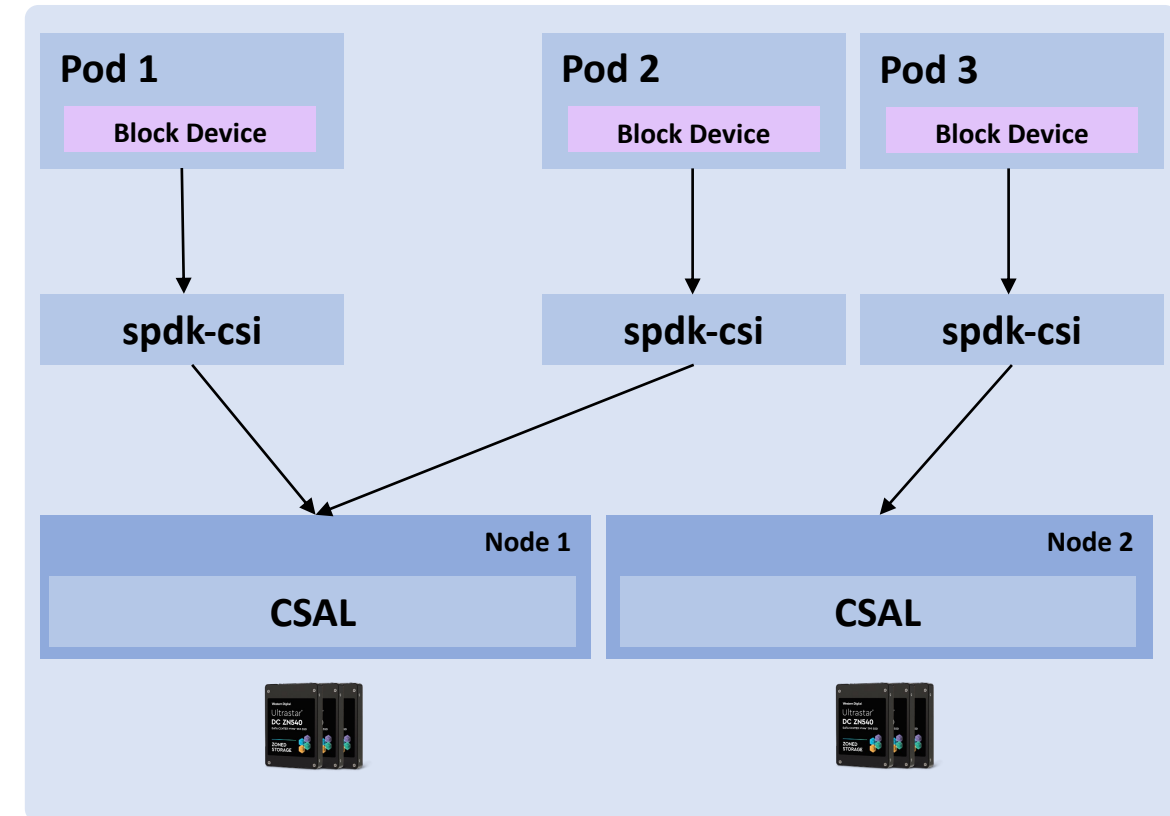
2) Mayastor – Zoned Cloud Native Storage

- OpenEBS - Mayastor (written in Rust ) is much like Longhorn a persistent highly available Container Attached Storage (CAS) solution
 - Utilizes SPDK → Avoiding storage related kernel modules for node uptime
 - Ongoing plumbing work in fixing up the Mayastor I/O path to natively handle ZNS NVMe block devices → First proposal with replica factor 1
 - Interesting challenges with zoned replication:
 - Keep zone write pointer across replicas in sync
 - What happens on I/O errors?
- CSI driver consumer will receive a zoned block device via NVMe-oF



3) SPDK CSAL – Zoned Cloud Native Storage

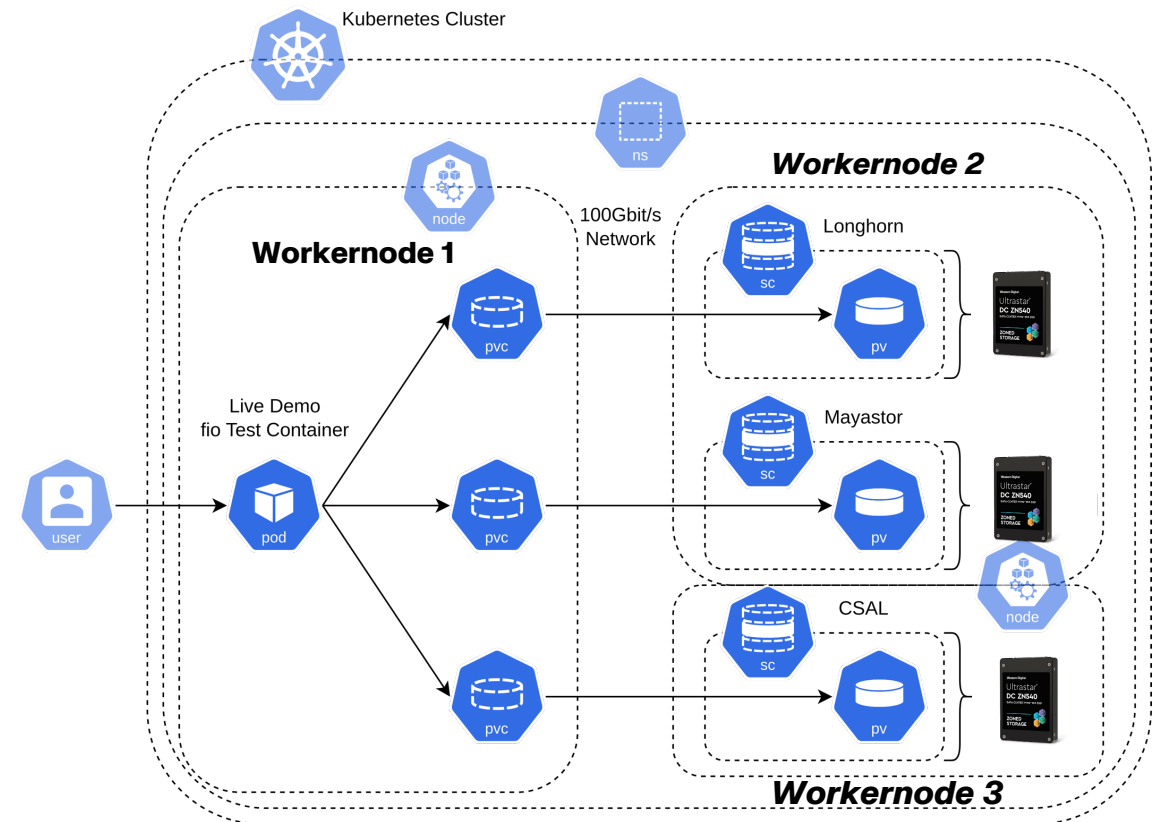
- The Cloud Storage Acceleration Layer (CSAL) which has WIP zoned storage support can be deployed as a CAS though the spdk-csi driver or Mayastor
 - Implements a caching and translation layer that transforms zoned storage to conventional storage
- CSAL uses a conventional (high-performance) block device for metadata and writes sequentially to the ZNS SSDs, thus hiding ZNS' sequential write constraint
- Exposed as a conventional block device over a NVMe-oF target



Showcase

Showcase Setup

- 4 Nodes (one control and three worker nodes)
 - Workernode 1
 - User workload - Container /w fio
 - Workernode 2
 - Longhorn - Exposes a generic file-system/block device backed by an Ultrastar® DC ZN540 SSD
 - Mayastor - Exposes an Ultrastar® DC ZN540 SSD natively to a container
 - Workernode 3
 - SPDK CSAL - Exposes a generic block device over NVMe-oF using an Ultrastar® DC ZN540 as its storage backend





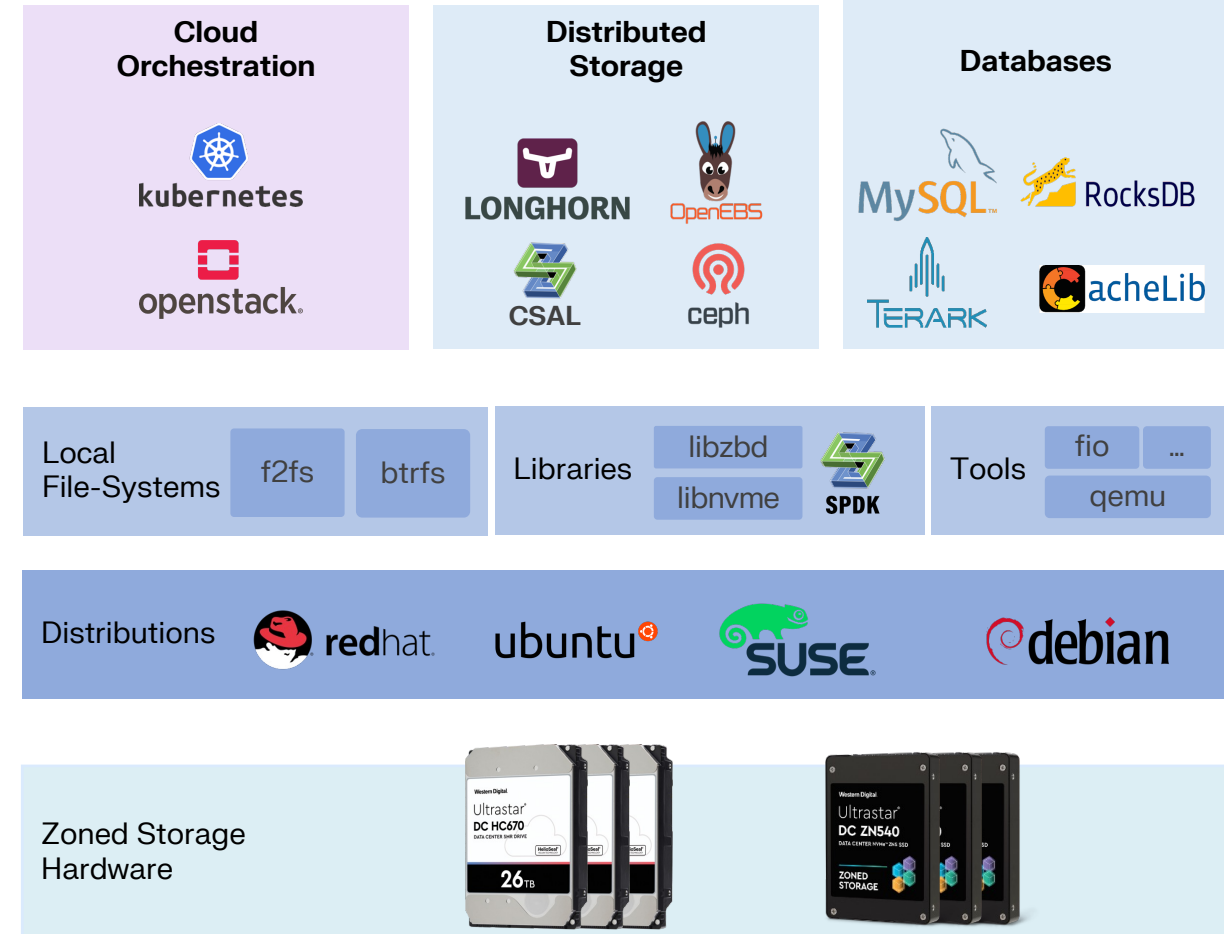
```
2 apiVersion: v1
3 kind: Pod
4 metadata:
5   labels:
6     app: percona-server
7     name: percona-server
8 spec:
9   initContainers:
10  - name: zenfs-setup
11    image: docker.io/library/percona:ps-8.0.33-25
12    command: ['/bin/bash', '-c', 'rm -rf /tmp-zenfs/zenfs-aux-test && zenfs mkfs --zbd=nvme0n1 --aux_path=/tmp-zenfs/zenfs-aux-test --force']
13    volumeMounts:
14    - name: mysql-config-directory
15      mountPath: /etc/my.cnf.d
16    - name: tmp-directory
17      mountPath: /tmp-zenfs
18    volumeDevices:
19    - name: zns-block-device
20      devicePath: /dev/nvme0n1
21  containers:
22  - name: percona-server
23    image: docker.io/library/percona:ps-8.0.33-25
24    ports:
25    - containerPort: 3306
26      name: mysql
27    volumeMounts:
28    - name: mysql-config-directory
29      mountPath: /etc/my.cnf.d
30    - name: tmp-directory
31      mountPath: /tmp-zenfs
32    volumeDevices:
33    - name: zns-block-device
34      devicePath: /dev/nvme0n1
35    env:
36    - name: MYSQL_ROOT_PASSWORD
37      value: "test"
38    - name: INIT_ROCKSDB
39      value: "yes"
40  volumes:
41  - name: mysql-config-directory
42    hostPath:
43      path: /home/debian/k8s-percona-mysql-config
44      type: Directory
45  - name: mysql-directory
46    hostPath:
47      path: /home/debian/k8s-percona-mysql
48      type: Directory
49  - name: tmp-directory
50    hostPath:
51      path: /home/debian/k8s-percona-tmp
52      type: Directory
53  - name: zns-block-device
54    persistentVolumeClaim:
55      claimName: block-device
```

Zoned Block Device

→ Native Zoned Storage applications integrate with Zoned Cloud Native Storage solutions

Summary

- The Linux® zoned storage ecosystem continues to mature
 - BTRFS, Longhorn, Ceph, OpenEBS - Mayastor, CSAL, end-to-end zoned applications
- Workloads orchestrated by Kubernetes are now able to seamlessly benefit from Zoned Storage devices as cost-effective and high-performance storage
- **No longer software stack changes required to deploy and use zoned storage at scale**



Thank You!

Western Digital and the Western Digital logo are registered trademarks or trademarks of Western Digital Corporation or its affiliates in the US and/or other countries. Linux® is the registered trademark of Linus Torvalds in the U.S. and other countries. The NVMe and NVMe-oF word marks are trademarks of NVM Express, Inc. All other marks are the property of their respective owners.