

# Zoned Storage MC

Hans Holmberg, Johannes Thumshirn



LINUX  
PLUMBERS  
CONFERENCE Vienna, Austria / Sept. 18-20, 2024

# Overview

- First half
  - Retrospective and current upstream state
  - Quick updates
    - Damien Le Moal - IO Stack & zonefs
    - Johannes Thumshirn - BTRFS
    - Dennis Maisenbacher - Cloud & virtualization
  - In-depth presentations of ongoing efforts
- Second half
  - BOF / discussions
  - Feel free to suggest topics
- Remote participants
  - Raise hand / turn on video for comments/questions

15:00	<b>Zoned Storage MC Intro</b> <i>Damien Le Moal et al.</i> <i>"Room 1.31-1.32", Austria Center</i> 15:00 - 15:30
	<b>Zoned storage support for QEMU</b> <i>Jia Li</i> <i>"Room 1.31-1.32", Austria Center</i> 15:30 - 15:50
16:00	<b>Zoned XFS Realtime Subvolumes</b> <i>Hans Holmberg</i> <i>"Room 1.31-1.32", Austria Center</i> 15:50 - 16:10
	<b>SSDFs: ZNS/FDP ready LFS file system saving your space and decreasing TCO cost</b> <i>Viacheslav Dubeyko</i> <i>"Room 1.31-1.32", Austria Center</i> 16:10 - 16:30
	<b>Break</b> <i>"Room 1.31-1.32", Austria Center</i> 16:30 - 17:00
17:00	<b>Flexible scheme of space management in ZNS SSD and/or SMR HDD storage pool for massive set of Virtual Machines (VMs)</b> <i>Viacheslav Dubeyko et al.</i>
	<b>Zoned Storage BOFs</b> <i>Hans Holmberg et al.</i>
18:00	<i>"Room 1.31-1.32", Austria Center</i> 17:20 - 18:30



# State of Zoned Storage



# The Evolution of the Zoned Storage Ecosystem

Research



Figure 6: SMR drive with the alternative rotation associated to each physical head assembly is visible parked at the inner diameter.

[2]

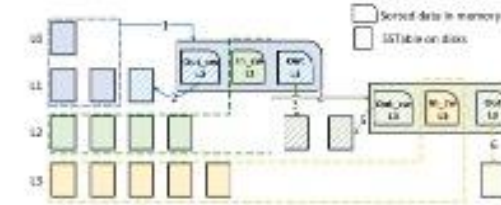


Figure 8: Process of gear compaction. The active compaction of  $L_0$  and  $L_1$  driver passive compactions in higher level. The resultant data of each compaction is divided into three parts according to its key range, including out of  $L_0$ 's compaction window ( $Out_{L_0}$ ), in  $L_1$ 's compaction window ( $In_{L_0}$ ), and out of  $L_1$ 's key range ( $Out_{L_1}$ ).

[3]

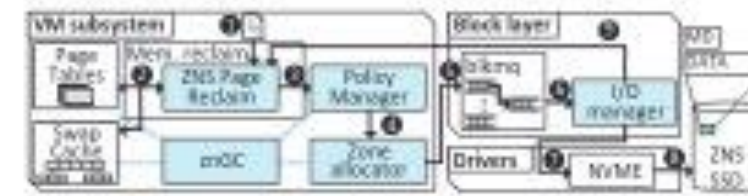
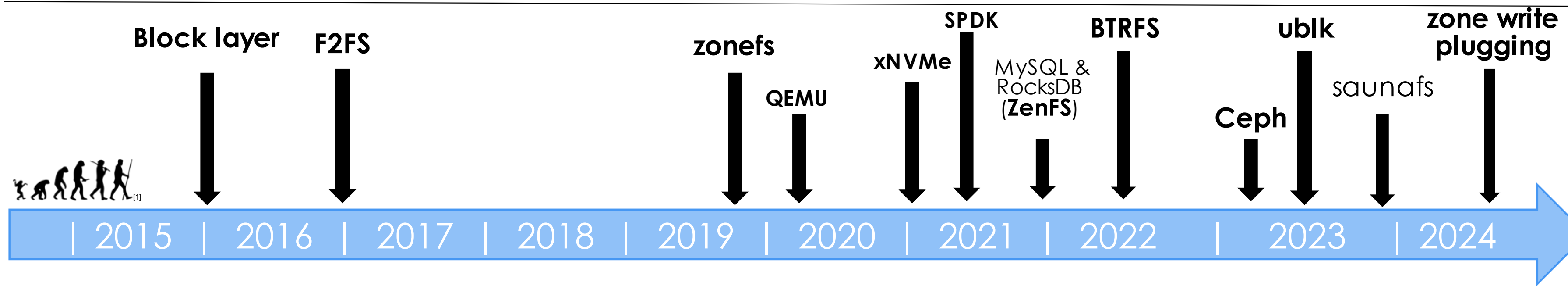


Figure 7: ZNSwap overview. Shaded shapes are internal ZN-Swap components.

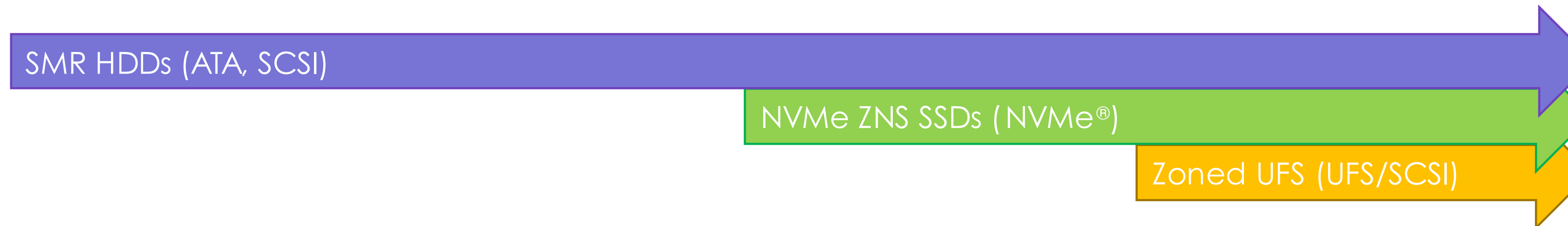
[4]

ZWAL: Rethinking Write-ahead Logs for ZNS SSDs with Zone Appends [5]

Linux (em) eco system



Standards



[1] [https://upload.wikimedia.org/wikipedia/commons/c/c2/Human\\_evolution\\_scheme.svg](https://upload.wikimedia.org/wikipedia/commons/c/c2/Human_evolution_scheme.svg)  
 [2] <https://www.usenix.org/system/files/conference/fast15/fast15-paper-aghayev.pdf>  
 [3] <https://www.usenix.org/system/files/fast19-yao.pdf>  
 [4] <https://www.usenix.org/system/files/atc22-bergman.pdf>  
 [5] <https://dl.acm.org/doi/abs/10.1145/3642963.3652203>

# Current state of the Zoned Storage Stack

## Library/Tools support

- Libzbd, libnvme, xNVME, SPDK, fio, qemu, blkzone, blktests,

## End-to-end Application Enablement

- MySQL, RocksDB, TerarkDB, ..

## Local File-system support

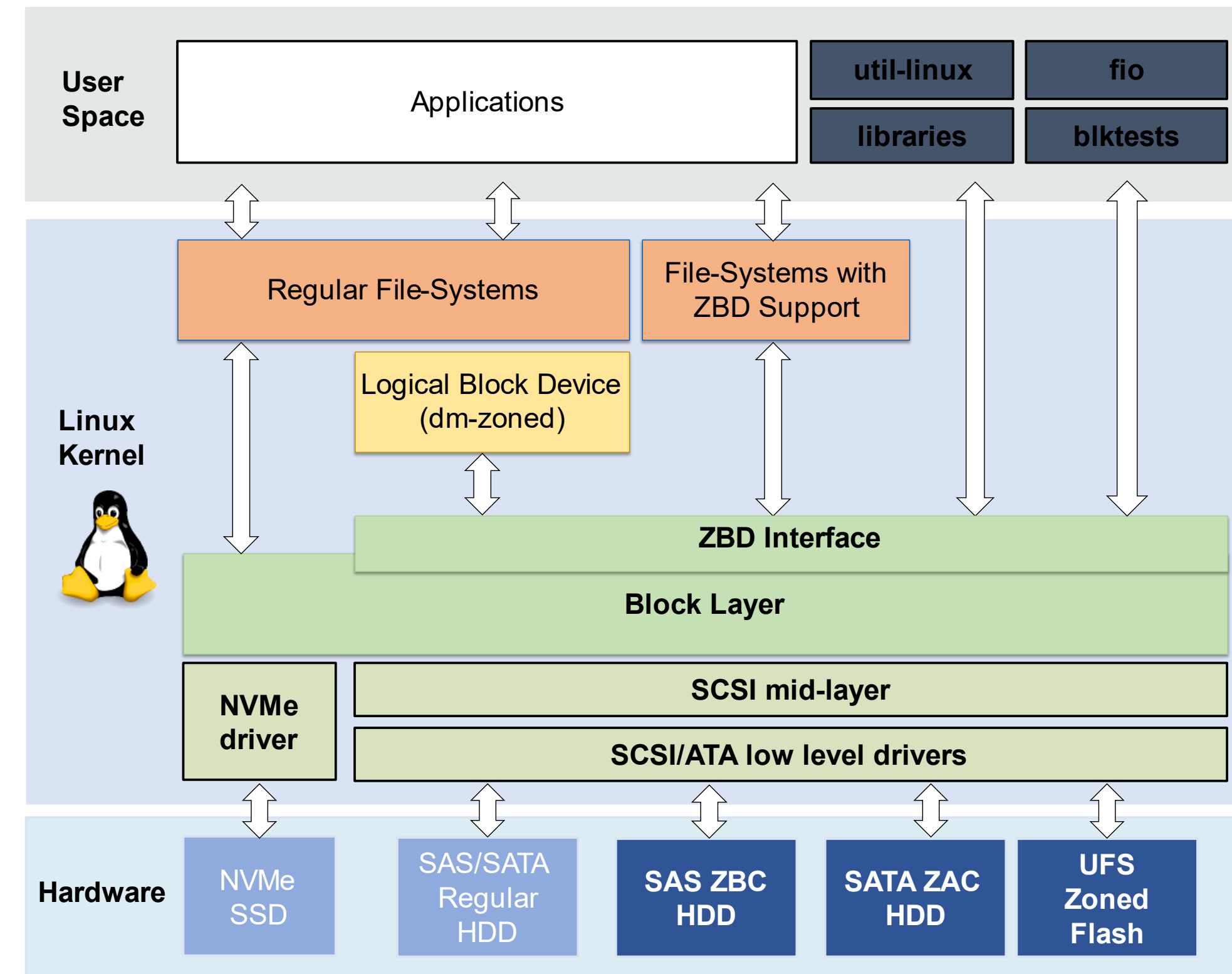
- f2fs (client – UFS)
- btrfs (enterprise – ZNS/SMR)
- xfs (enterprise, under development)

## Storage Systems

- Ceph, OpenEBS, Mayastor, Saunafs, SPDK's CSAL, ..

## Distributions

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



Mature, robust, and used in production by some of the biggest consumers of storage



# I/O Stack (Block Layer, DM, SCSI/NVMe)

- New per-zone sequential write ordering control with zone write plugging (v6.10)
  - Same "at most one write per zone in flight at any time" principle as with zone write locking (mq-deadline) but control is applied to BIOs instead of requests
- Advantages:
  - Remove dependency on mq-deadline (zone write locking): any block I/O scheduler can be used, included "none"
  - Generic zone append emulation implementation
    - Simplifies DM and SCSI sd code
  - Significant performance improvements in some cases, including for read operations
- Drawbacks
  - Not many so far (identified performance degradation with some SMR drives on zone boundary crossing)



# zonefs

- Some small changes
  - Error recovery bug fix (v6.8)
  - Conversion to new mount API (v6.9)
  - Large folio support added (v6.10)
- Started exploring zone append user interface through io\_uring
  - Changes to io\_uring and iomap needed
  - New write append operation will return the written offset to the user
    - Applicable to regular O\_APPEND writes







# BTRFS

- Initial support upstream v5.12
  - Garbage collection added in v5.13
  - NVMe ZNS support added in v5.16
  - (experimental) RAID0/1/10 support added in v6.7
- 
- Stabilization is still on-going
    - Early ENOSPC errors
      - Metadata overcommit
      - Space accounting
    - RAID support





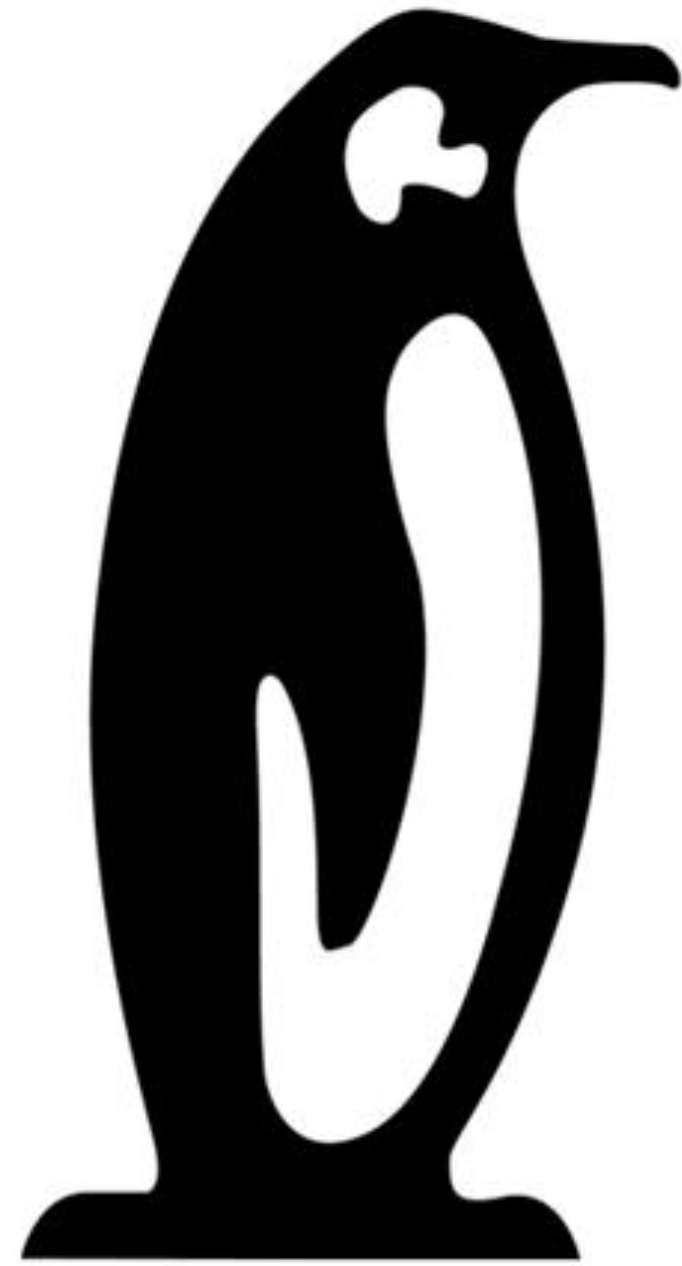
# Cloud storage & Virtualization

	 <b>SPDK</b>	 <b>ceph</b>	<b>Mayastor</b> 	 <b>QEMU</b>
<b>Upstream</b>	<ul style="list-style-type: none"> <li>Initial ZNS support since v20.10 with ongoing effort<sup>[1]</sup></li> </ul>	<ul style="list-style-type: none"> <li>Crimson – seastore ships with initial ZBD support (Reef stable release)<sup>[3][4]</sup></li> </ul>		<ul style="list-style-type: none"> <li>Initial ZNS emulation<sup>[6]</sup></li> <li>Attaching ZNS devices via vfiio PCI passthrough or virtio-blk<sup>[7]</sup></li> <li>Attaching SMR disks via virtio-scsi or vhost-scsi<sup>[8]</sup></li> </ul>
<b>In progress</b>	<ul style="list-style-type: none"> <li>CSAL - ZNS patches to be upstreamed<sup>[2]</sup></li> </ul>	<ul style="list-style-type: none"> <li>Fix for active zone resource exhaustion</li> </ul>	<ul style="list-style-type: none"> <li>Open PR for initial ZNS support under discussion<sup>[5]</sup></li> <li>CSAL Integration?</li> </ul>	<ul style="list-style-type: none"> <li>qcow2 support for zoned storage emulated devices<sup>[9]</sup></li> </ul>



[1] [https://spdk.io/release/2020/10/30/20.10\\_release/](https://spdk.io/release/2020/10/30/20.10_release/)  
 [2] <https://dl.acm.org/doi/abs/10.1145/3627703.3629566>  
 [3] <https://docs.ceph.com/en/reef/dev/zoned-storage/>  
 [4] <https://docs.ceph.com/en/latest/releases/reef/>  
 [5] <https://github.com/openebs/mayastor/pull/1298>

[6] <https://lists.nongnu.org/archive/html/qemu-block/2020-06/msg00720.html>  
 [7] <https://www.qemu.org/docs/master/devel/zoned-storage.html>  
 [8] <https://zonedstorage.io/docs/tools/qemu#qemu-virtio-scsi>  
 [9] <https://lore.kernel.org/all/20240122184830.40094-1-faithilikerun@gmail.com/>



# Linux Plumbers Conference

Vienna, Austria | September 18-20, 2024

