# Exploring Storage Stack on Modern NVMe Hardware

**(lots of graph and numbers coming up - apologies)**

**Animesh Trivedi**
Assistant Professor (tenured)
https://animeshtrivedi.github.io/
April 22nd, 2024

Invited talk at the 4th Workshop on Challenges and Opportunities of Efficient and Performant Storage Systems (CHEOPS'24)
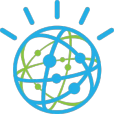
# Data is Essential to our Society



**Warehouse**

**Financial**

**AI/ML**

**Health**

**Mobility**

**Science**

Data is the new oil!

# A Minute on the Internet



# 1 Yottabytes

*(per year by 2030)*



1 byte = 1 grain

3

# Non-Volatile Memory (NVM) Storage to the Rescue…

# Rise of Domain-Specific Computing



**Stalled CPU-centric computing scaling**

*Rise of accelerator-centric computing*
+ Specialized hardware
+ Energy/Perf. gains over the CPU

5

# **Position:** Workload-Specialized Storage will Emerge

**Specialization**
- The nature of I/O operation
- Completion management
- Policies, scheduling, priority
- Performance

Classical Workloads

I/O APIs

File system

Block I/O

**OS**

Storage Devices (flash and NVM storage)

**Position:** Workload-Specialized Storage will Emerge

Classical Workloads

I/O APIs

File system

Block I/O

**OS**

Storage Devices (flash and NVM storage)

# Position: Workload-Specialized Storage will Emerge

Classical Workloads

| I/O APIs |
| File system |
| Block I/O | OS

KV Stores

Graphs

Trees

Tensors

Storage Devices (flash and NVM storage)

# Position: Workload-Specialized Storage will Emerge

**[Part - 1/2]** : **Study:** I/O Performance and Scheduling

**[Part - 2/2]** : Zone Namespace Devices (ZNS)

Classical Workloads

I/O APIs

File system

Block I/O

**OS**

Western Digital
Ultrastar
**DC ZN540**
DATA CENTER NVMe™ ZNS SSD

ZONED STORAGE

**NVMe Zone Namespace Interface**

Storage Devices (flash and NVM storage)

## [Part - 1/2] : **Study:** I/O Performance and Scheduling Overheads

Diego Didona, Jonas Pfefferle, Nikolas Ioannou, Bernard Metzler, and Animesh Trivedi. 2022. **Understanding modern storage APIs: a systematic study of libaio, SPDK, and io_uring**. In Proceedings of the 15th ACM International Conference on Systems and Storage (**SYSTOR '22**). Association for Computing Machinery, New York, NY, USA, 120–127. https://doi.org/10.1145/3534056.3534945

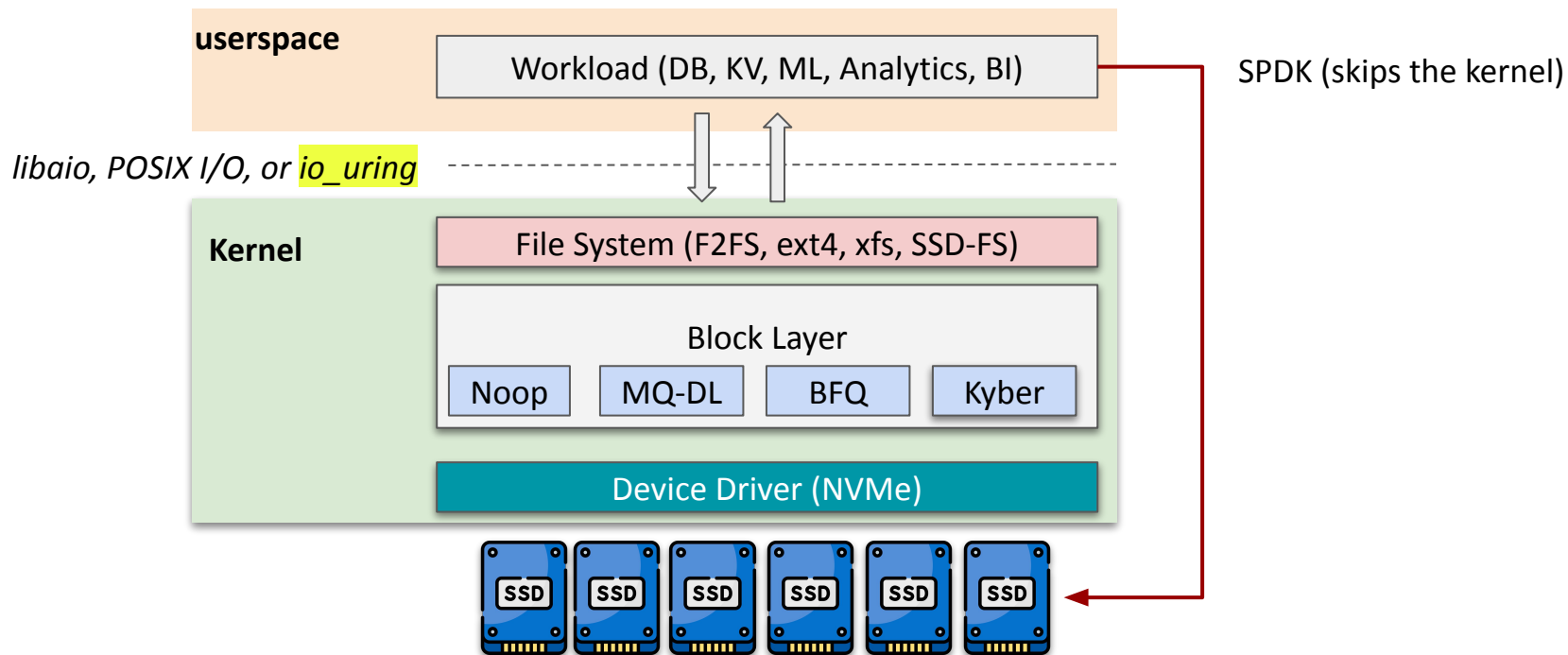Zebin Ren and Animesh Trivedi. 2023. **Performance Characterization of Modern Storage Stacks: POSIX I/O, libaio, SPDK, and io_uring.** In Proceedings of the 3rd Workshop on Challenges and Opportunities of Efficient and Performant Storage Systems (**CHEOPS '23**). Association for Computing Machinery, New York, NY, USA, 35–45. https://doi.org/10.1145/3578353.3589545

Zebin Ren, Krijn Doekemeijer, Nick Tehrany, Animesh Trivedi. 2024. BFQ, **Multiqueue-Deadline, or Kyber? Performance Characterization of Linux Storage Schedulers in the NVMe Era**, to appear in the 2024 ACM/SPEC International Conference on Performance Engineering (**ICPE '24**), London, UK.

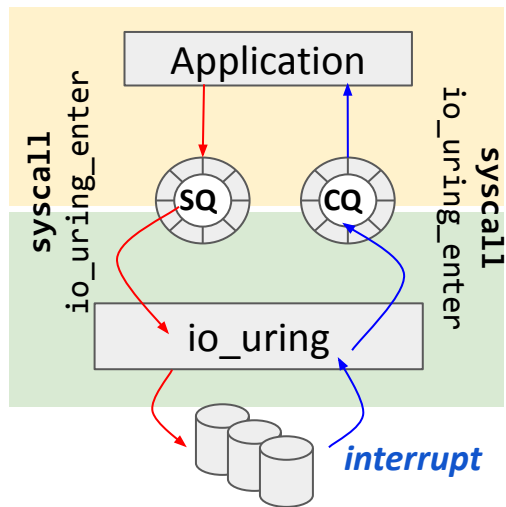## [Part - 2/2] : Zone Namespace Devices (ZNS) Performance Characterization

Krijn Doekemeijer, Nick Tehrany, Bala Chandrasekaran, Matias Bjørling and Animesh Trivedi. **Performance Characterization of NVMe Flash Devices with Zoned Namespaces (ZNS).** 2023 IEEE International Conference on Cluster Computing (**CLUSTER'23**), Santa Fe, NM, USA, 2023, pp. 118-131, doi: https://doi.org/10.1109/CLUSTER52292.2023.00018 .

# Workload-NVMe Interaction



userspace

Workload (DB, KV, ML, Analytics, BI)

SPDK (skips the kernel)

*libaio, POSIX I/O, or io_uring*

**Kernel**

File System (F2FS, ext4, xfs, SSD-FS)

Block Layer

Noop    MQ-DL    BFQ    Kyber

Device Driver (NVMe)

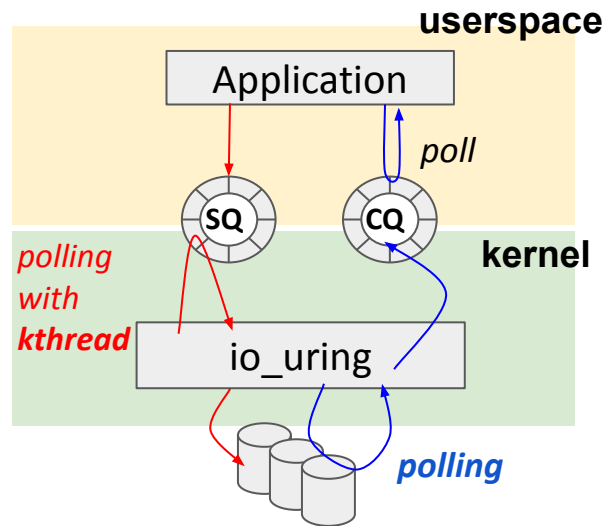SSD SSD SSD SSD SSD SSD

# Three Modes of io_uring API



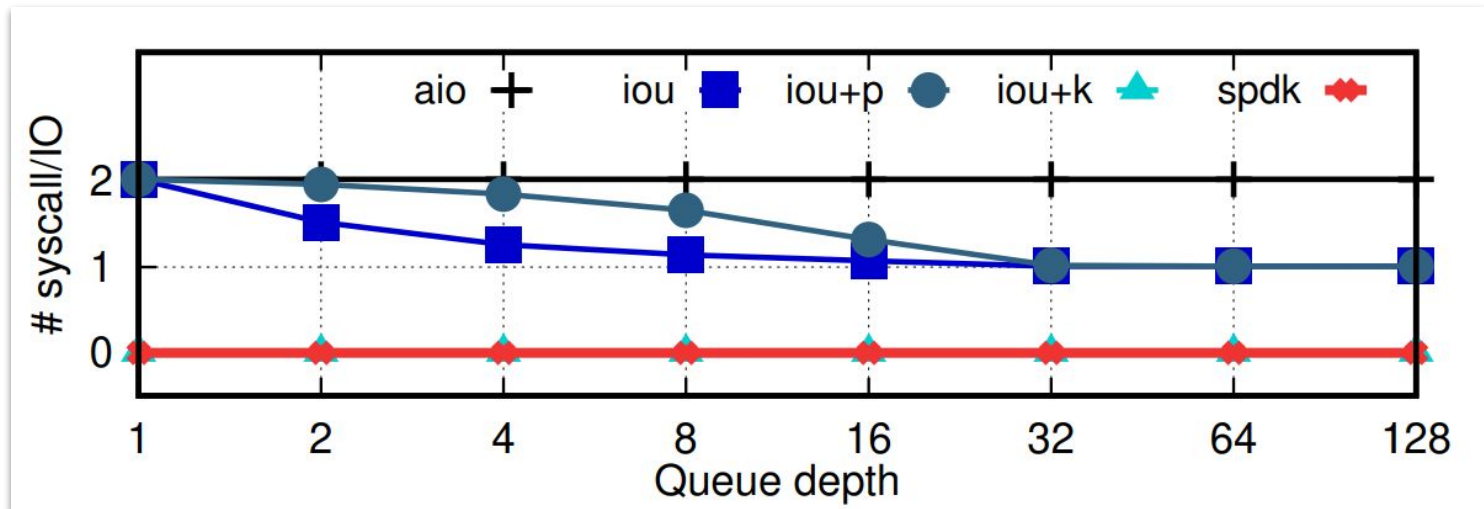**(a)** default with syscalls    **(b)** [**iou+p**] with completion polling    **(c)** [**iou+k**] with submission polling
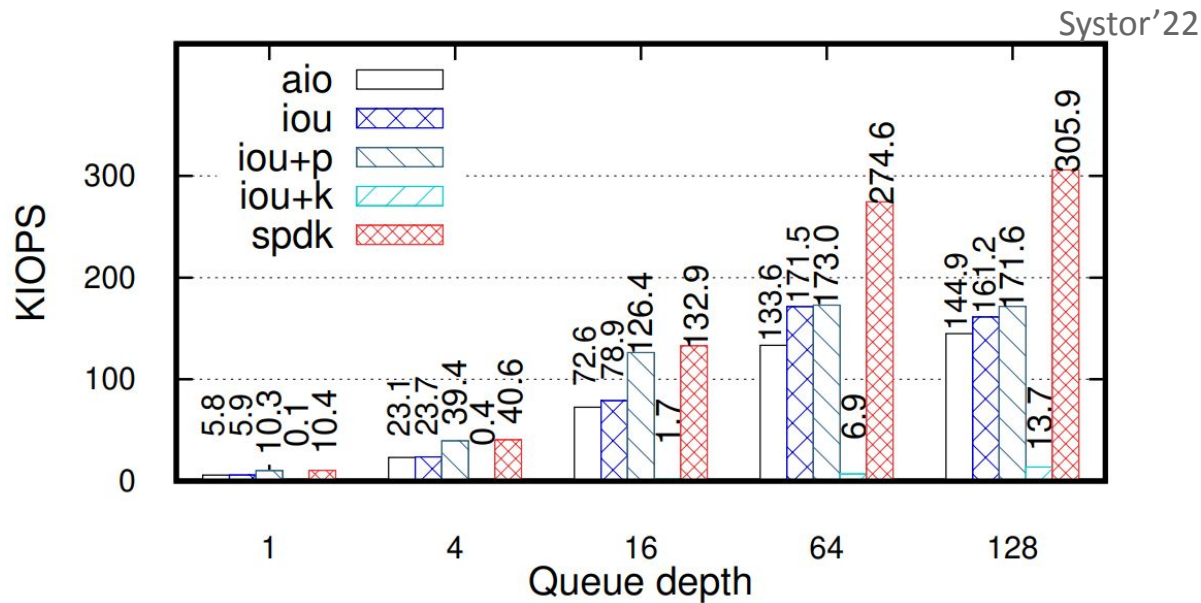
# io_uring: System call study



**Just like SPDK, io_uring can support a <u>pure polling based, ZERO system calls</u> I/O path!**
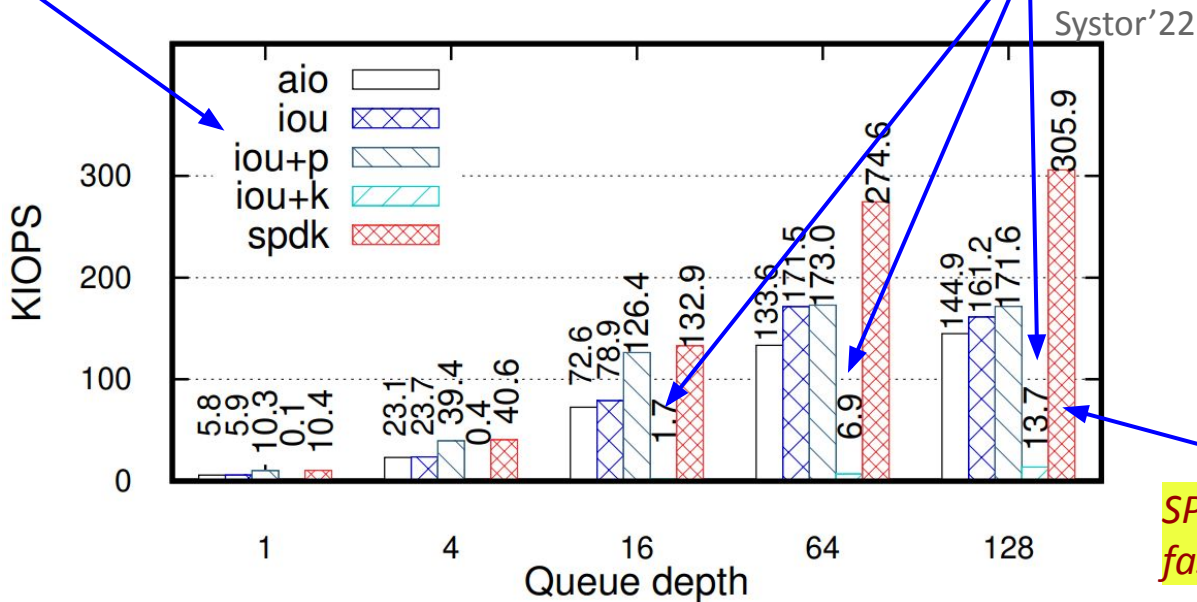
# Results: Efficiency (single CPU core)



Systor'22

# Results: Efficiency (single CPU core)

io_uring sits between libaio and SPDK

Performance collapses with the kernel polling

Systor'22



SPDK is the fastest API (still)!

16

# Analysis: CPU Profile

Systor'22

CHEOPS'23



50:50 CPU sharing with polling - *Careful*!

**SPDK** stack is still 5x more CPU efficient

# Results: Efficiency with <u>TWO</u> CPU cores

Systor'22



[ aio < iou < iou with polling < iou with kernel poll < SPDK ]

The "normal performance" order can be resumed (**but,** at the cost of 2x CPU cores)!

# **Pure** Performance Scaling



CHEOPS'23

# **Pure Performance Scaling**

CHEOPS'23



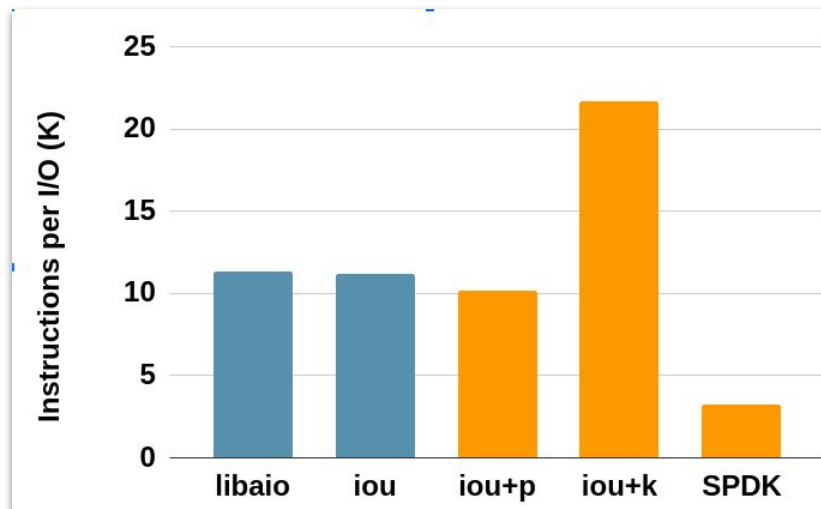- **There is a large gap (10x) in the CPU efficiency between SPDK and io_uring stacks**
- **In the Linux kernel, the block layer is the primary consumer of the CPU cycles**

# So, What's Wrong with SPDK?

Takes a pure performance-based approach

Highly CPU inefficient (only poll, 100% CPU utilization)
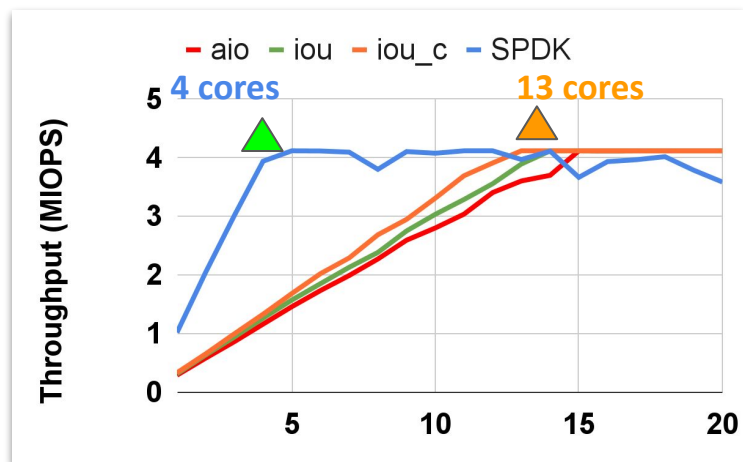
Fragile performance when polling on > #CPU cores

Does not have a file system

Does not have multi-tenancy (only single process)

No support for any other kind of devices except NVMe

No provision for the kernel supported services:

- Caching, buffering, security …
- **Importantly: Sharing and I/O Scheduling**

# What are the Scheduling Challenges

ICPE'24



(a) IOPS performance of schedulers;

High performance scaling with the none I/O scheduler
1.3 - 2.7x slowdown with other schedulers

Zebin Ren, Krijn Doekemeijer, Nick Tehrany, Animesh Trivedi. 2024. BFQ, **Multiqueue-Deadline, or Kyber? Performance Characterization of Linux Storage Schedulers in the NVMe Era**, to appear in the 2024 ACM/SPEC International Conference on Performance Engineering (ICPE '24), London, UK.

# What are the Scheduling Challenges

ICPE'24

*P95 latencies degradation*



(a) IOPS performance of schedulers;    Latency (P95) with background **(b) reads and (c) writes traffic**

- **No scheduling (NOOP) helps with pure performance scaling**
- **No scheduling (NOOP) has poor performance isolation with _interfering tasks_**

# The Tipping Point - the CPU bottleneck

(a) 1 L-app    (b) 16 L-apps    (c) 32 L-apps    (d) 64 L-apps    (e) 256 L-apps

**While the CPU is not the bottleneck, all I/O schedulers typically deliver the same performance**



24

# Can We Look at the SSD to Get Help for QoS Support?

# The Interference Control (or delivering Quality-of-Service)

*I/O Scheduling interference and overheads*

| P1 | P2 | P3 | P4 |

I/O Scheduling

Data Placement | Garbage collection

Wear-leveling (over provisioning)

**Flash** **Flash** **Flash** **Flash**

Ch#0    Ch#1    Ch#2    Ch#3

**Inside an SSD**
- Mixing of data (lifetime, workloads)
- I/O Scheduling
- Interference from GC
- Over provisioning
- Parallelism management
- …

## [Part - 1/2] : Study: I/O Performance and Scheduling Overheads

Diego Didona, Jonas Pfefferle, Nikolas Ioannou, Bernard Metzler, and Animesh Trivedi. 2022. **Understanding modern storage APIs: a systematic study of libaio, SPDK, and io_uring**. In Proceedings of the 15th ACM International Conference on Systems and Storage (SYSTOR '22). Association for Computing Machinery, New York, NY, USA, 120-127. https://doi.org/10.1145/3534056.3534945
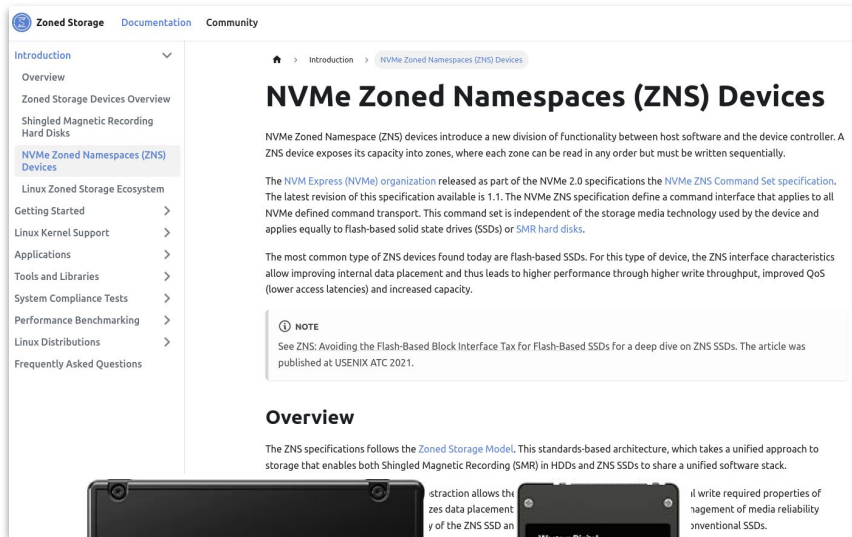
Zebin Ren and Animesh Trivedi. 2023. **Performance Characterization of Modern Storage Stacks: POSIX I/O, libaio, SPDK, and io_uring.** In Proceedings of the 3rd Workshop on Challenges and Opportunities of Efficient and Performant Storage Systems (CHEOPS '23). Association for Computing Machinery, New York, NY, USA, 35-45. https://doi.org/10.1145/3578353.3589545

Zebin Ren, Krijn Doekemeijer, Nick Tehrany, Animesh Trivedi. 2024. BFQ, **Multiqueue-Deadline, or Kyber? Performance Characterization of Linux Storage Schedulers in the NVMe Era**, to appear in the 2024 ACM/SPEC International Conference on Performance Engineering (ICPE '23), London, UK.
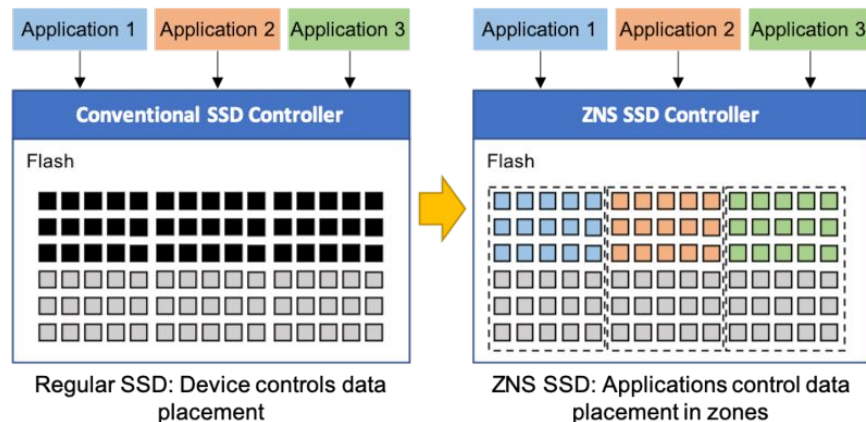
## [Part - 2/2] : Zone Namespace Devices (ZNS) Performance Characterization

Krijn Doekemeijer, Nick Tehrany, Bala Chandrasekaran, Matias Bjørling and Animesh Trivedi. **Performance Characterization of NVMe Flash Devices with Zoned Namespaces (ZNS).** 2023 IEEE International Conference on Cluster Computing (CLUSTER), Santa Fe, NM, USA, 2023, pp. 118-131, doi: https://doi.org/10.1109/CLUSTER52292.2023.00018 .

# ZNS: The New Storage Interface and Capabilities



Regular SSD: Device controls data placement

ZNS SSD: Applications control data placement in zones

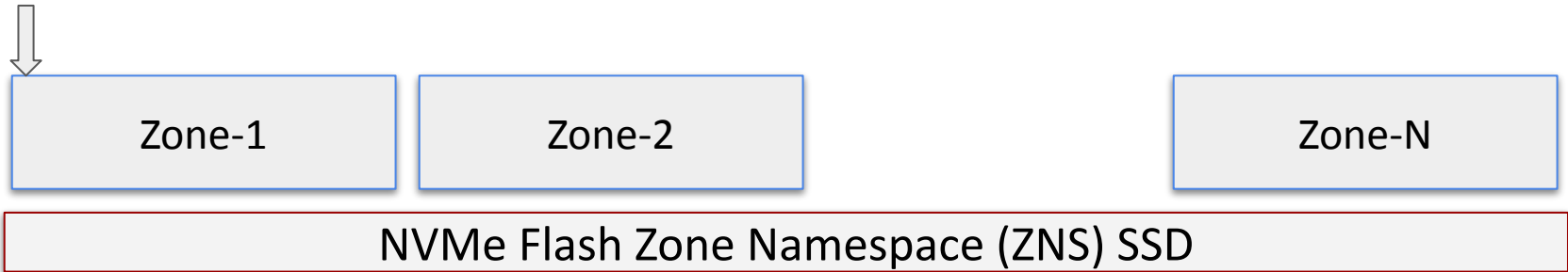https://zonedstorage.io/docs/introduction/zns

Standardized in the NVMe 1.4, July 2021

# Zone Namespace (ZNS) Devices : The Operational Model

A ZNS SSD is divided into Zones

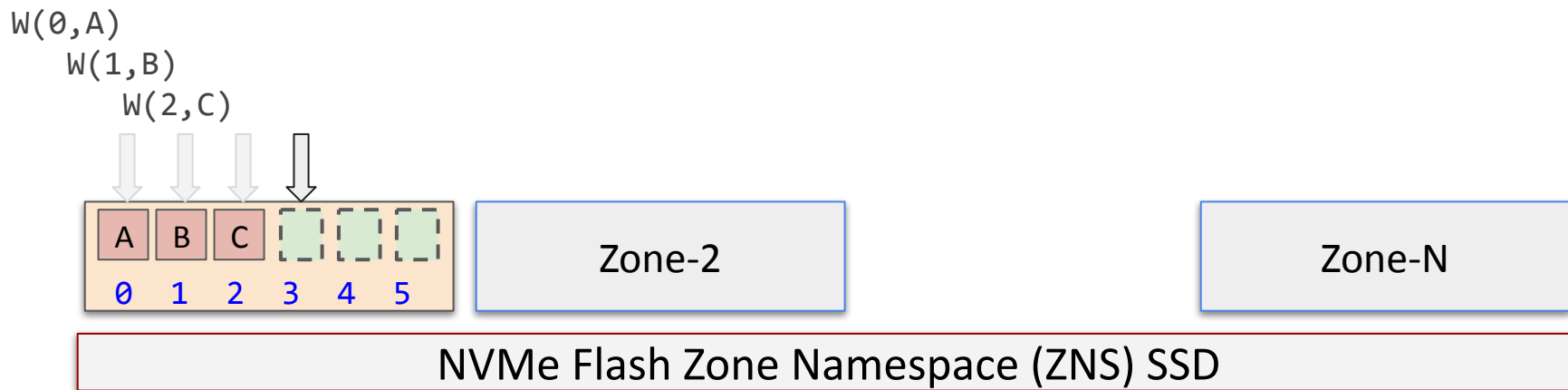Each zone has its size and a **write pointer**

Write pointer

| Zone-1 | Zone-2 | | Zone-N |
|--------|--------|--|--------|

NVMe Flash Zone Namespace (ZNS) SSD

# Zone Namespace (ZNS) Devices : The Operational Model

Each zone must be written sequentially

Limited **intra-zone** parallelism (only 1 write at a time)

```
W(0,A)
   W(1,B)
      W(2,C)
```

| A | B | C |   |   |   |
|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 |

Zone-2

Zone-N

NVMe Flash Zone Namespace (ZNS) SSD

# Zone Namespace (ZNS) Devices : The Operational Model

New I/O Command: **Append**
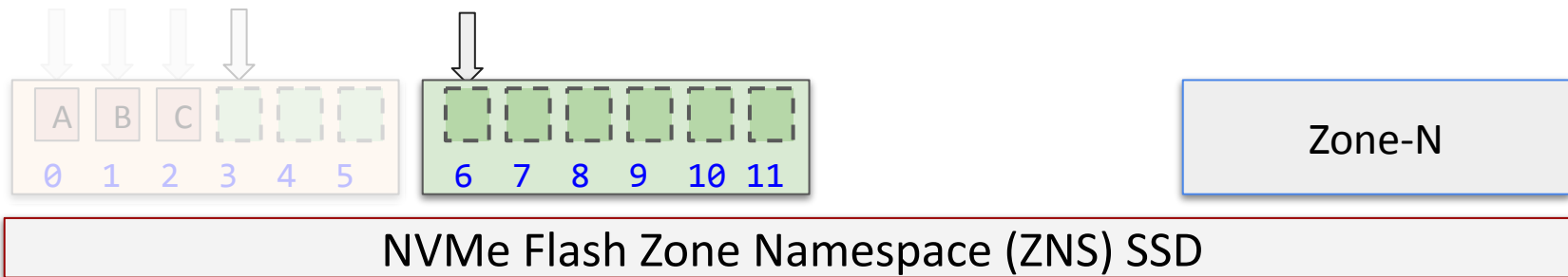
Multiple Append command can be issued to a zone (high **intra-zone** parallelism)

A(Z-2,M)
A(Z-2,N)                    *"Append M, N and O to Zone-2 (anywhere)"*
A(Z-2,O)



| A | B | C | | | |
|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 |

6  7  8  9  10  11

Zone-N

NVMe Flash Zone Namespace (ZNS) SSD

# Zone Namespace (ZNS) Devices : The Operational Model
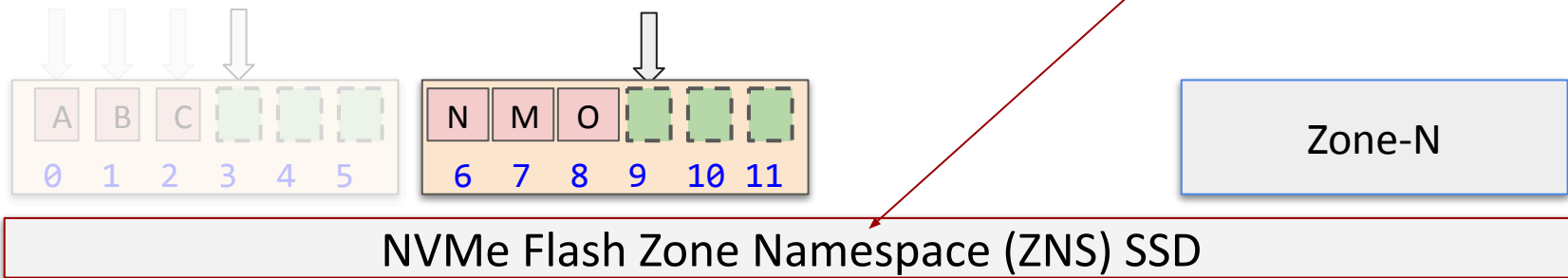
New I/O Command: **Append**

Multiple Append command can be issued to a zone (high **intra-zone** parallelism)

$$A(Z-2,M) \Rightarrow P7$$
$$A(Z-2,N) \Rightarrow P6$$
$$A(Z-2,O) \Rightarrow P8$$

ZNS SSD does **I/O scheduling** and **space allocation**

| A | B | C | | | |
|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 |

| N | M | O | | | |
|---|---|---|---|---|---|
| 6 | 7 | 8 | 9 | 10 | 11 |

Zone-N

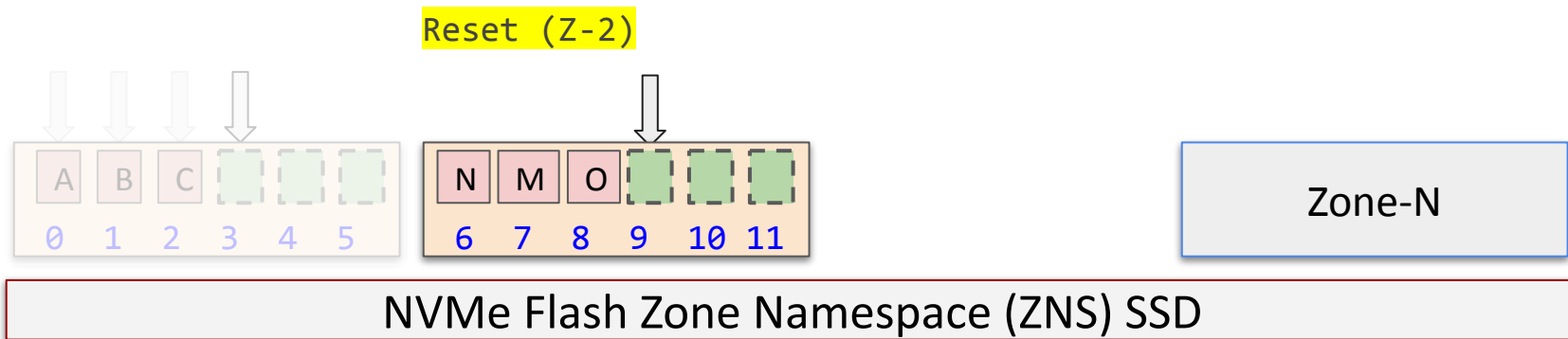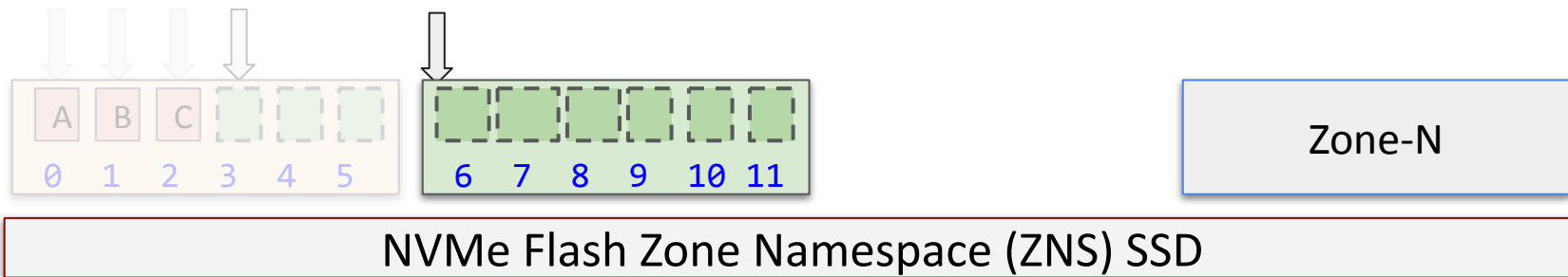NVMe Flash Zone Namespace (ZNS) SSD

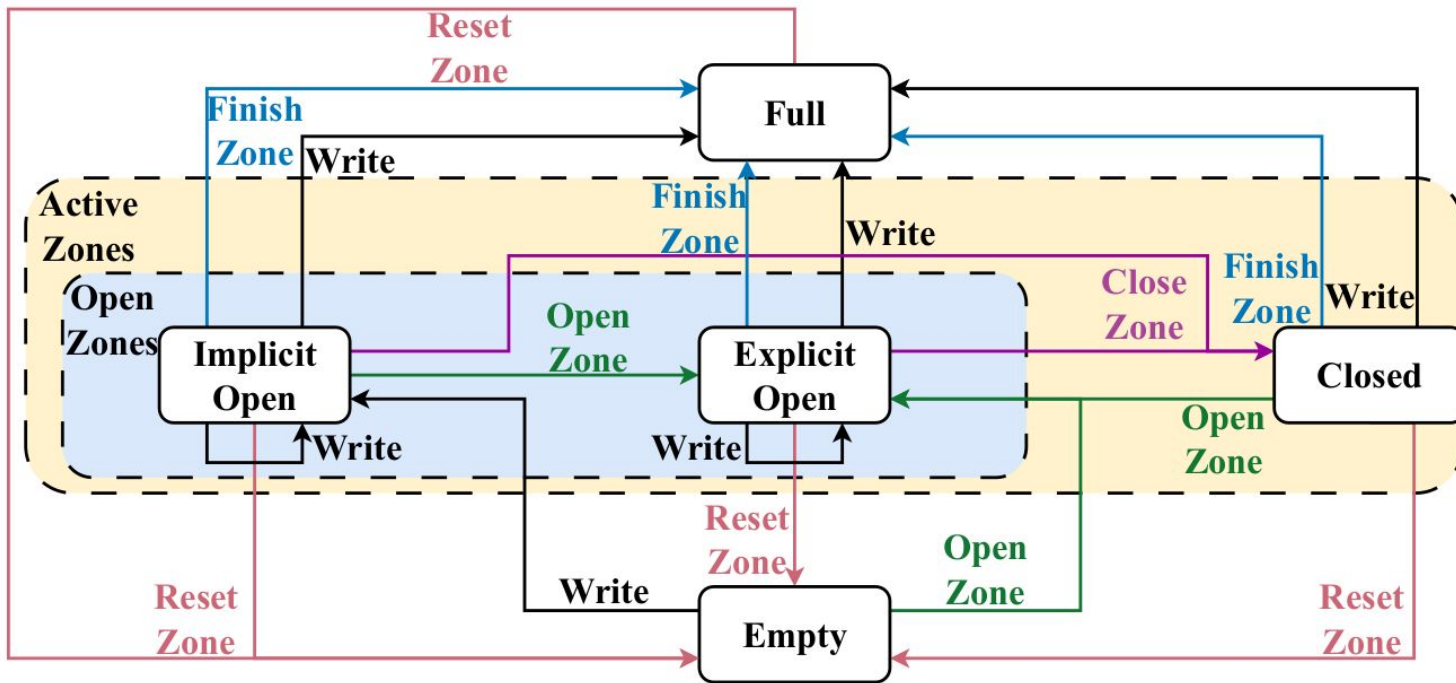# Zone Namespace (ZNS) Devices : The Operational Model

New zone-management commands: **Finish** and **Reset**

**Finish**: makes it read-only (release write resources)

**Reset**: garbage collect the zone

Reset (Z-2)

| A | B | C | | | |
|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 |

| N | M | O | | | |
|---|---|---|---|---|---|
| 6 | 7 | 8 | 9 | 10 | 11 |

Zone-N

NVMe Flash Zone Namespace (ZNS) SSD

# Zone Namespace (ZNS) Devices : The Operational Model

New zone-management commands: **<u>Finish</u>** and **<u>Reset</u>**

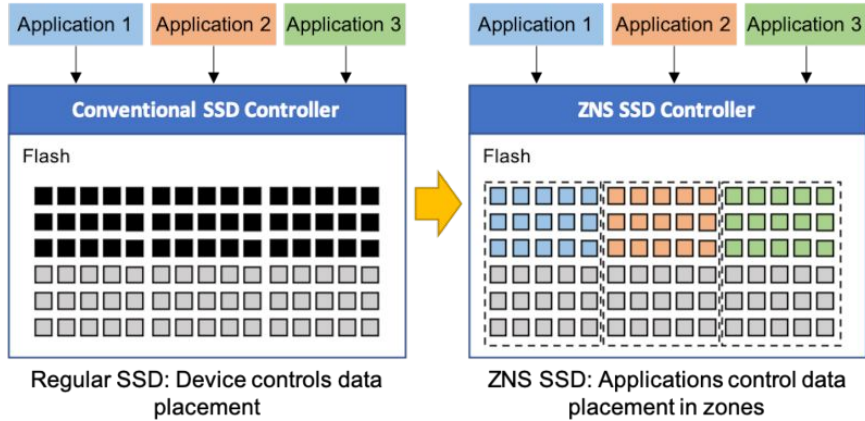    **Finish**: makes it read-only (release write resources)
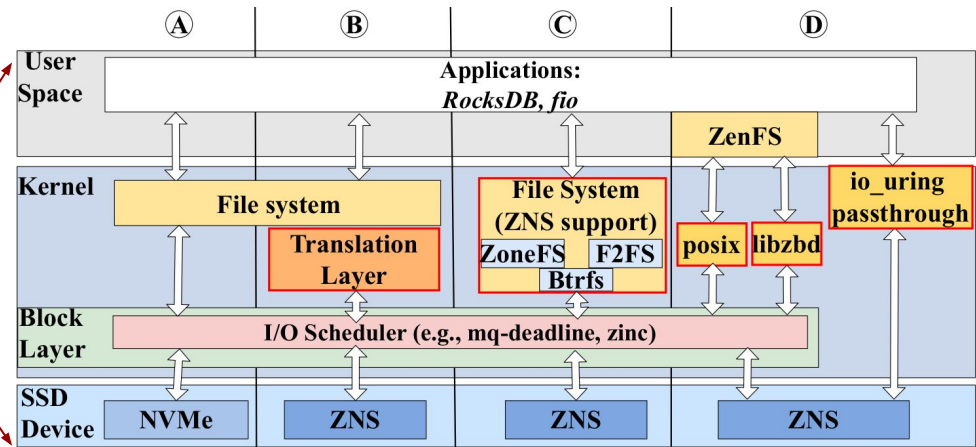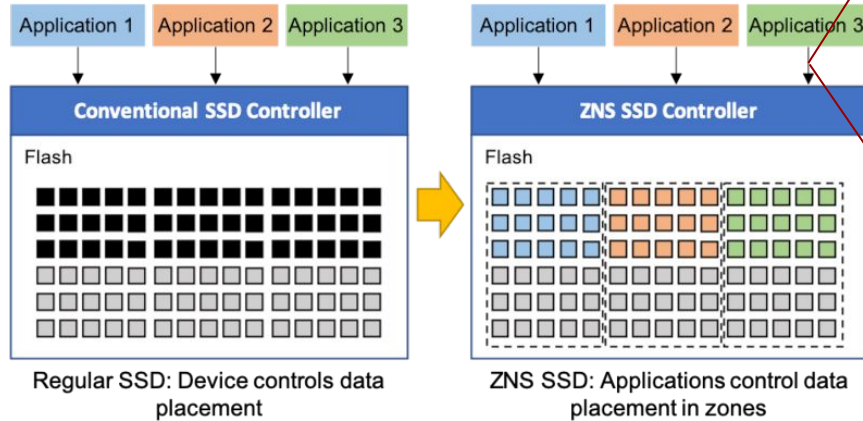
    **Reset**: garbage collect the zone

# Zone Namespace (ZNS) Devices: The State Machine

# State of the ZNS Software



Regular SSD: Device controls data placement

ZNS SSD: Applications control data placement in zones

# State of the ZNS Software



Regular SSD: Device controls data placement

ZNS SSD: Applications control data placement in zones



Understanding NVMe Zoned Namespace (ZNS) Flash SSD Storage Devices, Nick Tehrany, Animesh Trivedi, https://arxiv.org/abs/2206.01547 (2022).

**Idea**: Different zones helps to isolate workloads from each other and better Quality-of-Service (QoS)

**But:** There are multiple ways ZNS devices can be integrated

- Should I use **Append** or **Write**? How do I manage **parallelism**? Intra-zone or Inter-zone?
- What is the cost of **Reset** and **Finish**? And the state machine implementation
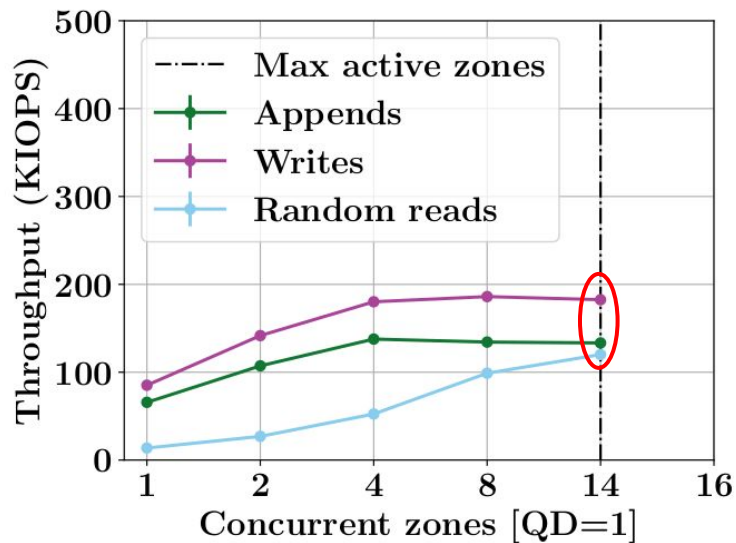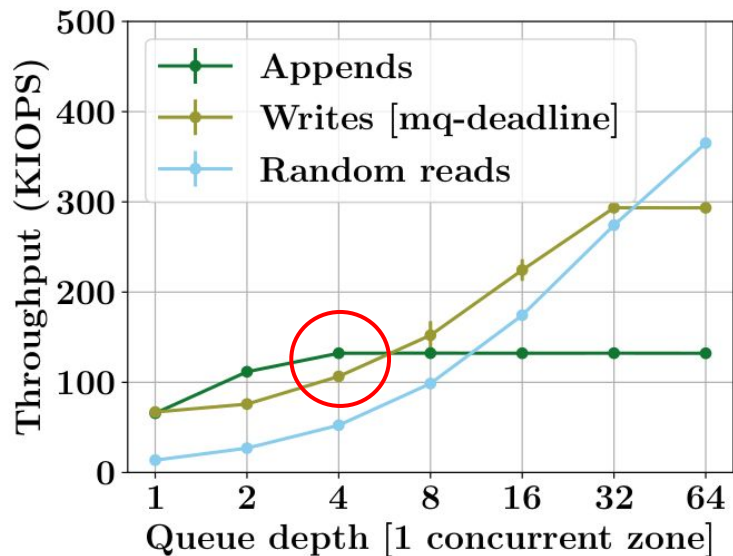- Do ZNS SSDs deliver **isolation**?

37

# Result [1 / 3]: Write vs Append Parallelism Management
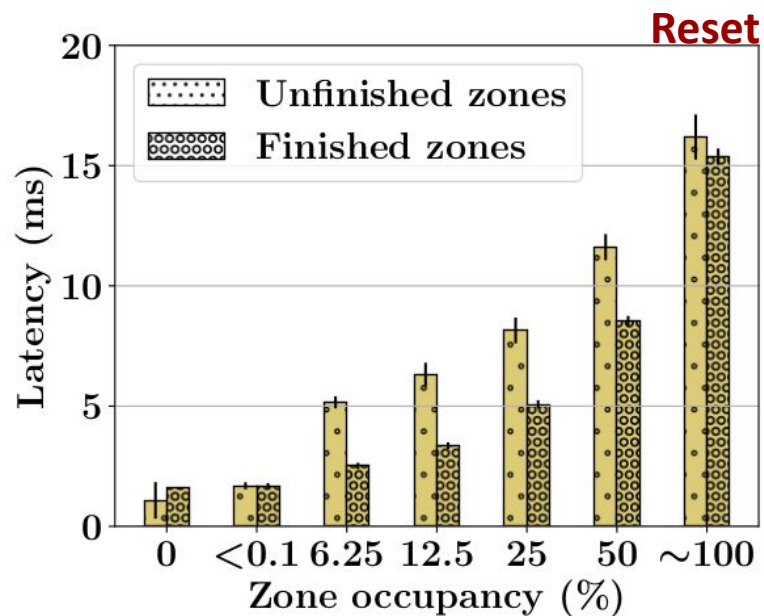


*mq-deadline scheduler merges adjacent writes*

**Single Zone Parallelism (intra-zone)**

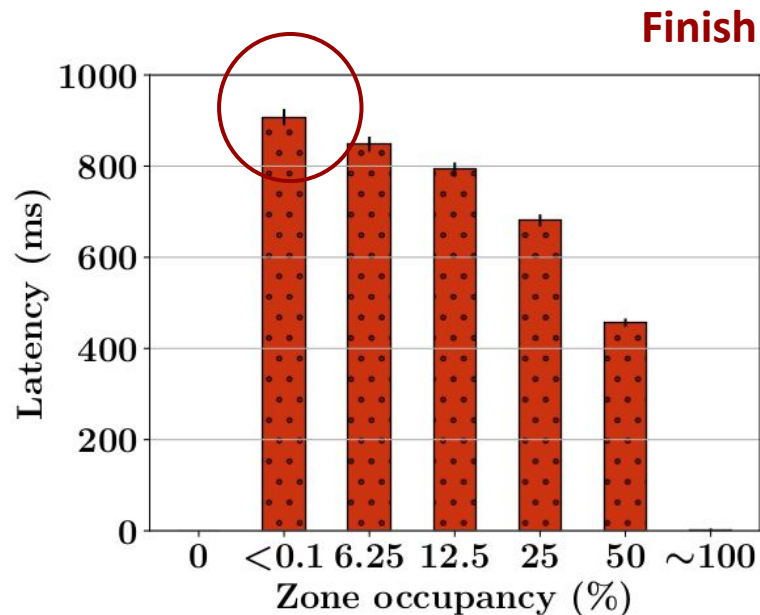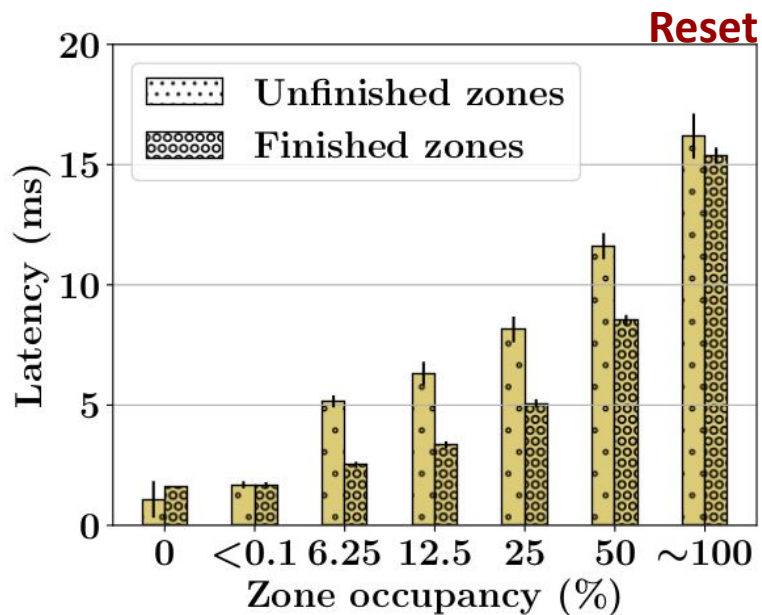# Result [1 / 3]: Write vs Append Parallelism Management



- **Intra-Zone parallelism has higher performance**
- **Writes have better performance scalability than Appends (!)**
- **Append scalability is independent of intra- or inter-zone, but limited in performance**

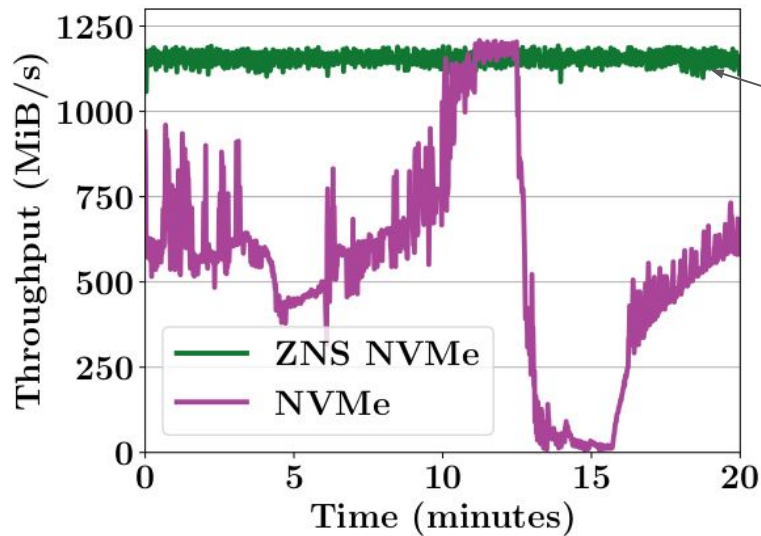# Result [2 / 3]: The Cost of Reset and Finish Operations

**Reset**

# Result [2 / 3]: The Cost of Reset and Finish Operations



- **The zone utilization --- Very important factor**
- **Finish is an extremely expensive operation (100 - 1,000s of milliseconds)**
- **Leverage intra-zone parallelism (*minimize half-written zones*)**
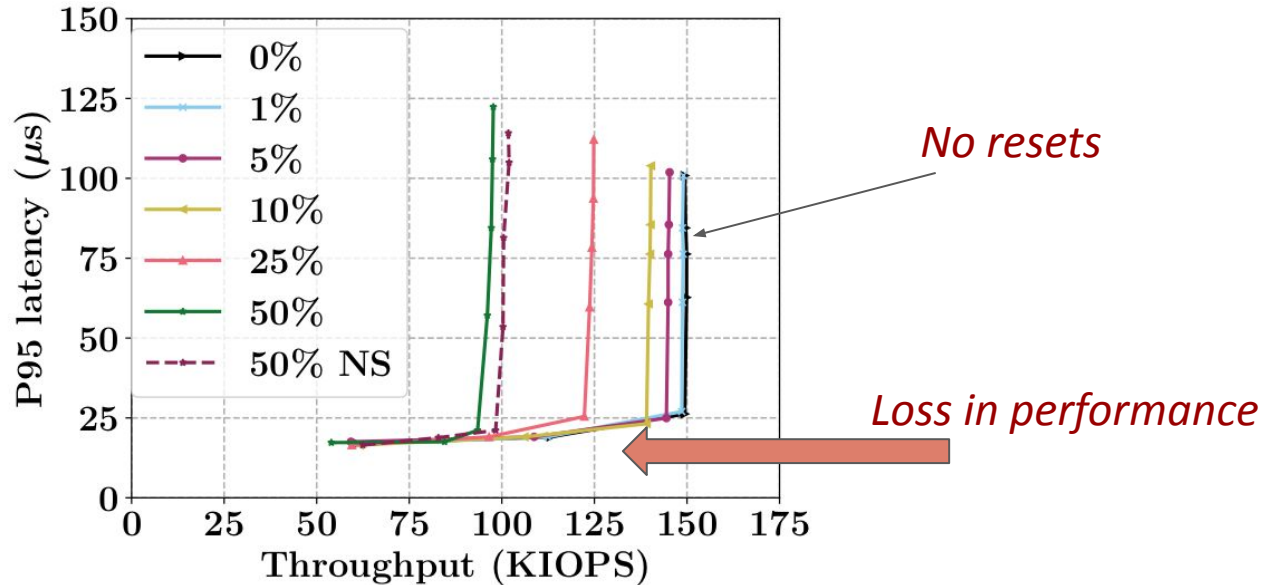
# Result [3 / 3]: <u>Read-Write</u> Isolation on ZNS



Stable performance

- **ZNS provides good read-write isolation when operating on multiple zones**
- **Stable performance (in comparison to NVMe)**

# Do Reset Commands Interfere with I/O Operations?



*No resets*

*Loss in performance*

**Initial results: Yes ... part of an active research now :)**

# [Part - 1/2] : Study: I/O Performance and Scheduling Overheads

Diego Didona, Jonas Pfefferle, Nikolas Ioannou, Bernard Metzler, and Animesh Trivedi. 2022. **Understanding modern storage APIs: a systematic study of libaio, SPDK, and io_uring**. In Proceedings of the 15th ACM International Conference on Systems and Storage (SYSTOR '22). Association for Computing Machinery, New York, NY, USA, 120–127. https://doi.org/10.1145/3534056.3534945

Zebin Ren and Animesh Trivedi. 2023. **Performance Characterization of Modern Storage Stacks: POSIX I/O, libaio, SPDK, and io_uring.** In Proceedings of the 3rd Workshop on Challenges and Opportunities of Efficient and Performant Storage Systems (CHEOPS '23). Association for Computing Machinery, New York, NY, USA, 35–45. https://doi.org/10.1145/3578353.3589545

Zebin Ren, Krijn Doekemeijer, Nick Tehrany, Animesh Trivedi. 2024. BFQ, **Multiqueue-Deadline, or Kyber? Performance Characterization of Linux Storage Schedulers in the NVMe Era**, to appear in the 2024 ACM/SPEC International Conference on Performance Engineering (ICPE '23), London, UK.
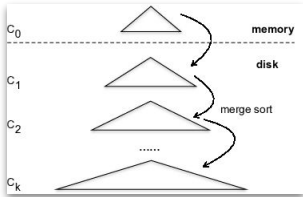
# [Part - 2/2] : Zone Namespace Devices (ZNS) Performance Characterization

Krijn Doekemeijer, Nick Tehrany, Bala Chandrasekaran, Matias Bjørling and Animesh Trivedi. **Performance Characterization of NVMe Flash Devices with Zoned Namespaces (ZNS).** 2023 IEEE International Conference on Cluster Computing (CLUSTER), Santa Fe, NM, USA, 2023, pp. 118-131, doi: https://doi.org/10.1109/CLUSTER52292.2023.00018 .
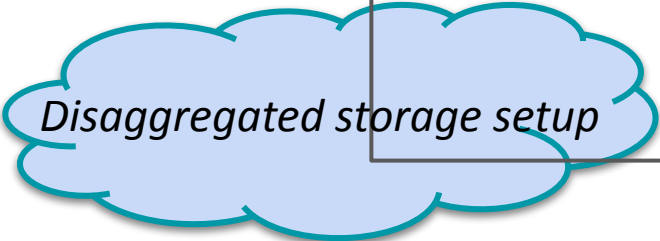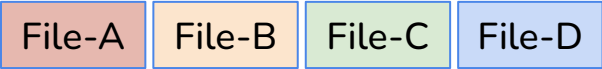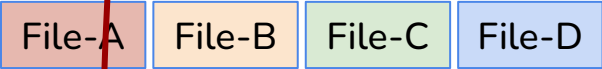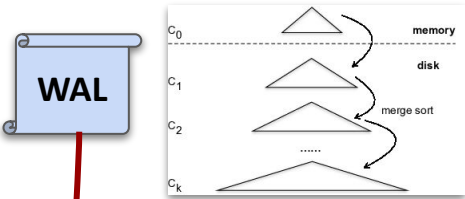
# Delivering QoS in a Distributed Setting

RocksDB

**WAL**

*Read, Write, Append, Reset, Finish, Close, Open,*

| File-A | File-B | File-C | File-D |

*Disaggregated storage setup*

# Delivering QoS in a Distributed Setting

# Delivering QoS in a Distributed Setting
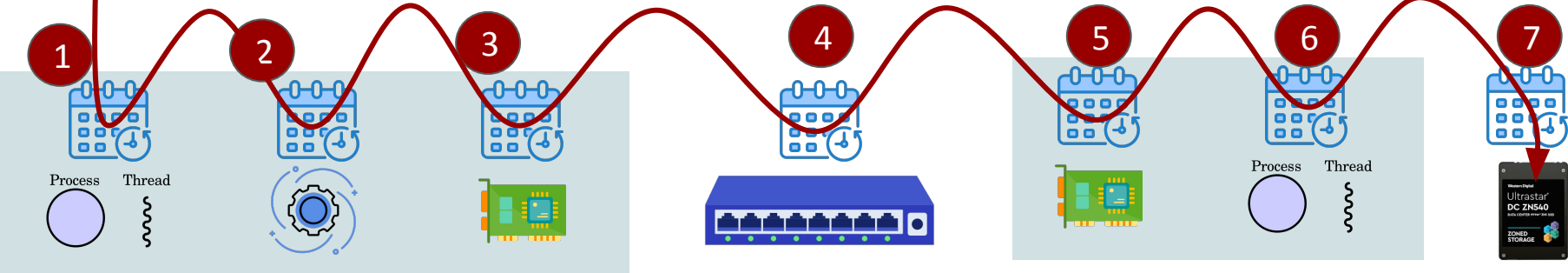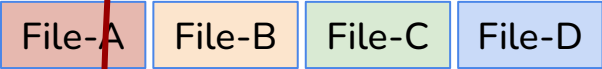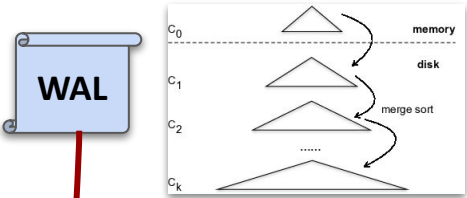


RocksDB

**WAL**

**End-to-End** abstraction for QoS:

**Co-design** workload-level storage-network data abstractions

**Co-schedule** them together (gang scheduling, co-flows)

File-A   File-B   File-C   File-D

1   2   3   4   5   6   7

Process   Thread

Process   Thread

# Conclusion

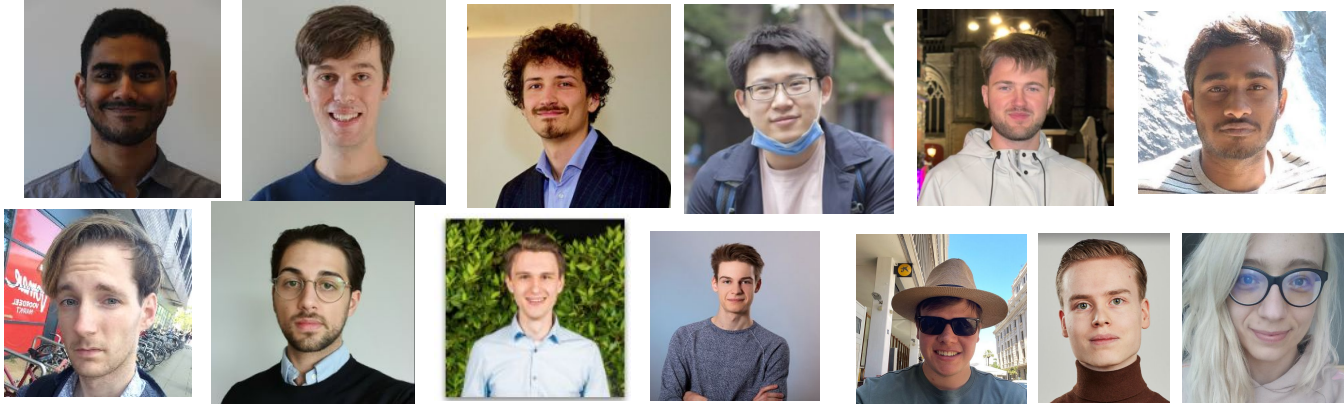**Vision: build your favorite workload-specialized data structure I/O stack!**

The era of <mark>workload-specialized storage stacks</mark> is here

We are exploring:

- Workload-specialized storage software abstractions
- Mapping software interfaces to the available hardware interfaces
  - NVMe ZNS, KV-SSD, CXL (emerging)

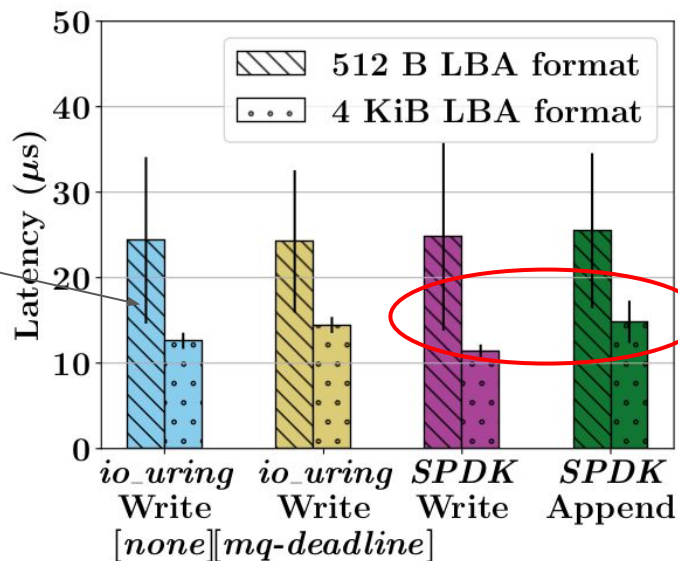**WiP: Co-scheduling (Network + ZNS Storage) ⇒ End-to-End QoS**

# Thank you!



*(past and present students)*

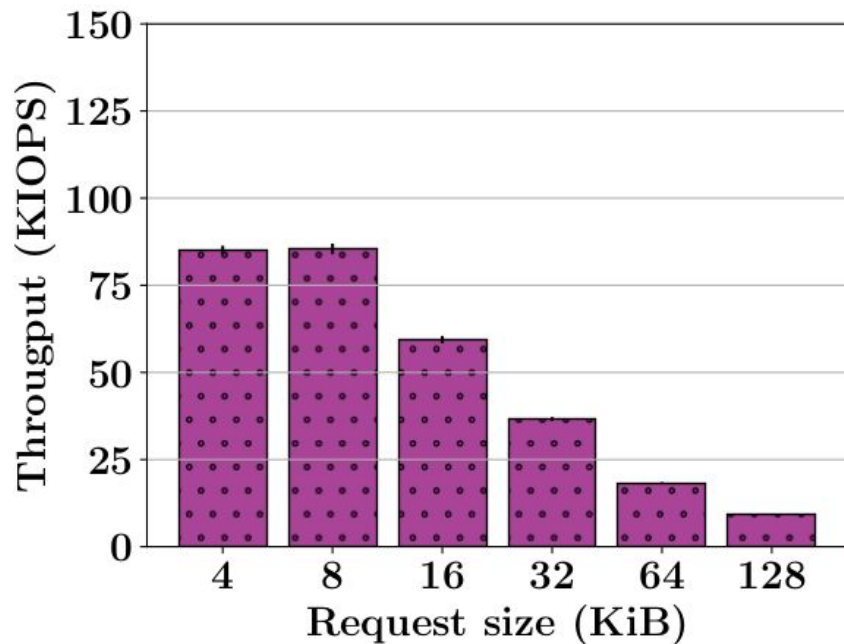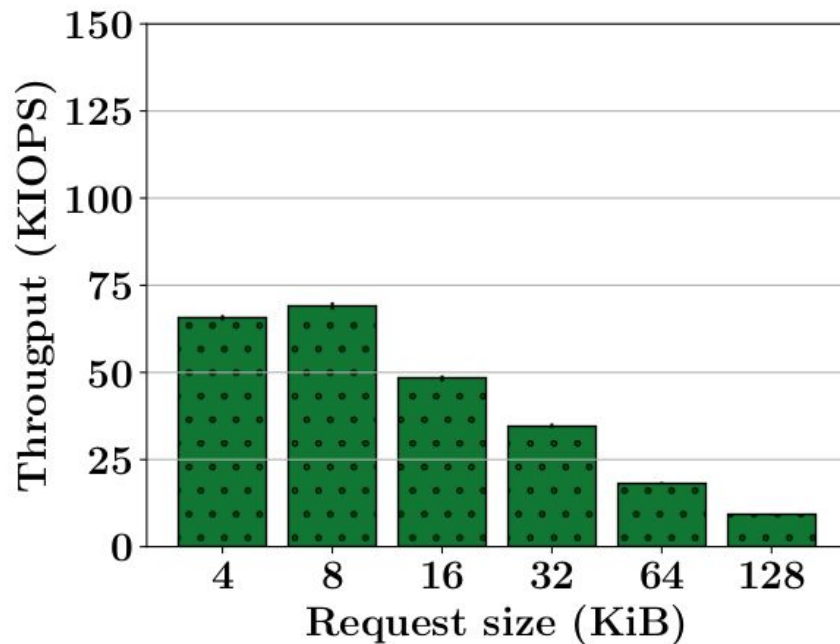# Result [1 / 4]: Write vs Append Latencies



Large gap in the LBA format

Writes lower than Append

- **4KiB block size has lower latencies (up to 2x)**
- **Writes have lower latencies than Append operations in our experiments**
- **SPDK has lower latencies than the Linux I/O stack (none, mq-deadline)**

# Write and Append: Bandwidth



(a) *write*

(b) *append*