# **BIZA: Design of Self-Governing Block-Interface ZNS AFA for Endurance and Performance**

### Shushu Yi, Shaocong Sun, Li Peng, Yingbo Sun Ming-Chang Yang, Zhichao Cao, Qiao Li, Myoungsoo Jung, Ke Zhou, Jie Zhang















### Why We Need All-Flash Array (AFA)?



#### All-flash array are widely adopted in diverse domains.





# **Open Question: I/O Interface Choices**



Zoned Namespace (ZNS) SSD



**EL2D** ZNS: Avoiding the Block Interface Tax for Flash-based SSDs (ATC'21) ZNS+: Advanced Zoned Namespace Interface for Supporting In-Storage Zone Compaction (OSDI'21)

**Conventional SSD** 

# **Open Question: I/O Interface Choices**



Zoned Namespace (ZNS) SSD



(ATC'21) ZNS: Avoiding the Block Interface Tax for Flash-based SSDs (OSDI'21) ZNS+: Advanced Zoned Namespace Interface for Supporting In-Storage Zone Compaction

**Conventional SSD** 

SELAD

# **Open Question: I/O Interface Choices**



Zoned Namespace (ZNS) SSD



(ATC'21) ZNS: Avoiding the Block Interface Tax for Flash-based SSDs (OSDI'21) ZNS+: Advanced Zoned Namespace Interface for Supporting In-Storage Zone Compaction

**Conventional SSD** 

ELAD

### **Existing Designs: Rob Peter to Pay Paul**

- Block-interface AFA: mdraid (Linux default)<sup>[1]</sup>, ScalaAFA (ATC'24)<sup>[2]</sup>
  - Random writes  $\rightarrow$  good compatibility, SSD-internal tasks  $\rightarrow$  bad performance & write amplification







### **Existing Designs: Rob Peter to Pay Paul**

- Block-interface AFA: mdraid (Linux default), ScalaAFA (ATC'24)
  - Random writes  $\rightarrow$  good compatibility, SSD-internal tasks  $\rightarrow$  bad performance & write amplification
- **ZNS-interface AFA:** RAIZN (ASPLOS'23)<sup>[1]</sup>

ASELAD

- host-managed tasks  $\rightarrow$  holistic designs for optimization, Sequential writes  $\rightarrow$  bad compatibility





### **Interface Adapter: A Simple Solution**

- Interface adapter: dm-zoned (Linux), dm-zap (Western Digital)
  - Converting ZNS interface to block interface
  - Maintaining mappings from block numbers to zone-related addresses



PEKING



## **Quantitative Analysis**

- Experimental setup: 3+1 (RAID 5) Western Digital ZN540 4TB SSDs
  - Peak throughput: 3265 MB/s and 2170 MB/s for read and write

### Challenge 1: Short endurance

- dmzap+RAIZN and mdraid+dmzap generate 33.3% and 54.6% more flash writes
- dm-zap muddles data with dfferent lifetimes in the same zones

### Challenge 2: Low throughput

- dmzap+RAIZN and mdraid+dmzap only achieve 47.7% and 18.4% of the ideal throughput
- dm-zap submits writes serially and only allows one in-flight write per zone
- RAIZN relies on a centralized zone to buffer the frequently updated metadata

### Challenge 3: High tail latency

- dmzap+RAIZN and mdraid+dmzap increase 99.99<sup>th</sup> tail latency by 10.3x and 2.2x after GC starts
- dm-zap handle GC and user I/O with the same I/O resources

Please refer to our paper for details.





#### WD ZN540 SSD



# **Our Solution: BIZA**



Key insight: new features & internal parallelism

• NVMe Technical Proposal 4076b: Zone Random Write Area (ZRWA)







#### • NVMe Technical Proposal 4076b: Zone Random Write Area (ZRWA)

- Break the strict sequential write constraint
- An efficient abstraction of the write buffer within SSDs
- Battery-backed DRAM, NVM, or FTL-mapped high-endurance flash (e.g., SLC)







#### • NVMe Technical Proposal 4076b: Zone Random Write Area (ZRWA)

- Break the strict sequential write constraint
- An efficient abstraction of the write buffer within SSDs
- Battery-backed DRAM, NVM, or FTL-mapped high-endurance flash (e.g., SLC)







#### • NVMe Technical Proposal 4076b: Zone Random Write Area (ZRWA)

- Break the strict sequential write constraint
- An efficient abstraction of the write buffer within SSDs
- Battery-backed DRAM, NVM, or FTL-mapped high-endurance flash (e.g., SLC)



Idea: absorb the frequently updated data within ZRWA







#### • NVMe Technical Proposal 4076b: Zone Random Write Area (ZRWA)

- Break the strict sequential write constraint
- An efficient abstraction of the write buffer within SSDs
- Battery-backed DRAM, NVM, or FTL-mapped high-endurance flash (e.g., SLC)



Idea: absorb the frequently updated data within ZRWA







#### [1] Understanding storage traffic characteristics on enterprise virtual desktop infrastructure. (SYSTOR'17)

### **ZNS SSD Exploration: Zone Random Write Area**

Idea: absorb the frequently updated data within ZRWA

However, the size of ZRWA is very limited per ZNS SSD

ZNS SSD	Zone Capacity	ZRWA size per open zone	Total ZRWA size
WD ZN540	1077 MB	1 MB	14 MB
DapuStor J5500Z	18144 MB	1 MB	16 MB
Inpur NS8600G	96 MB	1440 KB	11.25 MB

Lab















- I/O stack has no guarantee on the orders of request submissions
  - Reorder / merge I/O for high performance [1]
- Implicit I/O reorders may cause failures for parallel writes







- I/O stack has no guarantee on the orders of request submissions
  - Reorder / merge I/O for high performance [1]
- Implicit I/O reorders may cause failures for parallel writes







- I/O stack has no guarantee on the orders of request submissions
  - Reorder / merge I/O for high performance [1]
- Implicit I/O reorders may cause failures for parallel writes







- I/O stack has no guarantee on the orders of request submissions
  - Reorder / merge I/O for high performance [1]
- Implicit I/O reorders may cause failures for parallel writes







- I/O stack has no guarantee on the orders of request submissions
  - Reorder / merge I/O for high performance [1]
- Implicit I/O reorders may cause failures for parallel writes







- I/O stack has no guarantee on the orders of request submissions
  Reorder / merge I/O for high performance <sup>[1]</sup>
- Implicit I/O reorders may cause failures for parallel writes
- Only allow one-inflight write  $\rightarrow$  loses up to 65.3% throughput







- I/O stack has no guarantee on the orders of request submissions
  - Reorder / merge I/O for high performance [1]
- Implicit I/O reorders may cause failures for parallel writes







-

U



#### • ZNS SSDs have multiple parallel I/O resources

- We call them I/O channels (may consist of multiple flash channels or chips)
- Zones from with separated I/O channels can handle I/O requests in parallel



#### SSD Architecture





Write Tests

#### • ZNS SSDs have multiple parallel I/O resources

- We call them I/O channels (may consist of multiple flash channels or chips)
- Zones from with separated I/O channels can handle I/O requests in parallel
- Give up the centralized metadata zone
- Schedule user I/O and GC with parallel zones



#### SSD Architecture





Write Tests

#### • ZNS SSDs have multiple parallel I/O resources

- We call them I/O channels (may consist of multiple flash channels or chips)
- Zones from with separated I/O channels can handle I/O requests in parallel
- Give up the centralized metadata zone



#### SSD Architecture







### **BIZA Overview**

#### Zone group selector

- Isolate hot chunks from cold ones
- For better utilization of the scarce ZRWA
- A set of zones mapped with diff. channels

#### • GC avoidance mechanism

- Schedule I/O and GC with parallel zones
- Exhaust inter-zone parallelism

#### ZRWA-aware I/O scheduler

- Dispatches I/O requests in parallel
- Exhaust intra-zone parallelism

#### I/O channel detection

- Collect I/O completion latency
- Confirm the mappings between channels & zones



#### **Overview of BIZA**





### **BIZA Overview**

### • Zone group selector

- Isolate hot chunks from cold ones
- For better utilization of the scarce ZRWA
- A set of zones mapped with diff. channels
- GC avoidance mechanism\*
  - Schedule I/O and GC with parallel zones
  - Exhaust inter-zone parallelism

### • ZRWA-aware I/O scheduler\*

- Dispatches I/O requests in parallel
- Exhaust intra-zone parallelism

#### I/O channel detection\*

- Collect I/O completion latency
- Confirm the mappings between channels & zones

#### \*Please refer to our paper for more details.





#### **Overview of BIZA**



### **Zone Group Selector**

- ZRWA will shift right when multiple new writes arrives
- Isolate hot data with cold ones  $\rightarrow$  Keep hot data in ZRWA for a longer time
- Ghost-cache-based algorithm <sup>[1]</sup>
  - High revenue = the chunk will be updated multiple times
  - High profit = High revenue + Small cost (the chunk has short reuse distance)







### **Zone Group Selector**

- ZRWA will shift right when multiple new writes arrives
- Isolate hot data with cold ones  $\rightarrow$  Keep hot data in ZRWA for a longer time
- Ghost-cache-based algorithm <sup>[1]</sup>
  - High revenue = the chunk will be updated multiple times
  - High profit = High revenue + Small cost (the chunk has short reuse distance)

#### • Select different zone group for chunks in different cache

- High profit: can be updated multiple times within the limited ZRWA
- High revenue: hard to be absorbed in ZRWA but is the main source of GC





### **Experimental Setup**

- Implemented as a device mapper in Linux kernel
- Constructing AFA with 4 commodity SSDs as RAID 5
- **ZNS SSD:** Western Digital ZN540 SSD
  - Read: 3265 MB/s, Write: 2170 MB/s
  - Up to 14 open zones & 14 MB ZRWA
- Conventional SSD: Western Digital SN640 SSD
  - Read: 3331 MB/s, Write: 2250 MB/s
  - Developed on a similar hardware basis as ZN540 SSD

#### Comparison

- RAIZN: only support sequential write
- dmzap+RAIZN: stacking dm-zap on RAIZN for random writes
- mdraid+dmzap: stacking dm-zap on each ZNS SSD
- mdraid+ConvSSD: constructing AFA with conventional SSDs





#### WD SN640 SSD





### Microbenchmark

#### Bandwidth

- BIZA achieves 3.7x and 3.5x bandwidth than dmzap+RAIZN and mdraid+dmzap







### Microbenchmark

#### Bandwidth

- BIZA achieves 3.7x and 3.5x bandwidth than dmzap+RAIZN and mdraid+dmzap
- Even 1.4x bandwidth than mdraid+ConvSSD (software overhead of mdraid<sup>[1-2]</sup>)
- Almost exhaust the throughput of 4 ZNS SSDs (92.2% of the ideal)







[1] ScalaRAID: Optimizing linux software raid system for next-generation storage. (ATC'24)[2] stRAID: Stripe-threaded architecture for parity-based RAIDs with ultra-fast SSDs. (ATC'22)



### Microbenchmark

#### Bandwidth

- BIZA achieves 3.7x and 3.5x bandwidth than dmzap+RAIZN and mdraid+dmzap
- Even 1.4x bandwidth than mdraid+ConvSSD (software overhead of mdraid<sup>[1-2]</sup>)
- Almost exhaust the throughput of 4 ZNS SSDs (92.2% of the ideal)

#### • Average latency

Outperform RAIZN by 53.8% for sequential writes











### **Applications**

SELab

- Filesystem: F2FS + filebench
  - BIZA outperforms RAIZN by 26.6%, 24.9%, and 18.7% in randomwrite, fileserver, and oltp







### **Applications**

SELab

- Filesystem: F2FS + filebench
  - BIZA outperforms RAIZN by 26.6%, 24.9%, and 18.7% in randomwrite, fileserver, and oltp
  - Most (95.2%) requests in webserver are read requests







### **Applications**

- Filesystem: F2FS + filebench
  - BIZA outperforms RAIZN by 26.6%, 24.9%, and 18.7% in randomwrite, fileserver, and oltp
  - Most (95.2%) requests in webserver are read requests
- Key-value store: RocksDB + F2FS + db\_bench
  - 10.5% higher throughput than RAIZN







RAIZN

mdraid+dmzap

mdraid+ConvSSD

## Write Amplification Reduction

- Reducing **41.0%** flash writes (compared with mdraid+dmzap)
- Different sizes of ZRWA: larger ZRWA, less flash writes







### Conclusion

#### • Existing interface choices of AFA are all unsatisfactory

- Block-interface alone, ZNS-interface alone, and adapter
- Endurance, performance, and compatibility

#### • **BIZA**: benefit from the openness of ZNS whilst exposing block interface

- Fully exploiting the emerging ZRWA feature
- Exhausting the intra-zone and inter-zone parallelism
- Significantly improve I/O performance while mitigates write amplification
- Source code is accessible at:









# Thanks & QA

BIZA: Design of Self-Governing Block-Interface ZNS AFA for Endurance and Performance https://github.com/ChaseLab-PKU/BIZA

> Shushu Yi, Shaocong Sun, Li Peng, Yingbo Sun Ming-Chang Yang, Zhichao Cao, Qiao Li, Myoungsoo Jung, Ke Zhou, Jie Zhang















