

2024 USENIX Annual Technical Conference

ZMS: Zone Abstraction for Mobile Flash Storage

2024/7/10

Joo-Young Hwang¹, Seokhwan Kim¹, Daejun Park¹, Yong-Gil Song¹, Junyoung Han¹,
Seunghyun Choi¹, Sangyeun Cho¹ and Youjip Won²

¹ *Samsung Electronics Device Solutions*

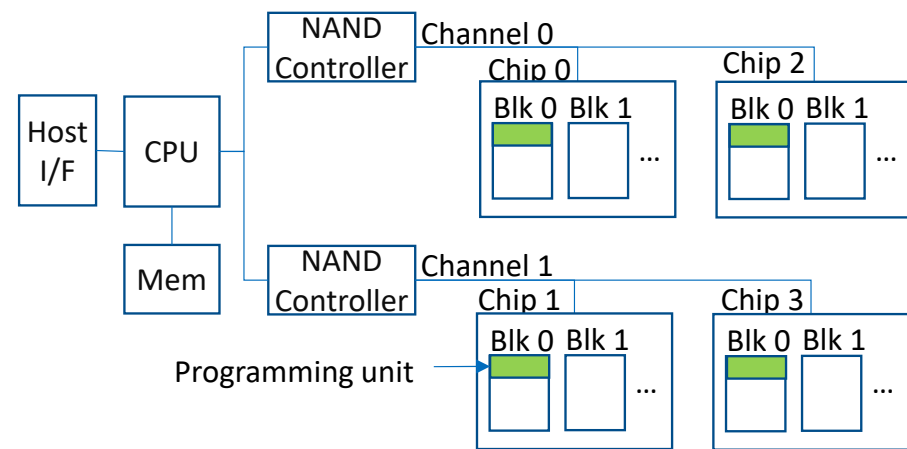
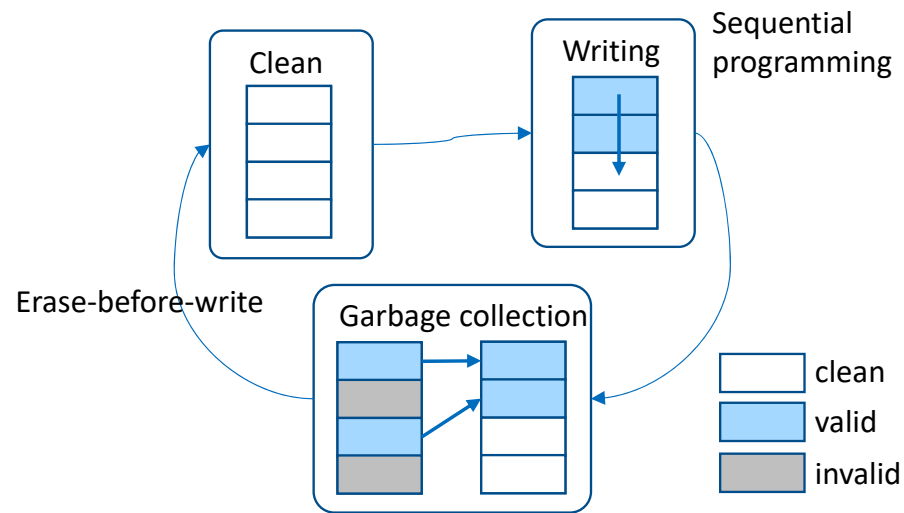
² *Korea Advanced Institute of Science and Technology*



Background

Conventional SSD (Block interface)

- Flash memory has sequential programming and “erase-before-write” characteristics.
- To provide block interface for flash memory, SSD performs logical-to-physical (L2P) address mapping and garbage collection.
- Data striping over multiple chips for parallel operation.



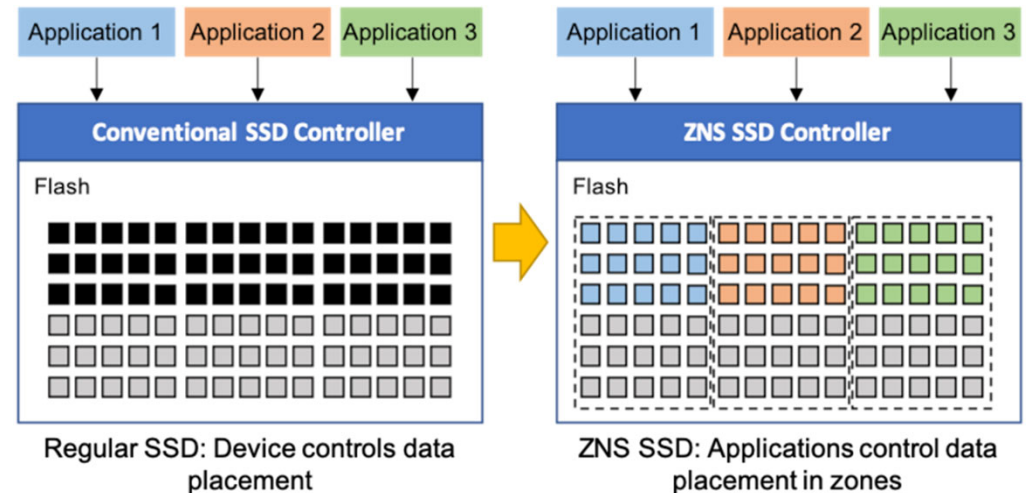
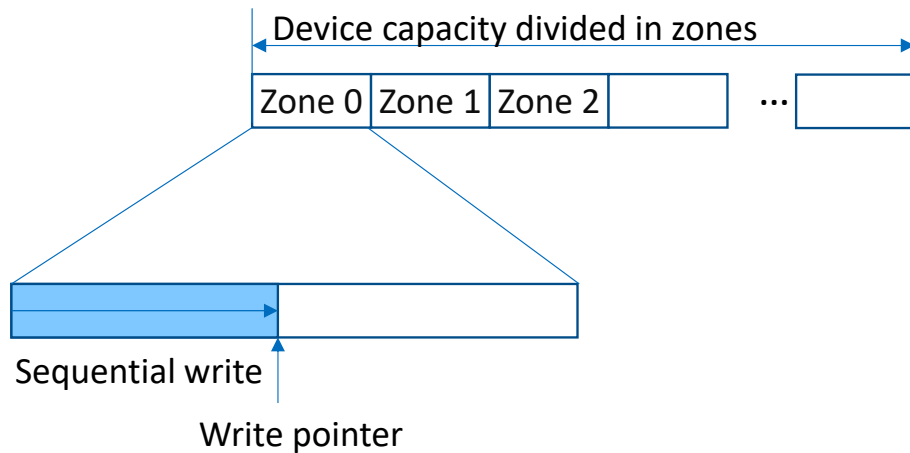
- **Superblock:** the set of flash blocks where data is striped
- **Superpage:** the set of programming units on the same offset of the superblock

• Issues

- Memory cost for L2P mapping table (← page mapping)
- Write amplification (← data with different lifetimes are mixed)

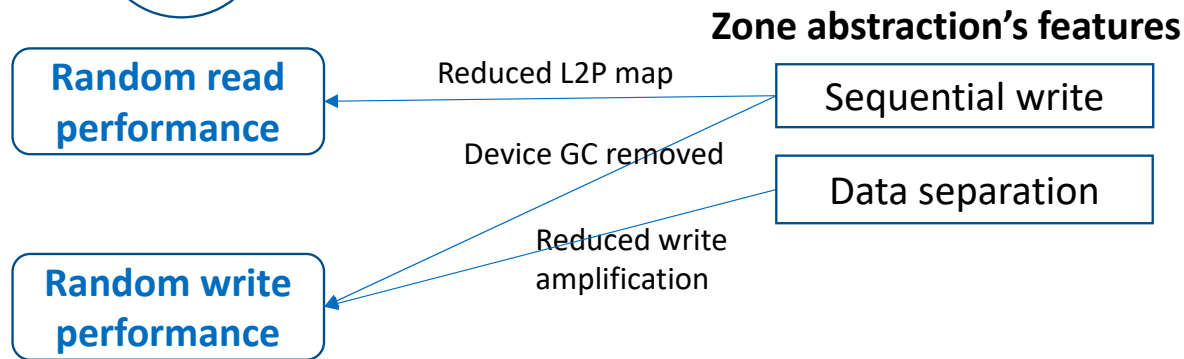
ZNS (Zoned Name Space) SSD

- Sequential write constraint → coarse-grain (zone) mapping → reduced L2P map
- Host controls data placement → removes device garbage collection
 - Host is responsible for data separation to reduce write amplification.



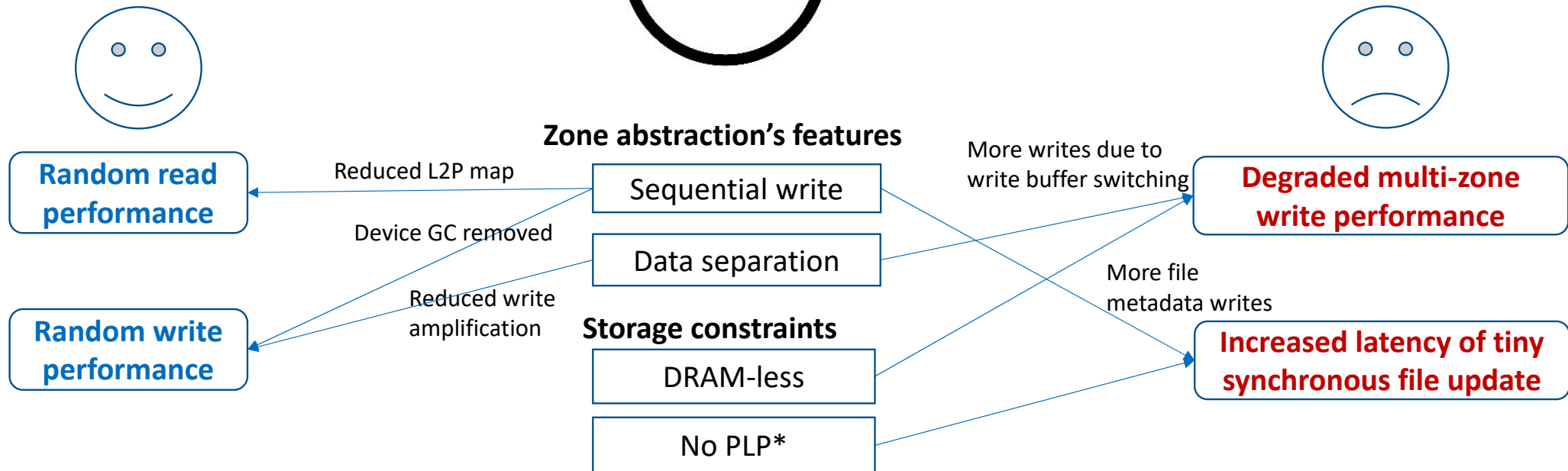
Does mobile storage benefit from zone abstraction?

Responsiveness is critical in mobile devices.



Does mobile storage benefit from zone abstraction?

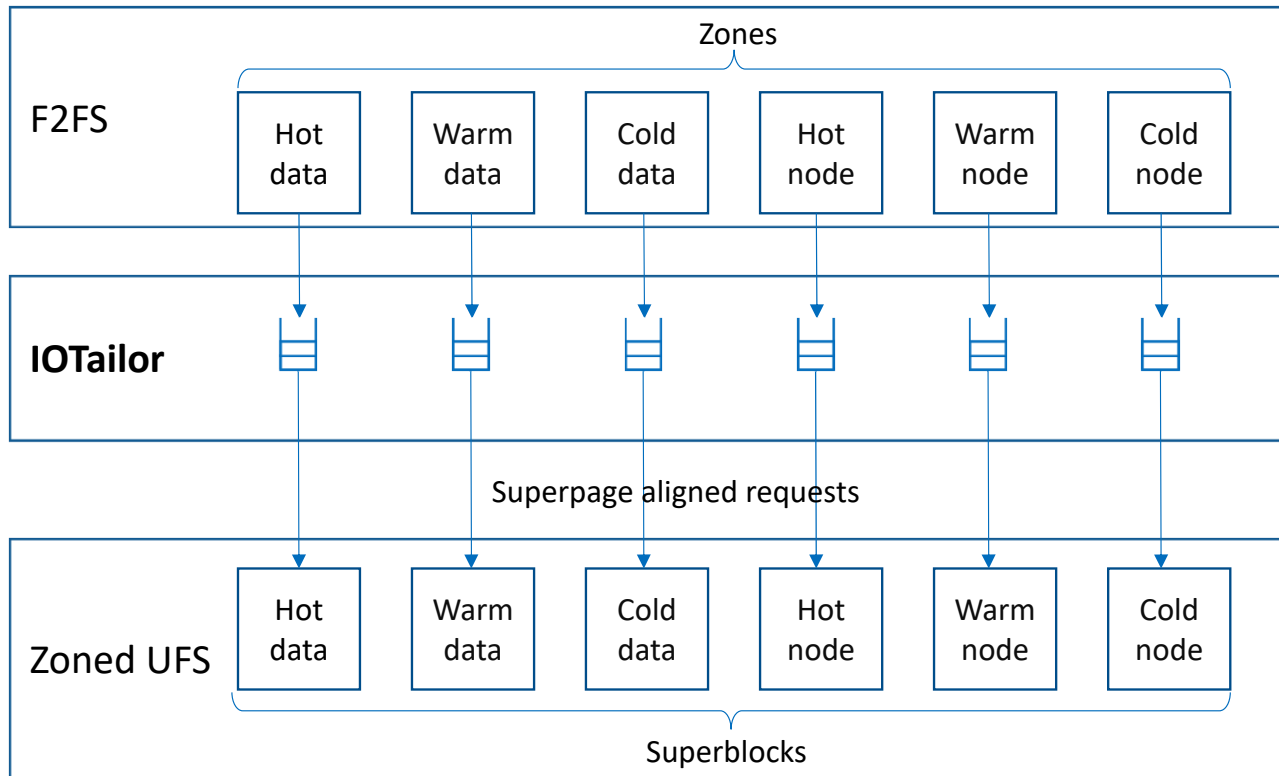
Responsiveness is critical in mobile devices.



*PLP: Power loss protection

ZMS Overview

- Utilizing F2FS*, data is separated according to six temperature types.
- Techniques to address the challenges: **IOTailor**, **budget-based in-place update**
- Optimization techniques: copy offloading, multi-granularity mapping (not covered in this talk)



- **Budget-based in-place update**
- Copy offloading to reduce GC cost

- Multi-granularity mapping

*F2FS: A new file system for flash storage, Lee et al. USENIX FAST '15

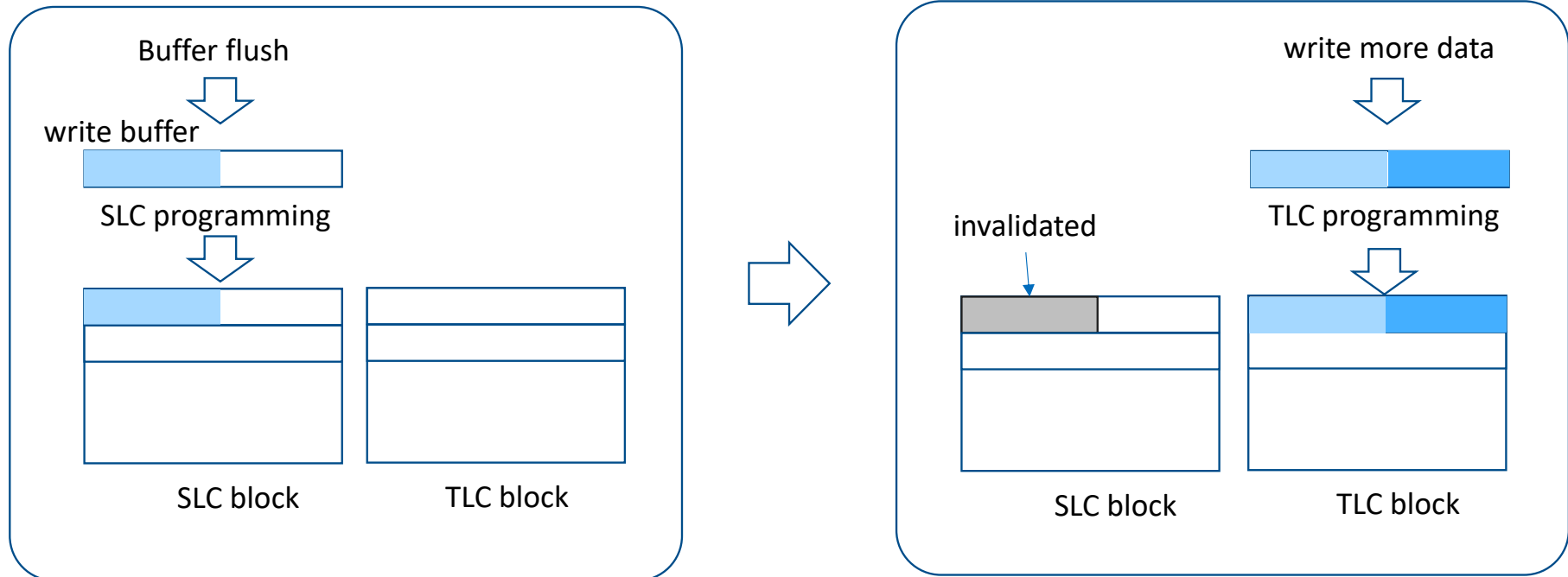
Talk Outline

- Challenge #1: Multi-zone write performance
- Challenge #2: Latency of Tiny Synchronous File Update
- Evaluation
- Conclusion

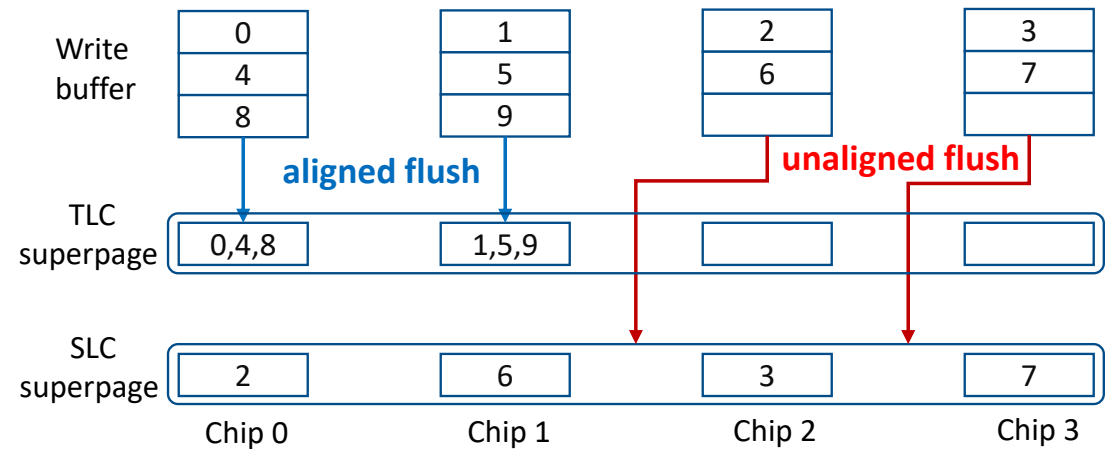
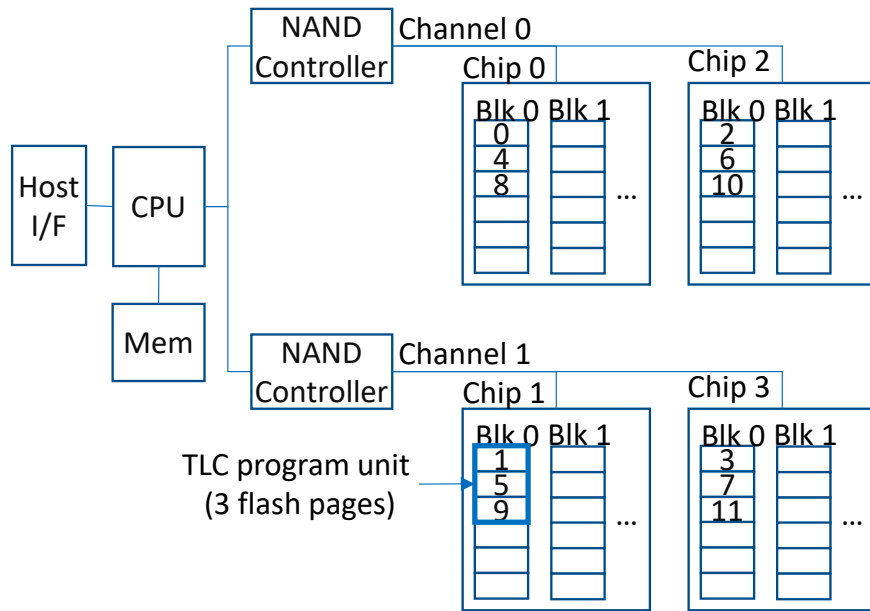
Challenge #1: Multi-zone Write Performance

SLC Buffering to Handle Unaligned Buffer Flushes

- Unaligned buffer flush: flush data that is smaller than TLC programming unit.
- Backup data to SLC, later migrate data to TLC
- Side-effect: double writes



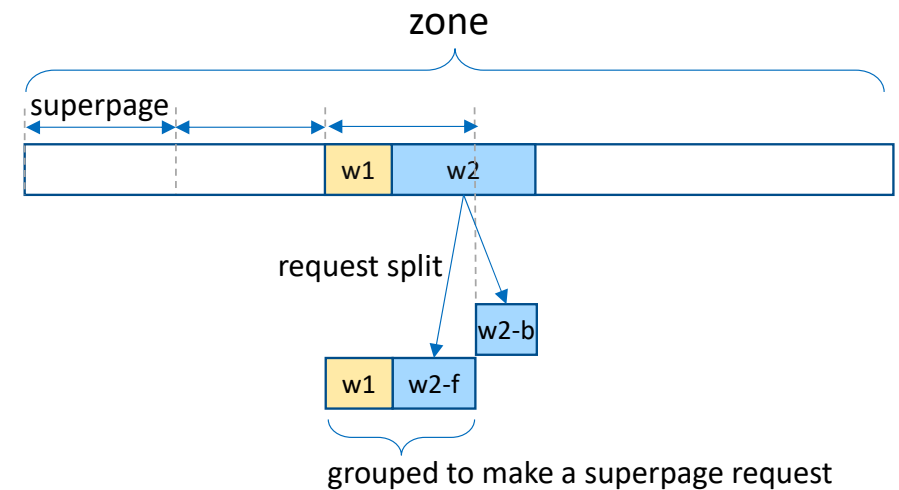
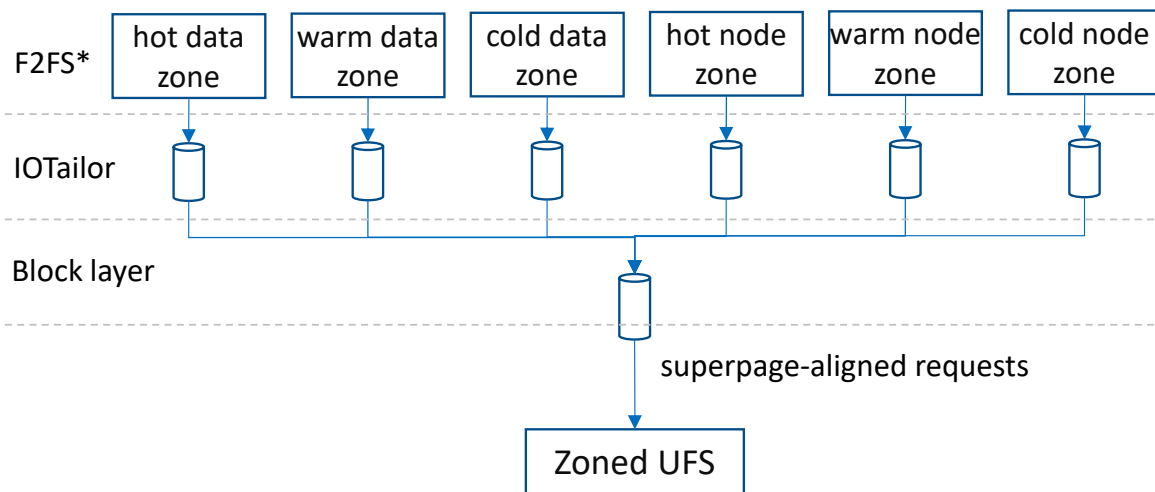
Example of Unaligned Buffer Flush Handling



Stripe unit: 32KiB (flash page)
 Superpage = 12 stripe units (384 KiB)

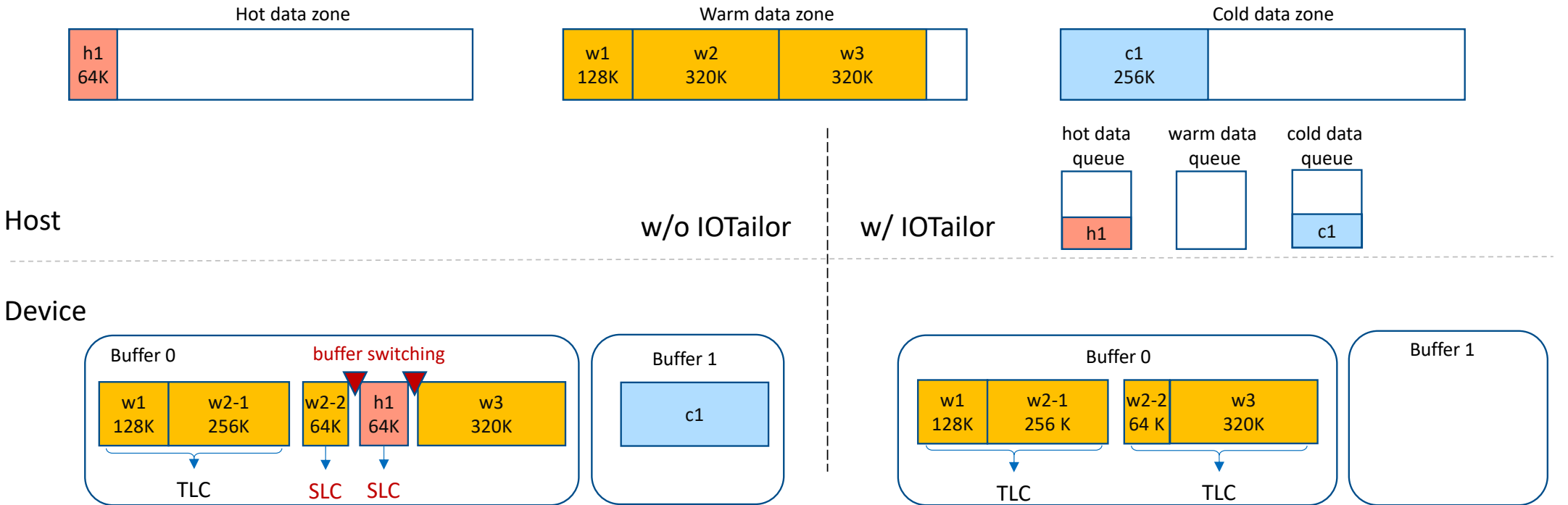
IOTailor

- Superpage-aligned request is good for parallelism and avoids unaligned buffer flushes.
- For each zone, IOTailor transforms requests to superpage-aligned requests
- Request split & request grouping in the per-zone queues



Example of Writing to Multiple Zones

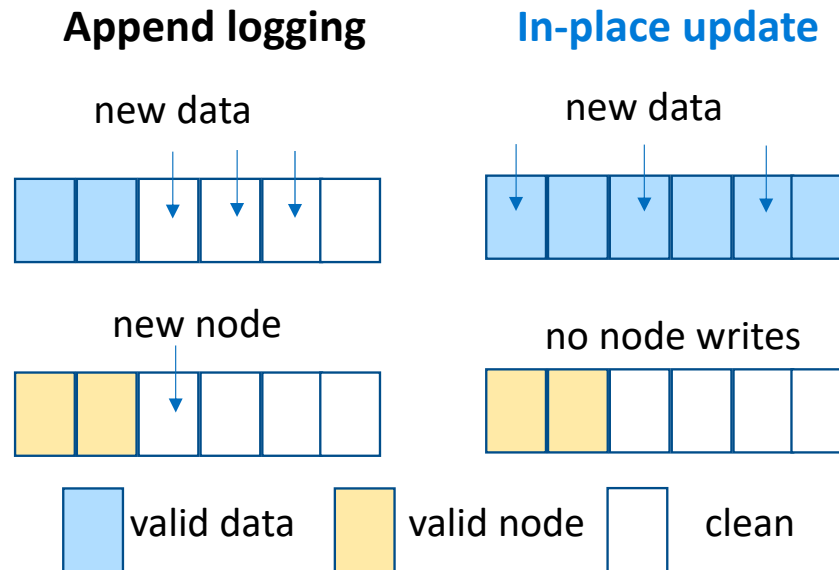
Command order: w1 – c1 – w2 – h1 – w3



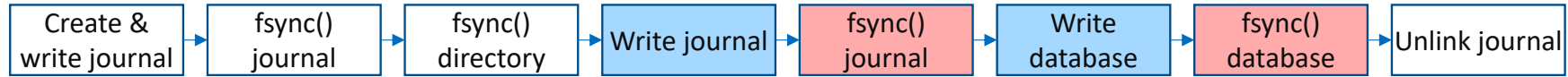
Challenge #2:
Latency of Tiny Synchronous
File Update

F2FS write optimization does not work for zoned device

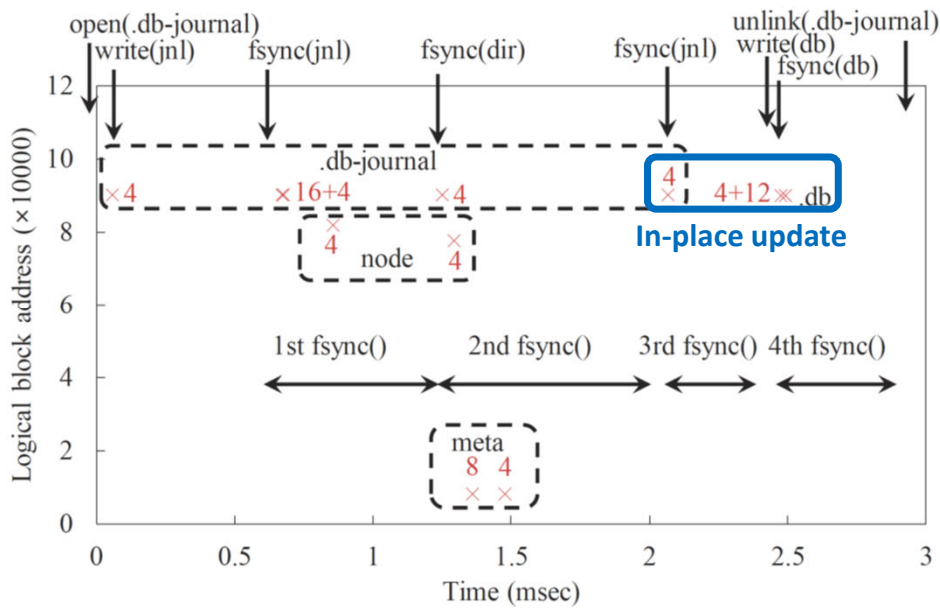
- Tiny synchronous file update is latency critical.
- For conventional SSDs, F2FS uses in-place update policy to reduce the latency of tiny (< 32KiB) synchronous file update.
- **In-place update policy cannot be used on zoned devices.**



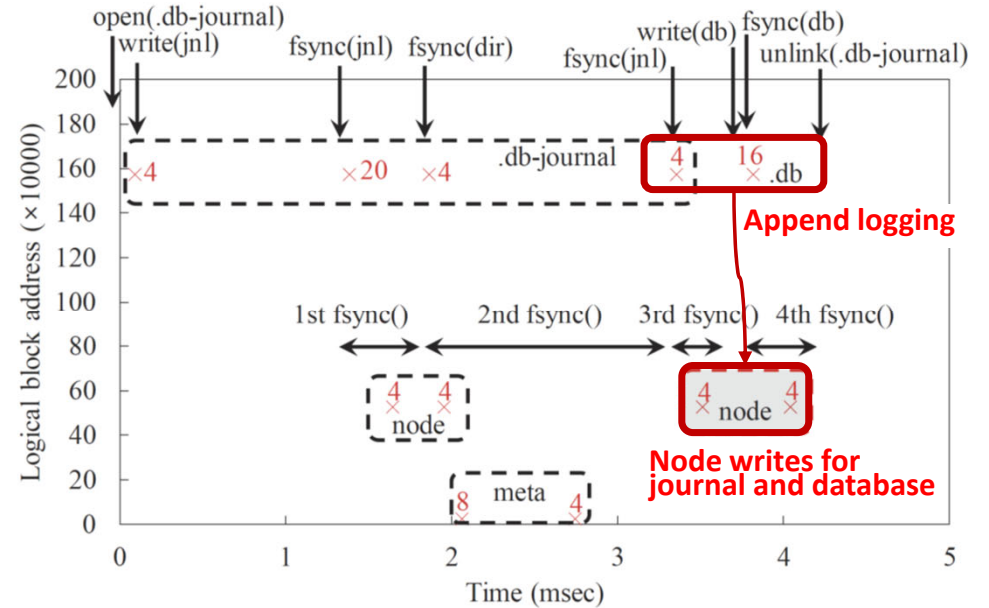
I/O Pattern of SQLite insert() Transaction



<Block device>

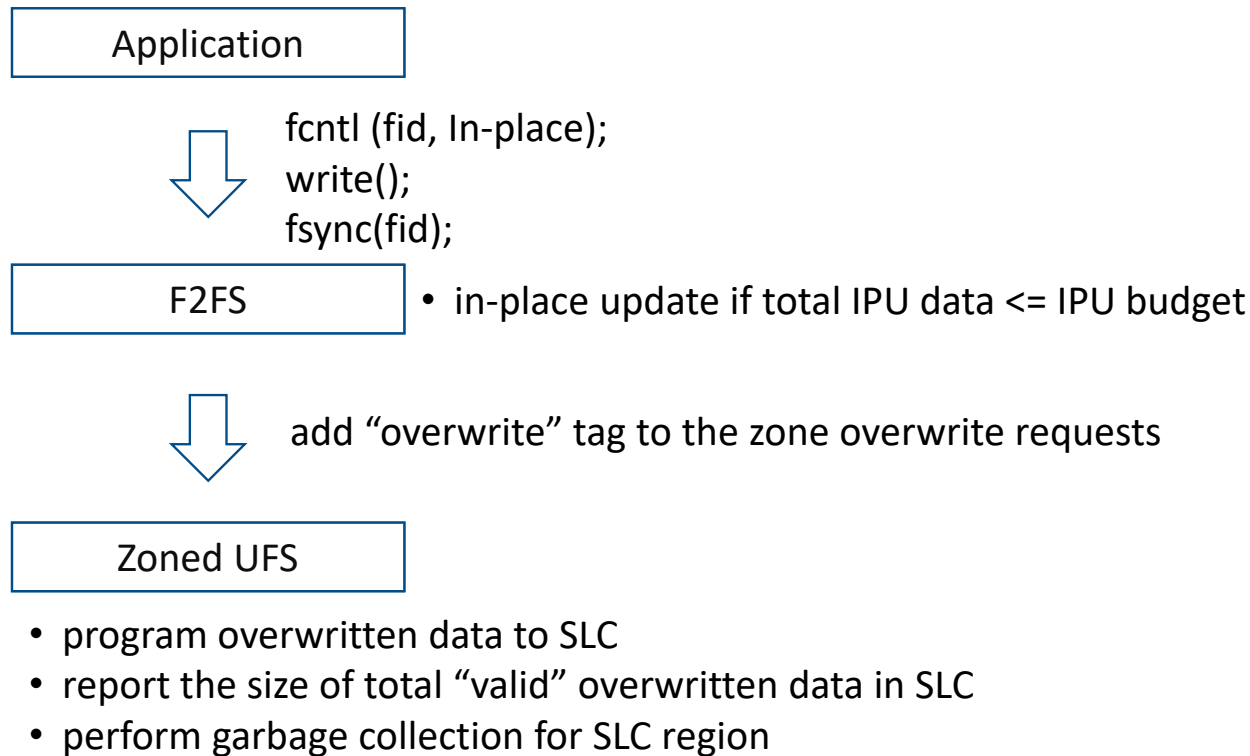


<Zoned device>



Budget-based In-Place Update

- Allow in-place update for select files as per the application request
- Device writes the In-place updated data into SLC blocks.
- Cap total valid data size for efficient garbage collection.

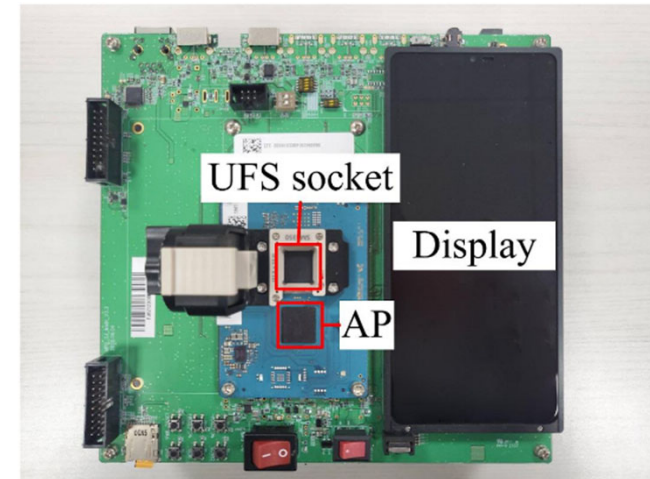


Evaluation

- How much are the benefits of zone abstraction for mobile storage?
 - are the challenges addressed well?

Evaluation Setup

- Host platform: SM8350 (8 cores), 12GiB DRAM, Android 11, Linux 5.4
- Zoned UFS: 128GiB, UFS 2.1
 - zone size: 138MiB
 - a conventional logical unit for F2FS meta area
- Baseline: the same device with a firmware that supports legacy block interface



| Workload | Configuration |
|--|--|
| Sequential read/write | Fio ¹ , 512 KiB IO size |
| Buffered random read/write | Fio, 4 KiB IO over 1 GiB file |
| Synchronous random write | Fio, 4 KiB write followed by fsync() |
| Wide range random read | Fio, 4 KiB read over 8 GiB file |
| Concurrent writing to multiple zones | Three concurrent Fio writing jobs, each writing to its own files |
| SQLite benchmark (Mobibench ²) | 1M insert(), 3.9 MiB WAL ³ file, 385 MiB database file |
| Application launch | Category (number of apps): basic (8), image (3), video (5), education (4), game (17) |

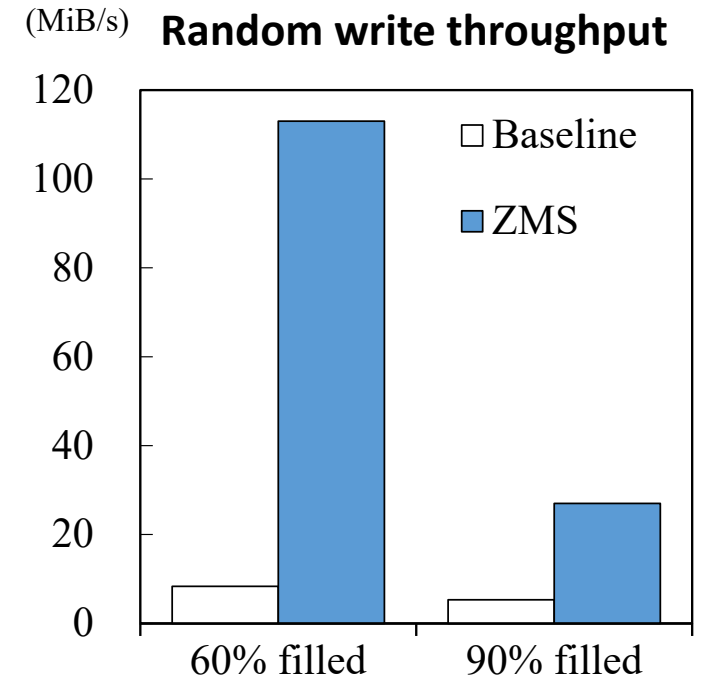
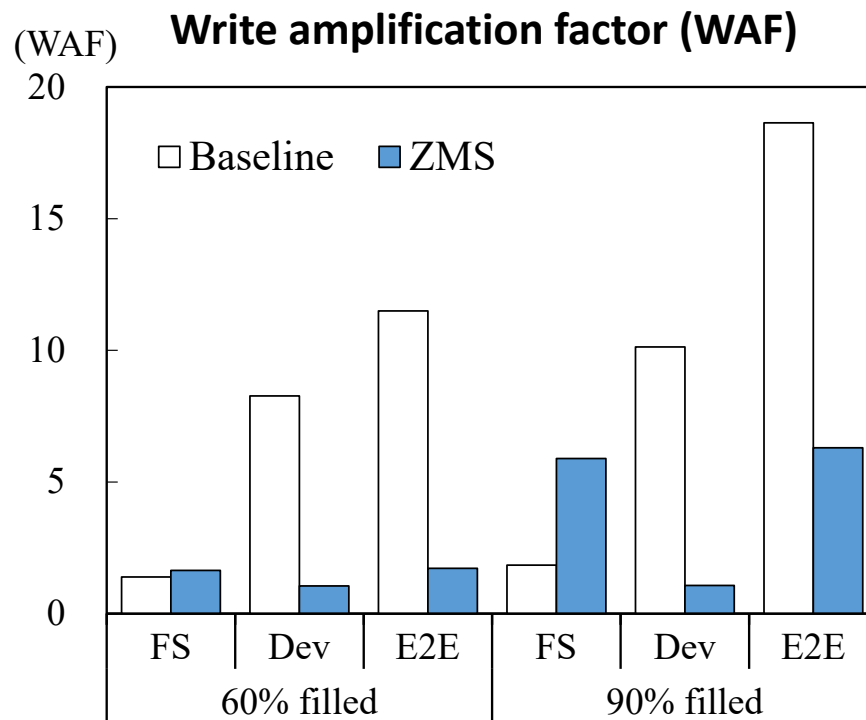
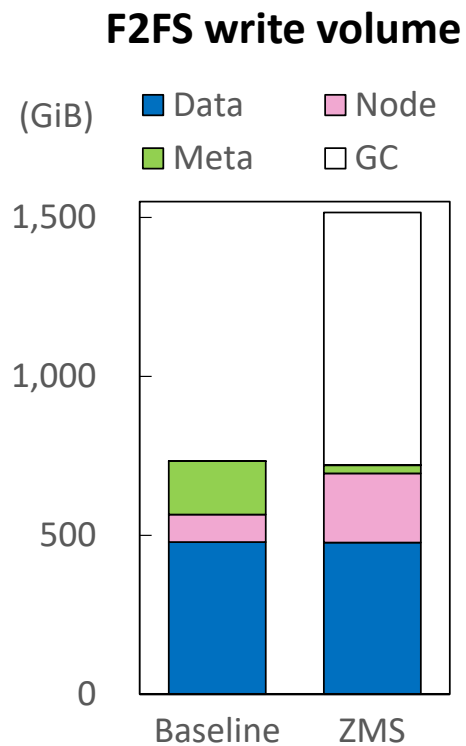
¹FIO: https://fio.readthedocs.io/en/latest/fio_doc.html

²Mobibench: <https://github.com/ESOS-Lab/Mobibench>

³WAL: Write ahead logging

Random Write Performance & Write Amplification

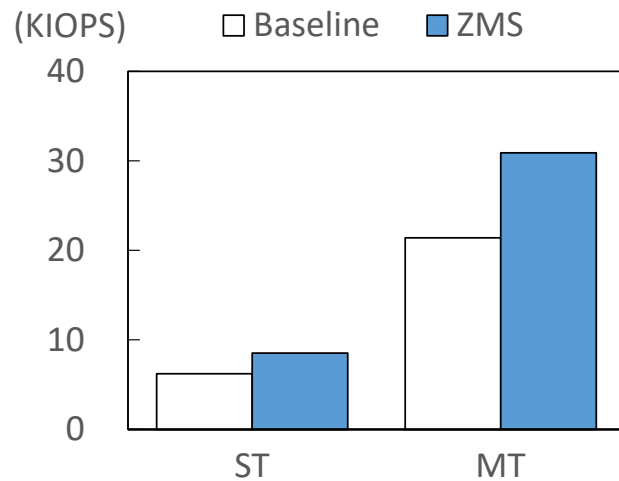
- **2.85x ~ 6.4x** lower write amplification
- **5x ~ 13.6x** higher random write throughput



Random Read Performance

- **37~44%** better random read performance
- Application launch time: **5.8 ~ 11.6%** reduction

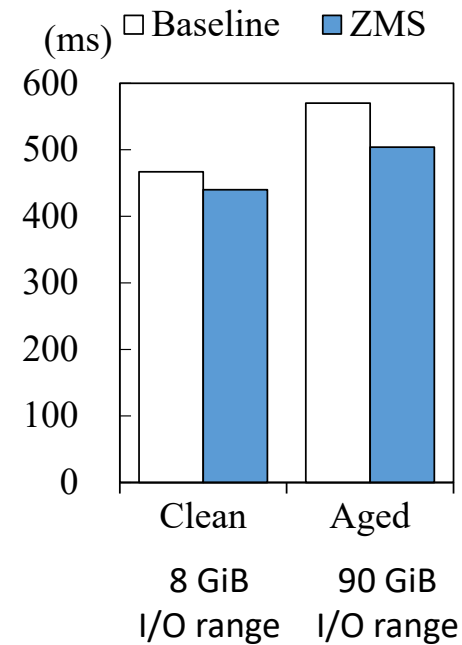
8GiB file random read



Baseline: map cache miss ratio: 27.1%

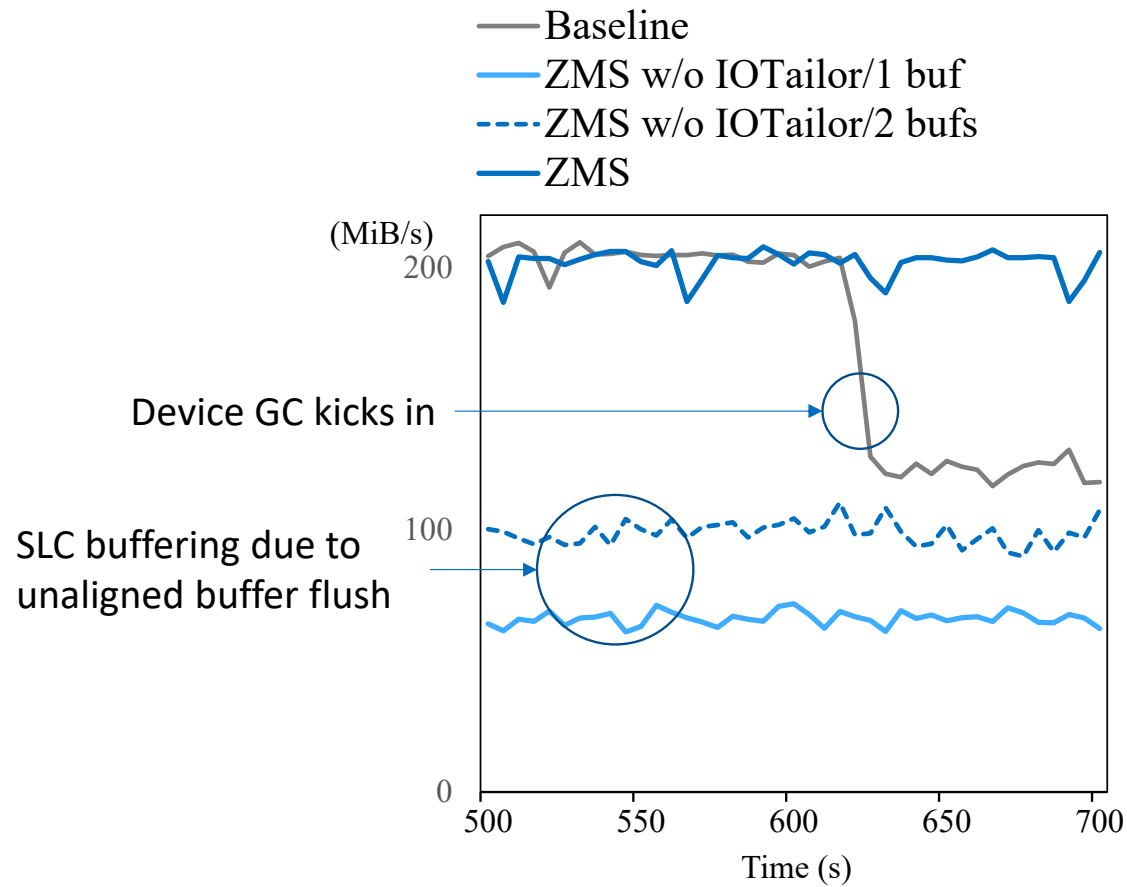
ZMS: no map cache miss

Application launch time



Performance of Writing to Multiple Zones

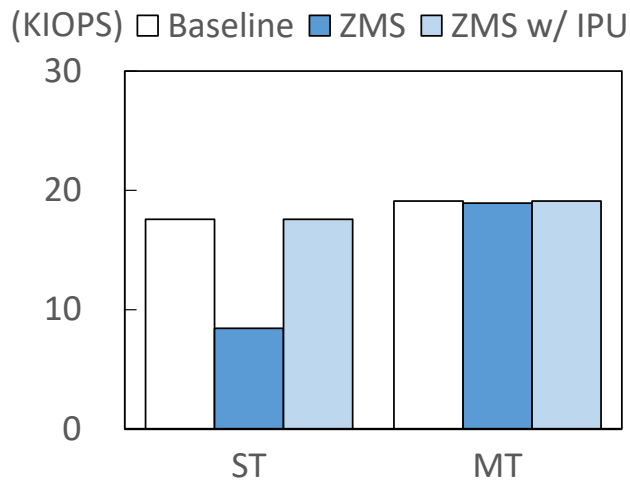
- IOTailor improves multi-zone write performance by reducing SLC buffering.



Synchronous Update Performance

- Using the budget-based in-place update, ZMS shows no performance degradation in tiny synchronous update and SQLite rollback journal mode.
- **60~100%** performance gain in write-ahead log (WAL) mode (append logging).

4KiB write + fsync()

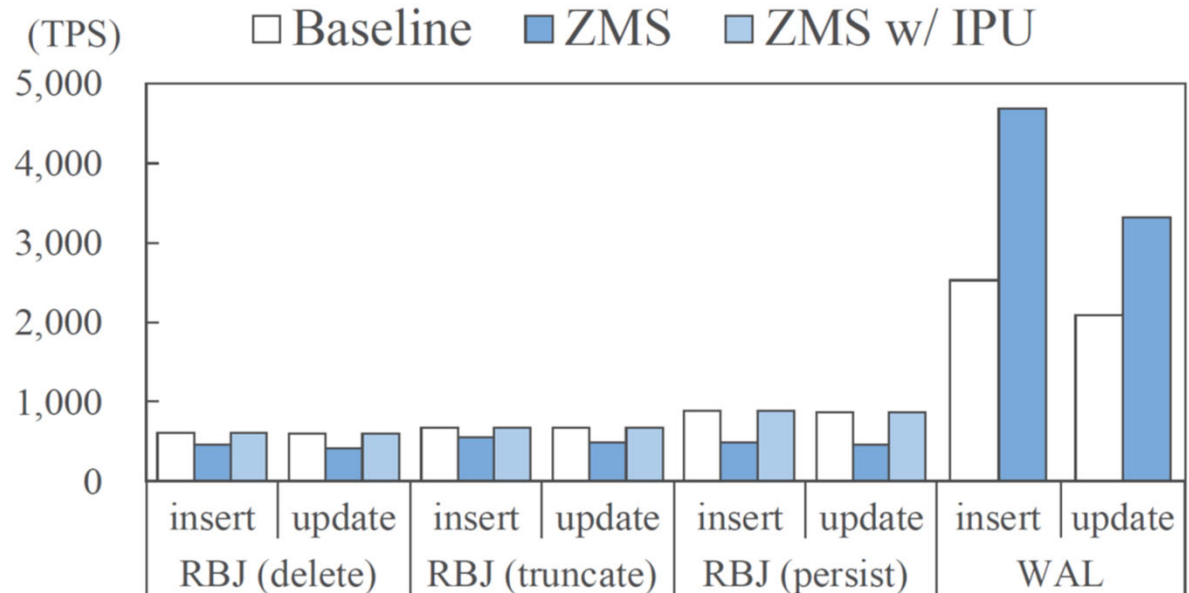


Baseline: In-place update

ZMS: append logging

ZMS w/ IPU: budget-based in-place update

SQLite throughput (Mobibench*)



Conclusion

- Zone abstraction is promising for enhancing responsiveness of mobile devices.
- Two challenges when employing zoned mobile storage
 - Degraded multi-zone write performance
 - Increased latency of tiny synchronous file update
- ZMS techniques to address the challenges
 - IOTailor improves performance of writing to multiple zones by avoiding unaligned buffer flushes due to buffer switching.
 - Budget-based in-place update improves synchronous update performance.
- ZMS improves random read/write performance and write amplification significantly.

Acknowledgement

- The authors express sincere gratitude to Jaegeuk Kim, Bart Van Assche and all other individuals who have contributed to zoned UFS standardization and Linux kernel support for zoned devices. Such collective efforts contributed to establishing ecosystem for zoned UFS.