# Compaction-Aware Zone Allocation for LSM based Key-Value Store on ZNS SSDs

**Hee-Rock Lee, Chang-Gyu Lee, Seungjin Lee, and Youngjae Kim**
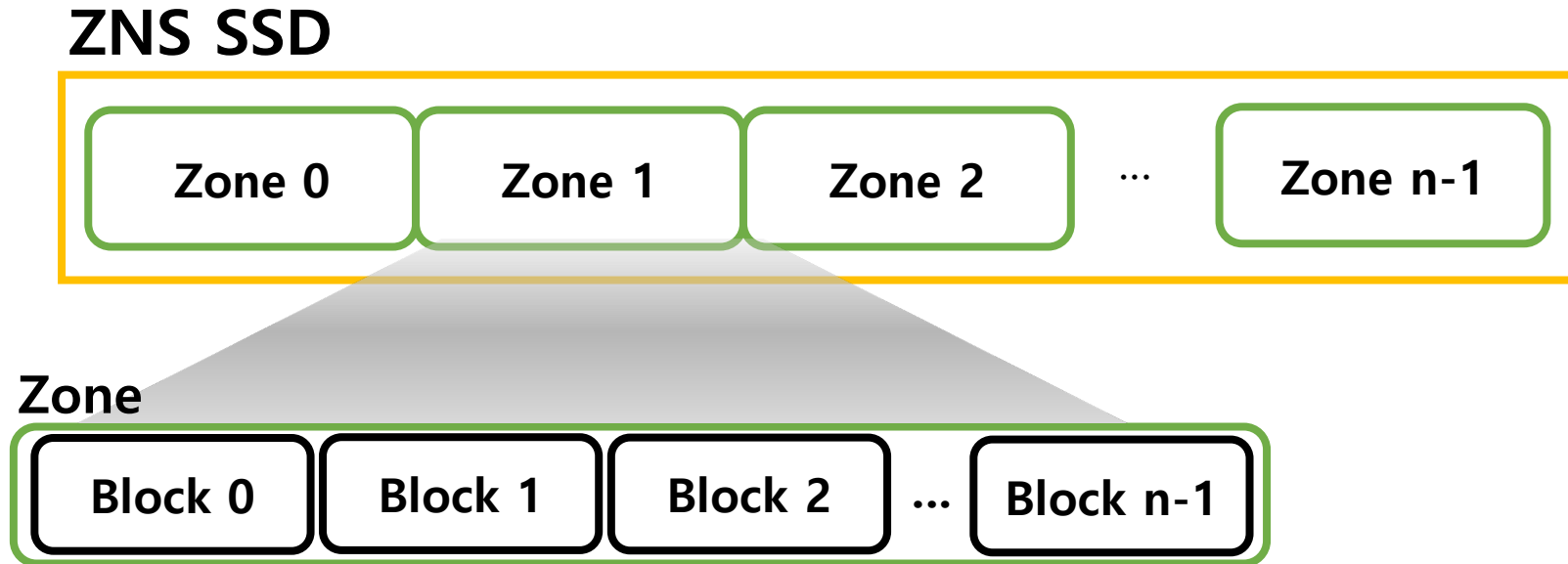
Sogang University, Seoul, South Korea

The 14th ACM Workshop on Hot Topics in Storage and File Systems
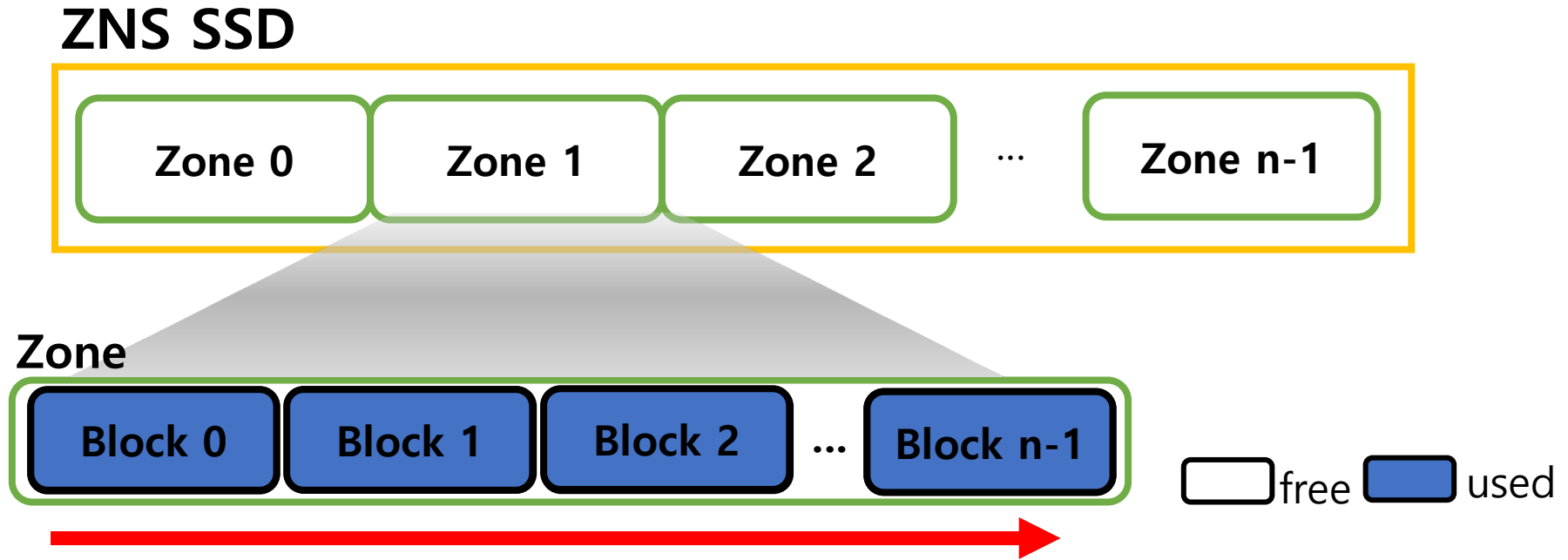**(HotStorage'22, June 27-28)**
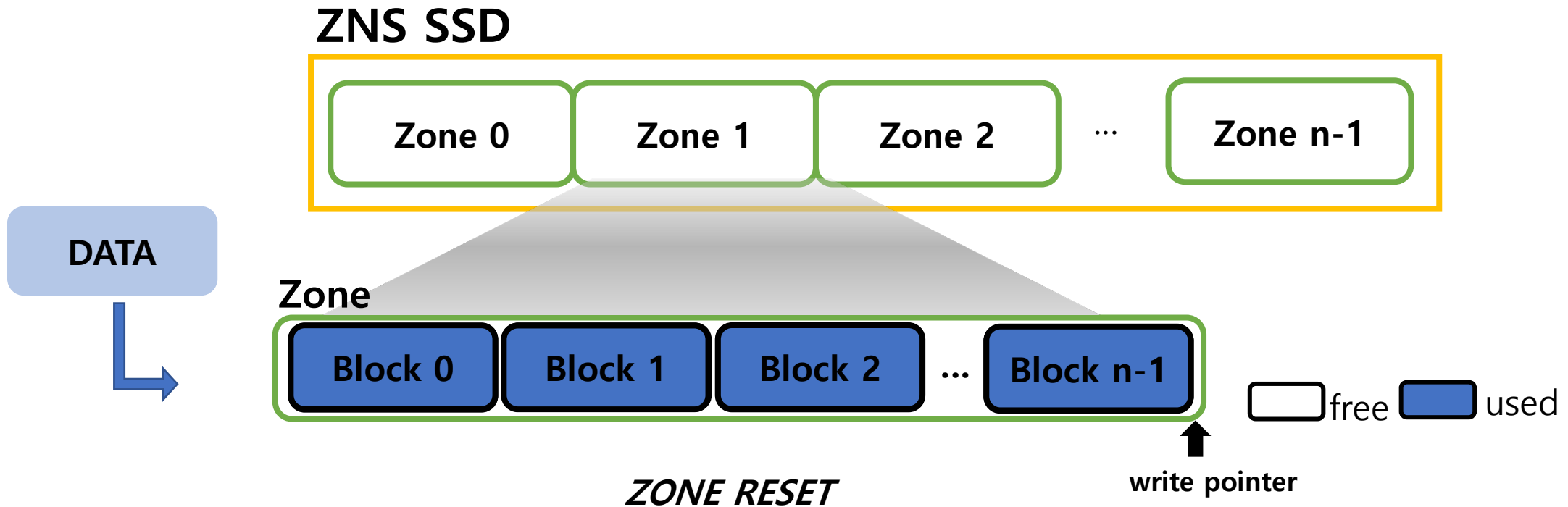
# Zoned Namespace SSD(ZNS SSD)

**ZNS SSD**

| | | | | |
|---|---|---|---|---|
| Zone 0 | Zone 1 | Zone 2 | ... | Zone n-1 |

**Zone**

| | | | |
|---|---|---|---|
| Block 0 | Block 1 | Block 2 | ... Block n-1 |

❑ Group multiple logical blocks into a **zone**

❑ Zone is an erase unit of ZNS SSD.

서강대학교
SOGANG UNIVERSITY

# Zoned Namespace SSD(ZNS SSD)

**ZNS SSD**

| Zone 0 | Zone 1 | Zone 2 | ... | Zone n-1 |

**Zone**

| Block 0 | Block 1 | Block 2 | ... | Block n-1 |

☐ free  ■ used

❑ Zone allows only <u>sequential writes</u>.

서강대학교
SOGANG UNIVERSITY

# Zoned Namespace SSD(ZNS SSD)

**ZNS SSD**

| | Zone 0 | Zone 1 | Zone 2 | ... | Zone n-1 |
|---|---|---|---|---|---|

**DATA**

**Zone**

| Block 0 | Block 1 | Block 2 | ... | Block n-1 |
|---|---|---|---|---|

free   used

↑ **write pointer**

*ZONE RESET*

❏ Zone allows only <u>sequential writes</u>.

❏ Zone disallows overwrite on same logical blocks.

4

서강대학교
SOGANG UNIVERSITY

# Zoned Namespace SSD(ZNS SSD)
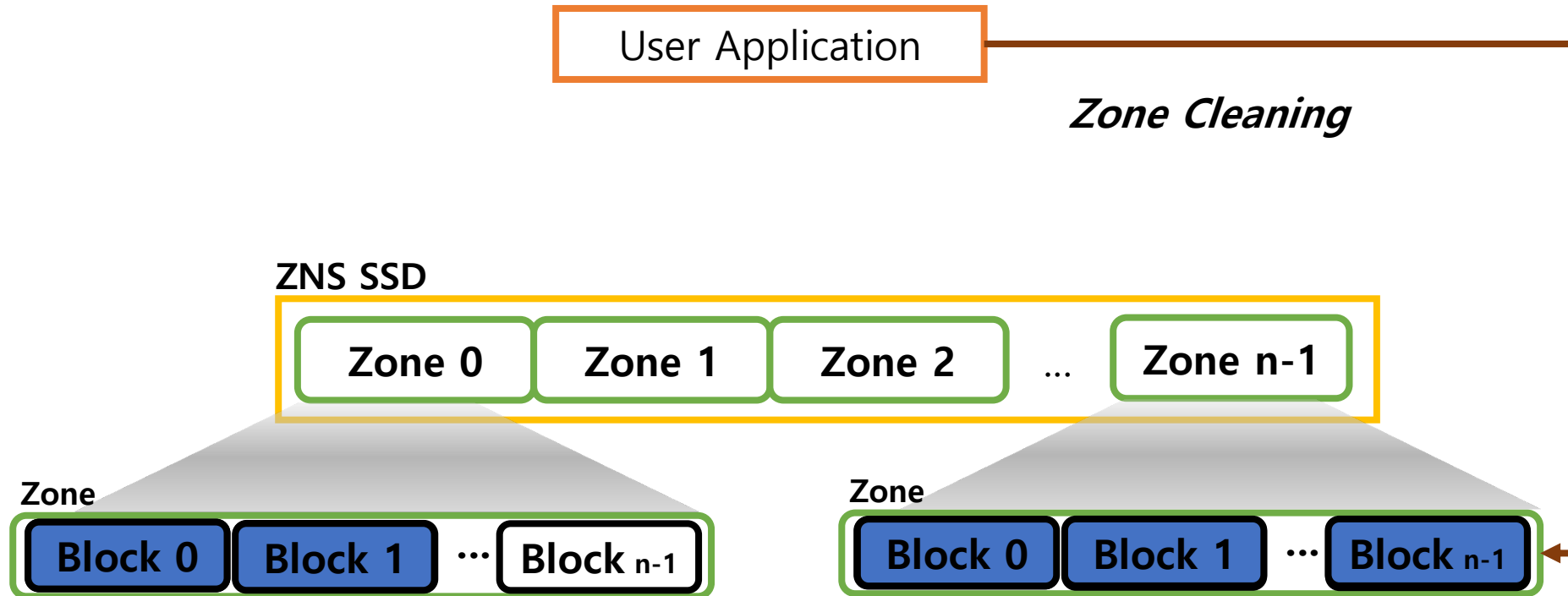
**ZNS SSD**

| Zone 0 | Zone 1 | Zone 2 | ... | Zone n-1 |

# *No More Garbage Collection in FTL*

- ❑ Zone allows only <u>sequential writes</u>.

- ❑ Zone disallows overwrite on same logical blocks.

서강대학교
SOGANG UNIVERSITY

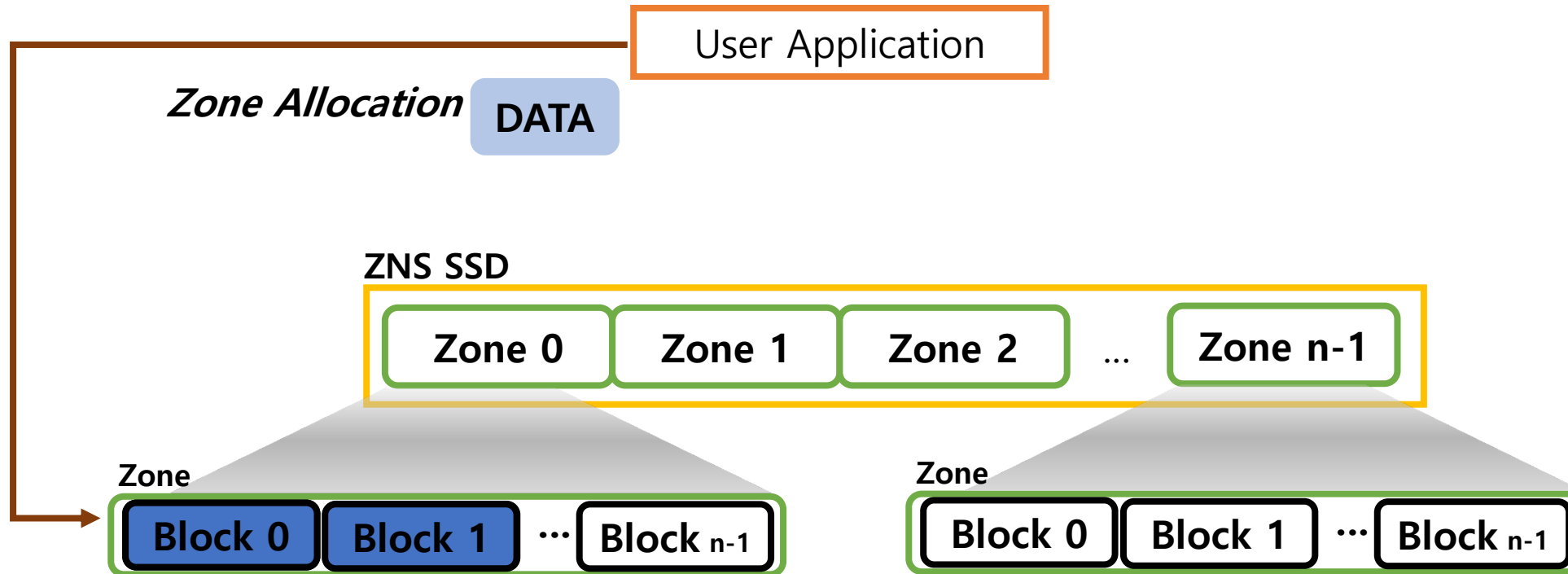# Zoned Namespace SSD(ZNS SSD)

## User applications take over responsibility for

- Free space reclamation (zone cleaning)

- Data placement (zone allocation)

# Zoned Namespace SSD(ZNS SSD)

## User applications take over responsibility for

- Free space reclamation (zone cleaning)

- Data placement (zone allocation)

User Application

*Zone Allocation*   DATA

ZNS SSD

| Zone 0 | Zone 1 | Zone 2 | ... | Zone n-1 |

Zone

| Block 0 | Block 1 | ... | Block n-1 |

Zone

| Block 0 | Block 1 | ... | Block n-1 |

서강대학교
SOGANG UNIVERSITY

# LSM-based Key-Value Store (LSM-KV)

LSM-KV is suitable for ZNS SSD.

❑ Log-structured merge-tree(LSM-tree)
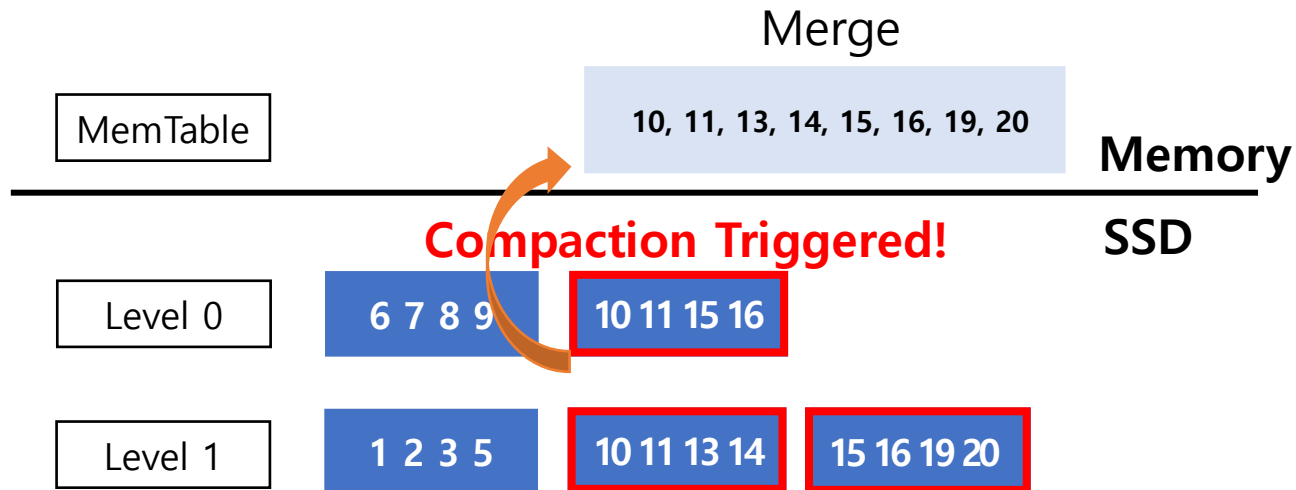
❑ Sequential I/O pattern

❑ Out-of-place update

❑ Application

# LSM-based Key-Value Store (LSM-KV)

## Data Update in LSM-tree

### Compaction
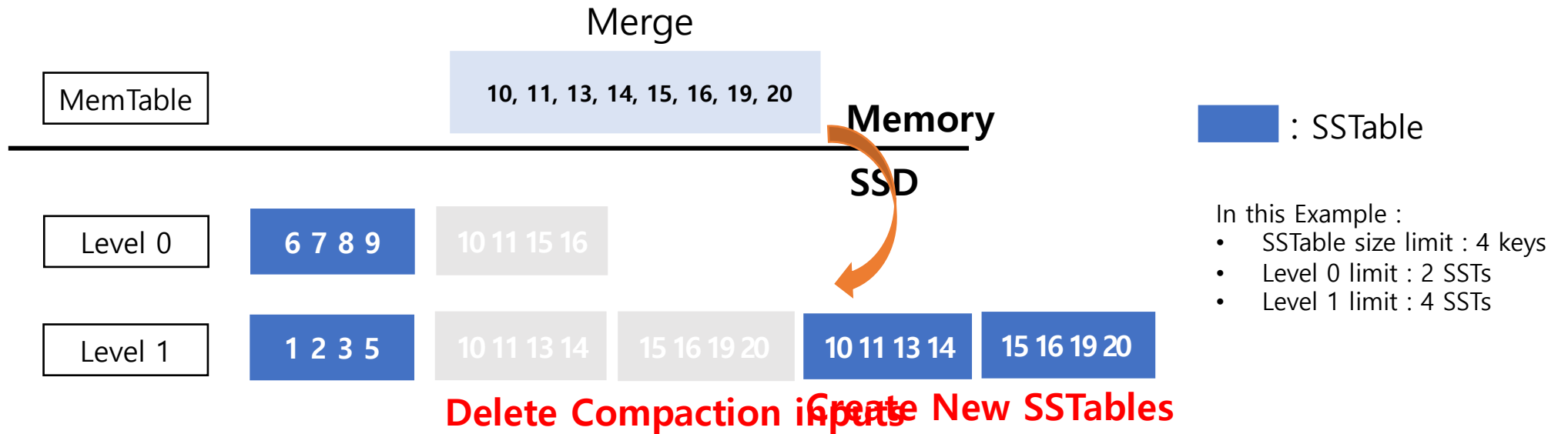
LSM-tree updates data file(SSTable) via compaction.



Merge

MemTable

| 10, 11, 13, 14, 15, 16, 19, 20 |

**Memory**

**SSD**

**Compaction Triggered!**

Level 0 | 6 7 8 9 | 10 11 15 16 |

Level 1 | 1 2 3 5 | 10 11 13 14 | 15 16 19 20 |

■ : SSTable

In this Example :
- SSTable size limit : 4 keys
- Level 0 limit : 2 SSTs
- Level 1 limit : 4 SSTs

서강대학교
SOGANG UNIVERSITY

# LSM-based Key-Value Store (LSM-KV)

## Data Update in LSM-tree

**Compaction**

LSM-tree updates data file(SSTable) via compaction.

Merge

| MemTable | | 10, 11, 13, 14, 15, 16, 19, 20 |

**Memory**

**SSD**

| Level 0 | 6 7 8 9 | 10 11 15 16 |

| Level 1 | 1 2 3 5 | 10 11 13 14 | 15 16 19 20 | 10 11 13 14 | 15 16 19 20 |

**Delete Compaction inputs**  **Create New SSTables**

: SSTable

In this Example :
- SSTable size limit : 4 keys
- Level 0 limit : 2 SSTs
- Level 1 limit : 4 SSTs

서강대학교
SOGANG UNIVERSITY

# ZenFS[1] [ATC'21]

❑ User-level file system for LSM-KV

❑ Backend module for RocksDB

❑ Responsible for

    - Zone Allocation of data in LSM-KV

    - Zone Cleaning for free space reclamation

1) Bjørling, Matias, et al. *ZNS: Avoiding the Block Interface Tax for Flash-based SSDs, ATC'21*

# LifeTime based Zone Allocation in ZenFS

❑ **LifeTime :** Reflect data hotness according to level in LSM-tree

# LifeTime based Zone Allocation in ZenFS

LifeTime - indicator of **<u>when an SSTable is invalidated</u>** in the device



Medium { L0 `0-99`

L1 `0-50` `53-99`

Long { L2 `0-27` `28-52` `53-80` `81-99`

L3

Extreme {

⋮ ⋮

*Hot*

*Cold*

# LifeTime based Zone Allocation in ZenFS

❑ **Zone Allocation using LifeTime**

   - ZenFS allocates SSTables with same LifeTime into the same zone.

*ZenFS fails to precisely predict which*

*SSTables will be deleted together*

zone 2

zone 3

서강대학교
SOGANG UNIVERSITY

# LifeTime based Zone Allocation in ZenFS

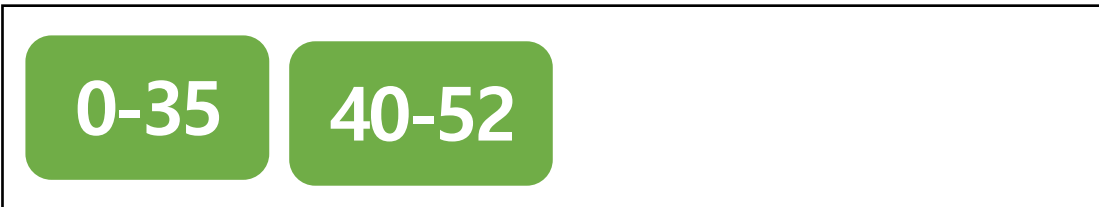**ZenFS fails to precisely predict which SSTables will be deleted together.**

① **Trigger Compaction in L1**  ② **Select Compaction inputs / Merge them**

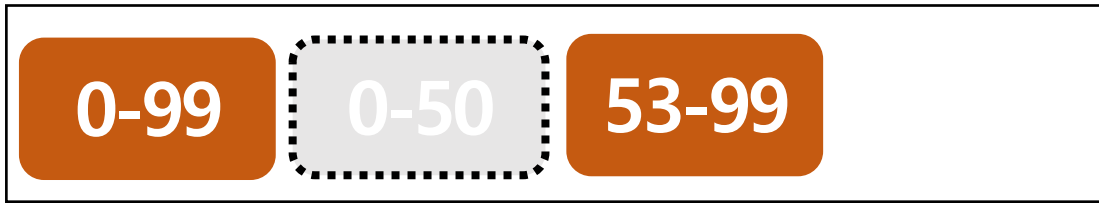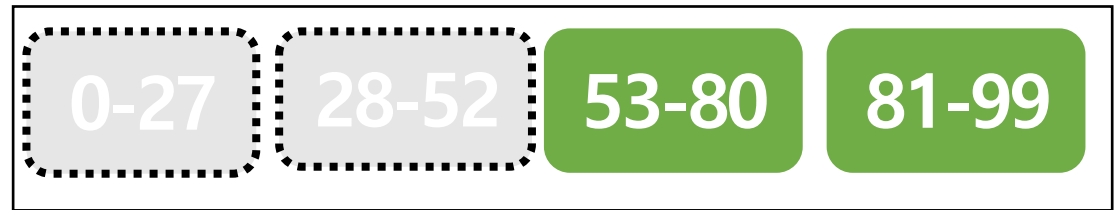③ **Create / Allocate new SSTables**  ④ **Delete Compaction inputs**

# LifeTime based Zone Allocation in ZenFS

**ZenFS fails to precisely predict which SSTables will be deleted together.**

**zone 0 (Medium)**

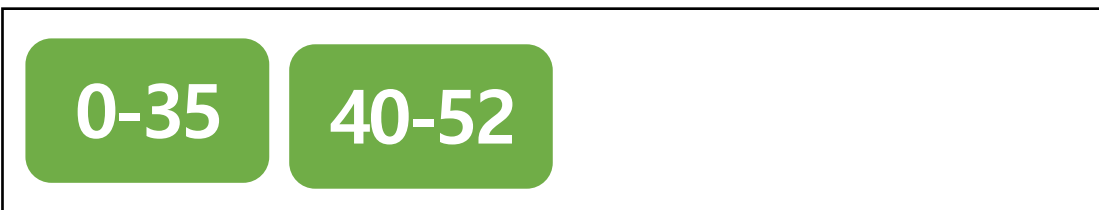| 0-99 | 0-50 | 53-99 |

**zone 1 (Long)**

| 0-27 | 28-52 | 53-80 | 81-99 |

**zone 2 (Long)**

| 0-35 | 40-52 |

**zone 3**

서강대학교
SOGANG UNIVERSITY

# LifeTime based Zone Allocation in ZenFS

**ZenFS fails to precisely predict which SSTables will be deleted together.**

**zone 0 (Medium)**

| 0-99 | 0-50 | 53-99 |
|------|------|-------|

**zone 1 (Long)**

| 0-27 | 28-52 | 53-80 | 81-99 |
|------|-------|-------|-------|

**zone 2 (Long)**

| 0-35 | 40-52 |
|------|-------|

**zone 3**

서강대학교
SOGANG UNIVERSITY

# Overhead of ZenFS Allocation

**Write Amplification**

Prediction failure of deletion time leads to valid data copying during zone cleaning.
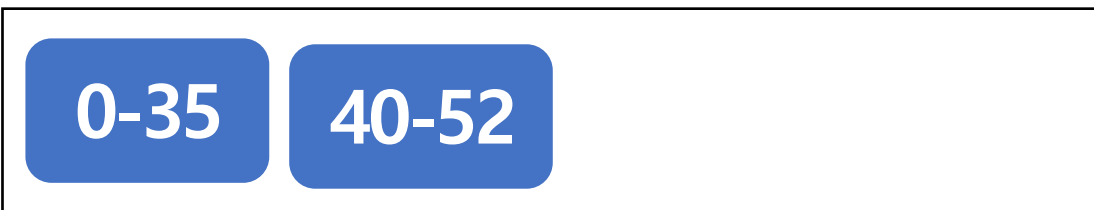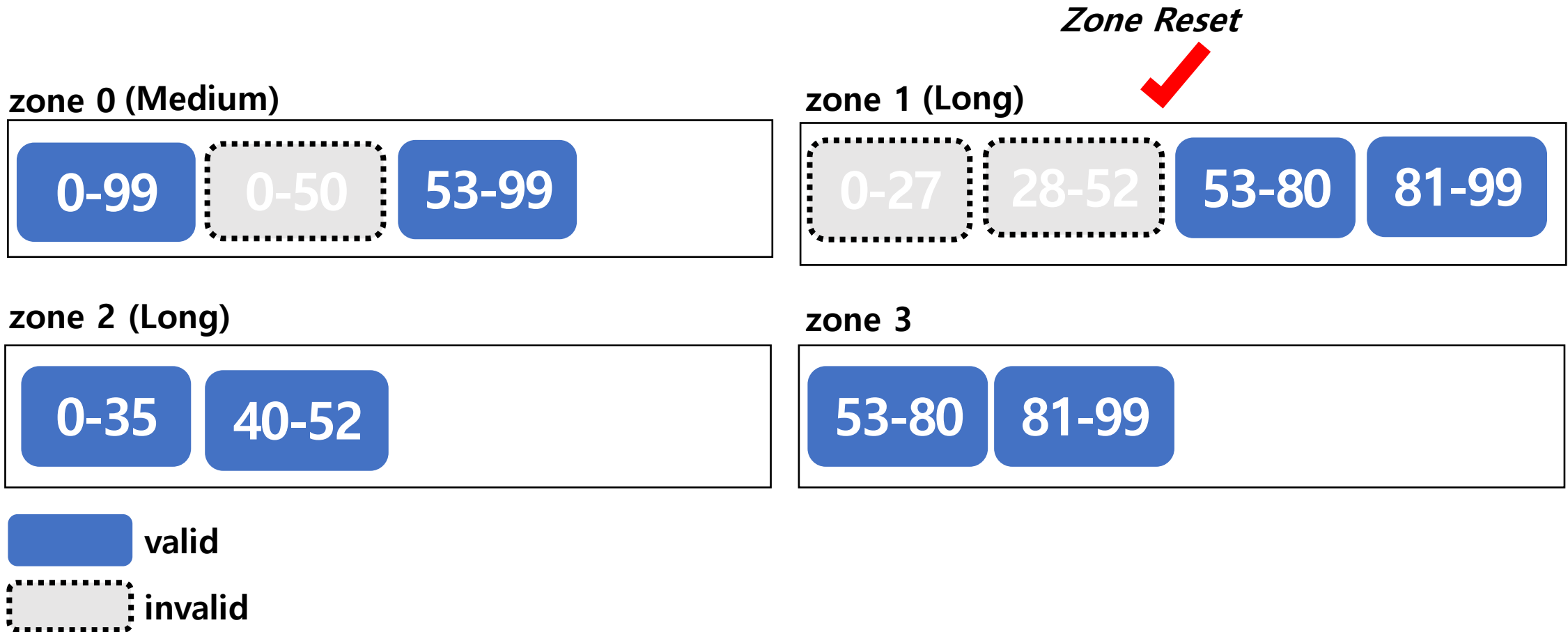


*Zone Cleaning Trigger*

zone 0 (Medium)
| 0-99 | 0-50 | 53-99 |

zone 1 (Long) ✔
| 0-27 | 28-52 | 53-80 | 81-99 |

zone 2 (Long)
| 0-35 | 40-52 |

zone 3

valid
invalid

# Overhead of ZenFS Allocation

**Write Amplification**

Prediction failure of deletion time leads to valid data copying during zone cleaning.

# Compaction-Aware Zone Allocation (CAZA)

## Main Idea

- ❑ We observed that SSTables to be compacted together are invalidated at the same time.

- ❑ We allocate SSTables to be compacted together into the same zone.

*Challenge : How can we predict which SSTables will be compacted together ?*

서강대학교
SOGANG UNIVERSITY

# Compaction-Aware Zone Allocation (CAZA)

**Conditions for SSTables to be compacted together**

*Condition 1: SSTables that are located in <u>adjacent levels</u>*

*Condition 2 : SSTables that have <u>overlapping key ranges</u>*

서강대학교
SOGANG UNIVERSITY

# Compaction-Aware Zone Allocation (CAZA)

## Compaction Input Picking in LSM-tree

**Compaction Triggered!**

A

$level_{(n)}$ ... **1-12** ...

B C D

$level_{(n+1)}$ **2-6** **7-11** **17-20** ...
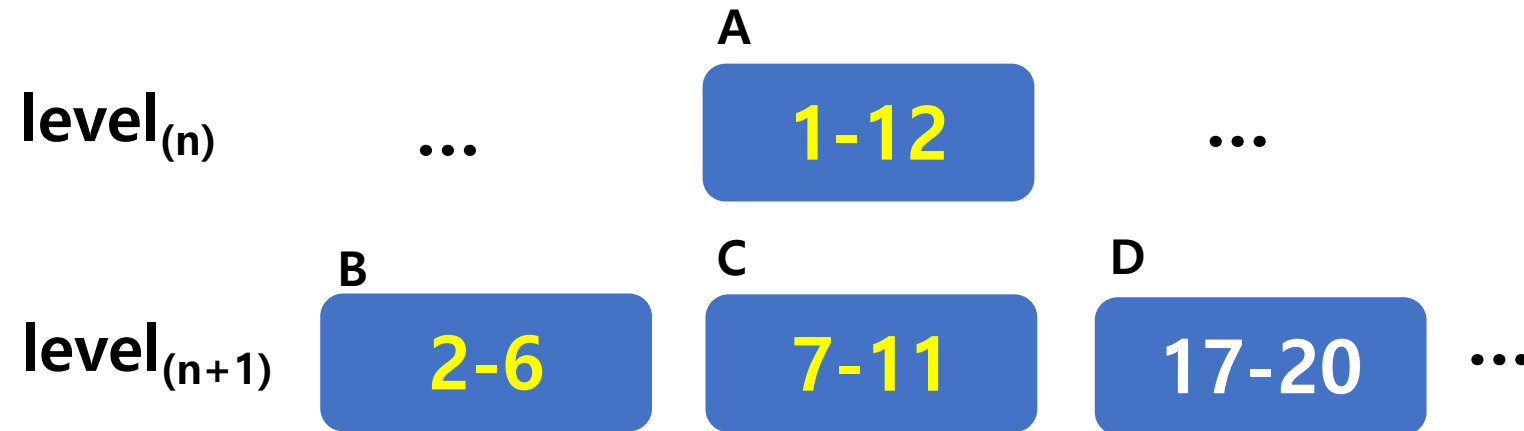
서강대학교
SOGANG UNIVERSITY

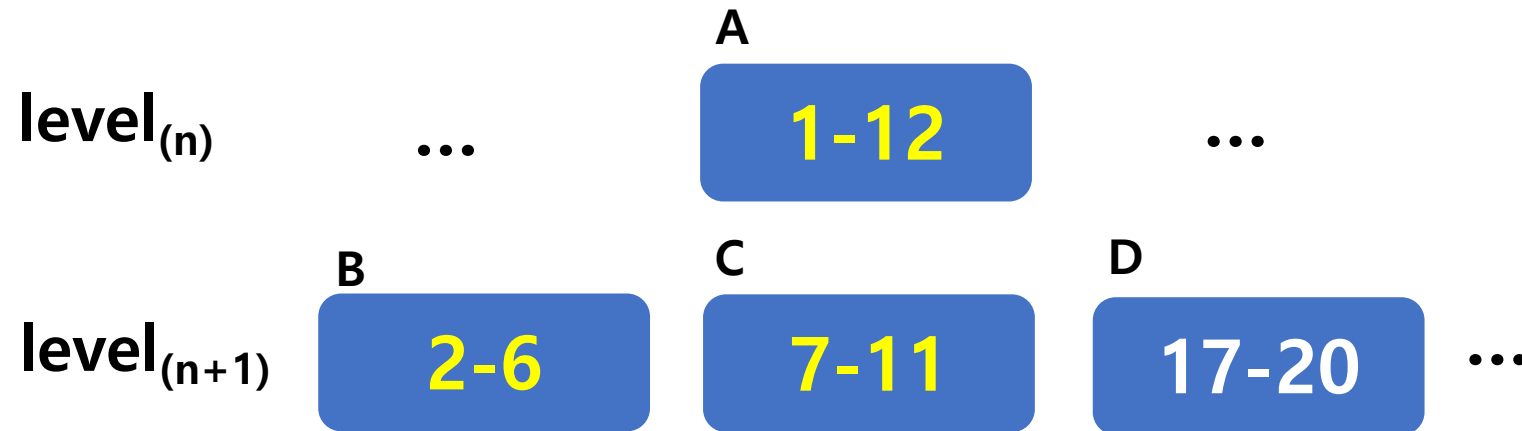# Compaction-Aware Zone Allocation (CAZA)

**Compaction Input Picking in LSM-tree**

# Compaction-Aware Zone Allocation (CAZA)

**Compaction Input Picking in LSM-tree**

# Compaction-Aware Zone Allocation (CAZA)

## Compaction Input Picking in LSM-tree

$level_{(n)}$    ...    **A** **1-12**    ...

$level_{(n+1)}$    **B** **2-6**    **C** **7-11**    **D** **17-20**    ...

서강대학교
SOGANG UNIVERSITY

# Compaction-Aware Zone Allocation (CAZA)



level$_{(n)}$  1-12

**Newly Created**

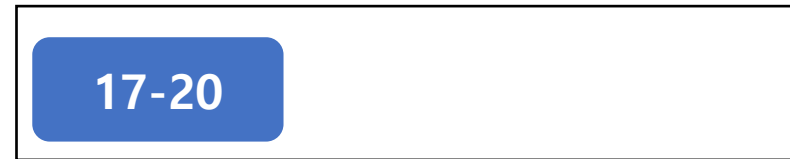level$_{(n+1)}$  2-6  7-11  17-20

zone$_k$
| 2 - 6 | 7-11 | 1-12 |

zone$_{k+1}$
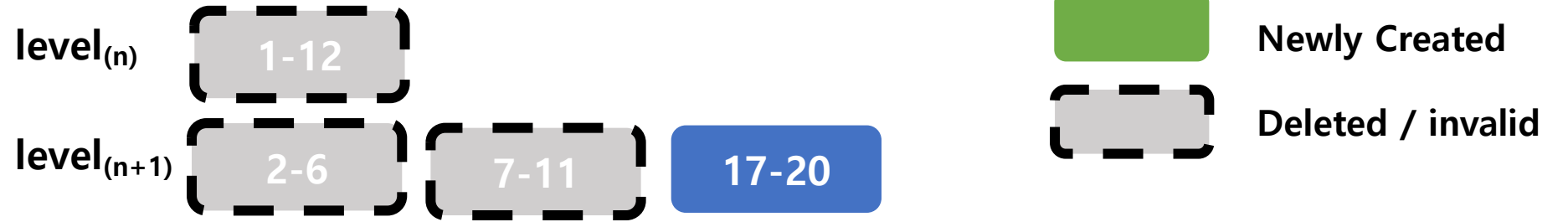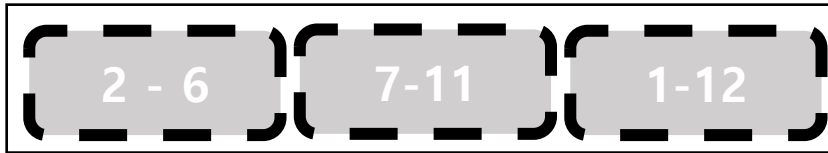| 17-20 |

# Compaction-Aware Zone Allocation (CAZA)
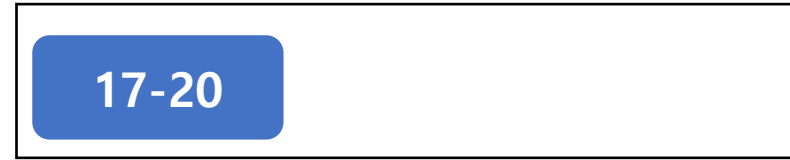
# Compaction-Aware Zone Allocation (CAZA)

# Compaction-Aware Zone Allocation (CAZA)

❑ Compaction-Aware Zone Allocation Rules

**Case 1)** SSTables with overlapping key ranges spread in several zones.

**Allocate the zone with more SSTables**

$level_{(n-1)}$ | 1-12

$level_{(n)}$ | 2-6 | 7-11 | 12-20

**$zone_k$**
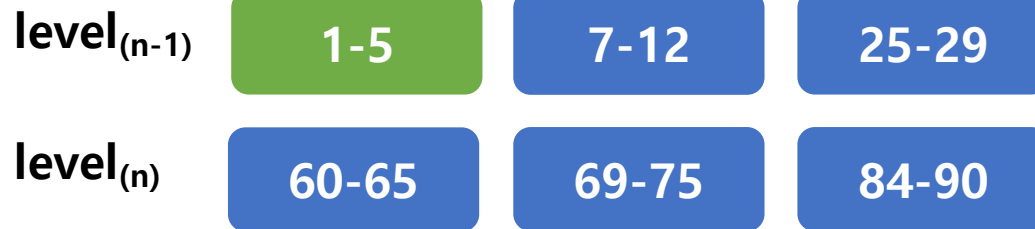
| 2 - 6 | 7-11 | 1-12 |

**$zone_{k+1}$**

| 12-20 |

# Compaction-Aware Zone Allocation (CAZA)

❑ Compaction-Aware Zone Allocation Rules

**Case 2)** No matching SSTables that meet the compaction condition

*Allocate Empty zone*

**level$_{(n-1)}$** | 1-5 | 7-12 | 25-29

**level$_{(n)}$** | 60-65 | 69-75 | 84-90

**zone$_k$**

7 - 12

**zone$_{k+1}$**

25-29

**zone$_{k+2}$**

1-15

**zone$_{k+3}$**
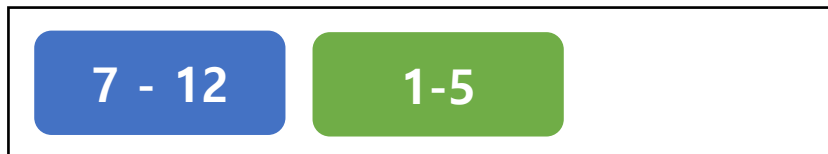
# Compaction-Aware Zone Allocation (CAZA)

❑ Compaction-Aware Zone Allocation Rules

**Case 2)** No matching SSTables that meet the compaction condition

*Allocate Zone with SSTables in*
*(1) Same level (2) Closest key range*

| | | | |
|---|---|---|---|
| **level$_{(n-1)}$** | 1-5 | 7-12 | 25-29 |
| **level$_{(n)}$** | 60-65 | 69-75 | 84-90 |

**zone$_k$**

| 7 - 12 | 1-5 |
|---|---|

**zone$_{k+1}$**

**zone$_{k+1}$**

**zone$_{k+3}$**

서강대학교
SOGANG UNIVERSITY
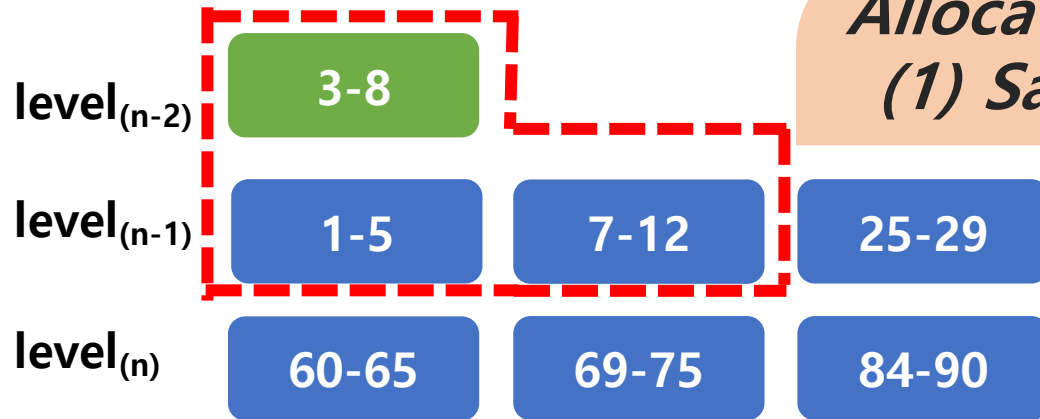
# Compaction-Aware Zone Allocation (CAZA)

❑ Compaction-Aware Zone Allocation Rules

**Case 2)** No matching SSTables that meet the compaction condition

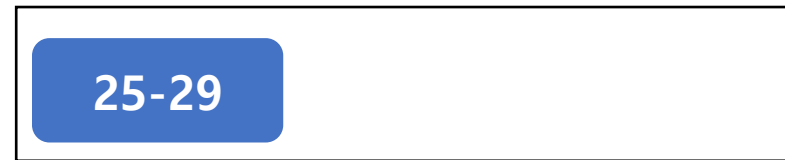*Allocate Zone with SSTables in (1) Same level (2) Closest key range*

level$_{(n-2)}$    **3-8**

level$_{(n-1)}$    **1-5**    **7-12**    **25-29**

level$_{(n)}$    **60-65**    **69-75**    **84-90**

**zone$_k$**

**7 - 12**    **1-5**    **3-8**

**zone$_{k+1}$**

**zone$_{k+1}$**    **25-29**

**zone$_{k+3}$**

서강대학교
SOGANG UNIVERSITY

# Compaction-Aware Zone Allocation (CAZA)

❑ Compaction-Aware Zone Allocation Rules

No matching case..

➔ Follow Zone Allocation in ZenFS

# Evaluation

## ❏ Comparison

- **LIZA** : LifeTime-based Zone Allocation in ZenFS

- **CAZA** : Compaction-Aware Zone Allocation

## ❏ Workload

- RocksDB micro-benchmark
    (*overwrite after fillseq*)
- 40GB KVs, 16B key & 128B Value

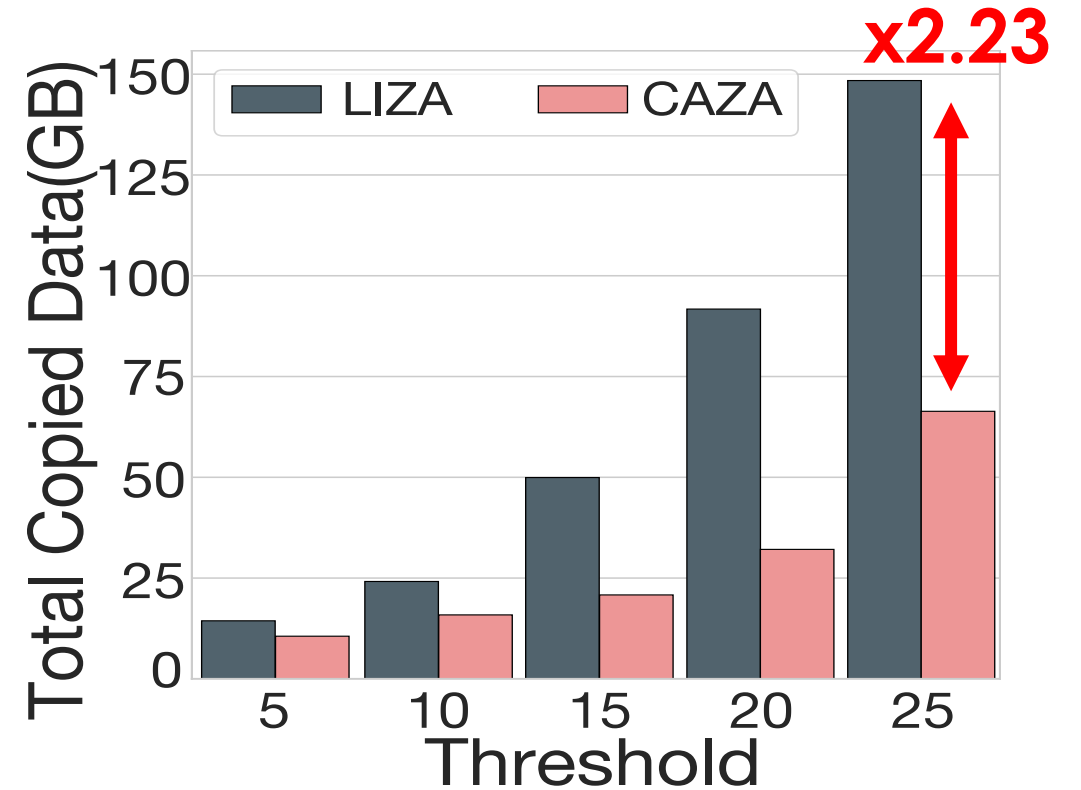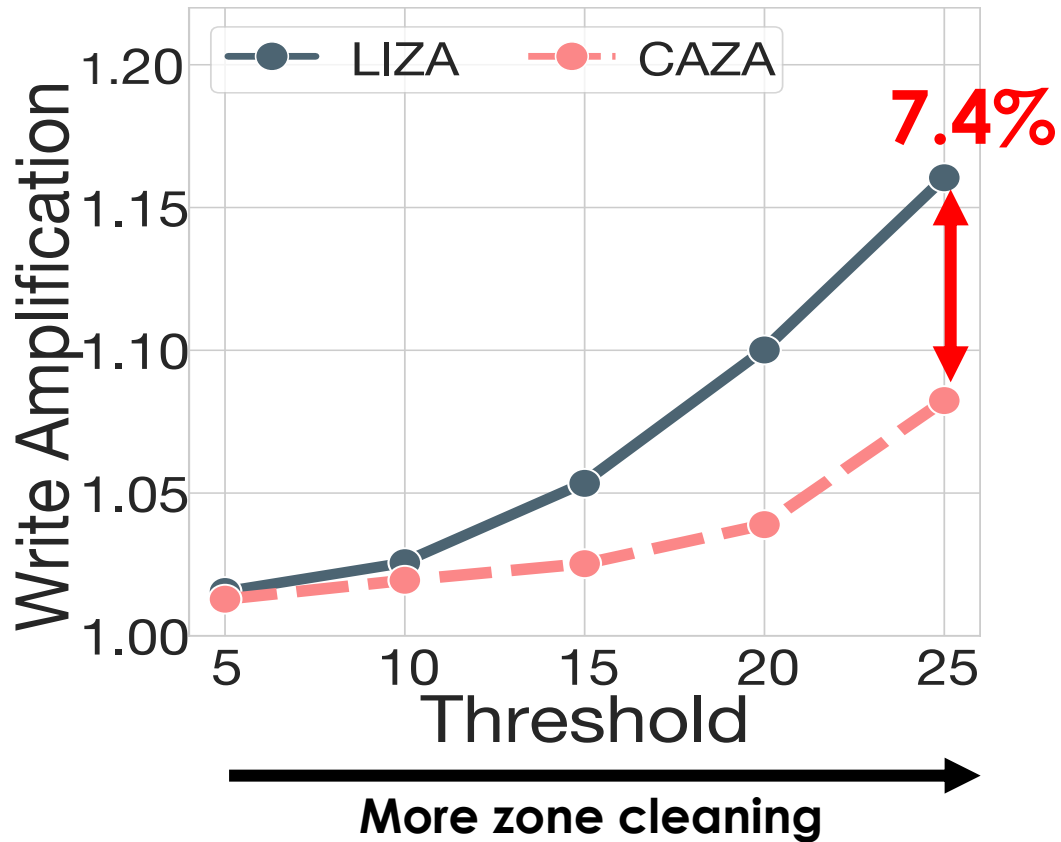| Testbed Specifications | |
|---|---|
| *ZNS SSD* | 100GB Capacity DRAM Emulation 1GB-sized zone (total 100 zones) |
| CPU | Intel Xeon E5-2640 |
| LSM-KV | RocksDB v6.13 |
| OS | Linux Kernel 5.10.13 |

## ❏ Zone Cleaning

- **Greedy Zone Cleaning** : Select zones with most invalid data for reset

- Varying threshold(5%, 10%, 15%, 20%, 25%) where zone cleaning must reclaim for free space

서강대학교
SOGANG UNIVERSITY

# Write Amplification(WA)

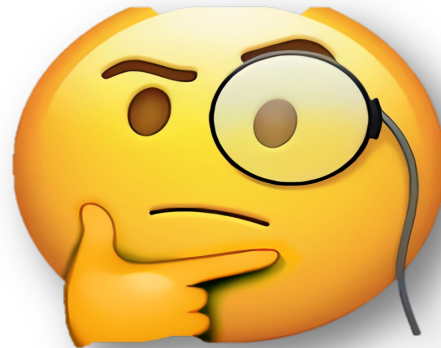$$WA = \frac{Written\ by\ LSM-KV +\ Copied\ by\ Zone\ cleaing}{Written\ by\ LSM-KV}$$

# Conclusion

❑ Compaction-Aware Zone Allocation(CAZA)

- We proposed novel data placement algorithm for LSM-KV on ZNS SSDs.

- CAZA precisely estimates lifetime of SSTables.

- CAZA offers 7.4% lower write amplification and 2x lower data copying during zone cleaning than LIZA.

서강대학교
SOGANG UNIVERSITY

# Thank you!

# Any Questions?

Hee-Rock : heeock@sogang.ac.kr / heerock1.lee@gmail.com

서강대학교
SOGANG UNIVERSITY