

# A New LSM-style Garbage Collection Scheme for ZNS SSDs

*Gunhee Choi, Kwanghee Lee, Myunghoon Oh and Jongmoo Choi*

*from Dankook University*

*Jhuyeong Jhin and Yongseok Oh*

*from SK hynix*

*12<sup>th</sup> USENIX Workshop on Hot Topic in Storage and File System (HotStorage 20), 2020*

2020. 07. 13

Presentation by Choi, Gunhee

choi\_gunhee@dankook.ac.kr

# Content

1. Zoned Namespace SSD
2. Motivation
3. LSM-ZGC Design
4. Evaluation
5. Conclusions
  - Appendix



- **What are Next Generation SSDs?**



Traditional SSD

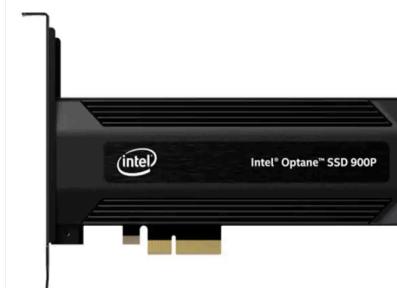
VS.



Open-Channel SSD



Key-Value SSD

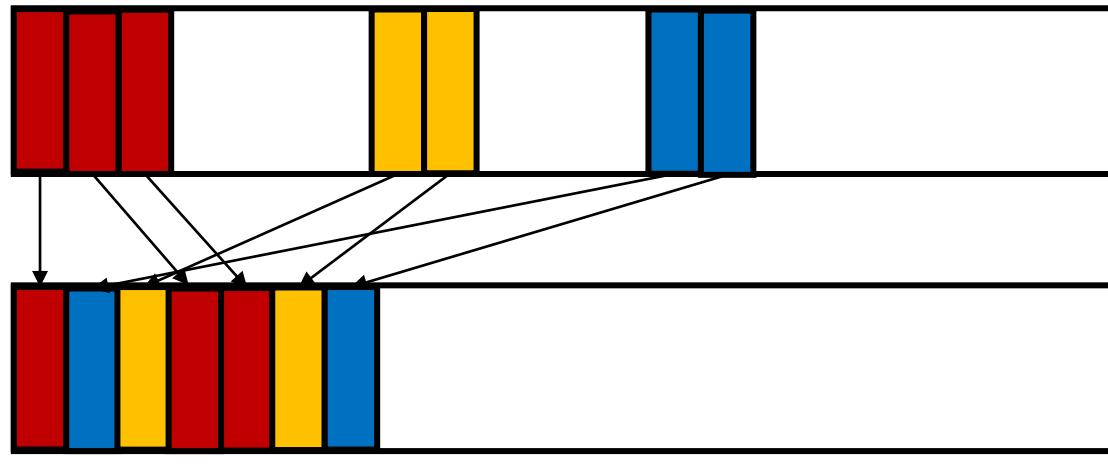


Optane SSD

- **What is ZNS SSD?**

## Traditional SSD

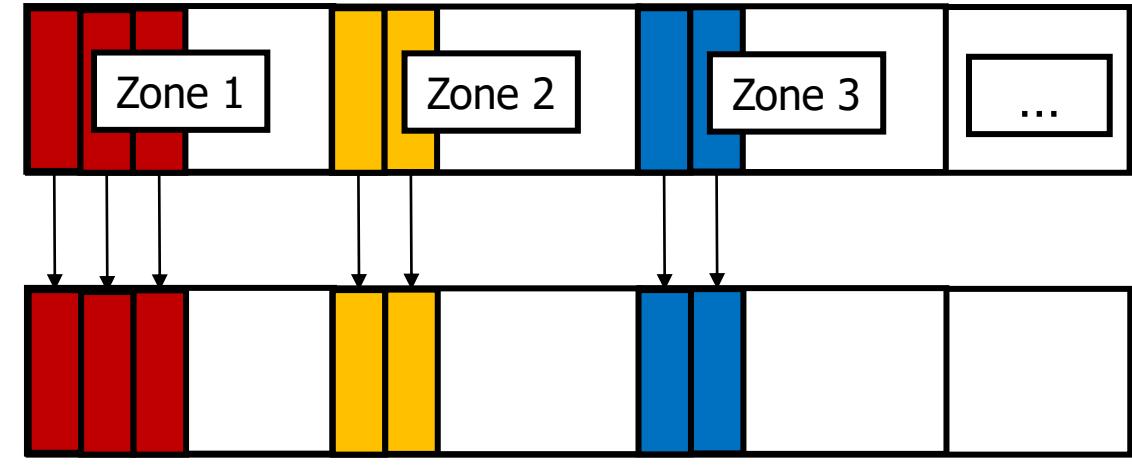
LBA space



NAND

## Zoned Namespace SSD

LBA space



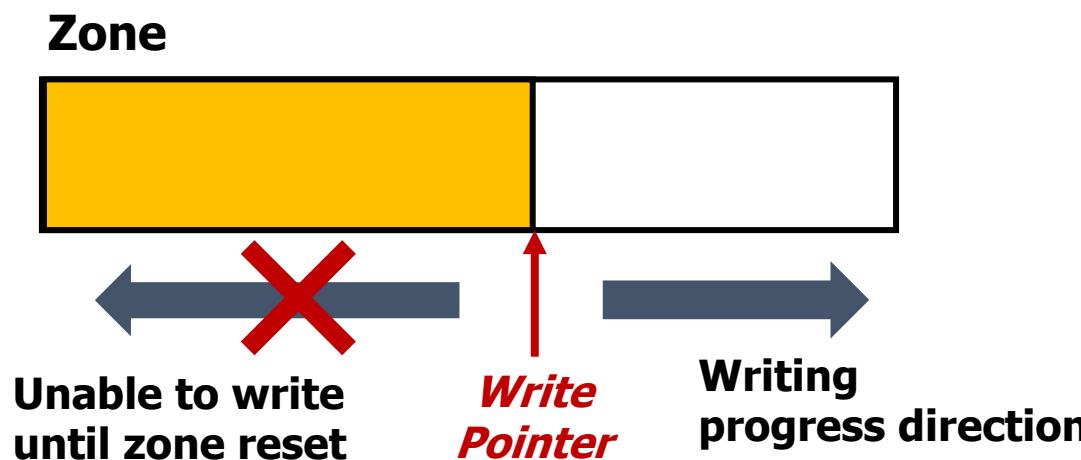
NAND

## ✓ Benefits

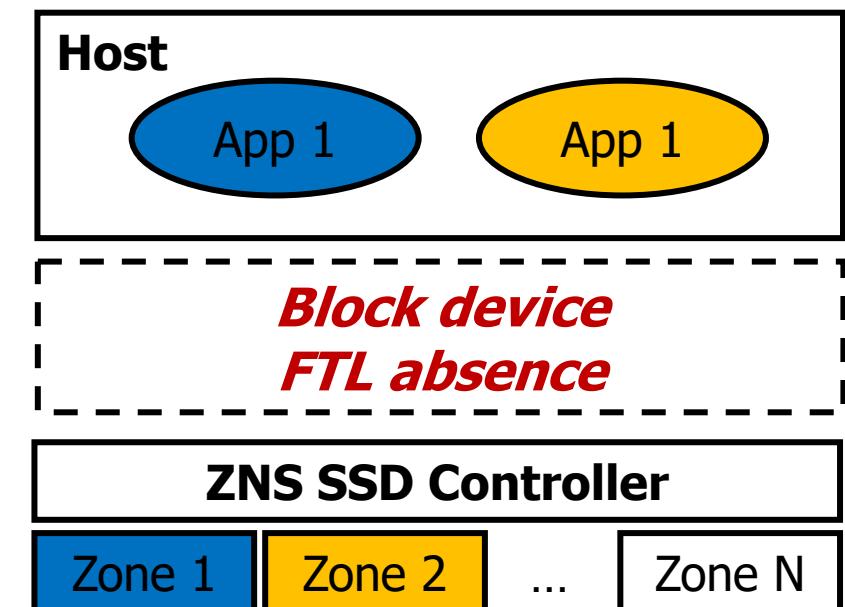
- Better performance and WAF by distributing different workloads into different zones
- Better isolation (IO Determinism)
- Reduce DRAM usage and Over-provisioning area in SSDs

- **What are the issues of ZNS SSD?**

- ✓ Sequential write constraint: writes need to be conducted in a sequential manner, like the SMR drives.
- ✓ Host needs to control zones directly such as zone open, close, reset and **zone garbage collection**.

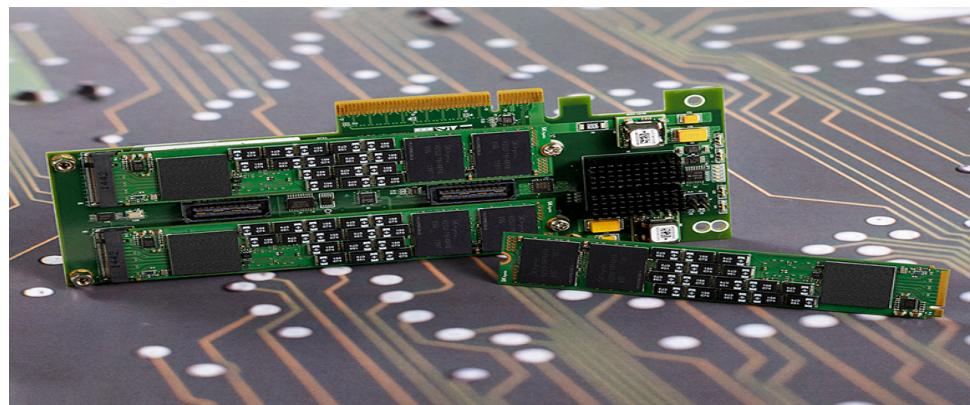


**Sequential write constraint**



**Host Needs to handle zone controls**

- **How much is the Zone Garbage Collection (hereafter ZGC) overhead?**
  - ✓ Using real ZNS SSD prototype
  - ✓ Zone size: **1GB** (note that the typical segment size in LFS is **2MB**)



***SK Hynix Prototype  
ZNS SSD***

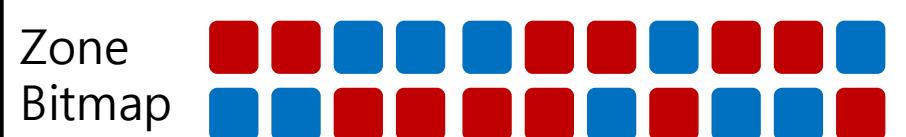
Table 1: ZNS SSD prototype information

Item	Specification
SSD Capacity	1TB
Size of a Zone	1GB
Number of Zones	1024
Interface	PCIe Gen3
Protocol	NVMe 1.2.1

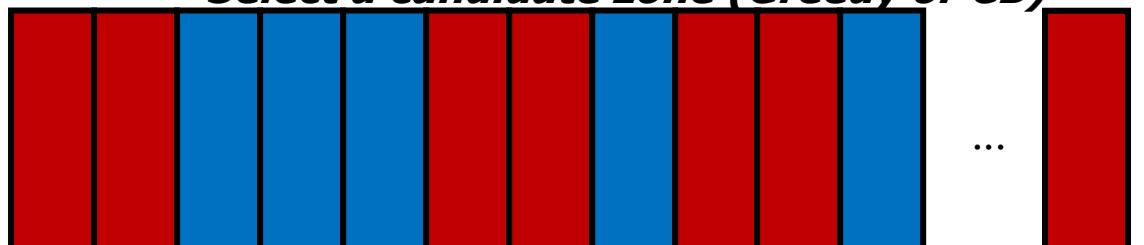
Ref : <https://news.skhynix.co.kr/1915>

## Basic Zone Garbage Collection (Basic\_ZGC)

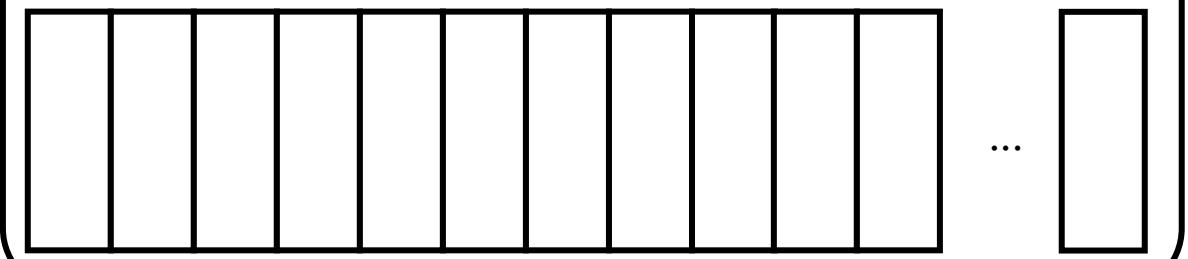
Memory



Zone 0    *Step 1.*  
*Select a candidate zone (Greedy or CB)*

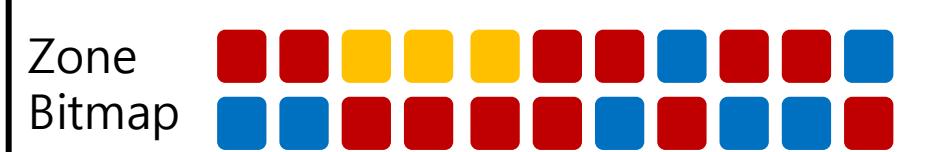


Zone 100



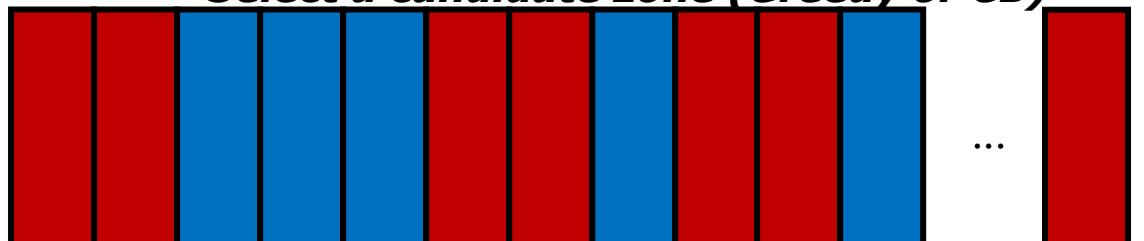
## Basic Zone Garbage Collection

Memory

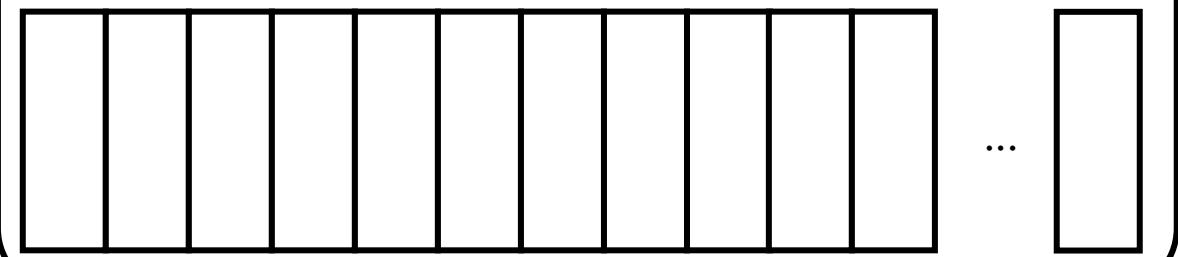


**Step 2.**  
*Find out valid blocks  
using a zone bitmap*

Zone 0    **Step 1.**  
*Select a candidate zone (Greedy or CB)*

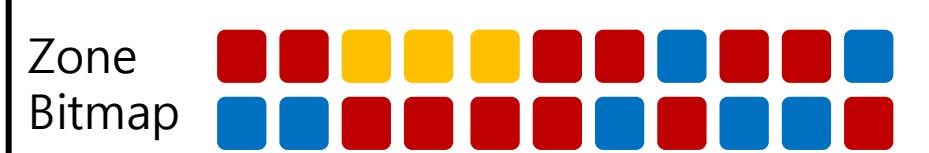


Zone 100

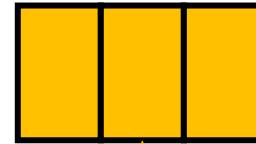


## Basic Zone Garbage Collection

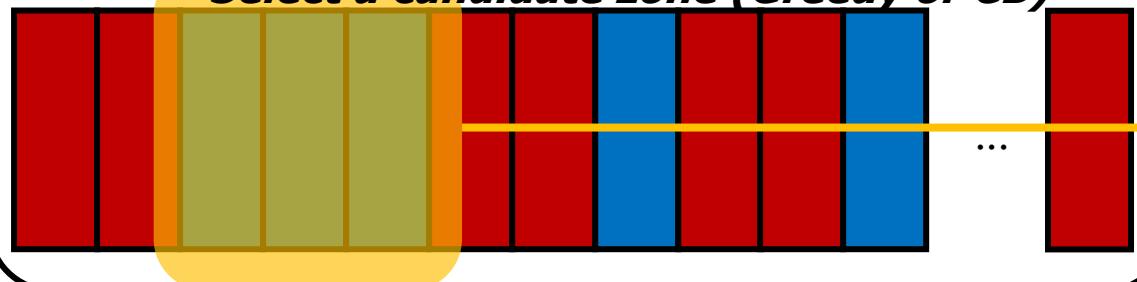
Memory



***Step 2.  
Find out valid blocks  
using a zone bitmap***

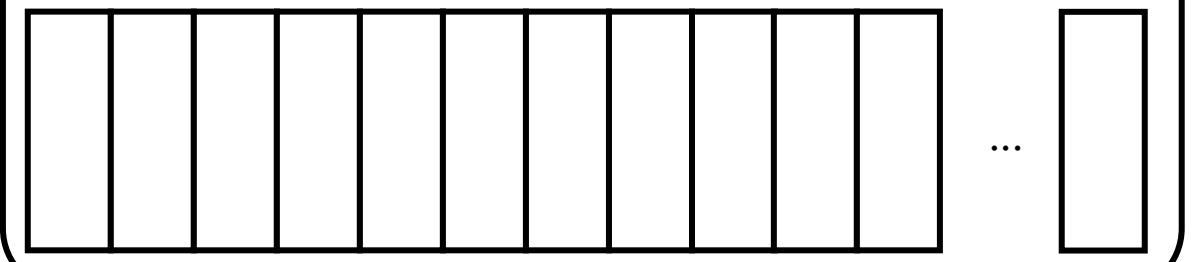


***Step 1.  
Select a candidate zone (Greedy or CB)***



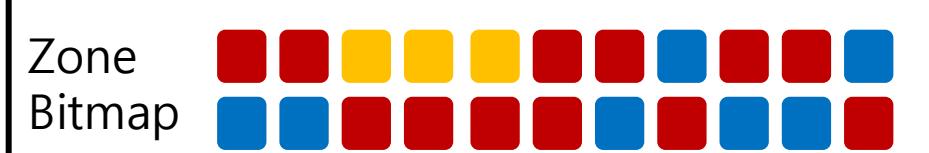
***Step 3.  
Read valid data in 4KB (or larger) I/O size***

Zone 100

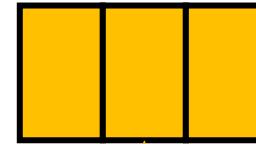


## Basic Zone Garbage Collection

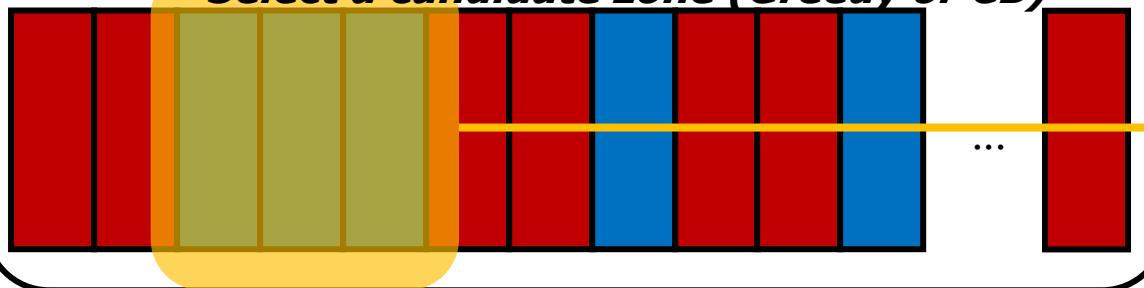
Memory



**Step 2.**  
*Find out valid blocks  
using a zone bitmap*



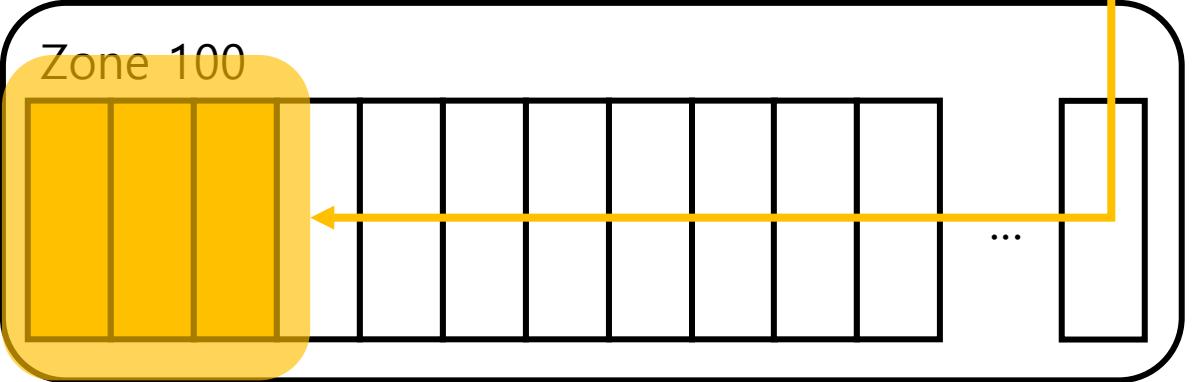
Zone 0  
**Step 1.**  
*Select a candidate zone (Greedy or CB)*



**Step 3.**  
*Read valid data in 4KB (or larger) I/O size*

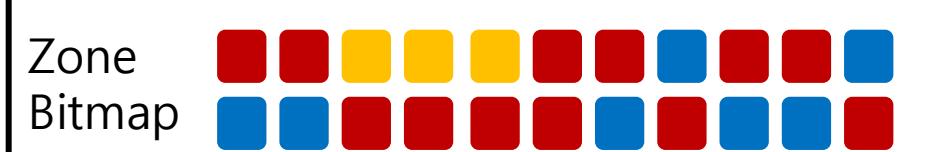
**Step 4.**  
*Write data in 4KB (or larger) I/O size*

Zone 100



## Basic Zone Garbage Collection

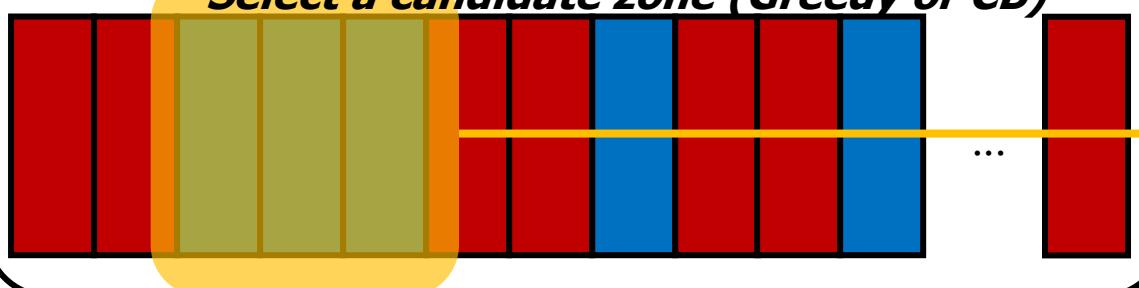
Memory



**Step 2.**  
*Find out valid blocks  
using a zone bitmap*



**Step 1.**  
*Select a candidate zone (Greedy or CB)*

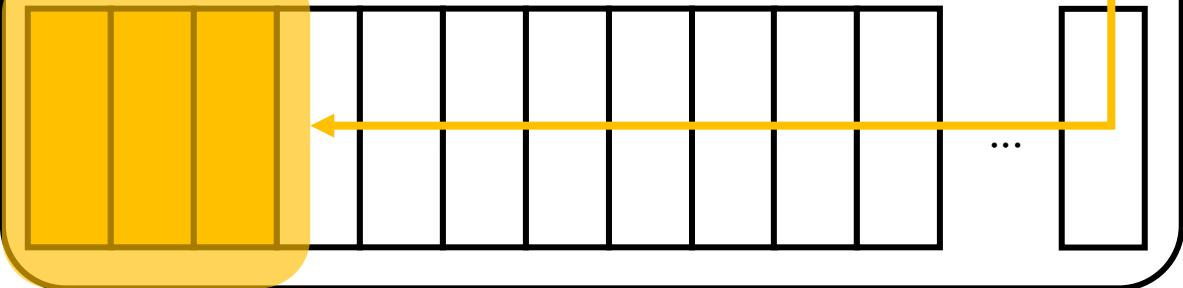


**Step 3.**  
*Read valid data in 4KB (or larger) I/O size*

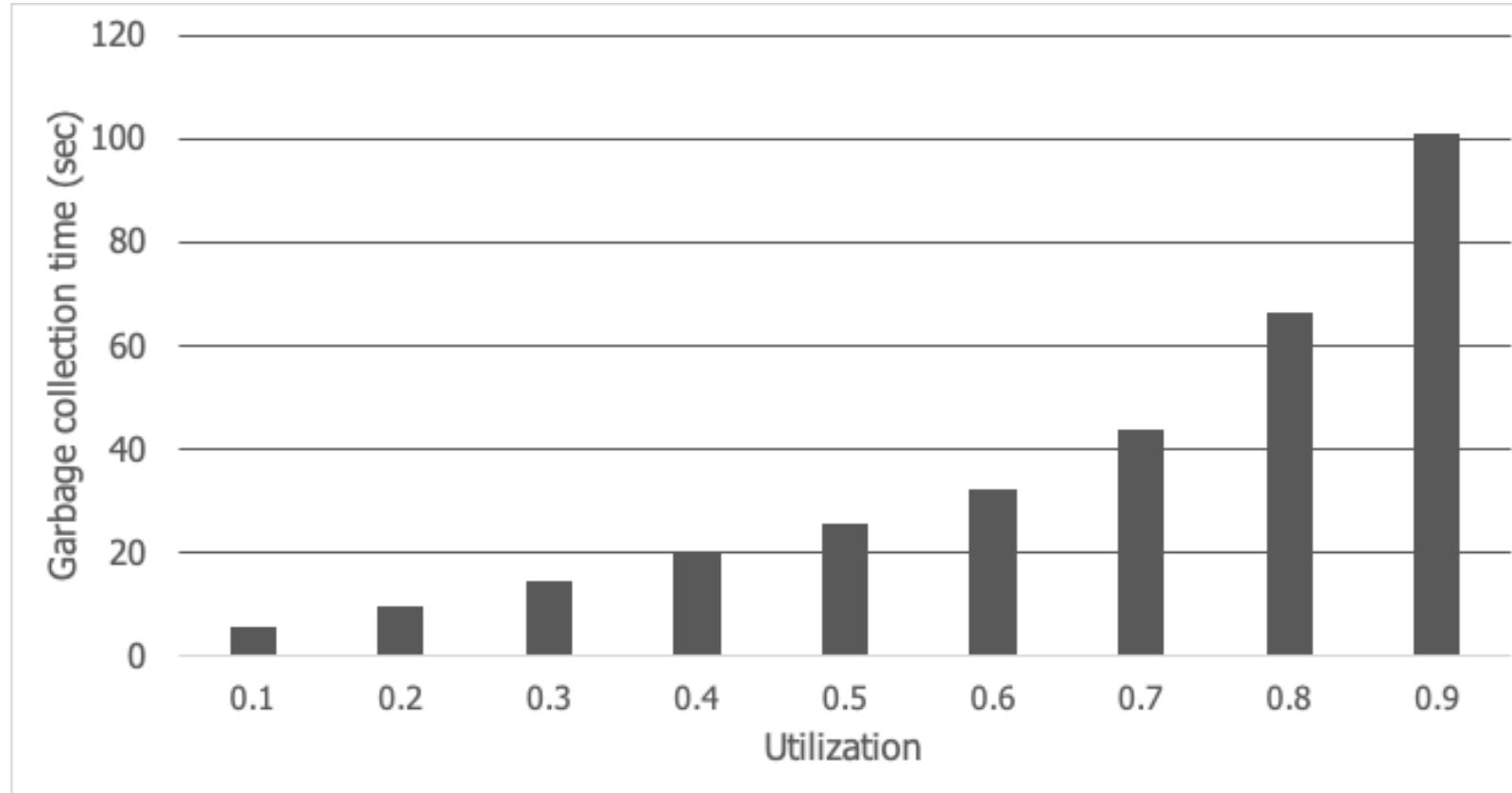
**Step 4.**  
*Write data in 4KB (or larger) I/O size*

**Step 5.**  
*Reset the selected zone*

Zone 100

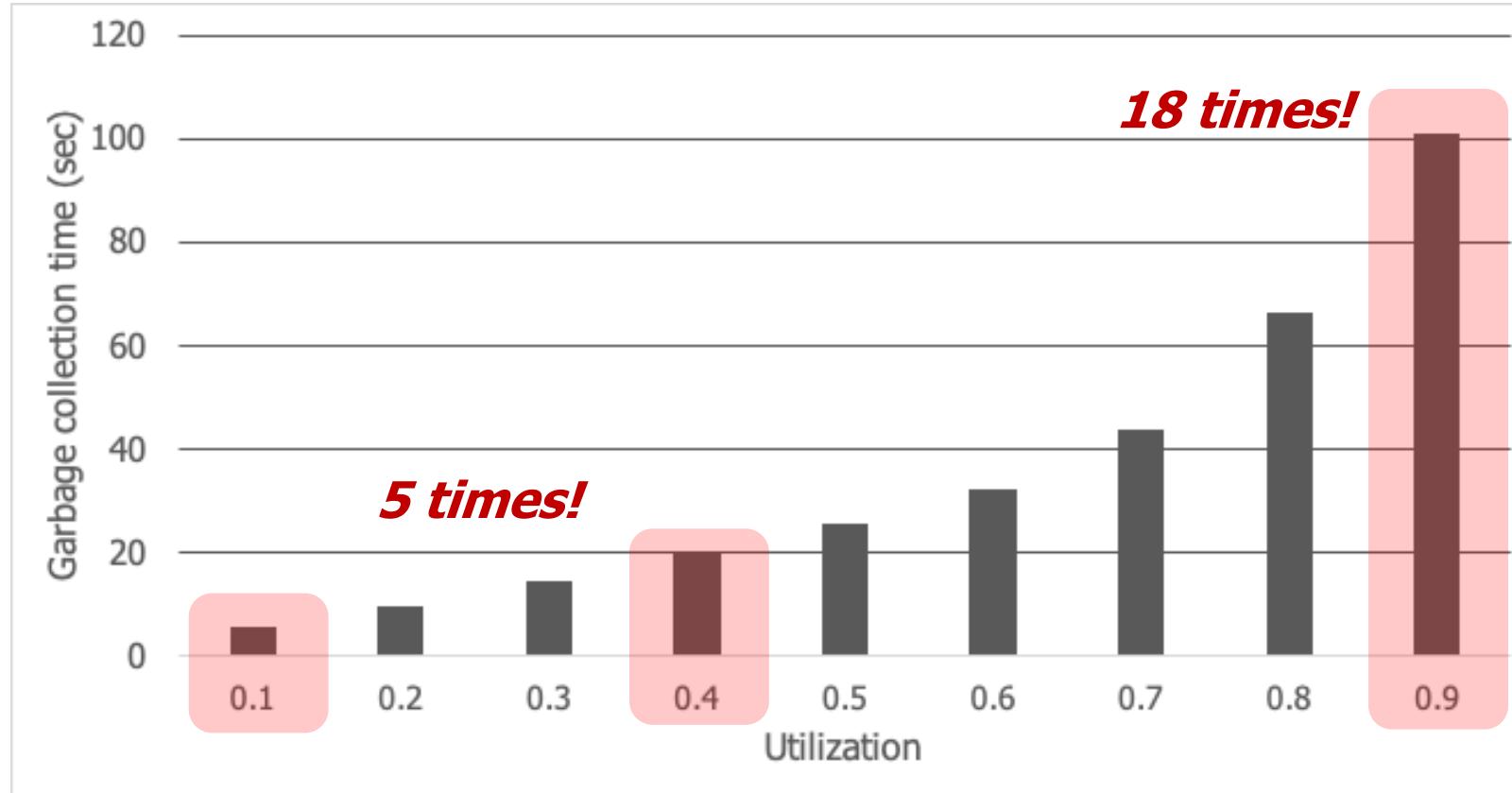


### Observation 1: Zone garbage collection overhead



- Zone : 1GB
- Block : 4KB

### Observation 1: Zone garbage collection overhead



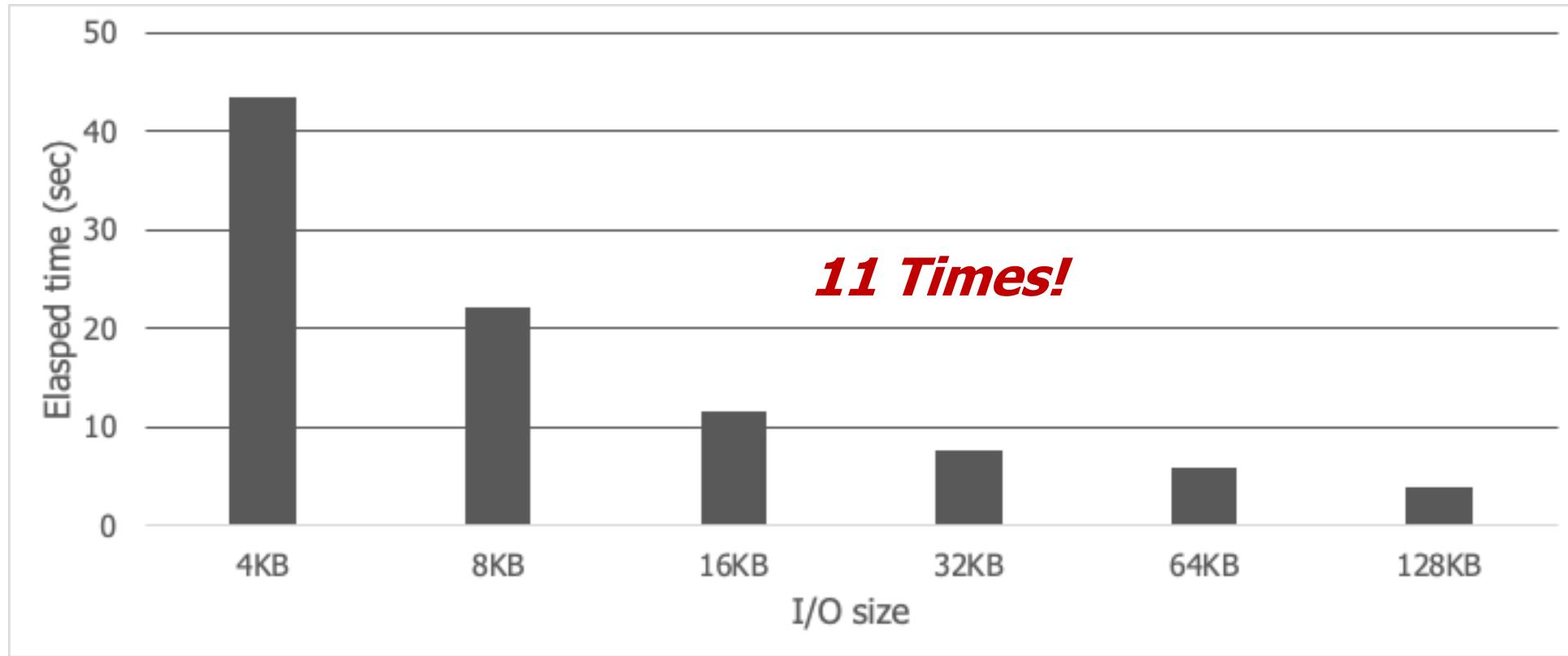
- Zone : 1GB
- Block : 4KB

☞ **Motivation 1: reducing utilization of a candidate zone is indispensable**

### Observation 2: I/O size for Read/Write

- **Another feature of ZNS SSD**
  - ✓ A zone is, in general, mapped into multiple channels/ways.
- **Then, how about read/write data in a larger I/O size (e.g. 128KB)?**



**Observation 2: I/O size for Read/Write**

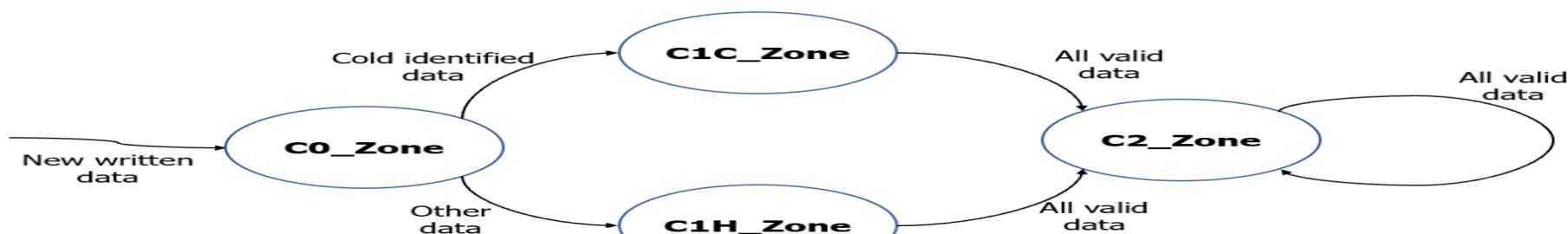
☞ **Motivation 2: accessing in a larger I/O size is beneficial in ZNS SSDs**

*So, Our ideas are*

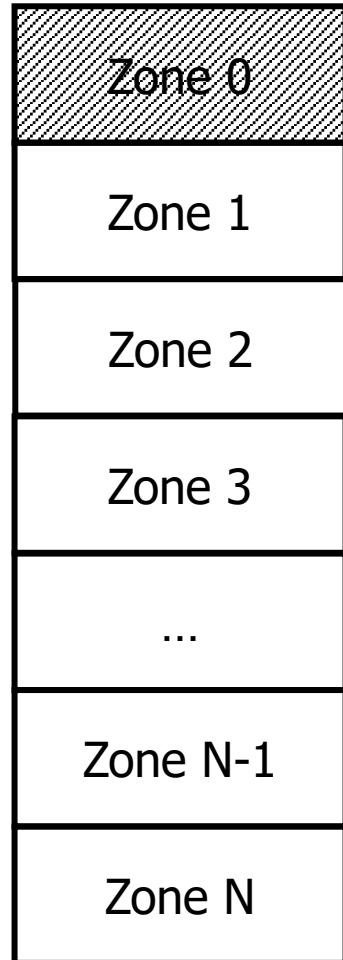
- 1) *Make the utilization of a candidate zone low***
- 2) *Access data in a larger I/O size***

- **How to access data in a larger I/O size?**
  - ✓ The coexistence of valid and invalid data makes it difficult
  - ✓ Read not only valid but also invalid data in a larger I/O size
- **How to make the utilization of a candidate zone low?**
  - ✓ Traditional hot/cold separation is not applicable in ZNS SSDs since zone is quite big
  - ✓ Employ the segment concept for finer-grained hot/cold separation

- Two management units
  - ✓ Zone: for garbage collection vs. Segment: for hot/cold separation
  - ✓ A zone is divided into multiple segments (1GB vs. 2MB in this study)
- Segment state and transition rule (refer to our paper for details)
  - ✓ New data → C0
  - ✓ During ZGC, survived data from C0
    - Data in a high utilized segment ( $> \text{threshold}_{\text{cold}}$ ): cold → C1C
    - Others: hot (or unknown) → C1H
    - Reasoning: spatial locality, also observed in previous studies such as F2FS (FAST'15), Multi-stream (FAST'19), Key-range locality (FAST'20)
  - ✓ During ZGC, survived data from C1C or C1H (second survived data) → C2

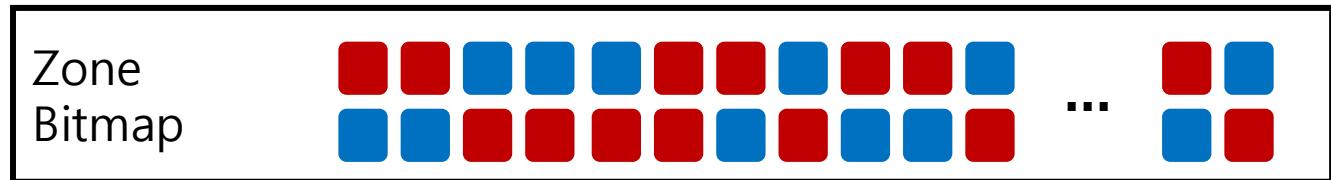


#### LSM(Log Structured Merge) Zone GC



*C0\_zone*

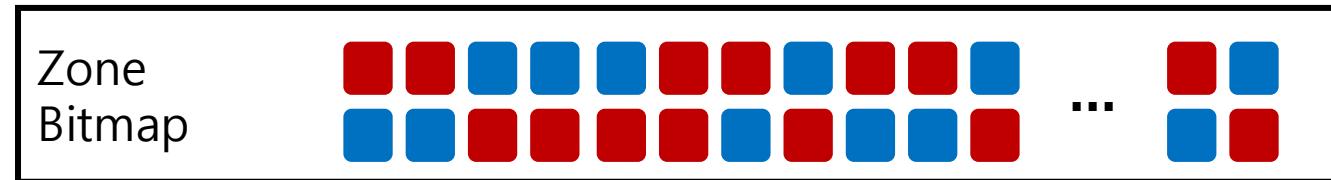
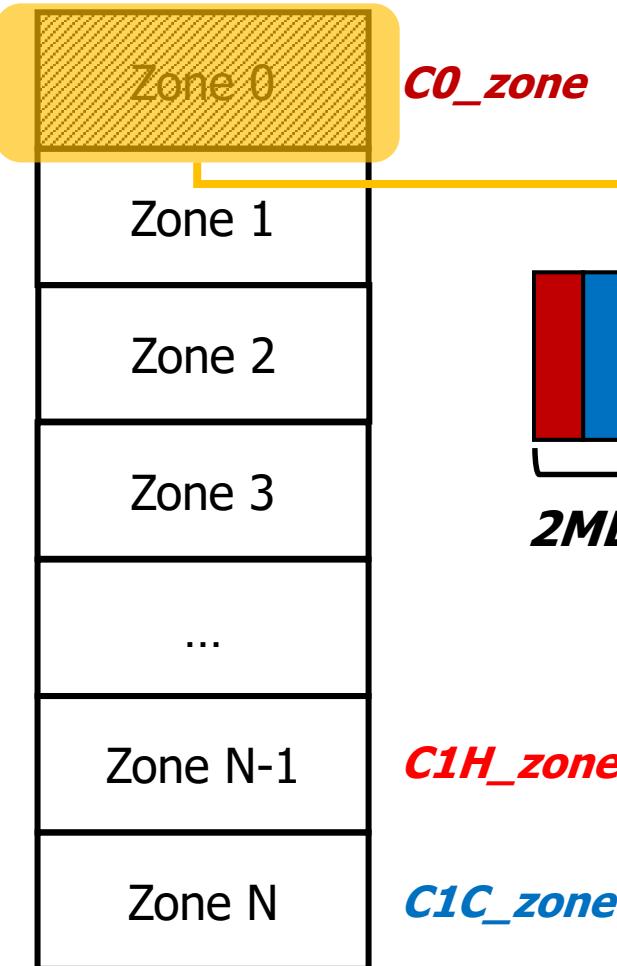
*Step 1.*  
*Select a candidate zone (or zones, Greedy or CB)*



*C1H\_zone*

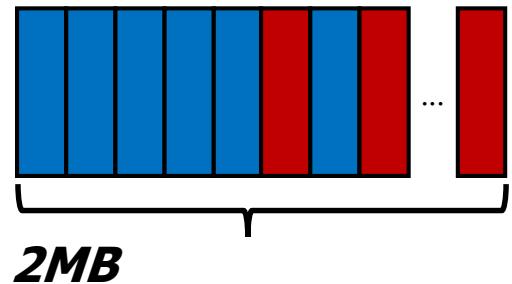
*C1C\_zone*

## LSM(Log Structured Merge) Zone GC



*Step 1.  
Select a candidate zone (or zones, Greedy or CB)*

*Step 2.  
Read all data in 128KB I/O size*

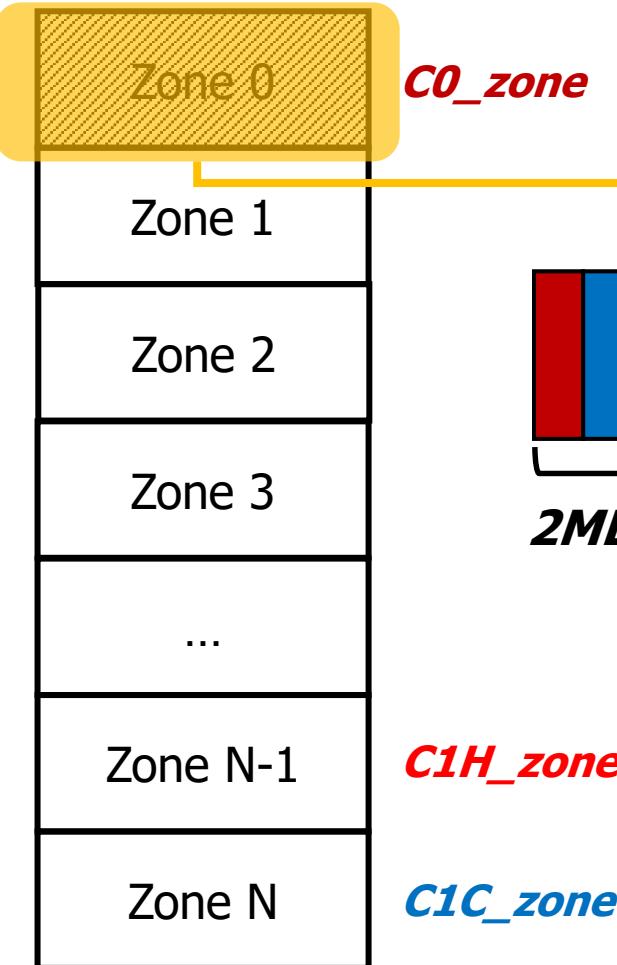


*2MB*

*2MB*

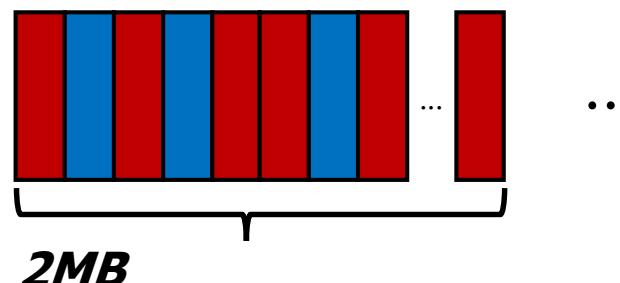
*2MB*

## LSM(Log Structured Merge) Zone GC



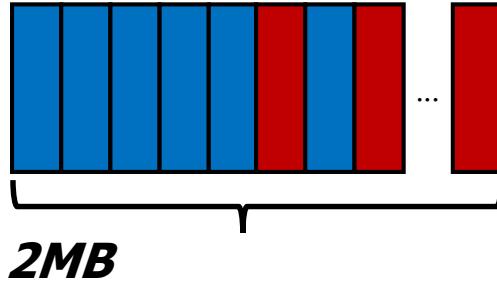
**Step 1.**  
*Select a candidate zone (or zones, Greedy or CB)*

**Step 2.**  
*Read all data in 128KB I/O size*

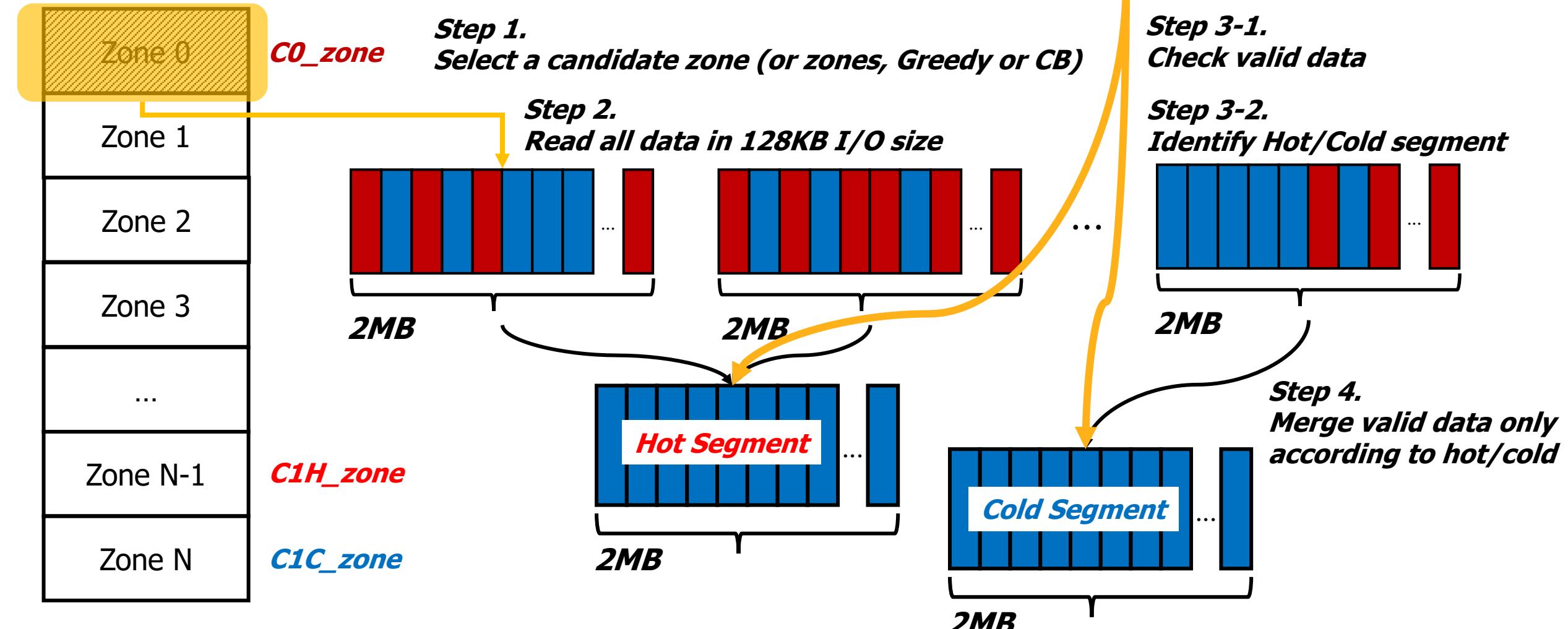


**Step 3-1.**  
*Check valid data*

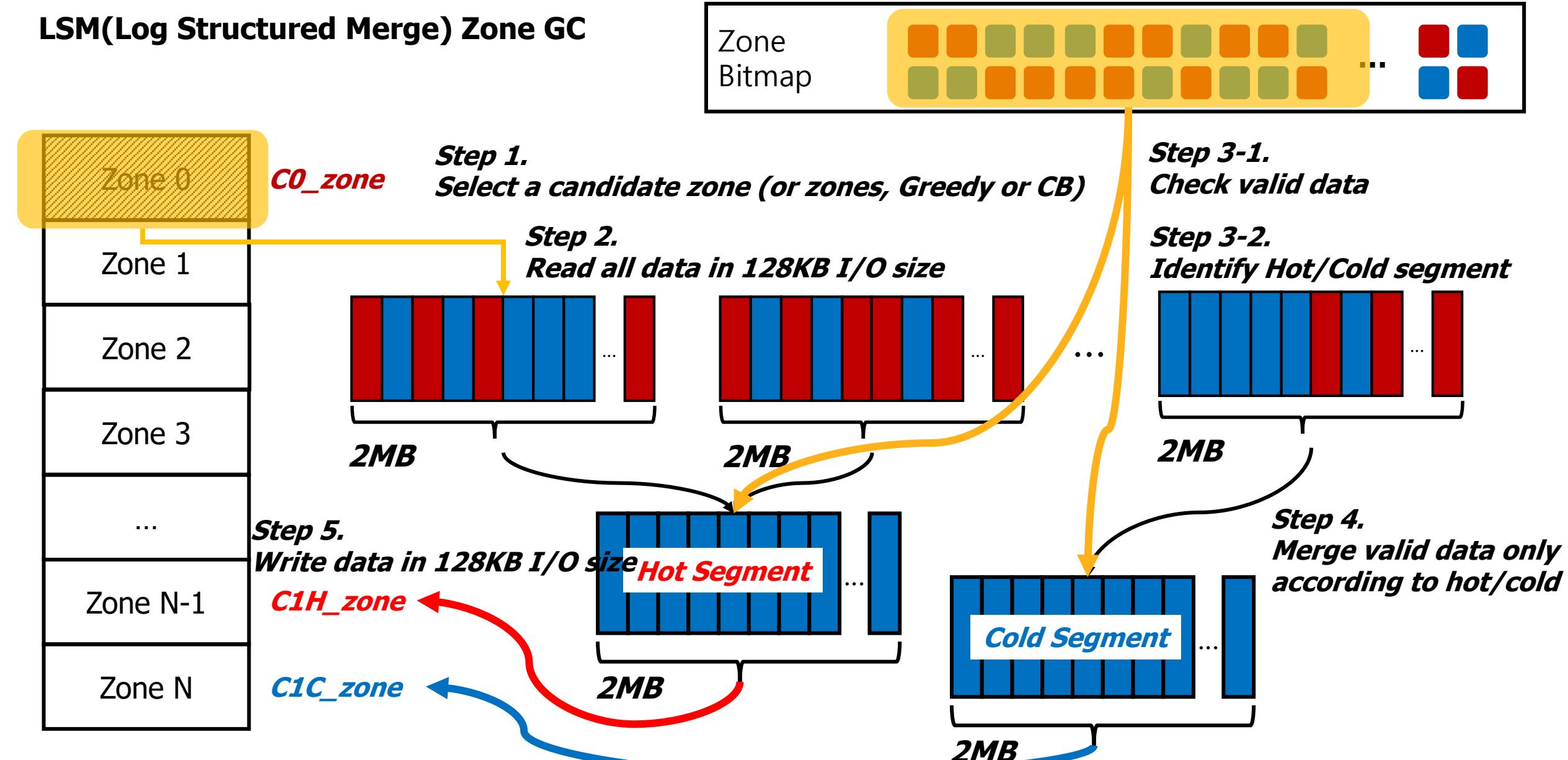
**Step 3-2.**  
*Identify Hot/Cold segment*



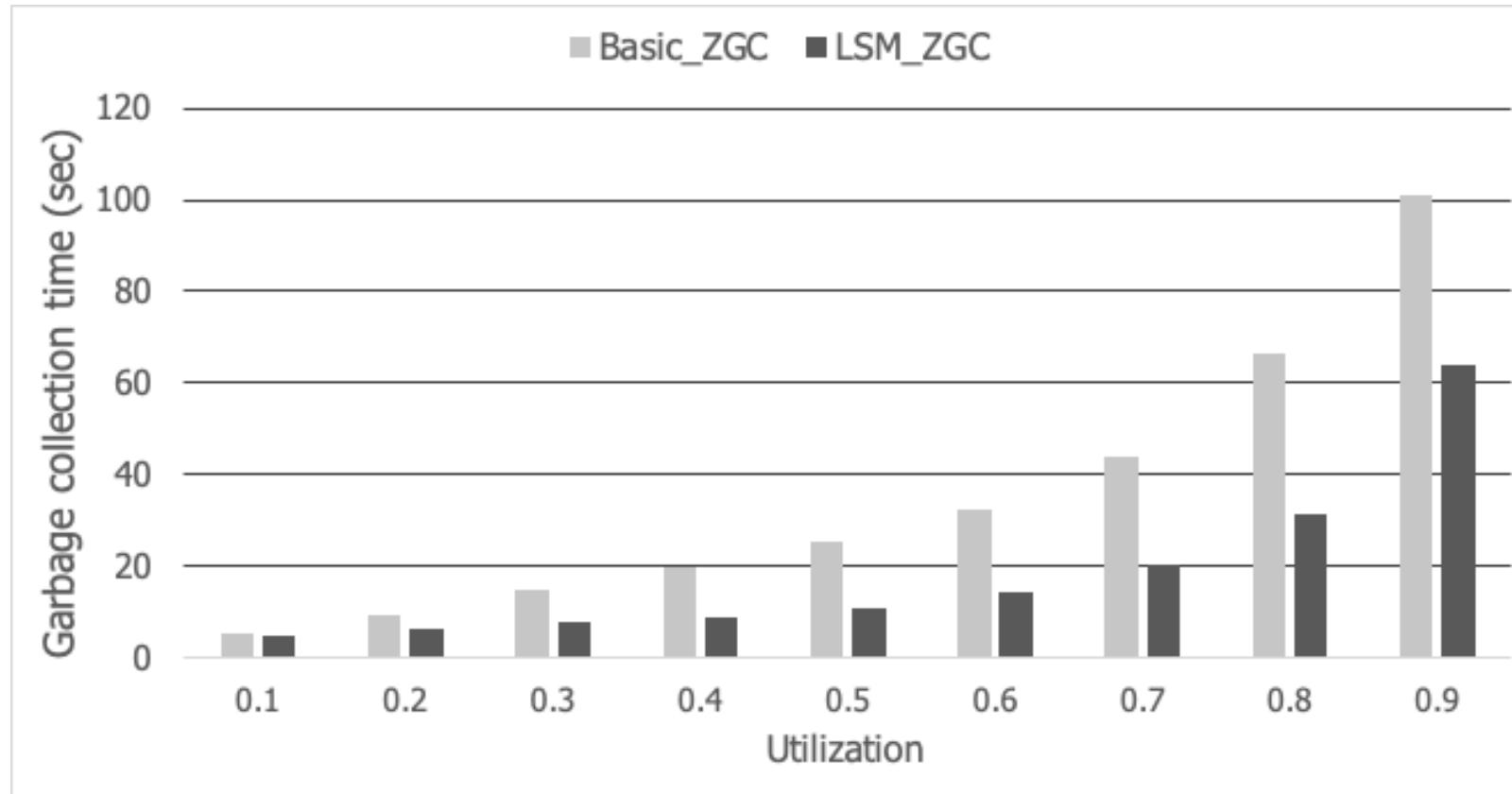
## LSM(Log Structured Merge) Zone GC



## LSM(Log Structured Merge) Zone GC



## Garbage collection overhead: uniform update pattern

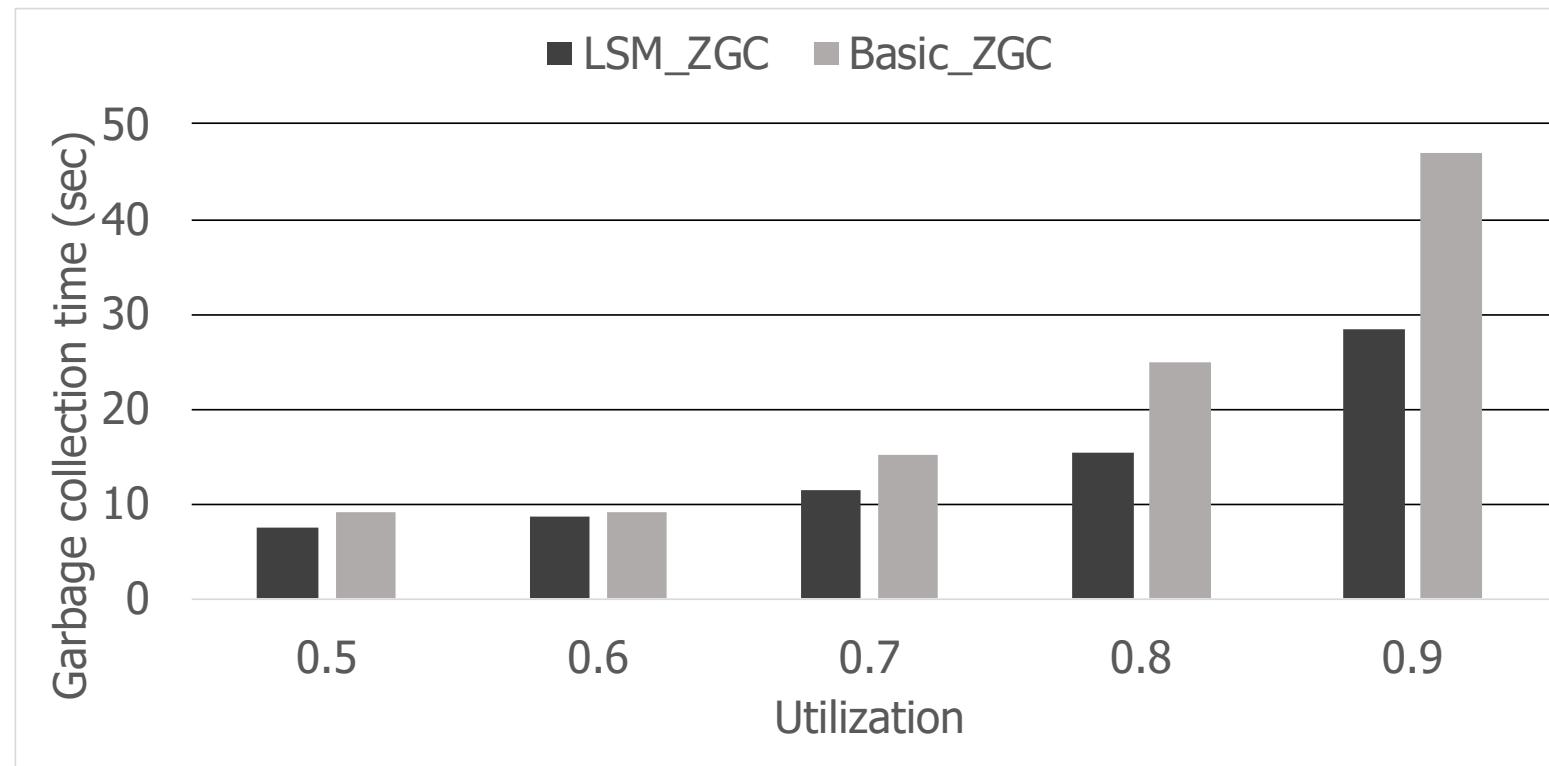


*Average of 1.9 times*  
*Max of 2.3 times*

### Experimental environment

- Intel Core i7 (8 core)
- 16GB DRAM
- 1TB ZNS SSD
- Size of Zone : 1GB

## Garbage collection overhead: skewed update pattern



**Average of 1.4 times**

**Max of 1.6 times**

### Parameters

- Workload: 70/30 hot/cold ratio
- Threahold<sub>cold</sub> : 0.8
- average utilization: x-axis

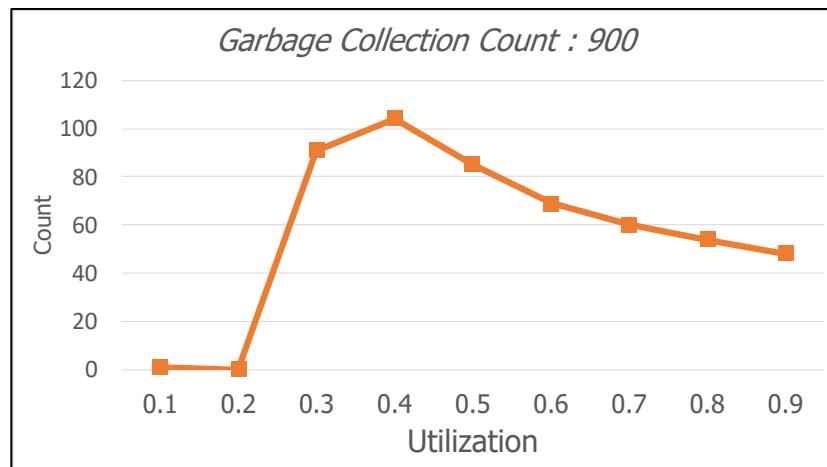
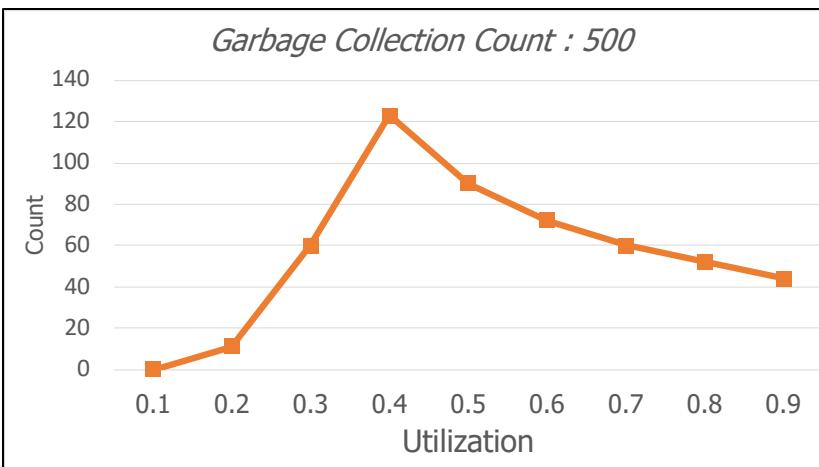
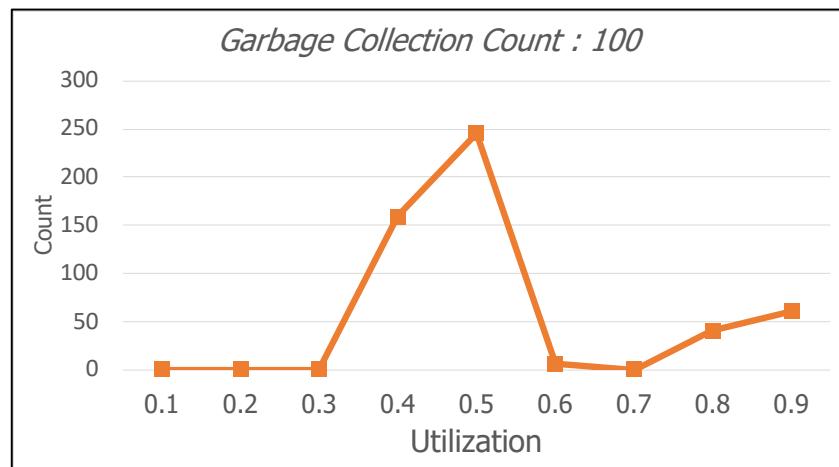
### Experimental environment

- Intel Core i7 (8 core)
- 16GB DRAM
- 1TB ZNS SSD
- Size of Zone : 1GB

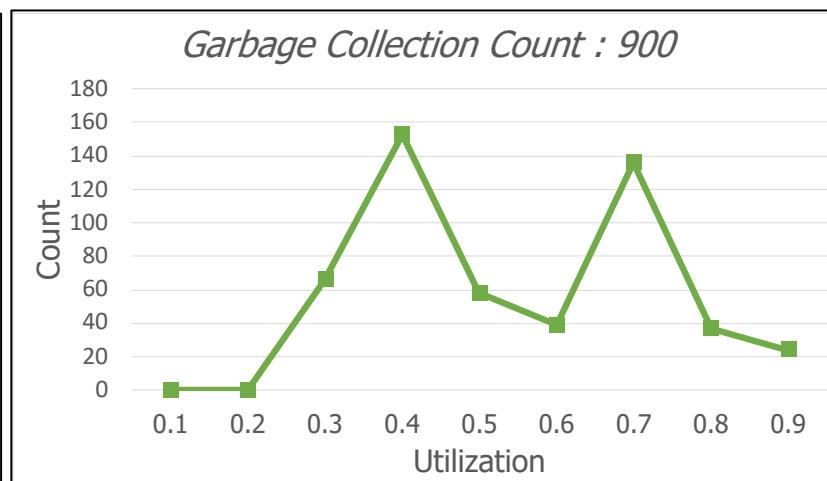
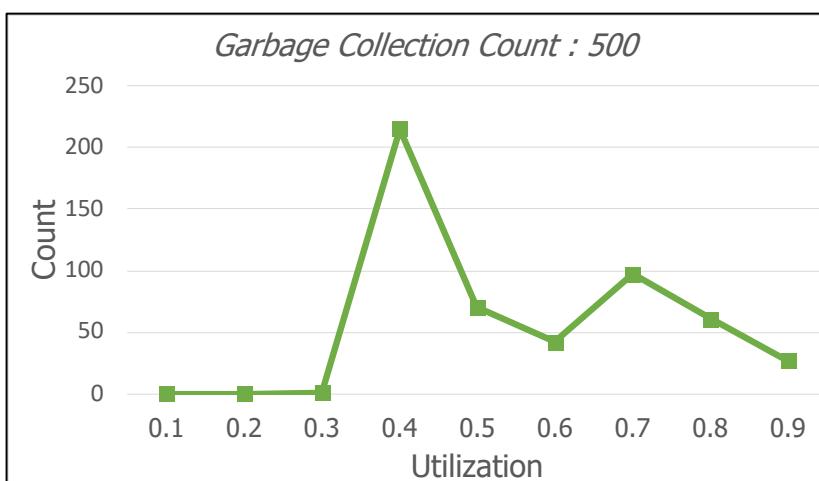
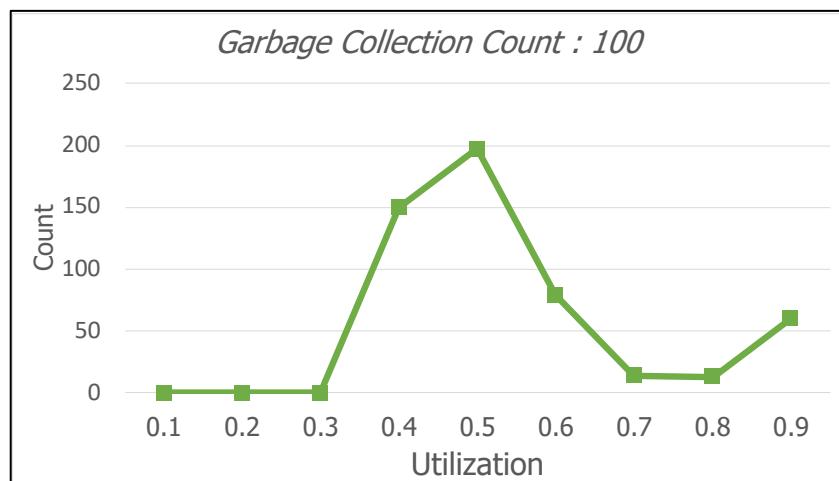
## 4. Evaluation

### Hot/Cold Separation

#### ✓ Without hot/cold separation



#### ✓ With hot/cold separation



### Parameters

- Workload: 70/30 hot/cold ratio
- Threahold<sub>cold</sub> : 0.8
- Average utilization: 0.6

- **Our contributions**

- *Observation: a zone garbage collection really matters*
- *Proposal: a new LSM-style zone garbage collection scheme*
- *Evaluation: real implementation based results*

- **Future work**

- *We are currently extending F2FS on our ZNS SSD prototype*
- *Also, evaluating LSM ZGC under diverse workloads with different hot /cold ratio, data size, initial placement and classification policies*

# A New LSM-style Garbage Collection Scheme for ZNS SSDs

*Gunhee Choi, Kwanghee Lee, Myunghoon Oh, Jhuyeong Jhin, Yongseok Oh, Jongmoo Choi*

*12<sup>th</sup> USENIX Workshop on Hot Topic in Storage and File System (HotStorage 20), 2020*

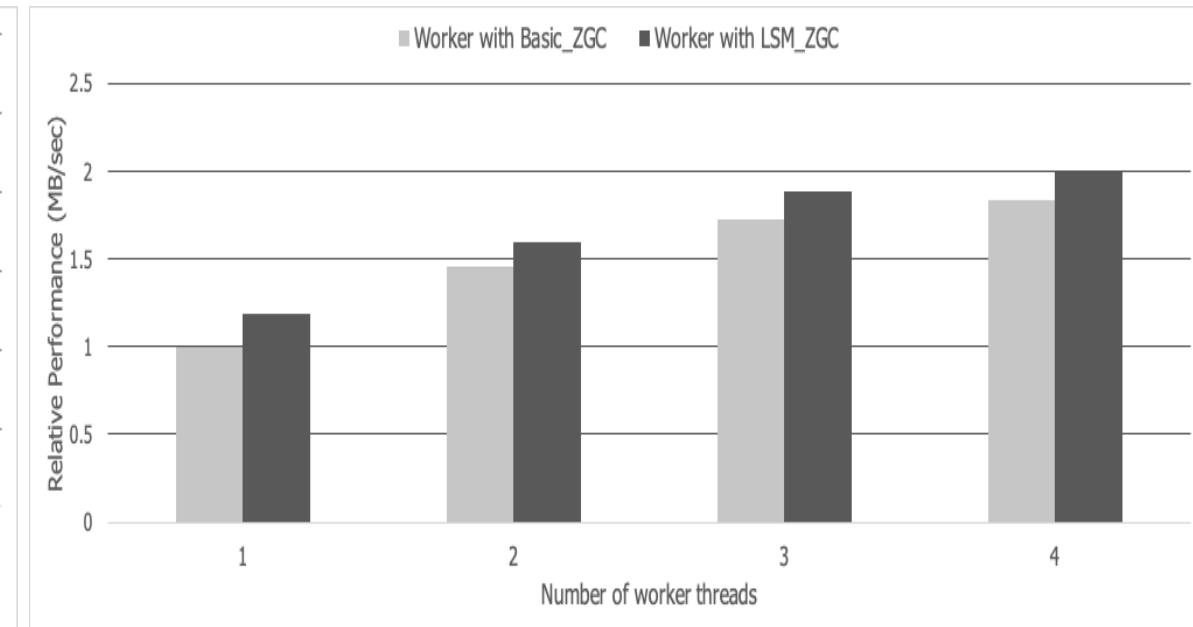
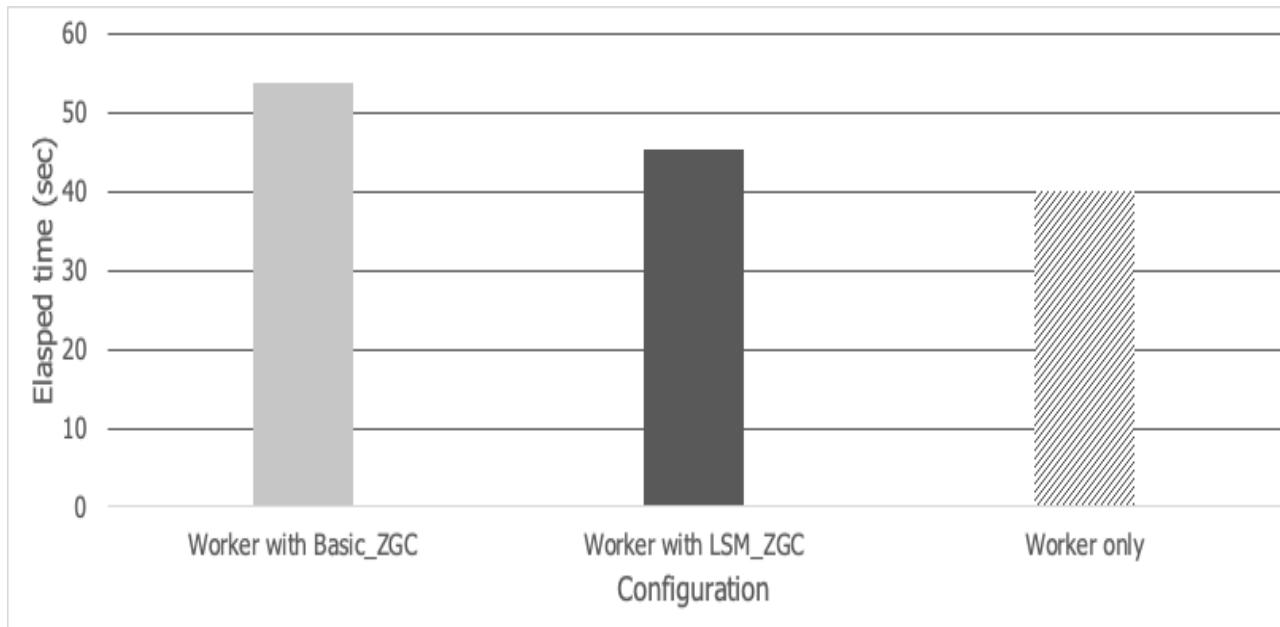
Thank You!  
Questions?

2020. 07. 13

Presentation by Choi, Gunhee

choi\_gunhee@dankook.ac.kr

## Performance comparison using multi-thread & Scalability

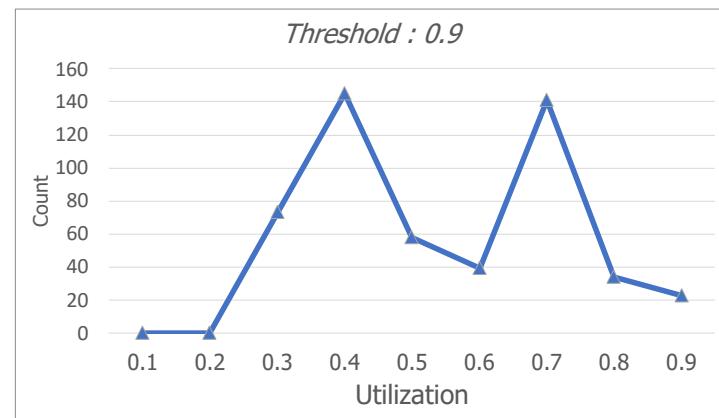
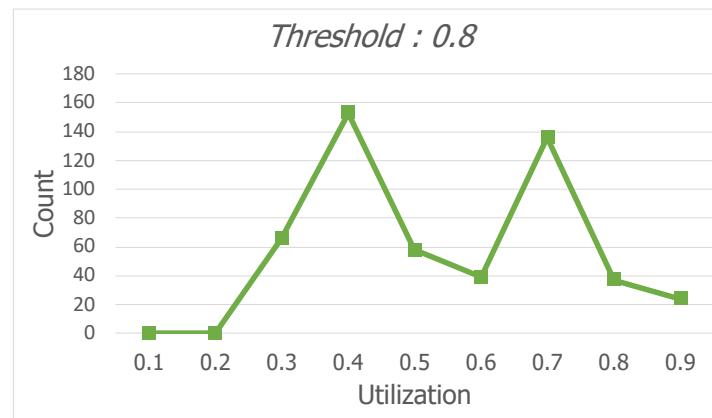
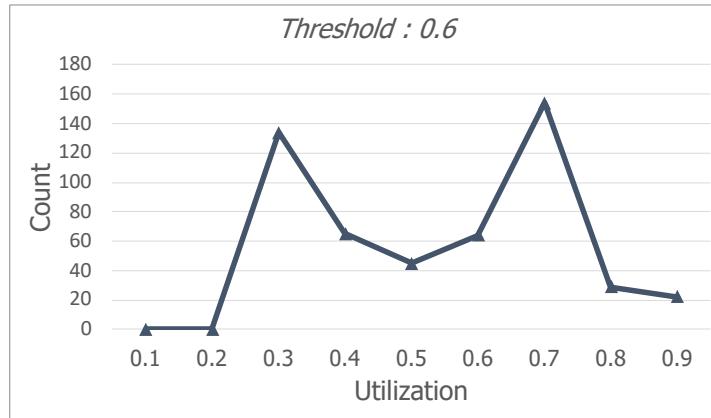


***Worker only: 40***  
***With LSM\_ZGC : 45***  
***With Basic\_ZGC : 53***

***Non-linear  
Scalability***

## Sensitive Analysis: various parameters

### ✓ Effect of threshold<sub>cold</sub> (initial utilization: 0.6)



### ✓ Effect of initial utilization of a zone (threshold<sub>cold</sub> : 0.8)

