# High-Level Summary

We analyze the problem of zoned namespaces (ZNS), by using two production ZNS SSDs

Problem: It is hard to exploit the internal parallelism of SSDs
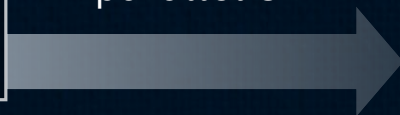(Reason: ZNS does not provide abstraction required to manage the parallelism)

We propose two simple modules

#1: Interference profiler

Get information which is required
to exploit the internal parallelism

information
related to the
parallelism

#2: Interference-aware I/O scheduler

Adjust the order of requests
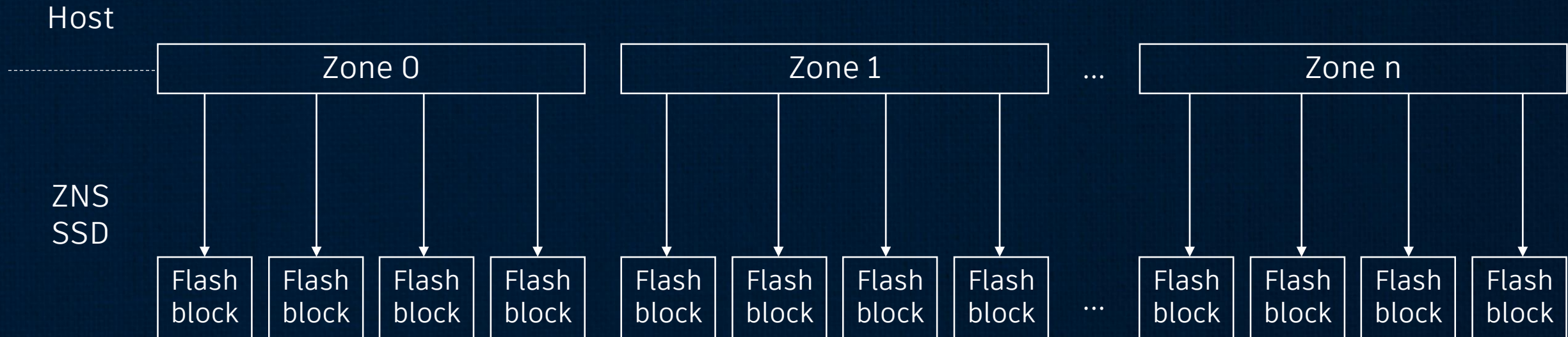to exploit the internal parallelism

# 1. Background – Zoned Namespaces

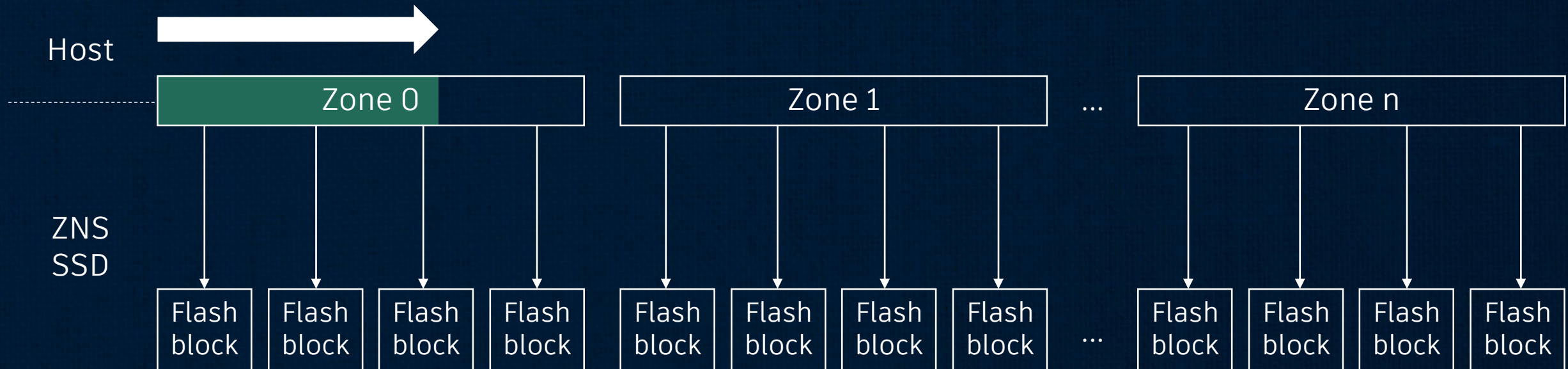2. Challenge

3. Solution

4. Evaluation

# What Is ZNS?

- Zoned namespaces (ZNS): Emerging storage interface
  - Divide logical address space into multiple zones
  - In general, each zone is mapped to one or more flash blocks

Host

| Zone 0 | Zone 1 | ... | Zone n |

ZNS
SSD

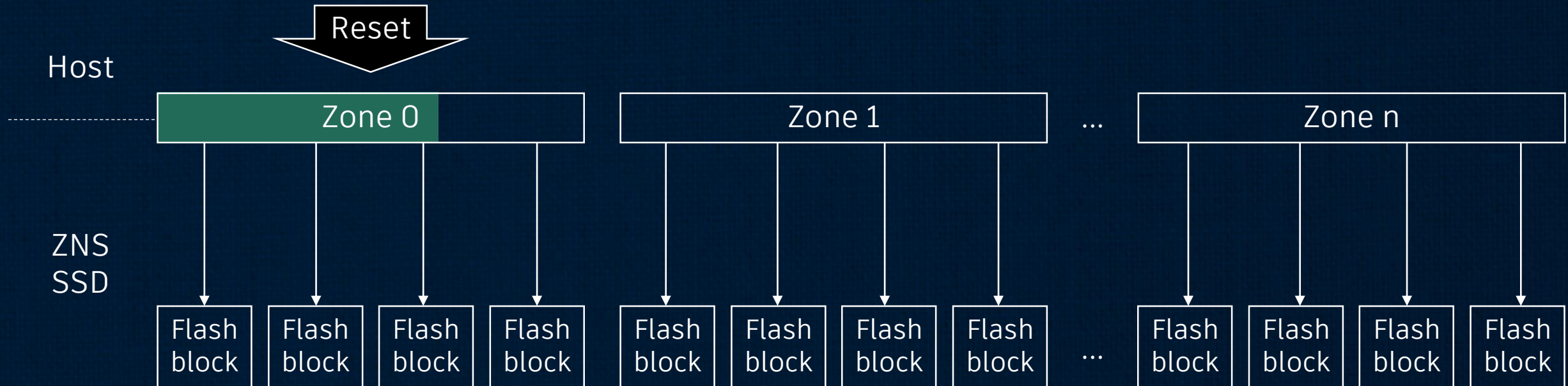| Flash block | Flash block | Flash block | Flash block | | Flash block | Flash block | Flash block | Flash block | ... | Flash block | Flash block | Flash block | Flash block |

# What Is ZNS?

- Two constraints on each zones
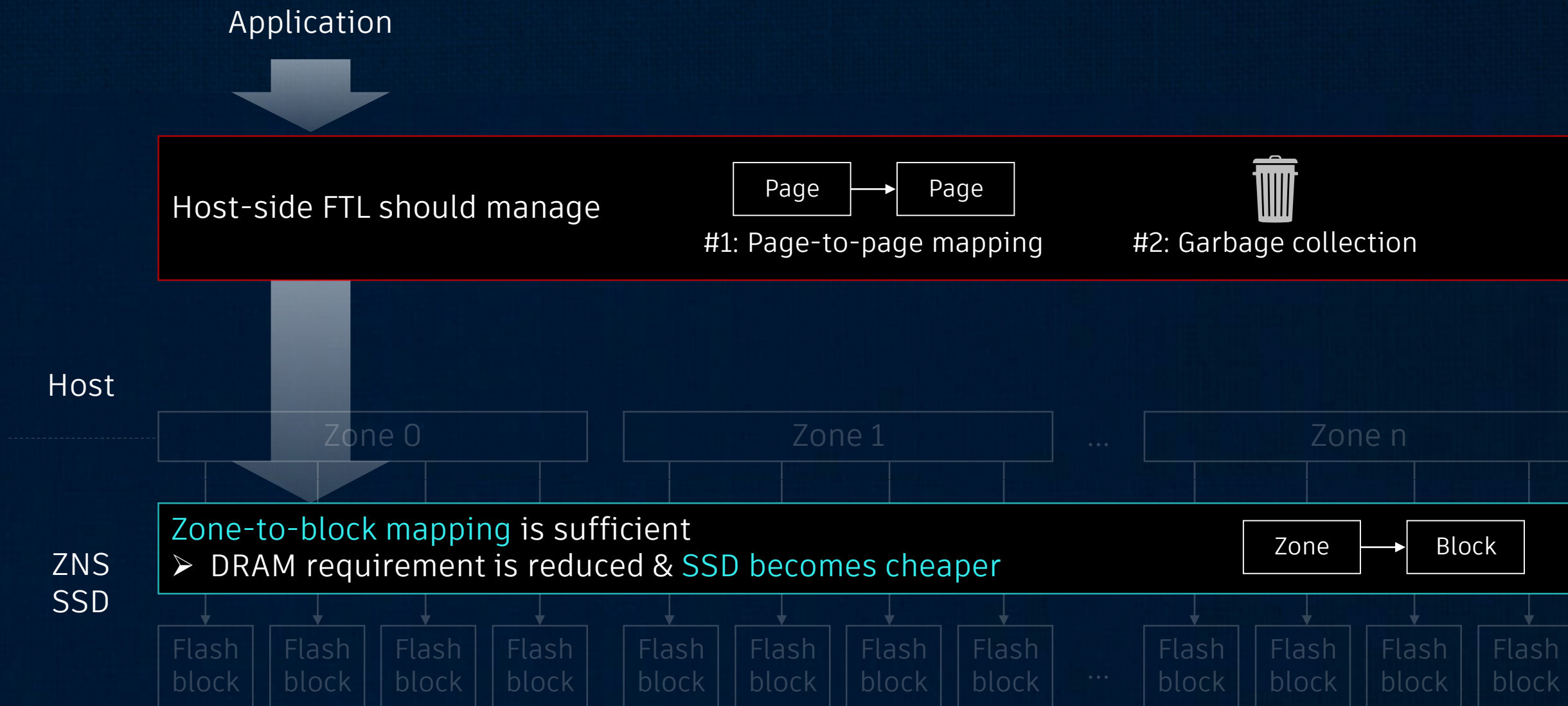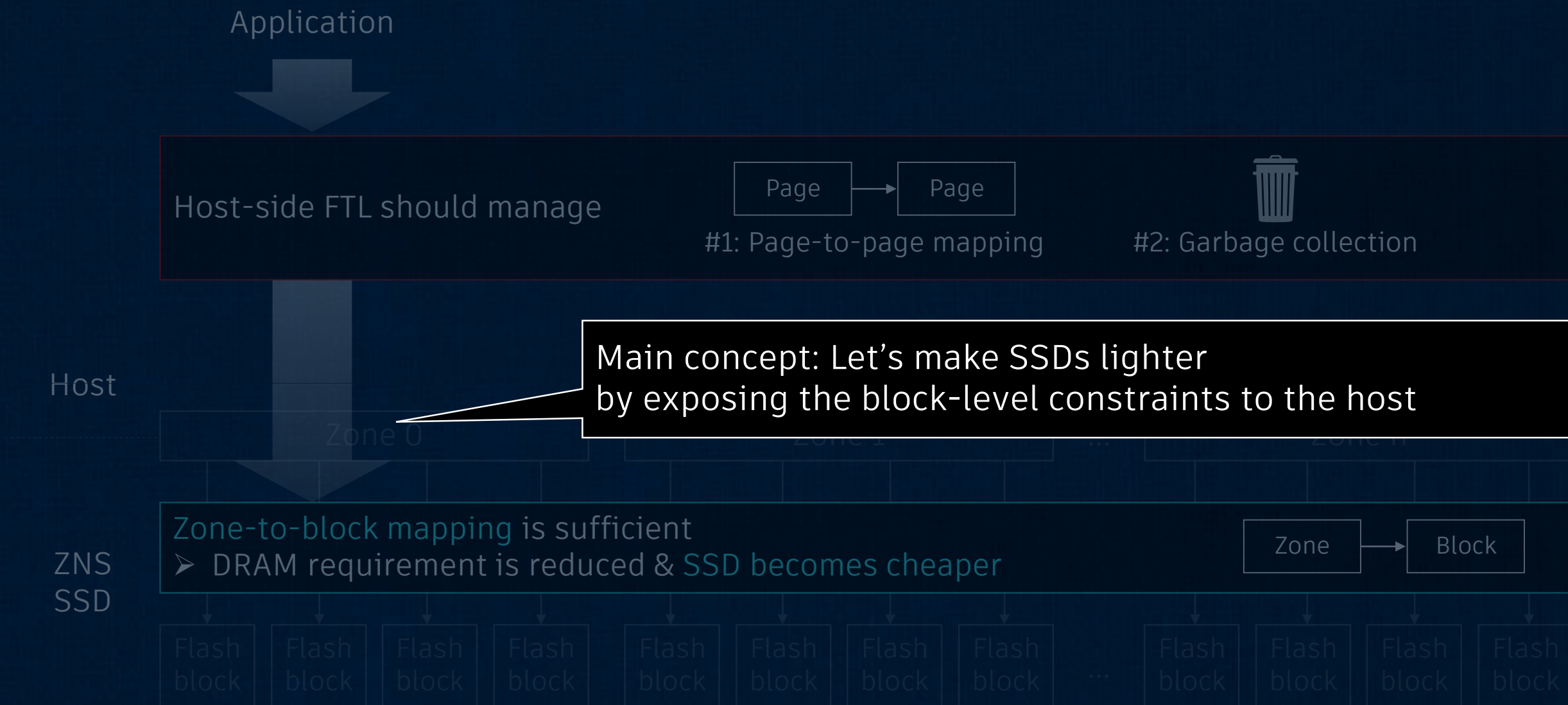    - Constraint #1: Sequential write

# What Is ZNS?

- Two constraints on each zones
    - Constraint #1: Sequential write
    - Constraint #2: Erase(reset)-before-write

# What Is ZNS?

Application

Host

**Host-side FTL should manage**

| Page | → | Page |

#1: Page-to-page mapping

#2: Garbage collection

Zone 0          Zone 1          ...          Zone n

## ZNS SSD

**Zone-to-block mapping** is sufficient
➤ DRAM requirement is reduced & **SSD becomes cheaper**

| Zone | → | Block |

Flash block | Flash block | Flash block | Flash block | Flash block | Flash block | Flash block | Flash block | ... | Flash block | Flash block | Flash block | Flash block

# What Is ZNS?

Application

Host-side FTL should manage

Page → Page

#1: Page-to-page mapping

#2: Garbage collection

Main concept: Let's make SSDs lighter
by exposing the block-level constraints to the host

Host

Zone 0          Zone 1          ...          Zone n

ZNS
SSD

Zone-to-block mapping is sufficient
➢ DRAM requirement is reduced & SSD becomes cheaper

Zone → Block

Flash block | Flash block | Flash block | Flash block | Flash block | Flash block | Flash block | Flash block | ... | Flash block | Flash block | Flash block | Flash block

# How Should Zones Be Configured?

- No constraint on zone size

**Large zone**



**Small zone**



** Above figure is different from
the configuration of the production ZNS SSDs that we used in this paper

# How Should Zones Be Configured?

- No constraint on zone size

Large zone

Small zone

Channel

Chip

| zone 0 | zone 0 | | zone 0 | zone 0 |
| zone 1 | zone 1 | | zone 1 | zone 1 |

| zone 0 | zone 0 | | zone 0 | zone 0 |
| zone 1 | zone 1 | | zone 1 | zone 1 |

➢ We advocate small zone devices

Channel

Chip

zone 0

zone 2

+) Higher degree of freedom for data placement

+) Less time required to mitigate
valid data in a zone for the host-level GC

zone 1

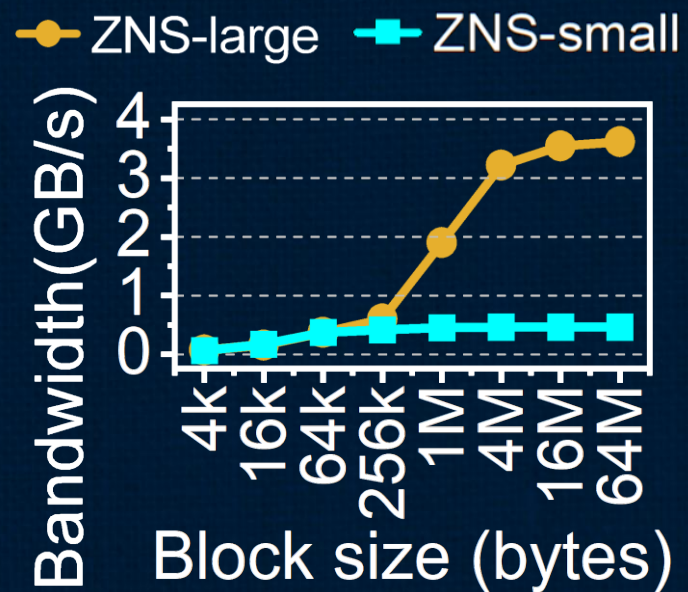zone 3

# Disadvantage of Small Zone Devices

Disadvantage:
They show much worse performance than large zone devices, when the request size is large



- Sequential read w/ single process
- Used two production ZNS SSDs
  - ZNS-large: ZNS SSD with large zones (2.18GB/zone)
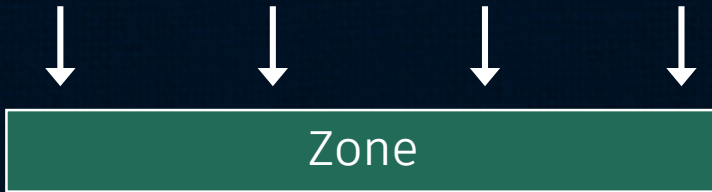  - ZNS-small: ZNS SSD with small zones (96MB/zone)

** Both SSDs utilize the same flash package
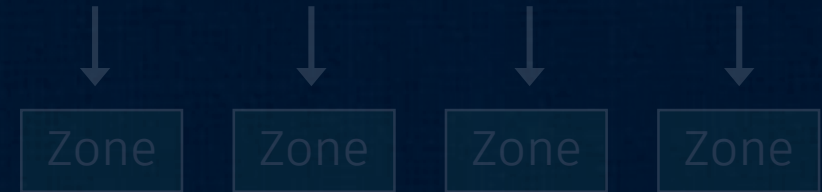** PCIe3.0 x4 (max. bandwidth = 3.94GB/s)

➢ Reason: Internal parallelism (especially, intra-zone parallelism)
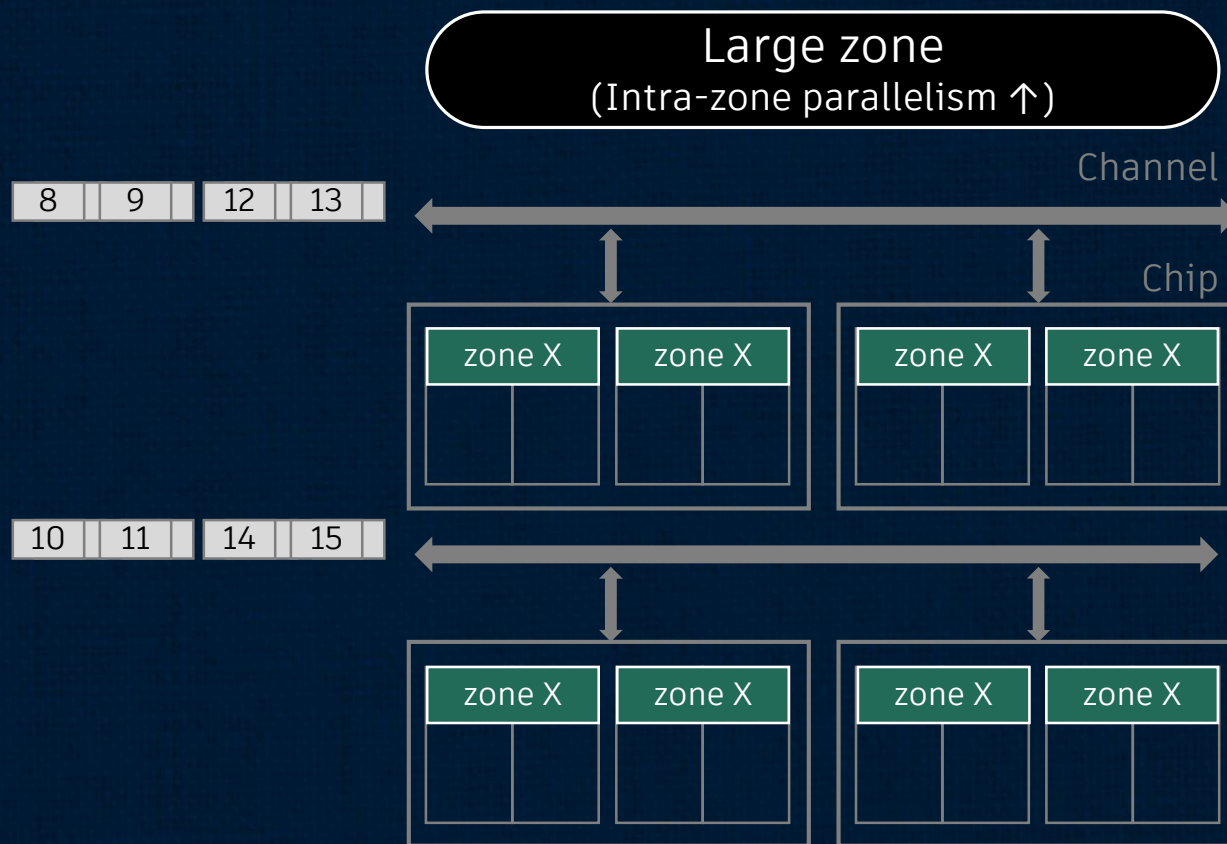
# Internal Parallelism in ZNS SSDs

Intra-zone parallelism
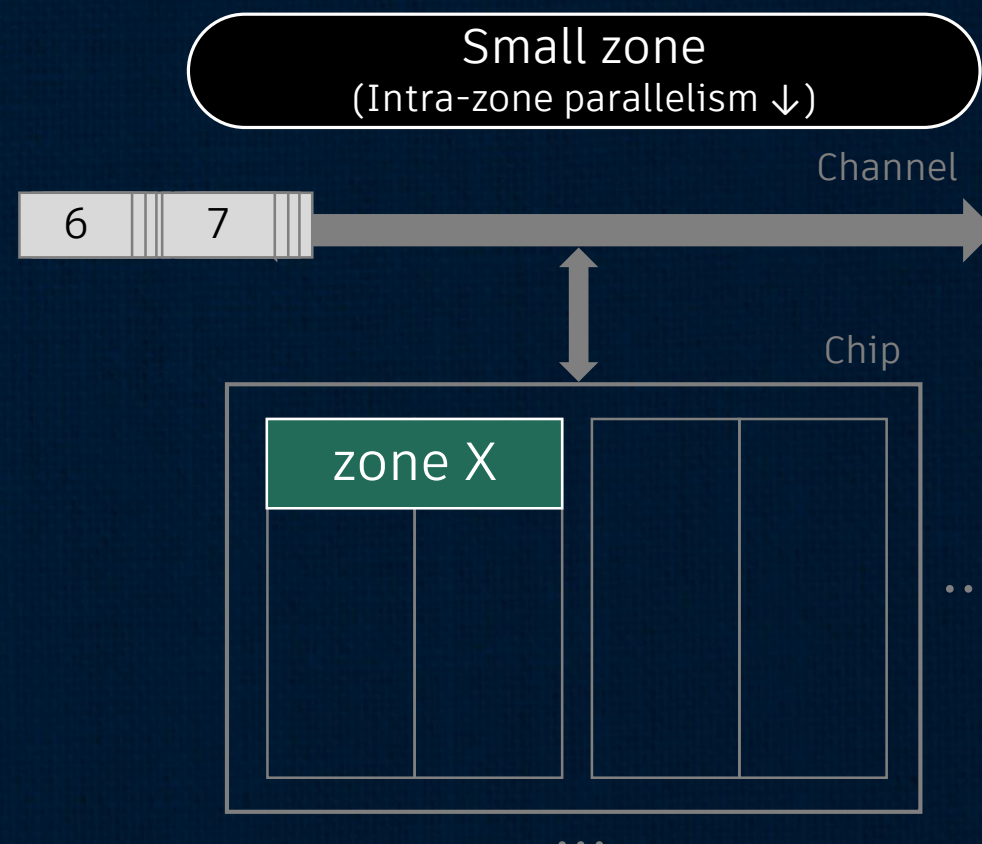
Zone

Inter-zone parallelism

Zone    Zone    Zone    Zone

# Intra-zone Parallelism

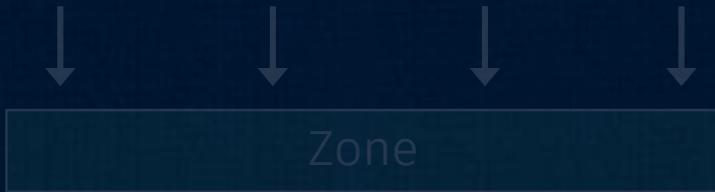Way to exploit the intra-zone parallelism
: Increase the request size

Request

| 0 | 1 | ... | 15 |

Sub-requests

**Large zone**
(Intra-zone parallelism ↑)

Channel

| 8 | 9 | 12 | 13 |

Chip

| zone X | zone X | | zone X | zone X |

| 10 | 11 | 14 | 15 |

| zone X | zone X | | zone X | zone X |

Able to get a high performance
by increasing the request size

**Small zone**
(Intra-zone parallelism ↓)

Channel

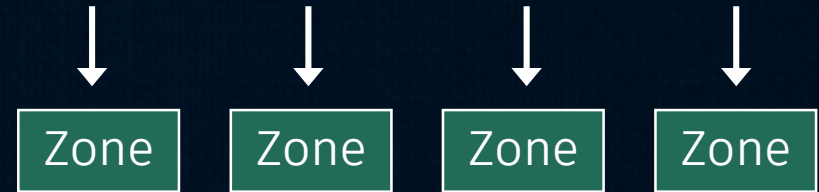| 6 | 7 |

Chip

| zone X | |

...

...

Unable to get a high performance
by increasing the request size

# Internal Parallelism in ZNS SSDs

Intra-zone parallelism

Zone

Inter-zone parallelism

Zone  Zone  Zone  Zone
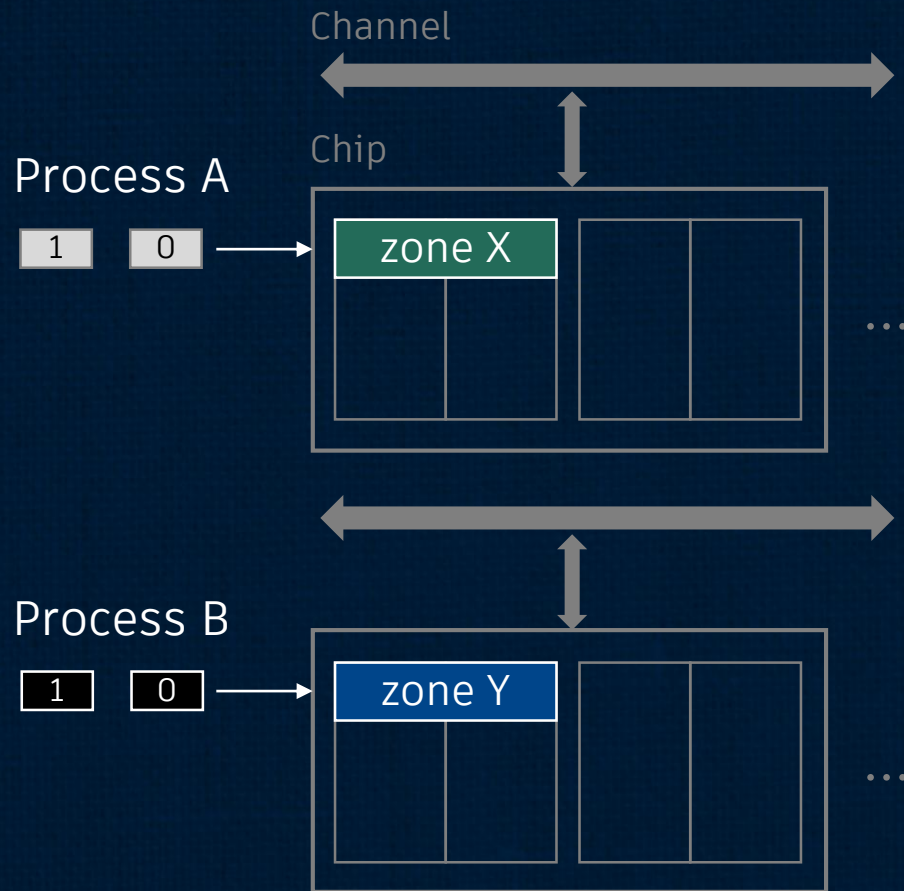
# Inter-zone Parallelism

Way to exploit the inter-zone parallelism
: Send requests to different zones at the same time

Channel

Chip

**Process A**

| 1 | 0 | → | zone X | | ... |

**Process B**

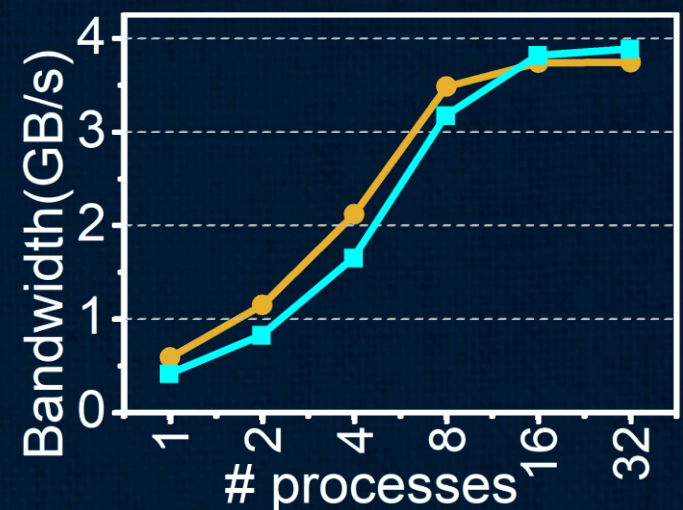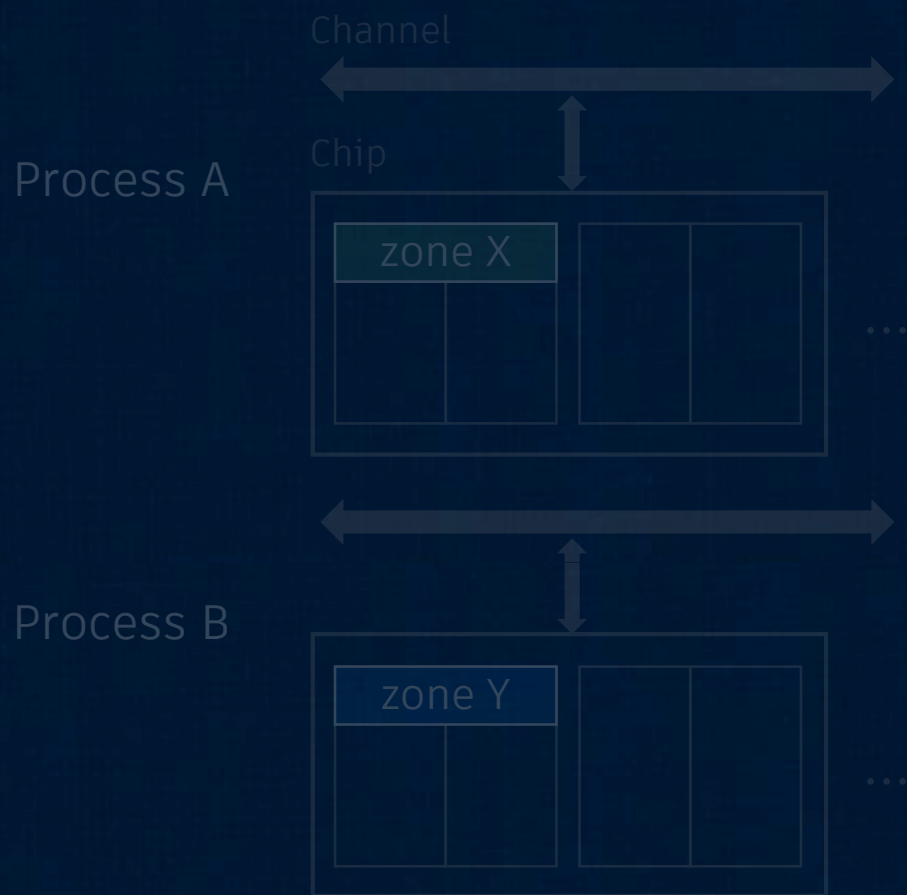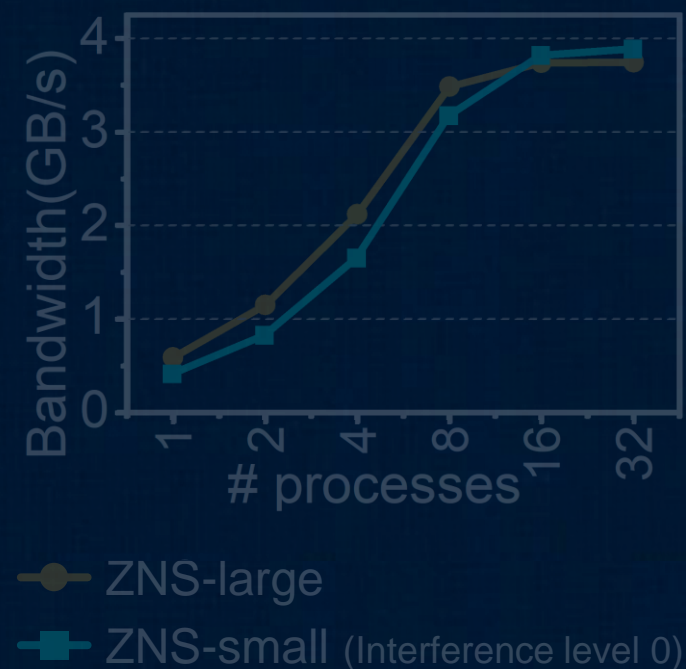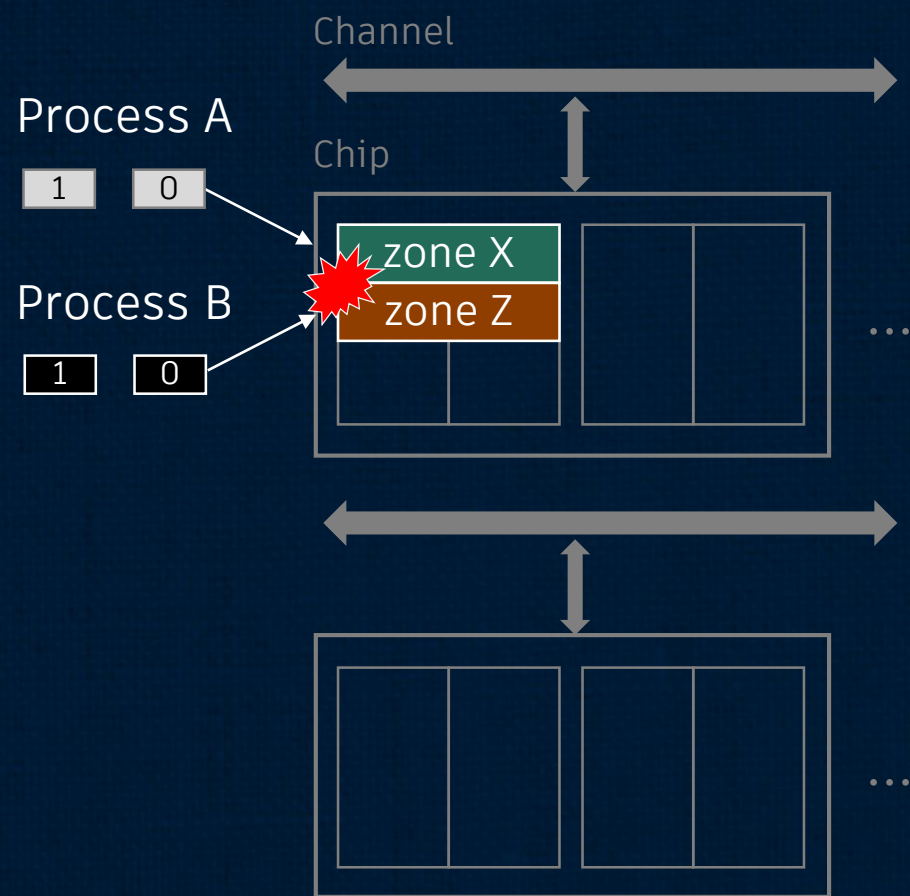| 1 | 0 | → | zone Y | | ... |

# Inter-zone Parallelism

Way to exploit the inter-zone parallelism
: Send requests to different zones at the same time

# Challenge: Inter-Zone Interference

Channel

Process A

`1`  `0`

Chip

zone X
zone Z

Process B

`1`  `0`

...

...

Bandwidth(GB/s)

# processes

ZNS-large

ZNS-small (Interference level 0)

CAMEL <18>

# Challenge: Inter-Zone Interference

The host **must be aware of the inter-zone interference**,
since it can cause a serious performance degradation

Channel

Process A

| 1 | 0 |

Chip

zone X
zone Z

Process B

| 1 | 0 |

...

...

2.84x

ZNS-large

ZNS-small (Interference level 0)

ZNS-small (Interference level 5)

ZNS-small (Interference level 6)

ZNS-small (Interference level 7)

# Requirements to Prevent the Interference

I/O requests

Host

zone X | zone Y | zone Z

Channel

Chip

zone X
zone Z

ZNS SSD

zone Y

#2: The host needs to schedule I/O requests appropriately, based on the information
(e..g., Prevent zone X and zone Z from being accessed simultaneously)

#1: The host must recognize
the zone-to-zone relationship
generating the interference

# Challenge: Not Enough Abstraction

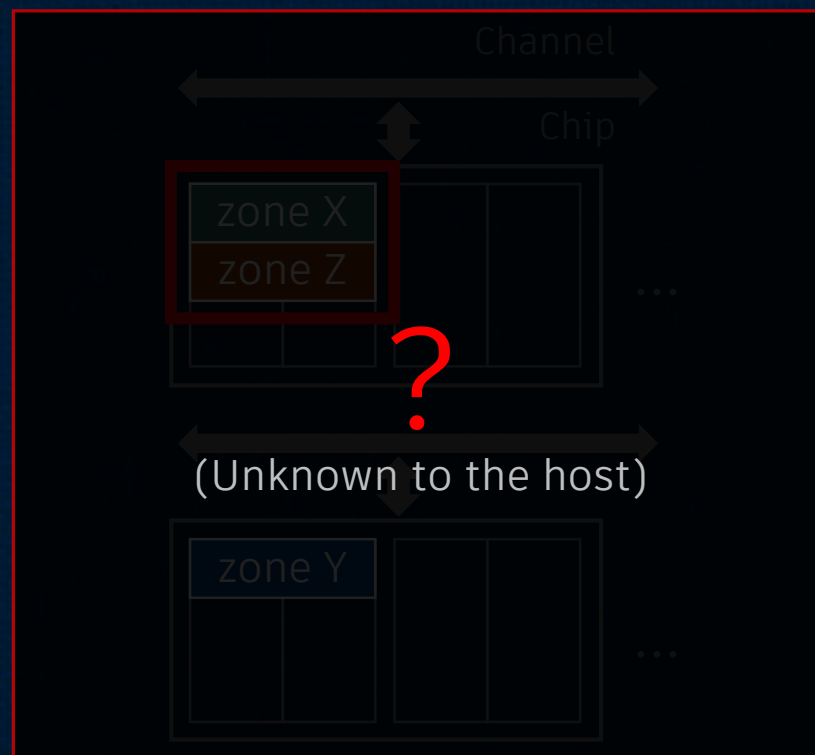I/O requests

Host

zone X  zone Y  zone Z

The host cannot prevent the interference since there is no information related to the interference/parallelism

ZNS SSD

Channel

Chip

zone X
zone Z

?

(Unknown to the host)

zone Y

...

Problem:
ZNS does not provide enough abstraction about hardware configuration
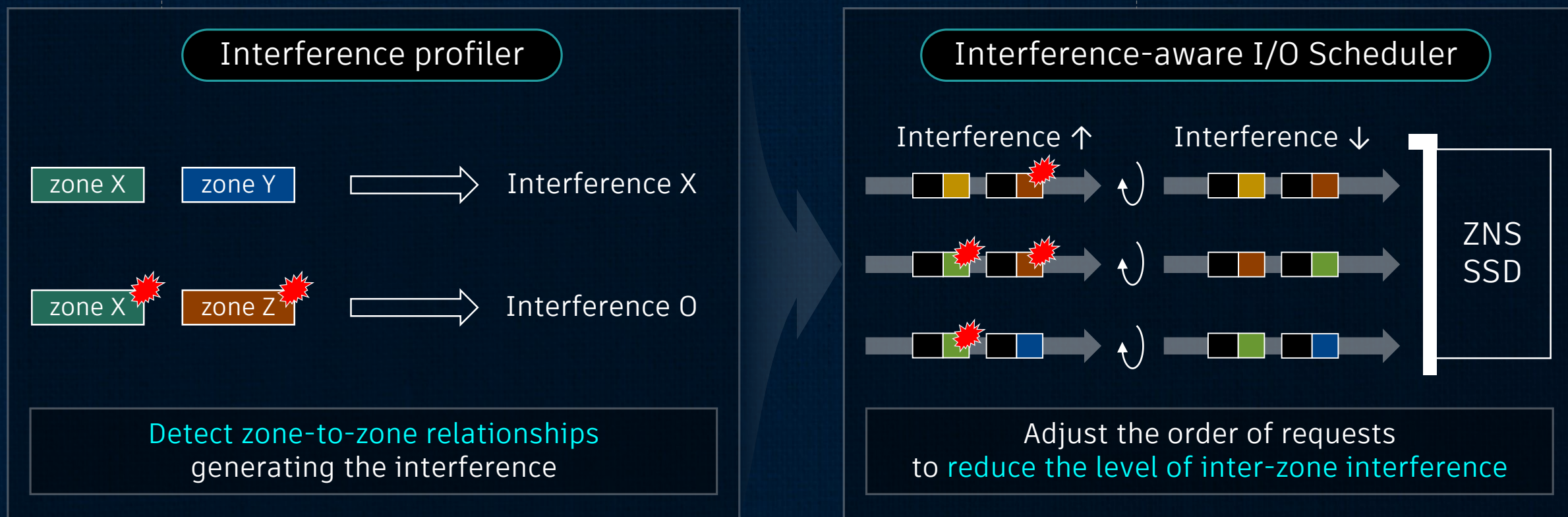
CAMEL <22>

# Overview

Problem:
Not enough abstraction about hardware ➜ Cannot prevent performance degradation due to the interference

## Interference profiler

| | | |
|---|---|---|
| zone X | zone Y | ➜ Interference X |
| zone X | zone Z | ➜ Interference O |

Detect zone-to-zone relationships
generating the interference

## Interference-aware I/O Scheduler

Interference ↑          Interference ↓

ZNS
SSD

Adjust the order of requests
to reduce the level of inter-zone interference

# Interference Profiler – Main Idea

# Interference Profiler

read                    read
⇩                       ⇩
zone X                  zone Y
baseline zone           target zone

High bandwidth          Low bandwidth
: No interference       : Interference

Low bandwidth: Interference

Baseline Zones (7 6 5 4 3 2 1 0)

Evaluation Target Zones (9 11 13 15 17 19 21 23 25 27 29 ... 2047)

Bandwidth (MB/s)

900
850
800
750
700
650
600

Conflict group 0    zone 0    zone 9    zone 16    zone 25  ...

Conflict group 1    zone 1    zone 8    zone 17    zone 24  ...

...

** Conflict group: Set of zones interfering with each other

CAMEL <25>

# Interference Profiler

In: List of zones to analyze

| zone 0 | ... | zone Z |

Interference profiler

Out: Zone-to-CG (Z2C) mapping table

Conflict group 0 | zone 0 | zone 9 | zone 16 | zone 25 | ...

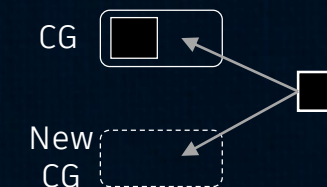Conflict group 1 | zone 1 | zone 8 | zone 17 | zone 24 | ...

...

Set the
threshold bandwidth

Initialize the first CG
with the first zone

CG 0 | zone A |

Classify the remaining zones
into conflict groups
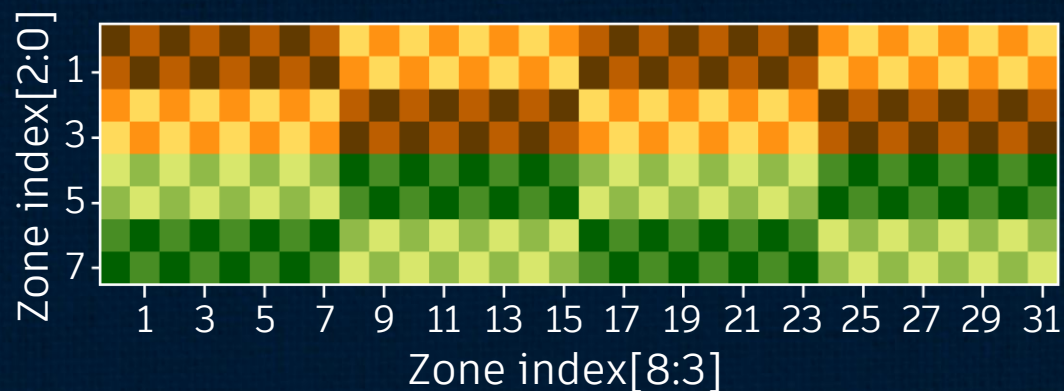Case 1: Add to the existing CG
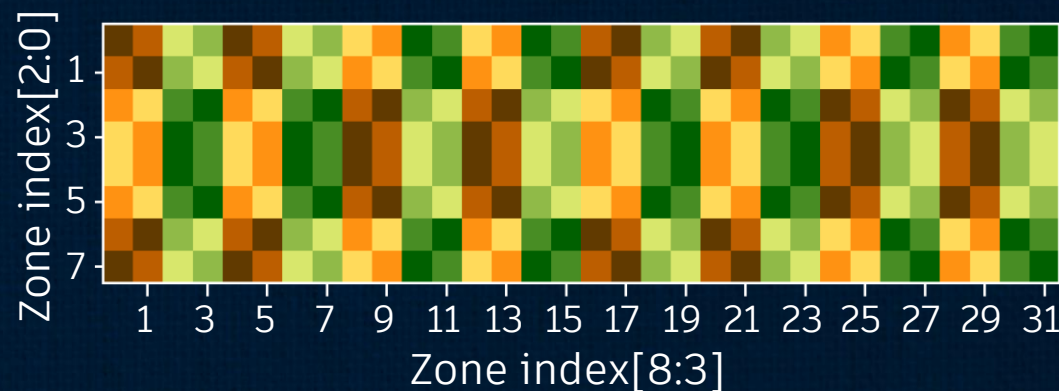Case 2: Create a new CG

CG

New
CG

# Interference Profiler - Results

Conflict group  0 1 2 3 4 5 6 7

Case #1

Case #2



**Zone index[2:0]** (y-axis, Case #1): 1, 3, 5, 7
**Zone index[8:3]** (x-axis, Case #1): 1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 31

**Zone index[2:0]** (y-axis, Case #2): 1, 3, 5, 7
**Zone index[8:3]** (x-axis, Case #2): 1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 31

...

\*\*Above patterns are repeated

- **Z2C mapping can vary** based on the order in which zones are written

# Interference Profiler - Results



Conflict group  0  1  2  3  4  5  6  7

Case #1                                    Case #2

Zone index[2:0]

Conflict group 0    | zone 0 | 💥 | zone 9 | 💥 | zone 16 | 💥 | zone 25 | ...

No interference

Conflict group 1    | zone 1 | 💥 | zone 8 | 💥 | zone 17 | 💥 | zone 24 | ...

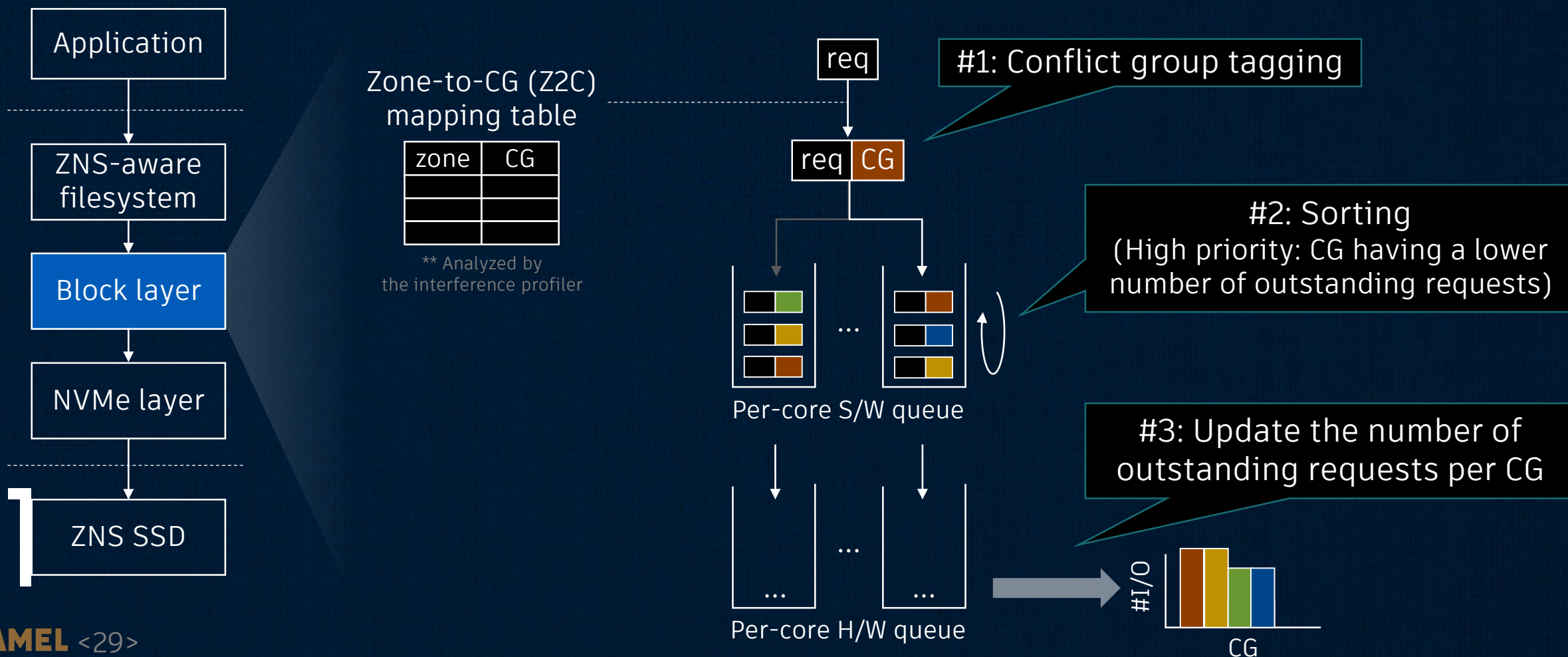1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 31

**Above patterns are repeated

We can exploit the paralleism by accessing zones from different CGs

- Z2C mapping can vary based on the order in which zones are written

CAMEL <28>

# Interference-aware I/O Scheduler

Goal: Schedule I/O requests coming from different CGs as many as possible
➢ Able to exploit the internal parallelism

Application

ZNS-aware filesystem

**Block layer**

NVMe layer

ZNS SSD

Zone-to-CG (Z2C)
mapping table

| zone | CG |
|------|-----|
|      |     |
|      |     |

** Analyzed by
the interference profiler

req

req CG

#1: Conflict group tagging

#2: Sorting
(High priority: CG having a lower
number of outstanding requests)

Per-core S/W queue

#3: Update the number of
outstanding requests per CG

Per-core H/W queue

#I/O

CG

# Evaluation Setting

## Schedulers

- blk_mq
  - Multi-queue I/O scheduler of Linux

- zns_mq
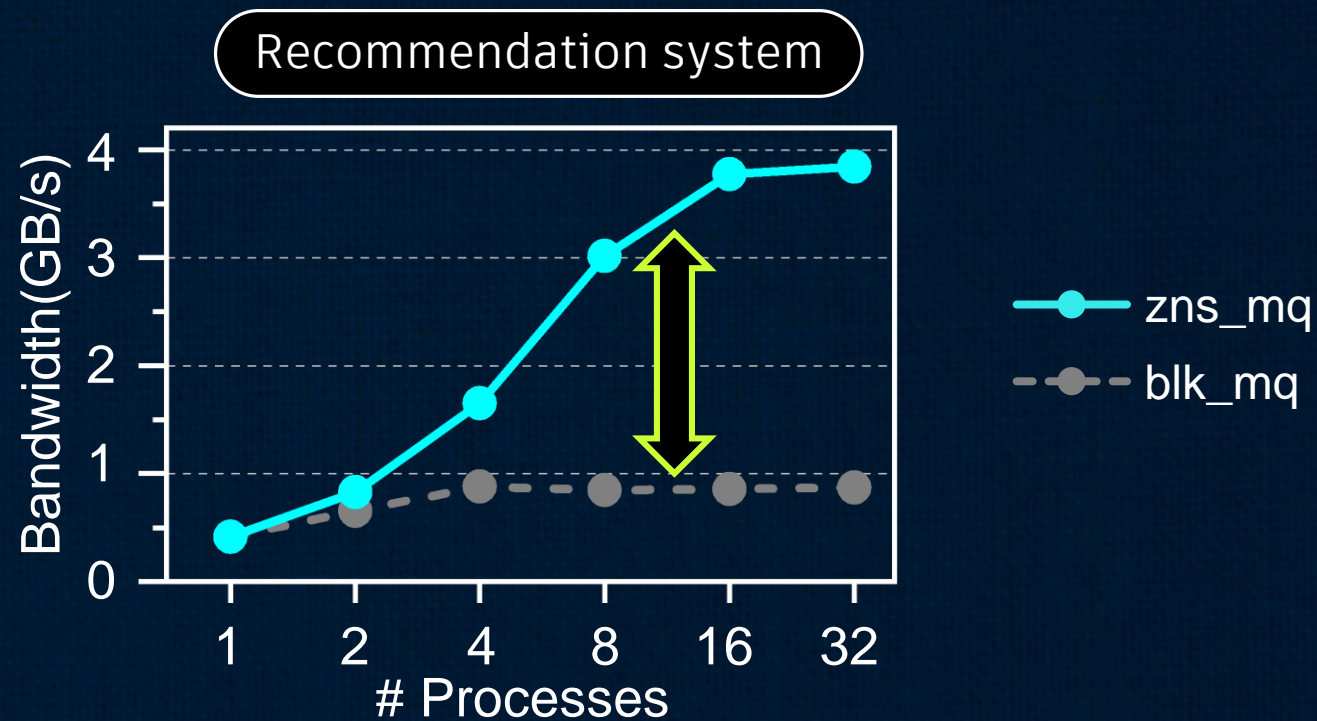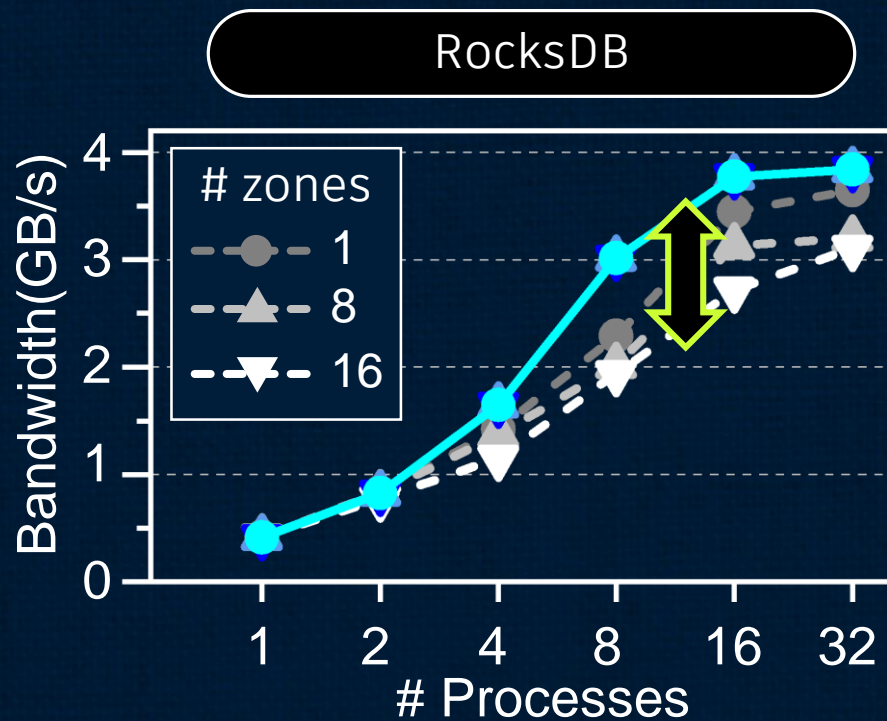  - Multi-queue I/O scheduler that utilizes our interference profiling information

## Workloads

- RocksDB
  - 1~16 zones / SST file

- Recommendation system
  - 128 zones / embedding table
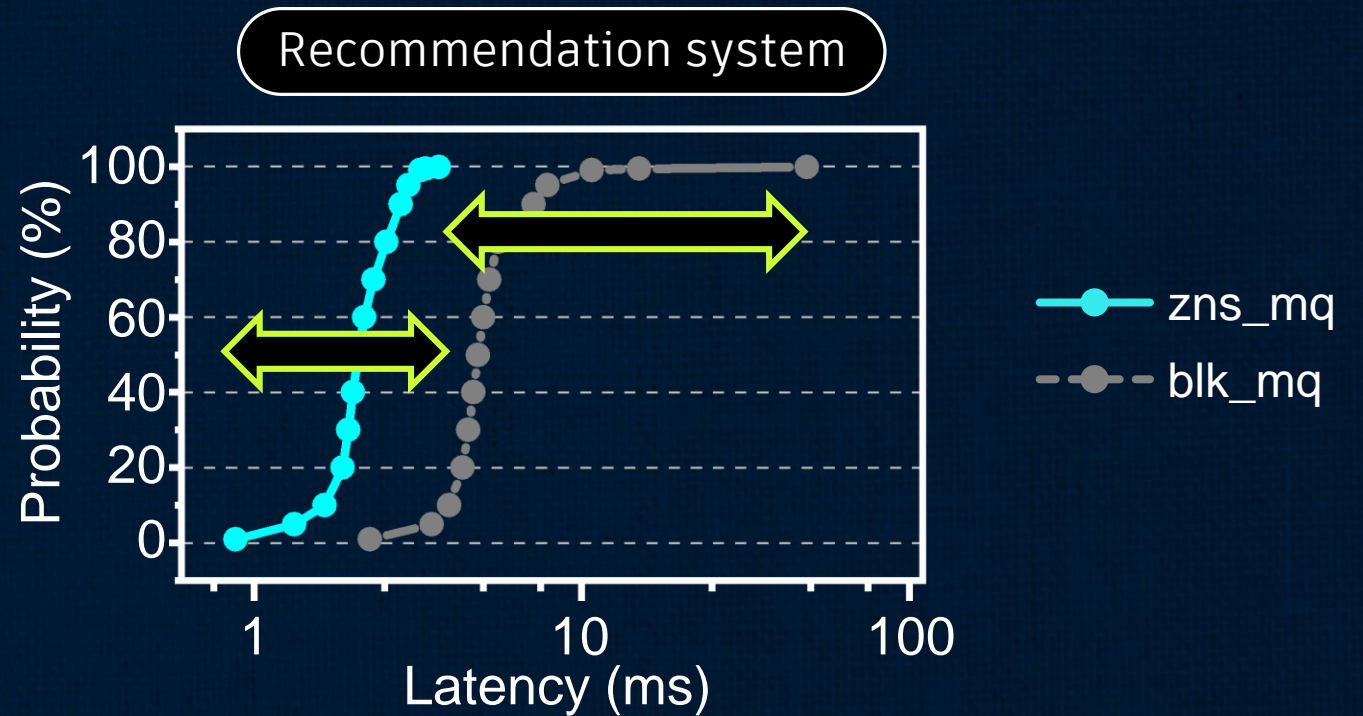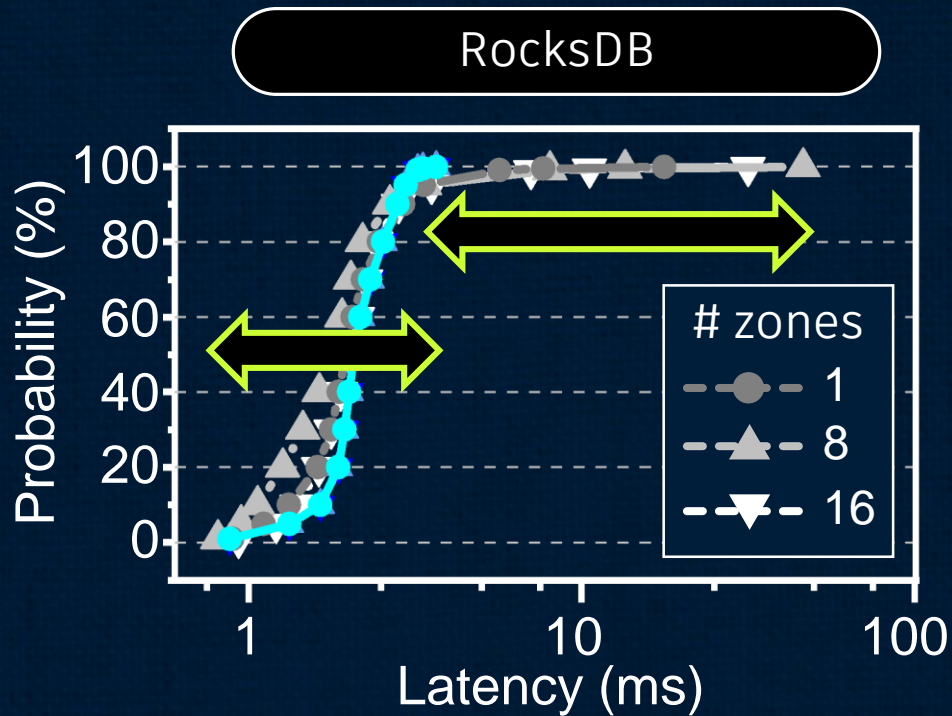  - Embedding table contains 50M indices with 64 dimensions

## Environments

- ZNS-small
  - production ZNS SSD
  - 96MB/zone
  - TLC-based flash

- Intel Xeon CPU
  - 2.3GHz, 20cores, 40 vCPUs

# Evaluation Results - Bandwidth



RocksDB

Recommendation system

- zns_mq improves the bandwidth of blk_mq by 1.98x, on average

# Evaluation Results – Tail Latency



- zns_mq shows narrower width of the distribution
  : All I/O requests experience similar interference levels
- zns_mq exhibits 11x shorter three nine (99.9%) tail latency, on average
- zns_mq exhibits 2.2x shorter latency, on average

# Conclusion

- By using two production ZNS SSDs, we quantitively analyze
  the performance degradation due to the inter-zone interference

- We propose two simple modules to exploit the internal parallelism of ZNS SSDs
  - Interference profiler
  - Interference-aware I/O scheduler

- Our evaluation results show that
  our mechanism can improve the bandwidth and latency, significantly

# Thank You

Contact: Hanyeoreum Bae (hyr.bae@camelab.org)