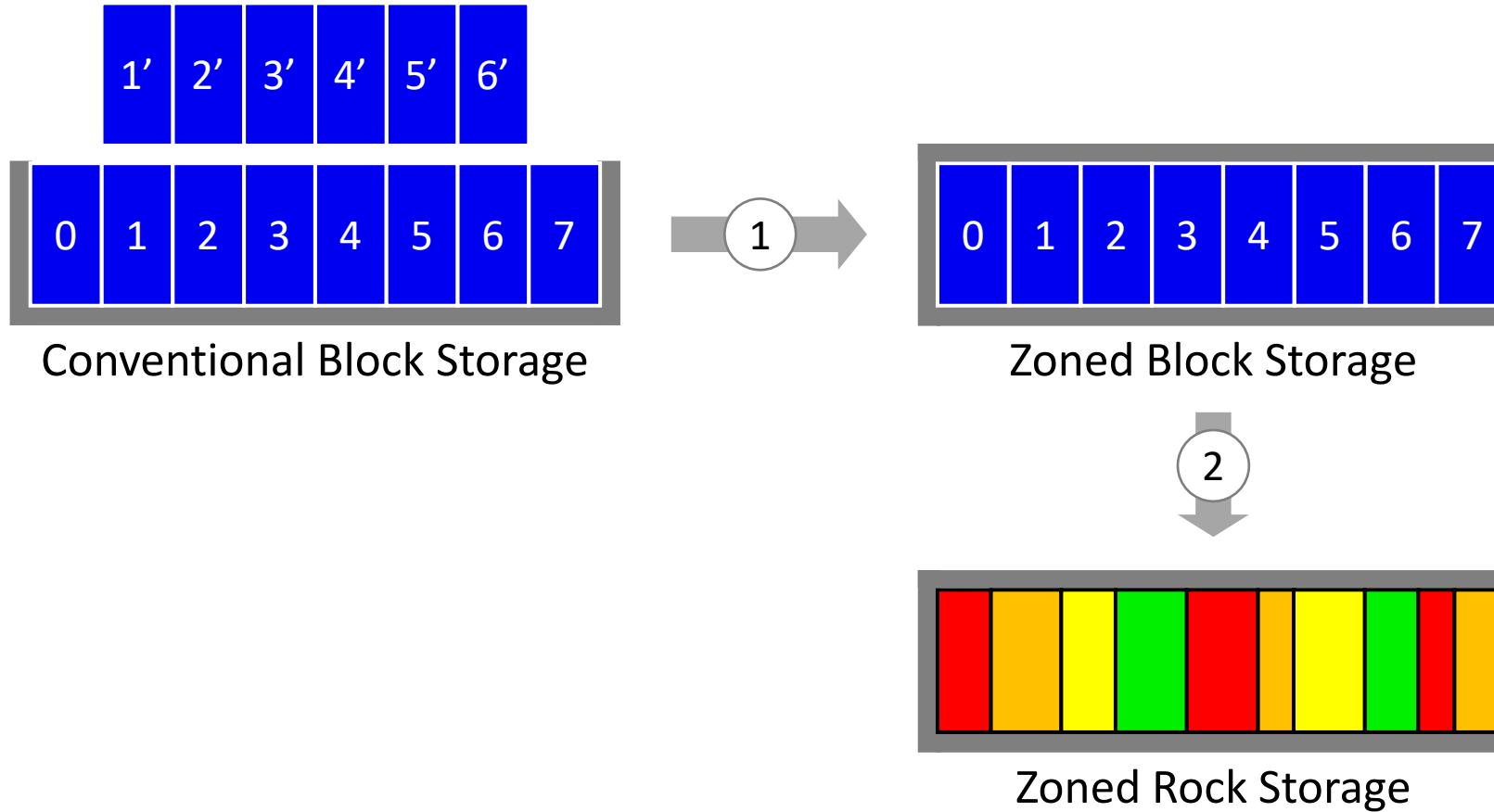# From Blocks to Rocks:
# A Natural Extension of Zoned Namespaces
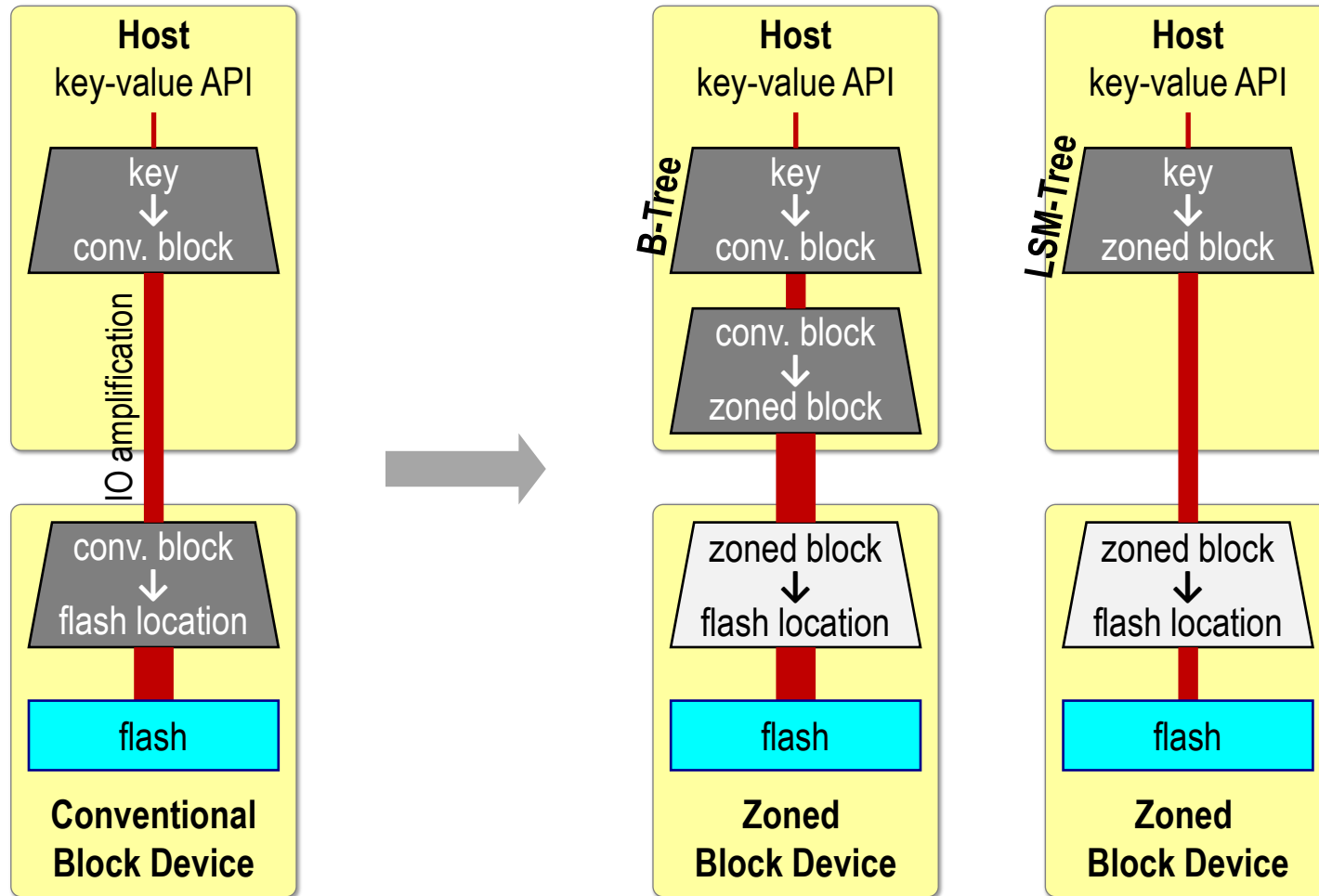
Umesh Maheshwari
Chiku Research

HotStorage 2021

# Storage Abstractions



Conventional Block Storage

1

Zoned Block Storage

2

Zoned Rock Storage

# Why Zoned Storage

# Why Rocks in Zoned Storage

## The Downside

Add little:

- complexity in specification

- overhead in implementation

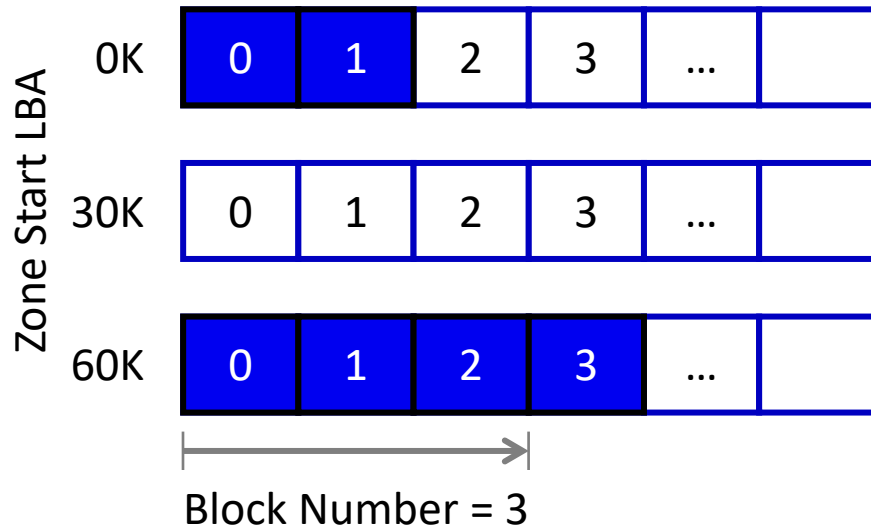## The Upside

Store small/variable-size data efficiently:
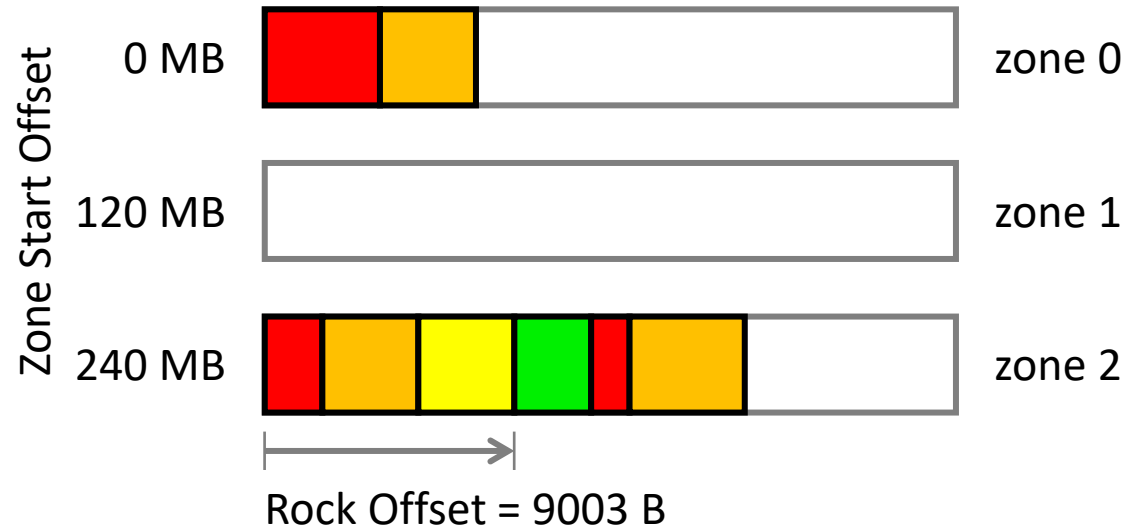
- compressed pages

- log records

# Potential Specification

## Zoned Block Namespace (ZBNS)

Zone Start LBA

| 0K | 0 | 1 | 2 | 3 | ... | |

| 30K | 0 | 1 | 2 | 3 | ... | |

| 60K | 0 | 1 | 2 | 3 | ... | |

Block Number = 3

Block Address = Zone Start LBA + Block Number = 60K+3

## Zoned Rock Namespace (ZRNS)

Zone Start Offset

| 0 MB | | | | zone 0 |

| 120 MB | | zone 1 |

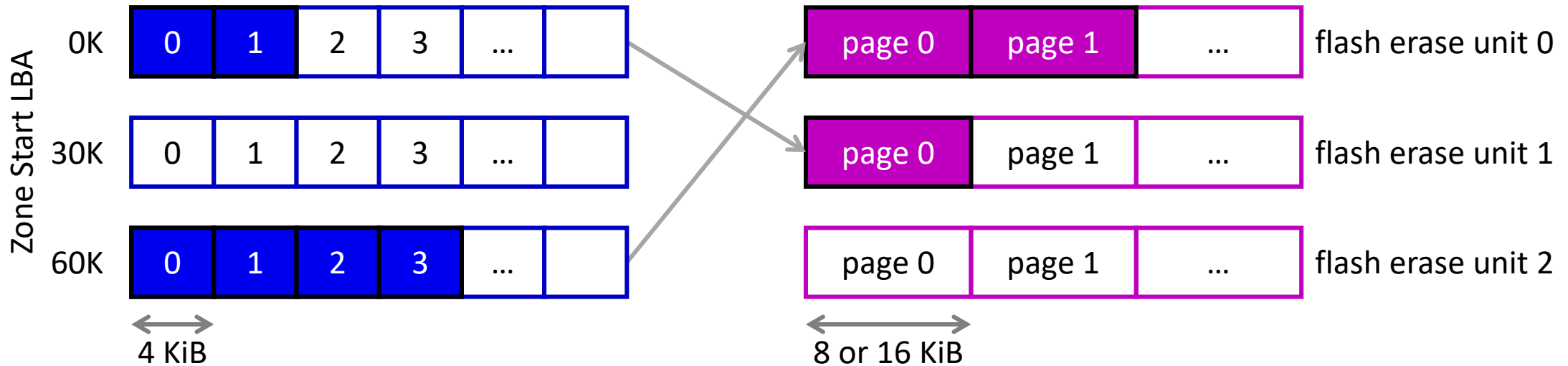| 240 MB | | | | | | | | zone 2 |

Rock Offset = 9003 B

Rock Address = Zone Start Offset + Rock Offset = 240M + 9003

Rock Address = (Zone Number, Rock Offset) = (2, 9003)

# Potential Implementation

## Zoned Block Namespace (ZBNS)
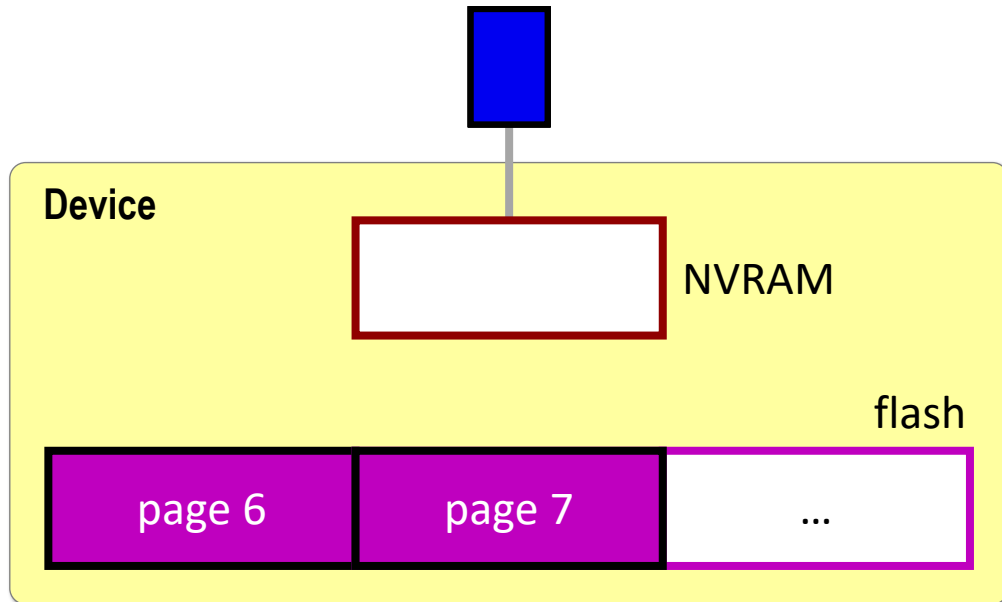
## Physical Locations in Flash



Zone Start LBA

0K | 0 | 1 | 2 | 3 | ... |

30K | 0 | 1 | 2 | 3 | ... |

60K | 0 | 1 | 2 | 3 | ... |

4 KiB

page 0 | page 1 | ... | flash erase unit 0

page 0 | page 1 | ... | flash erase unit 1

page 0 | page 1 | ... | flash erase unit 2

8 or 16 KiB

😊 Per-zone map, no per-block map.

☹️ Logical block size < flash page size

# Potential Implementation

## Zoned Block Namespace (ZBNS)

**Device**

NVRAM

flash

| page 6 | page 7 | ... |
|--------|--------|-----|

- 😊 One page-size NVRAM buffer per active zone
- 😊 No per-block map
- 😊 One command can read/write many blocks

## Zoned Rock Namespace (ZRNS)

**Device**

NVRAM

flash

| page 6 | page 7 | ... |
|--------|--------|-----|

- 😊 One page-size NVRAM buffer per active zone
- 😊 No per-rock map
- 😊 One command can read/write many rocks
- 😊 Can support rocks as small as 16 B

# Why Rocks in Zoned Storage

**The Downside**

Adds little:

- complexity in specification
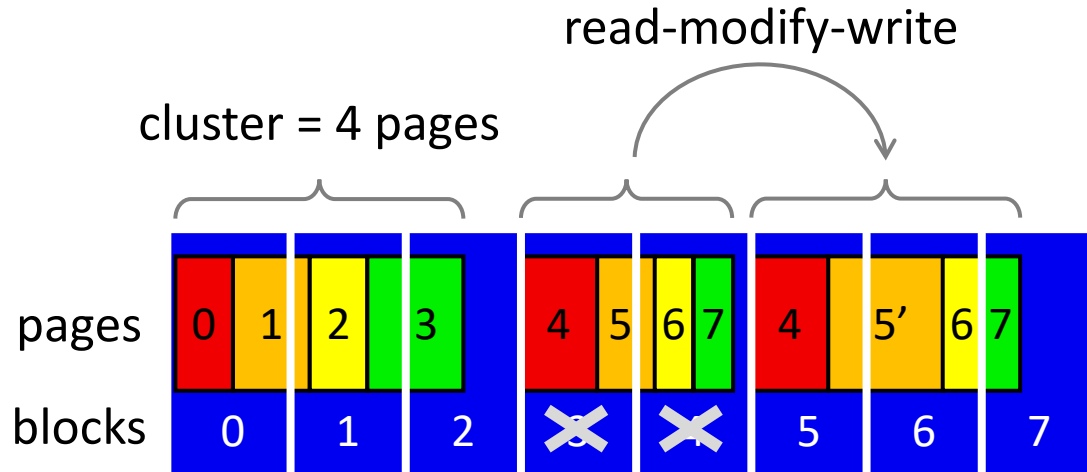
- overhead in implementation

**The Upside**

Store small/variable-size data efficiently:
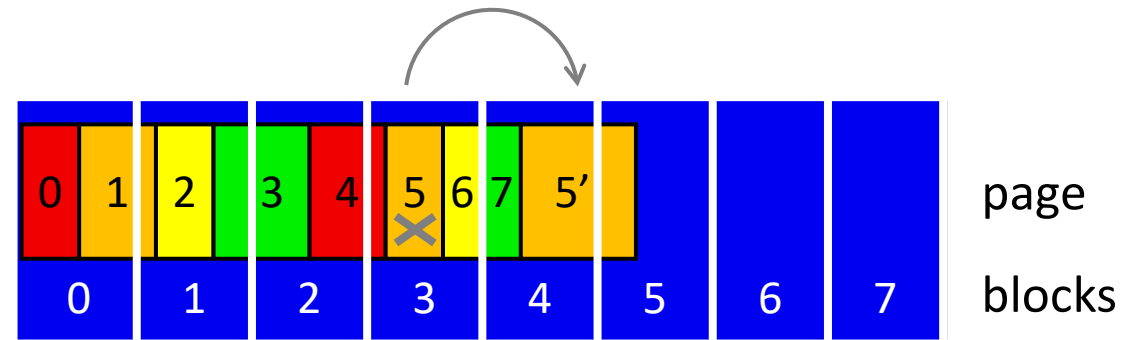
- compressed pages

- log records

# Transparent Compression

## 1. Block-Aligned Clusters

read-modify-write

cluster = 4 pages

pages

| 0 | 1 | 2 | 3 | | 4 | 5 | 6 | 7 | | 4 | 5' | 6 | 7 |

blocks: 0  1  2  ✗  ✗  5  6  7

Unit of garbage = block
e.g., WAFL®, Btrfs, F2FS

## 2. Soft Rocks Over Blocks

pages

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 5' | |

page

blocks: 0  1  2  3  4  5  6  7

blocks

Unit of garbage = rock
e.g., CASL®

# Transparent Compression

## 3. Device-Level Rocks

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 5' |

rocks

😊 Avoid reading extra bytes from device.

😊 Avoid redundant checksums on rocks and blocks.

😊 Offload to device: rock-level copy to optimize GC.

## 2. Soft Rocks Over Blocks

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 5' | | | page
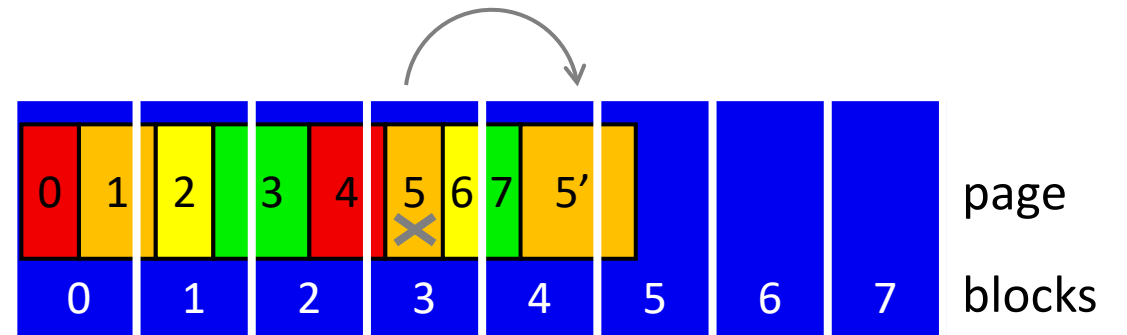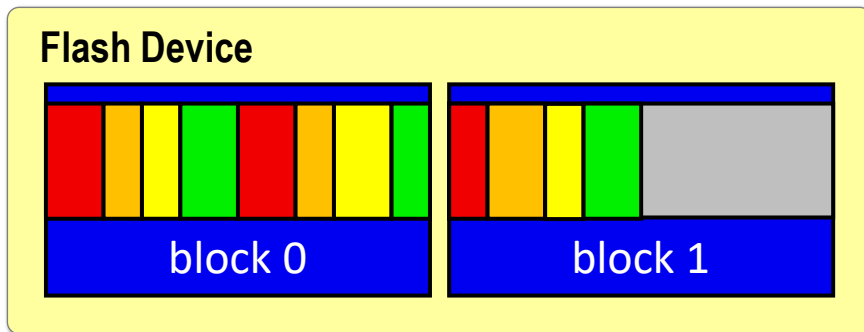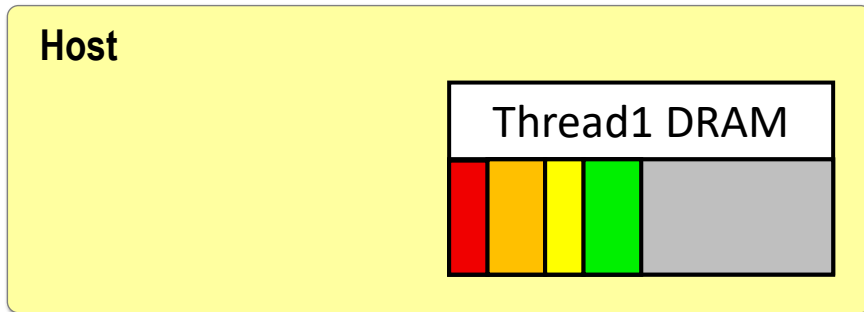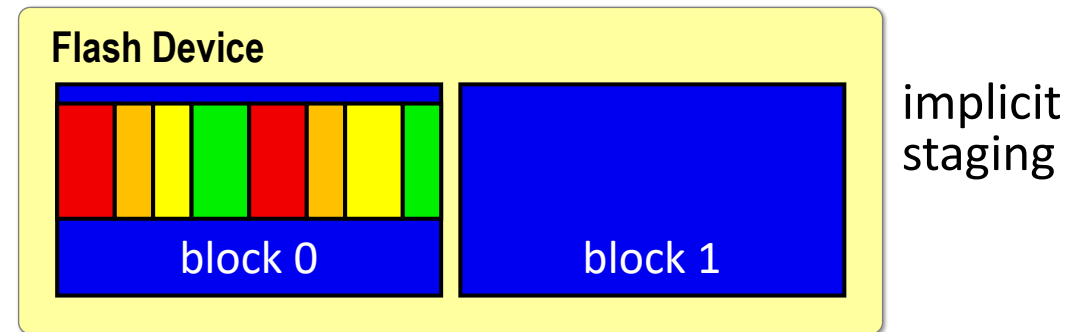
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | blocks

Unit of garbage = rock
e.g., CASL®

# Logging Change Records

## 1. Rewrite Last Block

**Host**

Thread1 DRAM

**Flash Device**

block 0

block 1

☹ Requires extension ZRWA, not yet standard.

☹ Can have at most 1 write pending to a block.

☹ Amplifies bytes transferred.

## 2. Stage in NVRAM

**Host**

NVRAM

explicit staging

**Flash Device**

block 0

block 1

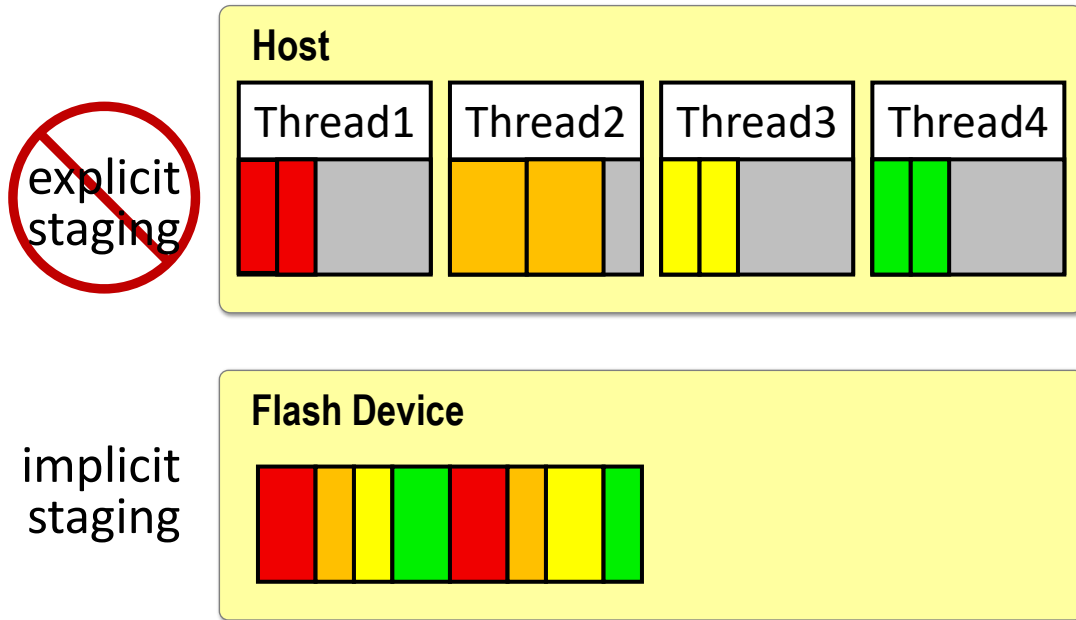implicit staging

☹ Explicit staging adds cost and complexity.

☹ Need for separate replication makes it worse.

# Logging Change Records

## 3. Device-Level Rocks



🚫 explicit staging

**Host**

| Thread1 | Thread2 | Thread3 | Thread4 |

implicit staging

**Flash Device**

Direct logging without (explicit) staging.

Concurrent appends by multiple threads.

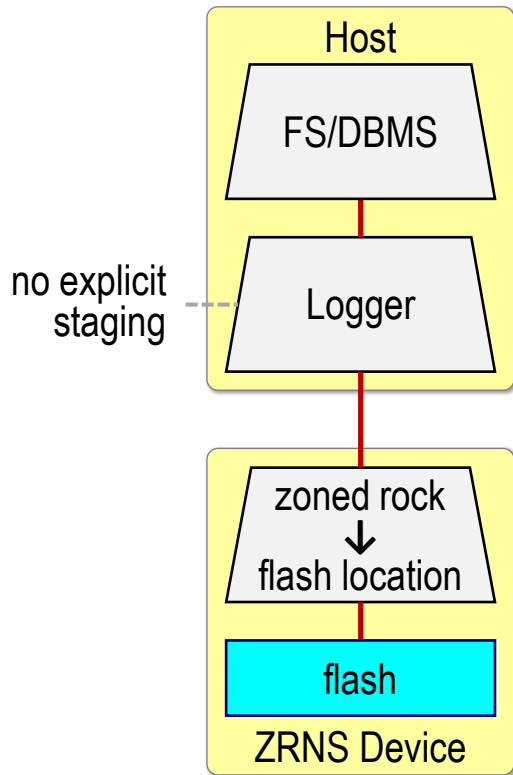As fast as explicit staging in PCIe-attached NVRAM.
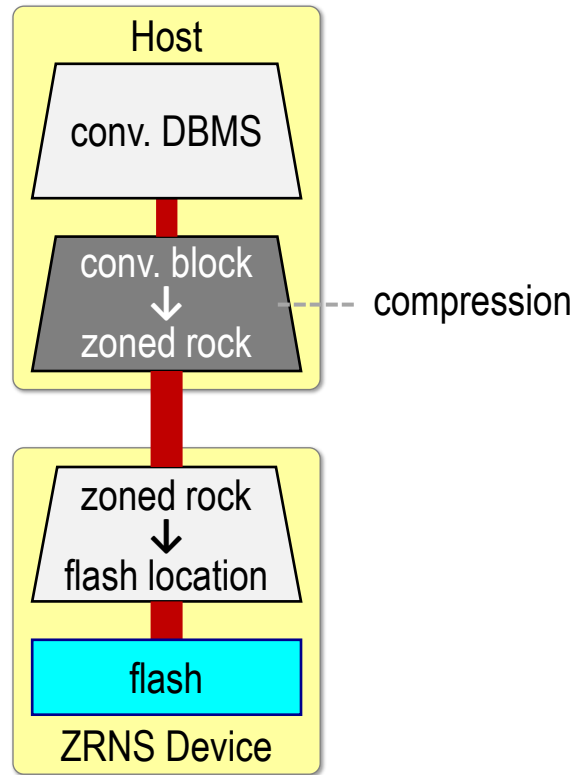
## 2. Stage in NVRAM

**Host**

NVRAM

explicit staging

**Flash Device**

implicit staging

block 0    block 1

😞 Explicit staging adds cost and complexity.

😞 Need for separate replication makes it worse.

# Future Directions

## 1. Direct Logging

| Host |
|---|
| FS/DBMS |
| Logger |

no explicit staging

| ZRNS Device |
|---|
| zoned rock ↓ flash location |
| flash |

## 2. Transparent Compression

| Host |
|---|
| conv. DBMS |
| conv. block ↓ zoned rock |

compression

| ZRNS Device |
|---|
| zoned rock ↓ flash location |
| flash |

## 3. Zoned KV Namespace (aka SSTables)

| Host |
|---|
| LSM Tree |
| zoned key ↓ zoned rock |

merge policy

zone compaction
key lookup within zone

| ZKVNS Device |
|---|
| zoned rock ↓ flash location |
| flash |

# Conclusions

1. ZNS has the potential to become a dominant abstraction:
   a. Helps avoid an un-necessary translation.
   b. Supports systems with different data layouts.

2. ZNS can be extended to support rocks (ZRNS) with little cost:
   a. Specification: command set similar to blocks.
   b. Implementation: needs same (small) amount of internal NVRAM.

3. ZRNS provides significant benefits:
   a. Store small/variable size data efficiently: inodes, small files, compressed data.
   b. Append log records concurrently without explicit staging in NVRAM.

4. ZRNS enables further extensions:
   a. Zoned key-value records for offloading merging in LSM Trees.
   b. Other domain-specific formats and functions?

Please send questions/suggestions to umesh at alum.mit.edu.