# PHISON

**Basic**

# FW Concepts

Lee Hao Zhi

PHISON

# AGENDA

Introduction

Block Base Mapping

Page Base Mapping

Garbage Collection

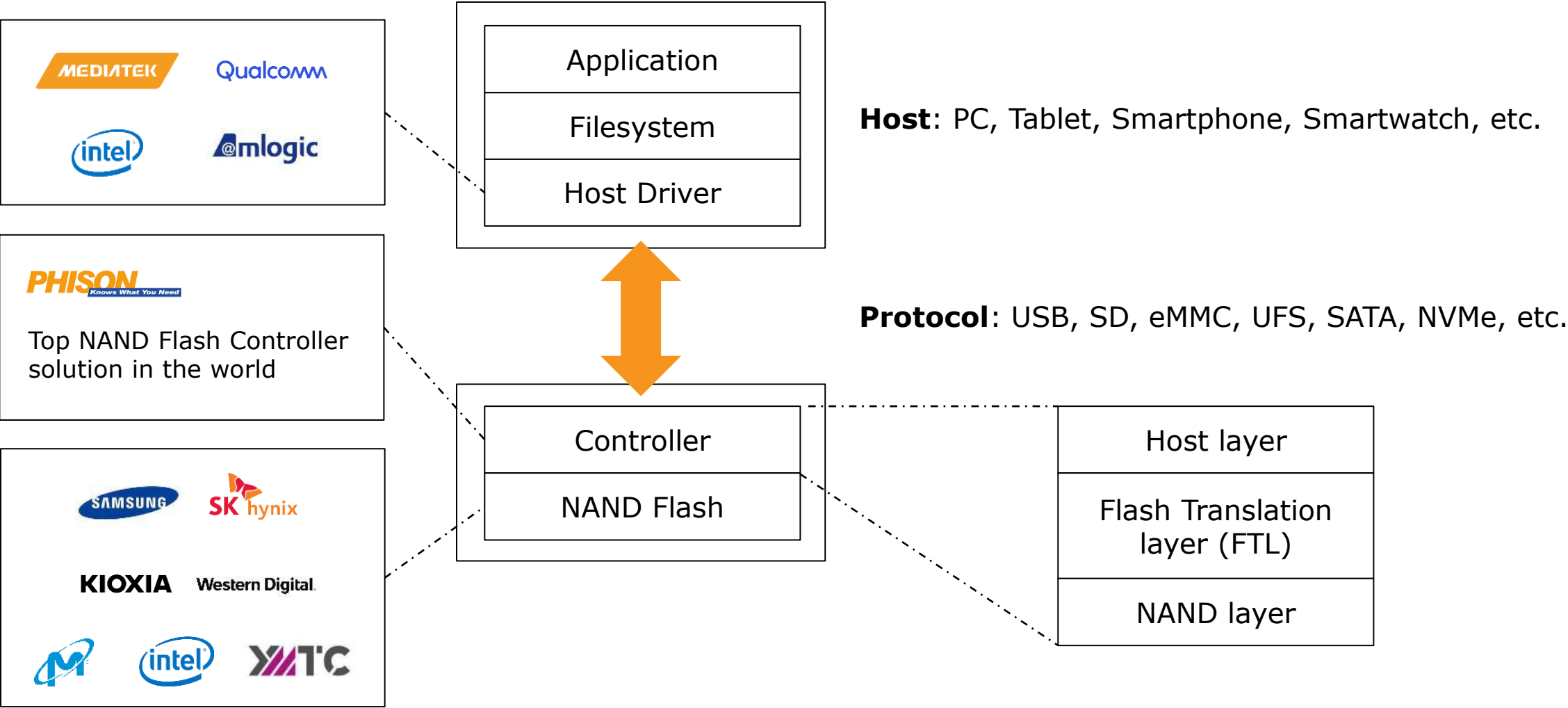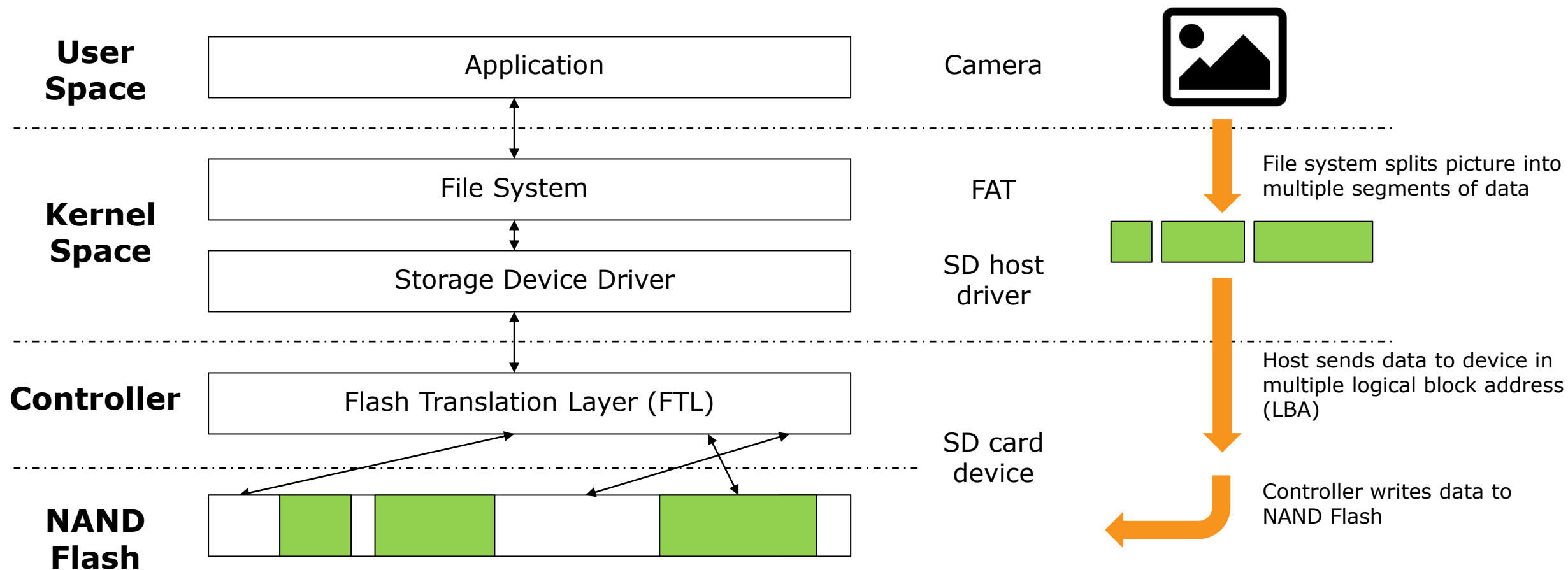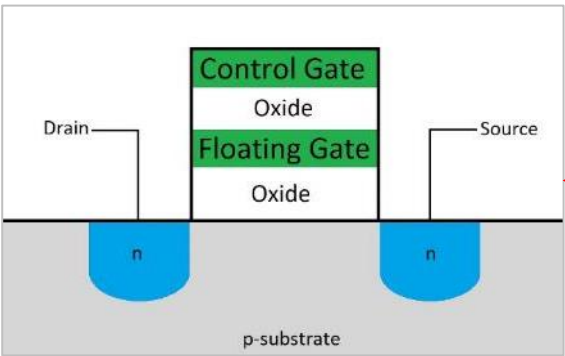Write Amplification Factor

Over Provision

Summary

# Introduction

PHISON

# Overview



**Host**: PC, Tablet, Smartphone, Smartwatch, etc.

**Protocol**: USB, SD, eMMC, UFS, SATA, NVMe, etc.

Top NAND Flash Controller solution in the world

| |
|---|
| Application |
| Filesystem |
| Host Driver |

| |
|---|
| Controller |
| NAND Flash |

| |
|---|
| Host layer |
| Flash Translation layer (FTL) |
| NAND layer |

# Application & Flash Storage Device

| User Space | Application | Camera |

| Kernel Space | File System | FAT |
| | Storage Device Driver | SD host driver |

| Controller | Flash Translation Layer (FTL) | SD card device |

| NAND Flash | | |

File system splits picture into multiple segments of data

Host sends data to device in multiple logical block address (LBA)

Controller writes data to NAND Flash

PHISON

# NAND Flash Structure

Device — Die(LUN) — Plane — Block — Page

User Area | Spare Area

Logical Unit (LUN)

18,592 bytes ←→ 18,592 bytes

Cache Registers | 16,384 | 2208 | 16,384 | 2208 | → DQ7 → DQ0

Data Registers | 16,384 | 2208 | 16,384 | 2208

504 blocks per plane

1008 blocks per LUN

1 Block | 1 Block

1 page = (16K + 2208 bytes)

1 block = (16K + 2208) bytes x 2304 pages
= (36,864K + 4968K) bytes

1 plane = (36,864K + 4968K) bytes x 504 blocks
= 168,666Mb

1 LUN = 168,666Mb x 2 planes
= 337,332Mb

Plane 0
(0, 2, 4, ..., 1006)

Plane 1
(1, 3, 5, ..., 1007)

Control Gate
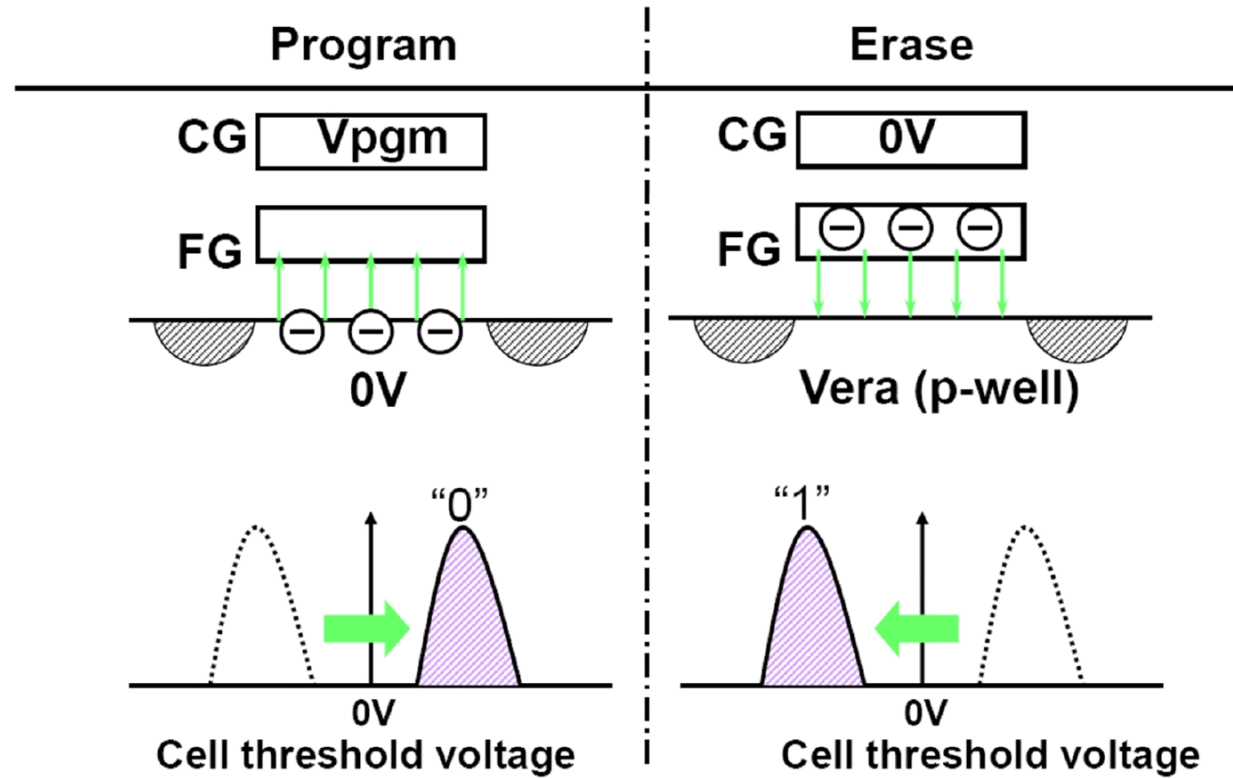Oxide
Floating Gate
Oxide

Drain — Source
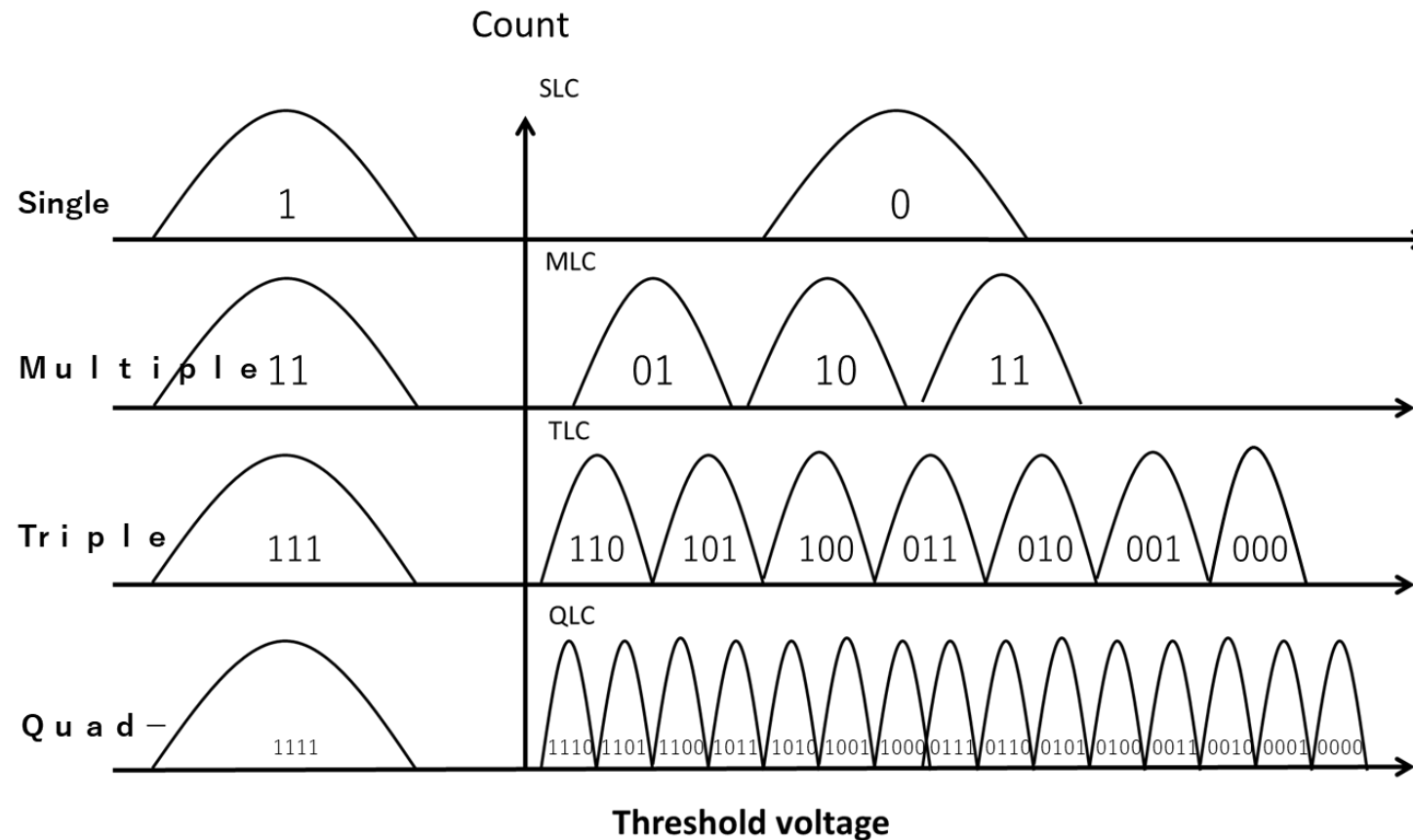
n | n

p-substrate

*Memory Cell*

PHISON

# Program & Erase (P/E) Operation



Principle of NAND memory programming and erasing

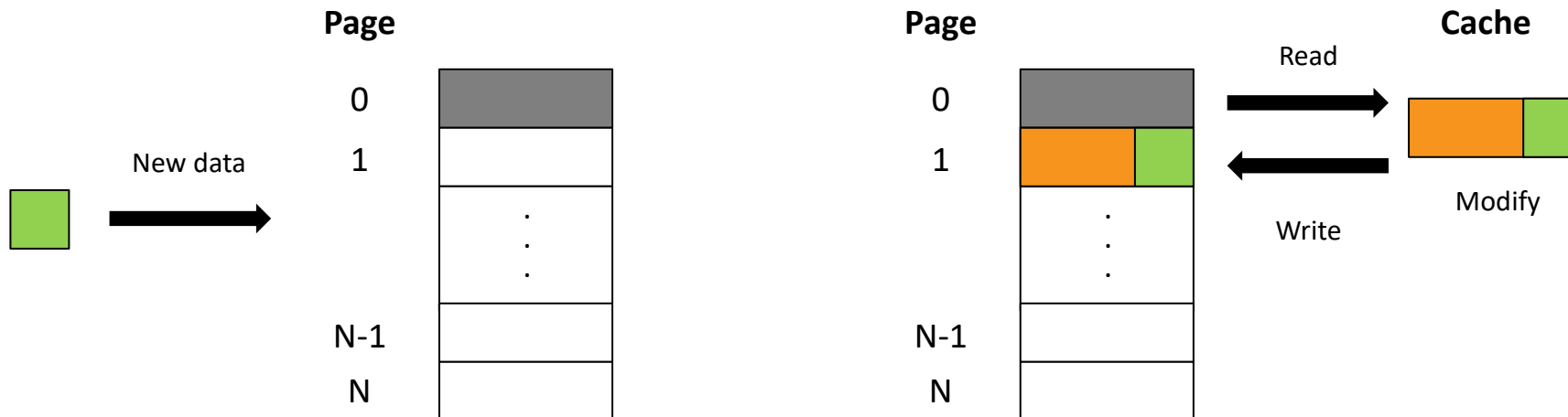# SLC vs MLC vs TLC

- Normal Voltage Threshold (VT) distribution
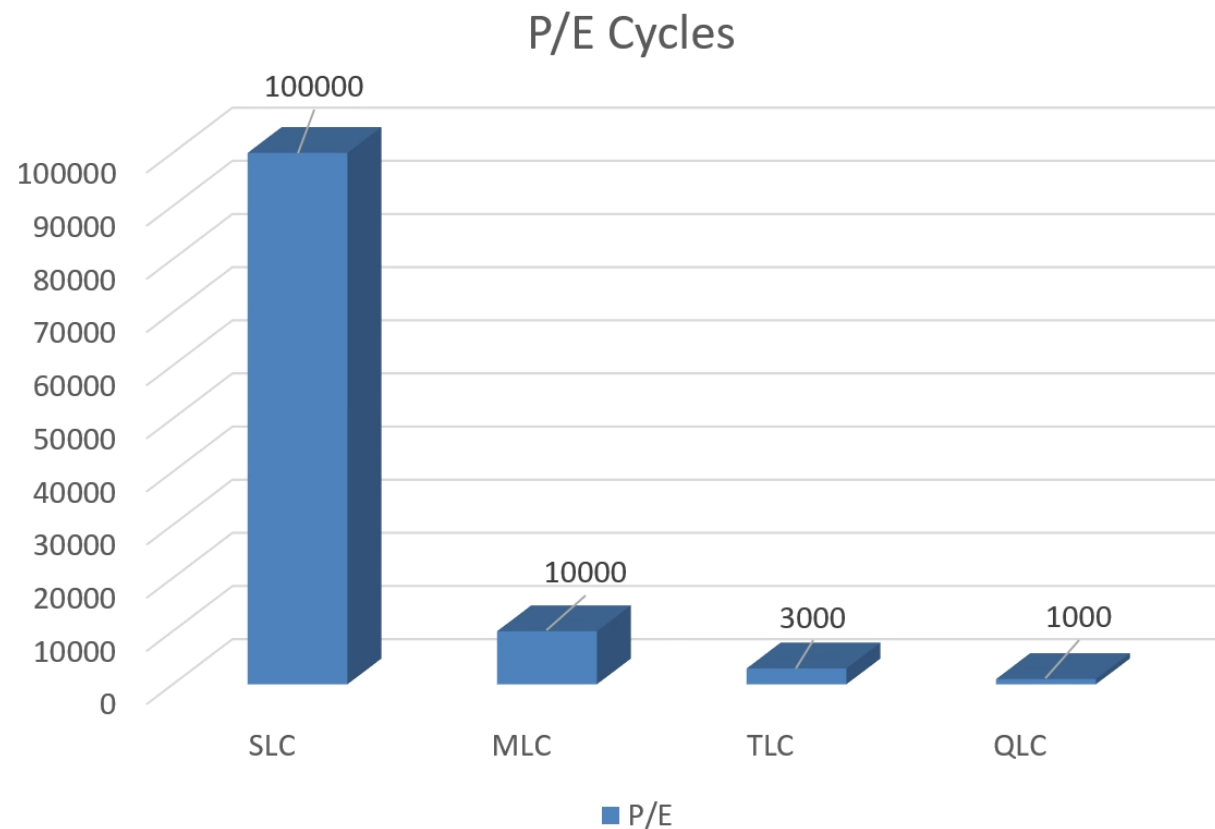
# NAND Flash Limitations

- **Page** is the smallest unit for read and program
- **Block** is the smallest unit for erase
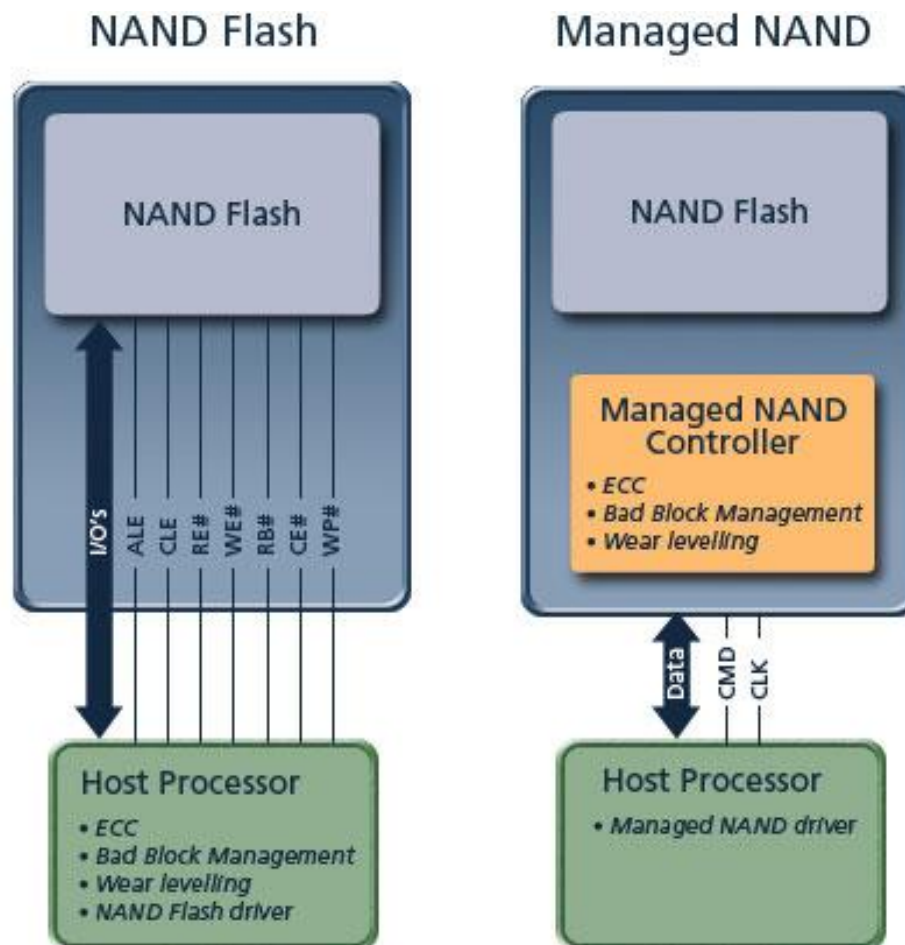- Must erase before program (cannot overwrite)



- Each read/program/erase operation has busy time to complete

# NAND Flash Limitations

- Program/Erase (PE) Cycle: SLC = 100k, MLC = 10k, TLC = 3k
- Initial and runtime bad blocks
- Read errors & disturbance

**P/E Cycles**

PHISON

# NAND Flash Controller
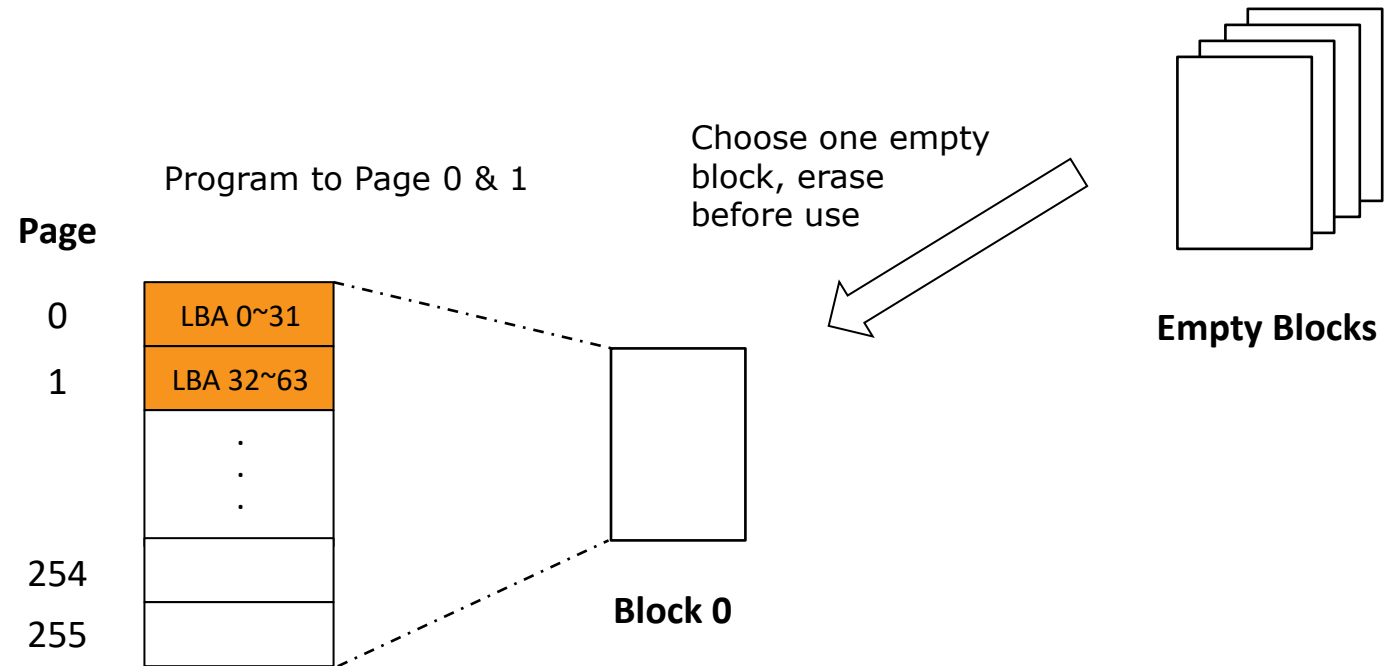
# Block Base Mapping

PHISON

# Definitions & Rules

- 8GB NAND Flash
  - 1 die = 2048 blocks
  - 1 block = 256 pages
  - 1 page = 16KB = 32 sectors
  - 1 LBA/sector = 512B

- Limitation
  - **Page** is the smallest unit for read and program
  - **Block** is the smallest unit for erase
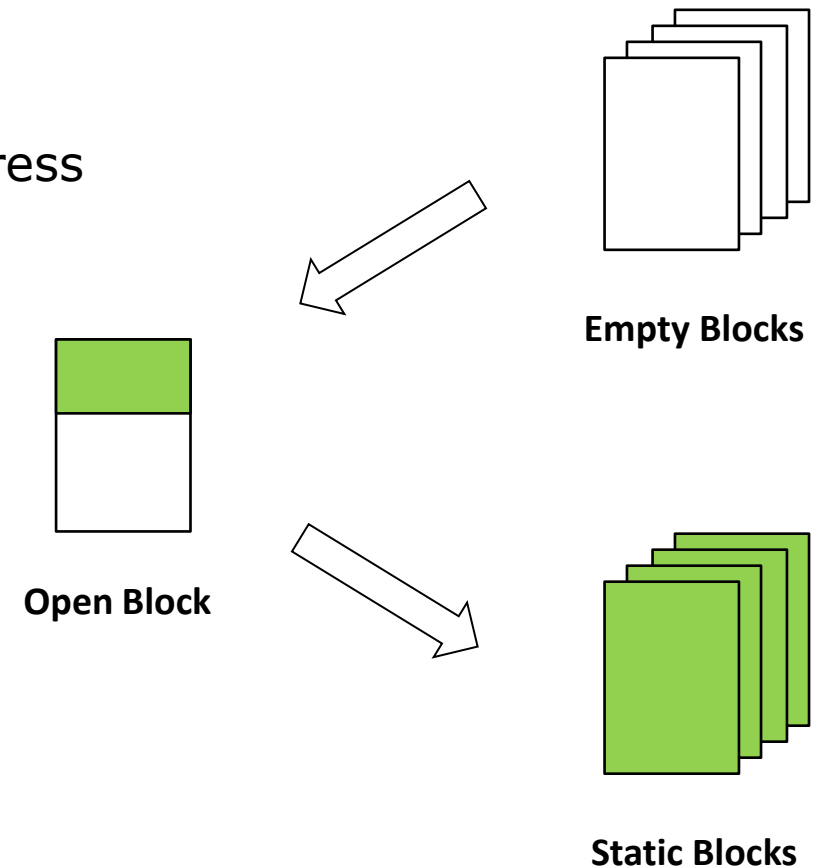  - Must erase before program (cannot overwrite)

# Write

- Host writes 2 chunk of data (16KB each)
  - LBA 0~31
  - LBA 32~63

Program to Page 0 & 1

Choose one empty block, erase before use

**Empty Blocks**

**Page**

| | |
|---|---|
| 0 | LBA 0~31 |
| 1 | LBA 32~63 |
| | . |
| | . |
| | . |
| 254 | |
| 255 | |

**Block 0**

# Write

- Host continues to write large chunk size data…
- Open block becomes static when full
- We record logical to physical block mapping
- Each Logical Block (LB) index maps to Block (PB) address

| LB | PB |
|------|--------|
| 0 | 0 |
| 1 | 1 |
| . | |
| . | |
| . | |
| 2046 | 0xFFFF |
| 2047 | 0xFFFF |

**Empty Blocks**

**Open Block**

**Static Blocks**

- Address = (PB)
- Might need 2B to store address

C O N F I D E N T I A L

# Read

- Host wants to read LBA 32~63
- How do we know where is the data located in NAND Flash?

# Read

- Logical to Physical (L2P) translation

    Step 1: Calculate Logical Block (LB) Index

    $$LB = \frac{LBA}{Total\ Physical\ Sector\ Number\ (in\ block)}$$

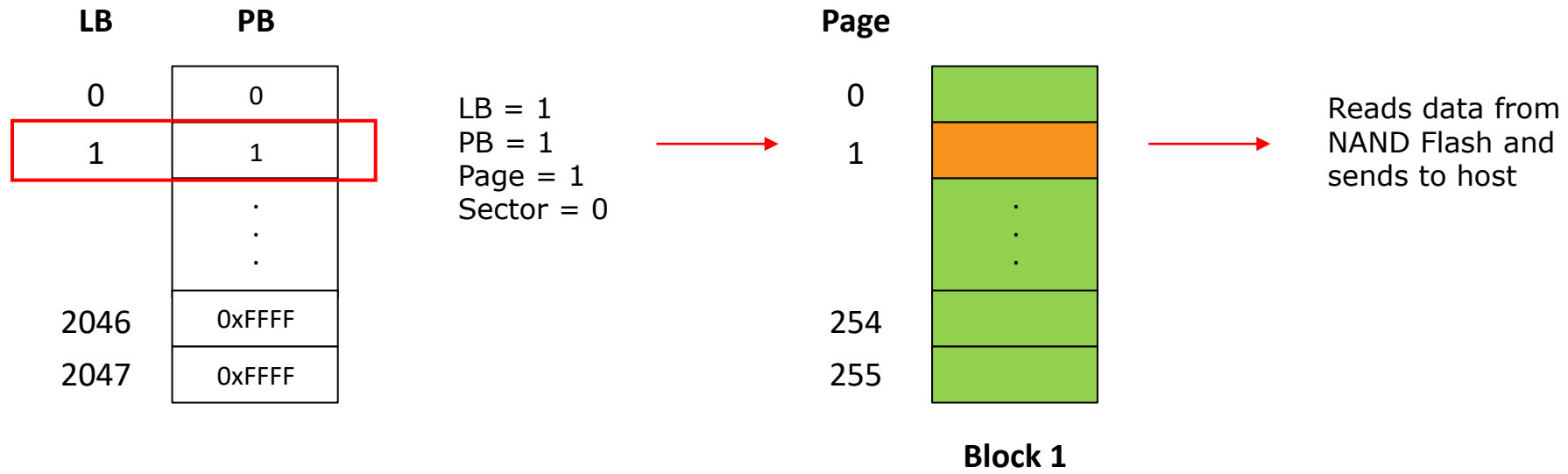    Step 2: Get Physical Block (PB) address from mapping table

    Step 3: Calculate Offset

    $$Page = \frac{LBA\ \%\ Total\ Physical\ Sector\ Number\ (in\ block)}{Total\ Physical\ Sector\ Number\ (in\ page)}$$

    $$Sector = LBA\ \%\ Total\ Physical\ Sector\ Number\ (in\ block)\ \%\ Total\ Physical\ Sector\ Number\ (in\ page)$$
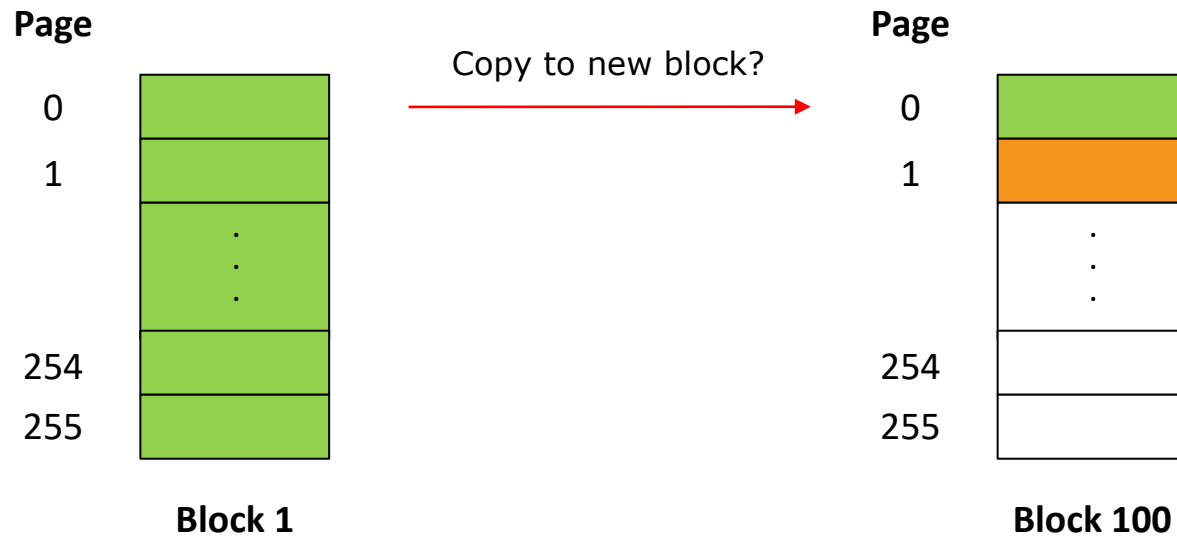
PHISON

# Read

- Host wants to read LBA 8224~8255

| LB | PB |
|----|------|
| 0 | 0 |
| 1 | 1 |
| . | . |
| . | . |
| . | . |
| 2046 | 0xFFFF |
| 2047 | 0xFFFF |

LB = 1
PB = 1
Page = 1
Sector = 0

→

| Page | Block 1 |
|------|---------|
| 0 | |
| 1 | |
| . | |
| . | |
| . | |
| 254 | |
| 255 | |

**Block 1**

→

Reads data from
NAND Flash and
sends to host

# Overwrite

- What if host overwrites same address LBA 8224~8255?
- Remember, we cannot overwrite a page that has been programmed

Page                          Copy to new block?                    Page

0                                                                    0
1                                                                    1

.                                                                    .
.                                                                    .
.                                                                    .

254                                                                  254
255                                                                  255

**Block 1**                                                          **Block 100**

- Not so efficient

**PHISON**
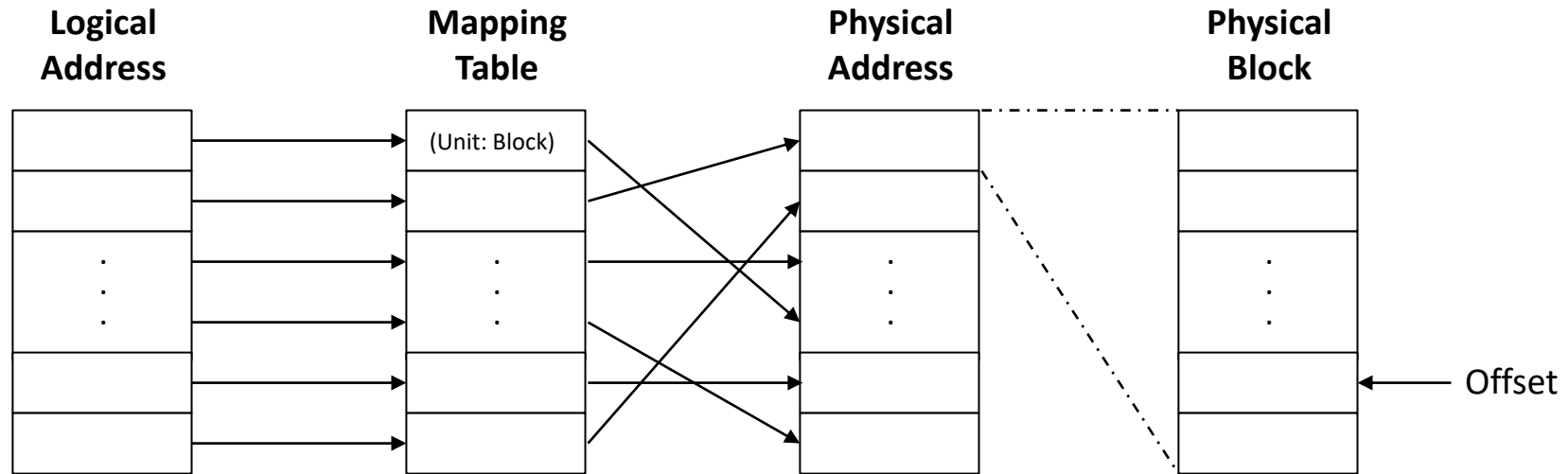
# Table Size (for Block Mapping)

- For 8GB NAND Flash, total table size required is 4KB
  - Table Size = 2048 blocks × 2B = 4KB
  - Table Entries = 2048

| LB | PB |
|----|----|
| 0 | 0 |
| 1 | 1 |
| | . . . |
| 2046 | 0xFFFF |
| 2047 | 0xFFFF |

# Block Base Mapping Table

- In a block level address mapping, a logical page address is made up of both a **logical block number** and an **offset**
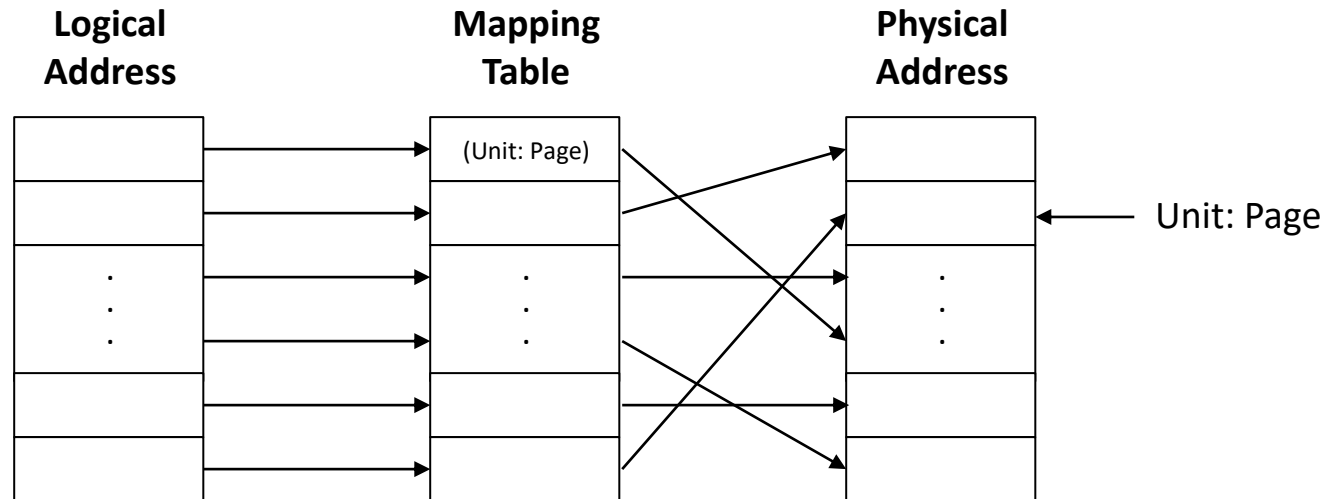
# Page Base Mapping

# Page Base Mapping Table

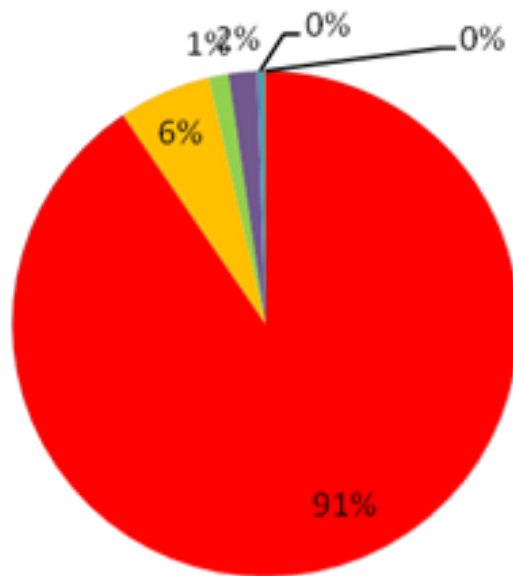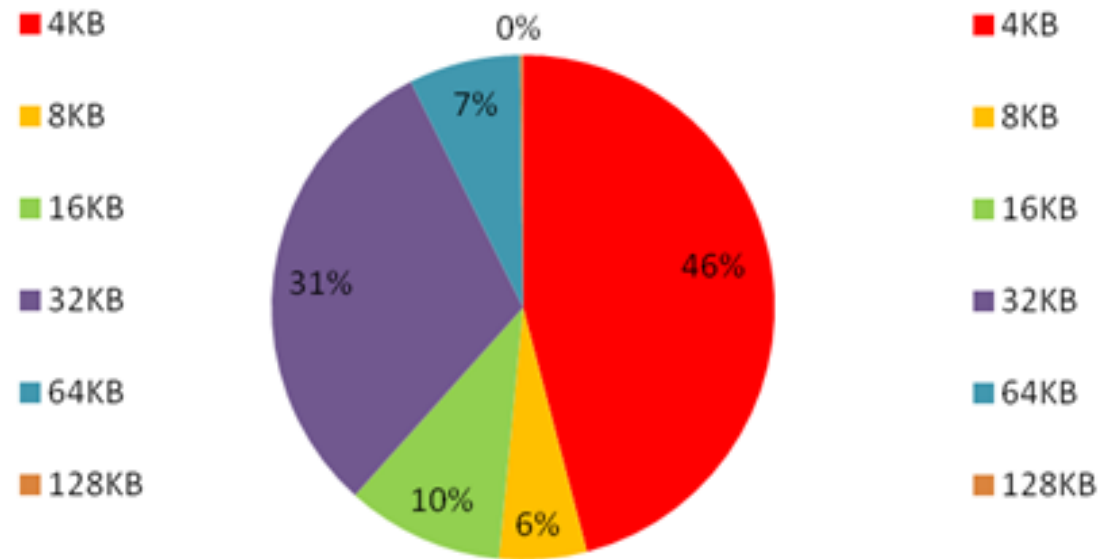- In a page level address mapping, a logical page can be **mapped into any physical page** in flash memory

# Page Base Mapping Table

- Data structure of host (eMMC) operations on WHCK performance test
- 4KB chunk size has higher percentage than other chunk size
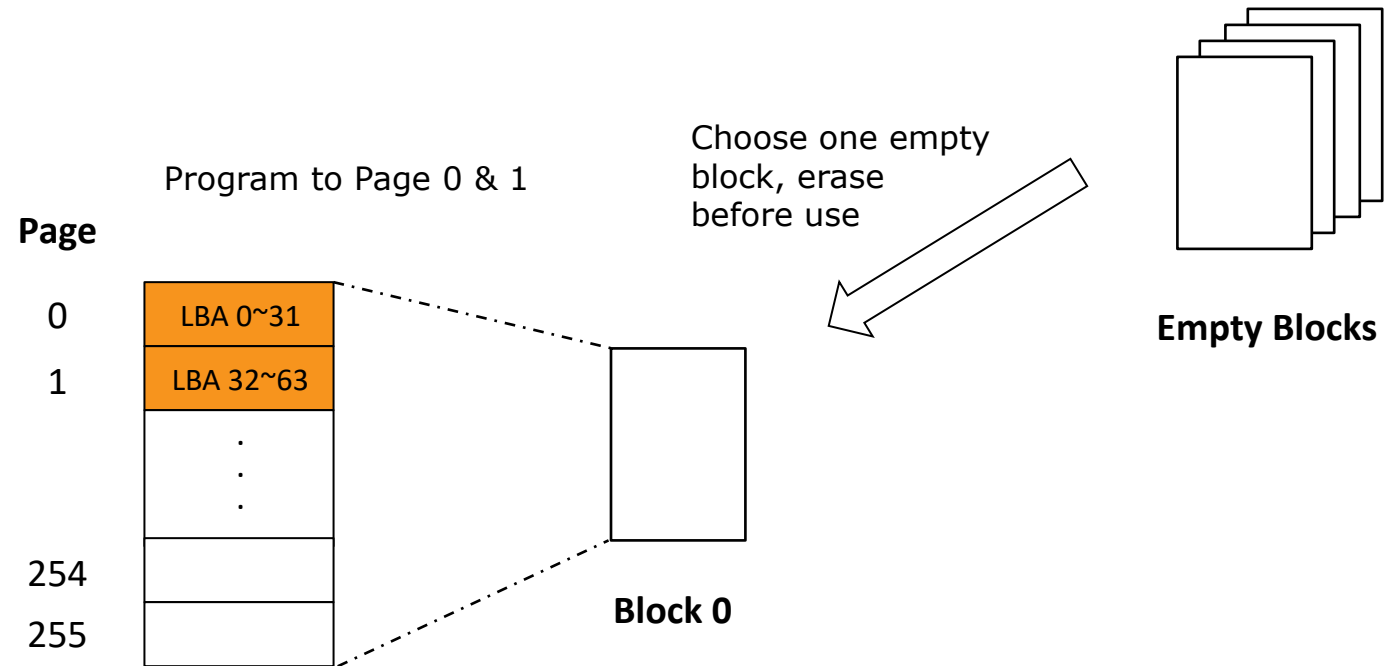


**Write Chunk Size**

1% 2% 0% 0%
6%
91%

Legend: 4KB, 8KB, 16KB, 32KB, 64KB, 128KB

**Read Chunk Size**

0%
7%
31%
46%
10%
6%

Legend: 4KB, 8KB, 16KB, 32KB, 64KB, 128KB

PHISON

# Write

- Host writes 2 chunk of data (16KB each)
  - LBA 0~31
  - LBA 32~63

Program to Page 0 & 1

Choose one empty block, erase before use

**Empty Blocks**

**Page**

| | |
|---|---|
| 0 | LBA 0~31 |
| 1 | LBA 32~63 |
| | . |
| | . |
| | . |
| 254 | |
| 255 | |

**Block 0**

**PHISON**

# Write

- Instead of block mapping, we now record logical to physical page mapping
- Each Logical Page (LP) index maps to Physical Page (PP) address

| LP | PP |
|---|---|
| 0 | 0, 0, 0 |
| 1 | 0, 1, 0 |
| . . . | |
| 542286 | 0xFFFFFFFF |
| 542287 | 0xFFFFFFFF |

- Address = (Block, Page, Sector)
- Might need 4B (instead of 2B) to store address

PHISON

# Write

- Host writes additional 2 chunk of data (16KB each)
  - LBA 64~95
  - LBA 32~63 (overwrite)

| Page | | | LP | PP |
|---|---|---|---|---|
| 0 | LBA 0~31 | | 0 | 0, 0, 0 |
| 1 | LBA 32~63 | No need to copy data to new block during overwrite, old data becomes invalid | 1 | 0, 3, 0 |
| 2 | LBA 64~95 | | 2 | 0, 2, 0 |
| 3 | LBA 32~63 | | 3 | 0xFFFFFFFF |
| | . . . | | | . . . |
| 254 | | | 542286 | 0xFFFFFFFF |
| 255 | | | 542287 | 0xFFFFFFFF |

**Block 0**

# Read

- Host wants to read LBA 32~63
- How do we know where is the data located in NAND Flash?

- Logical to Physical (L2P) translation

    Step 1: Calculate Logical Page (LP) Index

$$LP = \frac{LBA}{Total\ Physical\ Sector\ Number\ (in\ page)}$$

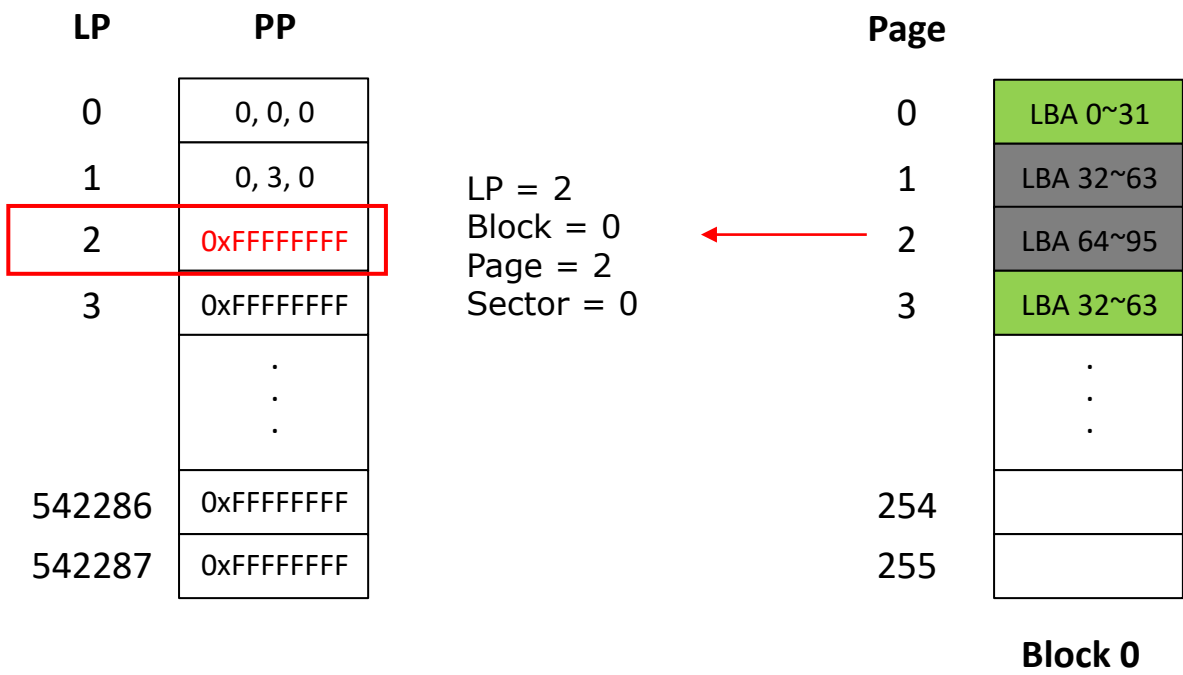    Step 2: Get Physical Page (PP) address from mapping table

# Read

- Host wants to read LBA 32~63

| LP | PP |
|---|---|
| 0 | 0, 0, 0 |
| 1 | 0, 3, 0 |
| 2 | 0, 2, 0 |
| 3 | 0xFFFFFFFF |
| . . . | . . . |
| 542286 | 0xFFFFFFFF |
| 542287 | 0xFFFFFFFF |

LP = 1
Block = 0
Page = 3
Sector = 0

| Page | |
|---|---|
| 0 | LBA 0~31 |
| 1 | LBA 32~63 |
| 2 | LBA 64~95 |
| 3 | LBA 32~63 |
| . . . | . . . |
| 254 | |
| 255 | |

**Block 0**

Reads data from NAND Flash and sends to host

# Erase

- Host wants to erase LBA 64~95

| LP | PP |
|---|---|
| 0 | 0, 0, 0 |
| 1 | 0, 3, 0 |
| 2 | 0xFFFFFFFF |
| 3 | 0xFFFFFFFF |
| . . . | . . . |
| 542286 | 0xFFFFFFFF |
| 542287 | 0xFFFFFFFF |

LP = 2
Block = 0
Page = 2
Sector = 0

| Page | |
|---|---|
| 0 | LBA 0~31 |
| 1 | LBA 32~63 |
| 2 | LBA 64~95 |
| 3 | LBA 32~63 |
| . . . | . . . |
| 254 | |
| 255 | |

**Block 0**

- Mark corresponding LP as invalid

**PHISON**

# Table Size (for 16KB Page Mapping)

- For 8GB NAND Flash, total table size required is 2MB
  - Table Size = 2048 blocks × 256 pages × 4B = 2MB
  - Table Entries = 2048 blocks × 256 pages = 524288

| LP | PP |
|----|----|
| 0 | 0xFFFFFFFF |
| 1 | 0xFFFFFFFF |
| . . . | |
| 542286 | 0xFFFFFFFF |
| 542287 | 0xFFFFFFFF |

# Table Size (for 4KB Page Mapping)

- For 8GB NAND Flash, total table size required is 8MB
  - Table Size = 2048 blocks × 256 pages × 4 nodes × 4B = 8MB
  - Table Entries = 2048 blocks × 256 pages × 4 nodes= 2097152

| LP | PP |
|---|---|
| 0 | 0xFFFFFFFF |
| 1 | 0xFFFFFFFF |
|  | . . . |
| 2097150 | 0xFFFFFFFF |
| 2097151 | 0xFFFFFFFF |

# Garbage Collection (GC)

PHISON

# Garbage Collection

- What is GC and why do we need to do GC?
- Host continues to write large chunk size data…
- Open block becomes static when full
- Some physical page become invalid

**Empty Blocks**

**Open Block**

**Static Blocks**

# Garbage Collection

- What if host keeps writing until empty block runs out?
- Before empty block runs out, we should do Garbage Collection (GC)

No more ??

Empty Blocks

Open Block

Static Blocks

# Garbage Collection

- We collect valid data to new block and reclaim blocks filled with invalid data
- So that we can erase the reclaimed blocks and use them for new data

Reclaimed blocks will
be erased and 're-used'

Merge valid data to new block

# Garbage Collection

- How do we pick static blocks as source of GC?

Valid cnt = 50       Valid cnt = 30       Valid cnt = 10

- Pick the block(s) with least valid count

Valid cnt = 10       Valid cnt = 30       Valid cnt = 50

# Garbage Collection

- How do we know which data is valid or invalid in the source blocks?
- Check the mapping table, look for the entries that point to these source blocks
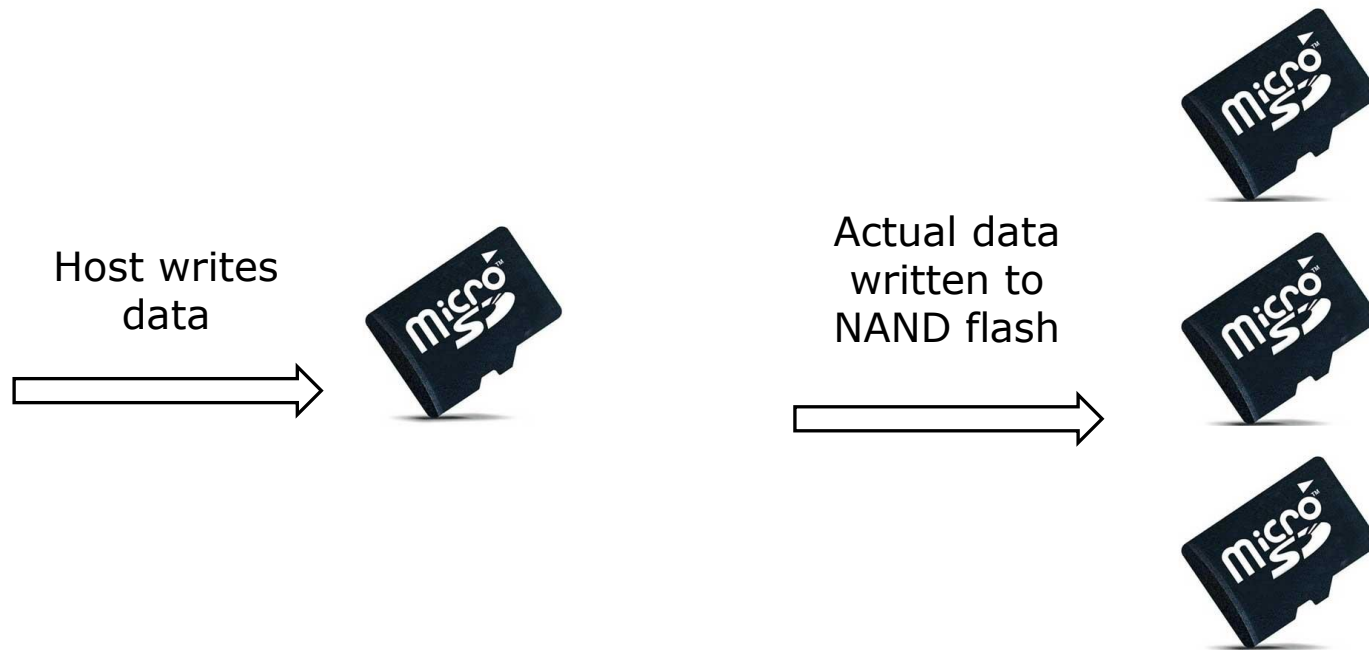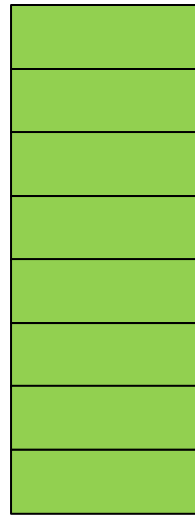
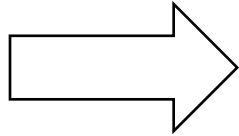# Write Amplification Factor (WAF)

# Definition

$$\text{Write Amplification Factor (WAF)} \quad = \quad \frac{\text{Data written to NAND Flash}}{\text{Data written by host}}$$

Host writes data

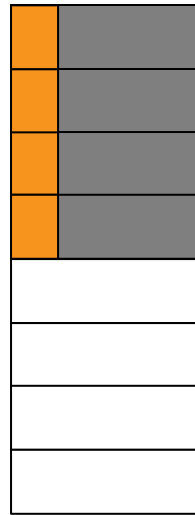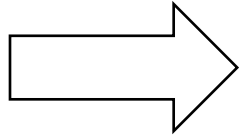Actual data written to NAND flash

PHISON

# Example

- WAF = 1

Host writes large
sequential data

# Example

- WAF = 4



Host writes 4 x 4KB

- Still remember what happens when host overwrites data of same LBA using block level mapping?

PHISON

# Summary

- Disadvantages of high WAF
    - Low performance
    - High erase count

**PHISON**

# Terabytes Written (TBW)

- Total amount of data that can be written to a storage device until it reaches its lifetime

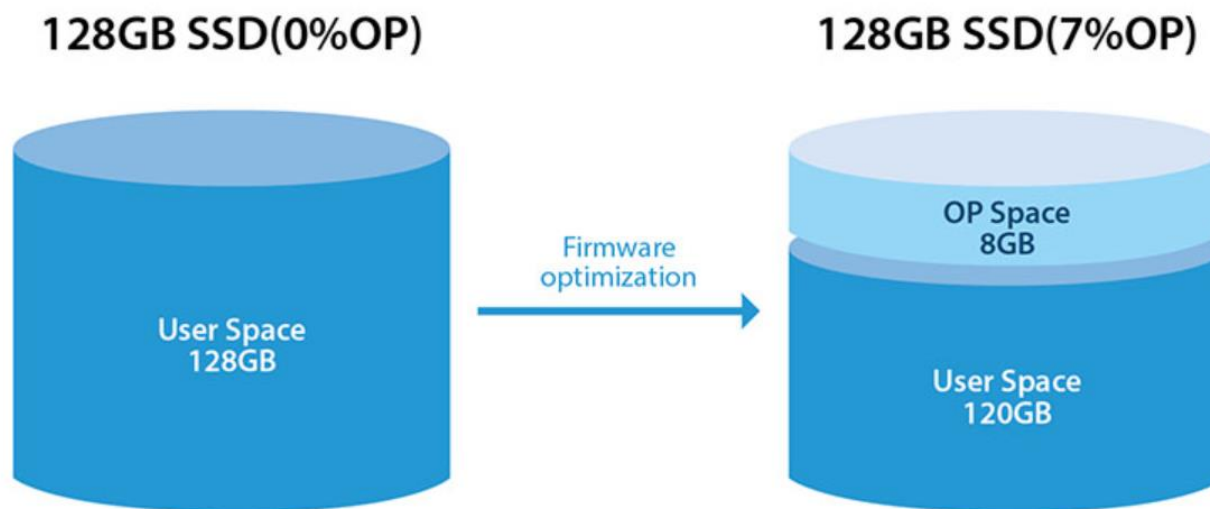$$Terabytes\ Written\ (TBW)\ =\ \frac{User\ Capacity\ (GB)\ \times NAND\ P/E\ Cycles}{WAF\ \times 1024}$$

PHISON

# Over Provision (OP)

PHISON

# Definition

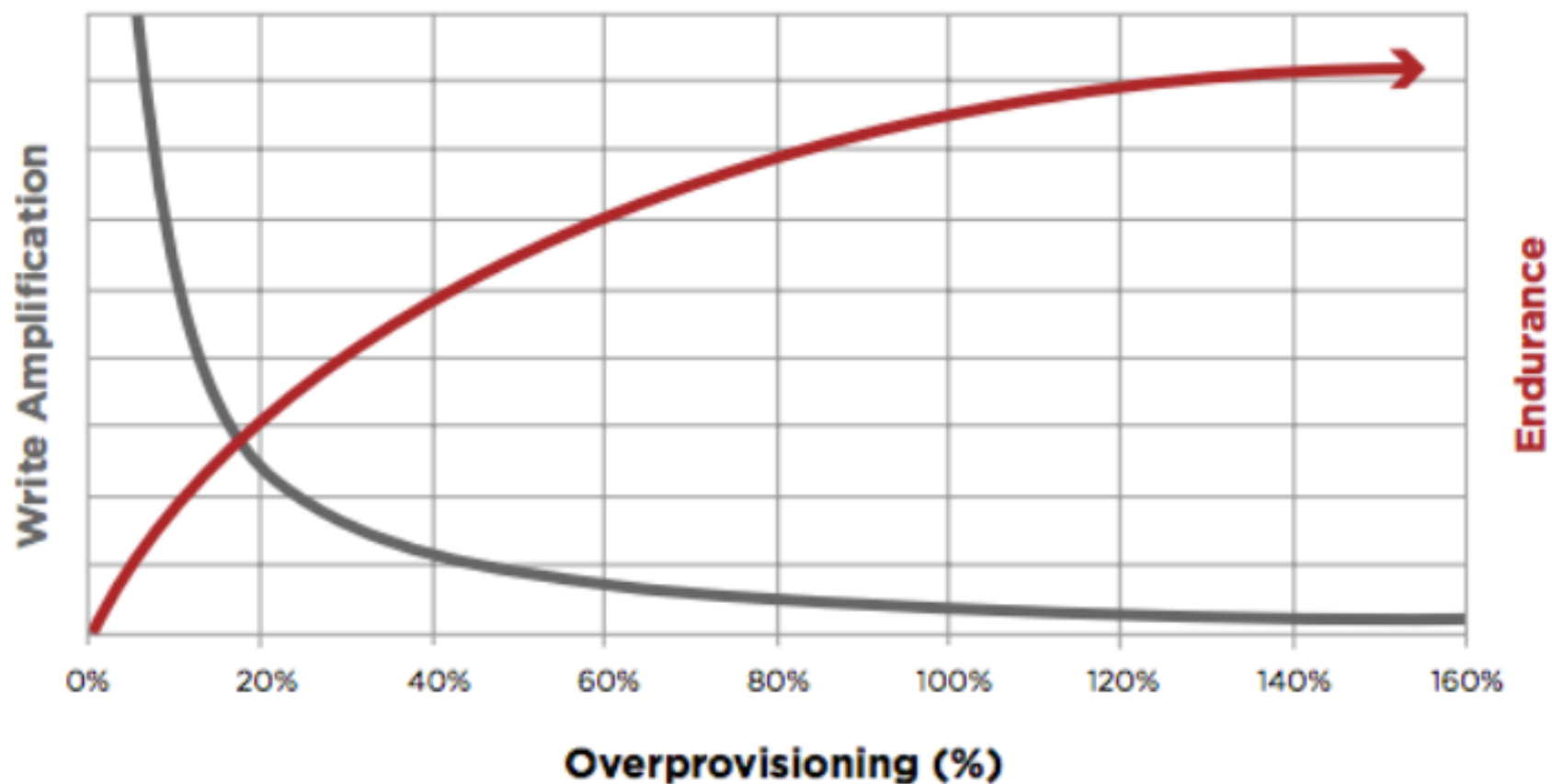$$Over\text{-}Provision \quad = \quad \frac{Physical\ capacity - User\ capacity}{User\ capacity}$$

# Example

- The larger the size of the spare area, the higher the operating efficiency, and the better the performance become

| 10 | 20 | 3  | 1  | 11 |
|----|----|----|----|----|
| 18 | 2  | 21 | 9  | 16 |
| 4  | 12 | 22 | 19 | 23 |
| 14 | 8  | 17 | 24 | 5  |
| 7  | 15 | 6  | 13 |    |

**VS**

| 10 | 20 | 3  | 1  | 11 |
|----|----|----|----|----|
| 18 | 2  | 21 | 9  | 16 |
| 4  | 12 | 5  | 19 | 13 |
| 14 | 8  | 17 |    |    |
| 7  | 15 | 6  |    |    |

The Sliding Puzzle

Figure 2: Write Amplification vs. Over Provisioning

# Summary

- Advantage(s) of higher OP due to less background data movement required
    - Higher performance
    - Better endurance (lower WAF)
- Disadvantage(s) of higher OP due to more reserved spare blocks
    - Less usable storage space

# **Summary**

# NAND Flash Limitation

- **Page** is the smallest unit for read and program
- **Block** is the smallest unit for erase
- Must erase before program (cannot overwrite)

# Page Base vs Block Base

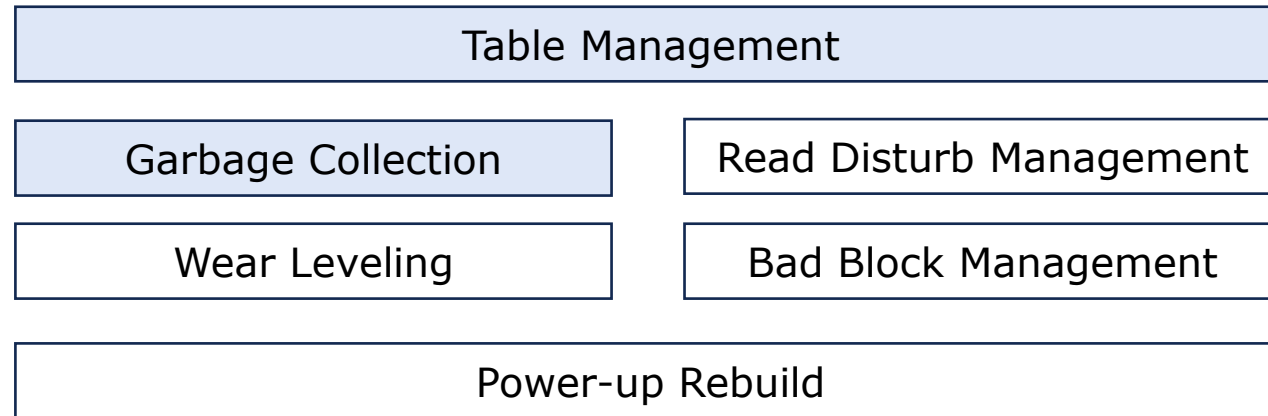| | Page Base | Block Base |
|---|---|---|
| **Mapping Unit** | Page (or 4KB) | Block |
| **Table Size** | Large (store in NAND Flash) | Small (store in SRAM) |
| **R/W Performance** | Excellent random write performance | Slow random write performance, but impressive read performance |
| **Garbage Collection** | Collect valid nodes when empty block becomes insufficient | Collect valid nodes when overwrite previous data or erase data |
| **WAF** | Low (efficient block utilization) | High (expensive merge operation) |

PHISON

# Cold Facts

- Do you know that your storage device with higher OP (less user space) has better performance and lifetime?

- Do you know that when you keep your storage usage full, the performance and lifetime become worst (due to frequent GC operation)?

- When purchasing storage device, consider performance vs lifetime (such as WAF or TBW)

# Advanced Questions

- What if power cycle occurs when we program NAND Flash?
- What if bad block occurs when we program NAND Flash?
- What if read error occurs when we read NAND Flash?
- What if certain blocks wear out quickly than other blocks?

| Table Management | |
|---|---|
| Garbage Collection | Read Disturb Management |
| Wear Leveling | Bad Block Management |
| Power-up Rebuild | |

THANK YOU!