# NAND Flash Controller

Lee Hao Zhi

2022/4/29

# 新竹交大

- Duration: 2h/week
- Time: 10am-12pm

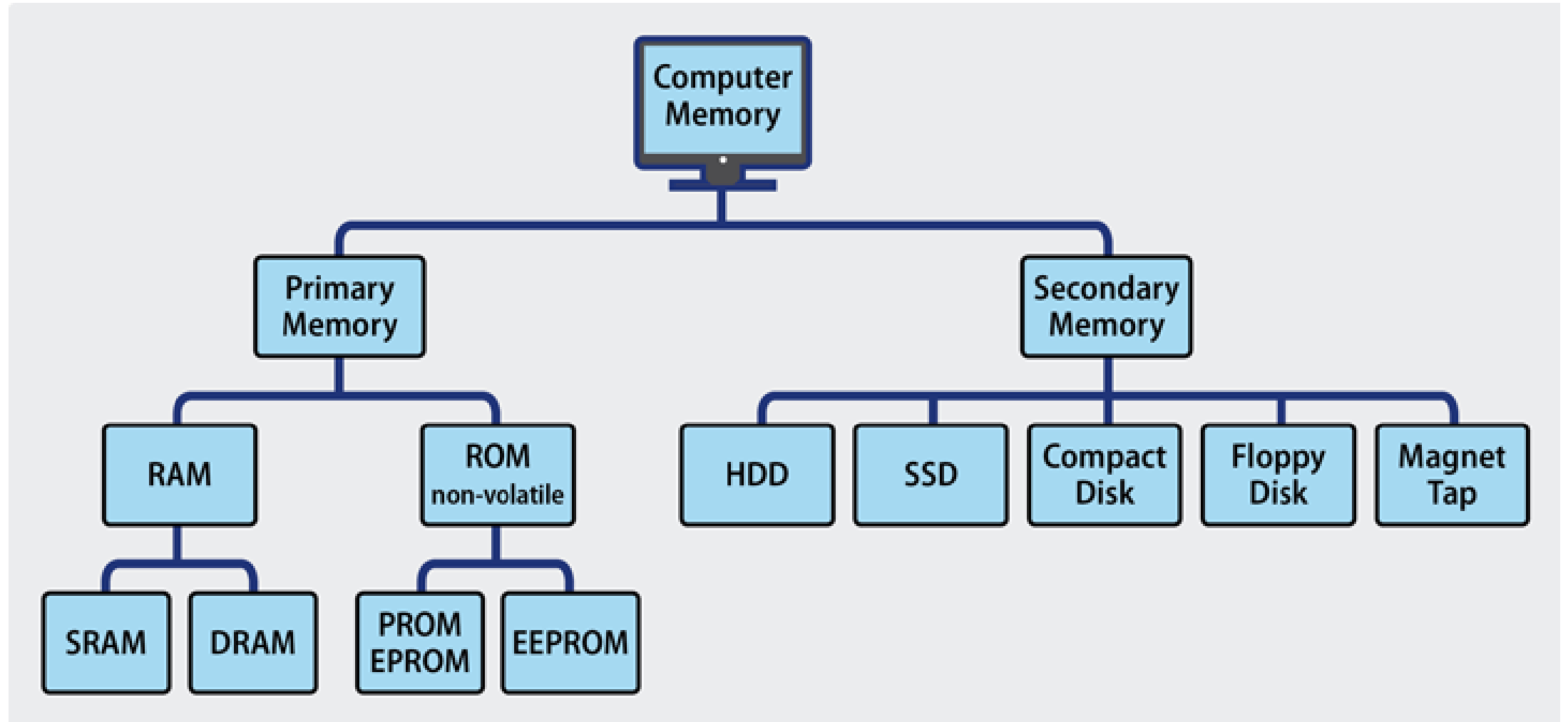| Date | Topic |
|------|-------|
| 4/29 | Lecture - Introduction to NAND Flash & Controller<br>Lecture - Introduction to FW concepts (basic)<br>Lab - FUSE environment set-up (before next lecture) |
| 5/6 | Lecture - Introduction to FW concepts (advanced)<br>Lab - Challenge |
| 6/B | Students provide report and source code |

# AGENDA

NAND Flash Market
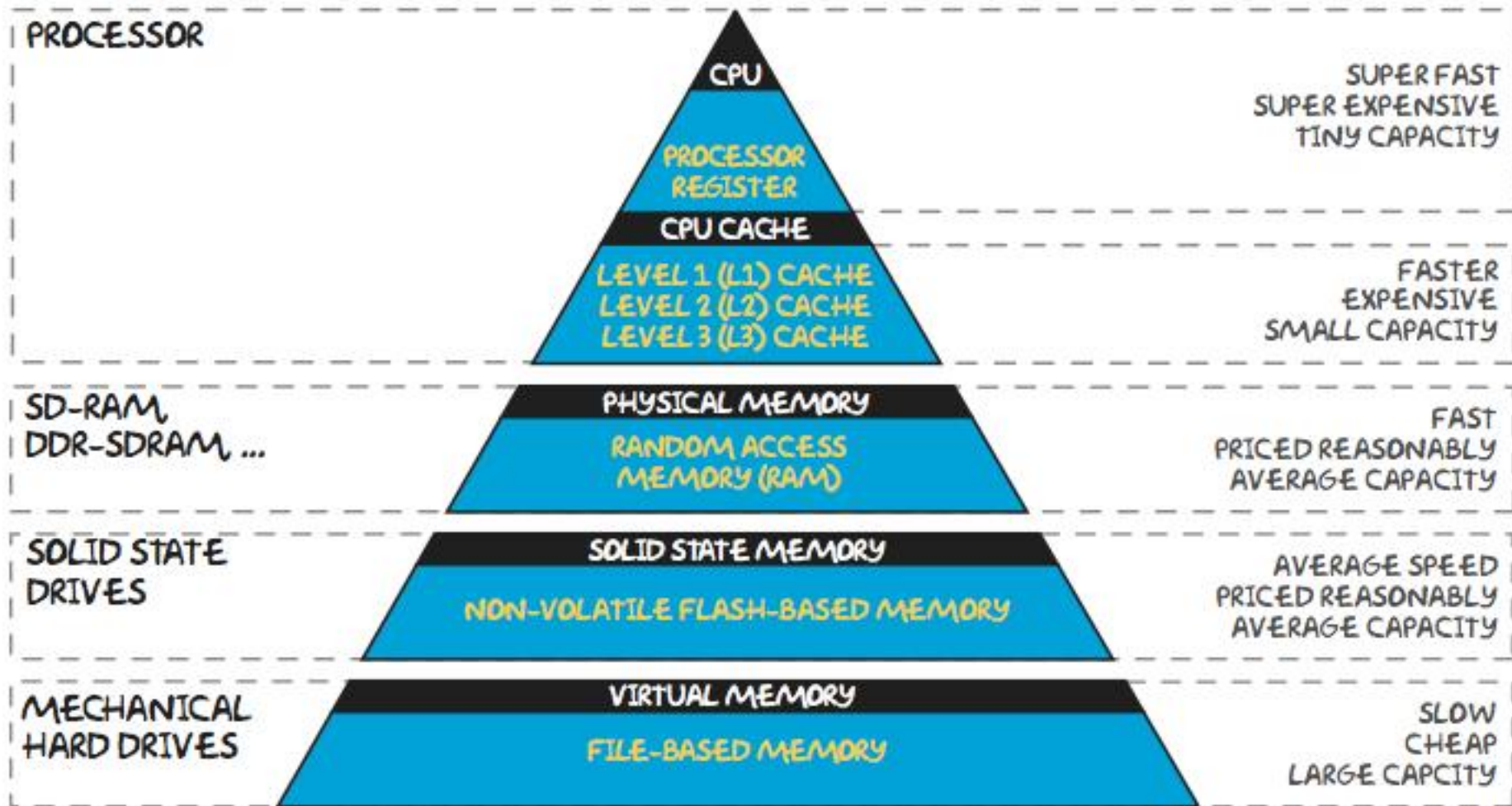
NAND Flash Device

FW Concepts (Basic)

Summary

# NAND Flash Market
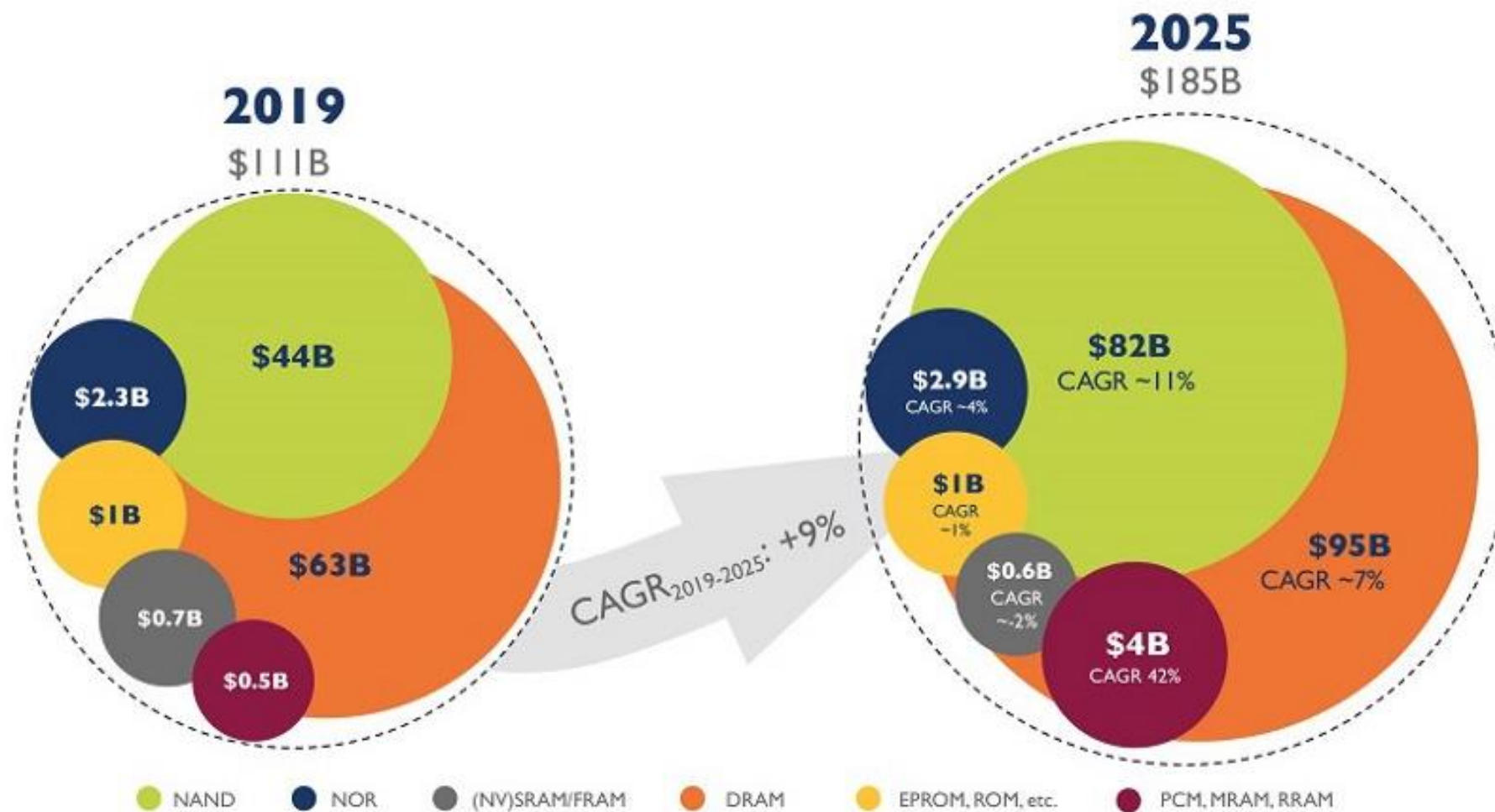
PHISON

# Memory Classification

# Memory Hierarchy

# Question

Which type of memory has the highest market cap growth in 5 years time?

Which type of memory has the highest CAGR in 5 years time?

PHISON

# Market of Memory



2019
$111B

2025
$185B

CAGR 2019-2025: +9%

**2019:**
- $2.3B (NOR)
- $44B (NAND)
- $1B (EPROM, ROM, etc.)
- $0.7B ((NV)SRAM/FRAM)
- $0.5B (PCM, MRAM, RRAM)
- $63B (DRAM)

**2025:**
- $2.9B CAGR ~4% (NOR)
- $82B CAGR ~11% (NAND)
- $1B CAGR ~1% (EPROM, ROM, etc.)
- $0.6B CAGR ~-2% ((NV)SRAM/FRAM)
- $4B CAGR 42% (PCM, MRAM, RRAM)
- $95B CAGR ~7% (DRAM)

Legend:
- NAND
- NOR
- (NV)SRAM/FRAM
- DRAM
- EPROM, ROM, etc.
- PCM, MRAM, RRAM

PHISON

# NAND Flash Is Everywhere



PC Computer

Netflix/FB/Amazon Servers

Smartphone

Game Consoles

Smart Speaker

**PHISON**

Google Glass

Smart Robots

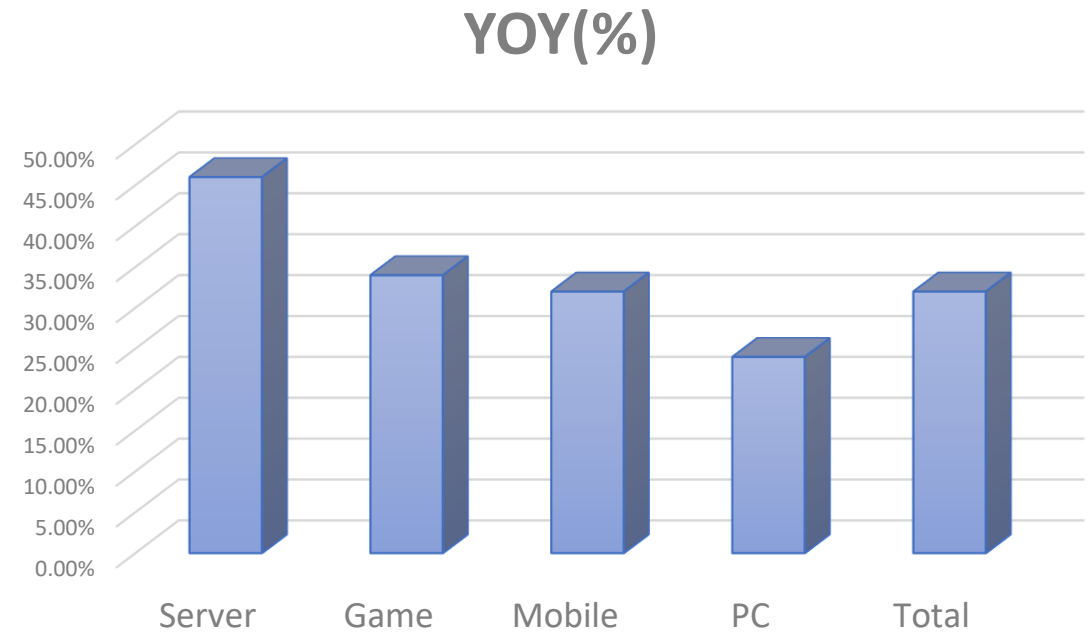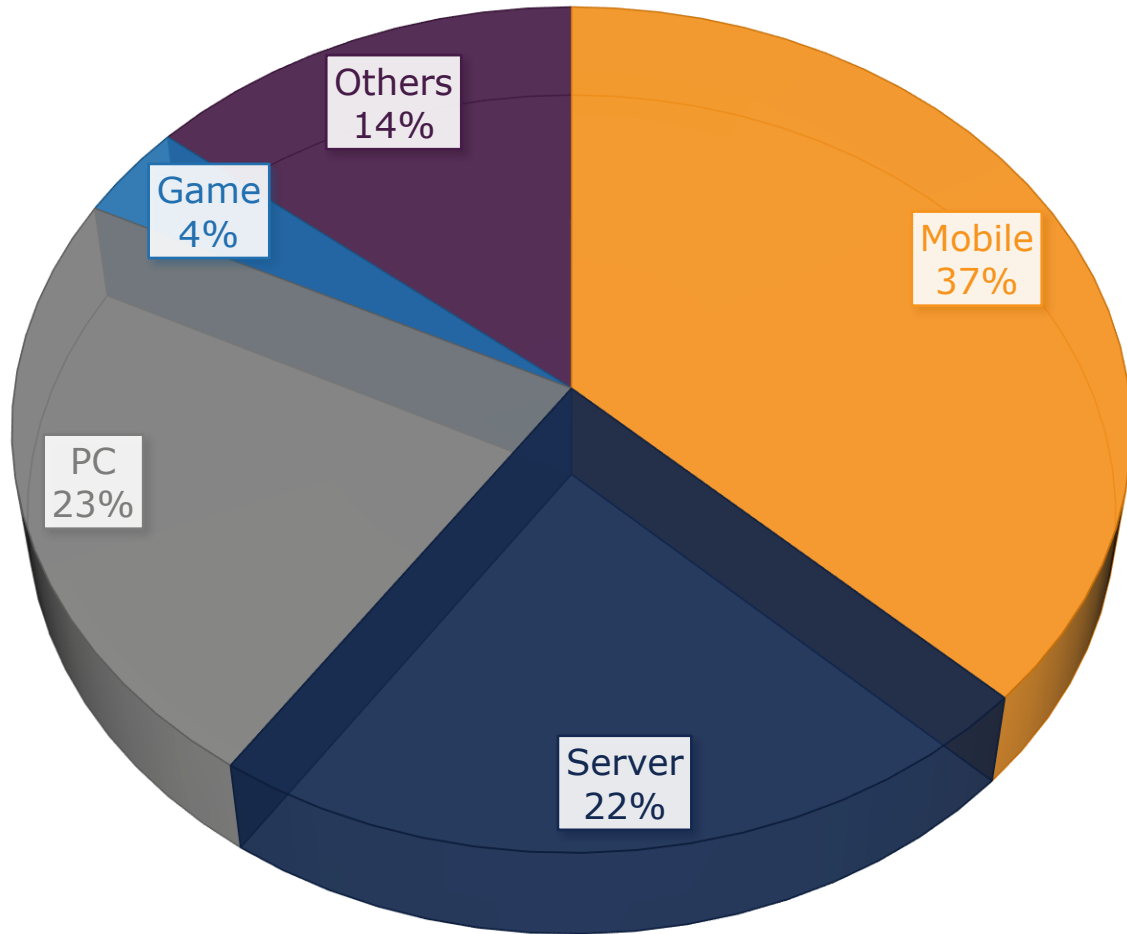Drones

Smart Treadmills

Smart Projectors

Electric Cars

**PHISON**

# Question

Which type of application used the most of NAND Flash?

Which type of application that uses NAND Flash has the highest growth?

**PHISON**

# Distribution and Trend by Applications

# Data Storage



**Floppy Disk**

1.44MB

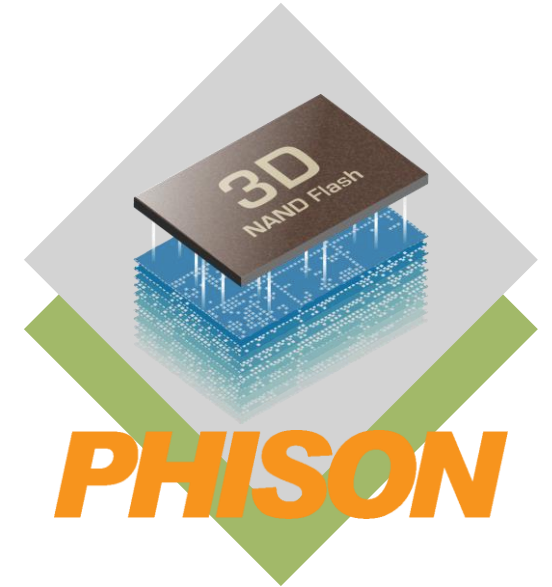**Optical Disc**

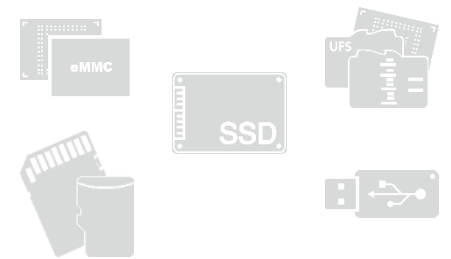CD: 700MB
DVD: 4.7GB
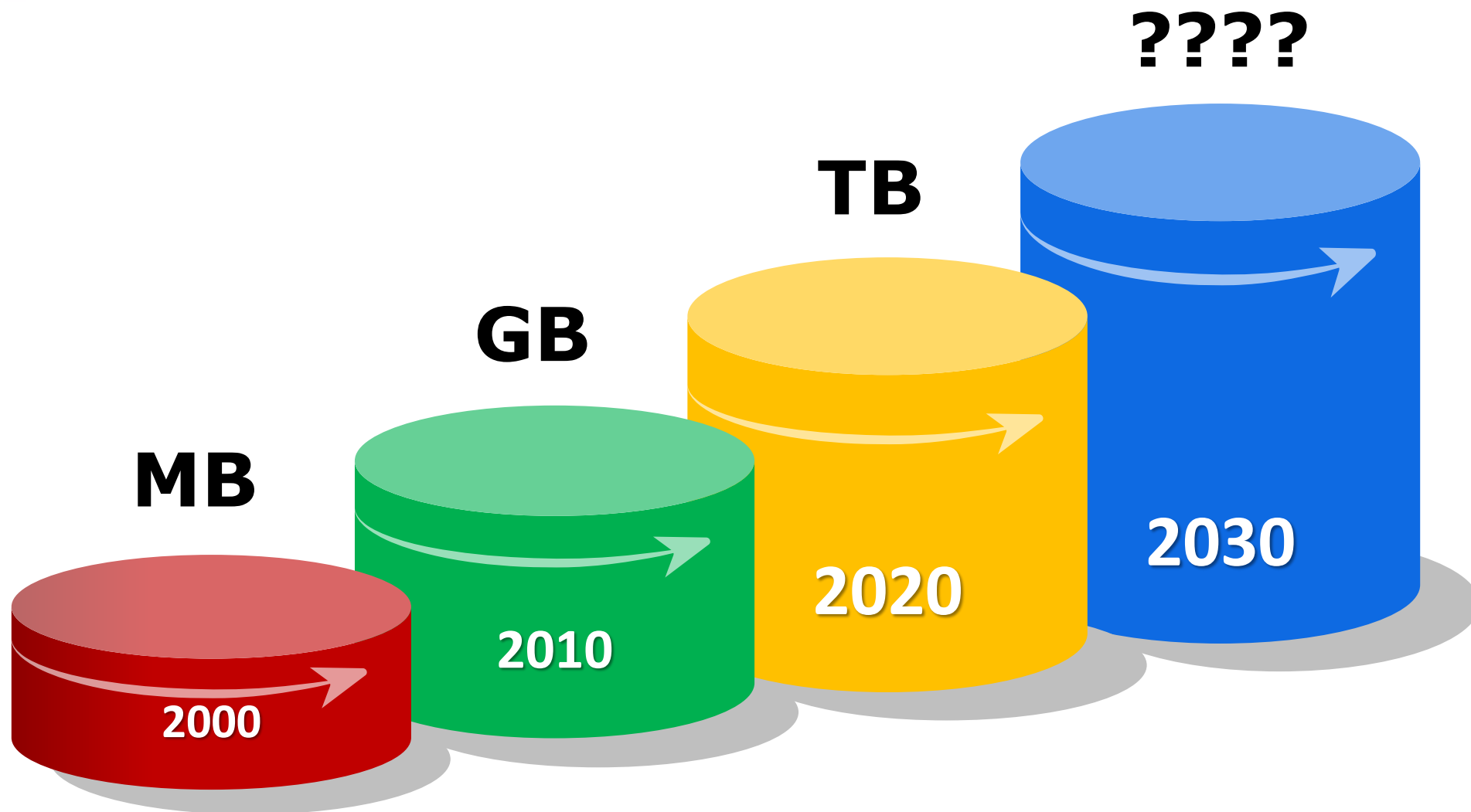Blu-ray Disc: 25GB~128GB

**HDD**

36GB~1TB~Beyond

**NAND Storage**
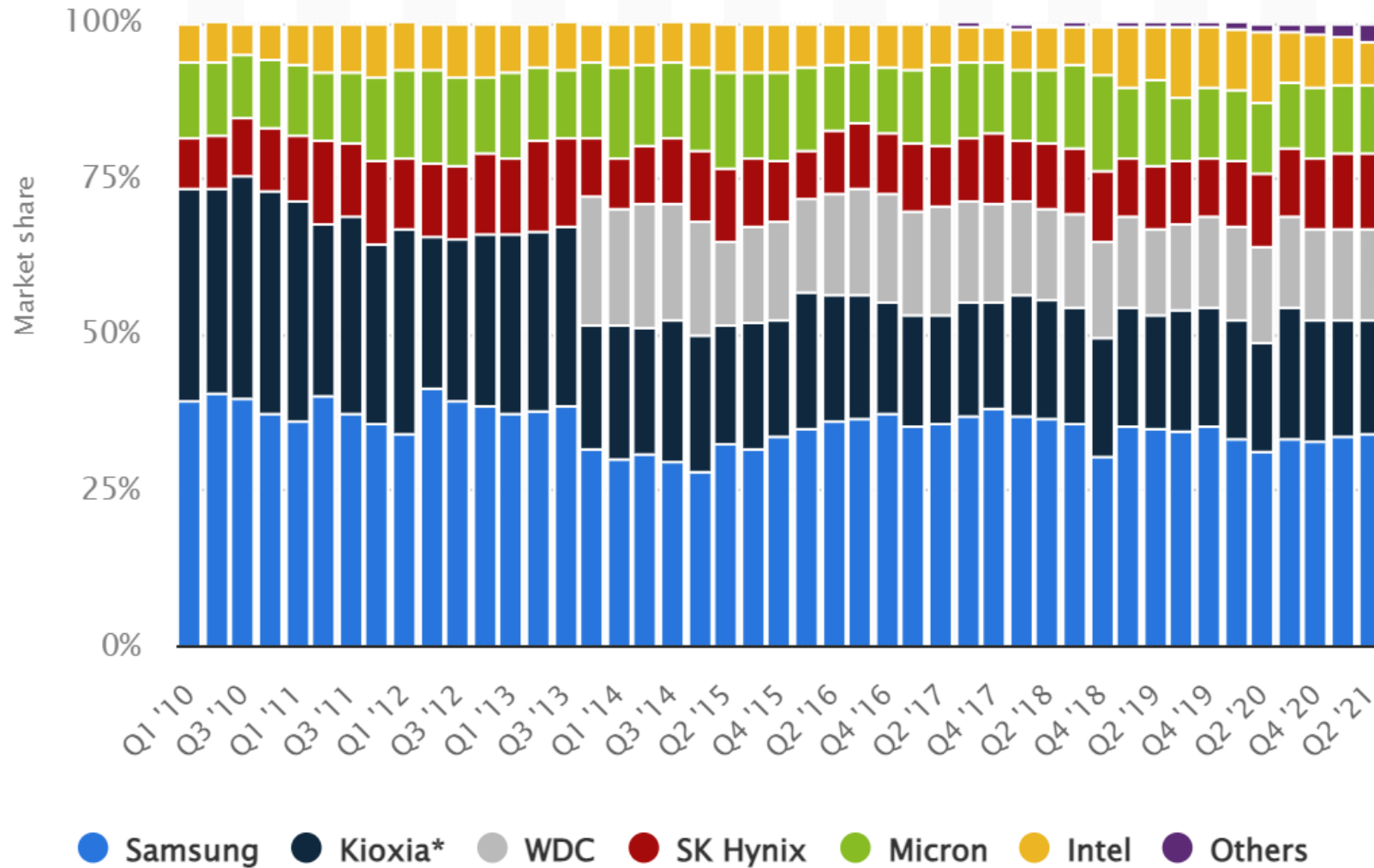
1GB~1TB~8TB~Beyond

# Demand Never Stop

# Question

Which company invented NAND Flash?


A NAND Flash company founder used to be a farmer, guess what did he used to grow?

**PHISON**

# NAND Flash Players

# NAND Flash Players



Legend: Samsung, Kioxia*, WDC, SK Hynix, Micron, Intel, Others

PHISON

# ECO System of NAND Flash Storage System
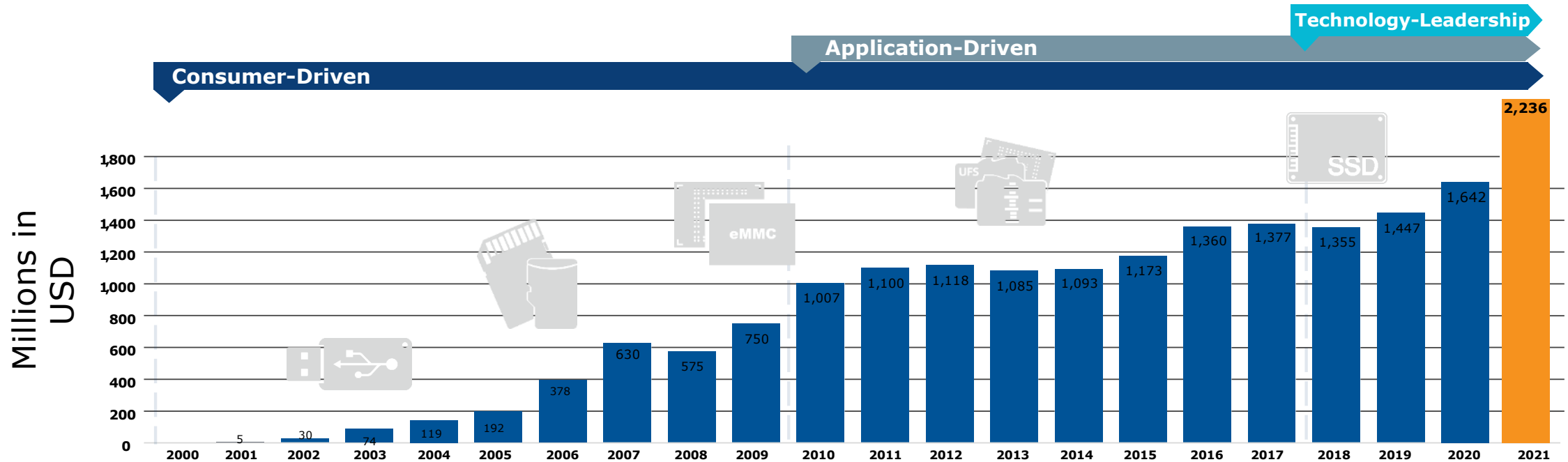
# About Phison

## World No.1
### NAND Flash Controller Solution

## Taiwan Top 4
### IC Design House

## 20%+ Worldwide
### SSD Controller Market Share

Gaming    Enterprise    Automotive    Industrial    Embedded ODM

Technology-Leadership

Application-Driven

Consumer-Driven

Millions in USD

| Year | Value |
|------|-------|
| 2000 | 5 |
| 2001 | 30 |
| 2002 | 74 |
| 2003 | 119 |
| 2004 | 192 |
| 2005 | 378 |
| 2006 | 630 |
| 2007 | 575 |
| 2008 | 750 |
| 2009 | 1,007 |
| 2010 | 1,100 |
| 2011 | 1,118 |
| 2012 | 1,085 |
| 2013 | 1,093 |
| 2014 | 1,173 |
| 2015 | 1,360 |
| 2016 | 1,377 |
| 2017 | 1,355 |
| 2018 | 1,447 |
| 2019 | 1,642 |
| 2020 | 2,236 |

**PHISON**

# Key Facts

# 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

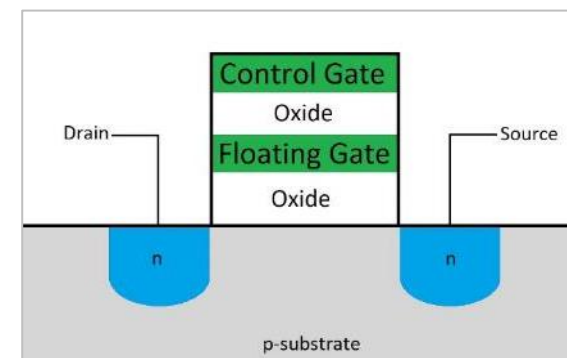Flash Translation layer (FTL)

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 Device

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            Plane 1
(0, 2, 4, ..., 1006)  (1, 3, 5, ..., 1007)



Control Gate
Oxide
Floating Gate
Oxide

Drain — Source

n         n

p-substrate

*Memory Cell*

# What's MOSFET

Gate

ゲート

Add voltage to gate

Create a channel in between

ソース ドレイン

Drain

N → チャネル N

Source

オン状態 オフ状態

e→ e→

伝導帯

価電子帯

ソース チャネル ドレイン ソース チャネル ドレイン

熱拡散による緩慢なオン・オフ特性

NMOS

PMOS

D. Kahng & Martin Atalla, 1960, Bell Lab

**PHISON**

# What's Threshold Voltage?



Id-Vg Characteristics

3D electron density for Vd=0.6

electron density for = 1-Vd=0.6V-Vg=0V

- As known as $V_{th}$, $V_t$ :a gate voltage which can "turn on" the channel of a transistor.

**PHISON**

# What's Floating Gate MOSFET



Control Gate

Floating Gate

Floating Gate (FG)-NVM cell, 1967 (Bell)

✓ **Electrically isolated**

✓ **Charge contained in it remains unchanged for long periods of time**


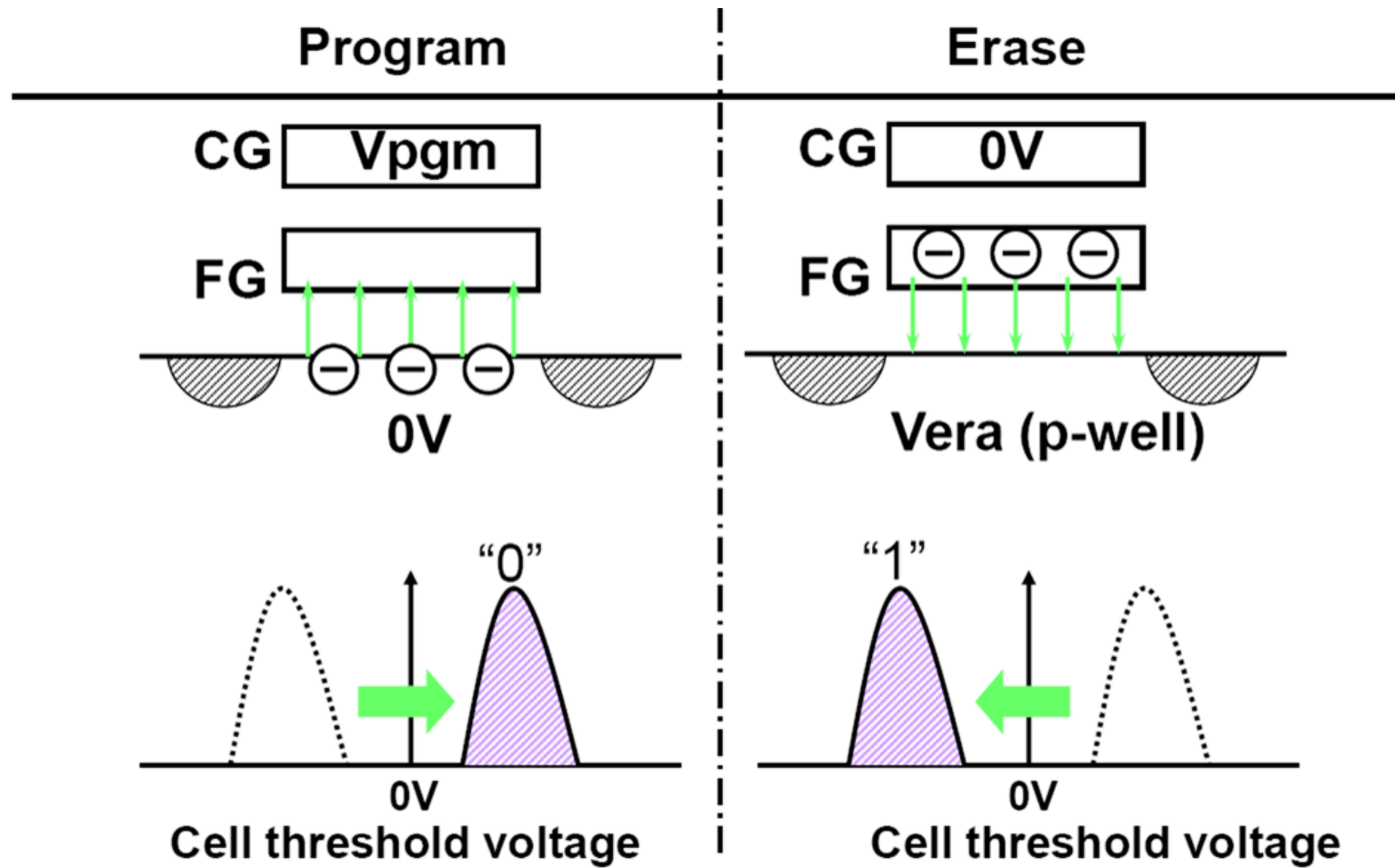
半導體之父 施敏教授研發的「浮閘記憶體」

# Threshold Voltage of Memory Cell

MOSFET

Logic Data "1"

Logic Data "0"

# Basic Operation of NAND Flash- P/E



Program operation

| | | |
|---|---|---|
| BL m-1 | BL m | BL m+1 |

Program

| | |
|---|---|
| | floating |
| SGD | 0V |
| SGD | 3V |
| WL31 | 10V (Vpass) |
| WL30 | 20V (Vpgm) |
| WL29 | 10V (Vpass) |
| WL1 | 10V (Vpass) |
| WL0 | 10 (Vpass) |
| SGS | 0V |
| Source | 3V |
| SGS | 0V |
| | floating |

Block n

3V    0V    3V

Non programm

Erase operation (Block erasing )

| | | |
|---|---|---|
| BL m-1 | BL m | BL m+1 |

| | |
|---|---|
| | floating |
| SGD | floating |
| SGD | floating |
| WL0 | 0V |
| WL1 | 0V |
| WL2 | 0V |
| WL30 | 0V |
| WL31 | 0V |
| SGS | floating |
| Source | floating |
| SGS | floating |
| | floating |

Block n

floating

*P-Well= ~20V Vera*

# Threshold Voltage Distribution

$V_{th}$ of a 128 bit NAND Array:



$V_{th}$ distribution of the array:

# Question

Why do we need XLC?

**PHISON**

# 2D (Planar) NAND Flash

# 3D NAND Array



**3D NAND Architecture**

Bit Line

Bit Line

BL Contact

BL Contact

Layer 96 SGD

String1   String2   String3   String4

WL

Layer 3 SGS
Layer 2
Layer 1

Memory Holes

Memory Cell

String

String

◆ Top View

Page Size = 16K Byte (16*1024*8)
1 WL = 4 strings = 4 Pages
1 Block = 4 * 96 WLs = 384 Pages (SLC)

**PHISON**

# Reliability Issues of NAND Flash



Program; FN Current from Inversion Layer

Si

Hole trappings in the tunnel oxide locally reduce barrier height

Anomalous increase of FN current

By hole detrapping, FN current shows normal value

FG

Weak Programming mode due to Vread stress

Vread=3.5 ~ 5V

0V

Vread = high      0V      Vread =high

0V

Bit Line

BL Contact

SGD

WL

SGS

Memory Holes      Source Plate

Memory Cell

Bit Line

Word Line

3D Charge Trap

WP n-1

WP n

WP n+1

Poly channel

= Loss through oxides
= Lateral migration

Source: G. Van Den Bosch, IMW 2014

# Question

Why do we need NAND Flash Controller?

What does NAND Flash Controller need to do?

**PHISON**

# 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)

Mapping Table &
Garbage Collection

**Page**

New data

0

1

.
.
.

N-1

N

**Page**

0

1

.
.
.

N-1

N

Read

Write

**Cache**

Modify

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

I/O scheduling

**PHISON**

# NAND Flash Limitations

- What if sudden power loss happens before VT reaches programmed state?



Read Level

Power Loss Rebuild

**(a) Erased State**

E $V_{TH}$

$V_{TARG}$ Power Loss

Electrons tunnel into FG

**(b) First programming pulse**

$V_{TH}$

$V_{TARG}$

**(c) N programming pulses**

$V_{TH}$

**ISPP Procedure**

START

Apply Programming Pulse

Verify most cells have $V_{TH}$ higher than $V_{TARG}$

FAIL

PASS

END

$t_{PROG}$~1500us

# NAND Flash Limitations

- Program/Erase (PE) Cycle: SLC = 100k, MLC = 10k, TLC = 3k  **Wear Leveling**
- Initial and runtime bad blocks  **Bad Block Management**



P/E Cycles

# NAND Flash Limitations

- Main sources of noise that would cause abnormal Vth distribution and read errors:
  - **Program/Erase cycling stress**: Program/erase pulses lead to degraded reliability of the underlying NAND flash cells.
  - **Cell-to-cell interference**: Threshold voltage of 'victim' cell is strongly affected by programming of neighboring 'aggressor' cells.
  - **Data retention**: Over time, electrons can escape from the programmed flash cells, causing a loss of threshold voltage.
  - **Read disturbance**: When reading a particular page in a block of NAND flash, a voltage is applied to all other WL in order to 'deselect' them. This applied voltage can affect the Vth distribution of the unselected WLs.

Error correction code &
Read disturb management

# NAND Flash Controller

# Key Departments in Controller House

**Electronics & Electrical, Computer Science, Physics, etc.**

## Firmware Engineer
- ✓ Skills :
  - ● C/C++
  - ● Python
  - ● Computer Architecture
- ✓ Firmware development
- ✓ FTL Algorithm design
- ✓ RMA Analysis

## Software Engineer
- ✓ Skills :
  - ● C/C++
  - ● Visual C++
  - ● Python
- ✓ V&V pattern design

## Field Application Engineer
- ✓ Skills :
  - • Electronics & Electrical fundamentals
- ✓ Circuitry design
- ✓ Customer support

## Analog IC Engineer
- ✓ Skills :
  - ● Analog IC Design
- ✓ Power analog circuit design
- ✓ Analog/Mixed signal design
- ✓ High Speed IO/transceiver circuit designs

## Digital IC Engineer
- ✓ Skills :
  - ● FPGA、Verilog
  - ● Digital Systems
  - ● Design & Simulation Tools
- ✓ IP design & validation
- ✓ CP/FT

## NAND Flash Engineer
- ✓ Skills :
  - ● C/C++
  - ● Semiconductor Device Physics
- ✓ Error correction algorithm
- ✓ NAND Flash characteristic analysis

**PHISON**

# 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

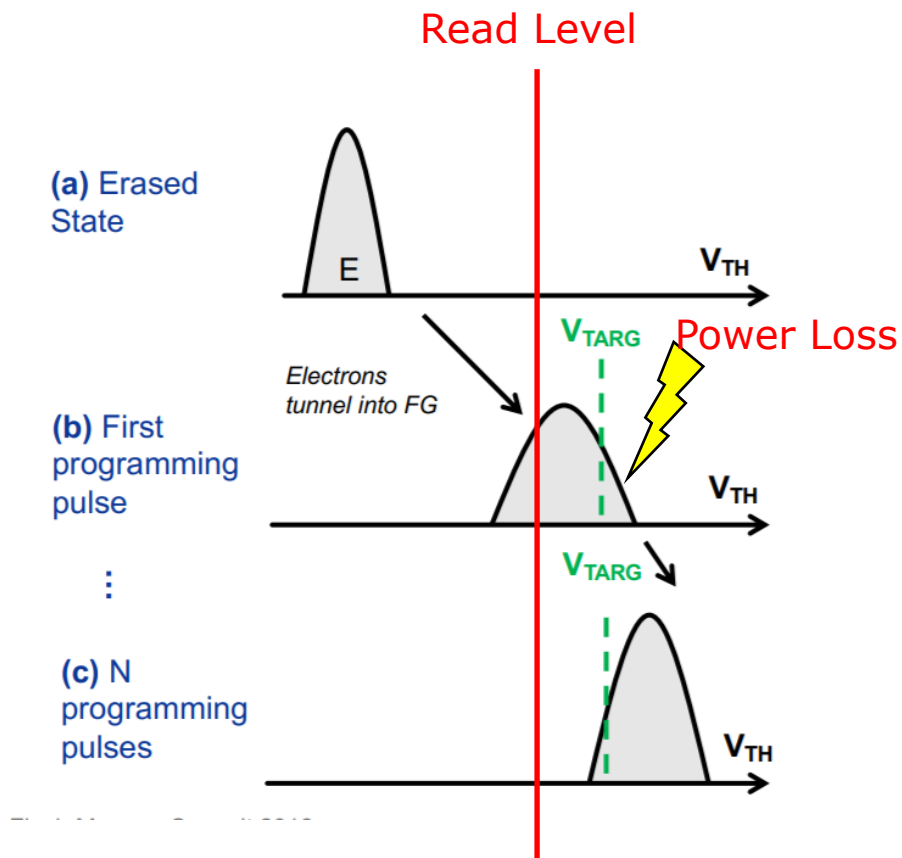- Host writes 2 chunk of data (16KB each)
  - LBA 0~31
  - LBA 32~63

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

**PHISON**

# 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 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

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

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

# 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 | |
|---|---|
| 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

| Logical Address | Mapping Table | Physical Address |
|---|---|---|

Unit: Page

# 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**

- 91% — 4KB
- 6% — 8KB
- 1% — 16KB
- 2% — 32KB
- 0% — 64KB
- 0% — 128KB

**Read Chunk Size**

- 46% — 4KB
- 6% — 8KB
- 10% — 16KB
- 31% — 32KB
- 7% — 64KB
- 0% — 128KB

# 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

- 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

# Write

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

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

**Block 0**

No need to copy data to new block during overwrite, old data becomes invalid

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

PHISON

# Read

- Host wants to read LBA 32~63
- How do we know where is the data located in NAND Flash?
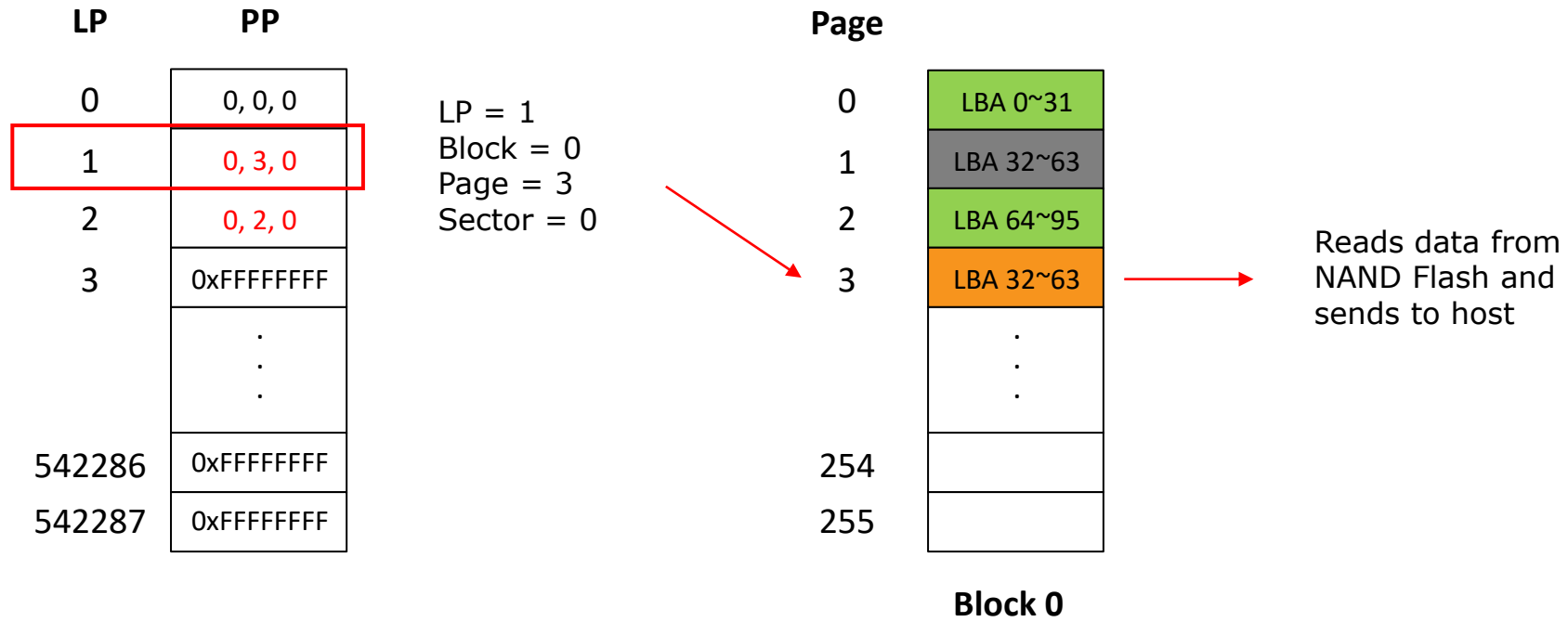
- Logical to Physical (L2P) translation

    Step 1: Calculate Logical Page (LP) Index

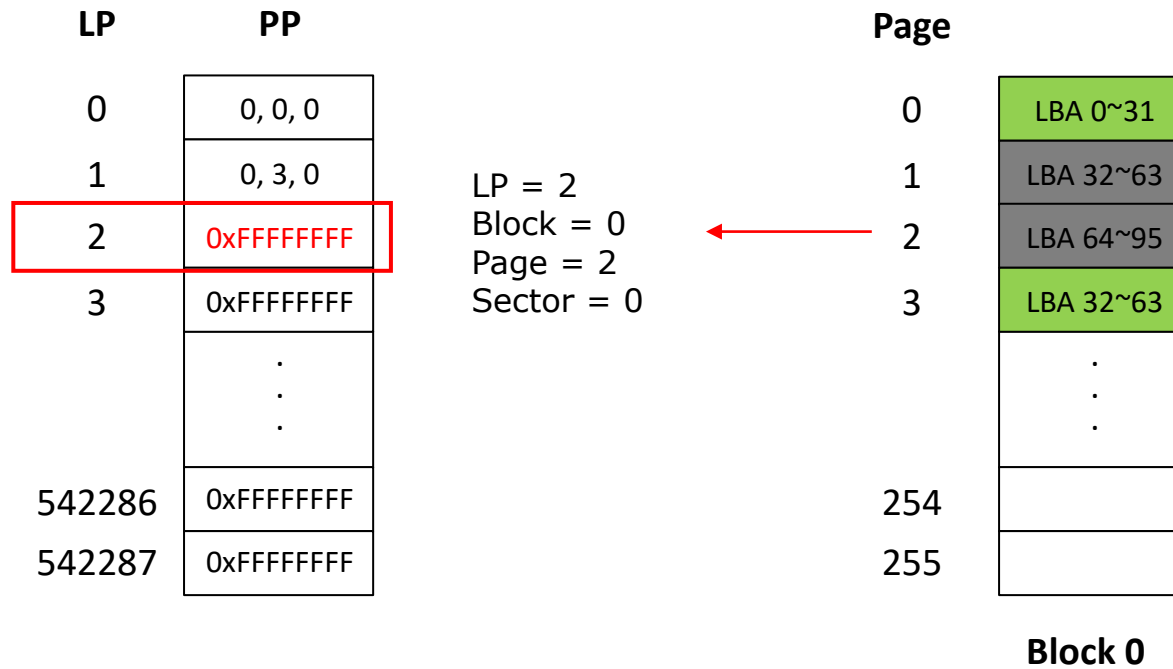    $$LP = \frac{LBA}{Total\ Physical\ Sector\ Number\ (in\ page)}$$

    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
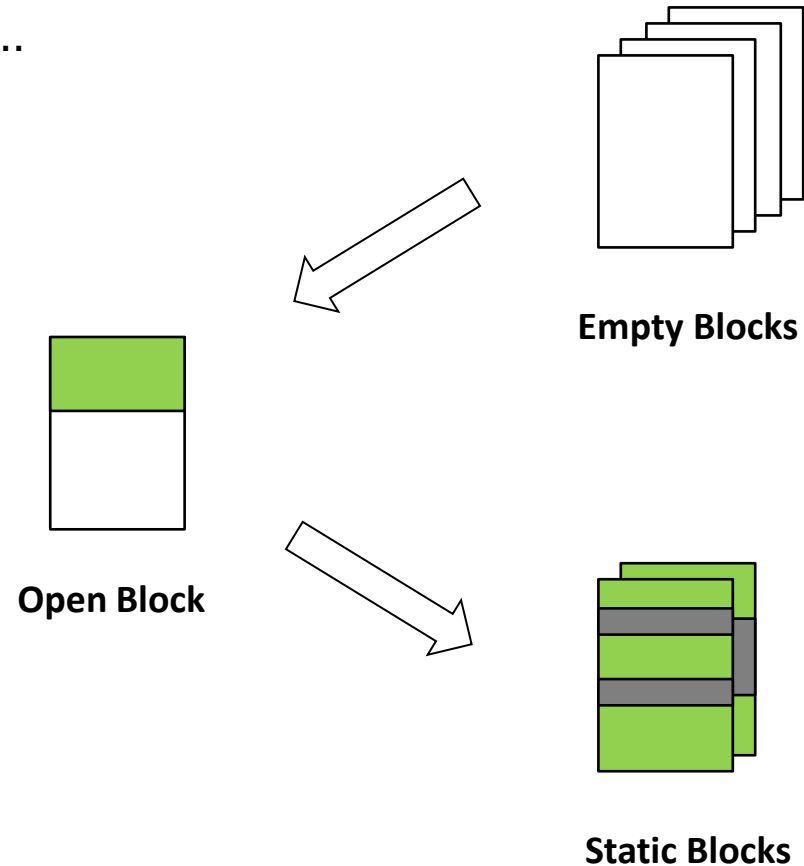  - Table Entries = 2048 blocks × 256 pages × 4 nodes= 2097152

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

PHISON

# Garbage Collection (GC)

# Garbage Collection

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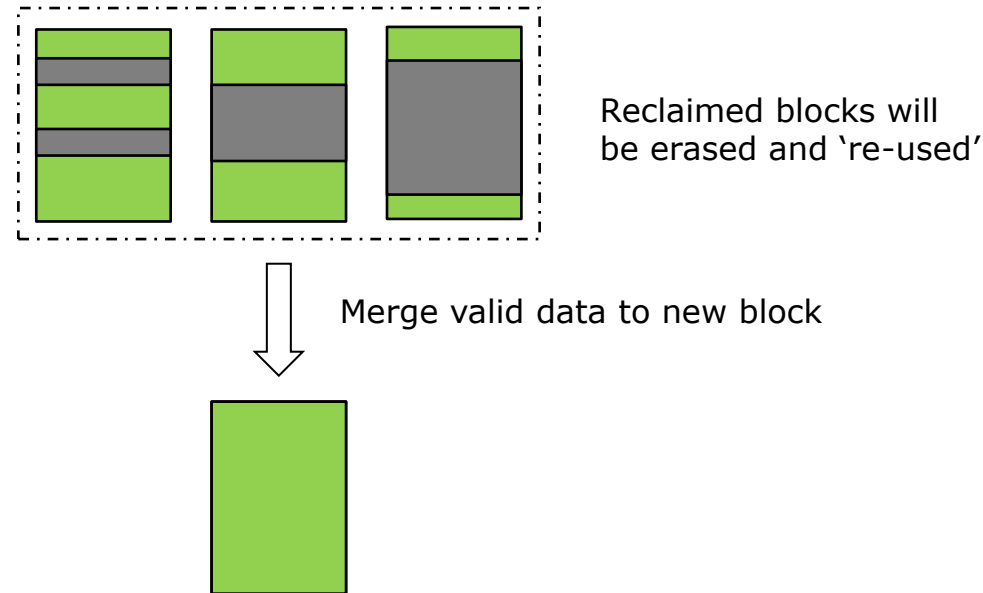- Open block becomes static when full
- Some physical page become invalid

**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

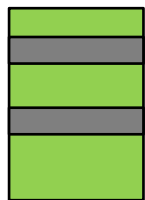- Pick the block(s) with least valid count
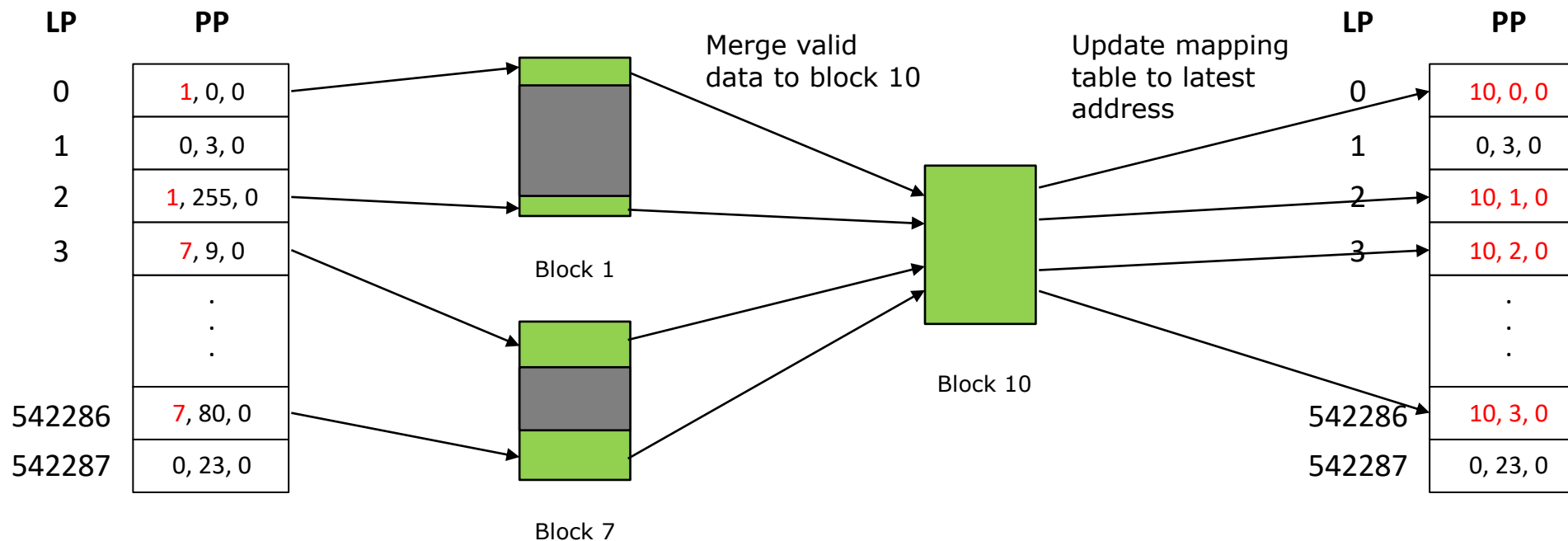
Valid cnt = 10    Valid cnt = 30    Valid cnt = 50

# Garbage Collection

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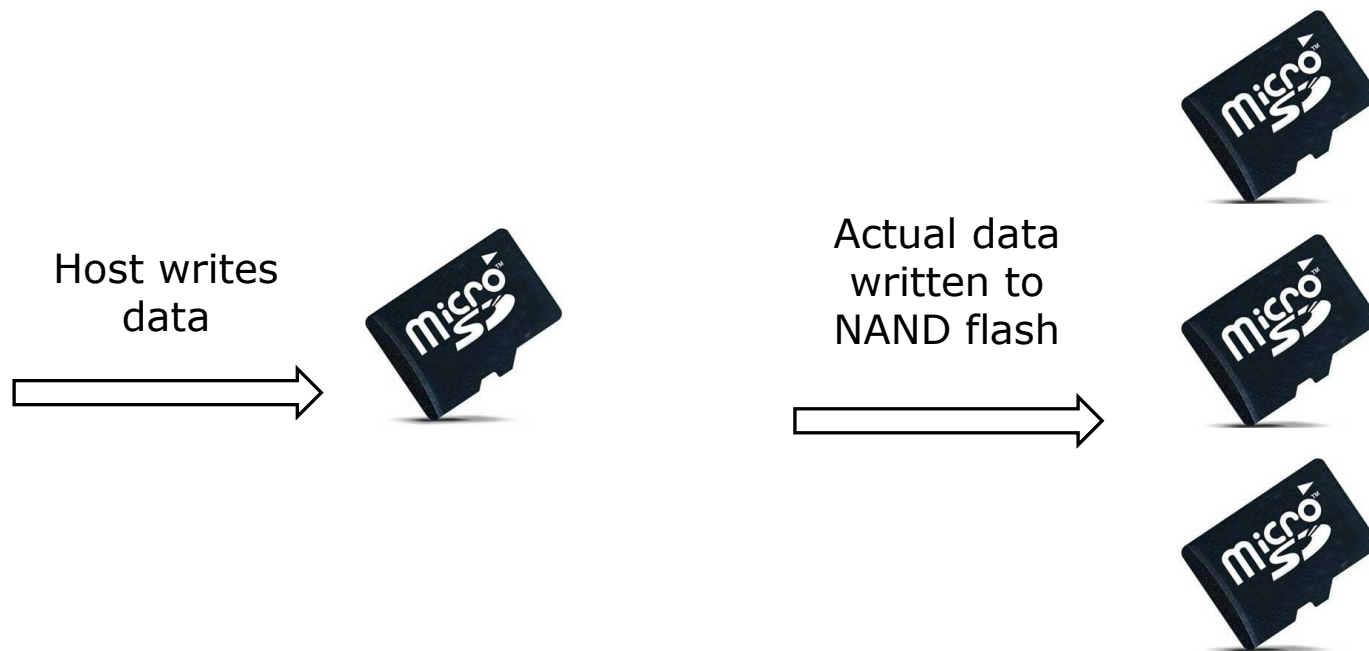- Check the mapping table, look for the entries that point to these source blocks
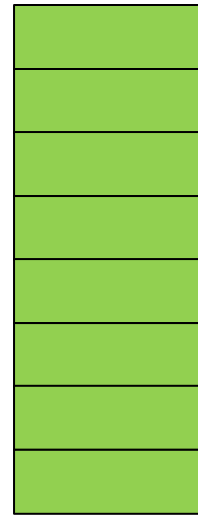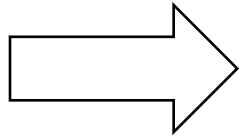
# Write Amplification Factor (WAF)

# Definition

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

Host writes data

Actual data written to NAND flash
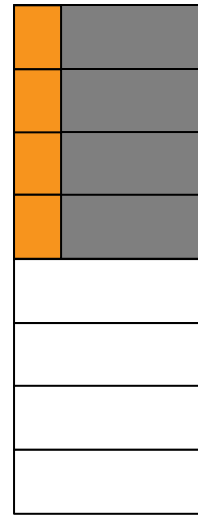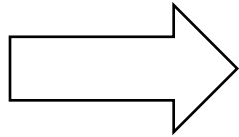
# Example

- 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?

# Summary

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

# 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

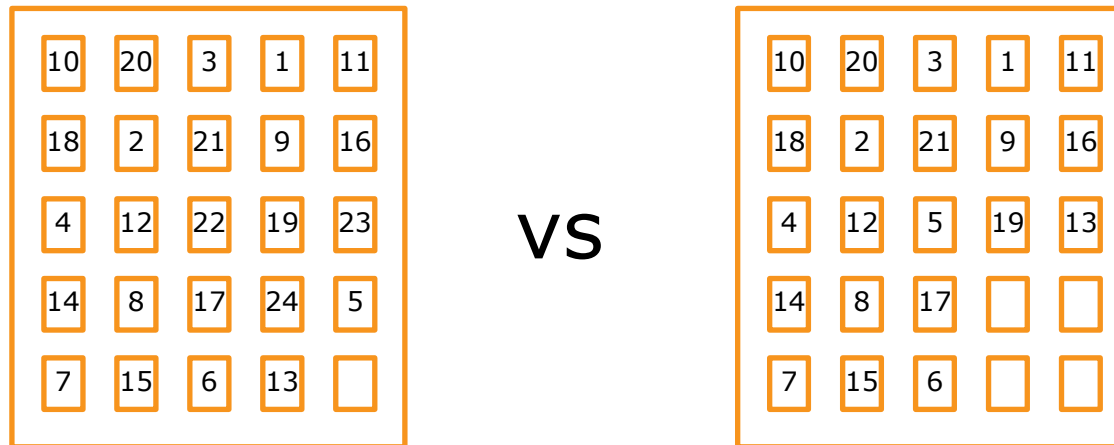# Definition

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

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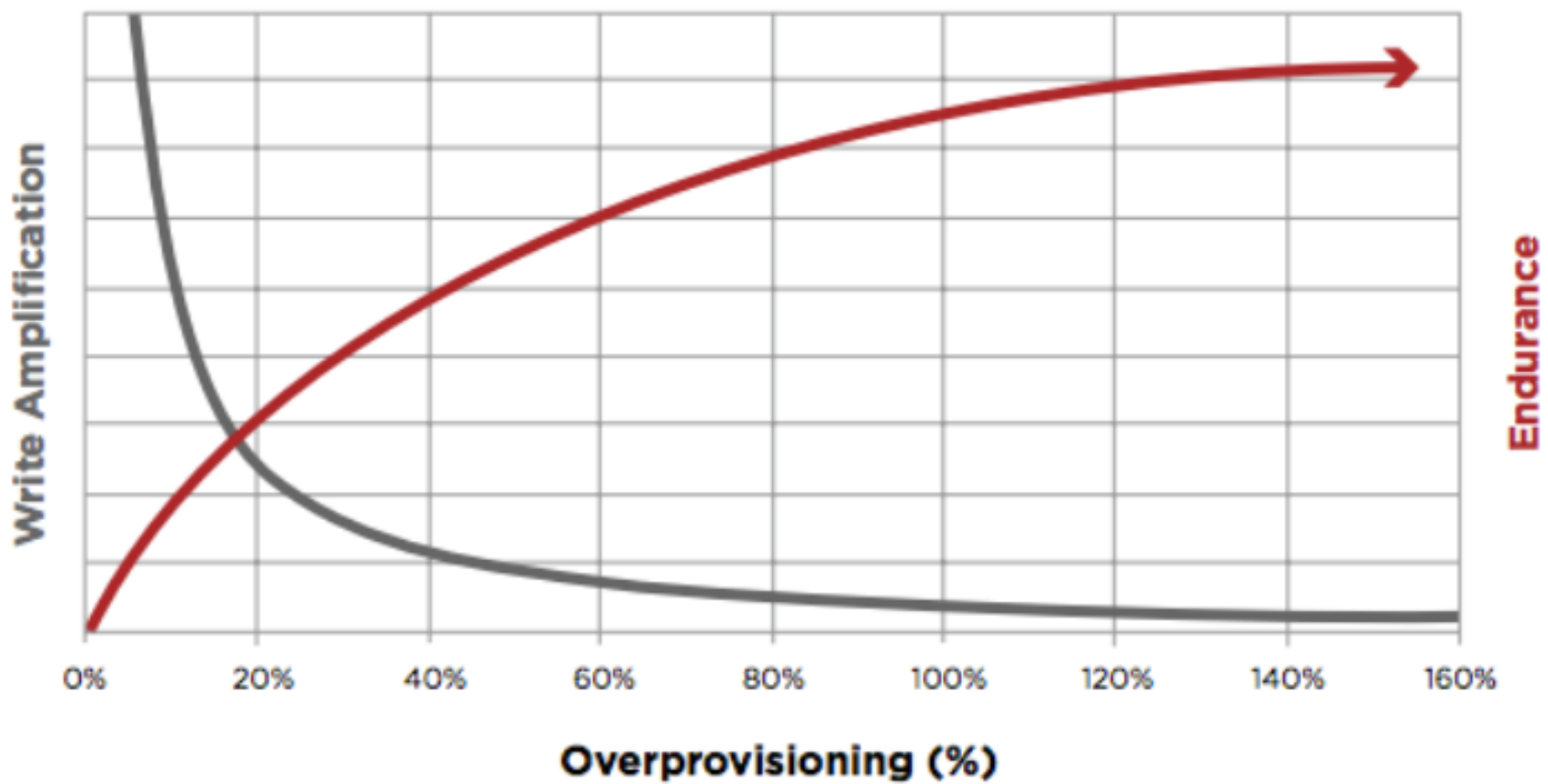
- The larger the size of the spare area, the higher the operating efficiency, and the better the performance become



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

PHISON

# 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)

**PHISON**

# Page Base vs Block Base
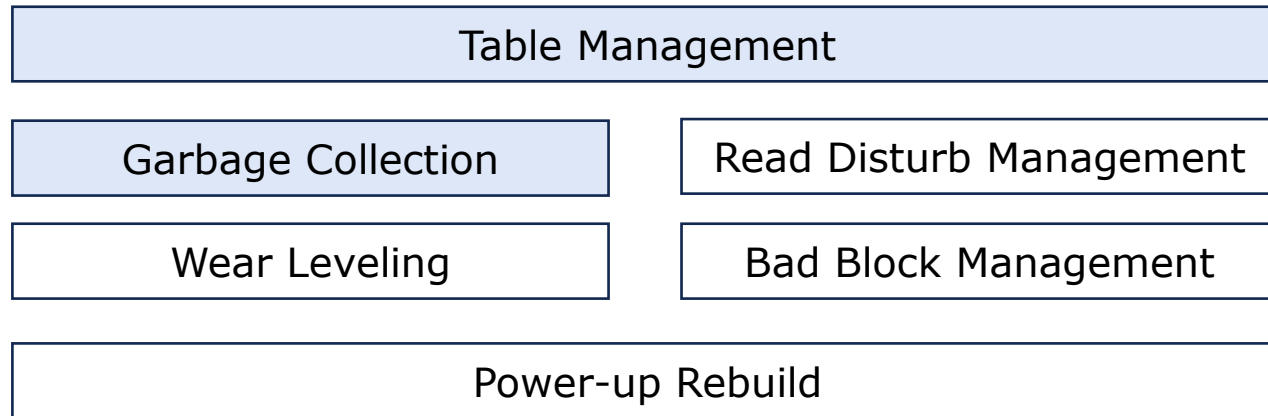
| | Page Base | Block Base |
|---|---|---|
| **Mapping Unit** | Page (or 4KB) | Block |
| **Table Size** | Large (store in NAND Flash) | Small (store in SRAM) |
| **R/W Performance** | Excellent random write performance | Slow random write performance, but impressive read performance |
| **Garbage Collection** | 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

**PHISON**

# THANK YOU!