

Overview of Flash and SSDs

Jihong Kim
Dept. of CSE, SNU

Based on:

M. Cornwell. Anatomy of a solid-state drive, ACM Queue 10(10), 2012.

Solid-State Drives (SSDs)

- A Solid-State Drive (SSD) is a data storage device that emulates a hard disk drive (HDD)
- Block devices
 - Small fixed contiguous segments of bytes as the addressable unit
 - Logical addressing to access data and abstract the physical media
- NAND Flash SSD's are essentially arrays of flash memory devices which include a controller that electrically and mechanically emulate, and are software compatible with magnetic HDD's



Traditional hard disk drive

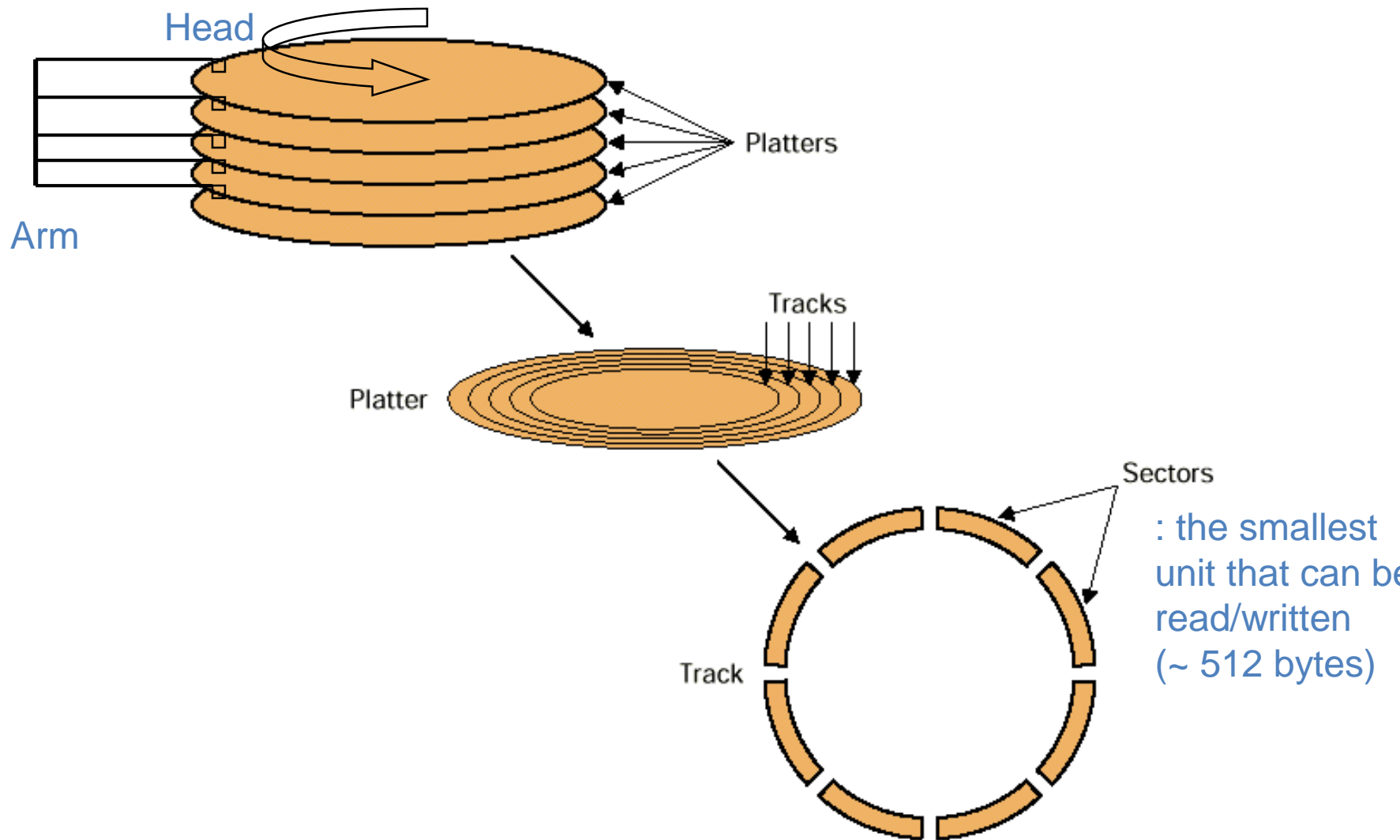


Solid state hard drive

SSD over HDD

- **No mechanical latency**
- **Erase-before-overwrite**
 - Out-place update
- **Asymmetric speed**
 - Program: $\sim 300 \mu\text{s}$ (100s of μs)
 - Read: $\sim 25 \mu\text{s}$ (10s of μs)
 - Erase: $\sim 2\text{ms}$ (a few ms)
- **Asymmetric unit**
 - Program & Read: page
 - Erase: block
- **Limited lifetime**

Magnetic Disks



Key Components of SSDs

- **Three Components**
 - **Storage media**
 - **Controller**
 - **Host interface**

Storage Media: NAND Flash Memory

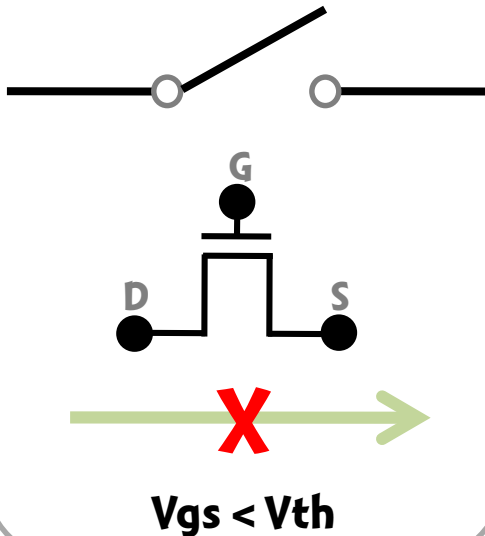
NAND Flash Memory

- **Flash memory** is a **non-volatile computer storage** chip that can be electrically erased and reprogrammed

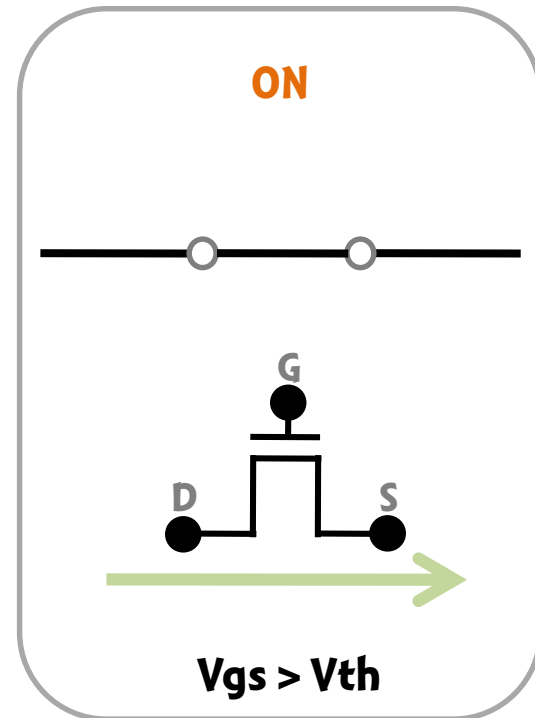
Transistor as an Electrical Switch

- Control gate voltage to turn on/off the switch.
 - Turn on → **Current flow**
 - Turn off → **No current flow**

OFF



ON



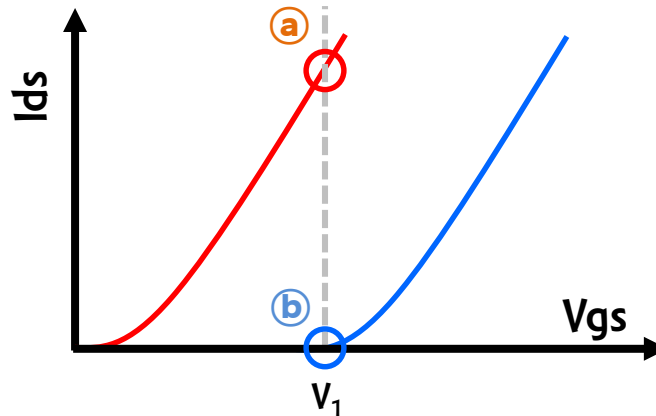
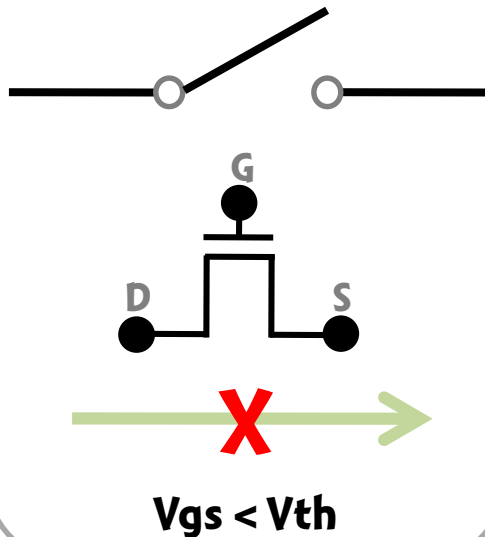
Q: Switch는 메모리는 아니나... 만약... 만약....

- 두 개의 상태를 구분은 하나 switch만으로는 새로운 데이터를 쓸 수가 없다.
- 해결책은?
 - 고정된 V_{th} 값에 대해 V_{gs} 를 조정하니.... On/Off 구분이 가능...
 - 만약 V_{th} 값을 변화시킬 수 있으면... 같은 V_{gs} 에 대해서... On/Off 구분이 가능...
 - 만약 다른 상태가 다른 V_{th} 값을 가지도록 할 수 있다면

Transistor as an Electrical Switch

- Control gate voltage to turn on/off the switch.
 - Turn on → **Current flow** → State 1
 - Turn off → **No current flow** → State 0

State "0"

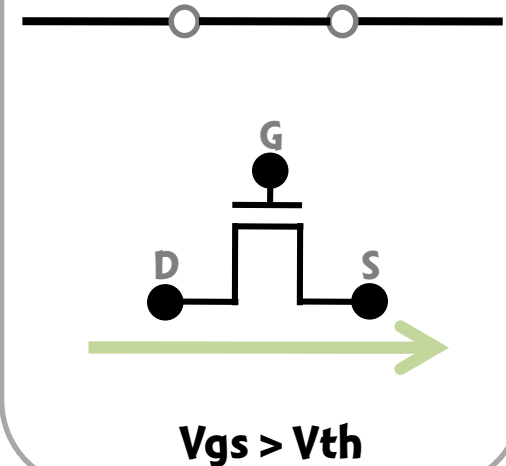


※ When V_1 is applied to gate,

Ⓐ: V_{th} is low → High current → State "1"

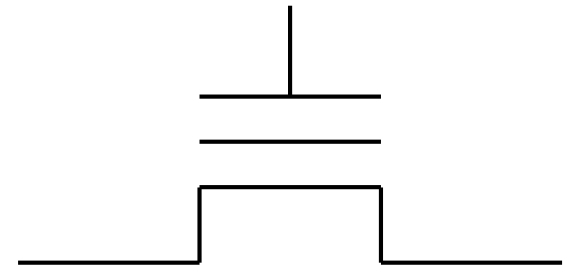
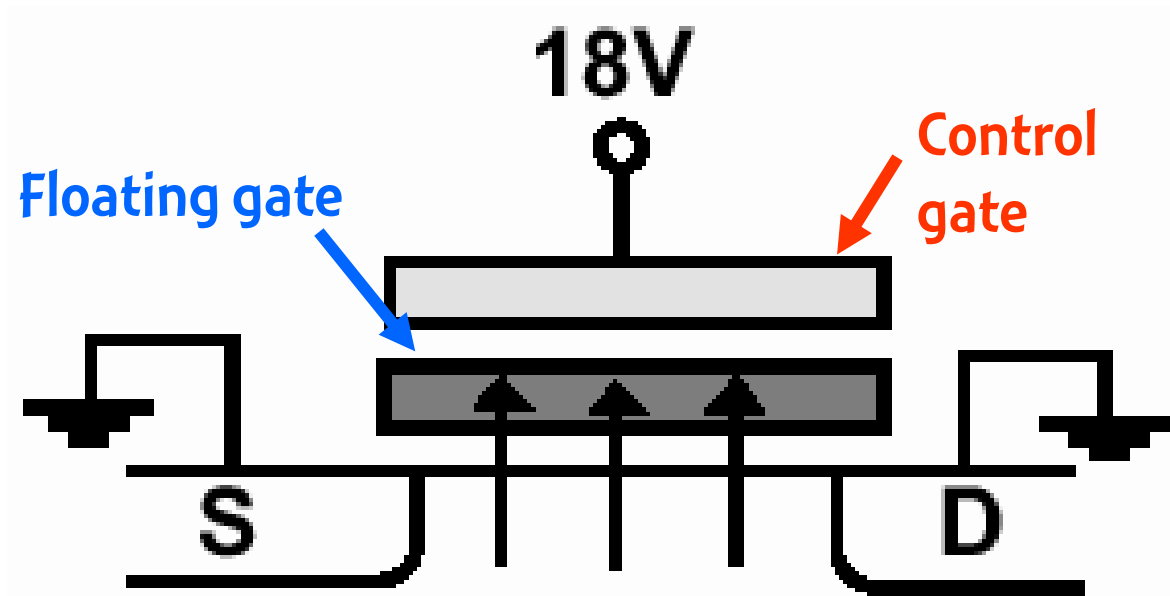
Ⓑ: V_{th} is high → No current → State "0"

State "1"



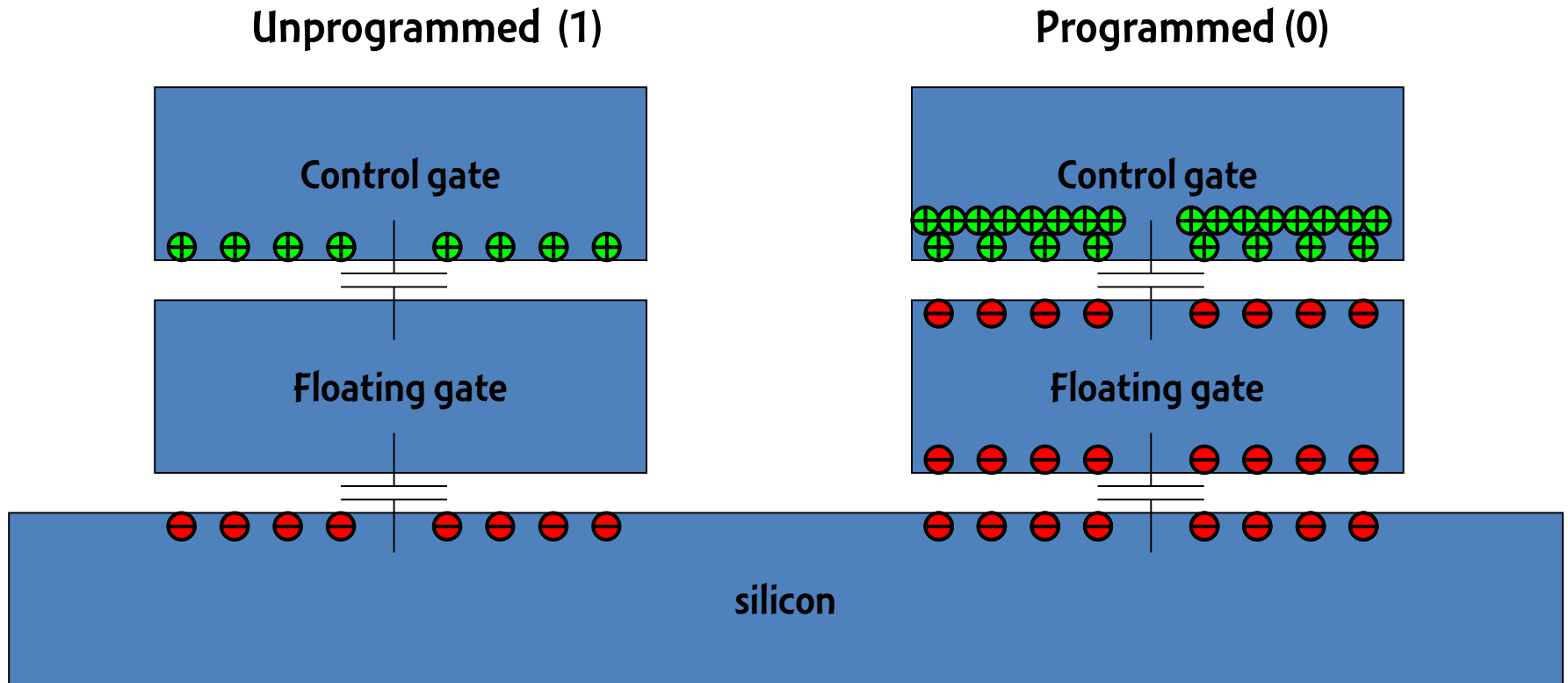
Floating Gate Transistor

- V_T is shifted by injecting electrons into the floating gate;
- It is shifted back by removing these electrons again.



- CMOS compatible technology!

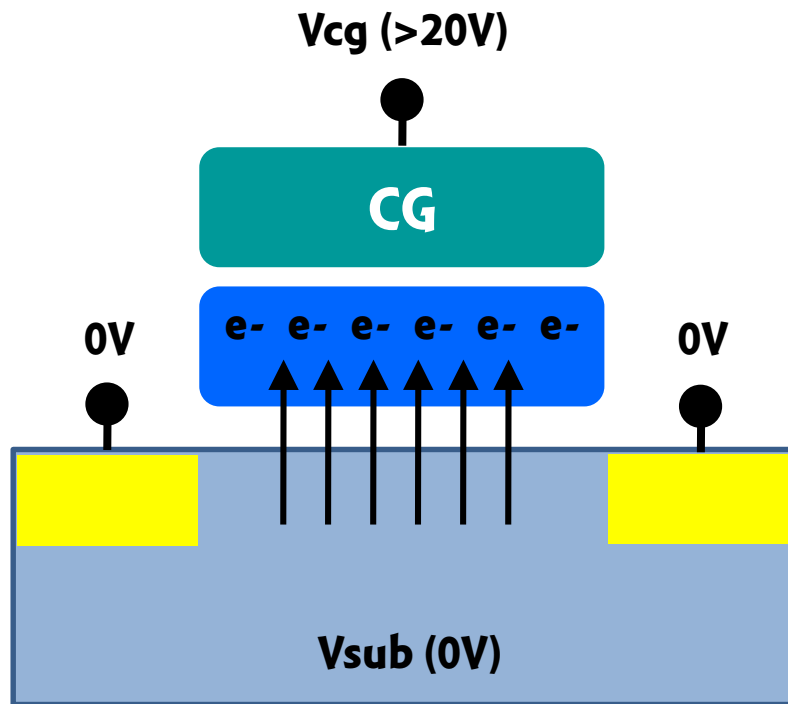
Channel Charge in Floating Gate Transistors



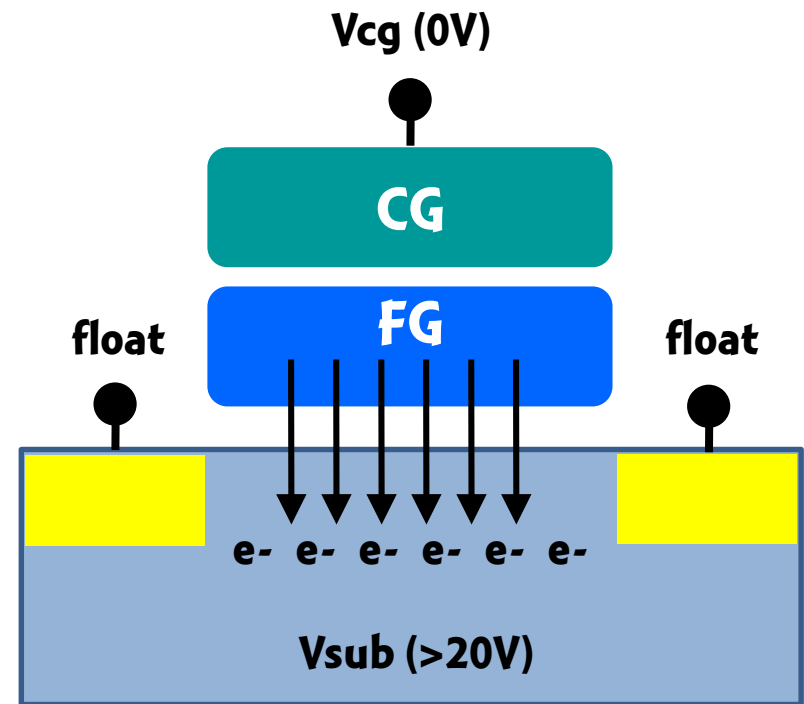
To obtain the same channel charge, the programmed gate needs a higher control-gate voltage than the unprogrammed gate

If we can control the V_{th} ?

- Increase V_{th} : Program
- Decrease V_{th} : Erase
- Determine V_{th} : Read



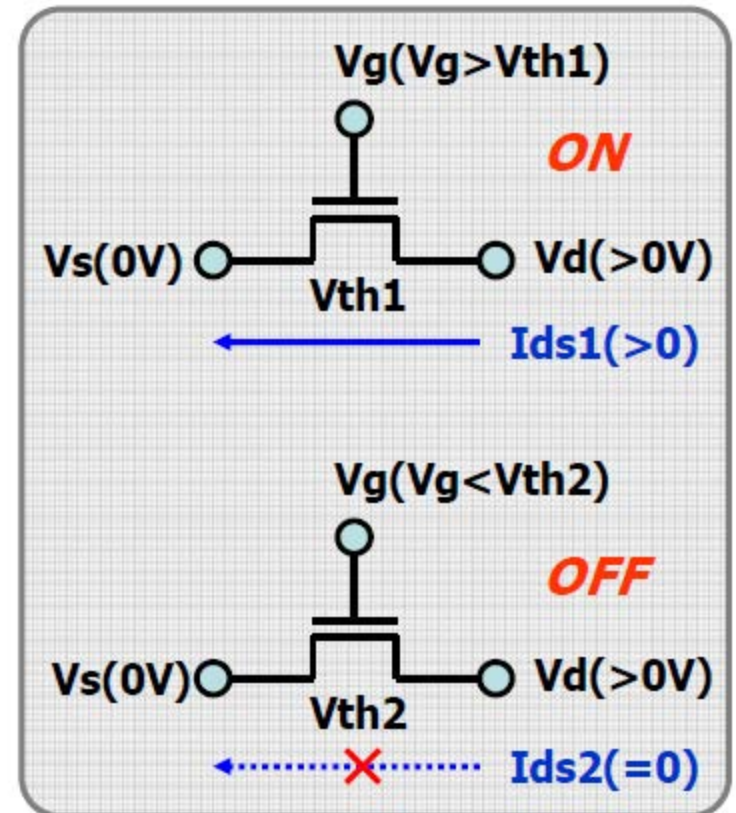
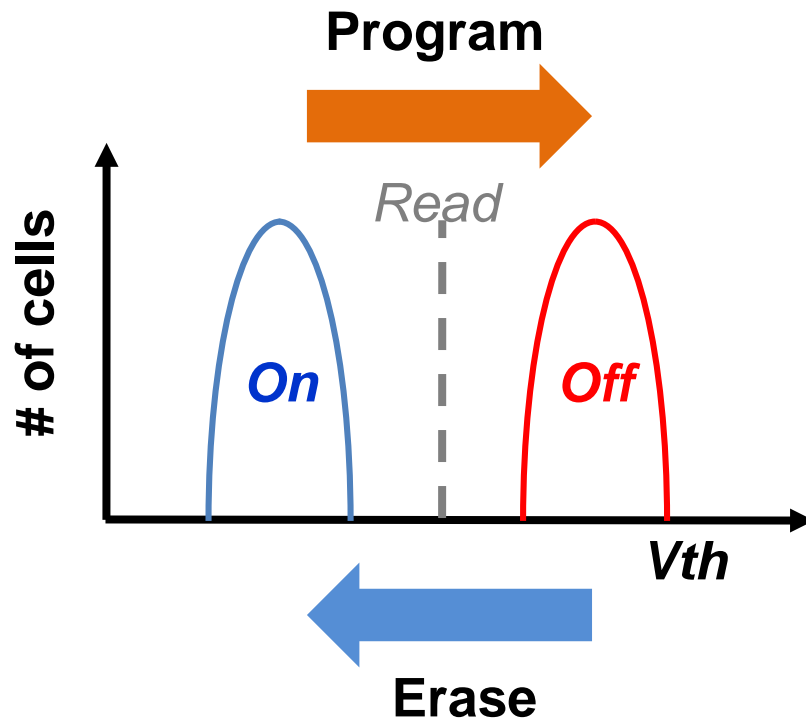
Program : FN tunneling



Erase : FN tunneling

NAND Operation – Overview

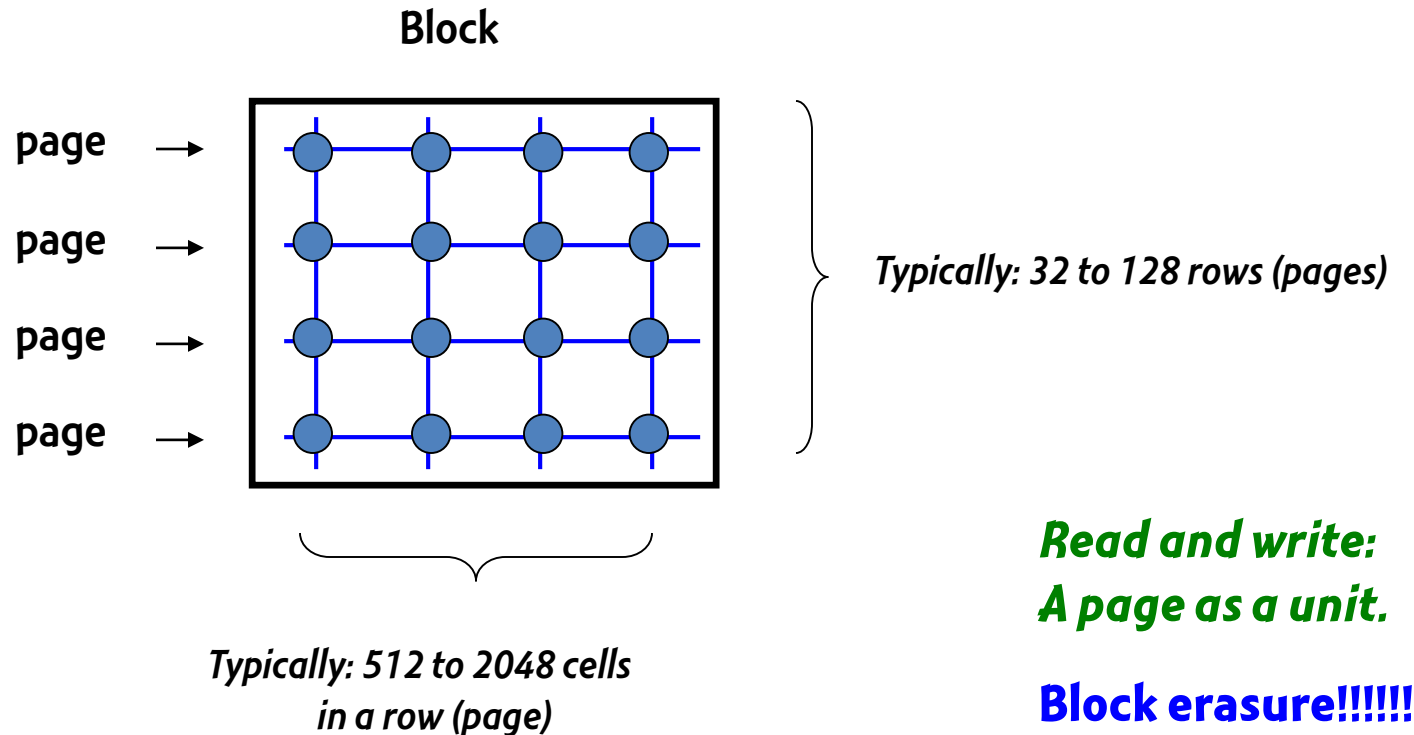
- **Write & Read binary data to a NAND flash cell**
 - Data “0” → Program → Shift cell V_{th} to high → Off state → No current flow
 - Data “1” → Erase → Shift cell V_{th} to low → On State → Current flow



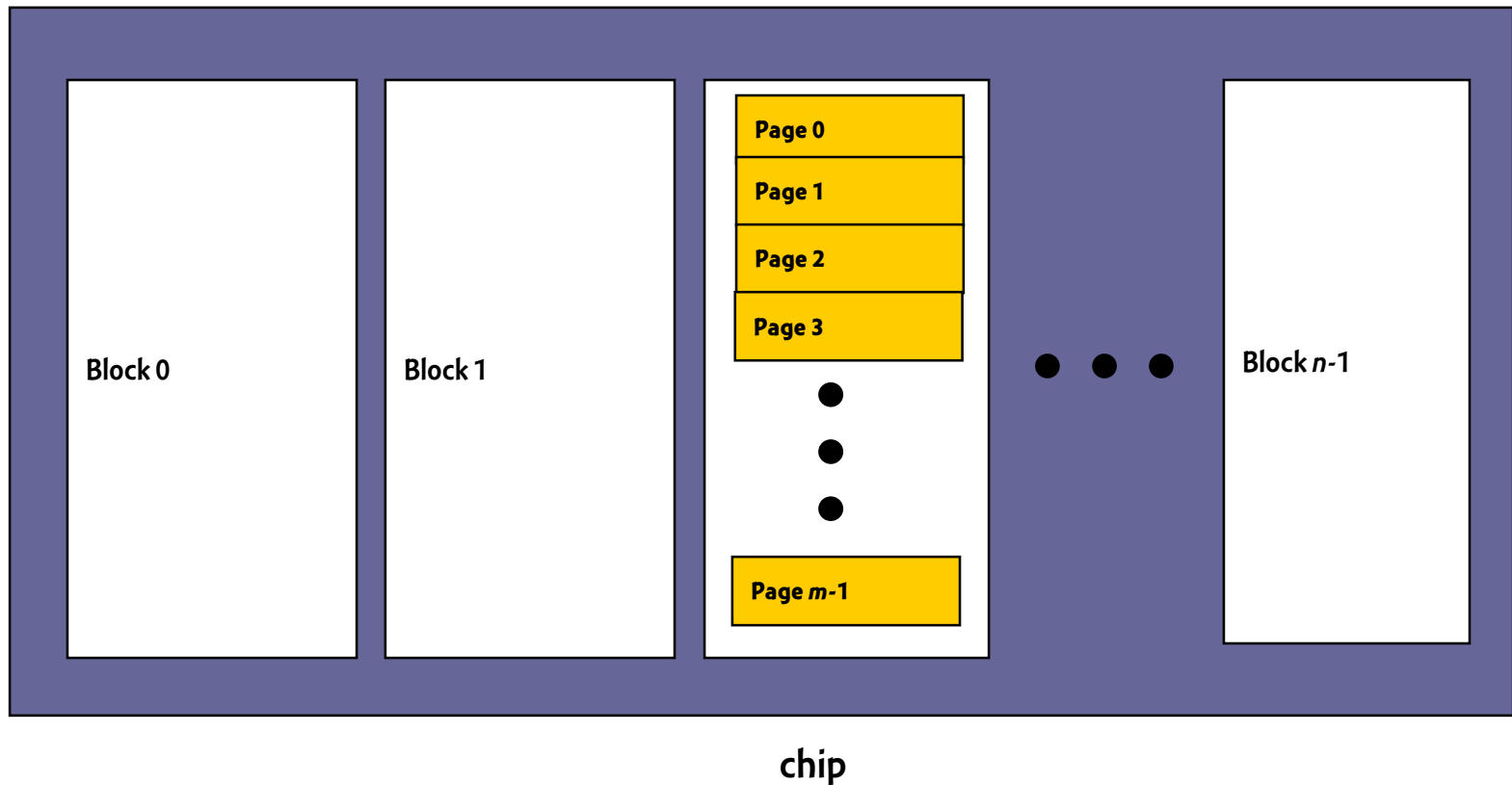
✓ **Read : Check the current flow**

Layout of Block & Page

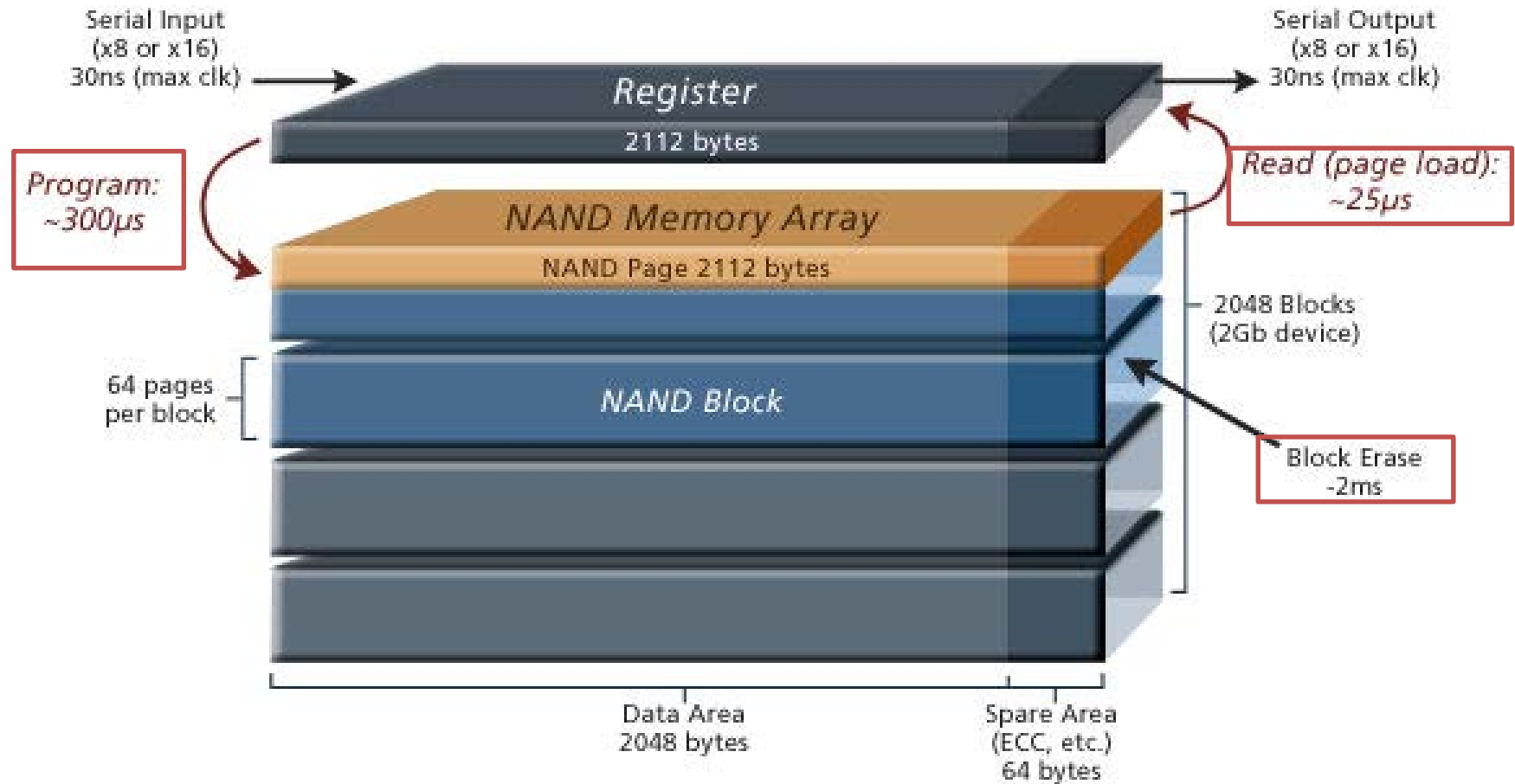
Cells form blocks. Every **block** is an array.
Every row is a **page**.



Organization of NAND Flash Memory

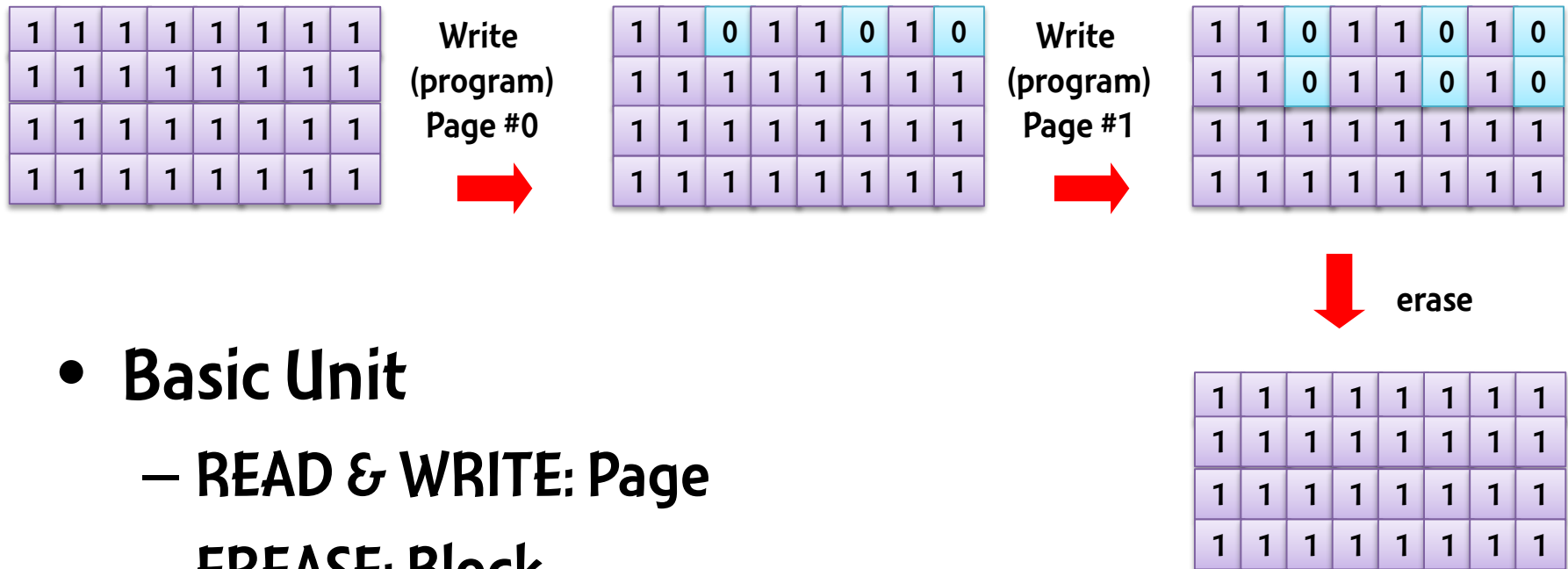


NAND Flash Operations



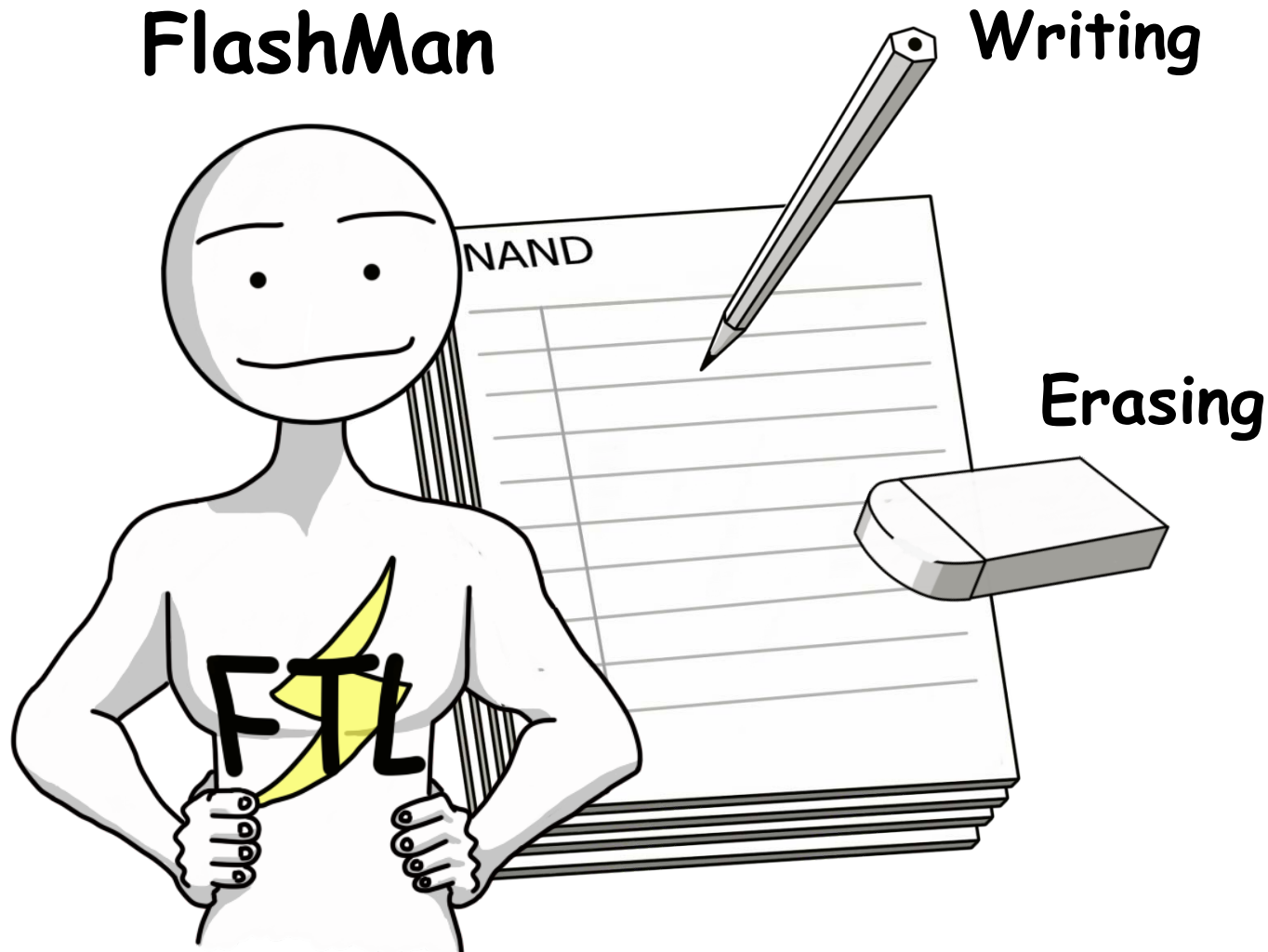
Source: Micron Technology, Inc.

Write & Erase



- **Basic Unit**
 - READ & WRITE: Page
 - ERASE: Block
- Program changes bit from '1' to '0'

NAND Flash Memory is like Sheets of Paper



Illustrated by Jisung Park, Seoul National University

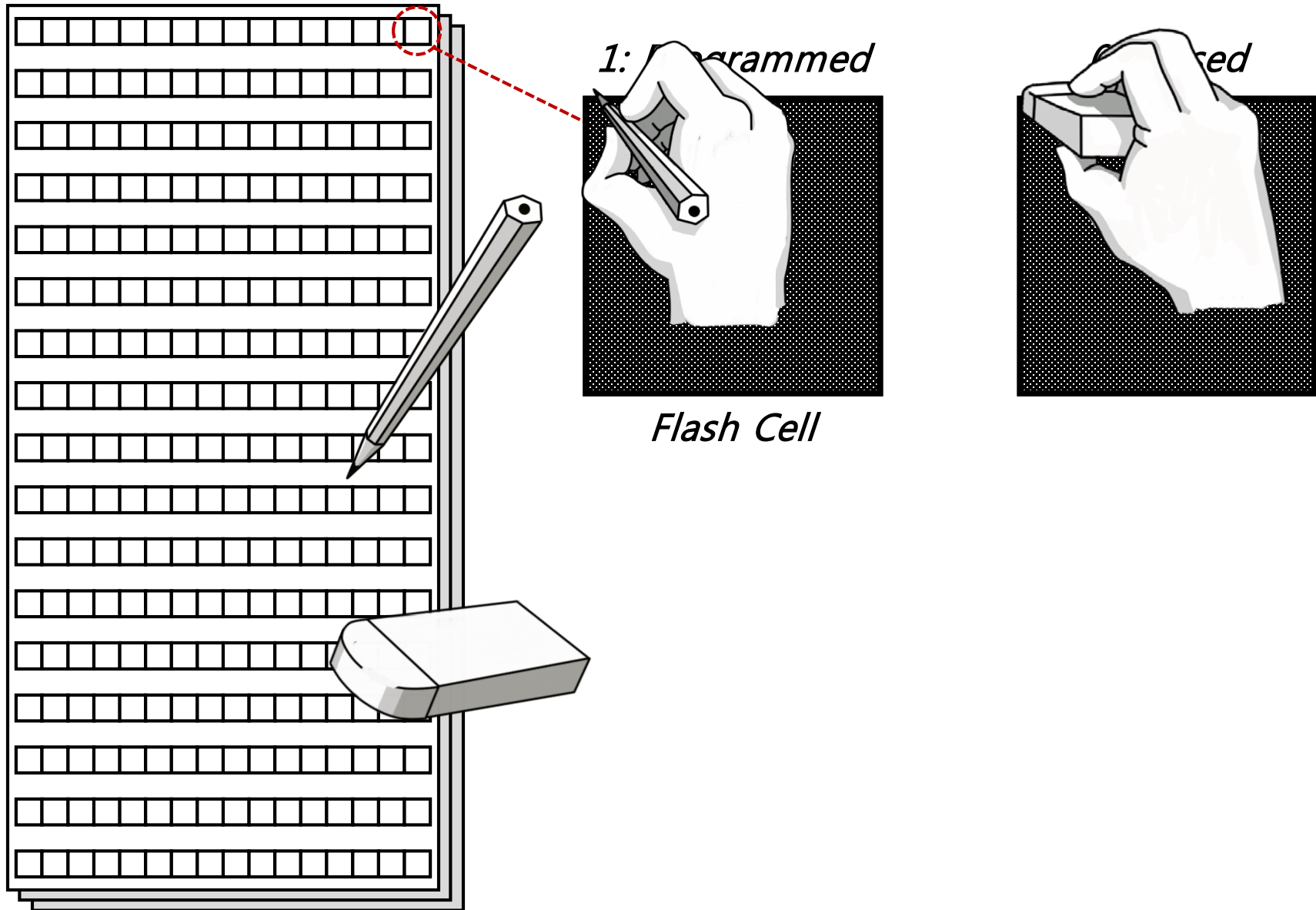
Writing Letters and Erasing Paper



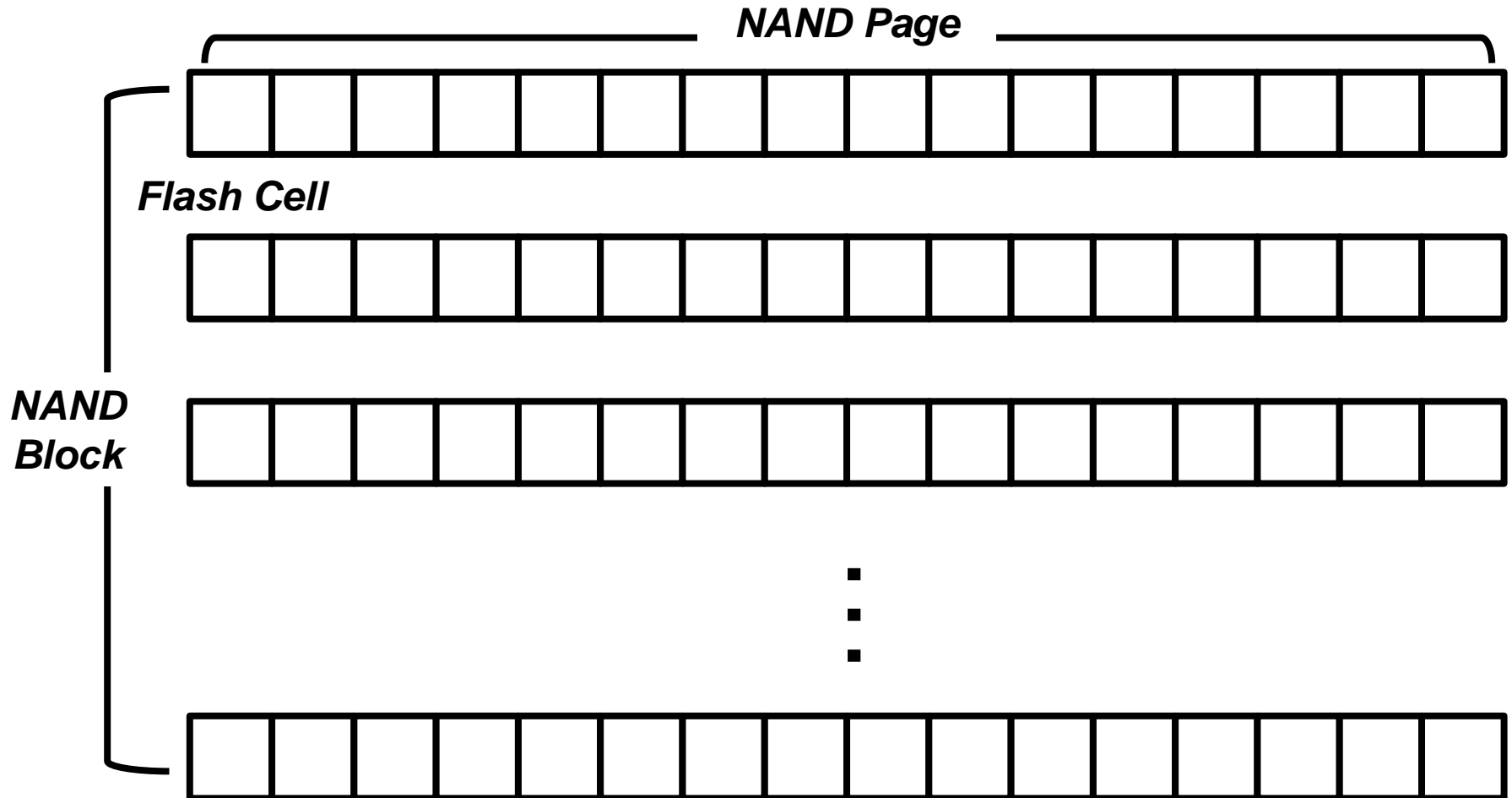
**Out-Place Update
Erase before Program**

(vs. In-Place Update)

NAND Flash Memory: Grid Paper



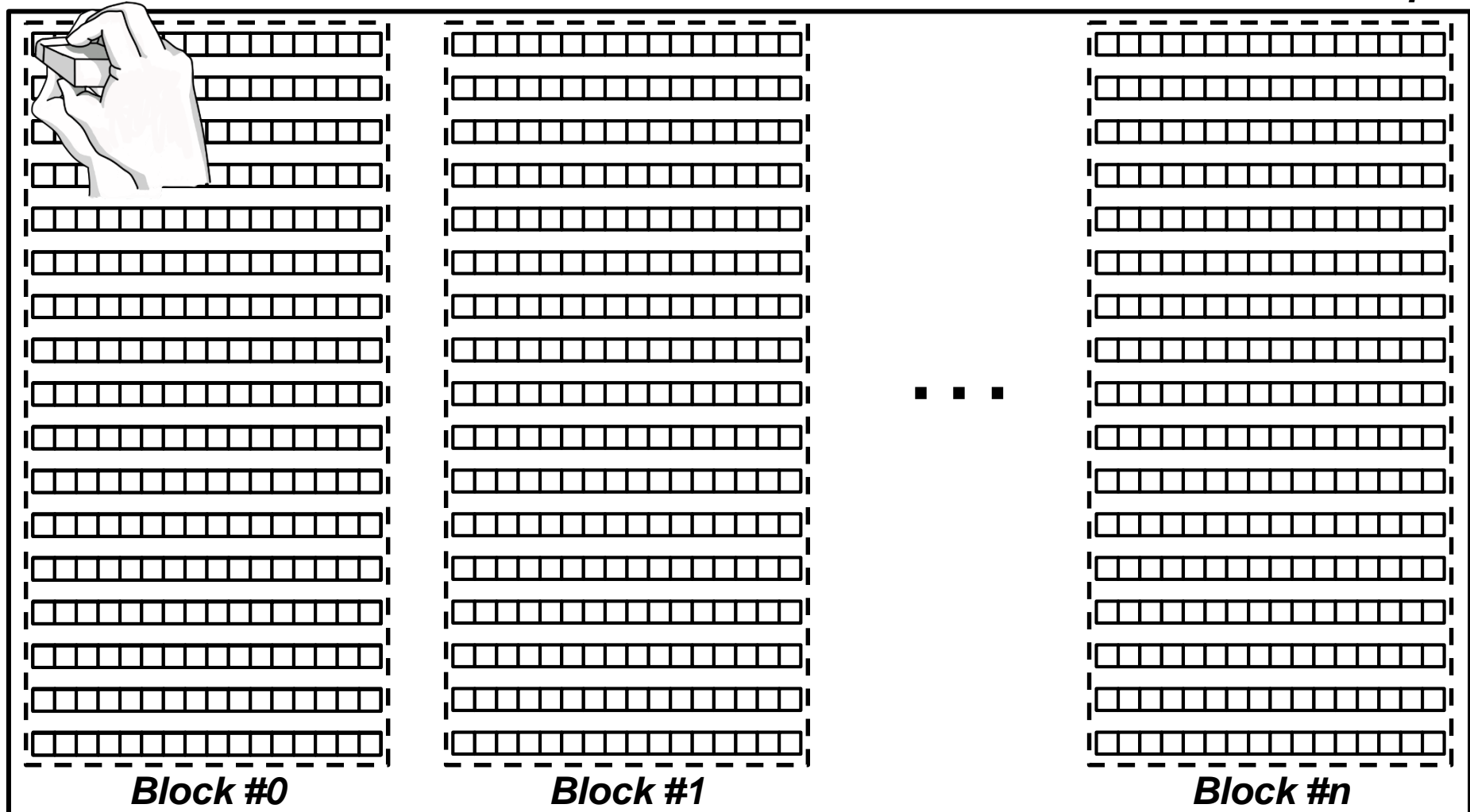
NAND Flash Architecture



NAND Flash Architecture (Cont'd)

Write unit: a page

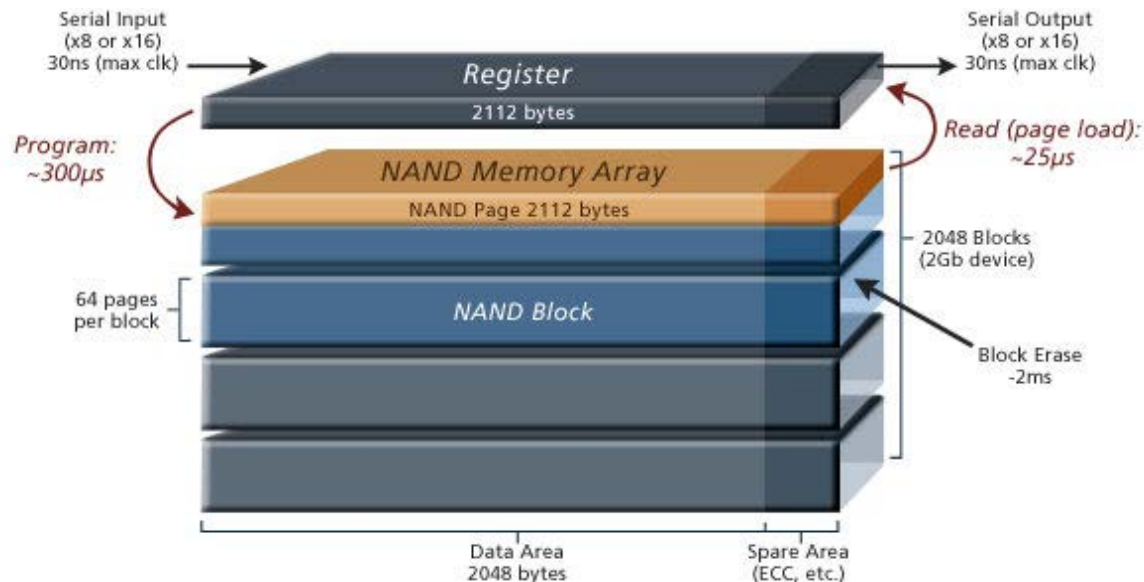
NAND Chip



Erase unit: a block

Unique Properties of NAND Flash Memory

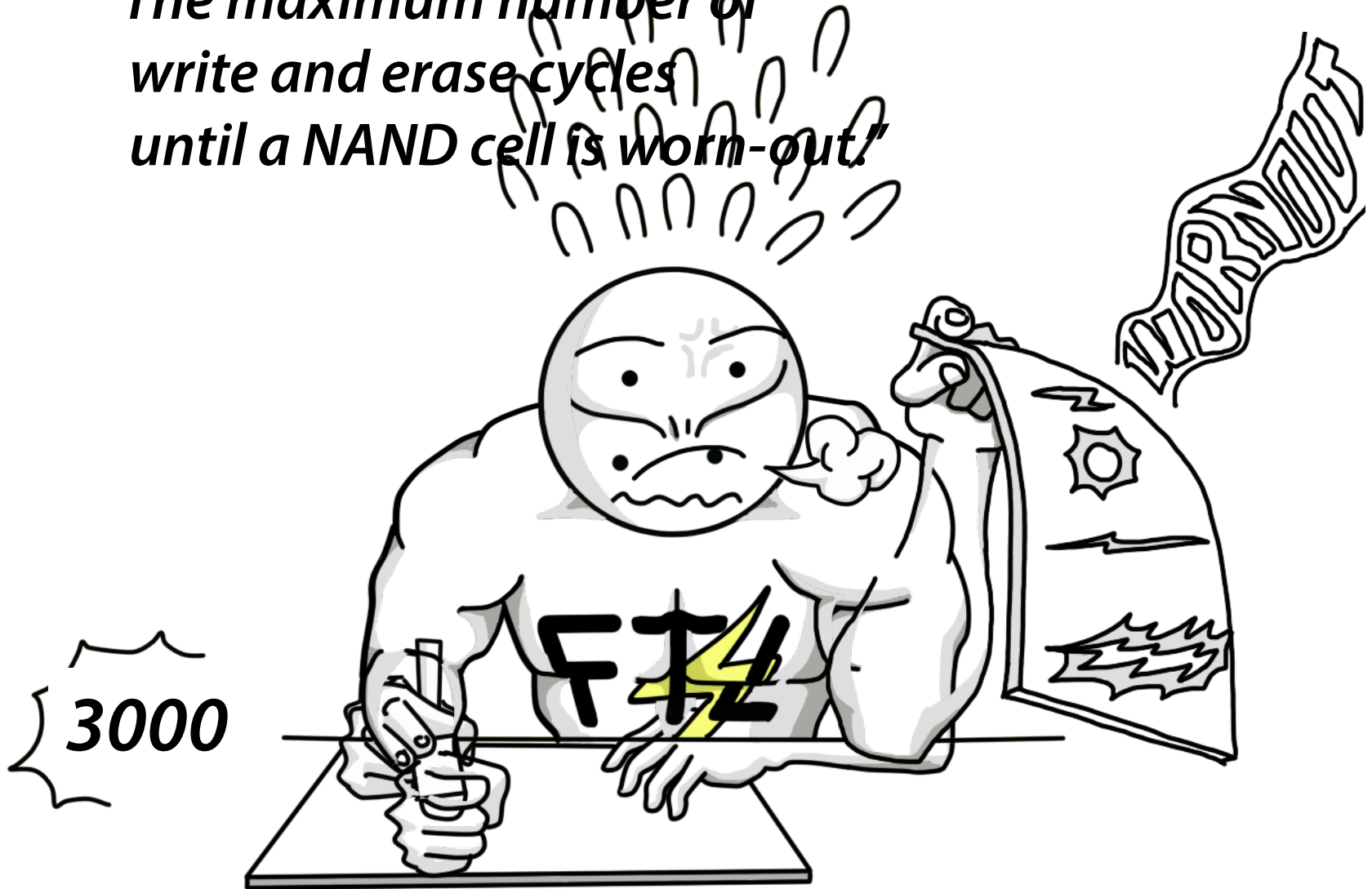
- No support for in-place update
 - **Erase-before-program** operation
- Asymmetric operation unit/performance
- Limited Endurance
 - A block becomes unreliable after a certain number of program/erase (P/E) cycles



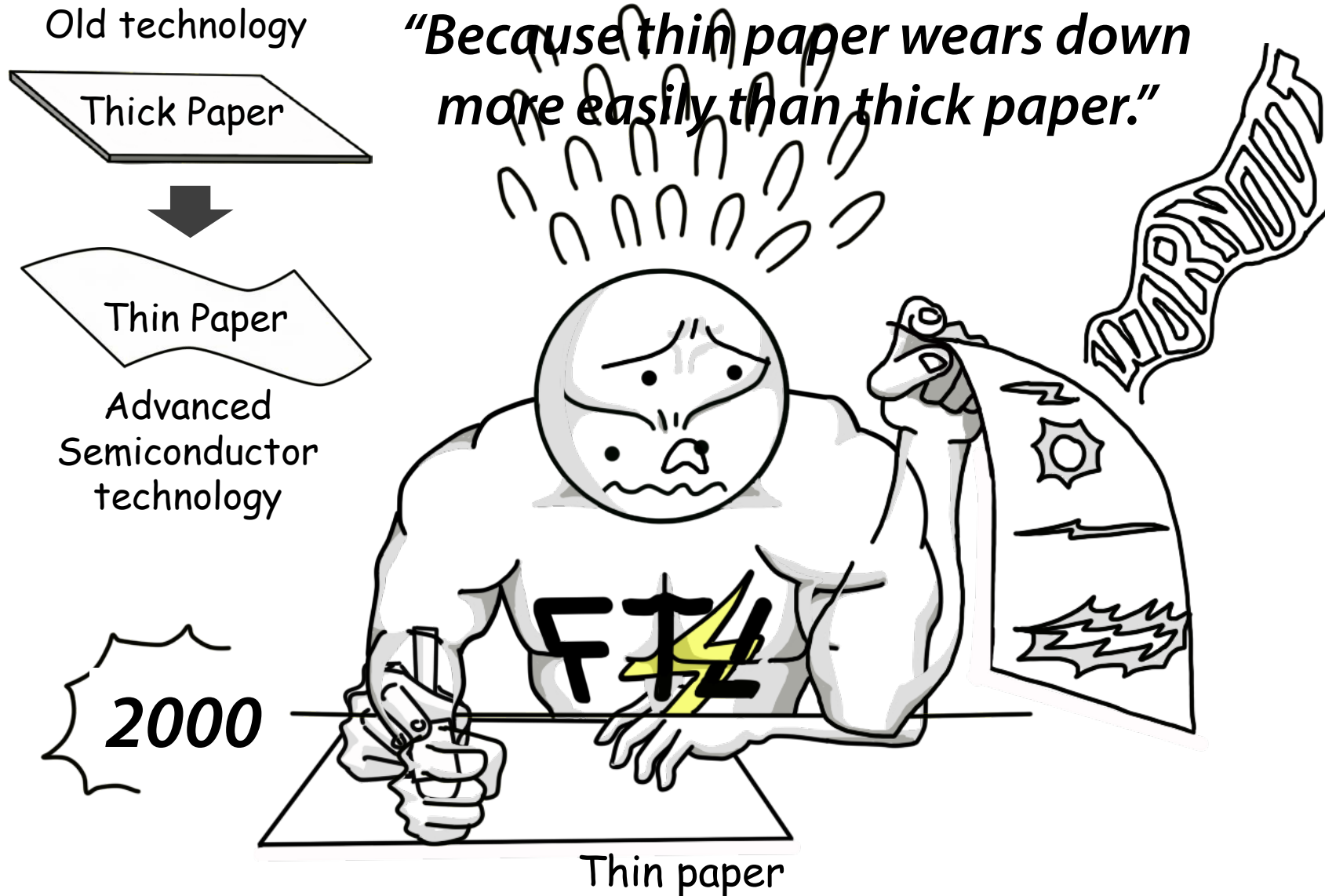
NAND flash memory organization

Limited NAND Endurance

"The maximum number of write and erase cycles until a NAND cell is worn-out."

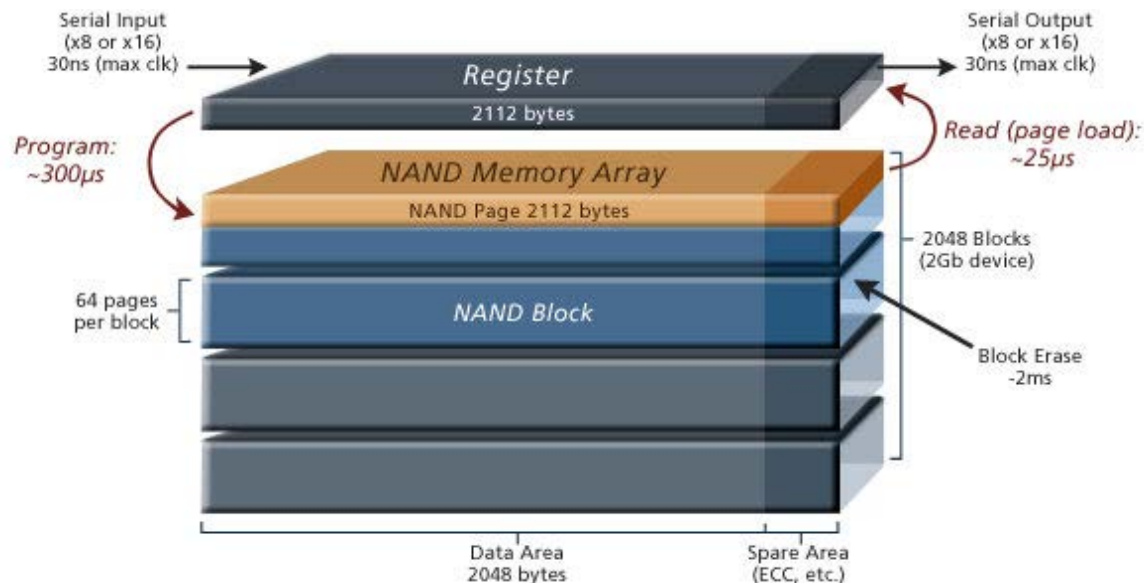


Why is the NAND Endurance Decreased ?



Unique Properties of NAND Flash Memory

- Asymmetric operation unit/performance
- No support for in-place update
 - **Erase-before-program** operation
- Limited Endurance
 - A block becomes unreliable after a certain number of program/erase (P/E) cycles



NAND flash memory organization

Like Papers ...

- Data stored in NAND flash **FADES!!**
 - **Not ideal for storage systems**
- The data dissipation speed depends on
 - The number of program/erase cycles
 - Exposure to high temperatures
 - The number of reads
 - Process technology
- Solution:
 - move data to a new location by a flash controller

Data Retention

[고객용]

농협중앙회

신협중앙회

서울대학교구내인경권

백준현

02-383-3000

서울 관악구 관악로 11111121

거래일시: 16/06/25 15:31:56

카드번호: 5461-1110-****-7641

승인번호: 30000400

카드종류: NH체크카드/신협중앙회

전표번호: 9000-0081-6727

부가세물품가액

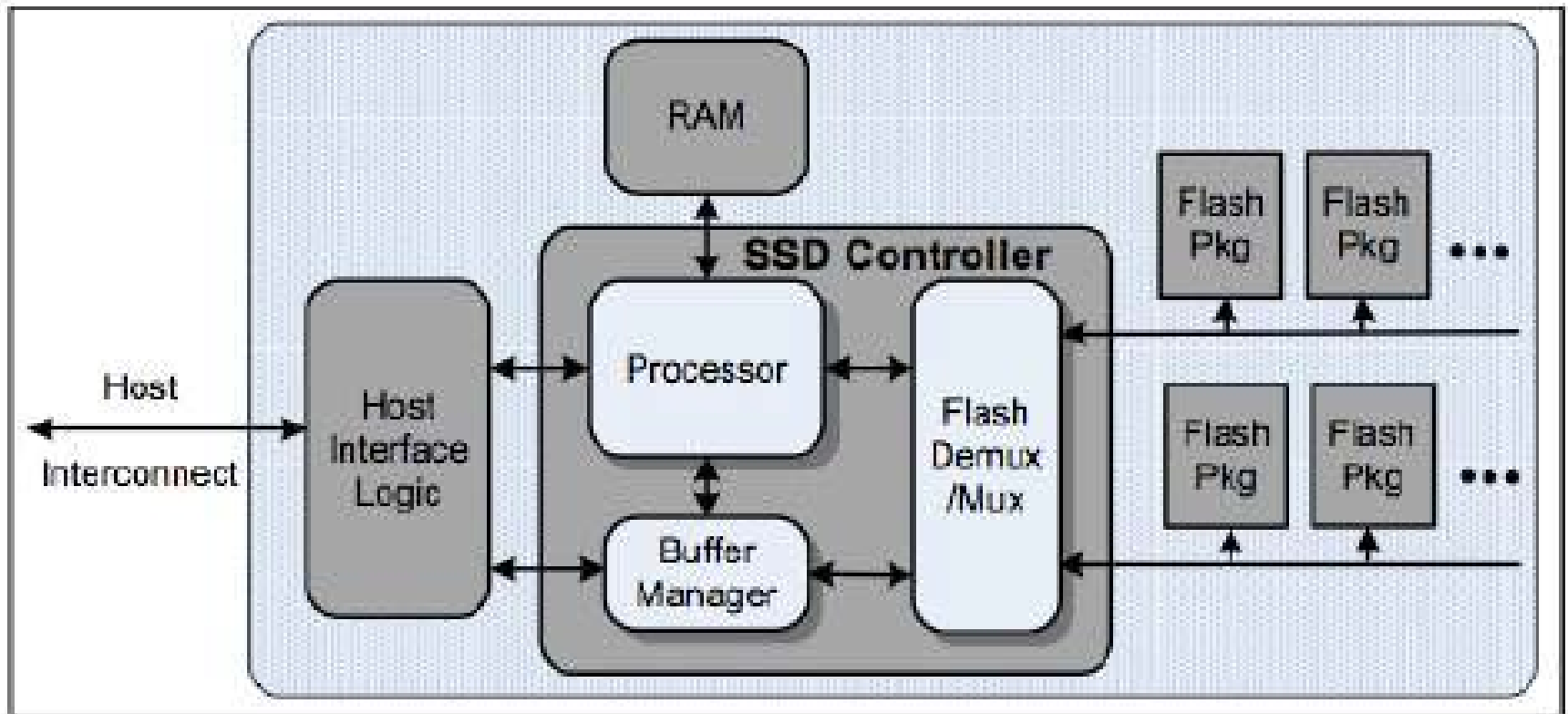
부 가 가 치 세 :

부 가 금 액 :

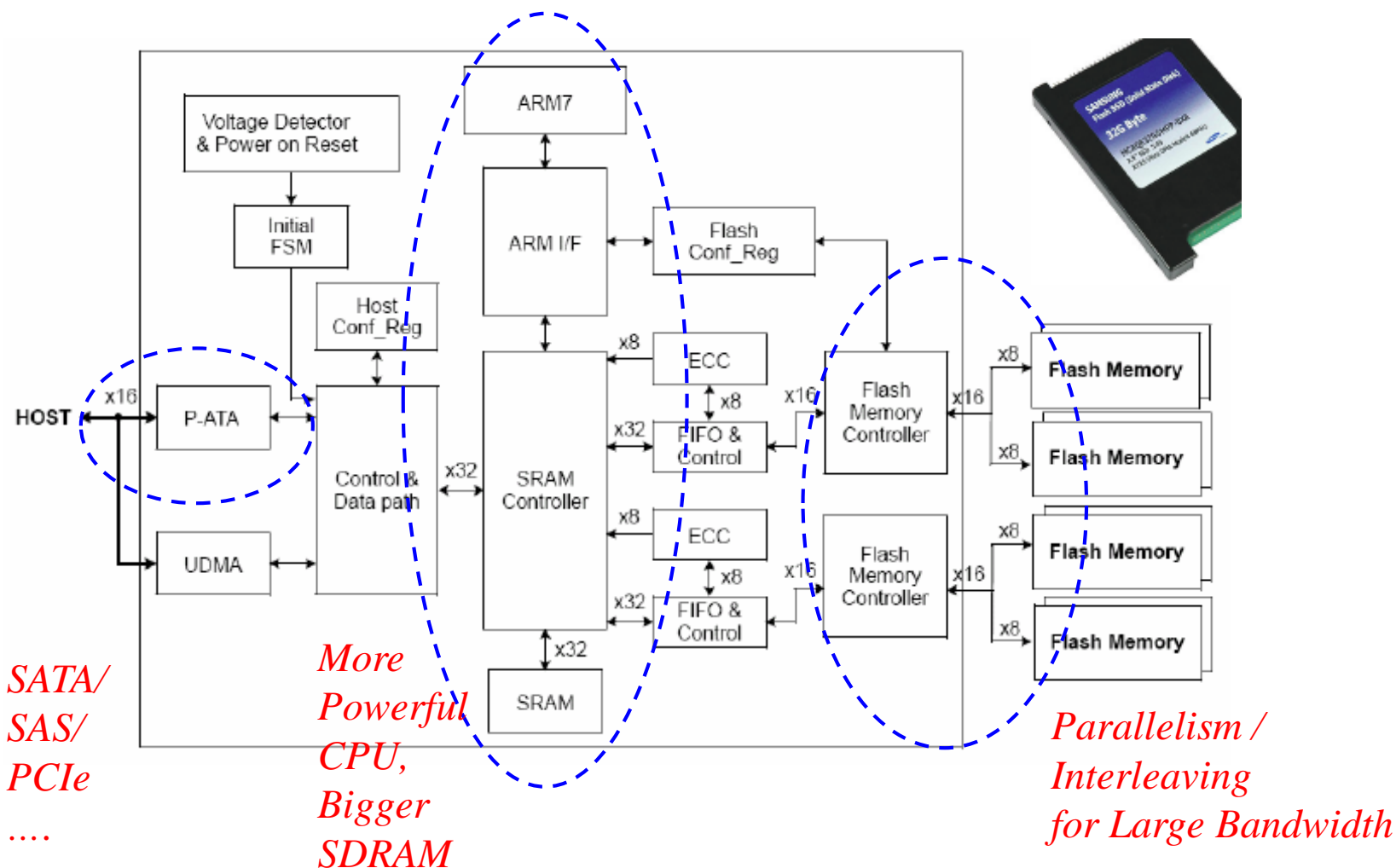
1 Year at 30 °C?

Flash Controller for Performance and Reliability

Overview of Controller



Flash Controller Diagram



Flash Controller

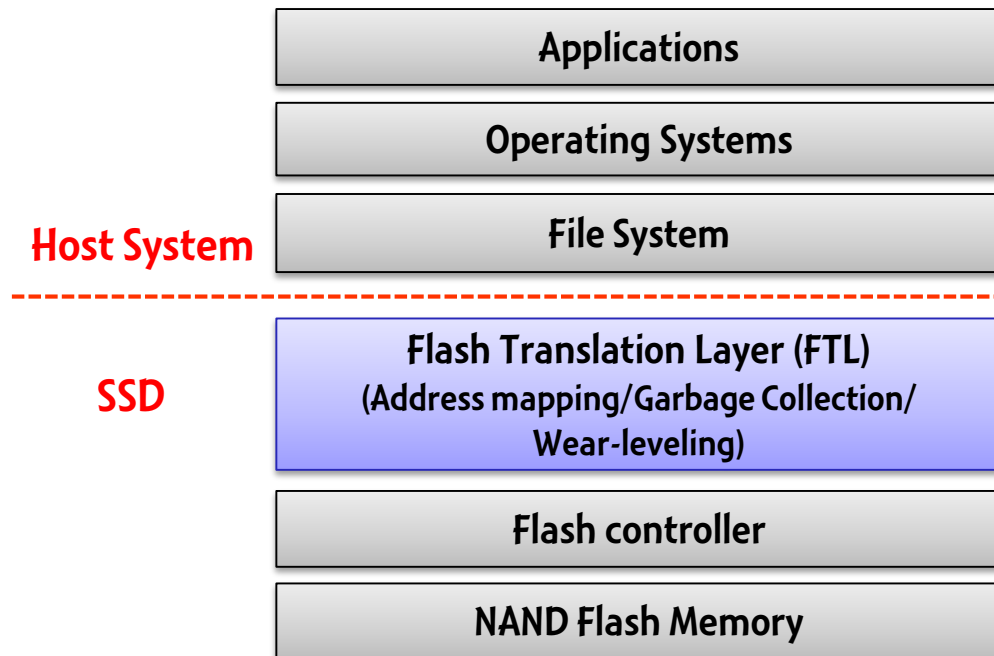
- **Make imperfect NAND flash robust and reliable**
- **Based on a single ASIC die**
 - **SRAM for firmware**
 - **DRAM for caching/buffering**
 - **Backup power system (e.g., batteries, capacitors) for sudden power off**
- **Parallelism support for high performance**
- **Various reliability functions for flash memory**

Error Correction

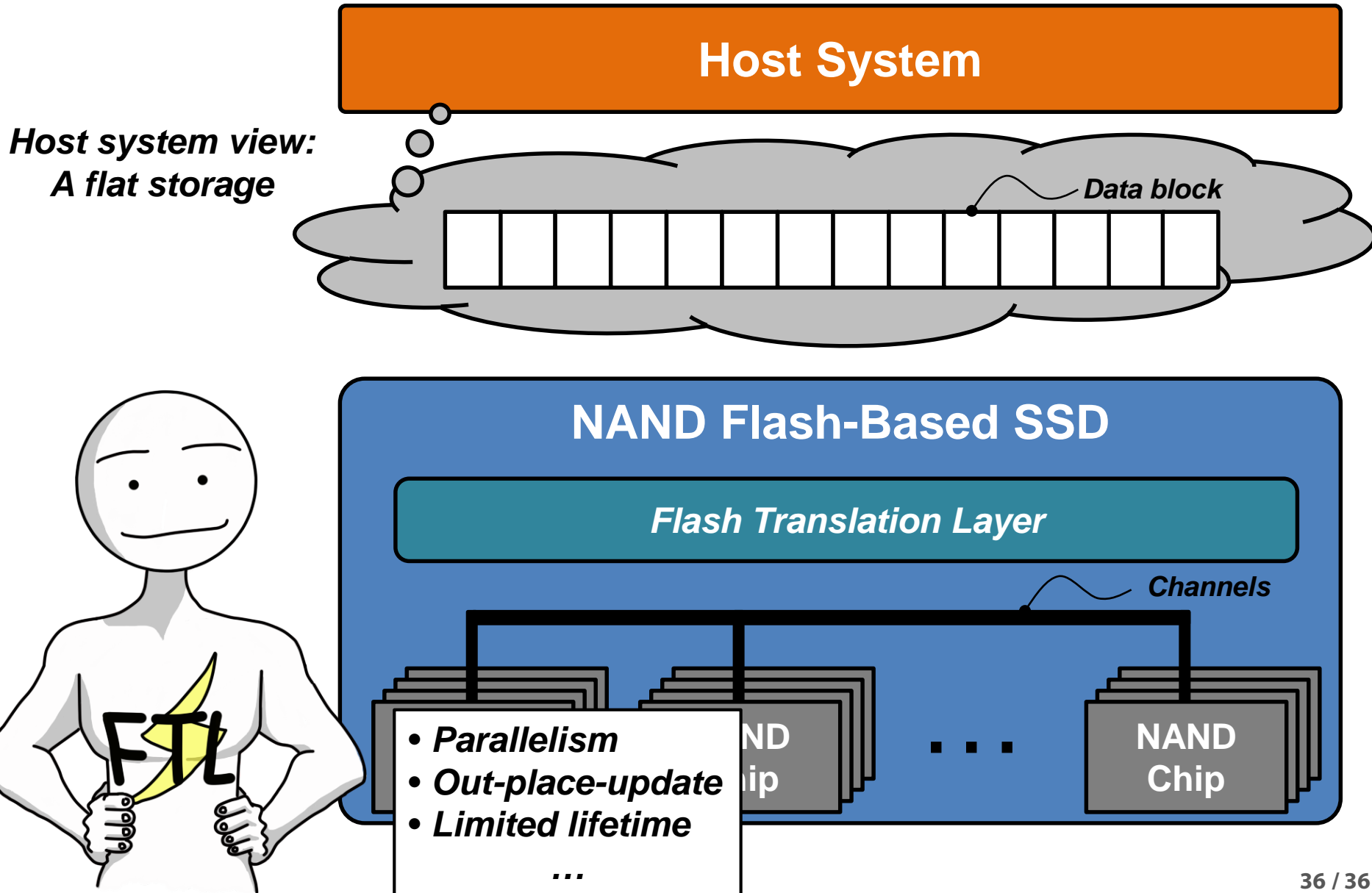
- **For all reads and writes from/to flash, error-correcting code is used.**
- **BCH**
- **LDPC**
- **XOR/Scramblers**

SSD Software

- Overcome the physical restrictions by employing system software called a flash translation layer (FTL)
 - Asymmetric operation unit/performance
 - No support for in-place update
 - Limited Endurance
- Address mapping/
Garbage collection
- Wear-leveling



Flash Translation Layer (FTL)



Address Translation

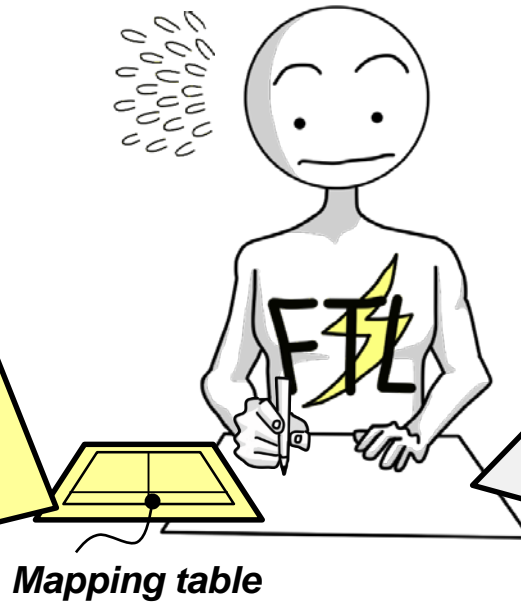
Write Request
Read Request
@ (logical) page 0
@ (logical) page 0

Host system






Flash-based SSD

Invalidation

logical page address (lpa)	physical page address (ppa)
0	12
1	11
2	
3	
4	
5	
...	...

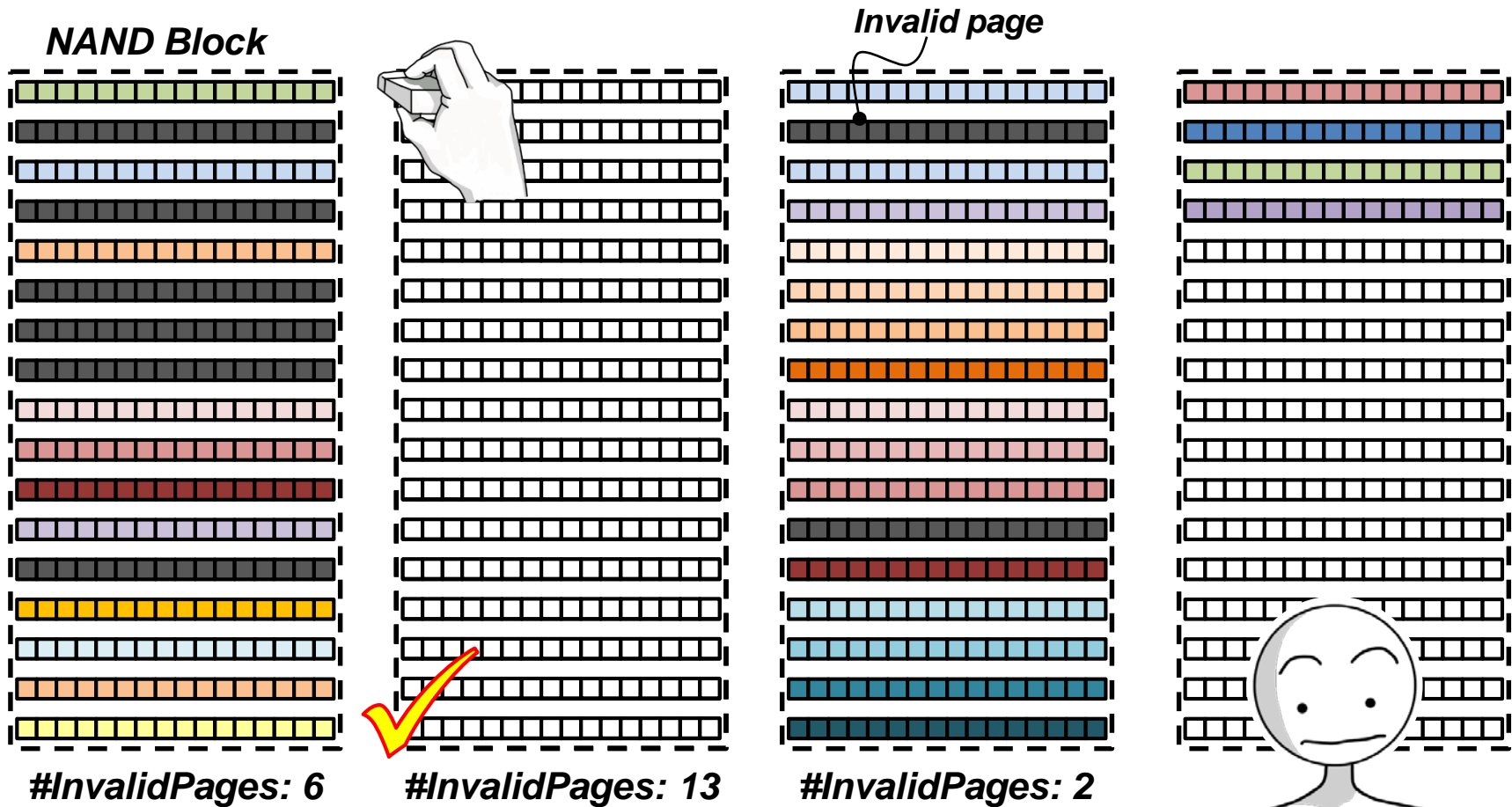


No In-Place-Update

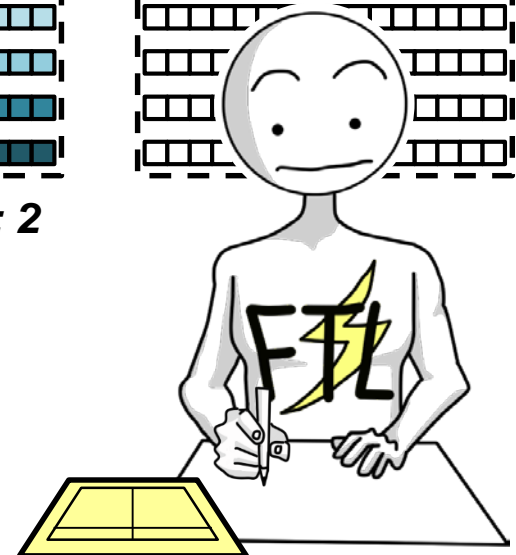
Page 10 
Page 11 
Page 12 
Page 13 
Page 14 

⋮

Garbage Collection

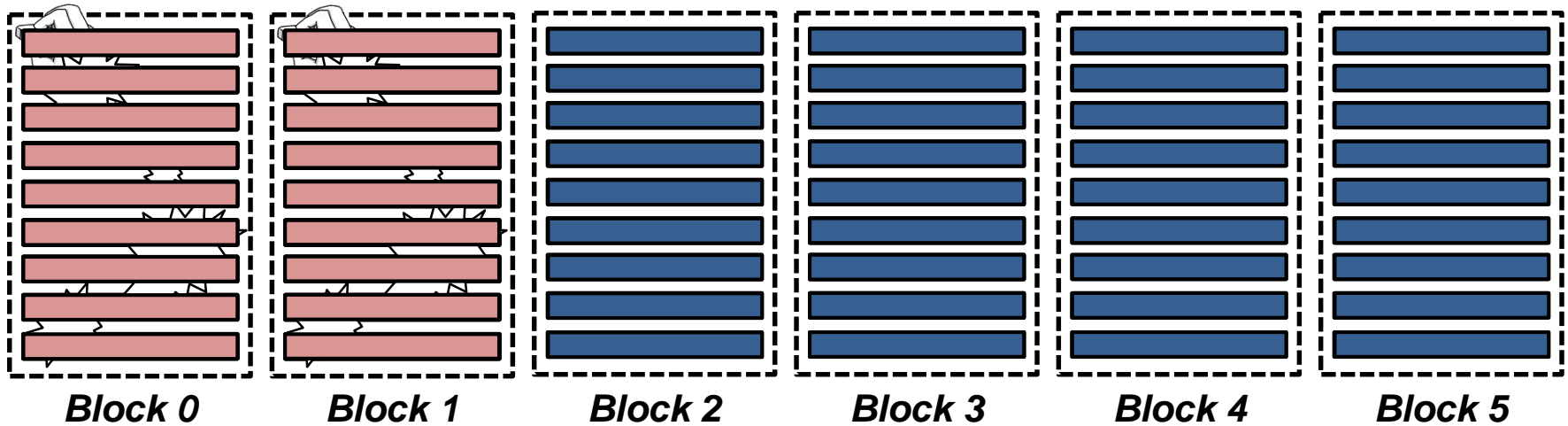


- 1) Choose a victim block
- 2) Copy valid pages and update mappings
- 3) Erase the victim block



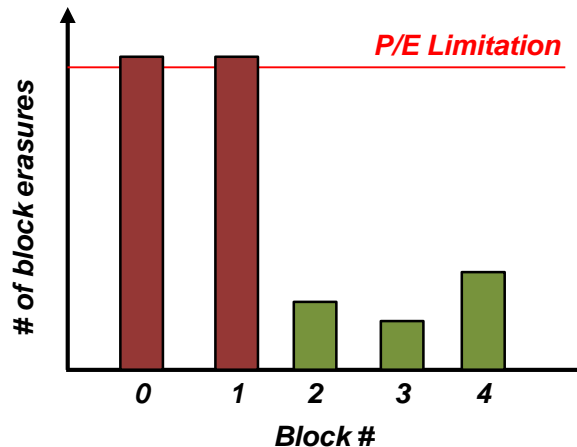
Wear Leveling

■ Invalid ■ Hot (frequently updated) ■ Cold (rarely updated)



Skewed block usage!

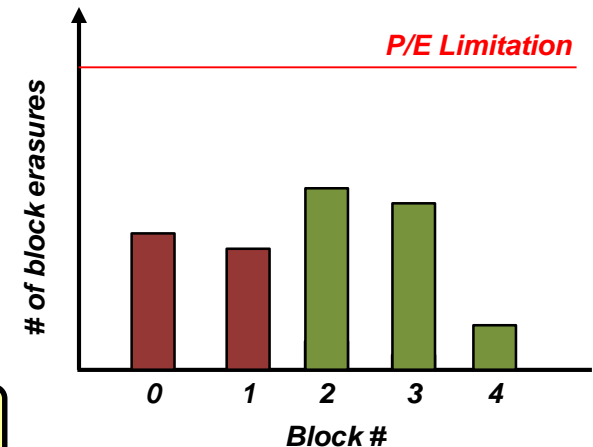
Short device lifetime!



Hot-cold swapping



Data Migration Department

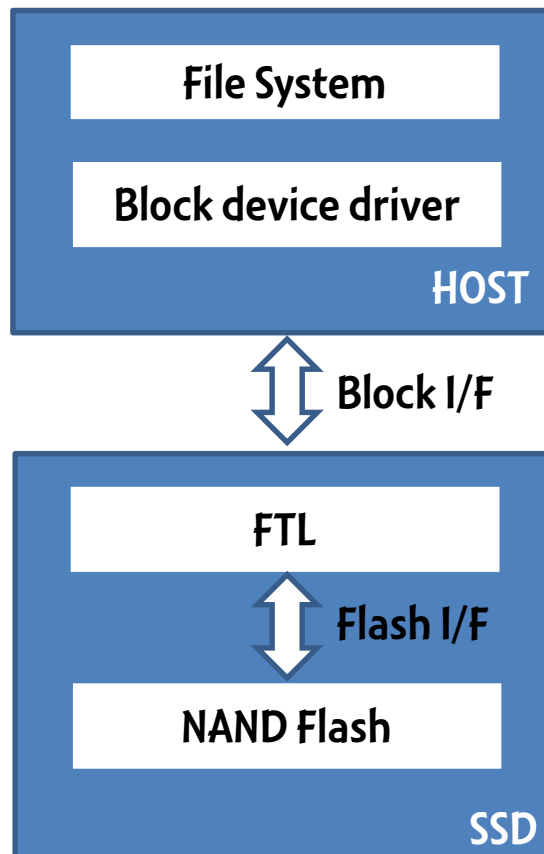


Controller Firmware

- **Flash Translation Layer (FTL)**
 - Mapping blocks
 - Garbage Collection
 - Wear Leveling

Flash Translation Layer

- A software layer to make NAND flash emulate traditional block devices (or disks)



Flash Management Tasks

- **Essential**

- **Address Translation**

- Avoid in-place update
 - Logical Block Address (LBA) -> Physical Block Address (PBA)

- **Garbage Collection**

- Reclaim invalid blocks -> Get new free blocks

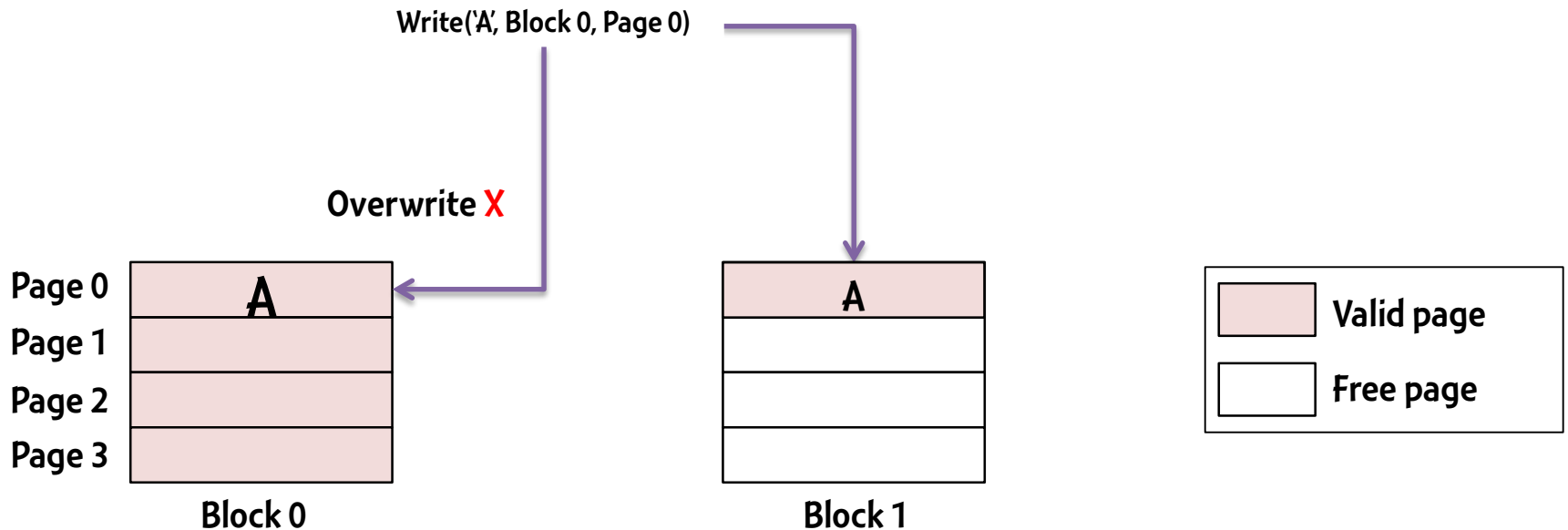
- **Optimization**

- **Wear Leveling**

- Erase all blocks evenly -> Extend life time

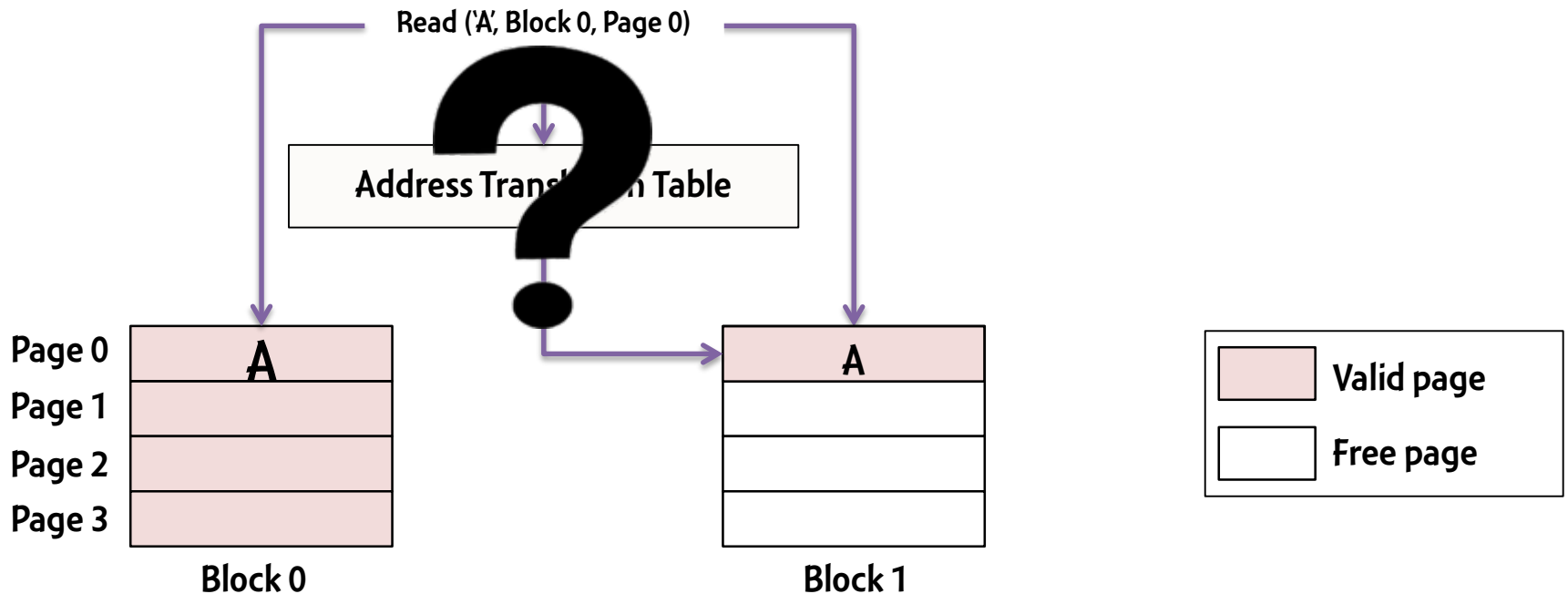
Out-place Update

File Systems



Address Translation

File Systems



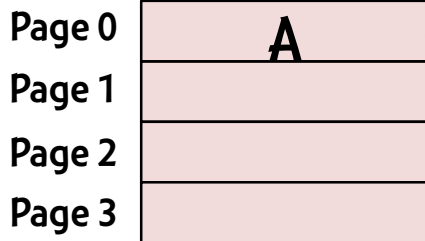
Garbage Collection

- NAND flash memory does not allow *in-place update*

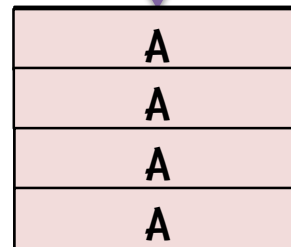
File Systems

Write('A', Block 0, Page 0) X4

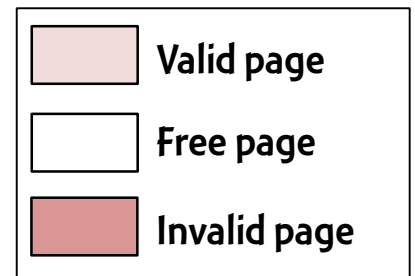
Overwrite **X**



Block 0



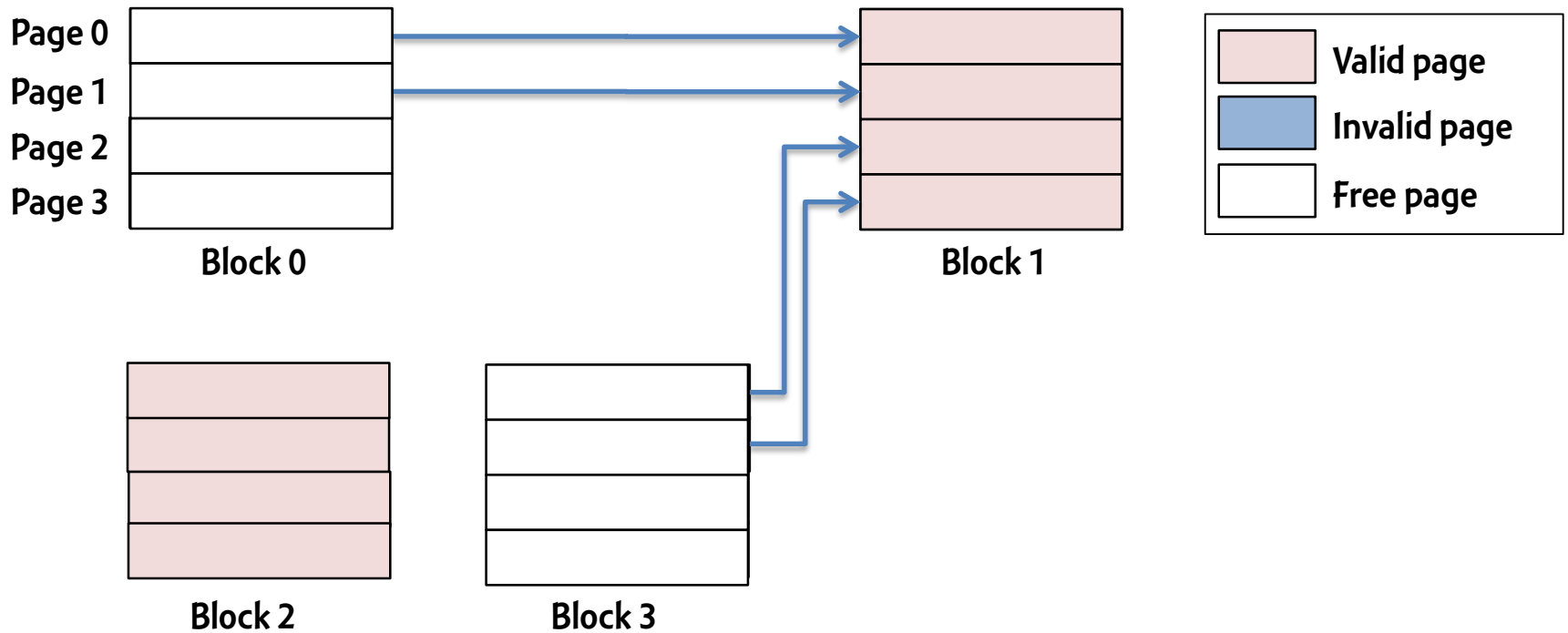
Free Block



- NAND flash memory will be full of invalid data...

Garbage Collection

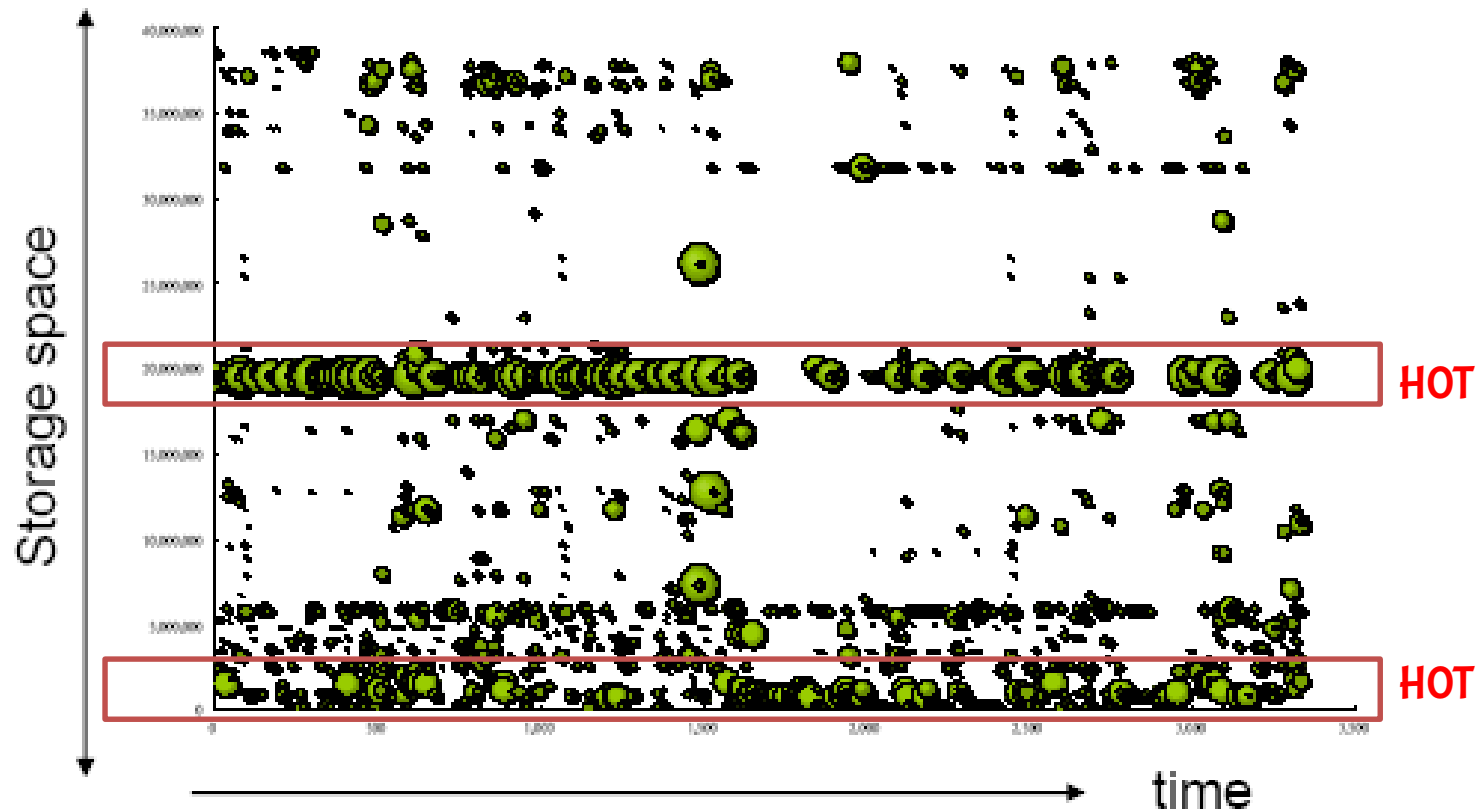
- Gather valid data -> Erase invalid data*



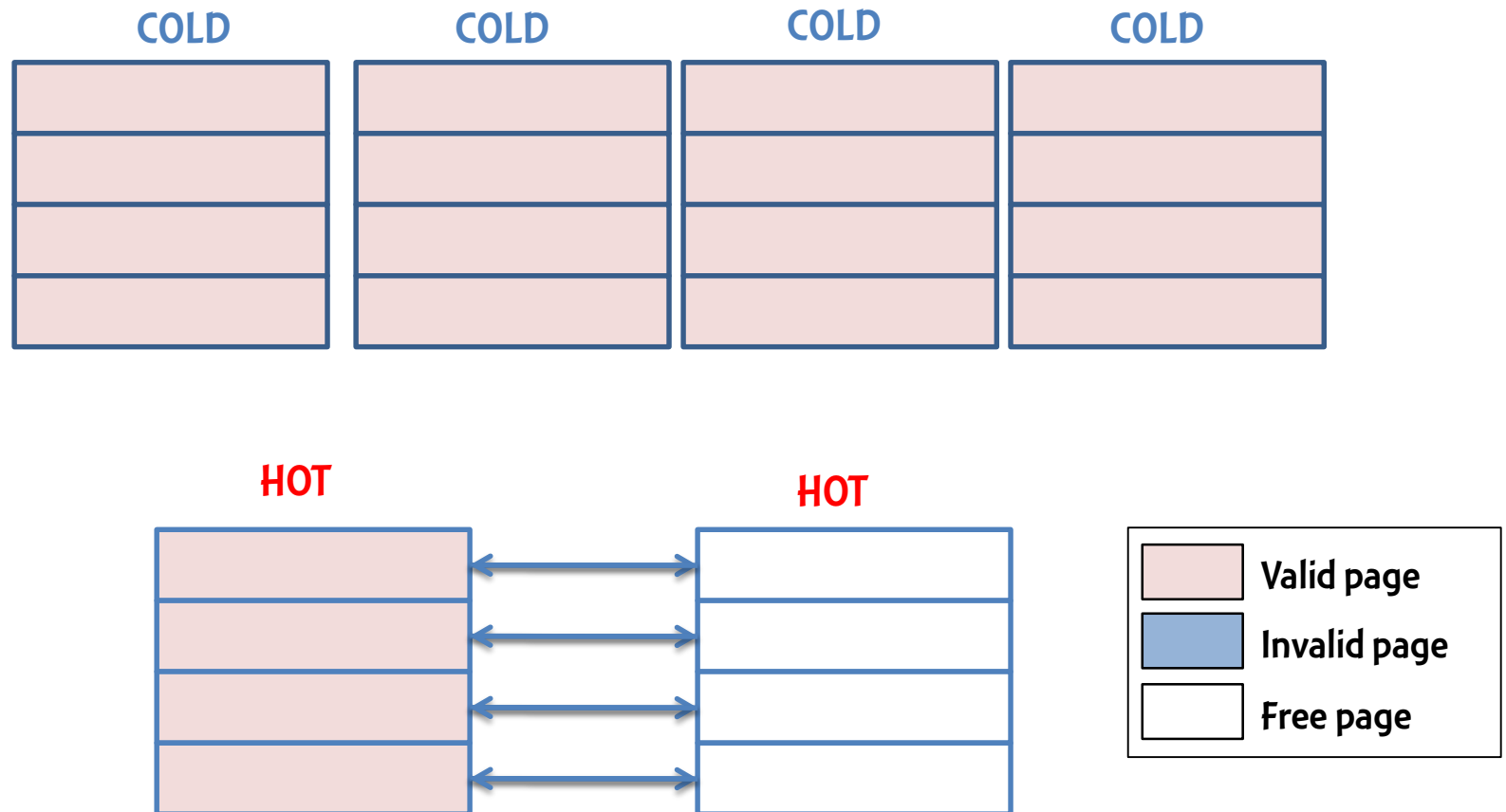
Flash Management Tasks

- **Essential**
 - **Address Translation**
 - Avoid in-place update
 - Logical Block Address (LBA) -> Physical Block Address (PBA)
 - **Garbage Collection**
 - Reclaim invalid blocks -> Get new free blocks
- **Optimization**
 - **Wear Leveling**
 - Erase all blocks evenly -> Extend life time

Spatial Locality of Write Request



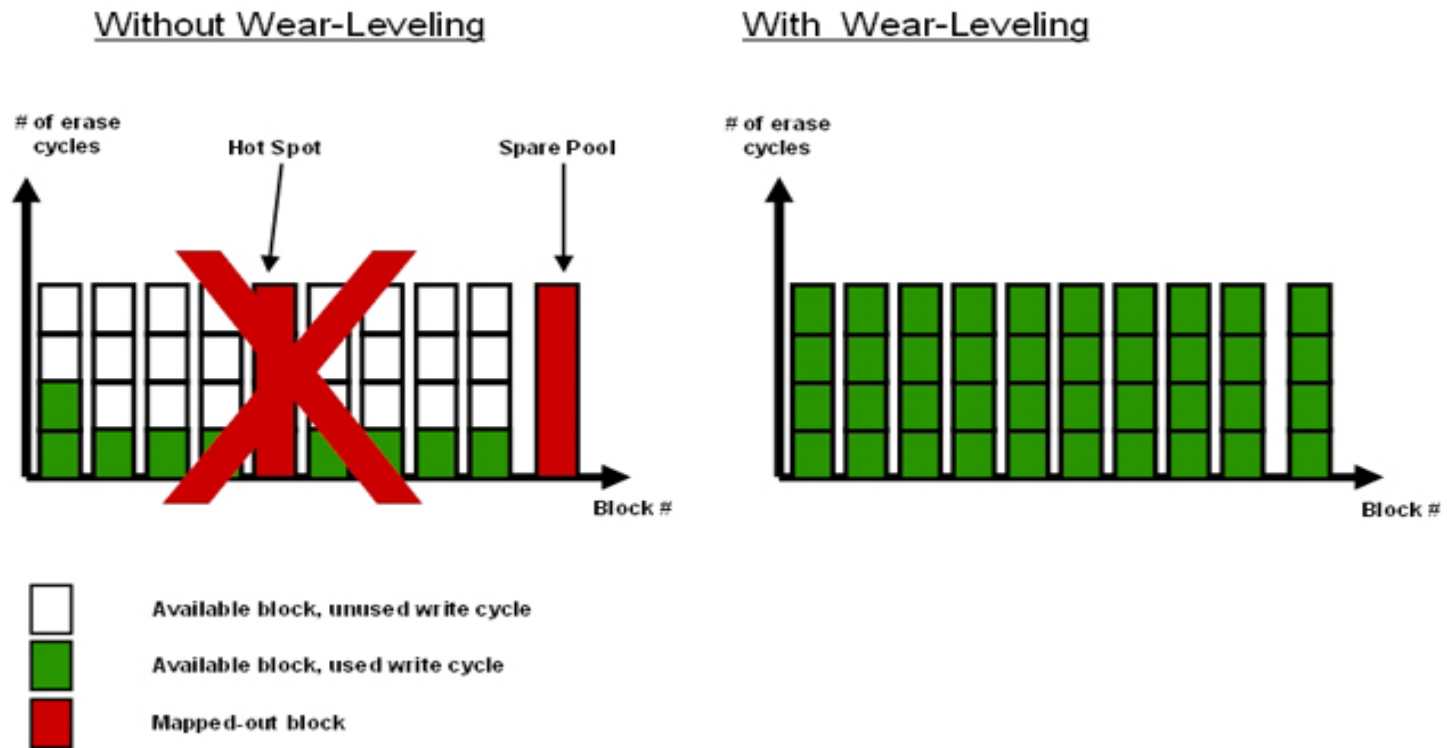
Wearing of Flash Memory Blocks



- Blocks with hot data are likely to be erased more

Wear Leveling

- Erase evenly all blocks



New Interface for SSDs

- **TRIM/UNMAP**
 - Let the SSD clear the LBA entries in the FTL, giving more free space to use
- **Scatter Gather**
 - Reduce command overhead of one command at a time
 - Gathers multiple noncontiguous requests into a single command, reducing overhead

Host Interface

Existing Interface: SATA, SAS

New interface: PCIe