

# I/O and Storage: Flash SSDs

*CS 571: Operating Systems (Spring 2020)*

Lecture 10a

Yue Cheng

# Disk vs. Flash

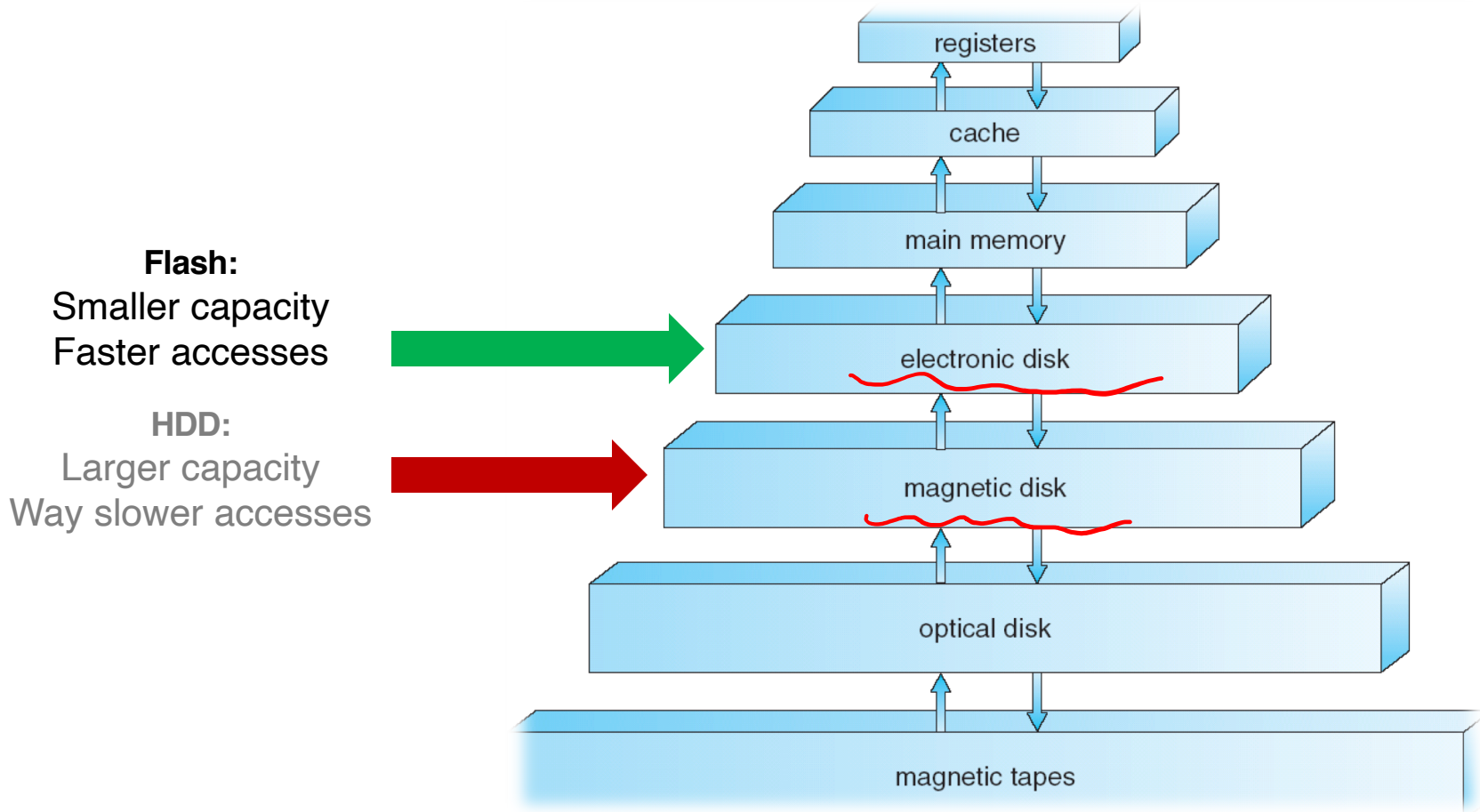
# Disk Overview

- I/O requires: seek, rotate, transfer
- Inherently:
  - Not parallel (only one head)
  - Slow (mechanical)
  - Poor random I/O (locality around disk head)
- Random requests each taking ~10+ms

# Flash Overview

- Hold charge in cells. No moving (mechanical) parts!
- Inherently parallel!
- No seeks!

# Storage Hierarchy Overview



# Disks vs. Flash: Performance

# data / time unit.

## • Throughput

- Disk: ~130MB/s (sequential)
- Flash: ~400MB/s

## • Latency

small, random I/Os.

block  
4KB

- Disk: ~10ms (one op)
- Flash:

1ms = 1k  $\mu$ s

- Read: 10-50us Fastest
- Program: 200-500us
- Erase: 2ms

Write ops. }

per-page

↑  
coarse-grained  
per-block

~~in-situ overwrite~~

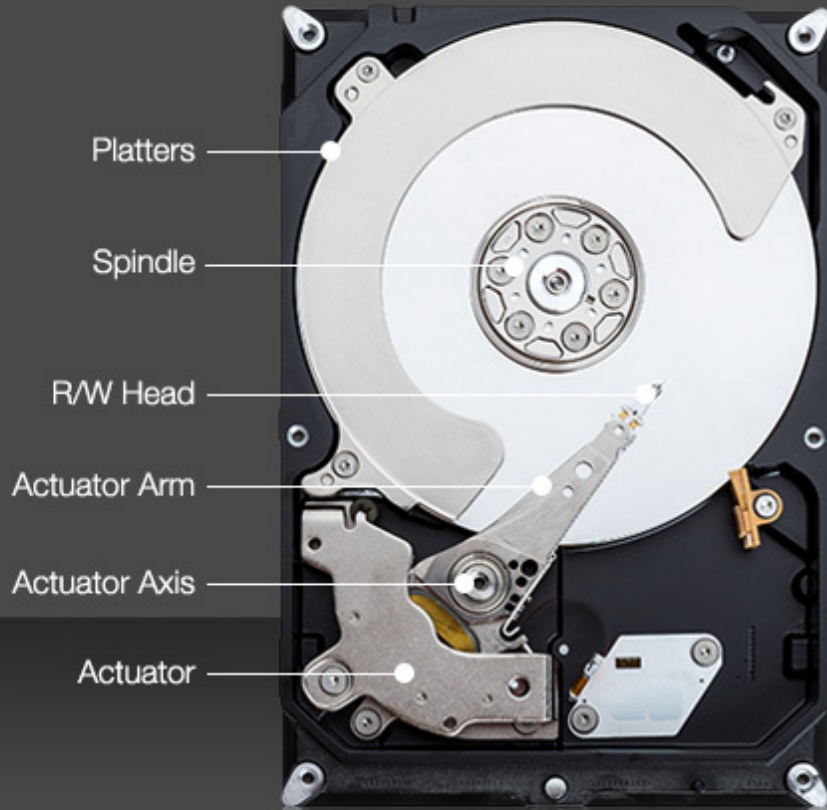
# Disks vs. Flash: Performance

- Throughput
  - Disk: ~130MB/s (sequential)
  - Flash: ~400MB/s
- Latency
  - Disk: ~10ms (one op)
  - Flash:
    - Read: 10-50us
    - Program: 200-500us
    - Erase: 2ms

Types of write, more later...

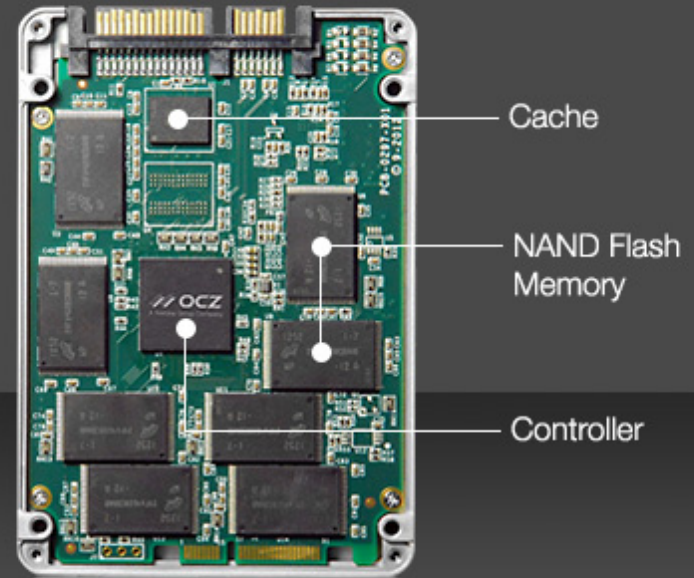
# Disks vs. Flash: Internal

**HDD**  
3.5"



Shock resistant up to 350g/2ms

**SSD**  
2.5"



Shock resistant up to 1500g/0.5ms



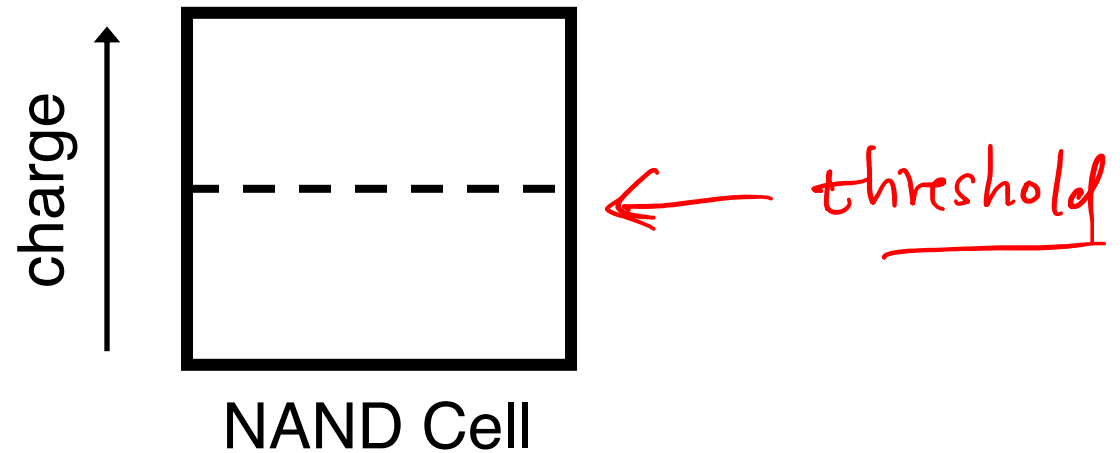
# Disks vs. Flash: Summary

	SSD	HDD
<b>Price</b>	\$0.25-\$0.27 per GB average	\$0.2-\$0.03 per GB average
<b>Lifespan</b>	30-80% test developed bad block in their lifetime	3.5% developed bad sectors comparatively
<b>Ideal for</b>	High performance processing Residing in APA or Tier 0/1 media in hybrid arrays	High capacity nearline tiers Long-term retained data
<b>Read/write speeds</b>	200 MB/s to 2500 MB/s	up to 200 MB/s
<b>Benefits</b>	Higher performance for faster read/write operations and fast load times	Less expensive Mature technology and massive installed user base
<b>Drawbacks</b>	May not be as durable/reliable as HDDs Not good for long-term archival data	Mechanical components take longer to read-write than SSDs

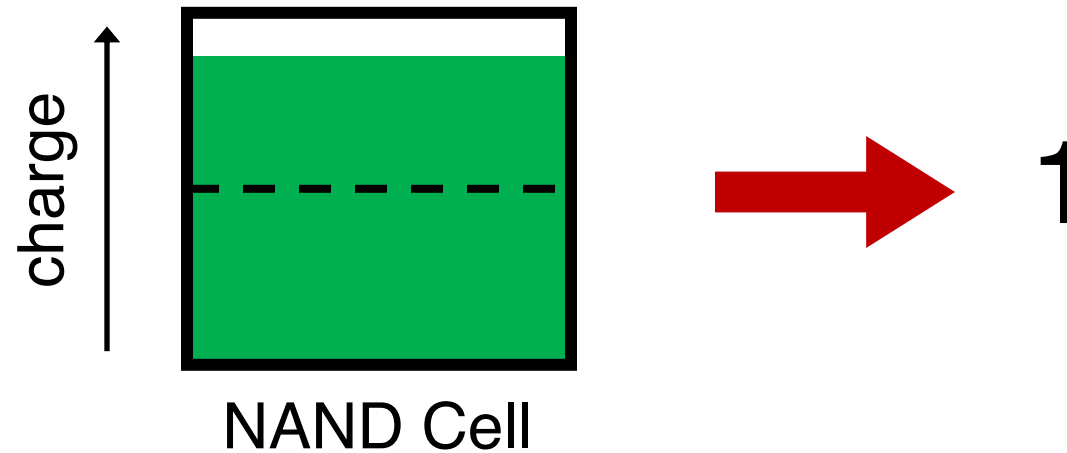
\* <https://www.enterprisestorageforum.com/storage-hardware/ssd-vs-hdd.html>

# Flash Architecture

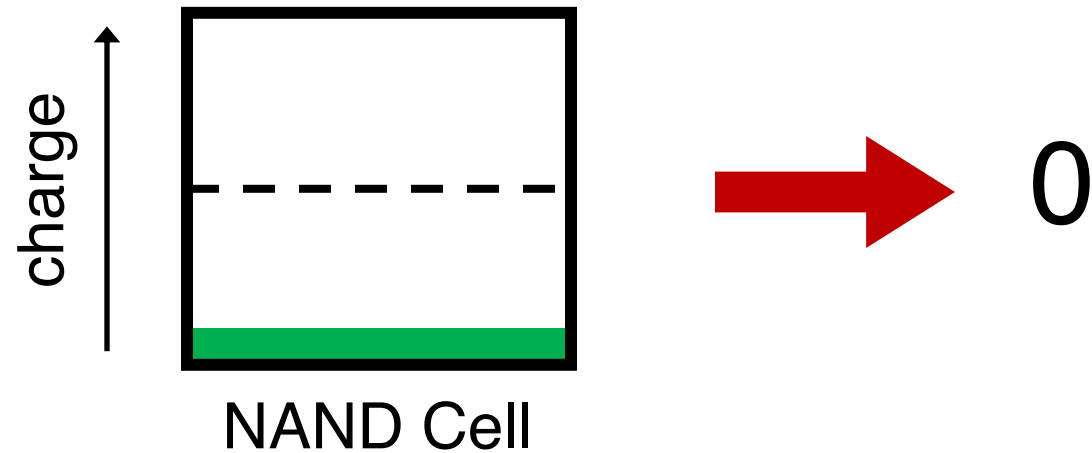
# SLC: Single-Level Cell



# SLC: Single-Level Cell



# SLC: Single-Level Cell

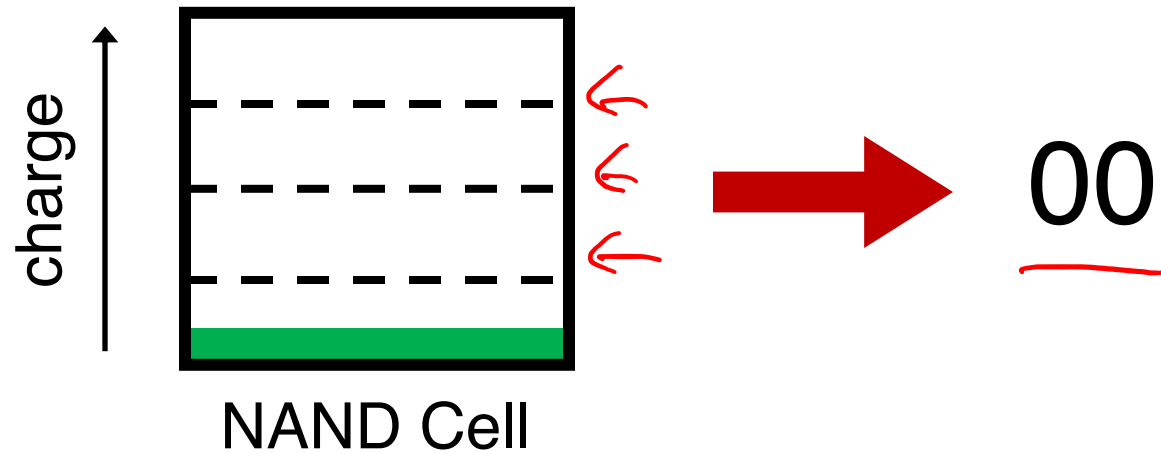


# MLC: Multi-Level Cell

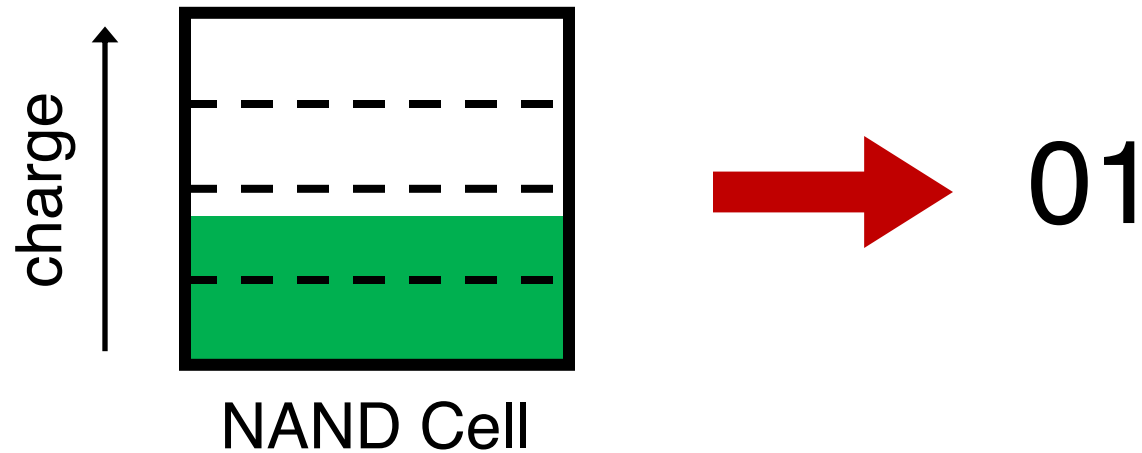
→ two bits

4 states.

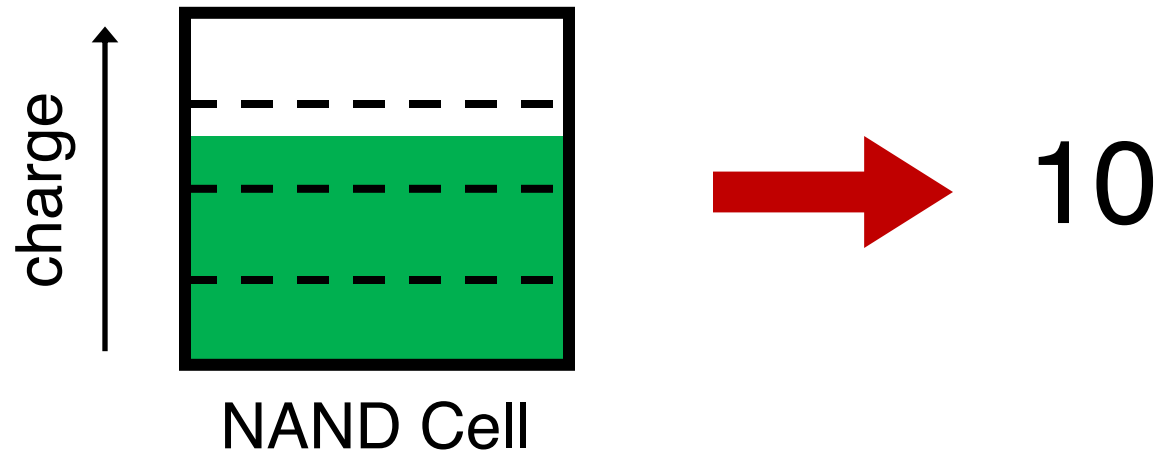
00
01
10
11



# MLC: Multi-Level Cell

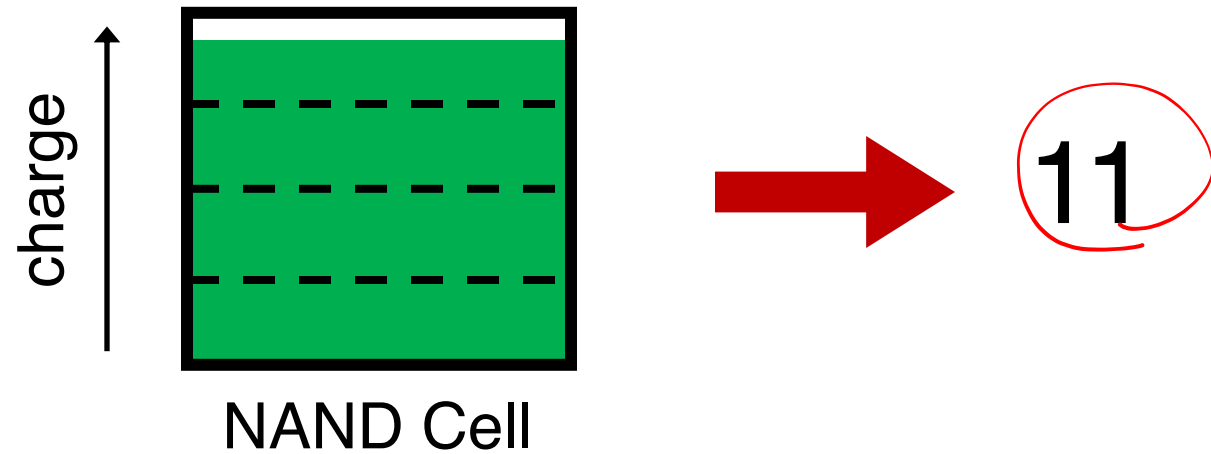


# MLC: Multi-Level Cell

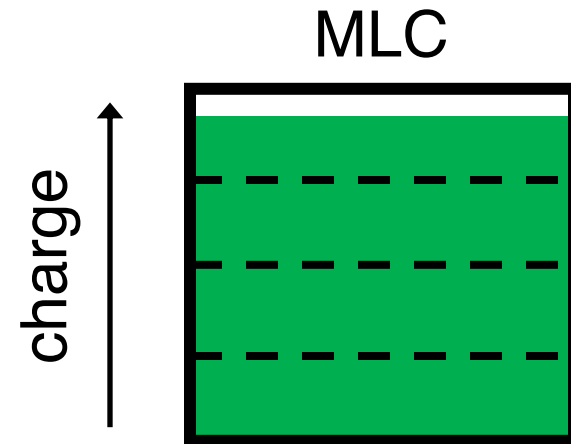
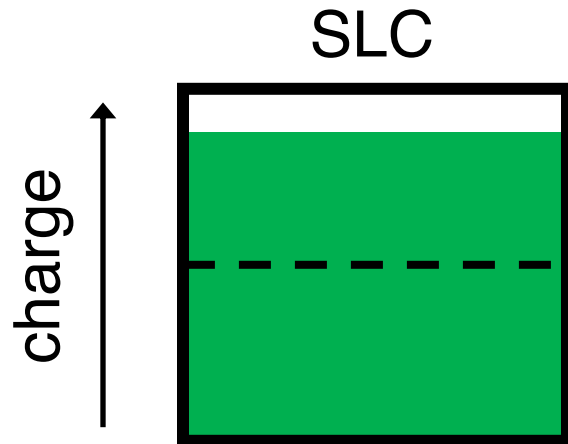




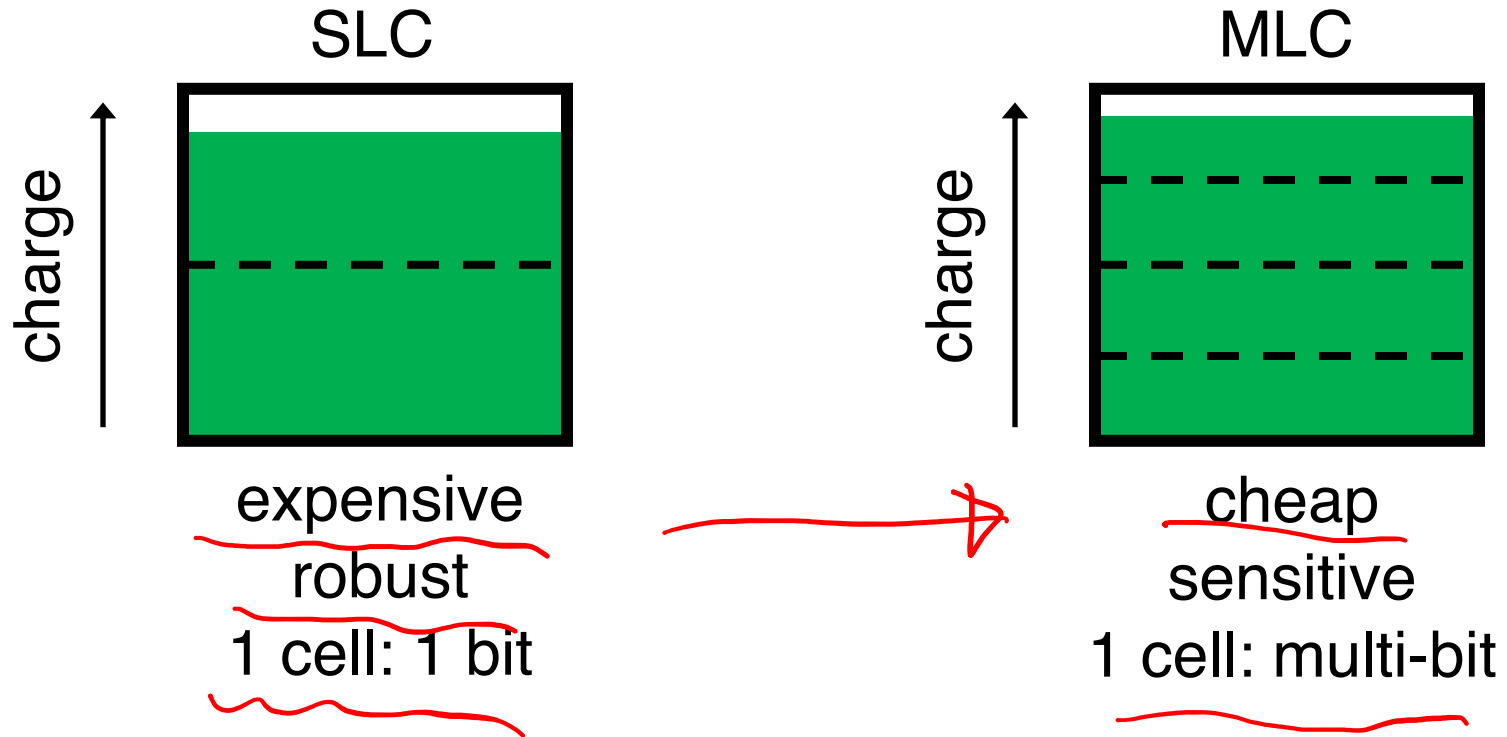
# MLC: Multi-Level Cell



# Single- vs. Multi-Level Cell



# Single- vs. Multi-Level Cell



# Wearout

no overwrite.

erase.

flash-block



overwrite.

- Problem: flash cells wear out after being erased too many times

endurance.

- MLC: ~10K times
- SLC: ~100K times
- Usage strategy: ???

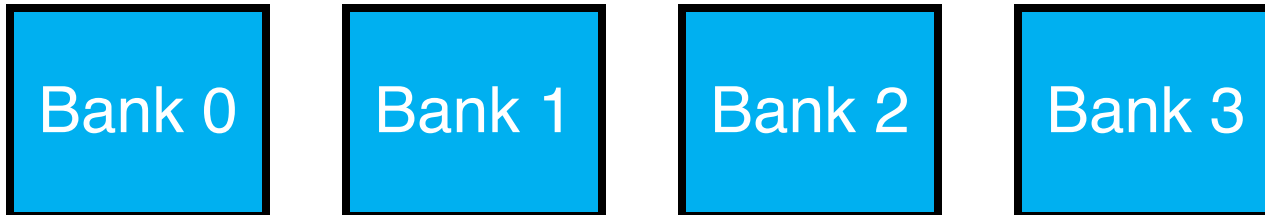
How  
maximize?

# Wearout

- Problem: flash cells wear out after being erased too many times
  - MLC: ~10K times
  - SLC: ~100K times
  - Usage strategy: wear leveling
    - Prevents some cells from being wornout while others still fresh
- write balancing*

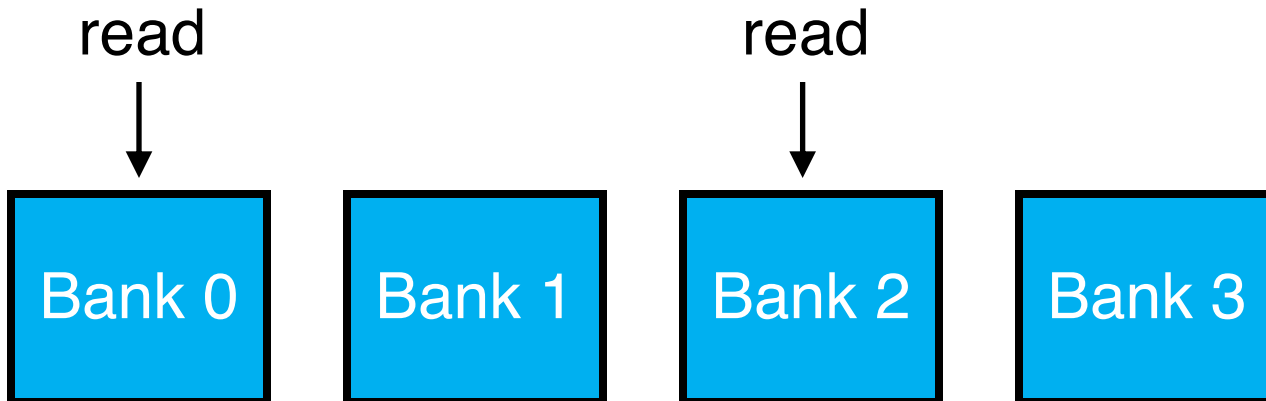
# Banks

- Flash devices are divided into banks (aka. planes)
- Banks can be accessed in parallel



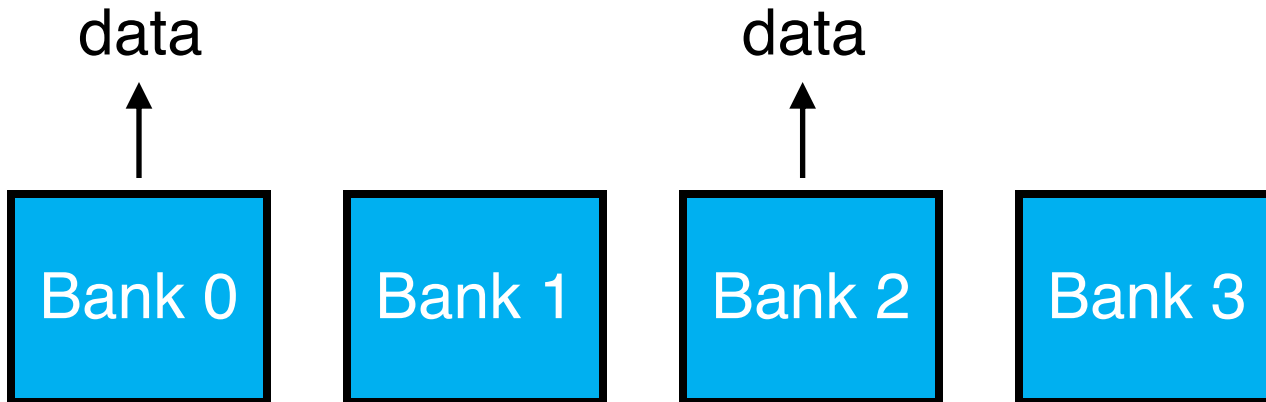
# Banks

- Flash devices are divided into banks (aka. planes)
- Banks can be accessed in parallel



# Banks

- Flash devices are divided into banks (aka. planes)
- Banks can be accessed in parallel







# Flash Writes

- Writing 0's → data.
  - Fast, fine-grained
- Writing 1's → fresh state.
  - Slow, coarse-grained

# Flash Writes

- Writing 0's
  - Fast, fine-grained
  - called “program”
- Writing 1's
  - Slow, coarse-grained
  - called “erase”

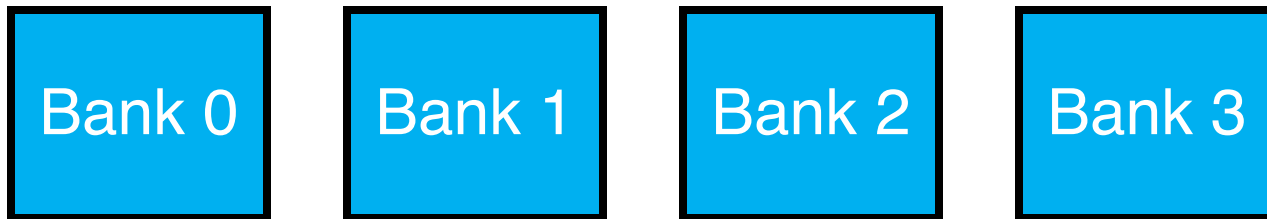
# Flash Writes

- Writing 0's
  - Fast, fine-grained [page-level]
  - called “**program**” 
- Writing 1's
  - Slow, coarse-grained [block-level]
  - called “**erase**” 

# Flash Writes

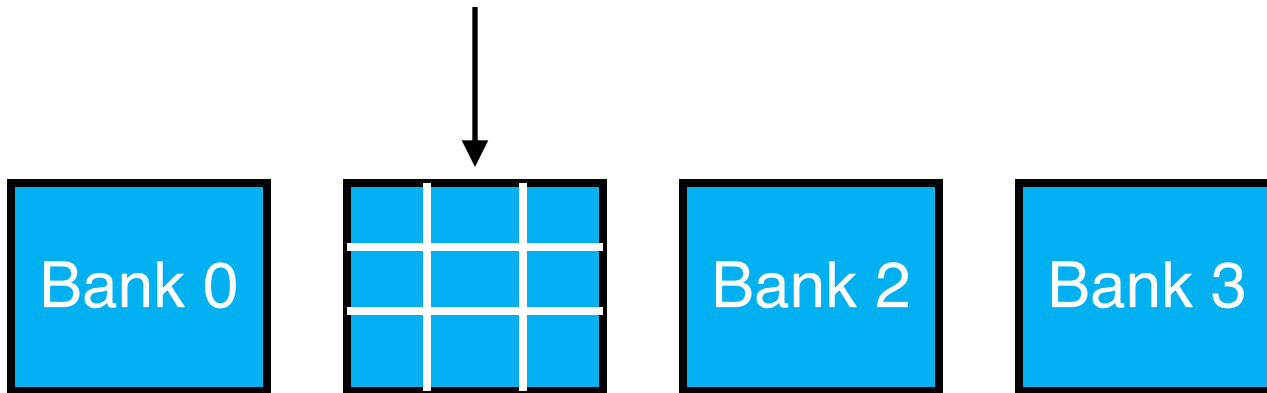
- Writing 0's
  - Fast, fine-grained [page-level]
  - called “**program**”
- Writing 1's
  - Slow, coarse-grained [block-level]
  - called “**erase**”
- Flash can only “write” (program) into **clean** pages
  - “**clean**”: pages containing all 1's (pages that have been erased)
  - Flash does not support in-place overwrite!

# Banks and Blocks

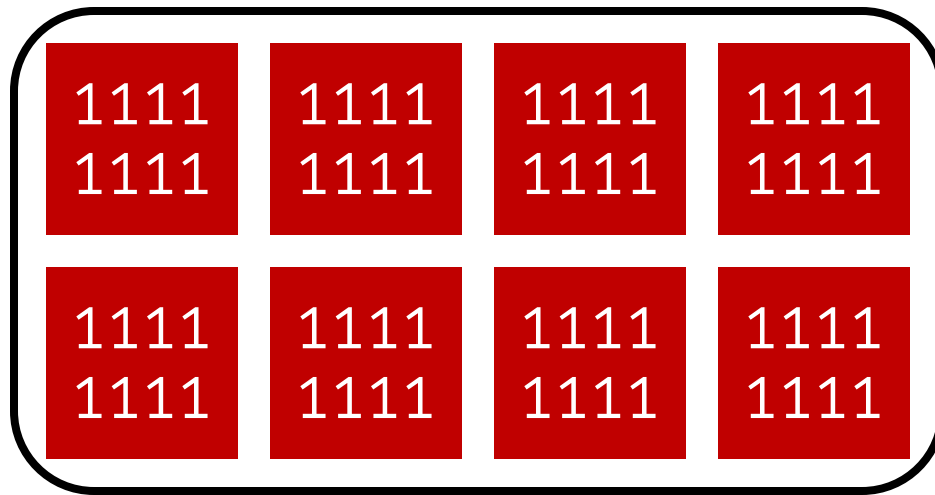


# Banks and Blocks

Each bank contains  
many “blocks”



# Block and Pages

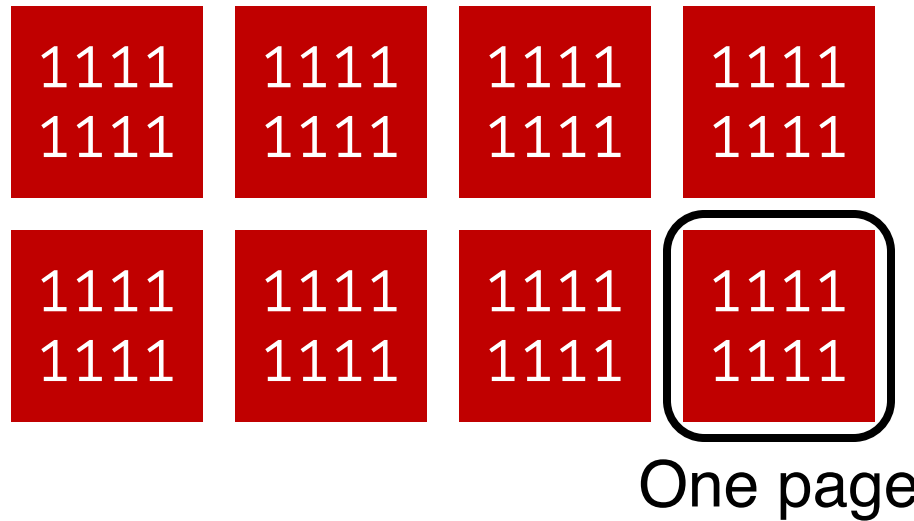


One block

# Block and Pages

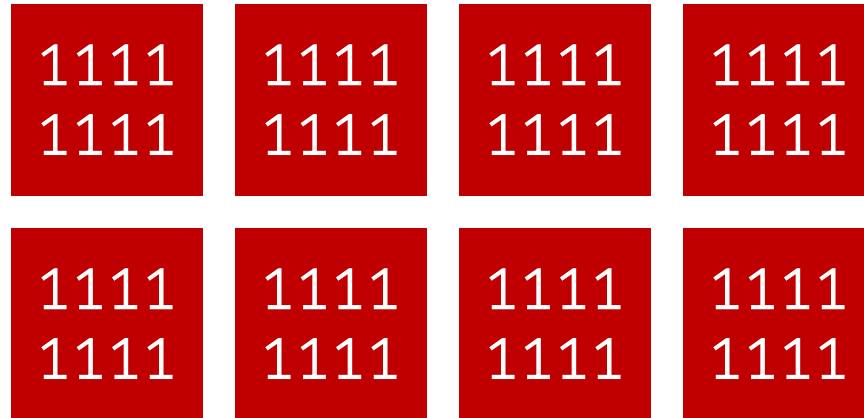
8 pages

basic unit  
of 40



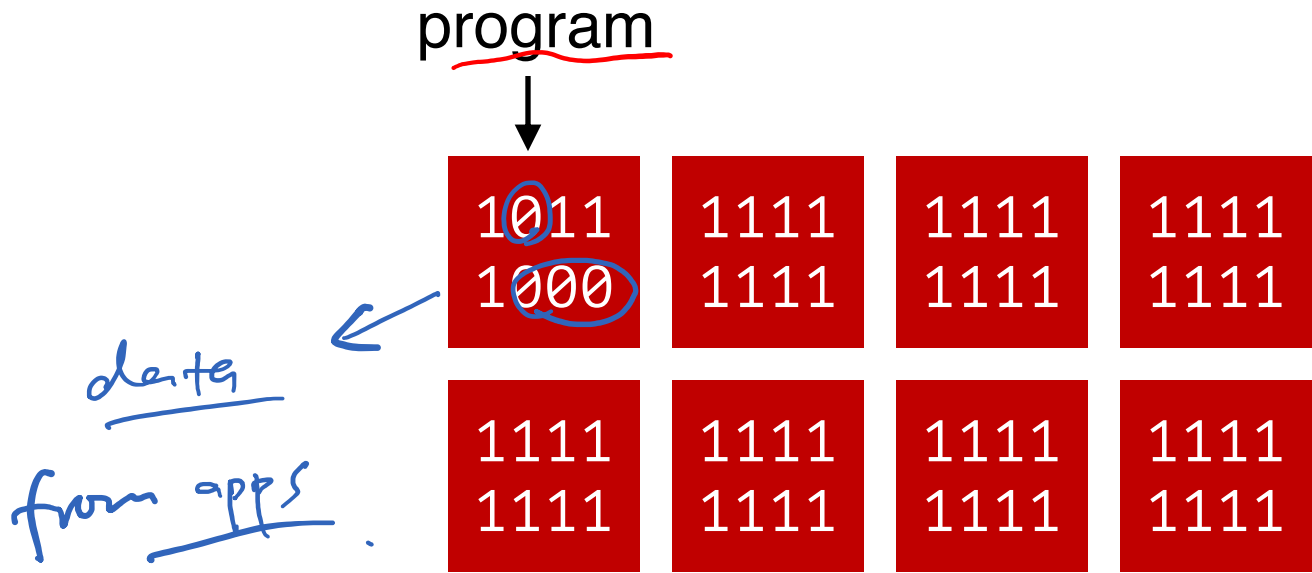


# Block and Pages

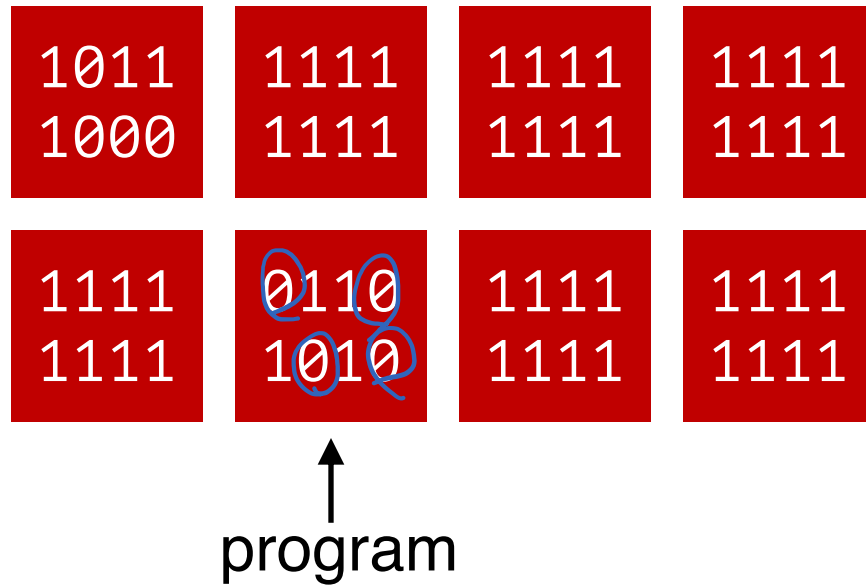


All pages are clean  
("programmable")

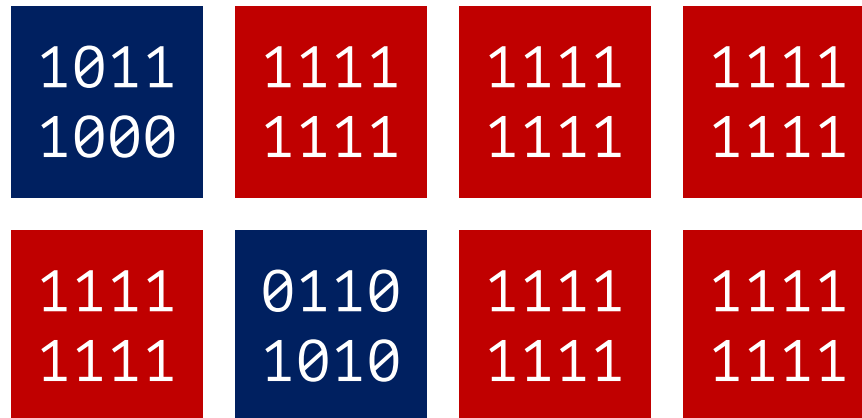
# Block



# Block



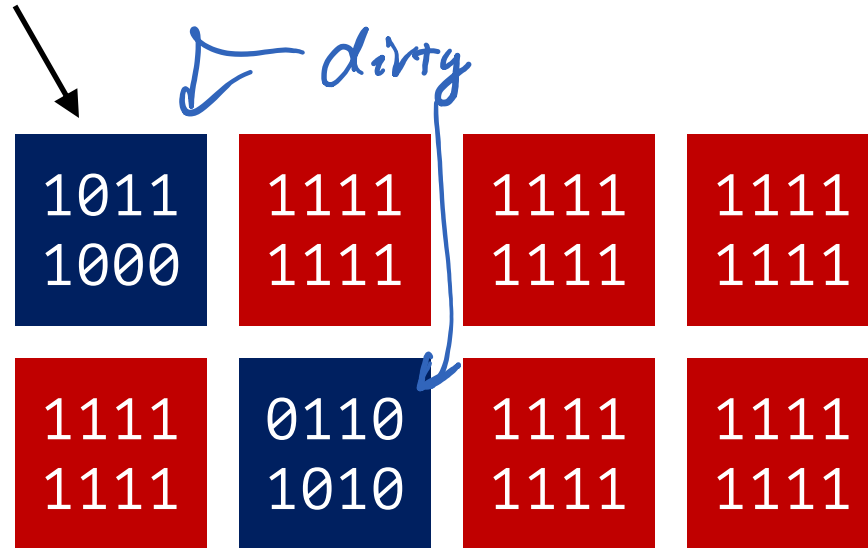
# Block



Two pages hold data  
**(cannot be overwritten)**

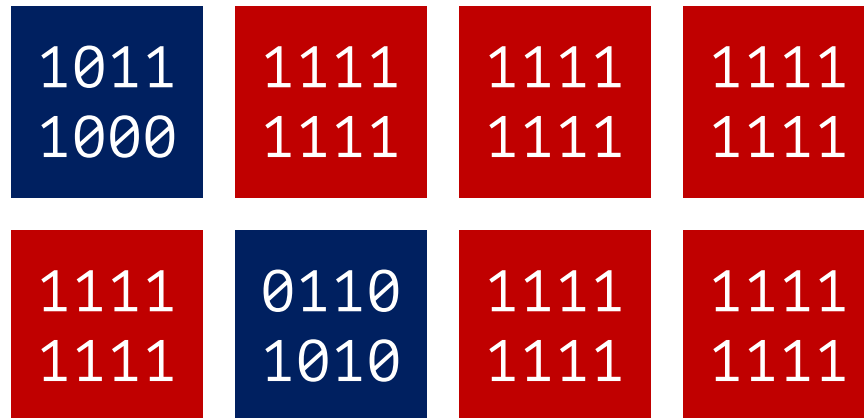
# Block

still want to write data into this page???



Two pages hold data  
**(cannot be overwritten)**

# Block

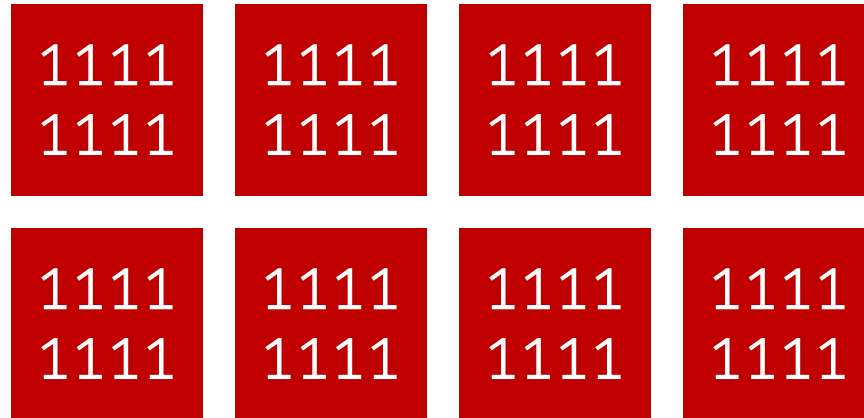


erase

the whole block!!

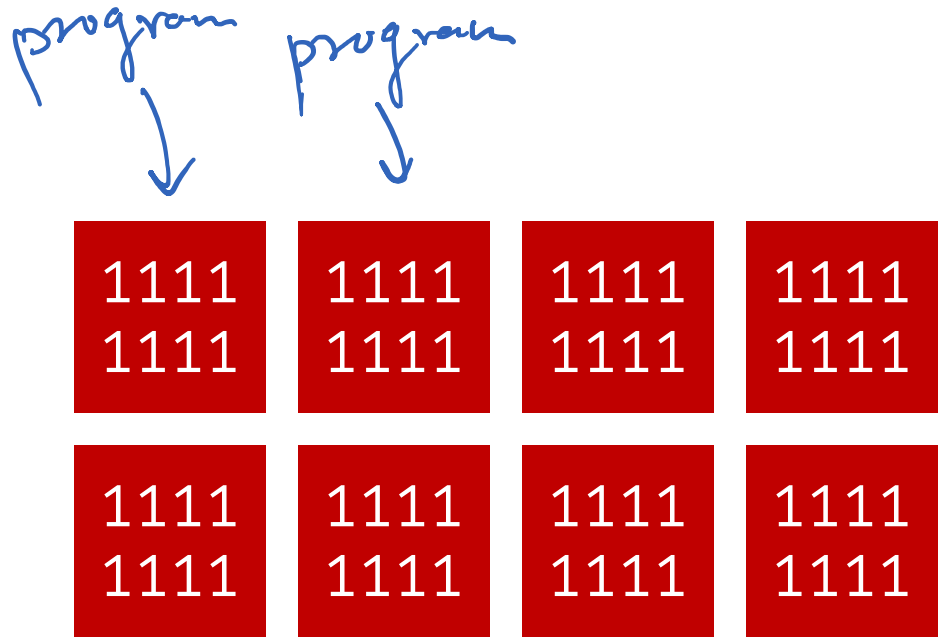
# Block

Flip 0's into 1's.



erase  
(the whole block)

# Block



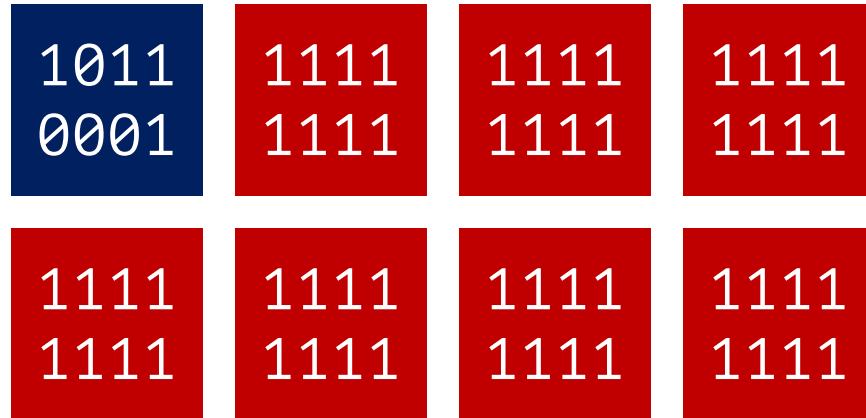
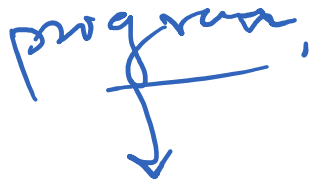
After erase, again, **free state**  
(can write new data in any page)

program



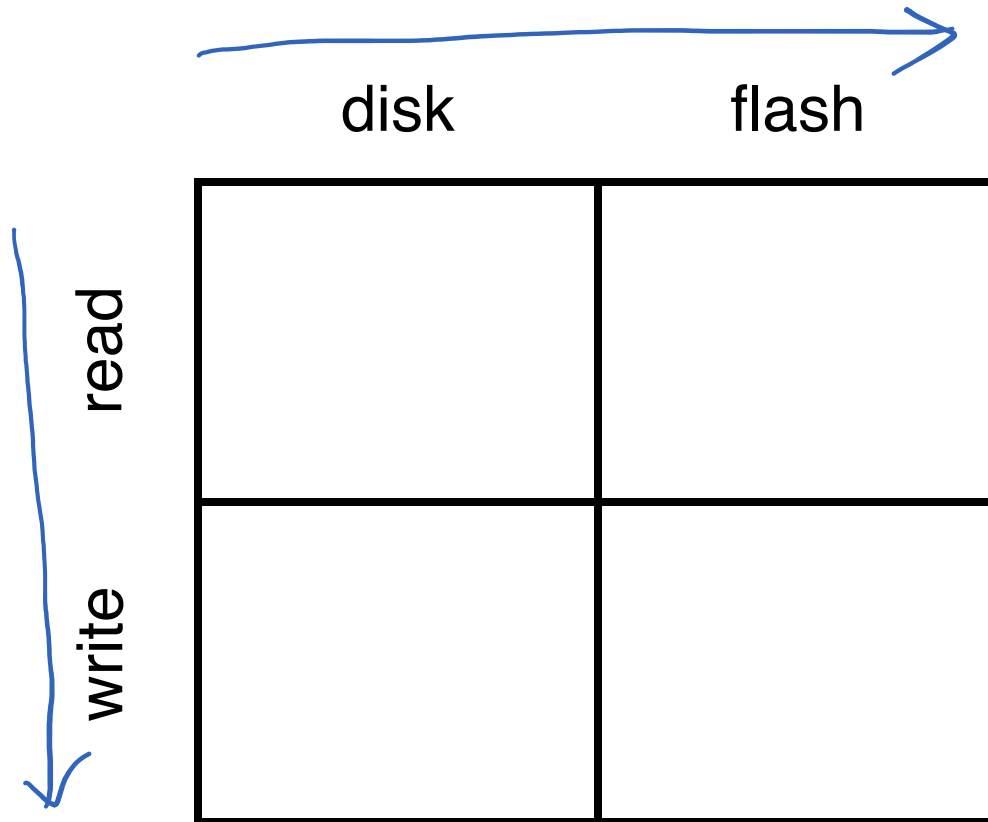
# Block

*program*

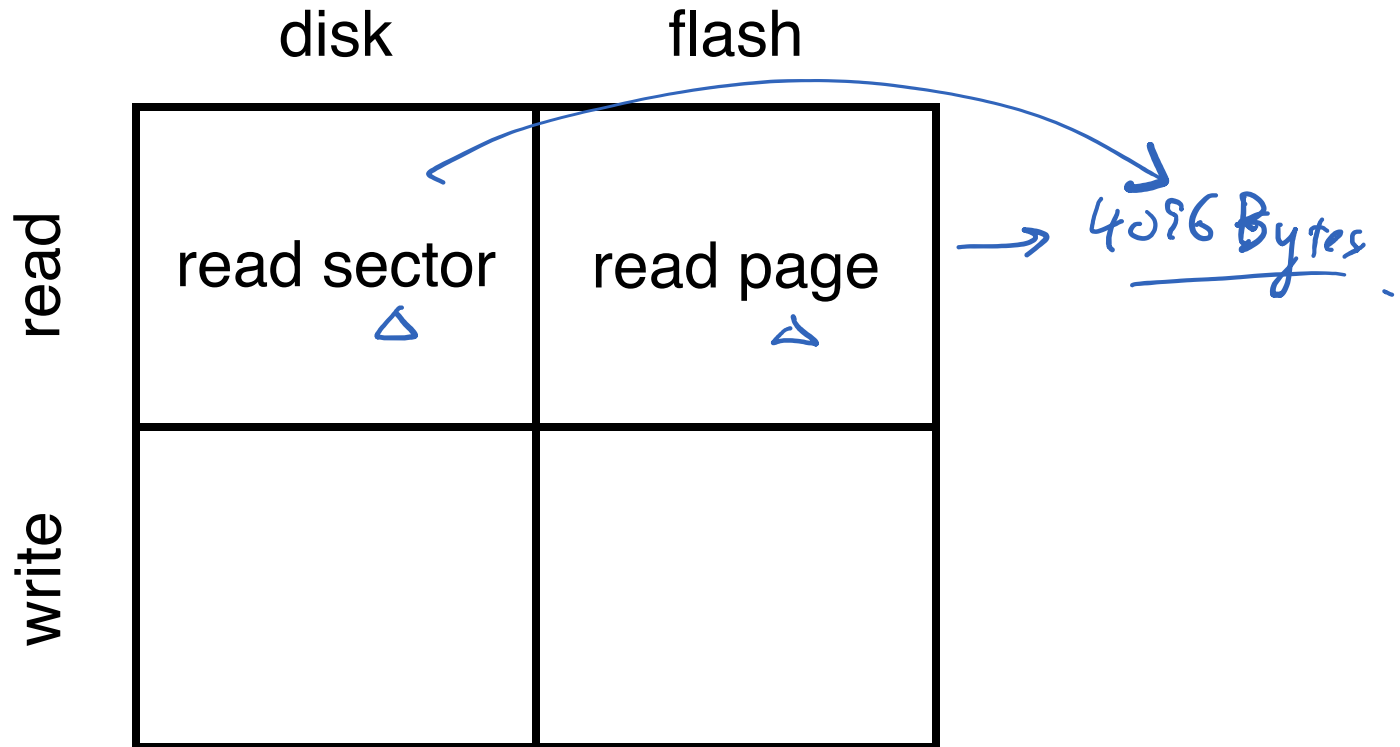


This blue page holds data

# Flash vs. Disks: APIs



# Flash vs. Disks: APIs



# Flash vs. Disks: APIs

	disk	flash
read	read sector	read page
write	write sector	<b>program</b> page → (0's) <b>erase</b> block ⤵ (1's)

program ops  
1 → 0  
erase op.  
0 → 1 ←

# Flash Architecture

- **Bank/plane**: 1024 to 4096 blocks
  - Banks accessed in parallel

- **Block**: 64 to 256 pages
  - Unit of erase

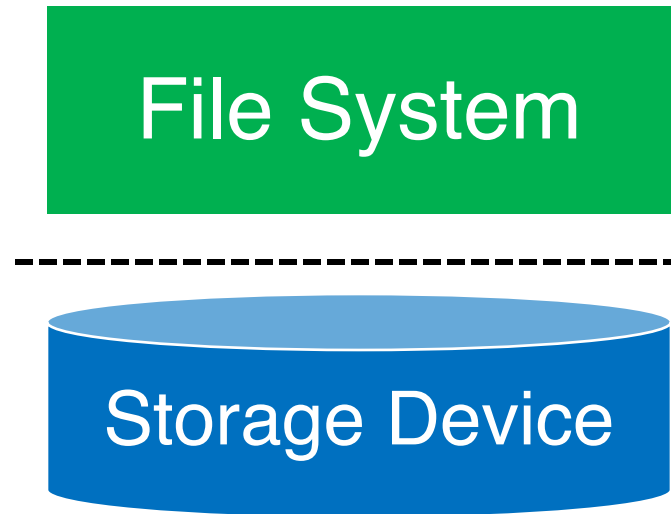
- **Page**: 2 to 8 KB
  - Unit of read and program

# Disks vs. Flash: Performance

- Throughput
  - Disk: ~130MB/s (sequential)
  - Flash: ~400MB/s
- Latency
  - Disk: ~10ms (one op)
  - Flash:
    - **Read**: 10-50us
    - **Program**: 200-500us
    - **Erase**: 2ms

# Working with File System

# Traditional File Systems

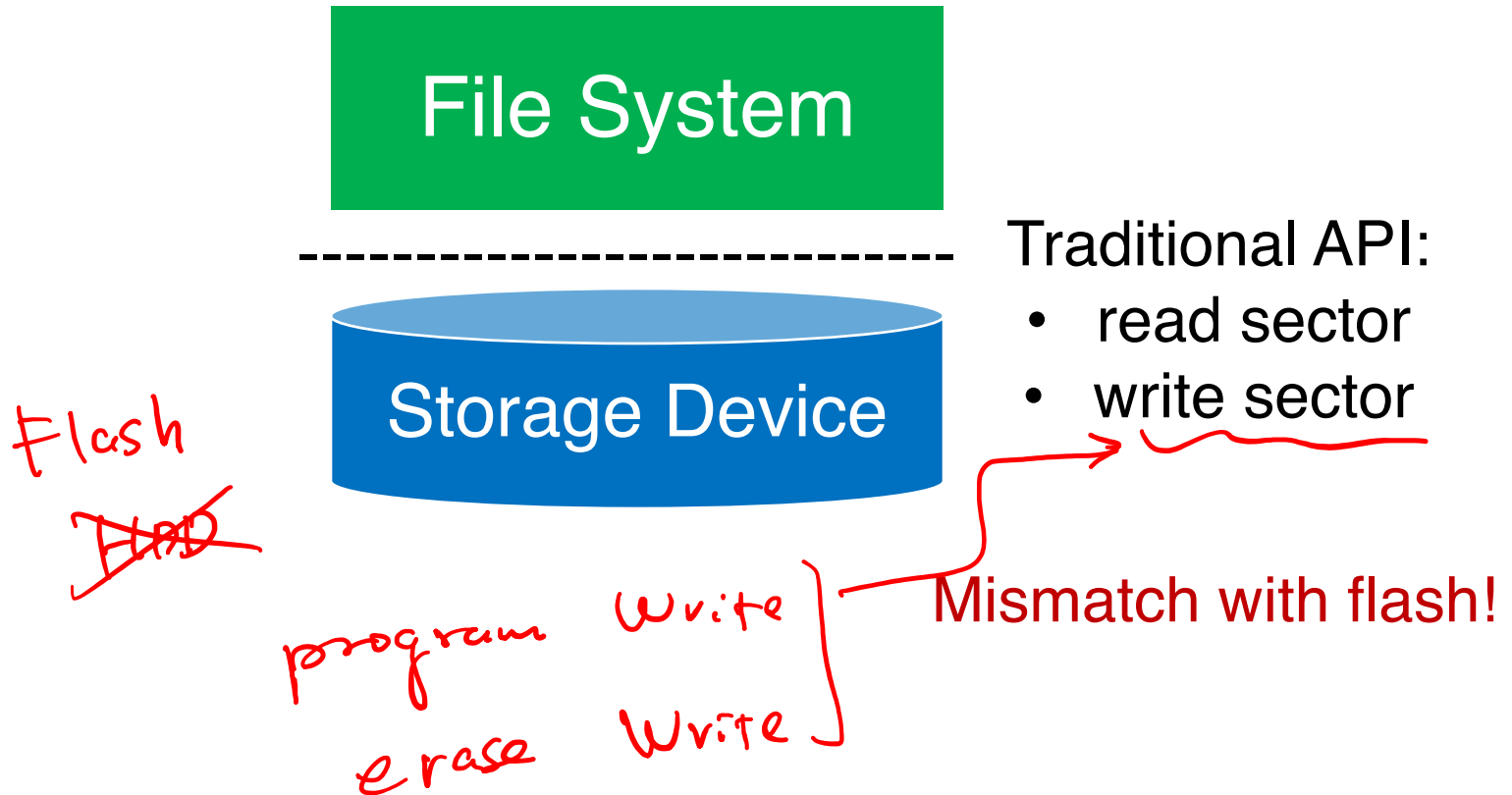


Traditional API:

- read sector
- write sector



# Traditional File Systems



# Traditional APIs wrapping around Flash APIs

read(addr):

*read flash-page.*

return flash\_read(addr)

△

~~\*~~  
↳ write(addr, data):

*HDD-friendly FS semantics.*

(mem) ① block\_copy = flash\_read(all pages of block)

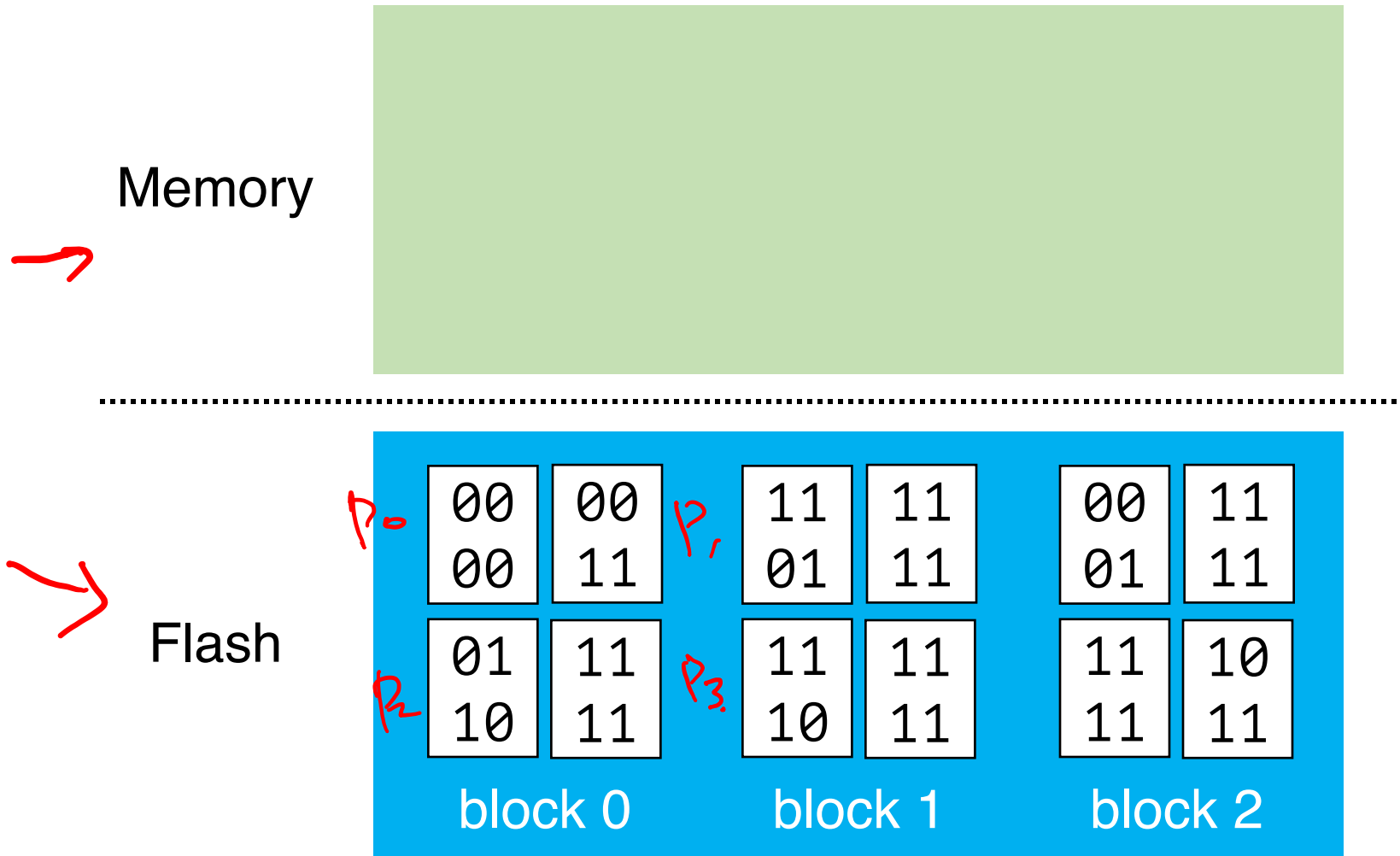
freely ② modify block\_copy with data → *in-mem data / state.*

③ flash\_erase(block of addr)

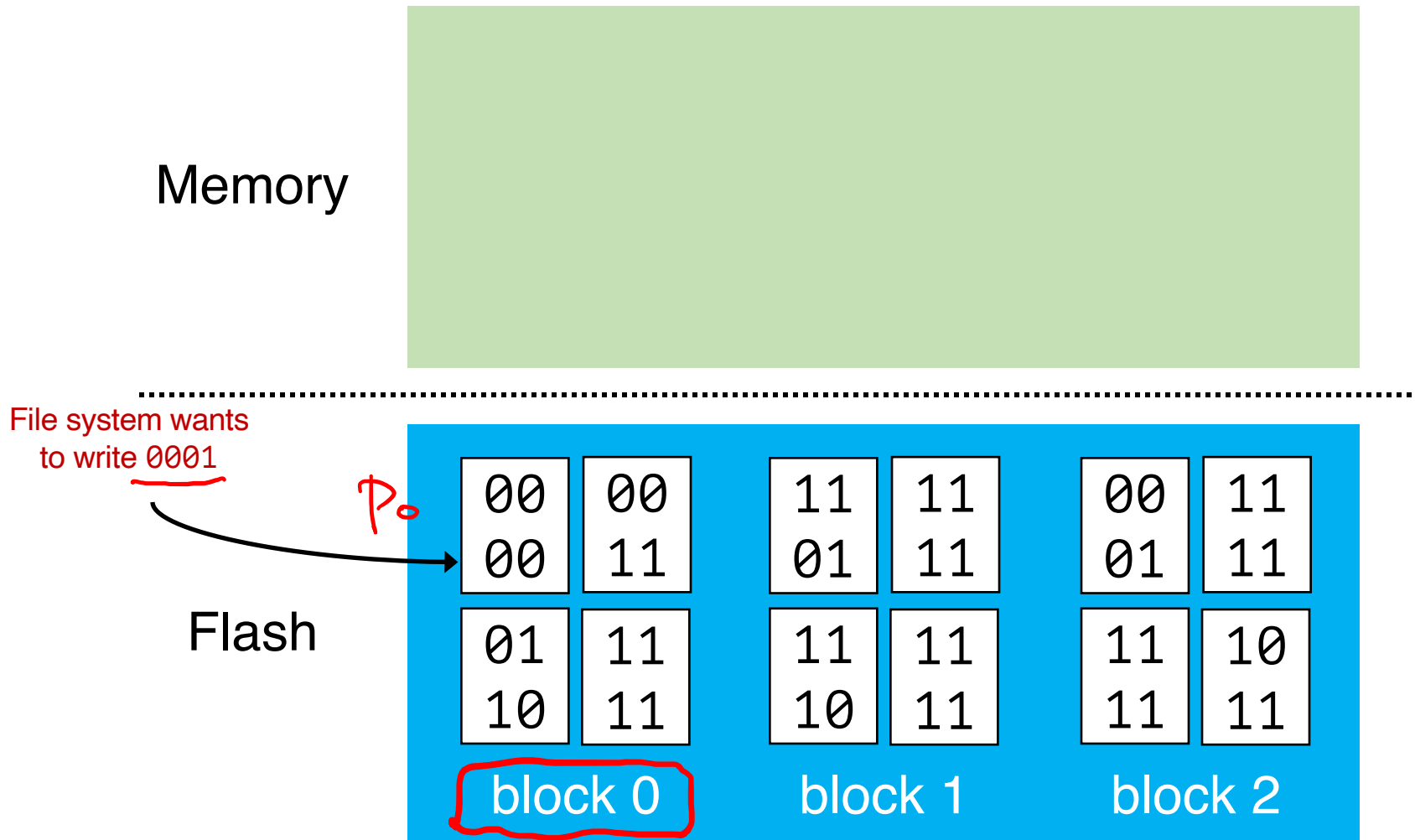
④ flash\_program(all pages of block\_copy)

finer-grained. △

# Awkward Flash Write

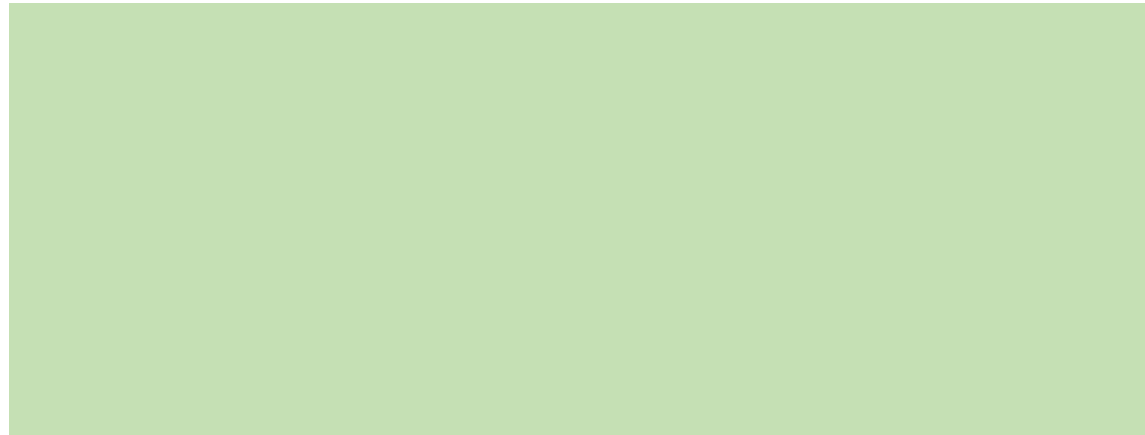


# Awkward Flash Write

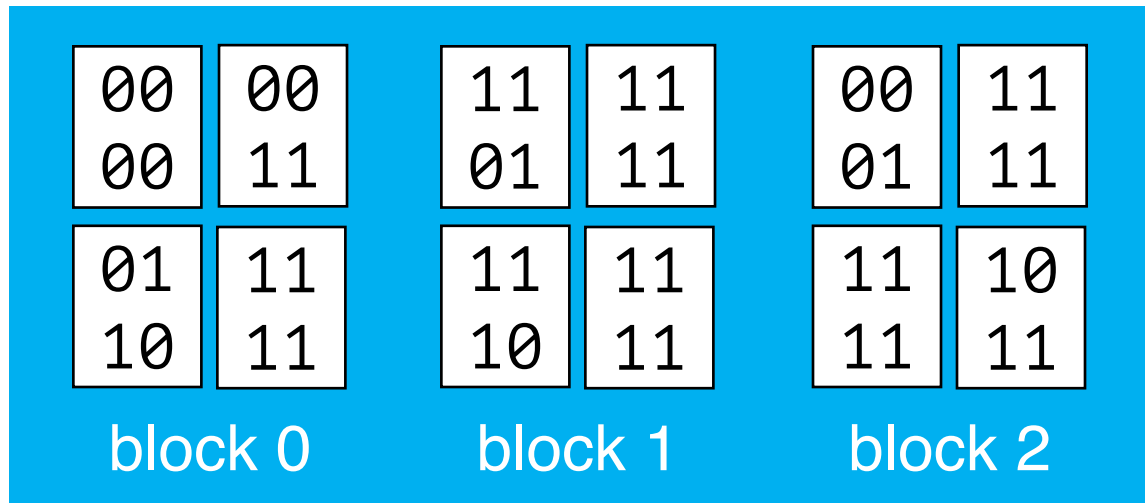


# Awkward Flash Write

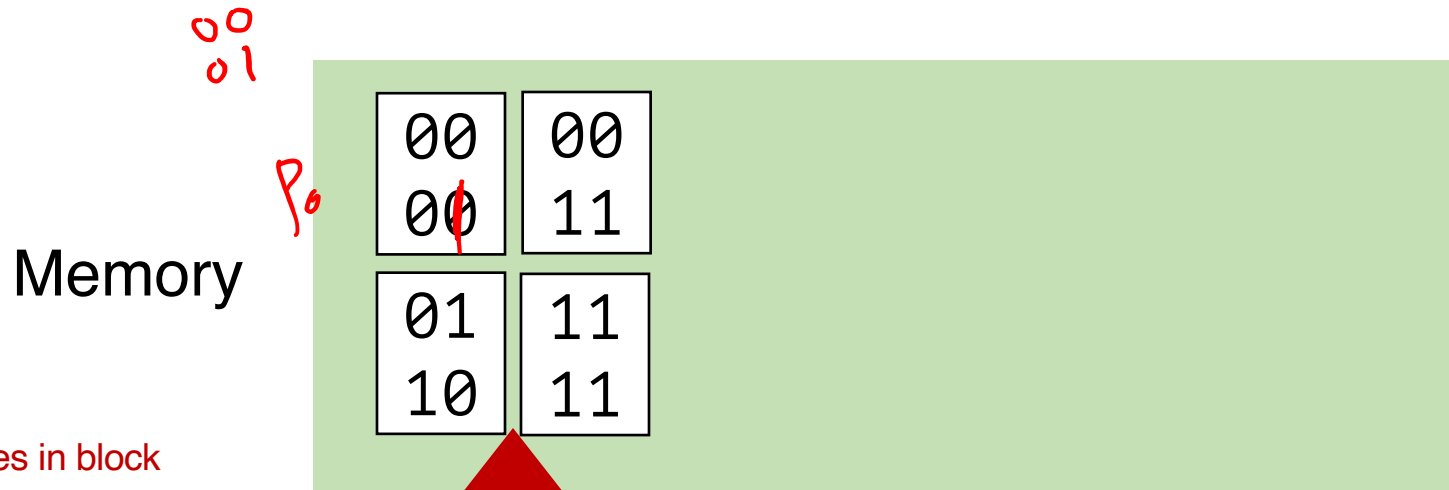
Memory



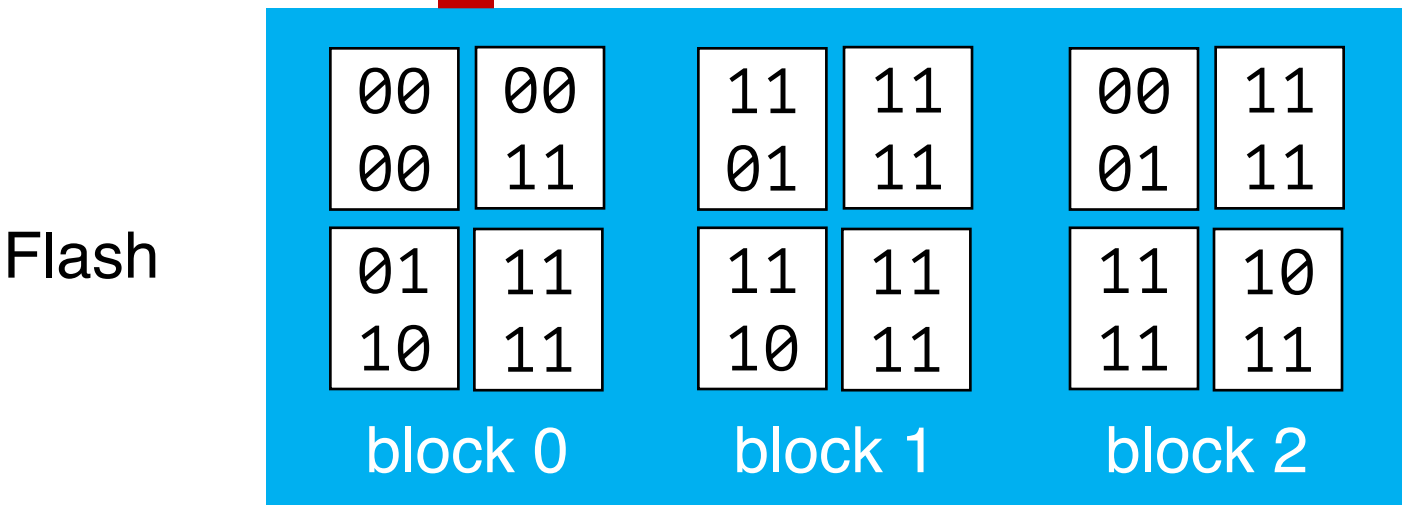
Flash



# Awkward Flash Write



Read all pages in block



# Awkward Flash Write

Memory

00	00
00	11
01	11
10	11

Flash

00	00	11	11	00	11
00	11	01	11	01	11
01	11	11	11	11	10
10	11	10	11	11	11
block 0		block 1		block 2	

# Awkward Flash Write

Modify target page  
in memory

Memory

00	00
01	11
01	11
10	11

Flash

00	00	11	11	00	11
00	11	01	11	01	11
01	11	11	11	11	10
10	11	10	11	11	11
block 0		block 1		block 2	



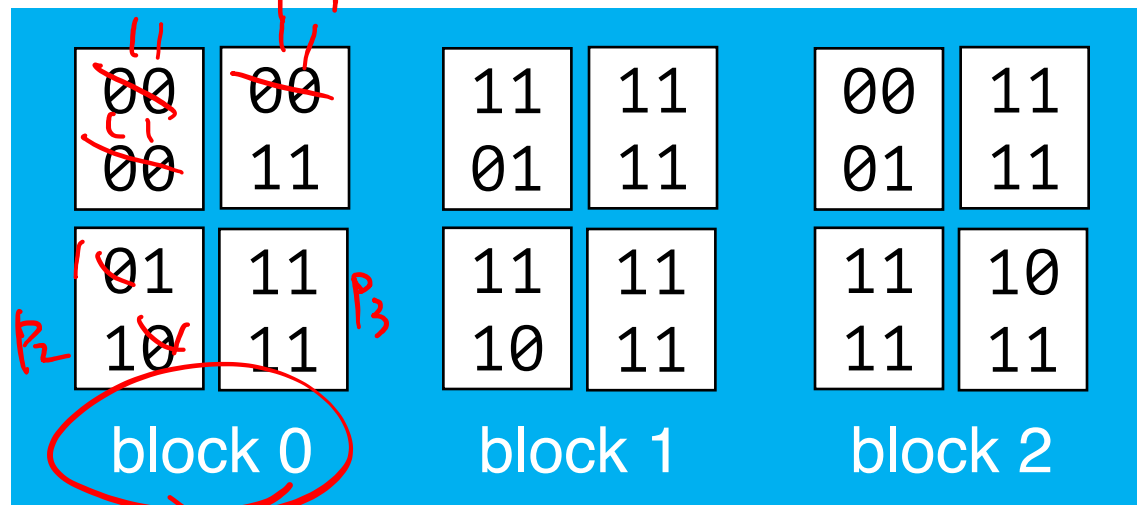
# Awkward Flash Write

Memory

00	00
01	11
01	11
10	11

~~no  
overwrite~~  
erase

Flash



# Awkward Flash Write

Memory

00	00
01	11
01	11
10	11

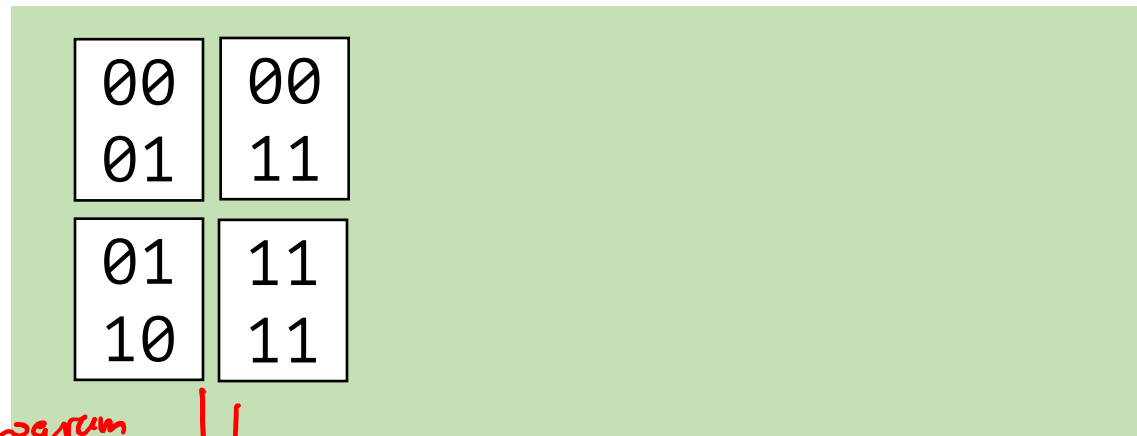
Erase whole block

Flash

<b>11</b>	<b>11</b>	11	11	00	11
<b>11</b>	<b>11</b>	01	11	01	11
<b>11</b>	<b>11</b>	11	11	11	10
<b>11</b>	<b>11</b>	10	11	11	11
block 0	block 1	block 2			

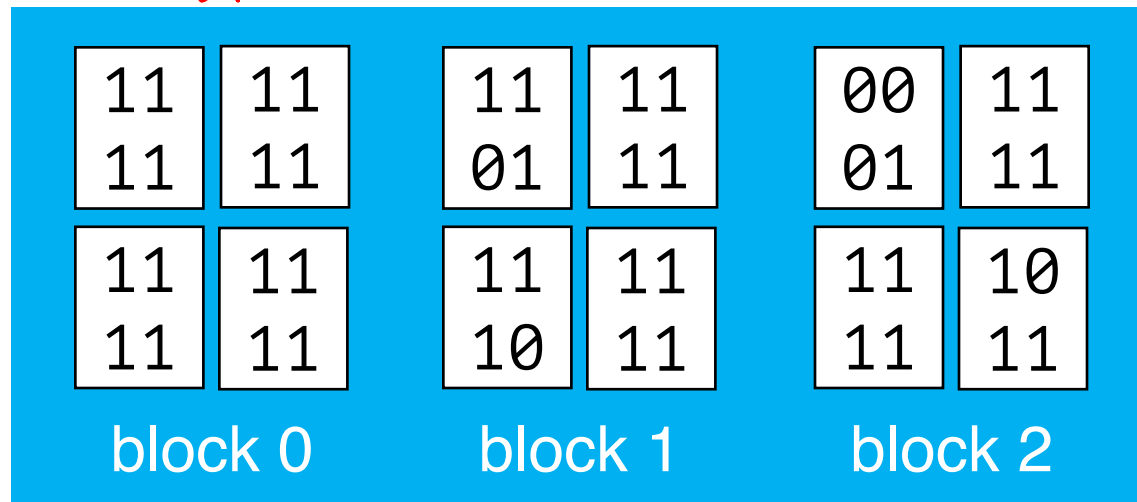
# Awkward Flash Write

Memory

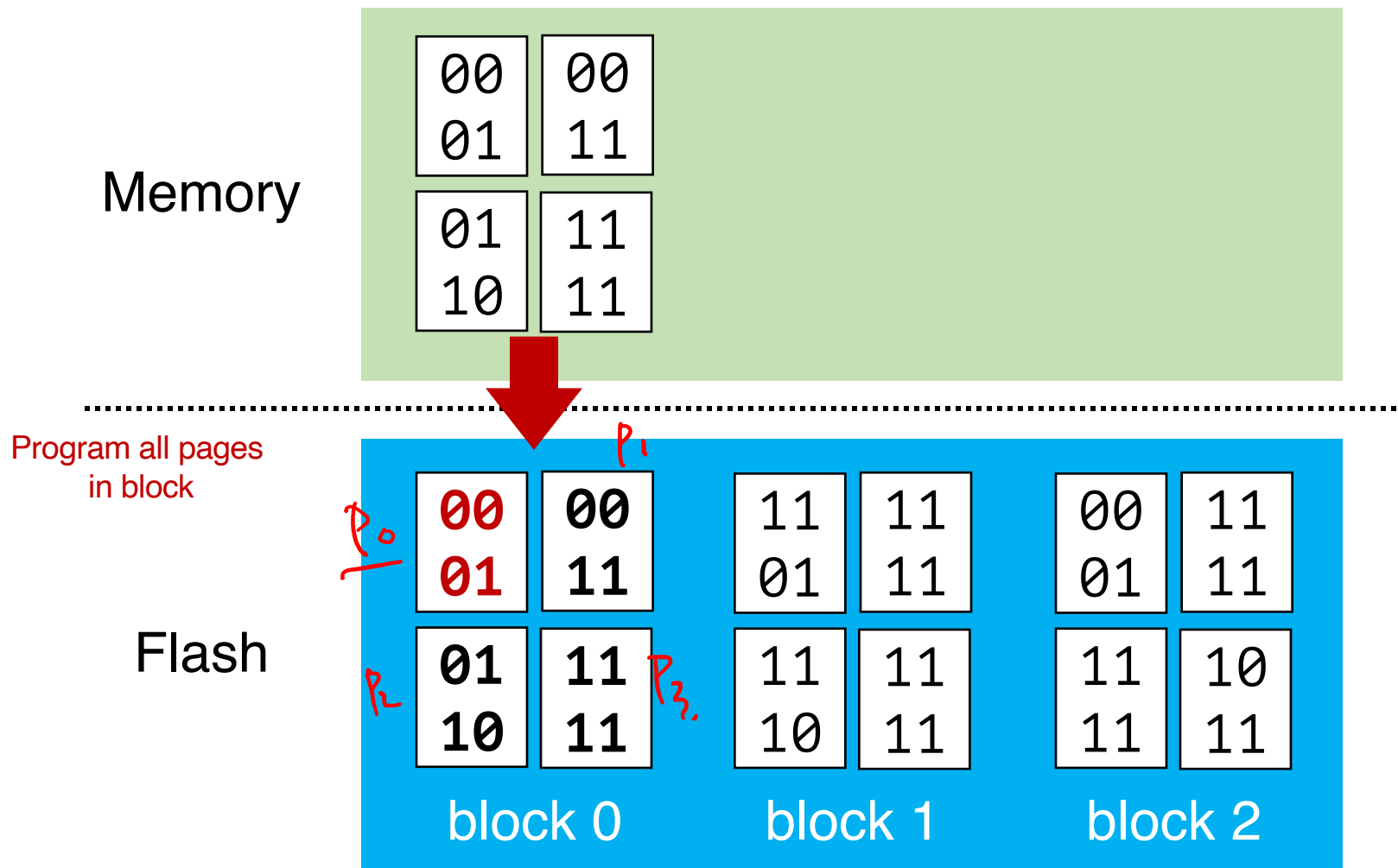


*Program* ↓ ↓

Flash

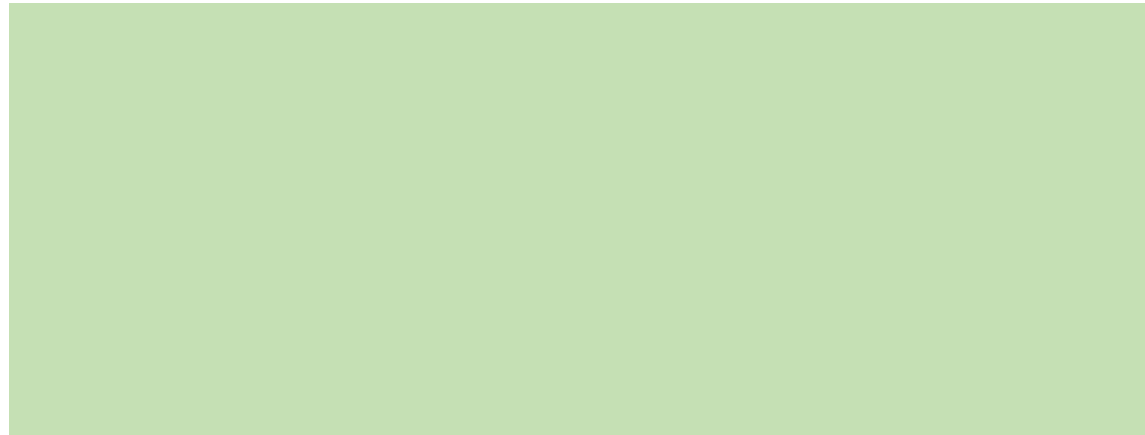


# Awkward Flash Write

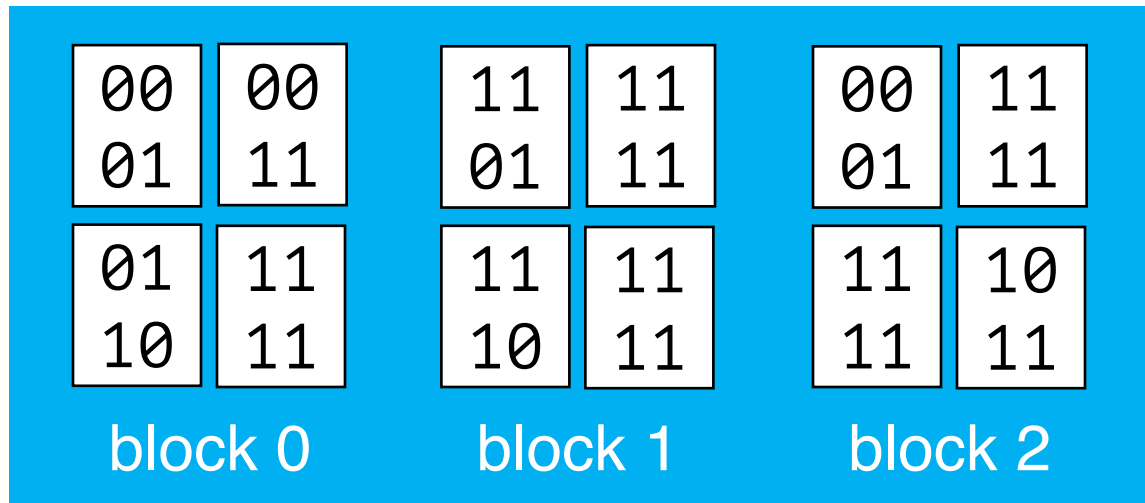


# Awkward Flash Write

Memory



Flash



# Issue: Write Amplification

factor  $\rightarrow \frac{256 \text{ KB}}{4 \text{ KB}} = \underline{64 \times}$

- Random writes are expensive for flash!  
△
- Writing one 4KB page may cause:
  - read, erase, and program of the whole 256KB block

Random I/Os → Sequential Appends.

By adding one more abstraction layer!

# Flash Translation Layer

# Flash Translation Layer (FTL)

Page table.  
VPA → PPA.

- Add an address translation layer between upper-level file system and lower-level flash

mem  
Virtualization.

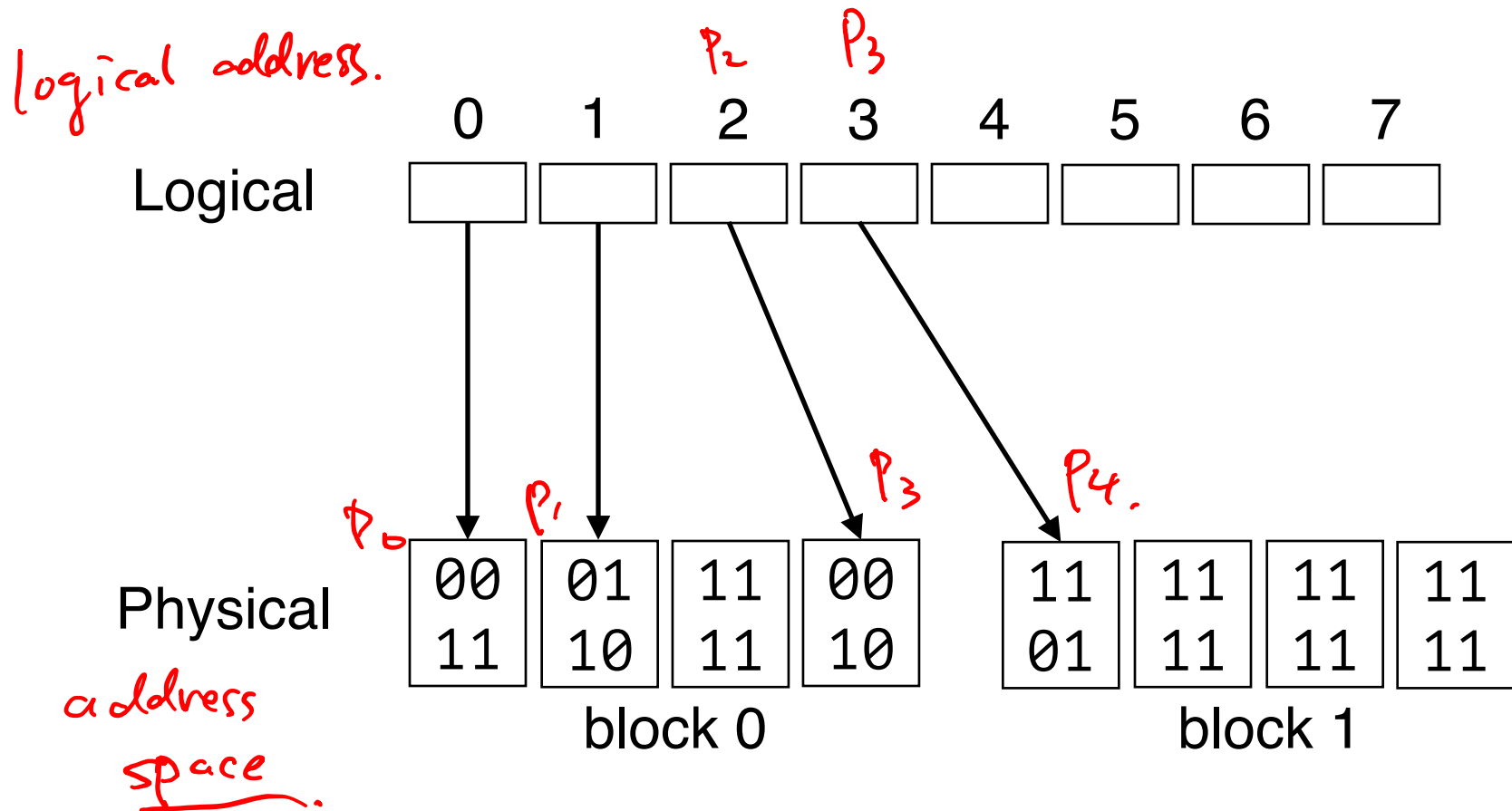
- Translate logical device addresses to physical addresses

- Goal →
- Convert in-place write into append-write
  - Essentially, a virtualization & optimization layer

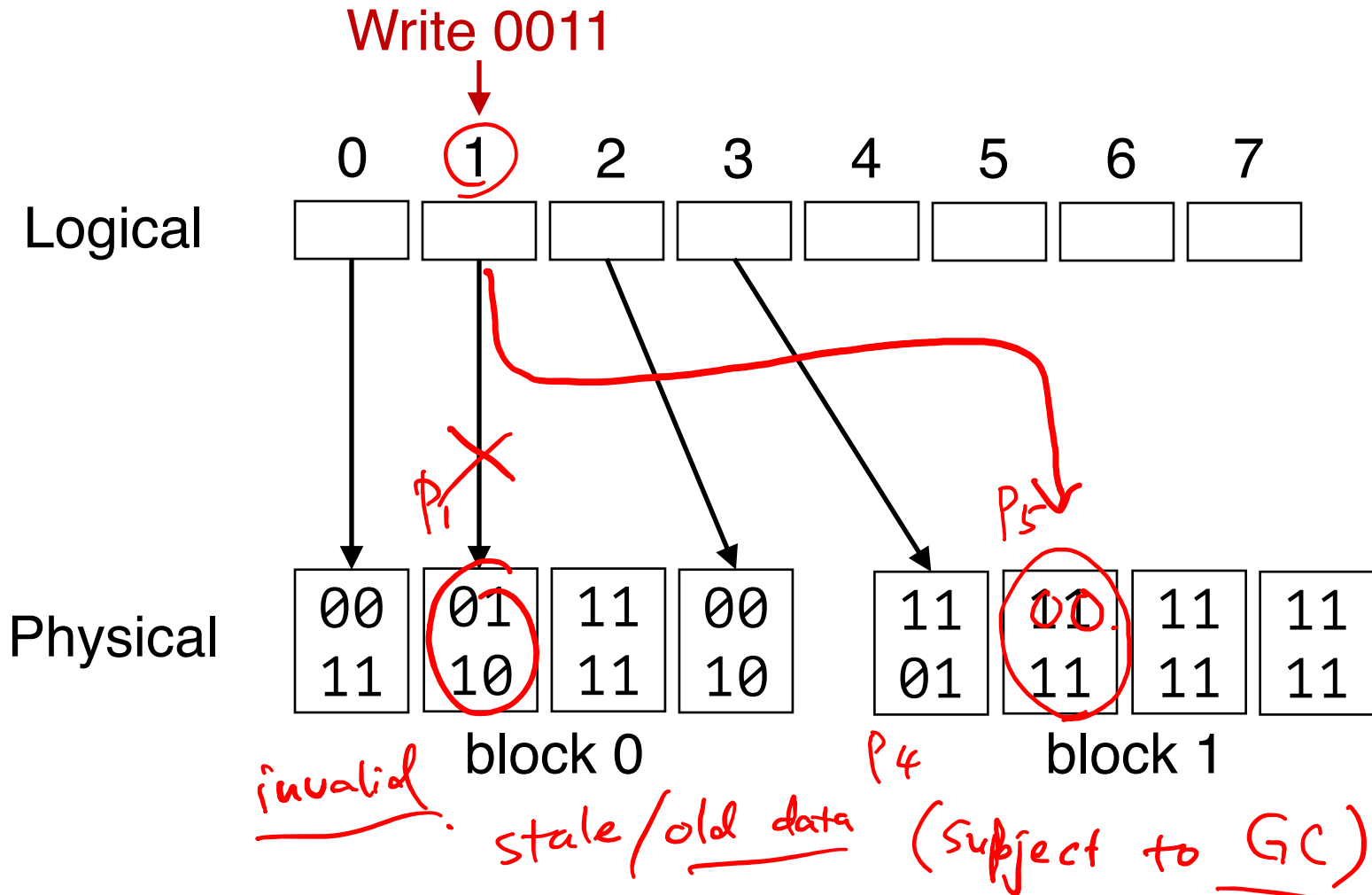
FTL → Flash  
space  
Virtualization.



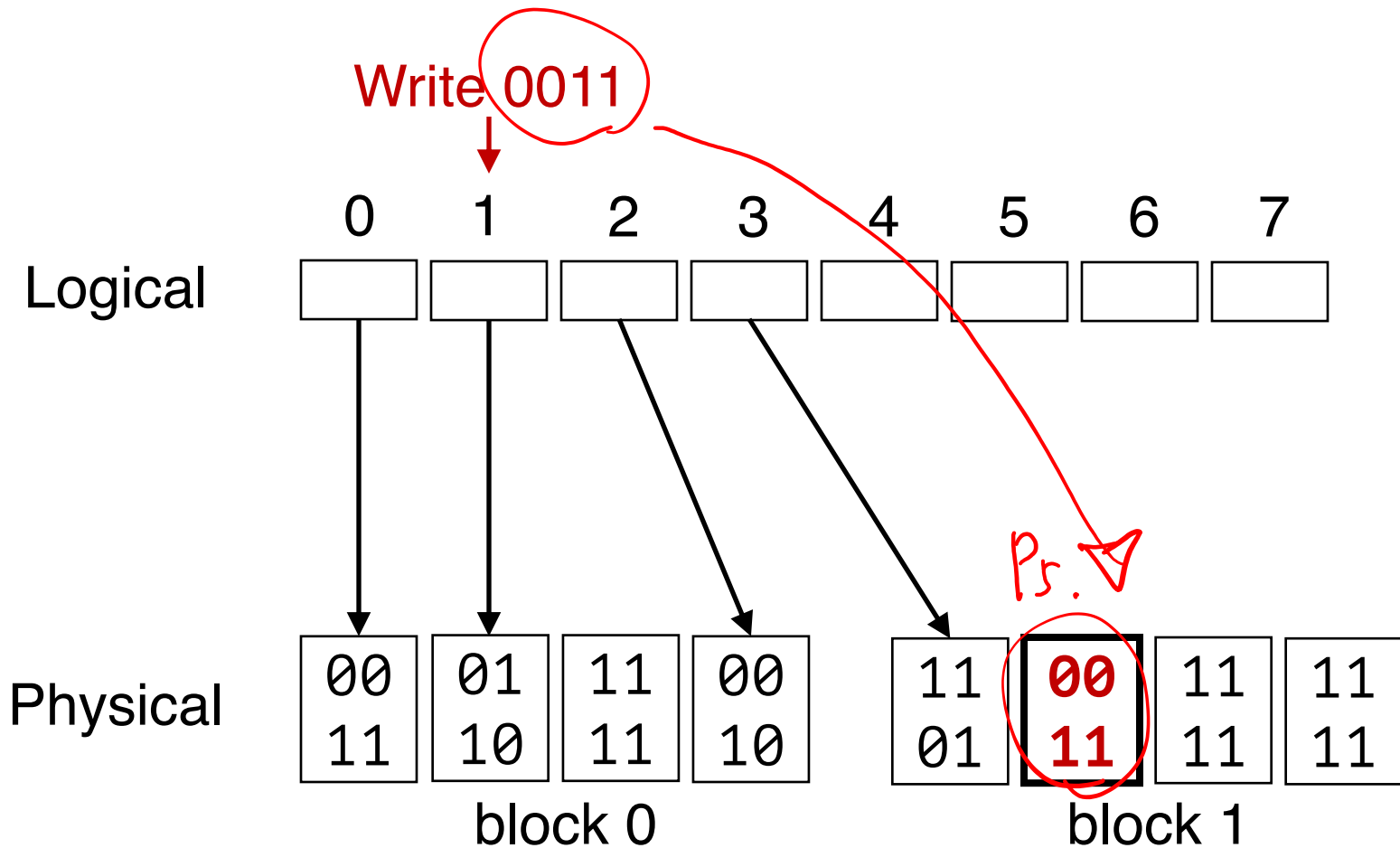
# Flash Translation Layer (FTL)



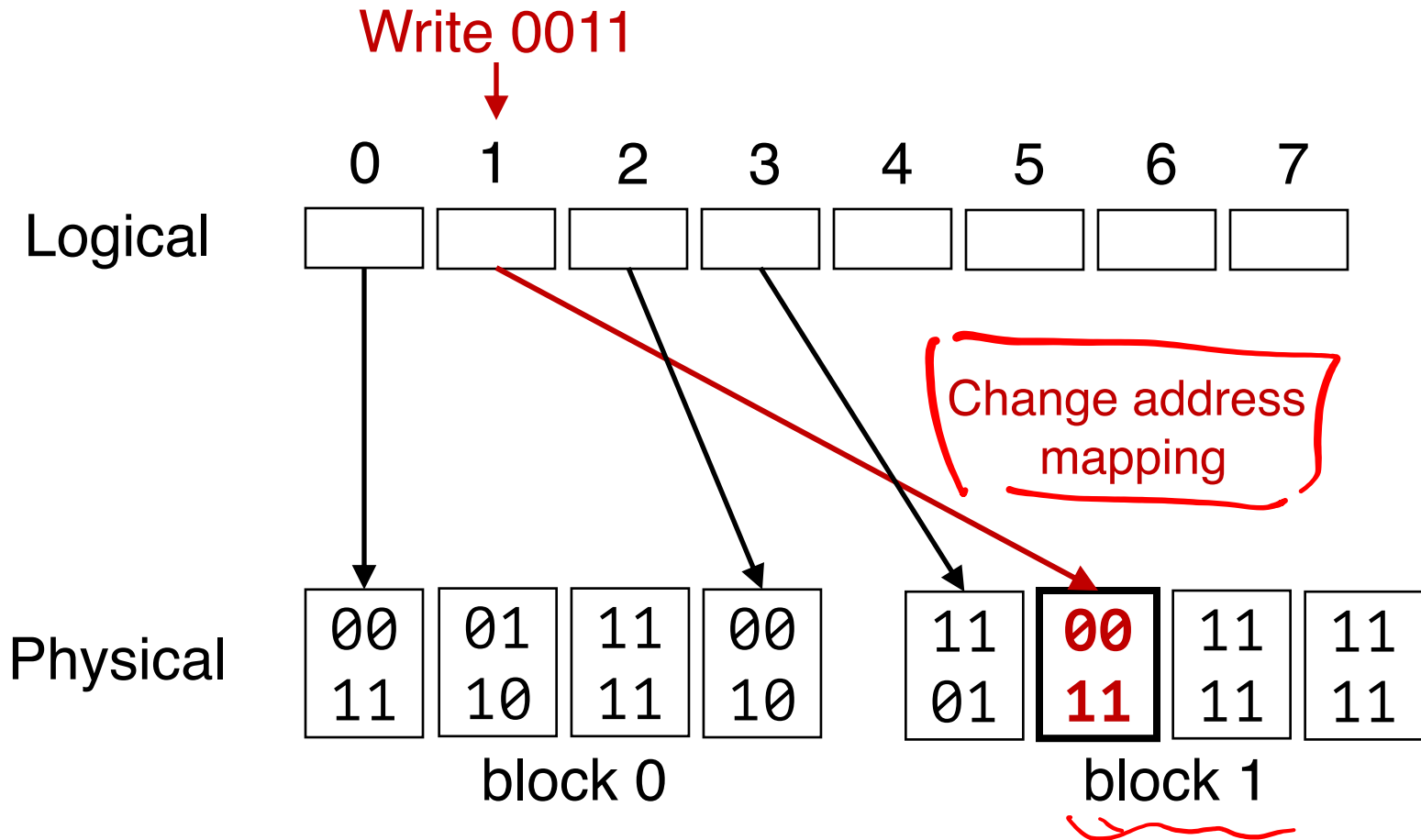
# Flash Translation Layer (FTL)



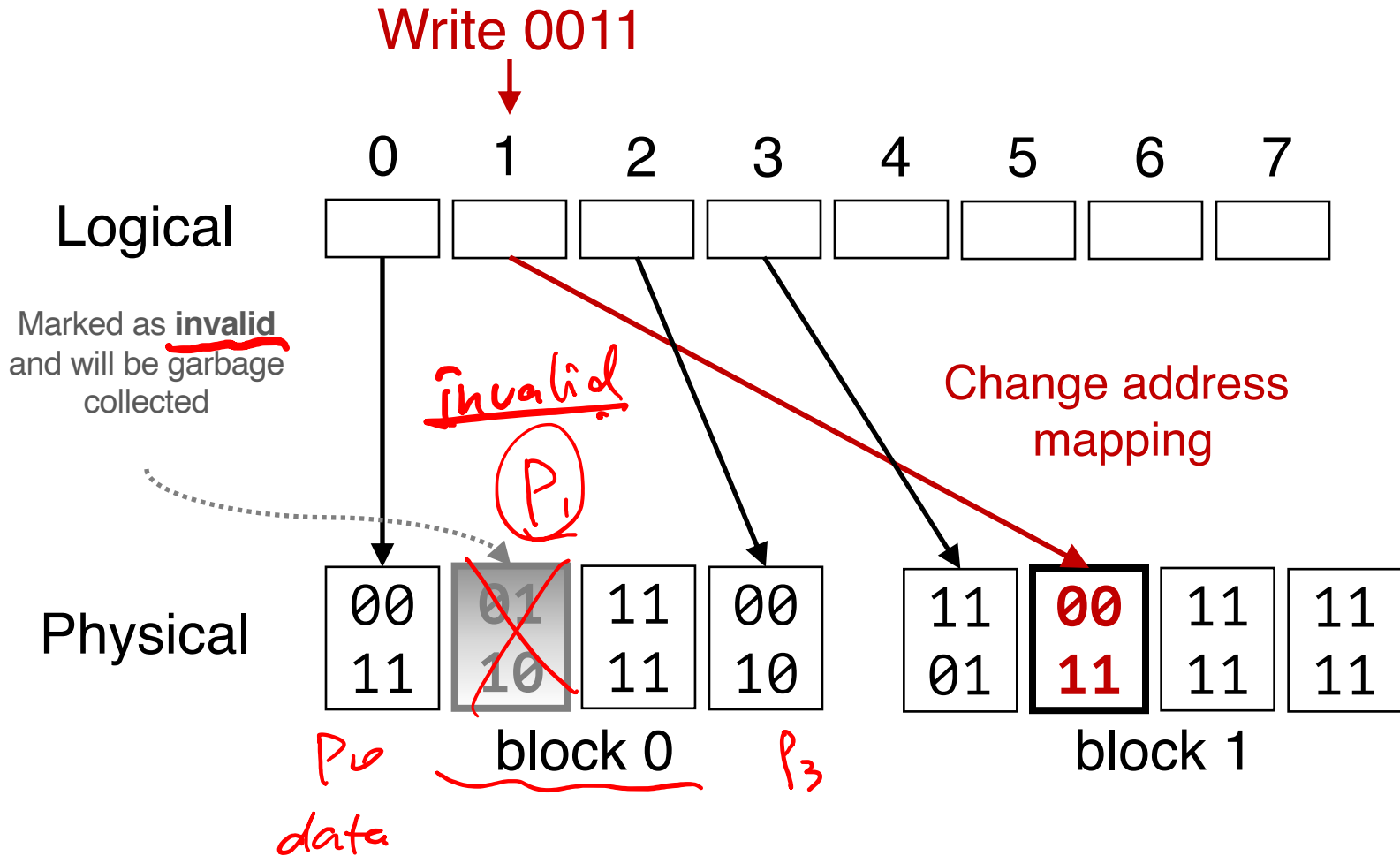
# Flash Translation Layer (FTL)



# Flash Translation Layer (FTL)

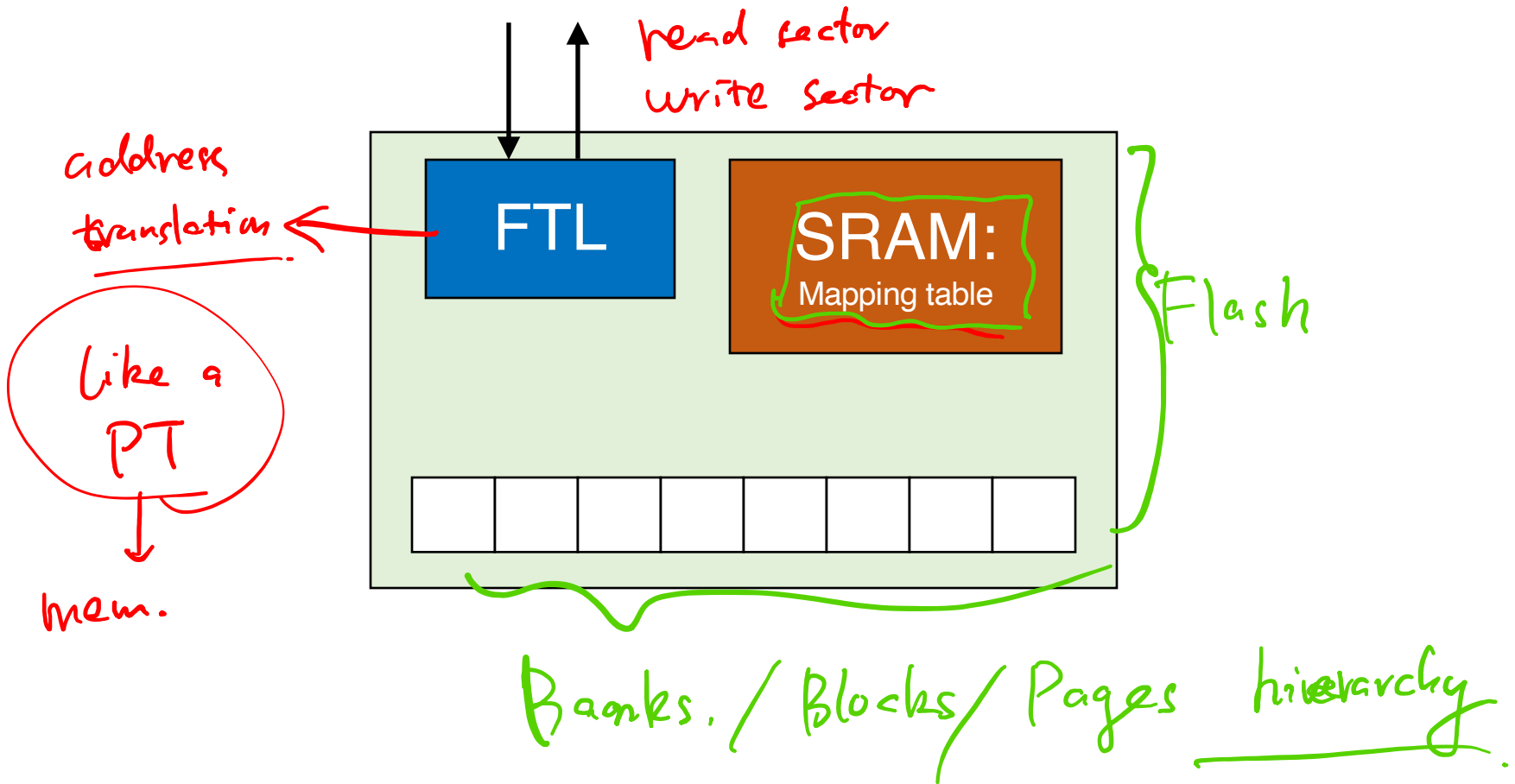


# Flash Translation Layer (FTL)



# SSD Architecture with FTL

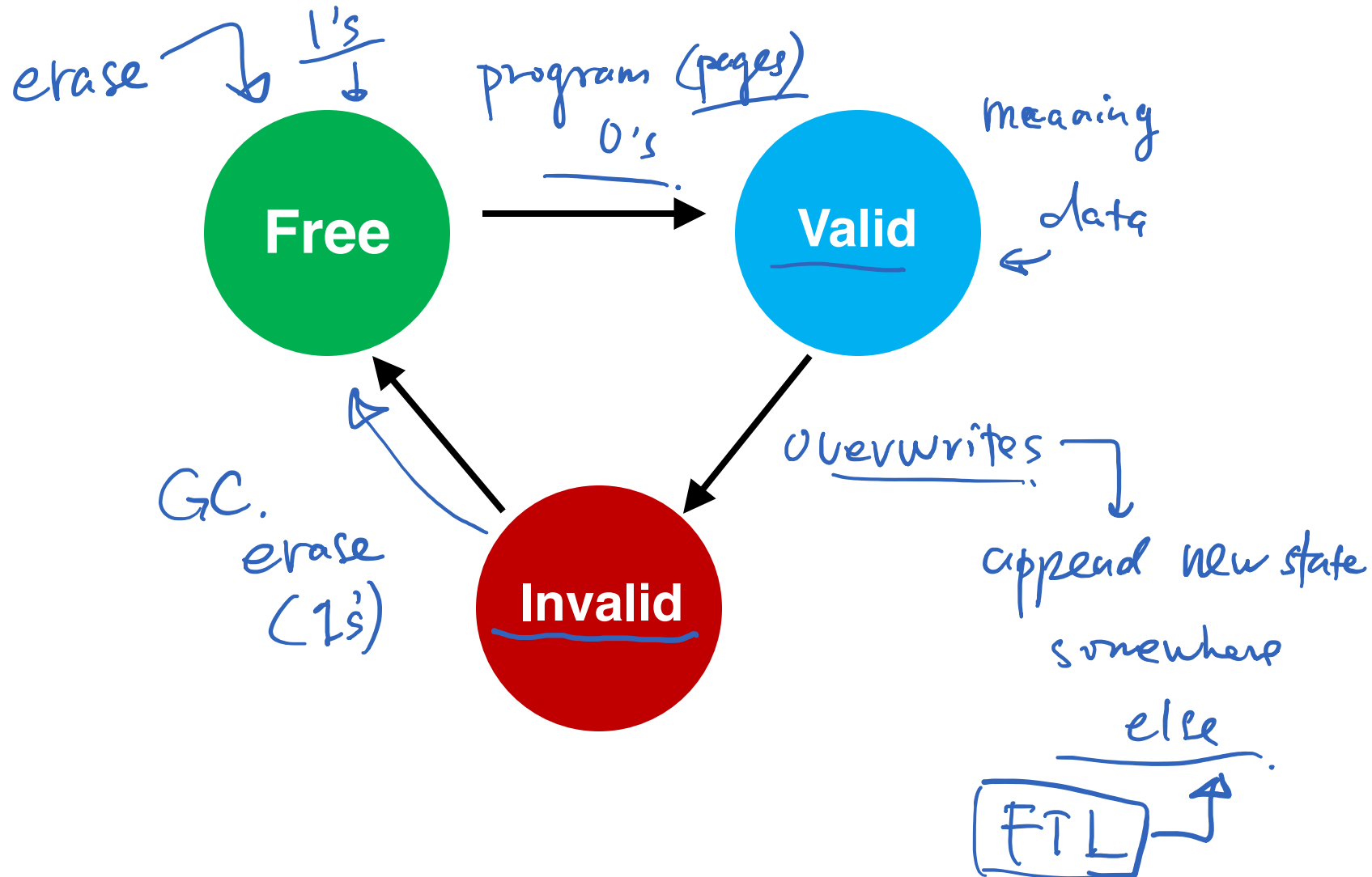
SSD provides disk-like interface



# Flash Translation Layer (FTL)

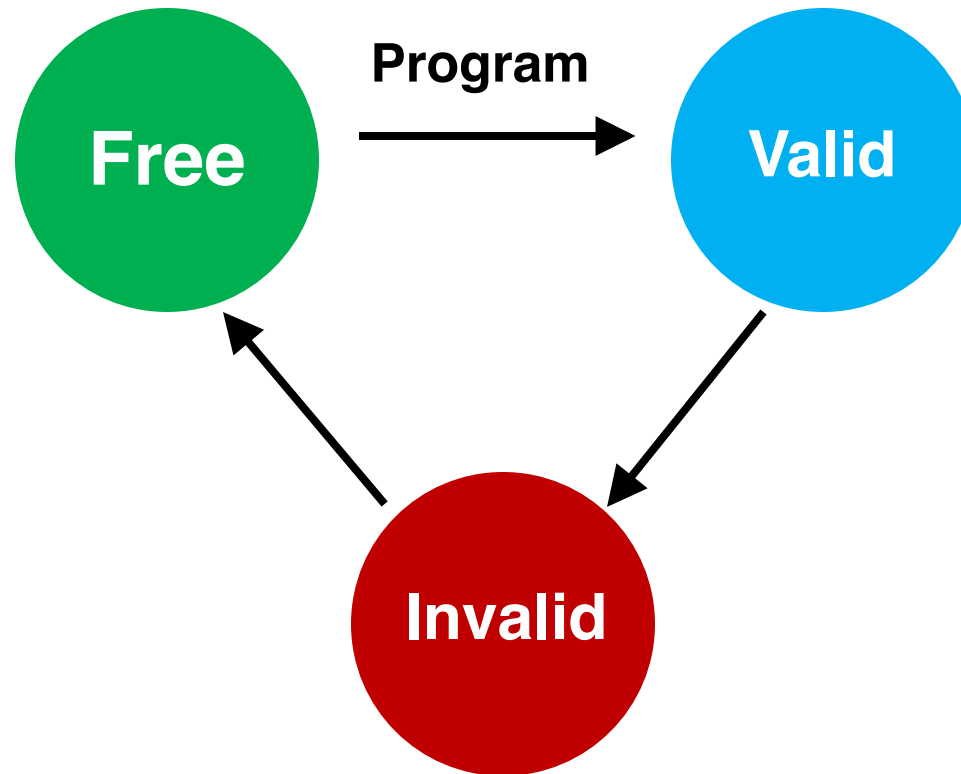
- Usually implemented in flash device's firmware
- Where to store mapping?
  - SRAM
- Physical pages can be in three states
  - valid, invalid, free

# State Transition of Physical Pages

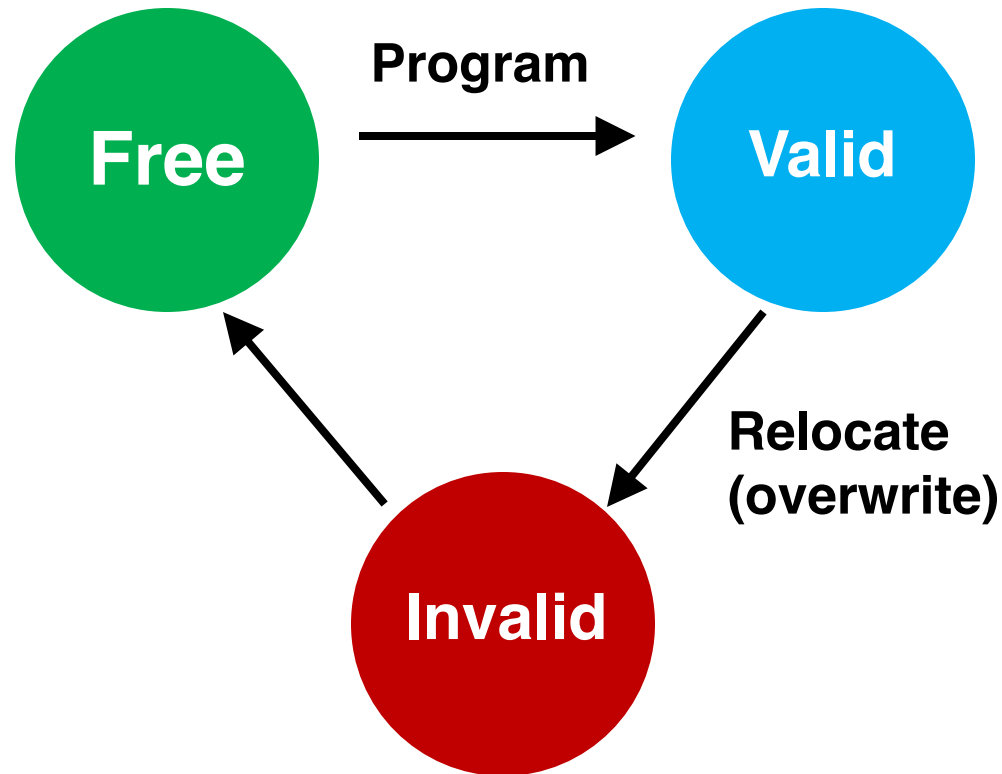




# State Transition of Physical Pages



# State Transition of Physical Pages



# State Transition of Physical Pages

