

# P&S Modern SSDs

Address Translation & Garbage Collection

Dr. Jisung Park

Prof. Onur Mutlu

ETH Zürich

Spring 2021

19 May 2021

# Recap: What We Have Discussed So Far

---

## ■ SSD Organization

- SSD controller
  - Multicore CPU running firmware (Flash Translation Layer, FTL)
  - Per-channel flash controller for ECC and data randomization
- DRAM: Metadata store (0.1% of SSD capacity)
- NAND flash packages (channel – die – plane – block – page)

## ■ NAND flash organization & operation

- Cell → NAND string & Page → Block → Plane → Die
- **Unique characteristics:** Erase-before-write, limited endurance, performance & operation-unit asymmetries, ...
- Basic operations: Page Read/Page Program/Block Erasure
- **Advanced operations:** Sub-page read, cache read, multi-plane, and program/erase suspend/resume commands

# Flash Translation Layer: Overview

---

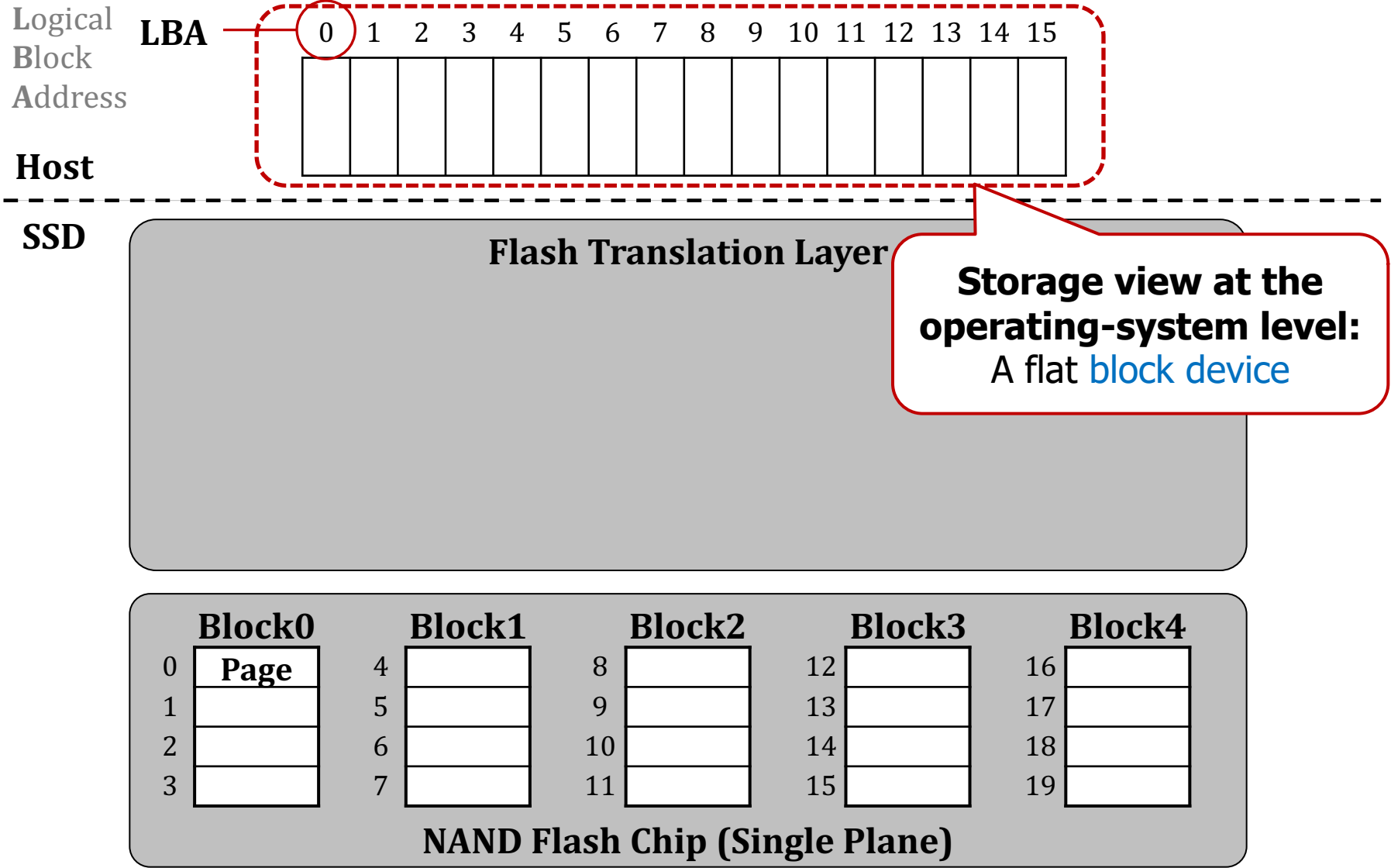
- SSD firmware (often referred to as SSD controller)
  - Provides **backward compatibility** with traditional HDDs
  - By **hiding unique characteristics** of NAND flash memory
- Responsible for many important **SSD-management tasks**
  - Address translation + garbage collection
    - Performs **out-of-place writes** due to erase-before-write property
  - Wear leveling
    - To prolong SSD lifetime by **evenly distributing** P/E cycles
  - Data refresh
    - Resets transient errors by **copying data** to a new page(s)
  - I/O scheduling
    - To take full advantage of **SSD internal parallelism**

# Flash Translation Layer: Overview

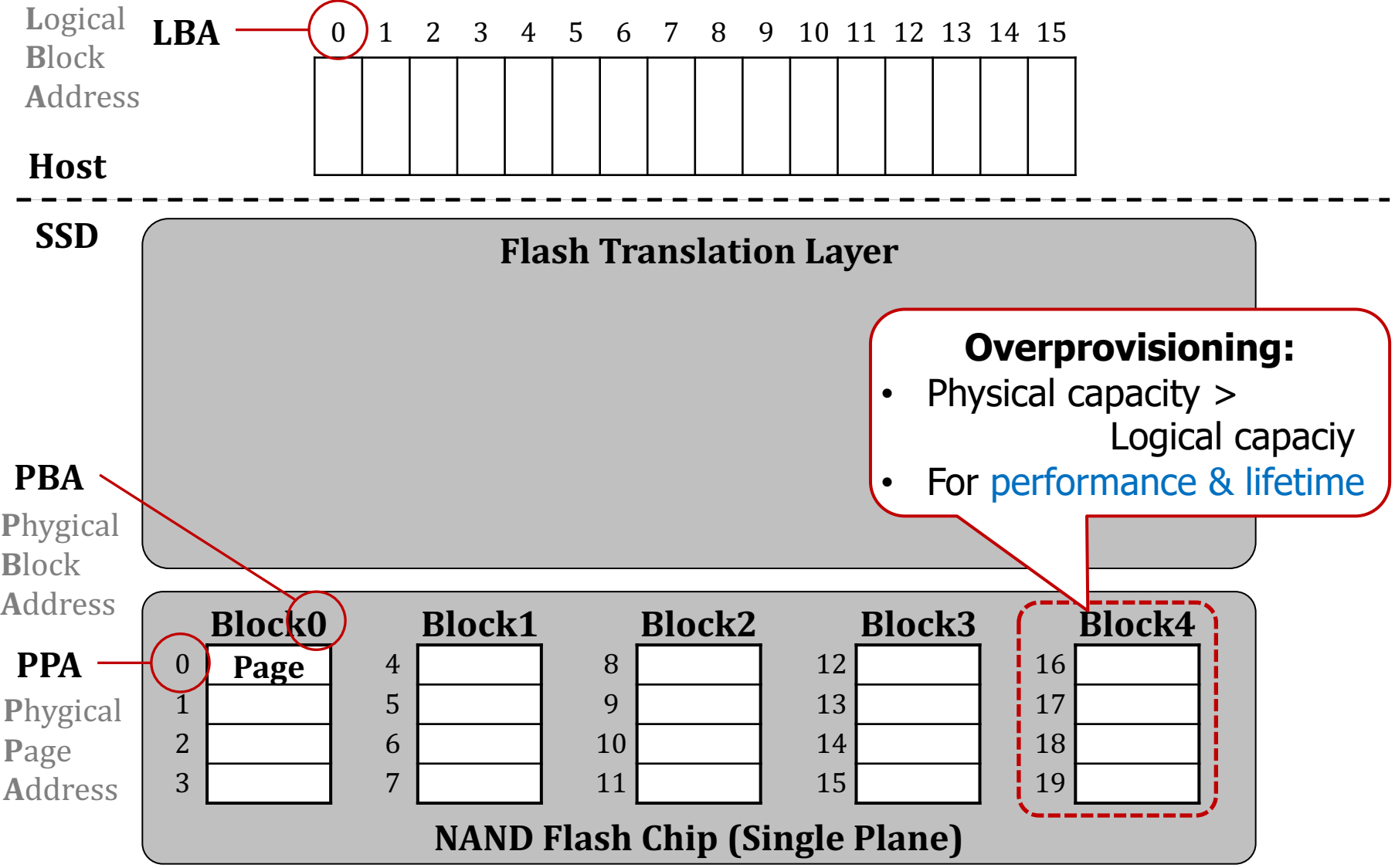
---

- SSD firmware (often referred to as SSD controller)
  - Provides **backward compatibility** with traditional HDDs
  - By **hiding unique characteristics** of NAND flash memory
- Responsible for many important **SSD-management tasks**
  - Address translation + garbage collection
    - Performs **out-of-place writes** due to erase-before-write property
  - Wear leveling
    - To prolong SSD lifetime by **evenly distributing** P/E cycles
  - Data refresh
    - Resets transient errors by **copying data** to a new page(s)
  - I/O scheduling
    - To take full advantage of **SSD internal parallelism**

# Simple SSD Architecture

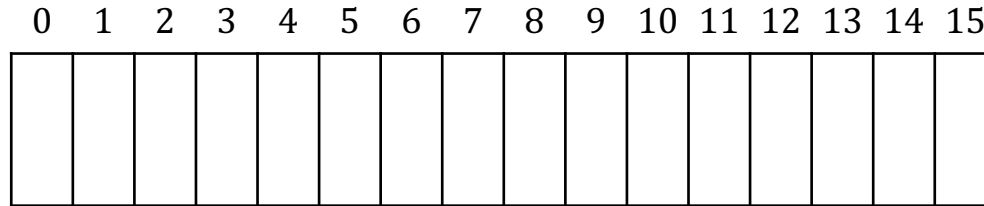


# Simple SSD Architecture



# Write Request Handling: Page Write

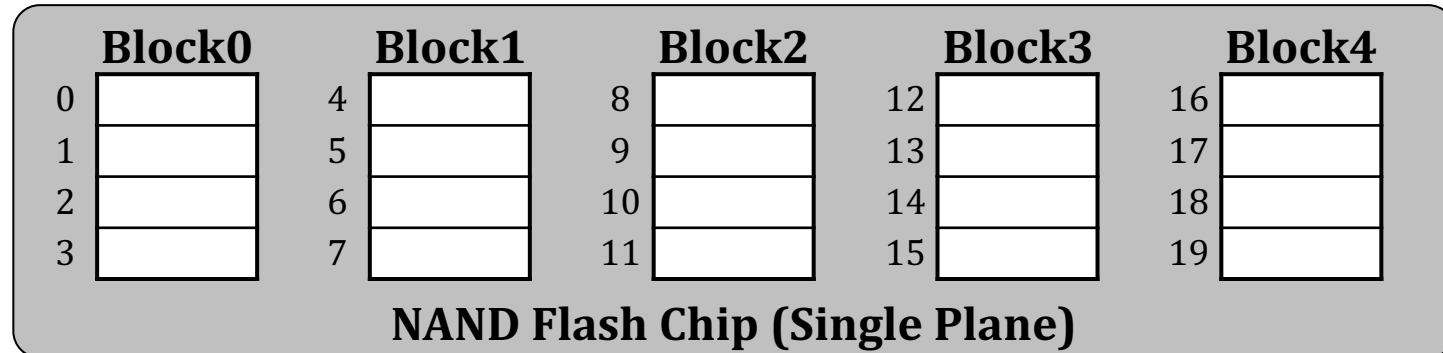
---



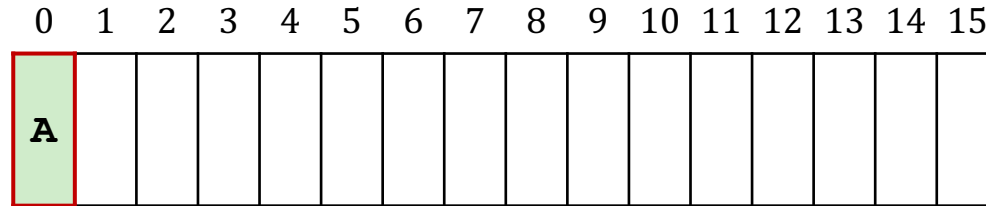
Host

SSD

**Flash Translation Layer**



# Write Request Handling: Page Write



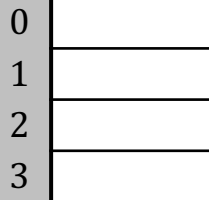
Host

SSD

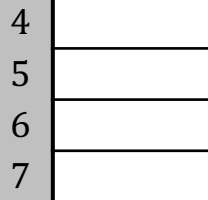
Flash Translation Layer

Req (LBA: 0, Size: 1, DIR: W, A)

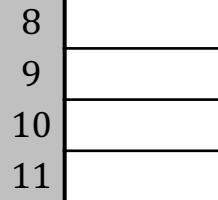
Block0



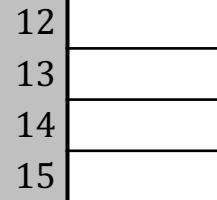
Block1



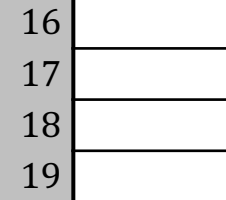
Block2



Block3



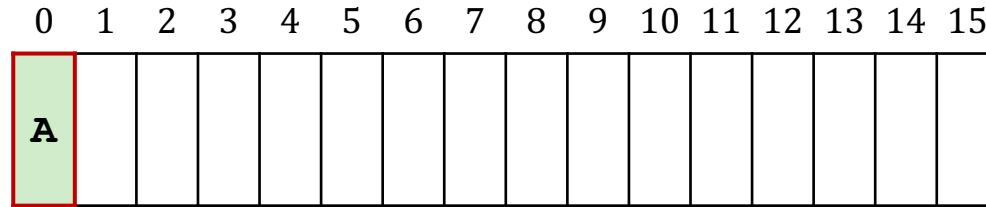
Block4



NAND Flash Chip (Single Plane)



# Write Request Handling: Page Write

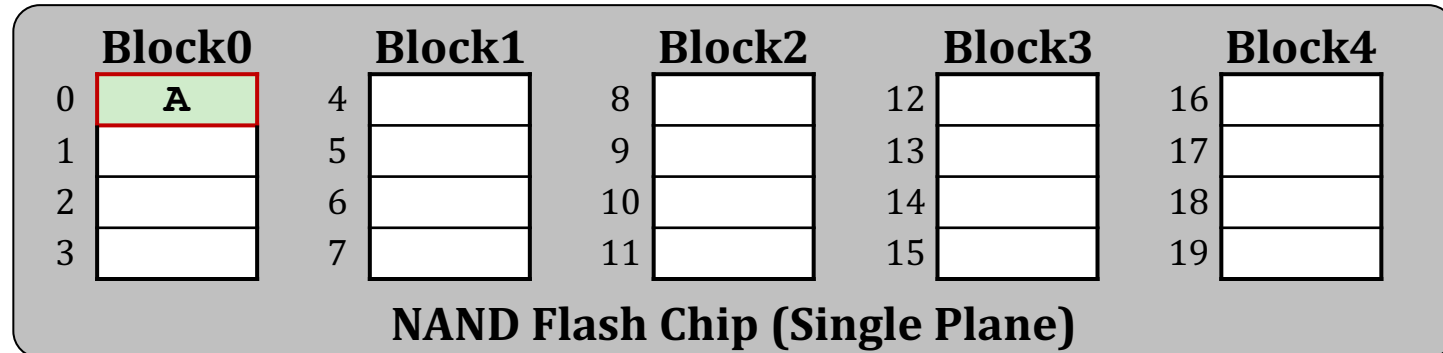


Host

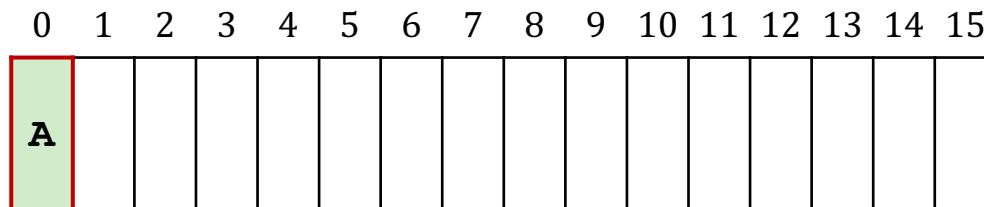
SSD

Flash Translation Layer

Req (LBA: 0, Size: 1, DIR: W, A)



# Write Request Handling: Page Write



Host

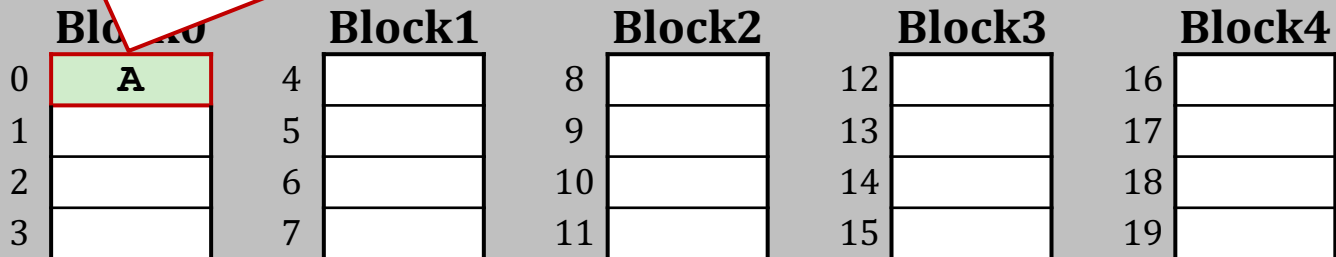
SSD

## Flash Translation Layer

Req (LBA: 0, Size: 1, DIR: W, A)

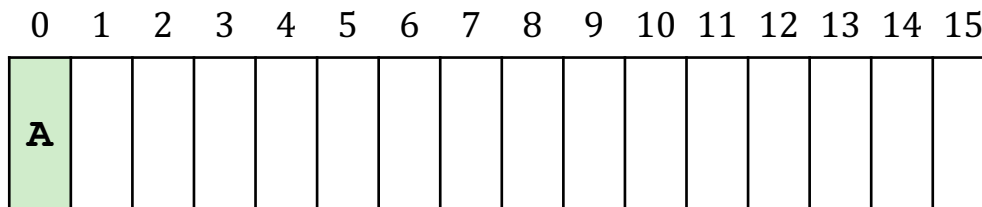
### Note:

- We are assuming that logical block size = physical page size
- LB size = 4 KiB, PP size = 16 KiB



NAND Flash Chip (Single Plane)

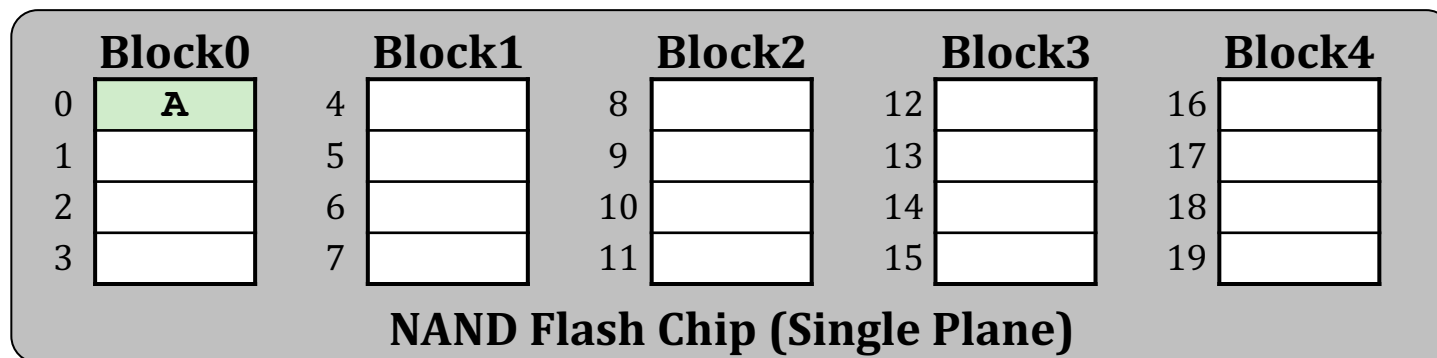
# Write Request Handling: Sequential Write



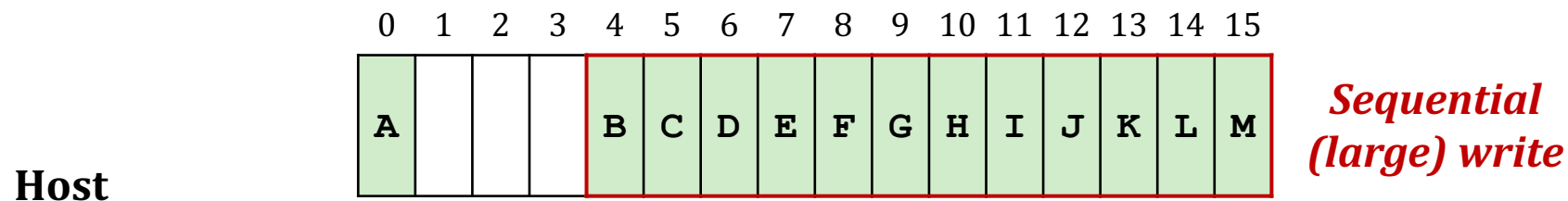
Host

SSD

Flash Translation Layer



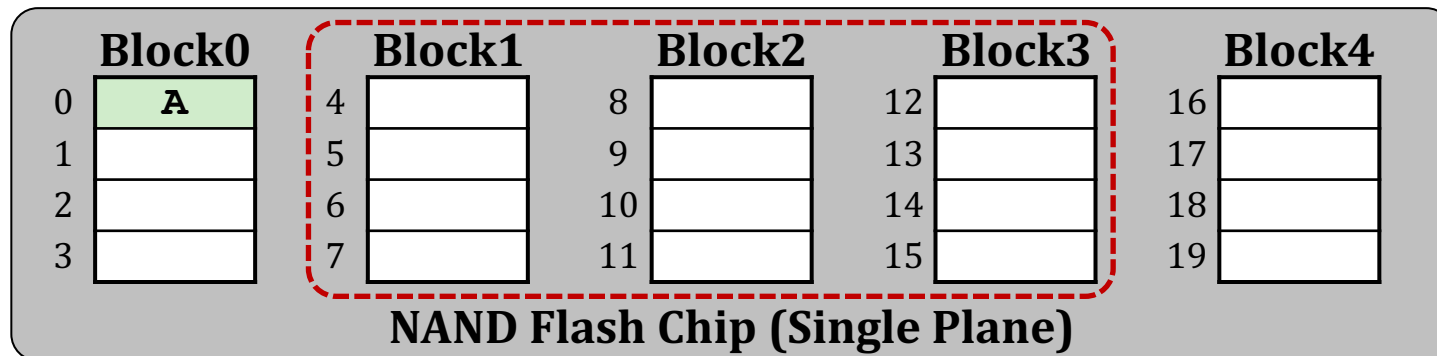
# Write Request Handling: Sequential Write



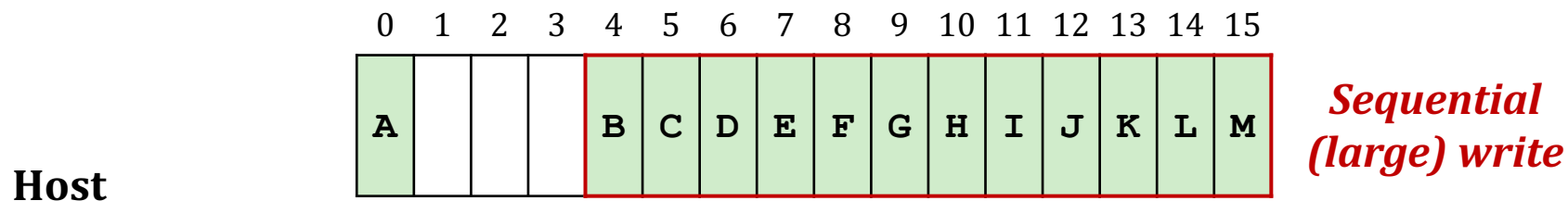
SSD

Flash Translation Layer

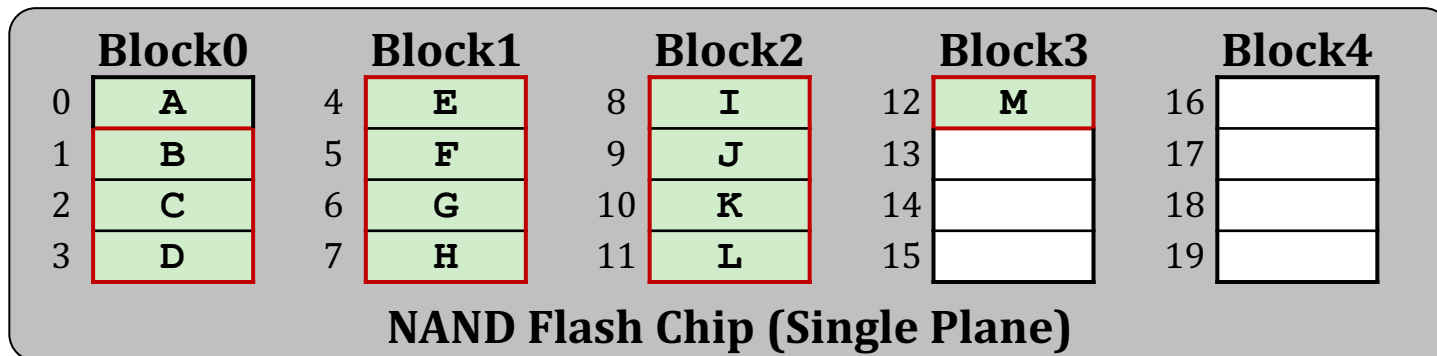
Req (LBA: 4, Size: 12, DIR: W, B ... M)



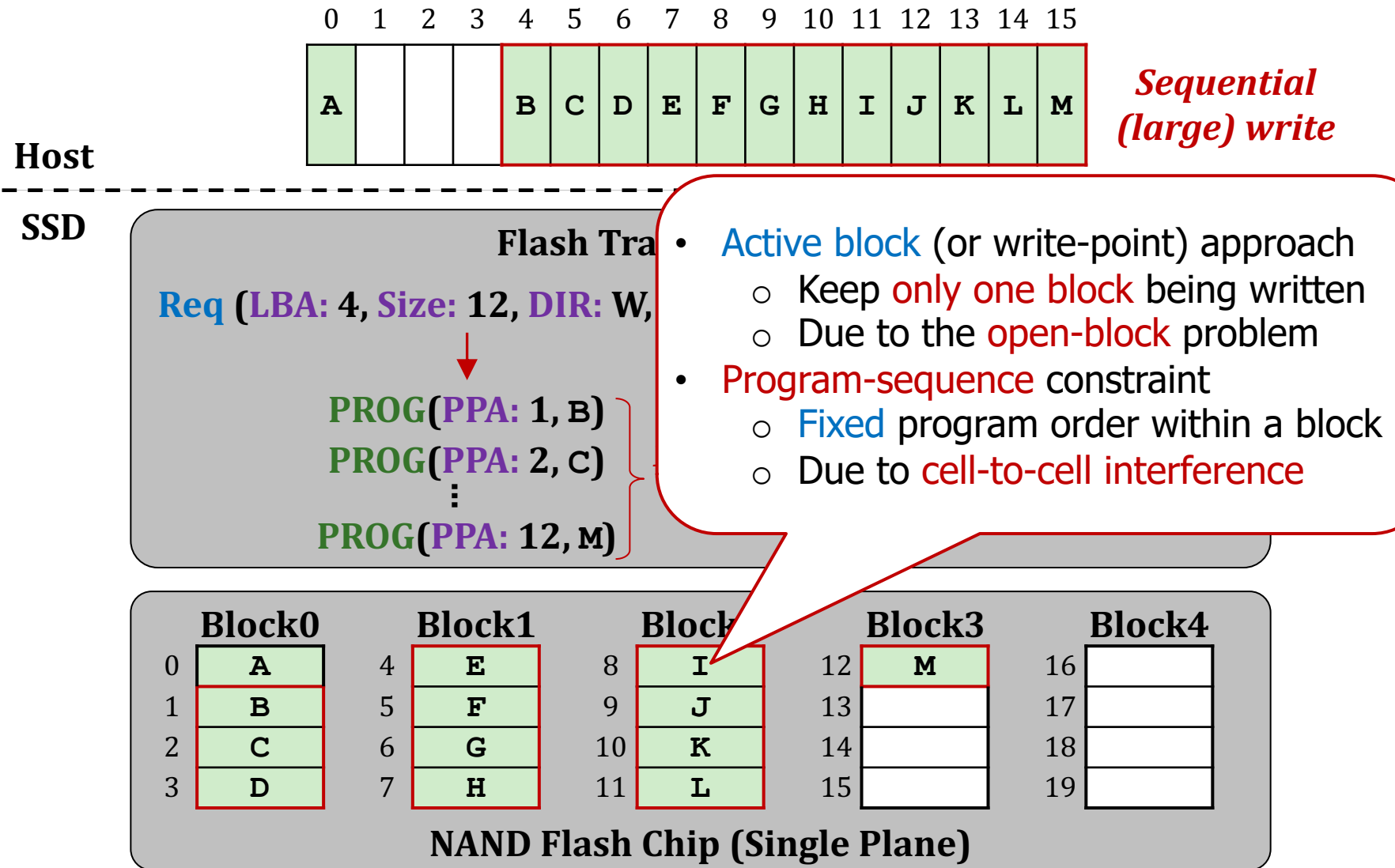
# Write Request Handling: Sequential Write



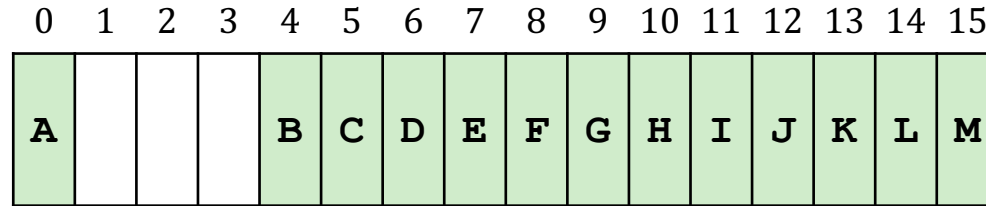
SSD



# Write Request Handling: Sequential Write



# Write Request Handling: Address Mapping

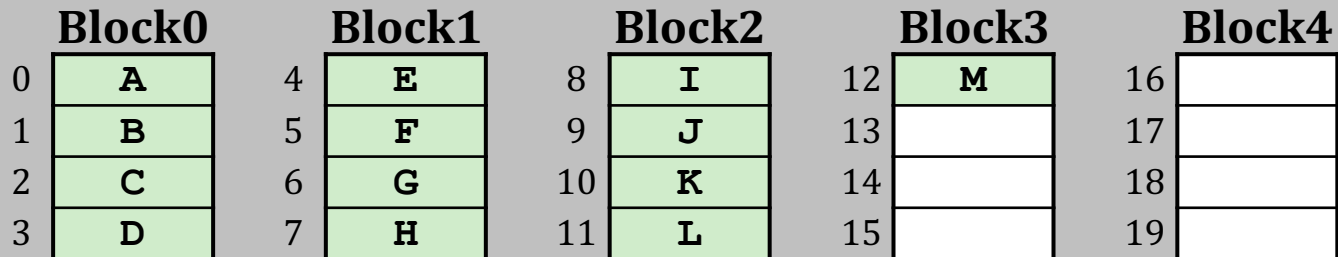


Host

SSD

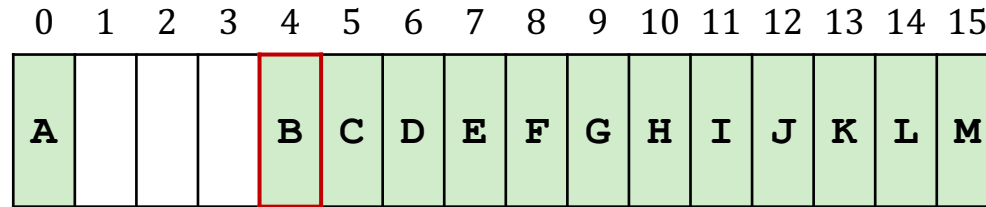
## Flash Translation Layer

**Problem: LBA (or LPA) does not match PPA!**



**NAND Flash Chip (Single Plane)**

# Write Request Handling: Address Mapping



Host

SSD

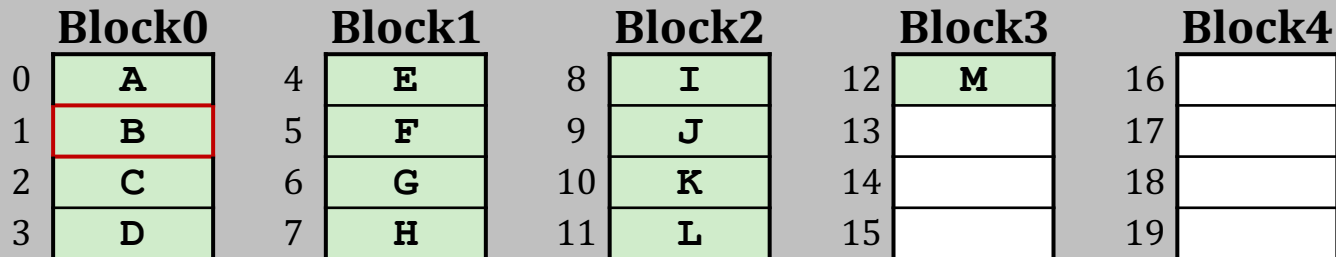
## Flash Translation Layer

**Problem: LBA (or LPA) does not match PPA!**

Req (LBA: 4, Size: 1, DIR: R)

↓  
READ (PPA: ?)

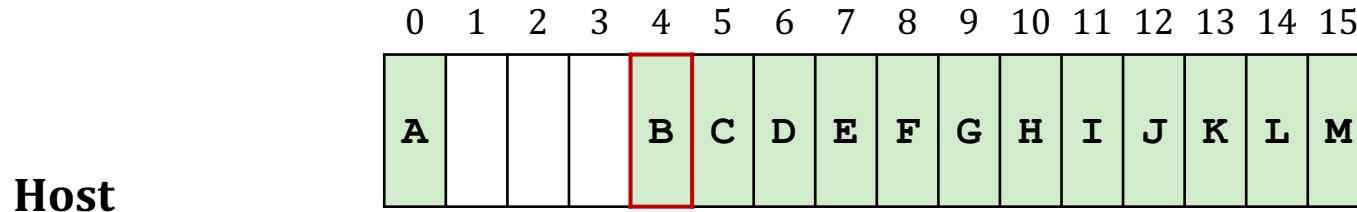
*Needs to maintain  
Address-mapping information*



NAND Flash Chip (Single Plane)



# Write Request Handling: Address Mapping



SSD

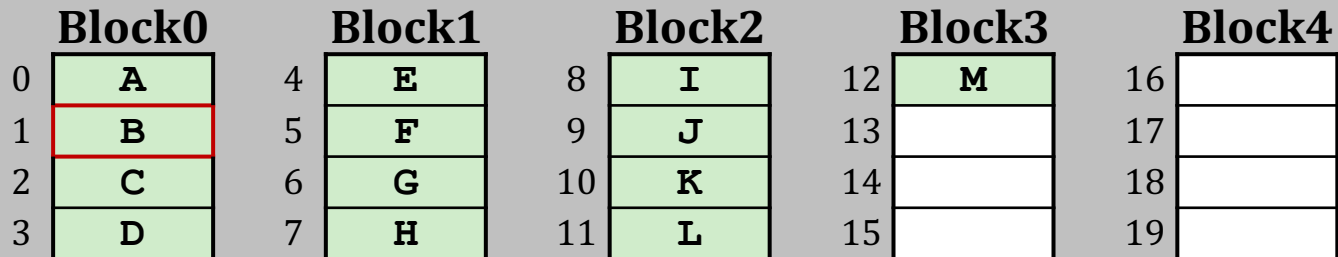
## Flash Translation Layer

Req (LBA: 4, Size: 1, DIR: R)

READ (PPA: ?)

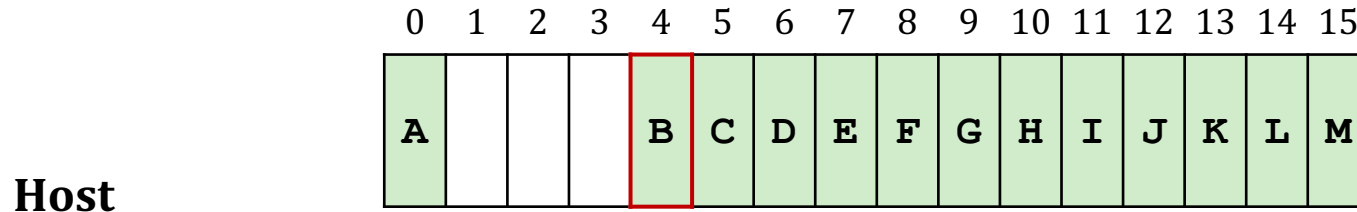
LPA	PPA
0	0
...	...
4	1
5	2
...	...

Mapping Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Address Mapping



SSD

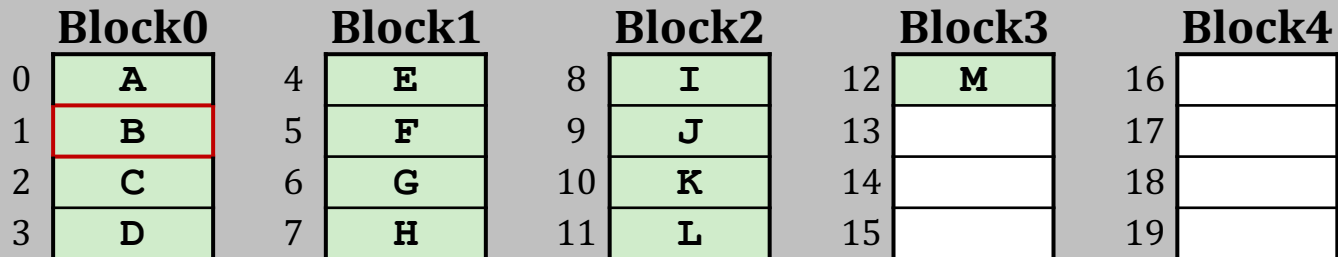
## Flash Translation Layer

Req (LBA: 4, Size: 1, DIR: R)

↓  
READ (PPA: 1)

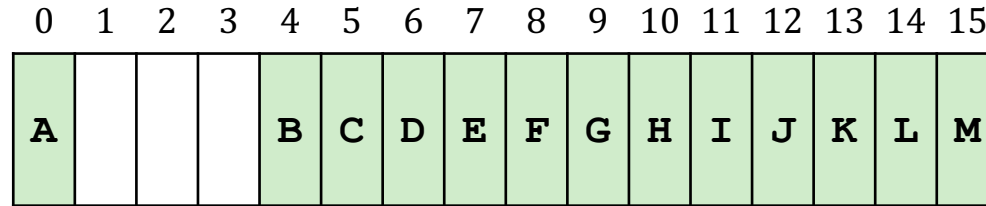
LPA	PPA
0	0
...	...
4	1
5	2
...	...

Mapping Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Update



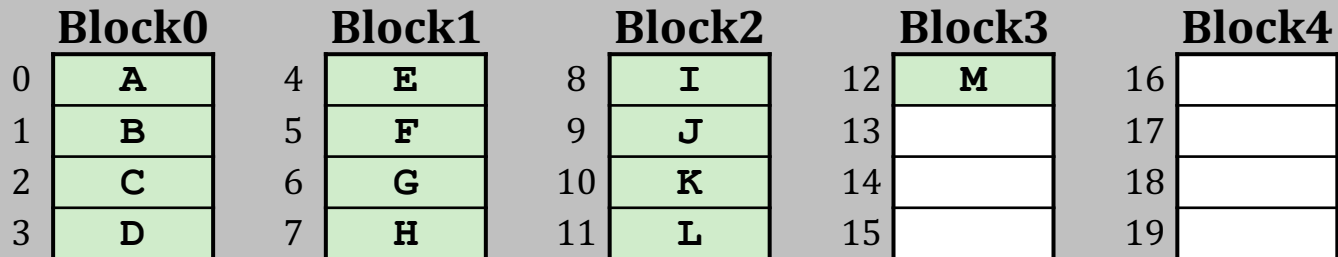
Host

SSD

## Flash Translation Layer

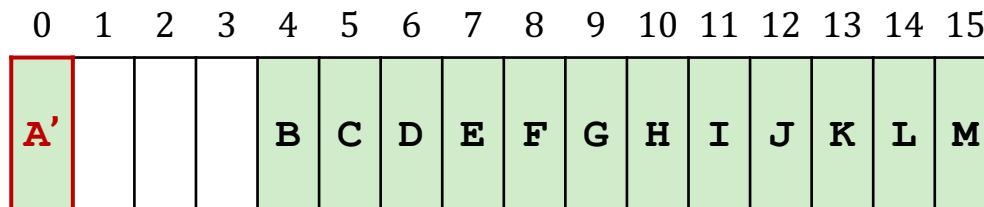
LPA	PPA
0	0
...	...
4	1
5	2
...	...

Mapping Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Update



Host

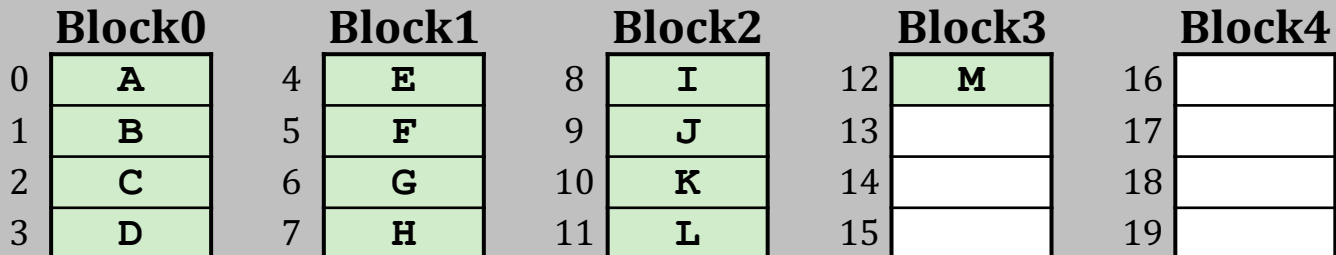
SSD

## Flash Translation Layer

Req (LBA: 0, Size: 1, DIR: W, A')

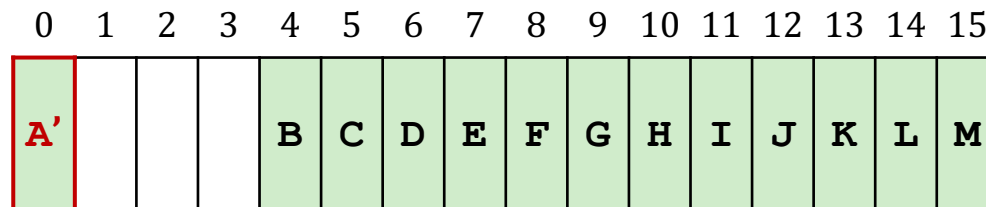
LPA	PPA
0	0
...	...
4	1
5	2
...	...

Mapping Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Update



Host

SSD

## Flash Translation Layer

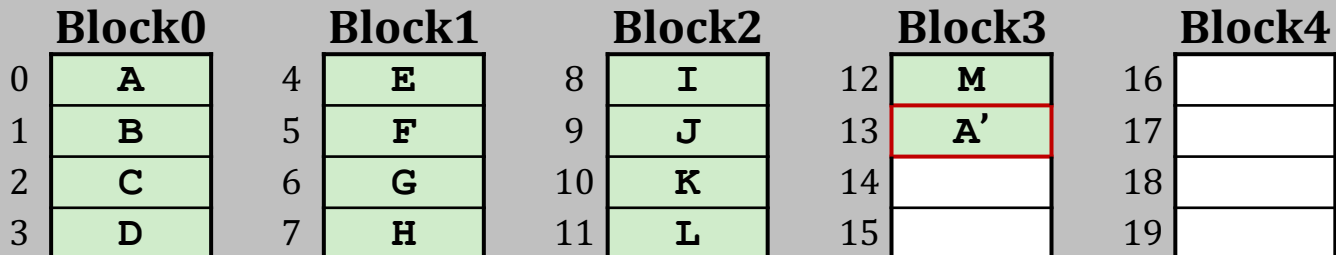
Req (LBA: 0, Size: 1, DIR: W, A')



READ (PPA: 13)

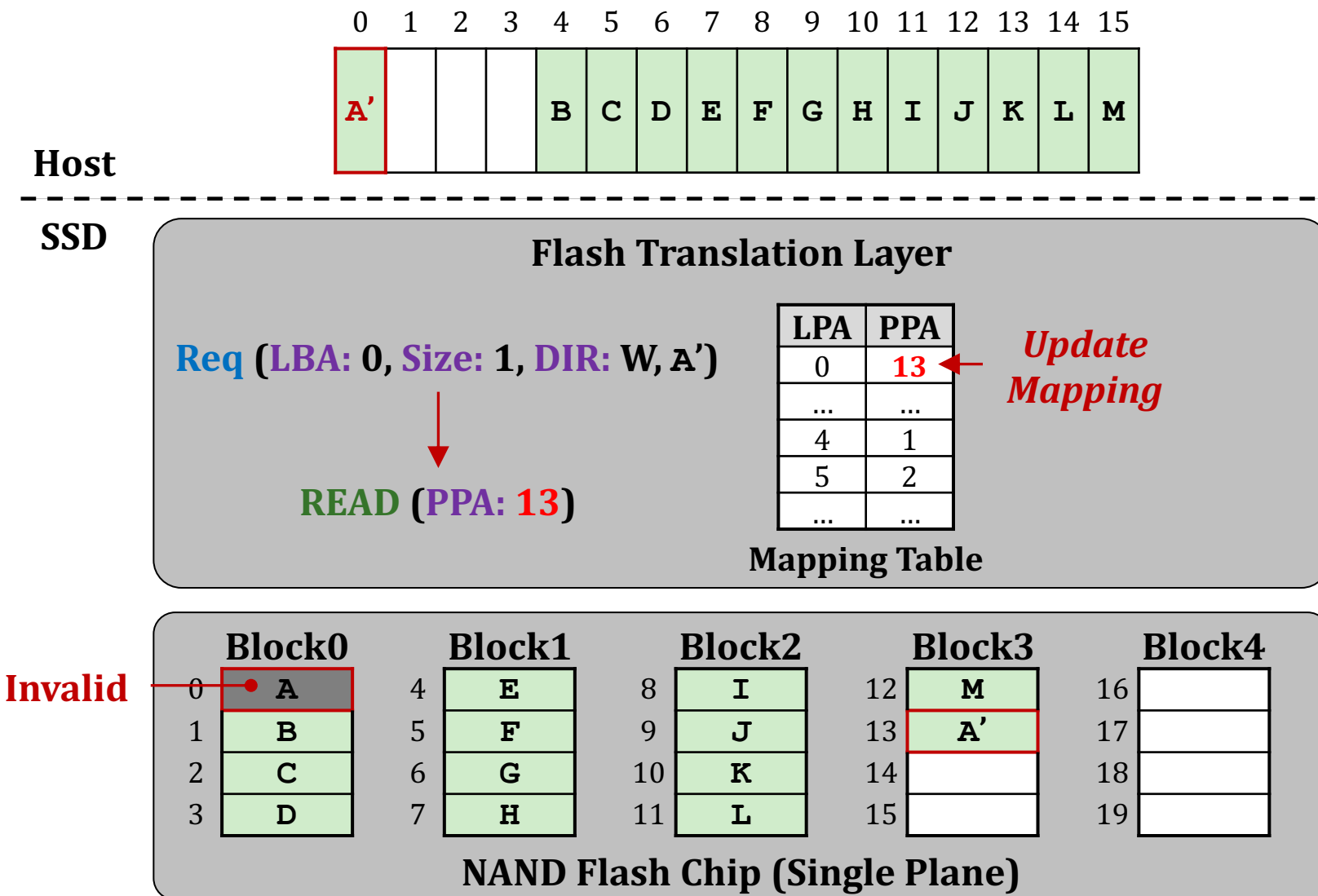
LPA	PPA
0	0
...	...
4	1
5	2
...	...

Mapping Table

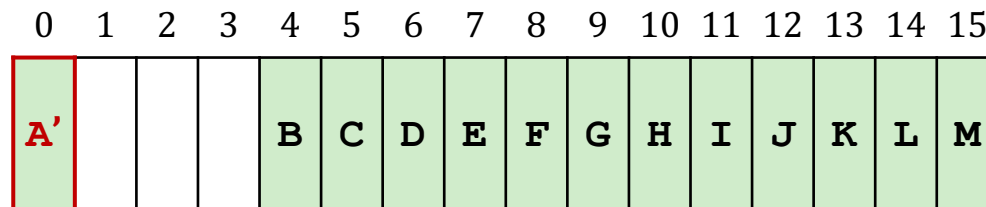


NAND Flash Chip (Single Plane)

# Write Request Handling: Update



# Write Request Handling: Update



Host

SSD

## Flash Translation Layer

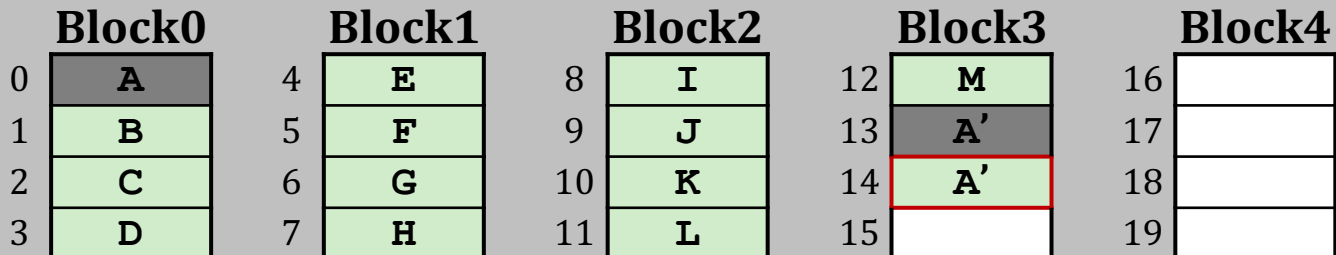
Req (LBA: 0, Size: 1, DIR: W, A')



READ (PPA: 14)

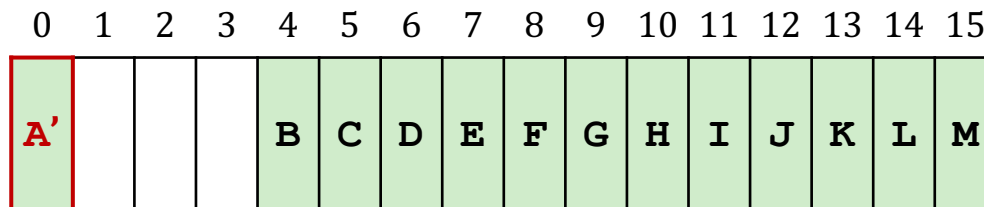
LPA	PPA
0	14
...	...
4	1
5	2
...	...

Mapping Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Update



Host

SSD

## Flash Translation Layer

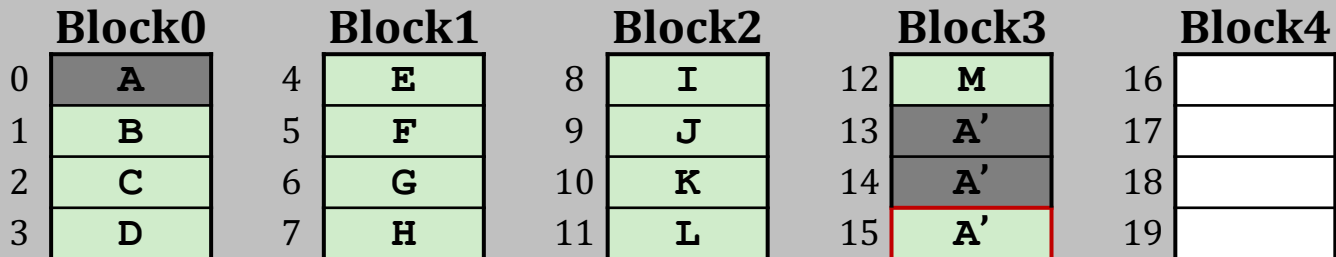
Req (LBA: 0, Size: 1, DIR: W, A')



READ (PPA: 15)

LPA	PPA
0	15
...	...
4	1
5	2
...	...

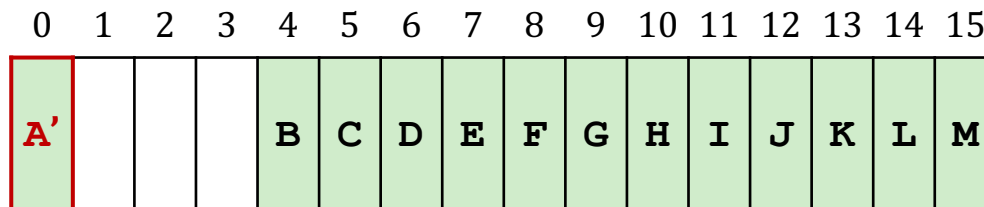
Mapping Table



NAND Flash Chip (Single Plane)

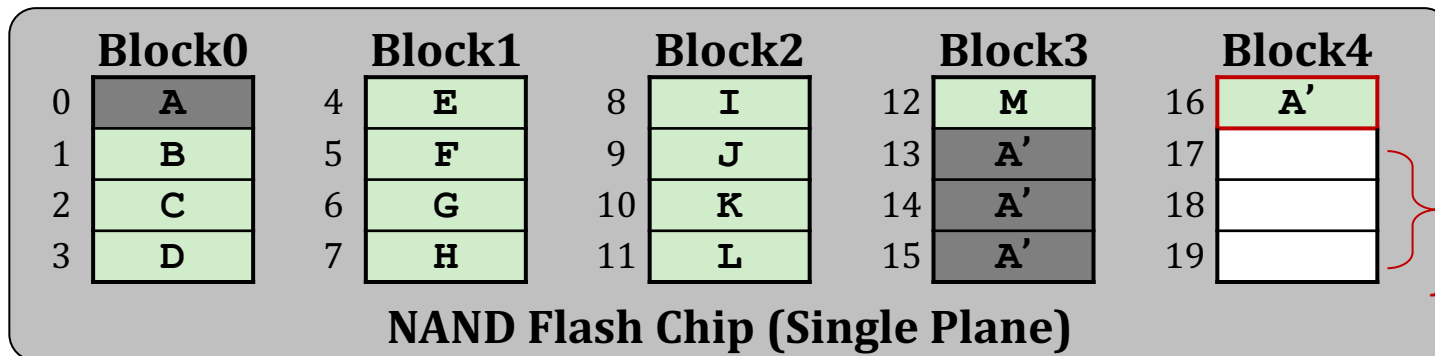
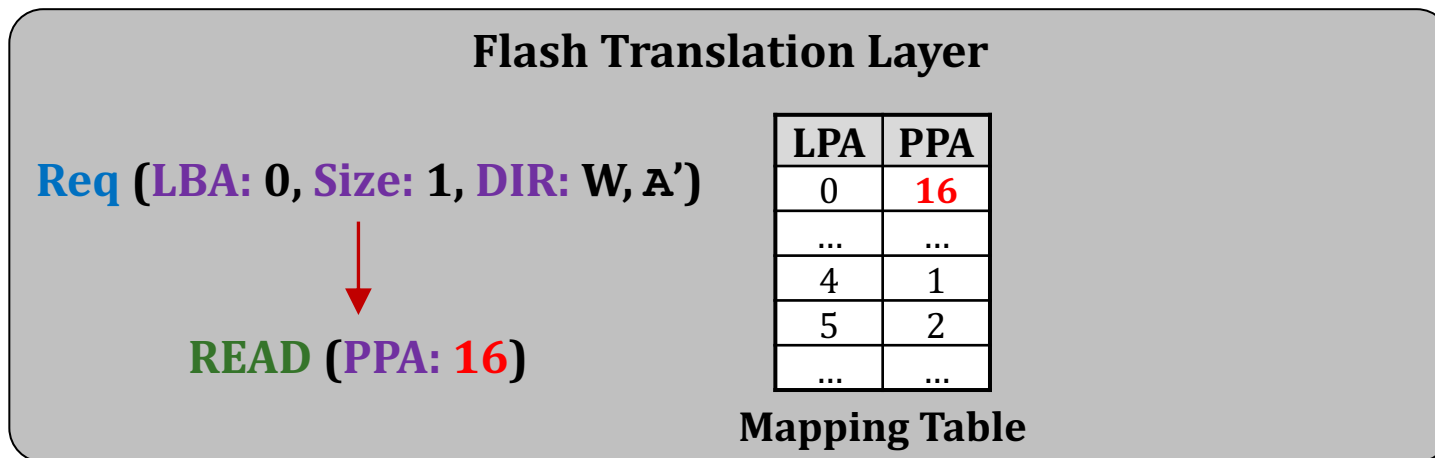


# Write Request Handling: Update



Host

SSD

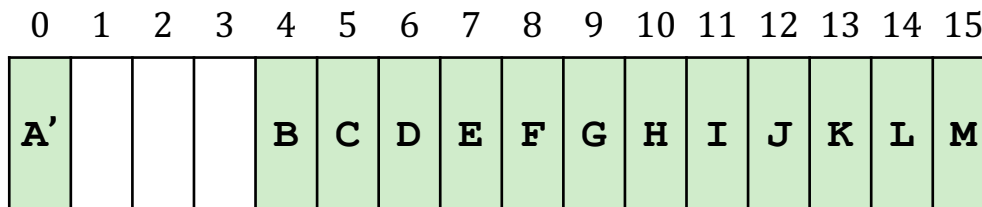


# Garbage Collection

---

- Reclaims **free pages** by erasing **invalid** pages
  - Erase unit: **block**
  - If a victim block (to erase) has **valid pages**, all the valid pages **need to be copied** to other free pages
    - **Performance overhead:**  $(t_{\text{READ}} + t_{\text{PROG}}) \times \#$  of valid pages
    - **Lifetime overhead:** additional writes  $\rightarrow$  P/E-cycle increase
- **Greedy** victim-selection policy:
  - Erases the block with the **largest number** of invalid pages
  - Needs to maintain **# of invalid (or valid) pages** for each block

# Write Request Handling: Garbage Collection



Host

SSD

## Flash Translation Layer

F: free, V: valid, I: invalid

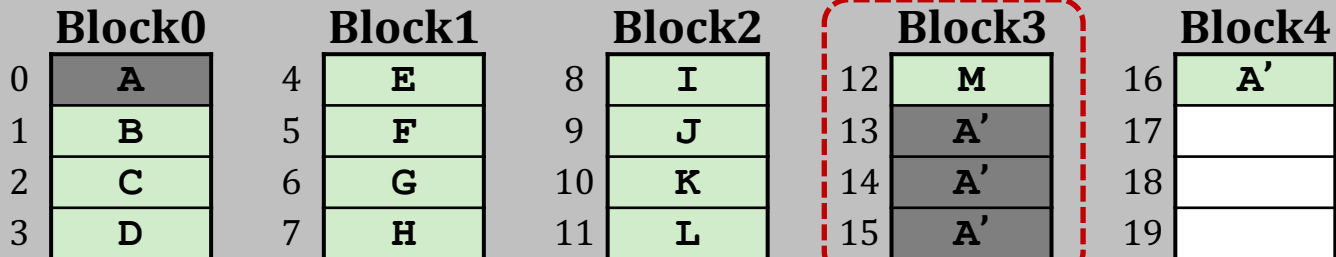
READ (PPA: 12)

LPA	PPA
0	16
...	...
4	1
5	2
...	...

Mapping Table

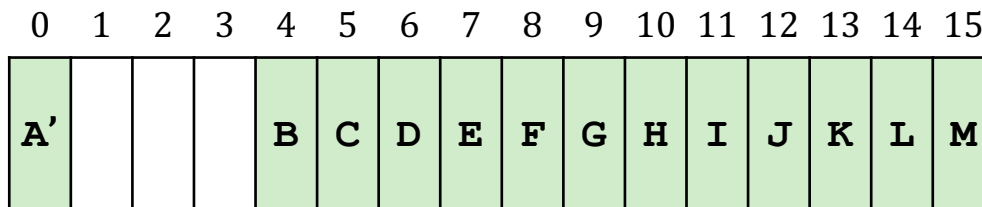
PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	VIII
4	VFFF

Status Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Garbage Collection



Host

SSD

## Flash Translation Layer

F: free, V: valid, I: invalid

READ (PPA: 12)

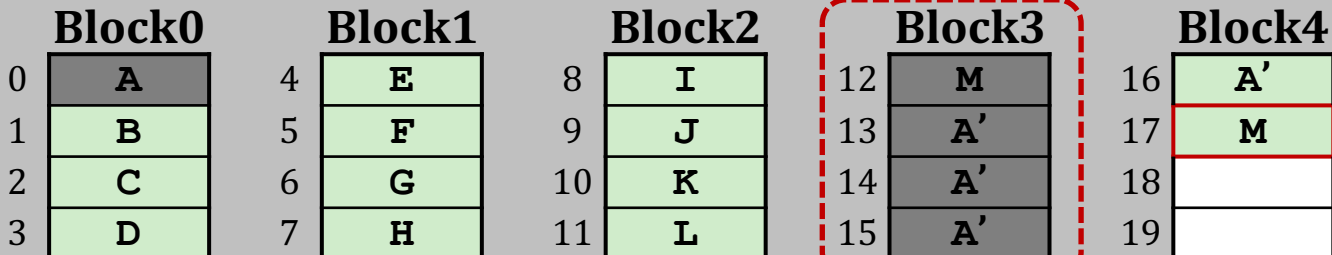
PROG (PPA: 17, M)

LPA	PPA
0	16
...	...
4	1
5	2
...	...

Mapping Table

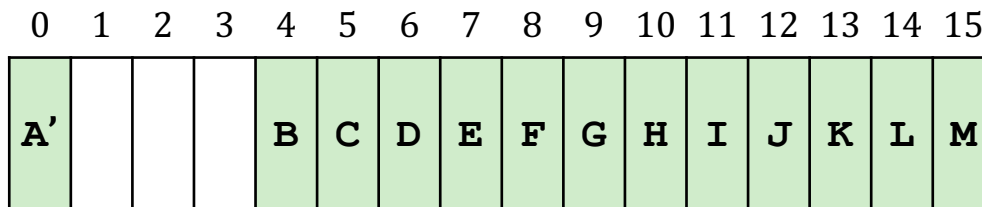
PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	VIII
4	VFFF

Status Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Garbage Collection



Host

SSD

## Flash Translation Layer

F: free, V: valid, I: invalid

READ (PPA: 12)  
 PROG (PPA: 17, M)

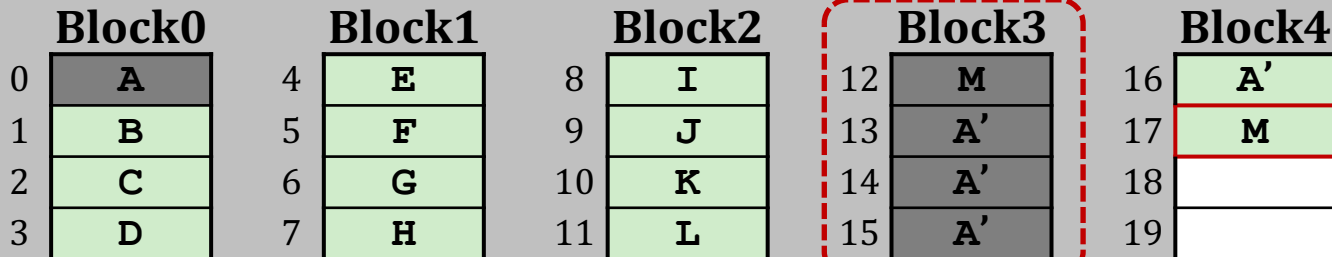
LPA	PPA
0	16
...	...
4	1
5	2
...	...

Mapping Table

PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	IIII
4	VVFF

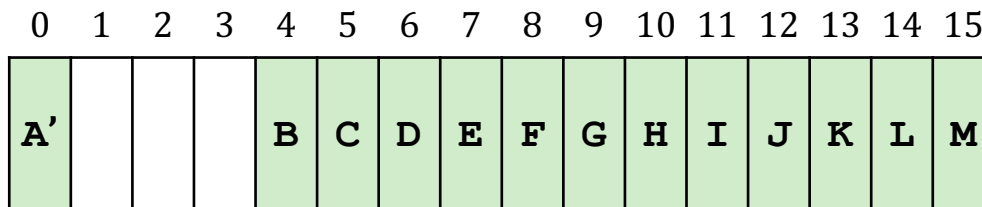
Status Table

*Update Status*

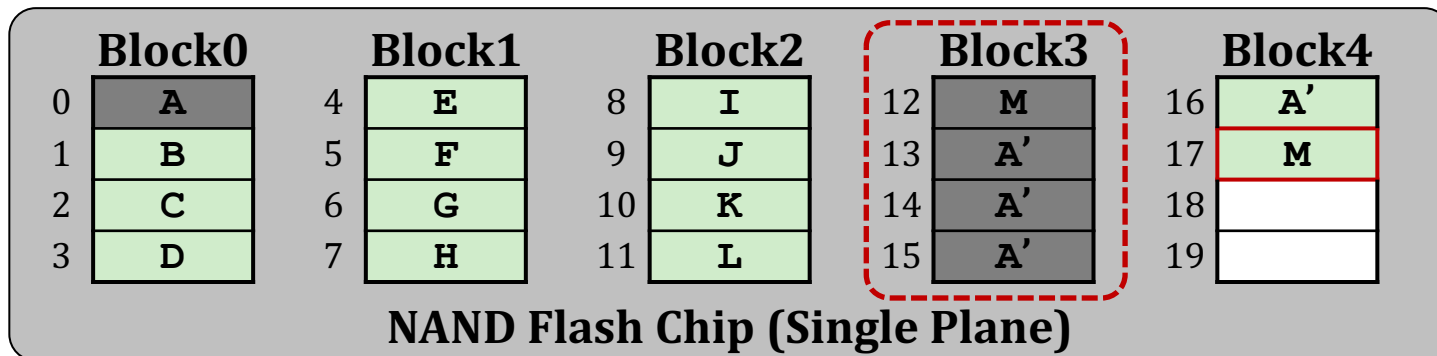
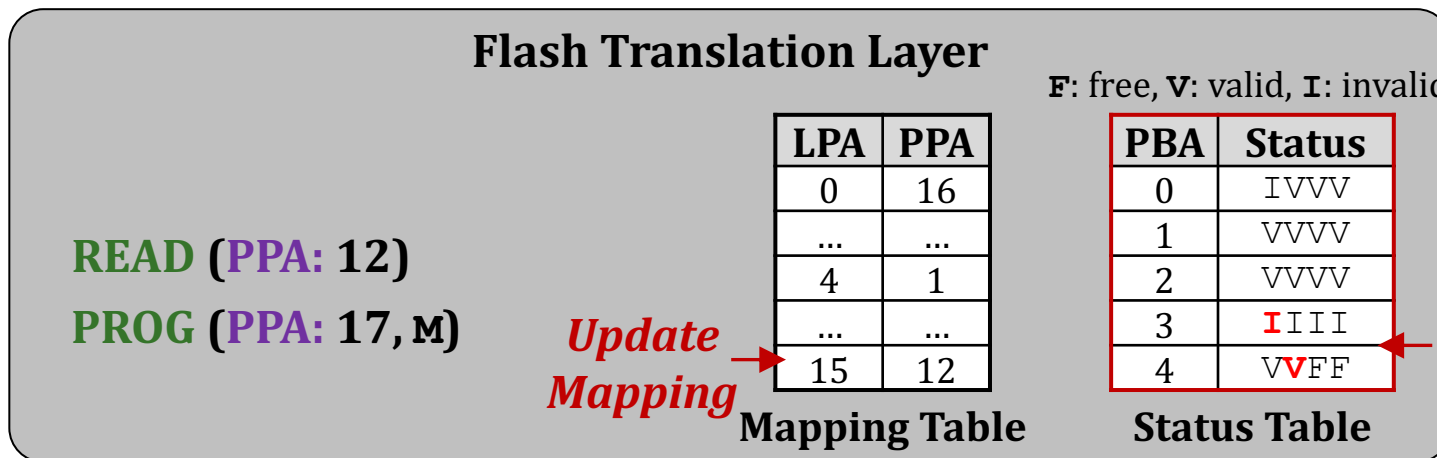


NAND Flash Chip (Single Plane)

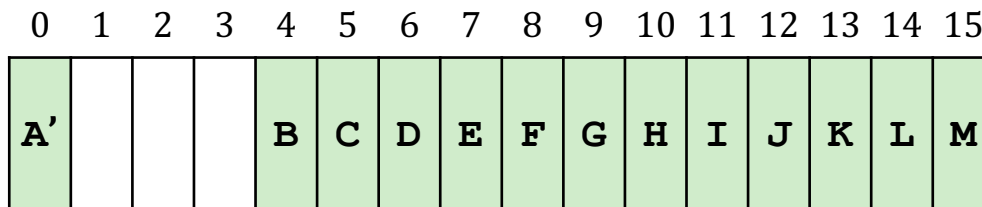
# Write Request Handling: Garbage Collection



SSD



# Write Request Handling: Garbage Collection



Host

SSD

## Flash Translation Layer

F: free, V: valid, I: invalid

READ (PPA: 12)  
 PROG (PPA: 17, M)

*Update Mapping* →

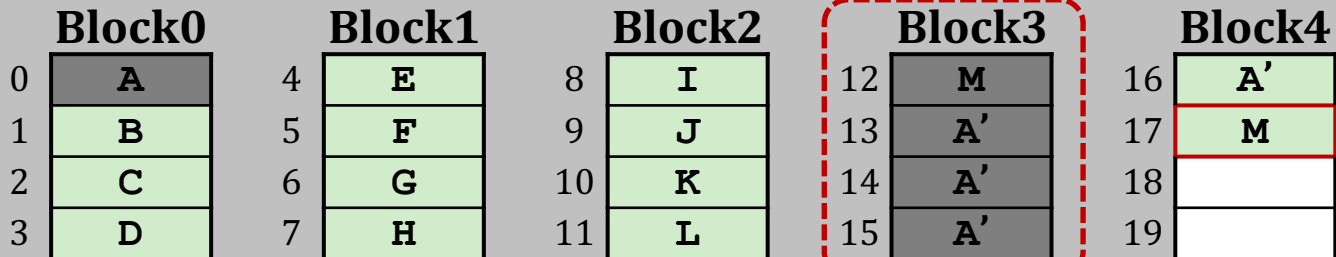
LPA	PPA
0	16
...	...
4	1
...	...
15	17

Mapping Table

PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	IIII
4	VVFF

Status Table

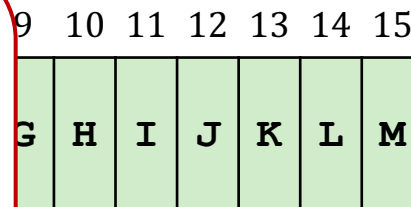
*Update Status* ←



NAND Flash Chip (Single Plane)

# Write Request Handling: Garbage Collection

- **Q:** How FTL knows PPA 12 (data **m**) is mapped to LPA 15?
  - Unless it maintains **P2L mappings**?
- **A:** P2L mapping is stored in each physical page's **OOB (Out-of-Band)** area
  - Cannot know P2L mapping **before reading the page**



**READ (PPA: 12)**  
**PROG (PPA: 17, M)**

*Update Mapping* →

**FTL Control Layer**  
 F: free, V: valid, I: invalid

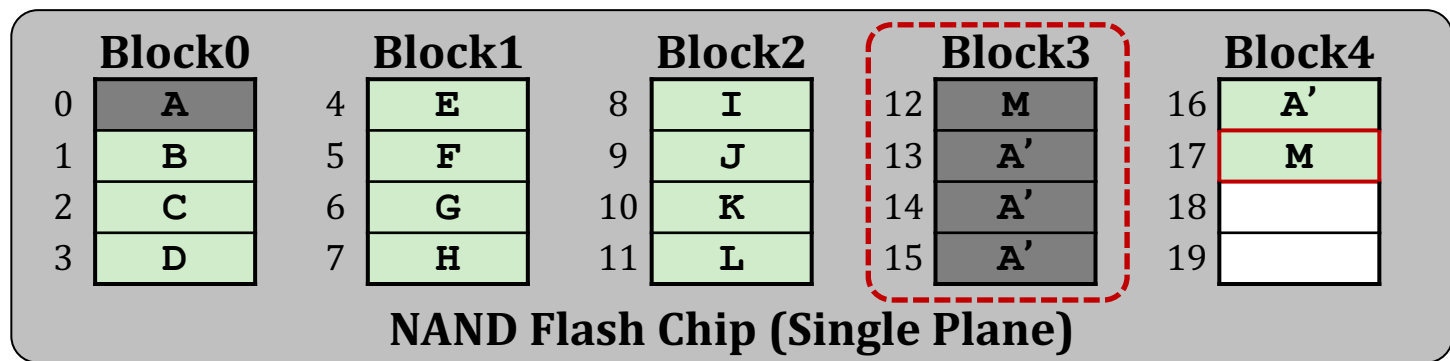
LPA	PPA
0	16
...	...
4	1
...	...
15	17

**Mapping Table**

PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	IIII
4	VVFF

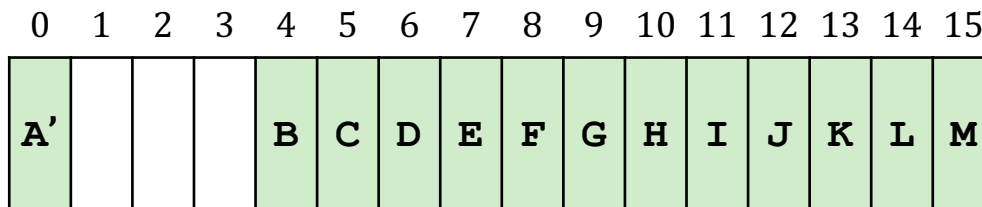
**Status Table**

*Update Status* ←





# Write Request Handling: Garbage Collection



Host

SSD

## Flash Translation Layer

F: free, V: valid, I: invalid

READ (PPA: 12)

PROG (PPA: 17, M)

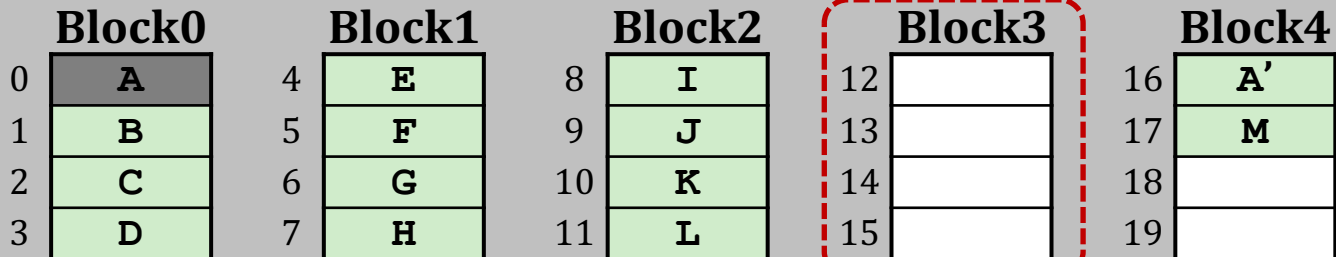
BERS (PBA: 3)

LPA	PPA
0	16
...	...
4	1
...	...
15	17

Mapping Table

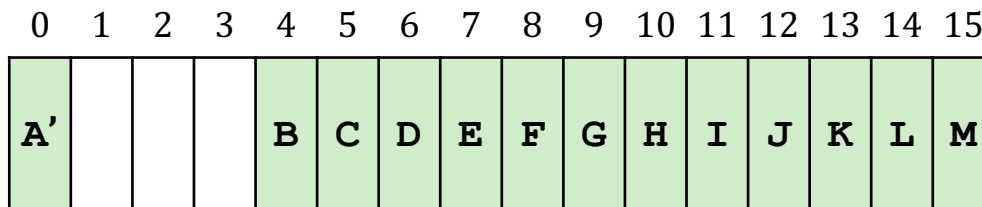
PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	IIII
4	VVFF

Status Table



NAND Flash Chip (Single Plane)

# Write Request Handling: Garbage Collection



Host

SSD

## Flash Translation Layer

F: free, V: valid, I: invalid

READ (PPA: 12)

PROG (PPA: 17, M)

BERS (PBA: 3)

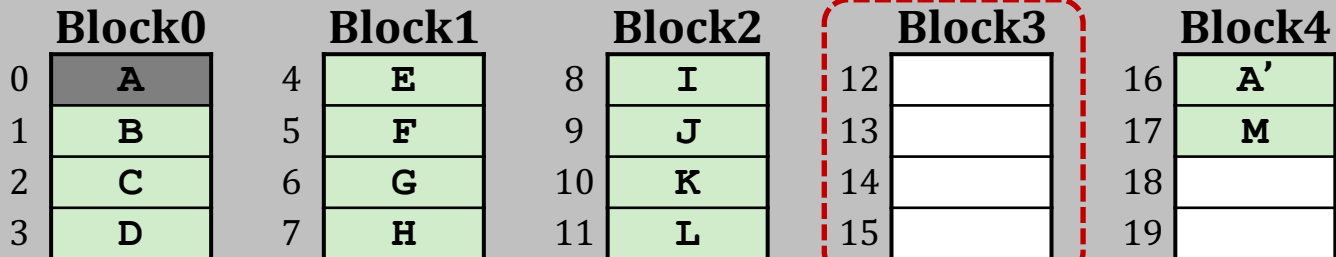
LPA	PPA
0	16
...	...
4	1
...	...
15	17

Mapping Table

PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	FFFF
4	VVFF

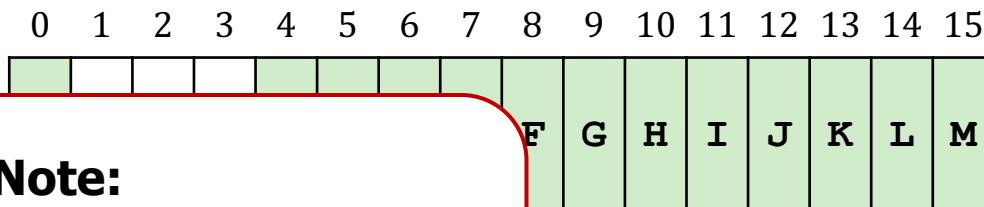
Status Table

*Update Status*



NAND Flash Chip (Single Plane)

# Write Request Handling: Garbage Collection



## Note:

- Block erasure (and status update) is done **just before** programming a new page to the block (i.e., **lazy erase**)
  - Due to the **open-block** problem

**F** (PPA: 12)  
**G** (PPA: 17, M)  
**ERS** (PBA: 3)

## Mapping Layer

F: free, V: valid, I: invalid

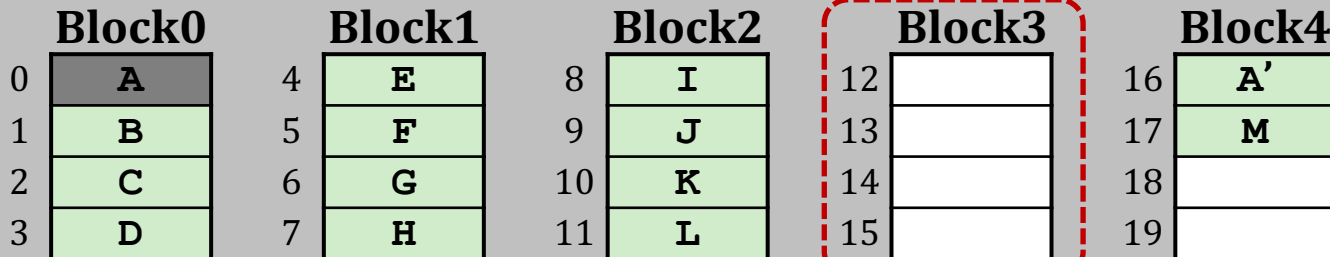
LPA	PPA
0	16
...	...
4	1
...	...
15	17

Mapping Table

PBA	Status
0	IVVV
1	VVVV
2	VVVV
3	<b>FFFF</b>
4	VVFF

Status Table

*Update Status*



NAND Flash Chip (Single Plane)

# Performance Issues

---

- Garbage collection **significantly affects** SSD performance
  - High latency: **Large block size** of modern NAND flash memory
    - Assume 1) a block contains **576** pages,  
2) only **5%** of the pages in the victim block are valid  
3)  $t_R = 100 \text{ us}$ ,  $t_{\text{PROG}} = 700 \text{ us}$ ,  $t_{\text{BERS}} = 5 \text{ ms}$ 
      - # of pages to copy =  $576 \times 0.05 = 28.8 \rightarrow 28$  pages
      - GC latency  $> 28 \times (t_R + t_{\text{PROG}}) + t_{\text{BERS}} = \mathbf{27,400 \text{ ms}}$
      - **Order(s) of magnitude larger** latency than  $t_R$  and  $t_{\text{PROG}}$
      - Copy operations are **the major contributor** (rather than  $t_{\text{BERS}}$ )
    - If FTL performs GC in an **atomic** manner,  
it **delays** user requests for a **significantly long time**
      - Long **tail latency** (performance fluctuation)
      - **Noisy neighbor**: a read-dominant workload's performance would be significantly affected when running with a write-intensive workload (+ performance fairness problem)

# Performance Issues: Mitigation

---

- **TRIM** (or **discard**) command
  - Informs FTL of **deletion/deallocation** of a logical block
  - Allows FTL to **skip copy** of obsolete (i.e., invalid) data
- **Background GC**: Exploits SSD idle time
  - Challenge: how to **accurately predict** SSD idle time
  - **Premature** GC: copied pages **could have been invalidated** by the host system
- **Progressive GC**: **Divide** GC process into subtasks
  - e.g., copying 28 pages  $\rightarrow$  (copying 1 page + **servicing user request**)  $\times$  28
  - Effective at decreasing tail latency

# P&S Modern SSDs

Address Translation & Garbage Collection

Dr. Jisung Park

Prof. Onur Mutlu

ETH Zürich

Spring 2021

19 May 2021