# P&S Modern SSDs

## Research Session 1:
## Data Sanitization and Read-Retry
## in Modern NAND Flash-Based SSDs

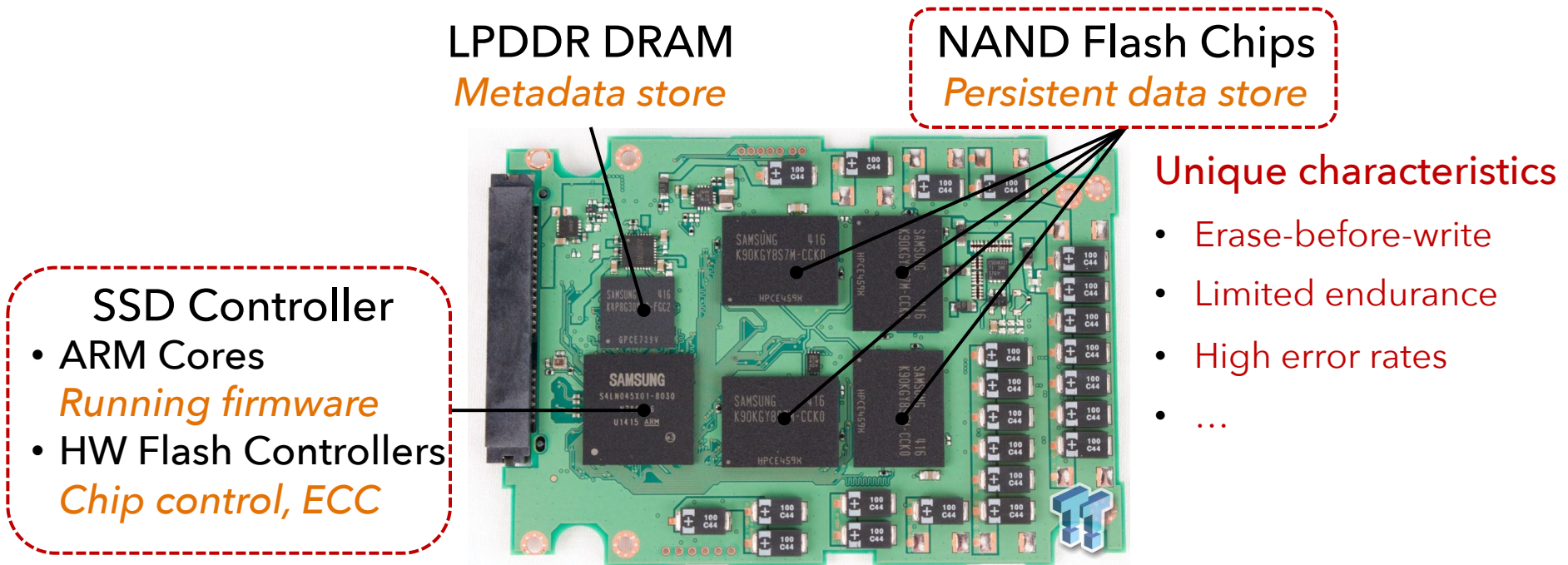Dr. Jisung Park

Prof. Onur Mutlu

ETH Zürich

Spring 2022

15 July 2022

# Outline

- NAND Flash Basics

- Read-Retry in Modern NAND Flash-Based SSDs

- Data Sanitization in Modern NAND Flash-Based SSDs
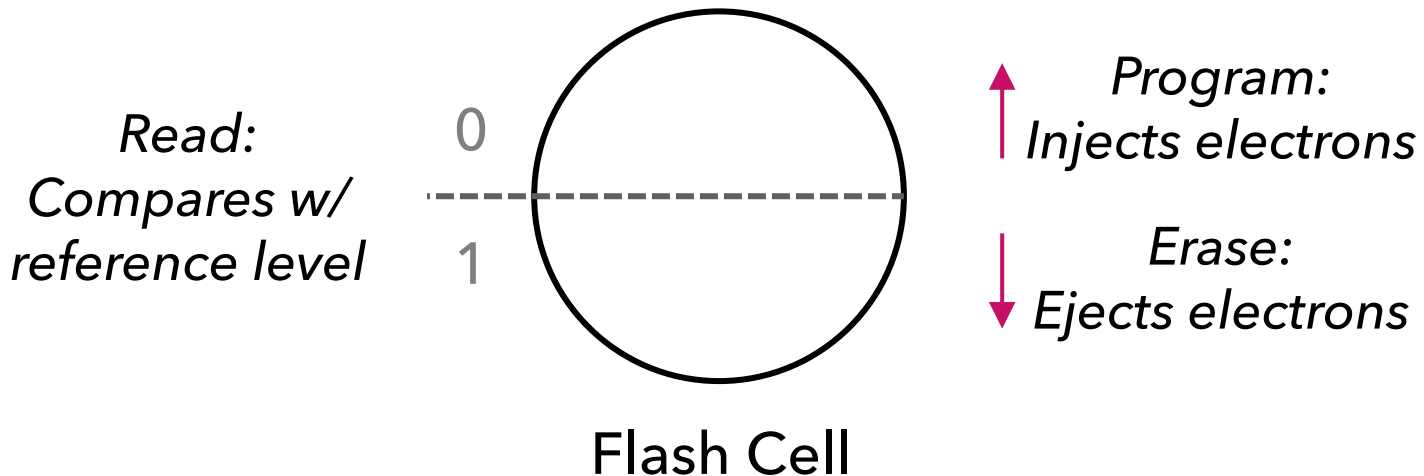
# NAND Flash-Based SSDs

- A complicated embedded system
    - Multiple cores, HW controllers, DRAM, and NAND flash chips

**LPDDR DRAM**
*Metadata store*

**NAND Flash Chips**
*Persistent data store*

**Unique characteristics**

- Erase-before-write
- Limited endurance
- High error rates
- ...

**SSD Controller**
- ARM Cores
  *Running firmware*
- HW Flash Controllers
  *Chip control, ECC*



Samsung PM853T 960GB Enterprise SSD (from https://www.tweaktown.com/reviews/6695/samsung-pm853t-960gb-enterprise-ssd-review/index.html)
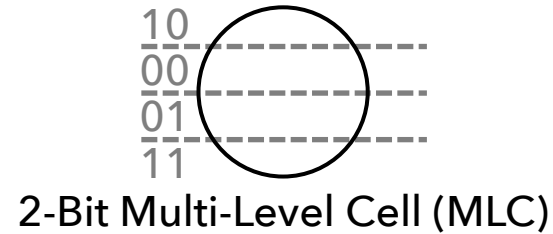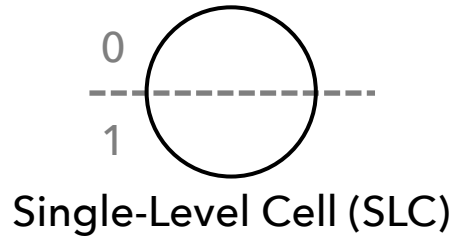
# Flash Cell

- Stores data using its threshold voltage ($V_{TH}$) level
  - Dictated by the amount of electrons (i.e., charge) in the cell
  - The more the electrons, the higher the $V_{TH}$ level

*Read:*
*Compares w/*
*reference level*

0

1

*Program:*
*Injects electrons*
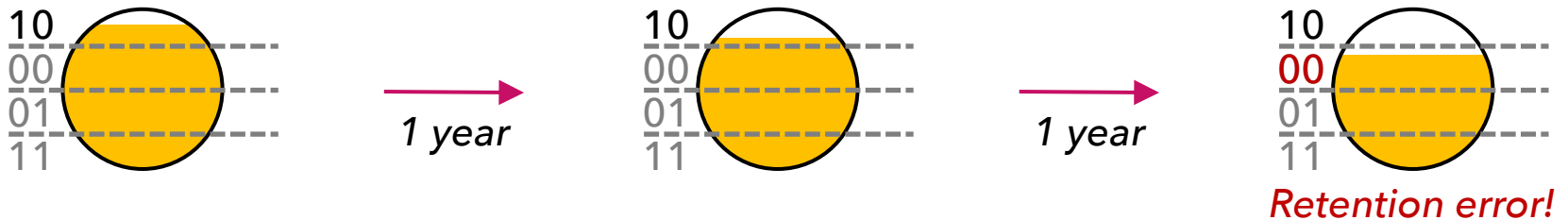
*Erase:*
*Ejects electrons*

Flash Cell

1. Can change $V_{TH}$ (=charge) in a non-volatile manner

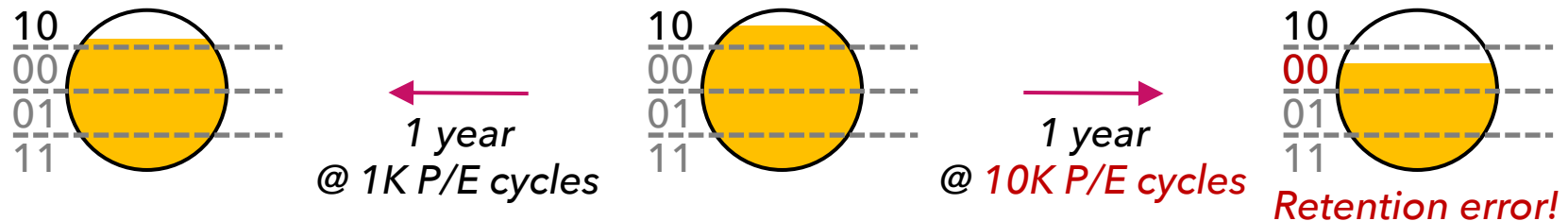2. Encodes bit data with different $V_{TH}$ ranges

# Flash Cell Characteristics

- Multi-leveling: A single flash cell can store multiple bits



Single-Level Cell (SLC)                2-Bit Multi-Level Cell (MLC)

- Retention loss: A cell leaks electrons over time



*1 year*          *1 year*          *Retention error!*

- Limited lifetime: A cell wears out after P/E cycling



*1 year*
*@ 1K P/E cycles*

*1 year*
*@ 10K P/E cycles*          *Retention error!*

# Page and Block

- A number of flash cells operate in parallel

*Cells in the same wordline (WL)*
*are concurrently read/programmed*

$WL_0$ — (1)(0)(0)(1) ... (0)(1)(1)(0)
$WL_1$
$WL_2$

<100,000 cells

$WL_{125}$
$WL_{126}$
$WL_{127}$

Block (128 WLs)

*Cells in the same block*
*are concurrently erased*

Long latency,
but high bandwidth

No overwrite of WLs
before erasing the block:
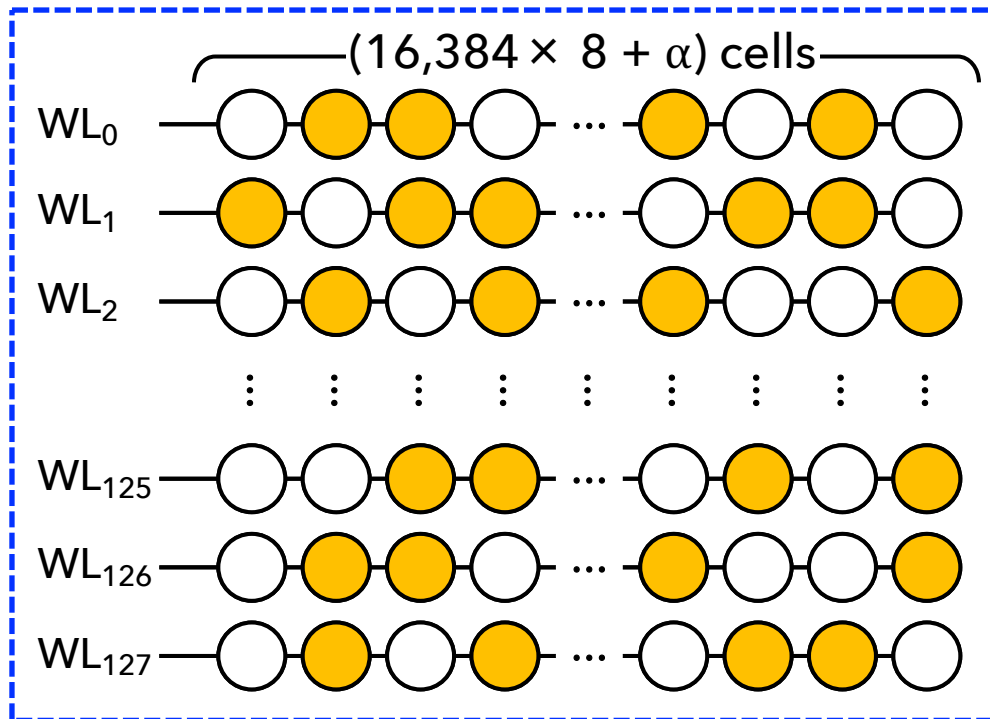*Erase-before-write*

# Page and Block

- A number of flash cells operate in parallel

*Cells in the same wordline (WL)*
*are concurrently read/programmed*

$(16{,}384 \times 8 + \alpha)$ cells

$WL_0$

$WL_1$

$WL_2$

$WL_{125}$

$WL_{126}$

$WL_{127}$

Block (128 WLs)

*Cells in the same block*
*are concurrently erased*

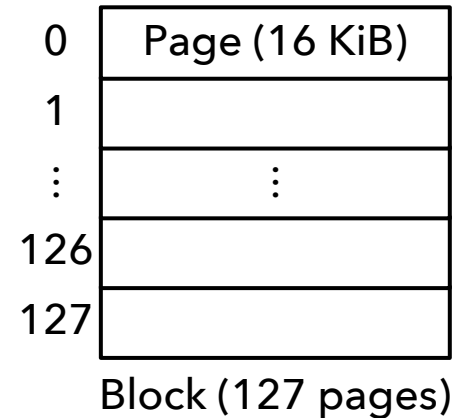| | Page (16 KiB) |
|---|---|
| 0 | |
| 1 | |
| ⋮ | ⋮ |
| 126 | |
| 127 | |

Block (127 pages)

# Page and Block

- A number of flash cells operate in parallel

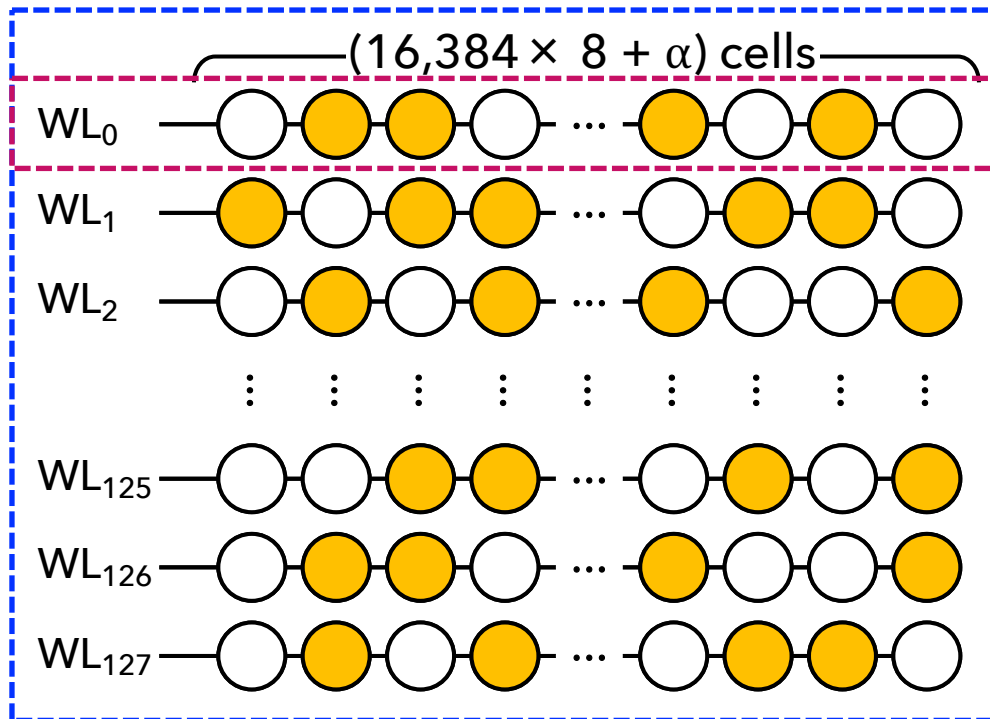*Cells in the same wordline (WL)*
*are concurrently read/programmed*



$(16,384 \times 8 + \alpha)$ cells

$WL_0$ ...
$WL_1$ ...
$WL_2$ ...
⋮
$WL_{125}$ ...
$WL_{126}$ ...
$WL_{127}$ ...

Block (128 WLs)

*Cells in the same block*
*are concurrently erased*

*Share the same cells (WL)*

| 0 | Page (16 KiB) |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| ⋮ | ⋮ |
| 378 | |
| 379 | |
| 380 | |
| 381 | |

Triple-Level Cell (TLC)
Block ($127 \times 3$ pages)

# Pipelined & Adaptive Read-Retry

## Reducing Solid-State Drive Read Latency by Optimizing Read-Retry

Jisung Park[1]   Myungsuk Kim[2,3]   Myoungjun Chun[2]   Lois Orosa[1]   Jihong Kim[2]   Onur Mutlu[1]

[1]ETH Zürich
Switzerland

[2]Seoul National University
Republic of Korea

[3]Kyungpook National University
Republic of Korea

The 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'21)

# Problem: Long, Non-Deterministic Latency

*Read request (4 KiB)* →

NAND Flash-Based SSD

← *Requested data*
*Expected latency: 100 μs*
*Actual latency: > 1 ms*

Degrades the quality of service
of read-intensive, latency-sensitive applications

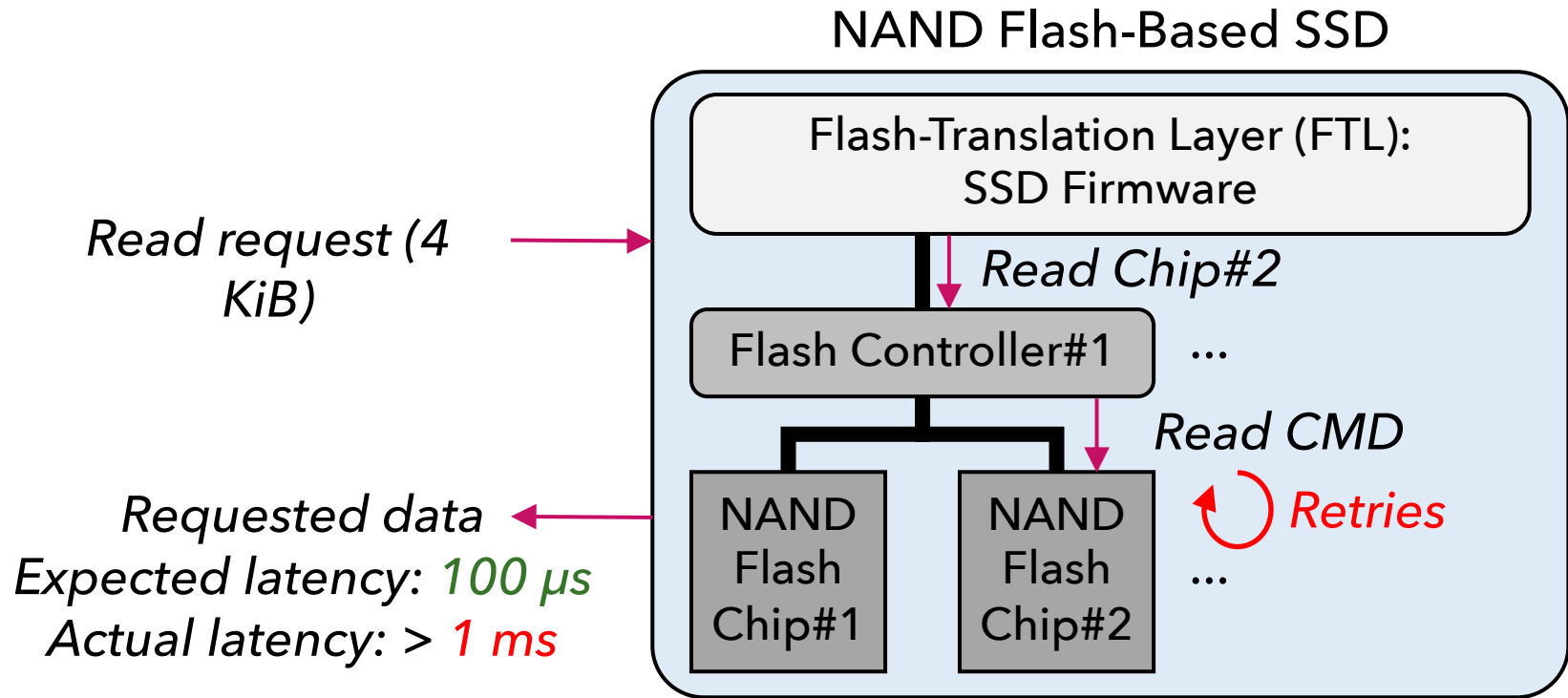# Problem: Long, Non-Deterministic Latency

NAND Flash-Based SSD



Read request (4 KiB)

Requested data
Expected latency: *100 µs*
Actual latency: > *1 ms*

Flash-Translation Layer (FTL): SSD Firmware

*Read Chip#2*

Flash Controller#1 ...

*Read CMD*

*Retries*

NAND Flash Chip#1

NAND Flash Chip#2 ...

## Internal Read-Retry Operations

# Errors in NAND Flash Memory

- NAND flash memory stores data by using cells' $V_{TH}$ levels

Read-Reference Voltage $V_{REF}$

Number of Cells

1
Erased
State

0
Programmed
State

Cell's $V_{TH}$ Level

<$V_{TH}$ Distribution of an SLC Page>

$V_{REF}$ — $\dfrac{0}{1}$

$V_{REF}$ — $\dfrac{0}{1}$

# Errors in NAND Flash Memory

- Various sources shift and widen programmed $V_{TH}$ states
  - Retention loss, program interference, read disturbance, etc.

Read–Reference Voltage $V_{REF}$

Number of Cells

*Interference*

*Retention loss*

1
Erased
State

0
Programmed
State

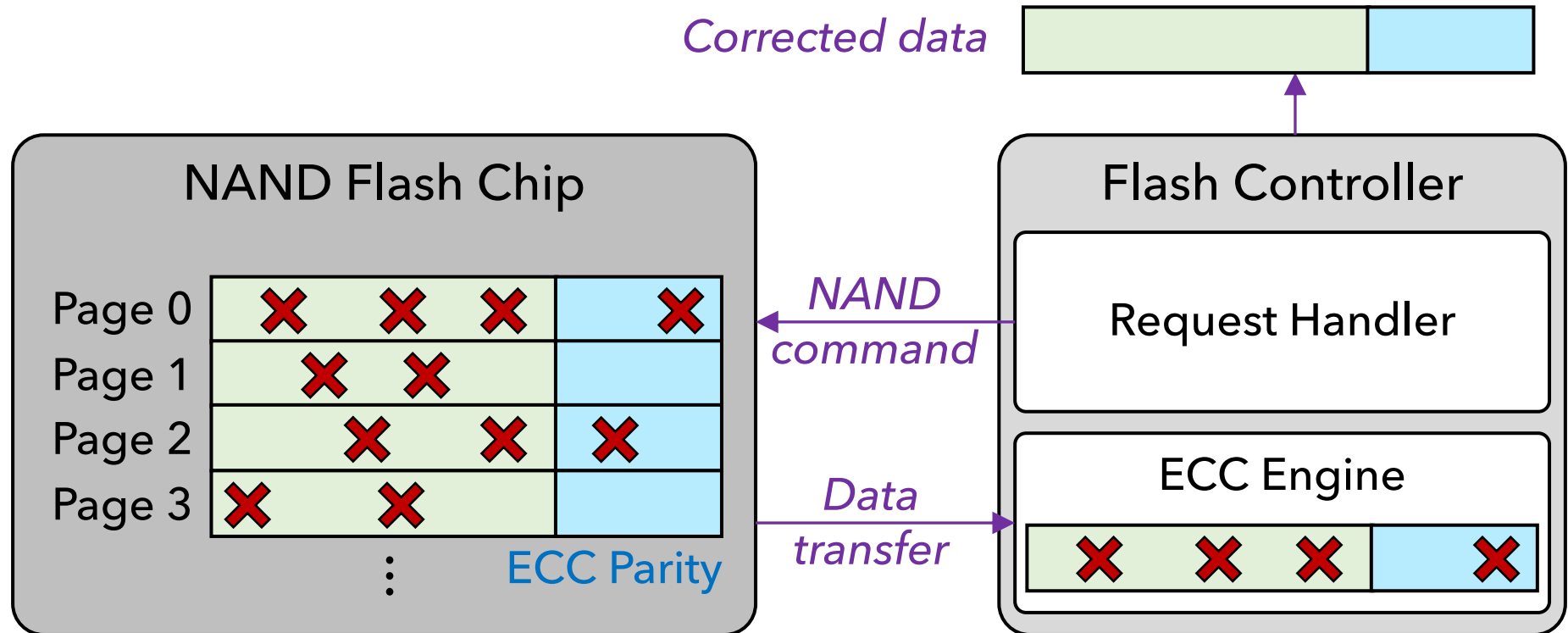Cell's $V_{TH}$ Level

Errors

<V_{TH} Distribution of an SLC Page>

# Error-Correcting Codes (ECC)

- Store redundant information (ECC parity)
  - To detect and correct row bit errors

# Error-Correcting Codes (ECC)

- Store redundant information (ECC parity)
  - To detect and correct row bit errors

# Errors in Modern NAND Flash Memory

- High row bit-error rates (RBER) in MLC NAND flash memory
  - Narrow margin b/w adjacent $V_{TH}$ states



$V_{REF}$

1
Erased (E)

0
Programmed

$V_{TH}$

# of cells

<$V_{TH}$ Distribution of an SLC Page>

$V_{REF0}$   $V_{REF1}$   $V_{REF2}$   $V_{REF3}$   $V_{REF4}$   $V_{REF5}$   $V_{REF6}$

CSB

MSB  LSB

111
E

110
P1

100
P2

000
P3

010
P4

011
P5

001
P6

101
P7

$V_{TH}$

# of cells

$V_{TH}$ Margin

<$V_{TH}$ Distribution of a TLC Page>

# Errors in Modern NAND Flash Memory

- High row bit-error rates (RBER) in MLC NAND flash memory
  - Narrow margin b/w adjacent $V_{TH}$ states

Strong ECC: Corrects ~80 bit errors per 1-KiB data

Not Scalable: Area, power, latency, …

What if RBER > ECC Capability?



$V_{REF0}$  $V_{REF1}$  $V_{REF2}$  $V_{REF3}$  $V_{REF4}$  $V_{REF5}$  $V_{REF6}$

CSB

MSB LSB

# of cells

111  110  100  000  010  011  001  101

E  P1  P2  P3  P4  P5  P6  P7  $V_{TH}$

$V_{TH}$ Margin

<$V_{TH}$ Distribution of a TLC Page>

# Read-Retry Operation

- Reads the page again with adjusted $V_{REF}$ values

# Read-Retry Operation

- Reads the page again with adjusted $V_{REF}$ values



Read-retry: Adjusting $V_{REF}$ values

$V_{REF(x-2)}$   $V_{REF(x-1)}$   $V_{REFx}$   $V_{REF(x+1)}$

$V'_{REF(x-2)}$   $V'_{REF(x-1)}$   $V'_{REFx}$   $V'_{REF(x+1)}$

Number of Cells

P(x-1) State    Px State    P(x+1) State

Erroneous cells

Cell's $V_{TH}$ Level

Read using properly-adjusted $V_{REF}$ values
→ Decreases # of raw bit errors
to be lower than the ECC capability

# Read-Retry: Performance Overhead

`tDMA:` *Data transfer*

`tR:` *Page sensing*          `tECC:` *ECC decoding*

READ A

$N_{ERR} = 32 < $ ECC capability $C_{ECC} = 72$

$\leftarrow$ `tREAD` $\rightarrow$

# Read-Retry: Performance Overhead



Read-retry increases the read latency
almost linearly with the number of retry steps

# P&AR$^2$: Outline

• Read-Retry in Modern NAND Flash-Based SSDs

• PR$^2$: Pipelined Read-Retry

• AR$^2$: Adaptive Read-Retry

• Evaluation Results

# Pipelined Read-Retry (PR$^2$): Key Idea



READ A — $N_{ERR} = 232 >$ ECC capability $C_{ECC} = 72$

RR1 — $N_{ERR} = 173$

RR2 — $N_{ERR} = 118$

RR($N - 1$) — $N_{ERR} = 87$

RR$N$ — $N_{ERR} = 23$

tR  tDMA  tECC

$\leftarrow$ tREAD $\rightarrow$  $\leftarrow$ tRETRY $\rightarrow$

$= N \times ($tR+tDMA+tECC$)$

In common cases, multiple (up to 25) retry steps occur

# Pipelined Read-Retry (PR$^2$): Key Idea



In common cases, multiple (up to 25) retry steps occur

→ Speculatively starts the next retry step

# Pipelined Read-Retry (PR$^2$): Key Idea



Removes tDMA & tECC
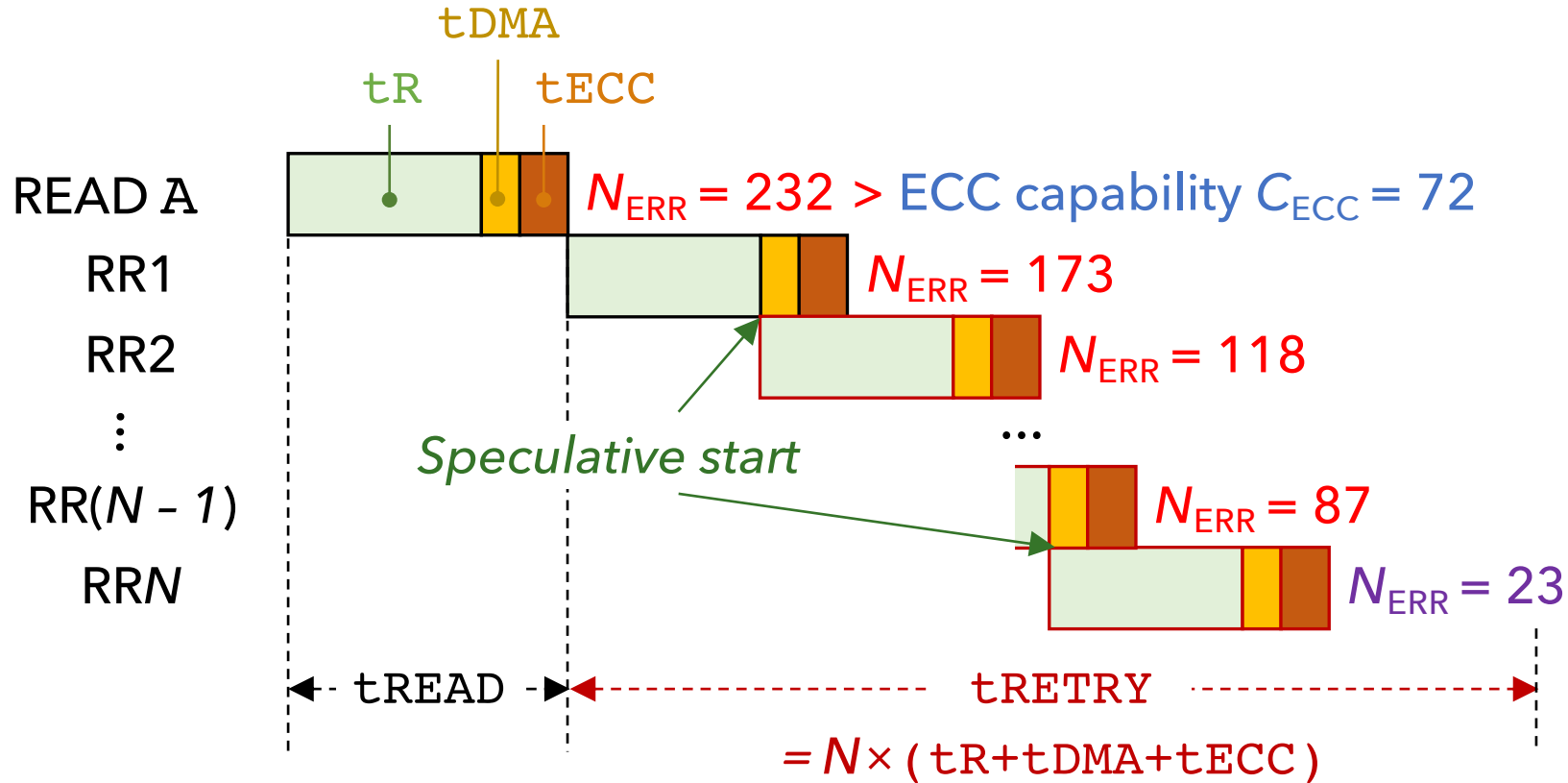(~30% of each retry step) from the critical path

# P&AR$^2$: Outline

- Read-Retry in Modern NAND Flash-Based SSDs

- PR$^2$: Pipelined Read-Retry

- AR$^2$: Adaptive Read-Retry

- Evaluation Results

# Adaptive Read-Retry (AR$^2$): Key Idea



READ A: $N_{ERR}$ = 232 > ECC capability $C_{ECC}$ = 72

RR1: $N_{ERR}$ = 173

RR2: $N_{ERR}$ = 118

RR($N$ – 1): $N_{ERR}$ = 87

RR$N$: $N_{ERR}$ = 23

tREAD

tRETRY = $N \times$ tR+tDMA+tECC

tR   tDMA   tECC

# Adaptive Read-Retry (AR$^2$): Key Idea



**READ A** — $N_{ERR} = 232 >$ ECC capability $C_{ECC} = 72$

**RR1** — $N_{ERR} = 173$

**RR2** — $N_{ERR} = 118$

**RR($N - 1$)** — $N_{ERR} = 87$

**RR$N$** — $N_{ERR} = 23$

tDMA
tR   tECC

tREAD   tRETRY
$= N \times$ tR+tDMA+tECC

A large ECC margin
in the final retry step when read-retry succeeds
→ Can we leverage this ECC (reliability) margin?

# Adaptive Read-Retry (AR$^2$): Key Idea



$tDMA$

$tR$        $tECC$

READ A        $N_{ERR}$ = 232 > ECC capability $C_{ECC}$ = 72

RR1        $N_{ERR}$ = 173 + $e_1$

RR2        $N_{ERR}$ = 118 + $e_2$

⋮        ...

$tR' = \alpha \times tR$        $N_{ERR}$ = 87 + $e_{N-1}$

RR($N$ – 1)

RR$N$        $N_{ERR}$ = 23 + $e_N$

← $tREAD$ → ← — — — — — — $tRETRY$ — — — — — →

= $N \times tR + tDMA + tECC$

Trading the large ECC margin to reduce $tR$

# Adaptive Read-Retry (AR$^2$): Key Idea



$\tt tDMA$

$\tt tECC$

$\tt tR$

READ A $\quad$ $N_{ERR} = 232 >$ ECC capability $C_{ECC} = 72$

RR1 $\quad$ $N_{ERR} = 173 + e_1$

RR2 $\quad$ $N_{ERR} = 118 + e_2$

$\vdots$ $\quad$ ...

RR($N - 1$) $\quad$ $N_{ERR} = 87 + e_{N-1}$

RR$N$ $\quad$ $N_{ERR} = 23 + e_N$

$\tt tR' = \alpha \times tR$

$\leftarrow$ $\tt tREAD$ $\rightarrow$ $\leftarrow$ $\tt tRETRY'$ $\rightarrow$ Latency Reduction

$\cancel{= N \times \tt tR + tDMA + tECC}$

$= N \times \alpha \times \tt tR + tDMA + tECC$

Trading the large ECC margin to reduce $\tt tR$
$\rightarrow$ Further reduction in the read-retry latency

# Adaptive Read-Retry (AR²): Key Idea



READ A — $N_{ERR} = 232 >$ ECC capability $C_{ECC} = 72$

RR1 — $N_{ERR} = 173 + e_1$

RR2 — $N_{ERR} = 118 + e_2$

...

RR($N - 1$) — $N_{ERR} = 87 + e_{N-1}$
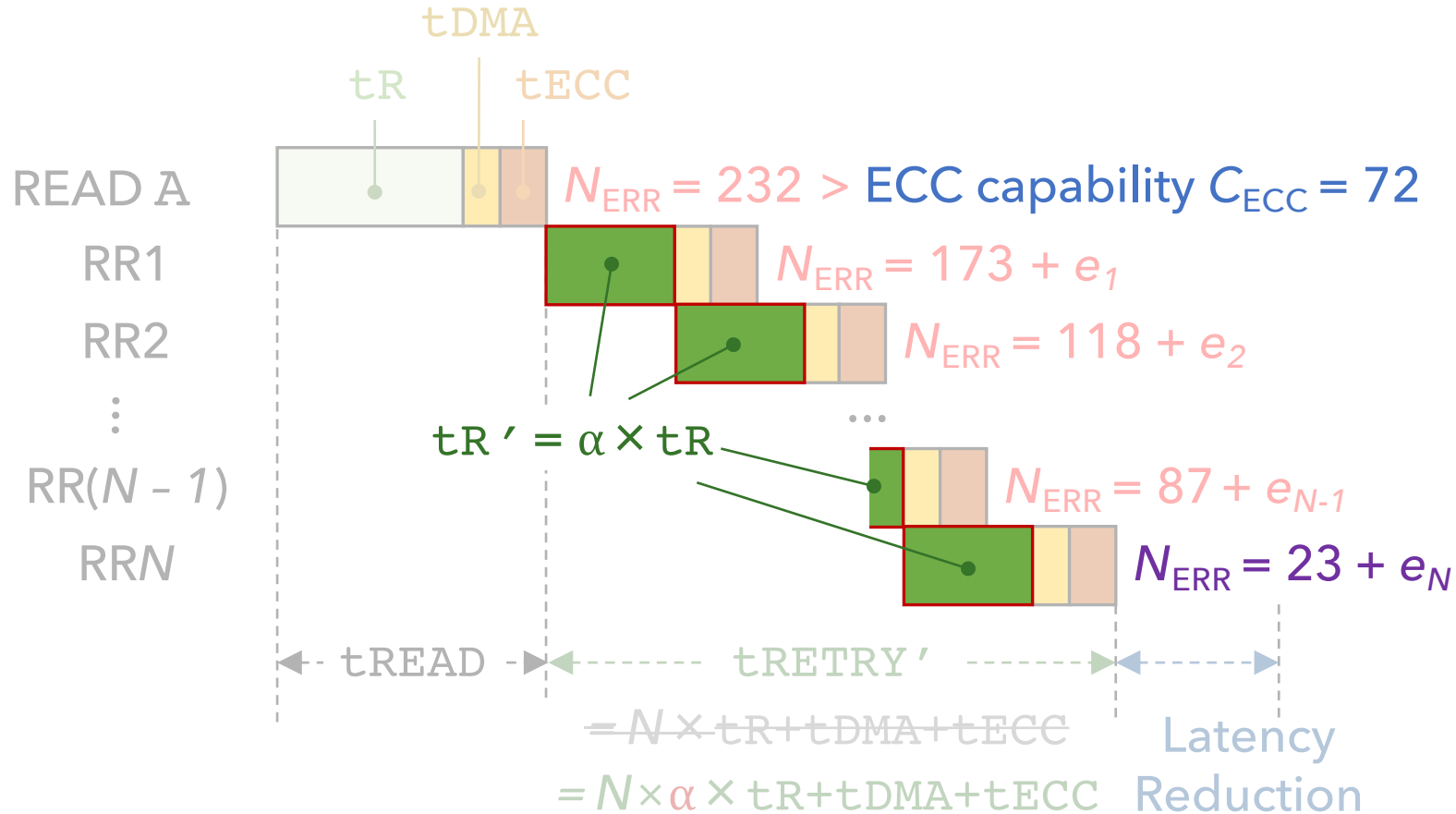
RRN — $N_{ERR} = 23 + e_N$

$\texttt{tR'} = \alpha \times \texttt{tR}$

$\texttt{tREAD}$ — $\texttt{tRETRY'}$ — Latency Reduction

$= N \times \texttt{tR} + \texttt{tDMA} + \texttt{tECC}$

$= N \times \alpha \times \texttt{tR} + \texttt{tDMA} + \texttt{tECC}$

## AR² reduces tR in every retry step

# Adaptive Read-Retry (AR$^2$): Key Idea



$N_{ERR}$ = 232 > ECC capability $C_{ECC}$ = 72

**READ A**

**RR1**    $N_{ERR}$ = 173 + $e_1$

**RR2**    $N_{ERR}$ = 118 + $e_2$

$\text{tR}' = \alpha \times \text{tR}$

**RR($N$ – 1)**    $N_{ERR}$ = 87 + $e_{N-1}$

**RR$N$**    $N_{ERR}$ = 23 + $e_N$

$\text{tR}$   $\text{tDMA}$   $\text{tECC}$

← tREAD → ← --------- tRETRY' --------- →  ← Latency

~~$= N \times \text{tR} + \text{tDMA} + \text{tECC}$~~

$= N \times \alpha \times \text{tR} + \text{tDMA} + \text{tECC}$   Reduction

## AR$^2$ reduces tR in every retry step
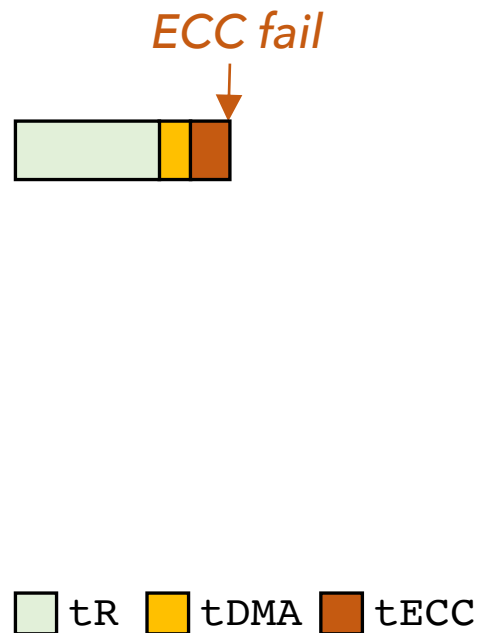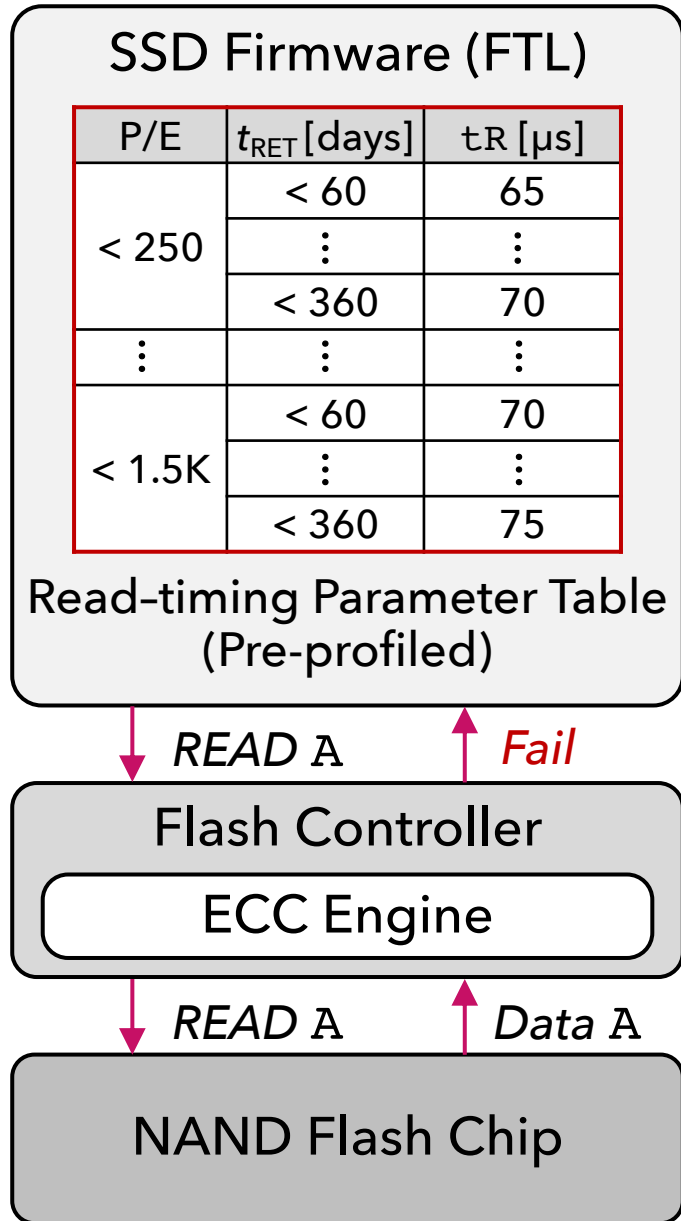## ensuring $N_{ERR} < C_{ECC}$ in the final retry step
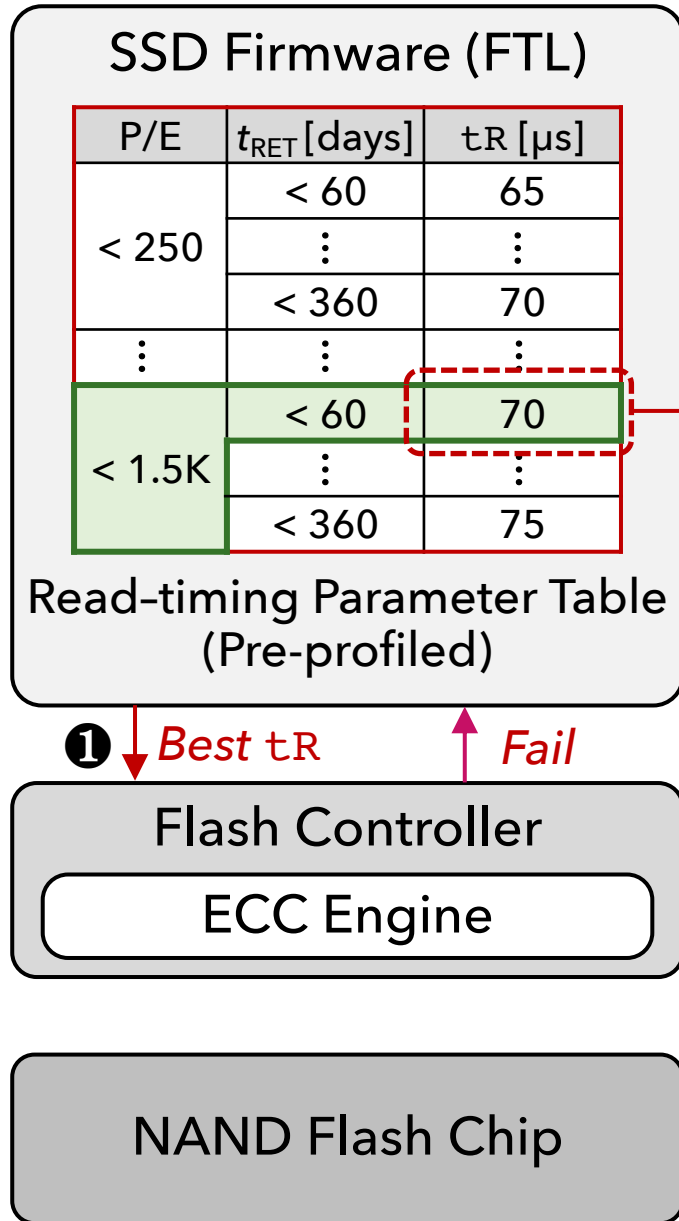
# Real-Device Characterization

- 160 real 48-layer TLC NAND flash chips

- Observation 1: A large ECC margin in the final retry step even under worst-case operating conditions
  - At most 40 errors per KiB under 1–year retention time @ 2K program and erase (P/E) cycles
  - Use of near-optimal $V_{REF}$ in the final retry step

- Observation 2: A large reliability margin incorporated in read-timing parameters
  - 25% `tR` reduction → At most 23 additional errors
  - Worst-case-based design due to process variations

## $AR^2$ can easily work in commodity NAND flash chips w/ at least 25% `tR` reduction
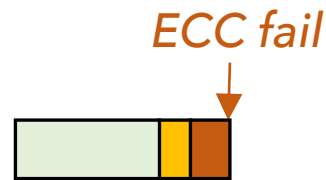
# P&AR$^2$: Design

## SSD Firmware (FTL)

| P/E | $t_{RET}$ [days] | tR [μs] |
|---|---|---|
| < 250 | < 60 | 65 |
| | ⋮ | ⋮ |
| | < 360 | 70 |
| ⋮ | ⋮ | ⋮ |
| < 1.5K | < 60 | 70 |
| | ⋮ | ⋮ |
| | < 360 | 75 |

**Read–timing Parameter Table (Pre-profiled)**

↓ *READ* A          ↑ *Fail*

## Flash Controller

### ECC Engine

↓ *READ* A          ↑ *Data* A

## NAND Flash Chip

*ECC fail*

□ tR  ■ tDMA  ■ tECC

# P&AR²: Design

## SSD Firmware (FTL)

| P/E | $t_{RET}$ [days] | tR [µs] |
|---|---|---|
| < 250 | < 60 | 65 |
|  | ⋮ | ⋮ |
|  | < 360 | 70 |
| ⋮ | ⋮ | ⋮ |
| < 1.5K | < 60 | 70 |
|  | ⋮ | ⋮ |
|  | < 360 | 75 |

Read–timing Parameter Table
(Pre-profiled)

*Best* tR *for the current operating conditions*

❶ *Best* tR          *Fail*

## Flash Controller

### ECC Engine

## NAND Flash Chip

*ECC fail*

☐ tR   ■ tDMA   ■ tECC

# P&AR²: Design

## SSD Firmware (FTL)

Read–timing Parameter Table (Pre-profiled)

| P/E | $t_{RET}$ [days] | tR [µs] |
|---|---|---|
| | < 60 | 65 |
| < 250 | ⋮ | ⋮ |
| | < 360 | 70 |
| ⋮ | ⋮ | ⋮ |
| | < 60 | 70 |
| < 1.5K | ⋮ | ⋮ |
| | < 360 | 75 |

*Best* tR *for the current operating conditions*

❶ *Best* tR    *Page* A

## Flash Controller

### ECC Engine

❷ *Read-retry w/ best* tR *+ speculative start*

## NAND Flash Chip

*ECC fail*

*ECC success*

...

☐ tR  ☐ tDMA  ☐ tECC  ☐ Best tR

# P&AR²: Design

SSD Firmware (FTL)

| PEC | $t_{RET}$ [days] | tR [μs] |
|-----|------------------|---------|
|     | < 60             | 65      |

## Key Takeaway

Strong ECC: to avoid read-retry as much as possible

→ Can provide high reliability margin
when read-retry occurs

→ Can be used to reduce the read-retry latency

② Read-retry w/ best tR + speculative start

NAND Flash Chip

☐ tR  ☐ tDMA  ☐ tECC  ☐ Best tR

37

# P&AR$^2$: Outline

- Read-Retry in Modern NAND Flash-Based SSDs

- PR$^2$: Pipelined Read-Retry

- AR$^2$: Adaptive Read-Retry

- Evaluation Results

# Evaluation Results

- Simulation using MQSim [Tavakkol+, FAST18] and 12 real-world workloads

- Our proposal improves SSD response time by
    - Up to 51% (35% on average) compared to a high-end SSD w/o read-retry mitigation
    - Up to 32% (17% on average) compared to a state-of-the-art read-retry mitigation technique [Shim+, MICRO19]

# Outline

- NAND Flash Basics

- Read-Retry in Modern NAND Flash-Based SSDs

- Data Sanitization in Modern NAND Flash-Based SSDs

# Access Control-Based Data Sanitization

## Evanesco: Architectural Support for Efficient Data Sanitization in Modern Flash-Based Storage Systems

**Myungsuk Kim***
morssola75@davinci.snu.ac.kr
Seoul National University

**Jisung Park***
jisung.park@inf.ethz.ch
ETH Zürich & Seoul National University

**Geonhee Cho**
ghcho@davinci.snu.ac.kr
Seoul National University

**Yoona Kim**
yoonakim@davinci.snu.ac.kr
Seoul National University

**Lois Orosa**
lois.orosa@inf.ethz.ch
ETH Zürich

**Onur Mutlu**
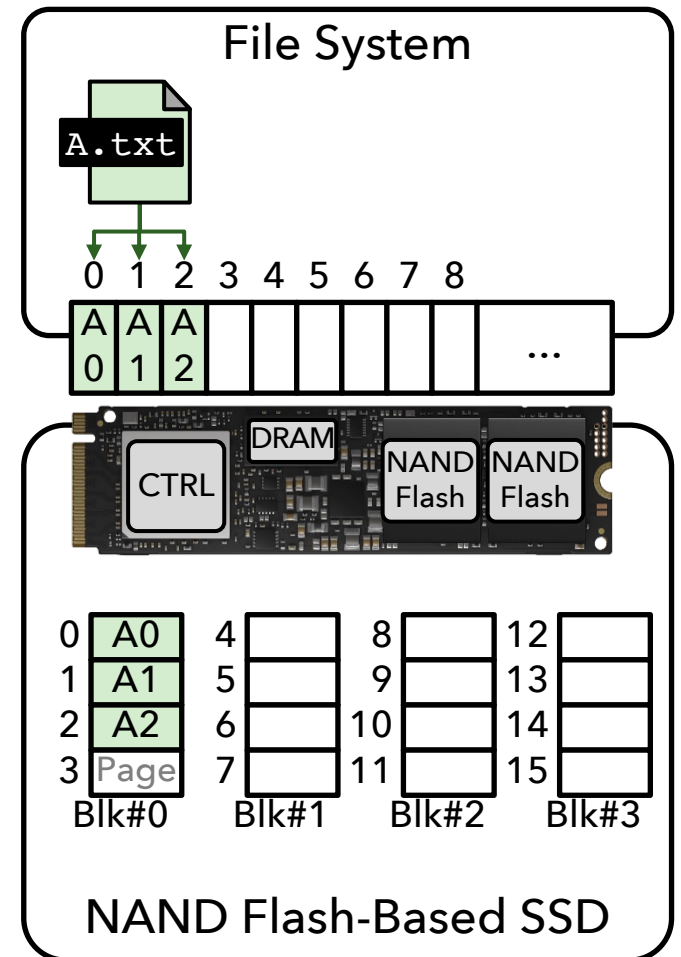omutlu@gmail.com
ETH Zürich

**Jihong Kim**
jihong@davinci.snu.ac.kr
Seoul National University

*The first two authors contributed equally to this research.

The 25th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'20)
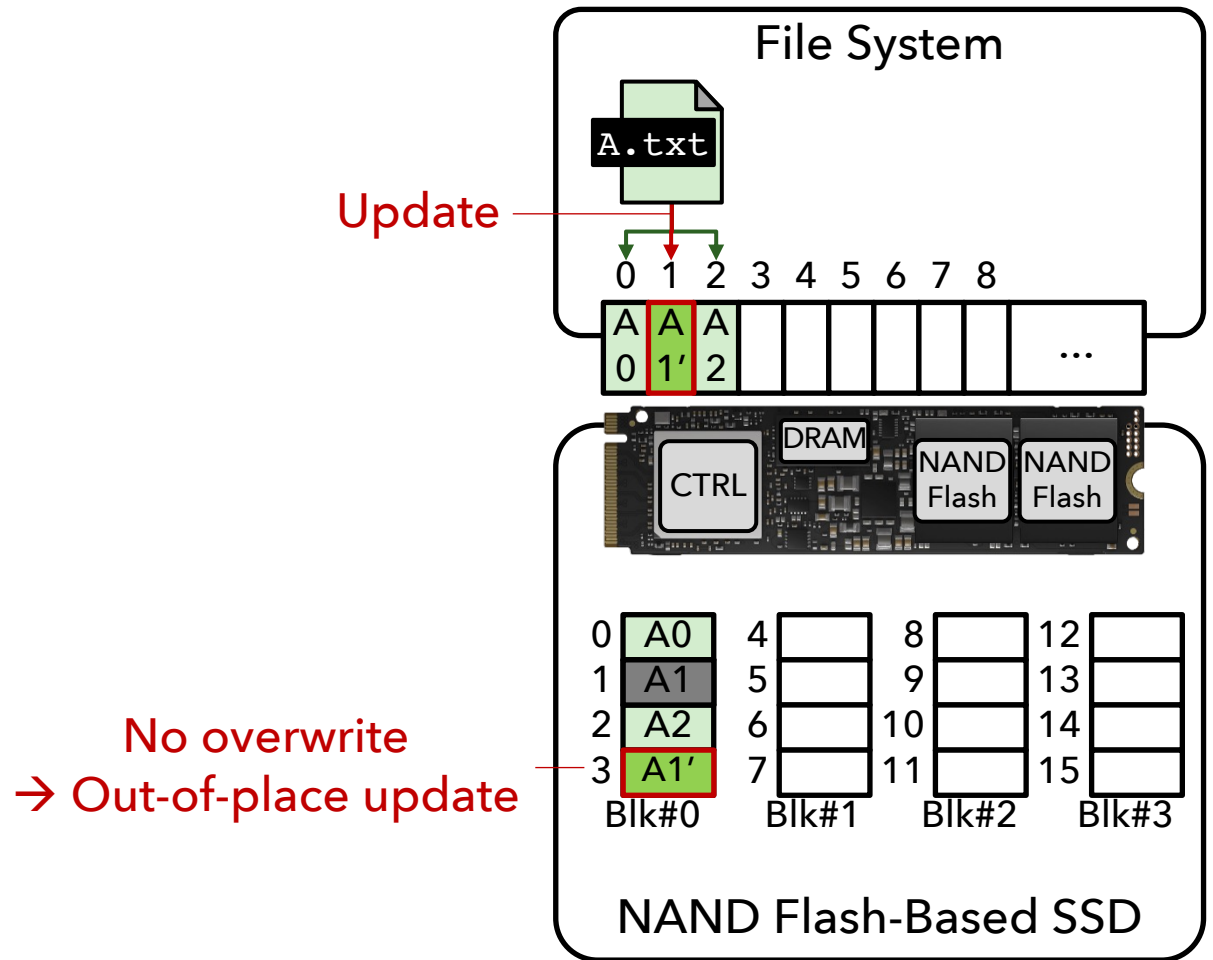
# Data-Remanence Problem in SSDs

- Deleted data can remain in SSDs for indefinite time



File System

A.txt

0 1 2 3 4 5 6 7 8

| A0 | A1 | A2 | | | | | | ... |

DRAM
CTRL
NAND Flash
NAND Flash

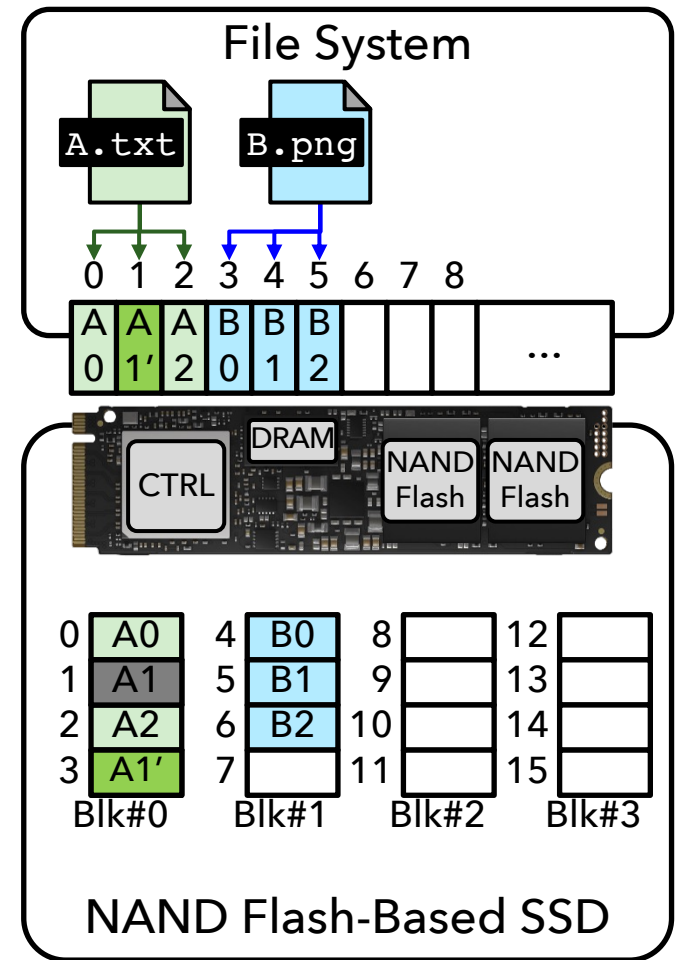| 0 | A0 | 4 | | 8 | | 12 | |
| 1 | A1 | 5 | | 9 | | 13 | |
| 2 | A2 | 6 | | 10 | | 14 | |
| 3 | Page | 7 | | 11 | | 15 | |
| Blk#0 | | Blk#1 | | Blk#2 | | Blk#3 | |

NAND Flash-Based SSD

# Data-Remanence Problem in SSDs

- Deleted data can remain in SSDs for indefinite time

# Data-Remanence Problem in SSDs

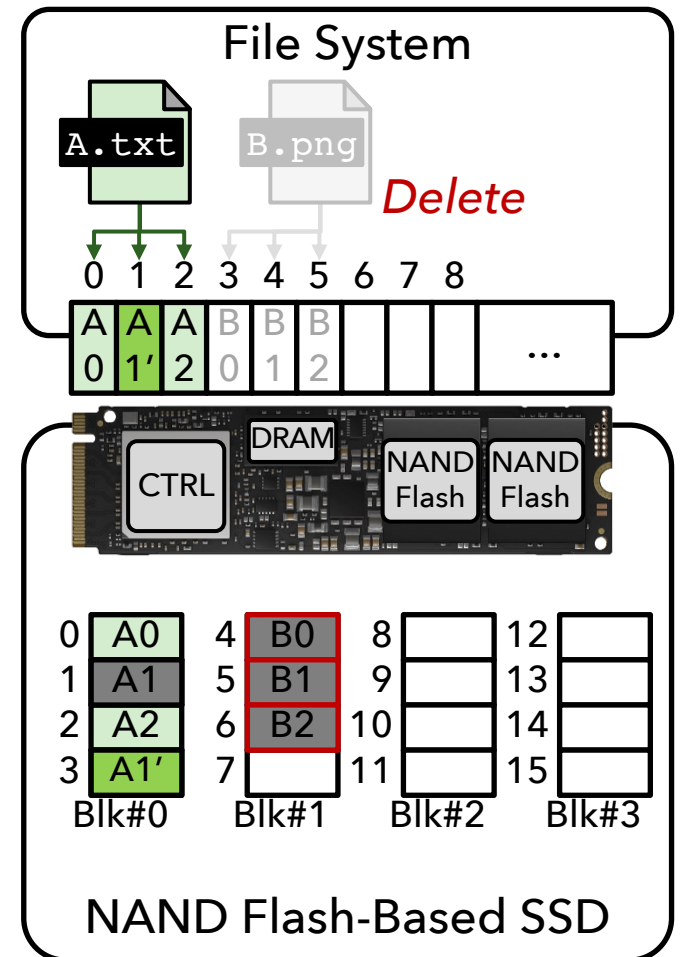- Deleted data can remain in SSDs for indefinite time

# Data-Remanence Problem in SSDs

- Deleted data can remain in SSDs for indefinite time

Q: When is a page erased?

A: Only in garbage collection
= when running out of free pages
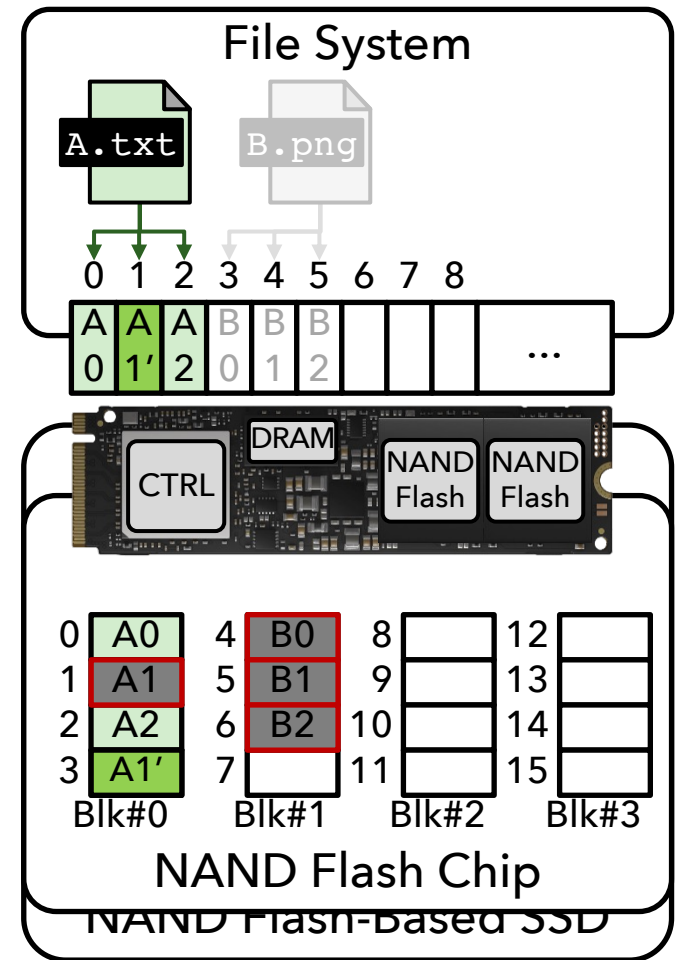
# Data-Recovery Attack

- System requirement: Obsolete data must be inaccessible

1. Detach the SSD

2. Detach the chip

3. Direct access to the chip

| A0 | A1 | A2 | B1 |

4. Run forensic tools

A.txt    B.png

ADVERSARY

## File System

A.txt    B.png

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| A 0 | A 1' | A 2 | B 0 | B 1 | B 2 | | | ... |

DRAM

CTRL    NAND Flash    NAND Flash

| 0 | A0 | 4 | B0 | 8 | | 12 | |
|---|----|---|----|---|---|----|---|
| 1 | A1 | 5 | B1 | 9 | | 13 | |
| 2 | A2 | 6 | B2 | 10 | | 14 | |
| 3 | A1' | 7 | | 11 | | 15 | |

Blk#0    Blk#1    Blk#2    Blk#3

NAND Flash Chip

NAND Flash-Based SSD

# Existing Solutions

- Why not immediately erase an invalid page?
  - Erase unit: a block (> 1,000 pages)



$$\texttt{tCOPY} = N_{\text{Valid}} \times (\texttt{tREAD} + \texttt{tPROG})$$

Target Page

*Copy $N_{\text{Valid}}$ Pages*

Free
Valid
Invalid

$N$ Pages
(>1,000)

$N_{\text{Valid}}$

*Wasted*

Target Block

Free Block

# Existing Solutions

- Why not immediately erase an invalid page?
  - Erase unit: a block (> 1,000 pages)



Target Page

$$\texttt{tCOPY}=N_{Valid}\times(\texttt{tREAD+tPROG})$$

Copy $N_{Valid}$ Pages

Free
Valid
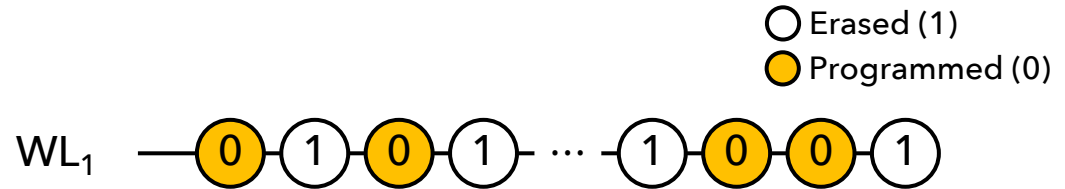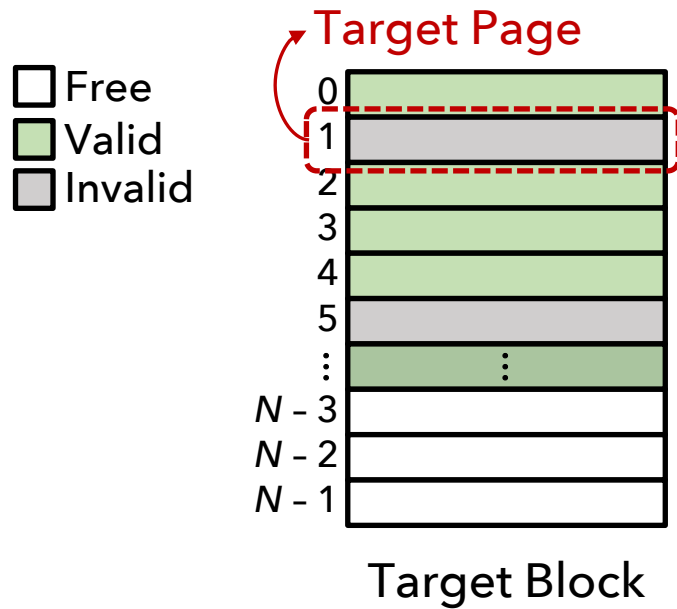Invalid

$N$ Pages
(>1,000)

Wasted

Target Block

Free Block

Immediate block erasure causes
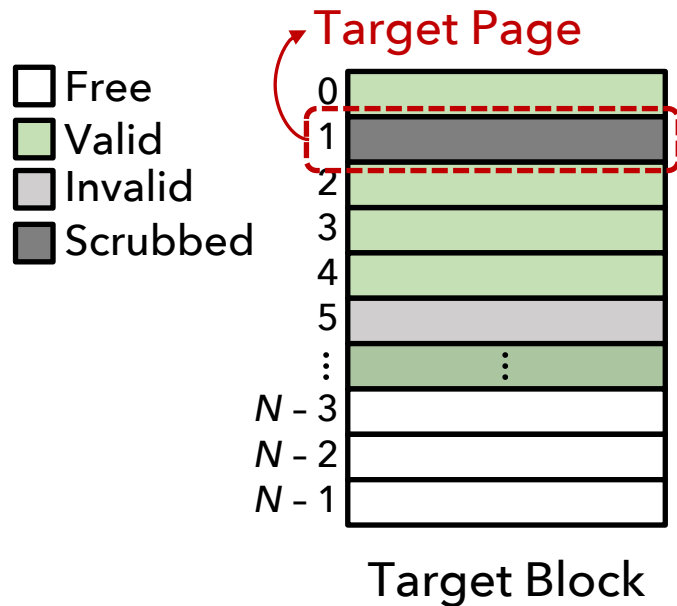prohibitive performance and lifetime overhead

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
  - Reprograms all the flash cells storing an invalid page

Target Page

Free
Valid
Invalid

0
1
2
3
4
5
...
N – 3
N – 2
N – 1

Target Block

○ Erased (1)
● Programmed (0)

$WL_1$ — 0 1 0 1 ... 1 0 0 1

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
  - Reprograms all the flash cells storing an invalid page

**Target Page**

Free
Valid
Invalid
Scrubbed

0
1
2
3
4
5
⋮
N – 3
N – 2
N – 1

**Target Block**

Erased (1)
Programmed (0)

$WL_1$ — 0 0 0 0 … 0 0 0 0
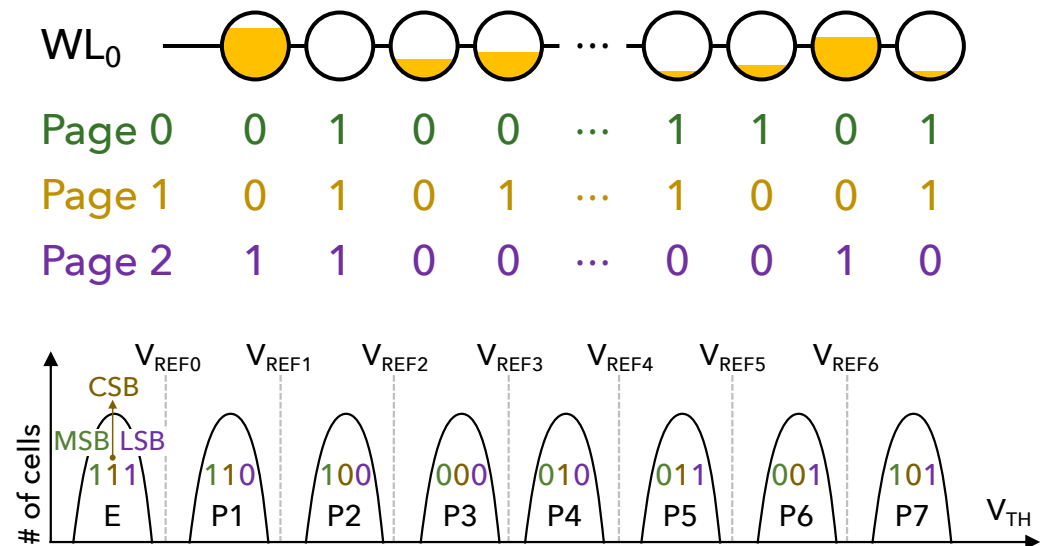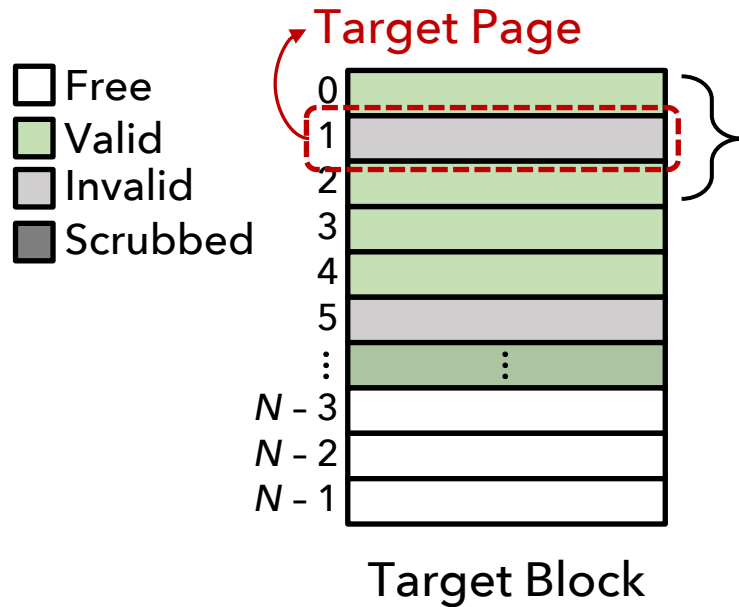
*Scrubbing: Reprogram the WL (to all '0's)*

*Overwriting is infeasible usually,
except when the data is all '0's*

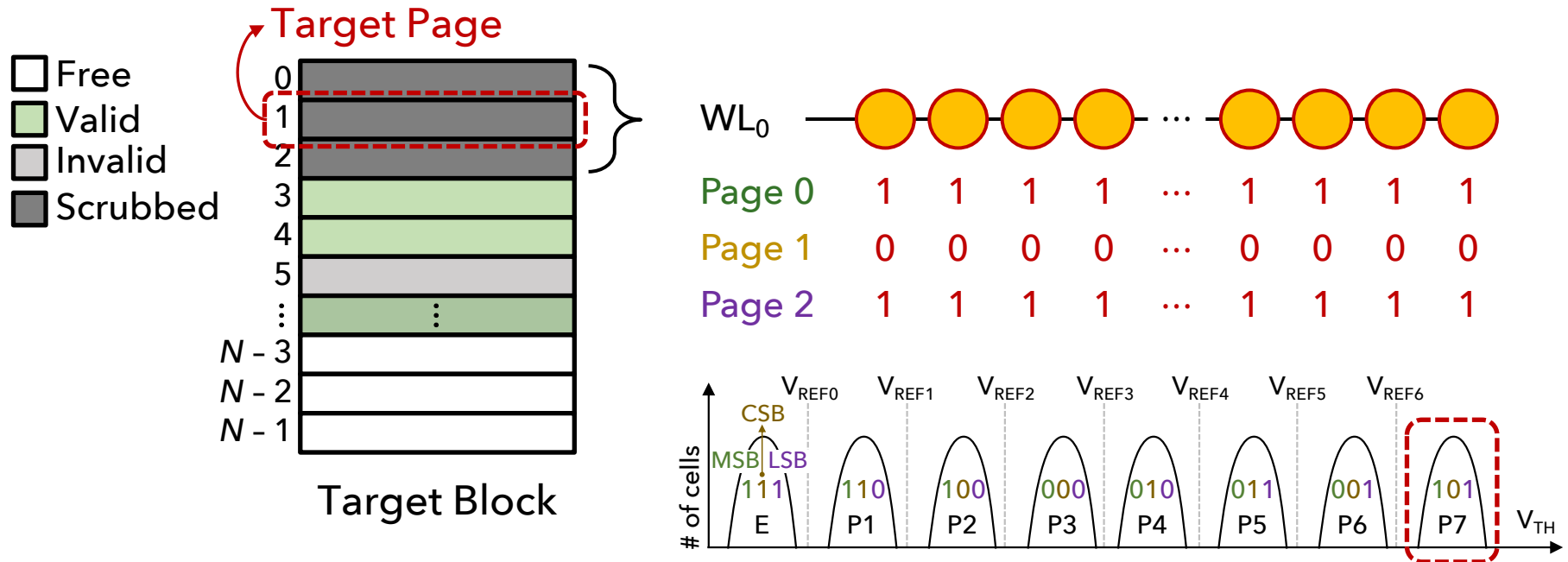*Destroys the page data w/o block erasure*

# Problem solved?

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
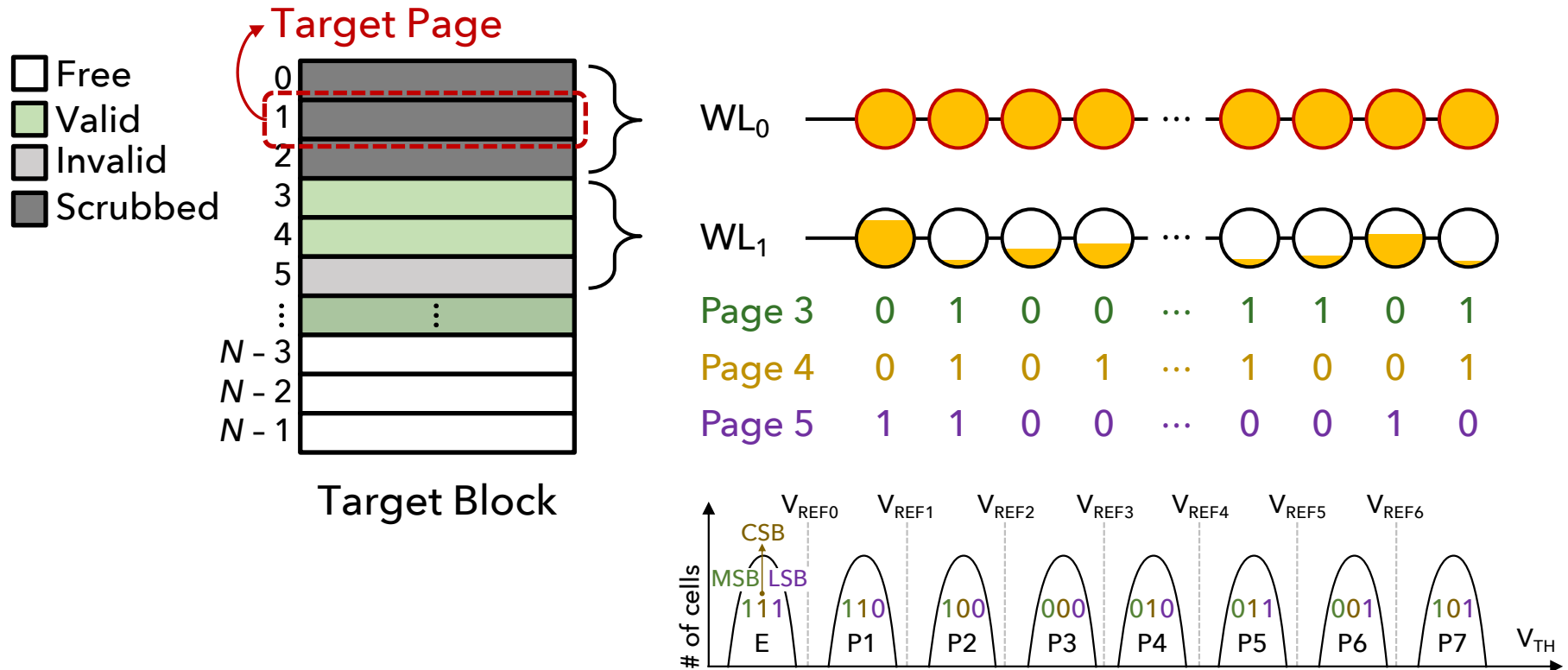    - Problem1: MLC NAND flash stores multiple pages in a WL

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
  - Problem1: MLC NAND flash stores multiple pages in a WL



Target Page

Free
Valid
Invalid
Scrubbed

Target Block

$WL_0$

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Page 0 | 1 | 1 | 1 | 1 | ⋯ | 1 | 1 | 1 | 1 |
| Page 1 | 0 | 0 | 0 | 0 | ⋯ | 0 | 0 | 0 | 0 |
| Page 2 | 1 | 1 | 1 | 1 | ⋯ | 1 | 1 | 1 | 1 |

$V_{REF0}$  $V_{REF1}$  $V_{REF2}$  $V_{REF3}$  $V_{REF4}$  $V_{REF5}$  $V_{REF6}$

CSB
MSB LSB
# of cells

111  110  100  000  010  011  001  101
E    P1   P2   P3   P4   P5   P6   P7   $V_{TH}$

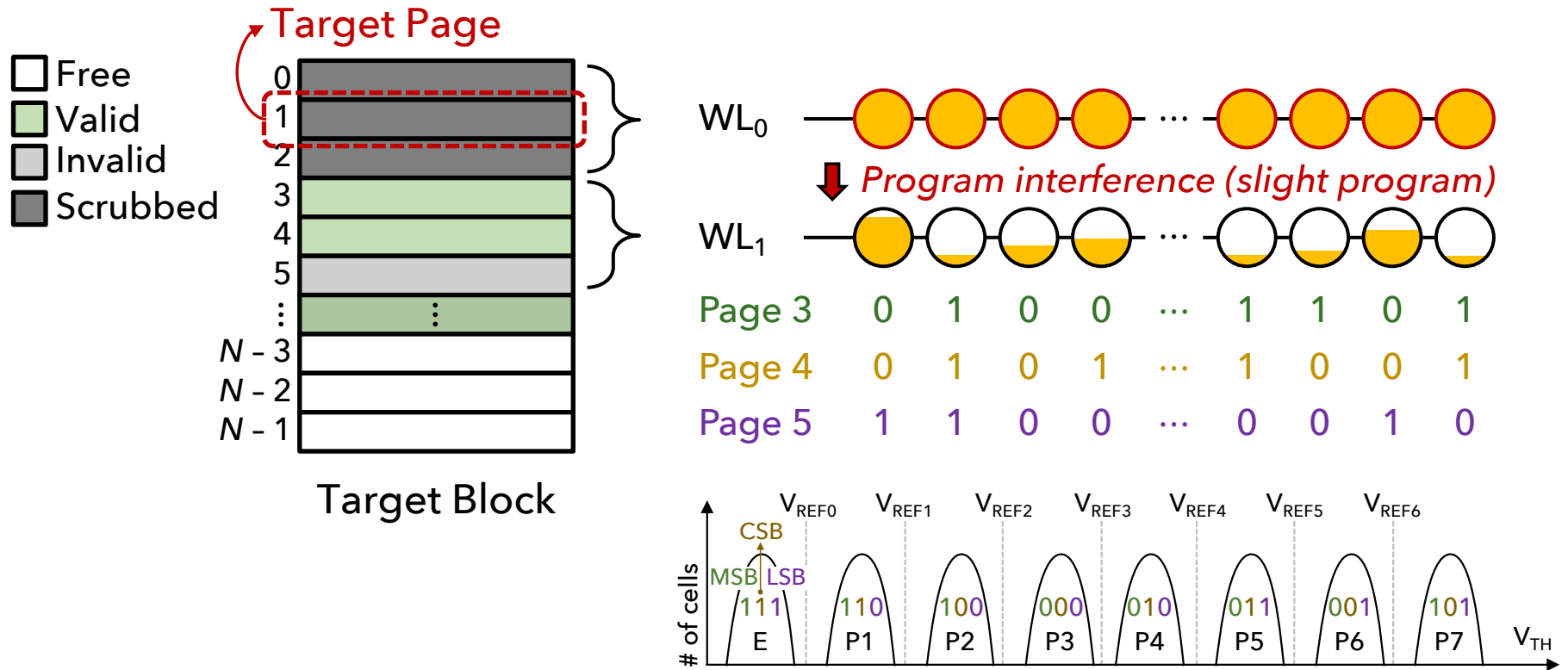## Scrubbing in MLC NAND flash memory → Destroys other valid pages → Copy overheads

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
  - Problem 1: MLC NAND flash stores multiple pages in a WL
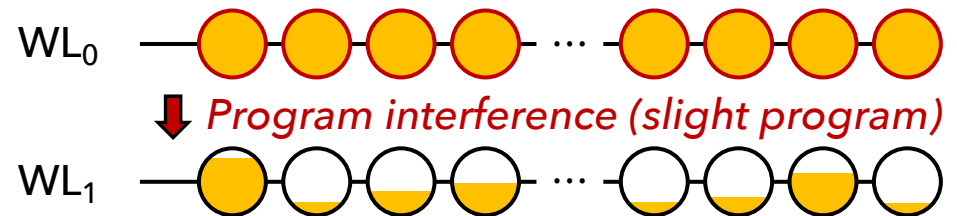  - Problem 2: Program interference



Target Page

Free
Valid
Invalid
Scrubbed

Target Block

$WL_0$

$WL_1$

| Page 3 | 0 | 1 | 0 | 0 | ⋯ | 1 | 1 | 0 | 1 |
| Page 4 | 0 | 1 | 0 | 1 | ⋯ | 1 | 0 | 0 | 1 |
| Page 5 | 1 | 1 | 0 | 0 | ⋯ | 0 | 0 | 1 | 0 |

$V_{REF0}$  $V_{REF1}$  $V_{REF2}$  $V_{REF3}$  $V_{REF4}$  $V_{REF5}$  $V_{REF6}$

CSB

MSB LSB

# of cells

111   110   100   000   010   011   001   101
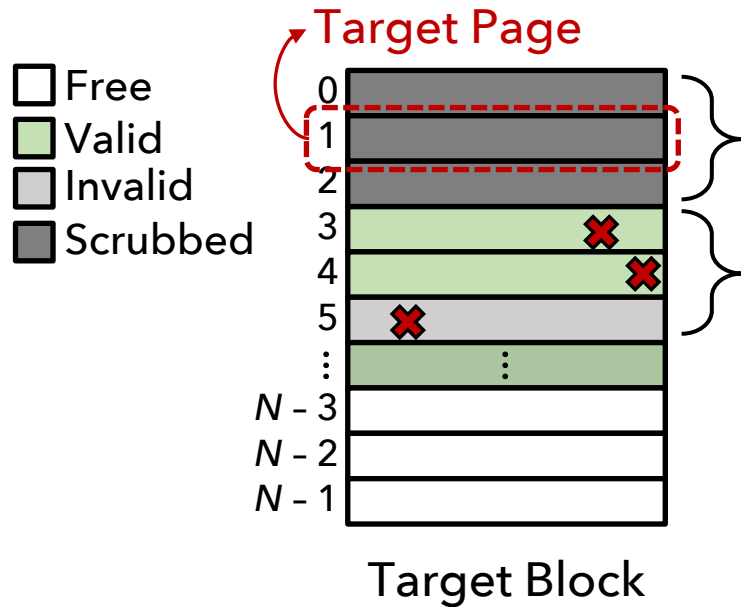E     P1    P2    P3    P4    P5    P6    P7

$V_{TH}$

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
    - Problem 1: MLC NAND flash stores multiple pages in a WL
    - Problem 2: Program interference
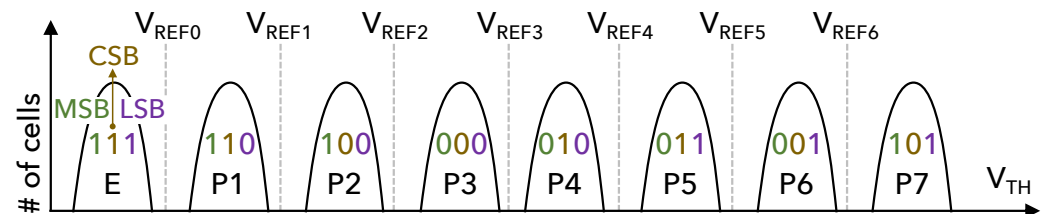


Target Page

Free
Valid
Invalid
Scrubbed

Target Block

$WL_0$

Program interference (slight program)

$WL_1$

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Page 3 | 0 | 1 | 0 | 0 | ⋯ | 1 | 1 | 0 | 1 |
| Page 4 | 0 | 1 | 0 | 1 | ⋯ | 1 | 0 | 0 | 1 |
| Page 5 | 1 | 1 | 0 | 0 | ⋯ | 0 | 0 | 1 | 0 |

$V_{REF0}$  $V_{REF1}$  $V_{REF2}$  $V_{REF3}$  $V_{REF4}$  $V_{REF5}$  $V_{REF6}$

# of cells

CSB
MSB LSB
111    110    100    000    010    011    001    101
E      P1     P2     P3     P4     P5     P6     P7     $V_{TH}$

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
  - Problem 1: MLC NAND flash stores multiple pages in a WL
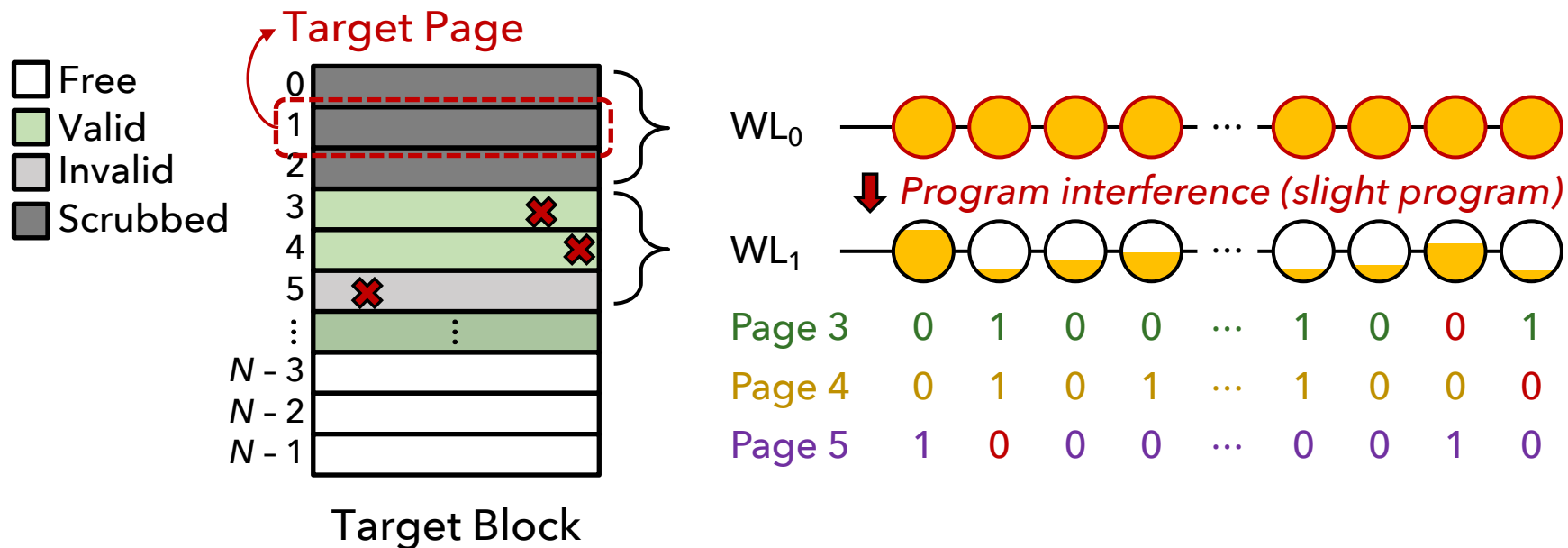  - Problem 2: Program interference

Target Page

Free
Valid
Invalid
Scrubbed

Target Block

$WL_0$

Program interference (slight program)

$WL_1$

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Page 3 | 0 | 1 | 0 | 0 | ⋯ | 1 | 0 | 0 | 1 |
| Page 4 | 0 | 1 | 0 | 1 | ⋯ | 1 | 0 | 0 | 0 |
| Page 5 | 1 | 0 | 0 | 0 | ⋯ | 0 | 0 | 1 | 0 |

# of cells

$V_{REF0}$  $V_{REF1}$  $V_{REF2}$  $V_{REF3}$  $V_{REF4}$  $V_{REF5}$  $V_{REF6}$

CSB

MSB LSB

111   110   100   000   010   011   001   101

E   P1   P2   P3   P4   P5   P6   P7

$V_{TH}$

# Existing Solutions

- Scrubbing [Wei+, FAST'11]
  - Problem 1: MLC NAND flash stores multiple pages in a WL
  - Problem 2: Program interference



Target Page

Free
Valid
Invalid
Scrubbed

Target Block

$WL_0$

Program interference (slight program)

$WL_1$

| Page 3 | 0 | 1 | 0 | 0 | ⋯ | 1 | 0 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|
| Page 4 | 0 | 1 | 0 | 1 | ⋯ | 1 | 0 | 0 | 0 |
| Page 5 | 1 | 0 | 0 | 0 | ⋯ | 0 | 0 | 1 | 0 |

Existing solutions incur
performance, lifetime, and reliability problems
in modern NAND flash memory

# Evanesco: Outline

- Data Remanace in NAND Flash-Based SSDs

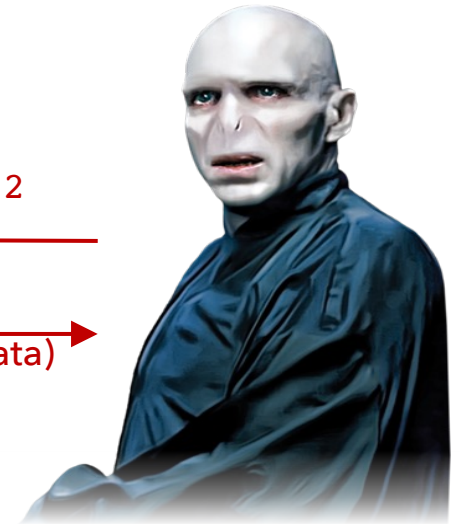- Evanesco: Access Control–Based Sanitization

- Evaluation Results
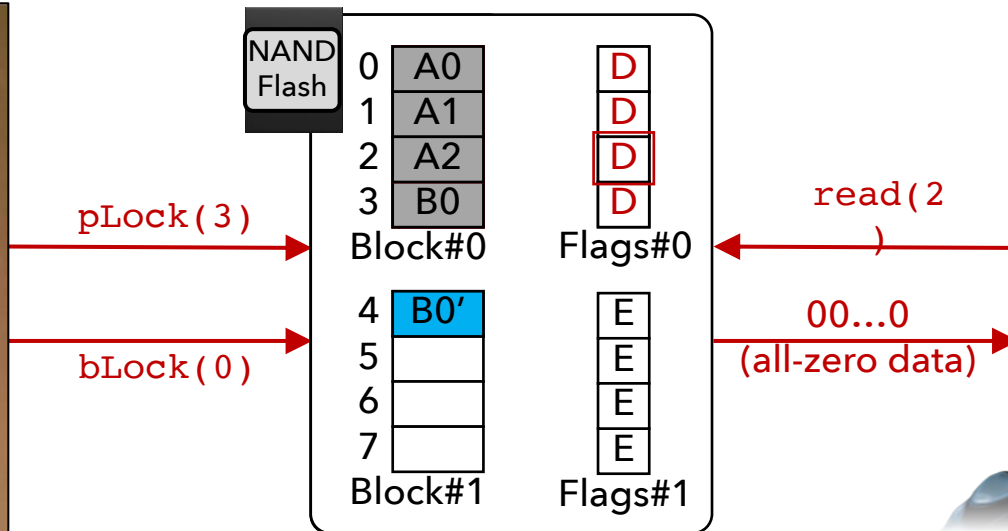
# Our Solution: Evanesco

- Allows a NAND flash chip to be aware of data validity
  - On-chip access control to avoid access to invalid data
  - Low overhead: No copy operations
  - High reliability: No program interference
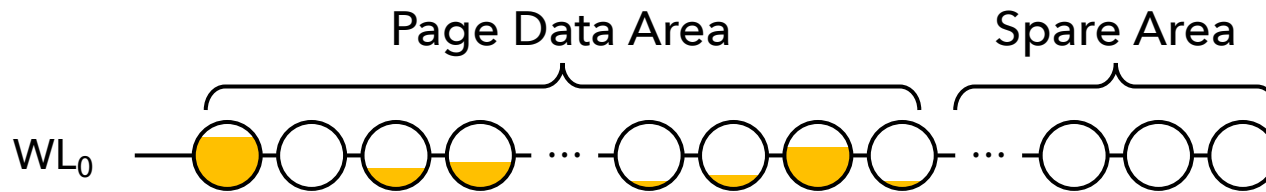
- Two new NAND commands: `pageLock` and `blockLock`

# Evanesco: Requirements

- Keep access-permission flags in a non-volatile manner

- Access-control logic inside a NAND flash chip

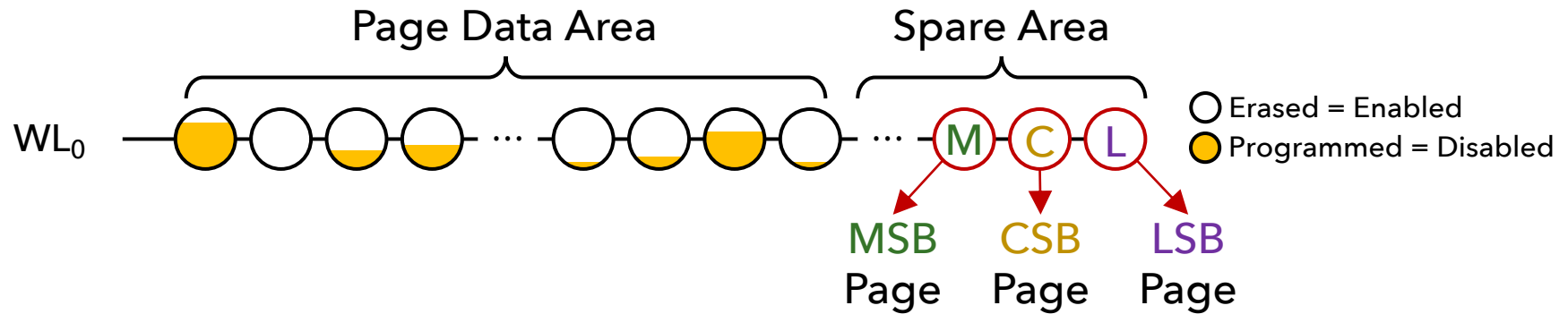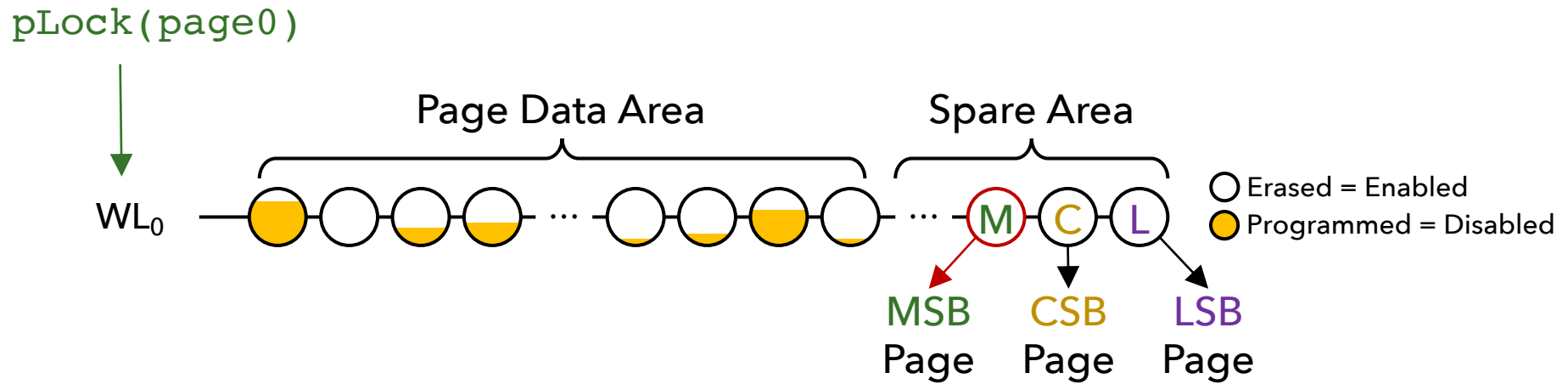- Minimal area overhead → High chip densitity is paramount
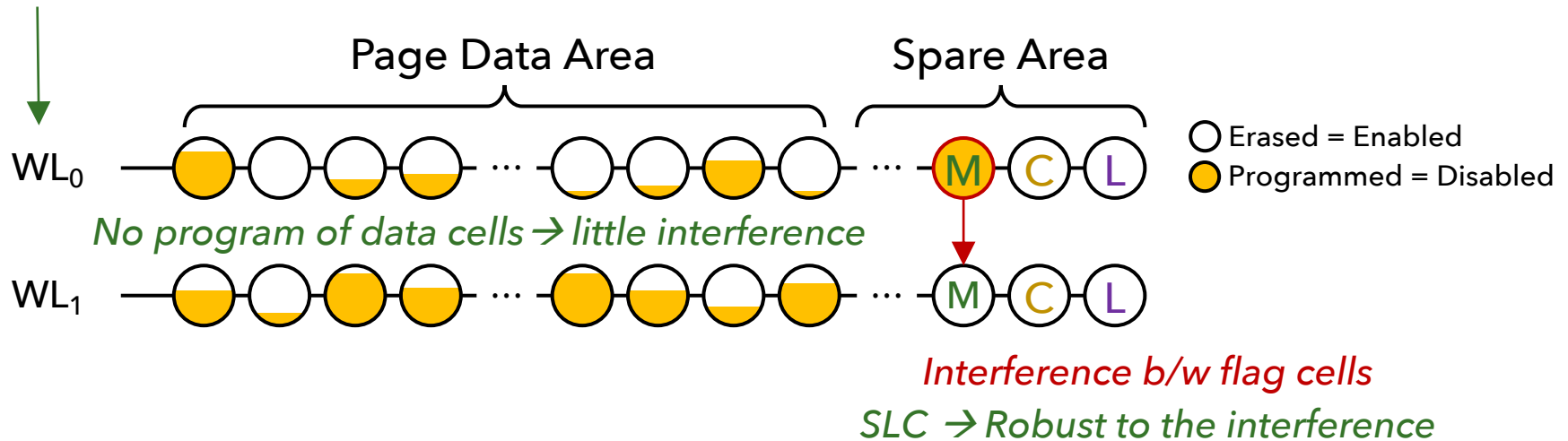
# pLock: Page-Level Data Sanitization

- On-chip access-permission flags: Spare cells in each WL

# pLock: Page-Level Data Sanitization

- On-chip access-permission flags: Spare cells in each WL



Page Data Area

Spare Area

$WL_0$

M  C  L

○ Erased = Enabled
● Programmed = Disabled

MSB      CSB      LSB
Page     Page     Page

# pLock: Page-Level Data Sanitization

- On-chip access-permission flags: Spare cells in each WL



pLock(page0)

Page Data Area  Spare Area

$WL_0$

MSB Page  CSB Page  LSB Page

○ Erased = Enabled
● Programmed = Disabled

# pLock: Page-Level Data Sanitization

- On-chip access-permission flags: <span style="color:orange">Spare cells</span> in each WL

`pLock(page0)`



Page Data Area      Spare Area

$WL_0$

*No program of data cells → little interference*

$WL_1$

○ Erased = Enabled
● Programmed = Disabled

*Interference b/w flag cells*

*SLC → Robust to the interference*

# pLock: Page-Level Data Sanitization

- On-chip access-permission flags: Spare cells in each WL
- On-chip access control logic: Small changes to data path

# pLock: Page-Level Data Sanitization
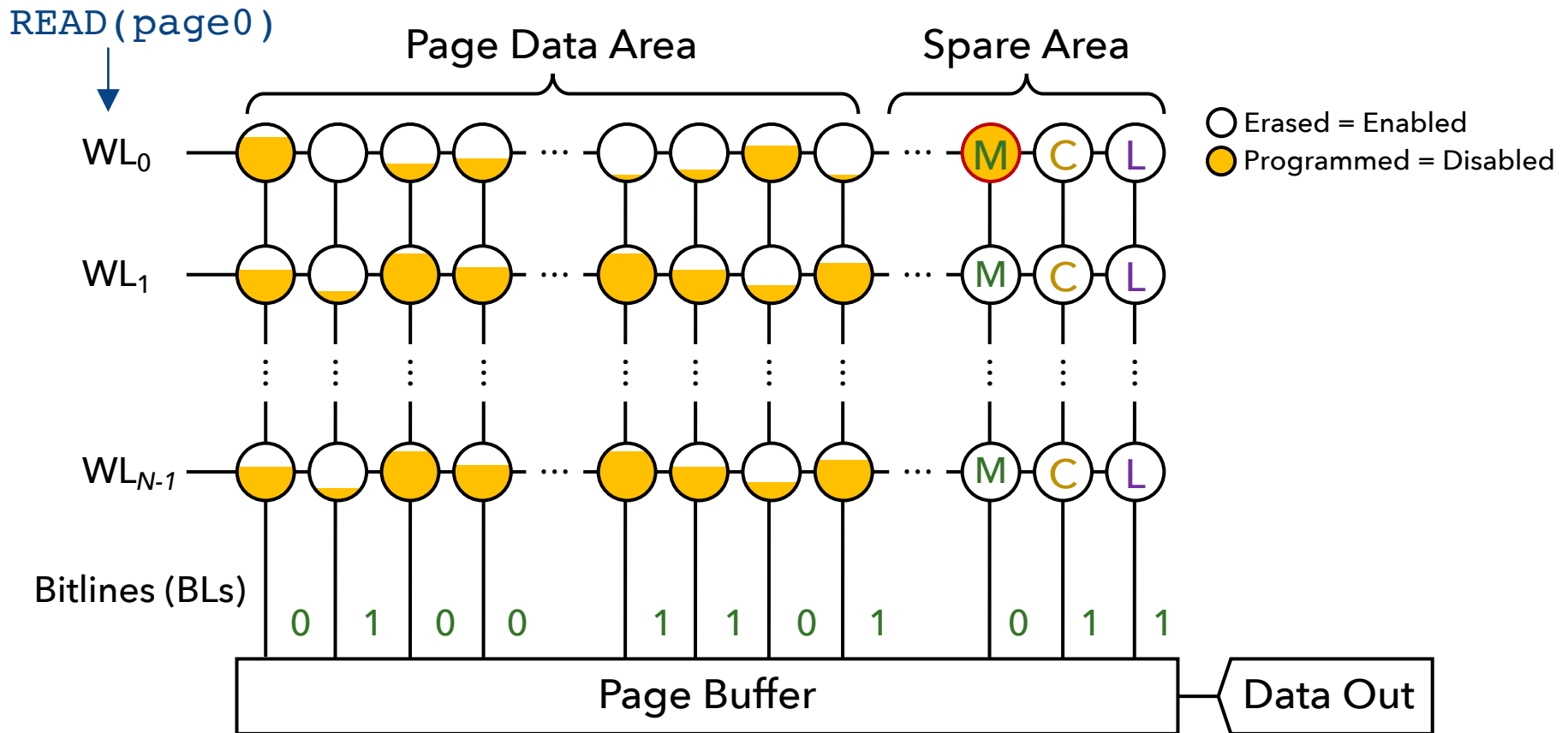
- On-chip access-permision flags: Spare cells in each WL
- On-chip access control logic: Small changes to data path

# pLock: Page-Level Data Sanitization

- On-chip access-permision flags: Spare cells in each WL
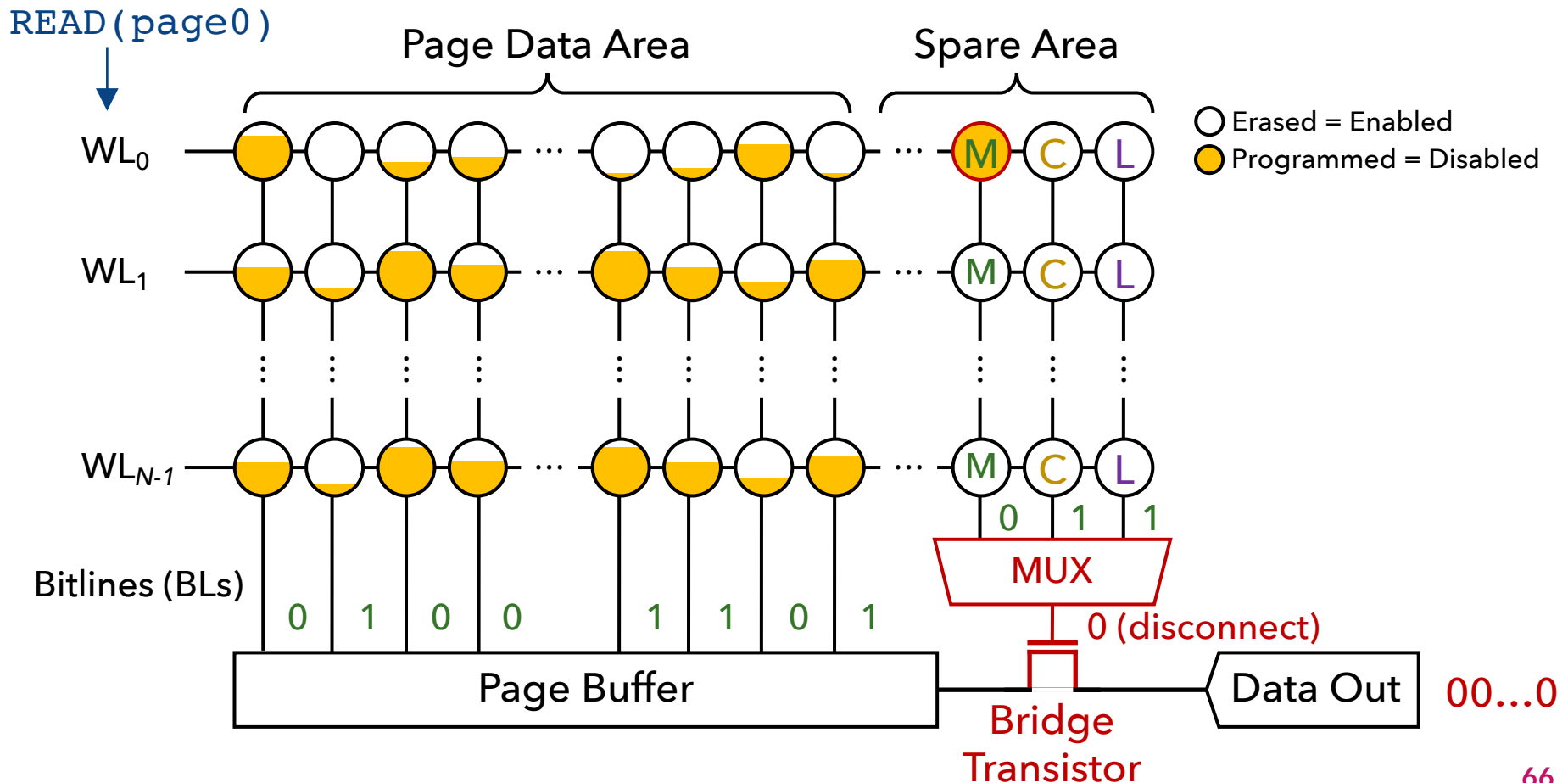- On-chip access control logic: Small changes to data path`

READ(page0)

Page Data Area      Spare Area
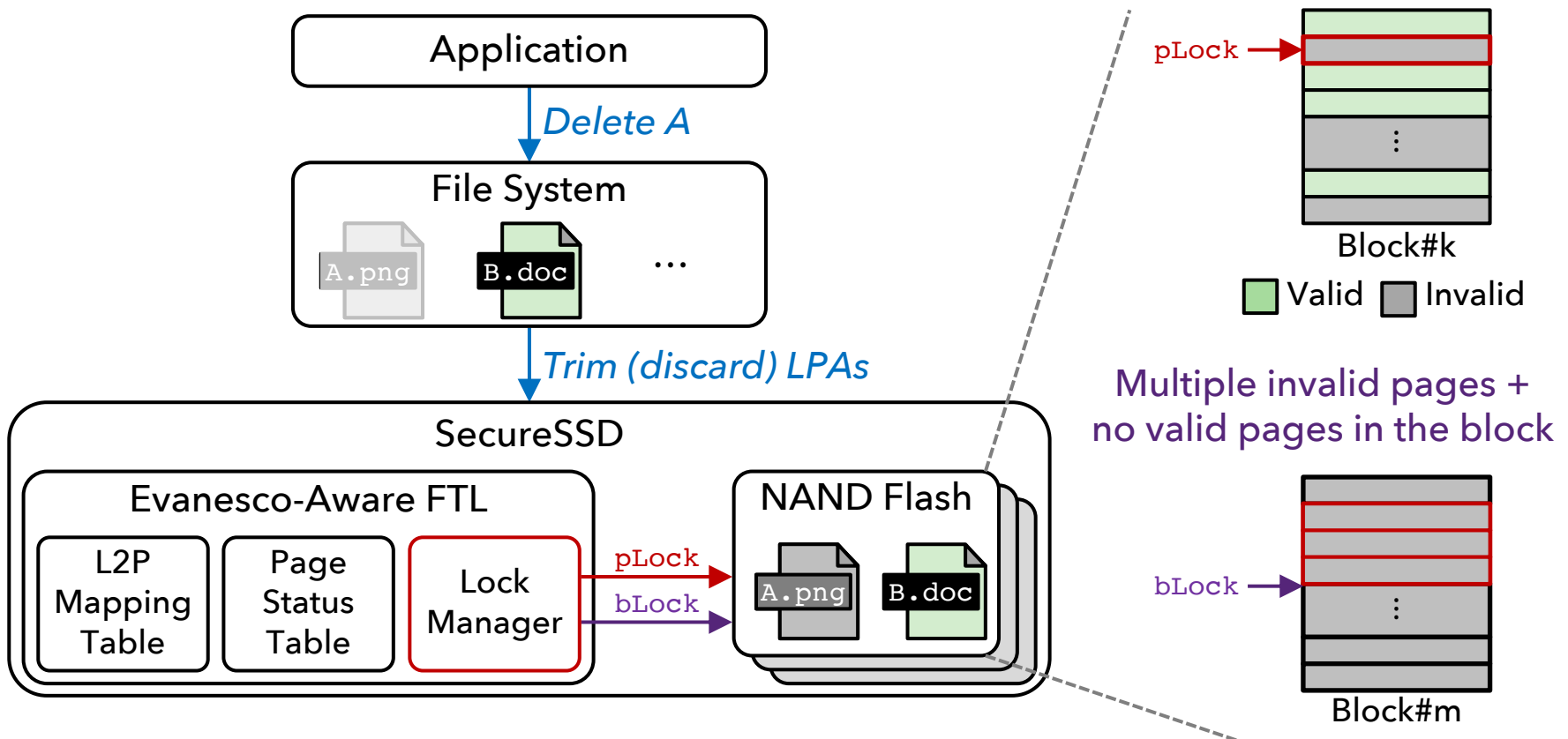
○ Erased = Enabled
● Programmed = Disabled

$WL_0$
$WL_1$
$WL_{N-1}$

M  C  L
0  1  1

MUX

Bitlines (BLs)

0  1  0  0    1  1  0  1

0 (disconnect)

Page Buffer

Bridge
Transistor

Data Out  00…0

# Real-Device Characterization

- Using 160 real 48-layer TLC NAND flash chips

- No reliability degradation for stored data

- `tPLOCK` = 100 μs, `tBLOCK` = 300 μs

Evanesco: No copy operation, no reliability issues
w/ minimal changes to NAND flash chip designs

But the performance overhead is not negligible

# SecureSSD: System-Level Optimization

- To minimize the performance overhead of data sanitization
  - Issues pLock and bLock commands depending on the status of the target block
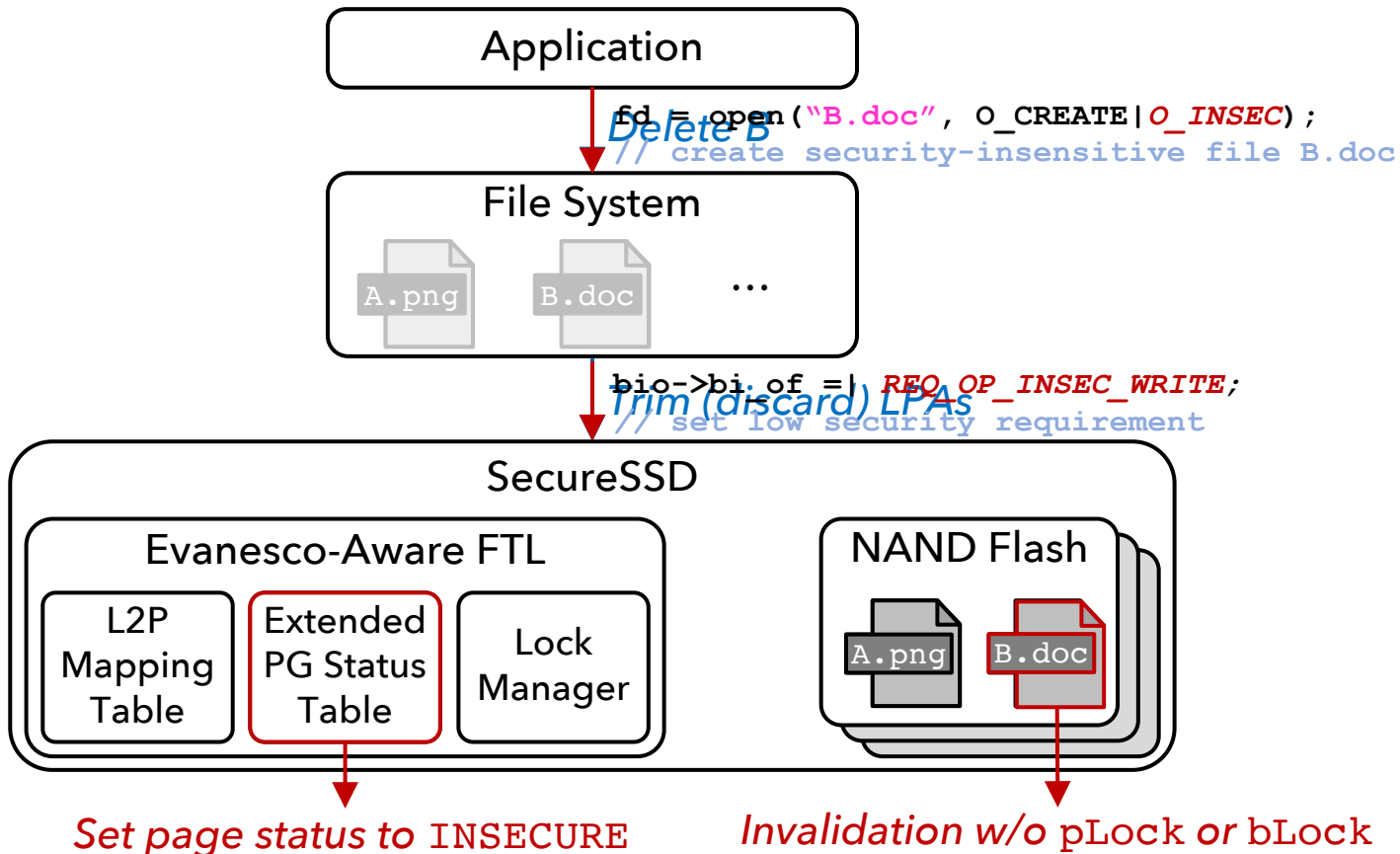
# SecureSSD: System-Level Optimization

- To minimize the performance overhead of data sanitization
  - Issues pLock and bLock commands depending on the status of the target block
  - Cross-layer interactions for selective data sanitization



```
fd = open("B.doc", O_CREATE|O_INSEC);
// create security-insensitive file B.doc
```

*Delete B*

```
bio->bi.of =| REQ_OP_INSEC_WRITE;
// set low security requirement
```

*Trim (discard) LPAs*

*Set page status to* `INSECURE`          *Invalidation w/o* `pLock` *or* `bLock`

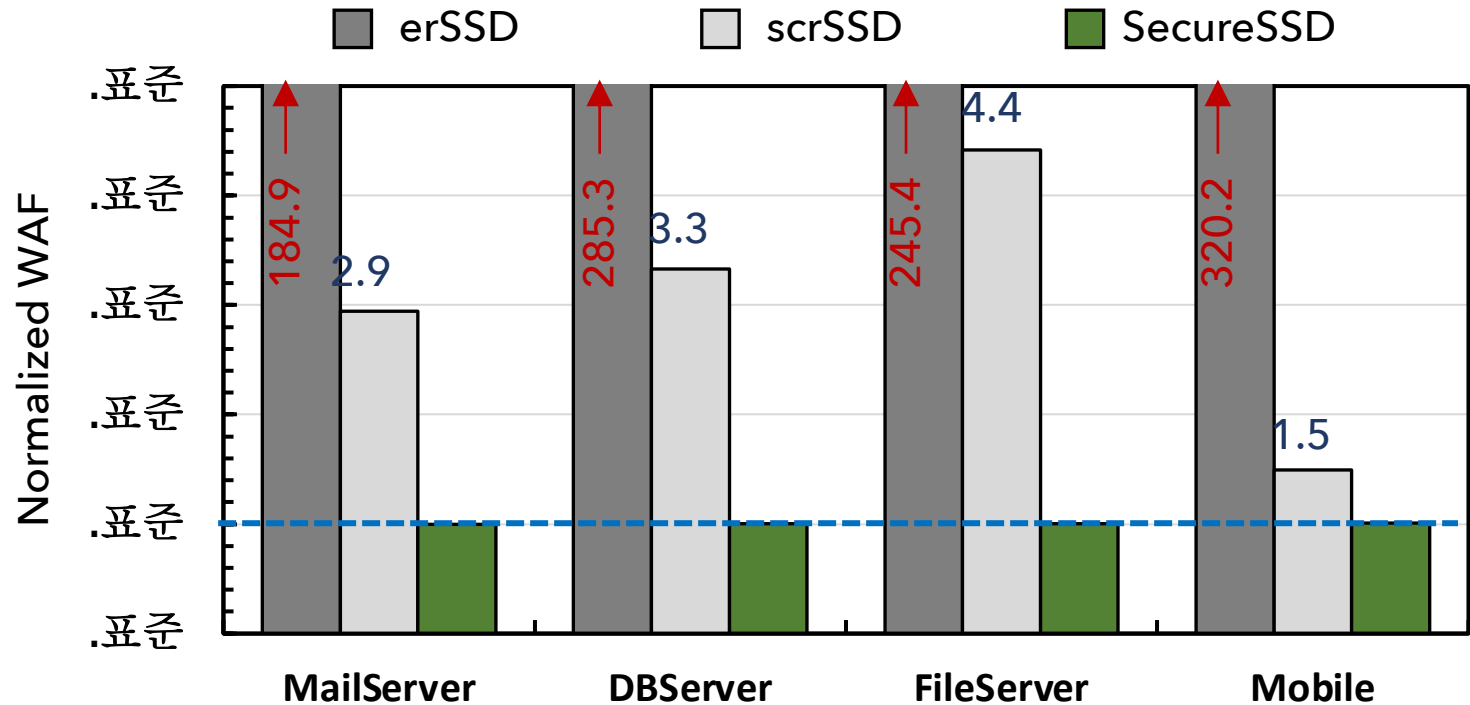# Evanesco: Outline

- Data Remanace in NAND Flash-Based SSDs

- Evanesco: Access Control–Based Sanitization

- Evaluation Results

# Results: Performance



Evanesco significantly reduces
performance overhead of data sanitization
(11% slowdown at most)

# Results: Lifetime



Legend: erSSD, scrSSD, SecureSSD

Normalized WAF

MailServer: 184.9, 2.9
DBServer: 285.3, 3.3
FileServer: 245.4, 4.4
Mobile: 320.2, 1.5

$$Write\ Amplification\ Factor\ (WAF) = \frac{\#\ of\ physical\ pages\ written\ by\ the\ SSD}{\#\ of\ logical\ pages\ written\ by\ the\ host\ system}$$

## No additional copy for data sanitization
## → No lifetime overhead

# Summary of Contribution

- **Problem1:** <span style="color:red">Long, non-deterministic</span> SSD read latency
  - Due to essential reliability management (read-retry)
  - Performance degradation of data-intensive applications

- Our solution: <span style="color:green">Pipelined & adaptive read-retry</span>
  - Leveraging <span style="color:blue">device characteristics</span> and <span style="color:blue">ECC margin</span>
  - Reducing read-retry latency
    → easy to combine with other optimizations

- **Problem2:** <span style="color:red">Data remanence</span> in NAND flash-based SSDs
  - Obsolete data remains intact in SSDs <span style="color:red">for an indefinite time</span>
  - Physical data destrubtion: <span style="color:red">prohibitive performance overheads</span>

- Our solution: <span style="color:green">Access control-based data sanitization</span>
  - Avoids <span style="color:blue">transfer of obsolete data</span> from NAND flash chips
  - Minimal performance and reliability overheads

# P&S Modern SSDs

## Research Session 1:
## Data Sanitization and Read-Retry
## in Modern NAND Flash-Based SSDs

Dr. Jisung Park

Prof. Onur Mutlu

ETH Zürich

Spring 2022

15 July 2022