



Mohamed bin Zayed  
University of  
Artificial Intelligence

# Unlocking the Full Potential of 3D NAND Flash

A Cross-Layer Approach to Reliability, Performance, and Lifetime

---

**Chun Jason Xue**

CHEOPS Workshop @ EuroSys 2026

## **Part I: Understanding Flash Reliability**

Physical characterization & key insights

## **Part II: Boosting Read Performance**

From LDPC optimization to near-zero read retry

## **Part III: Smart Data Encoding**

Encoding strategies for reliability & lifetime

## **Part IV: The Reprogram Revolution**

Repurposing program operations for SSD optimization

## **Part V: Future Challenges & Opportunities**

Where do we go from here?

# NAND Flash: The Backbone of Modern Storage

## Data Center



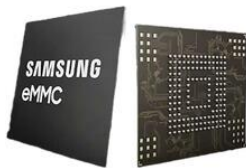
SATA eSSD, PCIe eSSD

## Automotive-grade Storage Products



eMMC

## Personal Mobile Phone



eMMC, UFS

## Consumer Mobile Storage Device



eMMC, UFS

## Computing Device



SATA cSSD, PCIe cSSD

Flash storage is **everywhere**: from mobile devices to data centers.

# NAND Flash: The Backbone of Modern Storage

**\$100B+**

Global NAND market  
(2025)

**60%+**

Data center storage is  
flash

**80%+**

Data is cold / rarely  
accessed

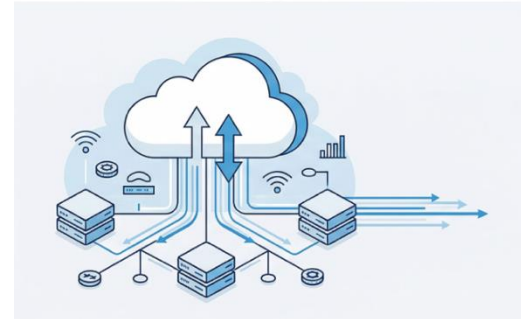
Understanding and optimizing NAND flash at every layer is critical to sustaining the storage industry's growth.



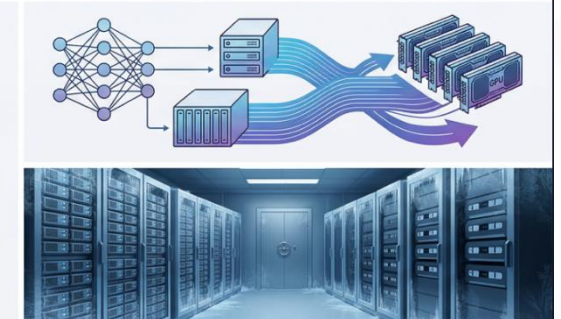
Consumer SSDs



Enterprise SSDs



Cloud Storage



AI/ML Training Data

ISSCC 2024

ISSCC 2026

ISSCC 2025

**13.3 A 280-Layer 1Tb 4b/cell 3D NAND Flash Memory with a 3.2GB/s High-Speed Write Throughput**  
H. Kim, W. Jung, D-B. Kim, T-H. Kim, T. Kim, H. Kim, G. Lee, J. Lee, J. Lee, T. Lee, B-K. Chun, T. Kim, Y. G. Lee, Y. Lee, S. Kim, J. Hwang, R. Song, M. Lee, S. Jo, C. H. Kim, J. C. Park, Samsung Electronics, Hwaseong, Korea

**13.7 A 1Tb Density 3b/Cell 3D-NAND Flash with a 3.2GB/s High-Speed Write Throughput**  
K. Kawai<sup>1</sup>, Y. Einaga<sup>1</sup>, Y. Oikawa<sup>1</sup>, A. D'alessandro<sup>3</sup>, E. Yu<sup>2</sup>, A. Murakami<sup>4</sup>, T. Ichikawa<sup>1</sup>, J. Yu<sup>2</sup>, G. Wang<sup>2</sup>, K. M. S. Bhushan<sup>2</sup>, D. Srinivasan<sup>2</sup>, H. Kuo<sup>5</sup>, C. Siau<sup>2</sup>, R. Ghodsji<sup>2</sup>, <sup>1</sup>Micron Technology, <sup>2</sup>Micron Technology, Avezzano, Italy, <sup>3</sup>Micron Technology, Avezzano, Italy, <sup>4</sup>Micron Technology, Hyderabad, India, <sup>5</sup>Micron Technology, Hyderabad, India

W. Cho, J. Jung, J. Kim, J. Ham, S. Lee, Y. Nam, H. Cho, J-S. Kim, C. Kwon, C. Park, H. Nam, T. Shin, J. Jang, J. Mun, J. Choi, H. Choi, S-V. H. Oh, H. Park, S. Shim, H. Huh, H. Choi, S. Lee, SK hynix, Icheon, Korea

**7.4 A 1Tb 3b/Cell 8th-Generation 3D-NAND Flash with a 2.4Gb/s Interface**  
M. Kim, S. W. Yun, J. Park, H. K. Park, J. Lee, K. Lee, B. Jeong, S. Kim, J. Park, C. A. Lee, J. Lee, J. Lee, J. Y. Chun, J. Jang, Y. Yang, S. H. M. Choi, W. Kim, J. Kim, S. Yoon, P. Kwak, M. Lee, R. Song, S. Kim, C. Yoon, D. Kang, J. J. Song, Samsung Electronics, Hwaseong, Korea

**15.1 A 2Tb 4b/Cell 6-Plane 3D-Flash Memory with 37.6Gb/mm<sup>2</sup> Bit Density and >85MB/s Write Throughput**  
J. M. Thimmaiah<sup>1</sup>, R. Yamashita<sup>2</sup>, I-S. Yoon<sup>3</sup>, J. Li<sup>3</sup>, C. Hsu<sup>3</sup>, T. Arikai<sup>2</sup>, N. Ookuma<sup>2</sup>, Y. Kato<sup>2</sup>, K. Hayashi<sup>2</sup>, K. Yamauchi<sup>2</sup>, I. K V<sup>1</sup>, M. Kano<sup>2</sup>, S. Bhamidipati<sup>1</sup>, S. Bhatia<sup>1</sup>, S. Malhotra<sup>1</sup>, N. Ojima<sup>2</sup>, E. Wu<sup>3</sup>, Z. Yang<sup>3</sup>, F. W. Tsai<sup>3</sup>, M. Bayle<sup>2</sup>, N. Minami<sup>2</sup>, Y. Fujihara<sup>2</sup>, K. Kitamura<sup>2</sup>, T. Kitani<sup>2</sup>, T. Kodama<sup>4</sup>, T. Handa<sup>4</sup>, N. Kanagawa<sup>4</sup>, Y. Ishizaki<sup>4</sup>, S. Fujimura<sup>4</sup>, Y. Suzuki<sup>4</sup>, M. Sako<sup>4</sup>, Y. Higashi<sup>4</sup>, Y. Watanabe<sup>4</sup>, T. Kouchi<sup>4</sup>, A. V<sup>1</sup>, C-Y. Chen<sup>3</sup>, X. Yang<sup>3</sup>, G. Liang<sup>3</sup>, J. Wang<sup>3</sup>

<sup>1</sup>Sandisk, Bengaluru, India; <sup>2</sup>Sandisk, Yokohama, Japan  
<sup>3</sup>Sandisk, Milpitas, CA; <sup>4</sup>KIOXIA, Yokohama, Japan

**30.6 A 16Mb 166.8TOPS/W Near-Memory Phase-Domain-Computing Ferroelectric NAND Flash for Approximate Nearest Neighbor Search on Edge Devices**  
W. Li<sup>\*1,2</sup>, B. Wang<sup>\*1,2</sup>, Z. Zhou<sup>1,2</sup>, J. Zhu<sup>1,2</sup>, Z. Li<sup>1,2</sup>, J. Shen<sup>1,2</sup>, W. Zha<sup>1,2</sup>, Z. Han<sup>1,2</sup>, Y. Wang<sup>1,2</sup>, L. Wang<sup>3</sup>, H. Hu<sup>1,2</sup>, Q. Luo<sup>1,2</sup>, C. Dou<sup>1,2</sup>, M. Liu<sup>1</sup>

<sup>1</sup>Institute of Microelectronics of the Chinese Academy of Sciences, Beijing, China  
<sup>2</sup>University of Chinese Academy of Sciences, Beijing, China  
<sup>3</sup>Columbia University, New York, NY  
\*Equally Credited Authors (ECAs)

**Cell WF-Bonding 3D-NAND Flash**  
G. Yu, C-H. Yu, H. Makoto, Y. Kwon, J-H. Park, Y. Kim, D-H. Kim, Y. Jo, H. Yoon, J. Park, H-S. Ku, J. Seo, J. Byun, S-H. Yun, K. Kang, Y. Ryu, H. Kim, W. Kim, H. Choi, J. Jeon, Lee, M-K. Lee, J-I. Son, J. Cho, M. Kim, S. Lee, K. Song, S-H. Hur

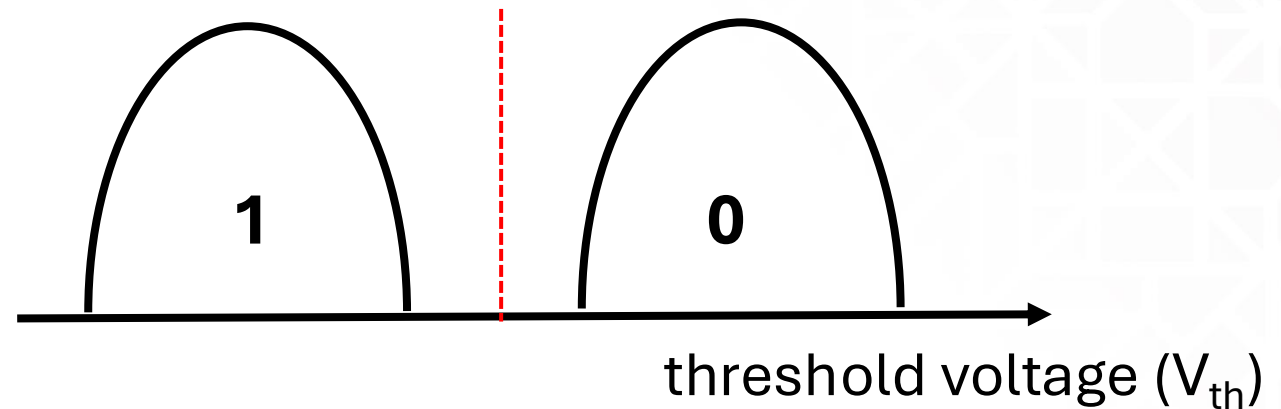
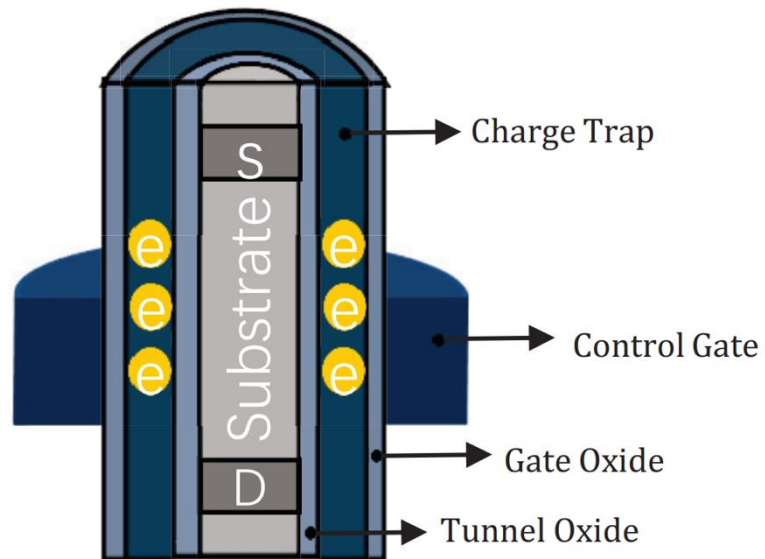
**M with a 29%-Improved-Energy- .8Gb/s Power-Isolated**  
Matsuno<sup>1</sup>, Y. Higashi<sup>1</sup>, Y. Shimizu<sup>1</sup>, Nakano<sup>1</sup>, Y. Ochi<sup>1</sup>, H. Hoshino<sup>1</sup>, T. Hioka<sup>1</sup>, Y. Kamata<sup>1</sup>, H. Chibvongodze<sup>2</sup>, N. Ojima<sup>2</sup>, Ukoshi<sup>2</sup>, R. Yamashita<sup>2</sup>, K. Abe<sup>2</sup>, Huh<sup>2</sup>, K. Htoo<sup>2</sup>, Y. Kato<sup>2</sup>, Y. Watanabe<sup>1</sup>, Western Digital, Milpitas, CA

# NAND Flash Structure

3D NAND flash package

Flash cell: charge trap based

Storing charges inside gates to represent data



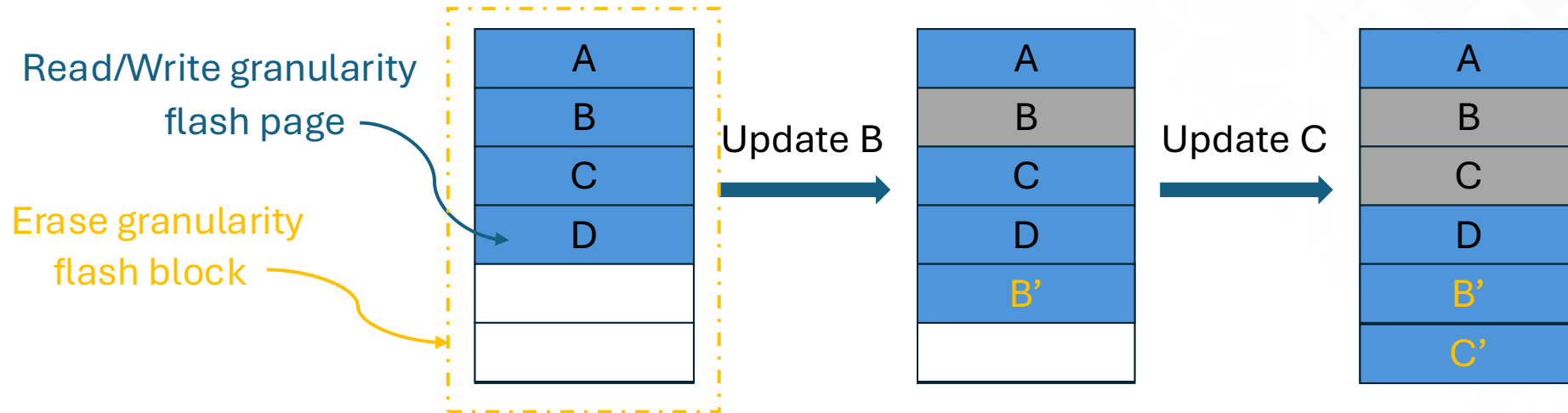
**Program (Write):** charge cells

**Read:** identify the voltage levels, low or high

**Erase:** discharge cells ← Unique

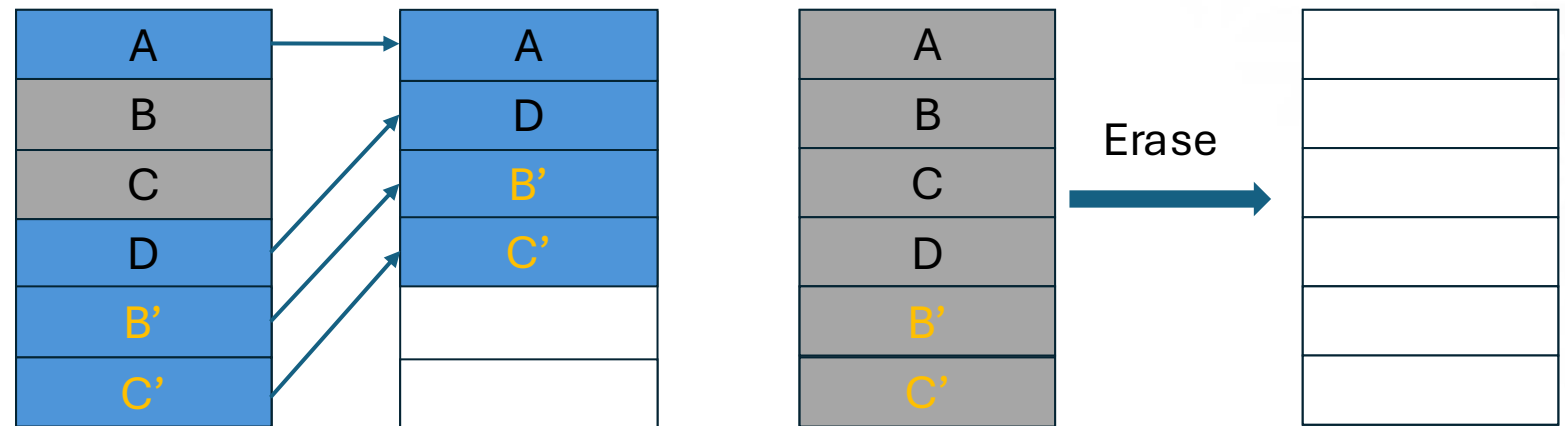
# Out-of-Place Update

Write updated data to a new free location, mark the original location as invalid.



## Garbage collection:

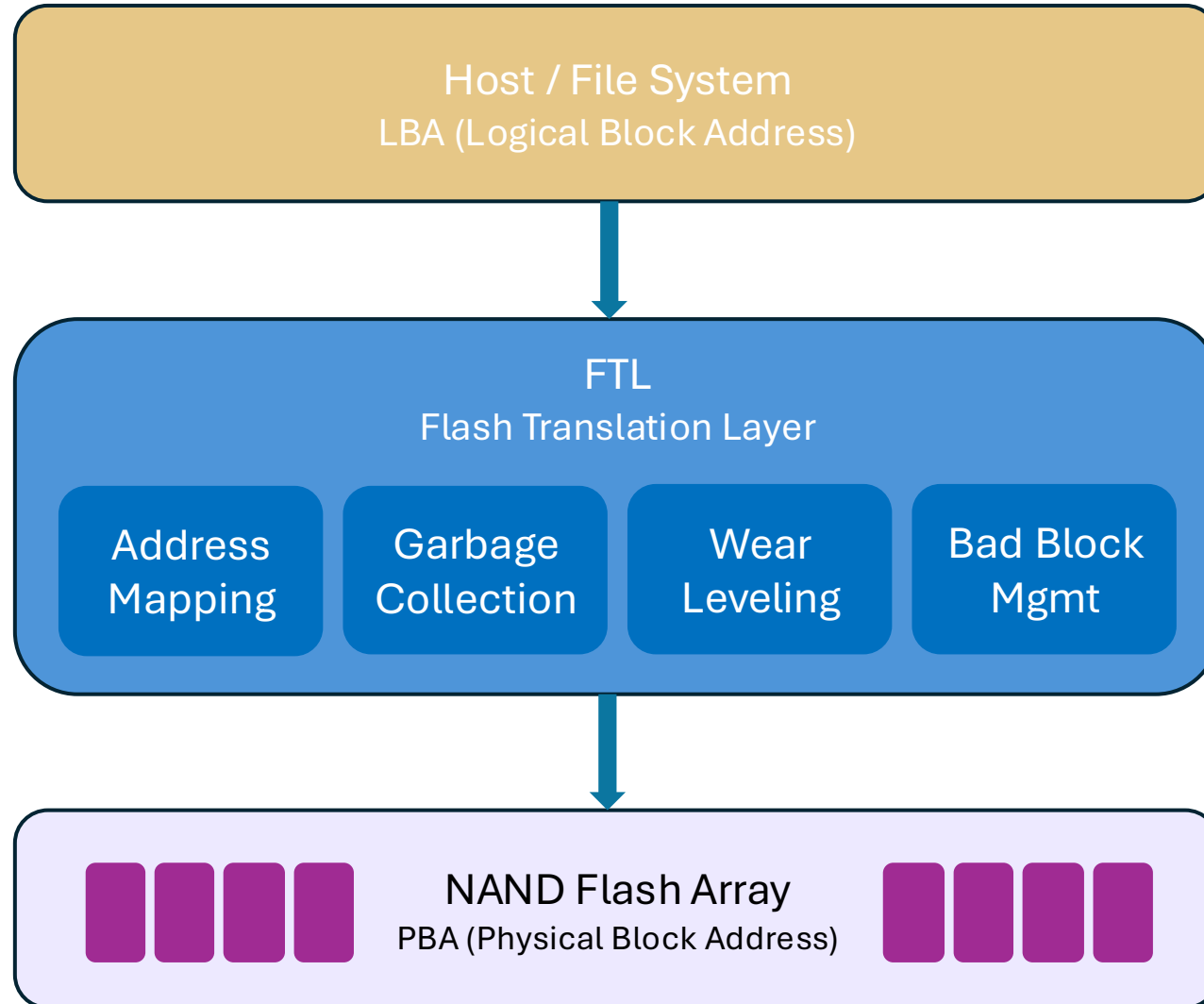
migrate valid data →  
erase the block →  
release free space.



migrate valid data

# Flash Translation Layer

Core firmware layer connecting the host with the NAND physical medium.



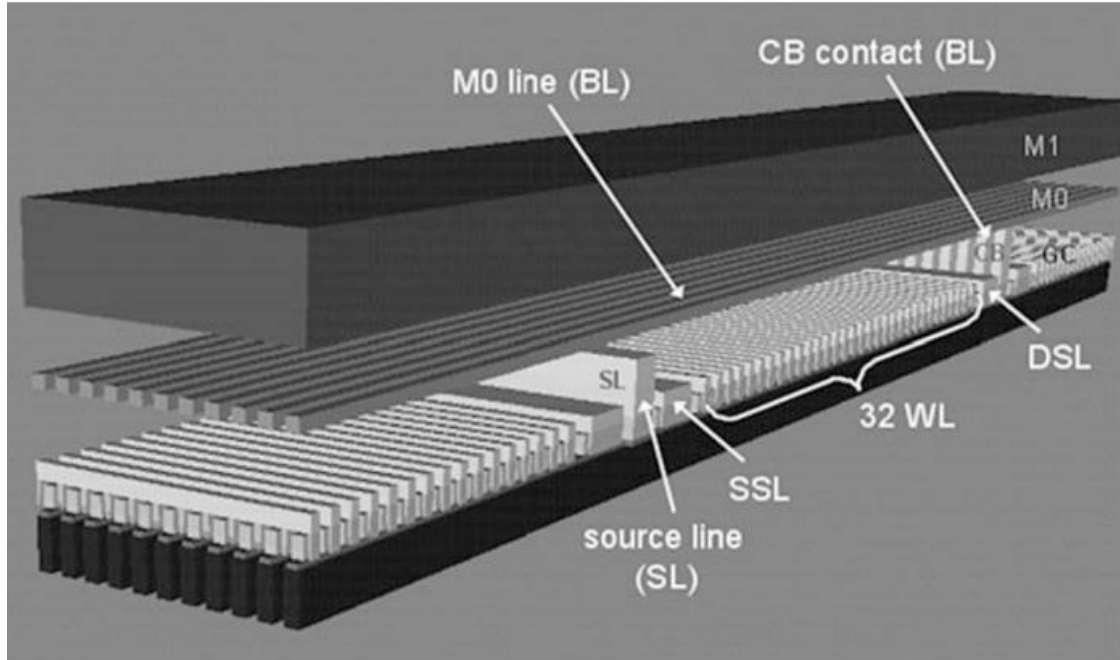
**Address Mapping:** Translates host LBAs to NAND PBAs and maintains the mapping table.

**Garbage Collection (GC):** Reclaims blocks containing invalid pages

**Wear Leveling:** Evenly distributes P/E cycles to prevent premature wear of specific blocks.

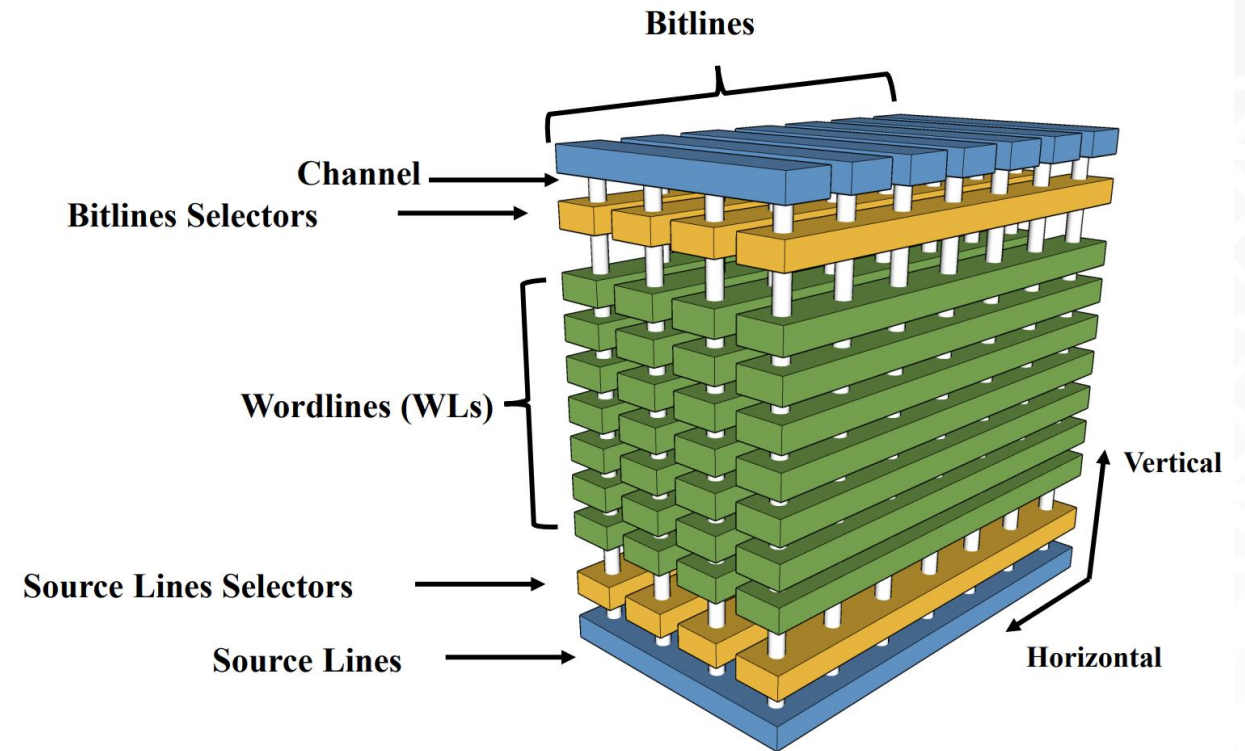
**Bad Block Management:** Detects and isolates factory-bad blocks and runtime-bad blocks, maintaining a bad block table.

## 2D NAND flash memory



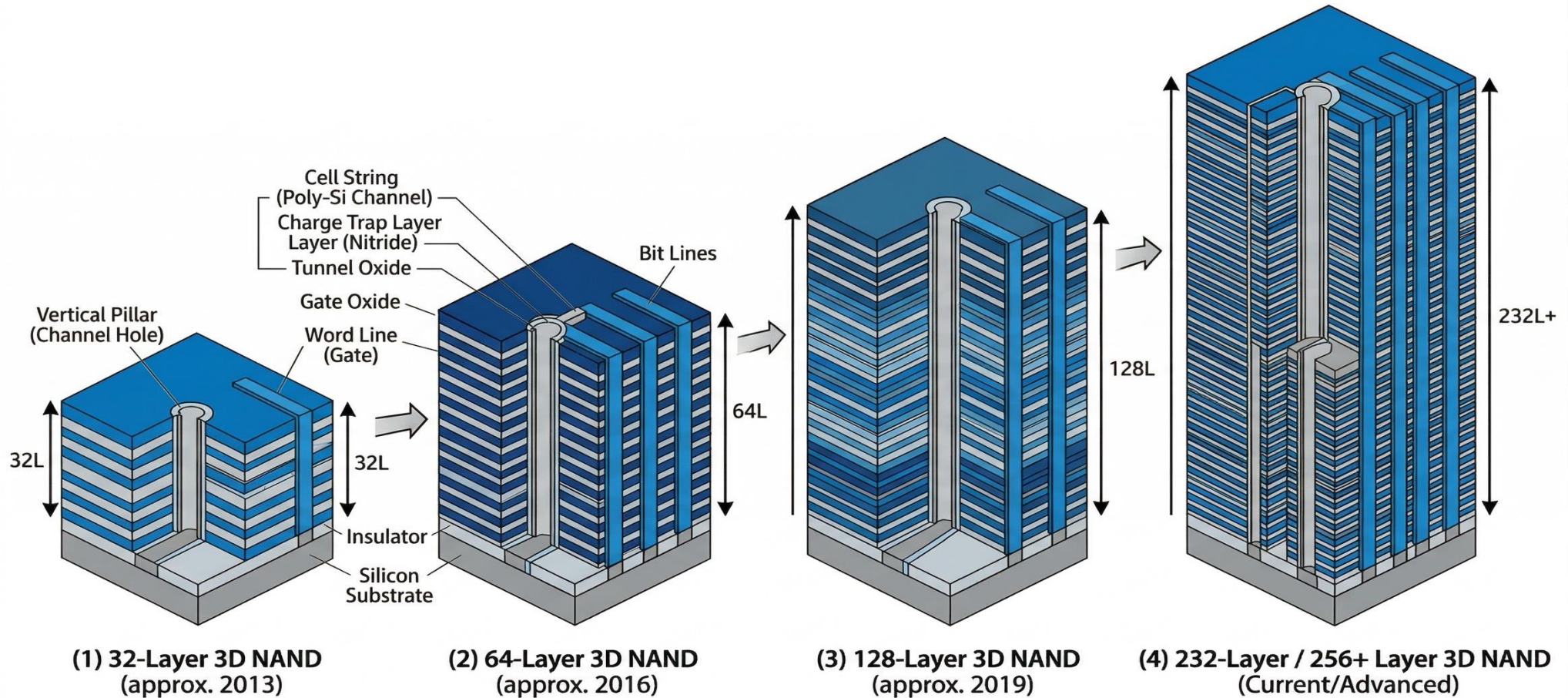
Scaling method: shrinking cell size  
Process node: 15–19 nm  
Cell-to-cell interference: severe

## 3D NAND flash memory



Scaling method: vertical layer stacking  
Layer count: 128 / 176 / 232+  
Cell-to-cell interference: significantly improved

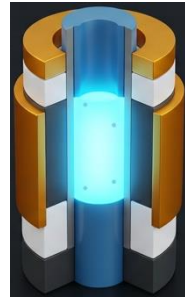
# 3D NAND Scaling



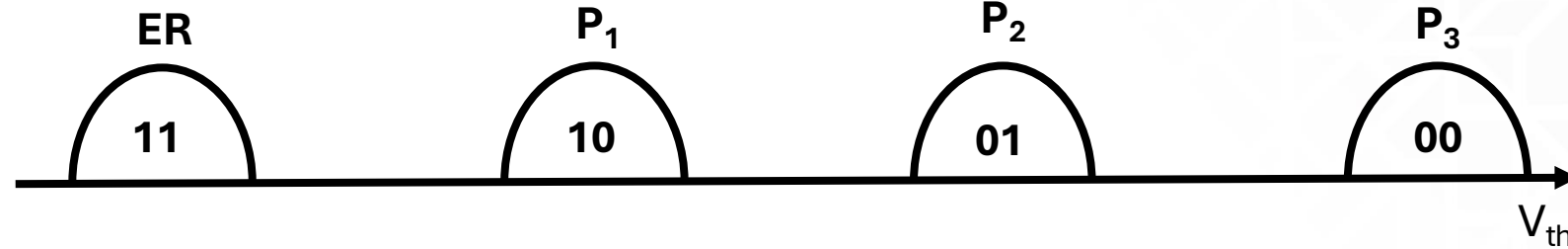
Number of layers inside 3D flash block in increasing

3D NAND flash tend to store more bits per cell

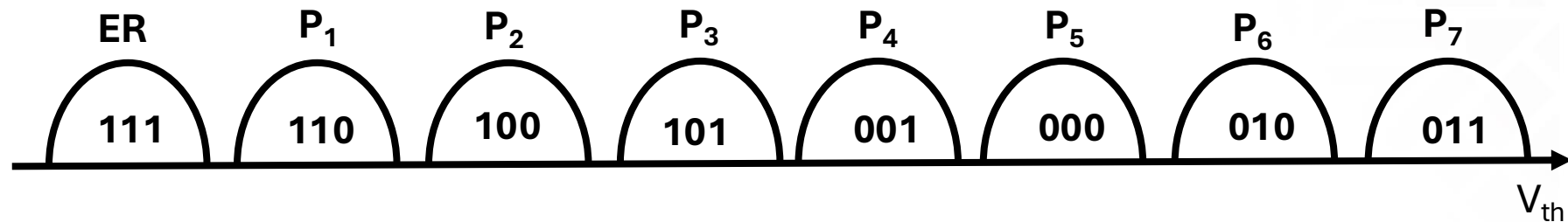
SLC  
1b/cell  
2 states



MLC  
2b/cell  
4 states



TLC  
3b/cell  
8 states

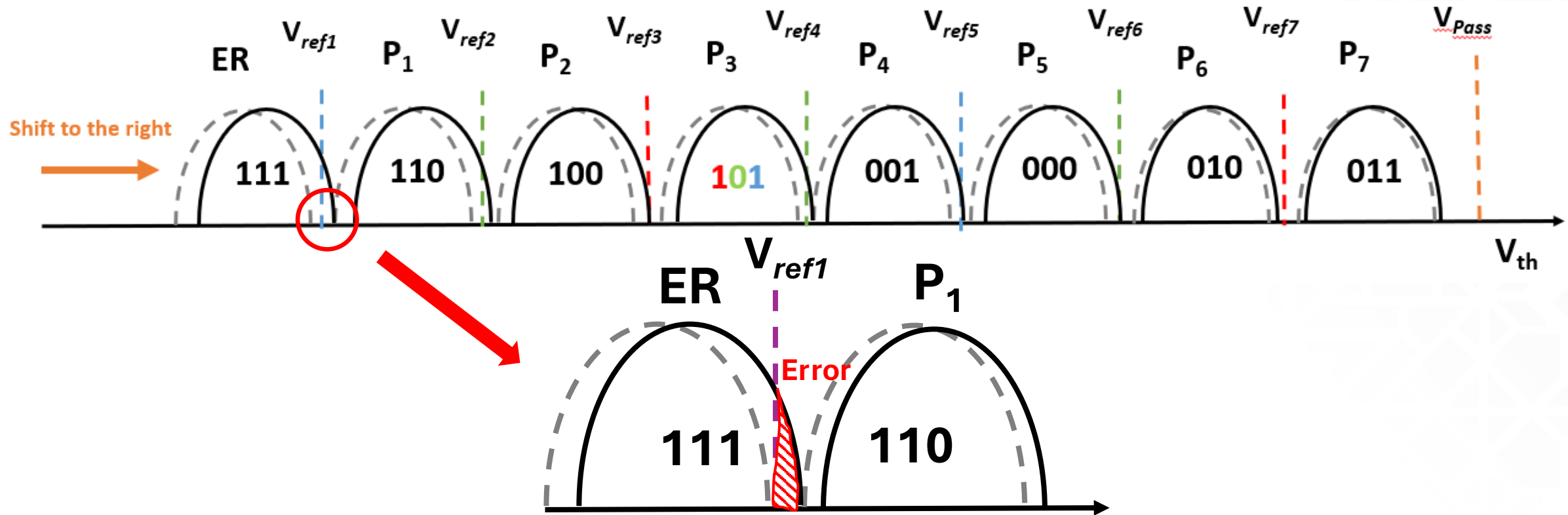


**SLC (1b/cell) → MLC (2b/cell) → TLC (3b/cell) → QLC (4b/cell)**

# 3D NAND Scaling

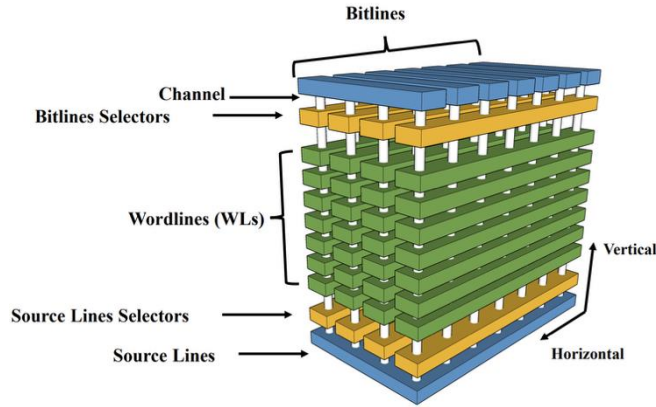
Suffer from multiple types of errors

More errors are introduced because of tighter voltage margins

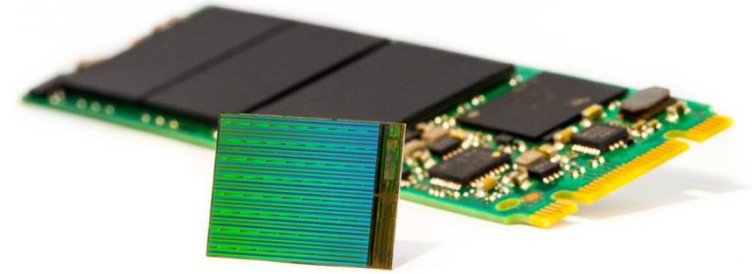


⚠ Higher density → tighter voltage margins → more errors → worse performance & shorter lifetime

# The Flash Scaling Trilemma



 **Reliability**



*Trade-offs*

 **Performance**

 **Lifetime**

As 3D NAND scales to higher density, all three attributes degrade simultaneously.

Our approach: Cross-layer understanding and optimization to push the Pareto frontier on all three dimensions simultaneously.

# A Decade of Flash Research (2015–2026)

- **2015–16** FAST '16: Access-guided R/W regulation
- **2017–18** ASP-DAC '17 (Best Paper Nominee) · ICCD '18: Selective compression
- **2019** DAC '19 · MICRO '19: Reprogramming TLC · TC '19: Asymmetric LDPC
- **2020** MICRO '20: Sentinel Cells · DAC '20: Partition for read perf.
- **2021–22** HotStorage '21 ★ **Best Paper** · TCAD '22: Open blocks
- **2023** HPCA '23: MGC · NVMSA '23 ★ **Best Paper** · TCAD '23: LDPC prediction
- **2024** ASPLOS '24: Near-zero retry · HPCA '24: Midas Touch · TACO '24
- **2026** ISCA '26: LOONG · EuroSys '26: ColdCode

Venues: ISCA, MICRO, HPCA, ASPLOS, EuroSys, FAST, DAC

# Part I

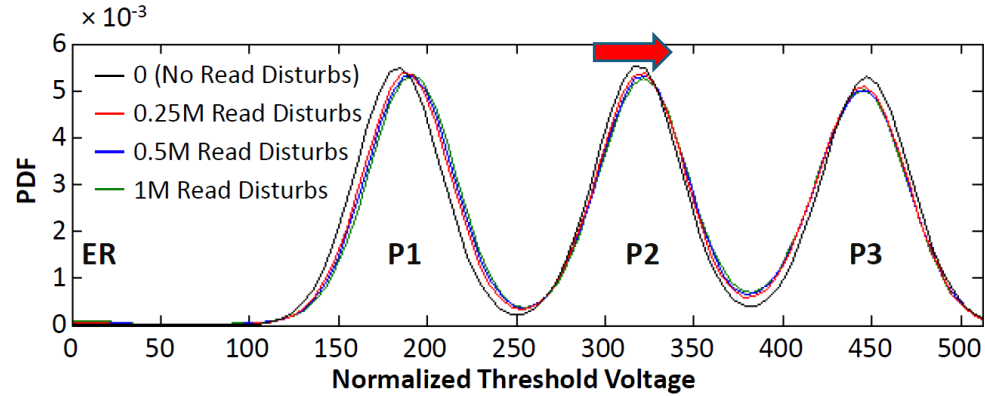
## Understanding Flash Reliability

Deep physical characterization as the foundation for all optimizations

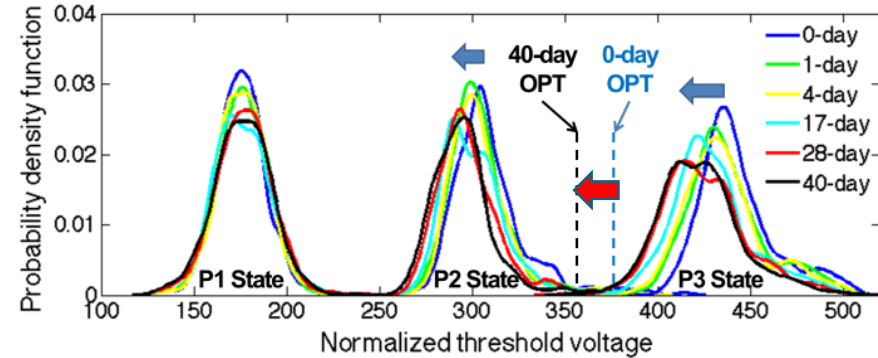
# NAND Flash Reliability

Different types of error sources impact flash reliability.

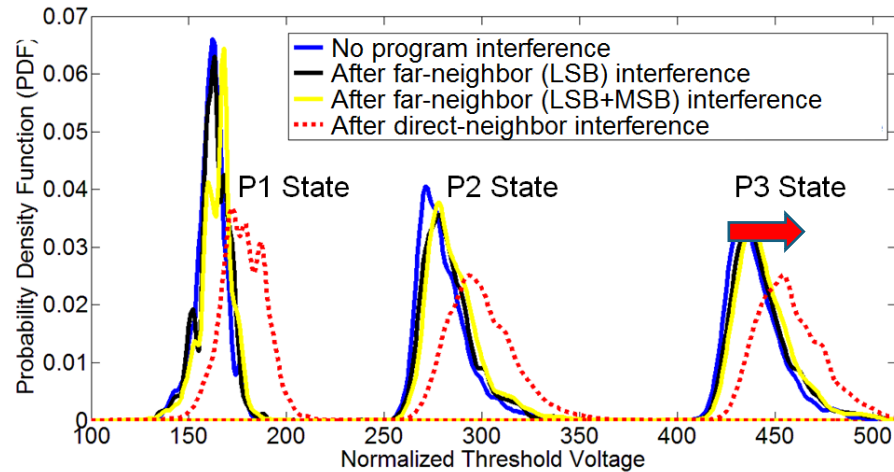
## Read disturb [Cai et al. DSN]



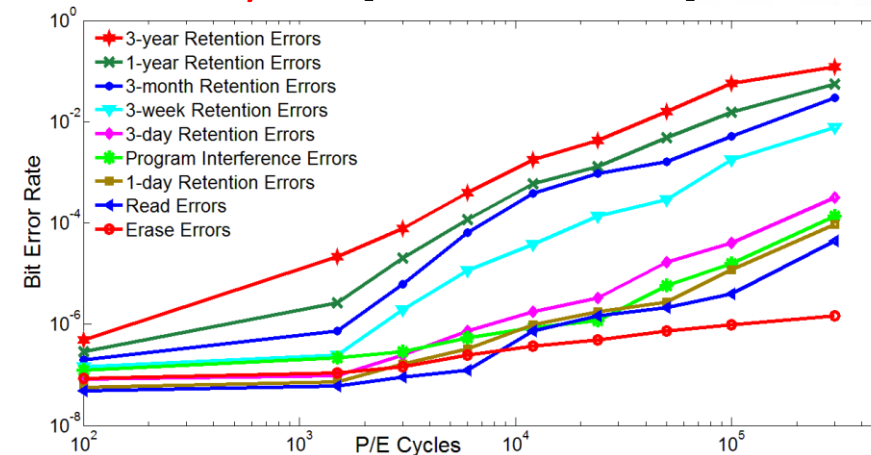
## Retention error [Cai et al. HPCA]



## Program Interference [Cai et al. ICCD]



## P/E Cycles [Cai et al. DATE]



# When Retention Testing Goes Wrong

Industry uses bake acceleration (high-temperature baking) to predict long-term retention degradation. But does this common method actually work for 3D NAND?

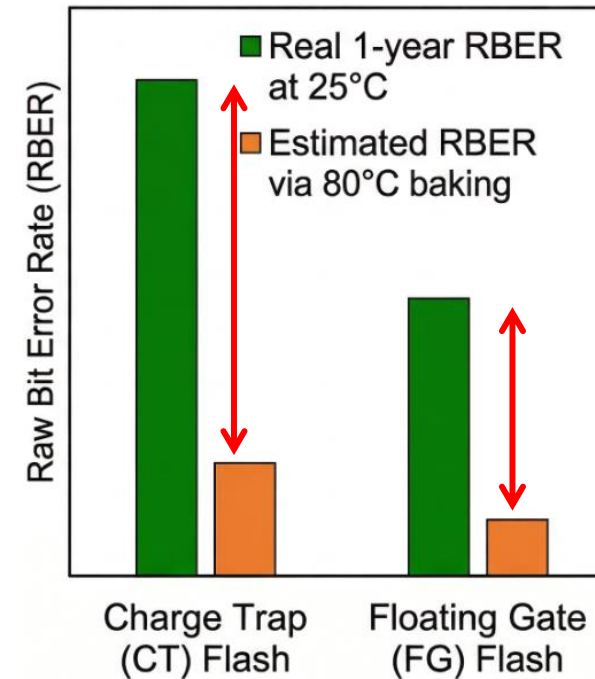
$$AF = \frac{t_1}{t_2} = \exp\left(\frac{E_a}{k} * \left(\frac{1}{T_1} - \frac{1}{T_2}\right)\right)$$

Retention time at  $T_1$

Retention time at  $T_2$

A high temperature

A low temperature



The answer is NO!

The baking underestimates real long-time errors

# When Retention Testing Goes Wrong

Industry uses bake acceleration (high-temperature baking) to predict long-term retention degradation. But does this common method actually work for 3D NAND?

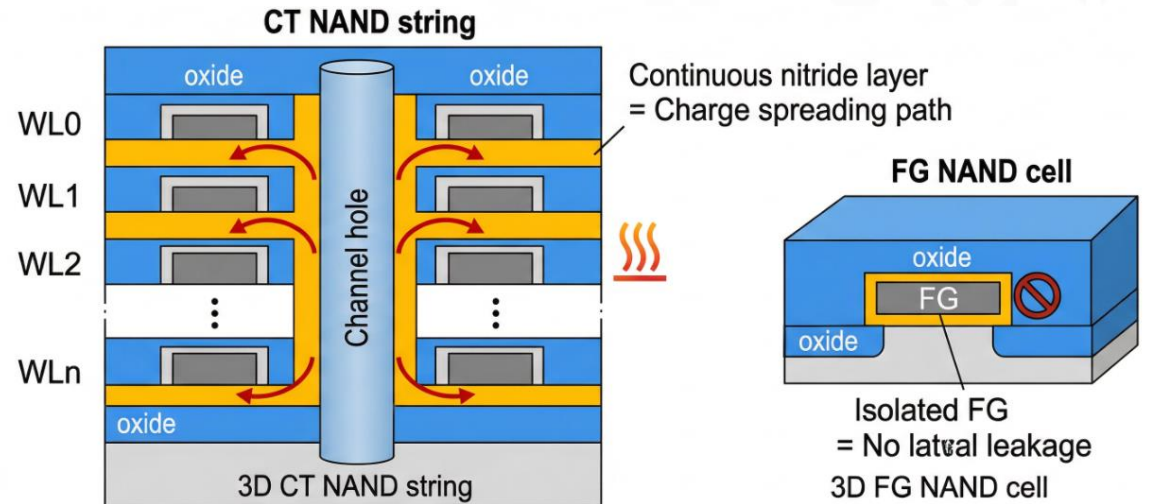
$$AF = \frac{t_1}{t_2} = \exp\left(\frac{E_a}{k} * \left(\frac{1}{T_1} - \frac{1}{T_2}\right)\right)$$

Retention time at  $T_1$

Retention time at  $T_2$

A high temperature

A low temperature



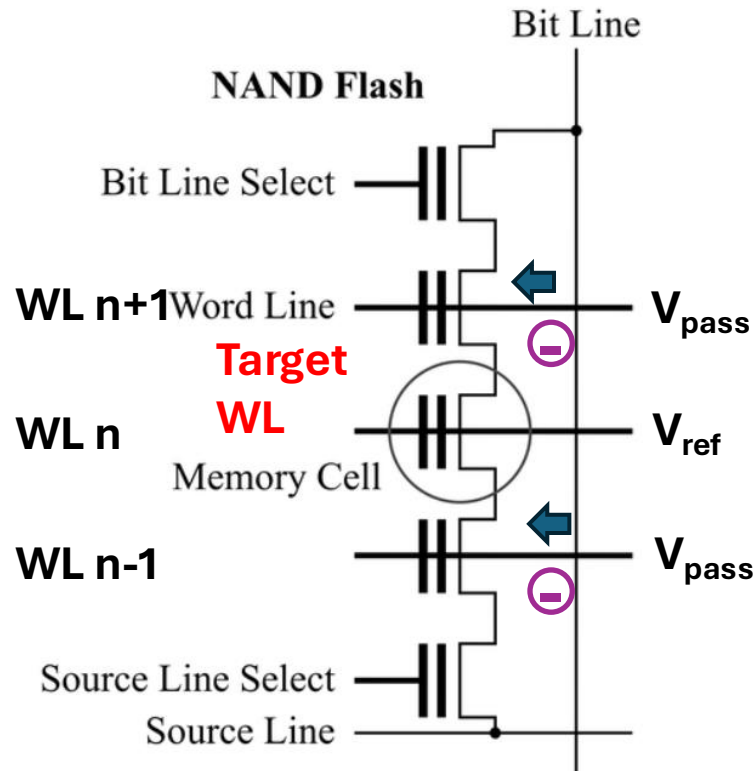
**CT Flash:** Fast lateral leakage → High real retention errors (underestimated by baking)

**FG Flash:** Isolated cells → Lower retention errors (better match with baking)

The answer is NO!

# Read Disturb: The Complete Story for 3D CT NAND

Read operations unintentionally disturb stored data in neighboring cells. This is a growing concern as layer counts increase.



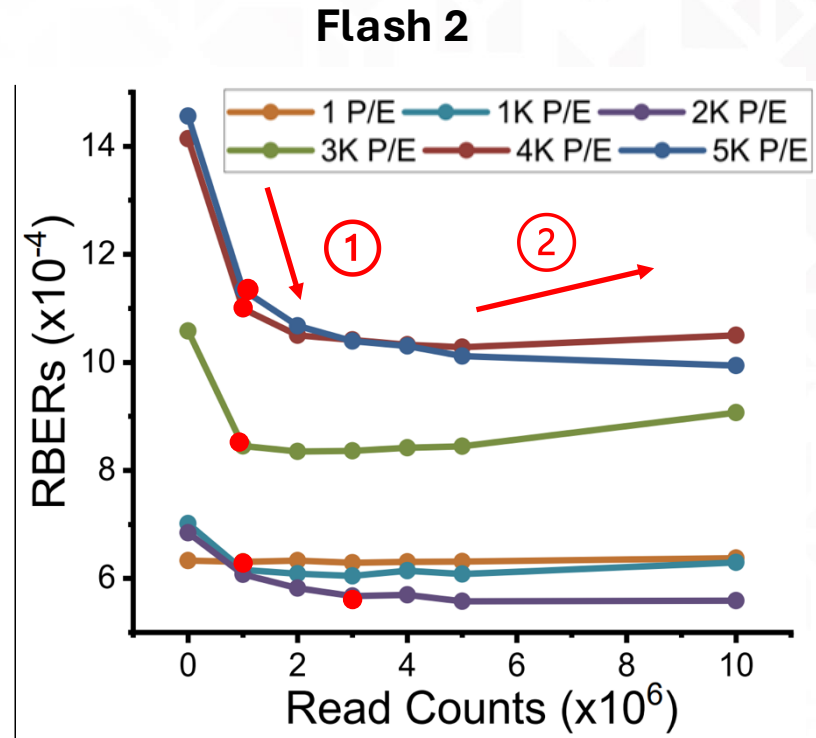
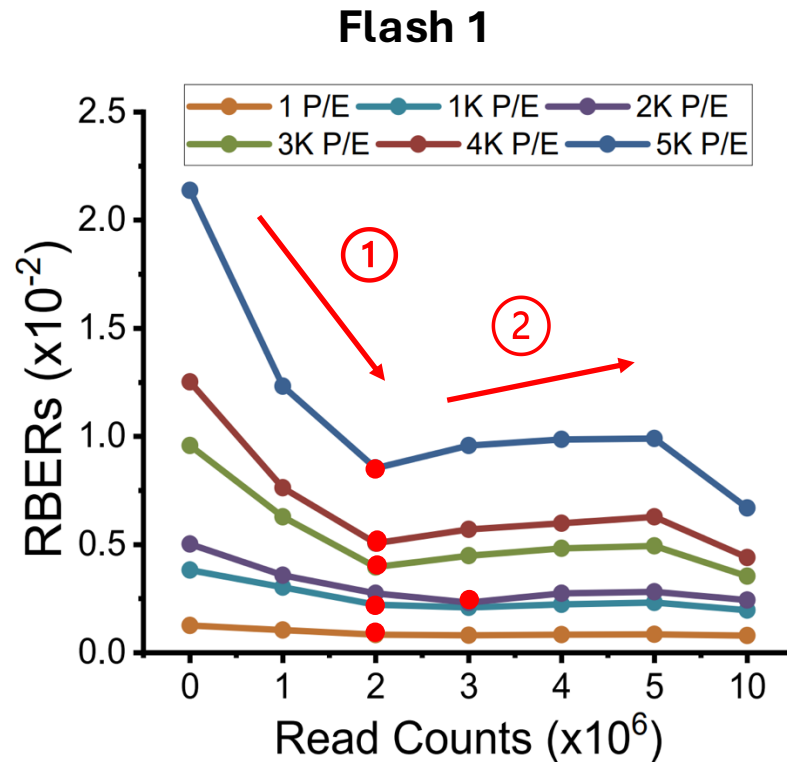
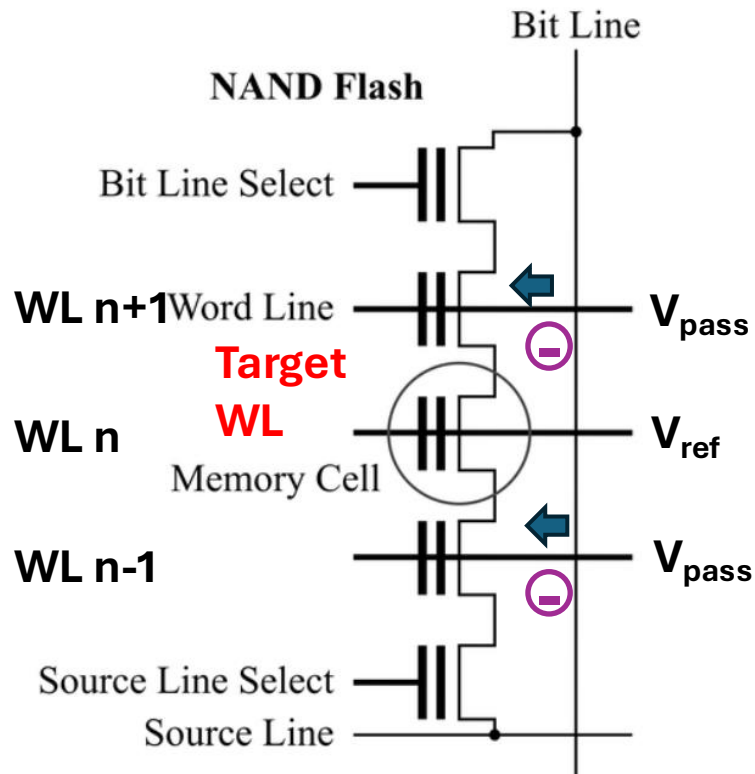
## New Insights

Read disturb behavior differs fundamentally from 2D NAND. It is location-dependent, asymmetric, and coupled with other error sources in complex ways.

Takeaway: Understanding read disturb's full picture in 3D CT NAND is essential for designing robust read management strategies.

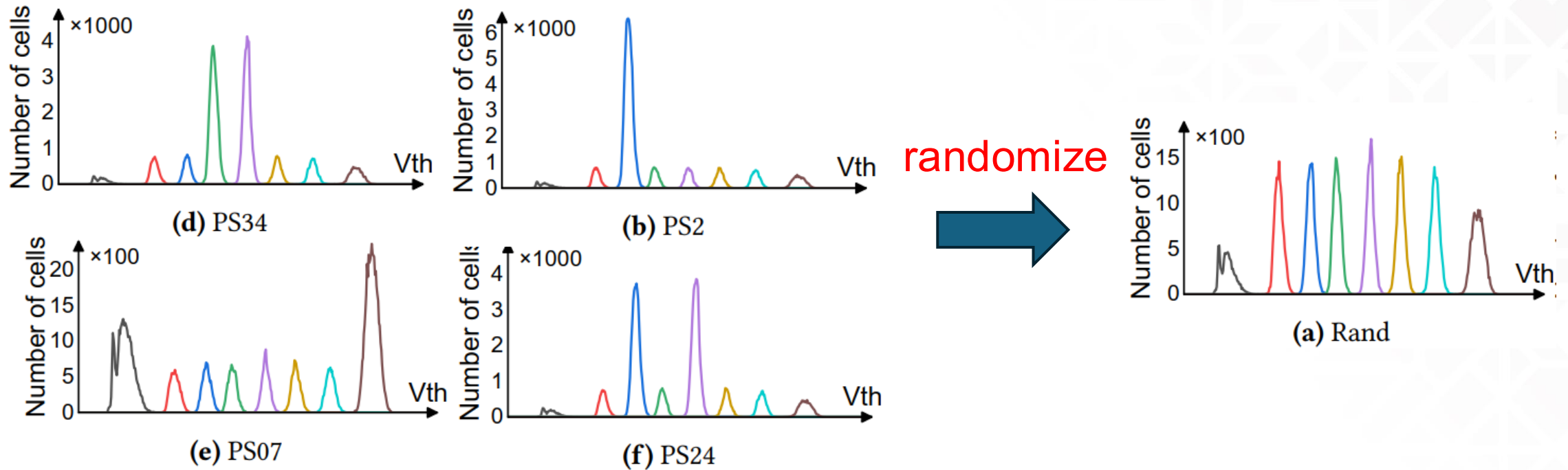
# Read Disturb: The Complete Story for 3D CT NAND

The read impacts as the **compensation** can **reduce** the RBERs at first, and than more reads will **increase** RBERs, the impacts of read will turn to be the **disturbance**.  
**Transformation from compensation to disturbance**



# Data Randomization

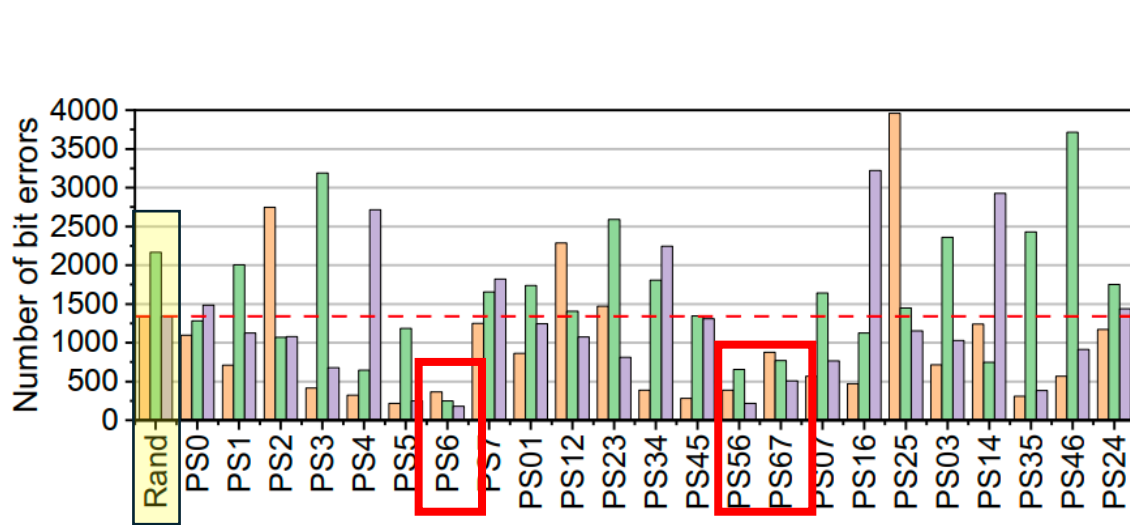
All NAND flash chips apply data **randomization** before programming to avoid extreme data patterns that generate worst-case RBER.



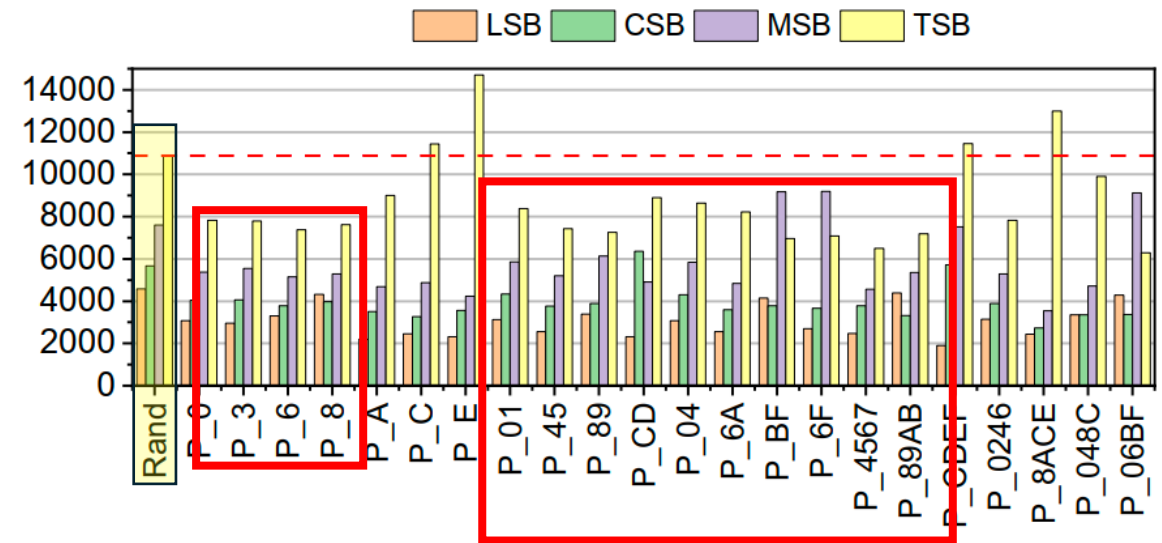
We tested the reliability of 24 different data patterns.

## Randomizer is suboptimal

Some data patterns outperform randomized data in RBER



(a) TLC Flash A (Retention)



(b) QLC Flash (Retention)

### ✓ What randomization does well

Prevents worst-case data patterns that cause extremely high raw bit error rate (RBER). Industry standard for decades.

### ✗ What randomization misses

Also eliminates best-case patterns that have very low RBER — a significant missed optimization opportunity for cold data.

## Part I — Key Insight

Deep physical characterization of 3D NAND flash reveals that many industry-standard assumptions are incomplete or incorrect. Every optimization must start from accurate hardware understanding.

*"Accurate understanding of the hardware is the prerequisite for effective optimization."*

Collectively these characterizations form the foundation of our optimization work

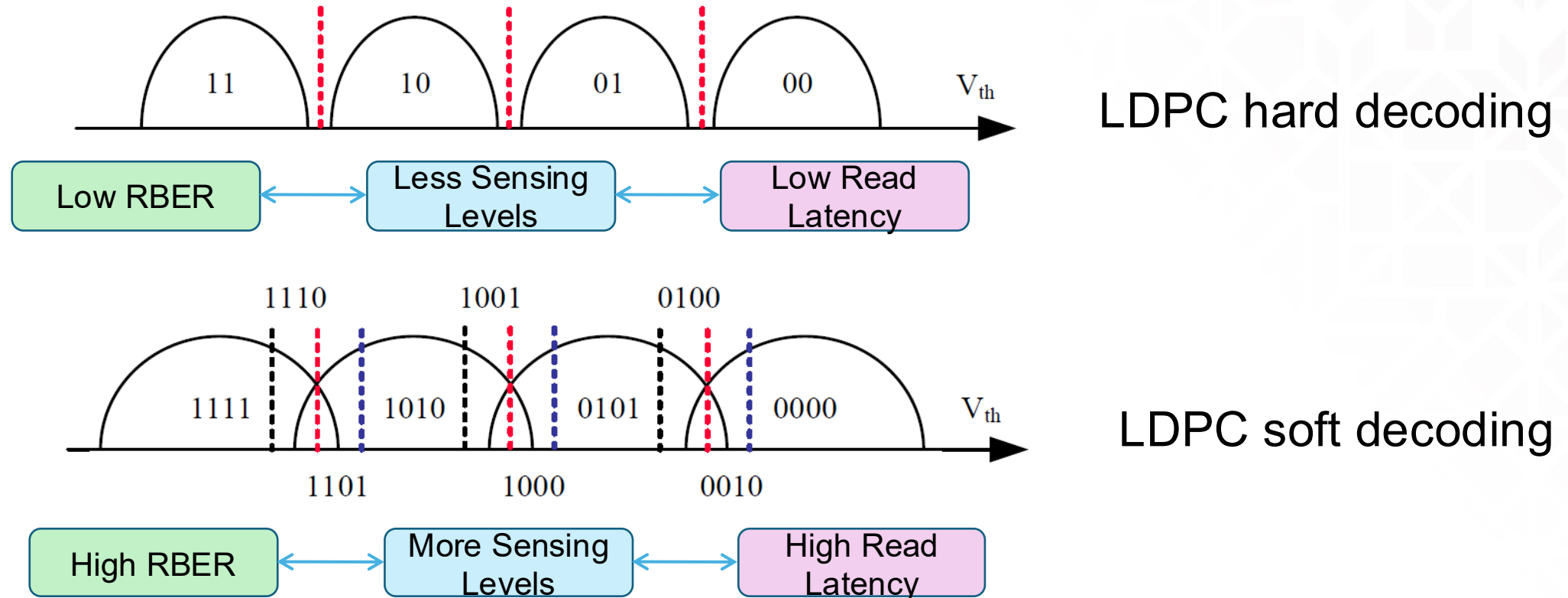
# Part II

## Boosting Read Performance

From LDPC optimization to achieving near-zero read retry

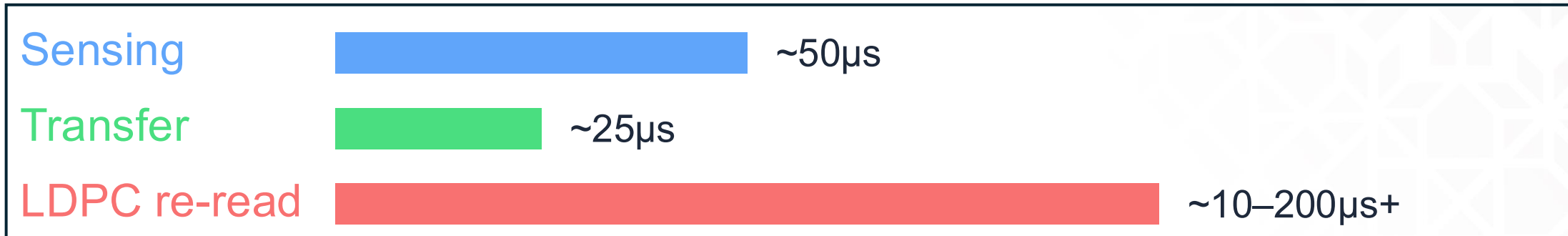
# LDPC Error Correction in NAND Flash

Modern high-density flash relies on LDPC (Low-Density Parity-Check) codes for error correction. But LDPC introduces significant read latency overhead.



Modern high-density flash relies on LDPC (Low-Density Parity-Check) codes for error correction. But LDPC introduces significant read latency overhead.

## Read Latency Breakdown:



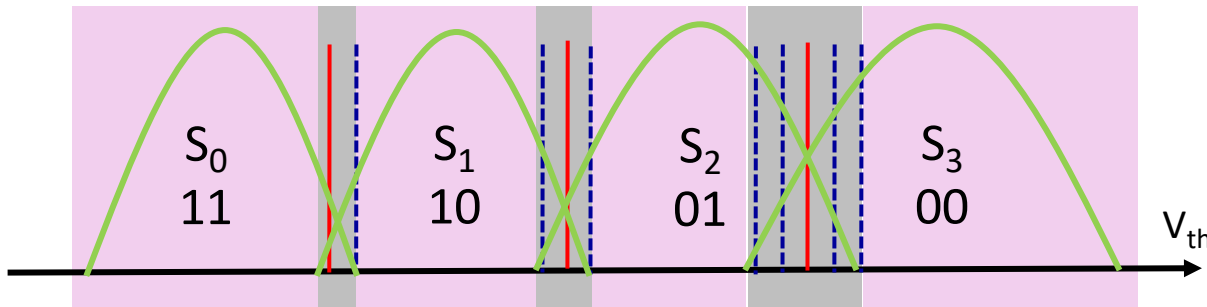
Problem: When LDPC fails at the default level, the SSD must perform read retries with adjusted voltage references — each retry adds 100µs+ of latency. This creates severe tail latency for storage systems.

*Our research tackles this from multiple angles: better sensing, better decoding, and ultimately eliminating retries.*

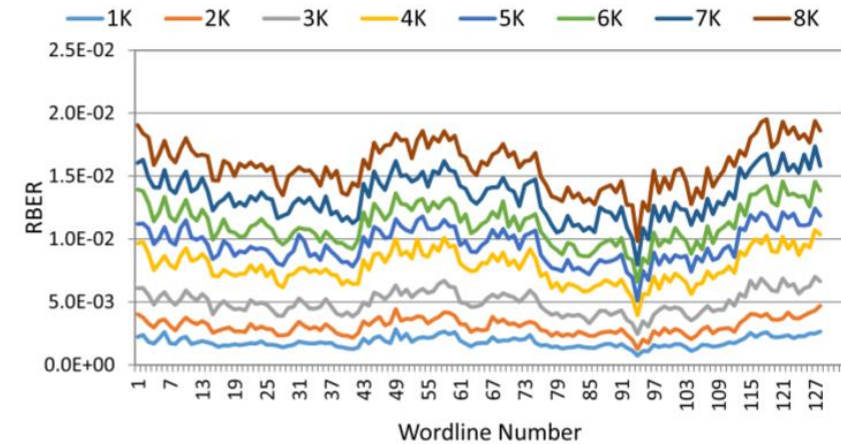
# Exploiting Error Asymmetry

Flash errors are not symmetric across voltage distributions.

By **placing sensing levels asymmetrically**, we minimize raw bit errors and dramatically improve LDPC success rate.



Asymmetry across states



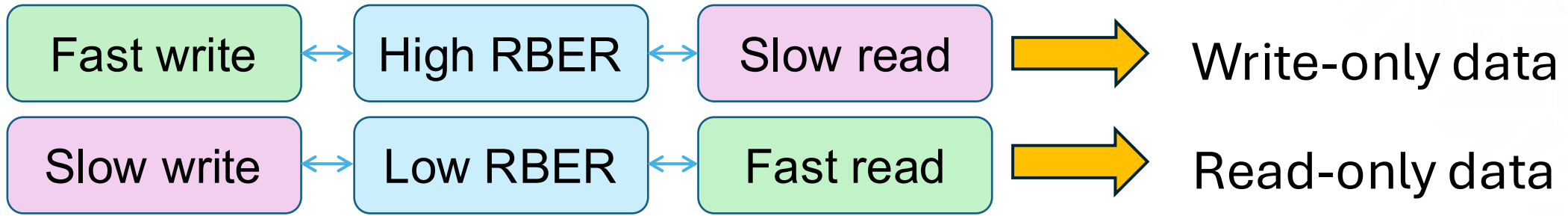
Asymmetry across layers

Key principles: exploiting physical error characteristics at the device level

# Exploiting Access Patterns



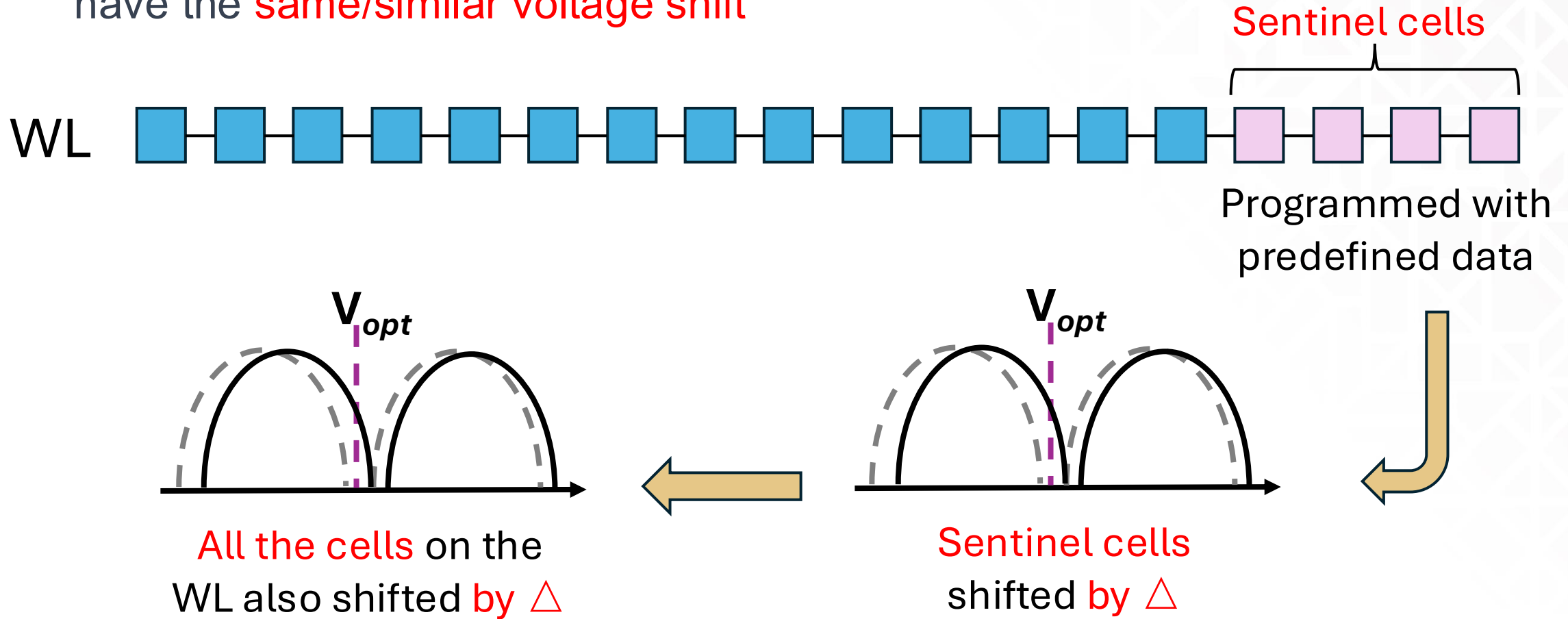
Most read requests occur on **read-only** data; Most write requests occur on **write-only** data



Key principles: exploiting data characteristics to guide flash optimization

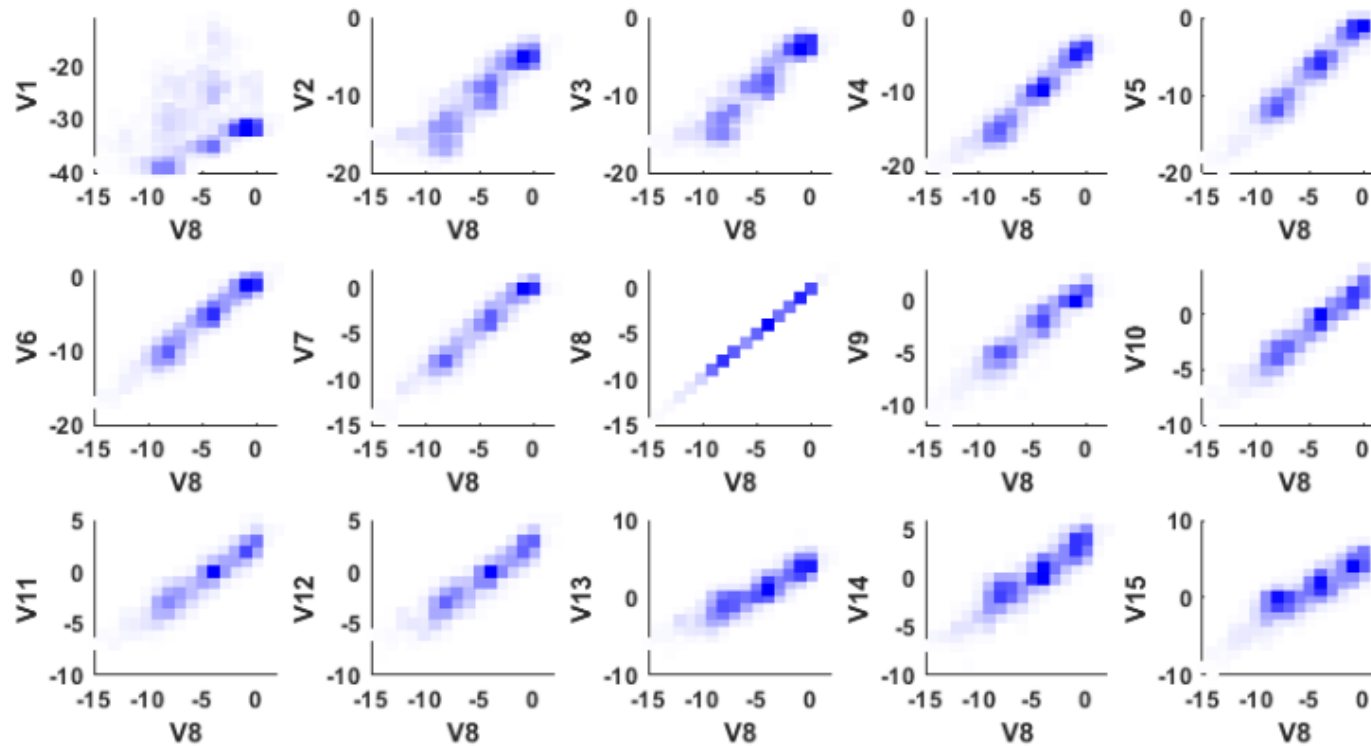
# Sentinel Cells: Shaving Retries for Fast Read

The cells on the same WL endure the same error conditions and tend to have the **same/similar voltage shift**



# Sentinel Cells: Shaving Retries for Fast Read

The optimal values of different read voltage have a linear relationship



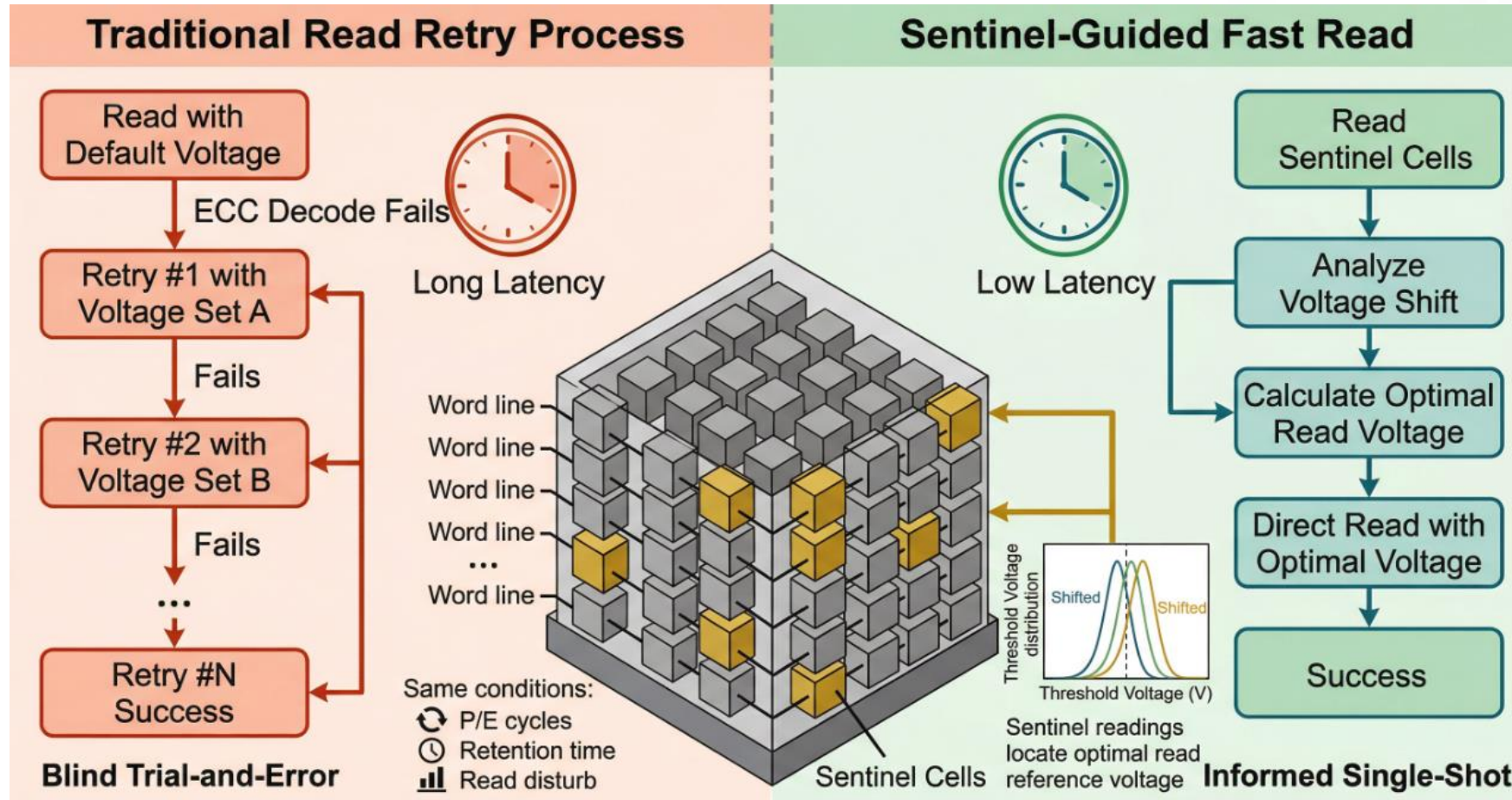
Obtaining the optimal value of one read voltage



Inferring the optimal value of other read voltages

# Sentinel Cells: Shaving Retries for Fast Read

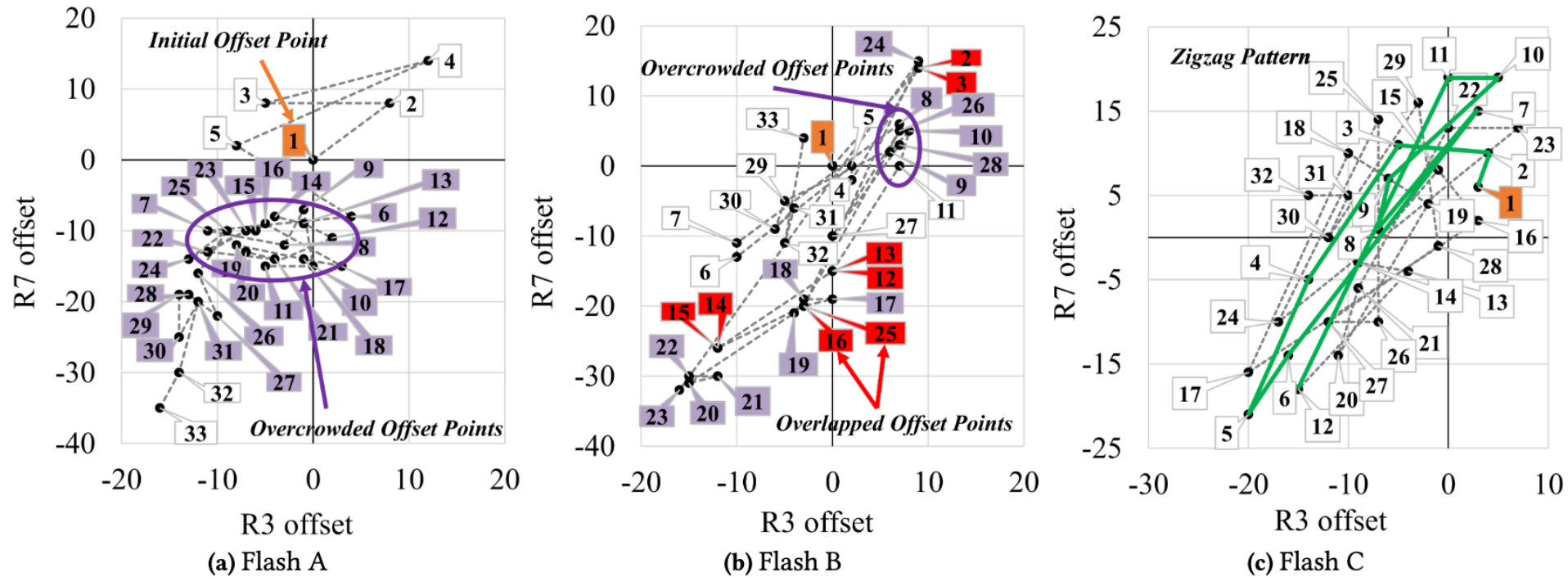
Core Idea: Designate a small set of flash cells as "sentinels" that serve as reliability indicators for the entire page.



81% reduction in read retries

# Achieving Near-Zero Read Retry

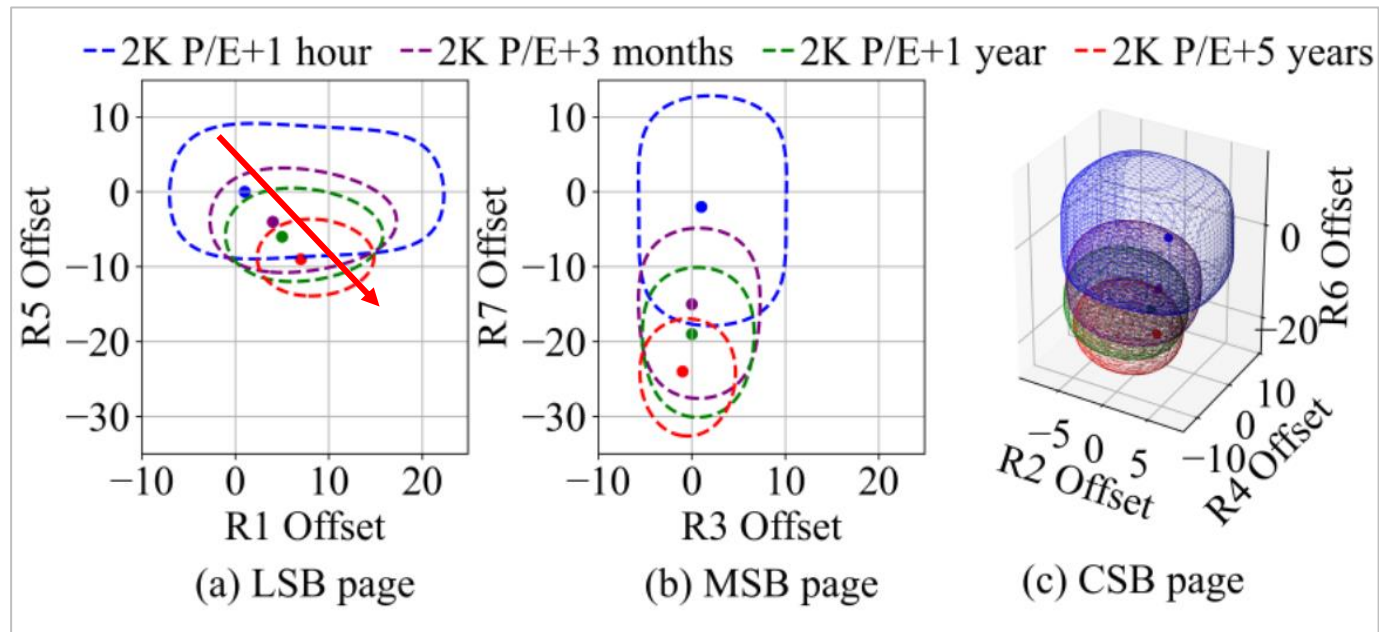
Building on sentinel work, this paper pushes toward the ultimate goal: virtually eliminating read retries in 3D NAND flash.



Current RRT in flash controller: **overcrowded, overlapped, zigzag**  
**RRT – Read Retry Table**

# Achieving Near-Zero Read Retry

Our work: **customize RRTs** based on the moving of RRCS  
First introduce **Read Retry Cover Space (RRCS)**



SOTA  
**15–30%**  
reads require retry



Our Approach  
**≈ 0%**  
near-zero read retry

Significance: Near-zero retry means consistent, predictable read latency — critical for latency-sensitive applications in data centers and cloud storage.

## Part II — Key Insight

Read performance evolution: from reducing LDPC iterations →  
sentinel-guided reads → near-zero retry.

Each layer of understanding enables the next breakthrough.

*" Each layer of understanding enables the next  
breakthrough in read performance."*



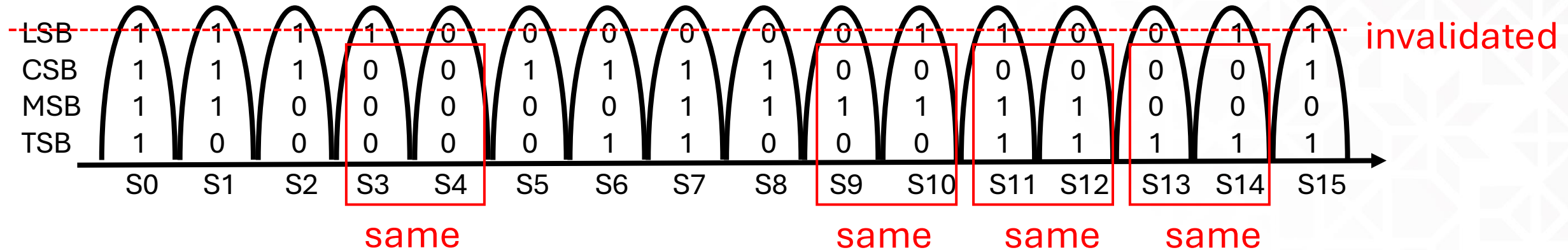
# Part III

## Smart Data Encoding

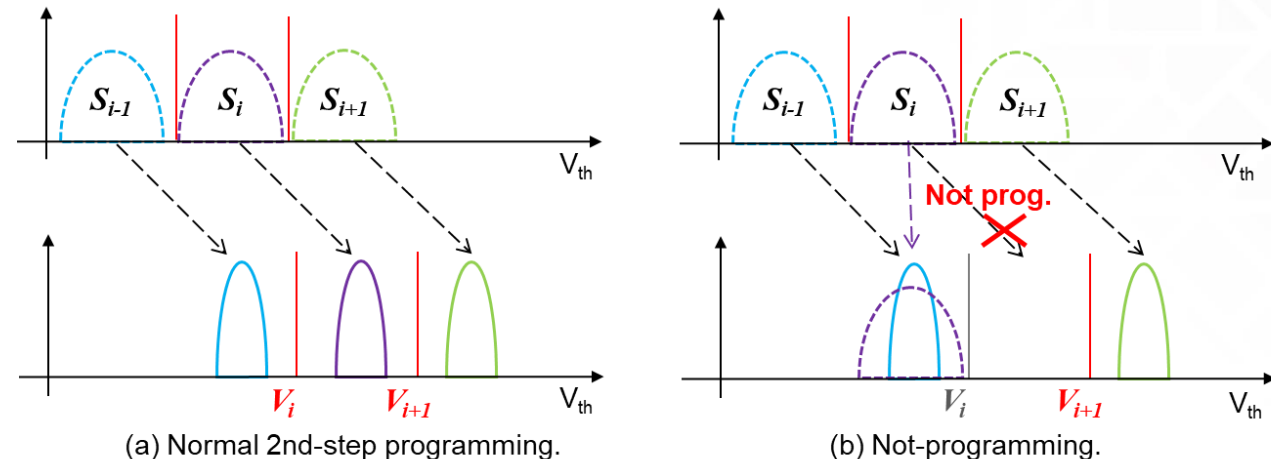
How data is encoded determines flash reliability and lifetime

# Midas Touch: Invalid Data Assists Flash

When data is logically invalidated (overwritten) before the second programming step, can we turn this "waste" into an advantage?

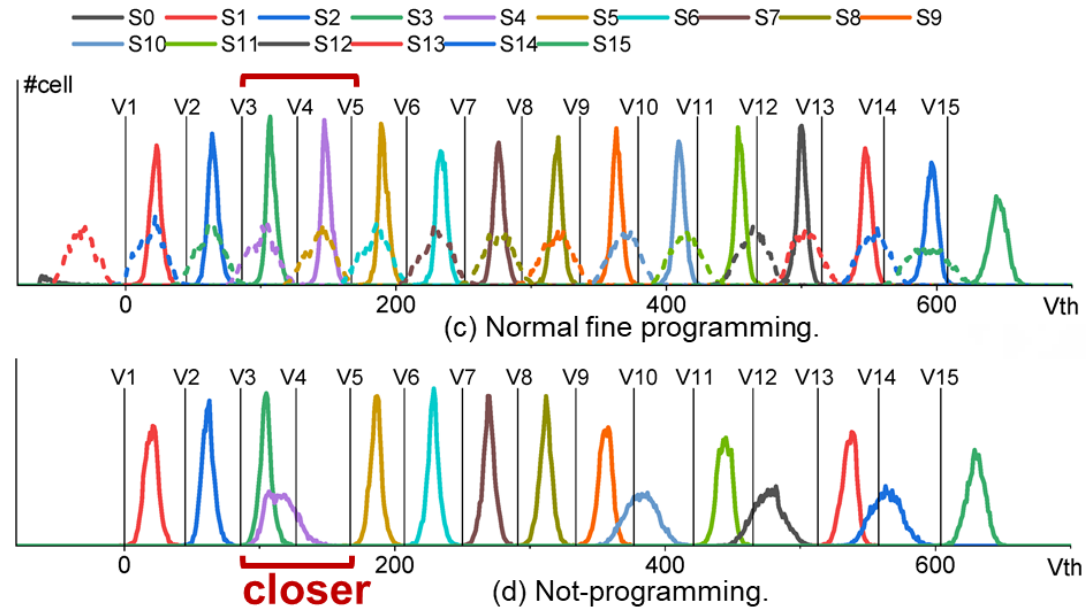


Key idea: **Enlarge the noise margin** between voltage states to **improve data reliability** of **valid pages**.



When data is logically invalidated (overwritten) before the second programming step, can we turn this "waste" into an advantage?

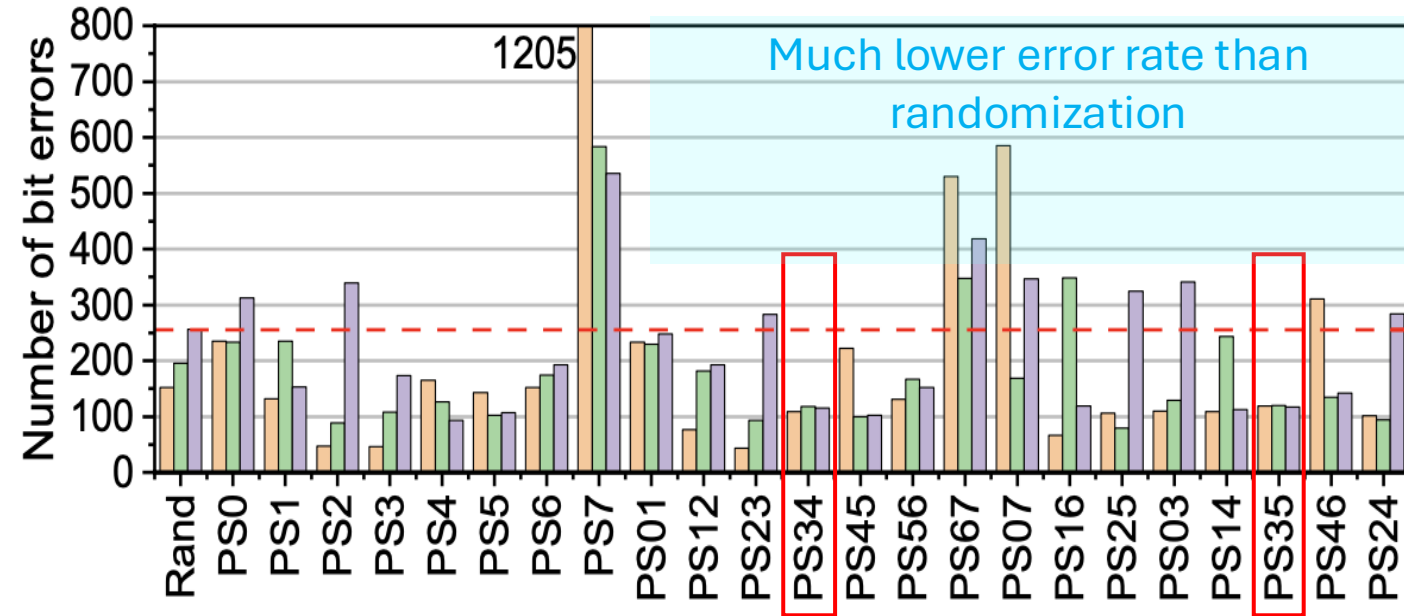
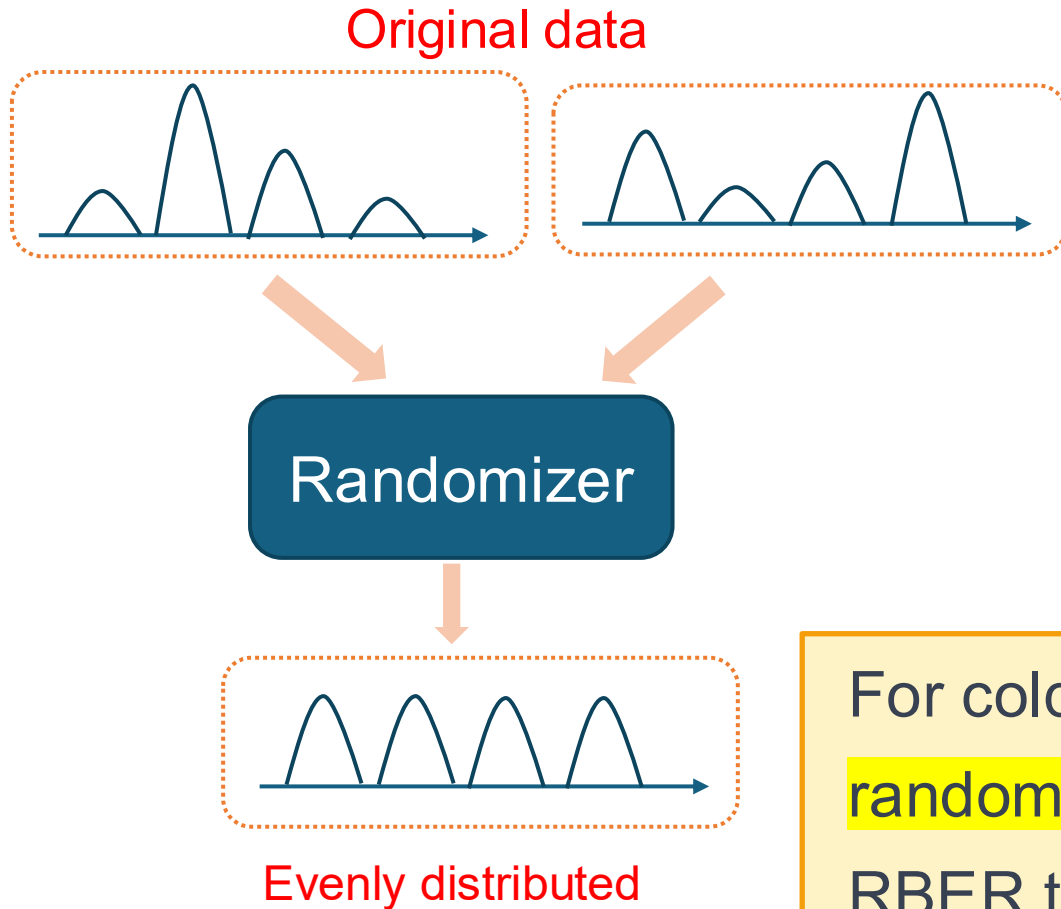
Real-device  
voltage  
distribution



Turns invalid data from a nuisance into an asset — "everything it touches turns to gold." Zero cost, pure gain.

# ColdCode: Data Encoding for Cold Storage

While **randomization** in NAND flash prevents worst-case RBER, it misses low-RBER patterns.



For cold data that won't be frequently rewritten, **replace randomization with carefully chosen encoding** to push RBER toward the best-case patterns.

# ColdCode: Data Encoding for Cold Storage

We propose **entropy-aware** encoding schemes for cold data: **skewed coding** and **reversed Huffman coding** for low- and high-entropy data, respectively.

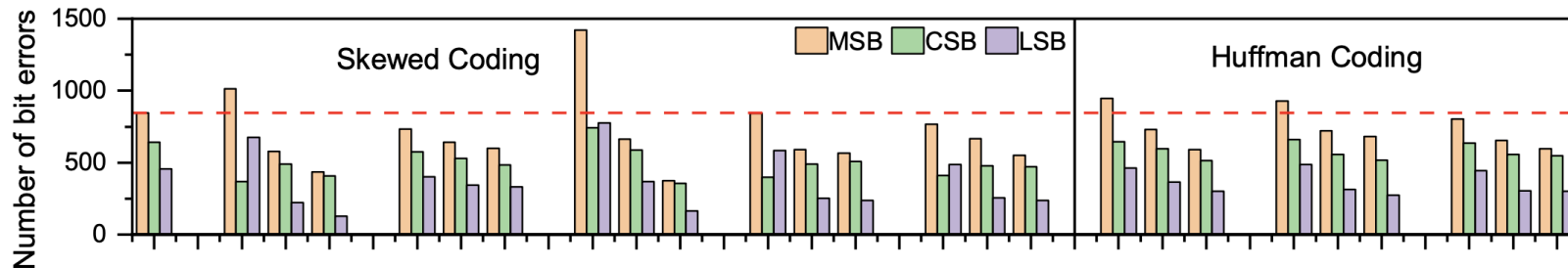
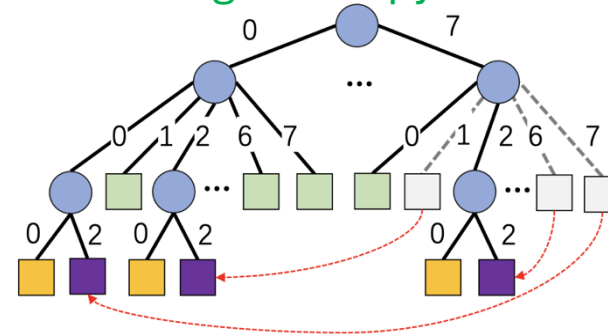
skewed coding for low-entropy data

01736415640173641640372564  
14420162740143516472016302

Frequent characters: 64 01

Desired characters: 20 02

Re-Huffman coding  
for high-entropy data



Evaluations on real-device demonstrate **42% RBER reduction** and **2.7X lifetime enhancement**.

## Part III — Key Insight

Data encoding is a powerful and under-explored lever. From invalid data patterns (Midas Touch) → cold data encoding (ColdCode) — each approach unlocks reliability and lifetime gains.

*"The bits you write determine how long your flash will last."*



# Part IV

## The Reprogramming Revolution

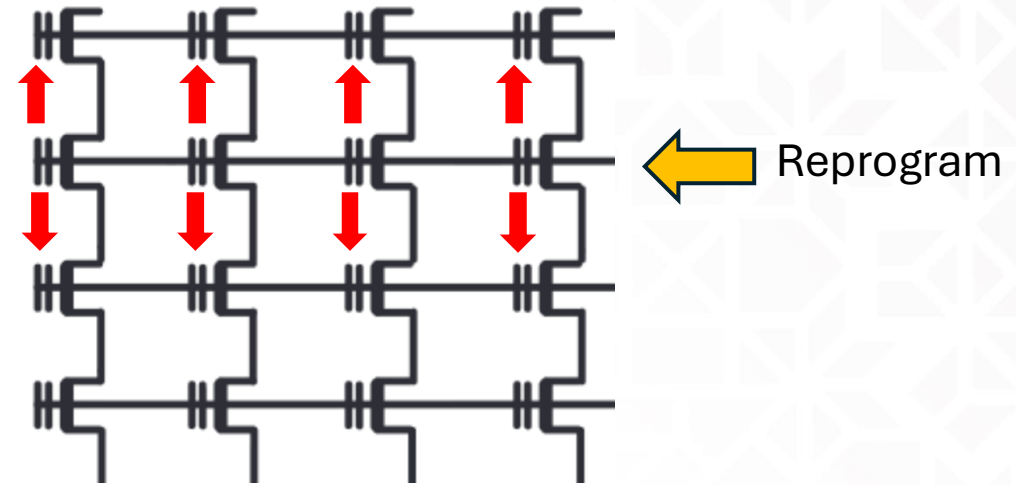
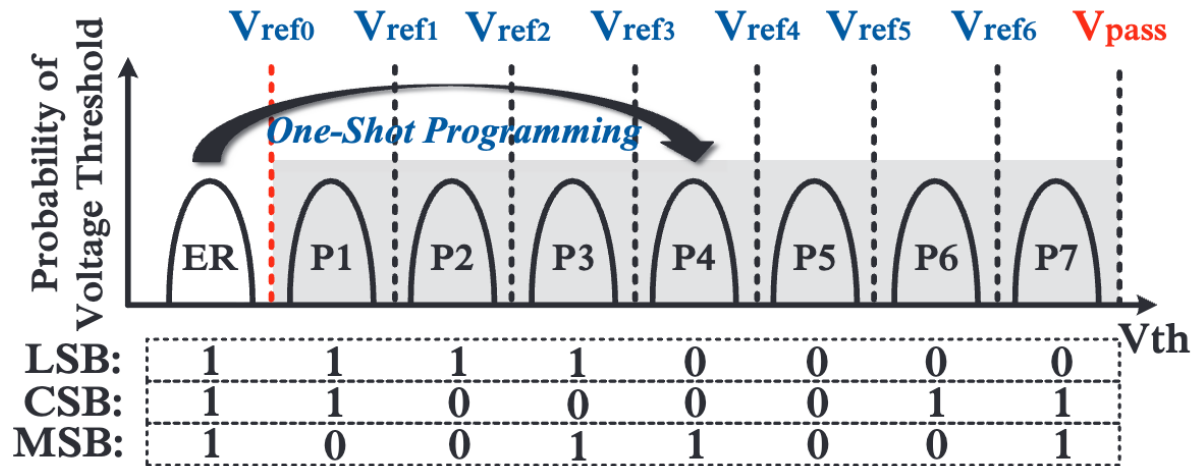
Breaking the sequential programming constraint of NAND  
flash

# ReSSD: Reprogramming 3D TLC Flash

NAND flash enforces strict sequential page programming.

Once a page is written, it typically cannot be modified.

But what if we could reprogram cells?

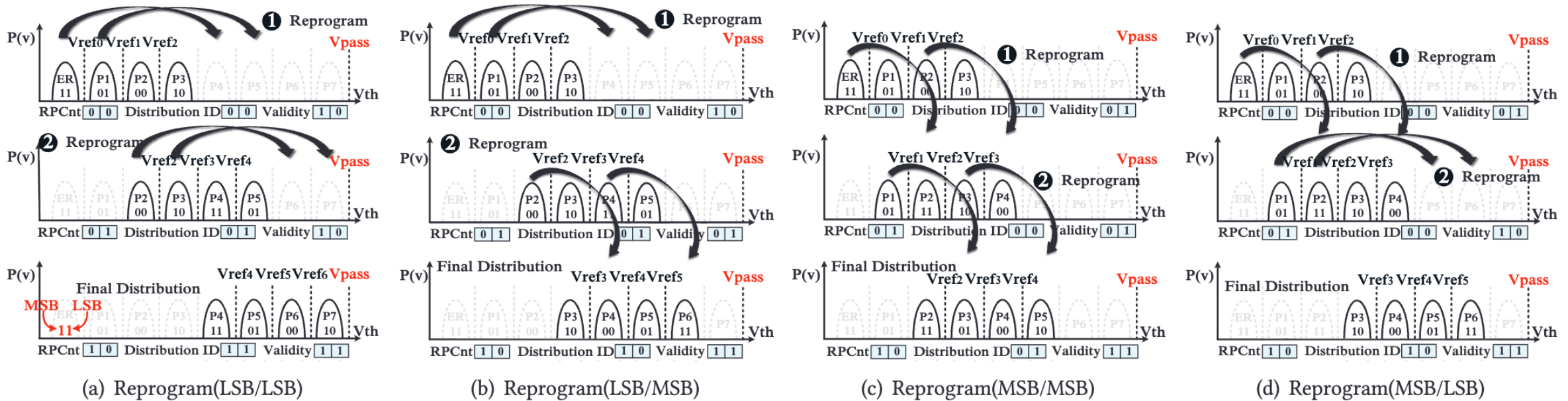


Severe interference to neighbors

For **update-intensive scenarios**, previous data is quickly invalidated in reprogrammable block

First program: use the first four states

Reprogram: one/two bits become invalid, use 1/2 more voltage states

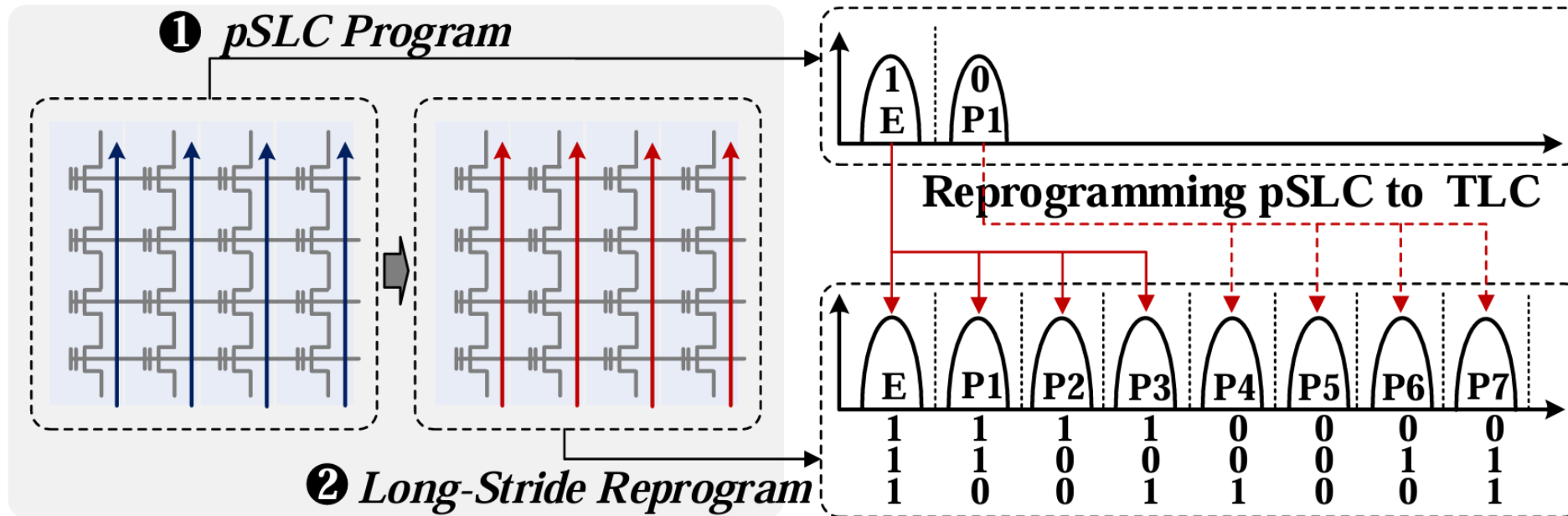


**Case study on RAID 5:** improve the endurance by 30.3%, boost write performance by 16.7%

# LOONG: Breaking the Stride Limitation

Prior reprogram work is limited: only a small subset of wordlines per block can be reprogrammed (small stride). LOONG expands this to entire blocks.

Key Innovation: **Spatially decouple the two program steps** across the entire block.



# LOONG: Breaking the Stride Limitation

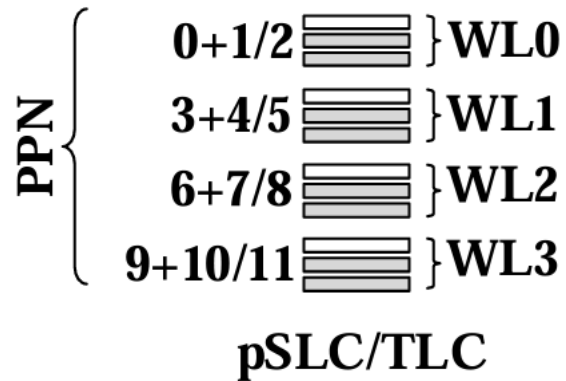
- 1

## Sequential SLC Program (All Wordlines)

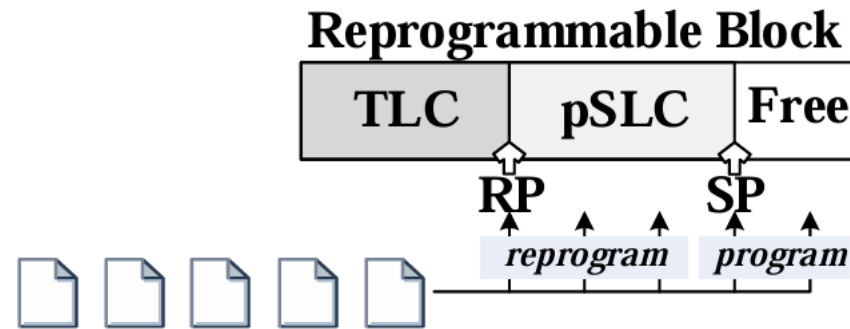
Program all WLs in a block in SLC mode (1 bit/cell). Fast and maximizing write throughput.
- 2

## Uniform TLC Reprogram (All Wordlines)

Reprogram all pages from SLC to TLC mode, preserving full storage capacity (3 bits/cell).



(a) PPN Assignment



(b) Active Page Pointers

SLC-speed writes → TLC-density storage. Best of both worlds.

Evaluated on two key SSD use cases

## Average Latency Reduction:

GC Acceleration:  **37.5%**

Program Opt.:  **18.1%**

### GC Acceleration

Valid page movement during garbage collection benefits from fast SLC writes + deferred TLC reprogram. Significant reduction in GC-induced latency spikes.

### Program Optimization

Normal write operations enjoy SLC-speed initial writes, with TLC conversion batched efficiently during idle periods.

## Part IV — Key Insight

Reprogram operations, once considered impractical, can harness Flash for major performance and lifetime gains: from small-stride reprogram → full-block reprogram.

*"Challenging fundamental constraints often yields the greatest rewards."*



# Part V

## Future Challenges & Opportunities

What's next for efficient and performant storage systems?

# The Cross-Layer Optimization Philosophy

**Our decade of work converges on a unified principle: effective flash optimization requires cross-layer co-design.**

**Application / File System**

Coldness tags, access patterns, data characteristics

**SSD Controller / FTL**

Encoding, scheduling, prediction, LDPC management

**NAND Flash Chip**

Program modes, error characteristics, reprogram operations

**Device Physics**

Charge trap behavior, process variation, self-healing effects

↕ **Information flows across all layers** ↕

Every breakthrough came from understanding one layer deeply and connecting it with another.

## QLC/PLC Scaling

4–5 bits/cell demand exponentially tighter margins. New encoding, error correction, and management strategies are essential for continued scaling.

## AI-Driven Storage Management

Machine learning for LDPC prediction, wear leveling, and workload-adaptive management. Moving from heuristics to intelligent, self-optimizing systems.

## CXL & Composable Storage

New memory/storage interfaces change the optimization landscape. Co-optimizing flash with CXL-attached memory opens new design spaces.


## Sustainability & Green Storage


Extending flash lifetime directly reduces e-waste and energy consumption. Encoding and self-healing contribute to environmentally sustainable storage.


## Workload-Specific Flash Optimization


AI training data, genomics, financial data — each has unique access patterns and reliability needs. One-size-fits-all approaches are leaving significant performance and lifetime on the table. The CHEOPS community is uniquely positioned to address this.


# Key Takeaways

 **Deep physical characterization is the foundation — standard assumptions often mislead.**

 **Read performance can be pushed to near-zero retry through sentinel-guided approaches.**

 **Smart data encoding (ColdCode) can extend lifetime by up to 3× with no hardware changes.**

 **Full-block reprogram (LOONG) reduces latency by up to 37.5%, deployable via firmware.**

 **Cross-layer co-design is the unifying principle — every optimization bridges layers.**



Mohamed bin Zayed  
University of  
Artificial Intelligence

# Thank You!

Questions & Discussion

---

**Chun Jason Xue**