**Technion**

Electronics
Computers
Communications

# NAND Flash Architectures Reducing Write Amplification Through Multi-Write Codes

**Saher Odeh**

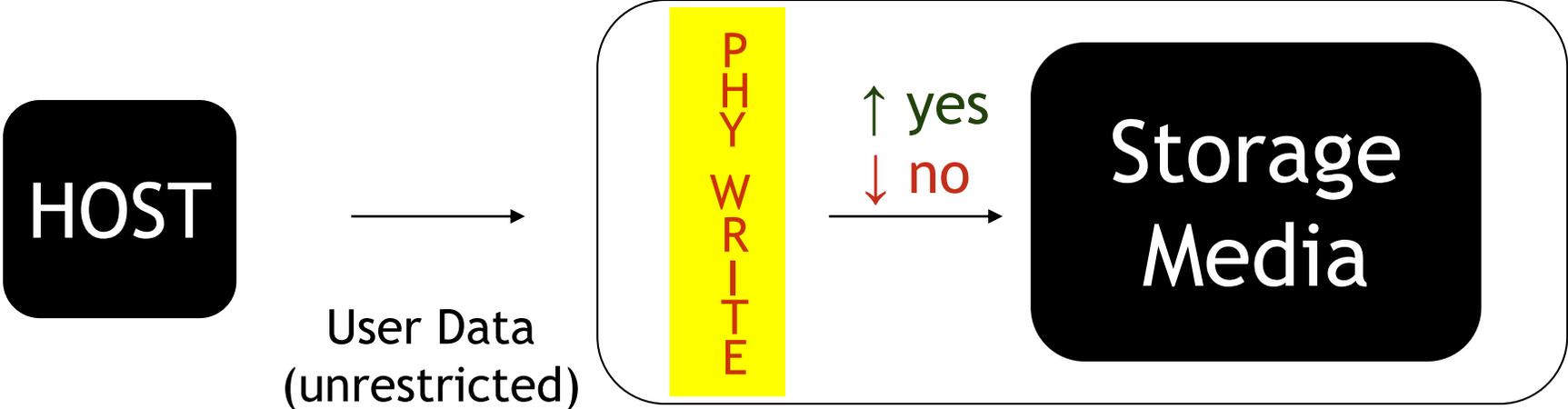Yuval Cassuto

Technion-EE

MSST 2014

June 6th

# Main Objective

► Improve device performance - speed, power and lifetime

► Reduce device cost per bit (silicon space)
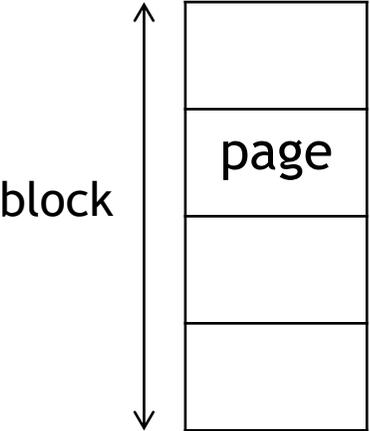
► Require minimal / no changes to NAND storage unit

Method:

► Propose architectures that cleverly use **Multi Write Codes**

# Problem: No In-Place Writes

HOST

User Data
(unrestricted)

PHY WRITE

↑ yes
↓ no

Storage
Media

- Page - write / read unit

- Block – erasure unit

block

page

# Write Amplification (WA)

- **Malady:** Write-Amplification
  - **Adding** spare blocks **reduces** WA (over-provisioning)
- Example:
  - Writes:

| 0 | 4 | 8 | 12 | |
|---|---|----|----|---|
| 1 | 5 | 9 | 13 | |
| 2 | 6 | 10 | 14 | |
| 3 | 7 | 11 | 15 | |

[Agarwal et. al 2010, Desnoyers 2012,...]

# Write Amplification (WA)

- **Malady:** Write-Amplification
  - **Adding** spare blocks **reduces** WA
- Example:
  - Writes: 0,4,8,12,

| ~~0~~ | ~~4~~ | ~~8~~ | ~~12~~ | 0 |
|---|---|---|---|---|
| 1 | 5 | 9 | 13 | 4 |
| 2 | 6 | 10 | 14 | 8 |
| 3 | 7 | 11 | 15 | 12 |

[Agarwal et. al 2010, Desnoyers 2012,...]

# Write Amplification (WA)

- **Malady:** Write-Amplification
  - **Adding** spare blocks **reduces** WA
- Example:
  - Writes: 0,4,8,12,1

| 2 | ~~4~~ | ~~8~~ | ~~12~~ | 0 |
|---|---|---|---|---|
| 3 | 5 | 9 | 13 | 4 |
| 1 | 6 | 10 | 14 | 8 |
| | 7 | 11 | 15 | 12 |

[Agarwal et. al 2010, Desnoyers 2012,...]

# Multi-Write Codes

- Remedy: Multi-Write Codes
  - Up to **t > 1** writes before erase!
- Example t = 2

| 0 | 4 | 8 | 12 | |
|---|---|---|---|---|
| 1 | 5 | 9 | 13 | |
| 2 | 6 | 10 | 14 | |
| 3 | 7 | 11 | 15 | |

# Multi-Write Codes

▶ **Remedy:** Multi-Write Codes

  ▶ Up to **t** writes before erase!

▶ Example t = 2

  ▶ Writes: 0,1,2,4,5,6,8,9,10,12

| 0 | 4 | 8 | 12 | |
|---|---|---|----|---|
| 1 | 5 | 9 | 13 | |
| 2 | 6 | 10 | 14 | |
| 3 | 7 | 11 | 15 | |

▶ Invalid

▶ Written twice

5

# Multi-Write Codes

▶ Remedy: Multi-Write Codes

    ▶ Up to **t** writes before erase!

▶ Example t = 2

    ▶ Writes: 0,1,2,4,5,6,8,9,10,12,0,4,8,12

| | | | | | |
|---|---|---|---|---|---|
| ~~0~~ | ~~4~~ | ~~8~~ | ~~12~~ | | 0 |
| 1 | 5 | 9 | 13 | | 4 |
| 2 | 6 | 10 | 14 | | 8 |
| 3 | 7 | 11 | 15 | | 12 |

| | |
|---|---|
| ~~4~~ | ▶ Invalid |
| 5 | ▶ Written twice |

# Multi-Write Codes

- **Remedy:** Multi-Write Codes
  - Up to **t** writes before erase!
- Example t = 2
  - Writes: 0,1,2,4,5,6,8,9,10,12,0,4,8,12,1
  - # Writes : 4 → 14

| | | | | |
|---|---|---|---|---|
| **2** | ~~4~~ | ~~8~~ | ~~12~~ | **0** |
| **3** | **5** | **9** | 13 | **4** |
| **1** | **6** | **10** | 14 | **8** |
| | 7 | 11 | 15 | **12** |

| | |
|---|---|
| ~~4~~ | ▶ Invalid |
| **5** | ▶ Written twice |

# The Tradeoff of Multi-Write Codes

▶ More in-place writes → higher code redundancy

▶ Code redundancy shrinks spare available for WA reduction

▶ Research question: can Multi-Write codes reduce overall WA?

# t=2 Multi-Write Codes

▶ Q: How big is the physical page for t=2 writes? (expansion-factor)

  ▶ A: capacity-achieving WOM codes:

$$|\text{physical page}| = \frac{2\log_2 q}{\log_2\binom{q+1}{2}} |\text{logical page}|$$

| q=2  SLC |
|----------|
| q=4  MLC |
| q=8  TLC |

▶ For t=2,q=2; expansion of data: ~1.26

▶ For t=2,q=8; expansion of data: ~1.16

Higher q → lower expansion

[Luojie, Kurkoski, Yaakobi, 2012]

# Prior Work

▶ Multi-write + compression

[Jagmohan, Franceschini, Lastras, 2010]

▶ Analysis of WA w/ multi-write codes

[Luojie, Kurkoski, Yaakobi, 2012]

# Known Results
### [Luojie, Kurkoski, Yaakobi, 2012]

- Multi Write codes reduce WA when
    1. Spare > 75% for SLC (q=2)
    2. Spare > 55% for MLC (q=4)
    3. Spare > 40% for TLC (q=8)

Cost likely too high for SSD deployment

# Partial Re-Write

▶ A natural idea:

t=2 →

t=1 ⟶

no re-write, no overhead

▶ <u>The challenge:</u>

Which writes should get re-write capabilities?

<u>The Answer</u>: the writes that will be rewritten before erasure.

# Re-Write Differentiation

▶ The Policy:

  ▶ Perform t=2 coding for all incoming user-writes

  ▶ Write without coding (t=1) all garbage-collection writes

▶ Observation:

  ▶ With high likelihood, user writes will be rewritten soon (temporal-locality)

# Double-Fronted Architecture



t=1

frntG

hotBlocksCount

t=2

User Writes

frntH

GC Write-backs

GC (Empty Pool)

Invalid

Hot State

Cold State

# Selective (Single-Fronted) Architecture

User Writes

frnt

GC Write-backs

GC (Empty Pool)

Invalid

Hot

Cold

# Evaluation

# Trace Results - TLC



PRXY q=8

# Trace Results - SLC
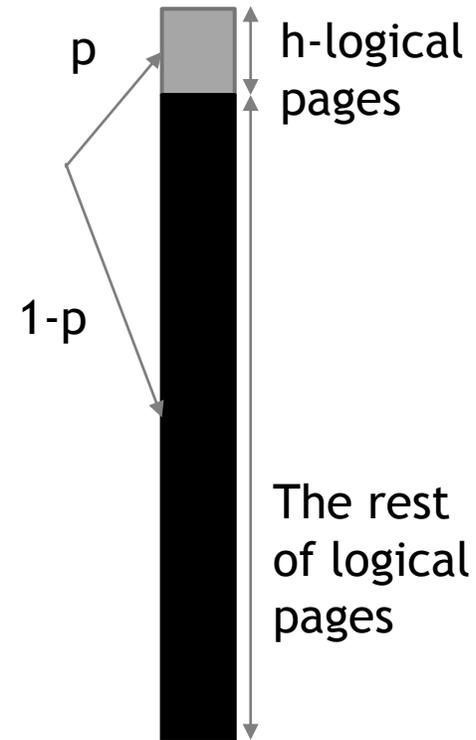


PRXY q=2

REG_RW1    REG_RW2    Selective_RW2    D_FRONT
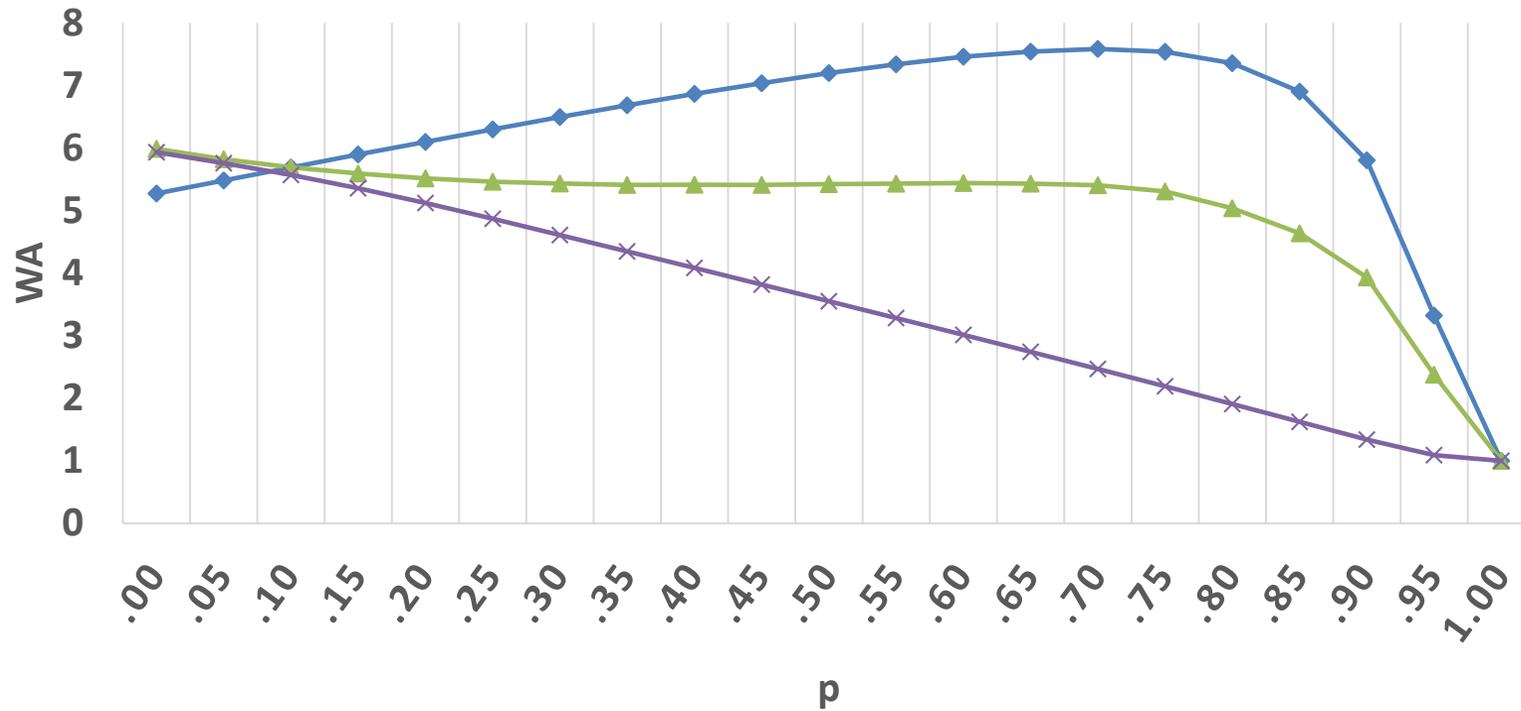
# Synthetic Workloads: p-Locality

- ▶ p – the probability for a "hot" write
  - ▶ Write is chosen uniformly from a pool of hot pages
  - ▶ h – the size of the hot pool

- ▶ 1-p – the probability for a "cold" write
  - ▶ Write will be chosen from the rest of the pages
  - ▶ The "cold" page is turned to "hot"



p

1-p

h-logical pages

The rest of logical pages

# p-Locality Results



SYNTHETIC q=8 SPARE=10%

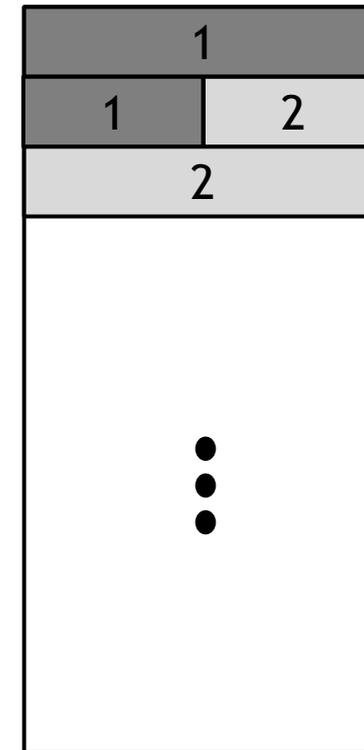# Real Implementation

- Problem: Variable Page Size

    - some pages use t=2 (larger).

    - some use t=1 (smaller).

- Solution:

    - w/o flash vendors support

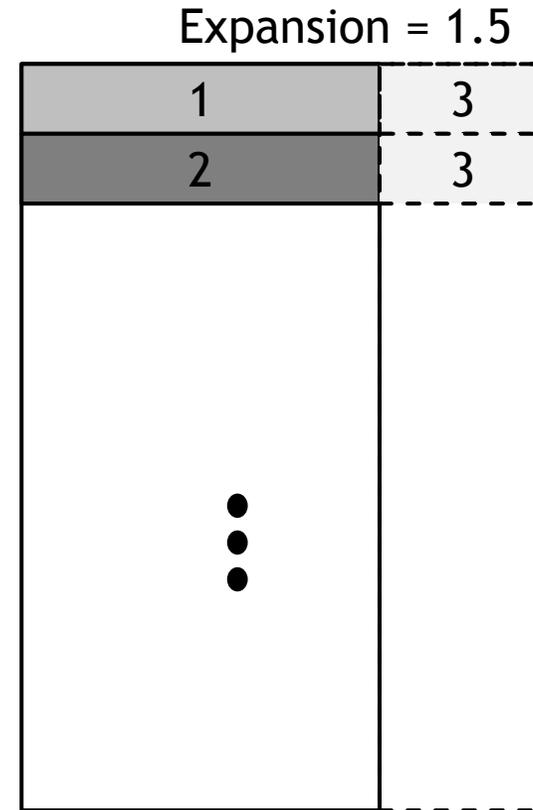    - w/   flash vendors support

# w/o Flash Vendors Support 1

Expansion = 1.5

- ▶ Normal Page Allocation
  - ▶ Normal pages – no change
  - ▶ Expanded pages – grouped together

- ▶ Read penalty
- ▶ In-place write → RMW

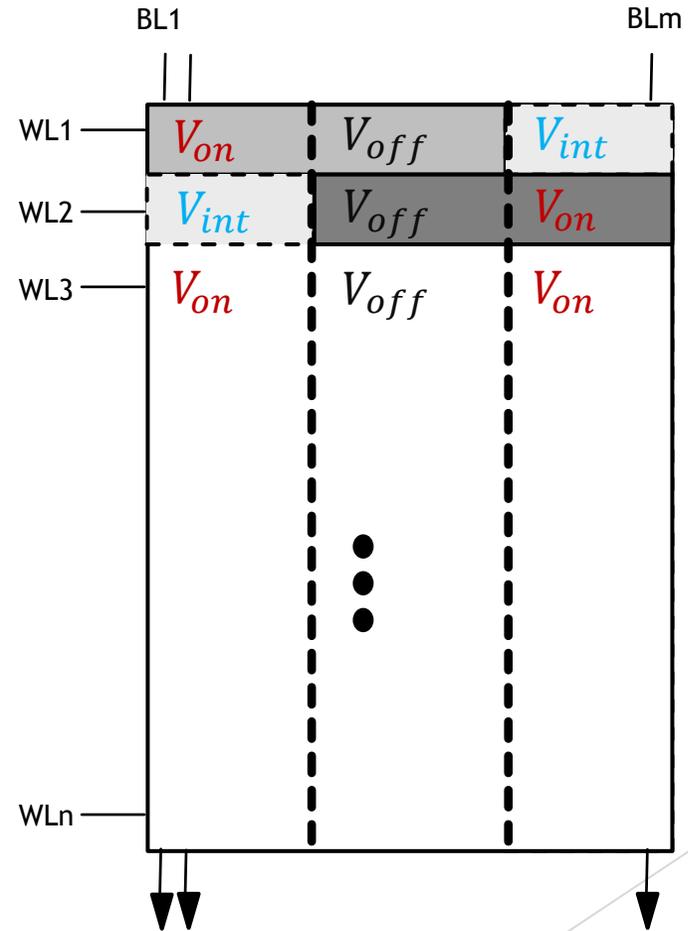# w/o Flash Vendors Support 2

Expansion = 1.5

- ▶ Expanded Page Allocation
  - ▶ Expanded pages – no change
  - ▶ Normal pages – grouped together

- ▶ GC pages are grouped and written together
  - ▶ Write (Scatter) - no penalty
  - ▶ Read (Gather) – read penalty
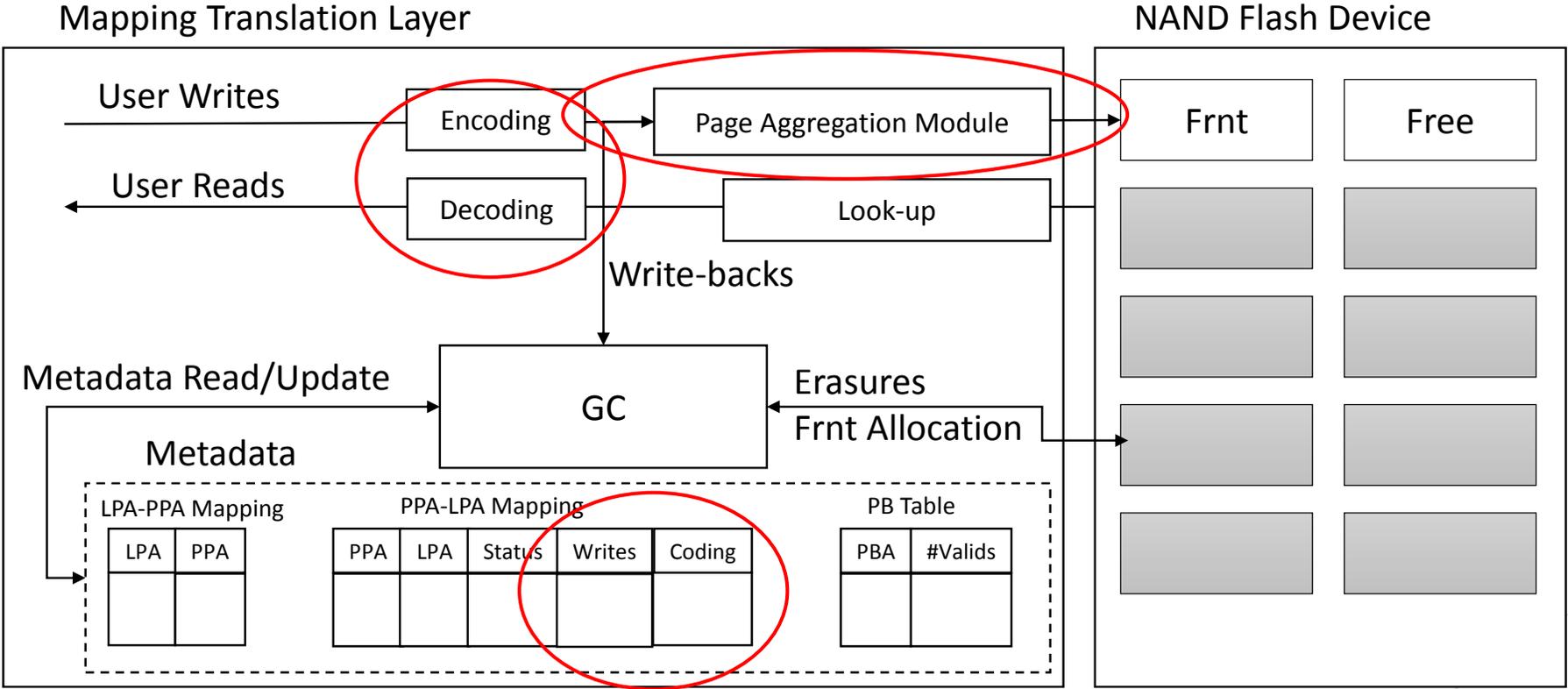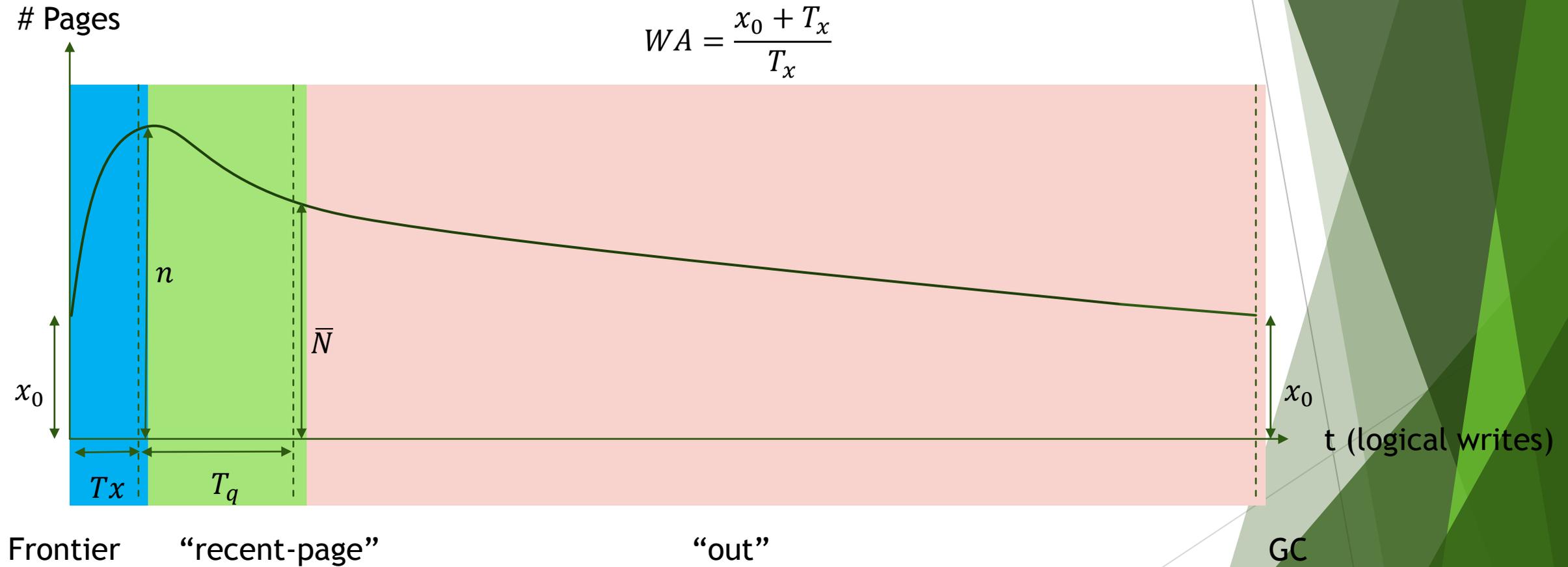
# w/ Flash Vendors Support

- Split pages are shifted
  - Parallel circuit read

- **No** penalty/ overhead

- Vendor support: Split word-line read

# Device Diagram

# The Block Life Cycle - Analysis

# Conclusion

▶ New scheme to improve WA with selective MW-writes

▶ Insensitive to specific workload properties

▶ Good potential for SSD deployment