# Predicting SSD Performance for Today's Dynamic Workloads

Shirish Bahirat, SSD Engineering System Architecture

**Micron**®

# Legal Disclaimer

- Presentation material is just conceptual representation of technicalities and for educational purpose only

    - Data represented in this presentation is not attributed to any product or application

    - Material contained in this presentation is copyright of Micron Technology

- Nothing in this presentation is considered as legal advice or official opinion of Micron Technology or presenter

    - Presenter do not assume any legal responsibility or liability arising due to use of any information in this presentation

    - No warranties, or no implied advice, use at your own risk

August 8, 2017     |   Micron Technology

# Contains

- Write amplification – basic terminology

  - Steady state write amplification

- Time domain performance

  - Why steady state write amplification analysis not enough?

- Components impacting dynamic write amplification

  - Logical saturation and degree of randomness – fresh out of box, streams, hot vs cold data, trim/unmap

  - Workload transitions – mix of sequential and random data

August 8, 2017    | Micron Technology

# Write Amplification

# Workload Visualization

**Initial state** → **Sequential Write**

| | | | | |
|---|---|---|---|---|
| B1 | B2 | B2 | B4 | B5 |
| B6 | B7 | B8 | B9 | B10 |
| B11 | B12 | B13 | B14 | B15 |
| B16 | B17 | B18 | B19 | B20 |
| B21 | B22 | B23 | B24 | B25 |
| OP | OP | OP | OP | OP |
| SP | SP | SP | SP | SP |

| | | | | |
|---|---|---|---|---|
| x | x | x | x | x |
| B6 | B7 | B8 | B9 | B10 |
| B11 | B12 | B13 | B14 | B15 |
| B16 | B17 | B18 | B19 | B20 |
| B21 | B22 | B23 | B24 | B25 |
| B1 | B2 | B3 | B4 | B5 |
| SP | SP | SP | SP | SP |

**Initial state** → **Random Write**

| | | | | |
|---|---|---|---|---|
| B5 | B18 | B10 | B21 | B14 |
| B15 | B25 | B3 | B6 | B19 |
| B24 | B20 | B12 | B2 | B9 |
| B17 | B1 | B7 | B13 | B23 |
| B4 | B11 | B22 | B16 | B8 |
| OP | OP | OP | OP | OP |
| SP | SP | SP | SP | SP |

| | | | | |
|---|---|---|---|---|
| B5 | x | B10 | B21 | B14 |
| B15 | B25 | B3 | x | B19 |
| B24 | B20 | B12 | B2 | B9 |
| x | B1 | B7 | B13 | x |
| B4 | x | B22 | B16 | B8 |
| B17 | B6 | B18 | B23 | B11 |
| SP | SP | SP | SP | SP |

NAND filled with sequential data, OP region empty

SP is spare region

Incoming sequential data written in OP region, entire block stripe validated

> Host Data Written = 1
> GC Data Written = 0

Drive filled with random data, OP region empty

Host random data written in OP region, 3 TU's need to be moved to create empty space

> Host Data Written = 1
> GC Data Written = 0.6

# Write Amplification (WA)

- Physical (Written to NAND) data is more than logical (Written by Host) data

  - NAND Flash memory pages must erased before re-written (Erase Unit - Block)

  - NAND Flash device can be programed only once after erase (Program Unit - Page)

- Lots of research and work to understand WA

  - Generally computed as a function of Overprovisioning (OP)

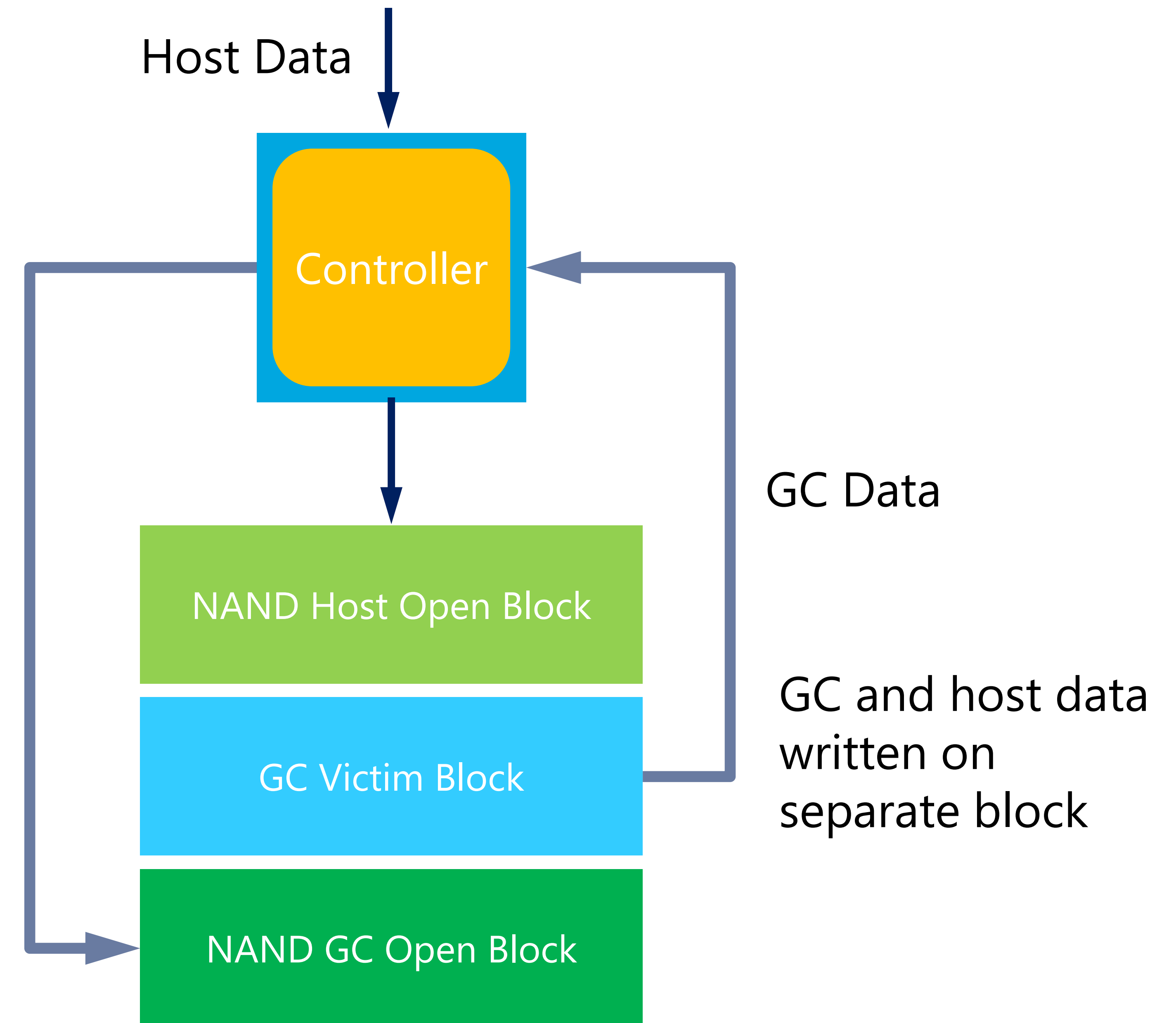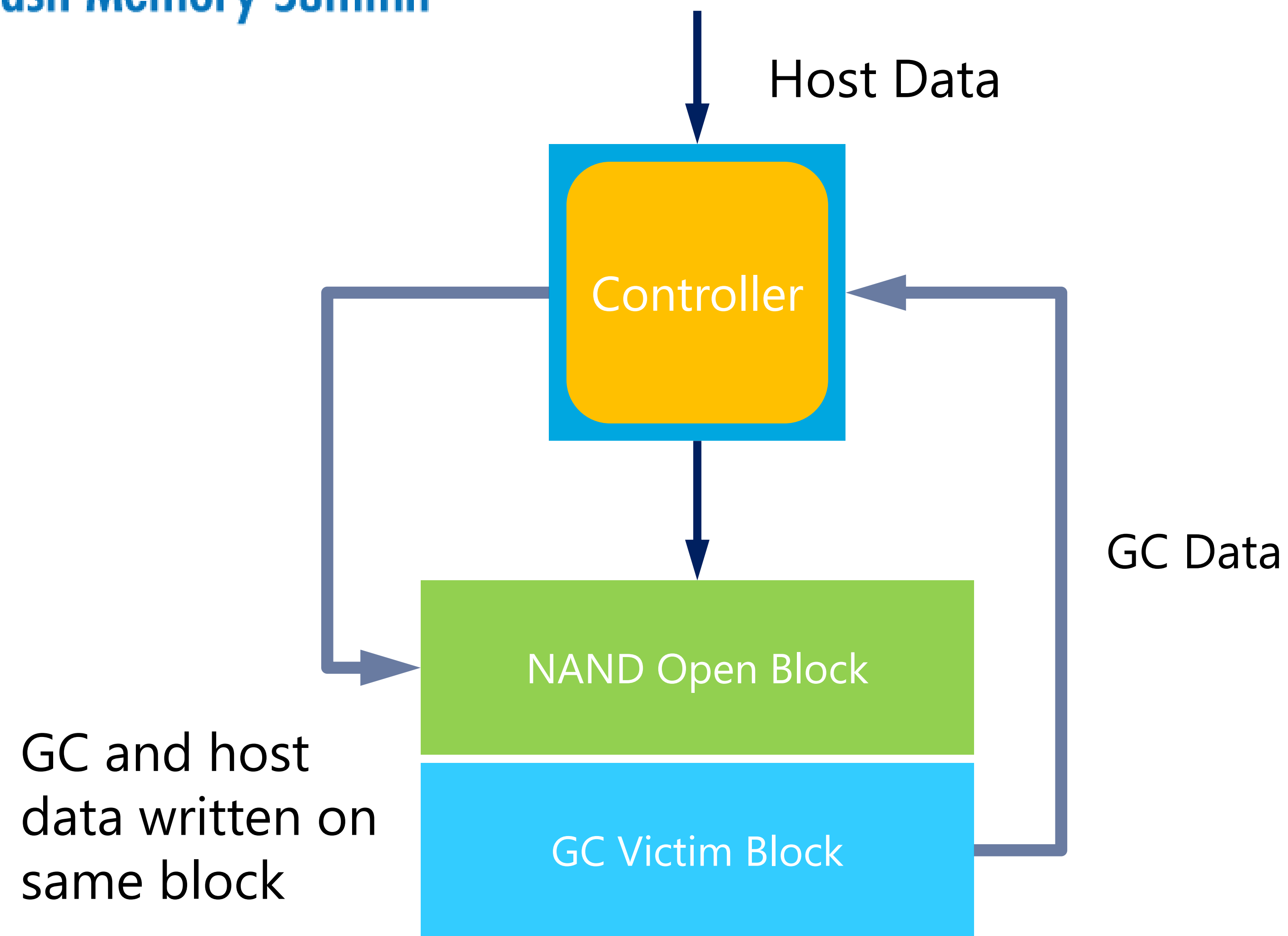  - Under steady state (degree of randomness is constant ) workloads, WA stabilizes (close to a) fixed number

$$WA = \frac{1}{2}\left(\frac{1 + \sigma}{\sigma}\right)$$

Where $\sigma$ denotes Overprovisioning factor

[Ref] Rajiv Agarwal and Marcus Marrow, "A closed-form expression for write amplification in NAND flash", in IEEE Globecom 2010 workshop on Applicat. of Commun. Theory of Emerging Memory Technologies, pp. 1908-1912   - also there is another paper on improved form of this equation

Real life workloads may not be steady state (or hard to predict)
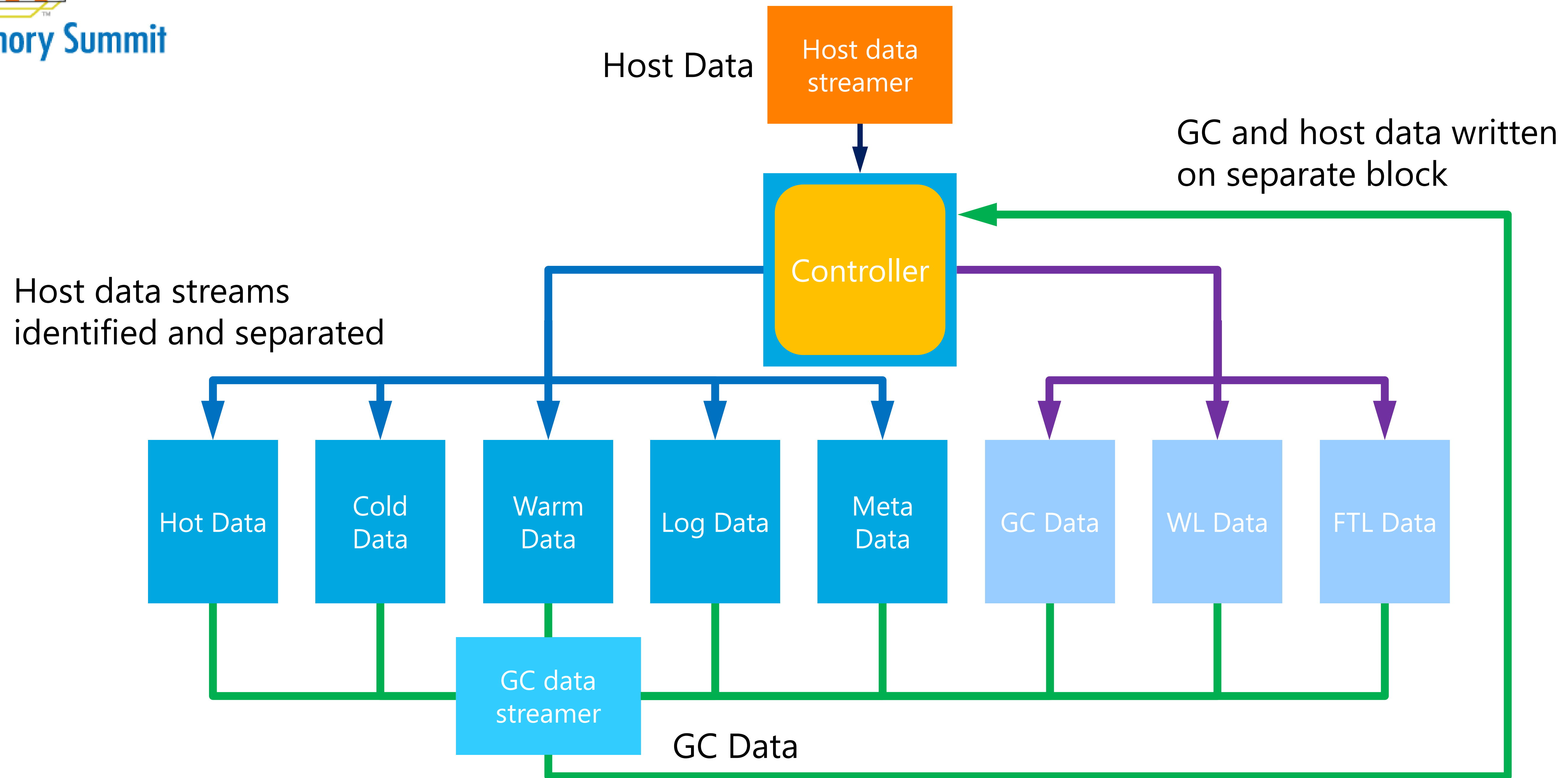Published equations fall apart under real life working conditions

August 8, 2017    |  Micron Technology

# Garbage Collection Process (GC)



Host Data

Controller

GC Data

NAND Open Block

GC Victim Block

GC and host data written on same block

Host Data

Controller

GC Data

NAND Host Open Block

GC Victim Block

NAND GC Open Block

GC and host data written on separate block

August 8, 2017   |  Micron Technology

# Separating incoming data based on expected life



Host Data

Host data streamer

GC and host data written on separate block

Controller

Host data streams identified and separated

| Hot Data | Cold Data | Warm Data | Log Data | Meta Data | GC Data | WL Data | FTL Data |

GC data streamer

GC Data

August 8, 2017    |  Micron Technology

# Time Domain Performance

# Typical SSD Performance – Bandwidth over time



- Bandwidth changes over time
  - No WA when drive is Fresh Out of Box
  - As drive is filled WA kicks in, dropping performance

- Steady state nature of WA is fully understood
  - Little efforts understanding performance consistency
  - Mechanism that impacts write amp under dynamic workload conditions

We will look at 2 key aspects for performance variance
1. logical saturation
2. random to sequential workload transition

August 8, 2017   | Micron Technology

# Logical Saturation and Degree of Randomness

# Logical Saturation – through simple geometrical mapping

$Op_{bs\ lsat}$

$x_{lsat}$

$x$

$T_{bs}$

*Random Data*

$Op_{bs}$

$N_{bs}$

Overprovisioned space  ≈ Invalided TU generated free space

- Logical Saturation
  - The portion of user logical block addresses (LBAs) that contain data

- Physical Saturation
  - The portion of physical NAND locations that contain data

$T_{bs}$ - *Translation units (TU) per blocks*

$N_{bs}$ - *Logical capacity blocks*

$Op_{bs}$ – *Blocks providing over provisioning*

$x$ – *valid Translation unit count on GC block*

$Op_{bs\ lsat}$ – *Blocks available due to lower logical saturation*

$x_{lsat}$ – *valid TU count due to lower logical saturation*

August 8, 2017    | Micron Technology

# Projecting Valid TU Count as a Function of Logical Saturation

$$WA \cong \frac{1}{1-x}$$

Legend:
- 1st drive fill
- 2nd drive fill
- 3rd drive fill
- 4th drive fill

$x_{i+n}$

Y-axis: Normalized Valid Block Count (1.0, 0.5)

X-axis: Block Number (100, 200, 300, 400, 500)

- Projecting $x$ as function of new blocks written

$$x_{i+1} = \left[1 - \frac{1}{Tu_{Logical}}\right]^{(Tb_i - Tb_{i-1}) * Gc_j}$$

- $x_{i+1}$ - Valid TU count on the GC victim block

- $Tu_{Logical}$ - Logical capacity of the drive

- $(Tb_i - Tb_{i-1})$ - New TUs written on last block, i.e. TUs per block – GC TUs moved $(x_i)$ for $j$ the block $(Gc)$

[Iterative form of ] Xiao-Yu Hu , Evangelos Eleftheriou , Robert Haas , Ilias Iliadis , Roman Pletka, Write amplification analysis in flash-based solid state drives, Proceedings of SYSTOR 2009: The Israeli Experimental Systems Conference, May 04-April 06, 2009, Haifa, Israel

Enables to projection of write amplification with any block range, any form of logical saturation, basically time domain WA as a function of number of random TUs written

August 8, 2017   | Micron Confidential

# Alternative (simple and less accurate) f



Dynamic Write Amplification

- Function of number of Overprovisioning block stripes and logical saturation – static equation

$$x = T_{bs} \left[ \frac{N_{bs} - Op_{bs}}{N_{bs} + Op_{bs}} \right]$$ Ignoring logical saturation

$$x = T_{bs} \left[ \frac{(N_{bs}\,(2 - l_{sat})) - Op_{bs}}{(l_{sat}\,N_{bs}) + Op_{bs}} \right]$$ Considering logical saturation

- Correlation with this methodology to transactional WA model

Iterative form we derived (1st of its kind) correlates WA within 1% to transactional model
Simpler form correlates within 5% at lower OP and within 20% to higher op (similar to some published work before)

August 8, 2017    | Micron Technology

# Intermixing of Sequential and Random Content

# Dynamic WA during Workload Transition

- After drive preconditioned with random data, subsequent sequential fills does not return the bandwidth to nominal state quickly
  - Nominal bandwidth correspond to sequentially preconditioned drive, executing sequential IOs

WA = 1

WA = 1

Non steady state WA

*Fill Random*

*Sequential Fills*

*Nominal Bandwidth*

1.0

0.5

*Normalized Bandwidth*

5000   10000   15000   20000   25000   30000   35000   40000

*Time (s)*

Understand how the performance recovery process works

# Modeling to Gain Insight – what we learned

- **Sequential BW with Sequential Preconditioning**
  - Incoming host data invalidates entire block stripe, no valid data left to move for garbage collection process

- **Sequential BW with Random Preconditioning**
  - Garbage collection intermix random data and sequential data for a single cursor drive



Transition workload
VTC plots look different

To understand performance during transition we need to understand VTC curves

August 8, 2017    |  Micron Technology

# 1st Sequential fill after random fills at 25% logical saturation



Sorted valid TU count per block stripe

Latest written blocks contain more valid data

This plot shows validity curve during sequential fill, when the drive was preconditioned with random data

August 8, 2017 | Micron Technology

# 50% sequential fill after random fills



2nd Sequential Fill after random fills, previous fill was sequential

Building pure sequential content at low VTC blocks

August 8, 2017   |  Micron Technolog...

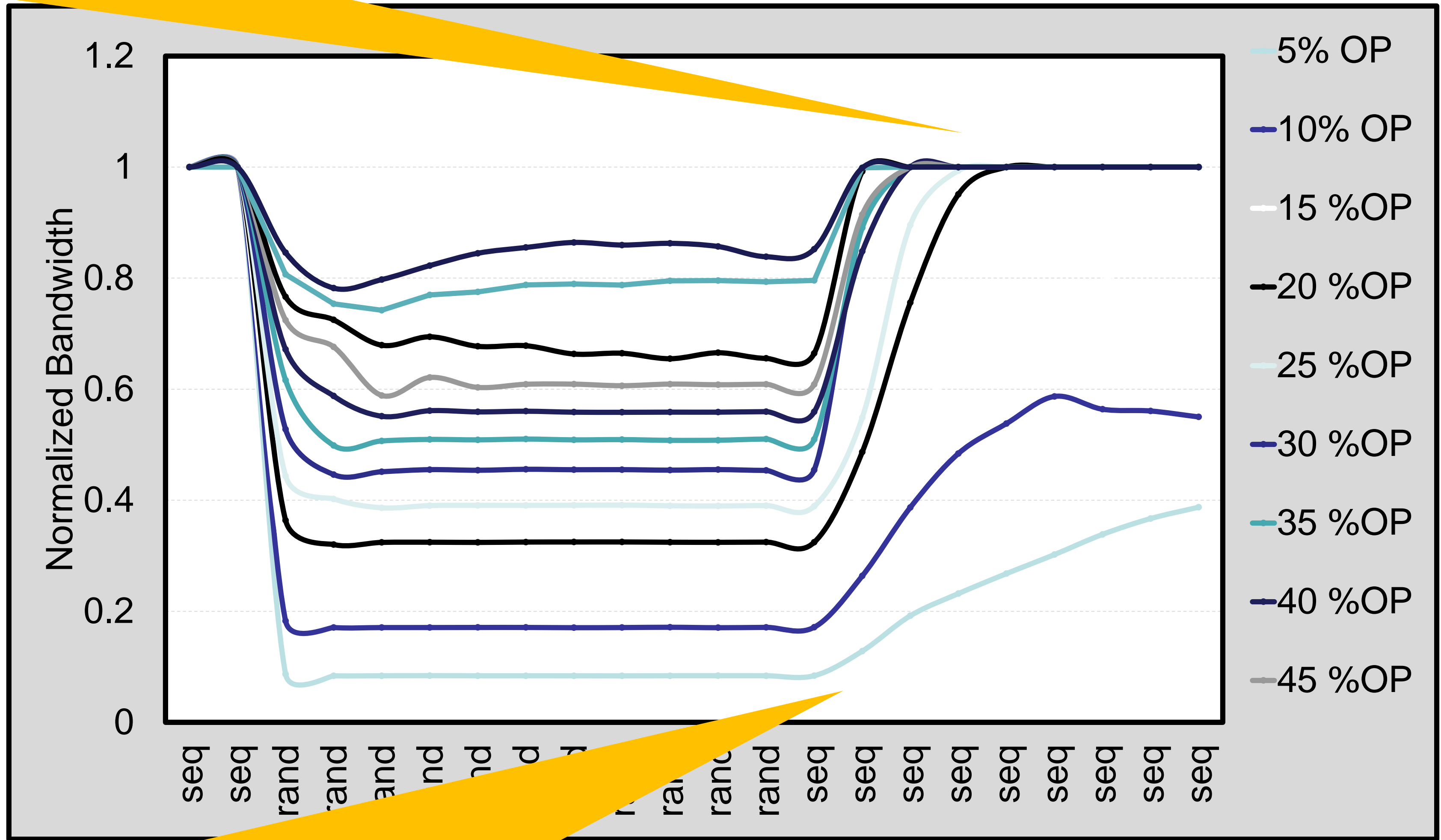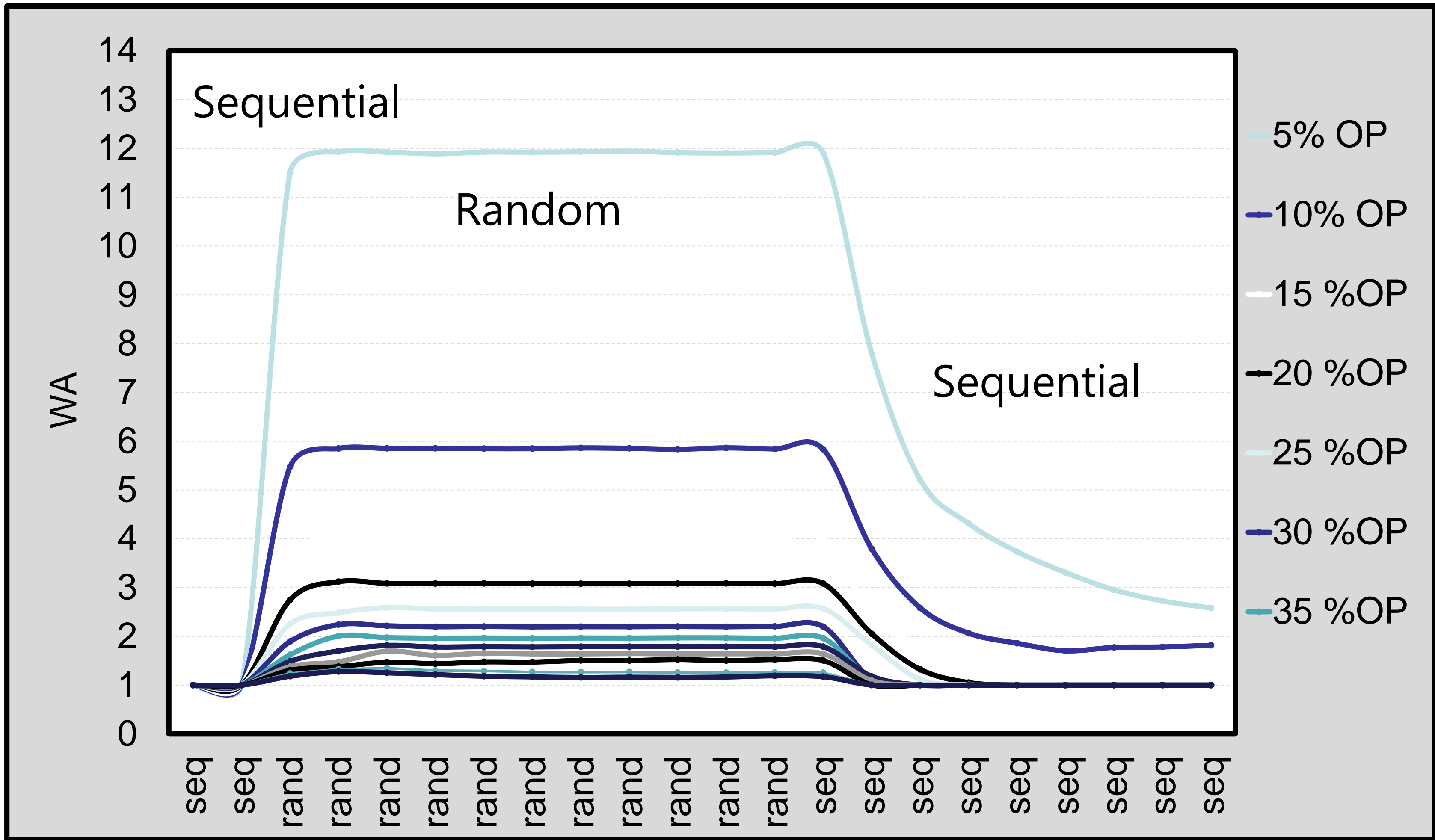# 100% sequential fill after random fills

August 8, 2017   |  Micron Technology

# Bandwidth and Write Amp Vs. OP during workload transitions



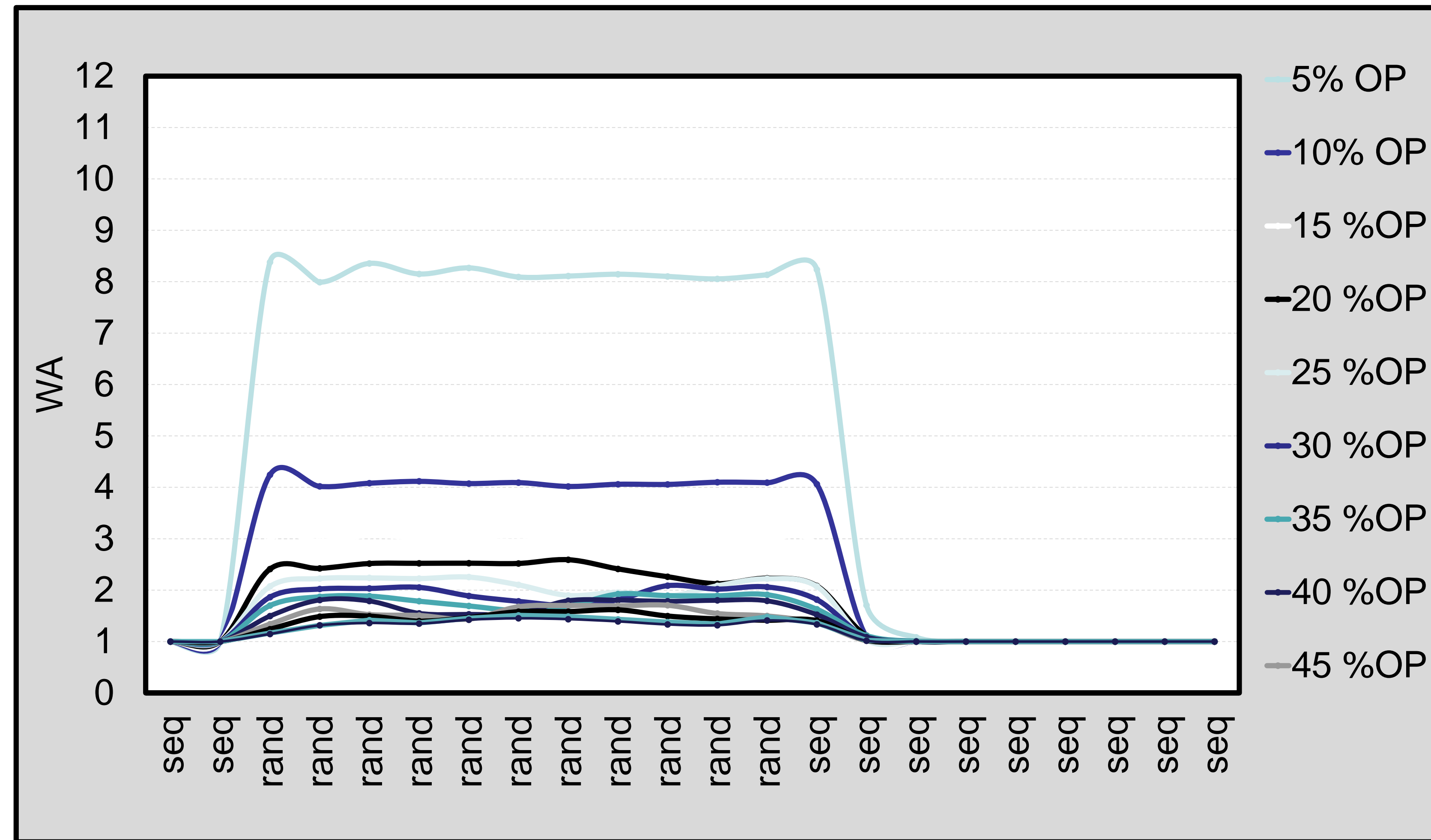On high OP, we get back to the Sequential Write performance after 2 drive fills

On low OP, it can take as much as 6 drive fills to back to the Sequential Write performance

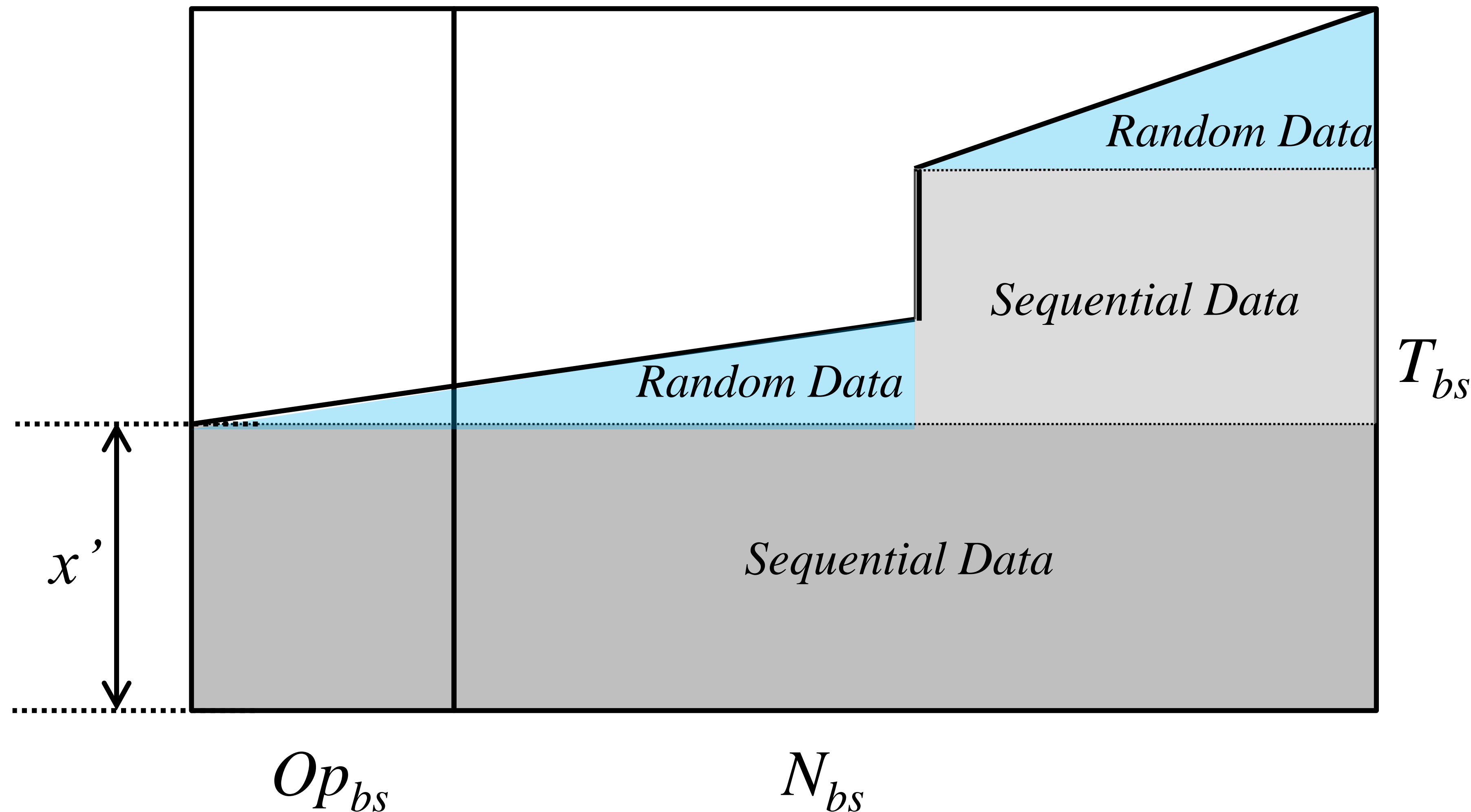| August 8, 2017    | Micron Technology

# Multi-cursor Recovery Process



- Drive can recover performance with 1 fill even for low OP configurations

- GC random data and Host sequential data is written to separate block stripes

- Avoiding intermixing of the sequential and random data

- However in real life there is always possibly even host sending intermixed sequential and random data

    August 8, 2017   |   Micron Technology

# Mathematical Formulation for Intermixed Workload



- Can be extended form of previous equations presented

- Number of discontinuous sections are function of Overprovisioning and WA

- Slope of VTC curve can be projected using iterative form presented previously

# Summary

# Dynamic WA is key to understand performance consistency

- Real life workloads can include number of variables impacting write bandwidth – trim/unmap, sequential/random, hot/cold etc

- Presented 1st of kind work (to best of our knowledge) to better understand how WA impacts performance consistency

- Demonstrated feasibility to characterize dynamic WA

August 8, 2017 | Micron Technology