# Flexible Data Placement using NVM Express®

*Implementation Perspective*

EMPOWERING OPEN.

OCP GLOBAL SUMMIT | OCTOBER 18-20, 2022 SAN JOSE, CA
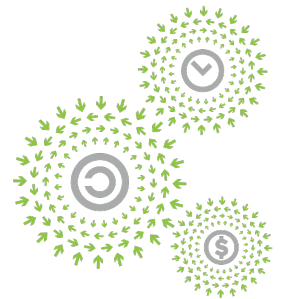
# Standards Perspective

STORAGE

Dave Landsman, Distinguished Engineer, Western Digital

OPEN PLATINUM™

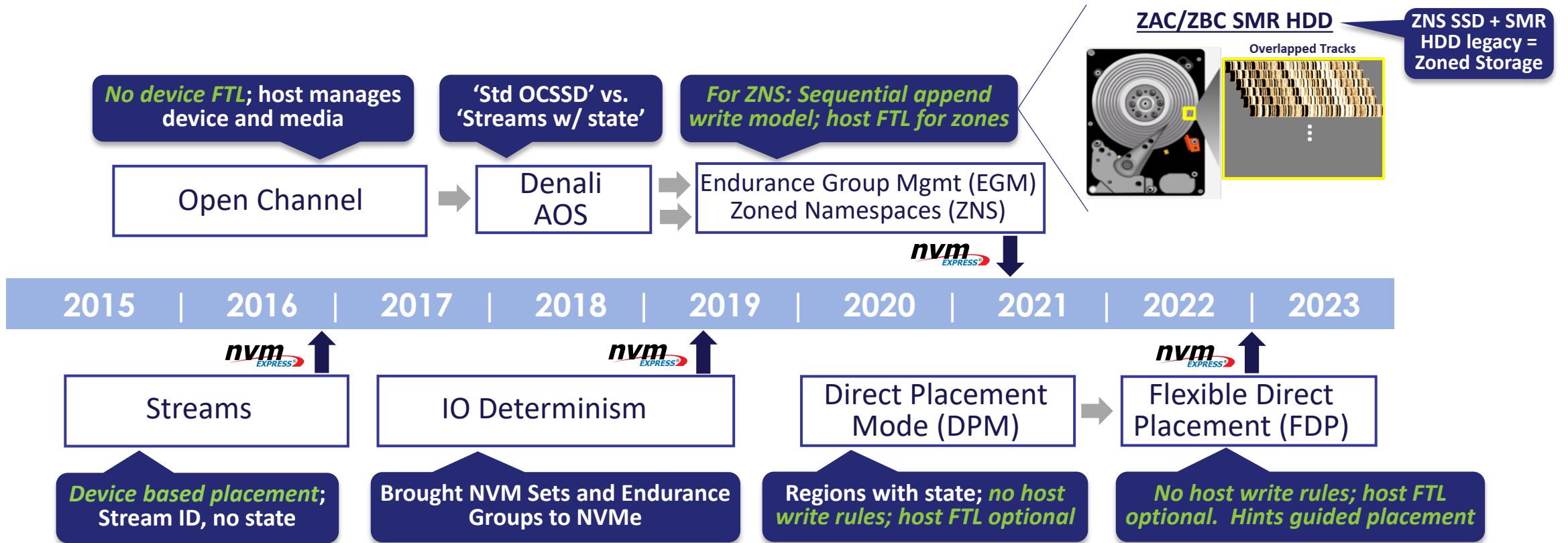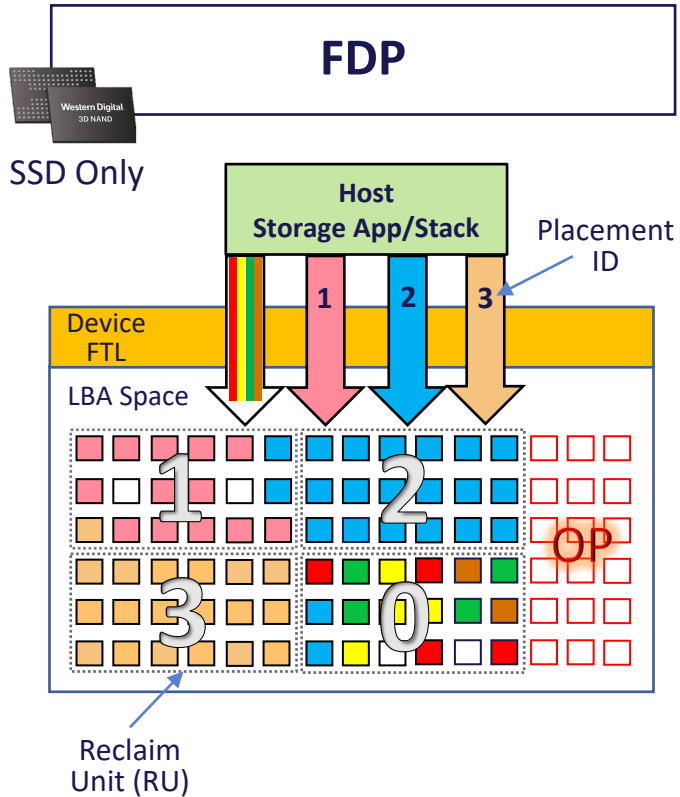# Data Placement Discussions In and Around NVMe
## Common Goal: Reduce Write Amp ➜ Reduce GC ➜ Better Endurance/Throughput/QoS

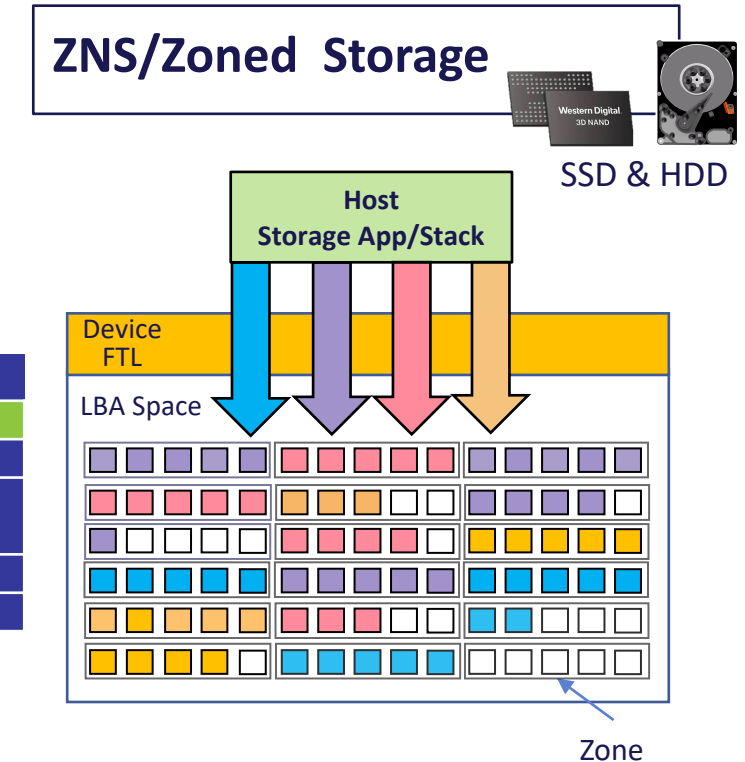**ZAC/ZBC SMR HDD**

Overlapped Tracks

**ZNS SSD + SMR HDD legacy = Zoned Storage**

**No device FTL**; host manages device and media

**'Std OCSSD' vs. 'Streams w/ state'**

**For ZNS: Sequential append write model; host FTL for zones**

| Open Channel | → | Denali AOS | → | Endurance Group Mgmt (EGM) Zoned Namespaces (ZNS) |
|---|---|---|---|---|

| 2015 | | 2016 | | 2017 | | 2018 | | 2019 | | 2020 | | 2021 | | 2022 | | 2023 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| Streams | | IO Determinism | | Direct Placement Mode (DPM) | → | Flexible Direct Placement (FDP) |
|---|---|---|---|---|---|---|

**Device based placement**; Stream ID, no state

**Brought NVM Sets and Endurance Groups to NVMe**

**Regions with state**; *no host write rules; host FTL optional*

**No host write rules; host FTL optional. Hints guided placement**

## All before ZNS and FDP abandoned or not widely adopted

# FDP and ZNS (Zoned Storage) have their own "Lanes"

**FDP**

SSD Only

Host Storage App/Stack

Placement ID

Device FTL

LBA Space

1  2  3  0

OP

Reclaim Unit (RU)

**ZNS/Zoned Storage**

SSD & HDD

Host Storage App/Stack

Device FTL

LBA Space

Zone

| Use Case Contrasts | |
|---|---|
| **FDP** | **ZNS** |
| Compute-centric | Capacity/Cost-centric |
| • Standard Block Device | • Zoned Block Device |
| • Mainstream NAND | • Highest capacity NAND |
| • Std OP%, Std DRAM | • ~0% OP, Reduced DRAM |

| Protocol Contrasts | |
|---|---|
| **FDP** | **ZNS** |
| Host writes to any LBA in any order | Host writes to append point in Zone |
| Host guides placement with hints; reduces device GC | Host manages zone GC; sequential writes in zones avoid device GC |
| Writes never rejected | Non-sequential writes rejected |
| Host FTL/Optimization optional | Host FTL/Optimization assumed |

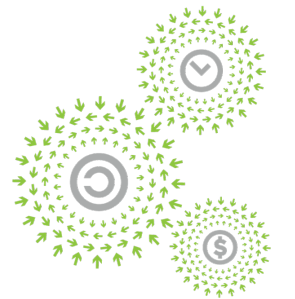## Let's build out FDP and ZNS; the ecosystem needs stability

STORAGE

# Implementation Perspective
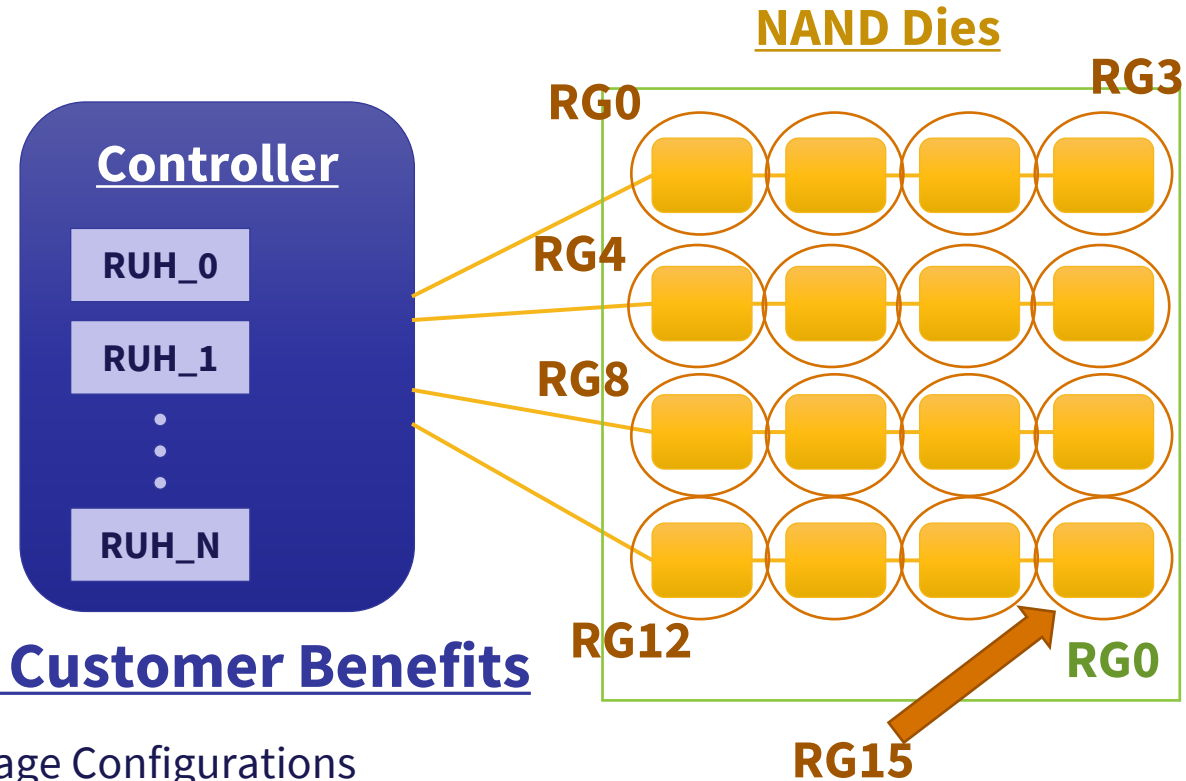
Dan Helmick, NVMe SSD Architect, Samsung

OPEN
PLATINUM™

EMPOWERING OPEN.

# FDP Implementation Externalizes Drive Design Options

**NAND Dies**

- NAND Layout per Configuration
    - Ex: **RG0** vs **RG0-RG15**
    - Ex: RU Formation

- Efficient Controller Resource Addressing and Utilization
    - RUH reporting
    - Power Fail Configurations

- Data Rerouting on Media Exceptions

- Reclaim Group selection for Legacy Users

- Extensible and Scalable for Future Feature Additions

**Controller**

RUH_0

RUH_1

⋮

RUH_N

RG0    RG3

RG4

RG8

RG12    RG0

RG15

## Highlighted Customer Benefits

- Server vs Storage Configurations

- Legacy Interoperability

- Command checking by Drive rather than Host Layers

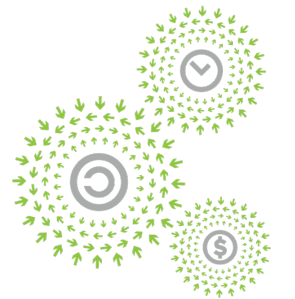- Event Logging rather than Error Interrupts

STORAGE

# Implementation Perspective

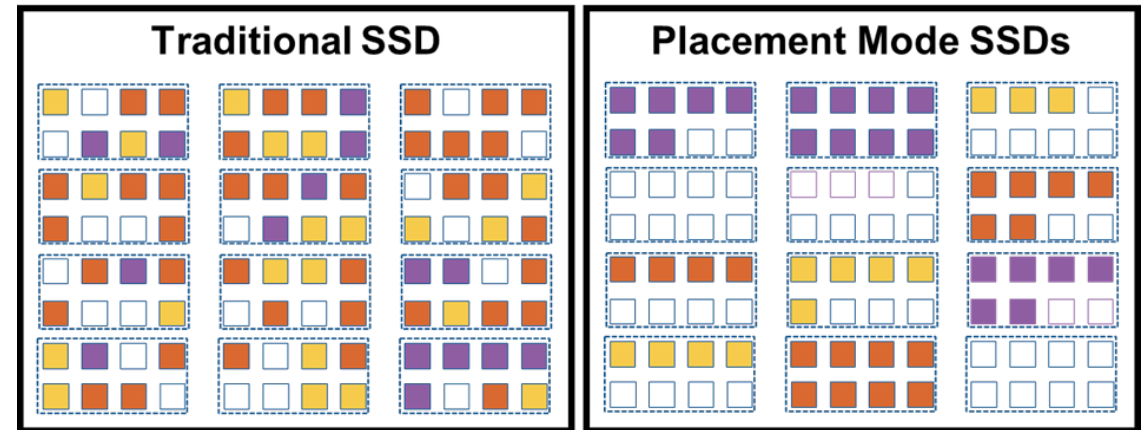John Rudelic, SSD Architect, Solidigm Technology

OPEN
PLATINUM™

# Optimization Relies on Partnership

## Key Host Side Considerations:

- Host responsibilities with TP-4146
  - Parallel IO scheduling for performance
  - NAND constraints/architecture
  - NAND features and benefits
  - SSD geometry / capacity
  - Garbage collection / housekeeping
  - New commands

- System Benefits – Better SSD Utilization
  - Significant performance (e.g. WA improvement)
  - QoS improvements
  - Write amplification improvement
  - Host placement granularity
  - Flexible use cases – (compute & storage)

## Controller Considerations & Industry Alignment on Use Cases:

- Number of supported configuration(s)
- Capability for "Default PID"
- Future NAND features
- NAND Evolution (divergent features)
- IU size with large capacity SSDs

# Implementation Perspective
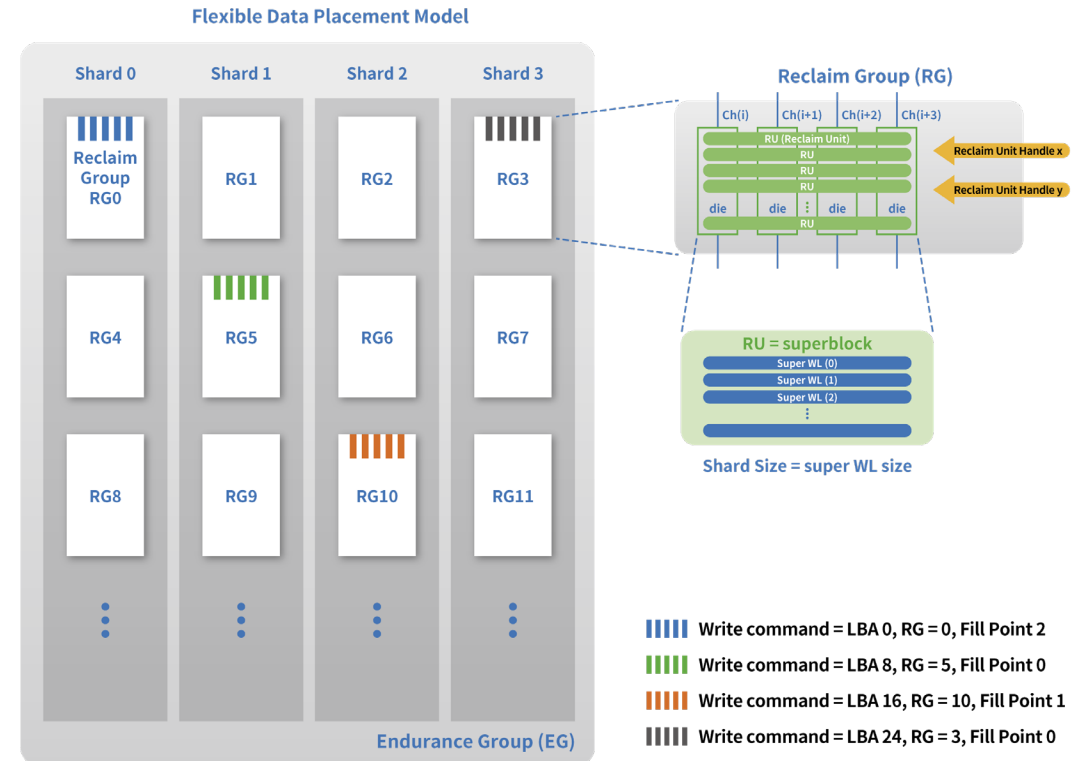
David Wang, SSD FW Architect, Silicon Motion Inc.

# FDP Implementation Difficulties

- FDP application usually requires a big number of open blocks. This implies SSD controller provide <u>a large size of high bandwidth write buffers</u>
  - Smaller size of RU (e.g., in 1GB level better for GC efficiency) means more Reclaim Groups (RGs) and open block number (i.e., concurrent write/fill pointers)
  - To deal with write-write collision and write-erase collision, also need the extra write buffer size to buffer host data during waiting time

- Using DRAM as write buffers is one feasible way but there are two challenges
  - <u>DRAM bandwidth</u>
  - <u>Supercap restriction</u>



**Flexible Data Placement Model**

Shard 0, Shard 1, Shard 2, Shard 3 — Reclaim Group RG0, RG1, RG2, RG3, RG4, RG5, RG6, RG7, RG8, RG9, RG10, RG11 — Endurance Group (EG)

Reclaim Group (RG) — Ch(i), Ch(i+1), Ch(i+2), Ch(i+3) — RU (Reclaim Unit), Reclaim Unit Handle x, Reclaim Unit Handle y, die

RU = superblock — Super WL (0), Super WL (1), Super WL (2) — Shard Size = super WL size

Write command = LBA 0, RG = 0, Fill Point 2
Write command = LBA 8, RG = 5, Fill Point 0
Write command = LBA 16, RG = 10, Fill Point 1
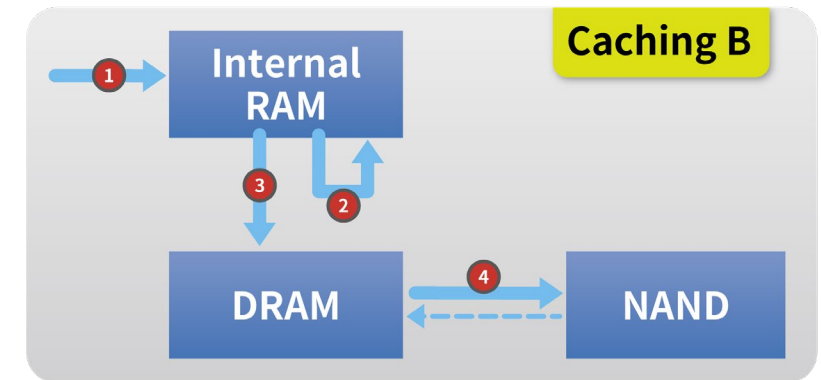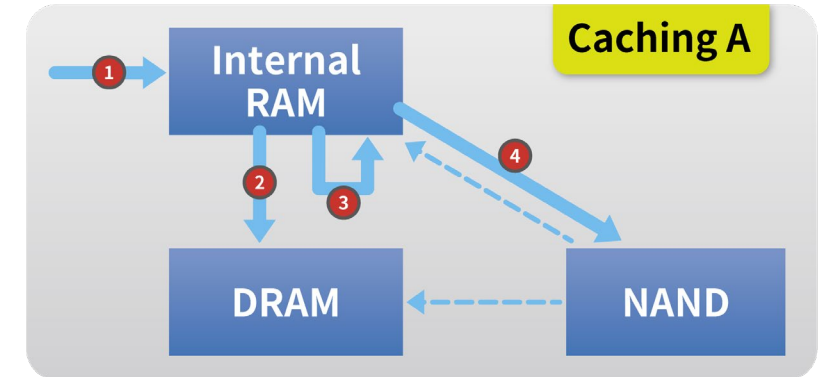Write command = LBA 24, RG = 3, Fill Point 0

# Gen5 NVMe Controller Solution

1. NVMe Controller may support high bandwidth DDR5 DRAM, that can be used as Write Buffers in Write IO path without performance loss
   - Dual Channel DDR5 bandwidth: 4.8GHz*8B*0.7 (DDR efficiency) = ~27GB/s, then max supported NAND Prog throughput is ~13.5GB/s

2. Use Flexible Write Cache to Configurable Open Block Number

| Open Block Number | Buffer (for Data Operation and NAND Prog) | Cache (for backup in case prog failure) | Mechanism |
|---|---|---|---|
| Small/medium num (e.g. <=32) | Internal RAM → NAND | DRAM | *A* |
| large num ( e.g. <=128) | Combined Internal RAM + DRAM → NAND | DRAM Limited by Supercap* | *B* |



* For example, Supercap capacitance 1500uF may only protect ~100MB DRAM data (write buffers and other FTL metadata) during power lose.