# STORAGE PERFORMANCE BENCHMARKING:
# PART 3 – BLOCK COMPONENTS

Ken Cantrell, NetApp

Mark Rogov, EMC

David Fair, SNIA ESF Chair, Intel

**March 8, 2016**

# SNIA Legal Notice

- The material contained in this tutorial is copyrighted by the SNIA unless otherwise noted.
- Member companies and individual members may use this material in presentations and literature under the following conditions:
  - Any slide or slides used must be reproduced in their entirety without modification
  - The SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of the SNIA Education Committee.
- Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.
- The information presented herein represents the author's personal opinion and current understanding of the relevant issues involved. The author, the presenter, and the SNIA do not assume any responsibility or liability for damages arising out of any reliance on or use of this information.

  NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.

# About The Speakers

## Ken Cantrell

NetApp

Manager Perf Engineering

@kencantrelljr
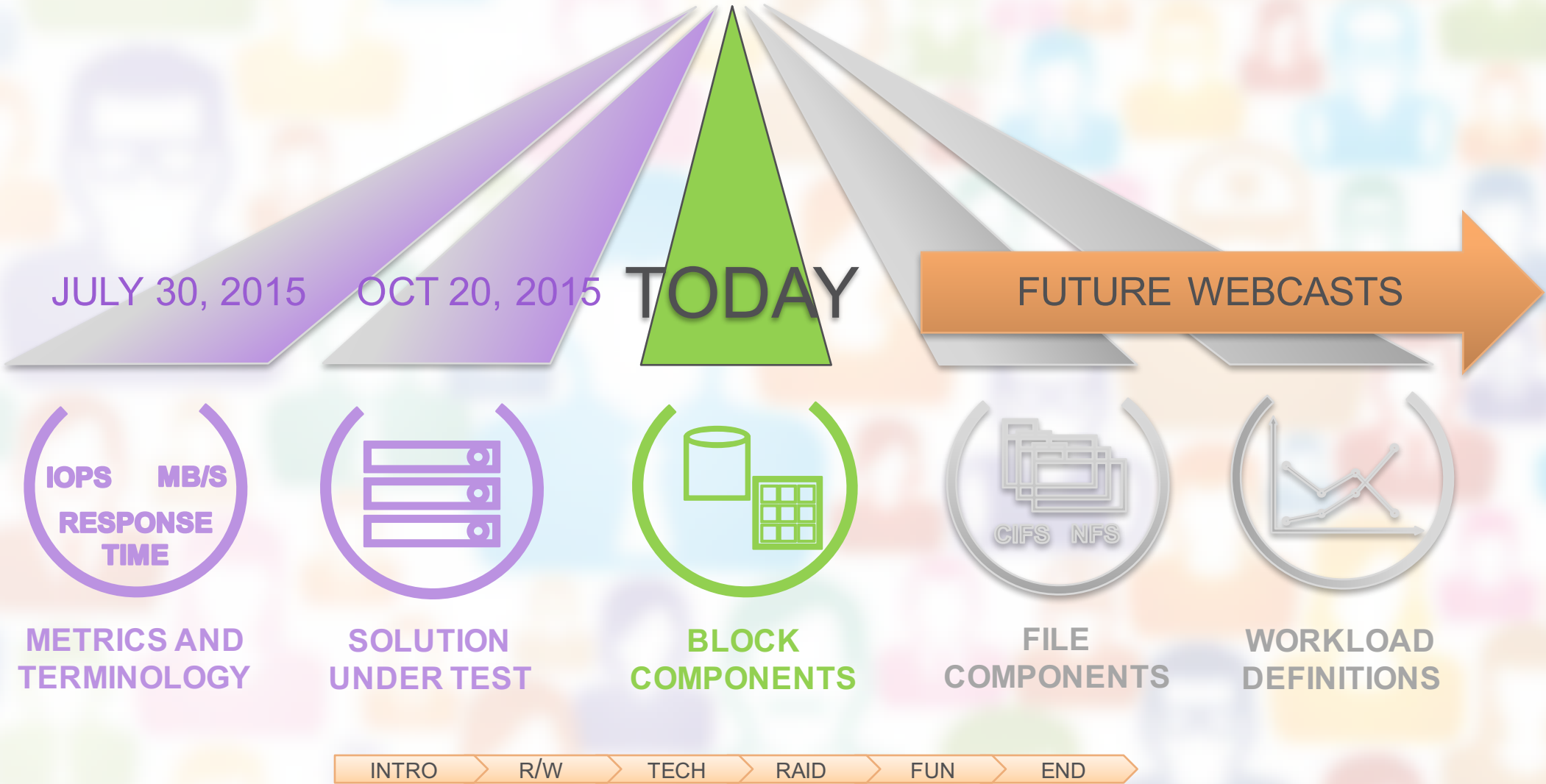
## Mark Rogov

EMC

Advisory Systems Engineer
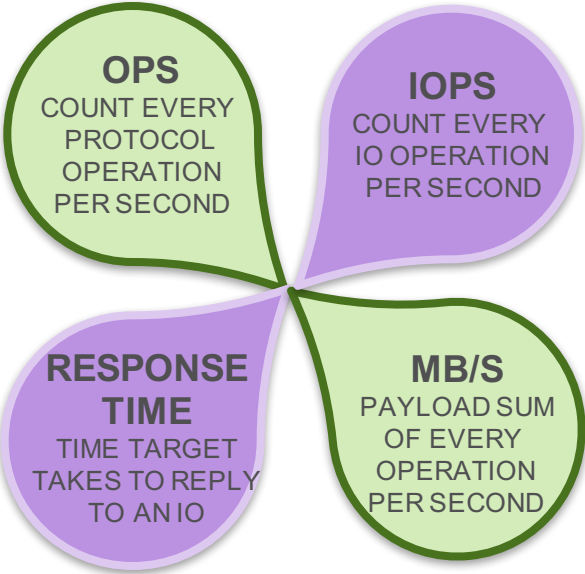
@rogovmark

## Dr. David Fair

SNIA ESF Chair
& Intel
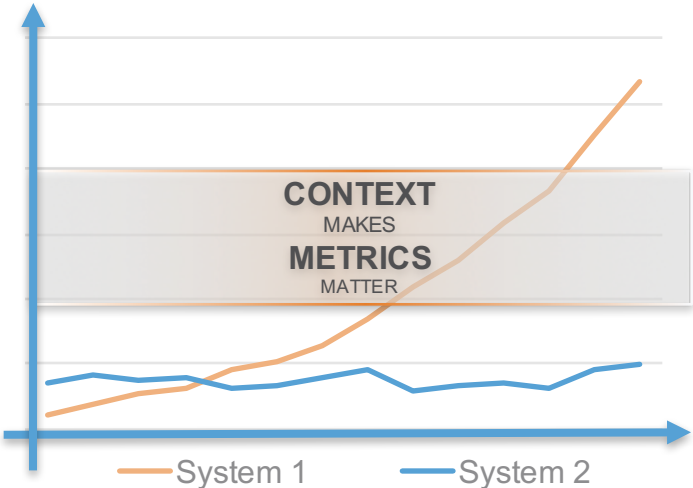Ethernet Networking
Marketing Manager

# Storage Performance Benchmarking

JULY 30, 2015     OCT 20, 2015     TODAY     FUTURE WEBCASTS

IOPS    MB/S
RESPONSE
TIME

CIFS   NFS

**METRICS AND TERMINOLOGY**

**SOLUTION UNDER TEST**

**BLOCK COMPONENTS**

**FILE COMPONENTS**

**WORKLOAD DEFINITIONS**

INTRO → R/W → TECH → RAID → FUN → END

# Session 1 – Terminology and Context

## TERMINOLOGY

**OPS**
COUNT EVERY PROTOCOL OPERATION PER SECOND

**IOPS**
COUNT EVERY IO OPERATION PER SECOND

**RESPONSE TIME**
TIME TARGET TAKES TO REPLY TO AN IO

**MB/S**
PAYLOAD SUM OF EVERY OPERATION PER SECOND

TERMINOLOGY

## GRAPH FUN

**CONTEXT**
MAKES
**METRICS**
MATTER

—— System 1      —— System 2

GRAPH FUN

INTRO | R/W | TECH | RAID | FUN | END

# Session 2 – The Slowest Component Matters Most

DISK BOUND

CLIENT BOUND

DO LESS WORK

DO WORK FASTER

INCREASE PARALLELISM

SLOW COMPONENT MATTERS MOST

BOTTLENECKS ALWAYS EXIST

3 PERFORMANCE PRINCIPLES

SNIA ESF | ETHERNET STORAGE

IDC
*Analyze the Future*

# 33.1 EXABYTES

**WORLD POPULATION**

7.4 BILLION

**SPLIT EQUALLY**

4.5 GIGABYTES PER INDIVIDUAL

**OR**

IOPS   MB/S   RESPONSE TIME

**METRICS AND TERMINOLOGY**

1541 COPIES OF OUR FIRST WEBCAST POWERPOINT

INTRO   R/W   TECH   RAID   FUN   END

# Eventually, All Data Goes To Block Storage

**BLOCK STORAGE**          **SOLUTIONS UNDER TEST**          **WORKLOADS**



HYPERVISOR

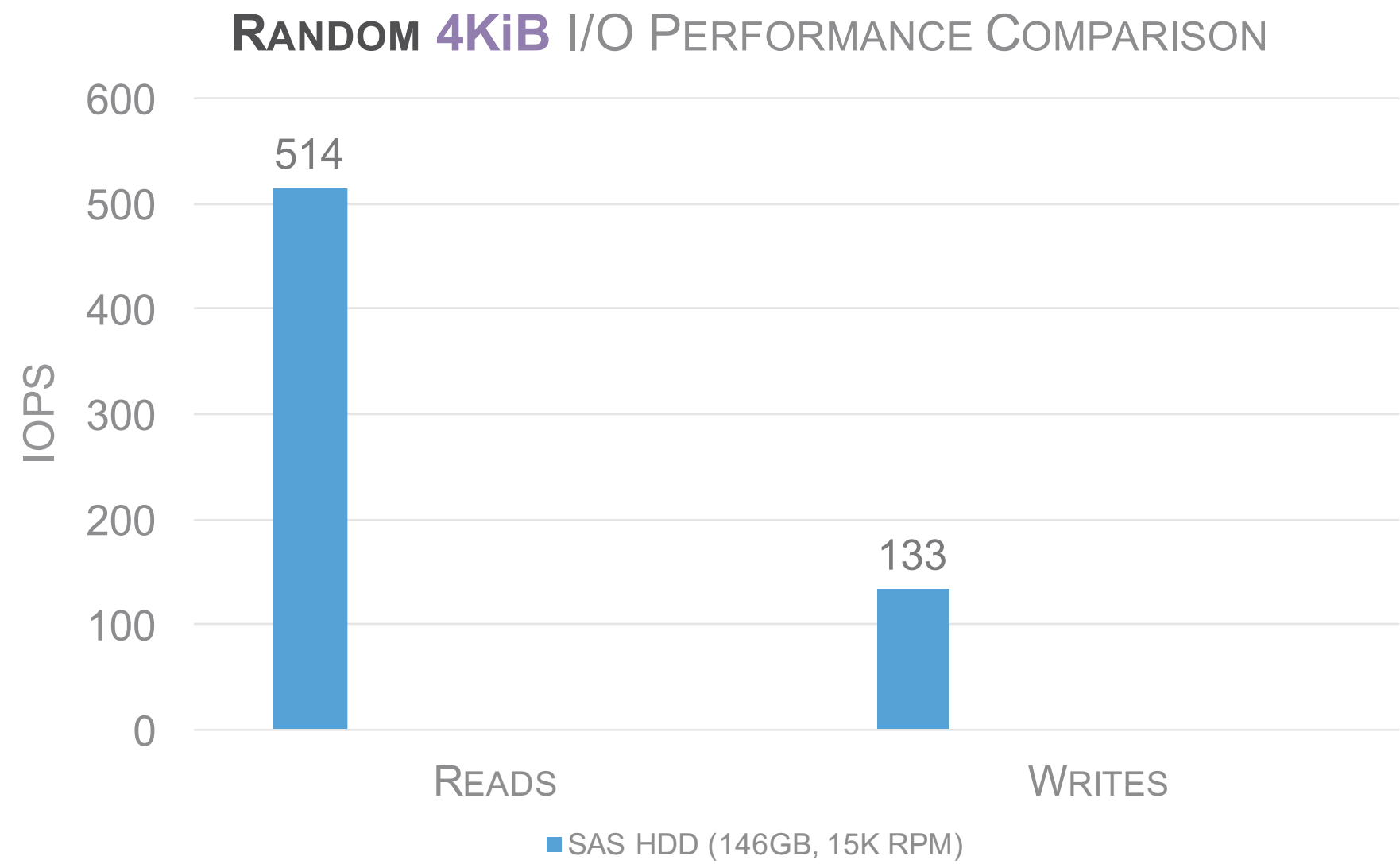HYPERVISOR

# Agenda

| | |
|---|---|
| INTRO | INTRODUCTION |
| R/W | READING, WRITING; WHAT IS THE DIFFERENCE? |
| TECH | HOW DOES THIS TECH WORK ANYWAY? |
| RAID | WHAT IF YOU NEED MORE THAN ONE? |
| FUN | PERFORMANCE? |
| END | SUMMARY |

# Let's Take A Drive… And Test It!

## RANDOM **4KiB** I/O PERFORMANCE COMPARISON



SAS HDD (146GB, 15K RPM)

# Detour! What Does "Random" Mean?

KEYS ARE ALL OVER

A QUICK BROWN FOX JUMPED OVER A LAZY DOG

IMAGINE THAT THE KEYBOARD IS A DISK DRIVE
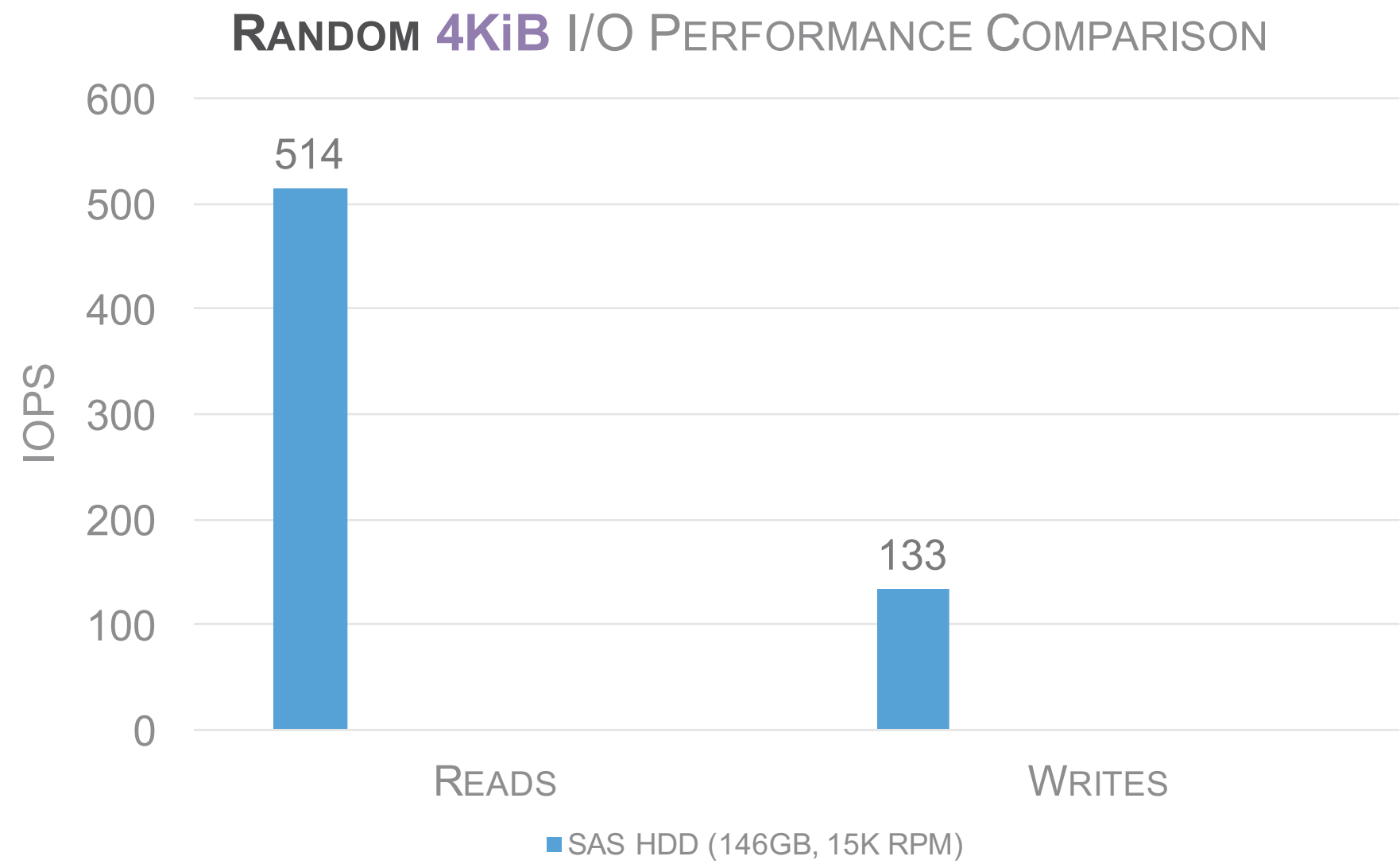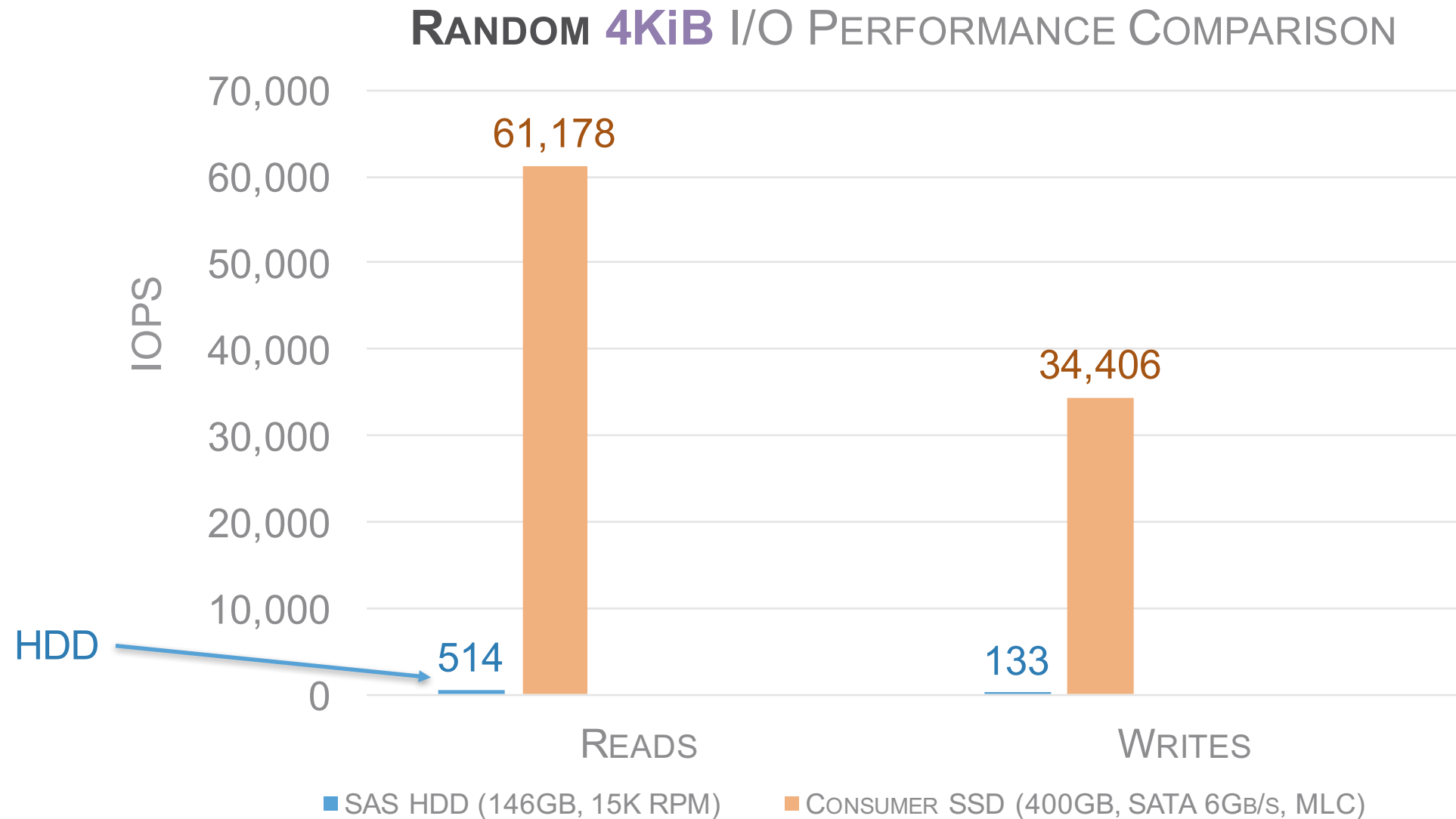
# What Does "Sequential" Mean?



EVERY KEY IS NEXT TO PREVIOUS

1 2 3 4 5 6

IMAGINE THAT THE KEYBOARD
IS A DISK DRIVE

INTRO    R/W    TECH    RAID    FUN    END

12

# "Sequential Read" Example

"SEQUENTIAL READ"

# Let's Take A Drive… And Test It!

## RANDOM 4KiB I/O PERFORMANCE COMPARISON



SAS HDD (146GB, 15K RPM)

INTRO  R/W  TECH  RAID  FUN  END

# Let's Take Two Drives… And Test Them!

## RANDOM 4KiB I/O PERFORMANCE COMPARISON



**IOPS** (y-axis: 0 to 70,000)

Reads:
- HDD (SAS HDD): 514
- Consumer SSD: 61,178

Writes:
- HDD (SAS HDD): 133
- Consumer SSD: 34,406

Legend:
- ■ SAS HDD (146GB, 15K RPM)
- ■ CONSUMER SSD (400GB, SATA 6GB/S, MLC)

INTRO | R/W | TECH | RAID | FUN | END

15

# And Add More SSDs

RANDOM 4KiB I/O PERFORMANCE COMPARISON



Chart: IOPS vs. Reads/Writes

SINGLE DRIVE

Legend:
- SAS HDD (146GB, 15K RPM)
- Consumer SSD (400GB, SATA 6Gb/s, MLC)
- Consumer SSD (128GB, M.2x4 AHCI, MLC)
- Enterprise SSD (1400GB, PCIe x8 AHCI, MLC)
- Enterprise SSD (1600GB, U.2 NVMe, MLC)

Reads:
- 514
- 61,178
- 172,881
- 750,846
- 585,884

Writes:
- 133
- 34,406
- 6,307
- 83,021
- 113,260

HDD

# Agenda

| | |
|---|---|
| **INTRO** | Introduction |
| **R/W** | Reading, Writing; What is the Difference? |
| **TECH** | How does this tech work anyway? |
| **RAID** | What if you need more than one? |
| **FUN** | Performance? |
| **END** | Summary |

# How Does This Tech Work?

FLASH

HDD OR DISK DRIVE

# Spinning Drives And Sectors



READ

WRITE

SEEK THE TRACK
SPIN TO THE SECTOR
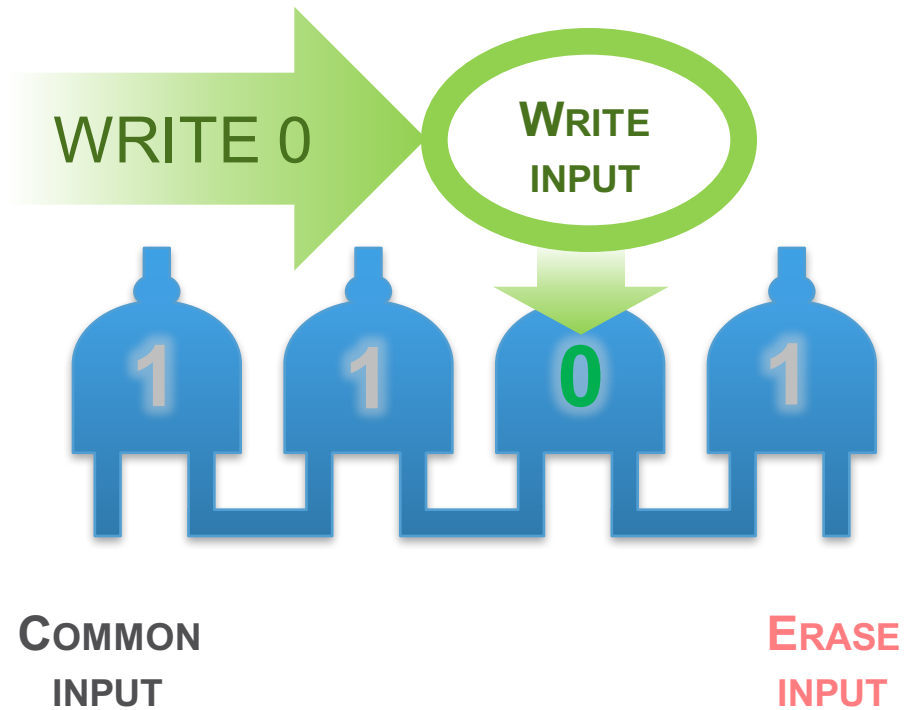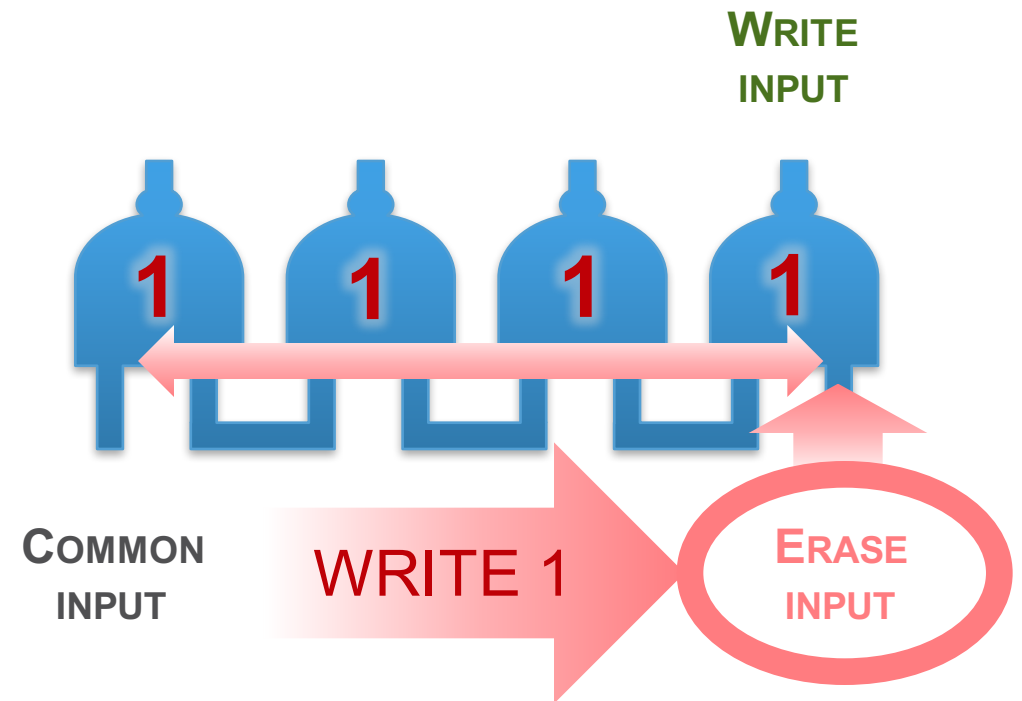
TRACK

SPIN

SECTOR

SEEK

# Flash And NAND Gates

EVERY NAND CAN BE SET TO 0 INDIVIDUALLY

TO SET BACK TO 1, AN ENTIRE GROUP NEEDS TO BE RESET



WRITE 0

WRITE INPUT

1  1  0  1

COMMON INPUT

ERASE INPUT

WRITE INPUT

1  1  1  1

COMMON INPUT

WRITE 1

ERASE INPUT

# Flash Construction

## FLASH BLOCK

PAGE

2KiB
4KiB
8KiB
…

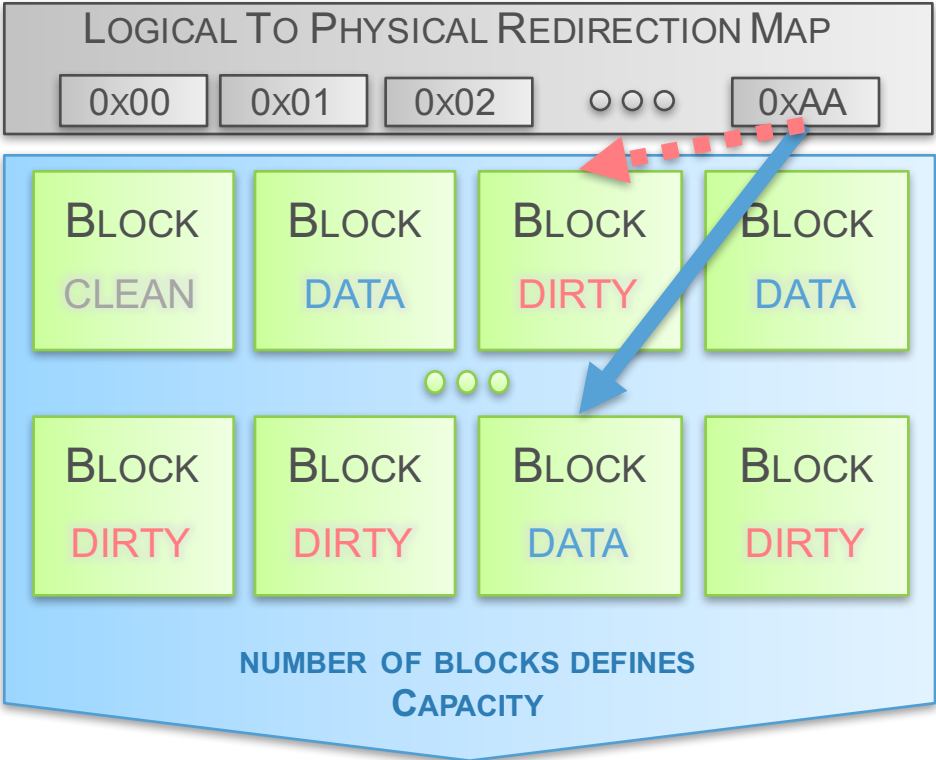**PAGES PER BLOCK (DEP ON MODEL)
128, 256, …, ETC.**

## MOST FLASH

WRITE—**1 PAGE** AT A TIME

## FLASH DEVICE

LOGICAL TO PHYSICAL REDIRECTION MAP

| 0x00 | 0x01 | 0x02 | ○ ○ ○ | 0xAA |

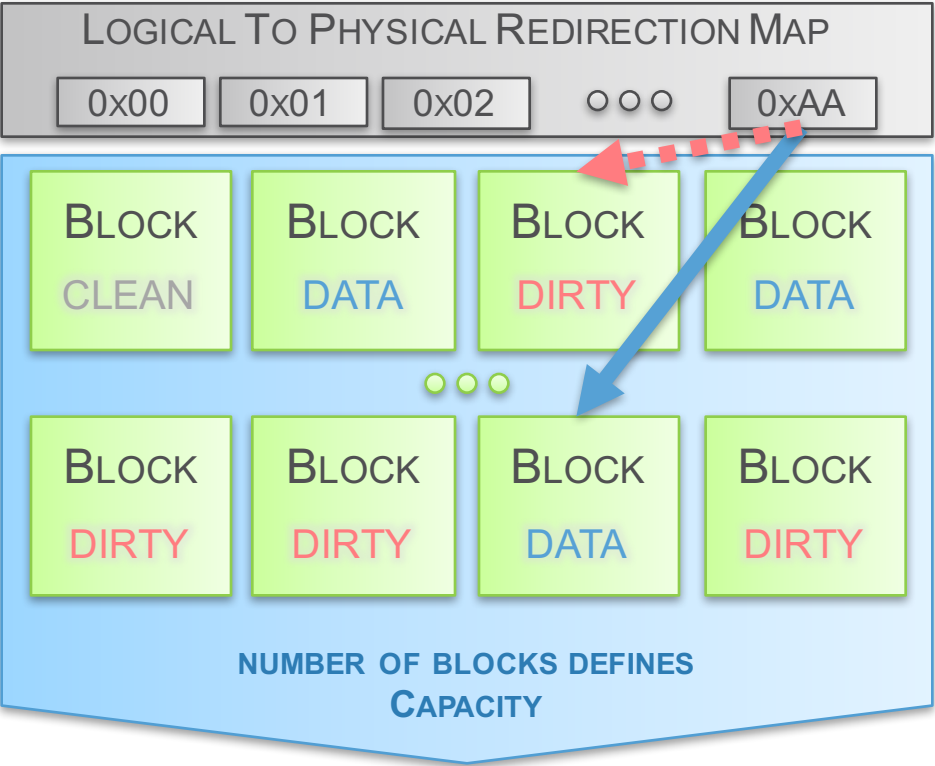| BLOCK CLEAN | BLOCK DATA | BLOCK DIRTY | BLOCK DATA |
|---|---|---|---|
| BLOCK DIRTY | BLOCK DIRTY | BLOCK DATA | BLOCK DIRTY |

**NUMBER OF BLOCKS DEFINES CAPACITY**

## REDIRECT ON OVER-WRITE

AN IO IS REDIRECTED TO A CLEAN BLOCK/PAGE
LEAVING OLD BLOCK/PAGE DIRTY

# Garbage Collection

## FLASH DEVICE

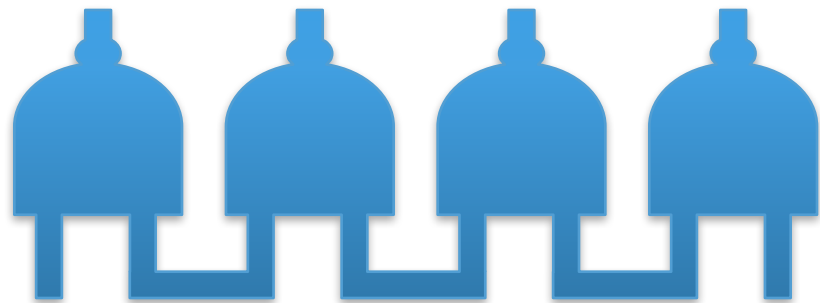LOGICAL TO PHYSICAL REDIRECTION MAP

| 0x00 | 0x01 | 0x02 | ○○○ | 0xAA |

| BLOCK CLEAN | BLOCK DATA | BLOCK DIRTY | BLOCK DATA |

○○○

| BLOCK DIRTY | BLOCK DIRTY | BLOCK DATA | BLOCK DIRTY |

**NUMBER OF BLOCKS DEFINES CAPACITY**

## GARBAGE COLLECTION

| BLOCK DIRTY | BLOCK DIRTY | BLOCK DIRTY | BLOCK DIRTY |

**ERASE**

| BLOCK CLEAN | BLOCK DIRTY | BLOCK DIRTY | BLOCK DIRTY |

**ERASE—1 DIRTY BLOCK AT A TIME**
(WHEN NUMBER OF CLEAN BLOCKS IS LOW)

# Sequential Vs. Random

## SSD or Flash

## HDD or Disk Drive



SPIN

SEEK

| EVERYTHING IS RANDOM IO FOR FLASH | WRITE | READ |
|---|---|---|
| | ERASE + WRITE | READ |

| | WRITE | READ |
|---|---|---|
| SEQUENTIAL | WRITE | READ |
| RANDOM | SEEK/SPIN + WRITE | SEEK/SPIN + READ |

**SLOWER PERFORMANCE**

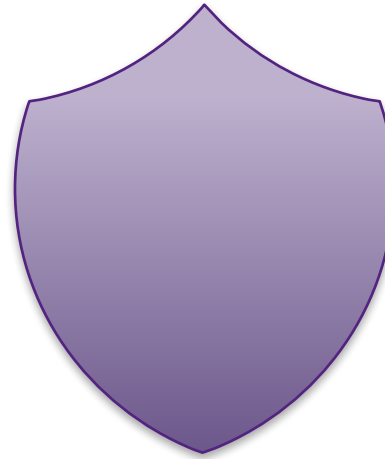INTRO  R/W  TECH  RAID  FUN  END

# Agenda

| | |
|---|---|
| INTRO | Introduction |
| R/W | Reading, Writing; What is the Difference? |
| TECH | How does this tech work anyway? |
| RAID | What if you need more than one? |
| FUN | Performance? |
| END | Summary |

# Just One?

**CAPACITY**

**PROTECTION**

**PERFORMANCE**

R

W

# RAID—**R**edundant **A**rray Of **I**nexpensive **D**isks

# RAID-0 (Striping Without Parity)

**CAPACITY**
100%

**PROTECTION**
NONE

**PERFORMANCE**
R 100%    W 100%



PHYSICAL STORAGE

VIRTUAL/ LOGICAL

CLIENTS / HOSTS

01

11

01  11

BACK END    FRONT END

# RAID-1 (Mirroring)

CAPACITY
**50%**

PROTECTION
**1 DRIVE**

PERFORMANCE
**R** 100%  **W** 50%

SNIA ESF | ETHERNET STORAGE

PHYSICAL STORAGE

VIRTUAL/ LOGICAL

CLIENTS / HOSTS

01 11

01 11

01 11

BACK END

FRONT END

# RAID-3, -4, -5 [-6, -DP]*

**Striping With Parity***

CAPACITY — N

PROTECTION — P

PERFORMANCE — N   "IT DEPENDS"

PHYSICAL STORAGE

VIRTUAL/ LOGICAL

**N** IS NUMBER OF DATA DRIVES
**P** IS NUMBER OF PARITY DRIVES

01

11

P*

01  11

BACK END
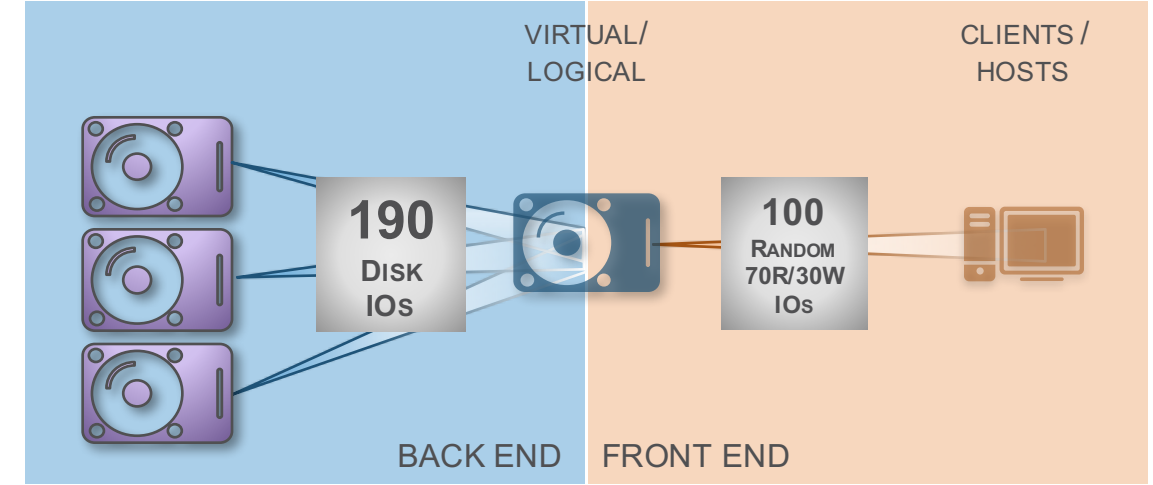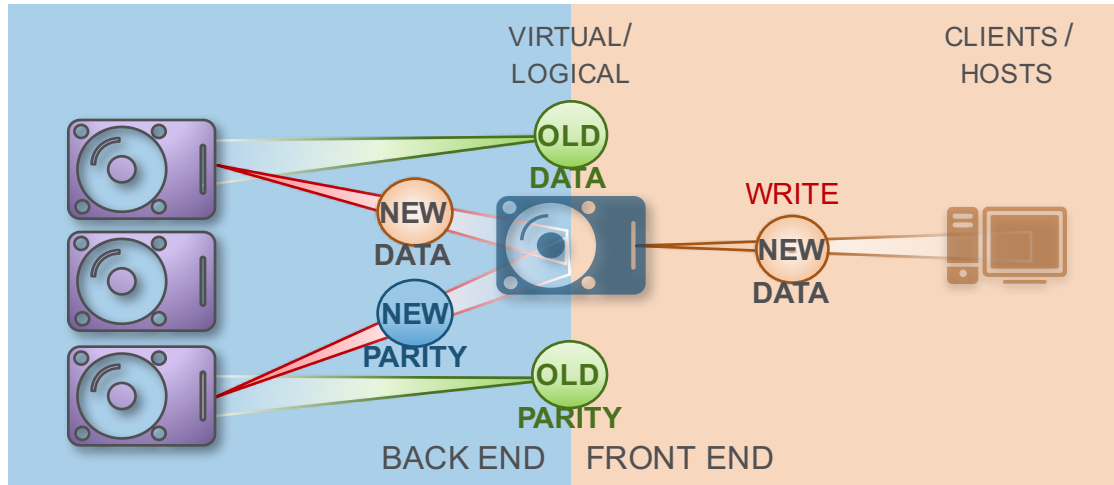
FRONT END

* RAID-6/-DP requires more than one parity

# RAID Partial Writes

**All Single Parity RAID: RAID-3, -4, -5, and etc.**

PERFORMANCE

R W

1 READ    2 READS
          2 WRITES

SNIA ESF | ETHERNET STORAGE

VIRTUAL/ LOGICAL

CLIENTS / HOSTS

OLD DATA

NEW DATA

NEW PARITY

OLD PARITY

WRITE

NEW DATA

BACK END    FRONT END

VIRTUAL/ LOGICAL

CLIENTS / HOSTS

190 DISK IOs

100 RANDOM 70R/30W IOs

BACK END    FRONT END

SINGLE PARTIAL WRITE:
- **READ** OLD DATA
- **READ** OLD PARITY
- CALCULATE NEW PARITY
- **WRITE** NEW DATA
- **WRITE** NEW PARITY

2 READS
2 WRITES

100 IOs 70R/30W = 70 READ + 30 WRITE IOs

BACKEND = (70R + 30 * (2W + 2R)) = 190 IOs

RAID PENALTY

# RAID Implementation



| PHYSICAL STORAGE | BACK-END CONNECT | STORAGE CONTROLLER | FRONT-END CONNECT | CLIENTS / HOSTS |
|---|---|---|---|---|

# Erasure Coding Implementation



| PHYSICAL STORAGE | BACK-END CONNECT | STORAGE CONTROLLERS | FRONT-END CONNECT | CLIENTS / HOSTS |
|---|---|---|---|---|

SCALE OUT

ERASURE CODING

# Erasure Coding

N + M = 2 + 1

SNIA | ETHERNET
ESF | STORAGE

**N** IS NUMBER OF DATA BLOCKS
**M** IS NUMBER OF PROTECTION BLOCKS

PHYSICAL
STORAGE

| 01 | 22 |

| 11 | P2 |

| P1 | 02 |

| 01 | 11 | 02 | 22 |

BACK END

FRONT END

# Agenda

| | |
|---|---|
| INTRO | Introduction |
| R/W | Reading, Writing; What is the Difference? |
| TECH | How does this tech work anyway? |
| RAID | What if you need more than one? |
| **FUN** | **Performance?** |
| END | Summary |

# What "Really" Happens With RAID-5?

## HDD POTENTIAL AGGREGATE 4KiB RANDOM WRITE PERFORMANCE (As Seen at Client)



IOPS

| | |
|---|---|
| 450 | |
| 400 | 106 |
| 350 | |
| 300 | 106 |
| 250 | |
| 200 | 106 |
| 150 | |
| 100 | 106 |
| 50 | |
| 0 | |

EACH DRIVE

■ DISK1　■ DISK2　■ DISK3　■ DISK4

# What "Really" Happens With RAID-5?

HDD POTENTIAL AGGREGATE 4KiB RANDOM WRITE PERFORMANCE (AS SEEN AT CLIENT)

# What "Really" Happens With RAID-5?

**HDD** Potential Aggregate **4KiB** Random Write Performance (As Seen at Client)

# What "Really" Happens With RAID-5?

**FLASH** POTENTIAL AGGREGATE **4KiB** RANDOM WRITE
PERFORMANCE (As SEEN AT CLIENT)



IOPS

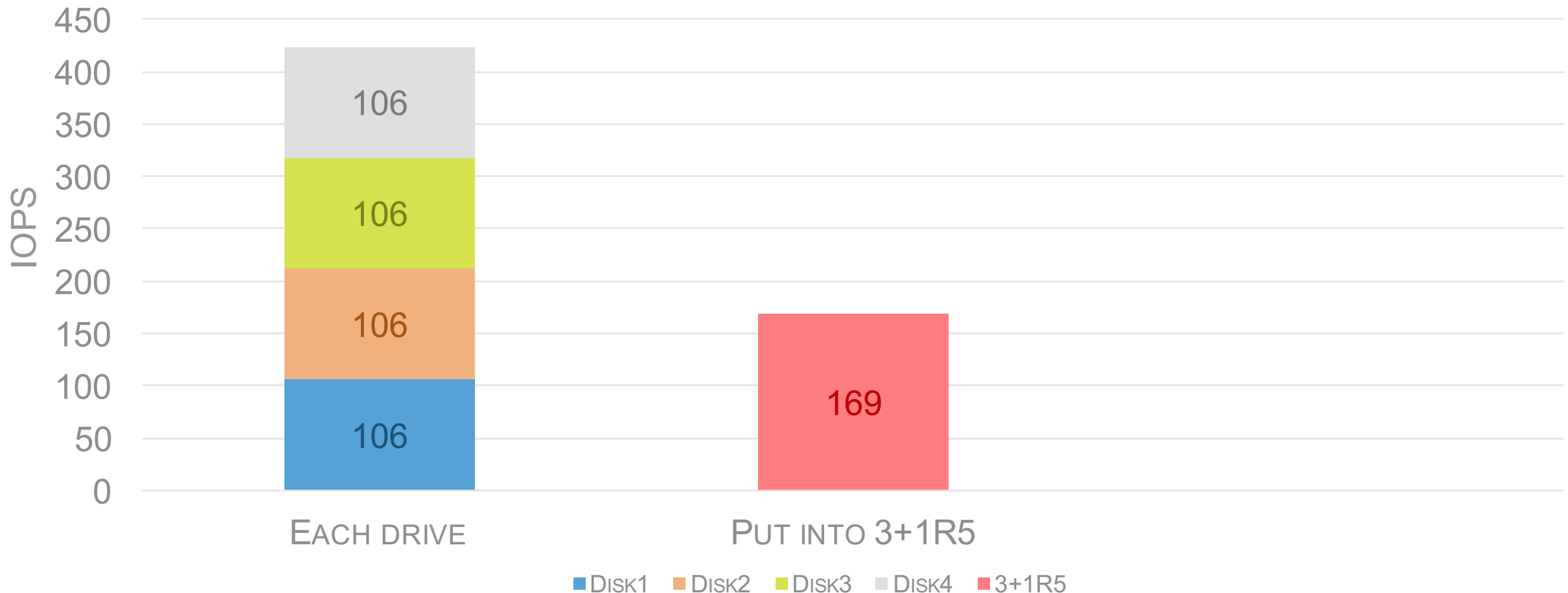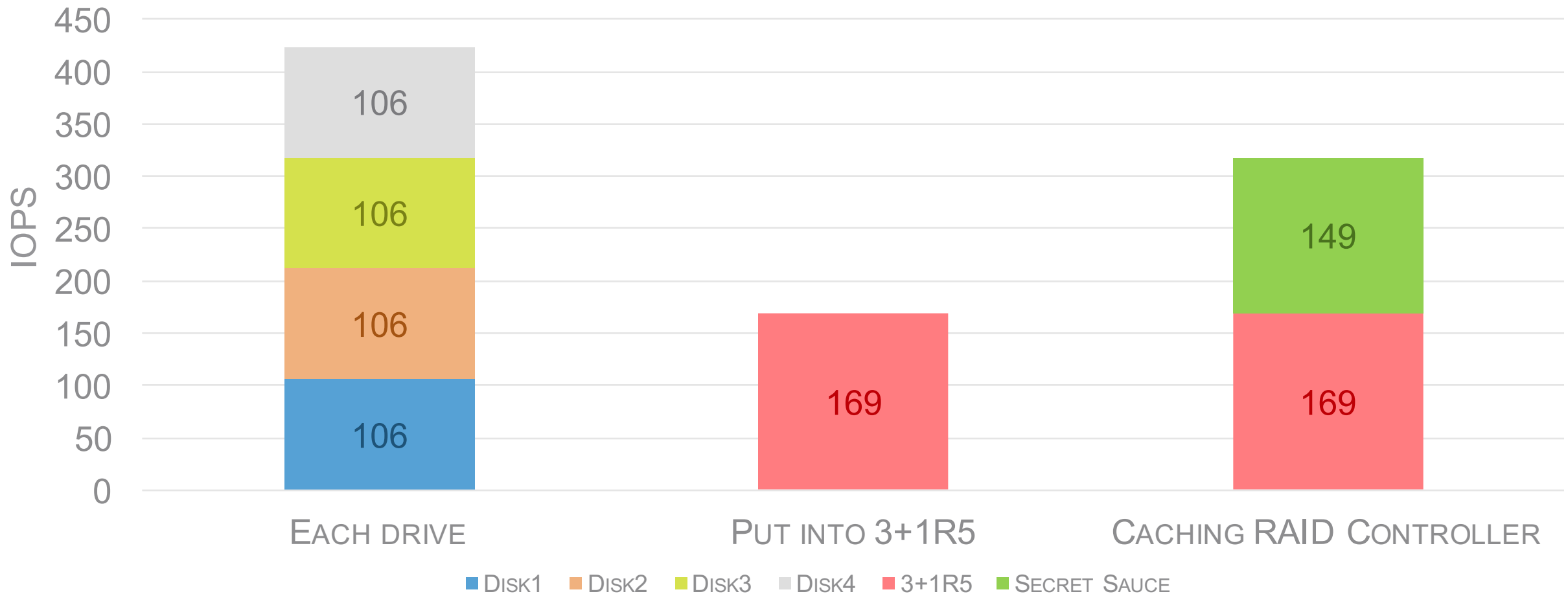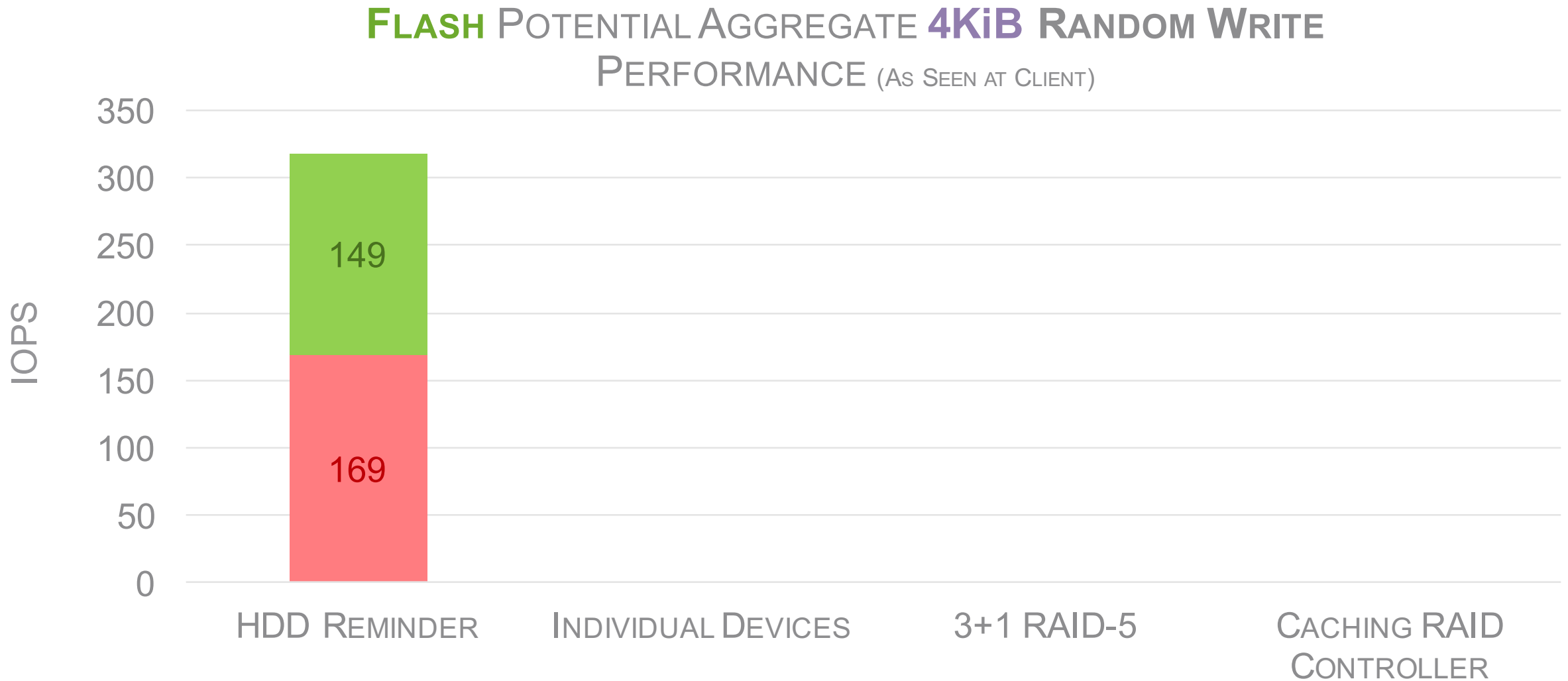| HDD REMINDER | INDIVIDUAL DEVICES | 3+1 RAID-5 | CACHING RAID CONTROLLER |

# What "Really" Happens With RAID-5?

FLASH POTENTIAL AGGREGATE 4KiB RANDOM WRITE PERFORMANCE (AS SEEN AT CLIENT)

# What "Really" Happens With RAID-5?

**FLASH** POTENTIAL AGGREGATE **4KiB** RANDOM WRITE PERFORMANCE (AS SEEN AT CLIENT)

IOPS

25,000
20,000
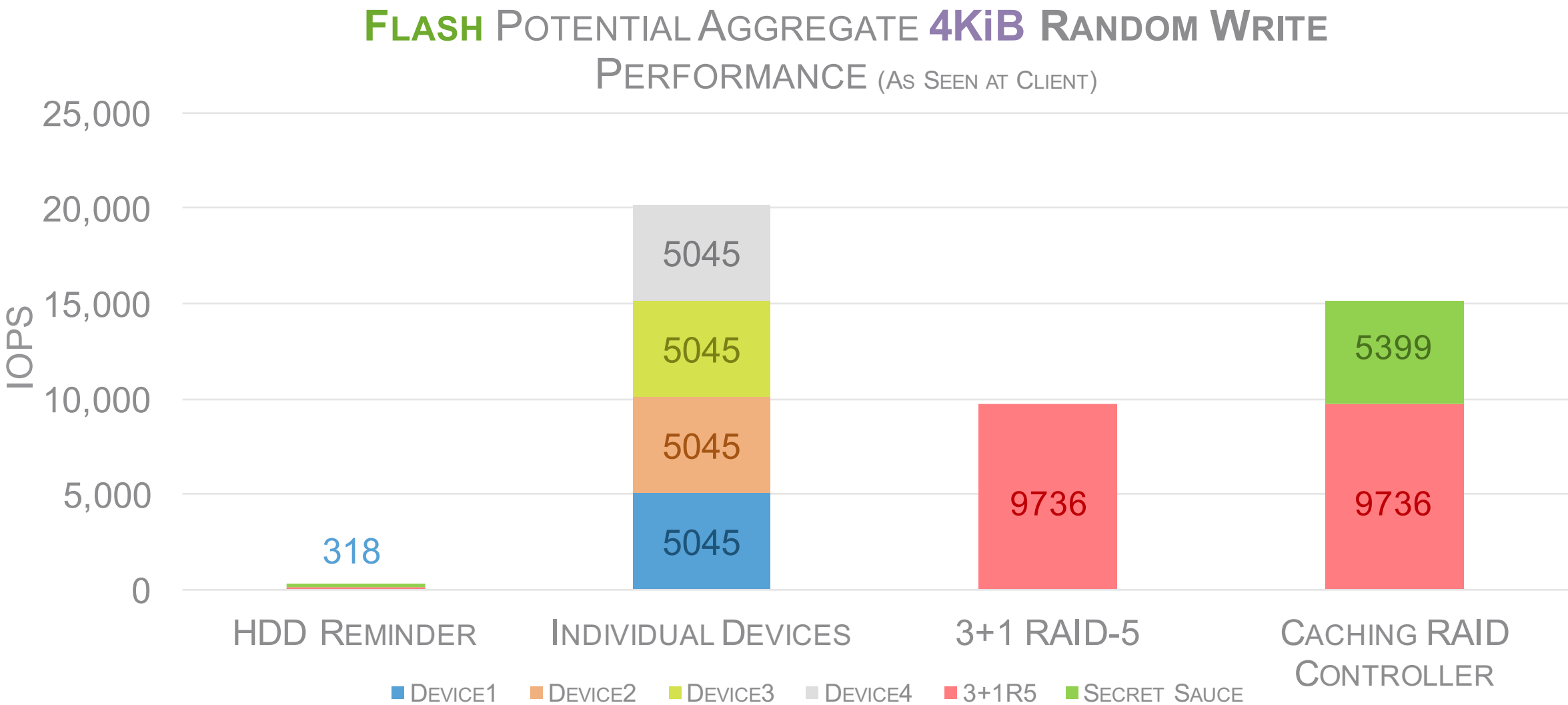15,000
10,000
5,000
0

HDD REMINDER — 318

INDIVIDUAL DEVICES — 5045 / 5045 / 5045 / 5045

3+1 RAID-5 — 9736

CACHING RAID CONTROLLER — 9736 / 5399

Legend: ■ DEVICE1 ■ DEVICE2 ■ DEVICE3 ■ DEVICE4 ■ 3+1R5 ■ SECRET SAUCE

# What "Really" Happens With RAID-5?

FLASH POTENTIAL AGGREGATE 4KiB RANDOM WRITE
(MiB/s, As Seen at Client)



MiB/s

| | |
|---|---|
| HDD REMINDER | 1 |
| INDIVIDUAL DEVICES | 79 |
| | 19.71 |
| | 19.71 |
| | 19.71 |
| | 19.71 |
| 3+1 RAID-5 | 38 |
| | 38.03 |
| CACHING RAID CONTROLLER | 59 |
| | 21.09 |
| | 38.03 |

Legend: ■ DEVICE1 ■ DEVICE2 ■ DEVICE3 ■ DEVICE4 ■ 3+1R5 ■ SECRET SAUCE

# Flash In The Real World

| PHYSICAL STORAGE | BACK-END CONNECT | STORAGE CONTROLLER | FRONT-END CONNECT | CLIENTS / HOSTS |
|---|---|---|---|---|
| 5685 MB/s | REQUIRES 10x 6Gb/s SAS | REQUIRES 6x PCIe 3.0 LANES | REQUIRES 4x 16Gb FC | HOW MANY HOSTS? |
| | 10 | 6 | 4 | ?? |

MB/s

# Flash In The Real World

| PHYSICAL STORAGE | BACK-END CONNECT | STORAGE CONTROLLER | FRONT-END CONNECT | CLIENTS / HOSTS |
|---|---|---|---|---|

NEW TECH

| | REQUIRES 10x 6GB/s | REQUIRES 6x PCIe | REQUIRES 4x 16 | HOW MANY HOSTS? |
|---|---|---|---|---|

5685 MB/s

MB/s

}10  }6  }4

**Flash In The Real World**

SNIA ESF | ETHERNET STORAGE

# BLOCK IS THE FOUNDATION
## OF STORAGE PERFORMANCE

# Agenda

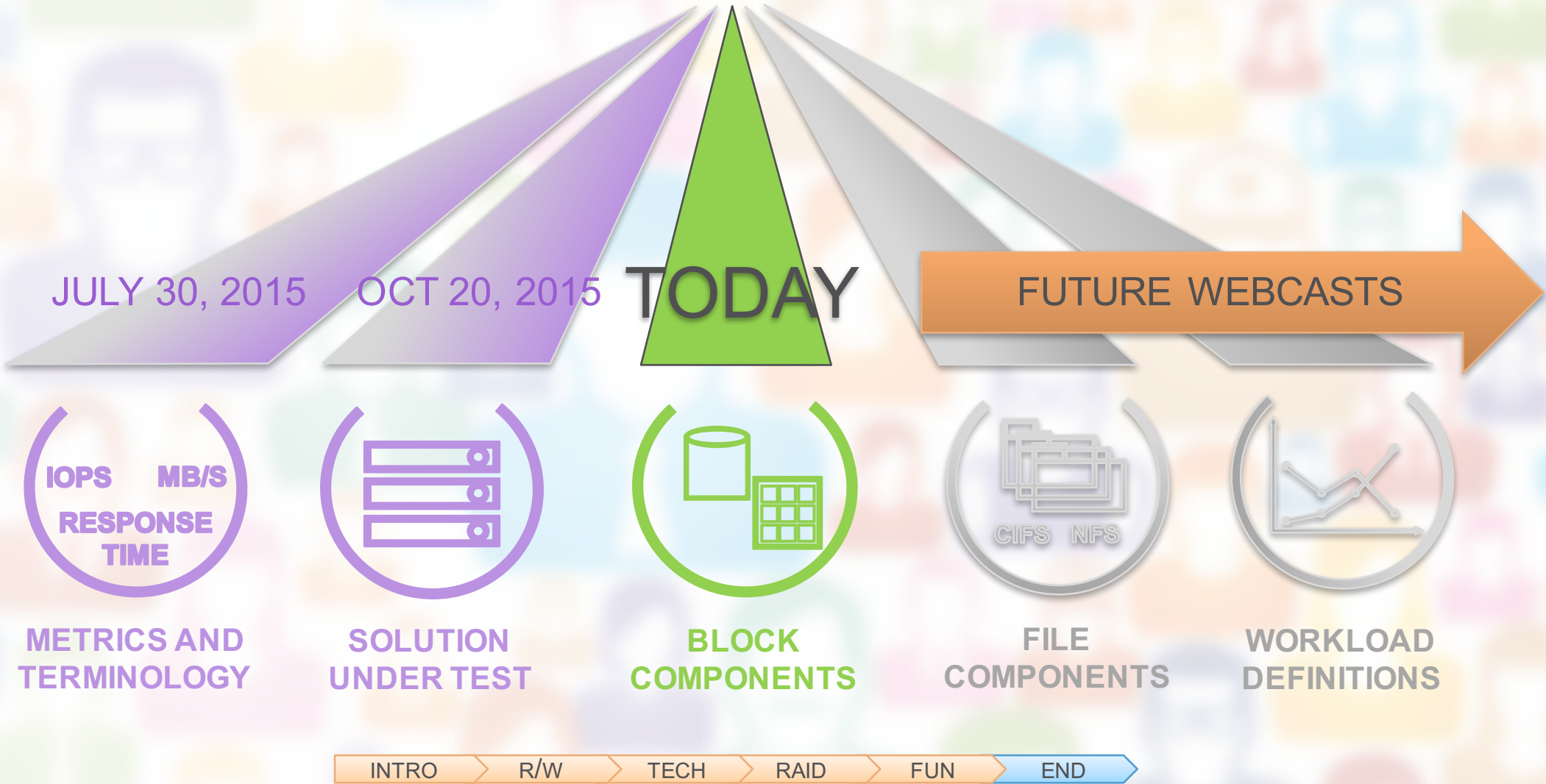| | |
|---|---|
| INTRO | INTRODUCTION |
| R/W | READING, WRITING; WHAT IS THE DIFFERENCE? |
| TECH | HOW DOES THIS TECH WORK ANYWAY? |
| RAID | WHAT IF YOU NEED MORE THAN ONE? |
| FUN | PERFORMANCE? |
| END | SUMMARY |

# Storage Performance Benchmarking

JULY 30, 2015   OCT 20, 2015   TODAY   FUTURE WEBCASTS

IOPS   MB/S
RESPONSE
TIME

**METRICS AND TERMINOLOGY**

**SOLUTION UNDER TEST**

**BLOCK COMPONENTS**

CIFS   NFS

**FILE COMPONENTS**

**WORKLOAD DEFINITIONS**

INTRO  >  R/W  >  TECH  >  RAID  >  FUN  >  END

# After This Webcast

- A PDF and a PPT of the slides for this and all previous parts of this Webcast series will be posted to the SNIA Ethernet Storage Forum (ESF) website and available on-demand
  - PPT and PDF: http://www.snia.org/forums/esf/knowledge/webcasts
  - Presentation Recording: https://www.brighttalk.com/webcast/663/164323
- A full Q&A from this webcast, including answers to questions we couldn't get to today, will be posted to the SNIA-ESF blog
  - http://sniaesfblog.org/

- Follow us on Twitter @SNIAESF, @RogovMark, @KenCantrellJr, @DrJMetz

- Next Webcast – Second Half of 2016
  - "Storage Performance Benchmarking: Part 4"

# Appendix – Additional Reading

# Appendix – More Reading

- SNIA S3 TWG Guide to SSD Performance: http://www.snia.org/sites/default/files/UnderstandingSSDPerformance.Jan12.web_.pdf
- SNIA S3 TWG SSD Performance Primer, 2013: http://www.snia.org/sites/default/files/SNIASSSI.SSDPerformance-APrimer2013.pdf
- Benchmarking methods for randomly sampling from a file, and why random seeks can (usually) hurt performance: http://simpsonlab.github.io/2015/05/19/io-performance/
- Excellent hard drive overview: https://www.backblaze.com/hard-drive.html
- SSD Performance results: http://www.tomshardware.com/charts/ssd-charts-2014/benchmarks,129.html
- SSD Performance results: http://www.anandtech.com/show/6433/intel-ssd-dc-s3700-200gb-review/3
- Intel Performance Benchmarking for PCIe* and NVMe* Enterprise Solid-State Drives: http://www.intel.com/content/dam/www/public/us/en/documents/white-papers/performance-pcie-nvme-enterprise-ssds-white-paper.pdf
- SSD M.2 Interface: http://arstechnica.com/gadgets/2015/02/understanding-m-2-the-interface-that-will-speed-up-your-next-ssd/
- More complete SSD interface article, covering NVMe, U.2 and M.2: http://blog.ocz.com/ssd-interfaces-sata-m2-u2-nvme/
- SSD vs HDD performance characteristics: http://www.tomshardware.com/reviews/ssd-gaming-performance,2991-3.html

- RAID
  - http://www.raid-recovery-guide.com/raid5-parity.aspx
  - http://rickardnobel.se/how-raid5-works/
  - http://igoro.com/archive/how-raid-6-dual-parity-calculation-works/
  - RAID Perf Calculator: http://wintelguy.com/raidperf.pl
  - RAID Reliability Calculator: http://wintelguy.com/raidmttdl.pl
  - RAID Failure Calculator: http://raid-failure.com/raid10-50-60-failure.aspx
  - RAID Survival Rate Simulation: https://linustechtips.com/main/topic/103179-lets-talk-about-raid-survival-rates/