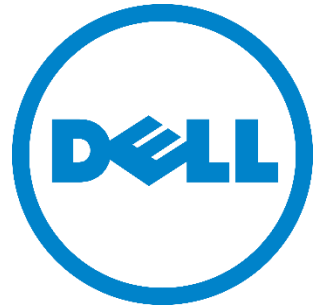# Something is changing in the Storage Panorama

Paolo Bianco
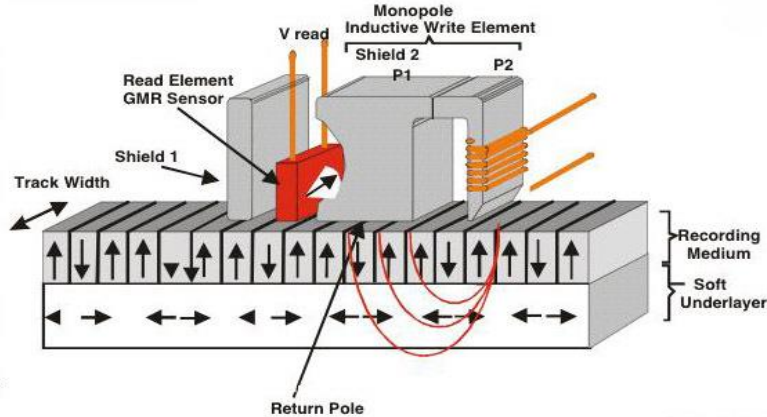Systems Engineer, Dell

# Agenda

- Nearline Capacity Drives Update

- SSDs and Performance Optimized Drives update

- Designing with SSDs

- Questions

# Nearline Capacity Drives Update

# Introducing PMR

- Data on conventional HDDs platters is written in circular, concentrical tracks (about 75nm wide), separed by guard spaces.

- Total track width is larger than necessary because write head poles needs to be large enough to generate sufficient coercitive force magnetization swap

- Effective read track width could be (and it is) smaller
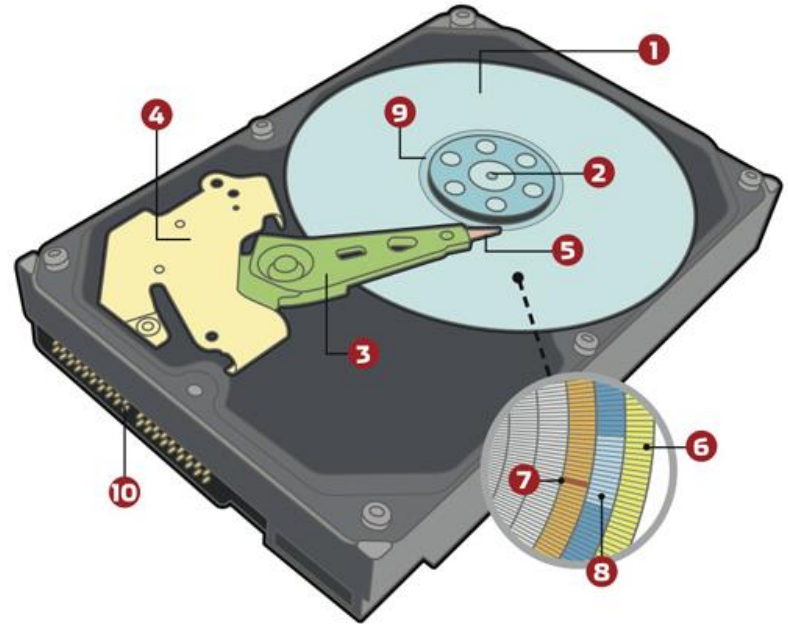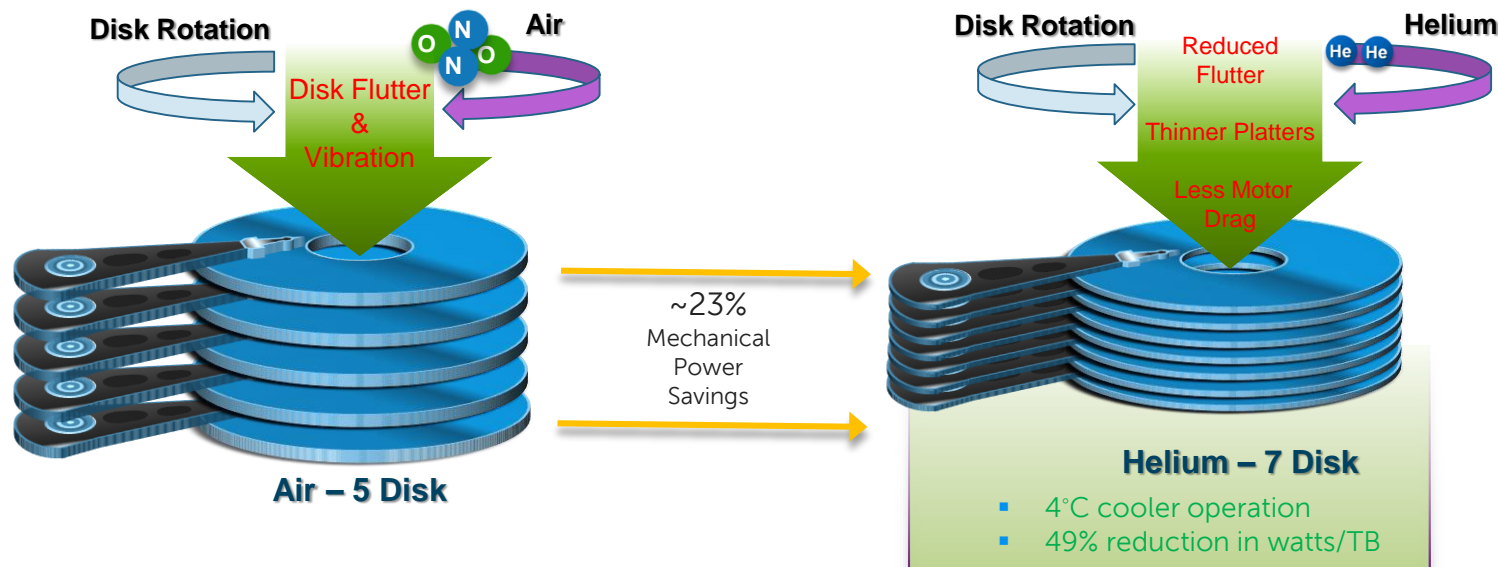
# PMR Technology has reached its own limits

- Diameter of platters: 3,5''
  - Total Surface 11 sq.in

- Useful Data Surface: about 5.5 sq.in

- Max Capacity per platter side: 0,68TB
  - @ 1Tb/sq.in

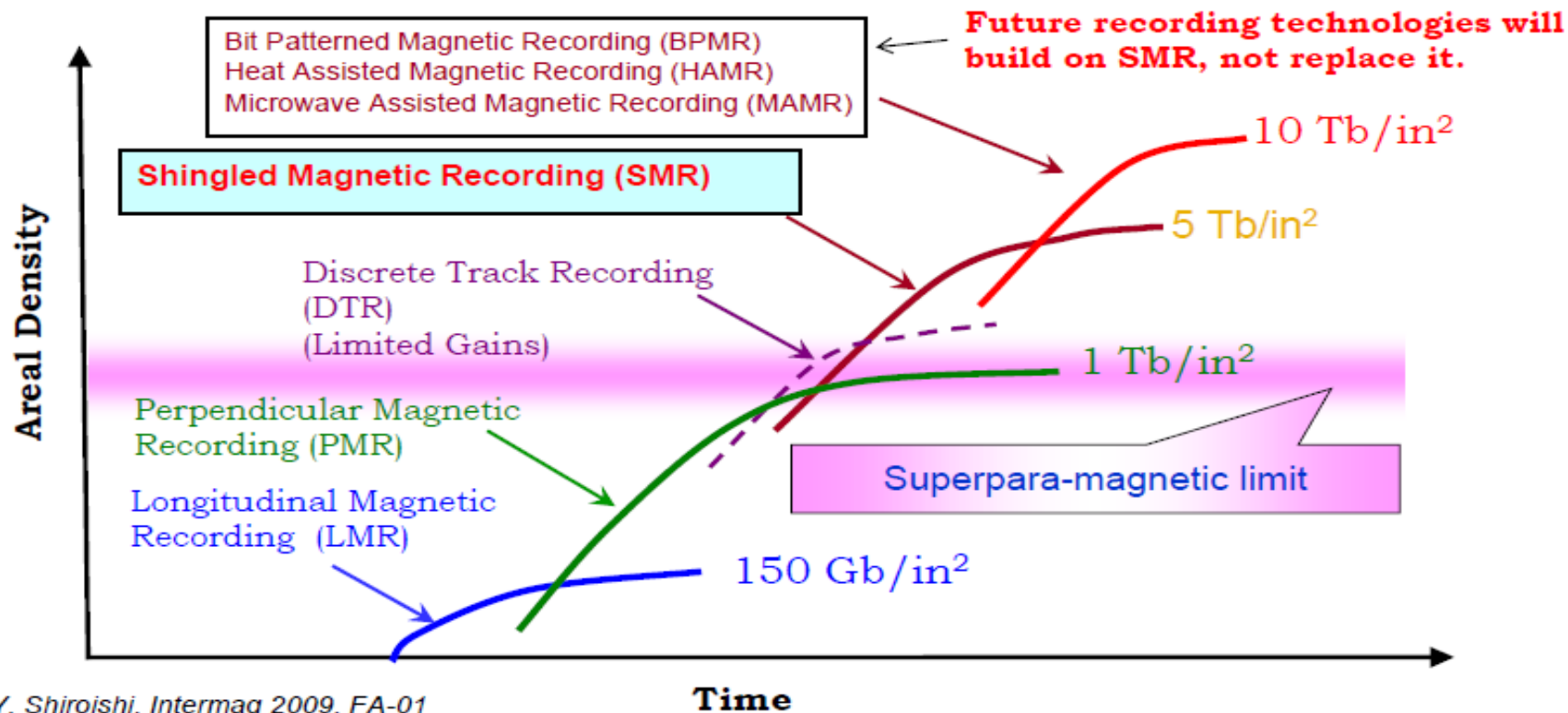- **Max Total Capacity (5 platters): 6,8TB**

# Helium-Filled: The last line of PMR Capacitive HDD

- Helium reduces mechanical power dissipated in air shear
- Allows platters to be placed closer together enabling more density
- 8TB He-Filled will probably be the last PMR-based cap.HDD generation on the market

# Magnetic Recording Technologies



Bit Patterned Magnetic Recording (BPMR)
Heat Assisted Magnetic Recording (HAMR)
Microwave Assisted Magnetic Recording (MAMR)

**Future recording technologies will build on SMR, not replace it.**

**Shingled Magnetic Recording (SMR)**

10 Tb/in$^2$

5 Tb/in$^2$

Discrete Track Recording (DTR) (Limited Gains)

1 Tb/in$^2$

Perpendicular Magnetic Recording (PMR)

Longitudinal Magnetic Recording (LMR)

Superpara-magnetic limit

150 Gb/in$^2$
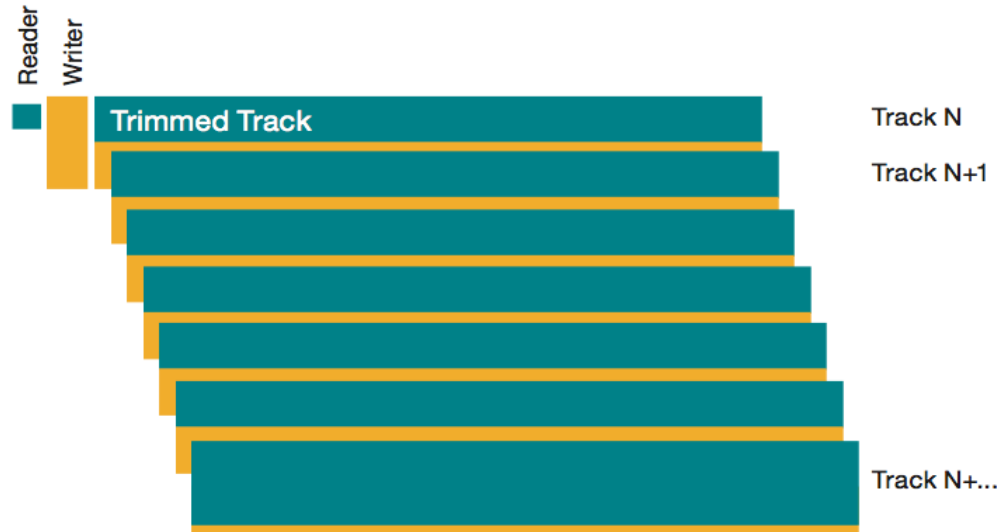
**Areal Density**

**Time**

Y. Shiroishi, Intermag 2009, FA-01

# Introducing SMR

- With Shingled Magnetical Recording  (SMR), clusters of tracks are superposed (just like «Roof Shingles») so that unnecessary track width space is recovered.

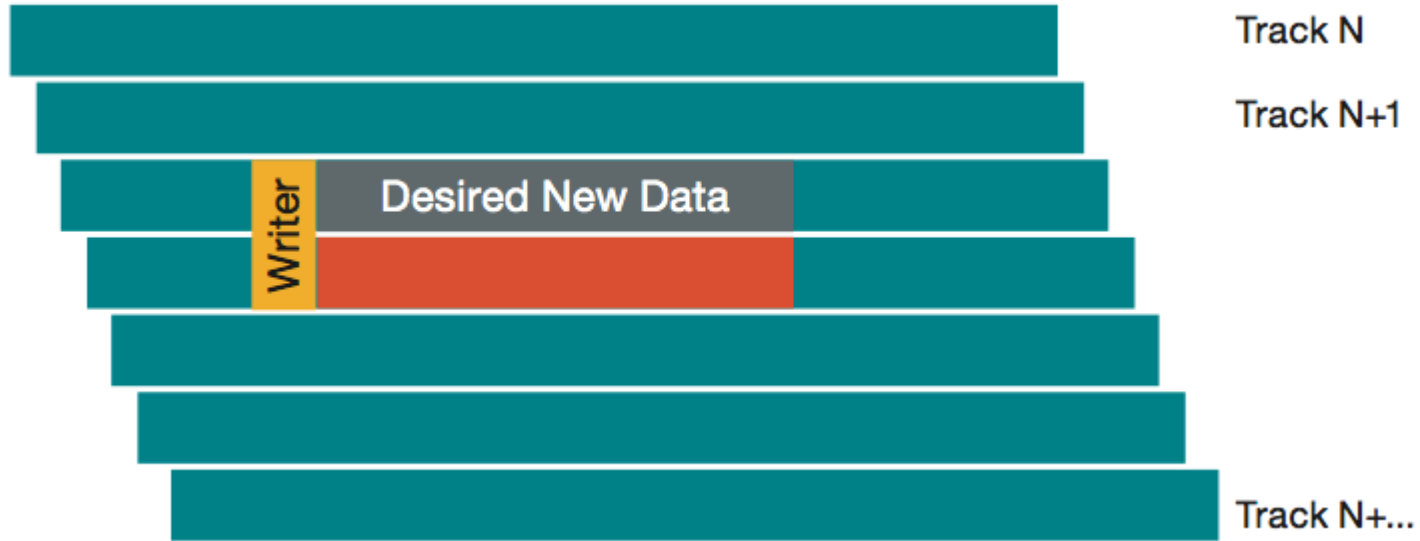**SMR Writes**

Reader

Writer

| Trimmed Track | Track N |
| | Track N+1 |

Track N+...

# SMR Disadvantages

- When a new data is written, new data track overwrites subsequent tracks...



Track N

Track N+1

Writer

Desired New Data

Track N+...

# SMR Disadvantages

- We need then to Load in a buffer all data following the new track in a cluster…



Desired New Data

Writer

Band N

- … and write down back again the cluster tracks starting from new data (aka R/M/W Penalty)
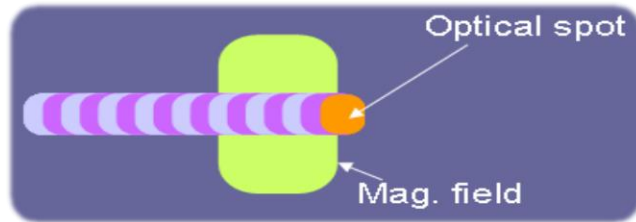
# SMR Challenges

- How to avoid performance loss (sustained data rate) due to Read-Modify-Write?
  - Onboard controller defragmentation (just like SSDs)
  - OS-aware SMR media management  (just like SSDs!)
  - Move BTL in Filesystem as exploring with FTL
  - T10 standards group working on this

- Short-term media capacity growing technology:
  - Move to HAMR in (probably) 3/5y

- Will all storage arrays manufacturers work on SMR-awareness?
  - Or will they try to just mitigate performance gaps with SSD caching?

# Near Future:Heat-Assisted Magnetic Recording

## HAMR : A Whole New Recording System

- Density growth limited by ability to make smaller bits thermally stable
- HAMR combines laser and magnetic field to write the media
- Allows for use of much higher coercivity media and hence enables higher densities



$$\frac{dH_{eff}}{dx} = \frac{dH_k}{dT} \cdot \frac{dT}{dx}$$



**Industry projecting the introduction of HAMR technology around 2018**

# HAMR is not too far....

- The right is a photo of an actual HAMR drive. You can tell is a HAMR drive because it has the Laser Warning Stick stuck in front of it
- Below is a picture of an integrated HAMR head including the laser (not the same head used in the drive)
- First fully-functional public HAMR Drive demo run in Sept. 2012 by Seagate.

# Key Takeaways

- Conventional Perpendicular Magnetic Recording technology  has reached its maximum areal density limit

- 8TB NL drives will be the latest capacity drives based on standard Magnetic Recording technology (PMR)

- The drive industry is introducing a new areal density enabling technology called Shingled Magnetic Recording (SMR).

- This technology will partially alter the throughput and response time behavior of IO, especially for random writes.

- SMR is a transition technology toward HAMR, which is expected to appear in the next 3-4years (if nothing changes in Solid State memory market...)

# SSDs and Performance Optimized Drives Update

# Performance Optimized Enterprise HDDs

# SSD predictions in 2008

Prediction 1: SSD will be bigger in capacity in 2015

# SSD predictions in 2008

Prediction 2: SSD will be more cost effective in 2009

# Performance HDDs replaced by SSDs



Projection 2015-2020 of 4-year Cost of Capacity Disk & NAND Flash

# NAND Memory Technologies

- Single Level Cell (SLC)
  - 1 bit/cell
  - Fastest
  - 100k P/E cycles

- e/HET Multi Level Cell (eMLC)
  - 2 bits/cell
  - Slightly Slower Writes
  - 30/40k P/E cycles

- Multi Level Cell (MLC)
  - 2 bits/cell
  - Slow
  - 10/20k P/E cycles

- Triple Level Cell (TLC)
  - 3 bits/cell
  - Slower
  - 3/5k P/E cycles

# SSDs: anatomy of a NAND Chip

## Asymmetrical access Storage Media

- **Asymmetrical Read/Write**
  - Read per page
  - Write per block (64 pages)
  - Block needs to be erased before a new write can occur

- **Read/Modify/Write Penalty**
  - Unavoidable!

- **P/E cycles wears out the media**
  - Electrical charges get trapped in the dielectric



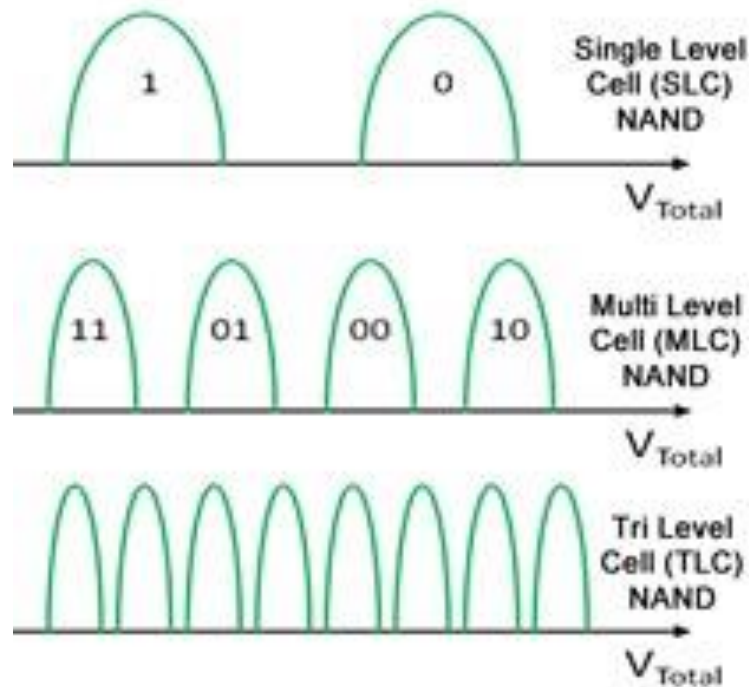Serial Input (x8 or x16) 30ns (max clk)

Serial Output (x8 or x16) 30ns (max clk)

Register
2112 bytes

Program: ~300µs

Read (page load): ~25µs

NAND Memory Array
NAND Page 2112 bytes

2048 Blocks (2Gb device)

64 pages per block

NAND Block

Block Erase -2ms

Data Area 2048 bytes

Spare Area (ECC, etc.) 64 bytes

# Dealing with Wear Out

- Spare Capacity (2002)
- Wear Leveling
  - Distribute data to even out the use of cells. (2003)
  - Background Read Data Refresh (HET, 2012)
- Error Correction Coding
  - BCH (2005), LDPC(2010), Polar (2012)
- Compression (2011)
  - Lempel-Ziv or derivative
  - Reduces the effective amount of stored data
  - Transparent to the host!!
- De-Duplication (2012)
  - Reduces the effective amount of stored data
  - Transparent to the host!!
- Endurance Coding , aka Data Shaping (2013)
  - Transform input data into shaped data having less "0"
  - Minimize the number of programmed cells per P/E cycle
- Increase in die/chipset capacity

# Steering away from SLC vs MLC discussion

## *Focusing on use profiles*

- Write Intensive SSDs
  - Mainly SLC
  - Highest longevity
  - Highest cost
- Read Intensive
  - MLC/eMLC
  - Lower longevity (but not affected by reads)
  - Lowest cost
- Multi-Use
  - Mainly eMLC
  - High performance and longevity
  - Medium cost

# Flash-Optimized SSD Comparison

| Storage Use | Write Intensive | | Read Intensive |
|---|---|---|---|
| Market Terminology | Write Intensive (WI) | Mixed Use (MU) | Read Intensive (RI) |
| Workload | Mainstream Applications Any usage | | Mostly Read 90/10 R/W Mix |
| Capacity | 200 / 400 GB | 800 GB | 480 / 1600 GB |
| Endurance (Full writes / day) | 30-10 | | 3 |
| Endurance (written PBs) | Up to 10 / 20 PB | | Up to 8 PB |
| Random Read IOPS (*) | Up to 20K+ | | 14K+ |
| Random Write IOPS (**) | 11K+ | 8K+ | 4K+ |
| Sustained Write Bandwidth (***) | 200 – 250 MB/s | 150-225 MB/s | 50 – 100 MB/s |
| List $/GB | Up to $20 | $11 | $4 |

**Managed NAND**

NAND Flash

Managed NAND Controller
• ECC
• Bad Block Management
• Wear levelling

# Dell 2.5″ SAS SSD Roadmap



**Write Intensive**

- 2.5″ SAS SED/FIPS TBD 6Gb/12Gb)
- 2.5″ 1600GB+ SAS 12Gb
- 2.5″ 200/400/800GB SAS 12Gb — 13G
- 2.5″ 200/400GB SAS 12Gb (running at 6Gb) — 12G

**Mix Use**

- 2.5″ 3.XTB+ SAS 12Gb
- 2.5″ 200/400/800/1600GB SAS 12Gb — 13G
- 2.5″ 200/400/800GB SAS 12Gb (running at 6Gbs) — 12G

**Read Intensive**

- 2.5″ 3.X-4TB+ SAS 12Gb
- 2.5″ 800/1600GB SAS 12Gb — 13G
- 2.5″ 800/1600GB SAS 12Gb (running at 6Gbs) — 12G

| CY14 | CY15 | | | | CY16 | | | |
|------|------|------|------|------|------|------|------|------|
| Q4 | Q1 | Q2 | Q3 | Q4 | Q1 | Q2 | Q3 | Q4 |

Investigate (Pre-APP) ▲   Develop (BC✓) ▲   Launch (RTS✓) ▲

Dell – Restricted – Confidential
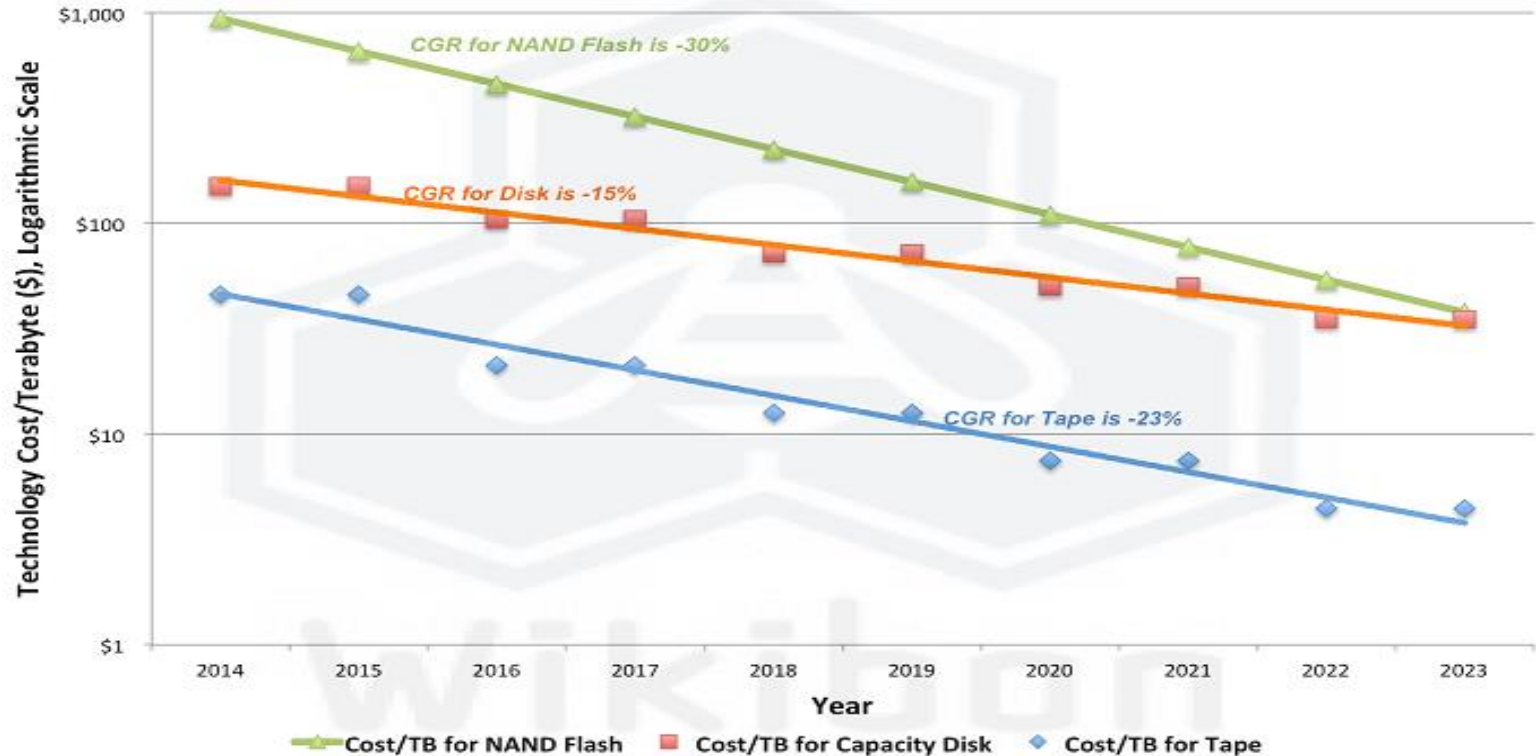
# SSDs and HDDs will still convive for a long time

# Designing with SSDs
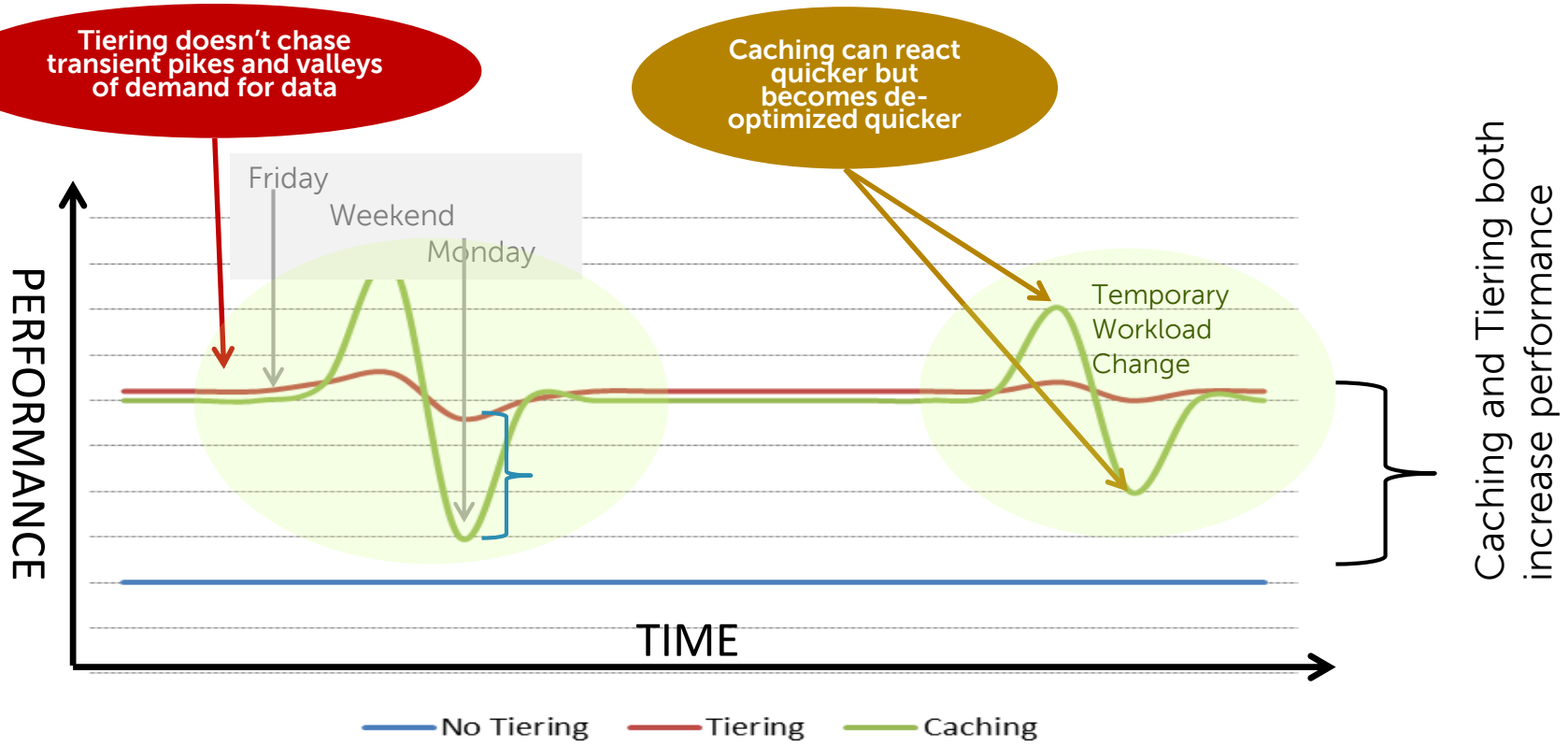
# SSD as Cache

**What is it?**

- SSD cache moves data from an HDD virtual disk to the SSDs following a host read or write.

- Subsequent host read of the same LBAs can be read directly from the SSDs with a <u>much lower response time</u> than re-reading the data from the HDD virtual disk.

- **All PV MD36XX/38XX acquired by INFN can do this**

**Workload characteristics that benefit from SSD Cache**

- Performance limited by HDD IOPs

- High percentage of Reads vs Writes / Large number of reads with intrinsic localty (repeated reads to the same or adjacent logical area of the LUN)

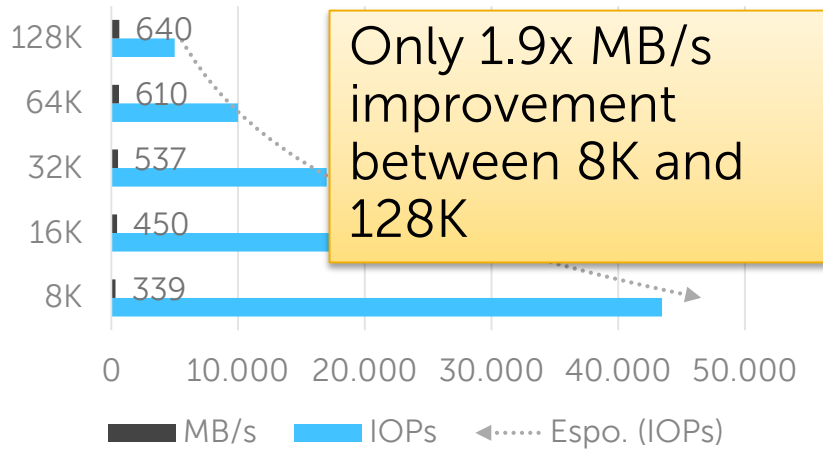- The working size set that is repeatedly accessed is smaller than the SSD cache capacity.

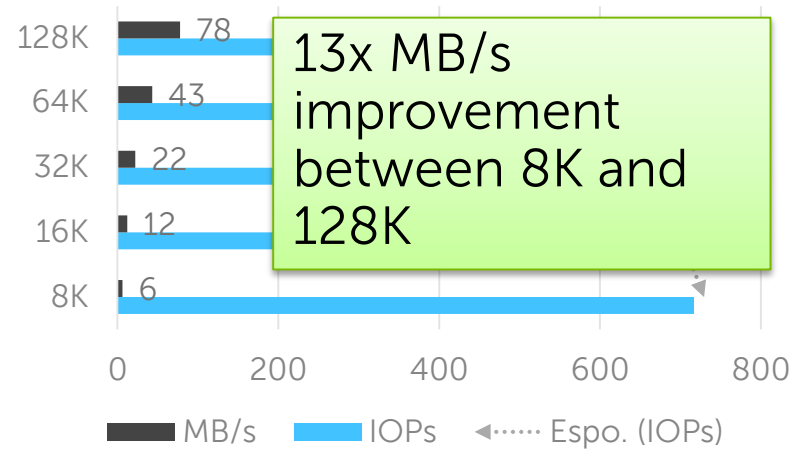# Caching is NOT Tiering (and vice-versa)...

# SSDs have very high Write potential throttled by MB/s

## 100% Random Writes Raid 10



### Write-Intensive SSDs

| Block | MB/s |
|-------|------|
| 128K  | 640  |
| 64K   | 610  |
| 32K   | 537  |
| 16K   | 450  |
| 8K    | 339  |

Only 1.9x MB/s improvement between 8K and 128K

Axis: 0, 10.000, 20.000, 30.000, 40.000, 50.000

Legend: ■ MB/s  ■ IOPs  ◄······ Espo. (IOPs)

### 15K HDDs

| Block | MB/s |
|-------|------|
| 128K  | 78   |
| 64K   | 43   |
| 32K   | 22   |
| 16K   | 12   |
| 8K    | 6    |

13x MB/s improvement between 8K and 128K

Axis: 0, 200, 400, 600, 800
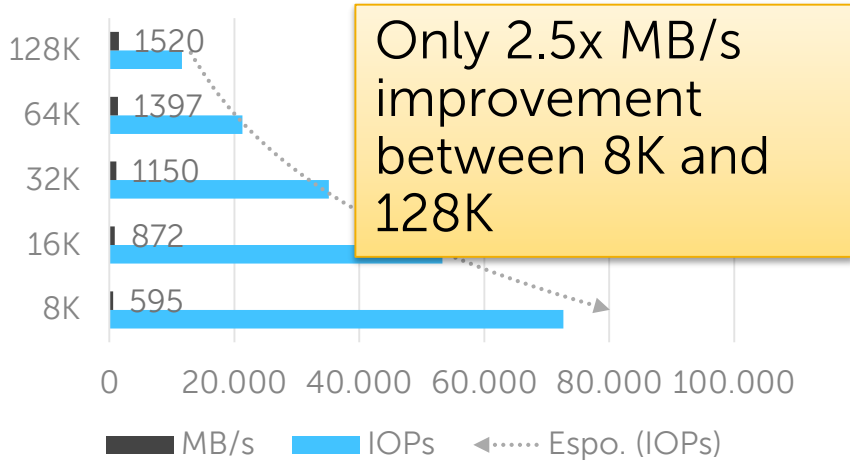
Legend: ■ MB/s  ■ IOPs  ◄······ Espo. (IOPs)

## Use caution when sizing to IOP/s per disk method with SSDs*

Expect about 16x MB/s improvement between 8K and 128K

# SSDs have very high Read potential throttled by MB/s
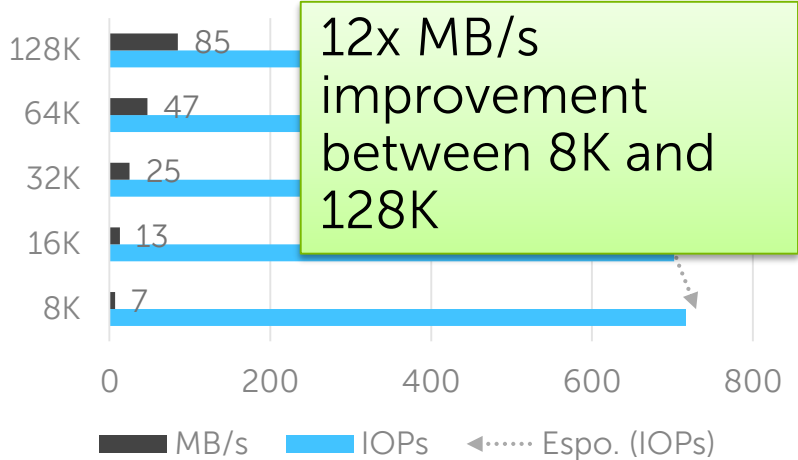
100% Random Reads Raid 5



Read-Intensive SSDs

| | |
|---|---|
| 128K | 1520 |
| 64K | 1397 |
| 32K | 1150 |
| 16K | 872 |
| 8K | 595 |

0    20.000   40.000   60.000   80.000  100.000

■ MB/s   ■ IOPs   ◀······ Espo. (IOPs)

Only 2.5x MB/s improvement between 8K and 128K

15K HDDs

| | |
|---|---|
| 128K | 85 |
| 64K | 47 |
| 32K | 25 |
| 16K | 13 |
| 8K | 7 |

0        200        400        600        800

■ MB/s   ■ IOPs   ◀······ Espo. (IOPs)

12x MB/s improvement between 8K and 128K

Sizing may be more appropriately based on MB/s for SSDs *
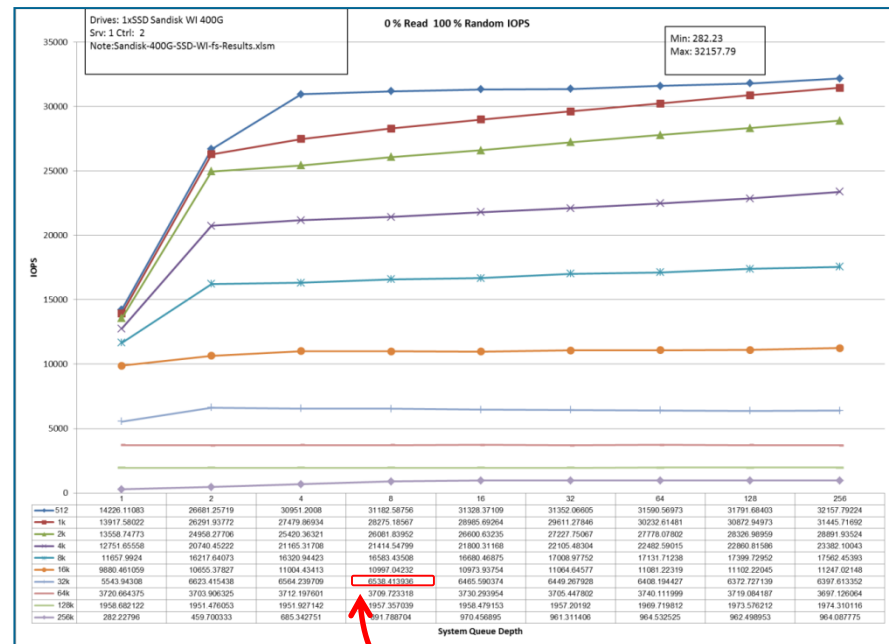
31

# Designing with SSD and Spinning Disk Drives

## Spinning Drives
- Need enough to provide needed IOPs

## SSD
- Need to provide enough throughput



BS=32KB ⇒ 215MB/s ⇒ 1075 MB/s (6x Pack) ⇒ /2 = 537.5 MB/s (Raid10)
Maximum IOPS = 537.5 * 1024 = 550.400 KB / 32 KB = 17.200 IOPS (~6538*5/2)

# Key Takeaways

- New Error Correction, Data Reduction and Cell Endurance algorythms can make Performance Drives replacement with SSDs a reality today.

- From now on, the drive industry will focus on SSDs more than 15k drives.

- 3,5" 15krpm drives will disappear shortly. Partial development will follow on 2.5" drives (lower seek time)

- Possible use of SSDs as Read Cache within INFN PV installed base. Read Cache, not Tiering!

- MB/s performance vs IO size scaling on SSDs not the same as HDDs. When dealing with SSD design, focus on MB/s performance and not on IOPS (If MB/s is a concern...)

# Questions?