# The Power of Memory

Tom Coughlin

President, Coughlin Associates,

[www.tomcoughlin.com](www.tomcoughlin.com)

President Elect, IEEE-USA

1

# Outline

- Digital Storage Drivers
- How Many IOPS are Enough?
- Digital Storage and Memory Technologies
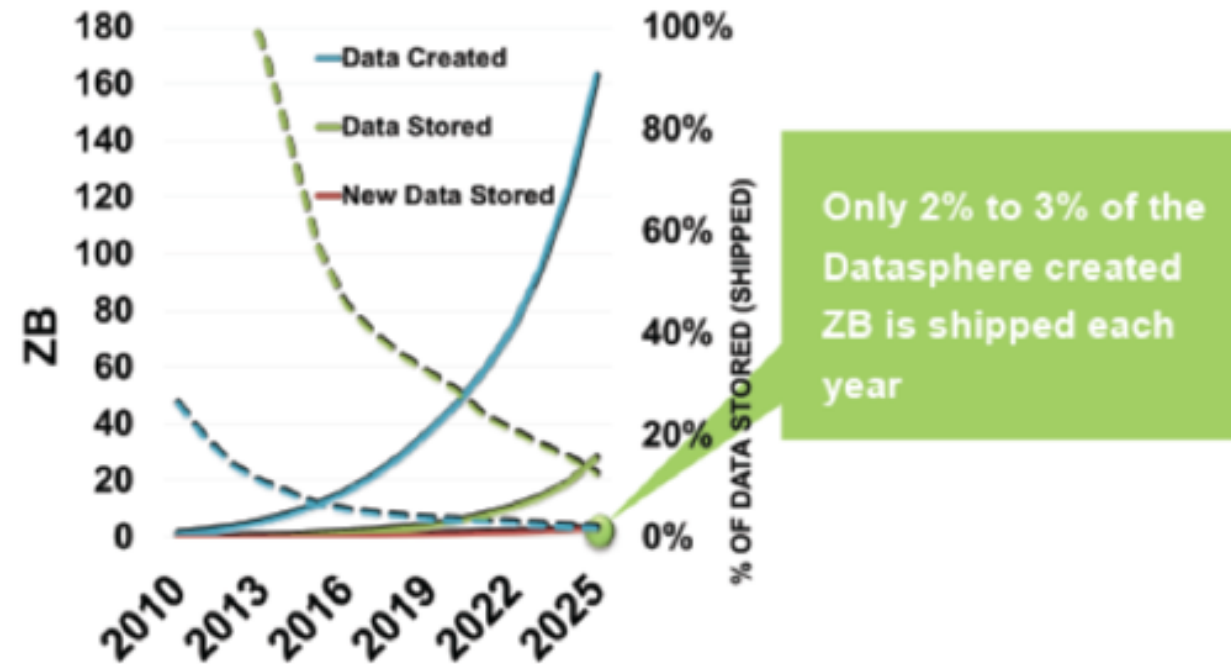- Conclusions
- References

# 50 minutes

- Abstract:  Advances in computing, communications, networking and electronics have enabled the generation of enormous amounts of data.  This is due to higher resolution data, captured more frequently, and the growing use of connected sensor data in IoT applications.  This data must be captured, processed and then stored.  Some data must be kept for a short time and processed quickly to be useful, e.g. in autonomous vehicles as well as robotics and automation.  The results of this processing and data generated having economic value (e.g. commercial video) may be kept for a long time.  This talk will examine the latest developments in the digital memory and storage hierarchy and how they will enable developments in AI, cloud computing, VR/AR, Smart Cities, robotics and automation and the preservation of culture, technology and history.

# Digital Storage and Memory Drivers

# Memory and storage needs are exploding

- Increasing storage demands—IDC 163 Zetabytes of data created by 2025 (16 ZB in 2016)

- New sources for unstructured data from media and entertainment, IoT, medicine, geo-science and big data

- Growth in local storage, storage at the edge (or the fog) and storage in large data centers (the cloud)

- There is a need for fast memory and storage to support processing and accessing this data and cheap storage to keep it for the long term
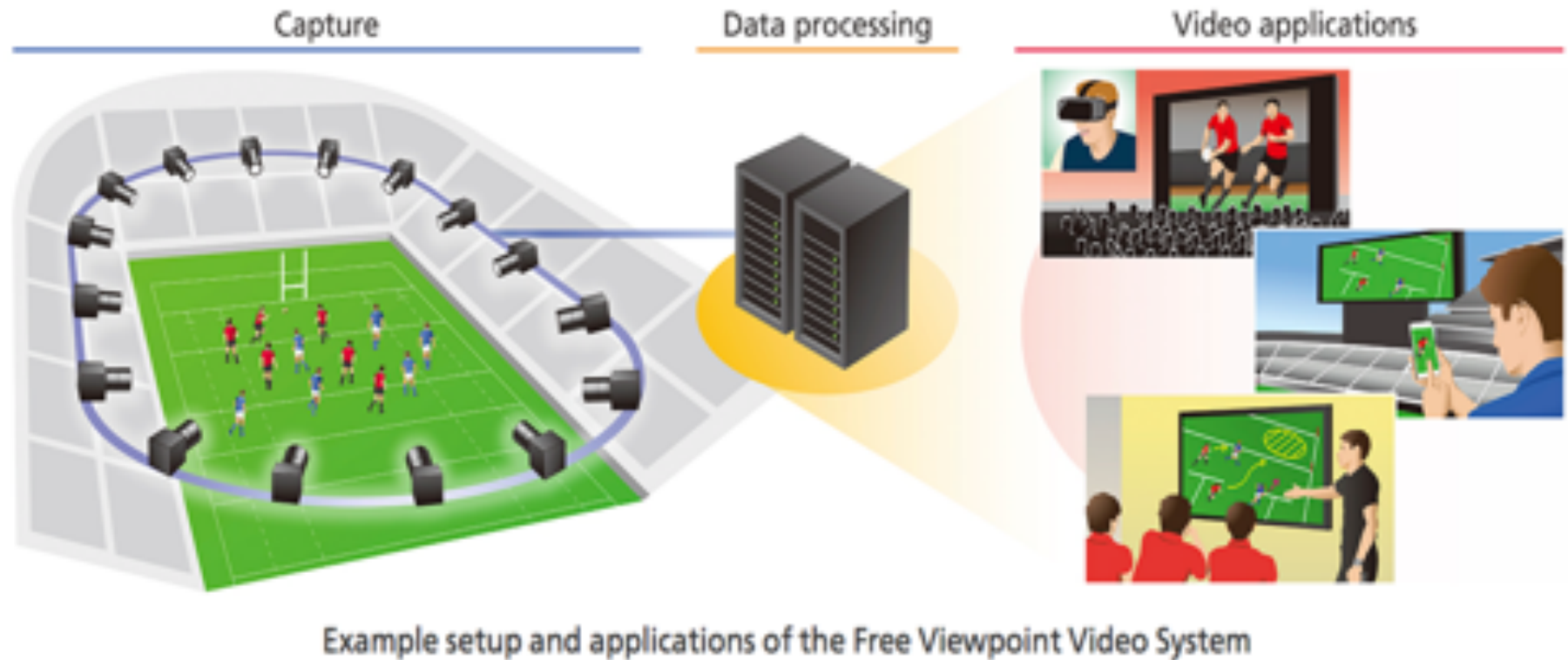


Only 2% to 3% of the Datasphere created ZB is shipped each year

# Example resolution, data rates and storage capacity requirements for professional media standards

| Format | Resolution (width X height) | Frame Rate (fps) | Data Rates (MB/s) | Storage Capacity GB/Hour |
|---|---|---|---|---|
| SDTV (NTSC, (8-bit) | 720 X 480 | ~30 | 31 | 112 |
| HDTV (1080p, 8-bit) RGB | 1920 X 1080 | 24 | 149 | 537 |
| UHD-1 4K (10-bit) RGB | 3840 X 2160 | 60 | 1,866 | 6,718 |
| UHD-2 8K (12–bit) RGB | 7680 X 4320 | 120 | 17,916 | 64,497 |
| Digital Cinema 2K (10-bit) YUV | 2048 X 1080 | 24 | 199 | 717 |
| Digital Cinema 4K (12-bit) YUV | 4096 X 2160 | 48 | 1,910 | 6,880 |
| Digital Cinema 8K (16 bit) YUV | 8192 X 4320 | 120 | 25,480 | 91,729 |

8K Ultra-HD may use more than 170X capacity of HD!

# Setup and application of Canon's free viewpoint video system (Volumetric video)

- 5-32 4K or higher resolution cameras surrounding an event
- Images are stitched together and allow rendering a view from anyway in the captured volume



Capture

Data processing

Video applications

Example setup and applications of the Free Viewpoint Video System
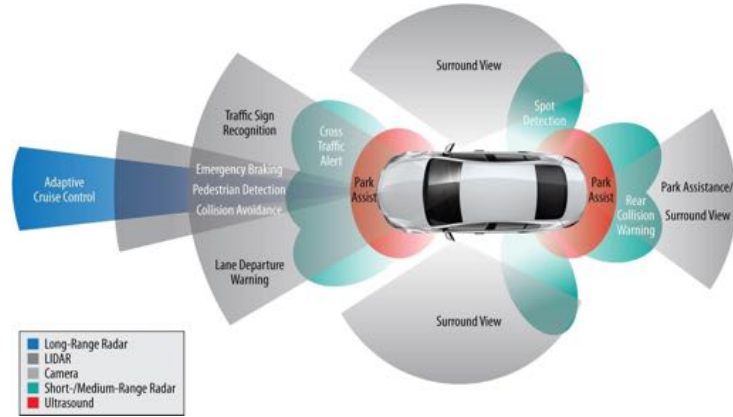
# 360° Video and VR Will Drive Content Growth



- VR and AR could drive 16k video (8K per eye)
- For a display as close as a pair of glasses, you can see the difference
- 8K per eye would give a more immersive experience

# Exabyte Video Projects Coming?

- Video at 16,000 X 8,000 pixel resolution, 24 bits/pixel, 300 fps raw video content could require **115 GB/s data rates and 414 TB/hour**. If 4 cameras were used to create data for a 360 degree presentation, the raw data would be **1.66 PB for an hour of content**

- Within 10 years we could have pro-video projects generating close to an exabyte of data

# Living the dream...



- Our personal devices are capturing more details of our lives than ever before and in greater details

- Our consumer electronics are connected to applications that can help us find our way, connect with others, maintain our health and make the world a better place

- In the future our cars may drive themselves, our homes make sure that we have what we need and our health and well being will be maintained by intelligent machines and applications

# Where is the data, where is the processing?

- Where data lives and where it is processed are dependent upon two important factors: latency to access data and speed and cost of processing

- Some data needs to be local where it is captured and processed because it is time critical (e.g. driver assistance) or important reference data

- Other data needs to be shared to be most valuable—such as local traffic information

- The network backbone has limited bandwidth, so local processing and intelligent compression can help manage traffic with the growing number of intelligent devices

- Thus we need a hierarchy of storage and processing, some local and some more remote
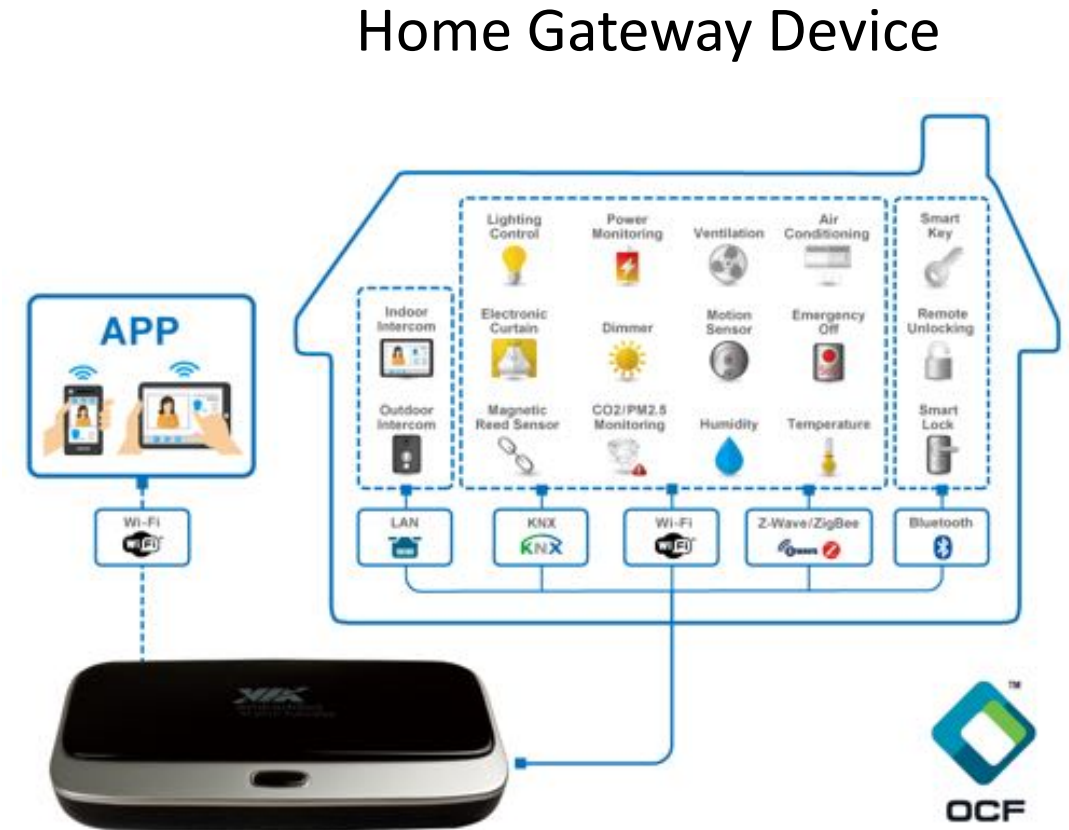
# The cloud and the fog





- The Cloud refers to compute resources, including storage, located in large data centers


- The Fog refers to local networks that connect thing (e.g. IoT) together
- Local fog networks may connect to the Internet

# In the device, in the fog or in the cloud?

- Data needing rapid access will be kept within a device

- Data having value from sharing and that can have someone longer latency may be in local networks (the fog)

- In a home IoT devices may be consolidated behind a gateway device (a local fog network)

- A home gateway may contain local process and storage in addition to connectivity

- In a metropolitan area the fog is part of a "smart city."

- Data for long term retention or that needs high performance processing can be in large data centers—the cloud

Home Gateway Device

# Mobile fog nodes



TechCrunch, October 28, 2016

- Around the world, wireless carriers are building all-new cellular networks for the Internet of Things.

- These new networks won't work with cell phones, they're made for IoT devices.

- Comcast, SoftBank, Orange, SKT, KPN, Swisscom and many others are building all-new nationwide IoT networks. Verizon and Vodafone are upgrading their networks, setting aside spectrum just for IoT. Cisco, Samsung, Nokia and Ericsson are selling equipment to make it work.

- New networks are necessary because cell phone networks fall short for IoT in three ways: battery life, cost and wireless coverage.

# 5G implementation



- In early 2018 Qualcomm demonstration 5G technology in San Francisco, CA and Frankfurt, Germany. In Frankfurt browsing jumped from 56 Mbps for a median 4G user connection to >490 Mbps for a median 5G user connection. In San Francisco these speeds increased from 71 Mbps to 1.4 Gbps.

- Significant investments in cell tower infrastructure will be required and telecom companies must agree upon standards for 5G networks.

- Will some cell towers be "data centers?"

# Fog computing visualization (Cisco)



Home Gateway

- Fog Data Services could include smart city services as well as commercial services
- Fog can include local networks between items in a home (home gateway) or other distributed products, such as V2V and V2X communication

# There may be many "edges" in the "fog"

- 5G is intimately tied to the Internet and other communications
- More layers of connectivity are envisioned with IoT devices resulting in multiple "edge" connectivity and possibly data storage and processing
- Eluv.io introduced the concept of a content fabric that uses blockchain authenticated peer to peer delivery of content versus today's highly centralized Content Delivery Networks (CDNs).

EXAMPLE - CDN VS CONTENT FABRIC

**LIVE STREAM**

8 versions / edits
3 formats
5 bitrates (4 sec seg) = 19.85 MB / seg

Master - 12 MB
20 Mbps - 10 MB
12 Mbps - 6 MB
6 Mbps - 3 MB
1.3 Mbps - 650 KB
400 Kbps - 200 KB

ALL BITRATES = 19.85 MB

**MANY BITRATES IN REGION**
(example: Japan)

only one version/edit
only push master

1 x 3 x 12 MB
36 MB per segment
32.4 GB per stream

**13x reduction**

**TYPICAL CDN PUSH**

8 x 3 x 19.85 MB = 476.4 MB per segment

900 segments per 1 h stream ==>

**428.76 GB per 1h stream**

**SINGLE BIT RATE REGION**
(example: India mobile)

only two versions/edits
only push low bw

2 x 3 x 850 KB
5.1 MB per segment
4.59 GB per stream

**93x reduction**

Michelle Munson, Eluvio presentation at the 2018 Creative Storage Conference (www.creativestorage.org)

# Autonomous Cars = Big Data

- Google autonomous car generates about 1 GB/sec.  Avg. US driver drives 600 hr/year.  This would generate 2 PB/car/year.

- The number of cars worldwide will surpass 1 billion

- So the potential size of data generated by autonomous cars is huge

- Most of the resulting long term storage is in the cloud, which also does the bulk of analytics

Increasing Storage Required

TODAY — NAVIGATION — ~10 m resolution

2020+ — AUTONOMOUS VEHICLE — ~10 cm resolution

https://company.here.com/en/newsroom/media-library/

ADAS ONLY — RADAR CAMERAS ULTRASONIC LIDAR — AUTONOMOUS VEHICLE — FULL AUTONOMOUS + REAL TIME DATA CAPTURE

Memory Development to Drive Autonomously, Kris Baxter, Micron, 2016 Flash Memory Summit

# Memory Trends in Automobiles

- The electric/electronic share of value added to a state-of-the-art vehicle is **already at 40 percent for traditional, internal combustion engine cars and jumps as high as 75 percent for electric or hybrid electric vehicles.**

- The total available market value for semiconductor memories in **automotive applications is expected to ….be well above the overall CAGR for the total memory semiconductor market.**

- **Automotives are about 5% of the memory market now but could grow to 10%** (2016 Micron Press Release)

Micron article in 2012 Embedded Computer Design

# Local Memory & Storage Needs

- 3D Maps
- Real Time Data Capture
- Black Box Data Recording
- On Board 4K Content
- **Greater than 1 TB capacity** will be common in new cars by 2020
- This is addition to storage in local and remote networks

# Location of Storage/Memory for Connected Cars

Fog Storage
Smart City

Cloud Storage

In-vehicle storage

Open Fog Consortium

# Connected Cars and Smart Cities

From: Tao Zhang, Cisco Distinguished Engineer, Co-Founder and Board Director of OpenFog Consortium

# Investments in infrastructure

- Faster data delivery speeds need faster equipment behind the radio delivery system

- Lower latency requirements will drive requirements for communication, processing and storage at the "edge."

- Changes to new distribution models, such as Eluvio's Content Fabric could change the nature of local connection nodes, including storage, processing and communication requirements

- In particular local content rendering, content caching and buffering will be key components in providing good content QoS

- This will lead to new investments in storage/memory architectures to support new ways to access, process and capture content and data

# Security in a connected age

- The data in our devices, in the fog and in the cloud can give a lot of information about us and our activities

- This can be used by people that mean us no good to rob or deceive us

- Encryption, good passwords and other forms of security must be important elements to protect personal data

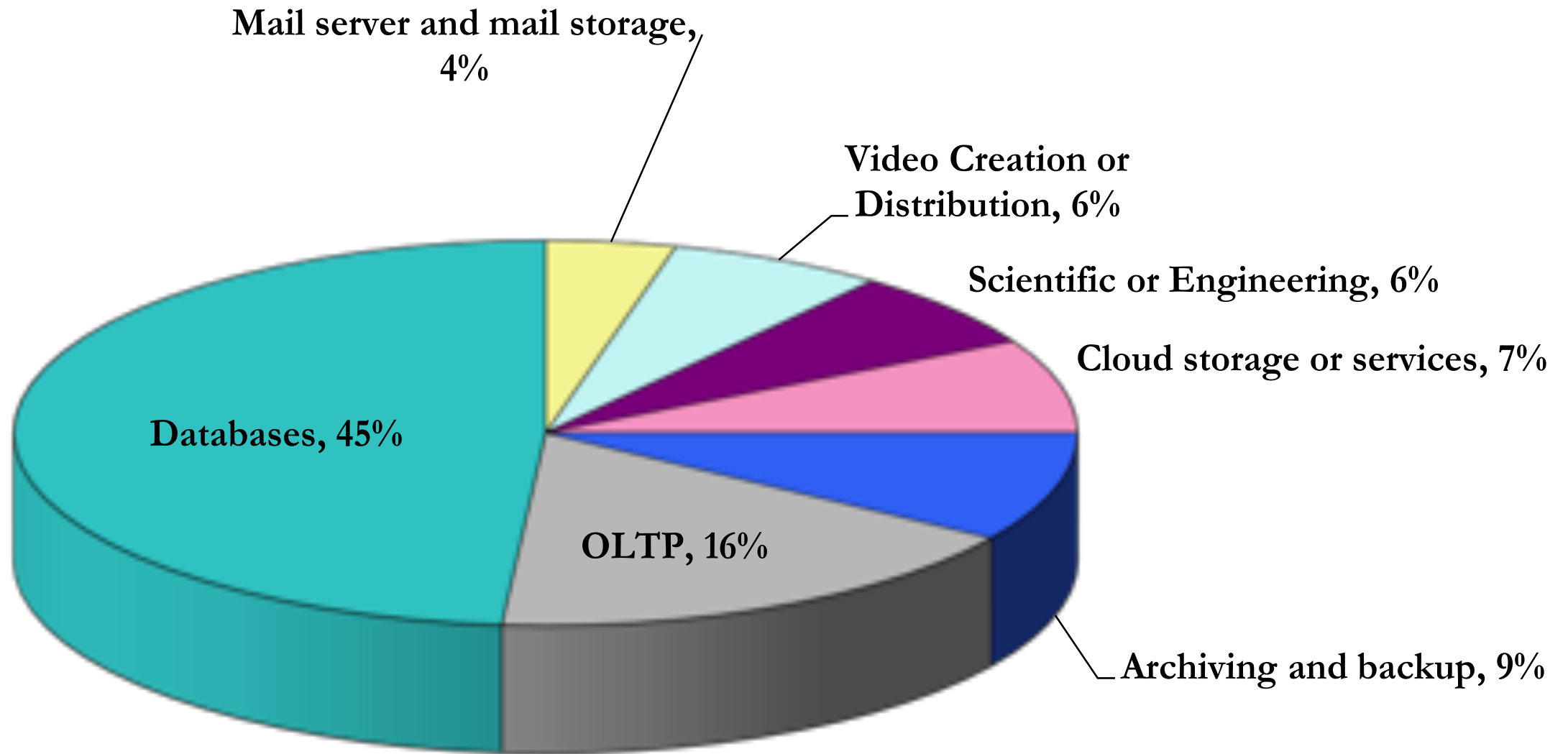- Anonymization may be the best way to share data and provide privacy

© 2018 Coughlin Associates

# How Many IOPS are Enough?

# Our Survey
(Objective Analysis and Coughlin Associates)

- Ongoing.  Take our survey at: http://TinyURL.com/IOPSsurvey

- IT participants participating

- Asks for IOPS, capacity and latency needs
  - Also their primary applications

- Some results are in a SNIA SSSI white paper

- We compared results from 2012 to those in in 2016

# Applications: 2012



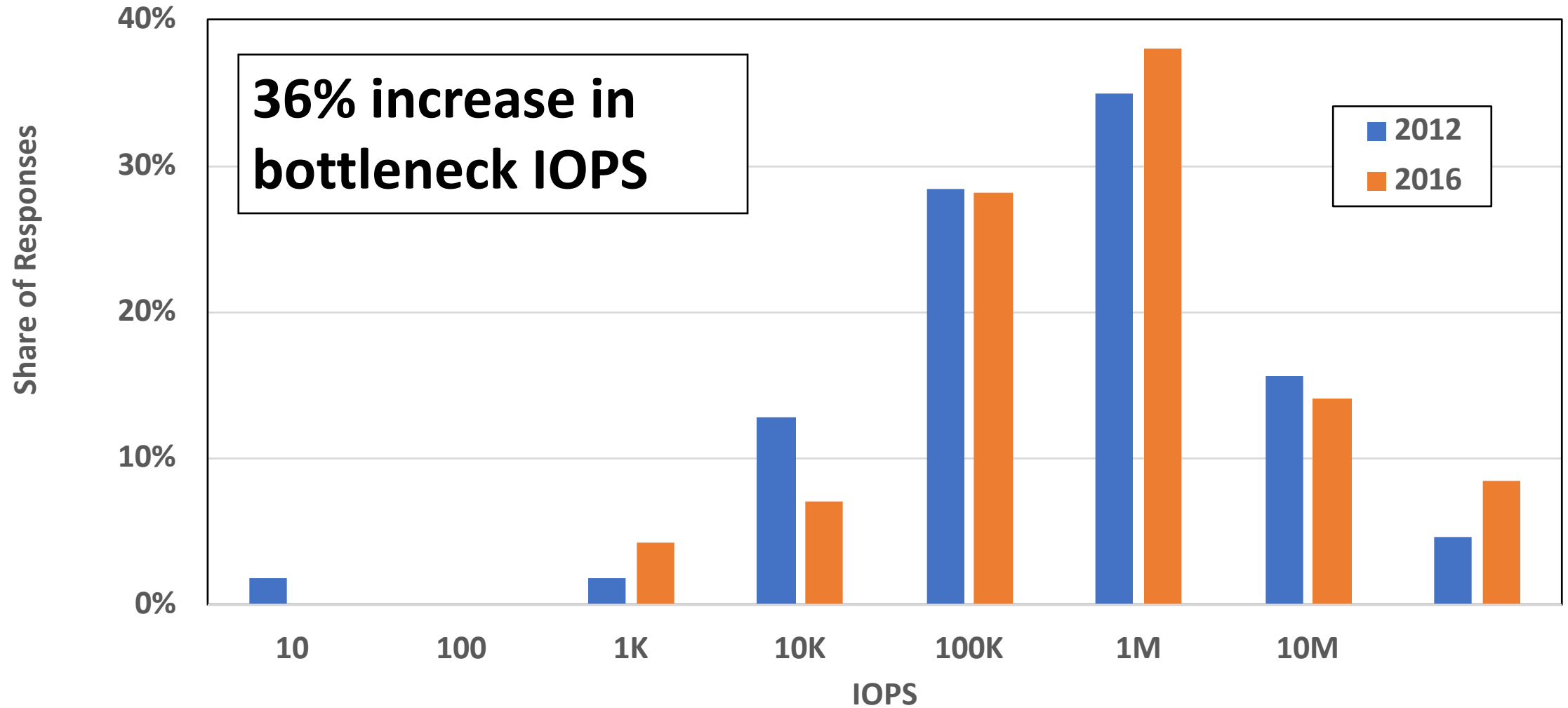Mail server and mail storage, 4%

Archiving and backup, 4%

Video Creation or Distribution, 7%

Cloud storage or services, 11%

Databases, 40%

OLTP, 24%

Scientific or Engineering, 10%

# Applications: 2016



Mail server and mail storage, 4%

Video Creation or Distribution, 6%

Scientific or Engineering, 6%

Cloud storage or services, 7%

Databases, 45%

OLTP, 16%

Archiving and backup, 9%

# IOPS Required for Dominant Application

# Capacity Required

# Other Hardware IOPS Bottleneck

**36% increase in bottleneck IOPS**

# Fastest Latency the System Can Use



73% decrease in mean latency

# IOPS by Form Factor

**HDD**

**SATA/SAS**

**NVMe/PCIe**

**Memory Channel**

$10^2$  $10^3$  $10^4$  $10^5$  $10^6$  $10^7$

# Digital Storage Capacity Projections



- The growth and processing of data will lead to the use of many types of digital storage

- SSDs will dominate for high performance storage and higher total revenue

- HDDs will be high capacity and used for colder storage

- Magnetic tape will be used by some organizations for the lowest cost (currently <1 cent/GB)

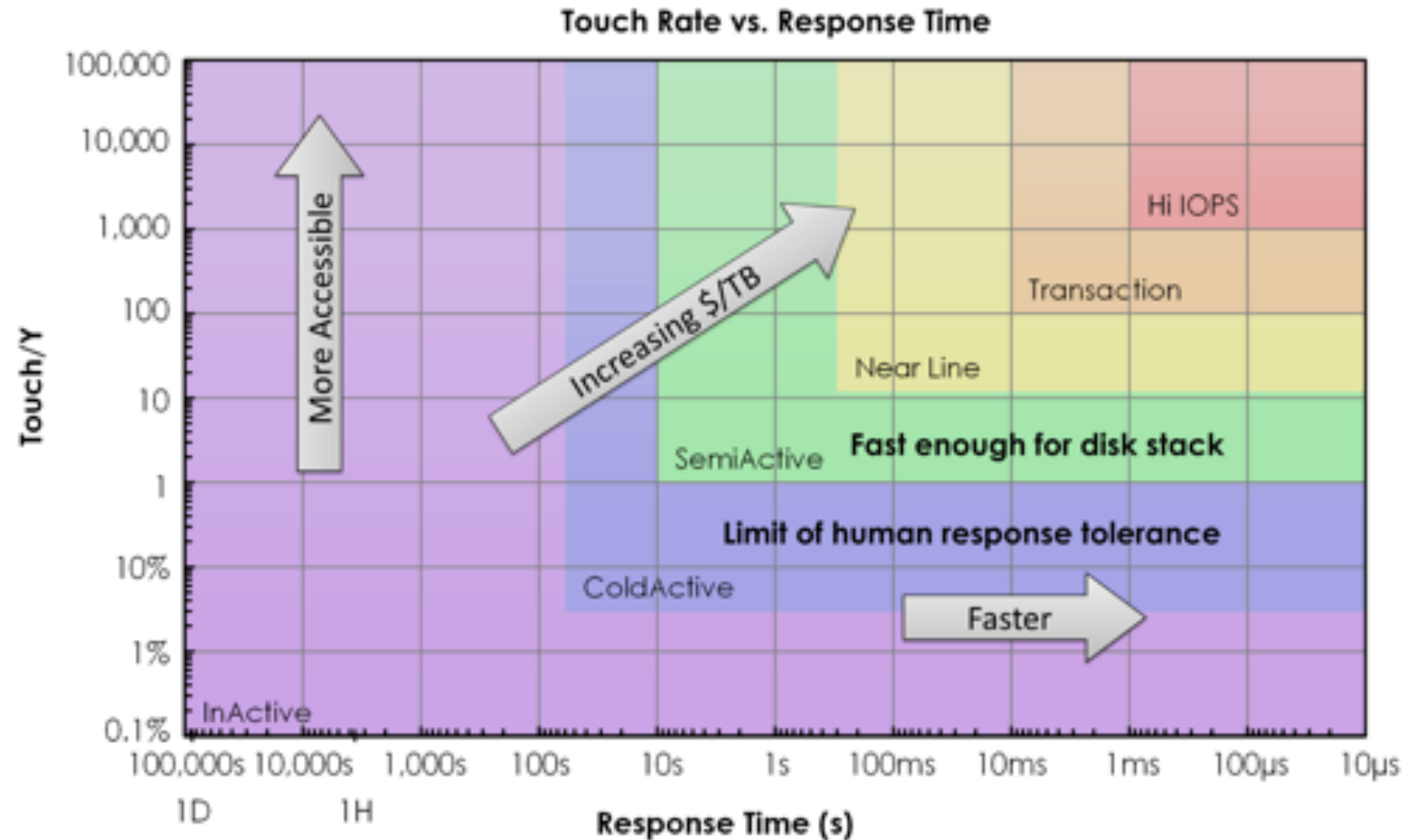# Digital Storage and Memory Technologies

# Digital Storage and Memory Tiering
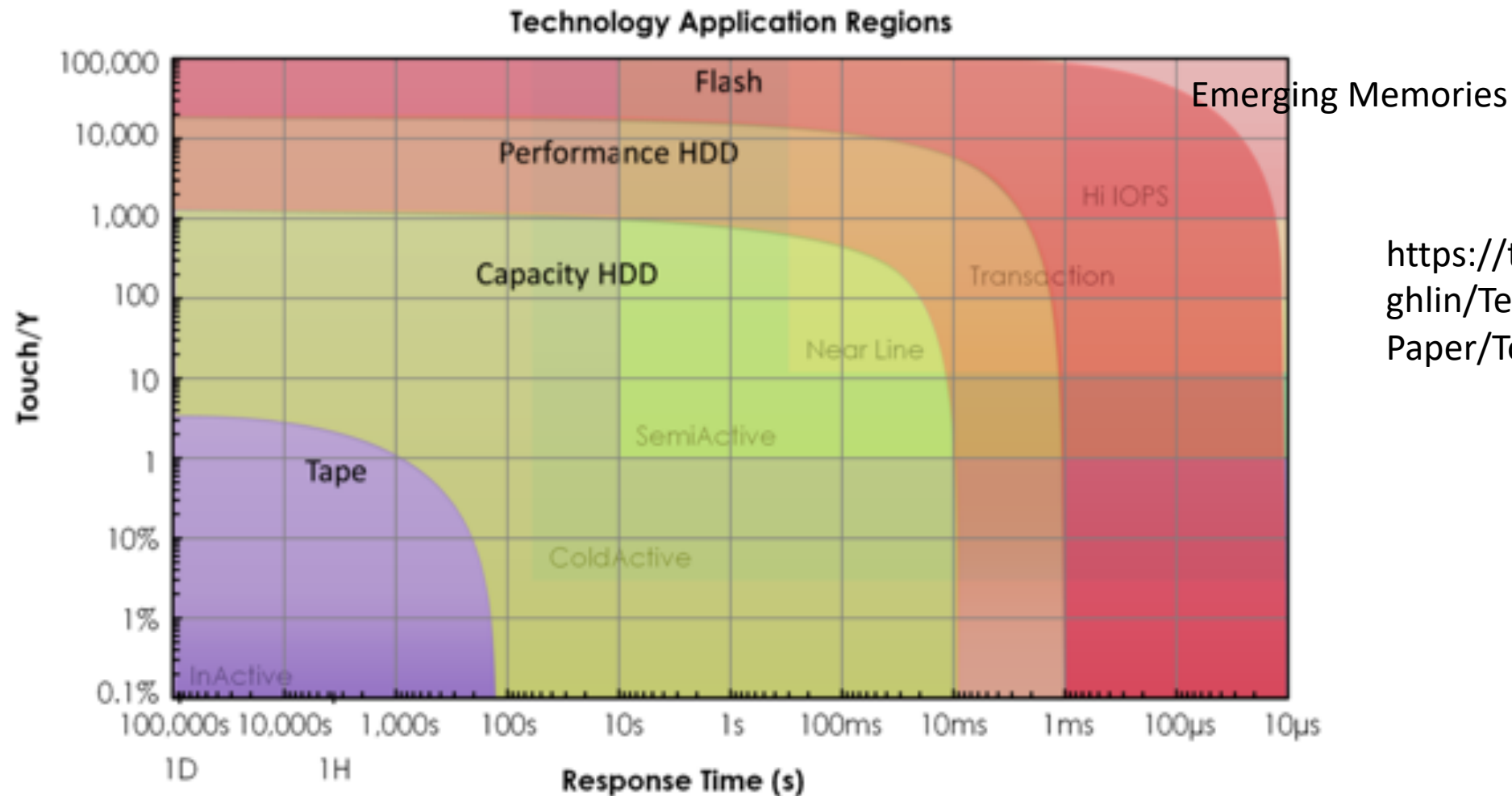


Intel Optane DIMM Announcement, June 2018

2018 Emerging Memories Report, Coughlin Associates, www.tomcoughlin.com

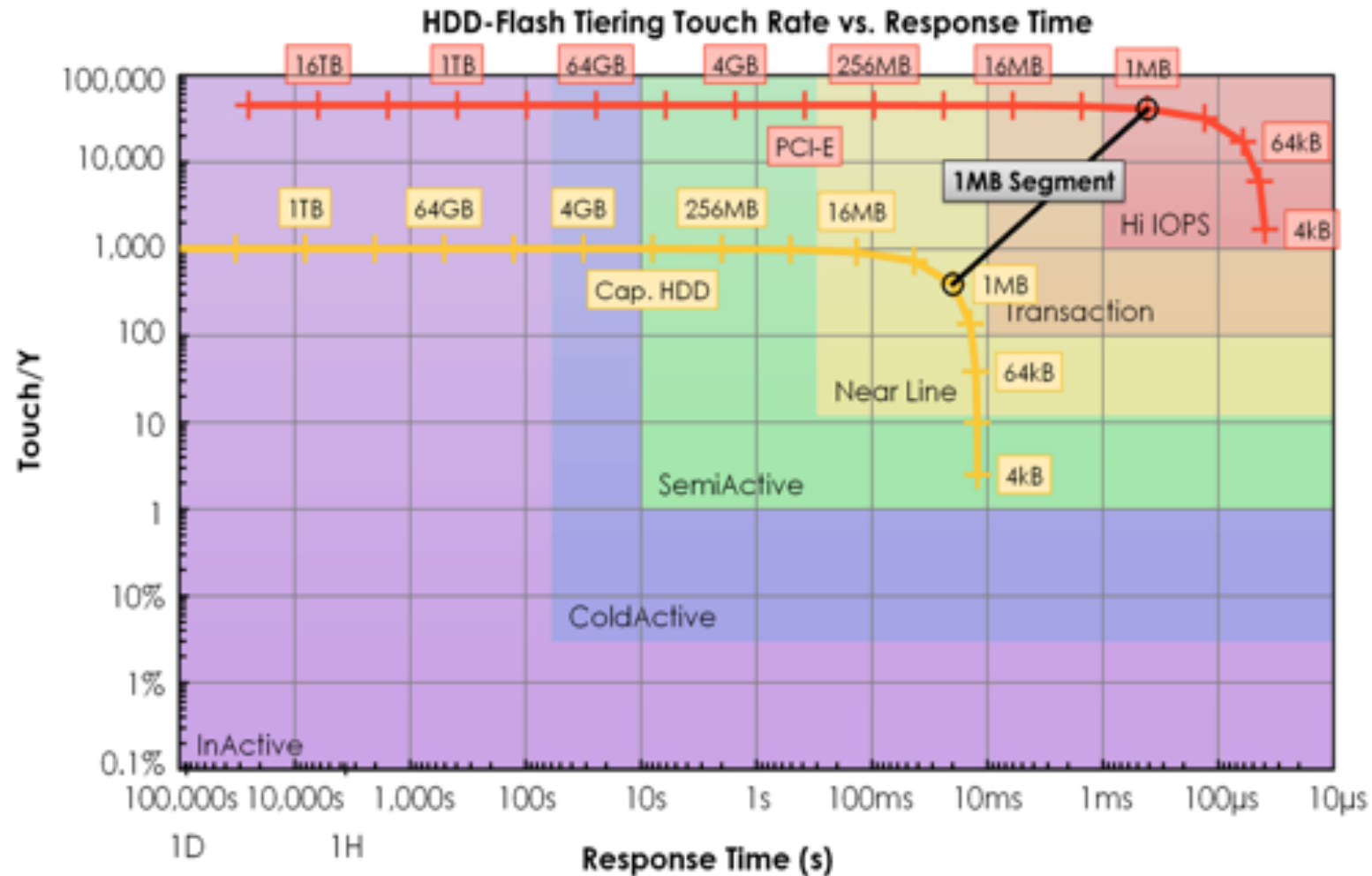# Touch rate versus response time indicating various types of uses



Touch Rate vs. Response Time

https://tomcoughlin.com/Coughlin/Techpapers/Hetzler%20Paper/Touch%20Rate.pdf

# Digital storage technologies regions overlaid on the Touch Rate/Response Time chart



## Technology Application Regions

Emerging Memories

https://tomcoughlin.com/Coughlin/Techpapers/Hetzler%20Paper/Touch%20Rate.pdf

40

# HDD-Flash tiering/caching touch rate chart



https://tomcoughlin.com/Coughlin/Techpapers/Hetzler%20Paper/Touch%20Rate.pdf
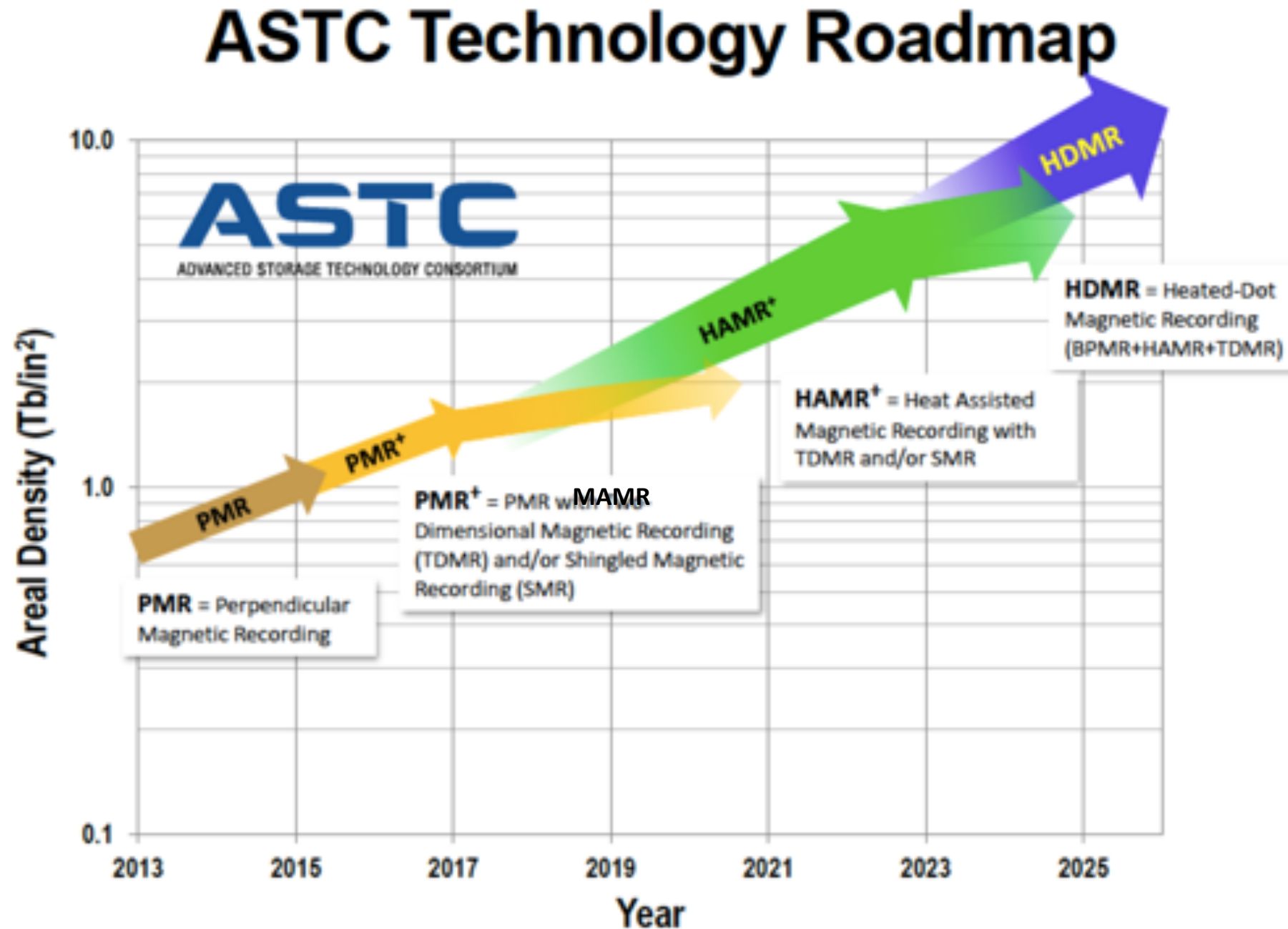
41

# Hard disk drives

- Currently shipping up to 14 TB with 20 + TB expected by 2019, 40 + TB by next decade

- Current HDD areal densities exceed 1 Gb per square inch

- He-filled HDDs provide greater energy efficiency and less cooling requirements

- Fastest growth in HDDs is for bulk cold storage in data centers

- Besides magnetic tape and optical discs, HDDs are the most cost effective storage medium for colder data
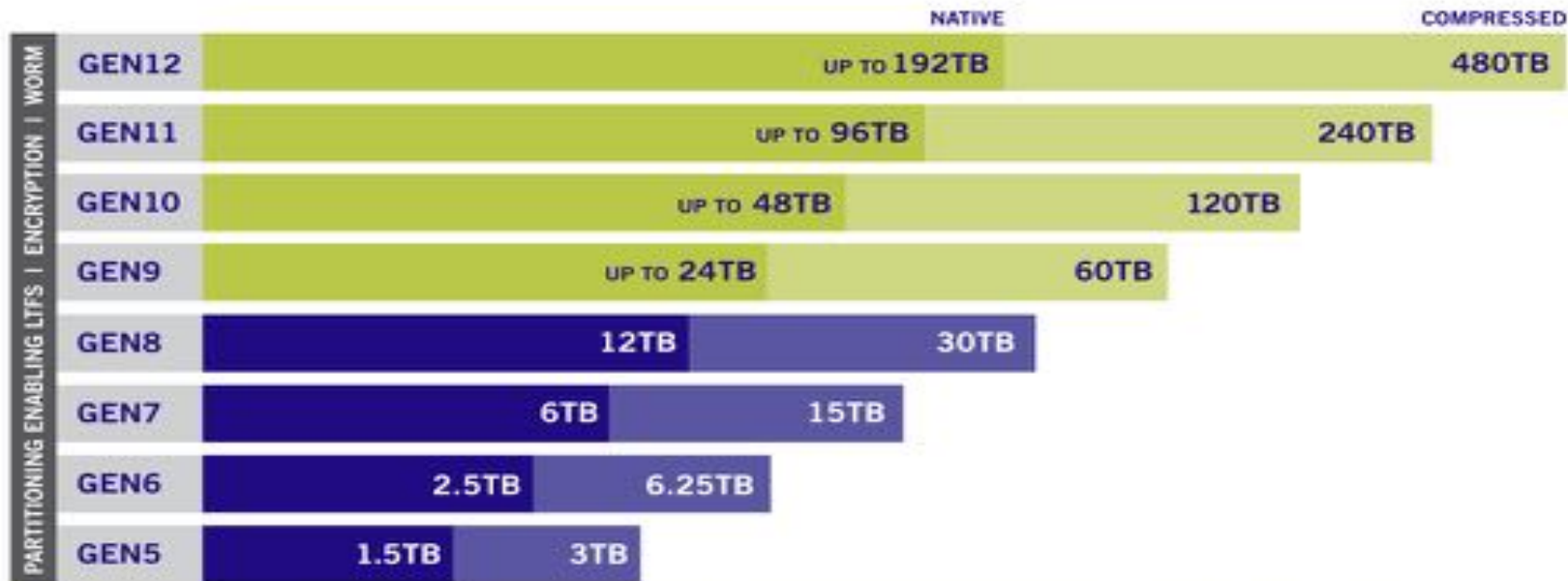
# HDD tech futures

- He and possibly "in vacuum" drives
- Thinner disks using glass
- HAMR/MAMR
- SMR and TDMR



## ASTC Technology Roadmap

**ASTC** — ADVANCED STORAGE TECHNOLOGY CONSORTIUM

**PMR** = Perpendicular Magnetic Recording

**PMR+** = PMR with Two Dimensional Magnetic Recording (TDMR) and/or Shingled Magnetic Recording (SMR)

**MAMR**

**HAMR+** = Heat Assisted Magnetic Recording with TDMR and/or SMR

**HDMR** = Heated-Dot Magnetic Recording (BPMR+HAMR+TDMR)

Areal Density (Tb/in²) vs Year (2013–2025)

# LTO Magnetic Tape Generations



**LTO ULTRIUM ROADMAP**
ADDRESSING YOUR STORAGE NEEDS

| | NATIVE | COMPRESSED |
|---|---|---|
| GEN12 | UP TO 192TB | 480TB |
| GEN11 | UP TO 96TB | 240TB |
| GEN10 | UP TO 48TB | 120TB |
| GEN9 | UP TO 24TB | 60TB |
| GEN8 | 12TB | 30TB |
| GEN7 | 6TB | 15TB |
| GEN6 | 2.5TB | 6.25TB |
| GEN5 | 1.5TB | 3TB |

PARTITIONING ENABLING LTFS | ENCRYPTION | WORM

**NOTE:** Compressed capacity for generation 5 assumes 2:1 compression. Compressed capacities for generations 6-12 assume 2.5:1 compression (achieved with larger compression history buffer).
**SOURCE:** The LTO Program. The LTO Ultrium roadmap is subject to change without notice and represents goals and objectives only. Linear Tape-Open, LTO, the LTO logo, Ultrium, and the Ultrium logo are registered trademarks of Hewlett Packard Enterprise, IBM and Quantum in the US and other countries.

- Current LTO gen 8 provides 15 TB native capacity, greater when compressed

- Laboratory demonstrations of tape technology capable of over 400 TB per cartridge

- Many generations of magnetic recording tech available for tape

# LTO-8 tape cartridge

- Current LTO gen 8 provides 12 TB native capacity, greater when compressed

- LTO consortium includes IBM, HP, Fujifilm and others

- Oracle and IBM also have competing "enterprise" tape technologies

# Sony/Panasonic optical archive roadmap

| | | | |
|---|---|---|---|
| **Capacity** | 300GB | 500GB | 1TB |
| **Signal Processing Technology** | | | High Linear Density (Multi Level Recording Technology) |
| | | High Linear Density (Inter Symbol Interference Cancellation Technology) | |
| | Narrow Track Pitch (Crosstalk Cancellation Technology) | | |
| **Basic Specification** | Double-Sided Disc Technology λ=405nm, NA=0.85, Layer Structure: 3Layers/side | | |

# Flash Memory

- Flash memory is increasing in storage capacity (density) and decreasing in $/GB pricing but still more expensive than HDDs

- Flash memory is winning more applications as its price ($/GB) drops

- Development of NVMe and NVMe-oF has enabled better access to the performance capabilities of flash memory

- In many data center applications, flash memory is now the primary storage

- Flash Memory can also handle more rugged environments, making this a favored storage media for remote location—such as for edge storage

# NAND Flash Expectations

- Flash Memory has moved from primarily planar to planar + 3D flash

- New 3D flash fab has reached parity with planar production in 2018, easing supply constraints

- This has resulted in a drop in flash prices in 2018

- Projections that 3D flash could go out many generations—

- The price reductions for 96 layer and higher will be less than going to 64 layer, because of slower process speeds
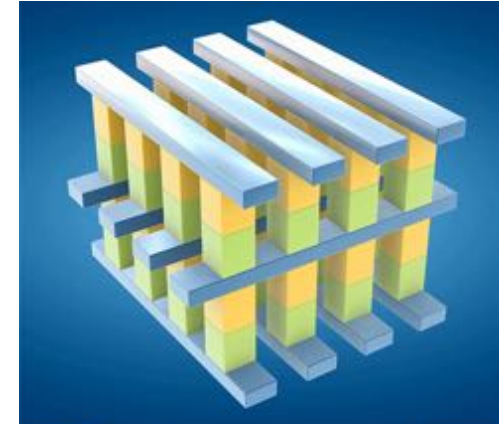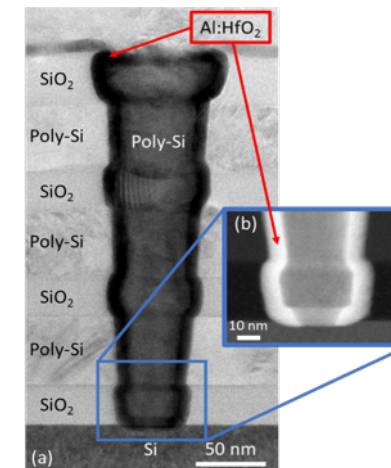


NAND Flash Technology Roadmap Update

# Persistent memory types

**MRAM**

**PCM**

**ReRAM**

**FRAM**

# Emerging Non-Volatile Memories

- There is intense effort to commercialize several non-volatile memories that could replace current volatile memories, such as DRAM and SRAM

- These technologies can be applied to stand along memory chips as well as in embedded memory

- This could reduce energy expenditure in battery and low power devices and also create more efficient data centers

- These NV memories will enable both IoT devices as well as data centers at the edge or in the cloud

- The memory technologies under consideration include magnetic random access memory (MRAM), resistive RAM (RRAM or ReRAM), phase change RAM (PRAM) and ferroelectric RAM (FRAM or FeRAM)

# MRAM and PRAM



- MRAM
  - Everspin shipped over 70 M MRAM Chips.  Company has partnership with Global Foundries, who is building 300 mm wafers and targeting embedded memory applications
  - Samsung and other foundries--plans to ship STT MRAM product samples by 2018.
  - IBM was showing an Everspin MRAM write cache for an SSD at the 2018 MRAM Developers Conference
- PRAM
  - Intel Optane NVMe products shipped in 2017.
  - Micron planning to introduce DIMM-based 3D XPoint product
  - Intel introduced their Optane DIMM products in June 2018
- Emerging NVM market could exceed $6B by 2023 (Emerging Memories are Poised to Explode, Coughlin Associates and Objective Analysis, http://www.tomcoughlin.com/techpapers.htm

# Memory Technology Comparison

| | Established Memory Types | | | | Emerging Memories | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | SRAM | DRAM | NOR Flash | NAND Flash | FRAM | ReRAM | Toggle MRAM | STT | PCM |
| Nonvolatile? | No | No | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Cell Size | 50-120f² | 6-10f² | 10f² | 4-5f² | 16-32f² | 4-6f² | 16-32f² | 5-7f² | 5-8f² |
| Read Time | 1-100ns | 30ns | 10ns | 50ns | 20-50ns | 10-20ns | 3-20ns | 3-15ns | 5-20ns |
| Write Time | 1-100ns | 50ns | $10^4$-$10^7$ns | $10^4$-$10^5$ns | 50ns | 20ns | 10-20ns | 3-15ns | >30ns |
| Endurance | ∞ | ∞ | $10^5$ | $10^3$ | $10^{12}$ | $10^5$ | $10^{15}$ | $10^{15}$ | $10^{12}$ |
| Write Energy | Low | Low | High | Med | Low | Low | Somewhat High | Low | Low |
| Write Voltage??? | None | 2 | 6-8 | 1.8 | 2-3 | 1.2 | 3 | 1.5 | 1.5-3 |

Emerging Memories Poised to Explore:  An Emerging Memory Report, Tom Coughlin and Jim Handy, Coughlin Associates, 2018

# Storage Systems for Edge and Data Center

- PCIe based NVMe storage interfaces will be the basis of future storage systems architectures using flash-based solid state drives

- This includes network storage capabilities including Remote Direct Memory Access (RDMA) and fabrics running the NVMe protocol, NVMe-oF

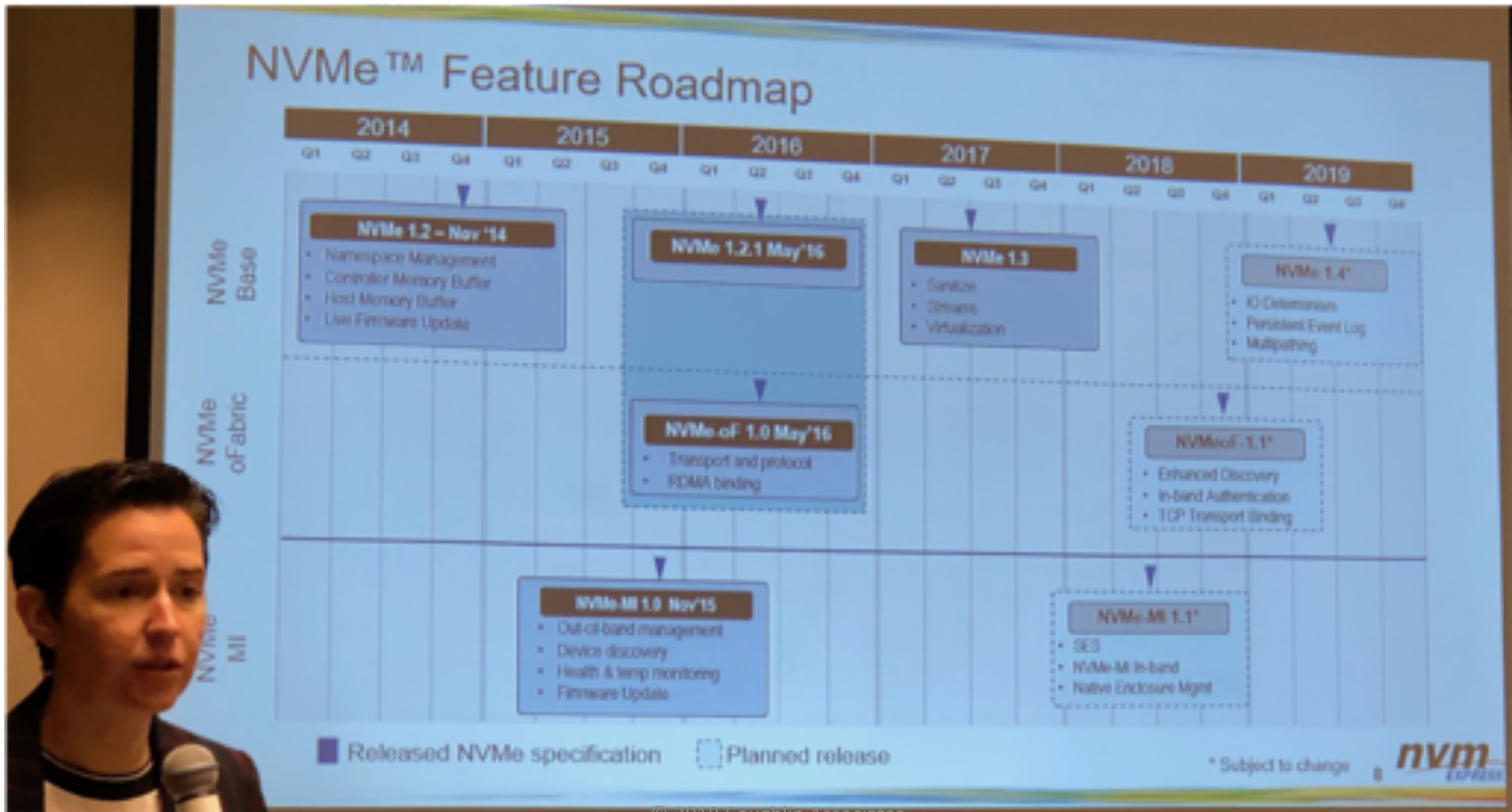- Currently products available running over fibre channel but TCP products announced

1 of 36 Hot-swap NGSFF Drives

# Storage Systems for Edge and Data Center (2)



- Solid state storage using server memory channels (DIMMs) will be the basis of non-volatile memory using flash memory and emerging technologies, such as Intel's Optane

- HDD object storage will be the backbone to mass storage for colder storage (behind solid state primary storage)

- Tape (or optical discs) may be used for longer term colder data

# NVMe Architectures and DIMMS

# Solarflare 2018 FMS presentation



- Computational storage (or memory centric computing)
- Consortiums such as Gen-Z are building networks that enable moving compute closer to the storage

# Solarflare 2018 FMS presentation

- NVMe networking can provide direct memory access

- For high speed storage devices

# PERSISTENT MEMORY OVER FABRICS (PMoF)
## FOR DATA REPLICATION WITH DIRECT LOAD/STORE ACCESS



INITIATOR

TARGET

FABRIC

# New Ways to Manage Storage



- AI methods, such as machine learning, will be used to harvest existing metadata and create new metadata from the essence in the content.

- Machine learning is being used to automate the management of data, including storage tier management, support of virtualization and resource optimization

# Conclusions

- 5G, IoT and AI will be a factor in the growth of digital content over the next decade and will increase requirements for more digital storage and memory solutions
- This could include real time rendering, buffering and caching solutions to improve content delivery QoS, perhaps using a peer to peer network
- The demand for lower latency and power efficiency will drive the use of solid state storage (including emerging memories) as primary storage in data centers and at the network edge
- HDDs and tape will continue to support longer term cold storage
- New storage technologies and architectures as well as traditional storage technologies will more efficient memory-centric computing

# References

- The Memory of Cars, Tom Coughlin, https://ieeetv.ieee.org/ieee-local-events/the-memory-of-cars-talk-by-tom-coughlin

- Michelle Munson, Eluv.io presentation at the 2018 Creative Storage Conference (www.creativestorage.org)

- Touch Rate:  A metric for analyzing storage system performance, Steve Hetzler and Tom Coughlin, 2015 https://tomcoughlin.com/Coughlin/Techpapers/Hetzler%20Paper/Touch%20Rate.pdf

- How Many IOPS Do you Really Need, https://tomcoughlin.com/product/how-many-iops-do-you-really-need-2016/

- Memory Development to Drive Autonomously, Kris Baxter, Micron,  2016 Flash Memory Summit

- 2018 Emerging Memories Poised to Explode:  Emerging Memory Report, Coughlin Associates, https://tomcoughlin.com/product/emerging-non-volatile-and-spin-logic-technology-and-manufacturing-report-2015/

- 2018 Digital Storage in Media and Entertainment Report, Coughlin Associates, https://tomcoughlin.com/tech-papers/

Thanks