

# System Impacts of HDD and Flash Reliability

Steven Hetzler  
IBM Fellow  
Manager Storage Architecture Research

# Outline

---

- **HDD**

- System impact of HDD capacity growth rate

- **Flash**

- Reliability and testing issues for IT applications

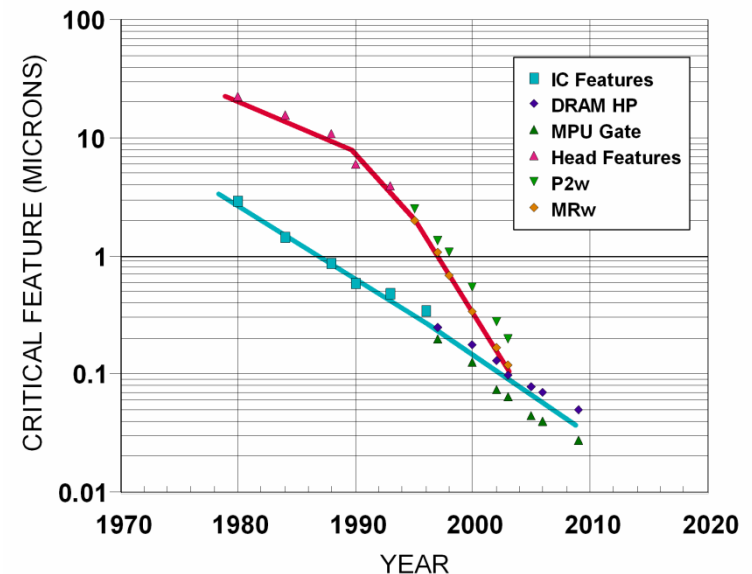
# Magnetic Recording Areal Density

- Capacity growth rates have slowed

- 1994 – 2003      90% CAGR
- 2003 – 2005      20% CAGR
- 2005 – ...      35-40% CAGR (SATA)
- 2003 – 2007      0-10% CAGR (Server)

- Key contributors

- Minimum feature
  - MR stripe width
- Silicon drives lithography tooling
  - Historical: 16%/year
  - NAND now driving Silicon density
- Constant BAR gives 35% CAGR



# HDD Growth Rate Effects

---

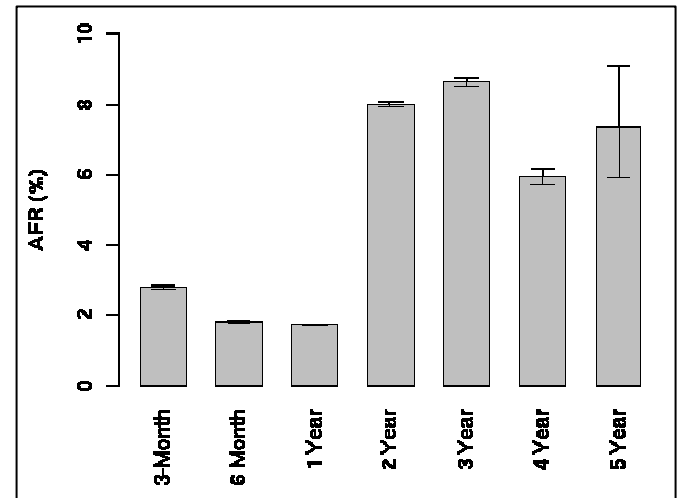
- **Required system field run time will increase**
- **100% CAGR market**
  - Had 10 years of 100% CAGR to entrench behaviors
  - Drove system replacement
  - 3 years = 8x capacity for same cost
    - Only 12% capacity hit to throw old system away
  - Disk longevity not a major issue here
- **35% CAGR changes economics**
  - 7 years = 8x capacity
  - Not economical to replace system @ 3 years
    - 40% capacity hit

## Take Away

**Mean age of disks in the field will grow...**

# Disk Longevity

- **Field data shows higher than spec failure rates**
  - Pinheiro et. al (Fast'07) may show indications of wearout behavior
    - AFRs in the 8% range @ 3+ years
  - Failures less likely to be independent
    - Rebuild activity will add stress
  - Unknown behavior in 5-7year range



Source: Pinheiro et. al, FAST '07

## Take Away

**We expect wearout at some point, increasing AFR**

# Impact of HDD Wearout

- 100TB System (user data)
- 300GB disks
- $8E-3$ /TB Hard error rate (enterprise HDD spec:  $1e15$  bits/err)
- 10 disk arrays (9+1 RAID 5, 8+2 RAID 6)

Disk AFR	1%	2%	4%	6%	8%
RAID 5					
Drive Loss/Y	4	8	15	23	30
Strip Loss/Y	$8E-02$	$2E-01$	$3E-01$	$5E-01$	$6E-01$
Array Loss/Y	$2E-03$	$7E-03$	$3E-02$	$7E-02$	$1E-01$
RAID 6 (DP)					
Drive Loss/Y	4	8	17	25	33
Strip Loss/Y	$4E-05$	$2E-04$	$6E-04$	$1E-03$	$2E-03$
Array Loss/Y	$5E-07$	$4E-06$	$3E-05$	$1E-04$	$2E-04$

- RAID 5 60x weaker to array kill at 8% AFR (8x to strip kill)
- RAID 6 500x weaker to array kill at 8% AFR (60x to strip kill)

# HDD End of Life Impacts

---

- **Increased exposure to data loss events**
  - Systems not designed for wearout AFRs
  - Wearout effects less likely to be independent
  - Disk longevity was less important during 100% CAGR days
  - Disk longevity not likely to be improved going forward
- **System impact**
  - Should account for disk aging in storage system design
- **Field impact**
  - Will be somewhat delayed due to system adoption rate

Take Away

**Slowed capacity growth impacts reliability**

# Flash Reliability in IT

---

- **Flash designed for consumer space**

- Extrapolating to IT applications requires caution
  - Consumer doesn't even approach low-end IT workload
- Digital photography
  - 100 images/card
  - 1000 write cycles (complete fills) = 100K images/card
  - Number of consumers with 100k images/card  $\ll 1\%$ 
    - Likely that camera or card are replaced long before 100k images
  - Failures not critical to consumer, unlikely to be reported
- MP3
  - Write stress even smaller

Take Away

**Consumer flash data of little value to IT**



# Flash IT Issues

---

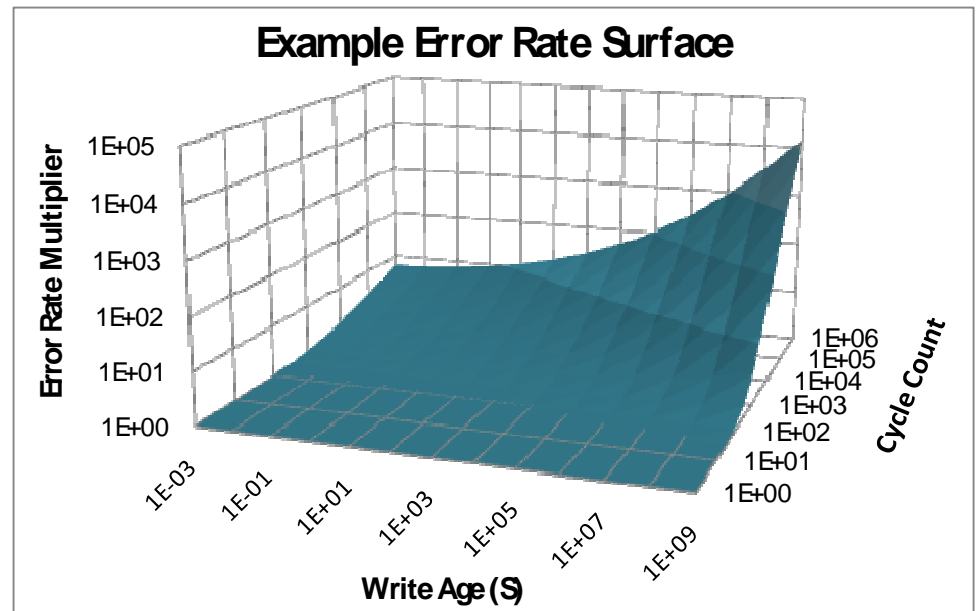
- **Flash has data integrity exposures**
  - Different than in HDD
- **Data Retention Limit *is* Data Corruption**
  - If data is not retained, it means *some* bits changed
    - We don't know which bits, therefore data is corrupted
  - ECC can help, if we know what the error rate is
    - Assuming that the power isn't exceeded
  - Checks (CRC) can help convert to hard errors
    - Hard errors not very desirable either (see strip loss on page 6)

Take Away

**Details critical to system design**

# Flash Error Rate Surface

- Flash has different contributors than HDD
  - Endurance
    - Bit error rate depends on cycle count
  - Retention
    - Bit error rate depends on time since last write
    - Also depends on cycle count!
  - Error rate is a surface
  - Expect BER Increase
    - w/ data age
    - w/ cycle count



# Flash System Design Issues

- **Set data integrity target**

- Start with system parameters

Flash unit IOPS	IOPS	of basic unit of flash
Field lifetime	10	Years
Field units	1,000,000	For full program
Duty cycle	80%	Percent of time at peak IOPS
Mean events/field/life	1	Data integrity target

- Target uncorrected error rate is  $1/(2.5e14 * \text{IOPS})$
- If IOPS = 10K @ 4kB, UER =  $1E-23/\text{bit}$

## Take Away

**Target takes 10M unit-years for mean of 1 event**

# Test Capability

- **Determine testing capability**

- Design test like HDD qualification

Test Duration	1,000	Hours
Test Units	1,000	Flash units
Duty cycle	80%	Need some time to read

- Test capability is  $1.1e-5$  of target IOP count
  - 5 orders of magnitude short
  - No way to get close to 10MUnit-Years
- Doesn't include measurements of error rate surface
  - Need more units and time to allow for data aging

## Take Away

**This is only a 100Unit-Year test (need 10MU-Y)**

# Recommendations

---

- **Design a bit error rate surface test**
  - Likely to be much longer than 1000 hours
    - Inability to accelerate SILC failures
      - Makes the data age test of long duration
  - Testing will be complex
  - BER variables
    - cycle #
    - data age
    - temperature
    - # read cycles
    - power cycles
    - virtual addressing (due to wear leveling)
    - ...
  - Understand supplier test procedures in detail
  - IDEMA could take the lead here

# Summary

---

- **HDD**

- Slowed growth rate will increase mean age of disks in systems
  - This will lead to reliability stress

- **Flash**

- Test, test, test
  - And then monitor in the field
  - Self test difficult due to finite write endurance

- **System vendor owns reliability**

- The reliability situation is changing
  - Customer experience with HDD will change
  - SSDs should be treated very differently than HDDs