# NAND Flash Solid State Storage
# Performance and Capability -- an In-depth Look

lance l. smith, Fusion-io

# SNIA Legal Notice

Tuesday, August 18, 2009

# Abstract

## NAND Flash Solid State Storage Performance and Capability

"This presentation provides an in-depth examination of the fundamental theoretical performance, capabilities, and limitations of NAND Flash-based Solid State Storage (SSS). The tutorial will explore the raw performance capabilities of NAND Flash, and limitations to performance imposed by mitigation of reliability issues, interfaces, protocols, and technology types. Best practices for system integration of SSS will be discussed. Performance achievements will be reviewed for various products and applications. "

Tuesday, August 18, 2009

# Moore's continues to beat Newton's Law

◆ **Mechanical Drives have hit their limits**

- ◆ Platter stability degrades at higher speeds
- ◆ Short-stroking reduces capacity for seek time
- ◆ Capacity is limited by smaller form factors

◆ **Solid State Storage continues to evolve**

- ◆ Greatest bit density (bits per cubic volume)
- ◆ Random IOPS are 250 times greater
- ◆ MLC increases capacity and lowers costs
- ◆ Advanced error correction improves reliability
- ◆ Performance and Capacity are intertwined

NAND Flash Solid State Storage Performance and Capability
© 2009 Storage Networking Industry Association. All Rights Reserved.

4

# Moore's continues to beat Newton's Law

- ## Mechanical Drives have hit their limits
    - Platter stability degrades at higher speeds
    - Short-stroking reduces capacity for seek time
    - Capacity is limited by smaller form factors

- ## Solid State Storage continues to evolve
    - Greatest bit density (bits per cubic volume)
    - Random IOPS are 250 times greater
    - MLC increases capacity and lowers costs
    - Advanced error correction improves reliability
    - Performance and Capacity are intertwined

NAND Flash Solid State Storage Performance and Capability
© 2009 Storage Networking Industry Association. All Rights Reserved.

4

**There can be no data integrity trade-off for performance**

Tuesday, August 18, 2009

# Media Reliability / Availability

◆ The GOOD

- No moving parts
- Post infant mortality (catastrophic) device failures are rare
- Predictable wear out

◆ The BAD

- Relatively high bit error rate, which increases with wear
- Higher density and MLC increases bit error rate
- Program and Read Disturbs

◆ The UGLY

- Partial Page Programming
- Data retention is poor at high temperature and wear
- Infant mortality is high (large number of parts…)

Tuesday, August 18, 2009

# Controller Reliability Management

◆ Wear leveling & Spare Capacity

◆ Read & Program Disturb control

◆ Data & Index Protection

  ◆ ECC Correction

  ◆ Internal RAID

  ◆ Data Integrity Field (DIF)

◆ Management

**Poor Media + Great Controller = Great SSS Solution**

Tuesday, August 18, 2009

Performance

**BAD**

**GOOD**

99%  99.9%  99.99%  99.999%  ...

Data Integrity

Tuesday, August 18, 2009

# Performance is about ROI

## Lower OpEx

- Less HW Maintenance
- Less SW Maintenance
- Greater Uptime
- Less Power/Cooling
- Fewer Diverse Skills

## Lower CapEx

- Fewer CPUs
- Less RAM
- Less Network Gear
- Fewer SW Licenses
- Less Space

HIGHER Productivity

# Media Performance

## The GOOD

- Performance is excellent (wrt HDDs)
- High performance per power (IOPS/Watt)
- Low pin count: shared command / data bus → good balance

## The BAD

- Not really a random access device
  - › Block oriented
  - › R/W access speed imbalance
  - Slow effective write (erase/transfer/program) latency
- Performance changes with wear

## The UGLY

- › Some controllers do read/erase/modify/write
- Others use inefficient garbage collection

Tuesday, August 18, 2009

# Performance Drivers – SSS Design

- Number of NAND Flash Chips (Die)
- Number of Channels (Real / Pipelined)
- Interconnect
- Data Protection (internal/external RAID; DIF; ECC…)
- SLC / MLC Flash Type
- Effective Block Size (LBA; Sector)
- Write Amplification Efficiency
- Garbage Collection (GC) Efficiency
- Bandwidth Throttling
- Buffer Capacity & Mgmt

Tuesday, August 18, 2009

# Simplified Theoretical Analysis

◆ Bandwidth Only (Not IOPS)

- Large Transfers (Data length = Integer times die count)
- Infinite Buffer
- Reads/Writes queued for maximum bandwidth
- No system latency

◆ Read/Write Ratio %'s fixed

- 100/0, 75/25, 50/50, 25/75, 0/100
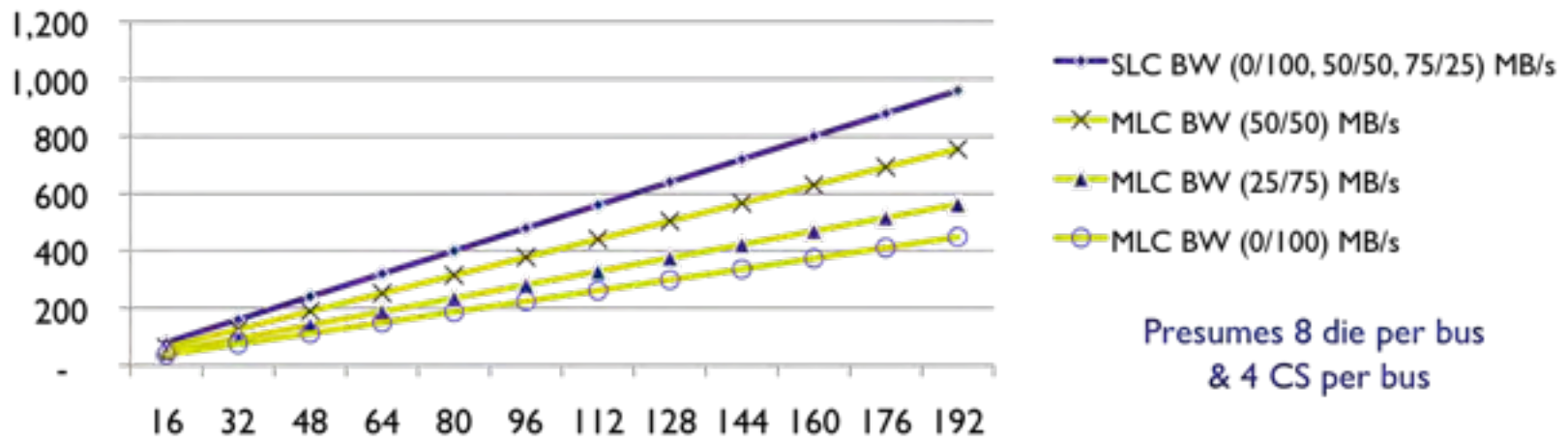- Steady State, 100% Efficient GC (EB erase/EB written = 1)

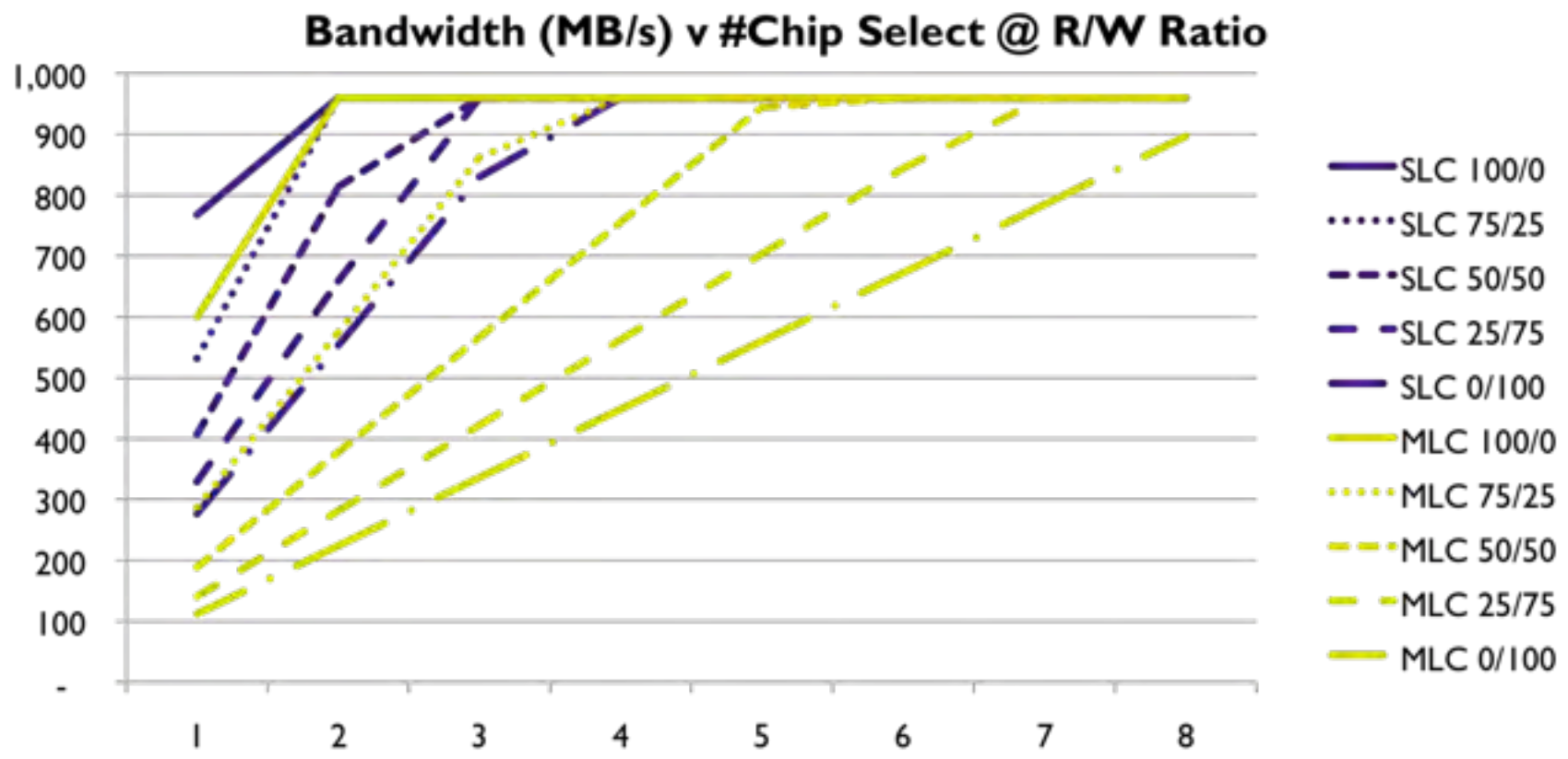◆ Maximum Total BW for SATA-II and PCI-e X4

- No overhead considered

Tuesday, August 18, 2009

| | | SLC | MLC |
|---|---|---|---|
| **Transfer Rate (MB/s)** | tRC & tWC | 400 | 400 |
| **Page Program (us)** | tProgram | 200 | 600 |
| **EB Erase (us)** | tErase | 3000 | 10,000 |
| **Load Page (us)** | tR (tRead) | 25 | 60 |
| **Capacity per die** | | 0.5 | 1.0 |



Theoretical BW (MB/s) v Number of Die (SLC, MLC)

Legend:
- SLC BW (0/100, 50/50, 75/25) MB/s
- MLC BW (50/50) MB/s
- MLC BW (25/75) MB/s
- MLC BW (0/100) MB/s

Presumes 8 die per bus
& 4 CS per bus

NAND Flash Solid State Storage Performance and Capability
© 2009 Storage Networking Industry Association. All Rights Reserved.

13

# Single-Level versus Multi-Level Cell



Bandwidth (MB/s) v #Chip Select @ R/W Ratio

Legend:
- SLC 100/0
- SLC 75/25
- SLC 50/50
- SLC 25/75
- SLC 0/100
- MLC 100/0
- MLC 75/25
- MLC 50/50
- MLC 25/75
- MLC 0/100

**Read / write performance imbalance closed with additional banks
Greater R/W imbalance in MLC requires more banks**

Tuesday, August 18, 2009

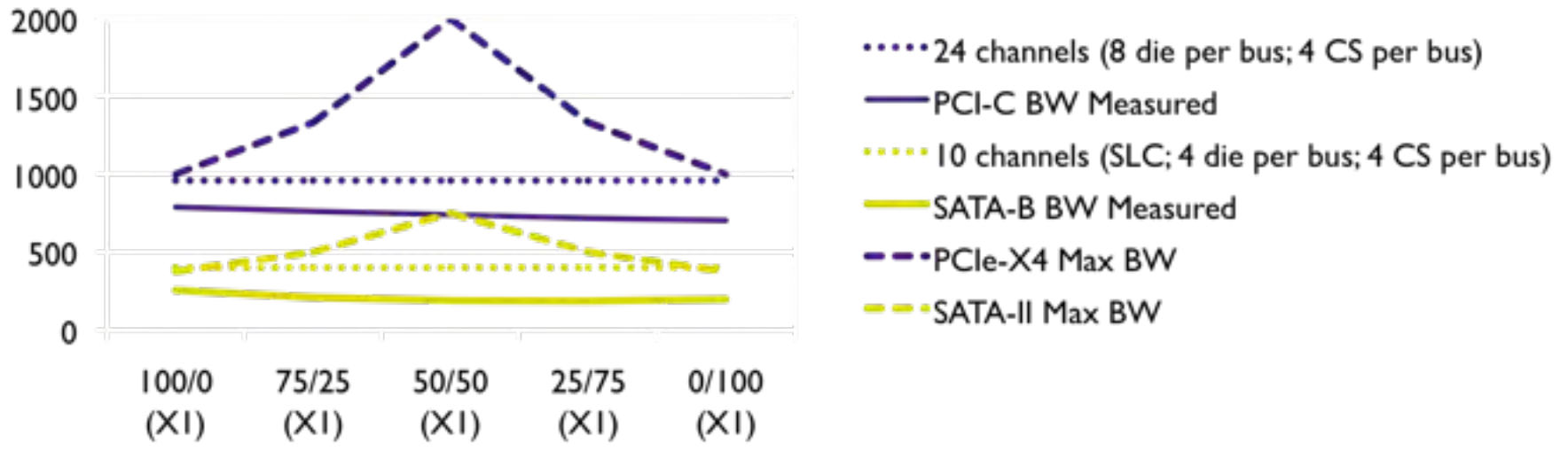# Product Comparison under Test

## Features directly affecting performance measurements

| | SATA (A) | SATA (B) | PCI (C) |
|---|---|---|---|
| Capacity (GB) | 32 | 32 | 160 |
| Bus/Link | SATA-II (3 Gb/s) | SATA-II (3 Gb/s) | PCI-E X4 1.1 |
| Memory Type | SLC | SLC | SLC |
| Adjustable Reserve Capacity | No | No | Yes |
| SSS Internal RAID -- Running during test | No N/A | No N/A | Yes Yes |
| K-IOPS (RMS) | 8 | 27 | 88 |
| K-IOPS (RMS) / WATT | 3 | ? | 7 |
| Bandwidth (RMS, MB/s) | 56 | 208 | 743 |
| ECC correction | 7 bits in 512B | 4 bits in ? | 11 bits in 240B |

Tuesday, August 18, 2009

Education
SNIA

## Measured versus Theoretical Max BW



- ····· 24 channels (8 die per bus; 4 CS per bus)
- —— PCI-C BW Measured
- ····· 10 channels (SLC; 4 die per bus; 4 CS per bus)
- —— SATA-B BW Measured
- – – PCIe-X4 Max BW
- – – SATA-II Max BW

Note: Theoretical Max BW with 24 channels (4 die per bus, 4 CS per bus) is identical to the PCI-C, 24 channel shown in these charts.

Capacity Multiplier:
   SATA-B: 1
   PCI-C: 2

## Measured BW as % of Theoretical Max



- —— PCI-C BW Actual/ Theoretical
- —— SATA-B BW Actual/ Theoretical

NAND Flash Solid State Storage Performance and Capability
© 2009 Storage Networking Industry Association. All Rights Reserved.

16

Tuesday, August 18, 2009

# Access Process (Physics Ignored)

## Read Access

- Address Chip / EB / Page
- Load Page into Register
- Transfer Data From Register 1-byte per cycle

> Typical NAND Flash Die:
> - 2000 Erase Blocks (EB)
> - 64 Pages per EB
> - 4000 Bytes per Page
> - 500 MByte Total Capacity

## Write Access

- Address Chip / EB
- Erase EB

### …some time later…

- Address Chip / EB / Page
- Transfer Data To Register 1-byte per cycle
- Program Register to Page

Tuesday, August 18, 2009

# Example 1: Read/Erase/Modify/Write

| Time = t1 | Time = t2 | Time = t3 |
|---|---|---|

**Starting State**

| Page | Erase Block 1 | | | |
|---|---|---|---|---|
| 0 | b | c | -- | -- |
| 1 | j | -- | k | l |
| 2 | m | -- | -- | -- |
| 3 | -- | -- | q | r |

**Write Buffer & W,X,Y**

| Page | Erase Block 1 | | | |
|---|---|---|---|---|
| 0 | b | c | W | X |
| 1 | j | Y | k | l |
| 2 | m | | | |
| 3 | | | q | r |

**Write Buffer & Z,A,B',C',R'**

| Page | Erase Block 1 | | | |
|---|---|---|---|---|
| 0 | B' | C' | w | x |
| 1 | j | y | k | l |
| 2 | m | Z | A | |
| 3 | | | q | R' |

Buffer holds data while EB-1 Erased

Buffer holds data while EB-1 Erased

| Page | Erase Block 1 | | | |
|---|---|---|---|---|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

| Page | Erase Block 1 | | | |
|---|---|---|---|---|
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 3 | | | | |

Tuesday, August 18, 2009

| Time = t1 | Time = t2 | Time = t3 |
|---|---|---|
| Starting State | Data to Buffer (not shown) | Data to Buffer (not shown) |
| | Erase EB-1 (not shown) | Erase EB-1 (not shown) |
| | Write Buffer & W, X, Y to EB-1 | Write Z, A & Replace b,c,r with B',C',R' & Write EB-1 |

**Erase Block 1**

| Page | | | | |
|---|---|---|---|---|
| 0 | b | c | -- | -- |
| 1 | j | -- | k | l |
| 2 | m | -- | -- | -- |
| 3 | -- | -- | q | r |

**Erase Block 2**

| Page | | | | |
|---|---|---|---|---|
| 0 | b | c | W | X |
| 1 | j | Y | k | l |
| 2 | m | | | |
| 3 | | | q | r |

**Erase Block 3**

| Page | | | | |
|---|---|---|---|---|
| 0 | B' | C' | w | x |
| 1 | j | y | k | l |
| 2 | m | Z | A | |
| 3 | | | q | R' |

Implicit wear leveling; EB-1 → EB-2 → EB-3
Presumes that destination EB-2 & EB-3 erased prior to transfer of data → higher performance (than previous "Read/Erase/Modify/Write" example)

Tuesday, August 18, 2009

# Example 3: Garbage Collection

| Time = t1 | Time = t2 | Time = t3 |
|---|---|---|
| **Start Garbage Collect EB-1** | **EB-1 GC'd to EB-2** <br> **W,X,Y added** | **EB-1 erase** <br> **b,c,r replaced by B',C',R'** |

# GC Performance Impact

- In this example,
  - COPIED DATA: {b, c, j, k, l, m, q, r} 8 blocks
  - NEW DATA {W, X, Y, B', C', Z, A, R'} 8 blocks
  - 50% (8 of 16) writes are user initiated
  - 50% (8 of 16) writes are internal movement (overhead)
- Important:
  - 50% of EB-1 was "invalid data"
  - What if only 10% had been "invalid data?"
  - GC efficiency is dependent upon % of reserve capacity

Tuesday, August 18, 2009

Want to do this in fewer moves?
Add more pegs!

Tuesday, August 18, 2009

SNIA Education

◆ If a high percentage of total storage capacity utilized

**AND**

◆ A High percentage of data has no correlation-in-time

**AND**

◆ Continuous writing (no recovery time for GC)

**THEN…**

*Efficiency of GC greatly diminished*

Tuesday, August 18, 2009

# Pathological Write Condition



**User Capacity Formatted of Total**

- 30GiB of 80G PCI-C
- 40GiB of 80G PCI-C
- 60GiB of 80G PCI-C
- 70GiB of 80G PCI-C
- 74GiB of 80G PCI-C
- 28GiB of 30G SATA-B

Tuesday, August 18, 2009

# Performance vs R/W Ratio

**IOPS @ 512 B**

**Bandwidth (MB/s) @ 128 KB**

Legend: PCI-C, SATA-A, SATA-B

X-axis labels: 100/0 (X1), 75/25 (X1), 50/50 (X1), 25/75 (X1), 0/100 (X1)

**Read/Write Collisions → Drop in Mixed Performance**

25

Tuesday, August 18, 2009

**R/W Ratio and Number of Devices in Parallel**

Tuesday, August 18, 2009

# Performance vs Block Size (75/25)

## 75/25 R/W IOPS

Block Size

SATA-A
SATA-B
PCI-C

## 75/25 R/W Bandwidth (MB/s)

Block Size

SATA-A
SATA-B
PCI-C

Tuesday, August 18, 2009

Tuesday, August 18, 2009

Tuesday, August 18, 2009

# System Level Considerations

- Data / Index Protection (RAID and DIF)
- Scalability
- Compare system- or data-center-level
  - Not device
- Best case: test on real application
  - Not benchmark
  - Plan to do tuning to reach top perf. / objectives
  - Applications may have contra-indicated optimizations
    - Keeping data in close physical proximity (short stroking)
    - Caching algorithms

Tuesday, August 18, 2009

Education
SNIA

☒ Bandwidth / IOPS at

   ☒ **Block size(s) you need**

   ◆ **R/W ratio you use**

   ◆ **Steady State / Burst**

   ☒ Reserve capacity used

   ◆ Data's temporal relationship

   ◆ Scalability

   ☒ RAIDing

   ◆ BOL / EOL

☒ Design impacts on data integrity; life; failures & perf.

   ☒ ECC robustness

   ◆ Write amplification / GC efficiency

   ◆ Internal RAID

   ☒ Bandwidth throttling

   ◆ Partial Page Programming

◆ Test Conditions

   ◆ Workload

   ◆ Temporal Relationships

   ☒ User capacity / reserve capacity

Tuesday, August 18, 2009

☒ Please send any questions or comments on this presentation to SNIA: : *__tracksolidstate@snia.org__*

**Many thanks to the following individuals
for their contributions to this tutorial.**

**- SNIA Education Committee**

**Jonathan Thatcher
Khaled Amer
Phil Mills
Rob Peglar
Marius Tudor**

Tuesday, August 18, 2009