# PCI Express Storage in Client Systems

John Carroll

Storage Architect

Intel

Flash Memory Summit 2013
Santa Clara, CA

## Why PCIe for Client Storage ?

- SSDs can be built that exceed SATA Gen3 (600 MB/s) today

- Enabling SATA beyond 600 MB/s is a long term development effort
  - Single lane scaling beyond ~ 8Gbps is challenging & requires trade-offs
  - Multi-lane SATA requires a new connector and modified chipset SATA controllers to make multi-lane software transparent

- To enable higher speed client SSDs in near term ('13 / '14), PCIe is the only choice
  - PCIe has bandwidth lead (1 GB/s with Gen3)
  - PCIe has multi-lane for scalability (x2, x4, ...)
  - Software compatible PCIe SSDs can be built as a single port AHCI device

**Interface Bandwidth** (Specification Intro Date)

**PCIe can deliver the performance client SSDs need today.**

Source: FMS '12 Amber Huffman

(intel)

## VAIO Pro 13 Ultrabook™

The world's lightest 13.3" touch Ultrabook[21].

Features:

- 4th gen Intel® Core™ i7 processor available
- Windows 8 Pro available
- Full HD TRILUMINOS IPS touchscreen (1920 x 1080)
- Super fast 512GB PCIe SSD available
- Ultra-light at just 2.34 lbs.

## Faster all-flash storage.
## Ready. Set. Done.

Flash storage in MacBook Air is now up to 45 percent faster than the previous generation. So everything you do is snappier and more responsive. MacBook Air even wakes up faster than ever, thanks to flash storage and the latest Intel Core processors. And now the 11-inch model comes standard with twice the capacity — 128GB — yet still starts at $999.

(intel)

# Agenda

- Client PCIe Form Factors

- SATA to PCIe Transition

  - Power Optimizations for PCIe SSDs

- Controller Interfaces

# Form Factor & Connector Landscape

Ultrabook™   Value Notebook   All-in-one   Desktop   WS   Server

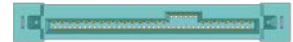M.2                                        2.5" SATA Express              SFF-8639

- ## M.2* was designed for the unique needs of Ultrabook™
  - However, M.2 is being used in a wide variety of devices
  - Cannot support HDDs or SSHDs

- ## 2.5" SATA Express* and SFF-8639 connectors provide flexibility to support HDDs, SSHDs, & SSDs on the same connector

# EMI Challenges for PCIe* Cabling

- A reference clock in the cable for 2.5" PCIe drives causes EMI issues for client PCs

- PCI-SIG solved this challenge with a new clocking mechanism (SRIS)

**Subpart A—General**
**§ 15.1   Scope of this part.**
(a)  This part sets out the regulations under which an intentional, unintentional, or incidental radiator may be operated without an individual license. It also contains the technical specifications, administrative requirements and other conditions relating to the marketing of part 15 devices.

(b)  The operation of an intentional or unintentional radiator that is not in accordance with the regulations in this part must be licensed pursuant to the provisions of section 301 of the Communications Act of 1934, as amended, unless otherwise exempted from the licensing requirements elsewhere in this chapter.

(c)  Unless specifically exempted, the operation or marketing of an intentional or unintentional radiator that is not in compliance with the administrative and technical provisions in this part, including prior Commission authorization or verification, as appropriate, is prohibited under section 302 of the Communications Act of 1934, as amended, and subpart I of part 2 of this chapter. The equipment authorization and verification procedures are detailed in subpart J of part 2 of this chapter.

**SRIS clocking enables cabling for PCIe* SSDs in client systems**

(intel)

# Agenda

- Client PCIe Form Factors
- **SATA to PCIe Transition**
  - **Power Optimizations for PCIe SSDs**
- Controller Interfaces

(intel)

# Transitioning from SATA* to PCIe*

- PCIe* storage devices incorporate controller functionality into SSD/SSHD

**SATA**

**PCIe**

Host

| SATA | PCIe |
|---|---|
| << upper-level OS & applications >> | << upper-level OS & applications >> |
| AHCI Driver | AHCI/NVMe Driver |
| OS driver PCI/PCIe | OS driver PCI/PCIe |
| AHCI Controller | |
| SATA interface | PCIe interface |

SSD/SSHD

| SATA | PCIe |
|---|---|
| SATA interface | PCIe interface |
| Device Controller | AHCI/NVMe Controller |
| NVM Storage | Device Controller |
| | NVM Storage |

(intel)

# Independent Power States

**SATA\***     **PCIe\***

- PCIe* separates link and device state into two independent states

- Link states defined by PCIe* spec

- Device states defined by controller interface

**Host**

| SATA* | PCIe* |
|---|---|
| << upper-level OS & applications >> | << upper-level OS & applications >> |
| AHCI* Driver | AHCI/NVMe* Driver |
| OS driver PCI/PCIe* | OS driver PCI/PCIe* |
| AHCI* Controller | PCIe* Link State |

**DEVSLP signal**

**SSD/SSHD**

SATA* interface ⇕ SATA* interface

PCIe* interface ⇕ PCIe* interface

**CLKREQ#**

Device Controller

NVM Storage

**SATA\* Power State**

AHCI/NVMe Controller

Device Controller

NVM Storage

**PCIe\* Device State**

(intel)

# PCIe* DevSleep Equivalent

The lowest non-zero power state for a PCIe* SSD

Host

<< upper-level OS & applications >>

AHCI/NVMe* Driver

OS driver PCI/PCIe*

PCIe* interface

Link in low power L1.2 state

PCIe* interface

SSD/SSHD

AHCI/NVMe* Controller

Device Controller

NVM Storage

Controller registers active for faster recovery

Rest of drive powered down

(intel)

# Resuming from DevSleep Equivalent

- For SATA*, AHCI controller in the host allows slower SSD recovery from DevSleep

- With PCIe*, register reads can stall the CPU consuming watts while waiting for controller to respond

- 2 step resume process for responsiveness and save power

```
PCIe* link in
L1.2 state
```

~150 µs

(1)  CLKREQ#

```
Links transition
to active state;
Register R/W
possible
```

30-50 ms

(2)  Drive loads context

```
Drive ready to
service I/O
```

**Link transitions first; drive catches up later**

(intel)

# Agenda

- Client PCIe Form Factors
- SATA to PCIe Transition
  - Power Optimizations for PCIe SSDs
- **Controller Interfaces**

# AHCI Lacks Scalability for the Future

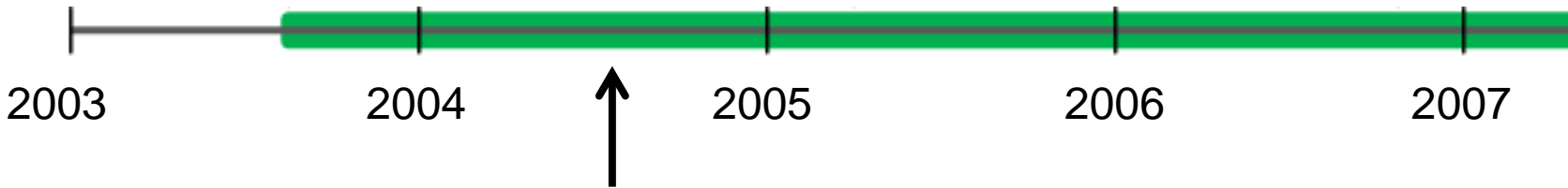A huge leap ahead of IDE, but still designed for hard drives…

| | **AHCI** | **nvm EXPRESS** |
|---|---|---|
| **Uncacheable Register Reads**<br>Each consumes 2000 CPU cycles | 4 per command<br>8000 cycles, ~ 2.5 µs | **0** per command |
| **MSI-X* and Interrupt Steering**<br>Ensures one core not IOPs bottleneck | No | Yes |
| **Parallelism & Multiple Threads**<br>Ensures one core not IOPs bottleneck | Requires synchronization lock to issue command | No locking, doorbell register per Queue |
| **Maximum Queue Depth**<br>Ensures one core not IOPs bottleneck | 32 | 64K Queues<br>64K Commands per Q |
| **Efficiency for 4KB Commands**<br>4KB critical in Client and Enterprise | Command parameters require two serialized host DRAM fetches | Command parameters in one 64B fetch |

**NVM Express* is architected from the ground up for NAND today and next gen NVM of tomorrow.**

# Another Controller Interface Transition

SATA* ships;
IDE* only
(Intel ICH5)

Microsoft* ships
AHCI* driver with
Windows* Vista*

2003　　　　　2004　　　　　2005　　　　　2006　　　　　2007

SATA* includes
AHCI* as an option;
(Intel ICH6)

# Summary

- Client PCIe* storage shipping now – primarily M.2

- Transitioning to PCIe* brings new tools to reduce storage power

- PCIe* AHCI* devices leverage existing AHCI* SW

- NVMe* the long term PCIe* solution; ecosystem establishing quickly