

Storage Interfaces



© 2021

Introduction

Tom Gardner



IEEE Silicon Valley Technology History Committee www.SiliconValleyHistory.com

IEEE Silicon Valley Technology History Committee

- Part of IEEE Santa Clara Valley Section
- Webinars on history of Silicon Valley history

www.SiliconValleyHistory.com

Storage Interfaces Agenda

Speaker	Subject
Tom Gardner	Introduction & History
Grant Saviers	Evolution of Storage Interfaces
Jai Menon	Evolution of Block Storage System Interfaces
Amber Huffman	Fast & Efficient Interfaces for the SSD Era
Questions & Answers	

Presentations are in two files Gardner and Saviers: Storage Interfaces 20210511 Part 1.pdf Menon and Huffman: Storage Interfaces 20210511 Part 2.pdf



Jai Menon

Evolution of Block Storage System Interfaces



Storage Interfaces

Silicon Valley Technology History Committee



Jai Menon

5/11/21



Talk Outline

- Mainframe block storage system interfaces
 - CKD, ECKD
 - CKD-Emulation on FB disks
- Block Storage systems interfaces for industry-standard servers
 - FC, iSCSI, NVMeOF

No coverage of file or object interfaces

Mainframes – CKD CKD refers to both disk format and commands (CCWs)



Introduced in 1964 with S/360

Synchronous operation

Gaps between fields must be large enough to allow channel to read, match, and turn around to issue next command

As disk track density increased, more and more space wasted in gaps

Each CKD disk cost IBM approx. \$100M to develop

Last CDK disk 3390-9 in 1993 11", 4.2 MB/s, 4200 rpm, 22.7 GB

Mainframes – CKD-Emulation

Possible because processors and memory became cheaper



Mainframes – ECKD processors and memory became cheaper



CCHHRKLDL

Introduced in 1985 with IBM 3880 Control Unit

Channel and CU/disk operate asynchronously

Longer distance possible from mainframe to disks

Better performance on hits (with caching)

Some access methods still used CKD

...

Data2

Mainframes – Escon Introduced in 1990 as part of S/390



Operated originally at 10 MB/sec Allowed greater distance Allowed switching Single CU to 8 mainframes

Mainframes – FICON, zHPF FICON Introduced in 1998 as part of 5th gen S/390



Industry-standard FB disks

FICON

- 1, 2, 4, 8, 16 Gbps
- FC switches and Directors
- Up to 100 kms

zHPF (high Performance FICON)

- Introduced in 2009
- Use FCP approach to transfer cmds, data & status
- Optional feature
- TCWs (Transfer Control Words) instead of CCWs
- 90% typical perf improvement
- Fewer channels needed

Mainframes today – DS8900 controller



Industry-standard servers – SCSI over FC



Industry-standard servers – iSCSI = SCSI on TCP/IP on Ethernet



- 1st draft IBM/Cisco 2000; ratified in 2003
- Lower cost than FC
- Longer distance than FC
- Runs on existing network infrastructure
- IBM Total Storage IP Storage 200i, which shipped in 2001, is the earliest iSCSI product in the industry.

Industry-standard servers – NVMe over Fabrics (NVMeoF) Focus on All Flash Storage Systems (Extends NVMe over PCIe to other transports)



Performance Comparison: Random Read



Open. Together.

iSCSI and NVMe/TCP AFA targets built from Intel Gold Servers

(Manoj Wadekar, FB and Anjaneya Chagam, Intel - 2019)



Minimal performance overhead with NVMe over TCP/IP



The Next Speaker

- Amber Huffman
- Fast and efficient interfaces for the SSD Era



NVMe® Technology Fast & Efficient Interfaces for the SSD Era

Amber Huffman

Fellow & Chief Technologist of IP Engineering Group, Intel Corporation President, NVM Express, Inc.





Intel Corporation



Memory and Storage Hierarchy Gaps



nvm

intel.

3

The Evolution of NVMe® Technology



ENABLE INNOVATION

SCALE OVER FABRICS

UNIFY PCIE* SSDs

2010

2020





Intel Corporation

4

Framing the Need Why Standard Interface for PCIe* SSDs

Enabling Broad Adoption of PCIe SSDs: Enterprise NVMHCI

- PCIe SSDs are attractive and popping out of the woodwork
 - Eliminates SAS infrastructure, plenty of PCIe lanes, bandwidth can be concentrated, lower latency
- PCIe SSDs **lack** standard OS support & driver infrastructure, since there is no standard host controller interface
- Impact of no standard infrastructure:
 - Requires SSDs vendor to provide drivers for every OS
 - OEM features, like error logs, are implemented in an inconsistent fashion
 - OEMs have to validate multiple drivers
- Enable broad adoption by extending NVMHCI (Non-Volatile Memory Host Controller Interface) to Enterprise
 - Leverage the client NVMHCI i/f, software infrastructure, and Workgroup to fill gap quickly with streamlined solution
 - http://www.intel.com/standards/nvmhci



450GB Flash storage
120,000 IOPS
700 MB/s random sustained external throughpu
ECC and RAID
Embedded CPU controller.

To order, please contact <u>Texas Memory Systems</u> Sales.

SMART Modular Technologies Announces Enterprise Class PCI Express Storage Solution



Delivering 140K random performance and 200x power reduction, the new 400GE XceedIOPS PCIe raises the bar in next-generation solid-state storage.

NEWARK, CA, June 1, 2009 - SMART Modular Technologies (WWH), Inc.





Building a Coalition Pull Together Key Stakeholders to Drive Change





6

Deliver the Baseline Specification And Select a Much Better Name!





Paint the Future Why It's Worth Crossing the Adoption Chasm

NVMe*: Efficient SSD Performance

	AHCI ¹	NV EXPRESS	
Uncacheable Register Reads Each consumes 2000 CPU cycles	4 per command 8000 cycles, ~ 2.5 μs	0 per command	
MSI-X and Interrupt Steering Ensures one core not IOPs bottleneck	No	Yes	
Parallelism & Multiple Threads Ensures one core not IOPs bottleneck	Requires synchronization lock to issue command	No locking, doorbell register per Queue	
Maximum Queue Depth Ensures one core not IOPs bottleneck	32	64K Queues 64K Commands per Q	
Efficiency for 4KB Commands 4KB critical in Client and Enterprise	Command parameters require two serialized host DRAM fetches	Command parameters in one 64B fetch	

NVMe designed for high parallelism and low latency

¹AHCI: Serial ATA programming interface. See http://www.intel.com/technology/serialata/ahci.htm

Scalability for Future NVM

NVMe* is defined to scale for future NVM

- Block layer: 2.8 µs, 9100 cycles - Traditional: 6.0 µs, 19500 cycles

Measurement taken on Intel® Core® I5-2500K 3.3GHz 6MB L3 Cache Quad-Core Desktop

- Host controller standards live for 10+ years
- Future NVM may have sub microsecond latencies

Linux * Storage Stack



Prototype

Measured LOPS

1 02

and SCSI layer while optimizing for NVMe Block layer attach reduces overhead > 50% Cores Used for 1M IOPs ∠.∪ µ∋∂CS USECS INTEL I



10

15

٠

Chatham NVMe Prototype

Processor using Linux RedHat EL6.0 2.6.32-71 Kernel

8

Provide Physical Infrastructure Optimized Form Factors for SSDs

SSD Form Factor Working Group Driving Enterprise SSD Infrastructure



An Optimized Caseless Form Factor

- Attributes to strive for in a new standardized caseless SSD FF:
 - Scalable from small to large capacity points
 - Support SATA Gen3 and two lanes of PCI Express* Gen3
 - Optimize for Z height (e.g. board edge connector, reduce PCB thickness)
 - Mounting strategy will limit board area and reduce fasteners
 - Optimize board size based on BGA NAND package & ensure efficient tiling







Enabling the Ecosystem **Drivers & BIOS Enabling Required**

Reference Drivers for Key OSs

- Linux*
 - Already accepted into the mainline kernel on kernel.org
 - Open source with GPL license
 - Refer to http://git.infradead.org/users/willy/linux-nvme.git
- Windows*
 - Baseline developed in collaboration by IDT*, Intel, and LSI* _
 - Open source with BSD license
 - Maintenance is collaboration by NVMe WG and Open Fabrics Alliance
 - Refer to https://www.openfabrics.org/resources/developer-tools/nvmewindows-development.html
- VMware*
 - Initial driver developed by Intel
 - Based on VMware advice, "vmk linux" driver based on Linux version
 - NVMe WG will collaborate with VMware on delivery/maintenance

Reference Drivers for Key OSs (cont.)

- Solaris*
 - There is a working driver prototype
 - Planned features include:
 - Fully implement and conform to 1.0c spec
 - Efficient block interfaces bypassing complex SCSI code path
 - NUMA optimized queue/interrupt allocation
 - Reliable with error detect and recovery fitting into Solaris^{*} FMA
 - Build ZFS with multiple sector sizes (512B, 1KB, 2KB, 4KB) on namespaces
 - Fit into all Solaris disk utilities and fwflash(1M) for firmware
 - Boot & install on SPARC and X86
 - Surprise removal support
 - Plan to validate against Oracle^{*} SSD partners
 - Plan to integration into S12 and a future S11 Update Release
- UEFI
 - The driver is under development
 - Plan to open source the driver in Q1 '13, including bug/patch process
 - Beta quality in Q1'13, production quality Q2'13

NVMe = NVM Express 23



24



Industry Thought Leaders Show What's Possible...

Designing with the Right Pieces



multi-core parallelism

- + PCI Express* Gen3 bandwidth
- + NUMA-aware software
- + NVMe flash
- = millions of IOPs...
- ...at microsecond latencies

new math for storage platforms



Proof Points in EMC* Products

Project Thunder: A networked non-volatile memory appliance

- Building block design center:
 - 10-20TB of flash capacity
 - Consistent 2.5M IOP throughput @ 150us
- NVMe ready:
 - Today: improved latency, reduced processor overhead
 - Tomorrow: ready for nano-second class NVM
 - On the showcase floor until 2pm



10 NVMe = NVM Express

NVMe = NVM Express



41

intel. ¹¹

Industry Thought Leaders **Inbox Driver Support by Microsoft**

Microsoft's Support of NVM Express

- The Natural Progression from SATA for NVM
 - Standardized PCI Express* Storage
 - First devices are enterprise-class
 - High-Density / High-Performance
 - Closing the latency gap with RAM
- Windows* Inbox Driver (StorNVMe.sys)
 - Windows Server 2012 R2 (enterprise)
 - Windows 8.1 (client)
 - Stable Base Driver
- The Storport Model
 - Reduced development cost
 - Offloads Basics: PnP, Power, Setup, Crash, Boot*
 - Mature / Hardened Model
 - Storport optimized for performance
 - **RAM-backed NVMe device**
 - > 1 million IOPS with < 20µs latencies</p>



StorNVMe Delivers a Great Solution

- StorNVMe Implementation Highlights
 - Uses hardened Enterprise Storage Stack
 - Strives for 1:1 manning of queues to processors



0.200

0.100

0,000

TI.CIDA

- With great IOPs 0.300
- And low latency

Random Read Workload, Tool: IOmeter





38

Crossing the Chasm Requires Patience and Perseverance

NVM Express* Deployment is Starting

- First plugfest held May 2013 with 11 companies participating
 - Three devices on Integrator's List
 - Next plugfest planned for Q4
- Samsung announced first NVM Express* (NVMe) product in July



Learn More in the NVMe Community

Check out the NVMe Community in the Showcase to see NVMe products in action

	Dell
NVM Express Community	Intel
	EMC
nvm	Micron
	SanDisk
	LSI
Fast and Efficient	SNIA
	PMC-Sierra
	Agilent
The function of the function o	Western Dig
	Teledyne Le
IDF13	Viking Techr
	Tektronix

Company	Booth #
Dell	726
Intel	727 & 734
EMC	728
Micron	729
SanDisk	730
LSI	731
SNIA	732
PMC-Sierra	733
Agilent	735
Western Digital	736
Teledyne LeCroy	737
Viking Technology	738
Tektronix	739

IDF13



55

Evangelize the Value Proposition Fast, Efficient, Simple

NVMe has ~ 1.5X better IOPs than SAS

NVM Express^{*} (NVMe) Delivers Best in Class IOPs

- 100% random reads: ٠
- NVMe has >3X better IOPs than SAS 12Gbps
- 70% random reads: NVMe has >2X better IOPs than SAS 12Gbps ٠
- 100% random writes: ٠



Note: PCI Express¹ (PCIe¹/NVM Express¹ (PCIe¹/NVMe) Measurements made on Intel®Core¹⁴/7-3770S system @ 3.1GHz and 4GB Mem running Windows¹ Server 2012 Standard O/S. Intel PCIe/NVMe SSDs. d IOmeter' tool. PCIe/IV/We SSD is under development. SAS Measurements from HGST Ultrastar' SSD800M 1000M (SAS) Solid State Drive Specification. SATA Measurements from Intel Solid State Drive Product Specification. Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark' and MobileMar usins specific comuter systems. components software, operations and functions. Any change to any of those factors may cause the results to avay. You should compute other and the software and workloads the results to avay. You should compute the system of the software and workloads the results to avay. in fully evaluating your contemplated purchases, including the performance of that product when combined with other products

The Efficiency of NVM Express^{*} (NVMe)

- CPU cycles in a Data Center are precious
- And, each CPU cycle required for an IO adds latency •



NVM Express^{*} (NVMe) takes less than half the CPU cycles per IO as SAS

With equivalent CPU cycles, NMMe delivers over 2X the IOPs of SAS!

Software and workloads used in performance tests may have been ontimized for performance only on Intel microprocessors. Performance tests such as SVSmark' and MohileMark', are measured using specific computer systems components software operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluation you contemplated purchases, including the performance of that product when combined with other products. For detailed configuration information, refer to "Setup for Efficiency and Latency Analysis" foil in Backup.



IDF14

Memory and Storage Hierarchy Gaps

nvm



Memory and Storage Hierarchy Gaps Solutions



The Evolution of NVMe® Technology



ENABLE INNOVATION

SCALE OVER FABRICS

UNIFY PCIE* SSDs

2010

2020





Intel Corporation

Enabling More Use Cases NVMe* over Fabrics

Why NVM Express* (NVMe) over Fabrics?

- Simplicity, Efficiency and End-to-End NVM Express* (NVMe) Model
 - NVMe has a single Admin queue pair with 10 required commands
 - NVMe supports up to 64 K I/O Queues with 3 required commands
 - Simplicity of protocol enables hardware automated I/O Queues transport bridge
 - No translation to or from another protocol like SCSI (in firmware/software)
 - Inherent parallelism of multiple I/O Queues is exposed
 - NVMe commands and structures are transferred end-to-end



<u>Goal:</u> Make remote NVMe equivalent to local NVMe, within ~ 10 μs latency.

Architectural Approach

- The NVM Express* (NVMe) Workgroup is starting the definition of NVMe over Fabrics
- The first fabric definition is the RDMA protocol – used with Ethernet and InfiniBand[™]
- A flexible transport abstraction layer is under definition, useful for many different fabrics





Enabling More Use Cases NVMe* over Fabrics

Commonality Between PCI Express®and Fabrics

- The vast majority of NVM Express™(NVMe) is leveraged as-is for Fabrics
 - NVM Subsystem, Namespaces, Commands, Registers/Properties, Power States, Asynchronous Events, Reservations, etc.
- Primary differences reside in enumeration and queuing mechanism

	Differences	PCI Express®(PCIe)	Fabrics
	ldentifier	Bus/Device/ Function	NVMe Qualified Name (NQN)
~ 90% Common Between	Discovery	Bus Enumeration	Discovery and Connect commands
PCIe and Fabrics	Queuing	Memory-based	Message-based
	Data Transfers	PRPs or SGLs	SGLs only, added Key
			IDF15

Enabling More Use Cases NVMe* over Fabrics

Flash Memory NVMe* over Fabrics

- Use NVMe* end-to-end to get the simplicity, efficiency and low latency benefits
- NVMe over Fabrics is a thin encapsulation of the base NVMe protocol across a fabric
 - No translation to another protocol (e.g., SCSI)
- NVMe over Fabrics 1.0 includes RDMA binding enabling Ethernet and InfiniBand[™]
 - INCITS T11 defining Fibre Channel binding



Flash Memory Summit 2016 Santa Clara, CA

*Other names and brands may be claimed as the property of others.

Intel Corporation

NVMe® Technology Roadmap Evolution





The Evolution of NVMe® Technology



SCALE OVER FABRICS

UNIFY PCIE* SSDs

2010

2020





Jeployments

Intel Corporation

NVMe® Technology Powers the Connected Universe



Unite (Ku)	2046	2047	2040	2040	2020*	2024*	Average Capacity (GB)			5)
Units (Ku)	2010	2017	2010	2019	2020"	2021		2016 2017	2018 2019 202	20* 2021*
							6000			
Enterprise	364	749	1,069	2,045	4,067	5,554				
							4000			
Cloud	2,051	3,861	10,369	12,276	18,982	21,999				
							2000			
Client	33 128	18 051	82 587	1/13 236	202 348	258 701				
Client	55,120	40,901	02,007	140,200	202,540	200,791	0			
* Projections p	provided by	Forward Ins	ights Q2'20)				Client	Enterprise	Cloud

- NVMe technology grew from 3 Petabytes to 29 PB shipped per year from 2016 to 2019
- For 2020, the projection is 54 PB
- NVMe technology demand projected to remain strong in a post COVID world



Driving Simplicity in a World of Complexity

- Focused on core values... Fast, Simple, Scalable
- Foster areas of innovation AND avoid impact to broadly deployed solutions
- Create an extensible infrastructure that will take us through the next decade of growth





Specification Families



- The core of NVMe and NVMe over Fabrics integrated into a base specification
- Modular command set specifications (Block, Zoned Namespaces, Key Value, etc)
- Modular transport layer specifications (PCI Express*, RDMA, TCP)
- Maintain Management Interface as separate modular specification

NVM Express Technology Specification



nvm

Intel Corporation



Q&A answers by:



I Dal Allan



Tom Gardner



Amber Huffman



Jai Menon



Grant Saviers

