

Reconfigurable Accelerators for Big Data and Cloud RAW 2016

(23rd Reconfigurable Architectures Workshop @ IPDPS)

H. Peter Hofstee, Ph.D.

Distinguished Research Staff Member, IBM Austin Research

Professor, Big Data Systems, TU Delft

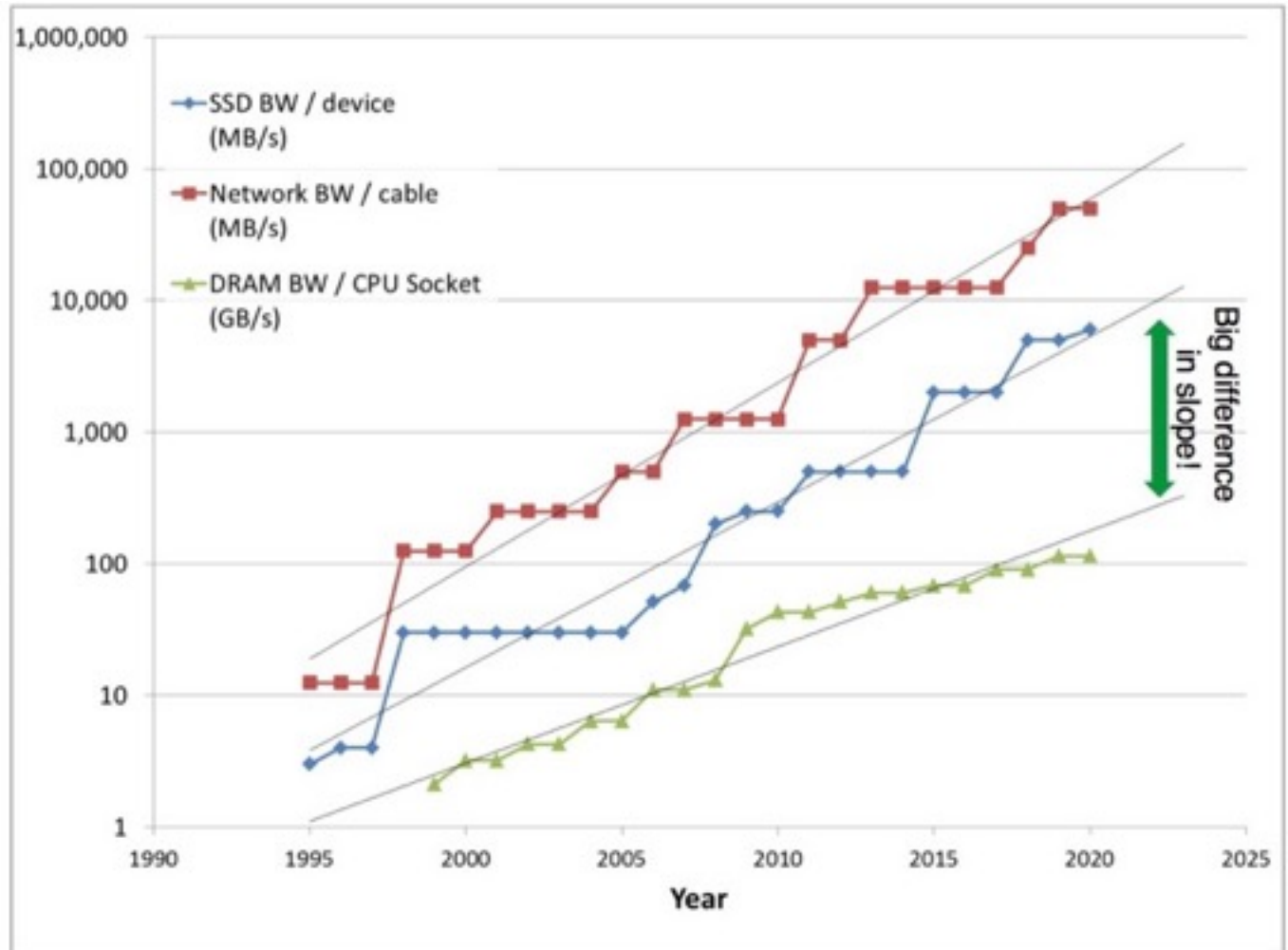
May 23, 2016



Network, Storage, & DRAM trends

Log scale

- Use DRAM Bandwidth as a proxy for CPU throughput
- Reasonable approximation for DMA and poor cache performance workloads (e.g. Storage)

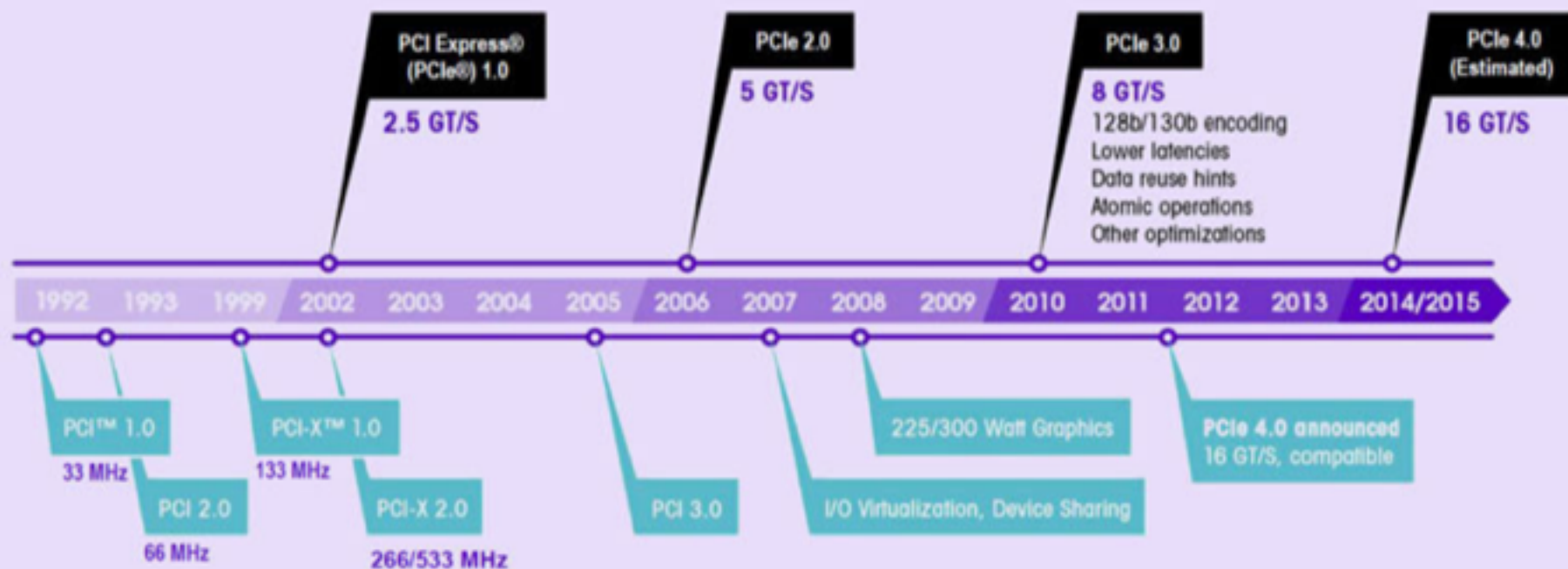


Source: Sandisk IT Blog



PCI & PCIe Technology Timeline

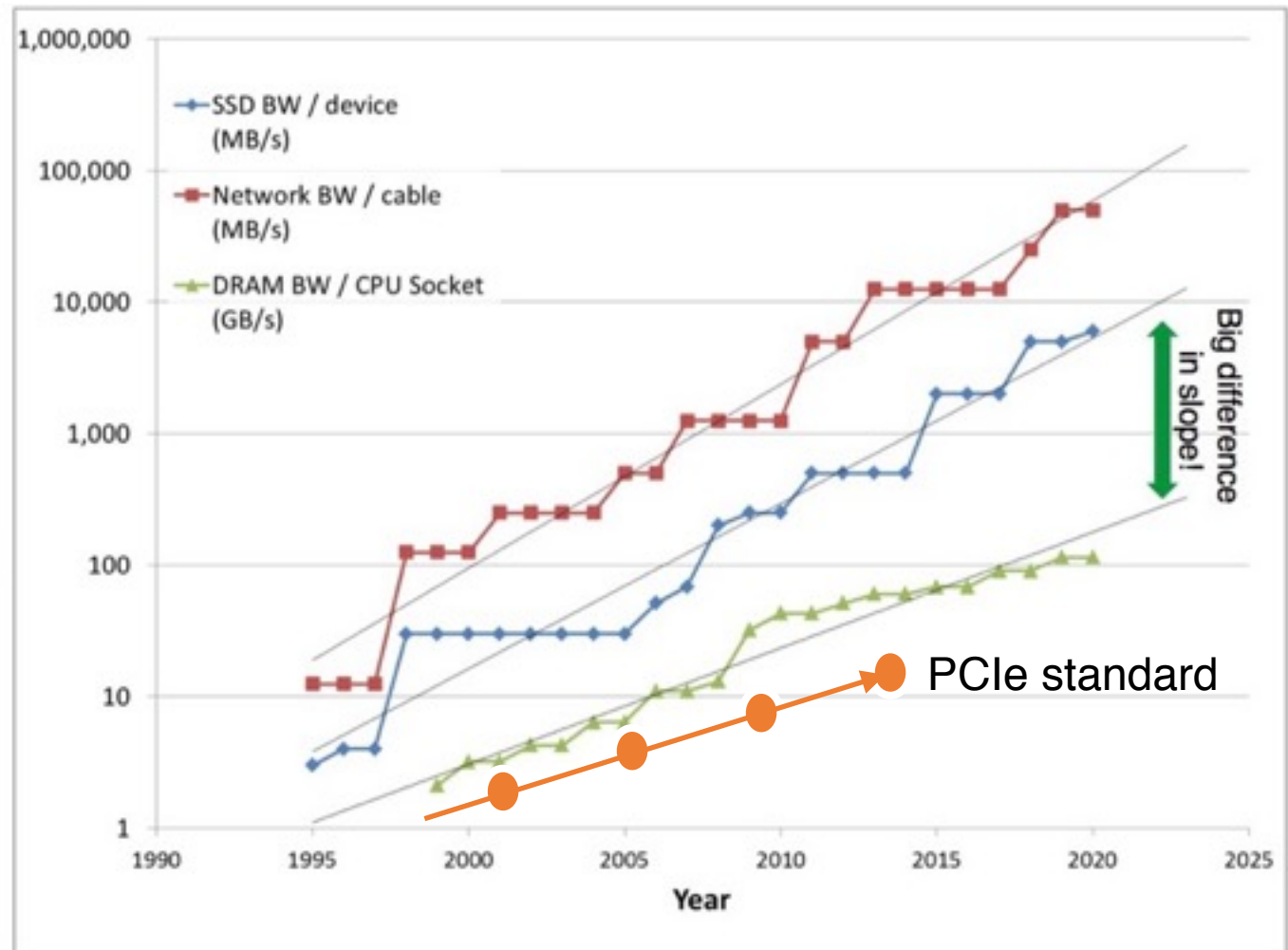
- PCI-SIG is the leader in I/O standard development and continues to evolve to meet existing and emerging usage models across the computing industry



Network, Storage, & DRAM trends

Log scale

- Use DRAM Bandwidth as a proxy for CPU throughput
- Reasonable approximation for DMA and poor cache performance workloads (e.g. Storage)

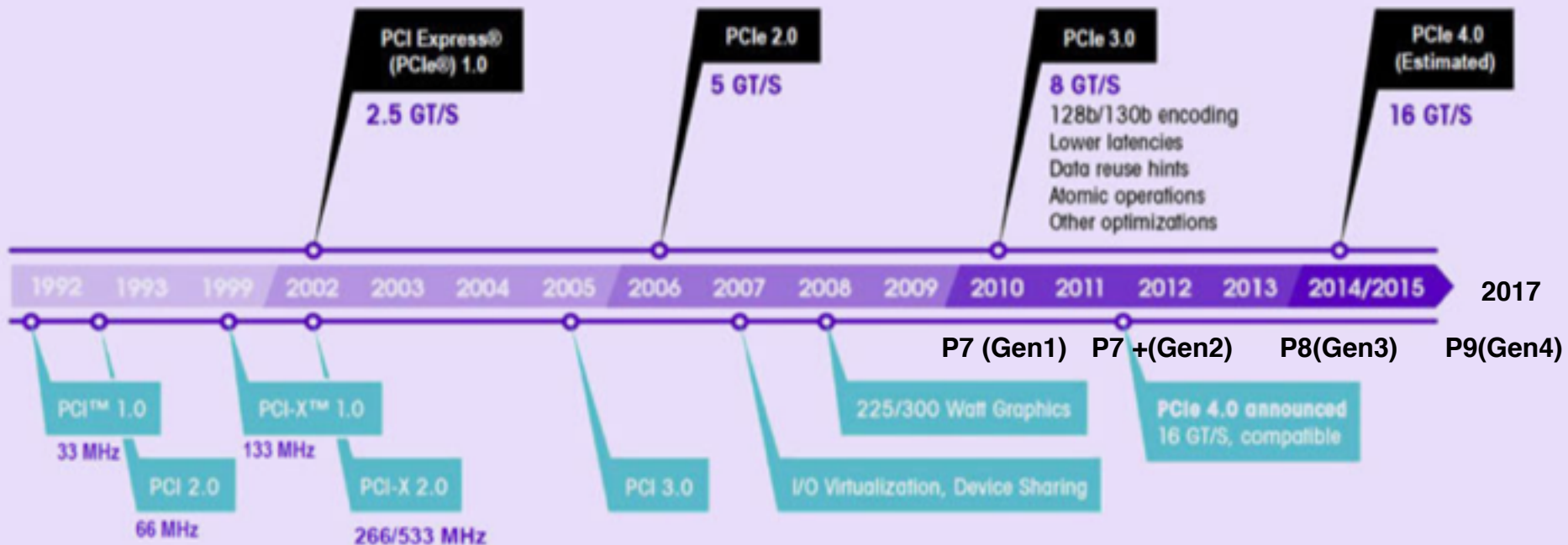


Source: Sandisk IT Blog



PCI & PCIe Technology Timeline

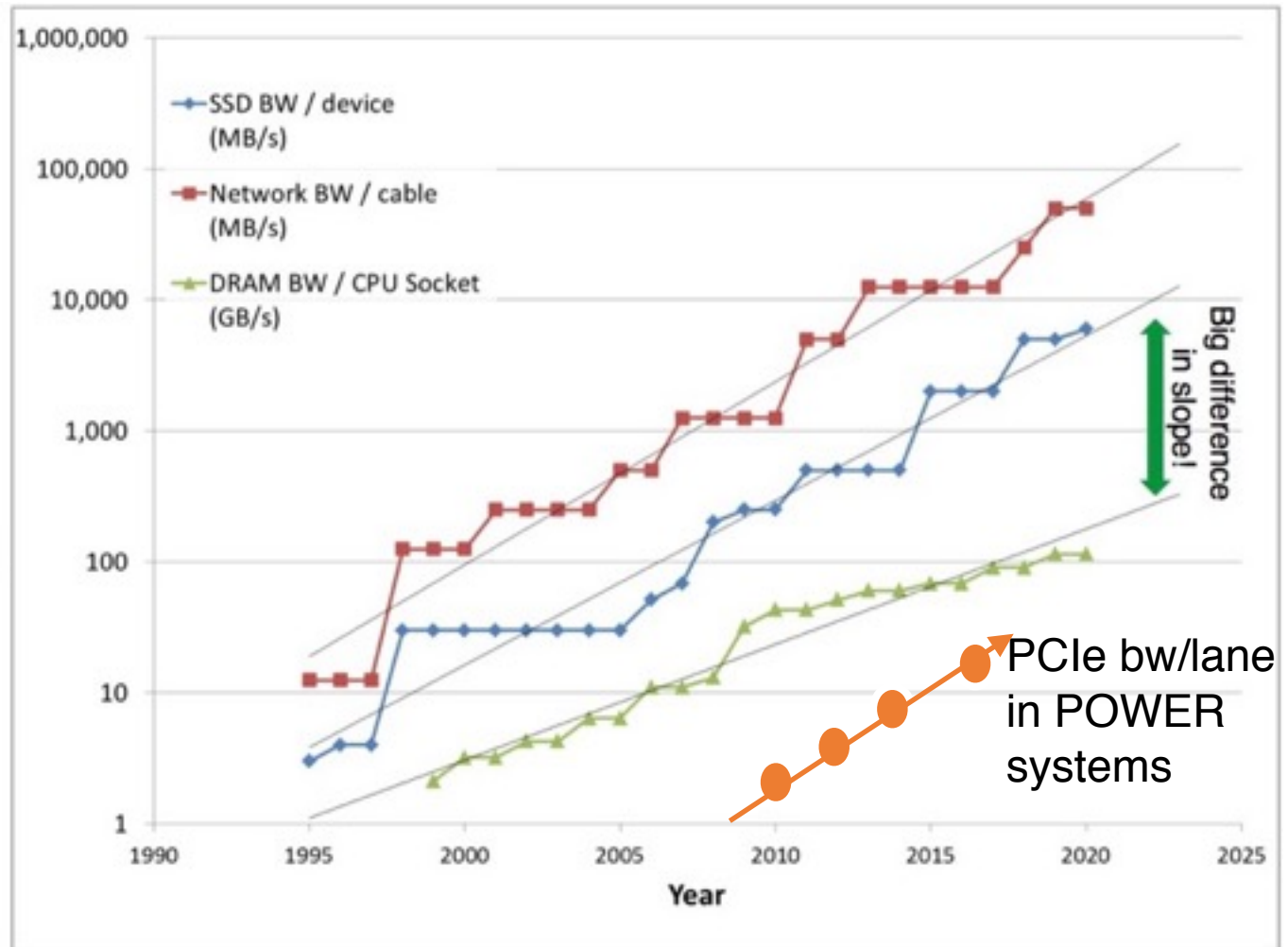
- PCI-SIG is the leader in I/O standard development and continues to evolve to meet existing and emerging usage models across the computing industry



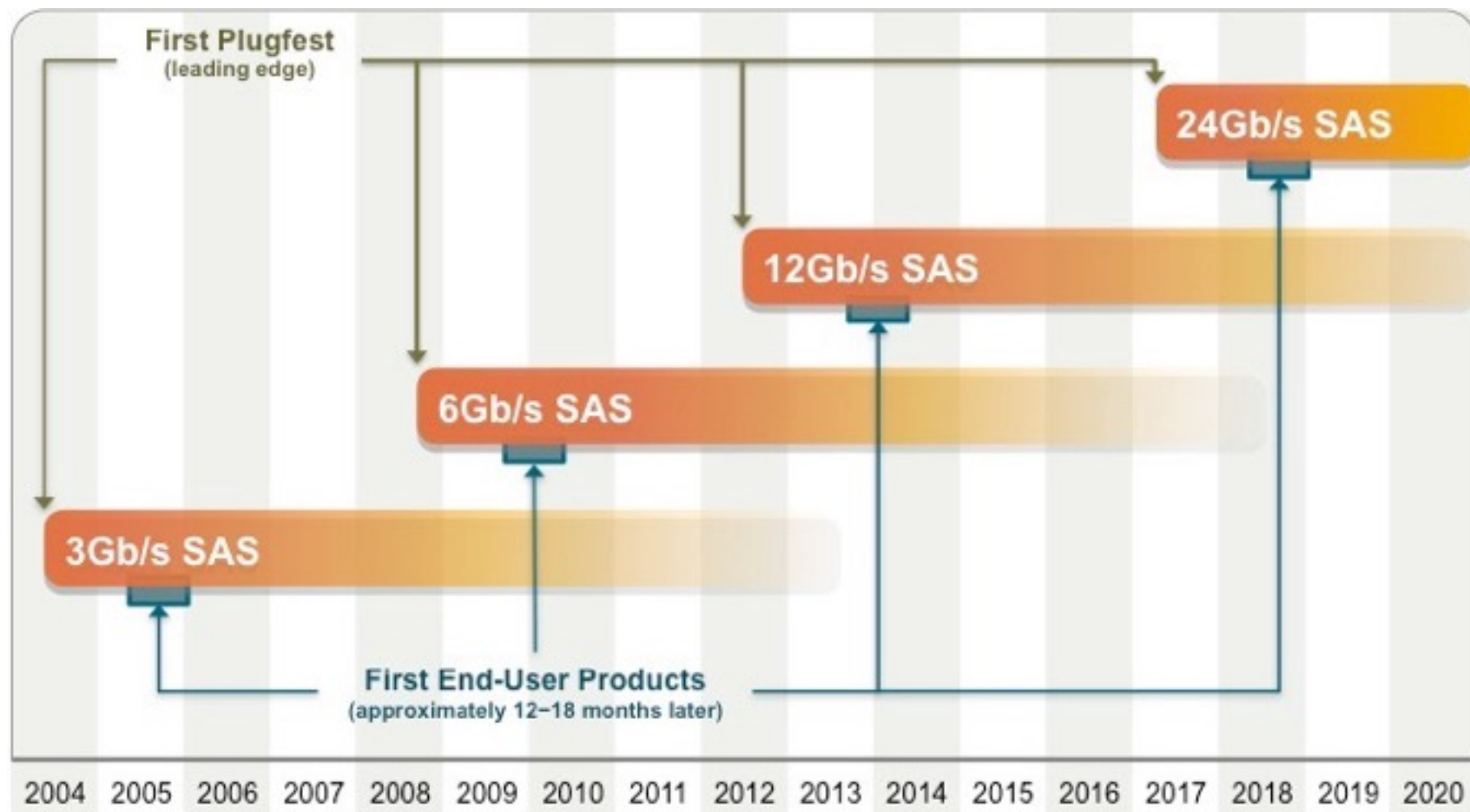
Network, Storage, & DRAM trends

Log scale

- Use DRAM Bandwidth as a proxy for CPU throughput
- Reasonable approximation for DMA and poor cache performance workloads (e.g. Storage)



Source: Sandisk IT Blog



Summary Performance Data – HDD, SSHD, SSD

Class	Type	FOB IOPS	IOPS (higher is better)				Throughput (larger is better)		Response Time (faster is better)	
			RND 4KIB 100%VW	RND 4KIB 100%VW	RND 4KIB 65:35 RVV	RND 4KIB 100% R	SEQ 1024KIB 100%VW	SEQ 1024KIB 100% R	RND 4KIB 100%VW AVE	RND 4KIB 100%VW MAX
HDD & SSHD										
7,200 RPM SATA Hybrid R30-4	2.5" SATA 500 GBWCD	125	147	150	135	97 MB/s	99 MB	15.55 msec	44.84 msec	
15,000 RPM SAS HDD IN-1117	2.5" SAS 80 GBWCD	350	340	398	401	84 MB/s	90 MB/s	5.39 msec	97.28 msec	
Client SSDs										
mSATA SSD R32-336	mSATA 32 GBWCD	18,000	838	1,318	52,793	79 MB/s	529 MB/s	1.39 msec	75.57 msec	
SATA3 SSD INB-1025	SATA3 256GBWCD	56,986	3,147	3,779	29,876	240 MB/s	400 MB/s	0.51 msec	1,218.45 msec	
SATA3 SSD R30-5148	SATA3 256GBWCE	60,090	60,302	41,045	40,686	249 MB/s	386 MB/s	0.35 msec	17.83 msec	
Enterprise SSDs										
Enterprise SAS SSD R1-2288	SAS 400GBWCD	61,929	24,848	29,863	53,942	393 MB/s	496 MB/s	0.05 msec	19.60 msec	
Server PCIe SSD INI-1727	PCIe 320GBWCD	133,560	73,008	53,797	54,327	663 MB/s	772 MB/s	0.05 msec	12.60 msec	
Server PCIe SSD IN24-1349	PCIe 700GBWCD	417,469	202,929	411,399	684,284	1,343 MB/s	2,053 MB/s	0.05 msec	0.58 msec	

Source: Calypso Testers, 2013

Today

- PCIe NVMe commoditized
 - Better BW/\$ than HDD !
- Typical enterprise NVMe
 - 3.2 GB/s seq. read
 - 3.2 TB
 - 800K IOP (4K)
 - 25W
 - 50-100usec
- 2U with 24 SFF NVMe
 - 76.8 GB/s (96 lanes of PCIe Gen 3)
 - all available lanes in a dual-socket POWER8
 - 76.8 TB
 - 19.2 Million IOPs
 - 600W(max) / 24 NVMe



SuperMicro

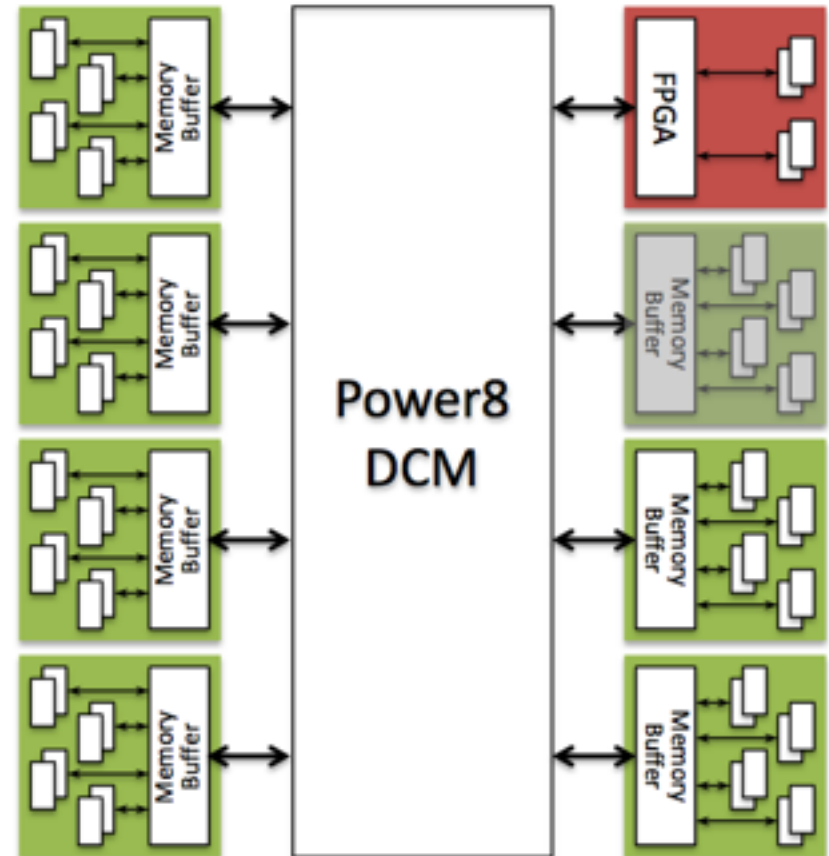
Memory-class bandwidth even before storage-class memory!

Three Ways to Attach

- As memory
 - Real address space
 - Any accelerator is highest privilege level only
 - Not attractive if latencies are high
- As I/O
 - I/O space
 - Offload (including DMA) to pinned memory only
 - Typically involves extra copy operation and coordination w.CPU
- On the SMP bus
 - Full access to all system resources
 - Best support for offload and near-memory compute
 - Most efficient synchronization mechanisms
 - Adapters can interact with other adapters w/o any CPU cycles

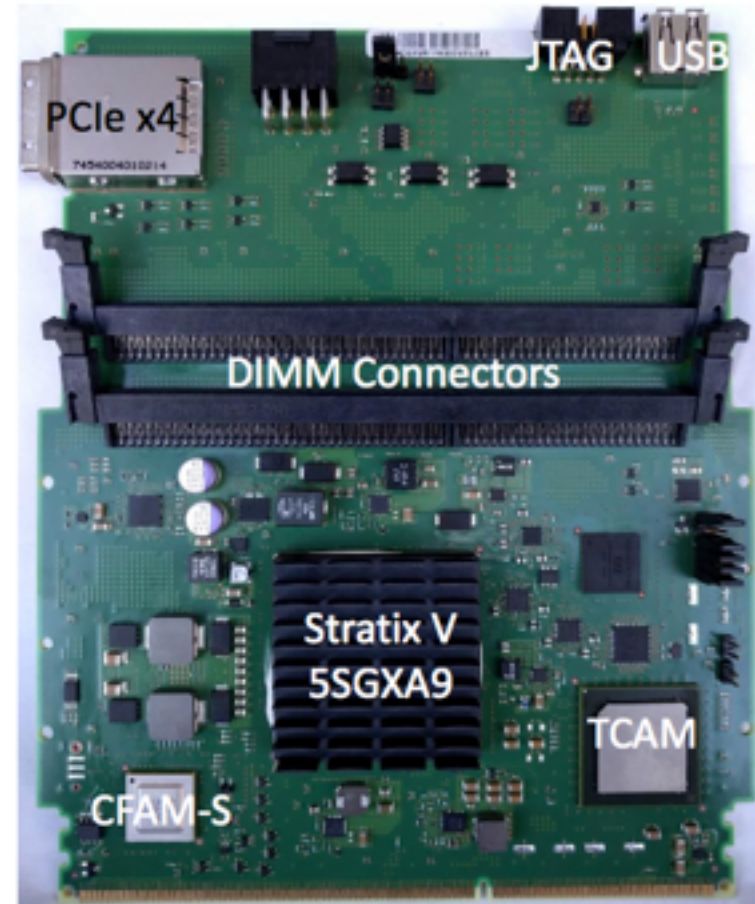
P8 Memory Sub-System with ConTutto

- Built an FPGA-based card that plugs into the DMI slot
- Enables regular system operation with any mix of CDIMMs and ConTutto cards populated
- Full compatibility with DMI protocol
- Memory controllers implemented in fabric logic and independent of DMI protocol logic
- Flexible system architecture enables easy implementation of additional features



ConTutto Card

- Intended to be an experimentation and prototyping vehicle
- Card characteristics
 - 10 signal layers
 - 10 power/ground layers
- Plug compatible with CDIMM, but 2.5" higher-- and DIMMs add width
- Large Altera FPGA with capacity for additional function incorporated
- CFAM-S (connection to service processor) enables system integration



DMI Connector
 OpenPOWER™

POWER8 CAPI Coherent Accelerator Processor Interface

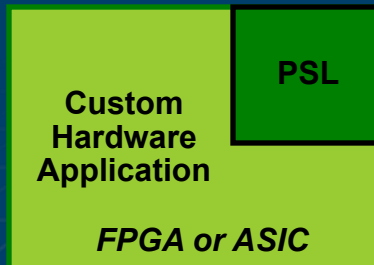
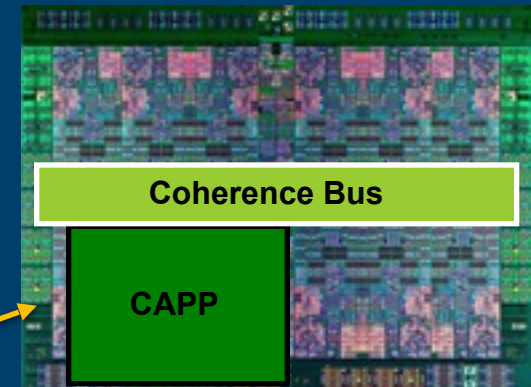
Virtual Addressing

- Accelerator can work with same memory addresses that the processors use
- Pointers de-referenced same as the host application
- Removes OS & device driver overhead

Hardware Managed Cache Coherence

- Enables the accelerator to participate in “Locks” as a normal thread
- Lowers Latency over IO communication model

POWER8



PCIe Gen 3

Transport for encapsulated messages

Processor Service Layer (PSL)

- Present robust, durable interfaces to applications
- Offload complexity / content from CAPP

Customizable Hardware Application Accelerator

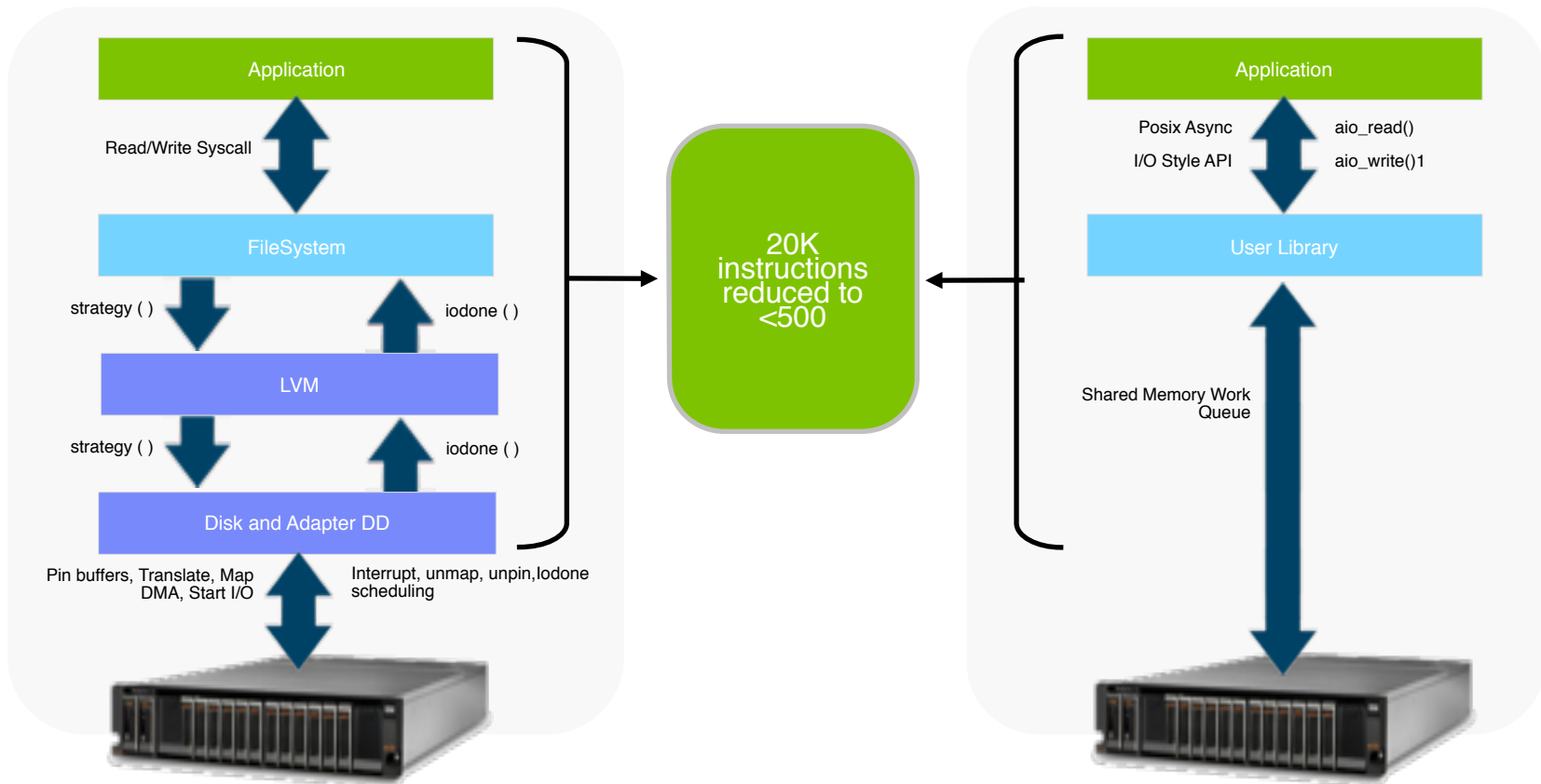
- Specific system SW, middleware, or user application
- Written to durable interface provided by PSL

Attach on the SMP bus: CAPI example

CAPI Attached Flash Optimization

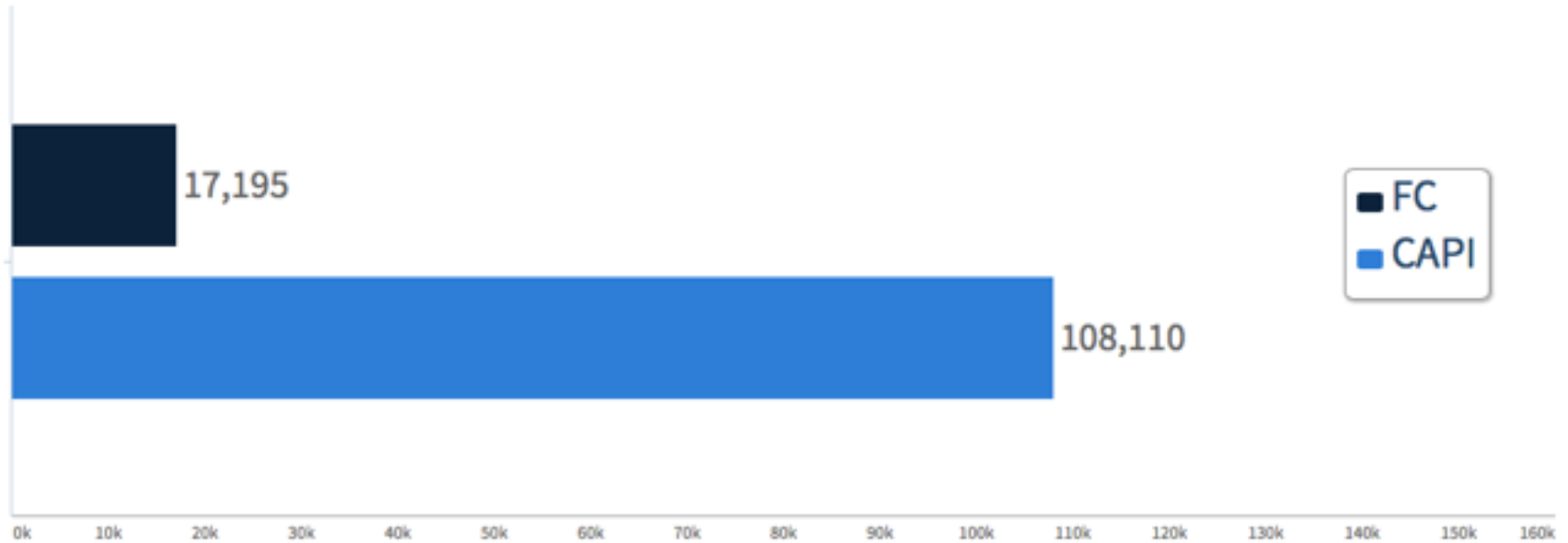


- Attach TMS Flash to POWER8 via CAPI coherent Attach
- Issues Read/Write Commands from applications to eliminate 97% of code pathlength
- Saves 20-30 cores per 1M IOPs

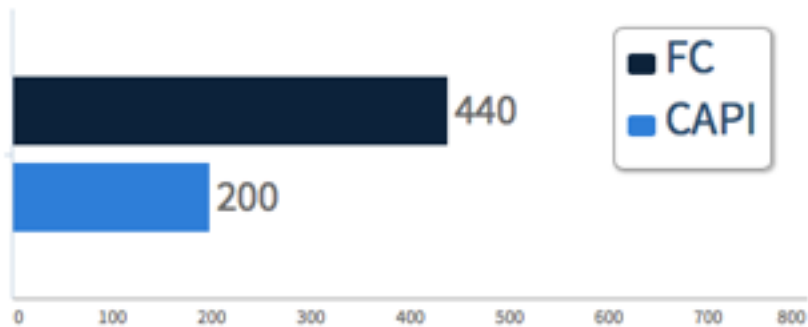


Attach on the SMP bus: CAPI example

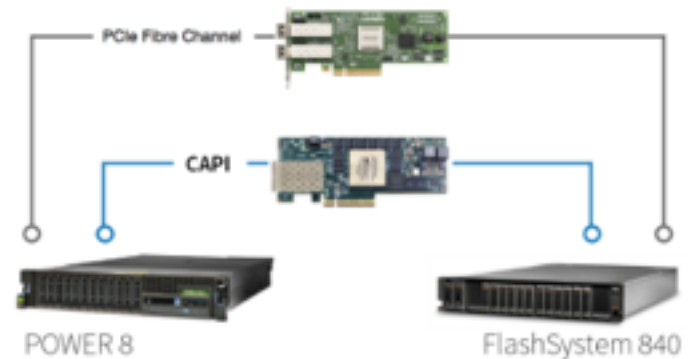
CAPI vs. Fibre Channel



IOPS per HW thread →

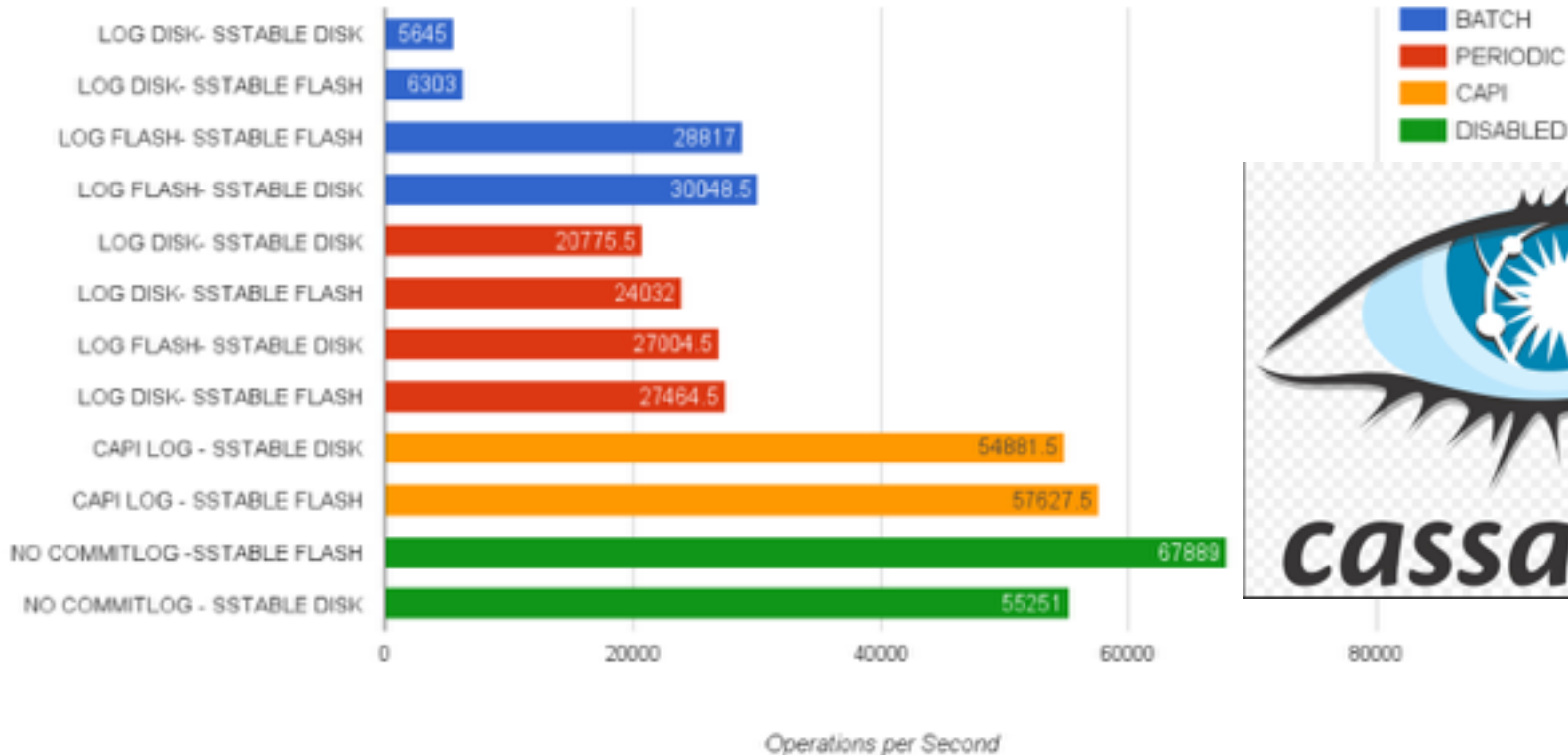


Latency [us] →



Attach on the SMP bus: CAPI example

YSCB 10M Records Insert - 64 Client Threads



CASSANDRA WRITES
w. Bedri Sendir



- **Periodic Commitlog:** Syncs every `commitlog_sync_period_in_ms`. *Writes are not durable.*
`commitlog_sync_period_in_ms` : 10000ms
- **Batch Commitlog:** Collects requests in batches `commitlog_sync_batch_window_in_ms` and sync's & acknowledge each add. *Durable writes.* `commitlog_sync_batch_window_in_ms` : 2ms
- **CAPI-Flash Commitlog:** Writes each mutation to the flash as soon as it is received on Cassandra server. It will guarantee that it writes on flash before acknowledging writes. *Fast Durable writes.*
- **Commitlog Disabled:** Writes go directly to the in memory data structure.



• BIG DATA • BIG ANALYTICS • BIG INSIGHTS •

Home About Whitepapers Events Subscribe

HOME

FEATURES ▾

SECTORS ▾

APPLICATIONS



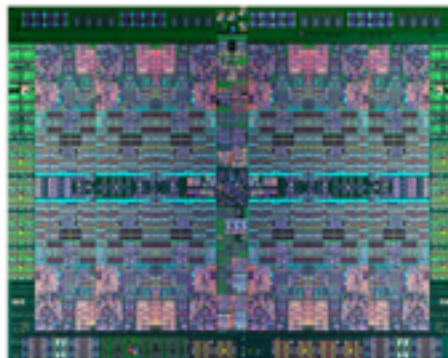
October 21, 2015

Neo4j Touts 10x Performance Boost of Graphs on IBM Power FGAs

Alex Woodie



Top News from
Leading Solution
Providers



Fraud Detection

Detecting fraud hidden in huge data sets should be easier with Neo4j running on Power8.

Neo Technology today announced that it's working with IBM to support its graph database on IBM Power8-based servers equipped with field programmable gate arrays (FPGAs). This will enable customers to run graph databases with hundreds of billions to trillions of edges, thereby tackling a new class of intractable big data problems in bioinformatics, fraud detection, and IoT analytics. In other news, Neo also open sourced a key piece of software.

Acceleration Use Case Example:

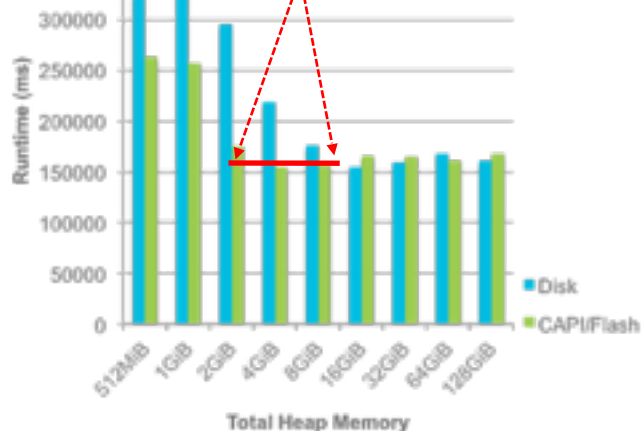


Fast and general engine for large-scale data processing

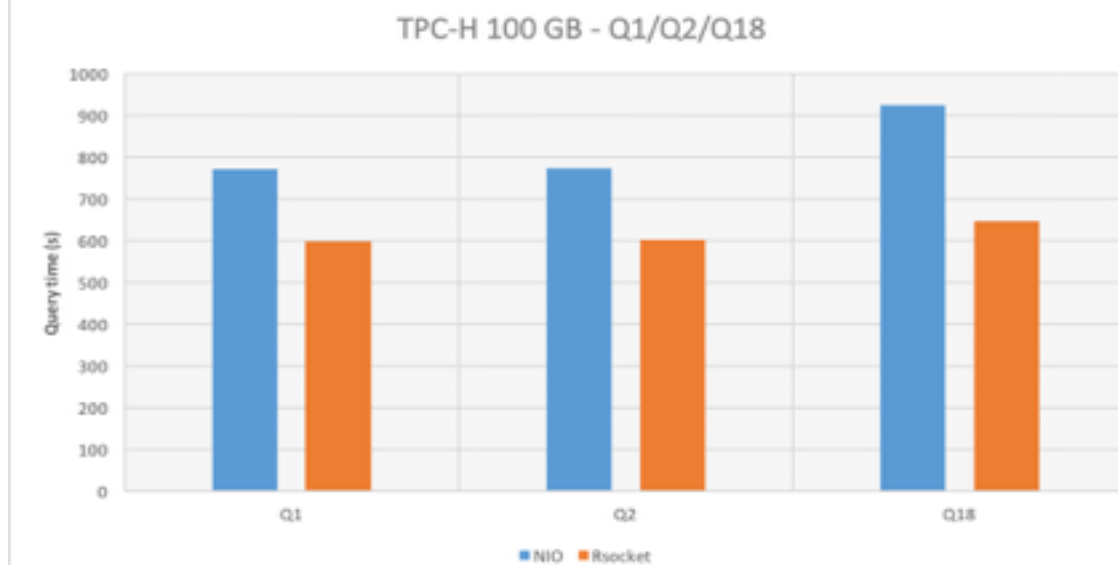
- CAPI Flash and RDMA Leveraged Transparently to Spark Applications
- Coming.... HDFS CAPI FPGA Erasure Code Acceleration, CAPI FPGA Compression Acceleration,

x Degrees of Separation on Spark

CAPI Flash for RDD Cache = 4X memory reduction at equal performance



RDMA for Spark Shuffle = 30% Better Response Time, Lower CPU Utilization



Attach on the SMP bus: CAPI example



In-the-box version of CAPI-attached Flash
(Uses standard M.2 formfactor flash)

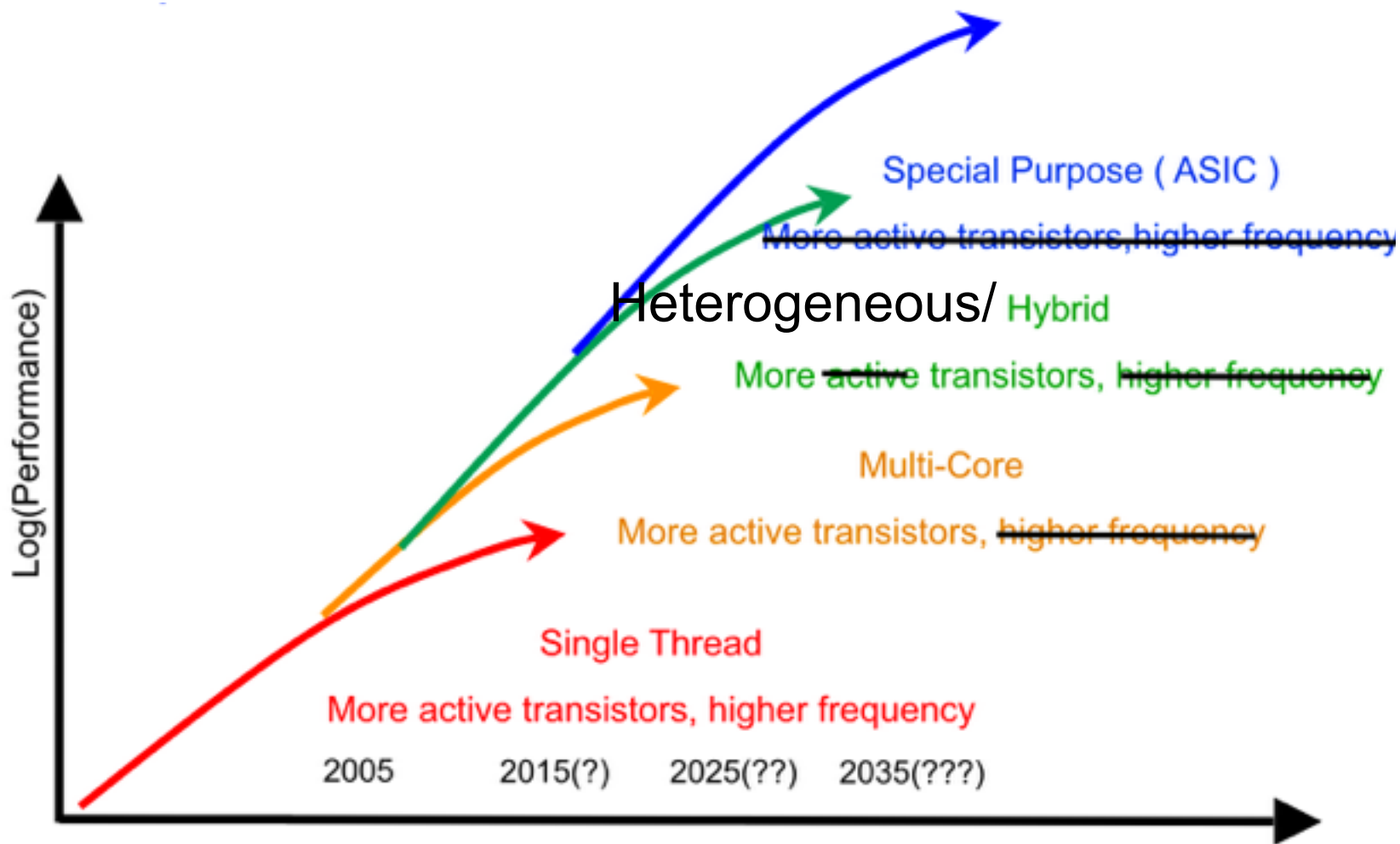


Nallatech CAPI Flash Accelerator based on a Xilinx Kintex UltraScale KU060 FPGA

Similar Story for Networking

- Fast transition from 1Gb/s to 10Gb/s to 100Gb/s
 - 25Gb, 50Gb really subsets of 100Gb
 - 40Gb looks less attractive (to me)
- Partly driven by in-memory paradigms like Spark
 - No longer just trying to match HDD speed
- Transition to optical (especially across racks)
 - Opens the door for much more bandwidth in future
- Coherent attach allows:
 - Most offload to network adapter
 - Network adapter to reach all resources (including storage) without touching the CPU

Usual Story for Compute



Reminder: Multi-Core Drives Memory Pressure

Conclusions so far ...

A LOT of pressure on bandwidth

Strong desire to be SMP-attached

So what are we doing about it?

1: This is a big challenge: don't go it alone ...



OpenPOWER

OpenPowerFoundation.org



POWER Processor Roadmap

Focus on Enterprise Technology and Performance Driven				Focus on Scale-Out and Enterprise Cost and Acceleration Driven				Future		
POWER6 Architecture		POWER7 Architecture		POWER8 Architecture		POWER9 Architecture		Partner Chip POWER8/9	POWER10	
2007 POWER6 2 cores 65nm New Micro-Architecture New Process Technology	2008 POWER6+ 2 cores 65nm+ Enhanced Micro-Architecture Enhanced Process Technology	2010 POWER7 8 cores 45nm New Micro-Architecture New Process Technology	2012 POWER7+ 8 cores 32nm Enhanced Micro-Architecture New Process Technology	2014 POWER8 12 cores 22nm New Micro-Architecture New Process Technology	2016 POWER8 w/ NVLink 12 cores 22nm Enhanced Micro-Architecture With NVLink New Process Technology	2017 P9 SO 24 cores 14nm New Micro-Architecture Direct attach memory New Process Technology	TBD P9 SU TBD cores 14nm Enhanced Micro-Architecture Buffered Memory	T B D	2018 - 20 P8/9 SO 10nm - 7nm Existing Micro-Architecture Foundry Technology	2020+ New Micro-Architecture New Technology
High Frequency Enhanced RAS Dynamic Energy Management		Large eDRAM L3 Cache Optimized VSX Enhanced Memory Subsystem		Optimized for Data-Centric Workloads Integrated PCIe CAPI Acceleration / I/O		Scale-Out Datacenter TCO Optimization Scale-up performance Optimization Acceleration Enhancements to CAPI and NVLINK Modularity for OpenPOWER			OpenPOWER Ecosystem Design Targeting Partner Markets & Systems Leveraging Modularity	New Features and Functions

Price, performance, feature and ecosystem innovation



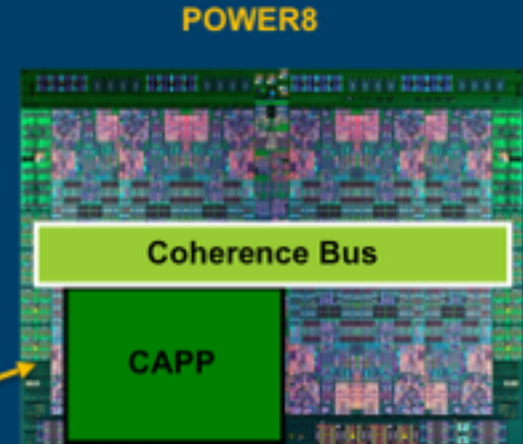
POWER8 CAPI Coherent Accelerator Processor Interface

Virtual Addressing

- Accelerator can work with same memory addresses that the processors use
- Pointers de-referenced same as the host application
- Removes OS & device driver overhead

Hardware Managed Cache Coherence

- Enables the accelerator to participate in "Locks" as a normal thread
- Lowers Latency over IO communication model



PCIe Gen 3
Transport for encapsulated messages

Customizable Hardware Application Accelerator

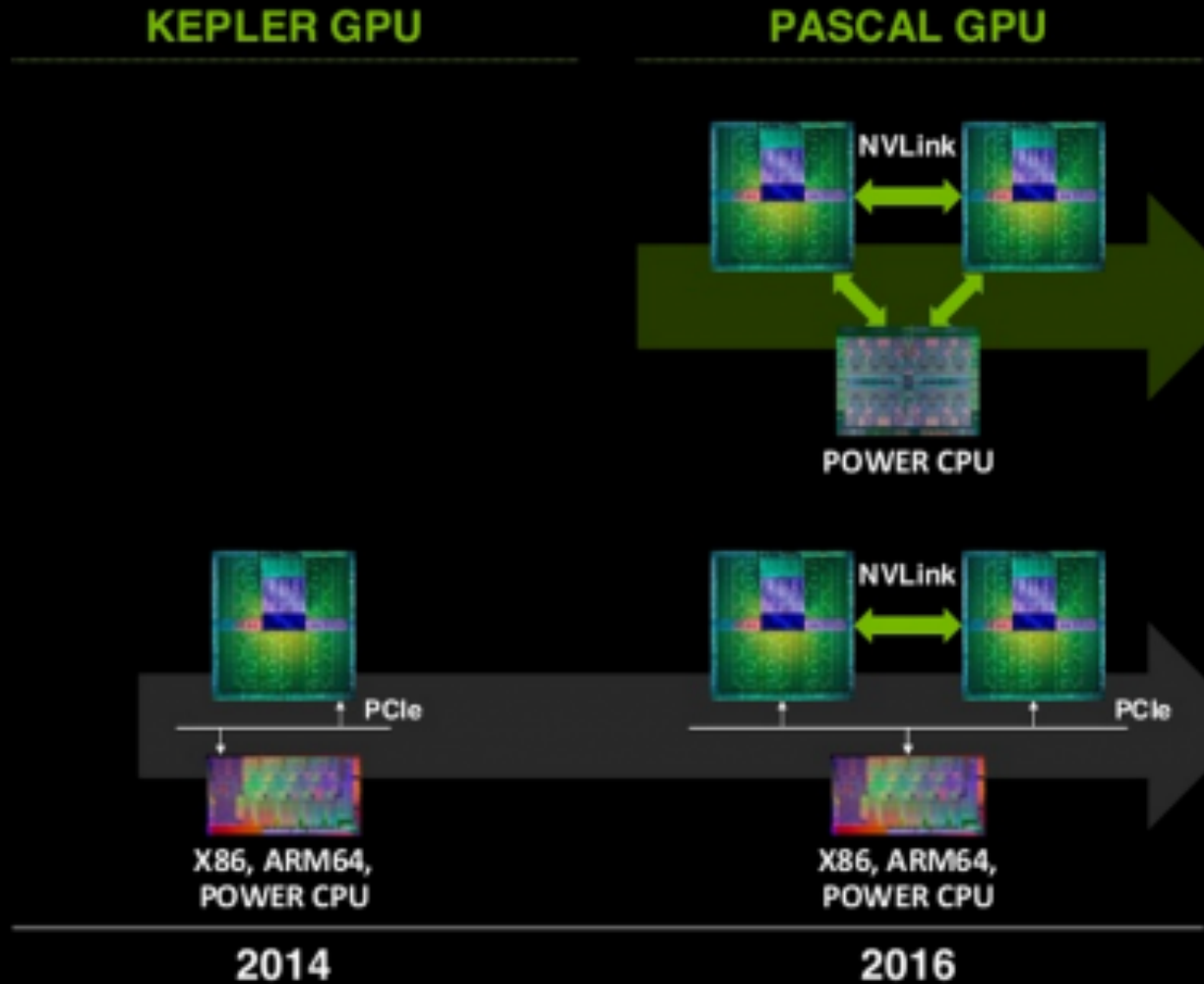
- Specific system SW, middleware, or user application
- Written to durable interface provided by PSL

Processor Service Layer (PSL)

- Present robust, durable interfaces to applications
- Offload complexity / content from CAPP

NVLink

High-Speed GPU Interconnect

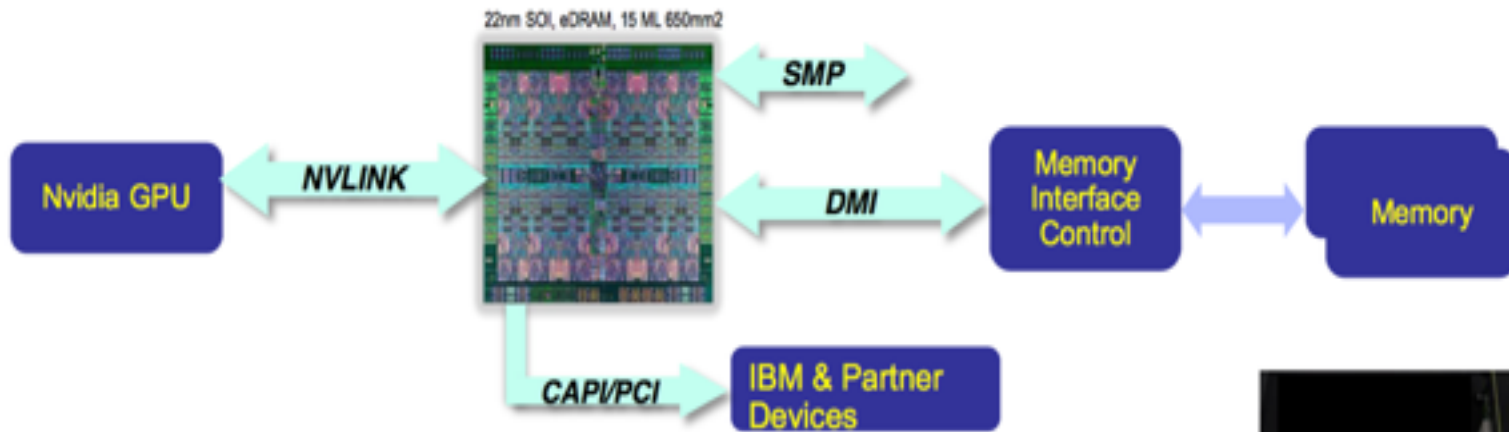


5-12x PCIe Gen3 x16

Source: NVIDIA

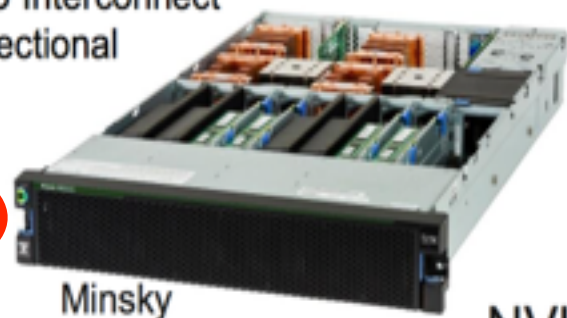
POWER8 + NVLink

POWER8 with NVLink



- NVLink High Speed CPU <-> GPU Interconnect
- 160+ GigaBytes per second bi-directional
- 5-12x faster than PCIe Gen3 x16

Nvlink Accelerator Lab
accellab@us.ibm.com



Minsky



Zoom

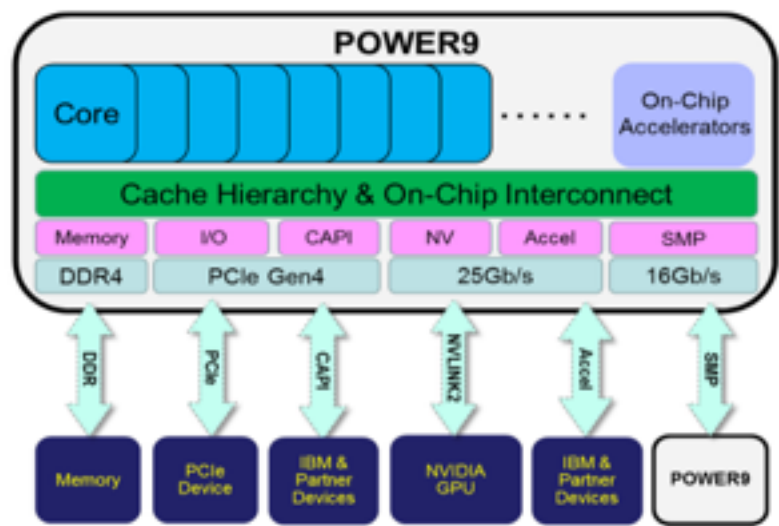


Inventec

NVLink POWER Systems



POWER9 SO with Advanced Accelerator Attach



24 newly designed POWER9 cores

- Leveraging execution slices for improved performance on cognitive, analytic, and big-data applications

Large, low-latency, eDRAM cache for big datasets

Global Foundries 14HP finFET technology with eDRAM

Cloud-focused innovation in Energy Efficiency, Security, and Quality of Service

State-of-the-art IO subsystem using PCIe Gen4

Leadership platforms for hardware acceleration

- High bandwidth, GPU interconnect (NV link2.0)
- Next-generation CAPI2.0 interface for coherent accelerator and storage attach
- On-chip compression & cryptography accelerators
- New 25Gb/s advanced accelerator attach bus

1st chip in POWER9 family

- Optimized for 2 socket scale out servers & hyperscale datacenters
- DDR4 direct attach memory channels

Full POWER9 family will address a broad range of scale out & enterprise servers

© 2016 IBM Corporation

FPGA Acceleration Service

3 use cases to experience the FPGA acceleration services.

- 🎮 Try FPGA acceleration service with a demo
- 📄 How to upload a private accelerator
- 📄 How to create a VM with accelerator

Try FPGA acceleration service with a demo

Play a demo with a pre-installed application and attached with an FPGA accelerator

- Visit SuperVessel Cloud (<http://www.ptopenlab.com>), click "Application Acceleration"
- Clone a dedicated virtual machine to run the Informix® demo. Get your host IP address.



- Connect VPN, go to Informix® dashboard and play the demo



How to upload private accelerator

To work with user's own accelerator, upload it to the cloud.

- Login to the FPGA Maker Zone



- Fill in accelerator information and upload.

Accelerator information

Required fields are marked with an asterisk (*) and must be filled in.

Name:

Accelerator Type: ①

Version:

Accelerator Logo: ②

Accelerator File: ③

- ① Upload CAPI accelerator or PCIe accelerator
- ② Choose a picture in local disk
- ③ Choose the accelerator package in local disk

- Check accelerator deployment status. All procedures should be green.

How to create VM with accelerator

Work in Linux environment with FPGA accelerator attached.

- Visit SuperVessel Cloud and click Apply VM->Launch Instance



- ① For PCIe, use KVM image. for CAPI, use docker
- ② Choose public or private accelerator
- ③ Launch

- Login the host with web console or SSH tool, start your development

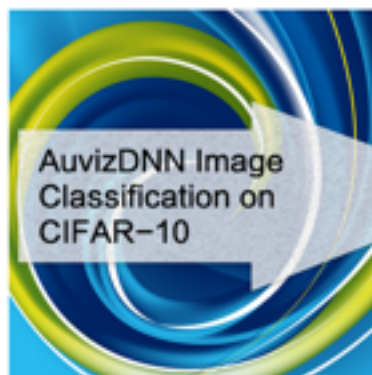
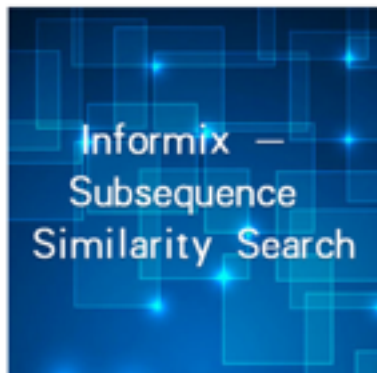
➤ Apply SuperVessel services from here: <http://www.ptopenlab.com>

Facebook: <http://www.facebook.com/supervesselcloud>

WeChat "SuperVessel"



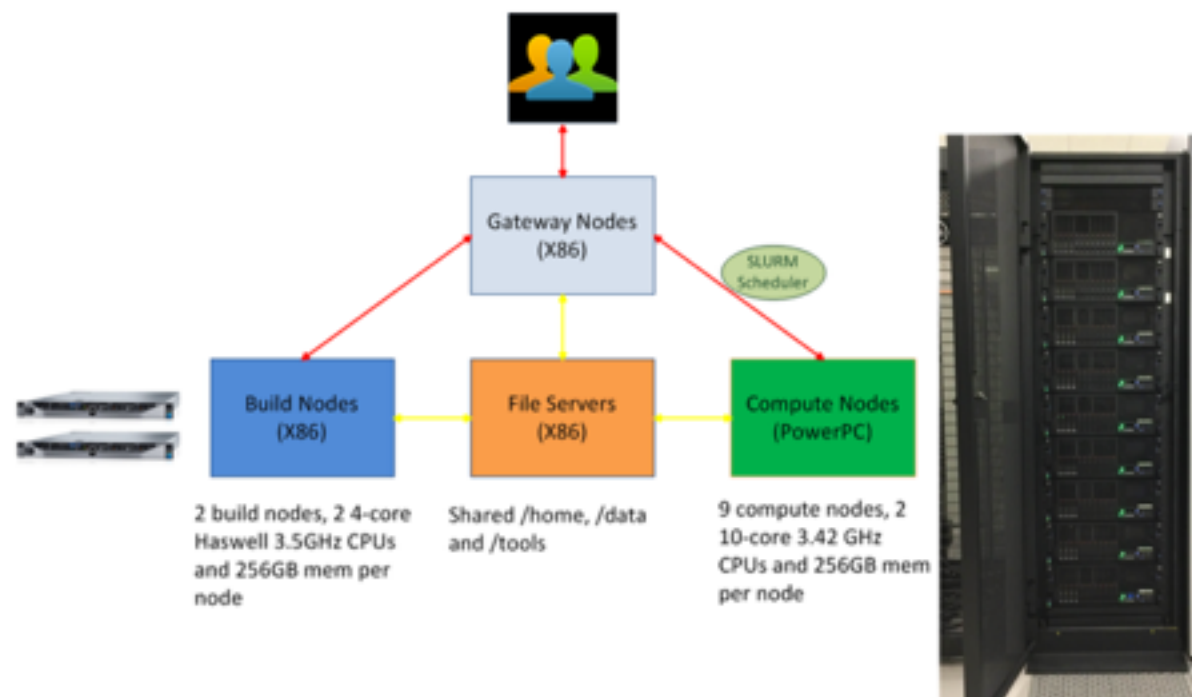
SuerVessel QQ group 344373069



FABRIC: FPGA Accelerator Research Infrastructure Cloud

IBM POWER8+CAPI cluster

The IBM POWER8+CAPI Cluster is a cluster of several x86 servers and nine POWER8 servers. Each POWER8 node is a heterogeneous platform capable of running GPGPU- and/or FPGA- accelerated applications.



Our current setup supports three accelerating devices:

- Nallatech 385 A7 Stratix V Altera-based FPGA adapter
- Alpha-data 7V3 Virtex7 Xilinx-based FPGA adapter
- NVIDIA Tesla K40m GPGPU card

<https://wikis.utexas.edu/display/fabric/Home>

OpenPOWER DEVELOPER CHALLENGE

Two tracks to challenge and win:

1. The Open Road Test
 - Port and optimize for OpenPOWER
 - Go faster with accelerators (*optional*)
2. The Spark Rally
 - Train an accelerated DNN and recognize objects with greater accuracy
 - Show you can scale with Spark



Key Dates

Register today

openpower.devpost.com

Register Now

Contest:

June 1 - Sep 1, 2016

**Grand prizes include a trip to
Supercomputing 2016**

Other prizes include iPads, Apple Watches

Join the conversation at #OpenPOWERSummit

© 2016 IBM Corporation

Thank You



Follow @OpenPOWERorg and use #OpenPOWER to share your thoughts and join the conversation!