

# HPC Strategy & OpenPOWER

# Our view of high-performance computing has evolved

### • The "Old View" of HPC

- The Value of an HPC System is Measured by FLOPS and TOP500 rank
- The Objective is to Make an Algorithm Run Faster
- HPC is a Special Category of Computing
- HPC Looks Only at the Cluster/Server
- Storage is an Afterthought

### • The IBM View of HPC

- Value is Measured by Application Performance
- The Objective is to make a workflow optimized
- HPC is another form of Analytics
- Influx of Large Data demands consideration of Data Management and Storage in HPC: We Must Look Beyond the Server. Performance and data availability are imperative

# Workflows Define HPC: Oil and Gas Example

•

•



The capability of any single piece of hardware is not what drives workflow.

# **Portfolio of HPC Solutions**

Processors & Systems	<ul> <li>High Performance Processors &amp; Systems</li> <li>Accelerator, networking, storage integration via NVLink &amp; CAPI</li> <li>Highest memory throughput</li> </ul>			
	<ul> <li>High speed interconnect / network fabric from Mellanox Technologies</li> </ul>			
High Speed Interconnect	MPI acceleration in the IB fabric, reducing CPU overhead Support for GPUDirect, NVMe over fabric			
High Performance File System & Storage	<ul> <li>Highest Performance HPC Storage: Elastic Storage Server</li> <li>High Performance Spectrum Scale (GPFS) Parallel File System</li> <li>Data centric design</li> </ul>			
	<ul> <li>Deployment tools, integrated management</li> </ul>			
HPC Software	<ul> <li>Compilers: gcc, IBM XLC, LLVM OpenMP4, PGI Fortran/C/C++, Java, OpenACC, OpenMP</li> </ul>			
	<ul> <li>Debuggers, Profilers, Math libraries, MPI &amp; HPC apps</li> </ul>			

# **OpenPOWER: Open Architecture for HPC & Analytics**

Processor IP Licensing

Licensing processor core to enable semiconductor partners like Suzhou Powercore to build POWER chips

Open Interfaces

Tight integration using CAPI & NVLink with Accelerators (NVIDIA, Xilinx), Networking (Mellanox), Storage (CAPI Flash)

Systems & Software Enabling System Partners to build POWER-based servers and Open Sourcing Software including Firmware & Hypervisor





Introducing the OpenPOWER Foundation... 5 Founding members in 2013



# 2016: 250+ Members



# Membership Options

Anyone may participate in OpenPOWER. Membership levels are designed for those that are investing to grow and enhance the OpenPOWER community and its proliferation within the industry.

	Membership Level	Annual Fee \$ USD	FTEs	Technical Steering Committee	Board / Voting position
	Platinum	\$100k	10	One seat per member not otherwise represented	Includes board position Includes TSC position
	Gold	\$60k	3	May be on TSC if Work group lead	Gold members may elect one board representative per three gold members
	Silver	<b>\$20k</b> \$5k if <300 employees	0	May be on TSC if Work group lead	Sliver members may elect one board representative for all silver members
w	Silver ISV	\$0 if ISV is <300 employees	0	May be on TSC if Work group lead	Sliver members may elect one board representative for all silver members
	Associate & Academic	\$0	0	May be on TSC if Work group lead	May be elected to one community observer, non-voting Board seat

Ne

# 2300+ Linux Applications on POWER





**Major Linux Distros** 



# **IBM Power Systems LC Line**

OpenPOWER servers for cloud and cluster deployments that are different by design





### Available now: Barreleye

In partnership with Avago, IBM, Mellanox, PMC & Samsung













# Zaius 1.25 OU

- 2 POWER9 CPUS
- 32 DDR4 DIMM SLOTS
- 2X G4 PCIE X16 FHFL SLOTS
- 1X G4 X16 HHHL SLOT
- 1X G4 X16 OCP MEZ
- 1X M.2 SATA PORT
- 1X SATA PORT
- 15X 2.5" SAS/SATA/NVME
   SLOTS
- BMC W/GBE LOM
- "DISKLESS" OPTION

### **OpenPOWER Innovation in the Design**

Power Systems S822LC for High Performance Computing (aka Minsky)



#### **NVIDIA:**

Tesla P100 GPU Accelerator with NVLink (GPU $\leftrightarrow$ GPU & GPU $\leftrightarrow$ CPU)

**Ubuntu by Canonical:** *Launch OS* supporting NVLink and Page Migration Engine

Wistron: Platform co-design

**Mellanox:** InfiniBand/Ethernet Connectivity in and out of server

**HGST:** Optional NVMe Adapters

**Broadcom:** Optional PCIe Adapters

**QLogic:** Optional Fiber Channel PCIe

Samsung: 2.5" SSDs

Hynix, Samsung, Micron: DDR4

**IBM:** POWER8 CPU with NVLink

# **IBM Strategy for HPC Systems**







High Performance Cores Fast & Large Memory System

Fast PowerAccel Interconnects for Accelerators

Faster Cores than x86 Larger Caches Per Core than x86 5x Faster Data Communication between POWER8 & GPUs

# **Roadmap for HPC / HPDA**

Mellanox	Connect-IB	ConnectX-4	ConnectX-5
Interconnect	FDR Infiniband	EDR Infiniband	Next-Gen Infiniband
Technology	PCle Gen3	CAPI over PCIe Gen3	Enhanced CAPI over PCIe Gen4
	Kanlan	Descal	) / alta
NVIDIA GPUs	Kepler		
	PCIe Gen3	INVLINK	NVLINK Next Gen

### POWER9

### **POWER8 with NVLink**

### IBM CPUs

# POWER8

OpenPower CAPI Interface



PowerAccel Interfaces: NVLink, CAPI, PCIe Gen3



PowerAccel: Enhanced CAPI, NVLink Next Gen, PCIe Gen4



### Why Accelerators and GPUs?



Shift back towards the Moore's Law prediction through:

- 1. IBM HPC Innovation (processor architecture enhancement, scalable filesystems, workflow management
- 2. Acceleration through partner ecosystem (e.g. NVIDIA GPUs deliver 2X perf/watt)

### **POWER8: Designed Memory Bandwidth**

### **IBM 22nm Technology**

- Silicon-on-Insulator, 15 metal layers,
- ~4.2 billion Transistors
- Deep trench eDRAM

### Compute

- 6/12 cores, ST/SMT2/SMT4/SMT8
- Enhanced, Auto balancing threads
- 8 dispatch/16 execution pipes/224 instructions in flight
- Transactional Memory/ Crypto & Crc instructions

### Cache

- 64KB L1 + 512KB L2 / core
- 96MB L3 + up to 128MB L4 / socket

### **System Interfaces**

- 230 GB/s memory bandwidth / socket
- Up to 48x Integrated PCI gen 3 / socket
- CAPI (over PCI gen 3)
- Robust, Large SMP Interconnect
- On chip Energy Mgmt, VRM / core





### **POWER8 Memory Organization (Max Config)**



### Memory bandwidth vs. most Xeon E5-2600v3 Configurations Based on STREAM Triad memory bandwidth

Deliver 79% greater memory bandwidth compared to Xeon E5-2600 v3 configurations with 2DPC

Deliver 60% greater memory bandwidth compared to Xeon E5-2600 v3 configurations with 1DPC

Only minor change vs Xeon E5-2600v4



• IBM Power System S822LC results are based on IBM internal measurements of STREAM Triad; 20 cores / 20 of 160 threads active, POWER8; 3.5GHz, up to 1TB memory,

Intel Xeon data is based on published data of Intel® Server System R2208WTTYS running STREAM Triad; 24 cores / 24 of 48 threads active, E5-2690 v3; 2.3GHz. For more details see <a href="http://www.intel.com/content/www/us/en/benchmarks/server/xeon-e5-2600-v3/xeon-e5-2600-v3-stream.html">http://www.intel.com/content/www/us/en/benchmarks/server/xeon-e5-2600-v3/xeon-e5-2600-v3-stream.html</a>

### What Does it Mean? Excellent CPU-Only Application Performance



Haswell Based

POWER8 – S822LC

# **Differentiated Acceleration - CAPI and NVLink**

### **CAPI-attached Accelerators**

#### POWER8



PSL

**FPGA or ASIC** 

#### New Ecosystems with CAPI

- Partners innovate, add value, gain revenue together w/IBM
- Technical and programming ease: virtual addressing, cache coherence
- Accelerator is hardware peer

### NVIDIA Tesla GPU with NVLink



#### Future, Innovative Systems with NVLink

- Faster GPU-GPU communication
- Breaks down barriers between CPU-GPU
- New system architectures

### **Power 8 CAPI – Coherent Accelerator Processor Interface**

### Virtual Addressing

• Accelerator can work with same memory addresses that the processors use

### Hardware Managed Cache Coherence

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

#### Customizable Hardware Application Accelerator

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



Transport for encapsulated messages

#### POWER8



#### Processor Service Layer (PSL)

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

#### **IBM** Confidential



# **OpenCAPI.org**

The OpenCAPI Consortium is an open forum to manage the OpenCAPI specification and ecosystem. OpenCAPI is a not-for-profit organization formed in October 2016 by OpenCAPI Board Members AMD, Google, IBM, Mellanox Technologies and Micron to create an open coherent high performance bus interface based on a new bus standard called Open Coherent Accelerator Processor Interface (OpenCAPI) and grow the ecosystem that utilizes this interface. This initiative is being driven by the emerging accelerated computing and advanced memory/storage solutions that have introduced significant system bottlenecks in today's current open bus protocols and requires a

technical solution that is openly available



Press Releases (1)

Archives

» November 2016 » October 2016

### **NVIDIA GPU Roadmap**





### **POWER8 with NVLink: 2.5x Faster CPU-GPU Connection**





GPUs Limited by PCIe Bandwidth From CPU-System Memory

NVLink Enables Fast Unified Memory Access between CPU & GPU Memories

### **Early Performance Results on Minsky**

### Speedup: NVIDIA P100 vs K80 GPU



### **Better Design: Flat and Fat**

Minksy is engineered both flat and fat

- Data flows freely across system
- Nearly as broad from CPU: GPU as System Memory: CPU
- Big pipes between GPUs on the same socket

# Addresses PCI-E Bottleneck for numerous usage models

- Burst at startup/teardown
- Stream data constantly Host-Device
- Constant Transfers between 2 GPUs
- Hidden Bus Transfers from Host-Device (due to insufficient BW)



### **Performance improvement with Power Architecture**

POWER8 with NVLink Platforms: up to a 4X performance uplift on Lattice QCD codes compared to their predecessors

x86 Alternatives: typically delivering 1.5-2.5X performance differentials on the same types of code





Minksy Performance Increase vs 2x Tesla K80 System: MILC/LQCD

## Page Migration Engine & POWER8 with NVLink





Far easier to create new applications on Tesla P100 + Minsky

- NVIDIA Page Migration Engine ensures unified memory space
  - Unified memory: address space spans CPU and GPU, 1TB+
  - Hardware managed transfers: eliminates explicit data transfers
  - Testing program implementing these advantages
- POWER8 with NVLink ensures speedy data throughput
  - 1TB memory space requires faster CPU:GPU data movement
  - Bus masks transfer times

### • Close code-base to parallel CPU code

![](_page_28_Picture_12.jpeg)

### **Application Potential Unlocked By Page Migration Engine and NVLink**

![](_page_29_Figure_1.jpeg)

### **IBM Technical Computing Software Portfolio**

![](_page_30_Figure_1.jpeg)

#### Spectrum Cluster Foundation Dynamic HPC Infrastructure Management

![](_page_30_Figure_3.jpeg)

![](_page_30_Figure_4.jpeg)

![](_page_30_Figure_5.jpeg)

### **Spectrum Scale – An High Performance Parallel File System**

![](_page_31_Figure_1.jpeg)

### **Power GPU Acceleration for HPC Compiler Roadmap**

![](_page_32_Figure_1.jpeg)

### **Summarizing our strategy**

- IBM remains committed to HPC
- We have a long term HPC roadmap already committed to multiple customers
- OpenPower is a broad play for entire HPC market, not just high end, and offers an alternative to the x86 monoculture
- Power outperforms x86 on key HPC apps
- We are actively attracting developers and ISVs to our platform
- We have differentiated solutions with accelerators and networking with CAPI and NVLink
- We have excellent storage solution for HPC (ESS)
- IBM Research is paving the way to exascale through innovation and collaboration