# IBM HPC Strategy & OpenPOWER

# 1st ASTERICS-OBELICS Workshop

12-14 December 2016, Rome, Italy.

H2020-Astronomy ESFRI and Research Infrastructure Cluster
(Grant Agreement number: 653477).
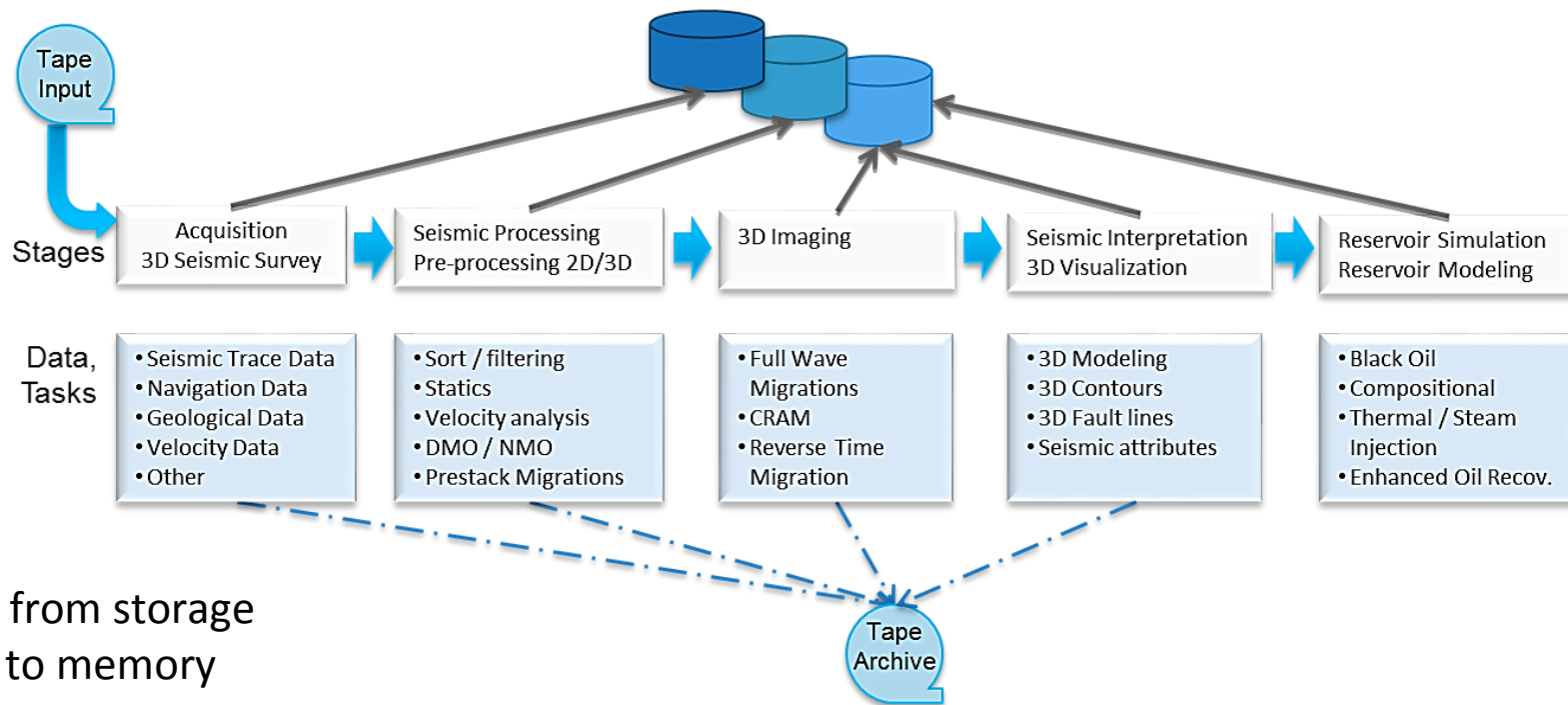
# View of HPC has changed

- **The "Old View"**
    - 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

- **Our View**
    - 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

# Oil and Gas Example



- Read from storage
- Load to memory
- Perform pre-processing
- Execute RTM algorithms
- Visualize and interpret
- Simulate and model

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

# Portfolio of HPC Solutions

**Processors & Systems**

- High Performance Processors & Systems
- Accelerator, networking, storage integration via NVLink & CAPI
- Highest memory throughput

**High Speed Interconnect**

- High speed interconnect / network fabric from Mellanox Technologies
- MPI acceleration in the IB fabric, reducing CPU overhead
- Support for GPUDirect, NVMe over fabric

**High Performance File System & Storage**

- Highest Performance HPC Storage: Elastic Storage Server
- High Performance Spectrum Scale (GPFS) Parallel File System
- Data centric design

**HPC Software**

- Deployment tools, integrated management
- Compilers: gcc, IBM XLC, LLVM OpenMP4, PGI Fortran/C/C++, Java, OpenACC, OpenMP
- Debuggers, Profilers, Math libraries, MPI & HPC apps

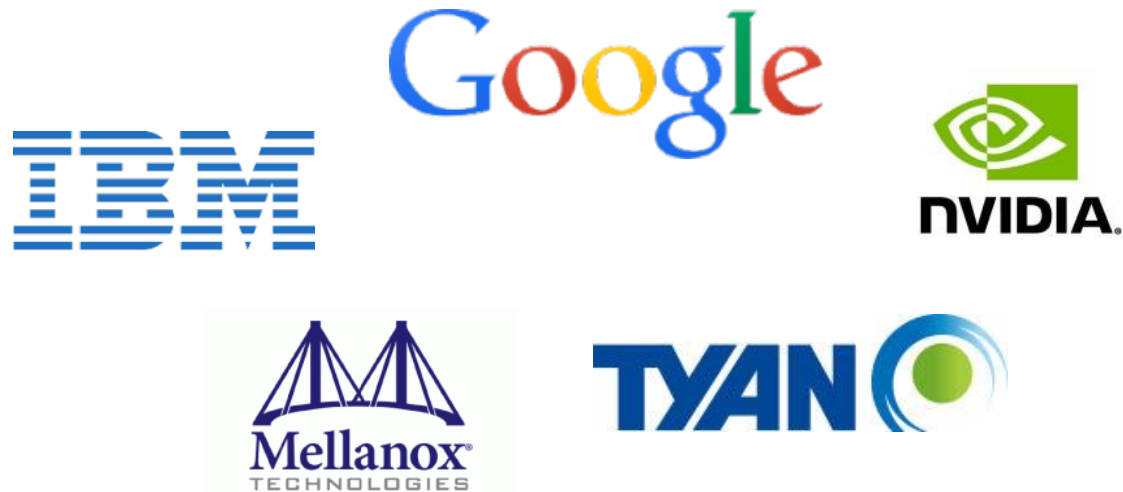# OpenPOWER: Open Architecture for HPC & Analytics

**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

**OpenPOWER™**

**Implementation / HPC / Research**

**Software**

**System / Integration**

**I/O / Storage / Acceleration**

**Boards / Systems**

**Chip / SOC**

# 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 | $20k $5k if <300 employees | 0 | May be on TSC if Work group lead | Sliver members may elect one board representative for all silver members |
| 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 |

New

[www.openpowerfoundation.org](www.openpowerfoundation.org)

# 2300+ Linux Applications on POWER

## HPC

| | |
|---|---|
| CHARMM | miniDFT |
| GROMACS | CTH |
| NAMD | BLAST |
| AMBER | Bowtie |
| RTM | BWA |
| GAMESS | FASTA |
| WRF | HMMER |
| HYCOM | GATK |
| HOMME | SOAP3 |
| LES | STAC-A2 |
| MiniGhost | SHOC |
| AMG2013 | Graph500 |
| OpenFOAM | Ilog |

## Cloud



## Big Data & Machine Learning



## Mobile Enterprise



**Major Linux Distros**

# IBM Power Systems LC Line

*High Performance Computing*

*Big Data*

*Compute Intensive*

**S822LC For Big Data**

**S822LC For High Performance Computing**

**S821LC**

**S812LC**

**S822LC**

1 socket, 2U

Storage rich for big data applications

Memory Intensive workloads

2 socket, 2U

Storage-centric and high through-put workloads

Big data acceleration with work CAPI and GPUs

2 socket, 2U

POWER8 with NVIDIA NVLink

Up to 4 integrated NVIDIA "Pascal" P100 GPUs

2 sockets, 1U

Dense computing

2 socket, 2U

Memory Intensive workloads

# Available now: Barreleye

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

# Google and Rackspace

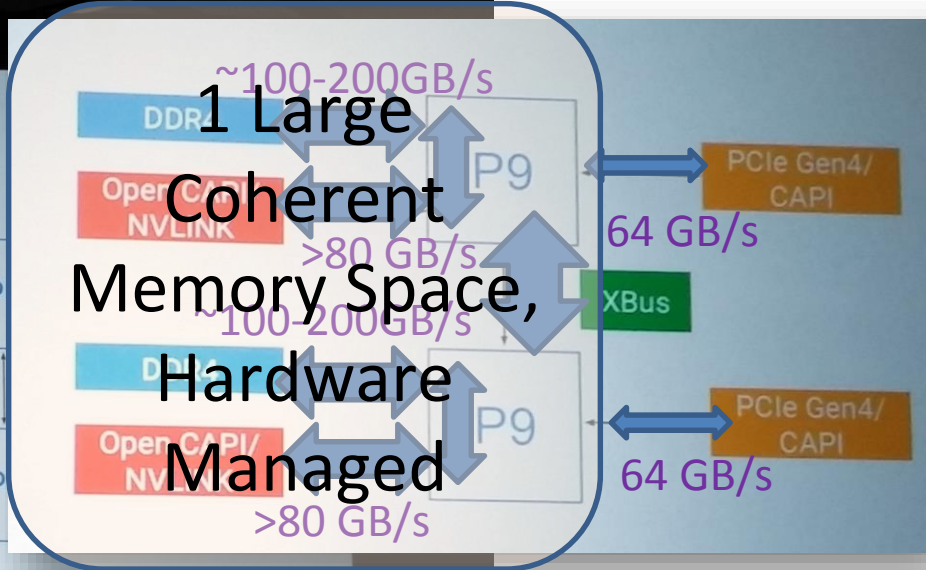# Zaius 1.25 OU

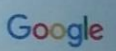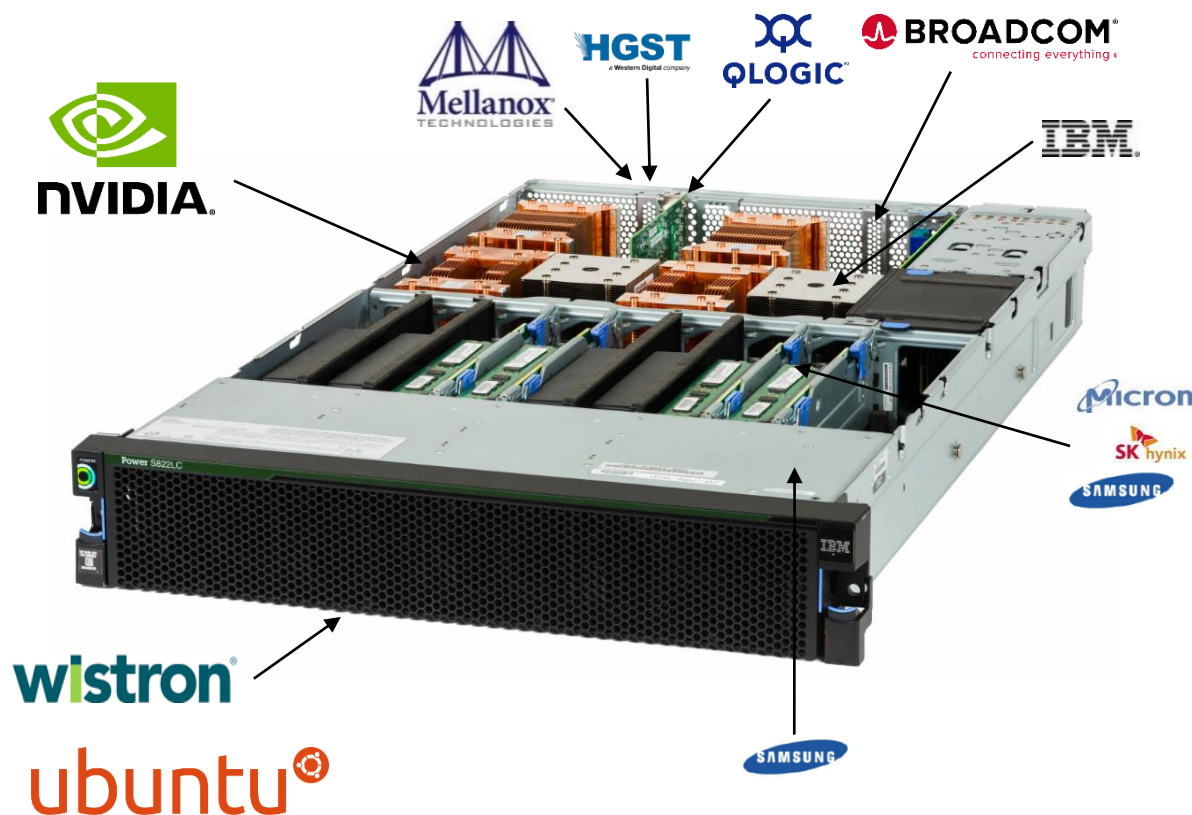- 2 POWER9 CPUS
- 32 DDR4 DIMM SLOTS
- 2X G4 PCIE X16 FHFL SLOTS
- 1X G4 X16 HHHL SLOT
- 1X G4 X16 OCP MEZ
- 1X M.2 SATA PORT
- 1X SATA PORT
- 15X 2.5" SAS/SATA/NVME SLOTS
- BMC W/GBE LOM
- "DISKLESS" OPTION

# Power S822LC HPC



**NVIDIA:**
Tesla P100 GPU Accelerator with NVLink (GPU↔GPU & GPU↔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

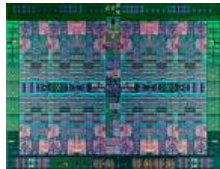**Broadcom:** Optional PCIe Adapters

**QLogic:** Optional Fiber Channel PCIe
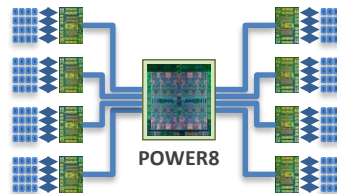
**Samsung:** 2.5" SSDs

**Hynix, Samsung, Micron:** DDR4

**IBM:** POWER8 CPU with NVLink

# IBM Strategy for HPC Systems

**High Performance
Cores**

**Faster Cores
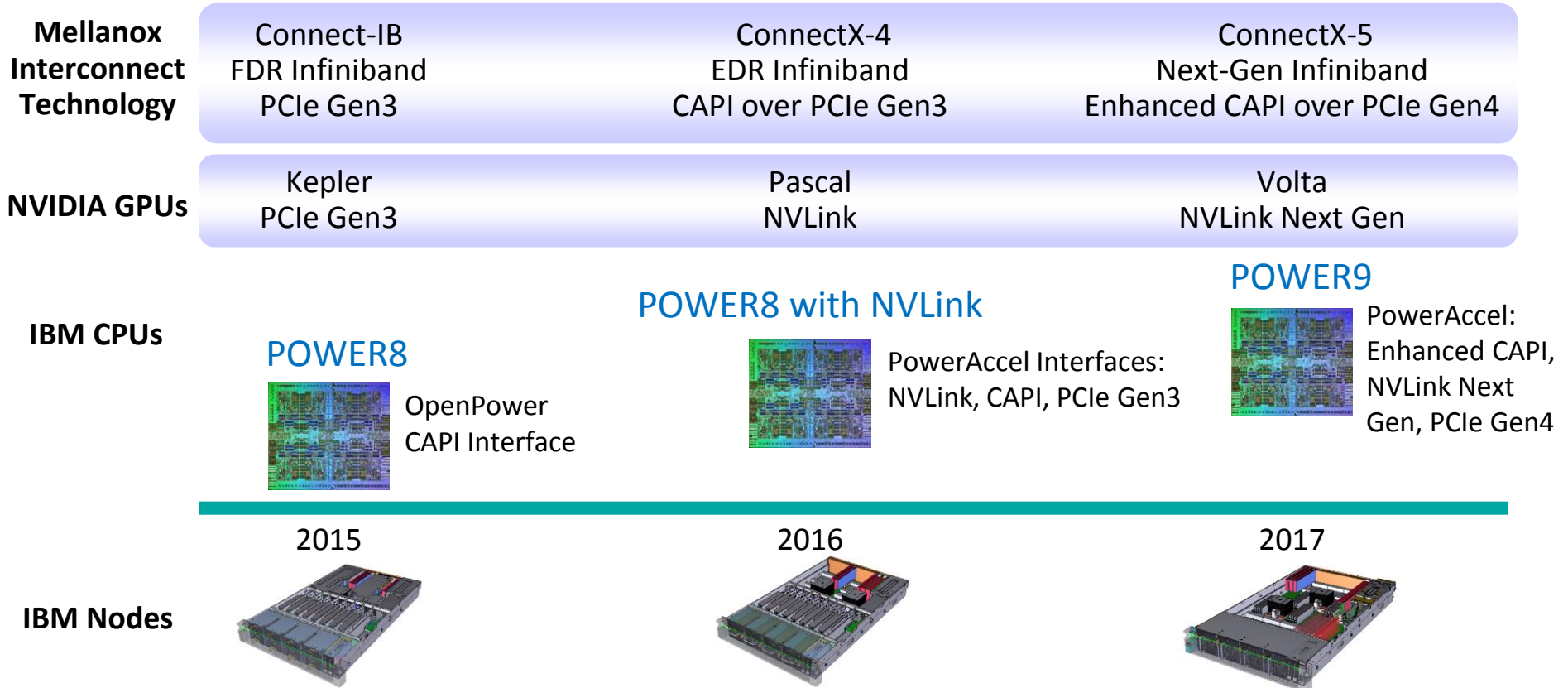than x86**

**Fast & Large
Memory System**

**Larger Caches
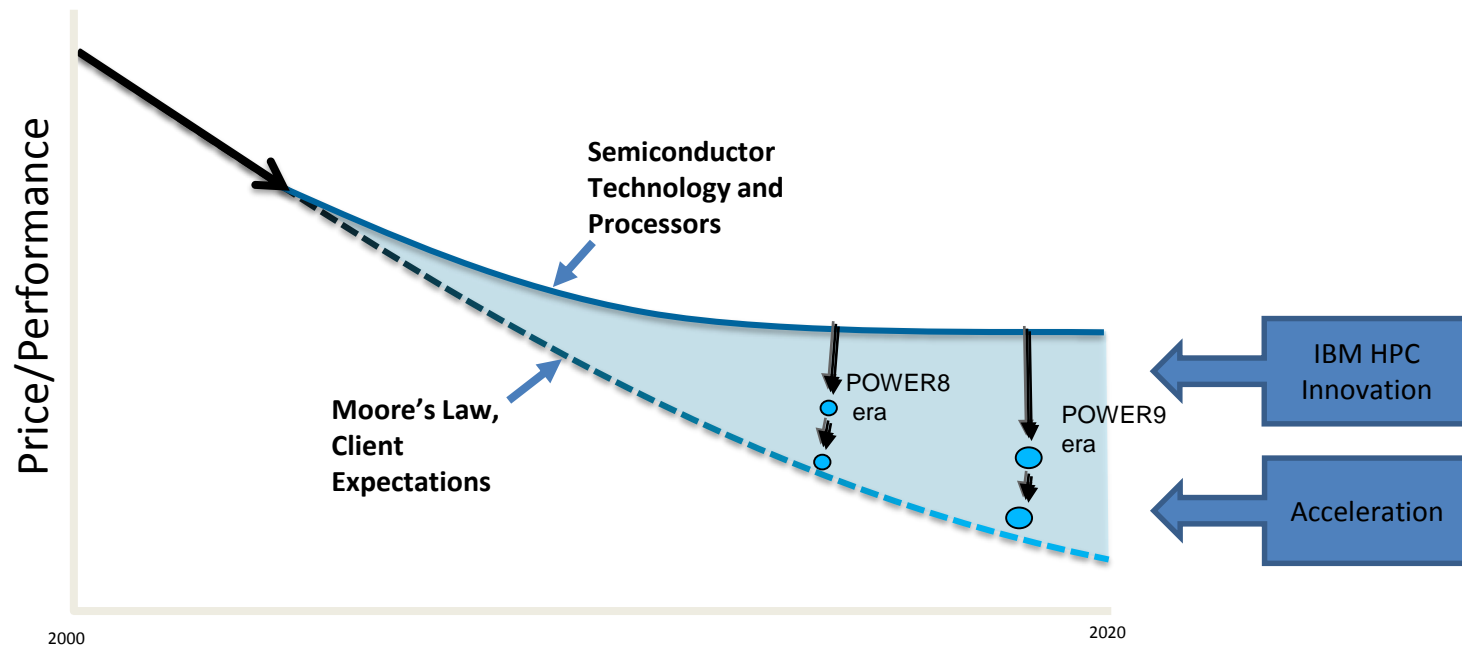Per Core than x86**

**Fast PowerAccel
Interconnects for
Accelerators**

**5x Faster Data
Communication between
POWER8 & GPUs**

# Roadmap for HPC / HPDA

| | | | |
|---|---|---|---|
| **Mellanox Interconnect Technology** | Connect-IB<br>FDR Infiniband<br>PCIe Gen3 | ConnectX-4<br>EDR Infiniband<br>CAPI over PCIe Gen3 | ConnectX-5<br>Next-Gen Infiniband<br>Enhanced CAPI over PCIe Gen4 |
| **NVIDIA GPUs** | Kepler<br>PCIe Gen3 | Pascal<br>NVLink | Volta<br>NVLink Next Gen |

**IBM CPUs**

POWER8 with NVLink

POWER9

POWER8

OpenPower
CAPI Interface

PowerAccel Interfaces:
NVLink, CAPI, PCIe Gen3

PowerAccel:
Enhanced CAPI,
NVLink Next
Gen, PCIe Gen4

| 2015 | 2016 | 2017 |
|---|---|---|

**IBM Nodes**

# 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,
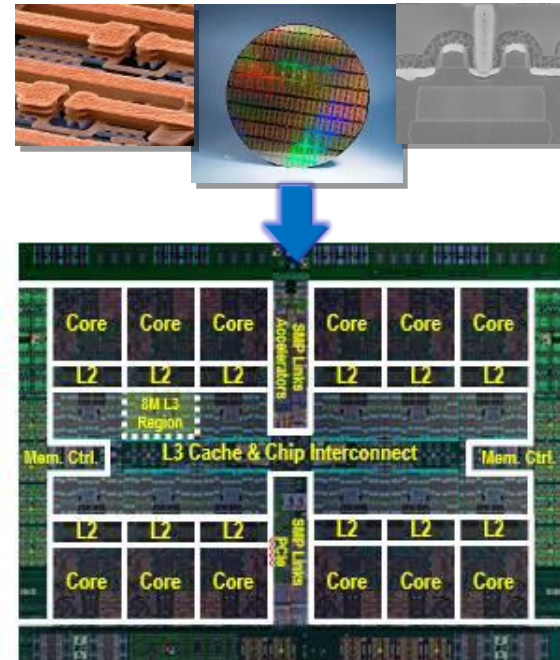- ~4.2 billion Transistors
- Deep trench eDRAM

## 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
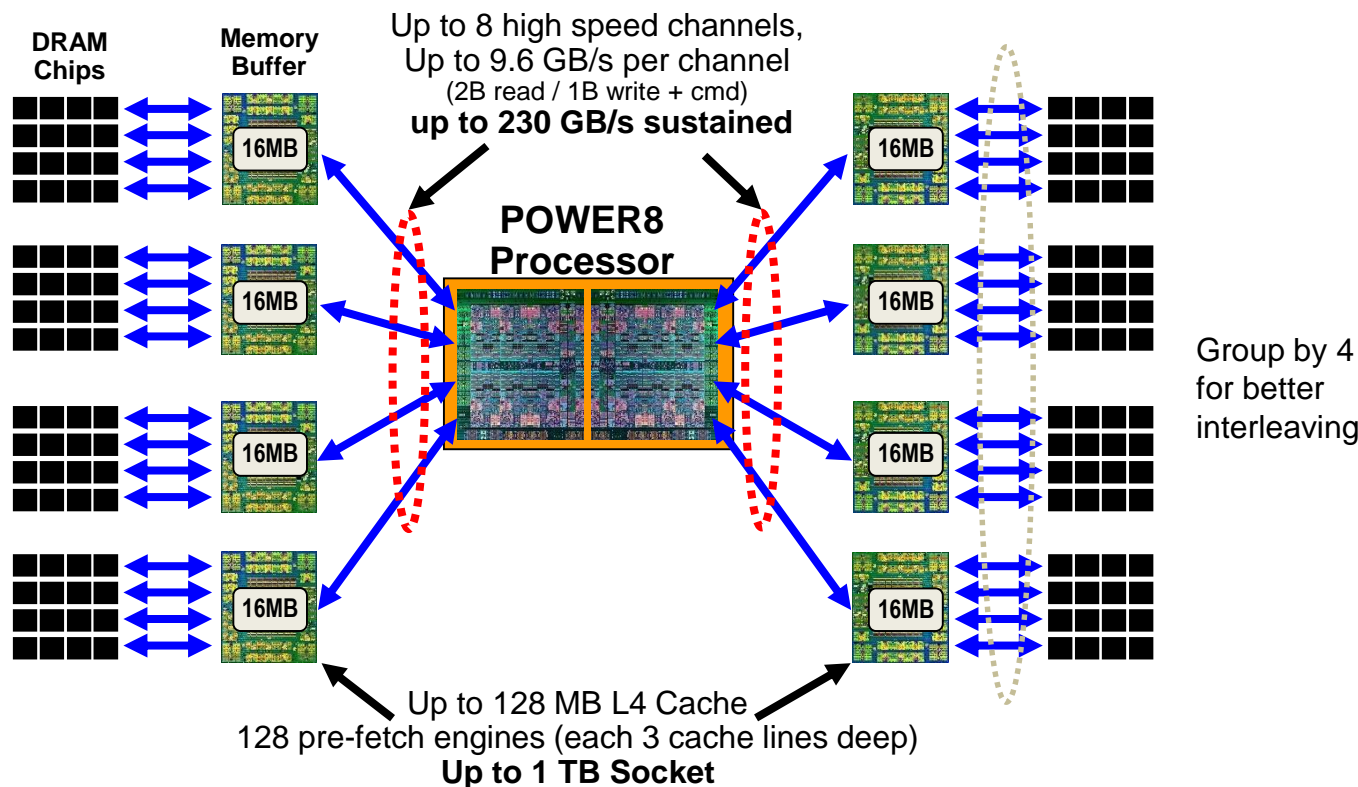
## Cache

- 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)
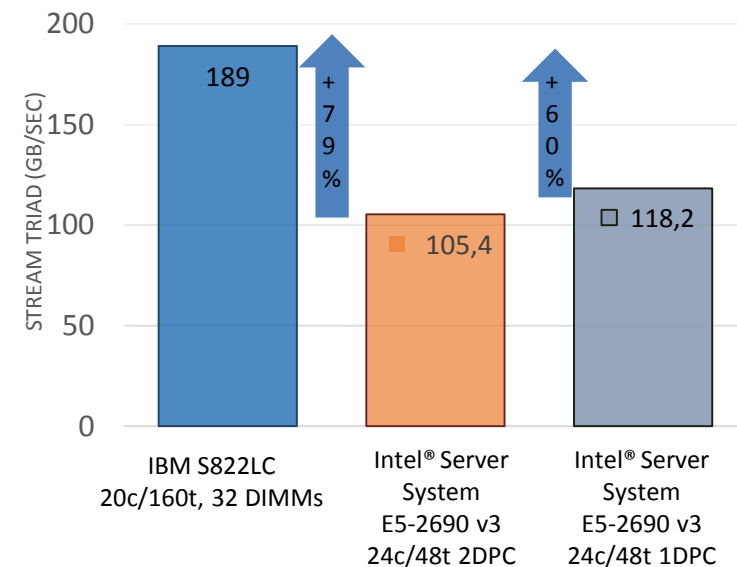- Robust, Large SMP Interconnect
- On chip Energy Mgmt, VRM / core



POWER8

# POWER8 Memory Organization
## (Max Config)

**DRAM Chips**

**Memory Buffer**

Up to 8 high speed channels,
Up to 9.6 GB/s per channel
(2B read / 1B write + cmd)
**up to 230 GB/s sustained**

16MB

**POWER8 Processor**

16MB

16MB

16MB

16MB

16MB

16MB

16MB

Group by 4 for better interleaving

Up to 128 MB L4 Cache
128 pre-fetch engines (each 3 cache lines deep)
**Up to 1 TB Socket**
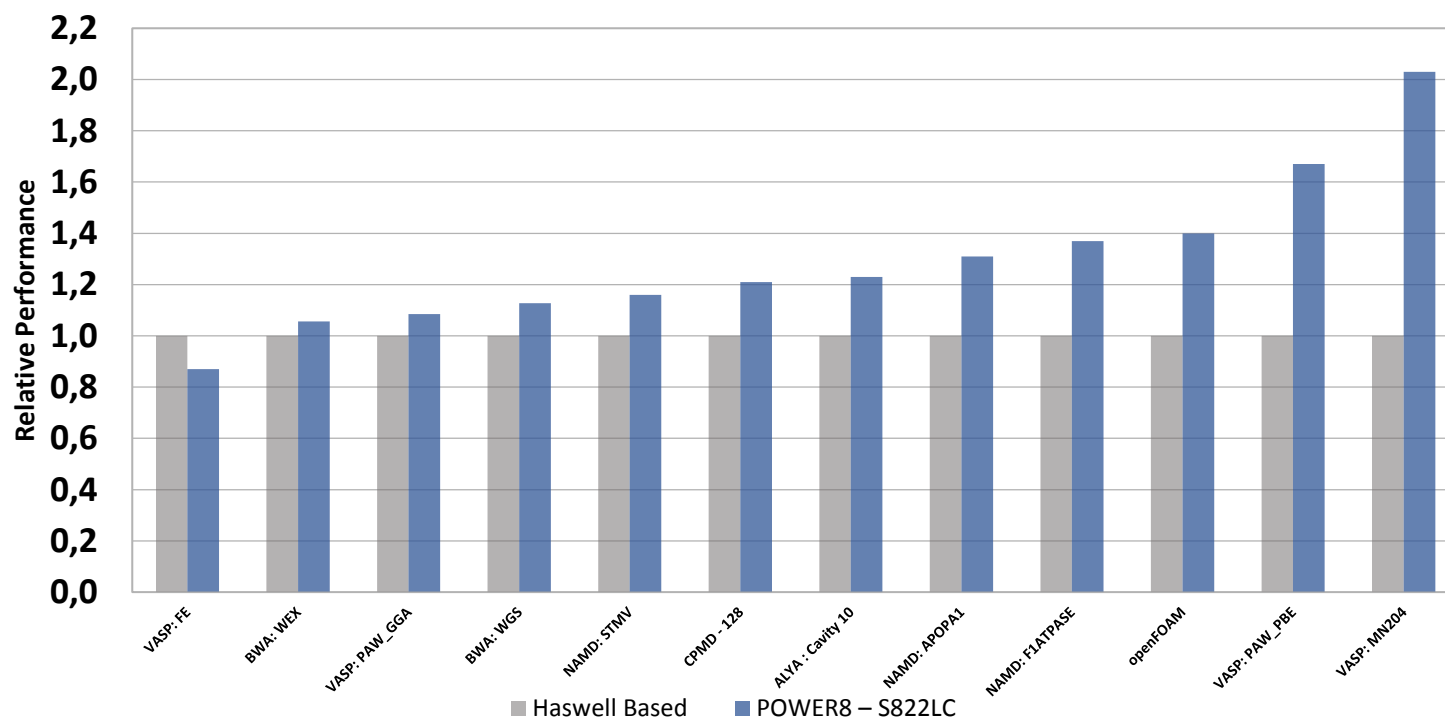
# Memory Bandwidth vs. Xeon E5-2600v3

- 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

Chart: STREAM TRIAD (GB/SEC)
- IBM S822LC 20c/160t, 32 DIMMs: 189
- Intel® Server System E5-2690 v3 24c/48t 2DPC: 105,4 (+79%)
- Intel® Server System E5-2690 v3 24c/48t 1DPC: 118,2 (+60%)

- 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 http://www.intel.com/content/www/us/en/benchmarks/server/xeon-e5-2600-v3/xeon-e5-2600-v3-stream.html
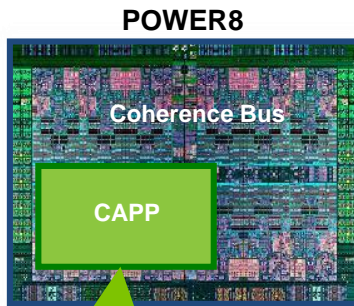
# What Does it Mean?
# Excellent CPU-Only Application Performance
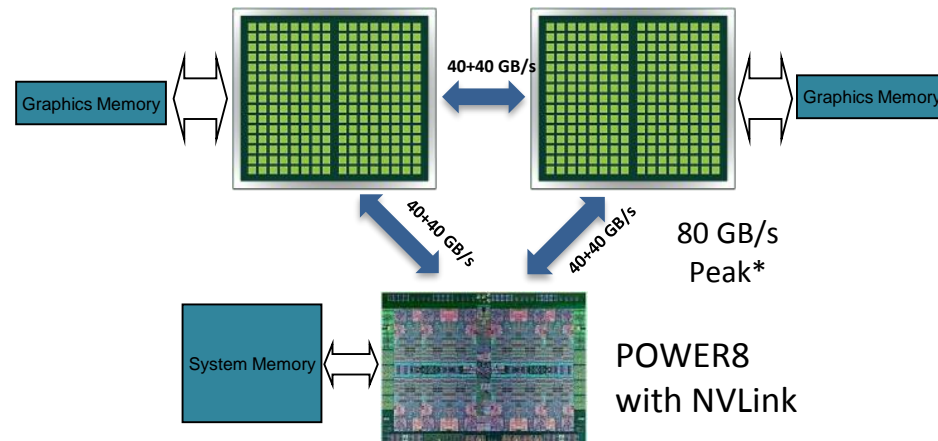
# Differentiated Acceleration
# CAPI and NVLink

## CAPI-attached Accelerators

**POWER8**

**Coherence Bus**

**CAPP**

**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

Graphics Memory

40+40 GB/s

Graphics Memory

40+40 GB/s

40+40 GB/s

80 GB/s Peak*

System Memory

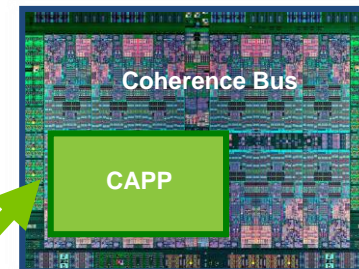POWER8 with NVLink

**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

**POWER8**



**Coherence Bus**

**CAPP**

PCIe Gen 3
Transport for encapsulated messages

**PSL**

**FPGA or ASIC**

**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

# 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

# NVIDIA GPU Roadmap

### Kepler
**CUDA 5.5 – 7.0**
**Unified Memory**

### Pascal
**CUDA 8**
**Full GPU Paging**

### Volta
**CUDA 9**
**Cache Coherent**

1 -2 GPU per board
235W – 300W
versions
1.5TF – 1.9TF std;
2.7TF (boost)
12GB @ 288GB/s
Or 24GB @ 480GB/s

Tesla
K40 – 2014
K80 – 2015

PCIe

POWER8

Buffered
Memory

Pascal
SXM2

4.0+TF std
16GB @ 1TB/s
SXM2 300W

**NVLink 1.0**

POWER8+

Volta
SXM2

7.0+ TF std
16GB @
1.2TB/s
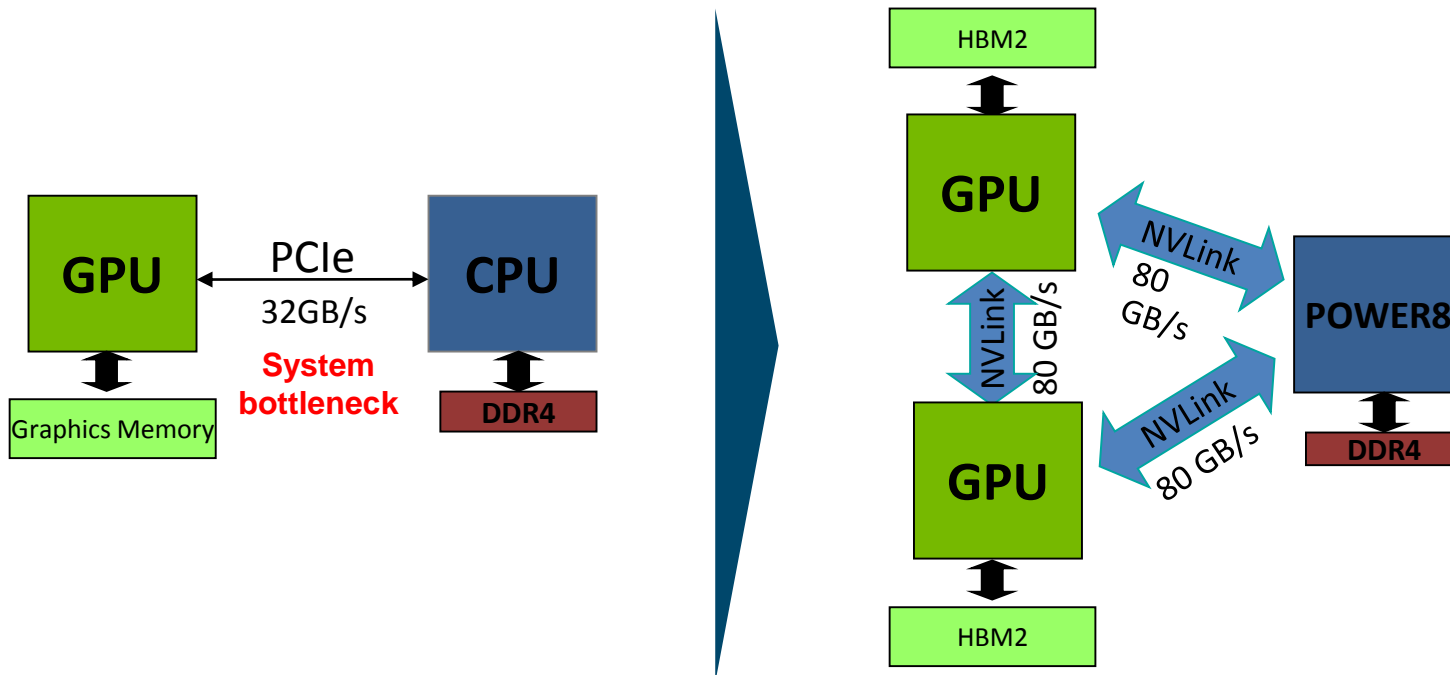SXM2 300W

**NVLink 2.0**

POWER9

Direct
attached

**2014-2015**

**2016**

**2017**

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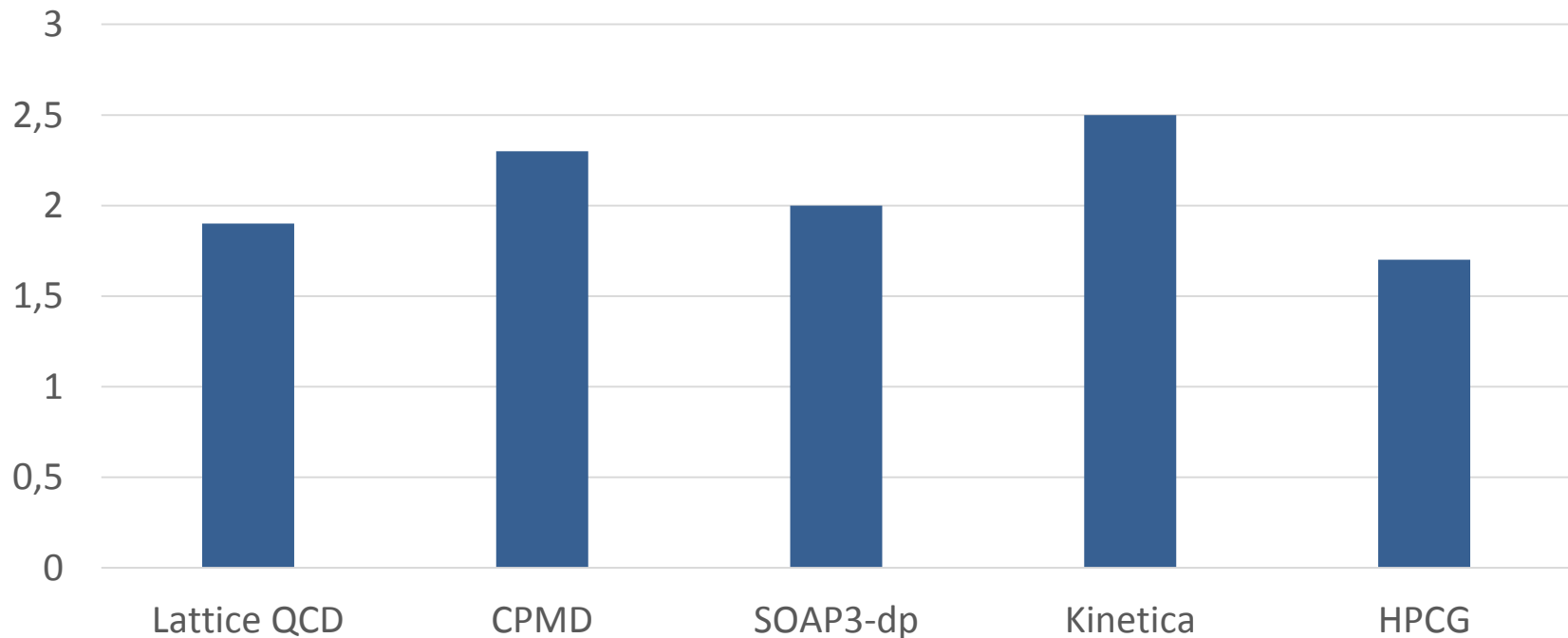# 2.5x Faster CPU-GPU Connection



GPUs Limited by PCIe Bandwidth
From CPU-System Memory

NVLink Enables Fast Unified Memory Access
between CPU & GPU Memories

# Early Performance Results on Minsky
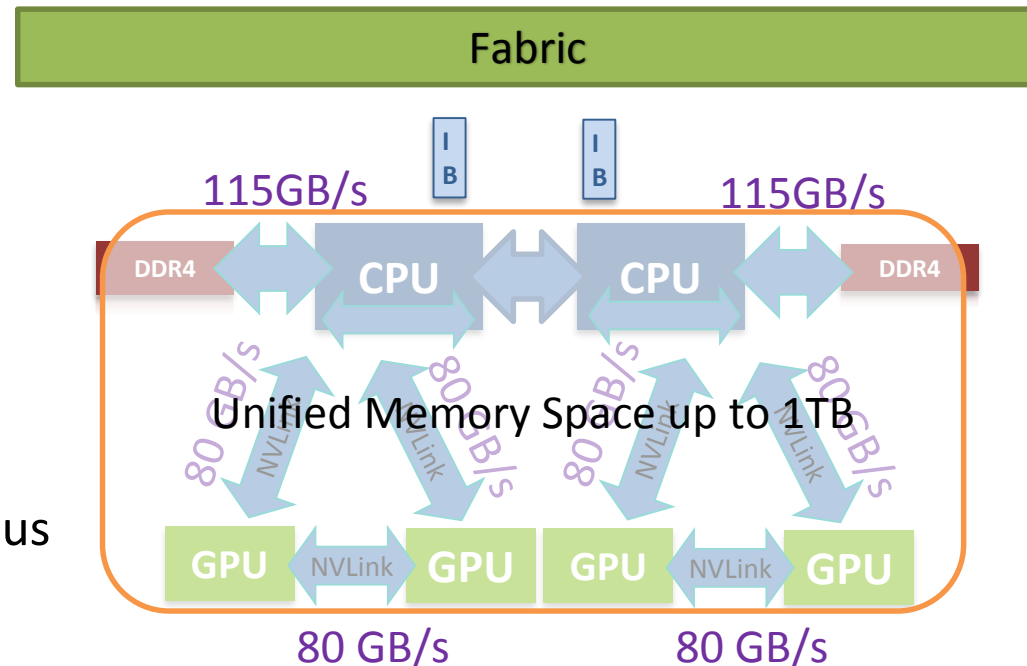
## Speedup: NVIDIA P100 vs K80 GPU

# 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

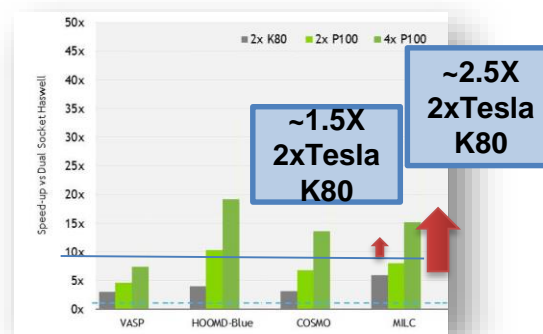Addresses PCI-E Bottleneck for numerous usage models

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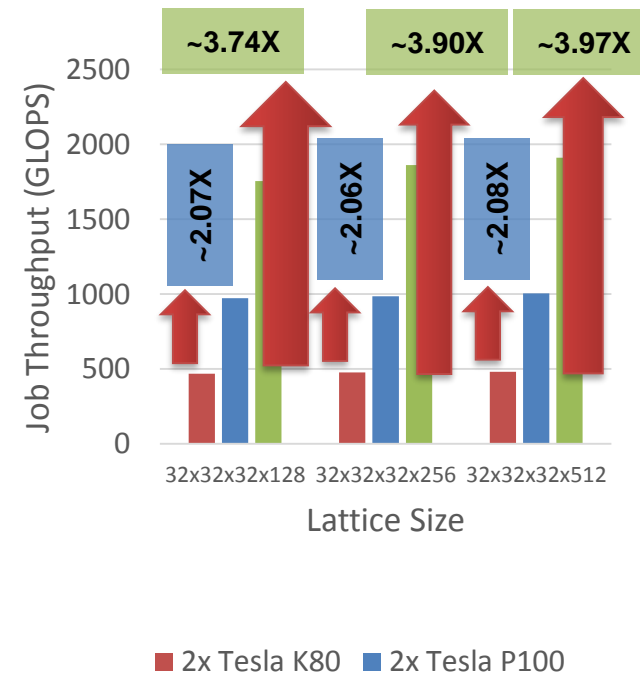- Hidden Bus Transfers from Host-Device (due to insufficient BW)

# 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
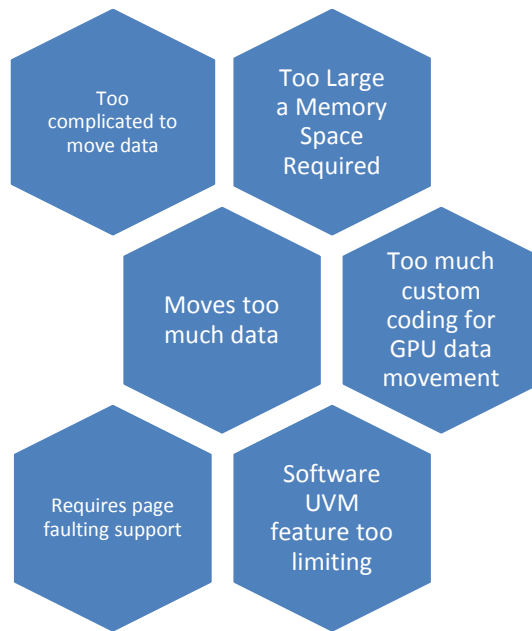


x86 Platform Speedup, vs CPU, 2x Tesla K80



Minksy Performance Increase
vs 2x Tesla K80 System: MILC/LQCD
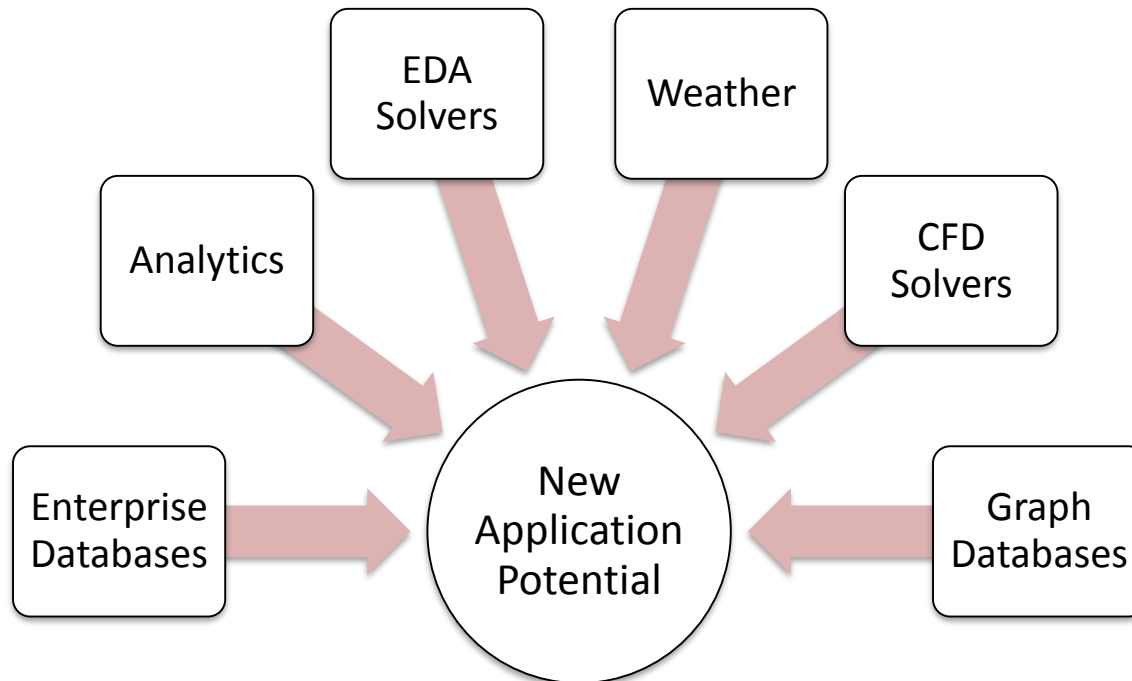
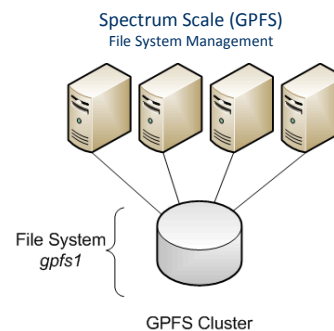# Page Migration Engine & POWER8 with NVLink

Barriers to Entry Removed

Too complicated to move data

Too Large a Memory Space Required

Moves too much data

Too much custom coding for GPU data movement

Requires page faulting support

Software UVM feature too limiting

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

# Application Potential Unlocked By Page Migration Engine and NVLink

# IBM Technical Computing Software Portfolio



**Spectrum LSF**
Workload Management

**Spectrum Symphony**
HPC Grid Services Management

**Parallel Environment (PE)**

**Spectrum Cluster Foundation**
Dynamic HPC Infrastructure Management

**Spectrum Cluster Foundation**
Systems Management and Provisioning

**Spectrum Scale (GPFS)**
File System Management

# Spectrum Scale
# An High Performance Parallel File System

# Power GPU Acceleration for HPC Compiler Roadmap

**CUDA**

**OpenACC**
Directives for Accelerators

**OpenMP**

Power Systems

NVIDIA

Mellanox
TECHNOLOGIES

**Power CUDA C/C++ GA**

**Power CUDA Fortran** Alpha  **CUDA Fortran** Beta

**Power CUDA C/C++/Fortran GA**

**PGI Power Acceleration Enabled Compiler CUDA, OpenACC C/C++/Fortran GA**

**PGI OpenACC C11** Alpha  **PGI OpenACC C++11** Alpha  **PGI OpenACC C/C++/Fortran** Beta

**Open Source OpenACC C/C++/Fortran**

**XL OpenMP 4 C/C++/Fortran Alpha**

**XL OpenMP 4 C/C++/Fortran GA**

**Open Source OpenMP 4 C/C++/Fortran GA**

P8 (4U) Tuleta
2 P8 + 2 GPU
PCIe Gen3

Power 2U
2 P8 + 2 GPU
PCIe Gen3

Power 2U
2 P8 Plus + 2/4 GPU
NVLink 1.0

Power 2U
2 P9 + 2/4 GPU
NVLink 2.0

**P8**  **P8**  **P8'**  **P9**

2015   2016   2017

# 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

**Text for acknowledgement Slide**

# Acknowledgement

- H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).