# Power and ML

**OpenPOWER™**

built with Scott Soutter & Mandie Quartly inputs

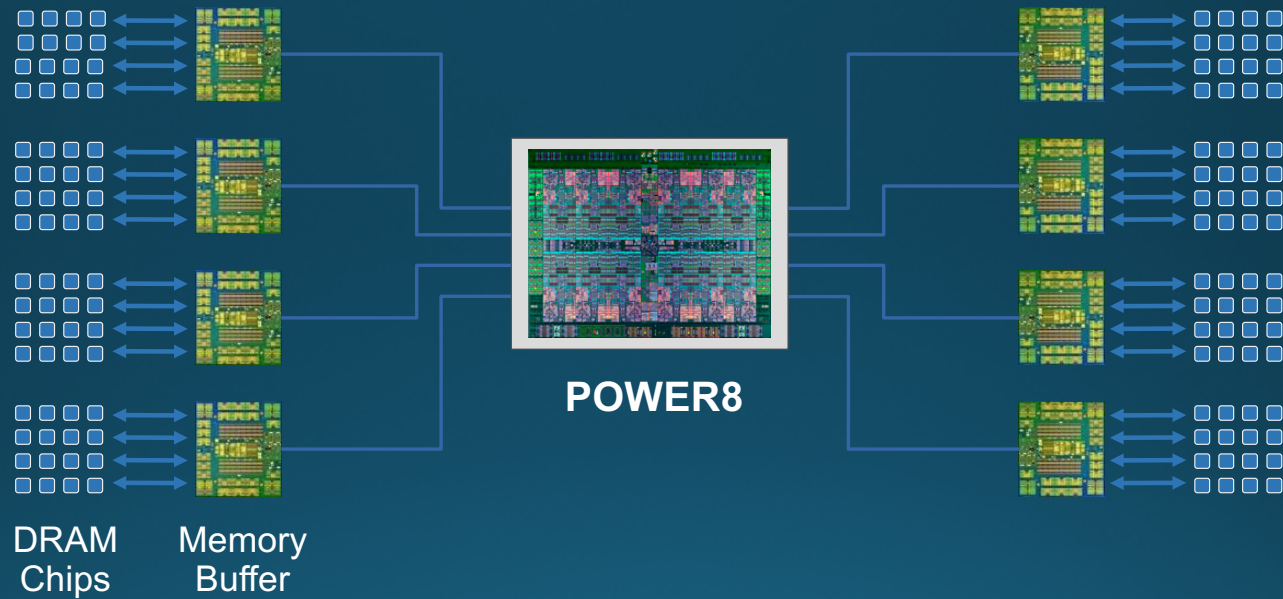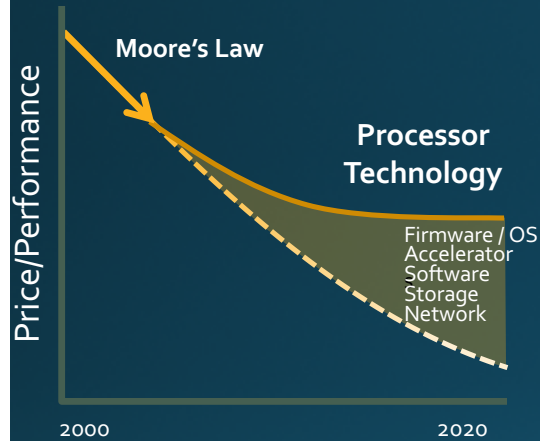Speed to innovation | Standards | Choice
Started in SW first

Leader in UNIX (firmware virtualization...)
Underdog in Linux

Then GOOGLE called…

# Fast Cores, Fast memory IO … and 8 SMT

IBM

Price/Performance

Moore's Law

Processor Technology

Firmware / OS
Accelerator
Software
Storage
Network

2000                    2020

DRAM Chips    Memory Buffer

**POWER8**

**Linux
OpenStack
…**

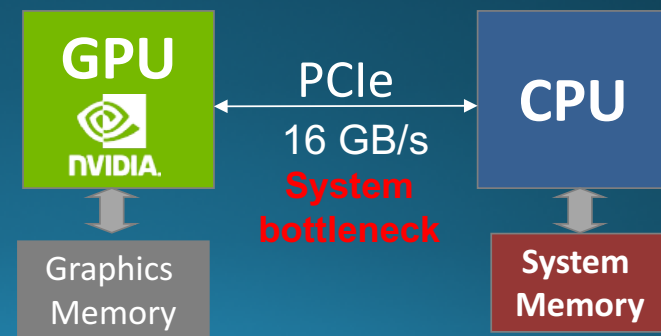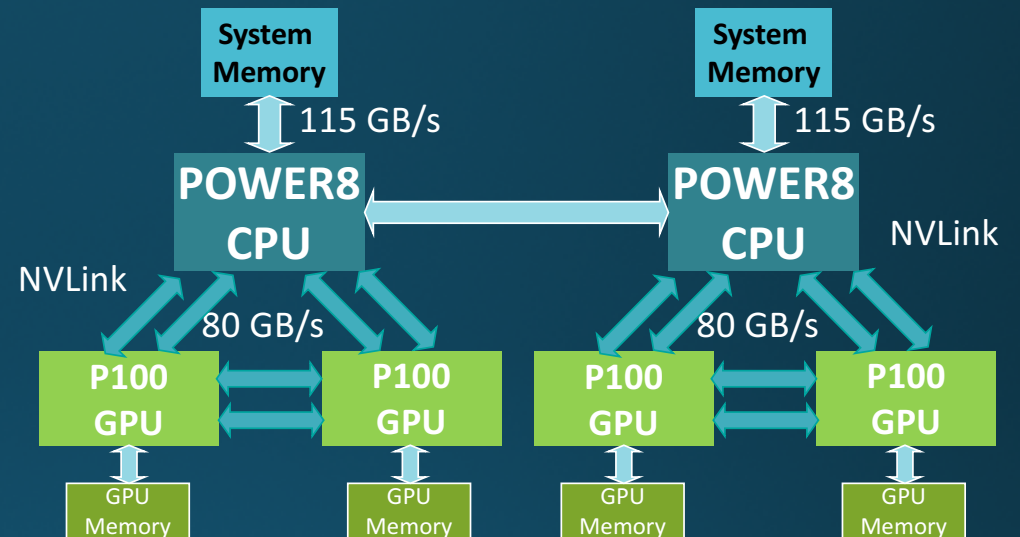**Power8: Up to 12 Cores, up to 96 Threads
L1, L2, L3 + L4 Caches
Up to 1 TB per socket
Up to 230 GB/s sustained memory
bandwidth**

https://www.ibm.com/blogs/systems/power-systems-openpower-enable-acceleration/

# ADD: CPU---P100 GPU NVLink

- Power NVLink between CPUs and GPUs to enable fast memory access to large data sets in system memory

- Two NVLink connections between each GPU and CPU-GPU leads to faster data exchange
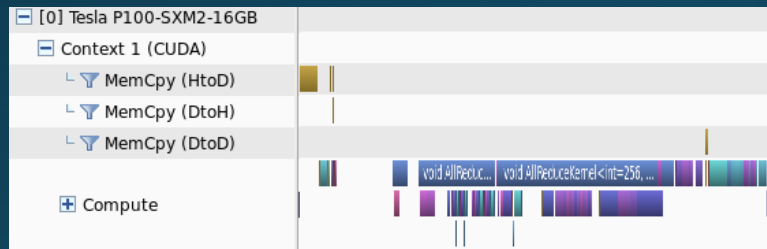
# NVLink and P100 advantage:
## reducing communication time, incorporating the fastest GPU for deep learning

IBM

- NVLink reduces communication time and overhead

- Data gets from GPU-GPU, Memory-GPU faster, for shorter training times
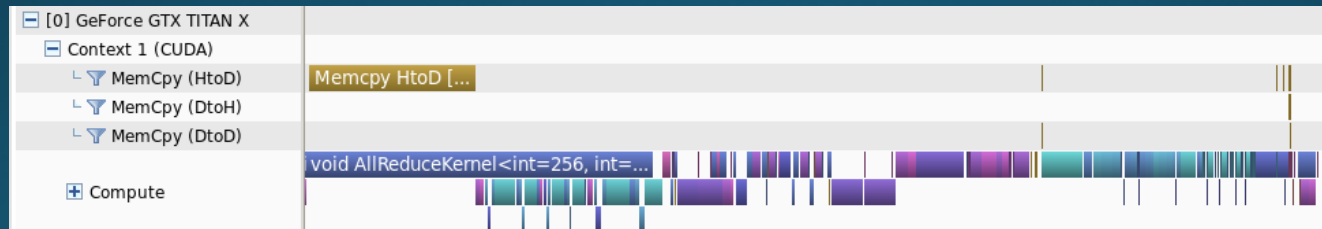
IBM advantage: data communication
and GPU performance

POWER8 +
Tesla
P100+NVLink



78 ms

x86 based
GPU system



170 ms

ImageNet / Alexnet: Minibatch size = 128

# ADD: Coherent Accelerator Processor Interface (CAPI)



**POWER8 Processor**

...FPGAs, networking, memory...

## Typical I/O Model Flow

DD Call → Copy or Pin Source Data → MMIO Notify Accelerator → Acceleration → Poll / Int Completion → Copy or Unpin Result Data → Ret. From DD Completion

## Flow with a Coherent Model

Shared Mem. Notify Accelerator → Acceleration → Shared Memory Completion

# POWER/OpenPower Processor evolution

| Focus on Enterprise Scale-Up 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 | **2017 P9 SO** 24 cores **14nm** New Micro-Architecture Direct attach memory New Process Technology | **TBD P9 SU** 12 cores **14nm** Enhanced Micro-Architecture Buffered Memory | **T B D** / **2018 - 20 P9 SO 10nm - 7nm** Existing Micro-Architecture Foundry Technology | **2020+** New Micro-Architecture New Technology |

# Community Accelerates innovation

- Over **2,500 Linux ISVs** developing on Power
- 50 IBM Innovation Centers
- Compelling PoCs
- Support for little endian applications

## 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

openstack cloud software · TOMCAT · Apache · Chef
STORIX SOFTWARE · REFLEXIS
DIASOFT the way it should be · OneView Commerce
POWERSENSE · IFS
KANA. · FIS
YUCHENG Technologies Limited · SKY solutions
betasystems · EXIGEN
ProjectWare · ALPHINAT

## Big Data & Machine Learning

hadoop · PostgreSQL
MariaDB · redislabs
SUGARCRM · gpudb
Datameer · Spark
neo4j · EDB ENTERPRISEDB
AsiaINFO
CROSSVIEW · INTERSYSTEMS
ZAO · POLARIS live your dream
Valogix

## Mobile Enterprise

zend · 亿阳信通 BOCO Inter-Telecom
php · CSC
RabbitMQ Open Source Enterprise Messaging · Kingdee 金蝶，企业管理专家
SYBASE | An SAP Company · FIS
NARI 南瑞集团公司 NARI GROUP CORPORATION · FINEOS
ALPHINAT · Comptel
Corent The MultiTenancy Company · TEMENOS The Banking Software Company
KUTIR

## Major Linux Distros

ubuntu · SUSE · redhat · freeBSD · Linux

Available now: Barreleye G1

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

CHASSIS: BARRELEYE G2
WITH CACHE COHERENT GPU

CHASSIS: BARRELEYE G2
WITH ALL PCIe SLOTS EXPOSED

ZAIUS MOTHERBOARD, USED IN BARRELEYE G2

FRONT OF BARRELEYE G2 2OU CHASSIS

http://blog.rackspace.com/first-look-zaius-server-platform-google-rackspace-collaboration

# POWER9 – dual memory subsystems

**Scale Out
Direct Attach Memory**

**Scale Up
Buffered Memory**

## 8 Direct DDR4 Ports

- Up to 120 GB/s of sustained bandwidth
- Low latency access
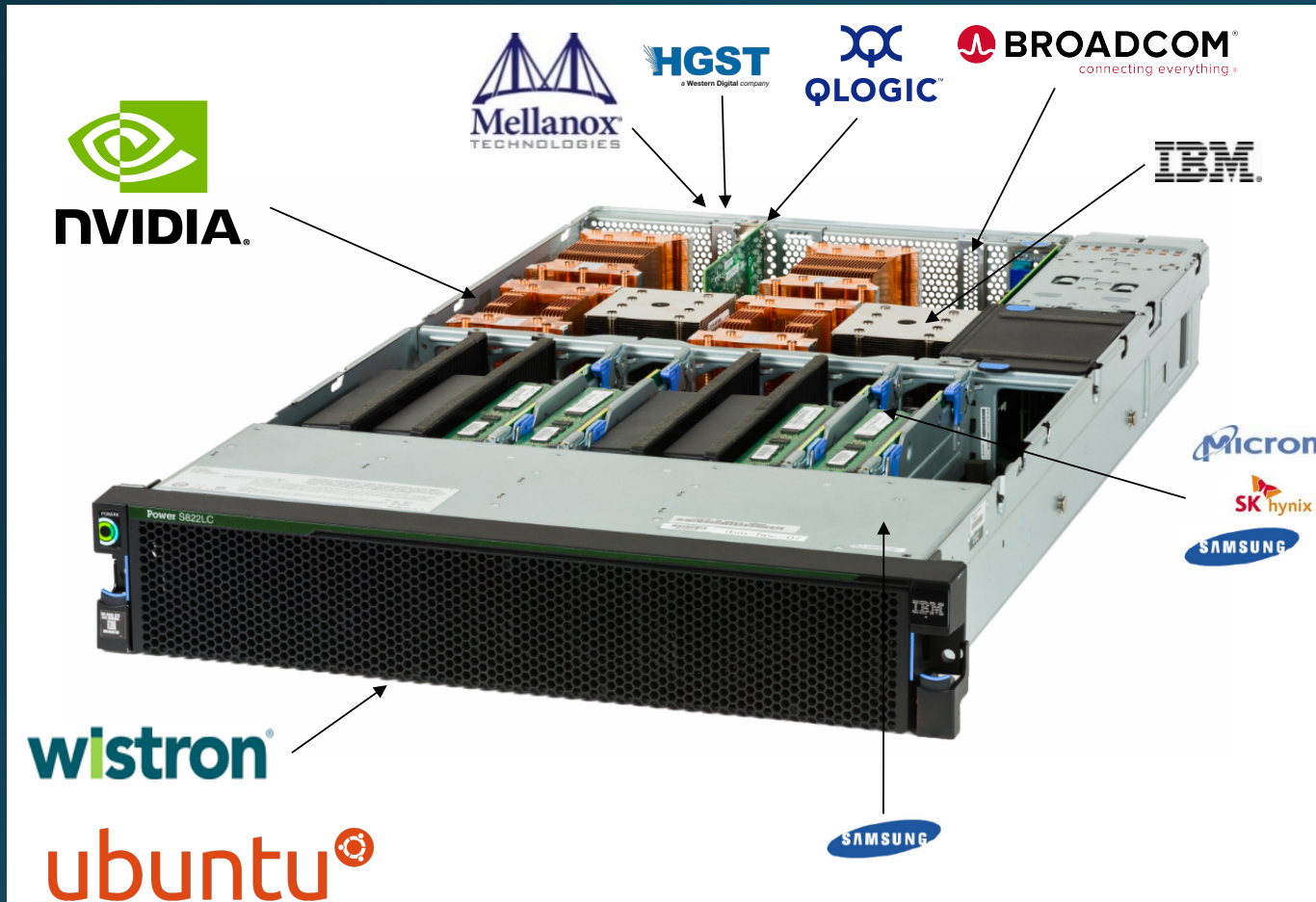- Commodity packaging form factor
- Adaptive 64B / 128B reads

## 8 Buffered Channels

- Up to 230GB/s of sustained bandwidth
- Extreme capacity – up to 8TB / socket
- Superior RAS with chip kill and lane sparing
- Compatible with POWER8 system memory
- Agnostic interface for alternate memory innovations

# Power Systems for High Performance Computing (aka Minsky)

**NVIDIA:**
Tesla P100 GPU Accelerator with NVLink

**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
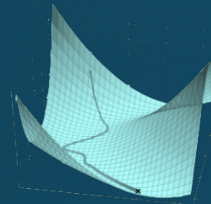
**Samsung:** 2.5" SSDs

**Hynix, Samsung, Micron:** DDR4

**IBM:** POWER8 CPU with NVLink

# PowerAI: Enterprise Deep Learning Distribution

**Enterprise Software Distribution**

Binary Package of Major Deep Learning Frameworks with Enterprise Support

**Tools for Ease of Development**

Graphical tools to Enhance Data Scientist Developer Experience

**Faster Training Times for Data Scientists**

Performance Optimized for Single Node & Distributed Computing Scaling

18

# PowerAI Deep Learning Software Distribution

**IBM**

**Deep Learning Frameworks**

| Caffe | NVCaffe | IBMCaffe | Torch |
|---|---|---|---|
| TensorFlow | Distributed TensorFlow | Theano | Chainer |

**Supporting Libraries**

| OpenBLAS | Bazel | Distributed Communications | NCCL | DIGITS |
|---|---|---|---|---|

**Accelerated Servers and Infrastructure for Scaling**

Cluster of NVLink Servers

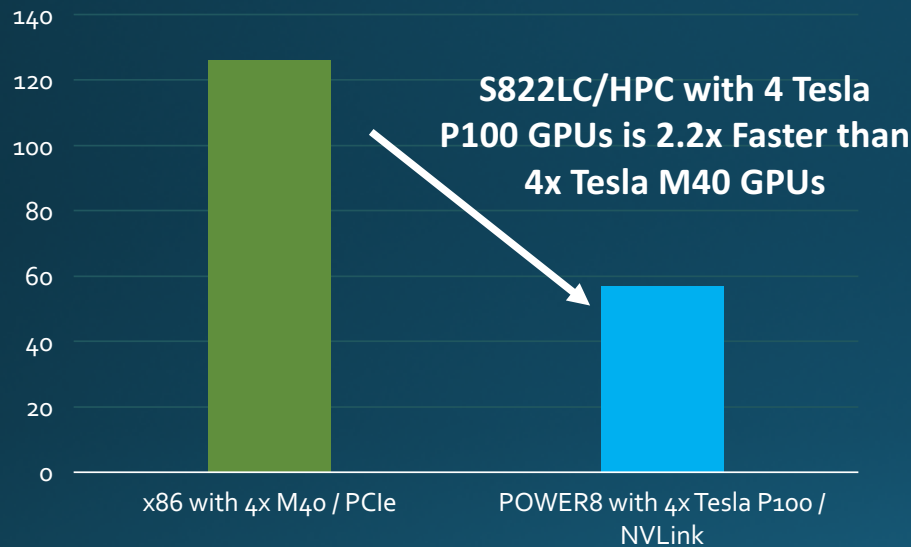Spectrum Scale: High-Speed Parallel File System

Scale to Cloud

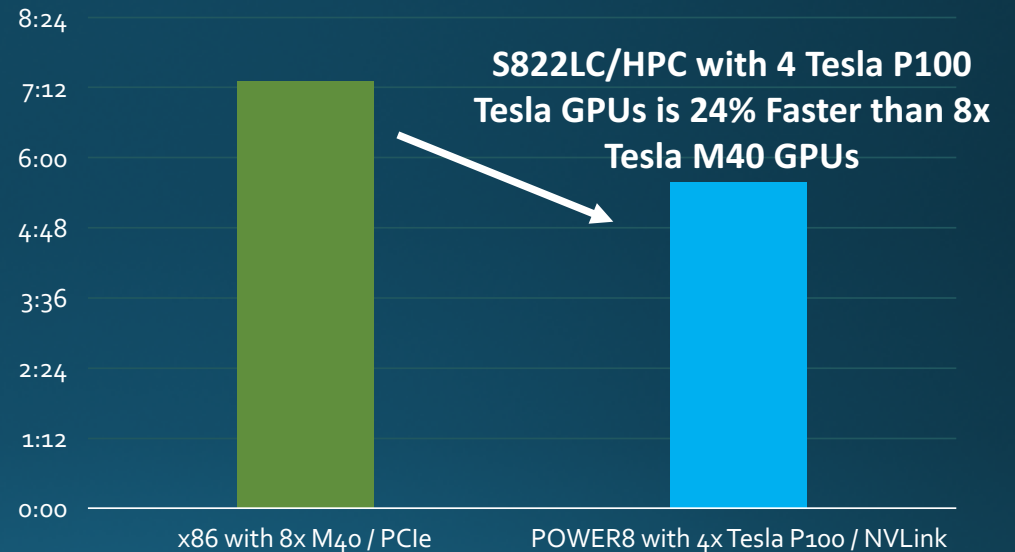# PowerAI Enterprise: Enhancing Developer Experience

IBM

Data Science Experience (DSX)

**PowerAI Enterprise** (Coming)

DL Developer Tools

**PowerAI**

DL Frameworks + Libraries
(TensorFlow, CAFFE, ..)

Distributed Computing
with Spark & MPI

Application Dev Services

IBM Enterprise Support

Cluster of NVLink Servers

Spectrum Scale High-Speed
File System via HDFS APIs

Scale to Cloud

# TensorFlow on Tesla P100: PowerAI is 30% faster
(larger is better)

Images Processed (Images/Sec)
(TensorFlow, Inception v3)

Minsky: 30% Faster

S822LC - Optimized
**Power8 "Minsky" Server**

E5-2640v4
**Intel x86-Based Server**

IBM S822LC 20-cores 2.86GHz 512GB memory / 4 NVIDIA Tesla P100 GPUs / Ubuntu 16.04 /
CUDA 8.0.44 / cuDNN 5.1 / TensorFlow 0.12.0 / Inception v3 Benchmark (64 image minbatch)

Intel Broadwell E5-2640v4 20-core 2.6 GHz 512GB memory / 4 NVIDIA Tesla P100 GPUs/ Ubuntu 16.04 /
CUDA 8.0.44 / cuDNN 5.1 / TensorFlow 0.12.0 / Inception v3 Benchmark (64 image minbatch)

# PowerAI Provides Latest DL Frameworks
## *No need to compile from open-source*

IBM

- Tested, binary builds of common Deep Learning frameworks for ease of implementation

- Simple, complete installation process documented on ibm.biz/powerai

- Future focus on optimizing specific packages for POWER: OpenBLAS, NVIDIA Caffe, TensorFlow, and Torch

| | PowerAI |
|---|---|
| OS | Ubuntu 16.04 |
| CUDA | 8.0 |
| cuDNN | 5.1 |
| Built w/ MASS | Yes |
| OpenBLAS | 0.2.19 |
| Caffe | 1.0 rc5 |
| NVIDIA Caffe | 0.14.5 + 0.15.14 |
| IBM Caffe | 1.0 rc3 |
| Chainer | 1.20.1 |
| NVIDIA DIGITS | 5 |
| Torch | 7 |
| Theano | 0.9 |
| TensorFlow | 1.0.0+ 0.12 |
| GPU | 4 x P100 |
| Base System | S822LC/HPC |

# Getting Started with PowerAI

- Install PowerAI on your existing S822LC for HPC server

  http://ibm.biz/powerai

- Don't have an S822LC for HPC?
    - Reference architecture / system requirements are available for the first system shipping with POWER8, NVLink, and Tesla P100 (next slide)
    - Visit IBM POWER HPC Cloud partners to test drive these frameworks on POWER8/P100 today
        - https://power.jarvice.com/  (Nimbix HPC Cloud)

# You can start small

https://www.hpcwire.com/2017/04/11/dutch-uni-builds-little-green-machine-ii-out-of-ibm-minsky-servers/

# Thank you and may the force be with you