



Partnership for Advanced Computing in Europe (PRACE) Training Event

November 2-4 2016, IDRIS-CNRS

Mastering GPU-Acceleration on OpenPOWER Platform for Optimal Application Performance

Stéphane CHAUVEAU	NVIDIA
François COURTEILLE	NVIDIA
Pascal VEZOLLE	IBM
Nicolas TALLET	IBM

schauveau@nvidia.com
fcourteille@nvidia.com
vezolle@fr.ibm.com
nicolas.tallet@fr.ibm.com

IBM Client Center



OpenPOWER Scientific Collaboration - GENCI & IBM/NVIDIA/Mellanox

▪ Objectives of the Collaboration

- OpenPower Technology Assessment for Multi-Petaflops Deployment in 2017
 - Analyze applications affinity and requirements
 - Port applications to OpenPower platform
 - Analyze applications performance
 - Organize workshops and exchanges with application developers
- Experiment OpenPower w/GPU Prototype Systems in 2015/2016
- Evaluate Programming Models (OpenMP, OpenACC)
- Provide Feedback To Development Teams

Workshop Agenda

Wednesday, November 2

- Introduction
- OpenPOWER Technical Architecture
- GPU Tesla P100 Technical Architecture
- HPC Software Stack
- Ouessant Platform
- Programming Model For GPU-Offload #1: CUDA
- Programming Model For GPU-Offload #2: OpenACC (Part 1)

Thursday, November 3

- Programming Model For GPU-Offload #2: OpenACC (Part 2)
- Programming Model For GPU-Offload #2: OpenMP (Part 1)

Friday, November 4

- Programming Model For GPU-Offload #2: OpenMP (Part 2)
- Conclusion

OpenPOWER Technical Architecture

Hardware Features

IBM OpenPOWER Accelerated Computing Roadmap

Mellanox
Interconnect
Technology

ConnectX-4
EDR Infiniband
PCIe Gen3

ConnectX-4
EDR Infiniband
CAPI over PCIe Gen3

ConnectX-5
HDR Infiniband
Enhanced CAPI over PCIe Gen4

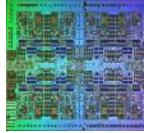
NVIDIA GPUs

Kepler
PCIe Gen3

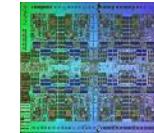
Pascal
NVLink
SXM2

Volta
Enhanced NVLink
SXM2

POWER8

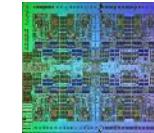


POWER8 w/NVLink



NVLink

POWER9



Enhanced
NVLink

IBM CPUs

2015



Firestone

2016



Minsky

2017



Witherspoon

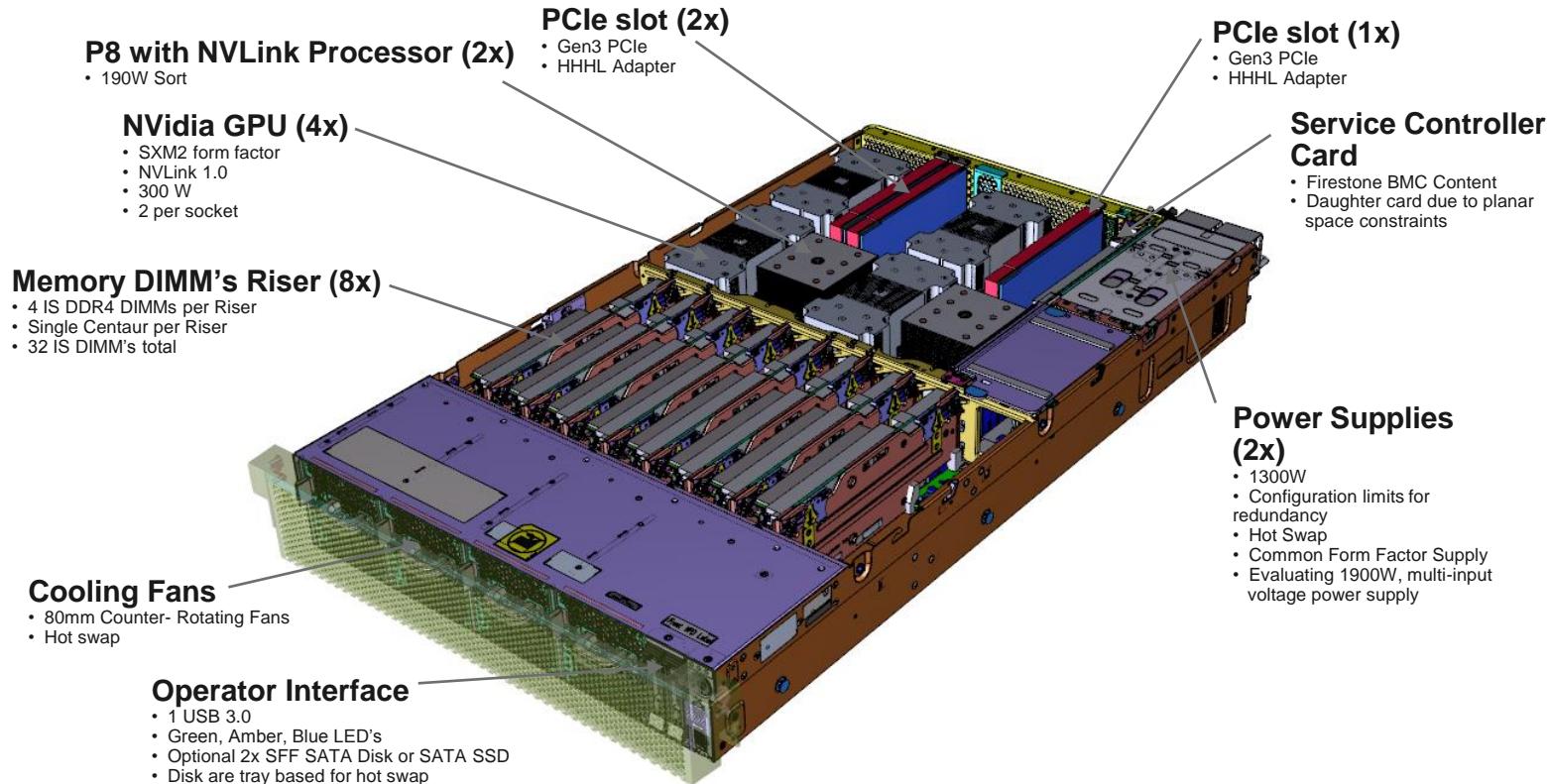
Server

IBM Power Systems S822LC 'Minsky' w/NVIDIA Tesla P100

- 2 POWER8 CPUs
- Up to 1TB DDR4 memory
- Up to 4 Tesla P100 GPUs
- 1st Server with POWER8 with NVLink Technology
- Only architecture with CPU:GPU NVlink

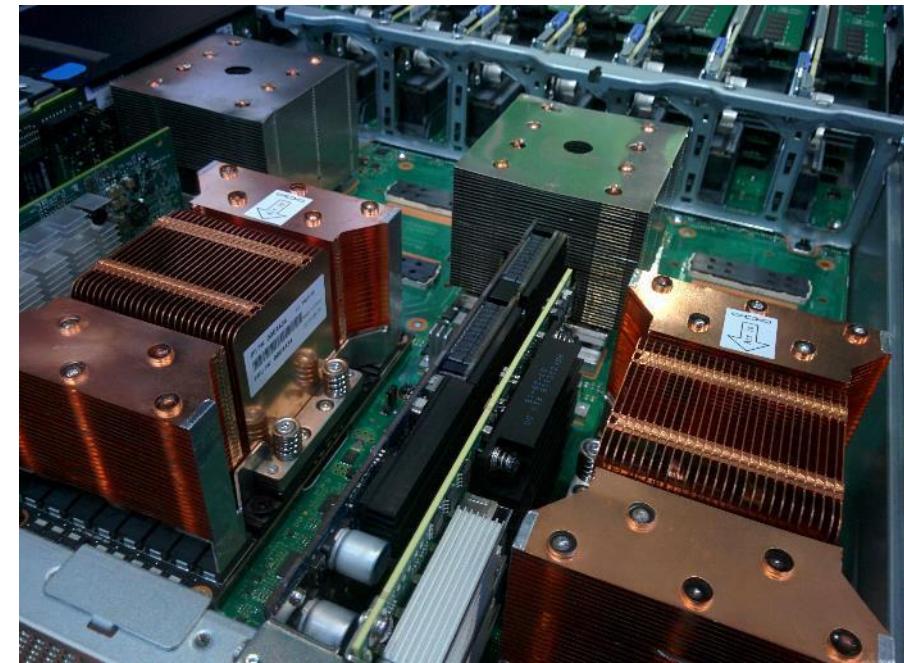


IBM Power System S822LC 'Minsky' (8335-GTB)



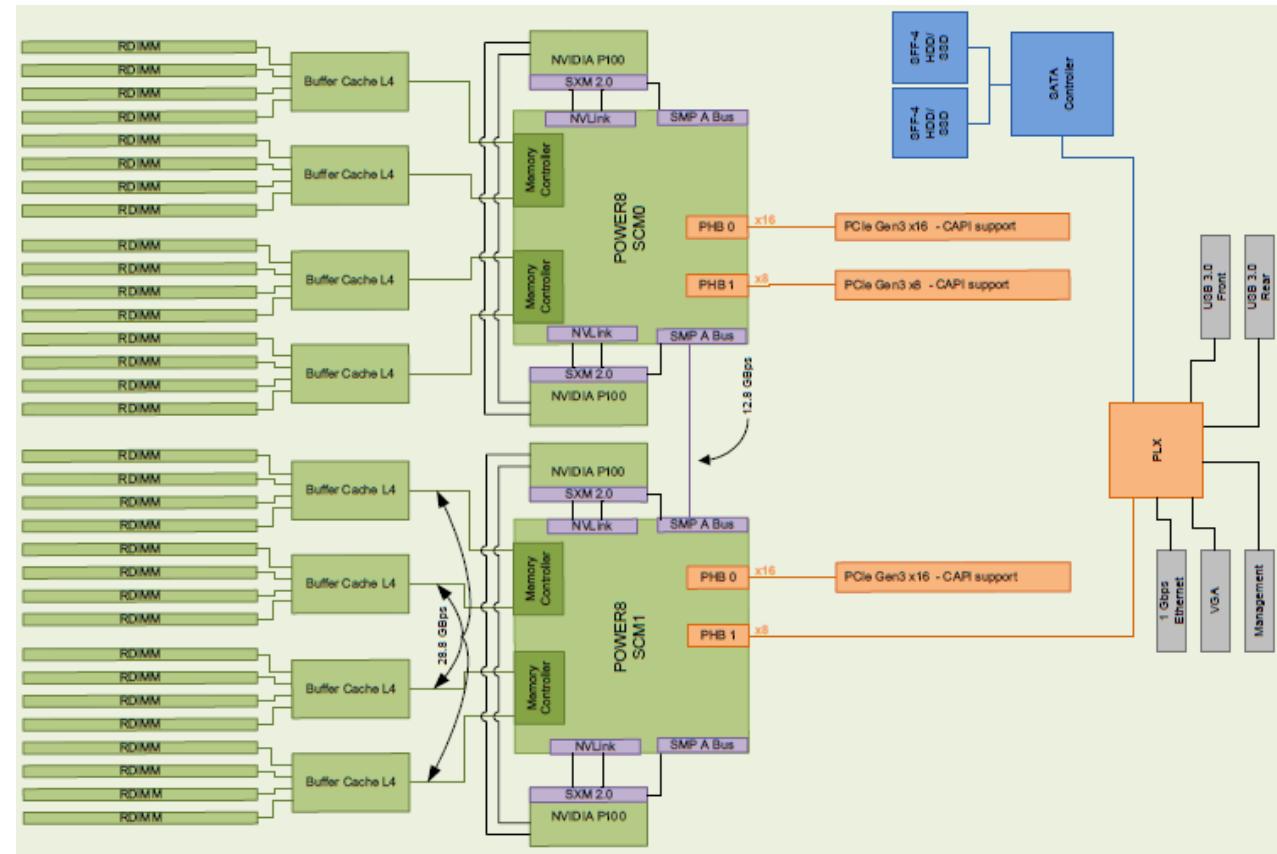
IBM Power System S822LC 'Minsky' (8335-GTB)

Sockets	2 x POWER8
Physical Cores	20
Hardware Threads (Logical Cores)	20 [SMT Off] 40 [SMT 2] 80 [SMT 4] 160 [SMT 8]
CPU Frequency	2.86 GHz [Nominal] 4.03 GHz [Turbo]
Memory Capacity	Up To 1 TB
Memory Bandwidth (Peak)	230 GB/s
DP Performance (Peak)	468 Gflops
GPUs	Up To 4 x NVIDIA Tesla P100
Link	NVLink
Link Bandwidth	10 GB/s
DP Performance (Peak)	4.9 TFlops



Logical System Diagram

- Two Single-Chip Modules (SCM)
- 4 Memory Riser Cards per SCM
 - Buffer Chips for L4 Cache
 - 4 RDIMM Slots
- Max. Capacity: 32 Memory DIMMs (1024 GB)
- 4 GPU Sockets (300 W Max.)
- 3 PCIe Gen3 Slots
- Dedicated PCI Bus:
 - Integrated SATA Controller: Up to 2 SATA Drives
 - Integrated Ethernet
 - Integrated USB Port

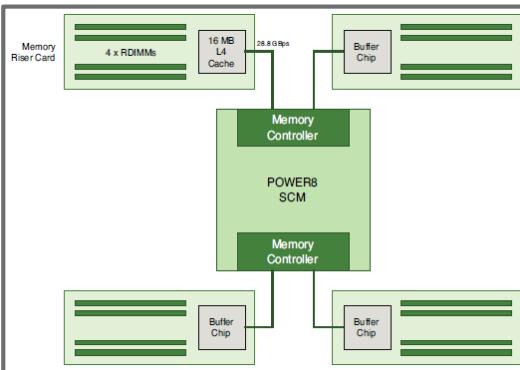
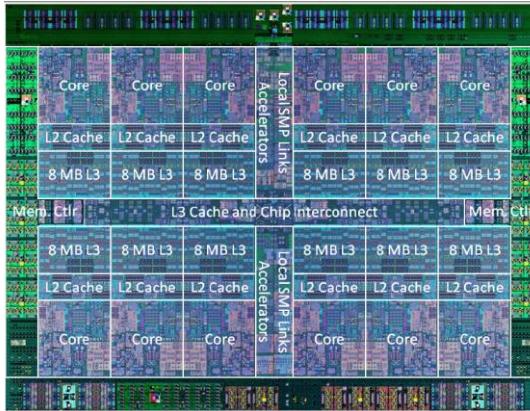


Simultaneous Multi-Threading (SMT)

- Allows a single physical processor core to dispatch simultaneously instructions from more than one hardware thread context
- With SMT, each POWER8 core can present up to 8 hardware threads
- Because there are multiple hardware threads per physical processor core, additional instructions can run at the same time
- SMT benefit highly depends on the workload
- Changing SMT mode does not require reboot

Mode	Logical Cores
Single Thread	20
SMT2	40
SMT4	80
SMT8	160

Memory Sub-System | L3 & L4 Cache



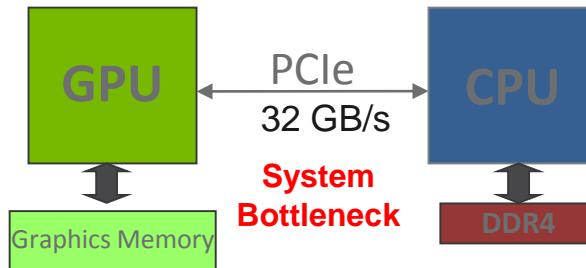
- L3 Cache

- eDRAM on the processor die
- Each processor core is associated with a fast 8-MB local region of L3 cache (FLR-L3)
- But can also access other L3 cache regions as shared L3 cache

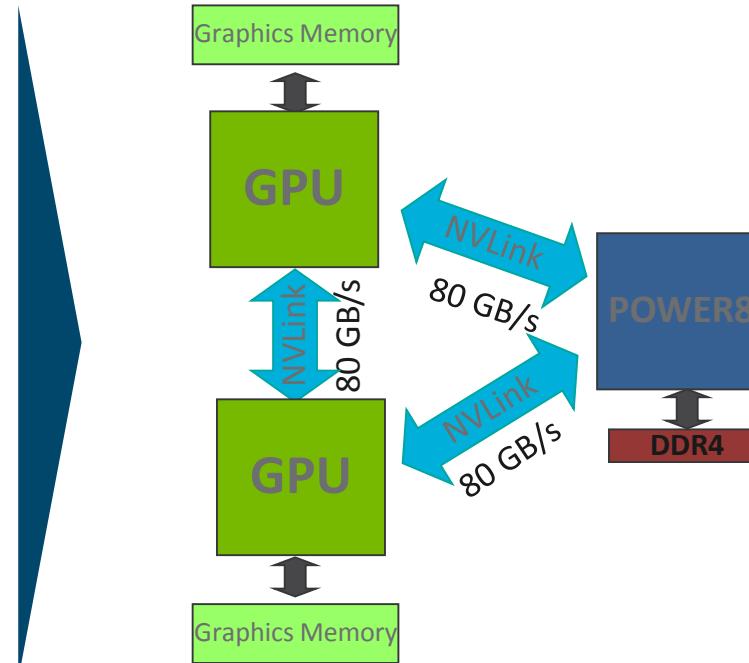
- L4 Cache

- Memory buffer on the memory riser cards
- Each memory buffer contains 16 MB of L4 cache
- Up to 128 MB of L4 cache by using all memory riser cards

POWER8 w/NVLink: 2.5x Faster CPU-To-GPU Connection

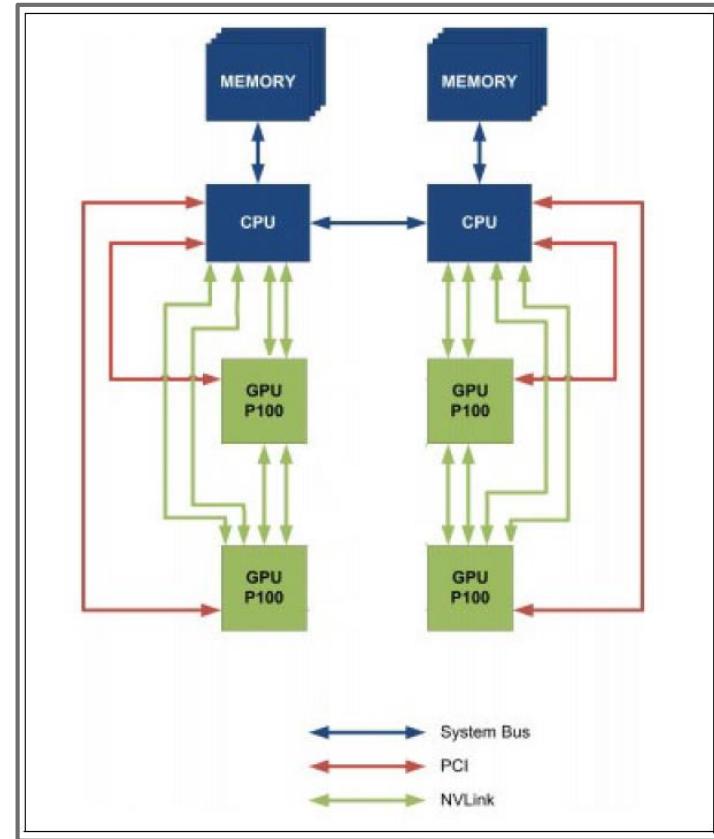
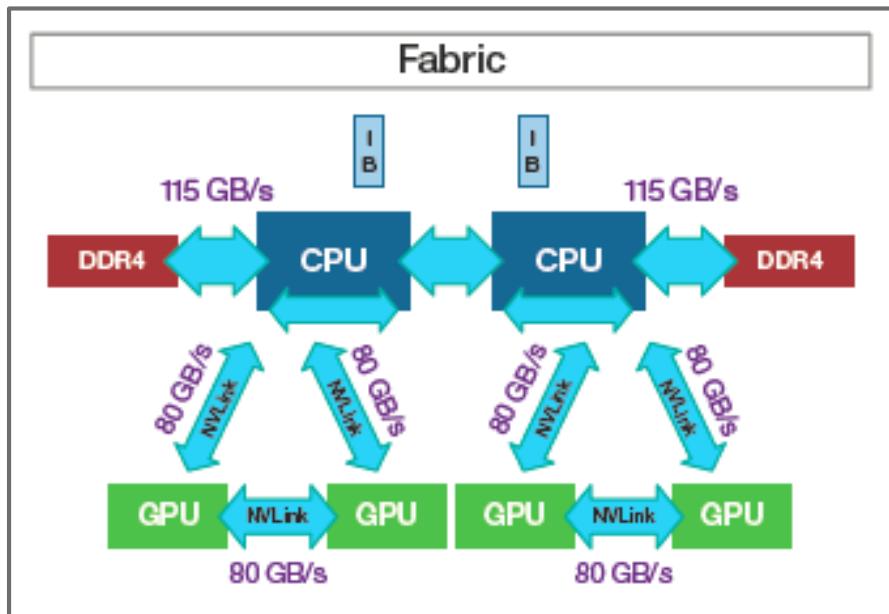


GPUs Limited by PCIe Bandwidth
From CPU-System Memory



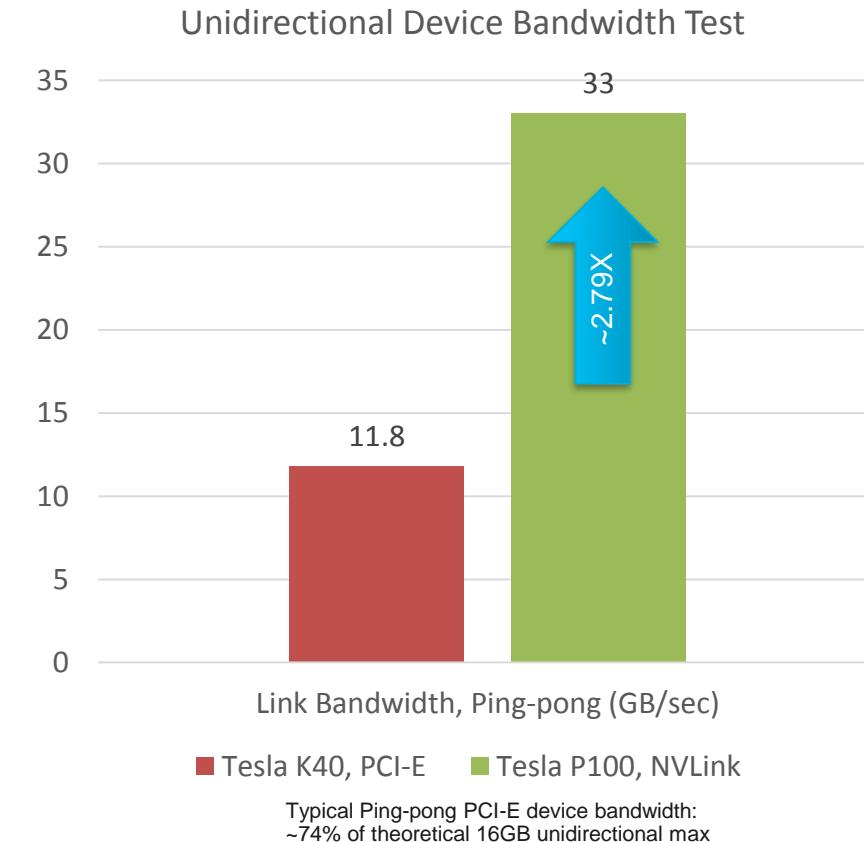
NVLink Enables Fast Unified Memory Access
between CPU & GPU Memories

Comprehensive Fabric Architecture



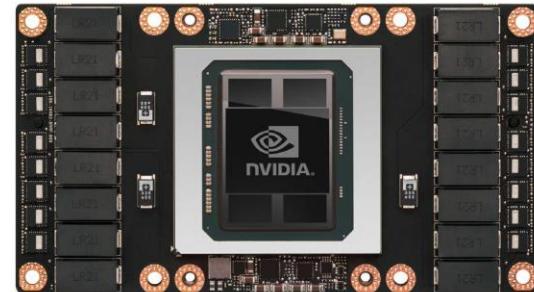
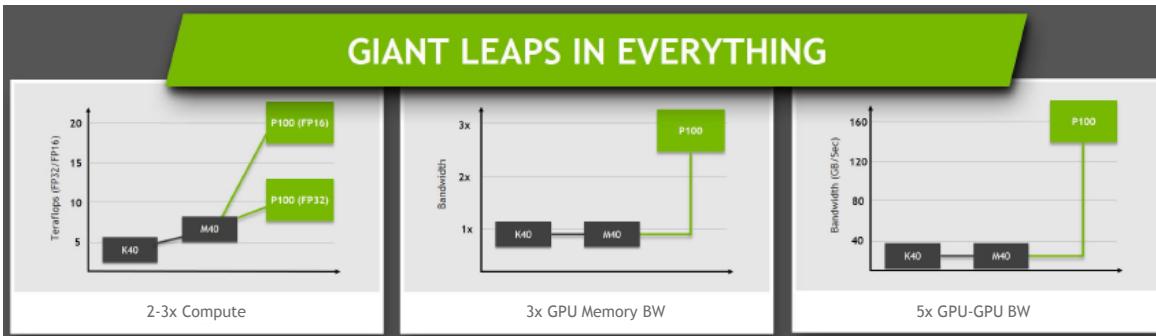
Unique Bandwidth Dividend: NVLink from CPU:GPU

- **POWER8 with NVLink: the only processor with NVLink from CPU to GPU**
- Delivering in excess of 2.5x bandwidth in testing **today**
- NVLink bus delivering higher efficiency than PCI-E links (82.5% vs 74% of peak)
- No code changes to start leveraging bus (CUDA 8.0 and go)
- **Platform for developers needing Bandwidth for the foreseeable future**
- POWER8 with NVLink ships in 2016
- Xeon E5-2600 Series CPUs retain PCI-E x16 3.0 connectivity through 2017

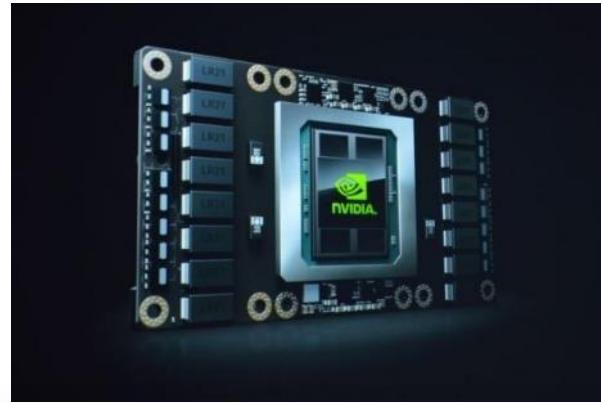


Accelerated Performance: NVIDIA Tesla P100 Pascal GPUs

	Tesla P100	Tesla K80	Tesla K40
DP TFLOPS	5.3 TFLOPS	2.91 TFLOPS (Max Boost)	1.4 TFLOPS
SP TFLOPS	10.6 TFLOPS (21.2 TFLOPS Half Precision)	8.74 TFLOPS (Max Boost)	4.3 TFLOPS
Memory Bandwidth	720 GB/sec	480 GB/sec (2x 240GB)	288 GB/sec
Memory Capacity	16 GB	24 GB (2x 12GB)	12 GB



NVIDIA Tesla P100 Architecture Details



INTRODUCING TESLA P100

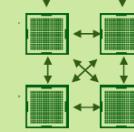
New GPU Architecture to Enable the World's Fastest Compute Node

Pascal Architecture



Highest Compute Performance

NVLink



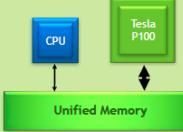
GPU Interconnect for Maximum Scalability

HBM2 Stacked Memory



Unifying Compute & Memory in Single Package

Page Migration Engine

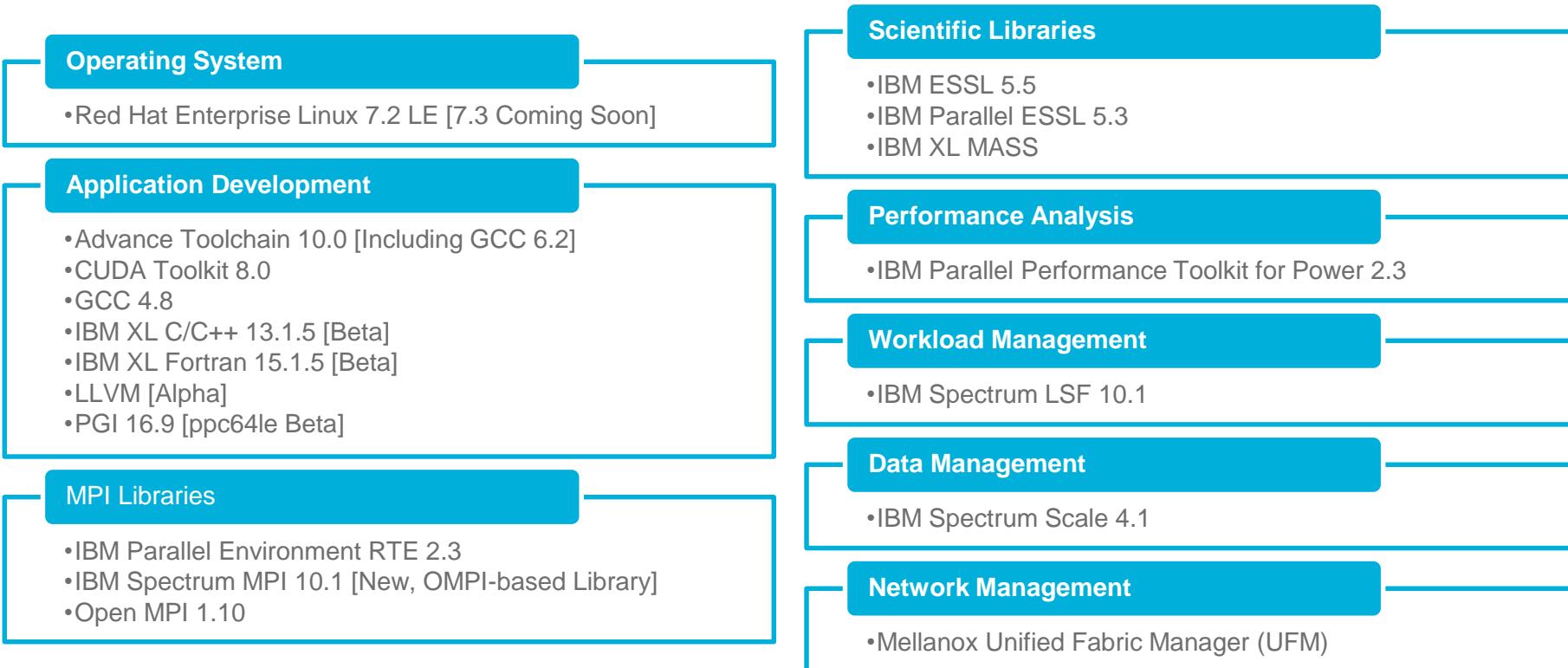


Simple Parallel Programming with 512 TB of Virtual Memory

HPC Software Stack

Software Components

Comprehensive HPC Software Stack



Engineering and Scientific Subroutine Library (ESSL)

ESSL

- Serial & SMP highly-tuned mathematical subroutines
- **30+ SMP CUDA BLAS3 subroutines**
 - POWER8 SMP / GPU / Hybrid (CPU+GPU) support
 - Leverage ESSL BLAS + NVIDIA cuBLAS
 - Support multiple GPUs
 - Support problem sizes larger than GPU Memory

Parallel ESSL

- 150+ SPMD highly-tuned mathematical subroutines
 - L2/L3 PBLAS, Linear Algebraic Equations, Fourier Transforms...

Acceleration Enabled Programming Models



CUDA

Key Features:

- Gives direct access to the GPU instruction set
- Supports C, C++ and Fortran
- Generally achieves best leverage of GPUs for best application performance
- PGI/NVIDIA Compiler
- CUDA C/C++ for Power via XL NVCC



Key Features:

- Designed to simplify Programming of heterogeneous CPU/GPU systems
- Directive based parallelization for accelerator device
- PGI/NVIDIA Compiler
- OpenACC/gcc



Key Features:

- OpenMP 4.0 introduces offloading and support for heterogeneous CPU/GPU
- Leverage existing OpenMP high level directives support
- IBM XL Compiler
- Open Source LLVM OpenMP Compiler

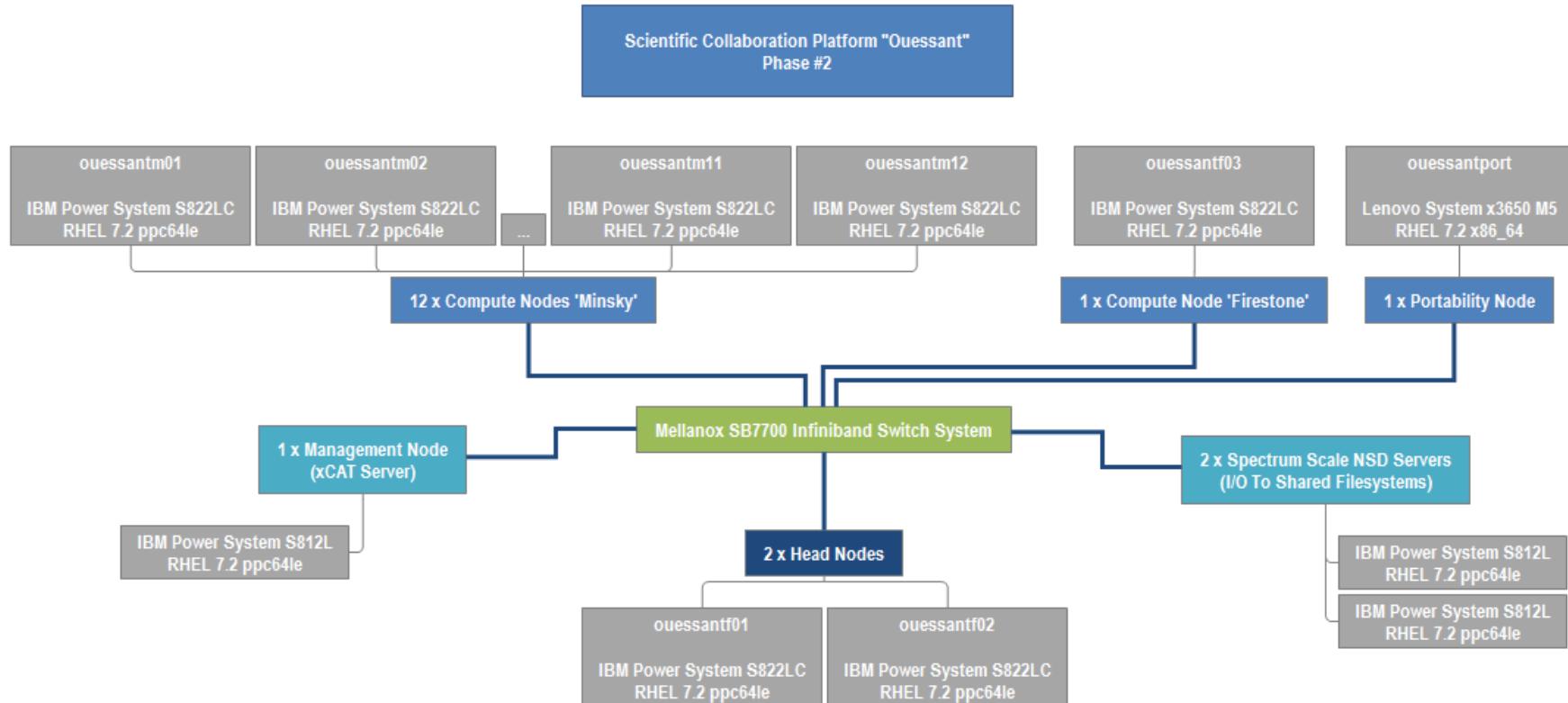
Supported GPU Offloading Features by Compiler Family (4Q 2016)

Compiler	OpenACC	OpenMP
GCC 4	-	4.0 w/o Offload
GCC 5	Experimental 2.0a No Offload on ppc64le	4.0 w/o Offload
GCC 6	Partial 2.0a	Partial 4.0
IBM XL	-	Partial 4.0+4.5
LLVM	-	Partial 4.0+4.5
PGI	Full 2.0a+2.5	-

Ouessant Platform

Scientific Collaboration Prototype

Platform Architecture



Connection

- **Hostname**
 - ouessant / ouessant.idris.fr
- **Users**
 - cfor001 To cfor032
 - Check on attendee list
 - Password: mastergpu
- **Job Submission Through IBM Platform LSF**

Job Submission: IBM Platform LSF

Typical Submission Script

```
#!/bin/bash

#BSUB -cwd ~/helloworld
#BSUB -e ~/helloworld/stderr
#BSUB -J HelloWorld
#BSUB -n 2
#BSUB -o ~/helloworld/stdout
#BSUB -q compute
#BSUB -R "affinity[core(1) : cpubind=core:distribute=balance] "
#BSUB -R "span[ptile=2]"
#BSUB -W 00:05

export OMP_NUM_THREADS=1

mpiexec ~/helloworld/helloworld.exe
```

Directives

Option	Value	Purpose
-cwd	<Path>	Execution Directory
-e	<File>	stderr File
-J	<Job Name>	Job Name
-n	<# MPI Tasks>	Total Number of MPI Tasks
-o	<File>	stdout File
-q	<Queue>	Target Queue
-R	"span[ptile=<ppn >]"	Resource Specification: Number of Tasks per Node
	"affinity[<level>(#):cpubind=<level>:distribute=<policy>]"	Resource Specification: Processor & Memory Affinity
-W	HH:MM	Run Limit
-X		Travail exclusif

User Commands

Command	Argument	Purpose
bjobs		List active jobs (waiting, in progress)
	-u { <User ID> all }	Restrict to specified user
	-X	List associated resources
bkill	<Job ID>	Kill job
bpeek	<Job ID>	Display job stdout/stderr
	-f	Refresh display in real time
bqueues		List existing queues
bstatus	<Job ID>	Display job status
bsub	< <Submission Script>	Submit job into queue

Task Placement

- **Default Policy: “Group Round Robin”**
 - ‘ptile’ MPI tasks per allocated compute node
 - One allocated compute node after the other until all MPI tasks have been placed
- **Alternative Policy**
 - Specified through environment variable ‘LSB_TASK_GEOMETRY’
 - Example:
 - `export LSB_TASK_GEOMETRY="{{(0,3)(1,4)(2,5)}«`
 - Tasks #0 & #3 placed on node #1
 - Tasks #1 & #4 placed on node #2
 - Tasks #2 & #5 placed on node #3

Processor Affinity

- Proper Processor Affinity is Mandatory for Performance
 - Especially when SMT is enabled
- Two Levels
 - CPU
 - GPU
- Requirement
 - Know and understand system topology

```
[user@host ~]$ ppc64_cpu --info
Core 0: 0*   1*   2*   3*   4*   5*   6*   7*
Core 1: 8*   9*   10*  11*  12*  13*  14*  15*
Core 2: 16*  17*  18*  19*  20*  21*  22*  23*
Core 3: 24*  25*  26*  27*  28*  29*  30*  31*
Core 4: 32*  33*  34*  35*  36*  37*  38*  39*
Core 5: 40*  41*  42*  43*  44*  45*  46*  47*
Core 6: 48*  49*  50*  51*  52*  53*  54*  55*
Core 7: 56*  57*  58*  59*  60*  61*  62*  63*
Core 8: 64*  65*  66*  67*  68*  69*  70*  71*
Core 9: 72*  73*  74*  75*  76*  77*  78*  79*
Core 10: 80*  81*  82*  83*  84*  85*  86*  87*
Core 11: 88*  89*  90*  91*  92*  93*  94*  95*
Core 12: 96*  97*  98*  99*  100* 101* 102* 103*
Core 13: 104* 105* 106* 107* 108* 109* 110* 111*
Core 14: 112* 113* 114* 115* 116* 117* 118* 119*
Core 15: 120* 121* 122* 123* 124* 125* 126* 127*
Core 16: 128* 129* 130* 131* 132* 133* 134* 135*
Core 17: 136* 137* 138* 139* 140* 141* 142* 143*
Core 18: 144* 145* 146* 147* 148* 149* 150* 151*
Core 19: 152* 153* 154* 155* 156* 157* 158* 159*
```

```
[user@host ~]$ nvidia-smi topo --matrix
      GPU0    GPU1    GPU2    GPU3    mlx5_0    mlx5_1    CPU Affinity
GPU0     X        NV2      SOC      SOC      SOC      SOC    0-79
GPU1      NV2     X        SOC      SOC      SOC      SOC    0-79
GPU2      SOC     SOC      X        NV2      SOC      SOC    80-159
GPU3      SOC     SOC      NV2      X        SOC      SOC    80-159
mlx5_0    SOC     SOC      SOC      SOC      X        PIX
mlx5_1    SOC     SOC      SOC      SOC      PIX      X
```

Legend:

X = Self
SOC = Connection traversing PCIe as well as the SMP link between CPU sockets (e.g. QPI)
PHB = Connection traversing PCIe as well as a PCIe Host Bridge (typically the CPU)
PXB = Connection traversing multiple PCIe switches (without traversing the PCIe Host Bridge)
PIX = Connection traversing a single PCIe switch
NV# = Connection traversing a bonded set of # NVLinks

Processor Affinity: CPU

- **Purpose**
 - Avoid as much as possible resource sharing
 - Avoid Linux Scheduler to move processes / threads between CPU cores
- **Management**
 - Manual
 - Through LSF Affinity String
 - Automated management is forthcoming
 - Complex Syntax
 - {core|thread}(n):cpubind={core|thread}:distribute={balance|pack}
 - <https://goo.gl/5v6Qmu>
- **Tips'n Tricks**
 - Checking Required!
 - ALWAYS check applied processor affinity with the help of monitoring tools
 - htop, nmon

Processor Affinity: CPU – A Few Useful Examples 1/2

thread(1, exclusive=(core, intask))*10:cpubind=thread:distribute=balance

CPUs/Cores	Core #01	Core #02	Core #03	Core #04	Core #05	Core #06	Core #07	Core #08	Core #09	Core #10	Core #11	Core #12	Core #13	Core #14	Core #15	Core #16	Core #17	Core #18	Core #19	Core #20
CPU #0	P00.t00	P00.t01	P00.t02	P00.t03	P00.t04	P00.t05	P00.t06	P00.t07	P00.t08	P00.t09	P01.t00	P01.t01	P01.t02	P01.t03	P01.t04	P01.t05	P01.t06	P01.t07	P01.t08	P01.t09
CPU #1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

thread(1):cpubind=thread:distribute=balance

CPUs/Cores	Core #01	Core #02	Core #03	Core #04	Core #05	Core #06	Core #07	Core #08	Core #09	Core #10	Core #11	Core #12	Core #13	Core #14	Core #15	Core #16	Core #17	Core #18	Core #19	Core #20
CPU #0	P00	P01	P02	P03	P04	-	-	-	-	-	P05	P06	P07	P08	P09	-	-	-	-	-
CPU #1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

thread(1,exclusive=(core,injob)):cpubind=thread:distribute=pack

CPUs/Cores	Core #01	Core #02	Core #03	Core #04	Core #05	Core #06	Core #07	Core #08	Core #09	Core #10	Core #11	Core #12	Core #13	Core #14	Core #15	Core #16	Core #17	Core #18	Core #19	Core #20
CPU #0	P00	P01	P02	P03	P04	P05	P06	P07	P08	P09	-	-	-	-	-	-	-	-	-	-
CPU #1	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

Processor Affinity: CPU – A Few Useful Examples 2/2

thread(2):cpubind=thread:distribute=balance

CPUs/Cores	Core #01	Core #02	Core #03	Core #04	Core #05	Core #06	Core #07	Core #08	Core #09	Core #10	Core #11	Core #12	Core #13	Core #14	Core #15	Core #16	Core #17	Core #18	Core #19	Core #20
CPU #0	P00.t00	P01.t00	P02.t00	P03.t00	P04.t00	-	-	-	-	-	P05.t00	P06.t00	P07.t00	P08.t00	P09.t00	-	-	-	-	-
CPU #1	P00.t01	P01.t01	P02.t01	P03.t01	P04.t01	-	-	-	-	-	P05.t01	P06.t01	P07.t01	P08.t01	P09.t01	-	-	-	-	-
CPU #2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

thread(2,same=core):cpubind=thread:distribute=balance

CPUs/Cores	Core #01	Core #02	Core #03	Core #04	Core #05	Core #06	Core #07	Core #08	Core #09	Core #10	Core #11	Core #12	Core #13	Core #14	Core #15	Core #16	Core #17	Core #18	Core #19	Core #20
CPU #0	P00.t00	P01.t00	P02.t00	P03.t00	P04.t00	P05.t00	P06.t00	P07.t00	P08.t00	P09.t00	P10.t00	P11.t00	P12.t00	P13.t00	P14.t00	P15.t00	P16.t00	P17.t00	P18.t00	P19.t00
CPU #1	P00.t01	P01.t01	P02.t01	P03.t01	P04.t01	P05.t01	P06.t01	P07.t01	P08.t01	P09.t01	P10.t01	P11.t01	P12.t01	P13.t01	P14.t01	P15.t01	P16.t01	P17.t01	P18.t01	P19.t01
CPU #2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

thread(1):cpubind=thread:distribute=balance

CPUs/Cores	Core #01	Core #02	Core #03	Core #04	Core #05	Core #06	Core #07	Core #08	Core #09	Core #10	Core #11	Core #12	Core #13	Core #14	Core #15	Core #16	Core #17	Core #18	Core #19	Core #20
CPU #0	P00	P02	P04	P06	P08	P10	P12	P14	P16	P18	P20	P22	P24	P26	P28	P30	P32	P34	P36	P38
CPU #1	P01	P03	P05	P07	P09	P11	P13	P15	P17	P19	P21	P23	P25	P27	P29	P31	P33	P35	P37	P39
CPU #2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
CPU #7	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

Processor Affinity: GPU

- **Purpose**
 - Avoid as much as possible resource sharing
 - Prefer GPU direct access (Avoid going through CPU-To-CPU link)
- **Management**
 - LSF-Automated
 - Supposed to be, at least ☺
 - Manual
 - Through environment variable setting
- **Tips'n Tricks**
 - Checking Required!
 - ALWAYS check applied processor affinity with the help of monitoring tools
 - gpustat, nvidia-smi pmon

GPU Compute Mode

Exclusive, Shared, or Multi-Process Service

GPU Compute Mode

- Available Compute Modes

- DEFAULT

- Equivalent to: Shared

- EXCLUSIVE_PROCESS

- Only one process allowed to run

- EXCLUSIVE_THREAD

- Only one thread allowed to run

- PROHIBITED

- Run not allowed

- GPU-Specific

- GPUs can have different Compute Modes

```
[user@host ~]$ nvidia-smi --query --display=COMPUTE
```

```
=====NVSMI LOG=====
```

Timestamp	:	Tue Nov 1
19:59:57 2016		
Driver Version	:	361.93.02

Attached GPUs	:	4
GPU 0002:01:00.0		
Compute Mode	:	Default

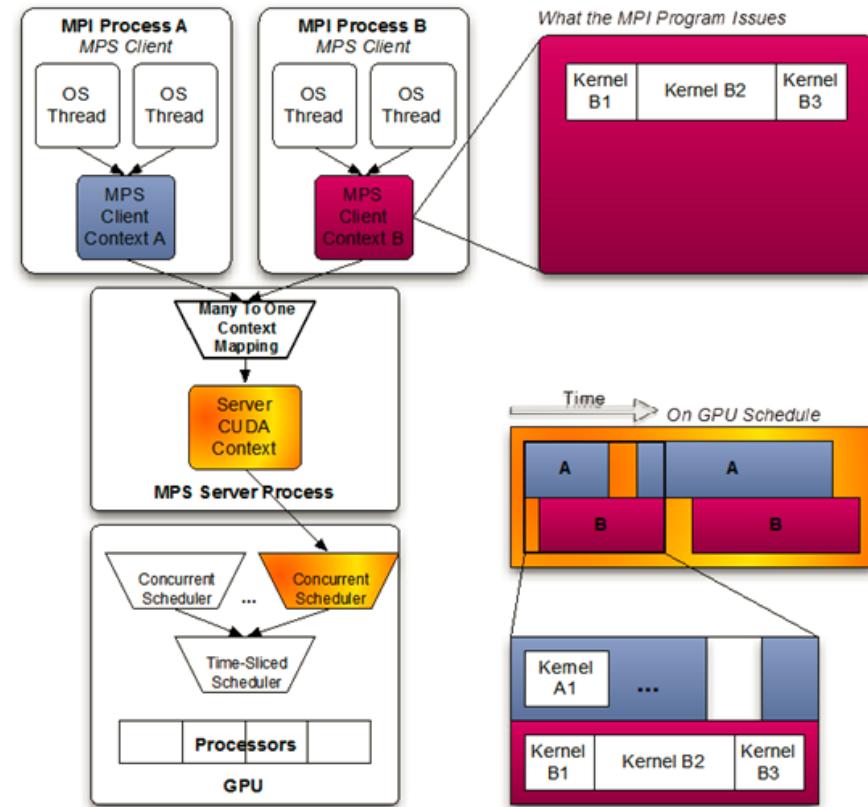
GPU 0003:01:00.0		
Compute Mode	:	Default

GPU 0006:01:00.0		
Compute Mode	:	Default

GPU 0007:01:00.0		
Compute Mode	:	Default

Multi-Process Service (MPS)

- MPS Benefits
 - GPU Utilization
 - w/o MPS: Single process may underutilize the compute and memory-bandwidth capacity
 - w/MPS: kernel and memcpy operations from different processes can overlap on the GPU, achieving higher utilization and shorter running times
 - Reduced On-GPU Context Storage
 - w/o MPS: each CUDA process using a GPU allocates separate storage and scheduling resources on the GPU
 - w/MPS: server allocates one copy of GPU storage and scheduling resources shared by all of its clients
 - Reduced GPU Context Switching
 - w/o MPS: when processes share the GPU their scheduling resources must be swapped on and off the GPU
 - w/MPS: server shares one set of scheduling resources between all of its clients, eliminating the overhead of swapping when the GPU is scheduling between those clients



End Of Presentation