

NVIDIA DGX-2

Haiduong Vo, DGX Product Management



NVIDIA DGX-2

Agenda

- NVIDIA DGX-2 Product Features and Benefits:
 - Integrated Hardware, including NVSwitch technology
 - Integrated Software
- DGX-2 Performance Results

DEMO: HIGH-RESOLUTION VIDEO GENERATION

NVIDIA DGX-2: Enabling New Use Case

- ▶ Input Video
- ▶ Before Video: With DGX-1
- ▶ After Video: With DGX-2

Show Videos

DGX-2: BUILT FOR THE MOST COMPLEX DL APPS

High-Resolution Video Generation from NV Research

- Generating 2048x1024 Video
- Custom Network Based on pix2pixHD Project
 - 6X the Size of Resnet152
- PyTorch Framework
- DGX-1 (V100/16GB): Training 4 Frames Simultaneously, 100GB GPU Memory usage
- DGX-2 (V100/32GB): Training **8+** Frames Simultaneously, 380GB+ Total GPU Memory usage
- “Everything just works” on DGX-2. No SW adaptation to run the code.

**THE WORLD'S FIRST
2 PETAFLOPS SYSTEM**



INTRODUCING NVIDIA DGX-2

THE WORLD'S MOST POWERFUL
AI SYSTEM FOR THE MOST COMPLEX
AI CHALLENGES

- DGX-2 is the newest addition to the DGX family, powered by DGX software
- Deliver accelerated AI-at-scale deployment and effortless operations
- Step up to DGX-2 for unrestricted model parallelism and faster time-to-solution

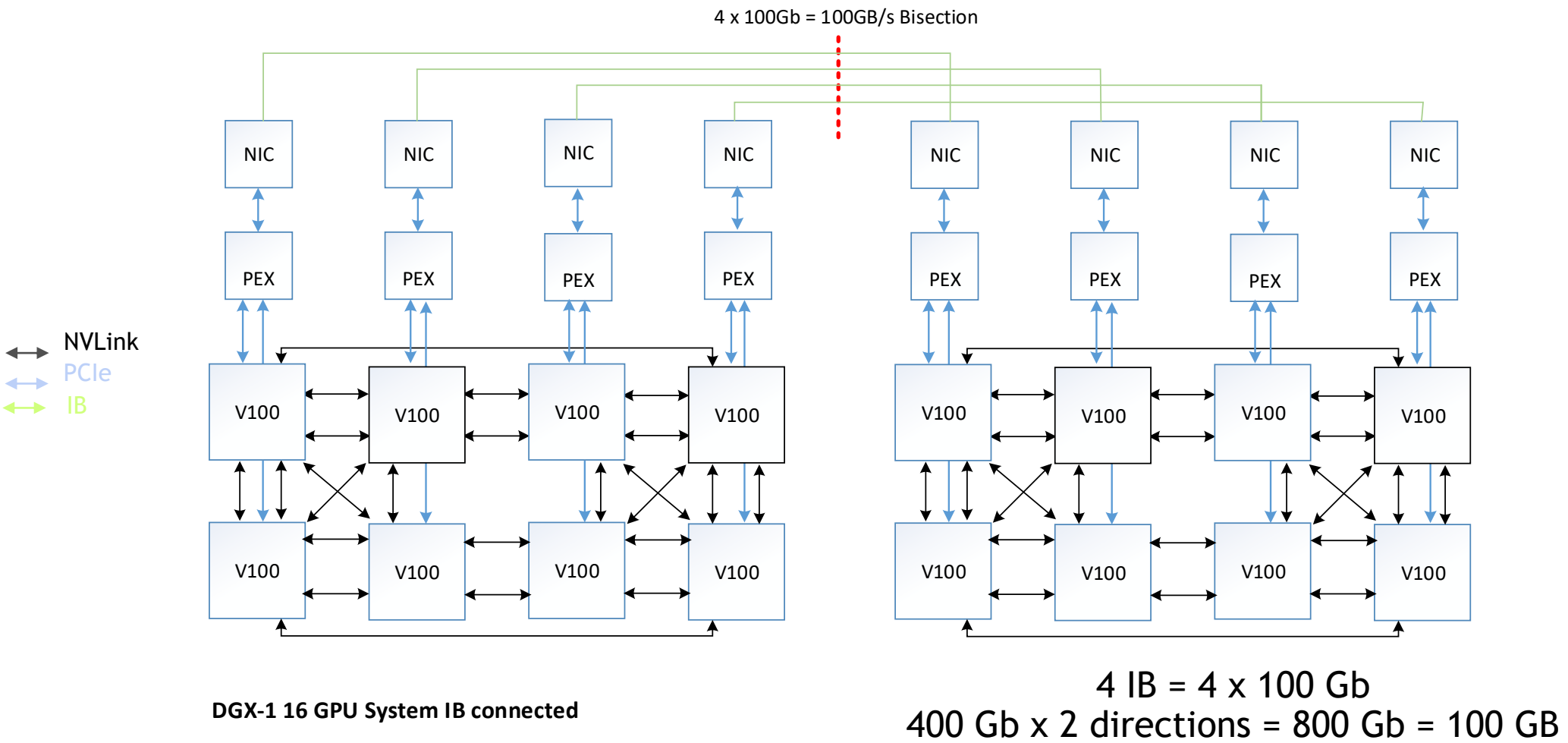


- 2.4 TB/s bisection bandwidth
- Equivalent to a PCIe bus with 1,200 lanes

NVSWITCH: THE REVOLUTIONARY AI NETWORK FABRIC

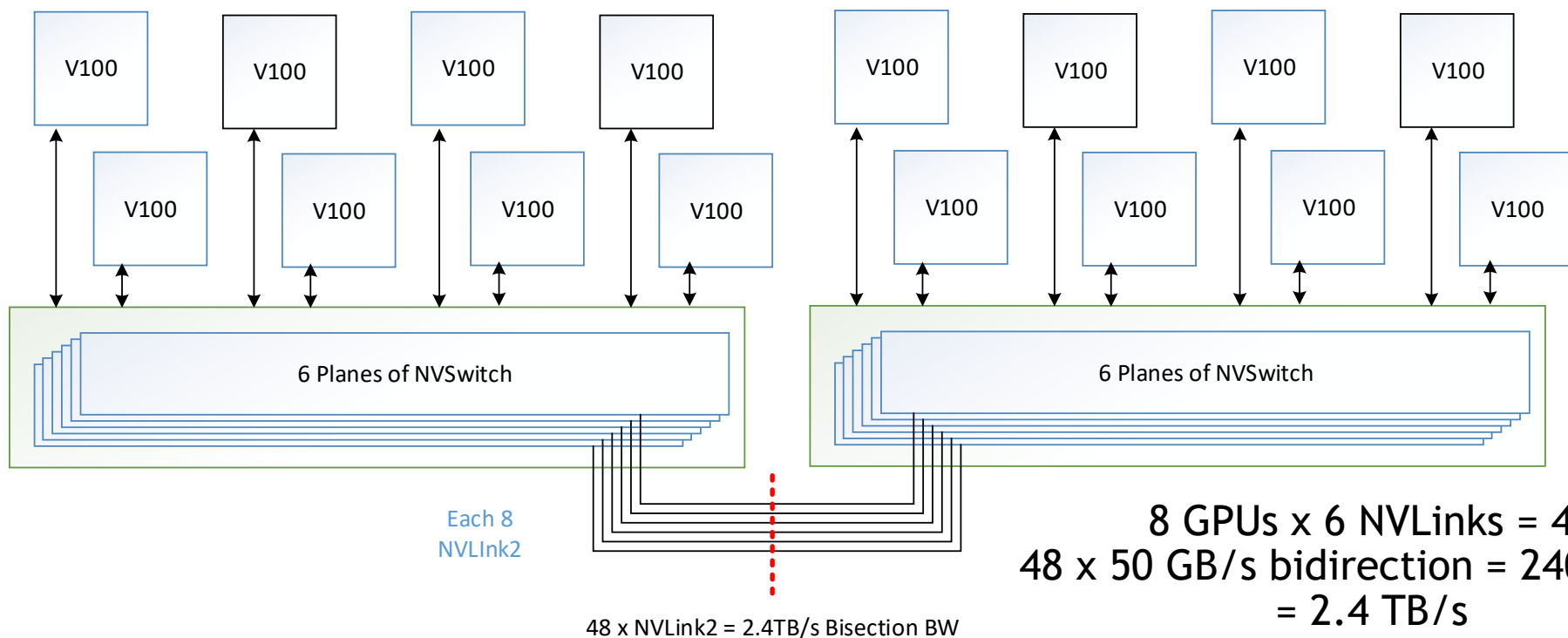
- Inspired by leading edge research that demands unrestricted model parallelism
- Like the evolution from dial-up to broadband, NVSwitch delivers a networking fabric for the future, today
- Delivering 2.4 TB/s bisection bandwidth, equivalent to a PCIe bus with 1,200 lanes
- NVSwitches on DGX-2 capable of downloading all of Netflix HD content in under a minute

100GB/S BISECTION B/W: USING IB ON TWO DGX-1



2.4 TB/S USING NVSWITCH PLANE ON DGX-2

24X bisection bandwidth



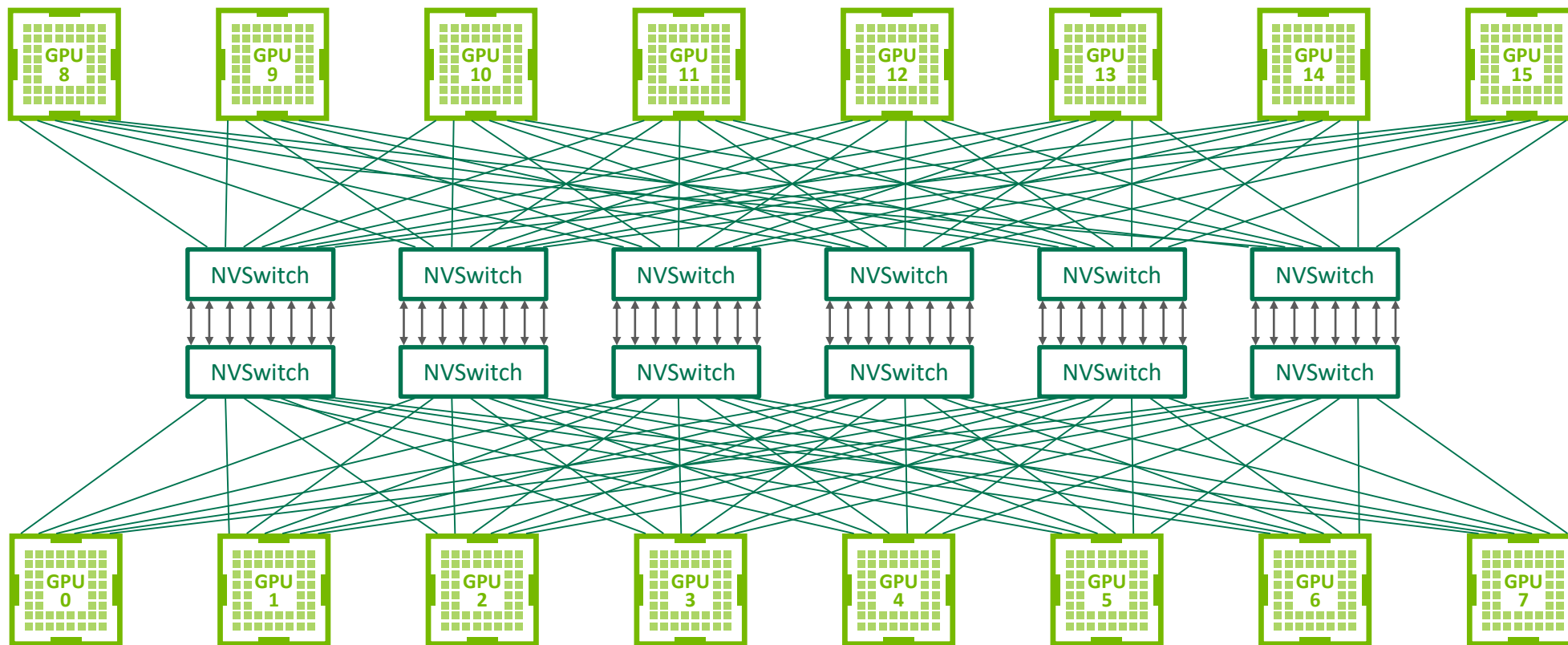
DGX-2 16 GPU System

Each 8
NVLink2

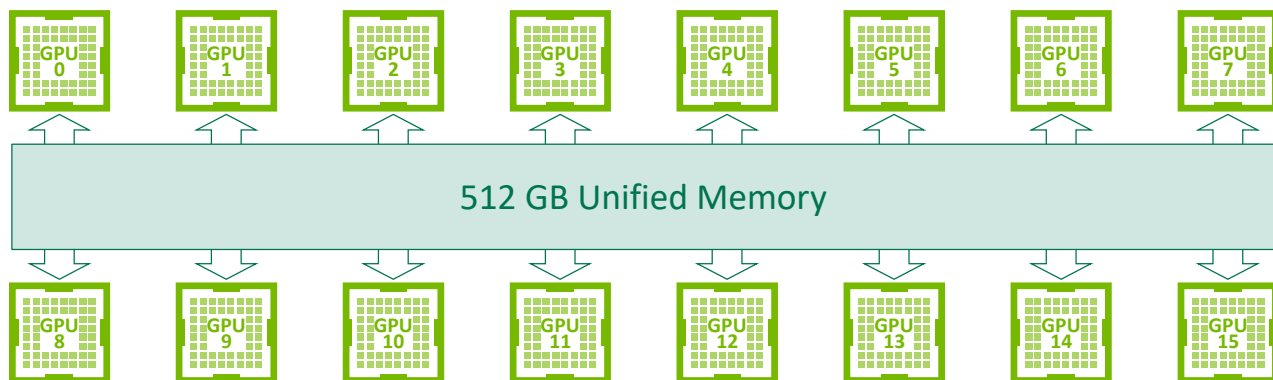
48 x NVLink2 = 2.4TB/s Bisection BW

8 GPUs x 6 NVLinks = 48
48 x 50 GB/s bidirection = 2400 GB/s
= 2.4 TB/s

FULL NON-BLOCKING BANDWIDTH



UNIFIED MEMORY + DGX-2



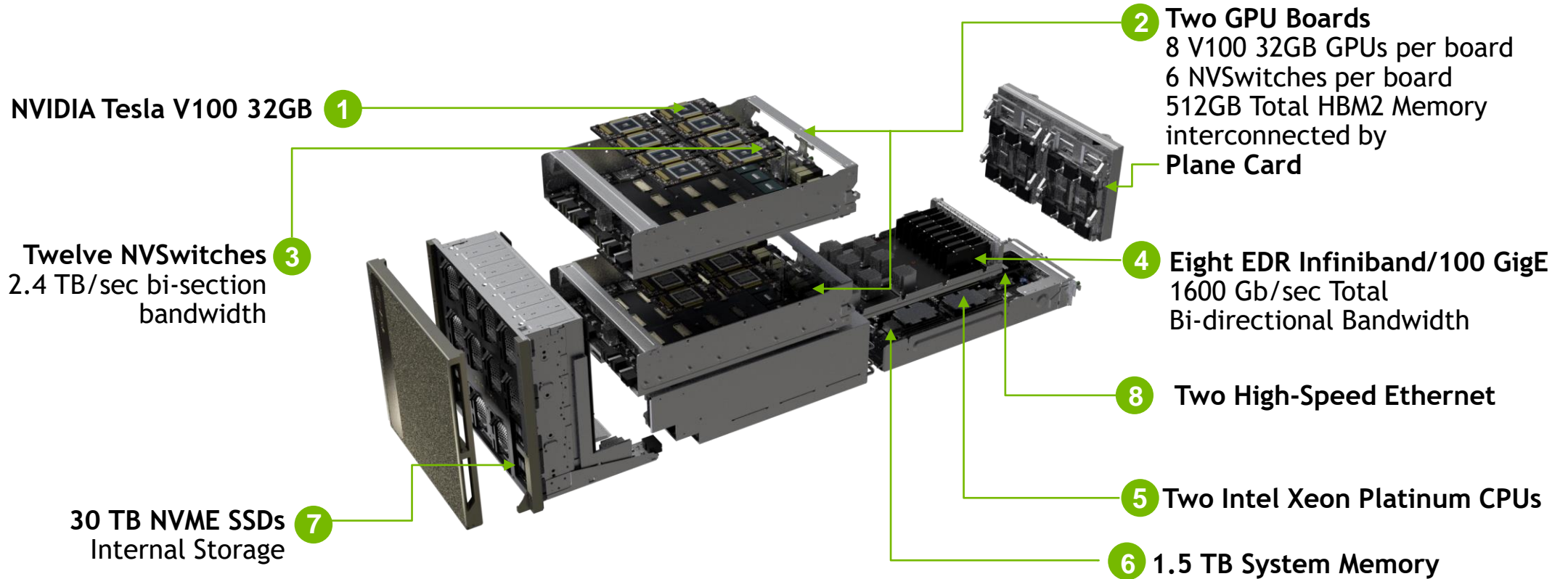
UNIFIED MEMORY PROVIDES

Single memory view
shared by all GPUs

Automatic migration of data
between GPUs

User control of data locality

DESIGNED TO TRAIN THE PREVIOUSLY IMPOSSIBLE



SYSTEM SPECS: DGX-2 AND DGX-1

App Focus Components: GPU AND CPU, NVSwitch

	NVIDIA DGX-2	NVIDIA DGX-1 (V100/32GB)
GPUs	16X NVIDIA Tesla V100	8X NVIDIA Tesla V100
GPU Memory	512 GB total and Nvswitch (closely resemble a large GPU)	256 GB total
NVIDIA NVSwitch	12 total	N/A
Performance	2 petaFLOPS (FP16)	1 petaFLOPS (FP16)
CUDA Cores/Tensor Cores	81920/10240	40960/5120
CPU	2X Intel Xeon Platinum 8168, 2.7 GHz, 24-cores	2X Intel Xeon E5-2698 v4, 2.2 GHz, 20-cores
System Memory	1.5 TB	512 GB
Network	8X 100 Gb/sec Infiniband/100GigE Dual PCIe slots for 10/25/40/100 Gb/sec Ethernet	4X 100 Gb/sec Infiniband/100GigE Dual 10 Gb/sec Ethernet
Storage	OS: 2 x 960GB NVME SSDs Internal Storage: 30TB (8 x 3.84TB) NVME SSDs	OS: 480 GB SAS SSDs Internal Storage: 7TB (4 x 1.92TB) SSDs
Software	Ubuntu Linux OS Same DGX SW stack	Ubuntu Linux OS Same DGX SW stack

SYSTEM SPECS: DGX-2 AND DGX-1

Power and Physical Dimensions

	NVIDIA DGX-2	NVIDIA DGX-1 (V100/32GB)
Maximum Power Usage	10 kW	3.5 kW
System Weight	340 lbs (154.2 Kgs)**	134 lbs
System Dimensions	10RU Height: 17.3 in (440.0 mm)** Width: 19.0 in (482.3 mm)** Length: 31.3 in (795.452 mm) ** - No Front Bezel 32.8 in (834.0 mm)** - With Front Bezel	3RU Height: 131 mm Width: 444 mm Length: 866 mm – No Front Bezel
Operating Temperature range	5 °C to 35 °C (41 °F to 95 °F)	5 °C to 35 °C
Cooling	Air	Air

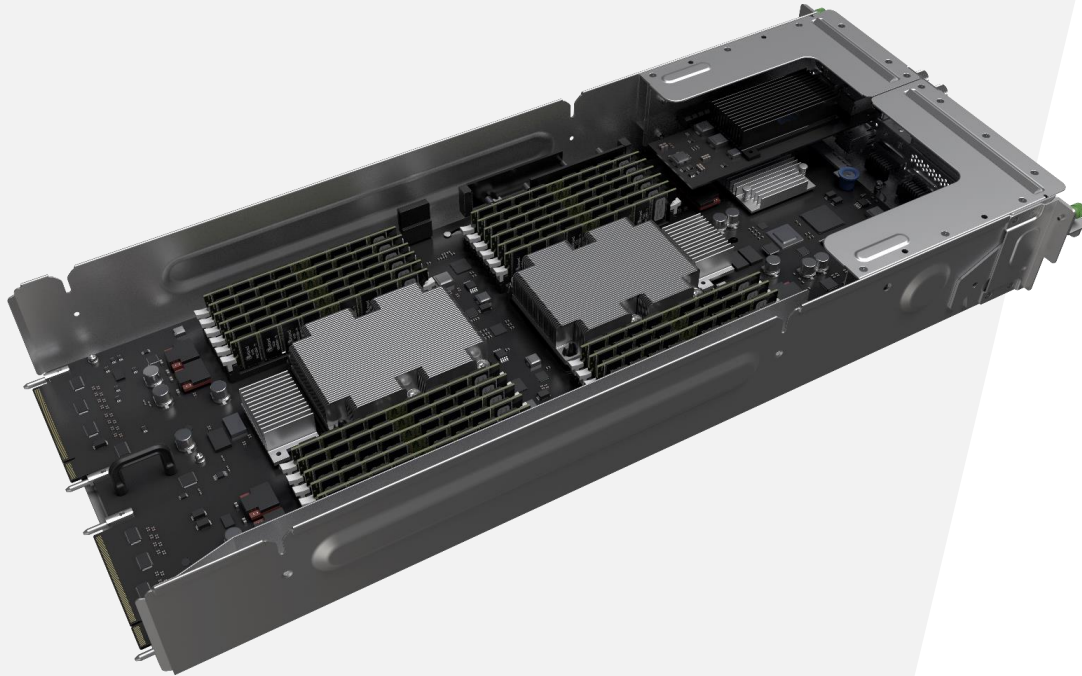
** Subject to Change

NVME SSD STORAGE

Rapidly ingest the largest datasets into cache

- Faster than SATA SSD, optimized for transferring huge datasets
- Dramatically larger user scratch space
- The protocol of choice for next-gen storage technologies
- 8 x 3.84TB NVMe in RAID0 (Data)
- 25.5 GB/sec Sequential Read bandwidth (vs. 2 GB/sec for 7TB of SAS SSDs on DGX-1)





LATEST GENERATION CPU AND 1.5TB SYSTEM MEMORY

Faster, more resilient, boot and storage management

- More system memory to handle larger DL and HPC applications
- 2 Intel Skylake Xeon Platinum 8168 - 2.7GHz, 24 cores
- 24 x 64GB DIMM System Memory

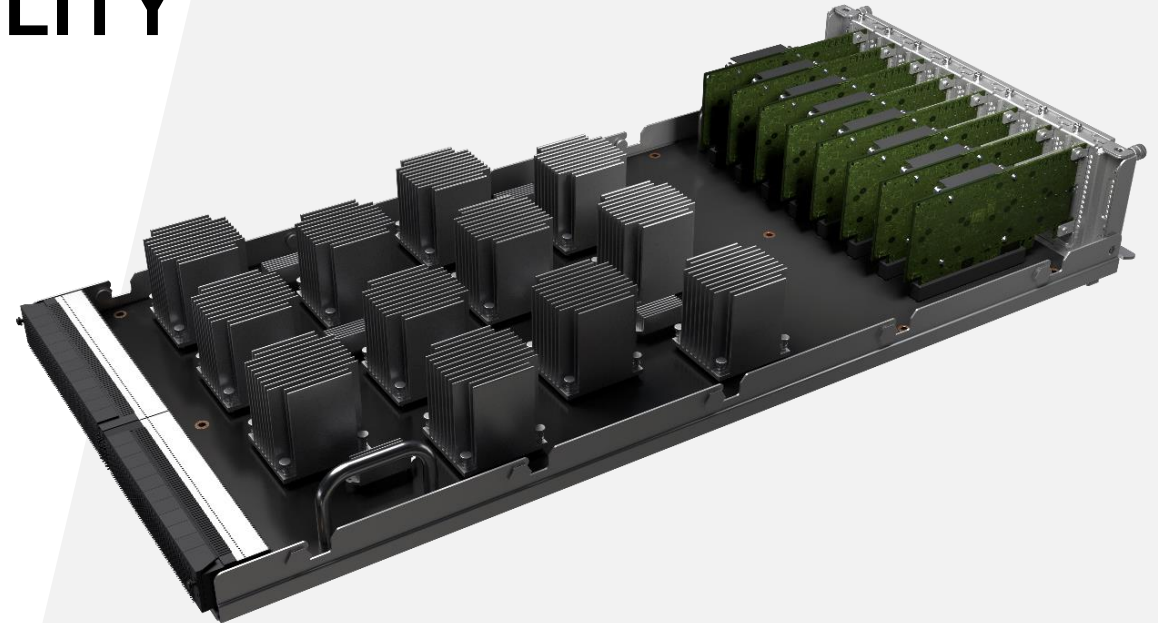
THE ULTIMATE IN NETWORKING FLEXIBILITY

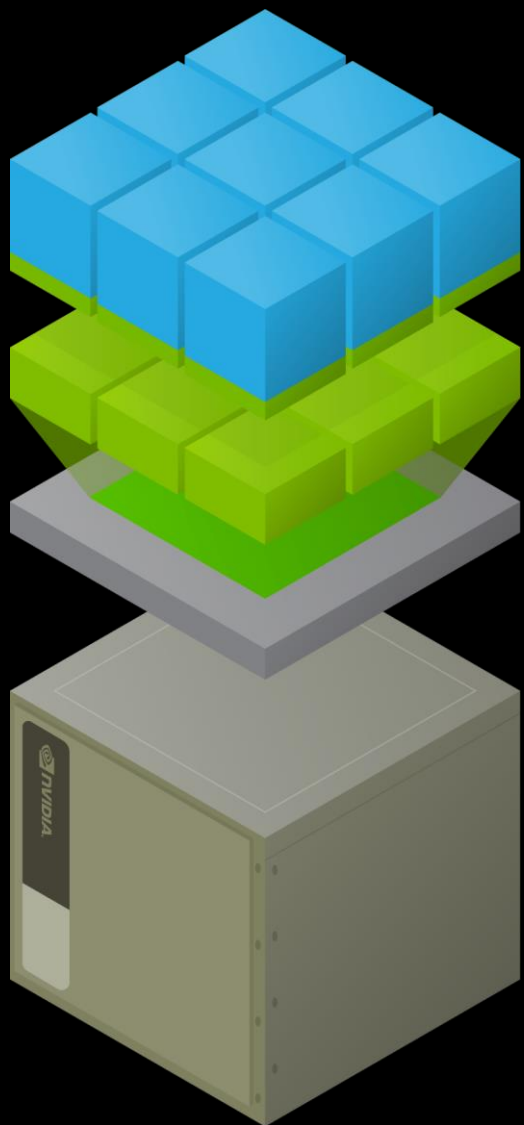
Grow your DL cluster effortlessly,
using the connectivity you prefer

- 8 EDR Infiniband / 100 GigE
- 1600 Gb/sec Total Bi-directional Bandwidth with low-latency
- Support for RDMA over Converged Ethernet (ROCE)

Also including dual-port Ethernet on CPU board

- Dual-port 10/25/40/56/100 GbE/sec





FLEXIBILITY WITH VIRTUALIZATION

Enable your own private DL Training Cloud for your Enterprise

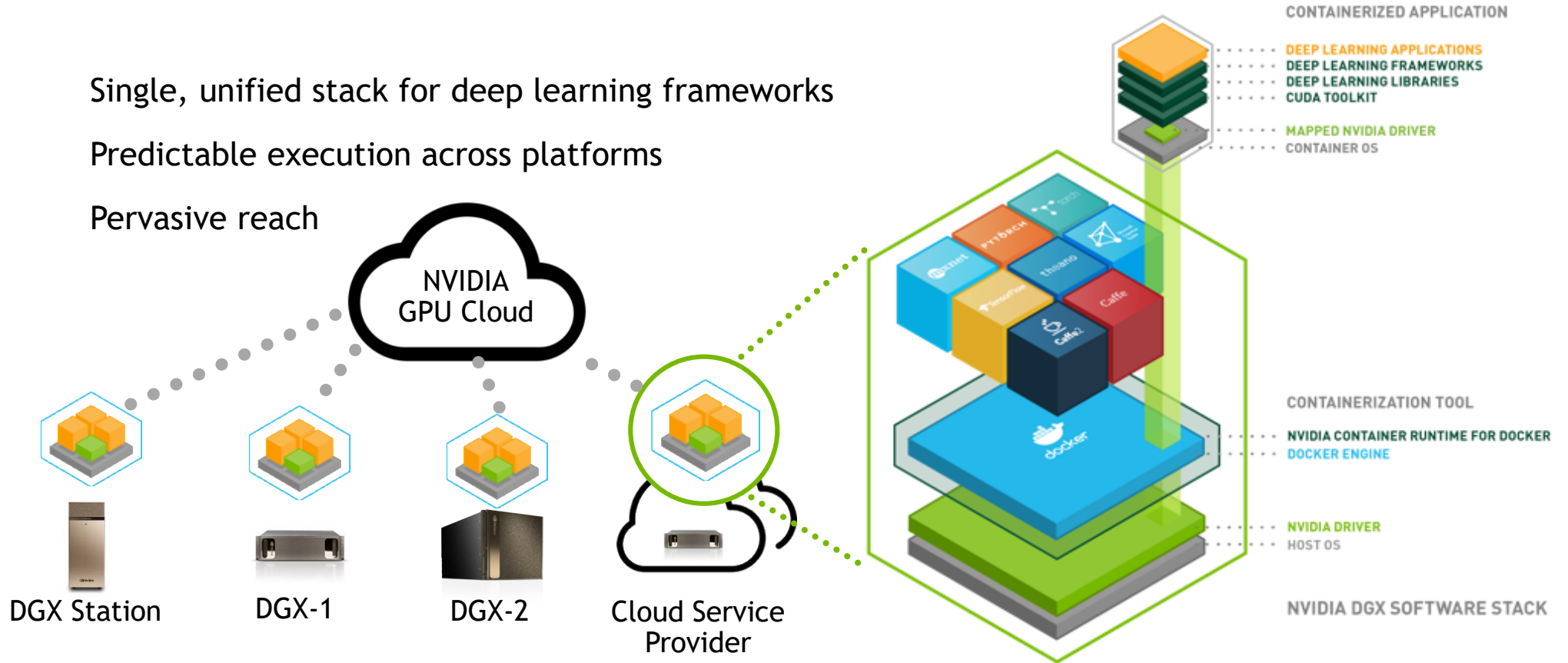
- KVM hypervisor for Ubuntu Linux
- Enable teams of developers to simultaneously access DGX-2
- Flexibly allocate GPU resources to each user and their experiments
- Full GPU's and NVSwitch access within VMs – either all GPU's or as few as 1

COMMON SOFTWARE STACK ACROSS DGX FAMILY

Single, unified stack for deep learning frameworks

Predictable execution across platforms

Pervasive reach



NGC REGISTRY

Simple access to a comprehensive catalog of GPU-accelerated software

Discover 30 GPU-Accelerated Containers

Deep learning, third-party managed HPC applications, NVIDIA HPC visualization tools, and partner applications

Innovate in Minutes, Not Weeks

Get up and running quickly and reduce complexity

Access from Anywhere

Use containers on PCs with NVIDIA Volta or Pascal™ architecture GPUs, NVIDIA DGX Systems, and supported cloud providers



NGC REGISTRY

25K User Registrations, 30+ Containers

DEEP LEARNING

caffe
caffe2
cntk
cuda
digits
mxnet
pytorch
tensorflow
tensorrt
theano
torch

HPC VIZ

paraview-holodeck
paraview-index
paraview-optix
index

HPC

bigdft
candle
gamess
gromacs
lammmps
lattice-microbes
milc
namd
relion
vmd

PARTNER

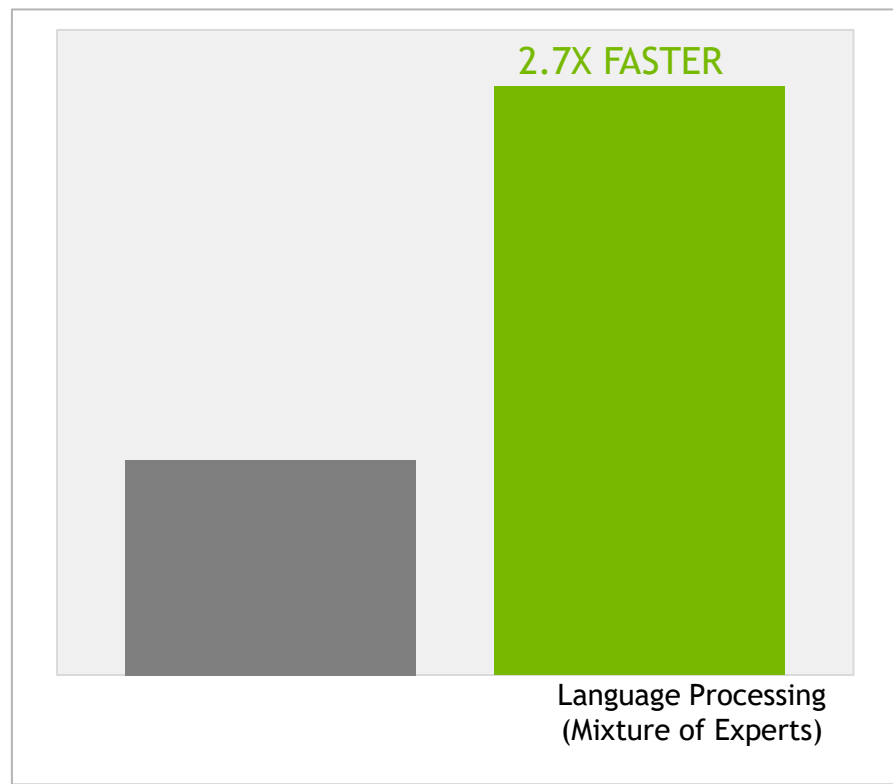
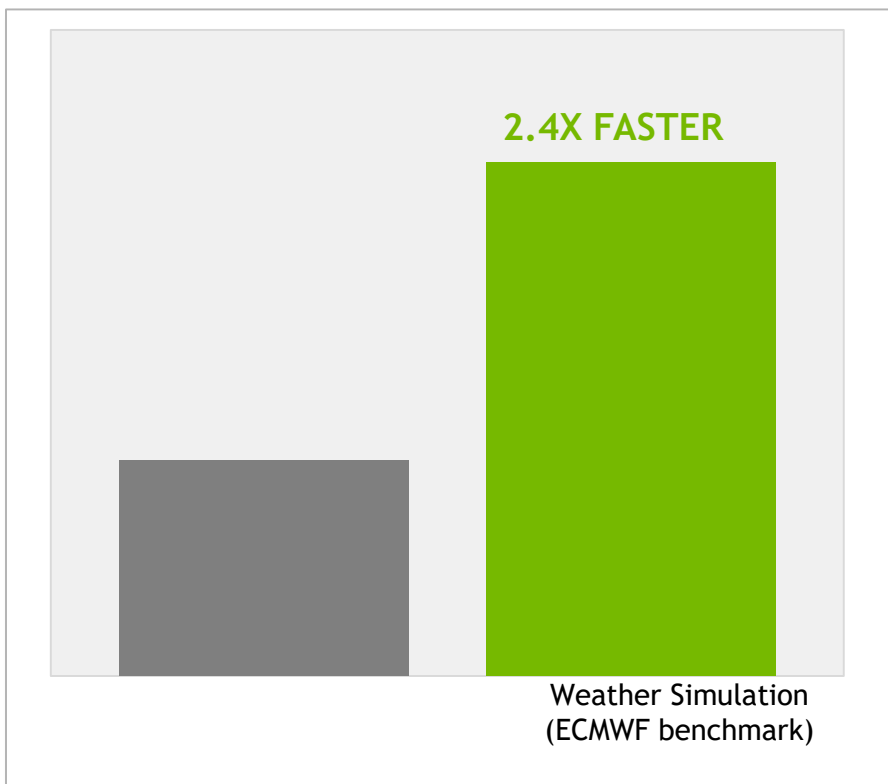
h2o
mapd
chainer
paddlepaddle
kinetica
matlab*

GUEST ACCESS

cuda
kubernetes*

*new

2X HIGHER PERFORMANCE WITH NVSWITCH

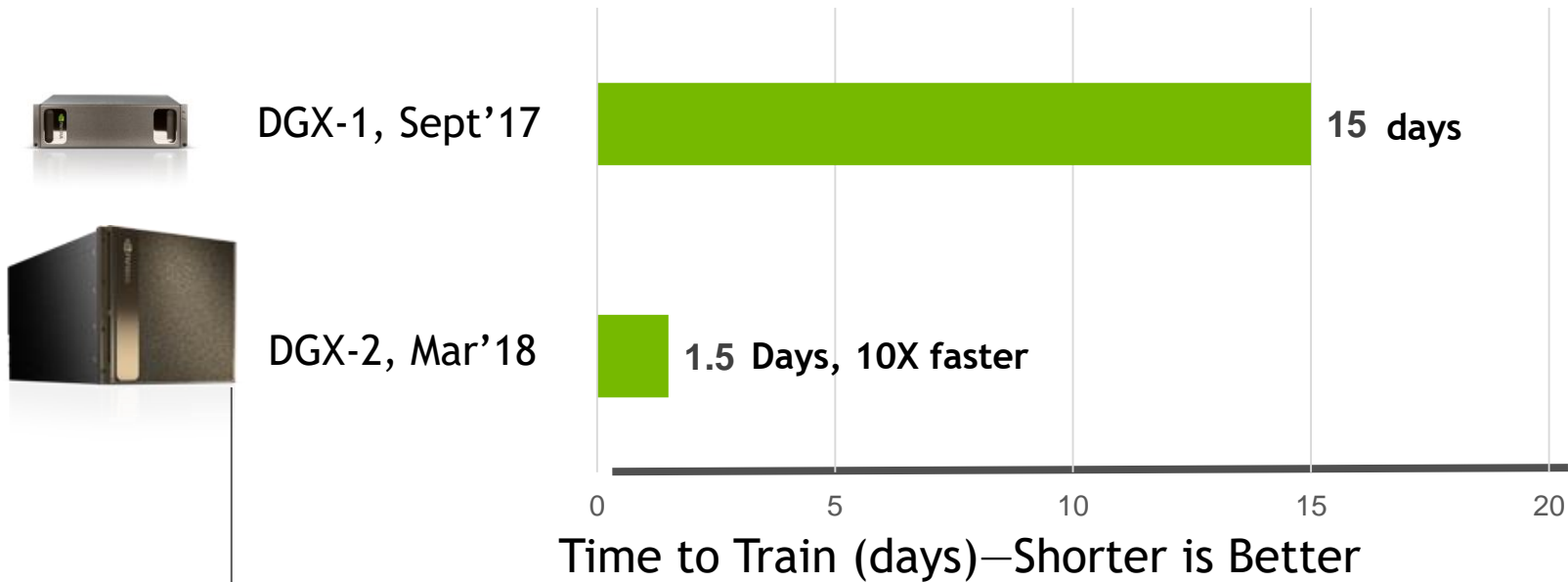


■ 2x DGX-1 (Volta)

■ DGX-2 with NVSwitch

2 DGX-1V servers have dual socket Xeon E5 2698v4 Processor. 8 x V100 GPUs. Servers connected via 4X 100Gb IB ports |
DGX-2 server has dual-socket Xeon Platinum 8168 Processor. 16 V100 GPUs

10X PERFORMANCE GAIN IN LESS THAN A YEAR




Performance gain through hardware and software improvements across the stack

Workload: FairSeq, 55 epochs to accuracy. PyTorch training performance.

“500X” IN 5 YEARS

2 GTX 580s – DEC ‘12

AlexNet

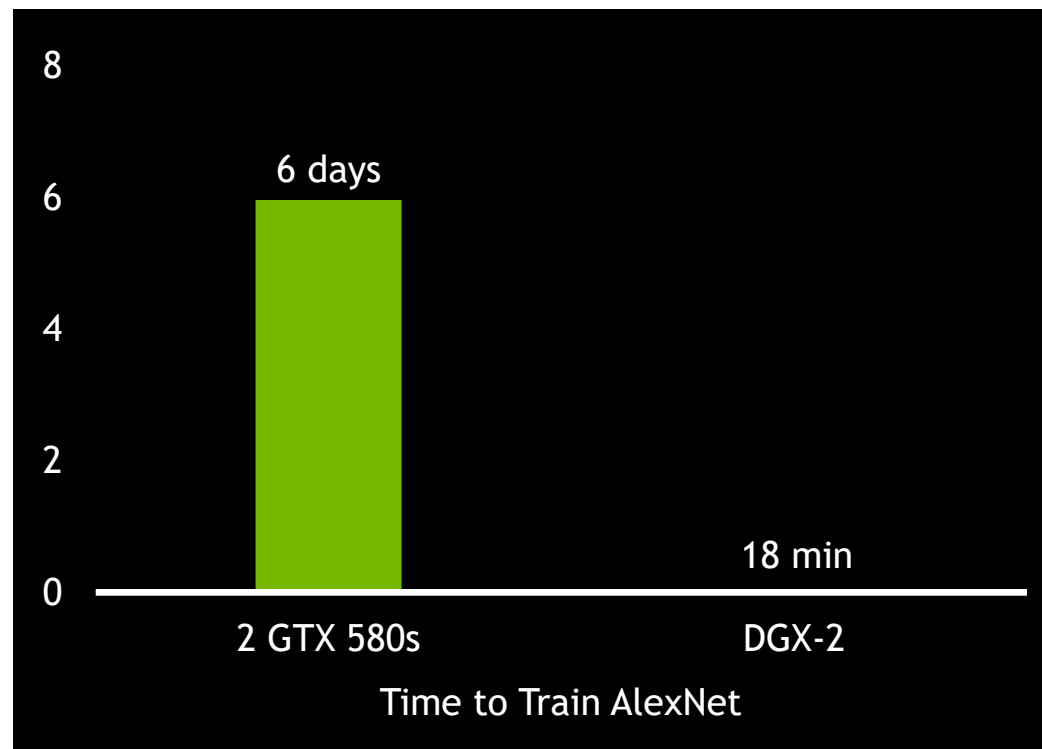


Framework	cuda-convnet	
System	NCCL	N/A
Software	cuDNN	N/A
Stack	cuBLAS	5.0
	cuFFT	5.0
	NPP	5.0
	CUDA	5.0
	Res Mgr	R304

DGX-2 – MAR ‘18



Framework	NV Caffe	0.17
System	NCCL	2.2
Software	cuDNN	7.1
Stack	cuBLAS	9.2
	cuFFT	9.2
	NPP	9.2
	CUDA	9.2
	Res Mgr	R396



THE PERFORMANCE OF 300 SKYLAKE SERVERS



300 Skylake Gold CPU Servers

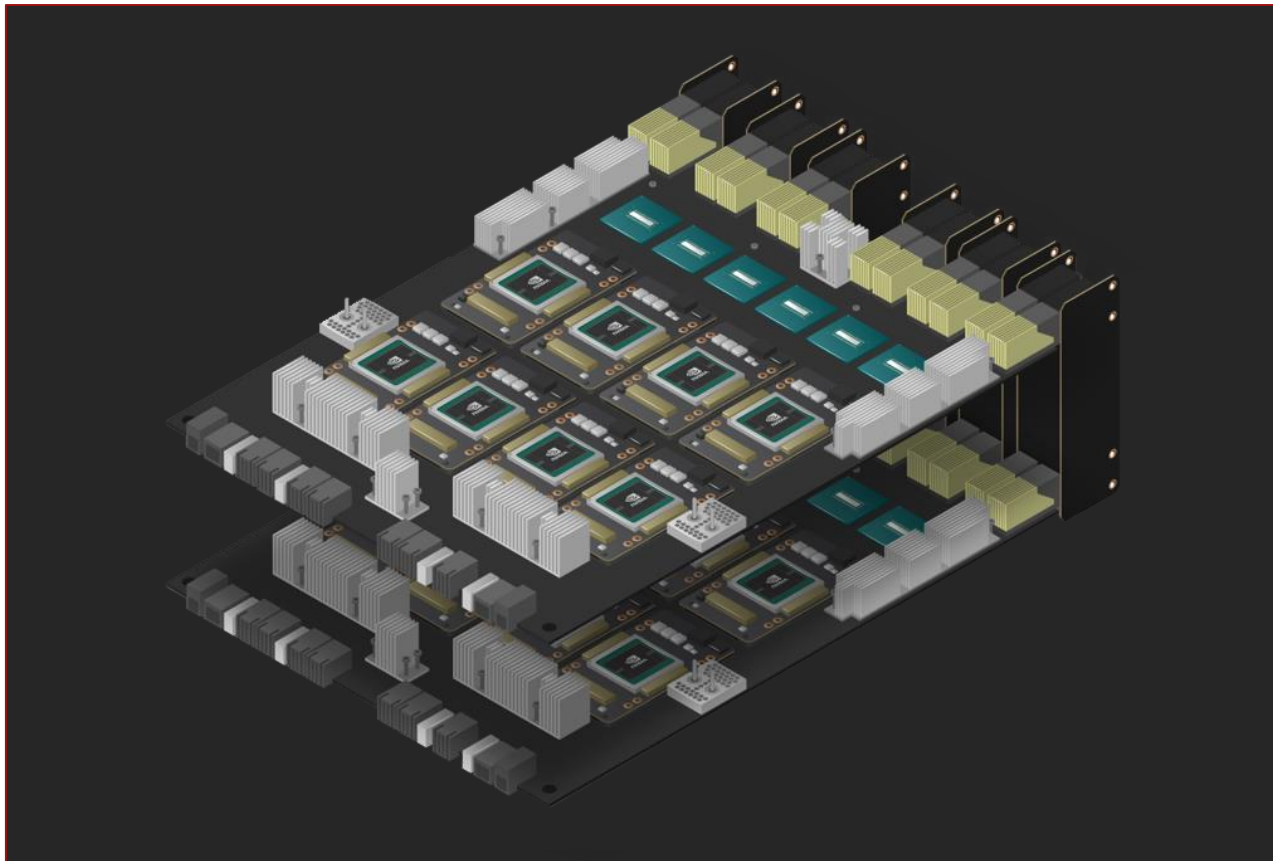
=



One DGX-2

NVSWITCH AVAILABILITY: NVIDIA HGX-2

- Eight GPU baseboard with six NVSwitches
- Two HGX-2 boards can be passively connected to realize 16-GPU systems
- ODM/OEM partners build servers utilizing NVIDIA HGX-2 GPU baseboards



GET EARLY ACCESS TO DGX-2

BE FIRST TO GET THE WORLD'S MOST POWERFUL DEEP LEARNING SYSTEM

Work with your NVIDIA team now:

- ✓ Review your DL capacity needs
- ✓ Submit your application for DGX-2 Early Access - Details Forthcoming
- ✓ Schedule site preparation
- ✓ Learn more about the DGX-2 - <https://www.nvidia.com/en-us/data-center/dgx-2/>



