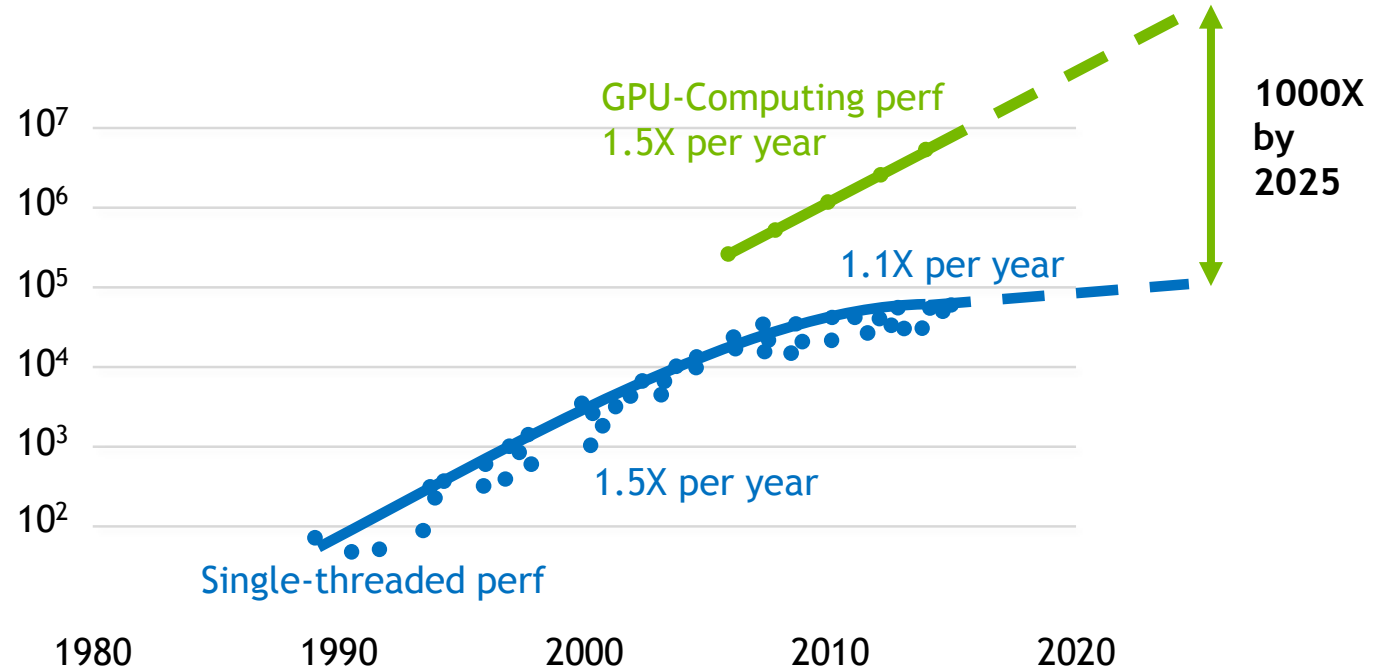# GPU COMPUTING

What's now ?

**NVIDIA.**

Guillaume BARAT
*gbarat@nvidia.com*
EMEA Higher-Education & Research

# RISE OF GPU COMPUTING

APPLICATIONS

ALGORITHMS

SYSTEMS

CUDA

CPU    GPU

ARCHITECTURE

$10^7$

$10^6$

$10^5$

$10^4$

$10^3$

$10^2$

GPU-Computing perf
1.5X per year

1000X
by
2025

1.1X per year

1.5X per year

Single-threaded perf

1980    1990    2000    2010    2020
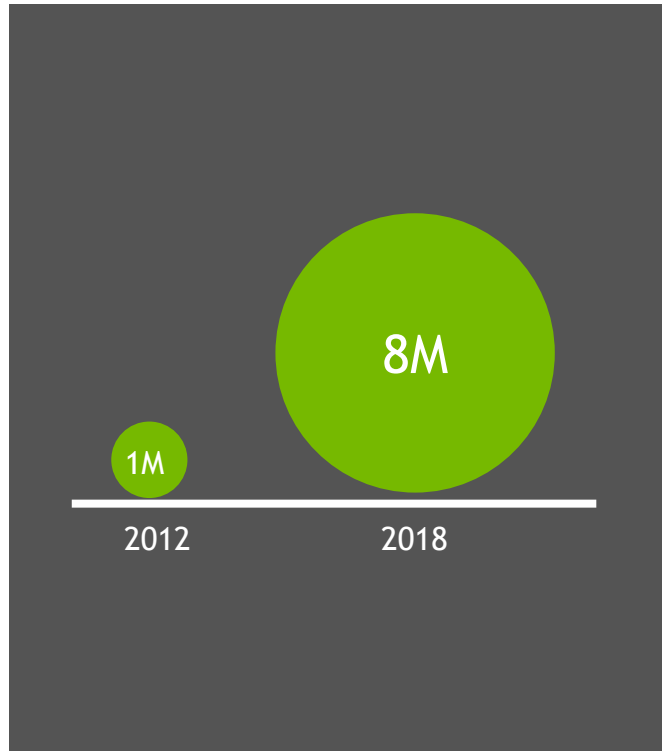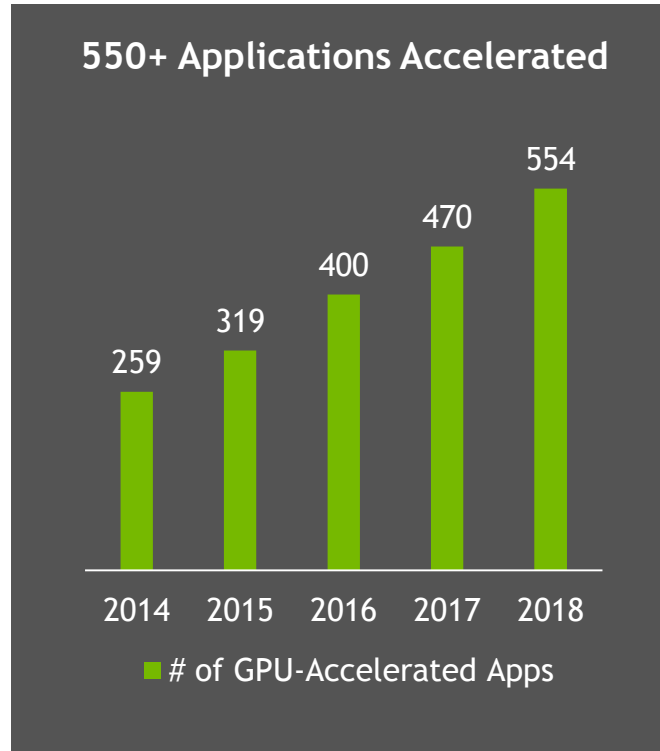
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

# MOST ADOPTED PLATFORM FOR ACCELERATING HPC



**8M**

**1M**

2012          2018

**8X CUDA DOWNLOADS**

---

## 550+ Applications Accelerated

554

470

400

319

259

2014    2015    2016    2017    2018

■ # of GPU-Accelerated Apps

**ALL TOP 15 APPLICATIONS ACCELERATED**

---



**OAK RIDGE SUMMIT**
US's next fastest supercomputer
200+ Petaflop HPC; 3+ Exaflop of AI

**ABCI Supercomputer (AIST)**
Japan's fastest AI supercomputer

**Piz Daint**
Europe's fastest supercomputer

**DEFINING THE NEXT GIANT WAVE IN HPC**

# GPU-ACCELERATED HPC APPLICATIONS

## 550+ APPLICATIONS

### LIFE SCIENCES

**50+ app**

Including:
- Gaussian
- VASP
- AMBER
- HOOMD-Blue
- GAMESS

### MFG, CAD, & CAE

**111 apps**

Including:
- Ansys Fluent
- Abaqus SIMULIA
- AutoCAD
- CST Studio Suite

### PHYSICS

**25 apps**

Including:
- QUDA
- MILC
- GTC-P

### OIL & GAS

**18 apps**

Including:
- RTM
- SPECFEM 3D

### CLIMATE & WEATHER

**3 apps**

Including:
- Cosmos
- Gales
- WRF

### DEEP LEARNING

**38 apps**

Including:
- Caffe2
- MXNet
- Tensorflow

### MEDIA & ENT.

**142 apps**

Including:
- DaVinci Resolve
- Premiere Pro CC
- Redshift Renderer

### FEDERAL & DEFENSE

**14 apps**

Including:
- ArcGIS Pro
- EVNI
- SocetGXP

### DATA SCI. & ANALYTICS

**23 apps**

Including:
- MapD
- Kinetica
- Graphistry

### SAFETY & SECURITY

**19 apps**

Including:
- Cyllance
- FaceControl
- Syndex Pro

### COMP. FINANCE

**16 apps**

Including:
- O-Quant Options Pricing
- MUREX
- MISYS

### TOOLS & MGMT.

**16 apps**

Including:
- Bright Cluster Manager
- HPCtoolkit
- Vampir

# 70% OF THE WORLD'S SUPERCOMPUTING WORKLOAD ACCELERATED

GROMACS
ANSYS Fluent
Gaussian
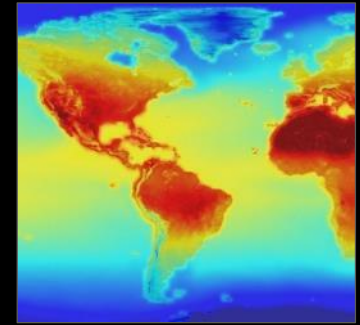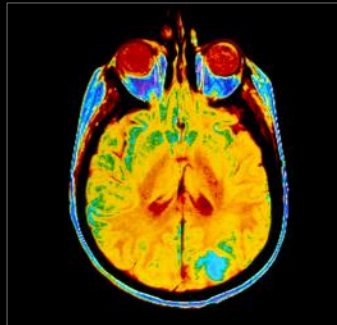VASP
NAMD
Simula Abaqus
WRF
OpenFOAM
ANSYS
LS-DYNA
NCBI-BLAST
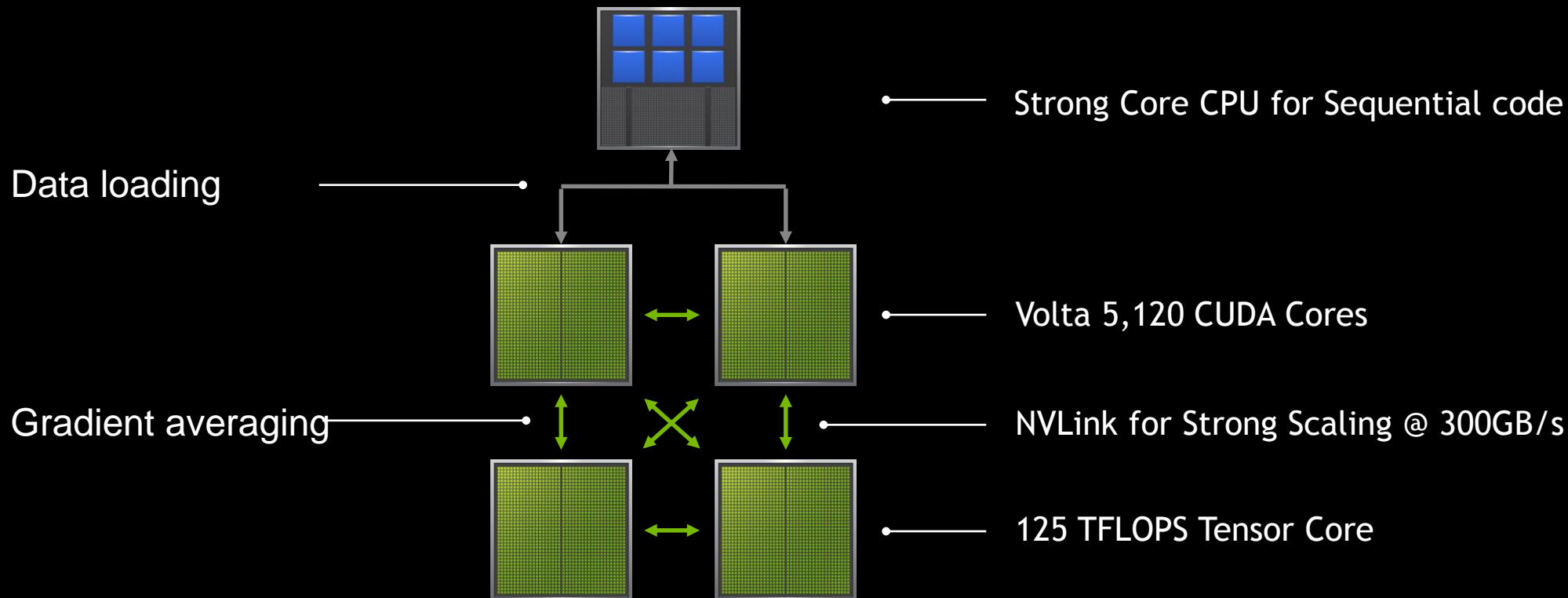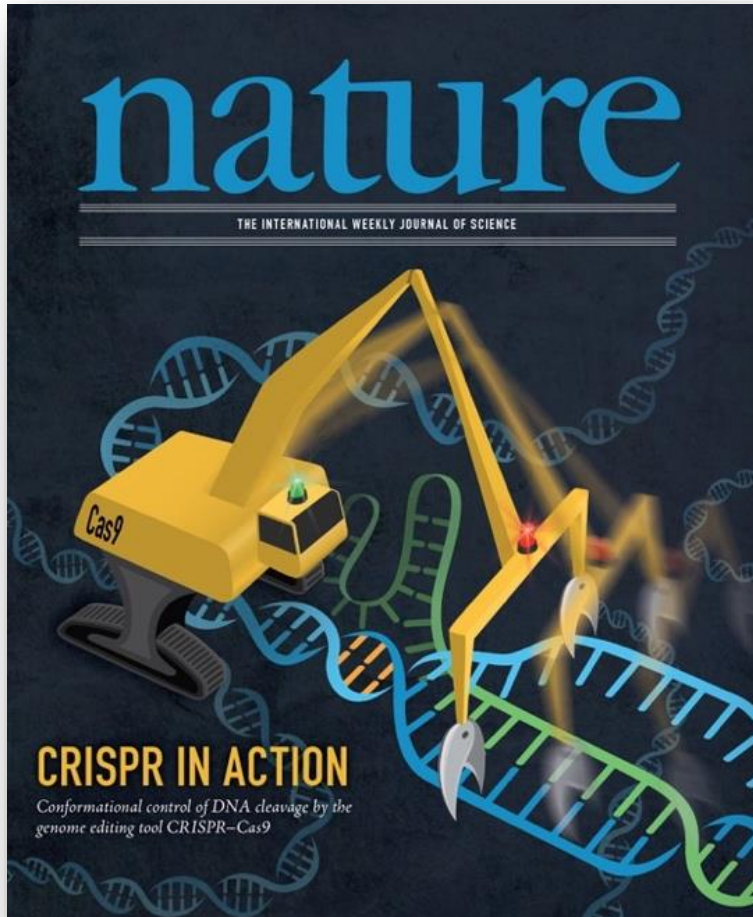LAMMPS
AMBER
Quantum Espresso
GAMESS

Top 15 HPC Applications

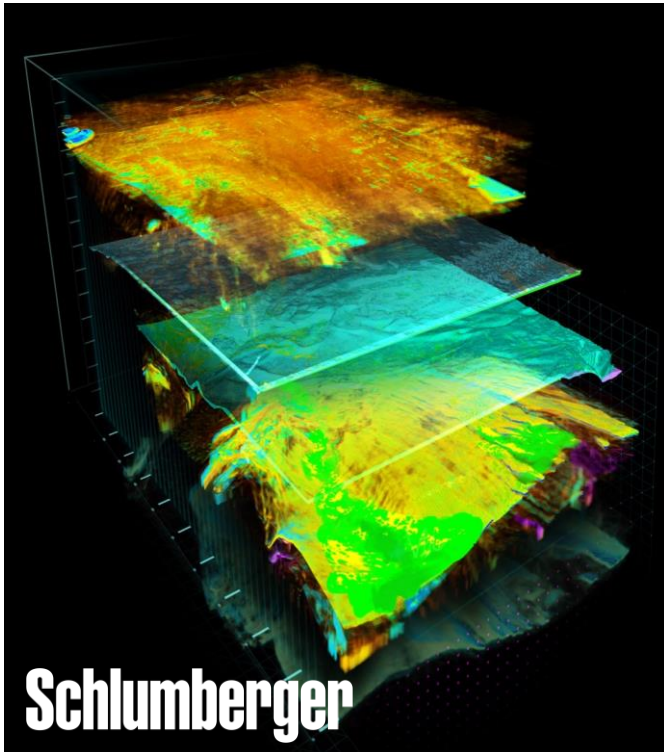500+ Accelerated Applications

# ARCHITECTING MODERN DATACENTERS
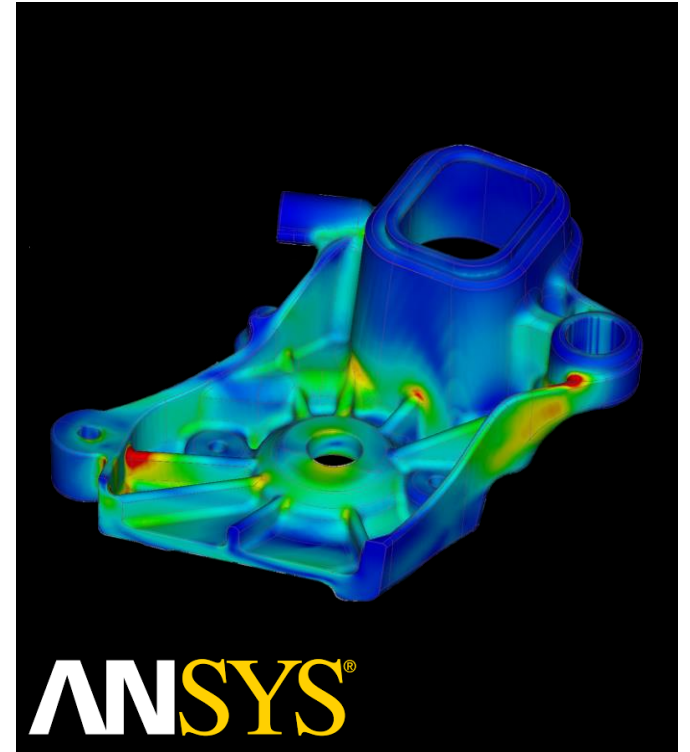
# THE POWER OF ACCELERATED COMPUTING

# INDUSTRY EMBRACING GPU SUPERCOMPUTING



**Schlumberger**

**OIL AND GAS DISCOVERY**
10X increase in data processing
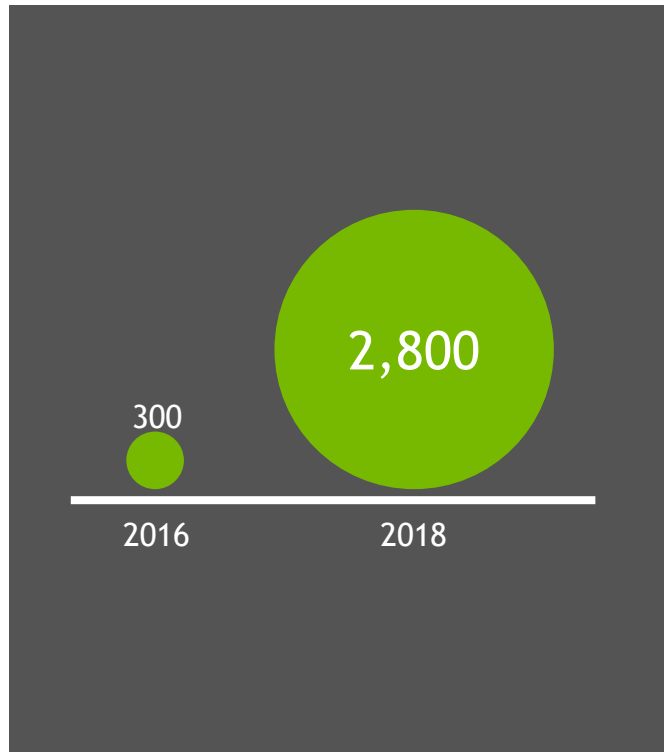
**UNITED STATES POSTAL SERVICE**

**REALTIME FLEET ANALYTICS**
Streamline routes to save >$28M

**ANSYS®**

**ENGINEERING DESIGN**
Accelerate from hours to minutes

# MOST ADOPTED PLATFORM FOR ACCELERATING AI



2016: 300
2018: 2,800

**9X STARTUPS ENGAGED VIA INCEPTION PROGRAM**

Caffe2 · Chainer · KALDI · Microsoft Cognitive Toolkit · mxnet · PaddlePaddle · PYTORCH · TensorFlow

**EVERY DEEP LEARNING FRAMEWORK ACCELERATED**

Alibaba Cloud aliyun.com · aws · Google Cloud · IBM Cloud · Microsoft Azure · Tencent Cloud

Cloud Services

DELL · Hewlett Packard Enterprise · IBM · inspur · Lenovo · SUPERMICRO

Systems

TITAN V

Desktops

**AVAILABLE EVERYWHERE**

# INTELLIGENT HPC

## DL Driving Future HPC Breakthroughs

- Trained networks as solvers
- Super-resolution of coarse simulations
- Low- and mixed-precision
- Simulation for training, network in production
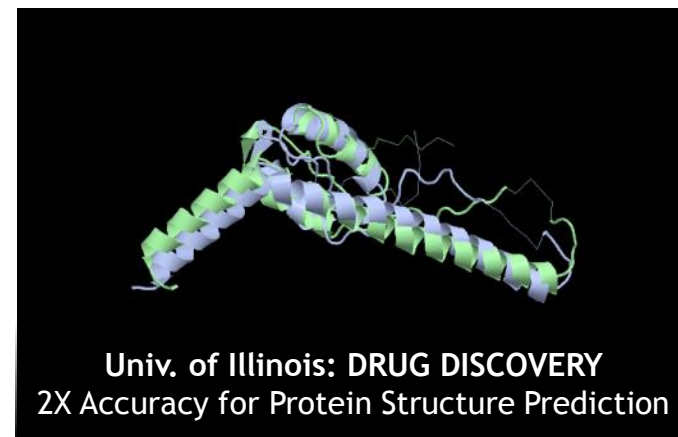
From calendar time to real time?

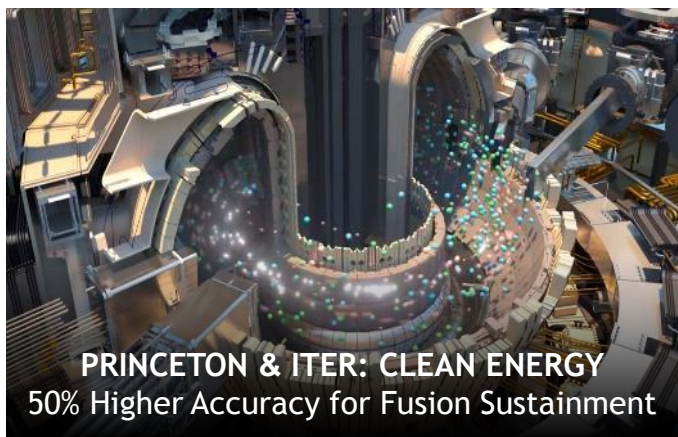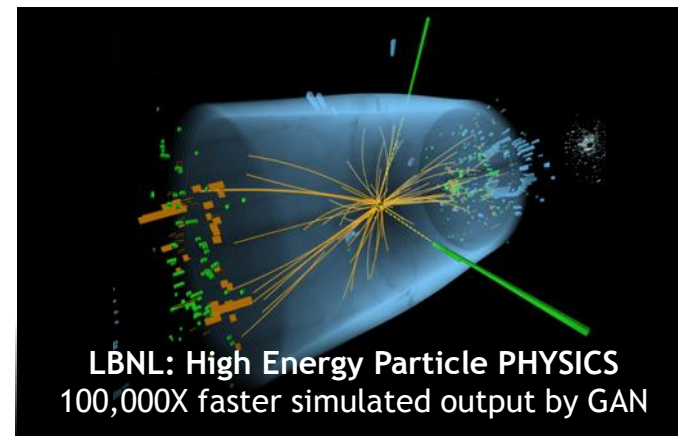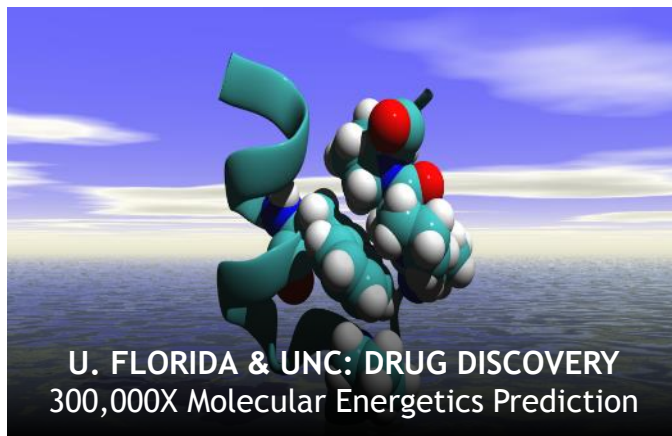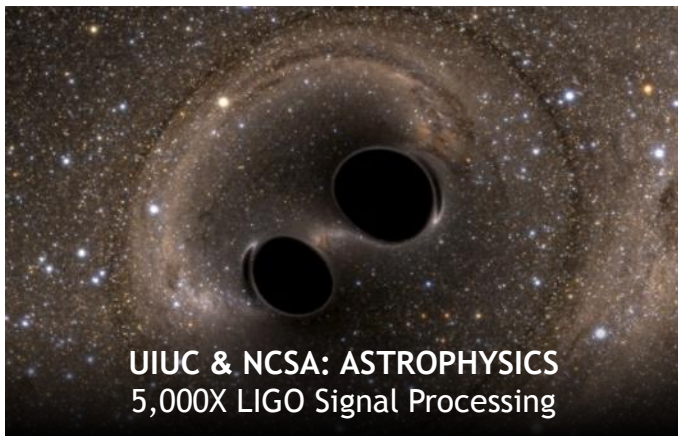**Pre-processing** → **Simulation** → **Post-processing**

- Select/classify/augment/distribute input data
- Control job parameters

- Analyze/reduce/augment output data
- Act on output data

NVIDIA

# DEEP LEARNING COMES TO HPC

## Accelerates Scientific Discovery



**UIUC & NCSA: ASTROPHYSICS**
5,000X LIGO Signal Processing

**U. FLORIDA & UNC: DRUG DISCOVERY**
300,000X Molecular Energetics Prediction

**LBNL: High Energy Particle PHYSICS**
100,000X faster simulated output by GAN

**PRINCETON & ITER: CLEAN ENERGY**
50% Higher Accuracy for Fusion Sustainment

**U.S. DoE: PARTICLE PHYSICS**
33% More Accurate Neutrino Detection

**Univ. of Illinois: DRUG DISCOVERY**
2X Accuracy for Protein Structure Prediction

# ONE PLATFORM BUILT FOR BOTH
# DATA SCIENCE & COMPUTATIONAL SCIENCE



CUDA

Tesla Platform



Accelerating AI



Accelerating HPC

# 4X BETTER HPC SYSTEM TCO

**Mixed Workload:**
Materials Science (VASP)
Life Sciences (AMBER)
Physics (MILC)
Deep Learning (ResNet-50)

160 Self-hosted Servers
96 KWatts

# 4X BETTER HPC SYSTEM TCO



**Mixed Workload:**
Materials Science (VASP)
Life Sciences (AMBER)
Physics (MILC)
Deep Learning (ResNet-50)

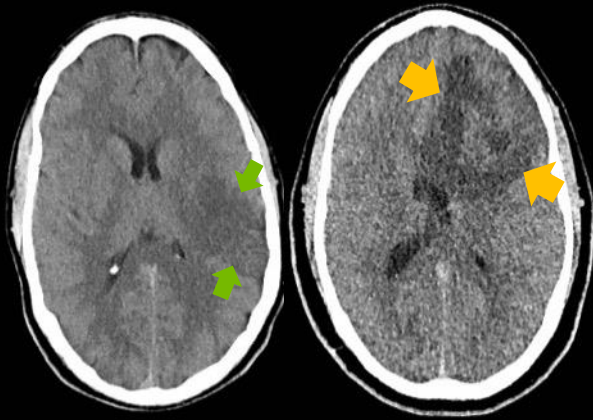12 Accelerated Servers w/4 V100 GPUs
20 KWatts

1/3 the Cost
1/4 the Space
1/5 the Power

# CUSTOMERS WANT MORE

# AI TO TRANSFORM EVERY INDUSTRY



**HEALTHCARE**

>80% Accuracy & Immediate Alert to Radiologists

**INFRASTRUCTURE**
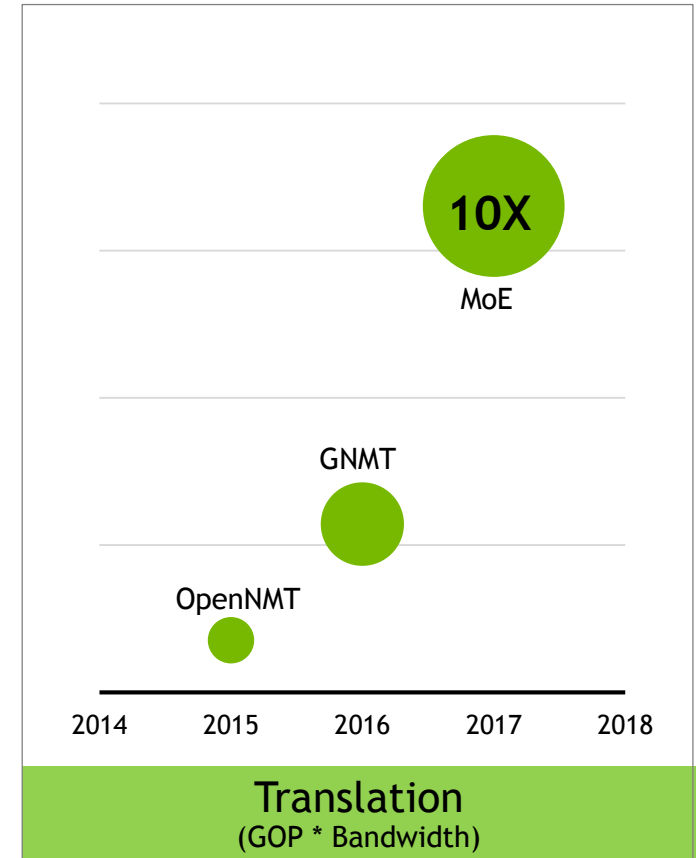
50% Reduction in Emergency Road Repair Costs

**IOT**

7HA GAS TURBINE

>$6M / Year Savings and Reduced Risk of Outage

# NEURAL NETWORK COMPLEXITY IS EXPLODING

## Bigger and More Compute Intensive



**Image**
(GOP * Bandwidth)

350X — Inception-v4
ResNet-50
GoogleNet
AlexNet
Inception-v2

2011 2012 2013 2014 2015 2016 2017

**Speech**
(GOP * Bandwidth)

30X — DeepSpeech 3
DeepSpeech 2
DeepSpeech

2013 2014 2015 2016 2017 2018

**Translation**
(GOP * Bandwidth)

10X — MoE
GNMT
OpenNMT

2014 2015 2016 2017 2018

# TESLA V100 32GB

**WORLD'S MOST ADVANCED DATA CENTER GPU**
**NOW WITH 2X THE MEMORY**

5,120 CUDA cores
**640 NEW** Tensor cores
7.8 FP64 TFLOPS | 15.7 FP32 TFLOPS | 125 Tensor TFLOPS
20MB SM RF | 16MB Cache
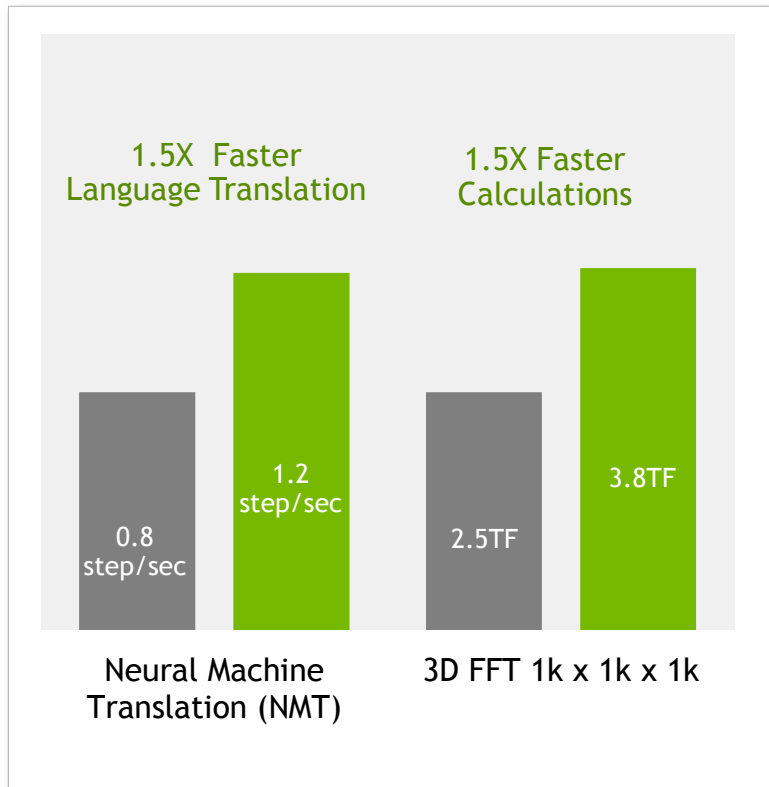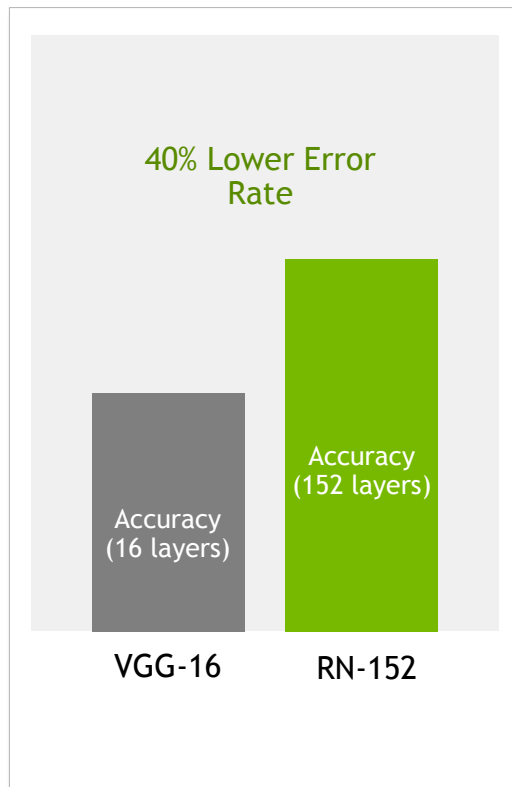**32GB** HBM2 @ 900GB/s | 300GB/s NVLink

# FASTER RESULTS ON COMPLEX DL AND HPC
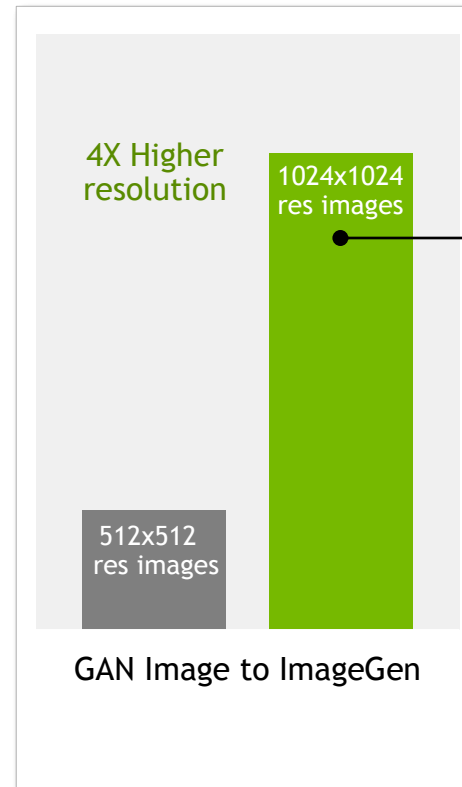
## Up to 50% Faster Results With 2x The Memory

### FASTER RESULTS

**1.5X Faster Language Translation**

1.2 step/sec

0.8 step/sec

Neural Machine Translation (NMT)

**1.5X Faster Calculations**

3.8TF

2.5TF

3D FFT 1k x 1k x 1k

### HIGHER ACCURACY

**40% Lower Error Rate**

Accuracy (152 layers)

Accuracy (16 layers)

VGG-16          RN-152

### HIGHER RESOLUTION

**4X Higher resolution**

1024x1024 res images

512x512 res images

GAN Image to ImageGen

Unsupervised Image Translation
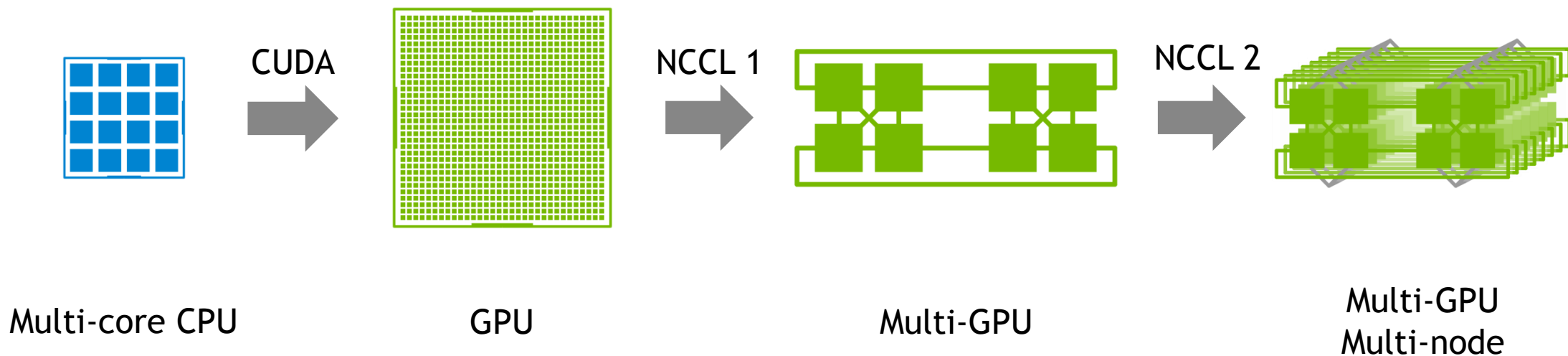
Input winter photo

AI converts it to summer

■ V100 16GB    ■ V100 32GB

⬅ NVIDIA.

# DEEP LEARNING ON GPUS

## Making DL training times shorter

Deeper neural networks, larger data sets ... training is a very, very long operation !

CUDA → NCCL 1 → NCCL 2 →

Multi-core CPU      GPU      Multi-GPU      Multi-GPU Multi-node

# NVLINK MULTI-GPU SCALING



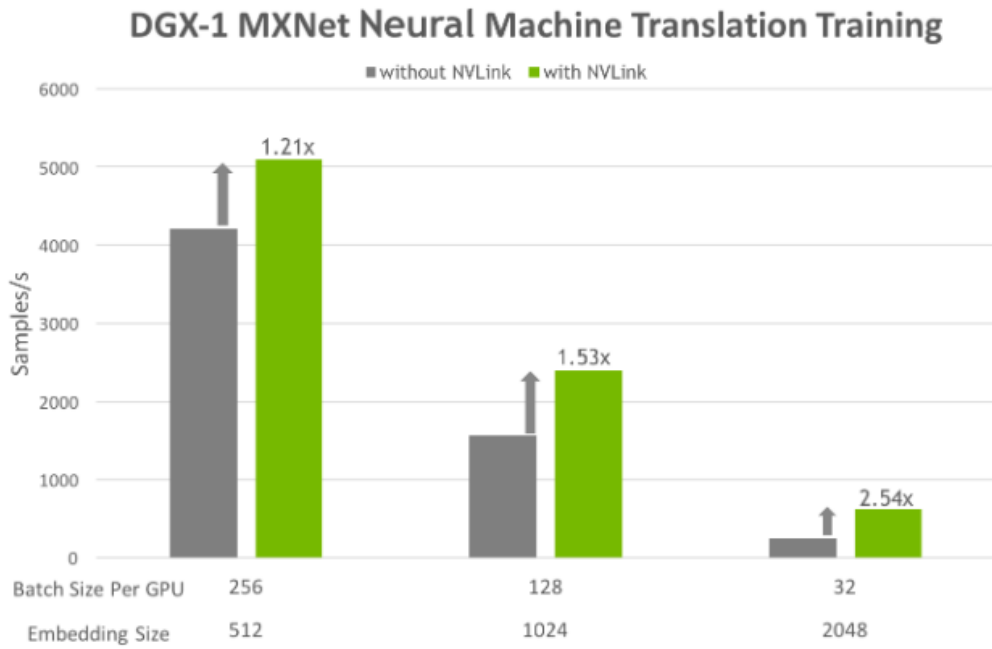**DGX-1 MXNet Neural Machine Translation Training**

Figure 7 Sockeye neural machine translation single-precision training with MXNet using MLP attention on DGX-1, demonstrating significant NVLink performance benefits. The bars present performance on eight Tesla V100 GPUs in a DGX-1 when using NVLink for communication (green), and when using PCIe for communication (gray). Performance benefits increase with the encoder/ decoder embedding size. Results are the average number of samples per second processed during a single epoch of training with the German to English dataset. Tests used NVIDIA DGX MXNet container version 17.11, processing real data with cuDNN 7.0.4, NCCL 2.1.2.
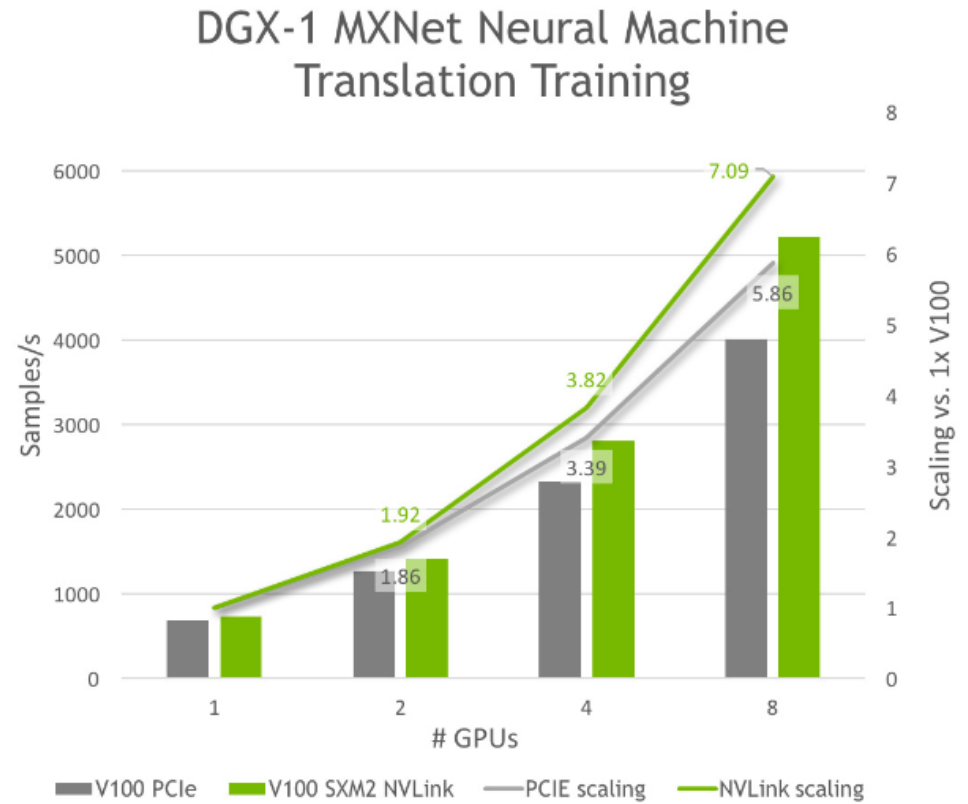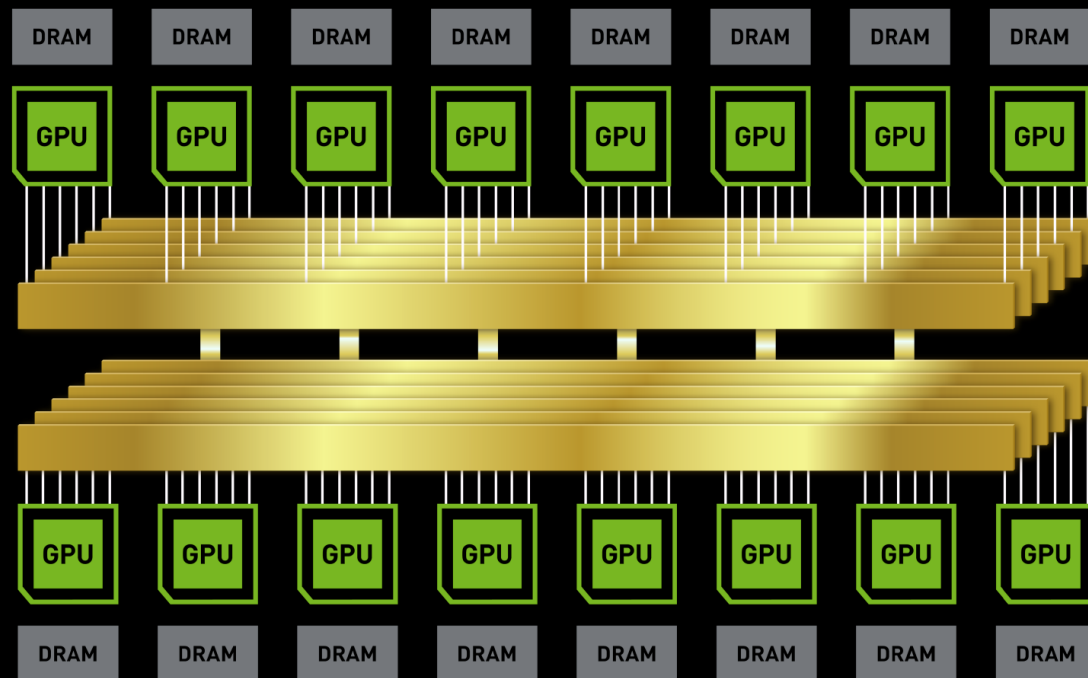
Figure 6 DGX-1 and V100 PCIe performance and scaling for single-precision training of a neural machine translation model with MLP attention and encoder/decoder embedding size of 512 and a batch size of 256 per GPU. The bars show performance on one, two, four, and eight GPUs, comparing an off-the-shelf system of eight Tesla V100 GPUs using PCIe for communication (gray) with eight Tesla V100 GPUs in a DGX-1 using NVLink for communication (green). The lines show the speedup compared to a single GPU. Tests used NVIDIA DGX containers version 17.11, processing real data with cuDNN 7.0.4, NCCL 2.1.2.
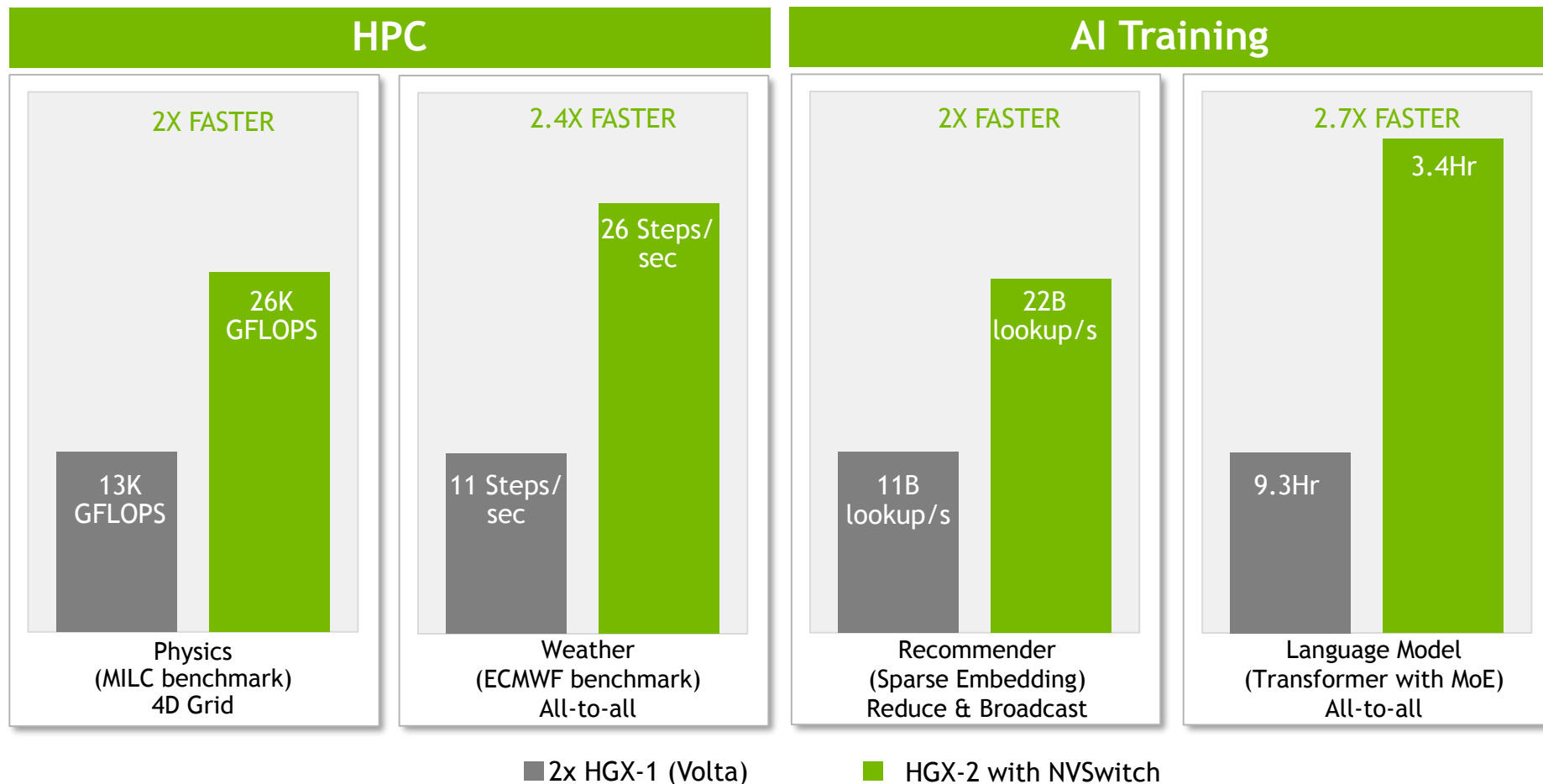
# NVSWITCH

ENABLES THE WORLD'S LARGEST GPU

16 Tesla V100 32GB Connected by New NVSwitch
2 petaFLOPS of DL Compute
Unified 512GB HBM2 GPU Memory Space
300GB/sec Every GPU-to-GPU
2.4TB/sec of Total Cross-section Bandwidth

# UP TO 3X HIGHER PERFORMANCE WITH NVSWITCH



| HPC | | AI Training | |
|---|---|---|---|

**2X FASTER** — Physics (MILC benchmark) 4D Grid: 13K GFLOPS → 26K GFLOPS

**2.4X FASTER** — Weather (ECMWF benchmark) All-to-all: 11 Steps/sec → 26 Steps/sec

**2X FASTER** — Recommender (Sparse Embedding) Reduce & Broadcast: 11B lookup/s → 22B lookup/s

**2.7X FASTER** — Language Model (Transformer with MoE) All-to-all: 9.3Hr → 3.4Hr

Legend: 2x HGX-1 (Volta) | HGX-2 with NVSwitch

2 HGX-1V servers have dual socket Xeon E5 2698v4 Processor. 8 x V100 GPUs. Servers connected via 4XIB ports | HGX-2 server has dual-socket Xeon Platinum 8168 Processor. 16 V100 GPUs
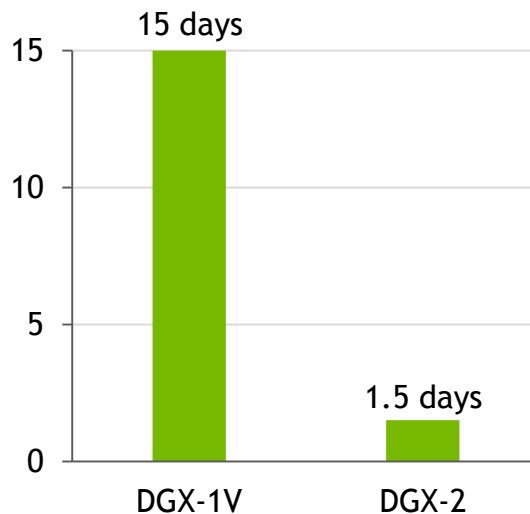
# THE WORLD'S FIRST 2 PETAFLOPS SYSTEM

# 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 simplified operations
- Step up to DGX-2 for unrestricted model parallelism and faster time-to-solution

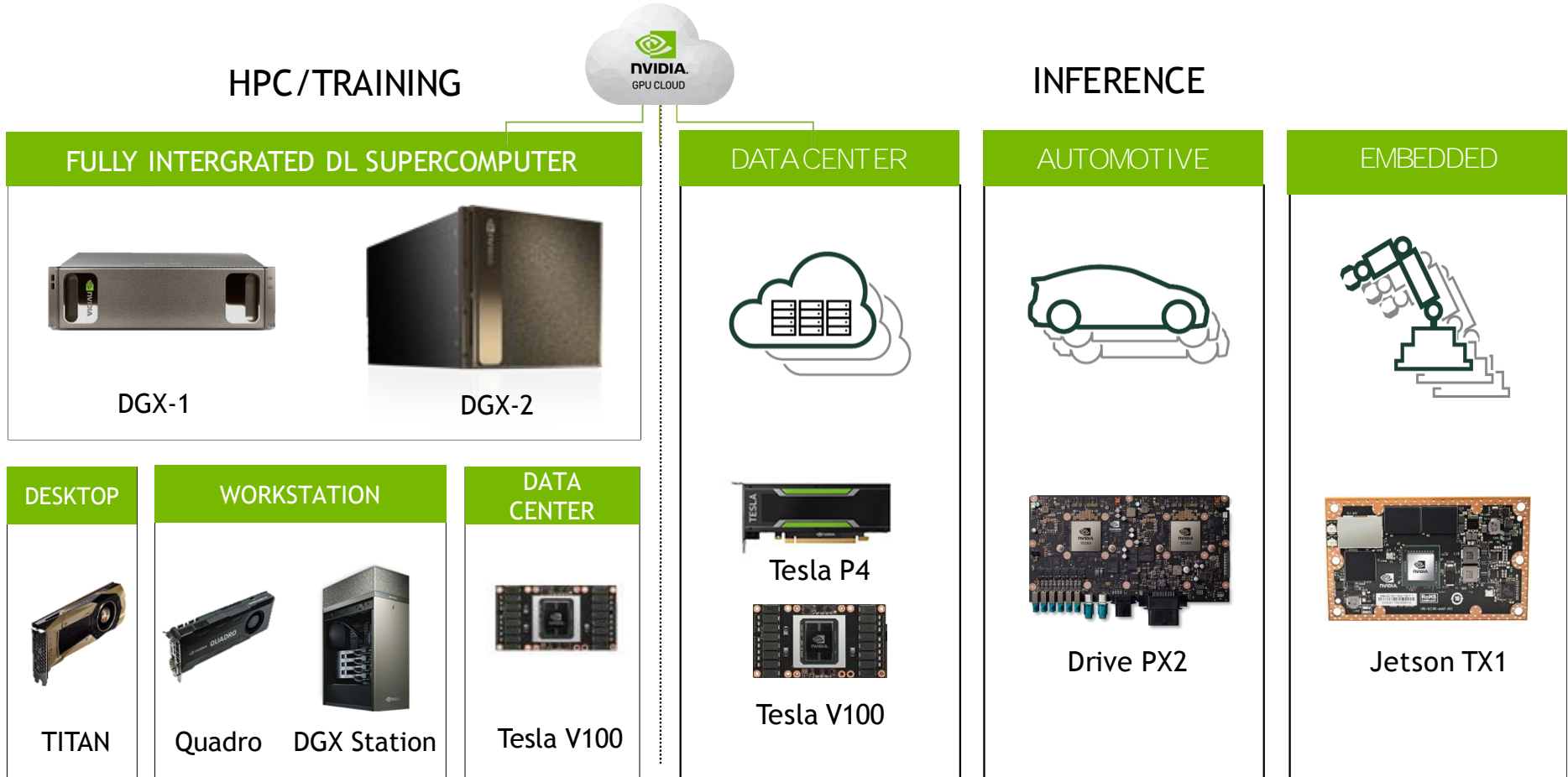# 10X PERFORMANCE GAIN LESS THAN A YEAR

**DGX-1, SEP'17**

**DGX-2, Q3'18**

15 days

1.5 days

DGX-1V

DGX-2

PyTorch Stack: Time to Train FAIRSEQ

software improvements across the stack including NCCL, cuDNN, etc.

# END-TO-END PRODUCT FAMILY

NVIDIA GPU CLOUD

HPC/TRAINING

INFERENCE

## FULLY INTERGRATED DL SUPERCOMPUTER

DGX-1

DGX-2

### DESKTOP

TITAN

### WORKSTATION

Quadro

DGX Station

### DATA CENTER

Tesla V100

## DATA CENTER

Tesla P4

Tesla V100

## AUTOMOTIVE

Drive PX2

## EMBEDDED

Jetson TX1

# TESLA STACK

## World's Leading Data Center Platform for Accelerating HPC and AI

**CUSTOMER USECASES**

Speech · Translate · Recommender

**CONSUMER INTERNET**

Healthcare · Manufacturing · Engineering

**ENTERPRISE APPLICATIONS**

Molecular Simulations · Weather Forecasting · Seismic Mapping

**SUPERCOMPUTING**

**INDUSTRY FRAMEWORKS & APPLICATIONS**

Caffe2 · Chainer · KALDI · Microsoft Cognitive Toolkit

mxnet · PaddlePaddle · PYTORCH · TensorFlow

Amber · ANSYS · CHROMA · GROMACS FAST. FLEXIBLE. FREE.

LAMMPS · NAMD · DS SIMULIA · VASP

+550 Applications

**NVIDIA SDK & LIBRARIES**

cuBLAS · cuDNN · cuFFT · cuRAND · cuSPARSE · DeepStream · NCCL · TensorRT · PGI OpenACC Directives for Accelerators

**CUDA**

**TESLA GPUs & SYSTEMS**

TESLA GPU · NVIDIA DGX STATION · NVIDIA DGX · NVIDIA HGX-1 · DELL · Hewlett Packard Enterprise · IBM · SYSTEM OEM · aws · Google Cloud Platform · Microsoft Azure · CLOUD

# NVIDIA GPU CLOUD
# SIMPLIFYING AI & HPC

**Cloud container registry for GPU accelerated apps**

Containerized in NVDocker

Optimized for GPU-accelerated Systems

Up-to-Date Containers

Available NOW

Sign up at nvidia.com/gpu-cloud

DEEP LEARNING          HPC APPS          HPC VIZ

# HPC APPS CONTAINERS ON NVIDIA GPU CLOUD
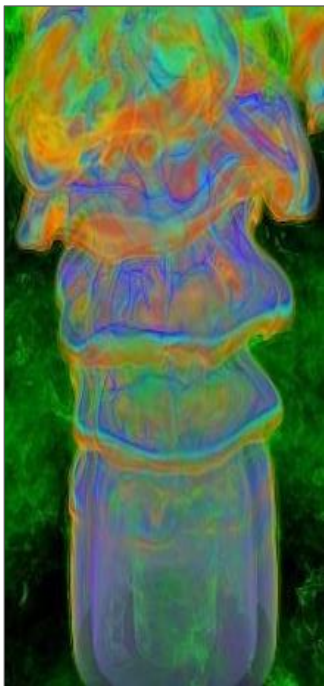


CANDLE

CHROMA*

MILC*

LATTICE MICROBES

GAMESS

GROMACS

LAMMPS

NAMD

RELION

DKRZ DEUTSCHES KLIMARECHENZENTRUM

JOHNS HOPKINS UNIVERSITY

KAUST

京都大学 KYOTO UNIVERSITY

MONASH University

Pfizer

PennState

Schlumberger

Stanford University

UNIVERSITY OF CAMBRIDGE

東京大学 THE UNIVERSITY OF TOKYO

YONSEI UNIVERSITY

**RAPID CONTAINER ADDITION**

**RAPID USER ADOPTION**

*Coming soon

# NVIDIA GPU CLOUD FOR HPC VISUALIZATION



ParaView with
NVIDIA IndeX

ParaView with
NVIDIA OptiX

ParaView with
NVIDIA Holodeck

VMD

IndeX

NEW CONTAINERS

# TESLA PLATFORM FOR DEVELOPERS

# NVIDIA SDK

The Essential Resource for GPU Developers

## DEEP LEARNING

**Deep Learning SDK**

High-performance tools and libraries for deep learning

## AUTONOMOUS VEHICLES

**NVIDIA DRIVE Platform**

Deep learning, HD mapping and supercomputing solutions, from ADAS to fully autonomous

## VIRTUAL REALITY

**NVIDIA VRWorks™**

A comprehensive SDK for VR headsets, games and professional applications

## GAME DEVELOPMENT

**NVIDIA GameWorks™**

Advanced simulation and rendering technology for game development

## ACCELERATED COMPUTING

**NVIDIA ComputeWorks**

Everything scientists and engineers need to build GPU-accelerated applications

## DESIGN & VISUALIZATION

**NVIDIA DesignWorks™**

Tools and technologies to create professional graphics and advanced rendering applications

## AUTONOMOUS MACHINES

**NVIDIA JetPack™**

Powering breakthroughs in autonomous machines, robotics and embedded computing

## SMART CITIES

**NVIDIA Metropolis**

Edge-to-cloud development platform for smart cities

# HOW GPU ACCELERATION WORKS

**Application Code**

Compute-Intensive Functions

5% of Code

Rest of Sequential
CPU Code

**GPU**

**CPU**

+

# HOW TO START WITH GPUS



1. Review available GPU-accelerated applications

2. Check for GPU-Accelerated applications and libraries

3. Add OpenACC Directives for quick acceleration results and portability

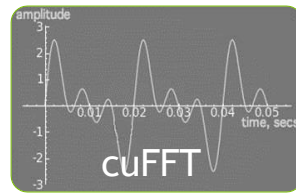4. Dive into CUDA for highest performance and flexibility

36

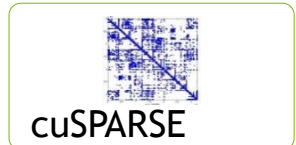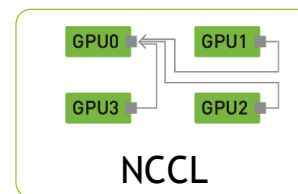# GPU ACCELERATED LIBRARIES

## "Drop-in" Acceleration for Your Applications

### DEEP LEARNING

cuDNN

TensorRT

DeepStream SDK

### SIGNAL, IMAGE & VIDEO

amplitude

cuFFT

NVIDIA NPP

CODEC SDK

### LINEAR ALGEBRA

cuBLAS

cuSPARSE

CUDA Math library

cuSOLVER

cuRAND

### PARALLEL ALGORITHMS

nvGRAPH

GPU0    GPU1
GPU3    GPU2

NCCL

Thrust

# WHAT IS OPENACC

## Programming model for an easy onramp to GPUs

Directives-based programming model for **parallel computing**

Add Simple Compiler Directive

```
main()
{
  <serial code>
  #pragma acc kernels
  {
    <parallel code>
  }
}
```

Designed for **performance portability** on CPUs and GPUs

**SIMPLE**

**POWERFUL & PORTABLE**

Read more at  www.openacc.org/about

OpenACC is an open specification developed by OpenACC.org consortium

NVIDIA.

# SINGLE CODE FOR MULTIPLE PLATFORMS
## OpenACC - Performance Portable Programming Model for HPC

OpenPOWER

Sunway

x86 CPU

x86 Xeon Phi

NVIDIA GPU

AMD GPU

PEZY-SC

AWE Hydrodynamics CloverLeaf mini-App, bm32 data set
http://uk-mac.github.io/CloverLeaf



Legend:
- PGI 18.1 OpenACC (green)
- Intel 2018 OpenMP (dark gray)

Y-axis: Speedup vs Single Haswell Core

Data points:
- Multicore Haswell: 7.6x, 7.9x
- Multicore Broadwell: 10x, 10x
- Multicore Skylake: 14.8x, 15x
- Kepler Pascal: 11x, 40x
- 1x Volta V100: 67x
- 2x Volta V100: 109x
- 4x Volta V100: 142x

NVIDIA

# OPENACC GROWING MOMENTUM

## Wide Adoption Across Key HPC Codes

### Over 100 Apps* Using OpenACC

| | |
|---|---|
| ANSYS Fluent | GTC |
| Gaussian | XGC |
| VASP | ACME |
| LSDalton | FLASH |
| MPAS | COSMO |
| GAMERA | Numeca |

**VASP**

*Top Quantum Chemistry and Material Science Code*

" For VASP, OpenACC is ***the*** way forward for GPU acceleration. Performance is similar to CUDA, and OpenACC dramatically decreases GPU development and maintenance efforts. We're excited to collaborate with NVIDIA and PGI as an early adopter of Unified Memory. "

*Prof. Georg Kresse*
*Computational Materials Physics*
*University of Vienna*

* Applications in production and development

NVIDIA.

# OPENACC.ORG RESOURCES

Guides ● Talks ● Tutorials ● Videos ● Books ● Spec ● Code Samples ● Teaching Materials ● Events ● Success Stories ● Courses ● Slack ● Stack Overflow

## OpenACC
## Now in GCC

https://www.openacc.org/community#slack

### Resources
https://www.openacc.org/resources

### Success Stories
https://www.openacc.org/success-stories

### Compilers and Tools
https://www.openacc.org/tools

### Events
https://www.openacc.org/events

# PGI COMPILERS FOR EVERYONE
## The PGI 17.10 Community Edition

| | **FREE** PGI Community EDITION | PGI Professional EDITION | PGI Enterprise EDITION |
|---|---|---|---|
| **PROGRAMMING MODELS** OpenACC, CUDA Fortran, OpenMP, C/C++/Fortran Compilers and Tools | ✔ | ✔ | ✔ |
| **PLATFORMS** X86, OpenPOWER, NVIDIA GPU | ✔ | ✔ | ✔ |
| **UPDATES** | 1-2 times a year | 6-9 times a year | 6-9 times a year |
| **SUPPORT** | User Forums | PGI Support | PGI Premier Services |
| **LICENSE** | Annual | Perpetual | Volume/Site |

**PGI**

⊙ **NVIDIA**

# NVIDIA ACADEMIC PROGRAMS

| EDUCATION | | GPU ACCESS | |
|---|---|---|---|
| **NVIDIA Teaching kits** | **DLI Ambassador Hands-on training** | **To start / teach GeForce, Jetson** | **To scale DGX, Tesla, GRID, Cloud** |



Fundamentals

Training — forward — "dog" / labels — =? — "human face"
Large N — backward — error
Inference
Smaller, varied N — forward — "human face"

Autonomous Vehicles      Healthcare      Robotics

**Complete course solutions:**
1. Deep Learning
2. Accelerated Computing
3. Robotics

Both Fundamentals and Industry-Specific Labs Available

Free for DLI Ambassador

Perfect platform to prototype

Accelerating your ideas
Enabling
new scientific publication

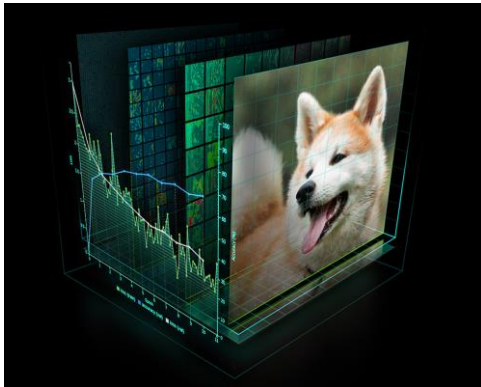Get an academic discount on Tesla, DGX, Jetson TX2 or GRID from our NPN partners

https://developer.nvidia.com/academia

# GPU TECHNOLOGY CONFERENCE

**Oct 9 -11, 2018 | Munich**
**www.gputechconf.eu** **#GTC18**
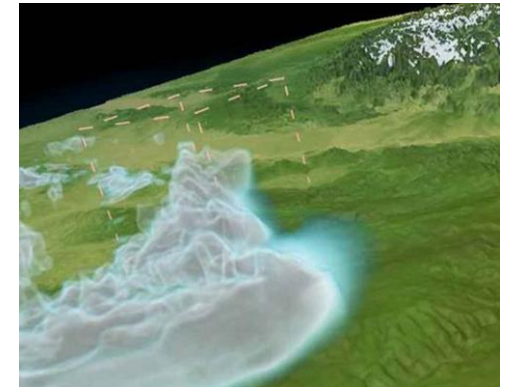
# EUROPE'S BRIGHTEST MINDS & BEST IDEAS

DEEP LEARNING &
ARTIFICIAL INTELLIGENCE

SELF-DRIVING CARS

VIRTUAL REALITY &
AUGMENTED REALITY

SUPERCOMPUTING & HPC

Join us in Munich and discover the latest breakthroughs in autonomous vehicles, HPC, smart cities, healthcare, big data, virtual reality, and more.

3 Days | 3000 Attendees | 80+ Exhibitors | 100+ Speakers | 10+ Tracks | 10+ Hands-on labs | 1-to-1 Meetings
-25% with the promo code NVGBARAT