# ALISON B LOWNDES

AI DevRel | EMEA

*@alisonblowndes*

**October 2018**

**NVIDIA.**

www.FrontierDevelopmentLab.org

NASA FRONTIER DEVELOPMENT LAB

# TESLA V100 32GB

5,120 CUDA cores
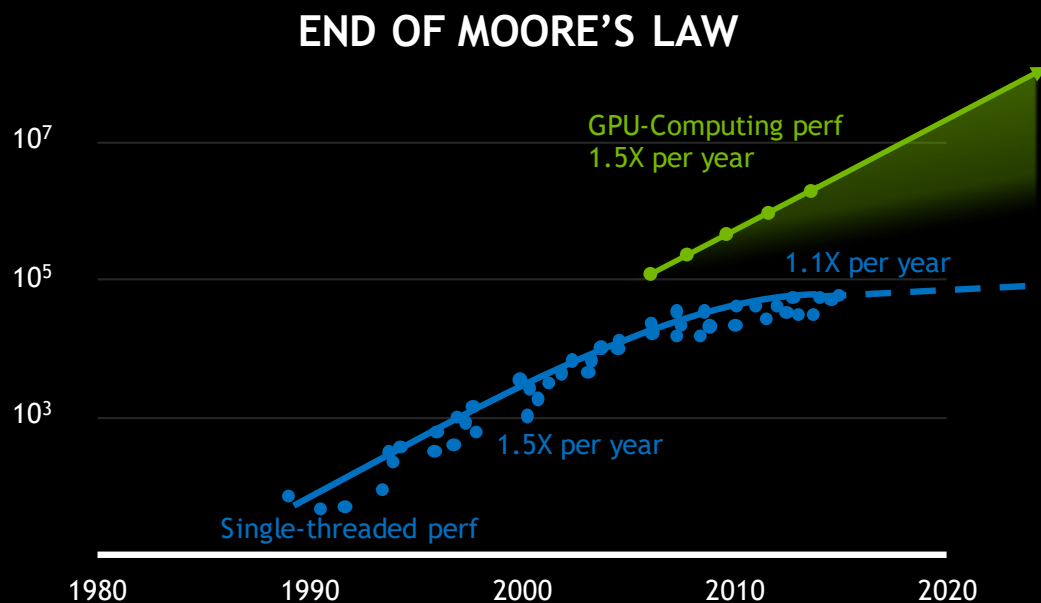**640 NEW** Tensor cores
7.5 FP64 TFLOPS  |  15 FP32 TFLOPS

## 120 Tensor TFLOP

20MB SM RF  |  16MB Cache  |  32GB HBM2 @ 900 GB/s
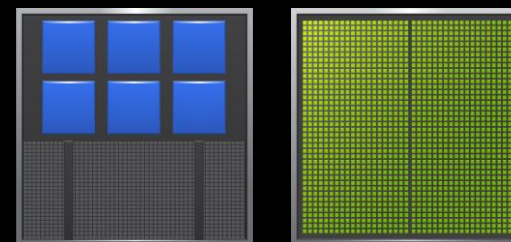300 GB/s NVLink

# THE RISE OF GPU COMPUTING
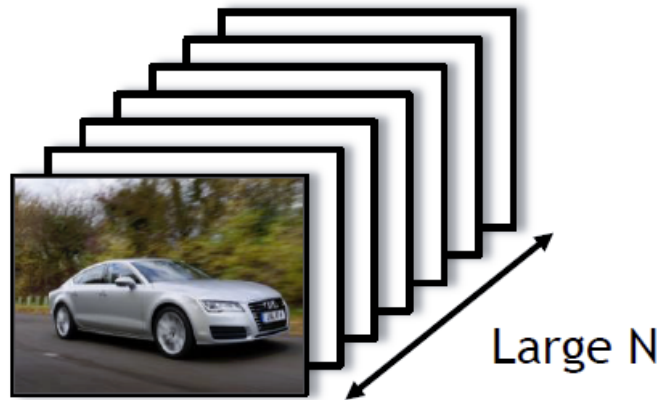
## Big Data Needs Algorithms and Compute That Scales

**END OF MOORE'S LAW**

**CPU vs. GPU**

$10^7$

GPU-Computing perf
1.5X per year

$10^5$

1.1X per year

$10^3$

1.5X per year

Single-threaded perf

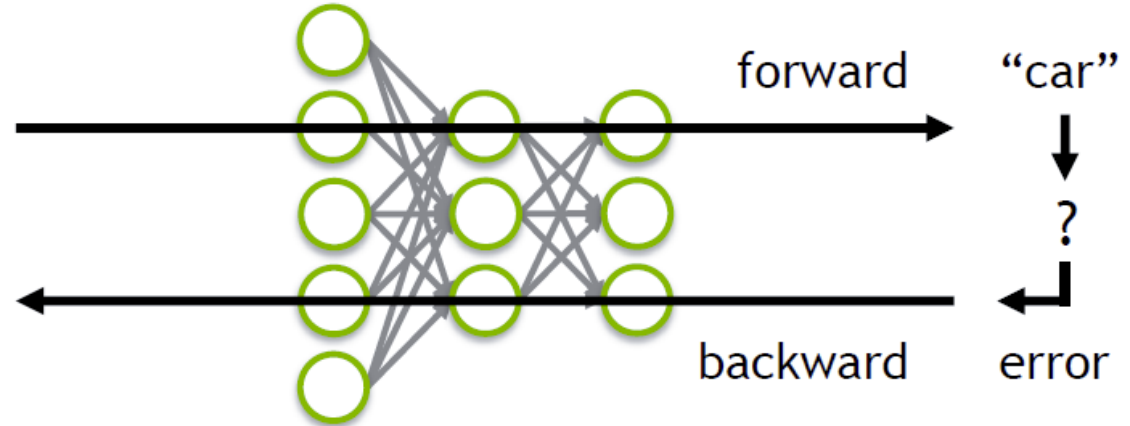1980    1990    2000    2010    2020

Original data up to the year 2010 collected and plotted by M. Horowitz,
F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten Newplot and data collected for 2010-2015 by K. Rupp
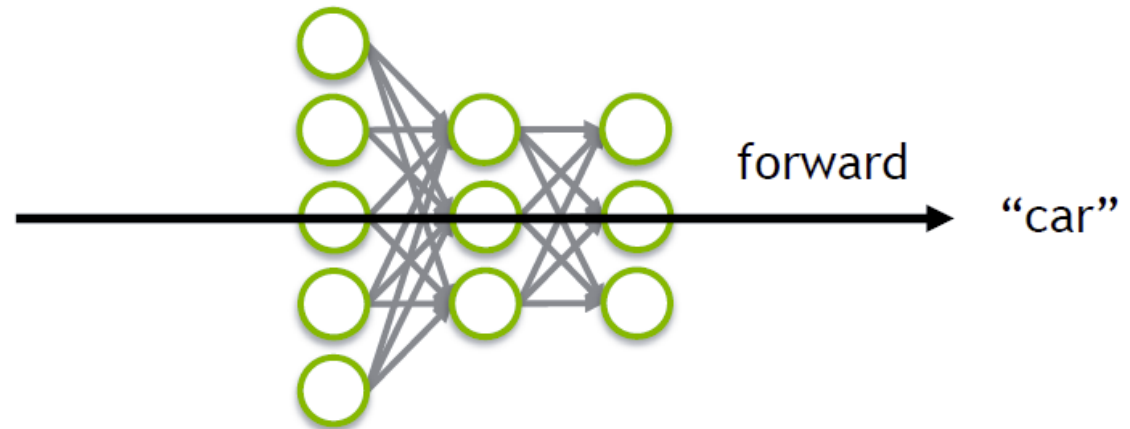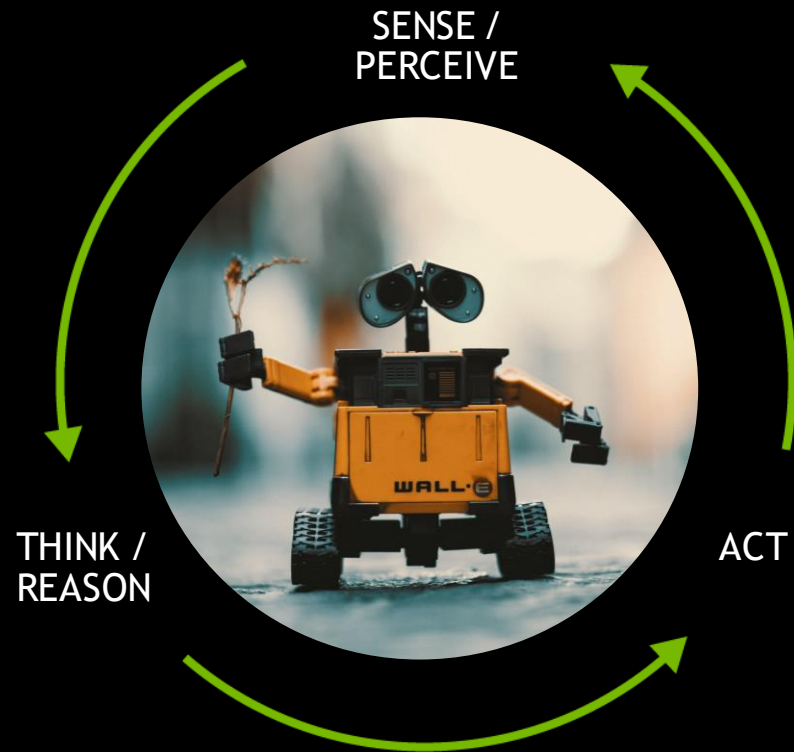
# TRAINING VS INFERENCE



Large N

**TRAINING**

forward "car"

?

backward error

Smaller, varied N

**INFERENCE**

forward "car"

NVIDIA.

# ROBOTS

RL CMU Humanoid
Rigid Terrain
RL Full Humanoid
RL Ant
RL Atlas Flagrun
RL Hard Flagrun
RL Fetch - Rigid
RL Fetch - Rope
RL Fetch - Cloth

Particle Count: 0
Diffuse Count: 0
Shape Match Count: 0
Rigid Body Count: 6500
Rigid Shape Count: 9500
Rigid Joint Count: 12000
Spring Count: 0
Tetra Count: 0
Num Substeps: 4
Num Iterations: 30

Device: TITAN X (Pascal)

Options

Global
☐ Emit particles
⦿ Pause

☐ Wireframe
⦿ Draw Points
☐ Draw Fluid
⦿ Draw Mesh
☐ Draw Basis
☐ Draw Springs
☐ Draw Contacts
☐ Draw Joints

Reset Scene

⦿ Jacobi
☐ LDLT
☐ PCG (CPU)
☐ PCG (GPU)
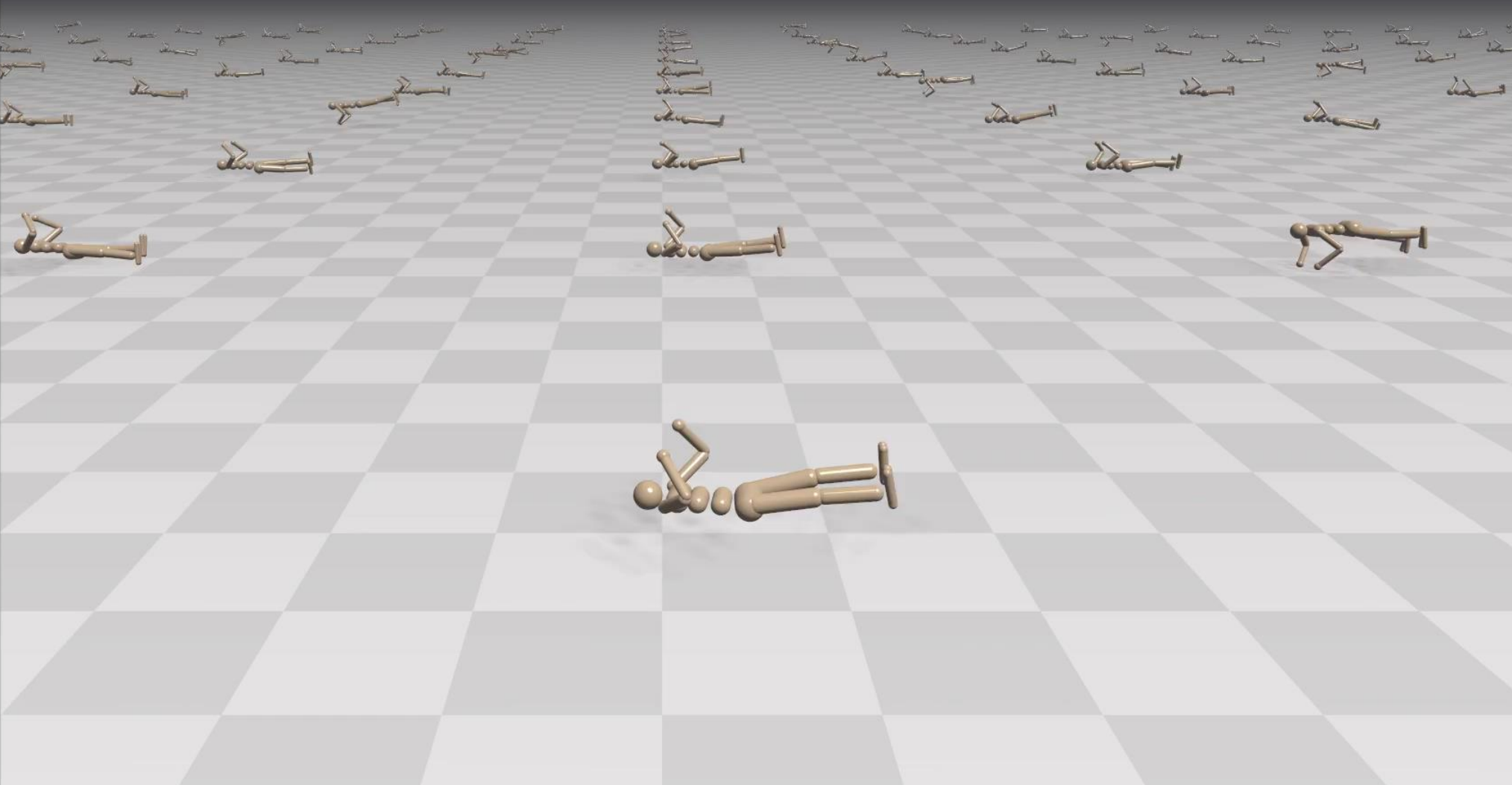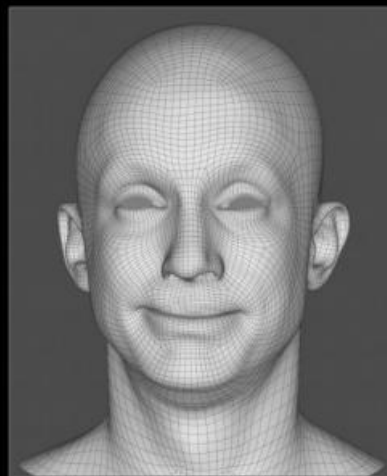
Num Substeps          4
Num Outer Iterations  30
Num Inner Iterations  20

Gravity X             0
Gravity Y           -10
Gravity Z             0

Radius             0.15
Solid Radius      0.150
Fluid Radius      0.000

SOR                1.00
Geometric Stiffness 1.000

# NVIDIA RESEARCH



NVIDIA Research
AI Autoencoder

NVIDIA Research / Remedy
Audio-driven Facial Animation
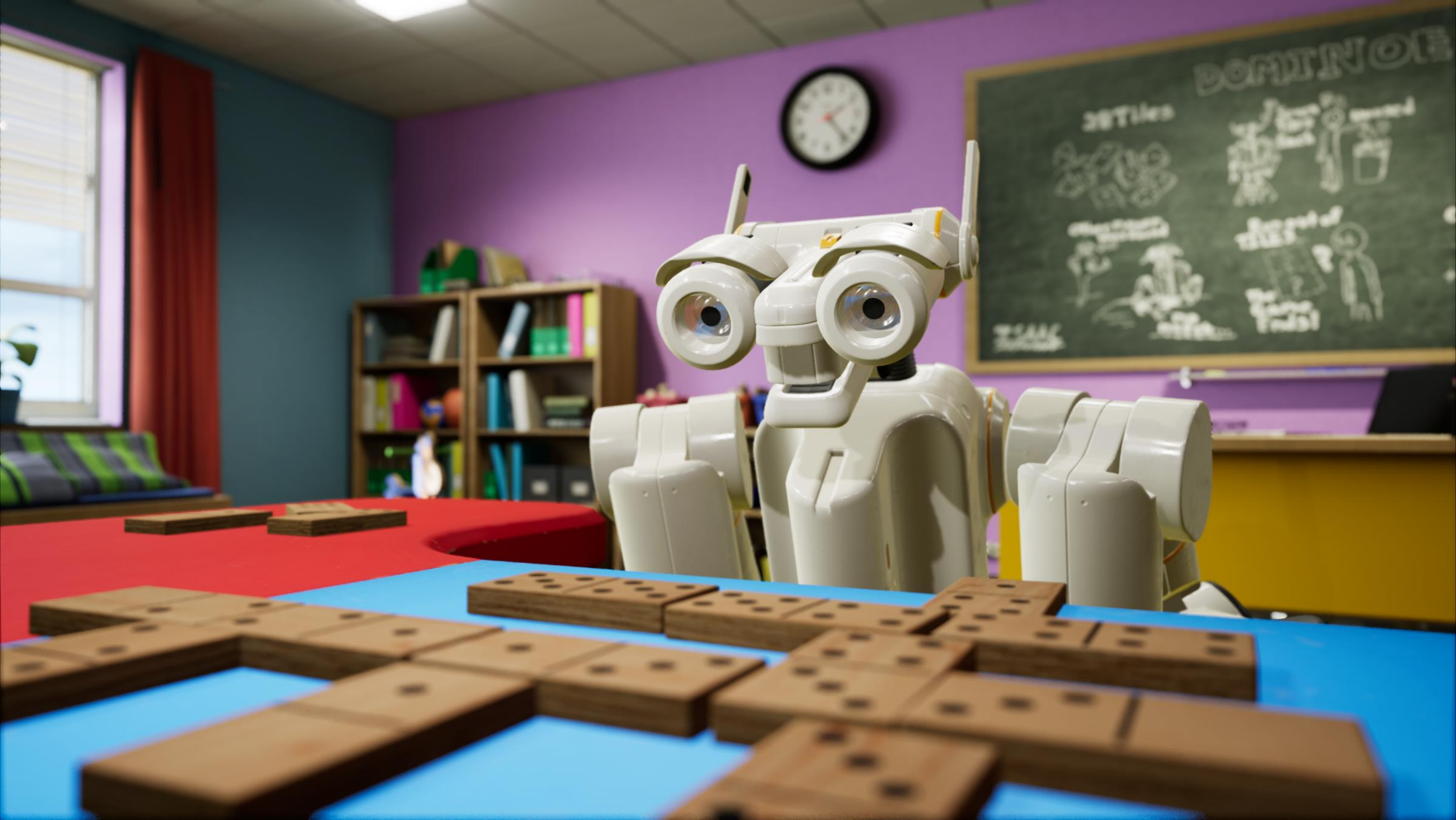
NVIDIA Research
Semantic Manipulation with GANs

NVIDIA Research
Progressive GAN

NVIDIA Research / AIVA
RNNs for Music

# ISAAC

**Simulate** → **Develop** → **Deploy**



| Navigation | Behaviors | Perception | Manipulation | Interactions with humans |



**World model**
Warehouse · Office · Store · Home

**Robot model**
Carter · URDF loader

**ML**
TensorRT · CUDA · Tensorflow · ...

**Gems**
Optimizers · Algebra · EKFs · Depth · ...

**Drivers**
Lidar · Camera · IMU · Robot Base · ...

**Jetson**
Fully integrated with X2 and Xavier

**Simulation Engine**
Photo-realistic Graphics · Physics · Soft bodies · · Procedural Generation · Massive parallelism · Unreal Engine 4 / Unity 3D

**Isaac Framework**
Codelets · Behaviors · 3D Poses · Distributed · Messaging · Synchronization · Record & Replay · Configuration · Visualization

**Unified Message API**
Use the same messages for simulation, actual hardware and across all apps

Virtual Sensors

Virtual Actuators

Sensor Processing

Actuator Control

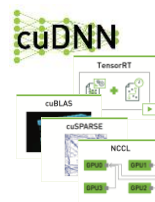HW Sensor

HW Actuator

NVIDIA.

# CUDA DEVELOPMENT ECOSYSTEM

*GPU Users*        *Domain Specialists*        *Problem Specialists*        *New Algorithm Developers and Optimization Experts*

CUDA-C++
CUDA Fortran

GROMACS
FAST. FLEXIBLE. FREE.

VASP    NAMD

ANSYS
v-ray    FLUENT    otoy

cuDNN
TensorRT
cuBLAS
cuSPARSE
NCCL

OpenACC
Directives for Accelerators

Thrust

**Applications**        **Frameworks**        **Libraries**        **Directives and Standard Languages**        **Extended Standard Languages**
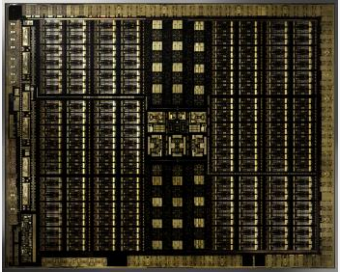
──────▶ *Ease of use* ──────────────────── *Specialized Performance* ──────▶
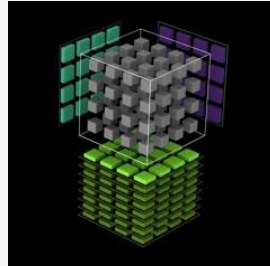
CUDA: Programming Model, GPU Architecture, System Architecture

# CUDA 10 - TURING

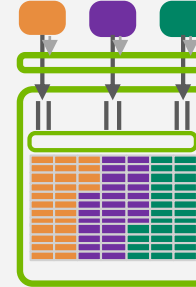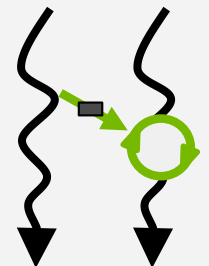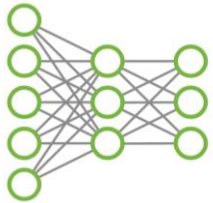| Turing Architecture | Multi-Precision Tensor Core | RT Core | Turing MPS | Independent Thread Scheduling |
|---|---|---|---|---|

**Inference Accelerated, Graphics Reinvented, Volta's Programmability**

# NVIDIA DEEP LEARNING SDK UPDATE

## GPU-accelerated DL Primitives



**cuDNN 7.3**

Support for Turing

Optimizations for RNNs

Leading frameworks support

## Multi-GPU & Multi-node
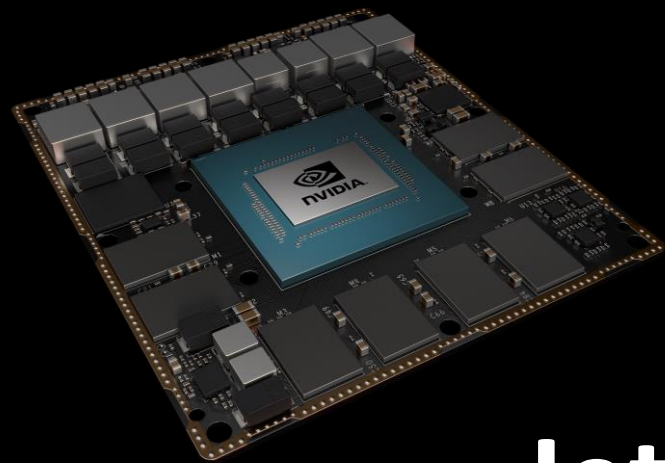


**NCCL 2**

Multi-node distributed training (multiple machines)

Leading frameworks support

## High-performance Inference Engine



**TensorRT 5**

TensorFlow model reader

Object detection

INT8 RNNs support

NVIDIA.
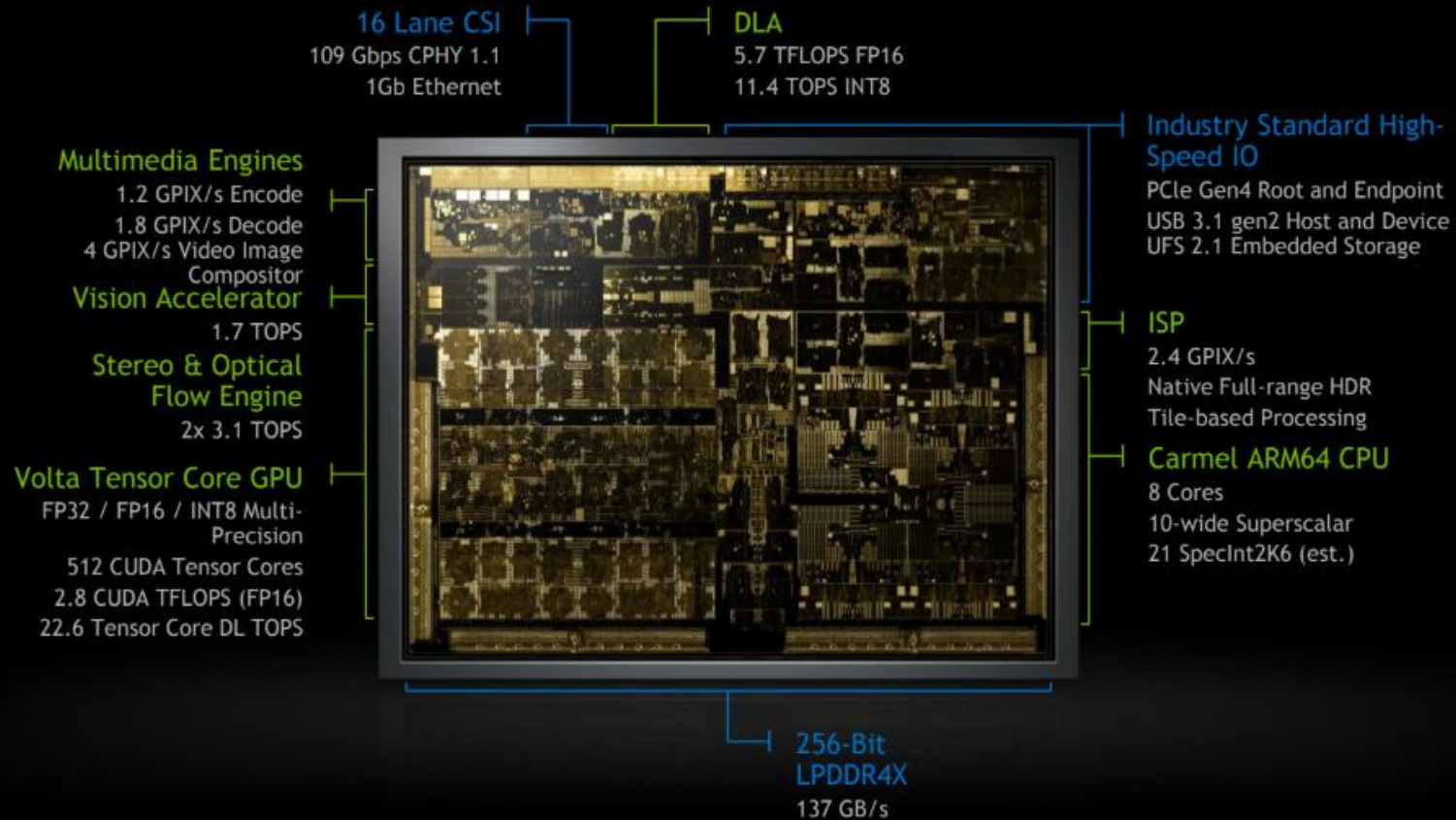
# Jetson
# Xavier

developer.nvidia.com/
jetson-xavier

30W • 15W • 10W
512 Volta CUDA Cores
8 core CPU
32 DL TOPS

NVIDIA.

# XAVIER

## World's First Autonomous Machines Processor



**16 Lane CSI**
109 Gbps CPHY 1.1
1Gb Ethernet

**DLA**
5.7 TFLOPS FP16
11.4 TOPS INT8

**Industry Standard High-Speed IO**
PCIe Gen4 Root and Endpoint
USB 3.1 gen2 Host and Device
UFS 2.1 Embedded Storage

**Multimedia Engines**
1.2 GPIX/s Encode
1.8 GPIX/s Decode
4 GPIX/s Video Image Compositor

**Vision Accelerator**
1.7 TOPS

**Stereo & Optical Flow Engine**
2x 3.1 TOPS

**Volta Tensor Core GPU**
FP32 / FP16 / INT8 Multi-Precision
512 CUDA Tensor Cores
2.8 CUDA TFLOPS (FP16)
22.6 Tensor Core DL TOPS

**ISP**
2.4 GPIX/s
Native Full-range HDR
Tile-based Processing

**Carmel ARM64 CPU**
8 Cores
10-wide Superscalar
21 SpecInt2K6 (est.)

**256-Bit LPDDR4X**
137 GB/s

Most Complex SOC Ever Made | 9 Billion Transistors, 350mm$^2$, 12FFN | ~8,000 Engineering Years
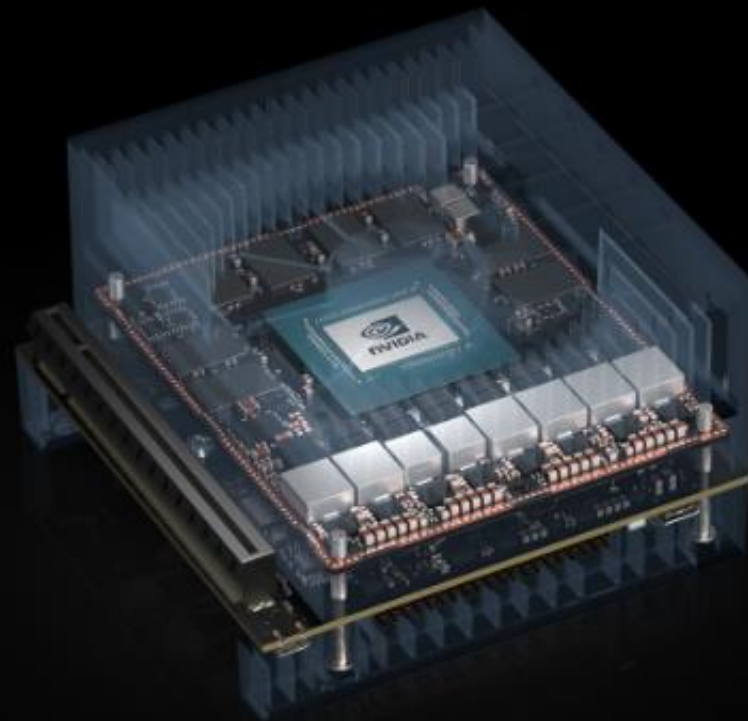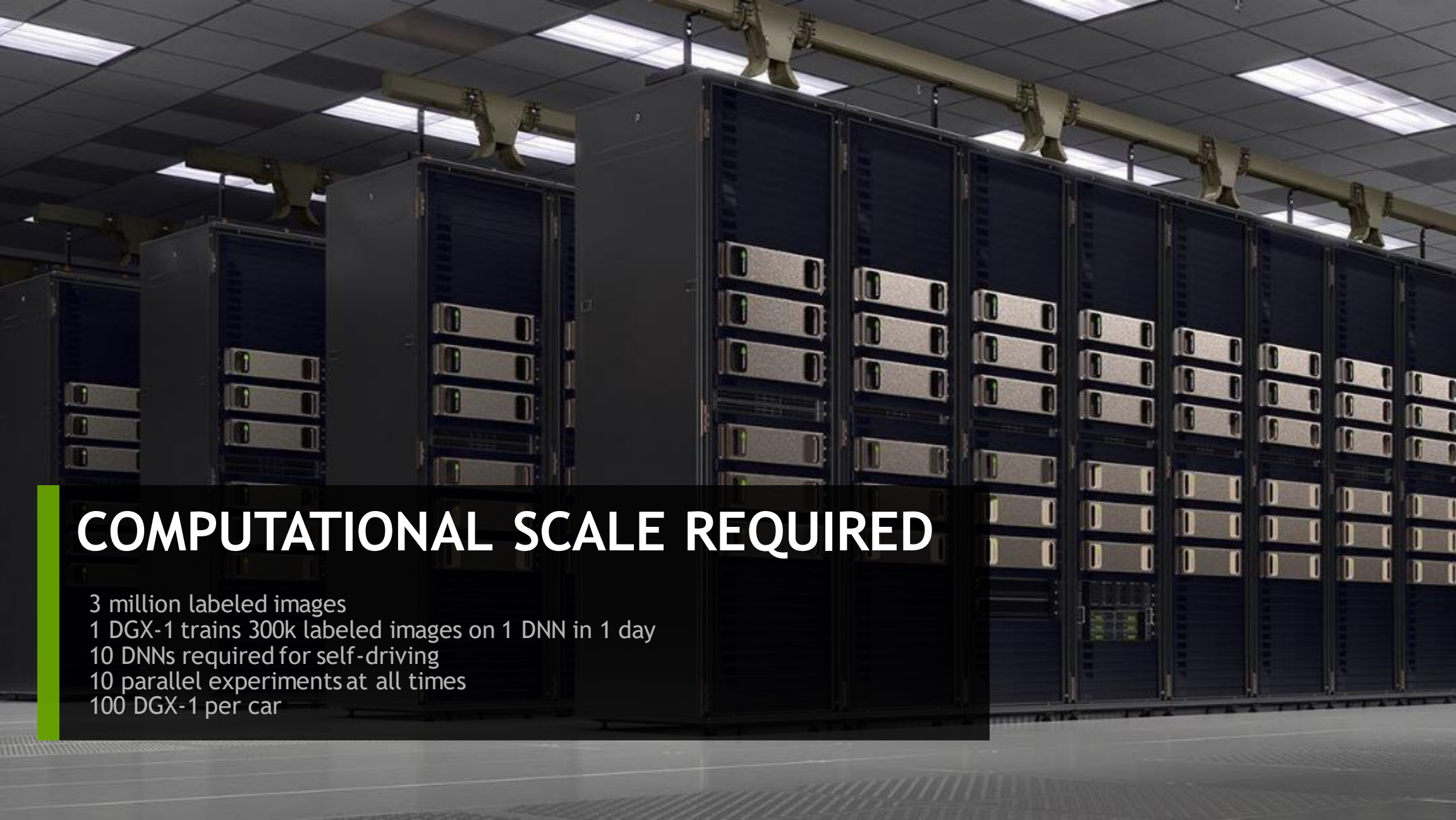
# JETSON XAVIER DEVELOPER KIT

$1299 (pre-order only)

Available from distributors WW

Early access August 2018

# JETSON XAVIER
# DEVELOPER KIT

# COMPUTATIONAL SCALE REQUIRED

3 million labeled images
1 DGX-1 trains 300k labeled images on 1 DNN in 1 day
10 DNNs required for self-driving
10 parallel experiments at all times
100 DGX-1 per car

# TensorRT SUPPORTS Xavier
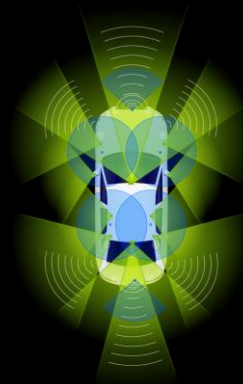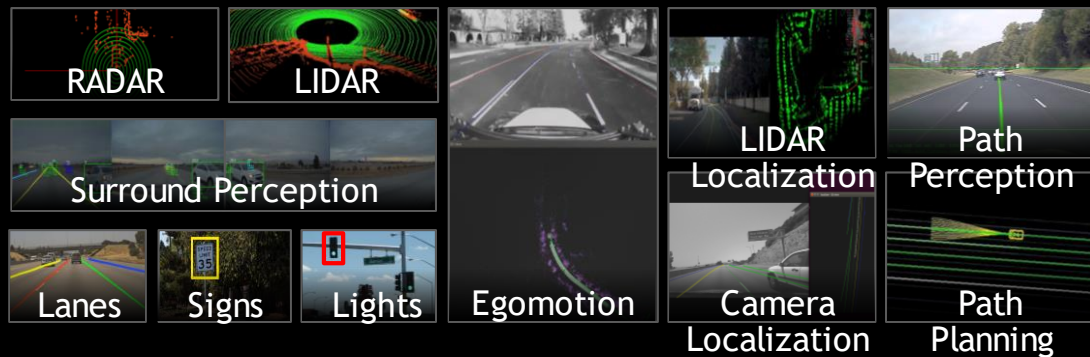
## Optimized Inference on the World's Most Powerful SoC

Deploy deep learning inference on Xavier platforms through NVIDIA DRIVE AI Platforms

Import models in any framework (including TensorFlow, Caffe and Torch) through ONNX, Universal Framework Format or custom C/C++ API

Optimize CNN, RNN and novel neural network layers and deploy reduced precision on Tensor Cores

Download for development or host environment today

developer.nvidia.com/tensorrt

# AI IS THE FUTURE OF INTELLIGENT INSTRUMENTS



AUTOMAP RECONSTRUCTION
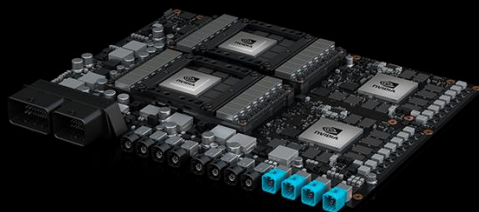
*Source: https://arxiv.org/pdf/1704.08841.pdf*
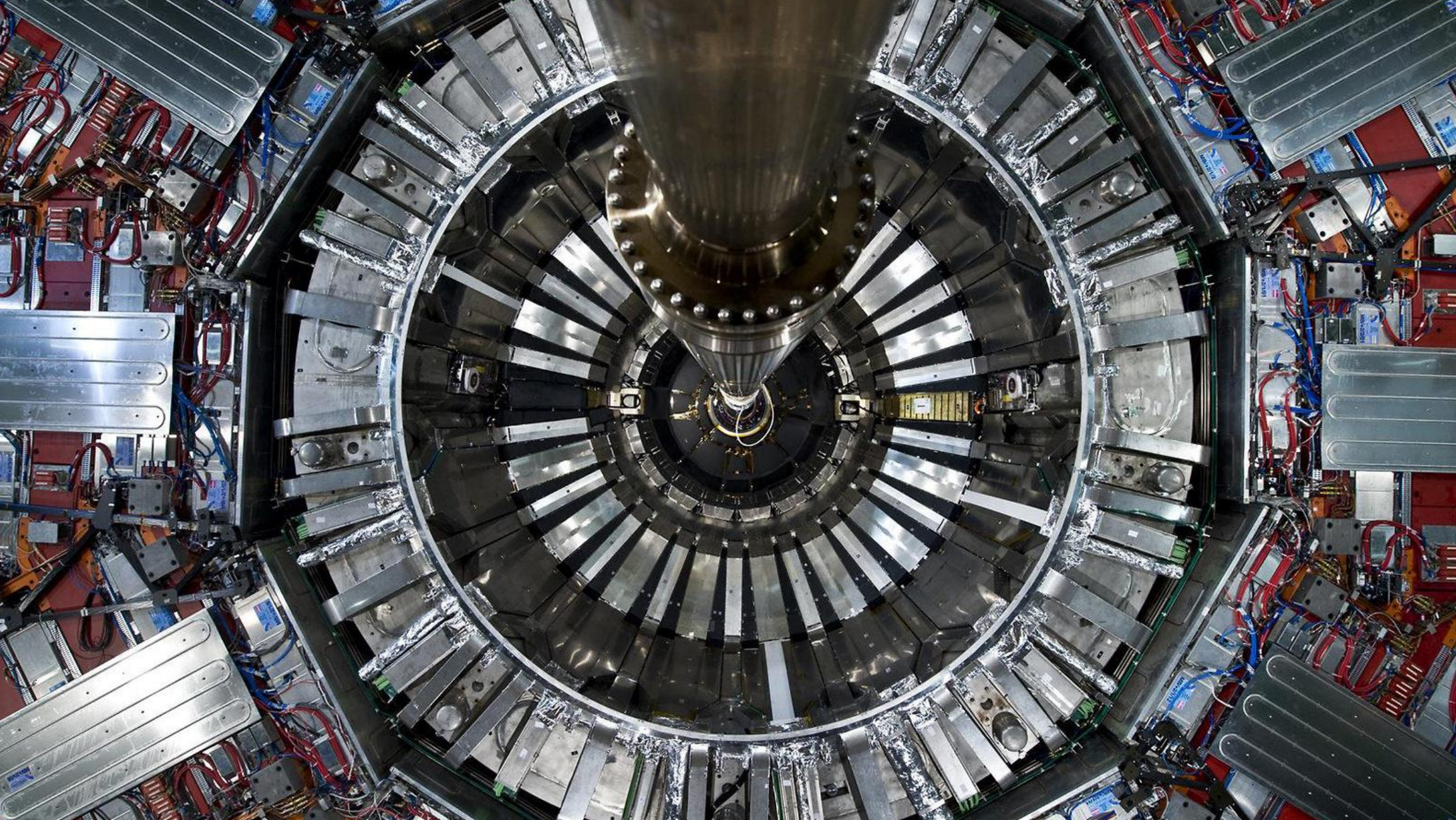


ANATOMY DETECTION

*Source: GE Healthcare*



INTELLIGENT RENDERING

*Source: GE Healthcare*

# NVIDIA POWERS FASTEST SUPERCOMPUTERS IN US, EUROPE, JAPAN, INDUSTRY

## 17 of World's 20 Most Energy-efficient Supercomputers

**ORNL Summit**
**World's Fastest**
**27,648 GPUs| 122 PF**

**LLNL Sierra**
**US 2nd Fastest**
**17,280 GPUs| 72 PF**

**ABCI**
**Japan's Fastest**
**4,352 GPUs| 20 PF**

**Piz Daint**
**Europe's Fastest**
**5,320 GPUs| 20 PF**

**ENI HPC4**
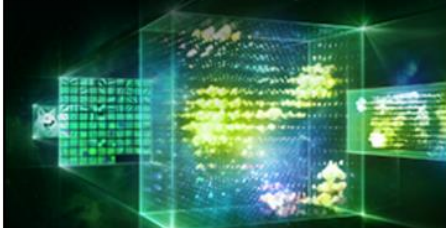**Fastest Industrial**
**3,200 GPUs| 12 PF**

# NVIDIA SDK

**The Essential Resource for GPU Developers**

developer.nvidia.com

## NVIDIA SDK

### DEEP LEARNING

**Deep Learning SDK**
High-performance tools and libraries for deep learning

### SELF-DRIVING CARS

**NVIDIA DriveWorks™**
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

### ADDITIONAL RESOURCES

More resources for GPU Developers

# WHY CONTAINERS?



CONTAINERIZED APPLICATION

DEEP LEARNING APPLICATIONS
DEEP LEARNING FRAMEWORKS
DEEP LEARNING LIBRARIES
CUDA TOOLKIT

MOUNTED NVIDIA DRIVER
CONTAINER OS

CONTAINERIZATION TOOL

NVIDIA CONTAINER RUNTIME FOR DOCKER
DOCKER ENGINE

NVIDIA DRIVER
HOST OS

NVIDIA GPU CLOUD SOFTWARE STACK

## Benefits of Containers:

Simplify deployment of
GPU-accelerated software, eliminating time-
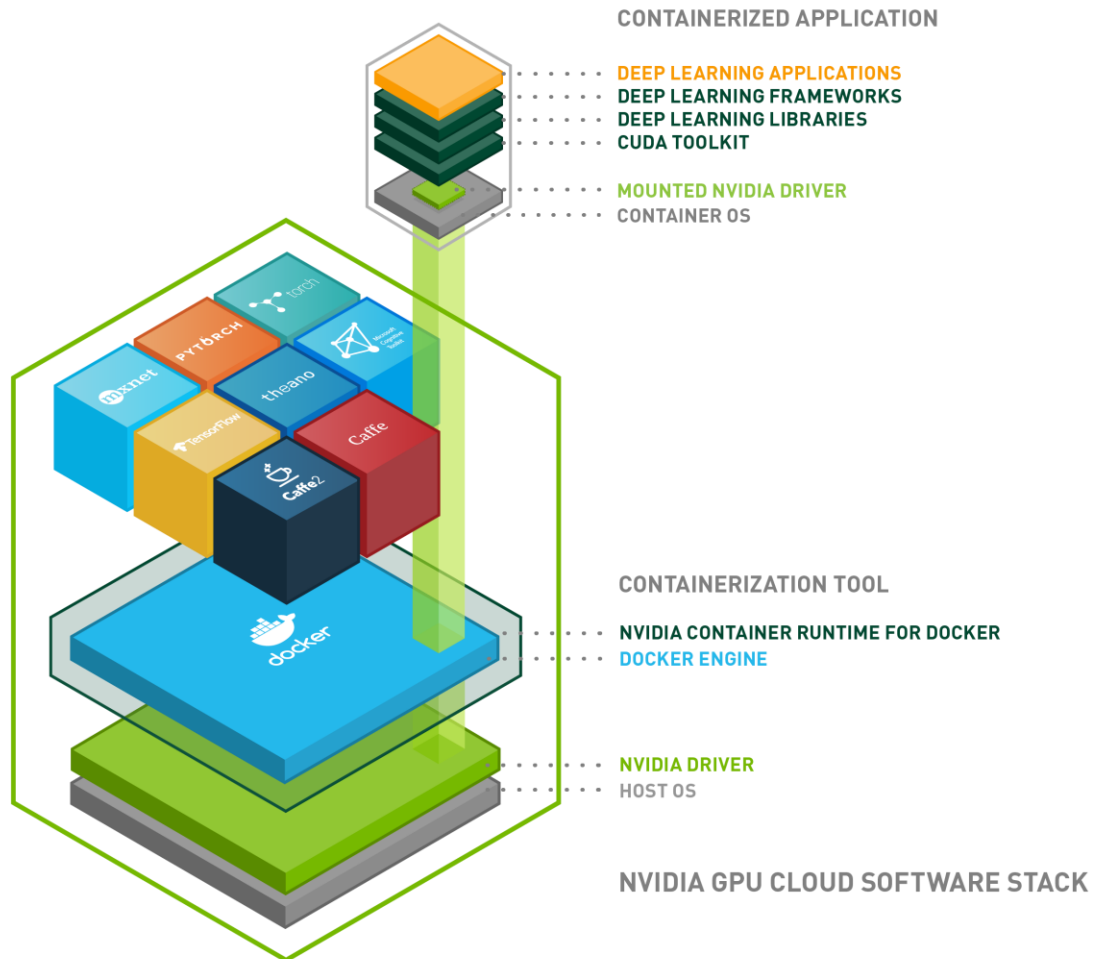consuming software integration work

Isolate individual deep learning frameworks
and applications

Share, collaborate,
and test applications across
different environments

To learn more:
## nvidia.com/ngc
To sign up:
## ngc.nvidia.com

NVIDIA.

# CUDA CONTAINERS ON NVIDIA GPU CLOUD

CUDA containers available from NGC Registry at **nvcr.io/nvidia/cuda**

Three different flavors:

## Base

Contains the minimum components required to run CUDA applications

## Runtime

Contains *base* + CUDA libraries (e.g. cuBLAS, cuFFT)

## Devel

Contains *runtime* + CUDA command line developer tools. Some *devel* tags also include cuDNN

# DGX POD ARCHITECTURE

a single data center rack containing up to 9x NVIDIA DGX-1 servers, storage, networking & NVIDIA AI software



Nine DGX-1 servers

12 storage servers

10 GbE (min) storage & management switch

Mellanox 100 Gpps intra-rack high speed network switches.

# NVIDIA DGX-2

## THE LARGEST GPU EVER CREATED

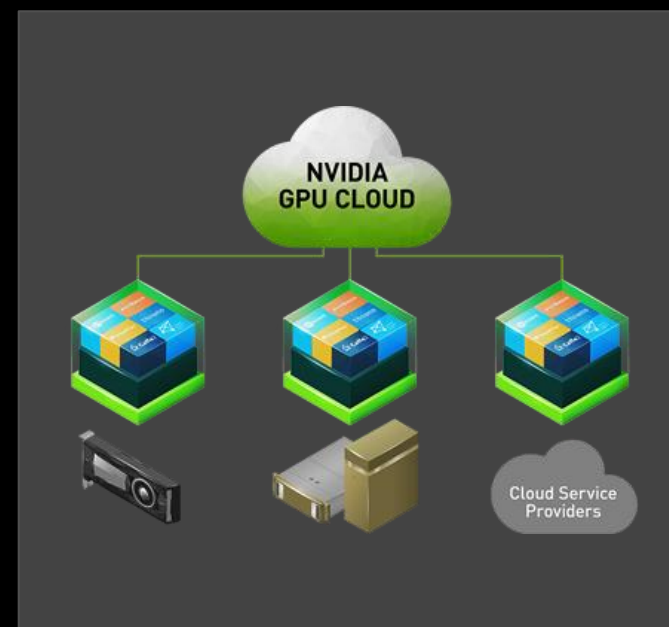2 PFLOPS | 512GB HBM2 | 10 kW | 350 lbs

# NEXT STEPS



**GTC Munich | October 9-11 2018**
www.nvidia.com/en-eu/gtc/
25% discount: NVALOWNDES

**NVIDIA Deep Learning Institute**
www.nvidia.com/en-us/deep-learning-ai/education

**NGC**
www.nvidia.com/en-us/gpu-cloud

# the esa earth observation φ-week

EO Open Science and FutureEO

12–16 November 2018 | ESA–ESRIN | Frascati (Rome), Italy