

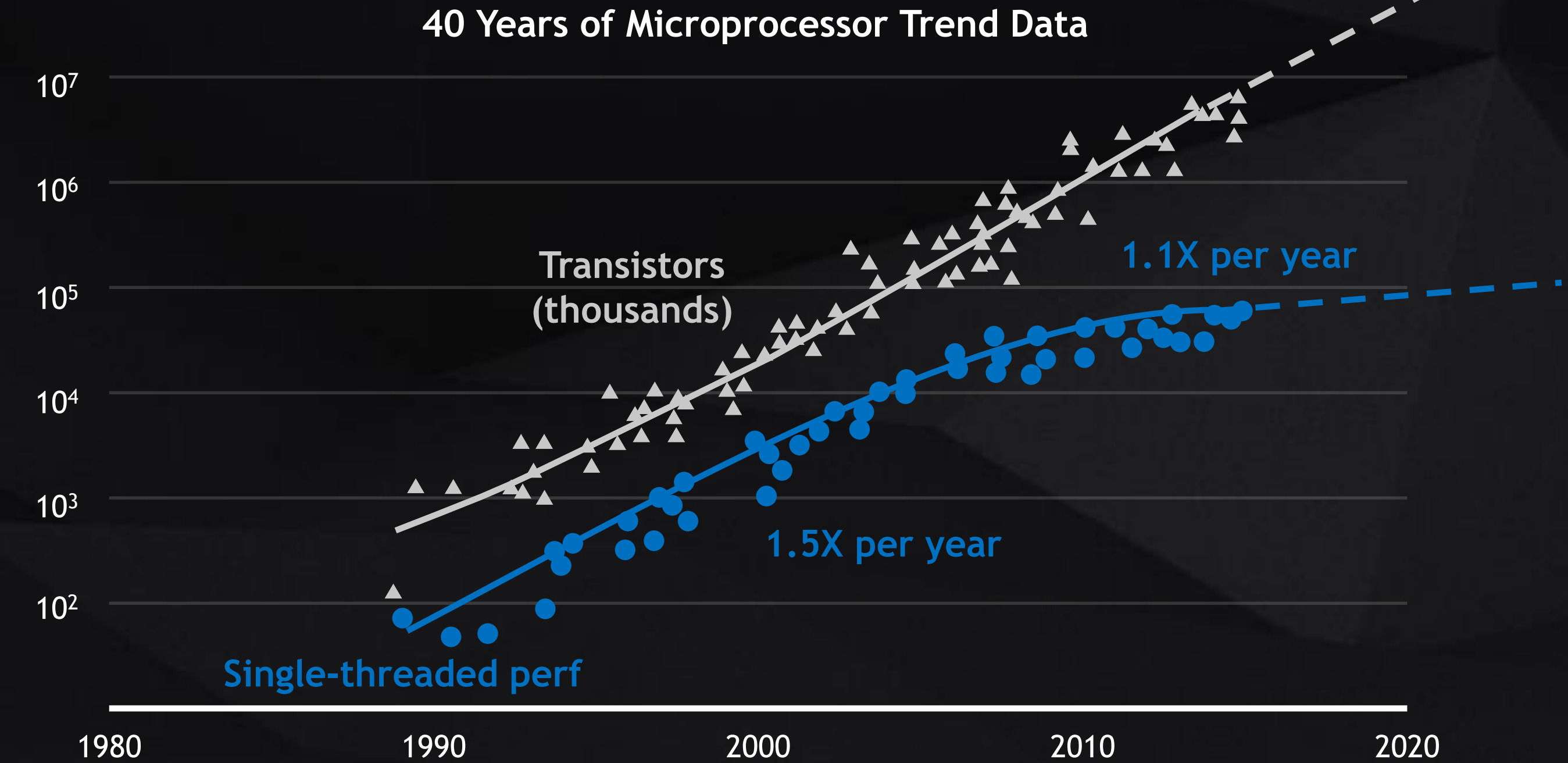
POWERING THE AI REVOLUTION

JENSEN HUANG, FOUNDER & CEO | GTC 2017

LIFE AFTER MOORE'S LAW

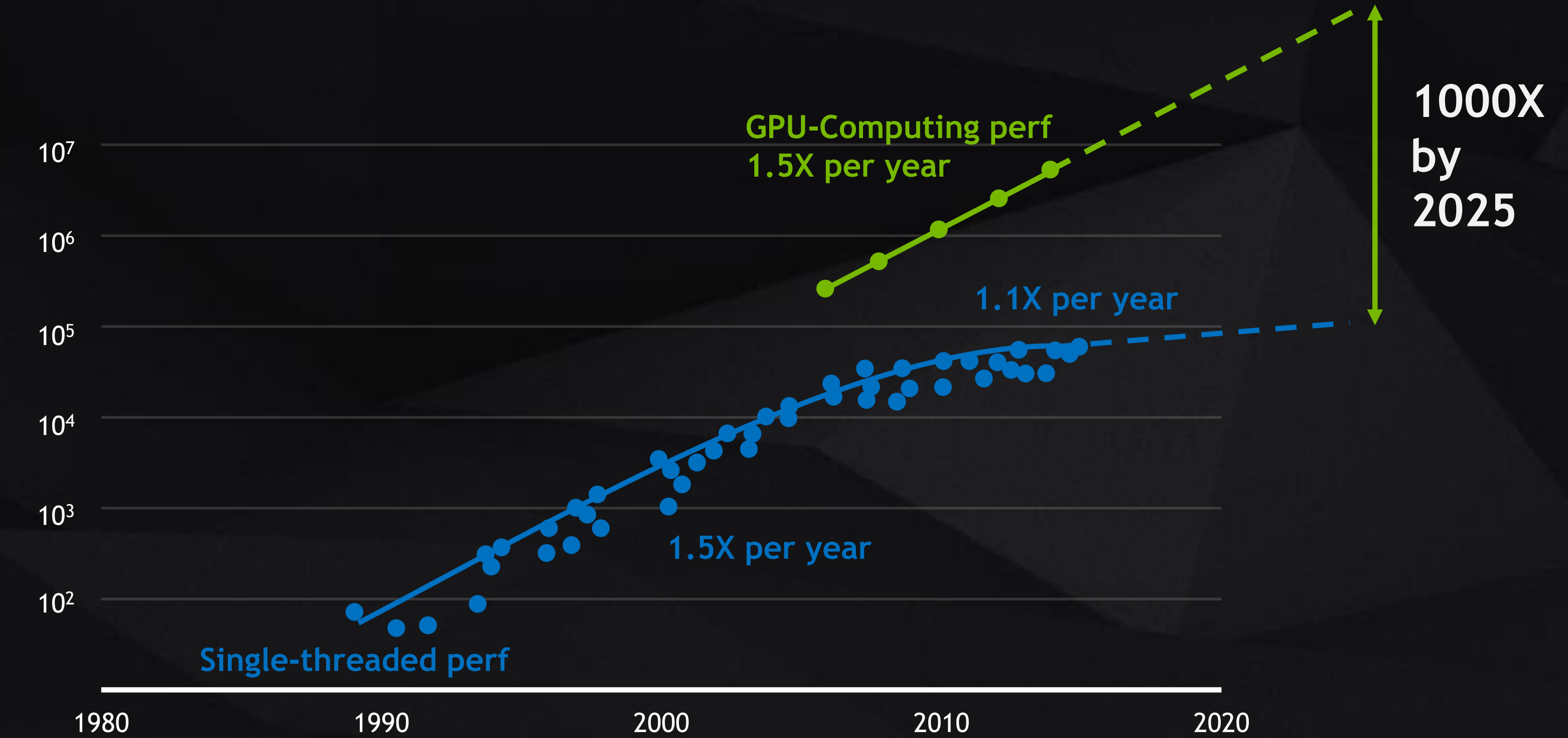
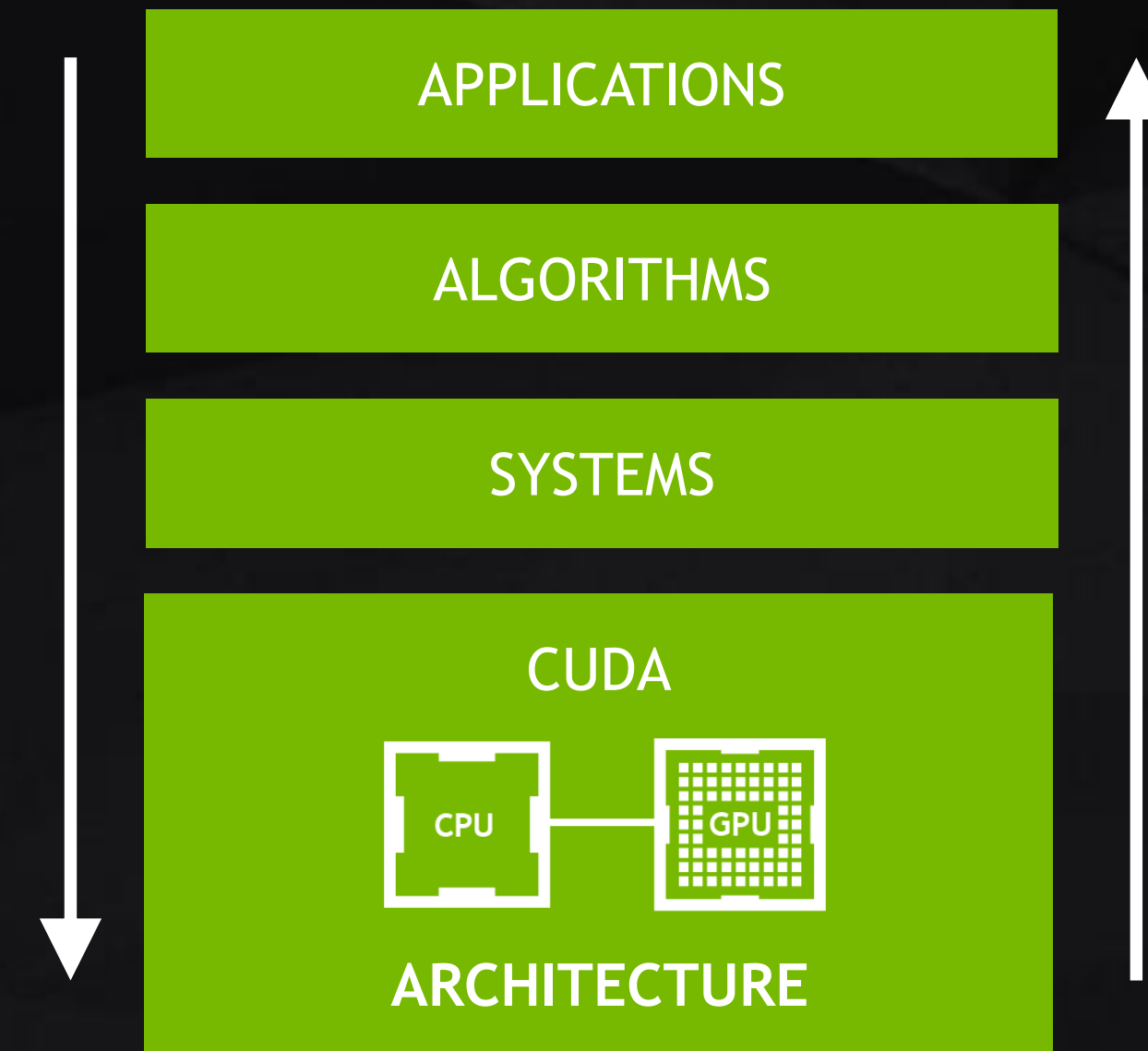
The End of Road for General Purpose Processors and the Future of Computing

John Hennessy
Stanford University
March 2017



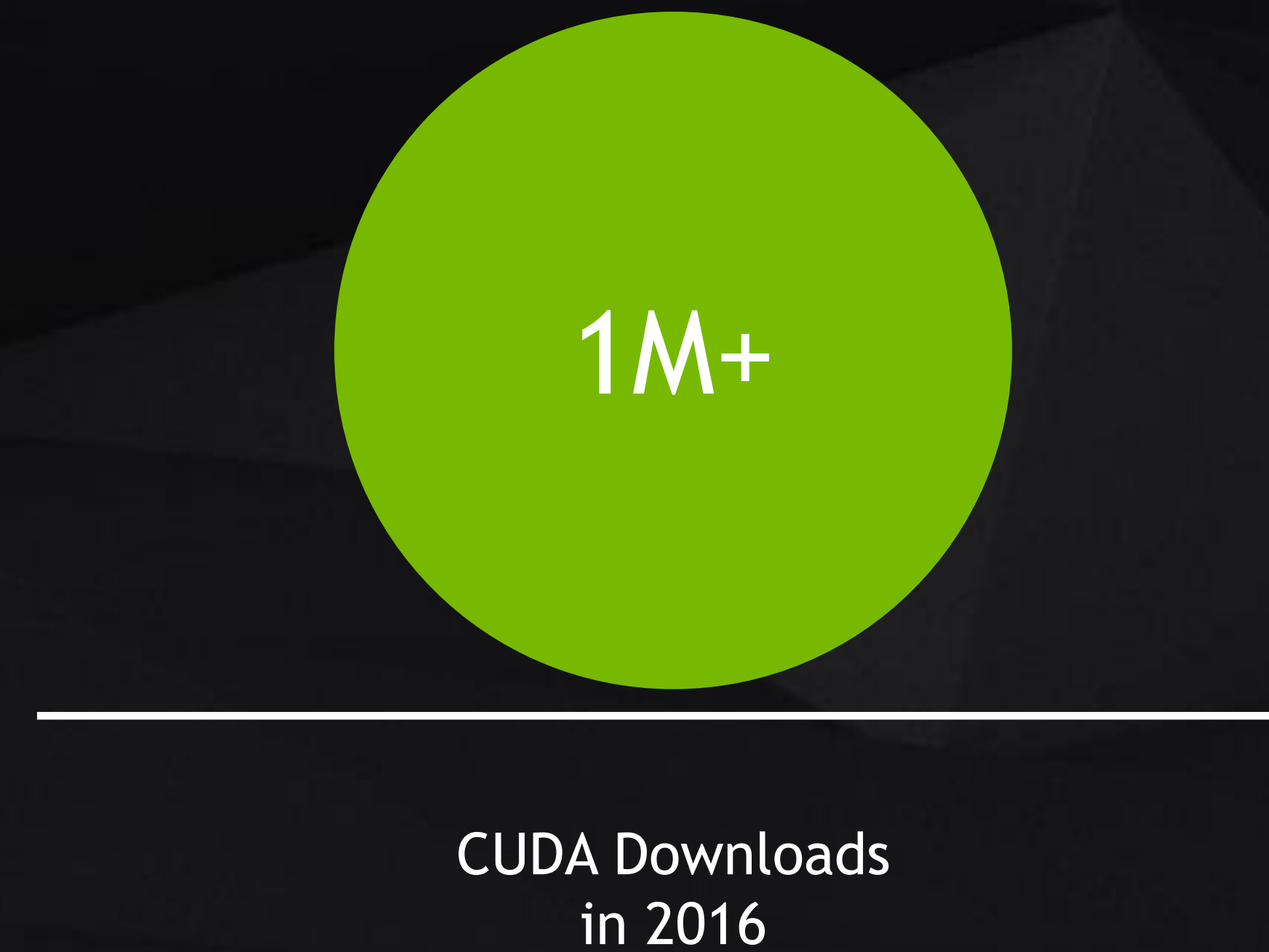
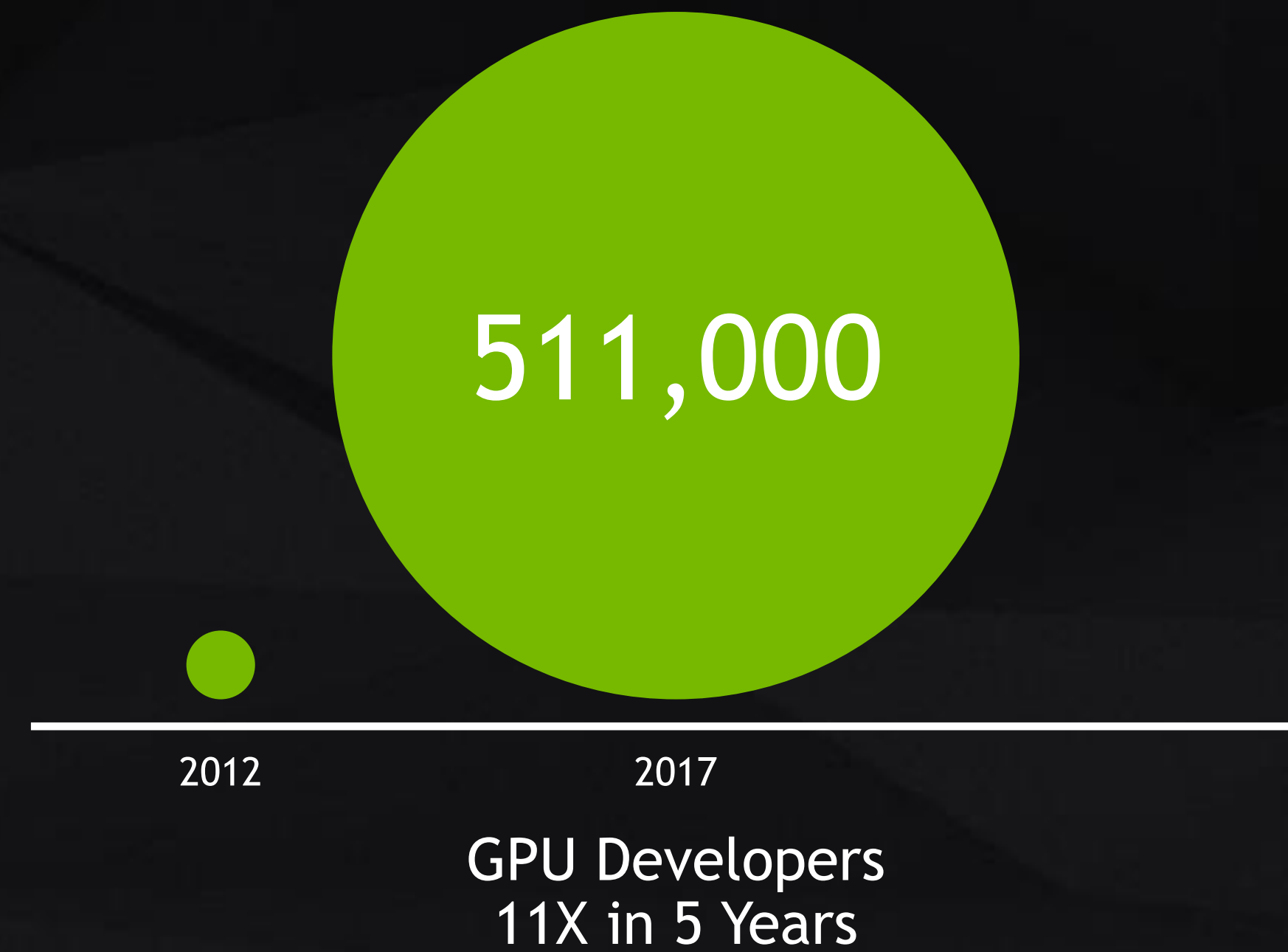
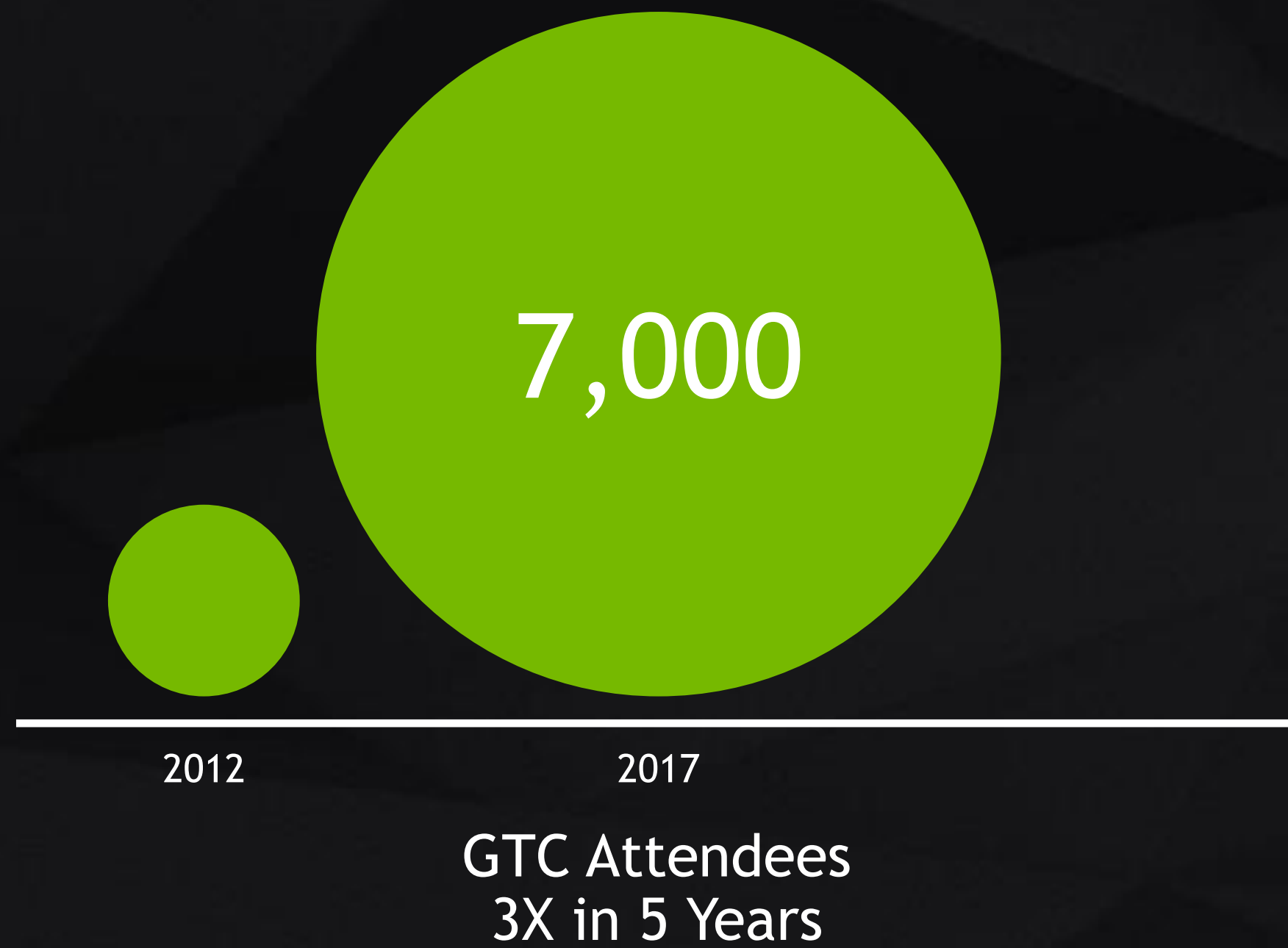
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

RISE OF GPU COMPUTING



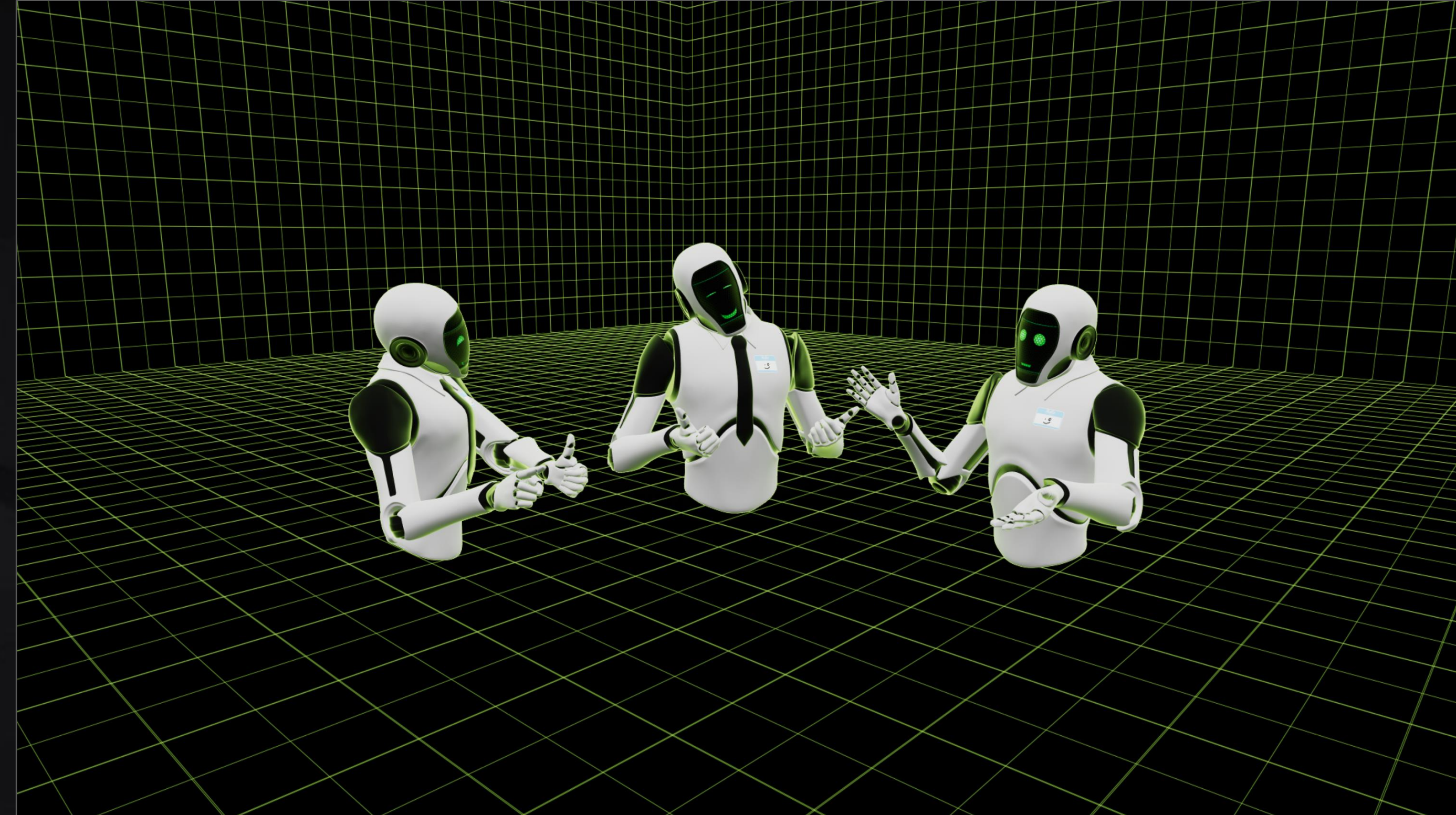
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

RISE OF GPU COMPUTING



ANNOUNCING PROJECT HOLODECK

Photorealistic models
Interactive physics
Collaboration



ANNOUNCING PROJECT HOLODECK

Photorealistic models

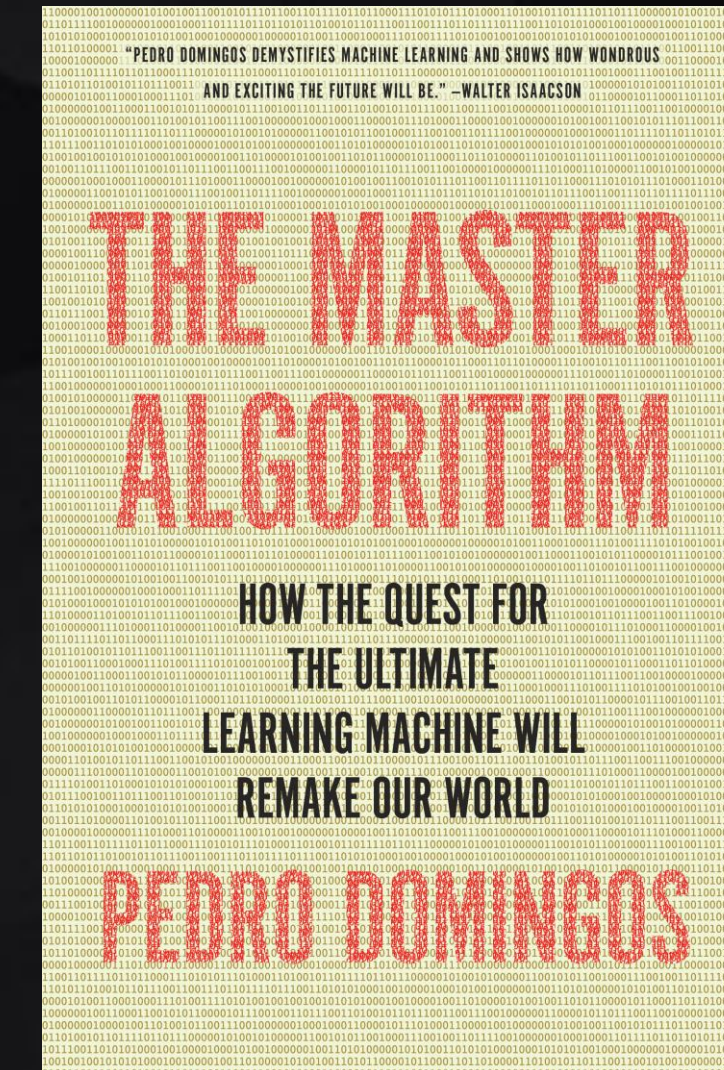
Interactive physics

Collaboration

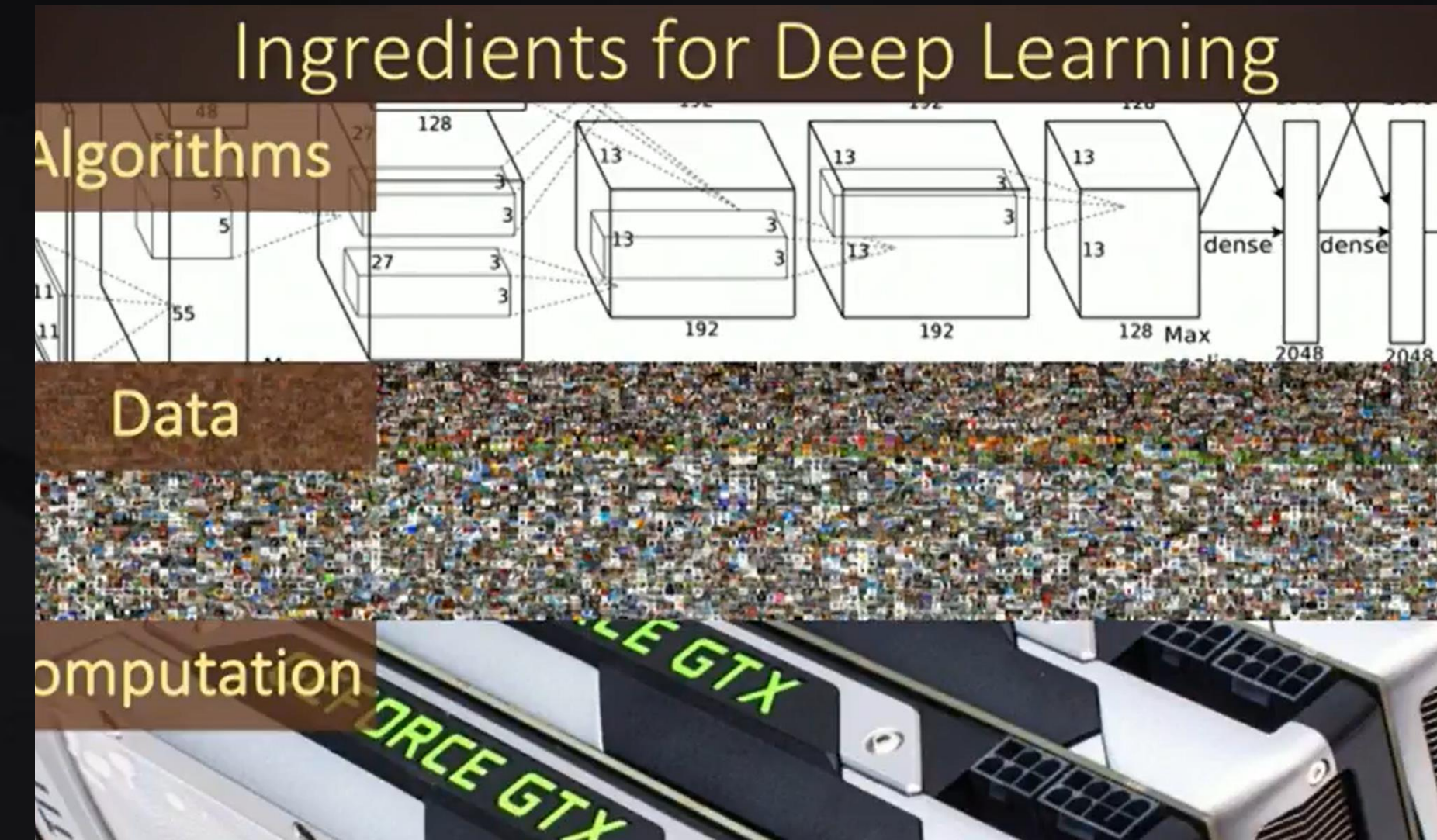
Early access in September



ERA OF MACHINE LEARNING

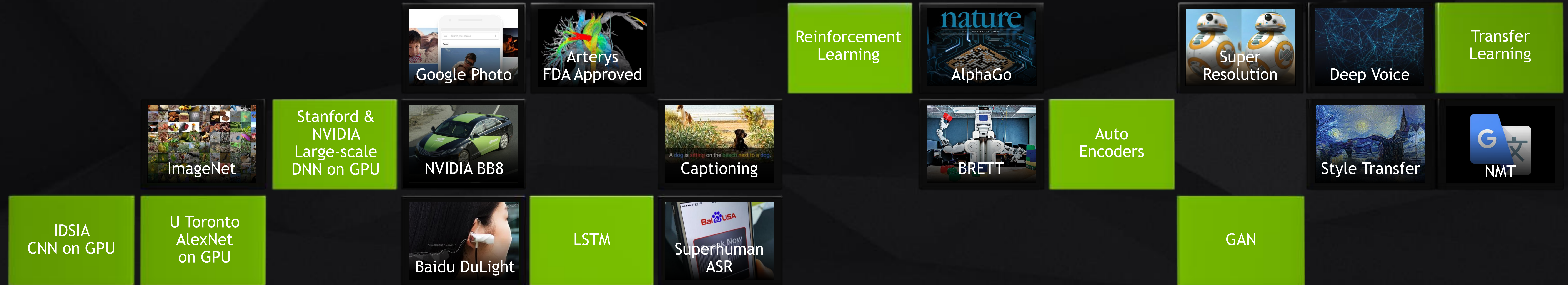


“The Master Algorithm”
— Pedro Domingos

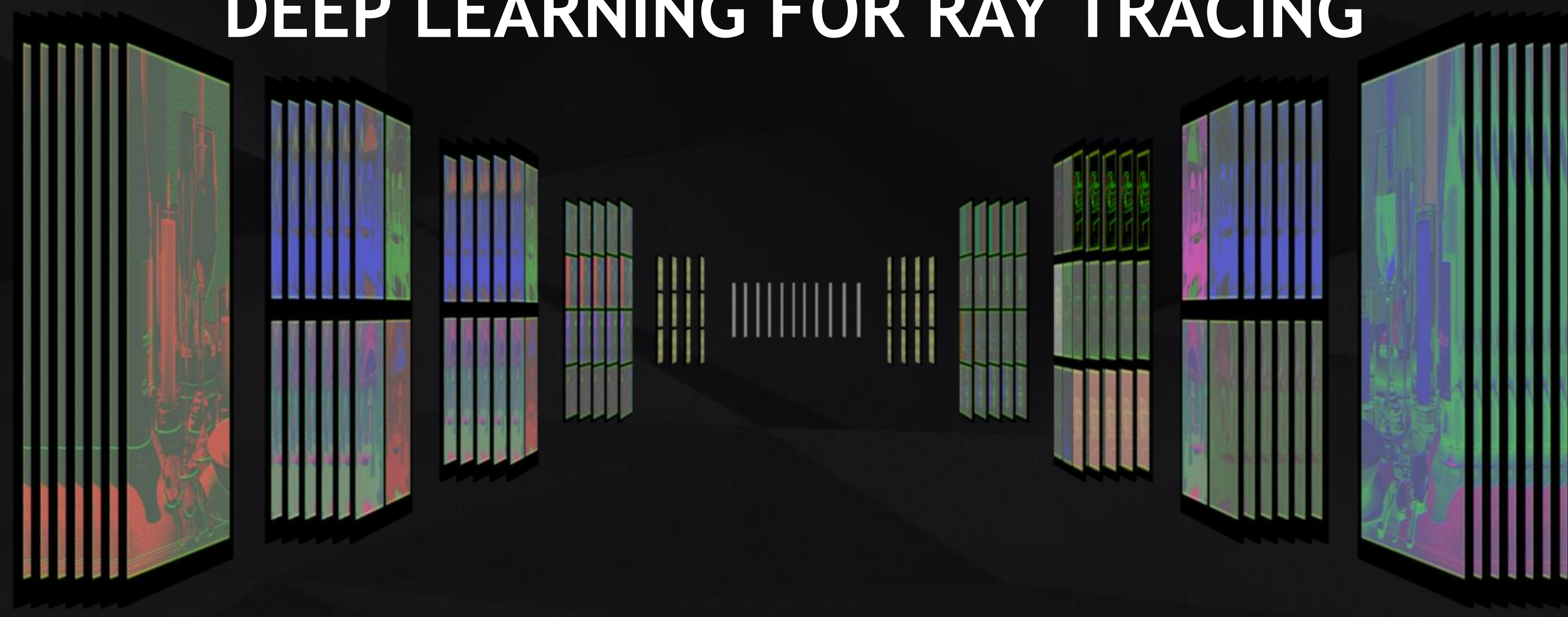


“A Quest for Intelligence”
— Fei-Fei Li

BIG BANG OF MODERN AI



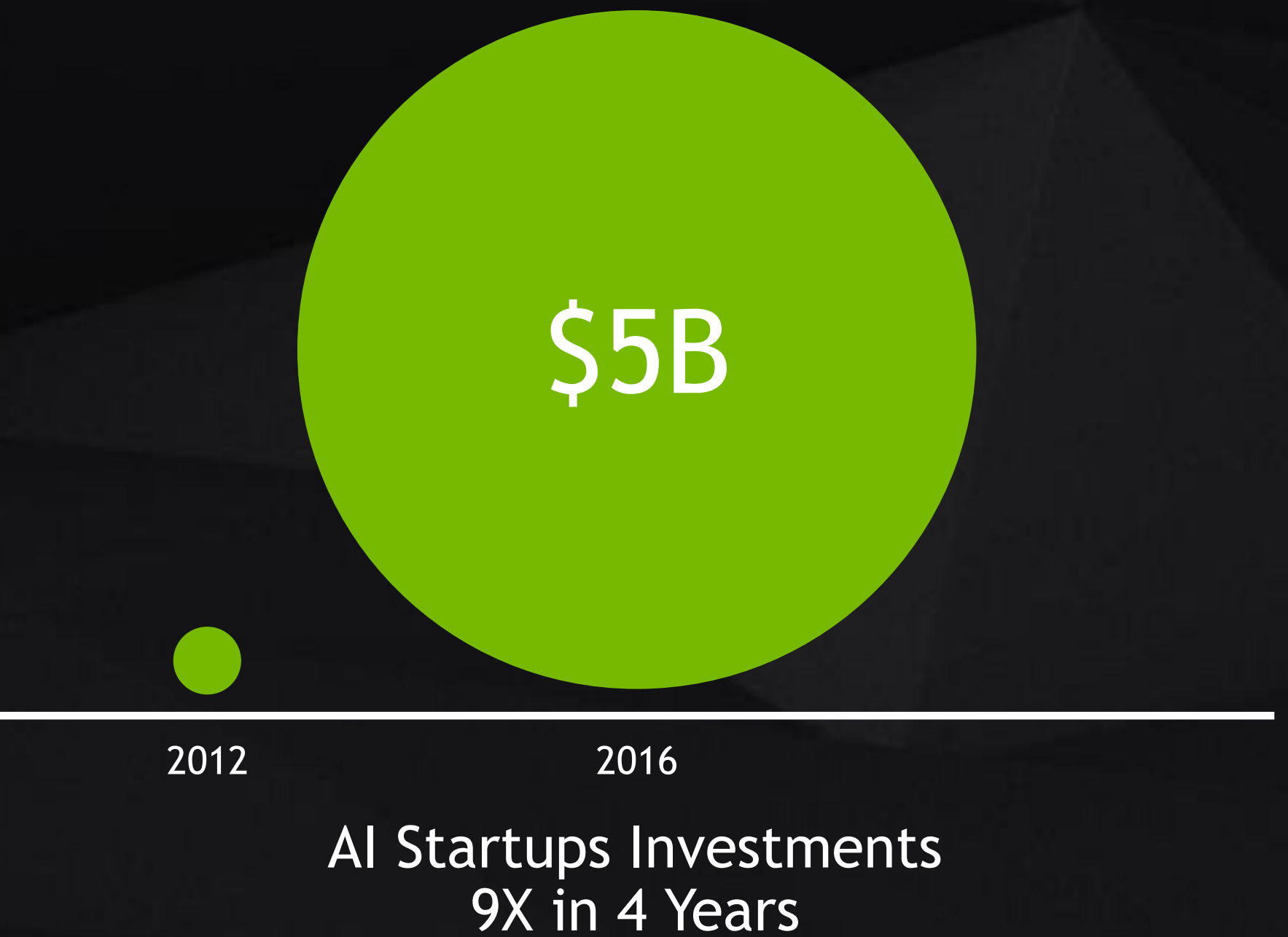
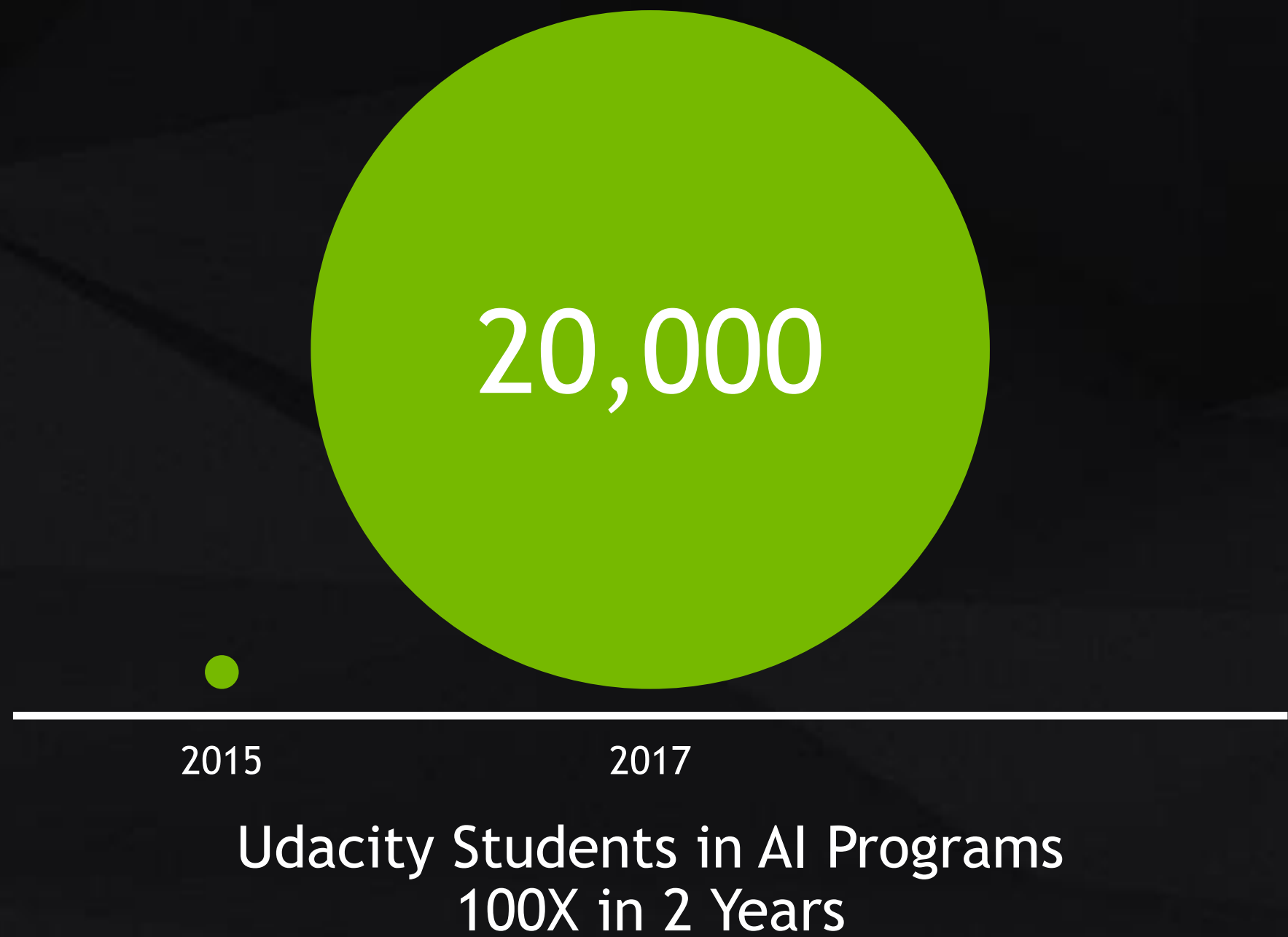
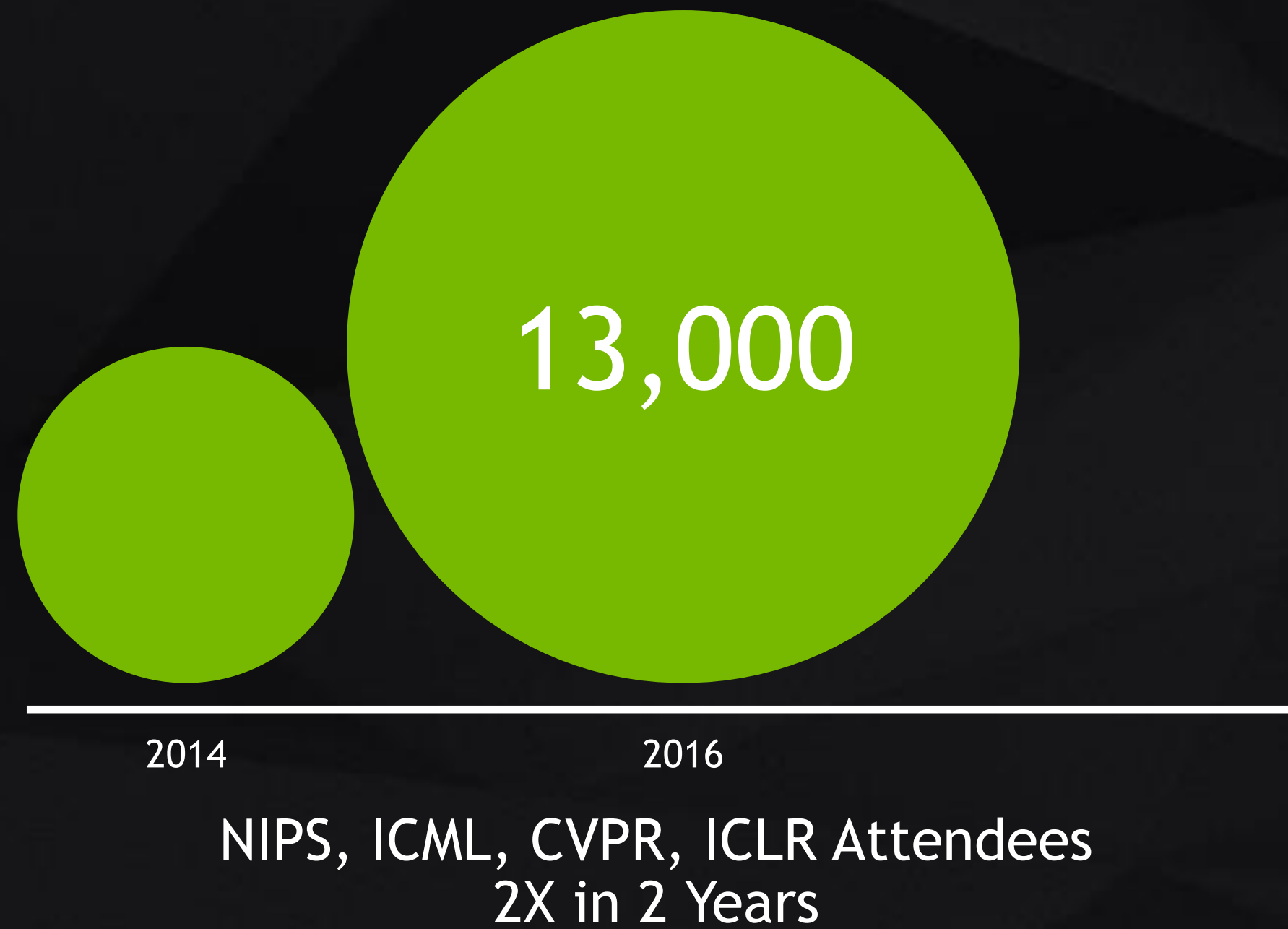
DEEP LEARNING FOR RAY TRACING



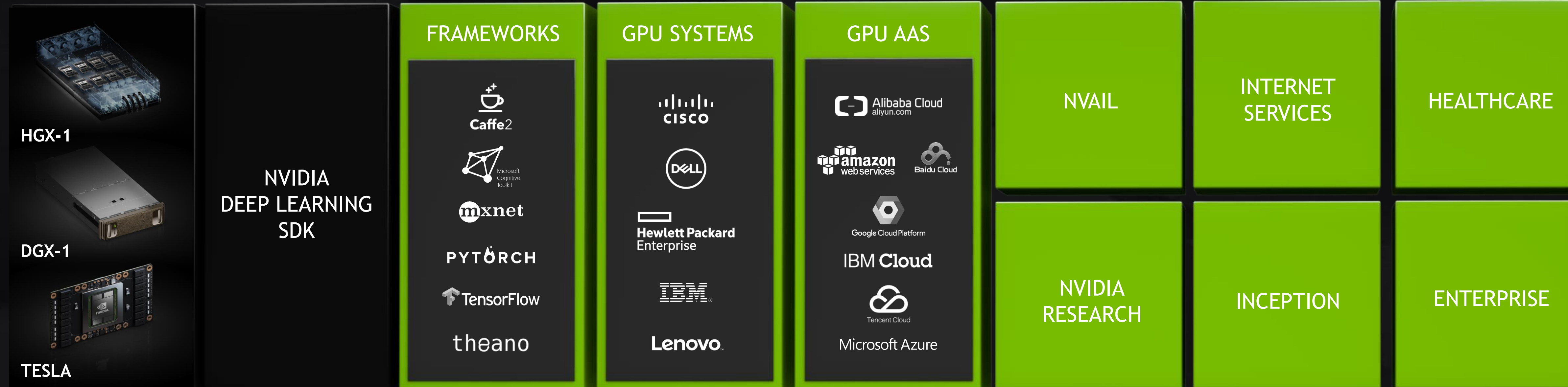
IRAY WITH DEEP LEARNING



BIG BANG OF MODERN AI



POWERING THE AI REVOLUTION



NVIDIA INCEPTION — 1,300 DEEP LEARNING STARTUPS

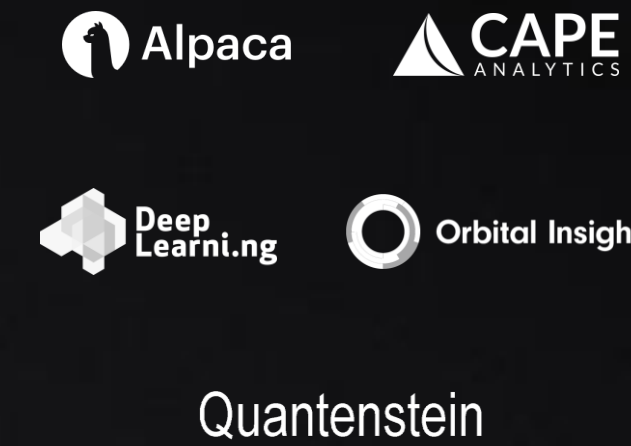
HEALTHCARE



RETAIL, ETAIL



FINANCIAL



SECURITY, IVA



PLATFORMS & APIs



DATA MANAGEMENT



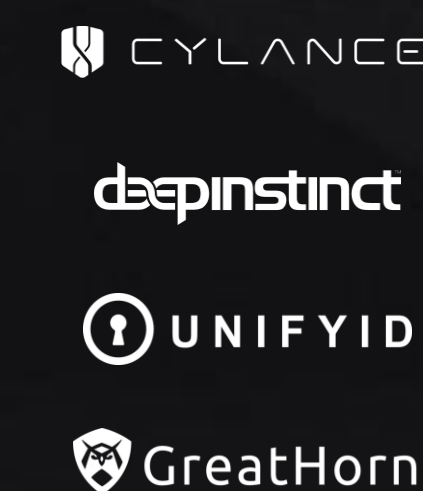
IOT & MANUFACTURING



AUTONOMOUS MACHINES



CYBER



AEC



DEVELOPMENT PLATFORMS

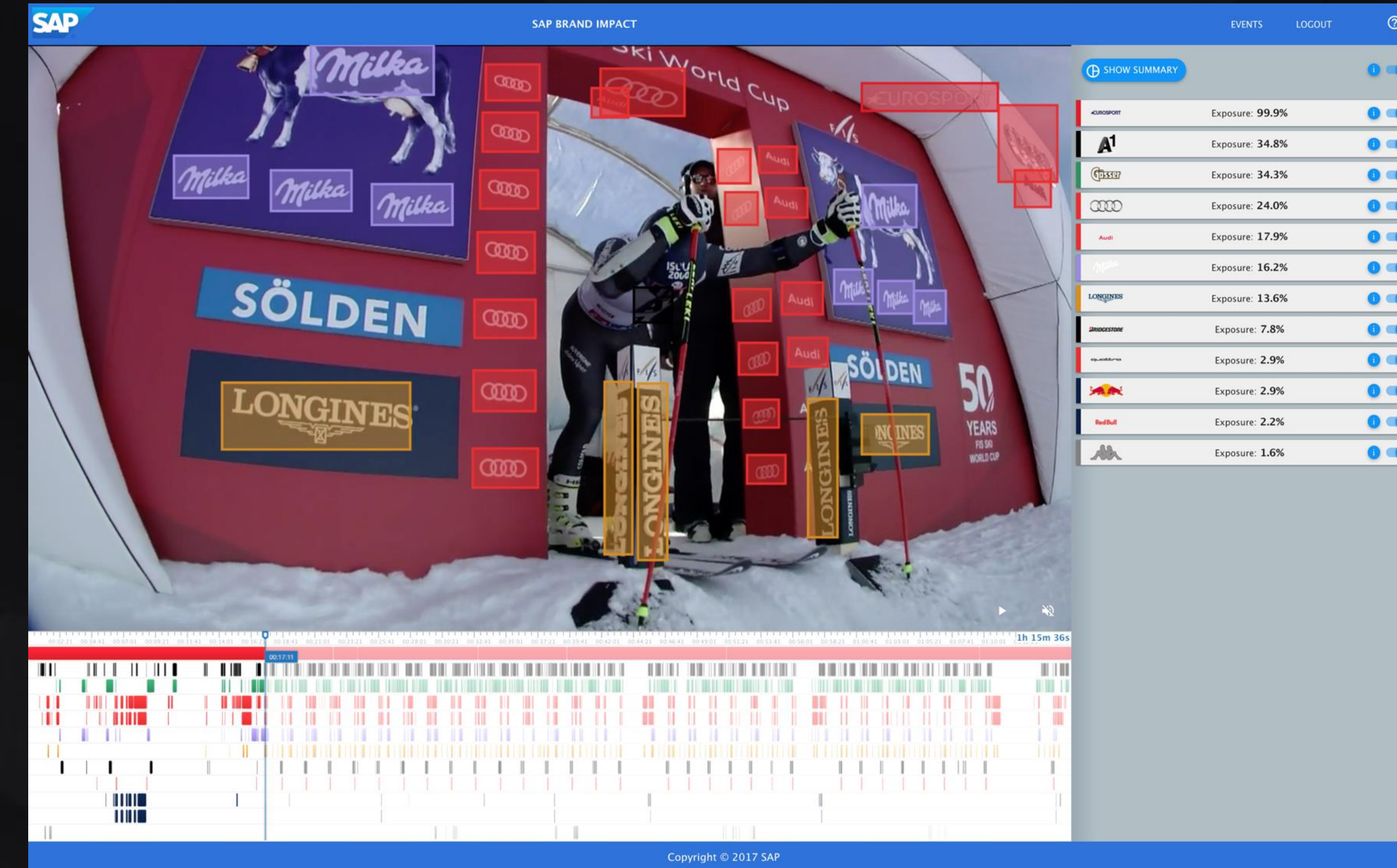


BUSINESS INTELLIGENCE & VISUALIZATION

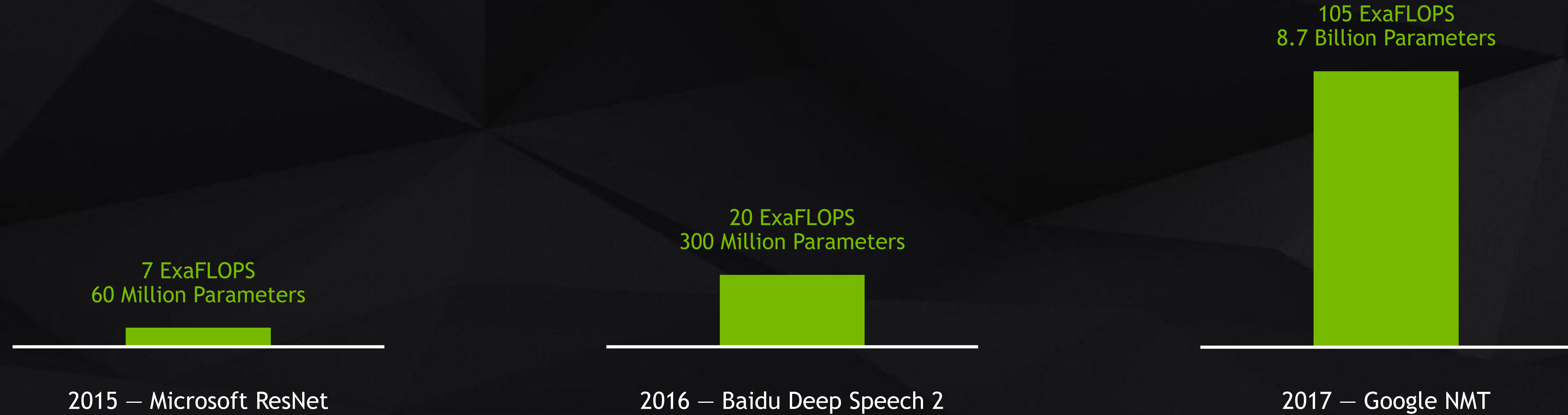


SAP AI FOR THE ENTERPRISE

First commercial AI offerings from SAP
Brand Impact, Service Ticketing, Invoice-to-Record applications
Powered by NVIDIA GPUs on DGX-1 and AWS



MODEL COMPLEXITY IS EXPLODING



ANNOUNCING TESLA V100

GIANT LEAP FOR AI & HPC
VOLTA WITH NEW TENSOR CORE

21B xtors | TSMC 12nm FFN | 815mm²

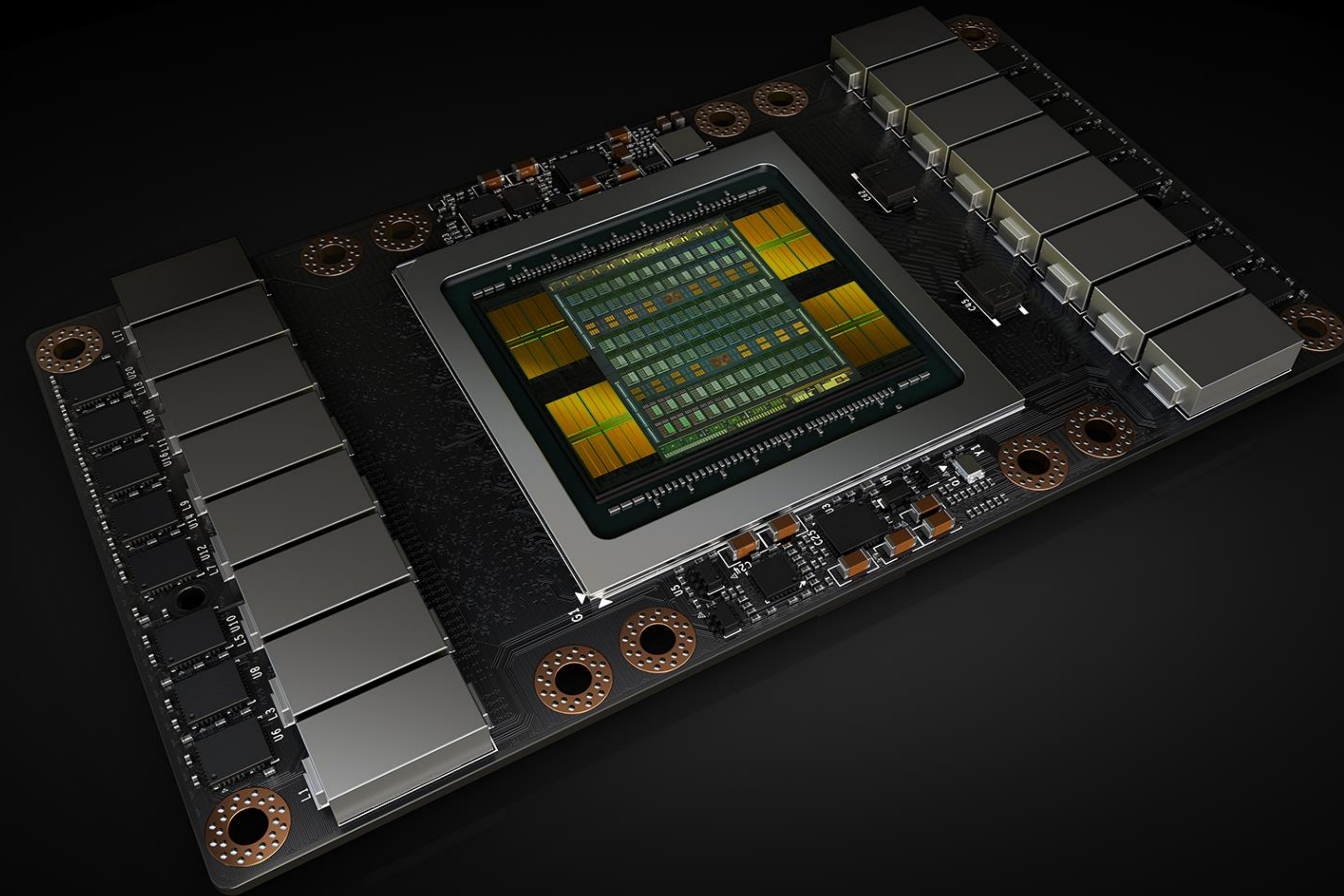
5,120 CUDA cores

7.5 FP64 TFLOPS | 15 FP32 TFLOPS

NEW 120 Tensor TFLOPS

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

300 GB/s NVLink



NEW TENSOR CORE

New CUDA TensorOp instructions & data formats

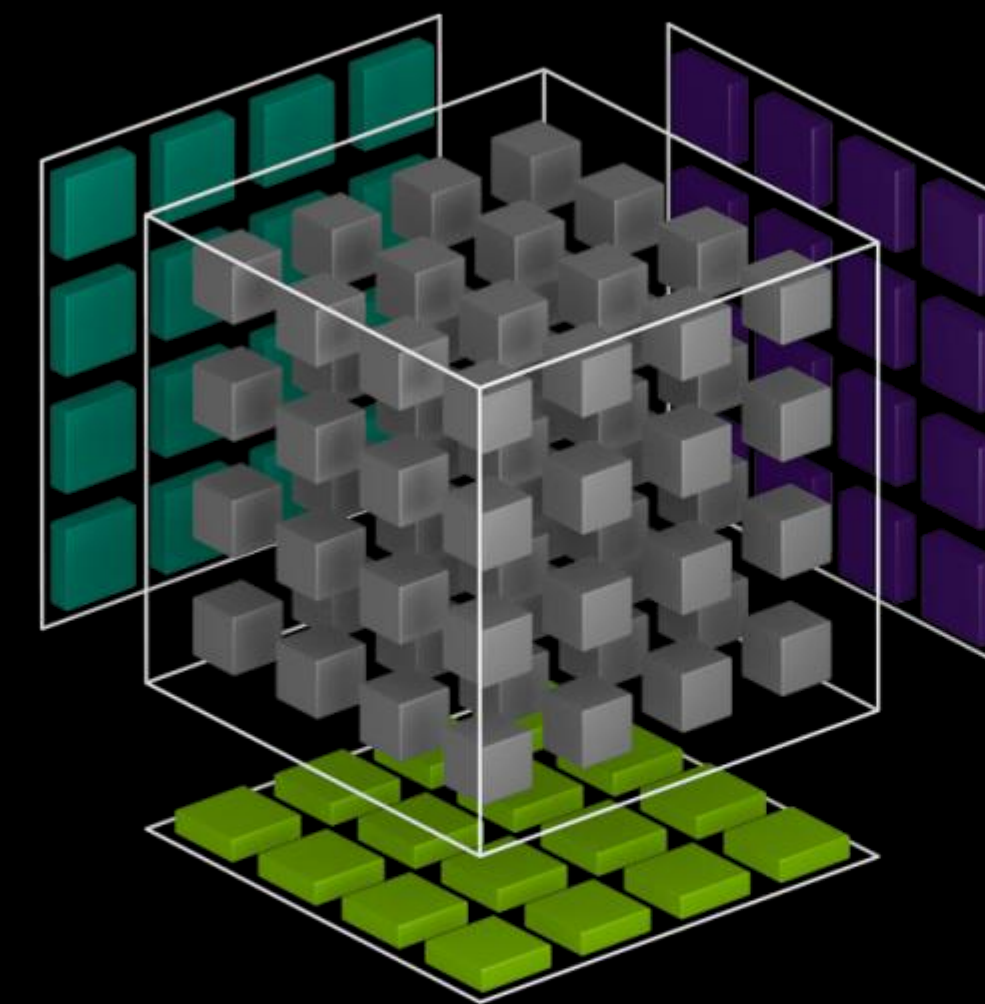
4x4 matrix processing array

$D[\text{FP32}] = A[\text{FP16}] * B[\text{FP16}] + C[\text{FP32}]$

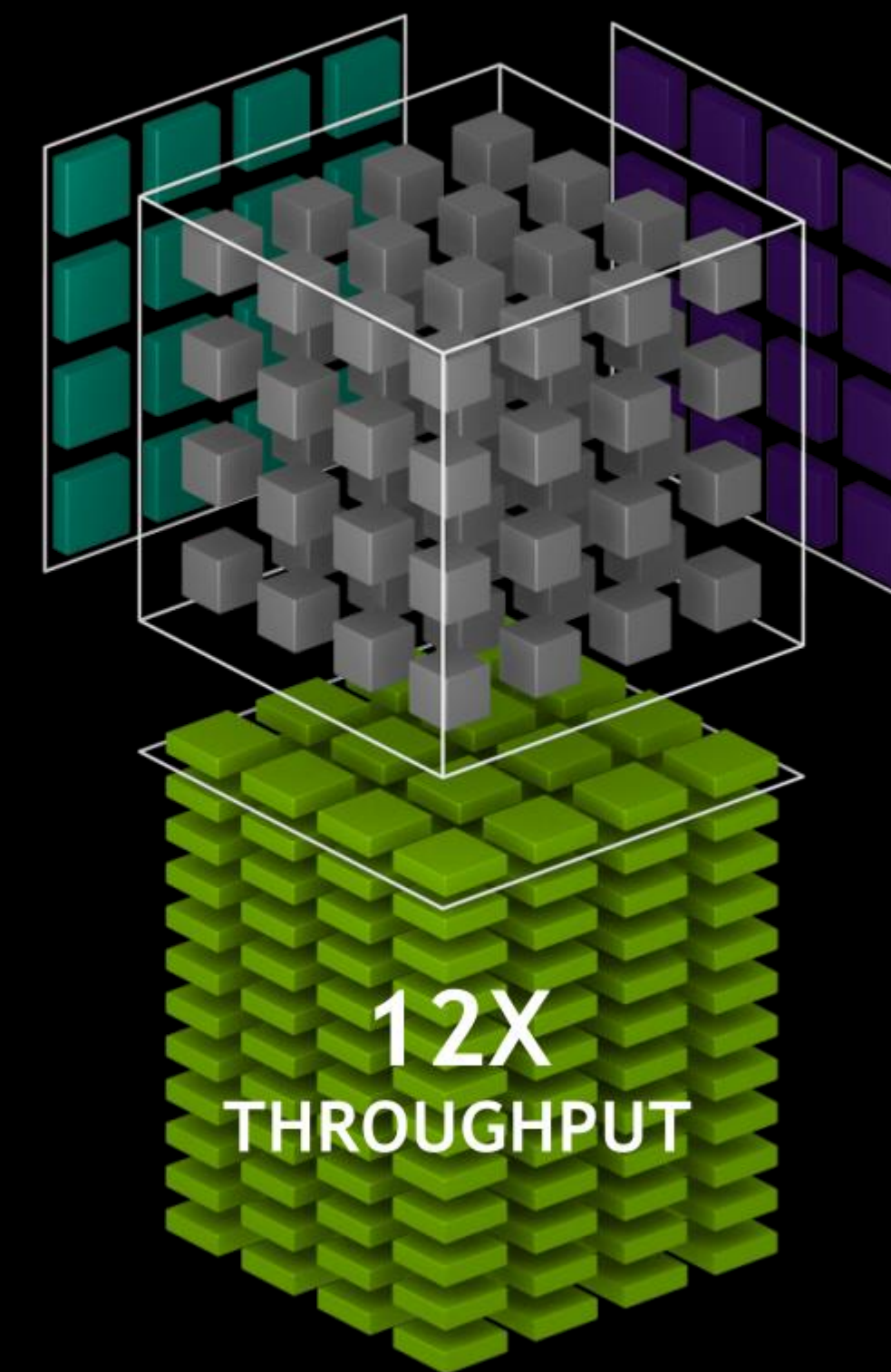
Optimized for deep learning

■ Activation Inputs ■ Weights Inputs ■ Output Results

PASCAL



VOLTA TENSOR CORES



ANNOUNCING TESLA V100

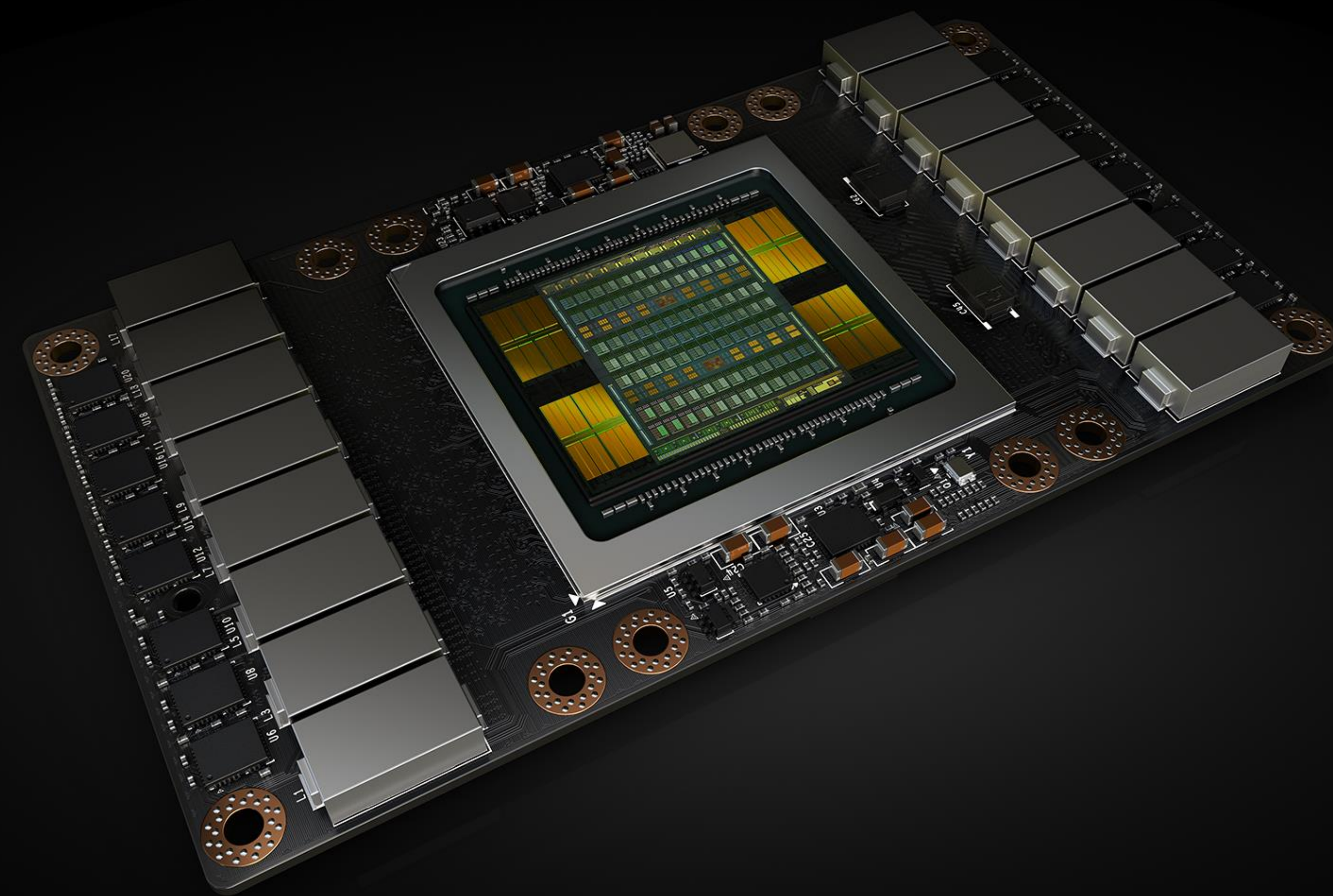
GIANT LEAP FOR AI & HPC
VOLTA WITH NEW TENSOR CORE

Compared to Pascal

1.5X General-purpose FLOPS for HPC

12X Tensor FLOPS for DL training

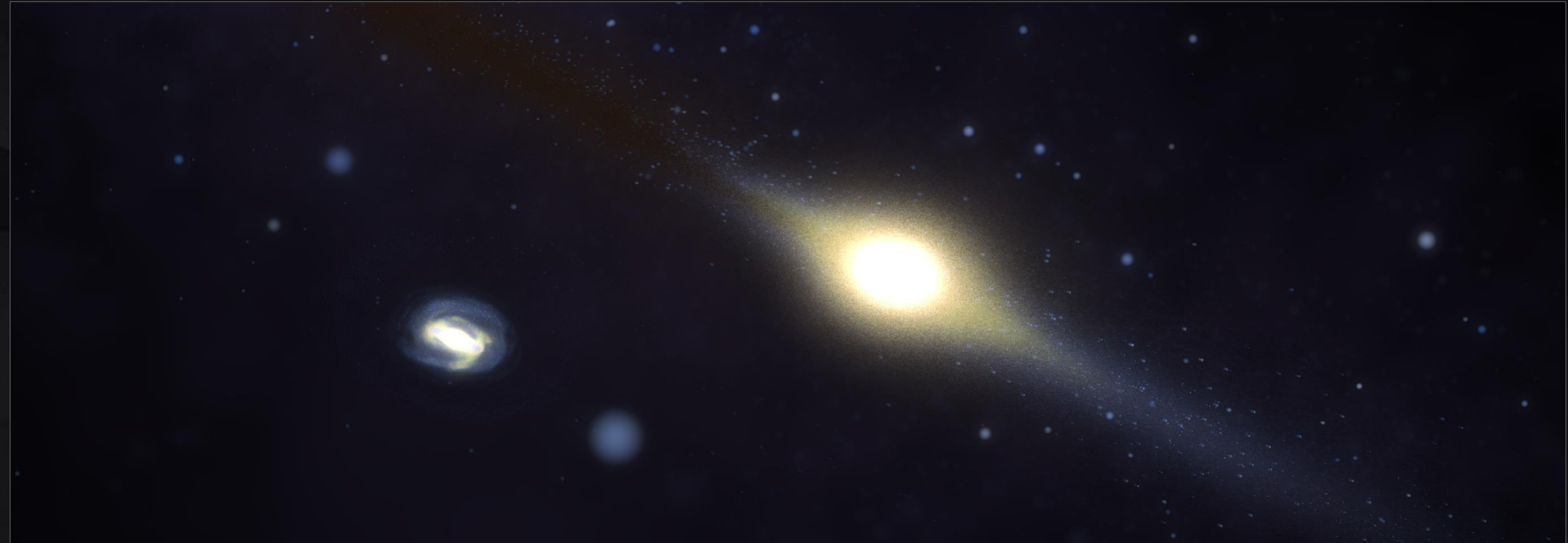
6X Tensor FLOPS for DL inferencing



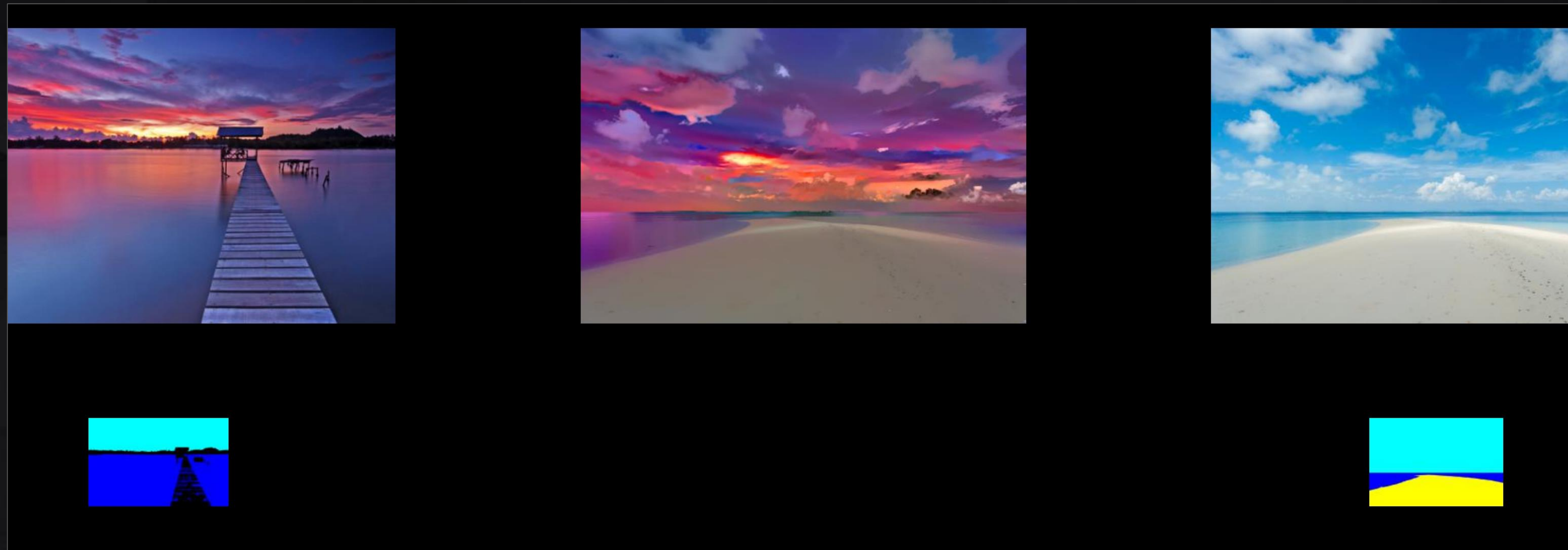


SQUARE ENIX®

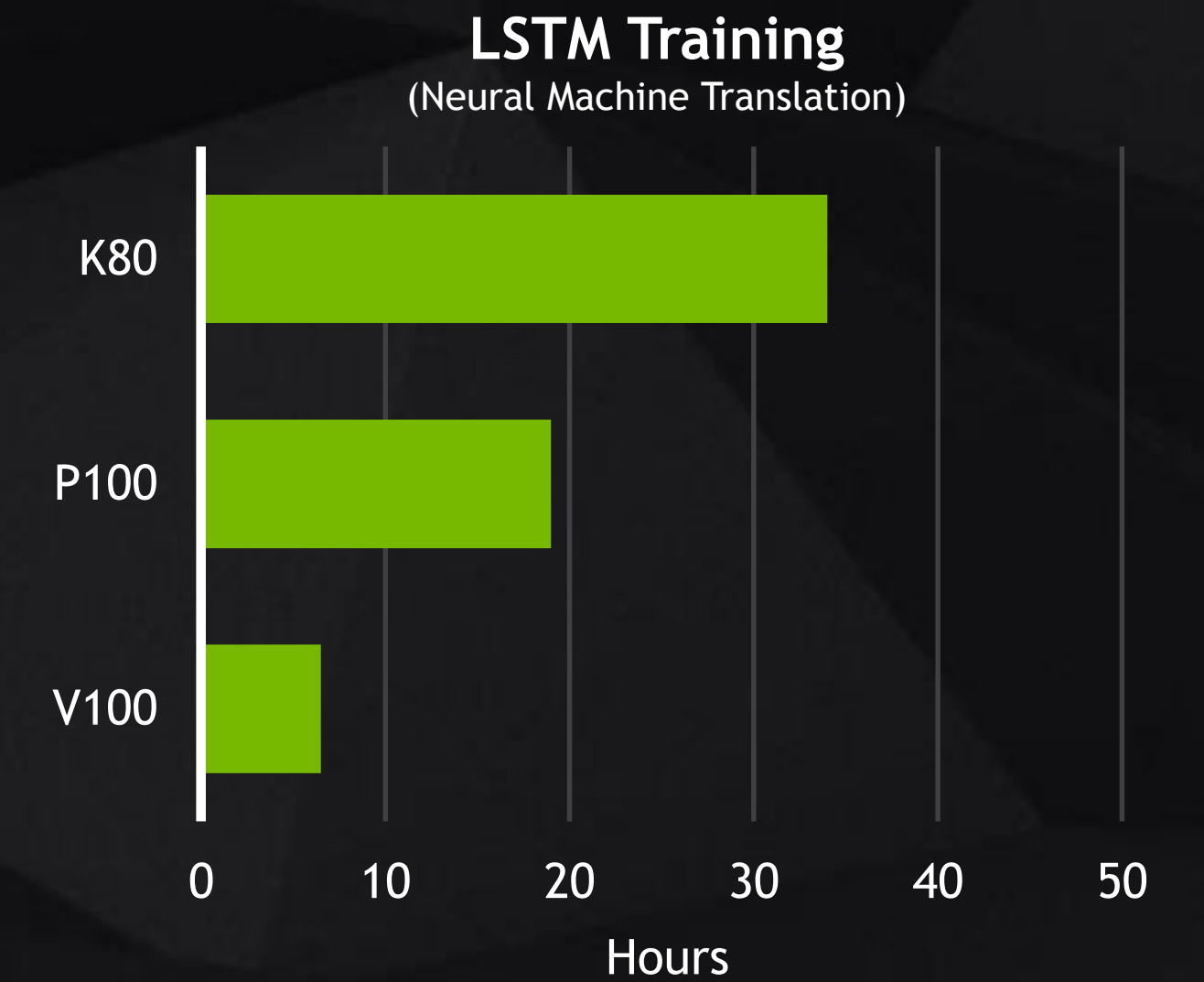
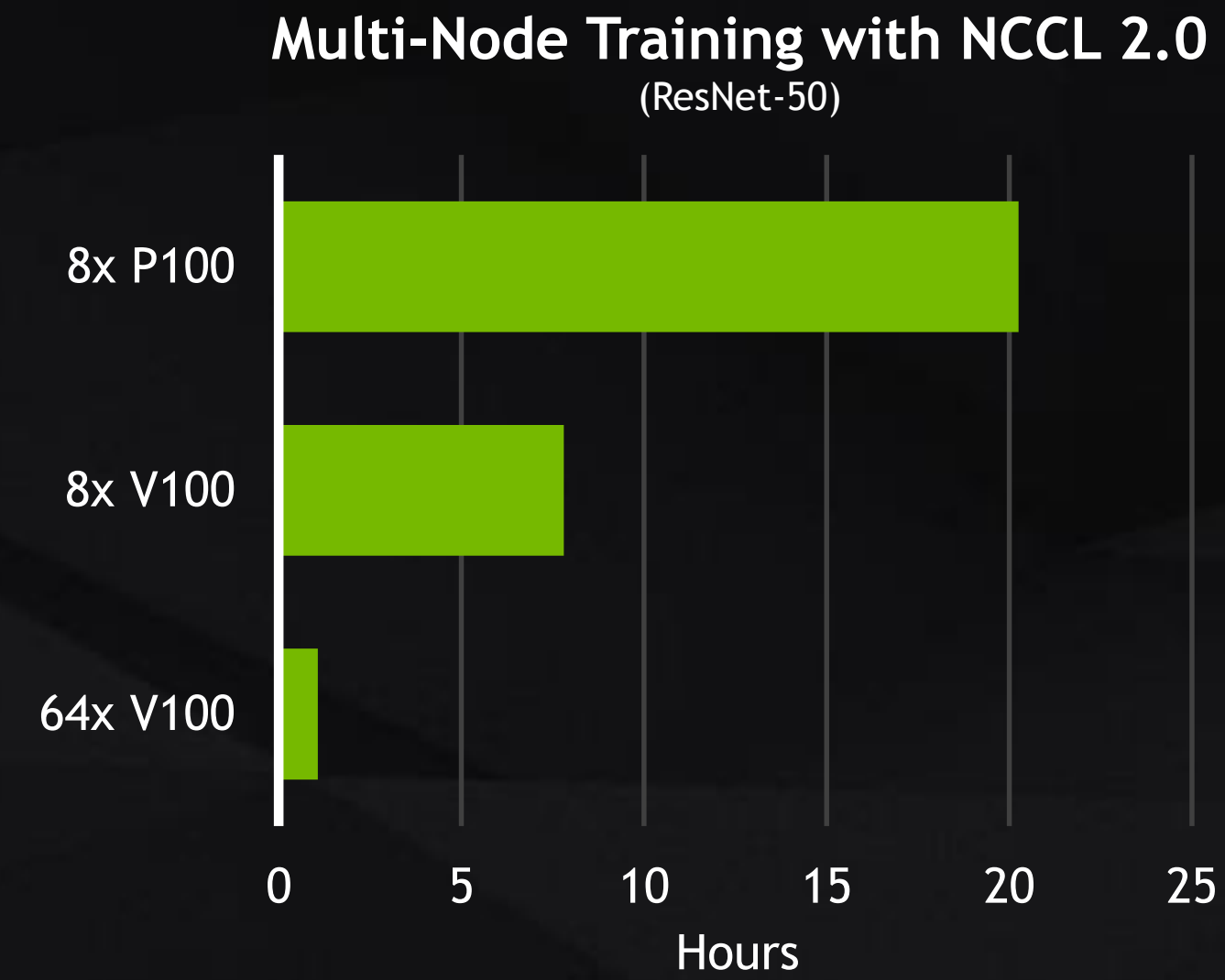
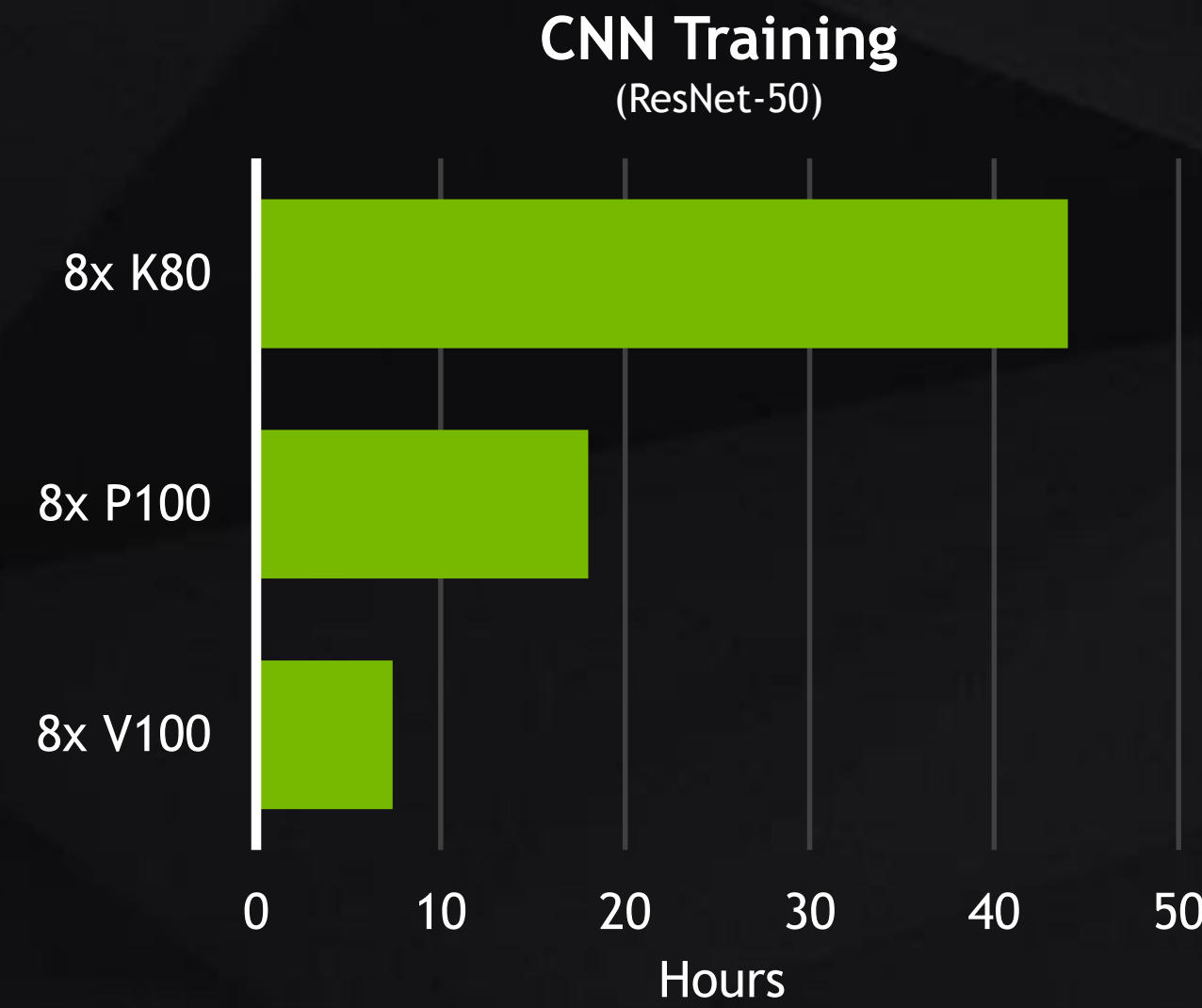
GALAXY



DEEP LEARNING FOR STYLE TRANSFER



ANNOUNCING NEW FRAMEWORK RELEASES FOR VOLTA





MATT WOOD

General Manager
Deep Learning and AI, AWS



ANNOUNCING NVIDIA DGX-1 WITH TESLA V100

ESSENTIAL INSTRUMENT OF AI RESEARCH

960 Tensor TFLOPS | 8x Tesla V100 | NVLink Hybrid Cube
From 8 days on TITAN X to 8 hours
400 servers in a box



ANNOUNCING NVIDIA DGX-1 WITH TESLA V100

ESSENTIAL INSTRUMENT OF AI RESEARCH

960 Tensor TFLOPS | 8x Tesla V100 | NVLink Hybrid Cube
From 8 days on TITAN X to 8 hours
400 servers in a box



ANNOUNCING NVIDIA DGX-1 WITH TESLA V100

ESSENTIAL INSTRUMENT OF AI RESEARCH

960 Tensor TFLOPS | 8x Tesla V100 | NVLink Hybrid Cube

From 8 days on TITAN X to 8 hours

400 servers in a box

\$149,000

Order today: nvidia.com/DGX-1



ANNOUNCING NVIDIA DGX STATION PERSONAL DGX

480 Tensor TFLOPS | 4x Tesla V100 16GB | NVLink Fully Connected
3x DisplayPort | 1500W | Water Cooled



ANNOUNCING NVIDIA DGX STATION

PERSONAL DGX

480 Tensor TFLOPS | 4x Tesla V100 16GB | NVLink Fully Connected
3x DisplayPort | 1500W | Water Cooled

\$69,000

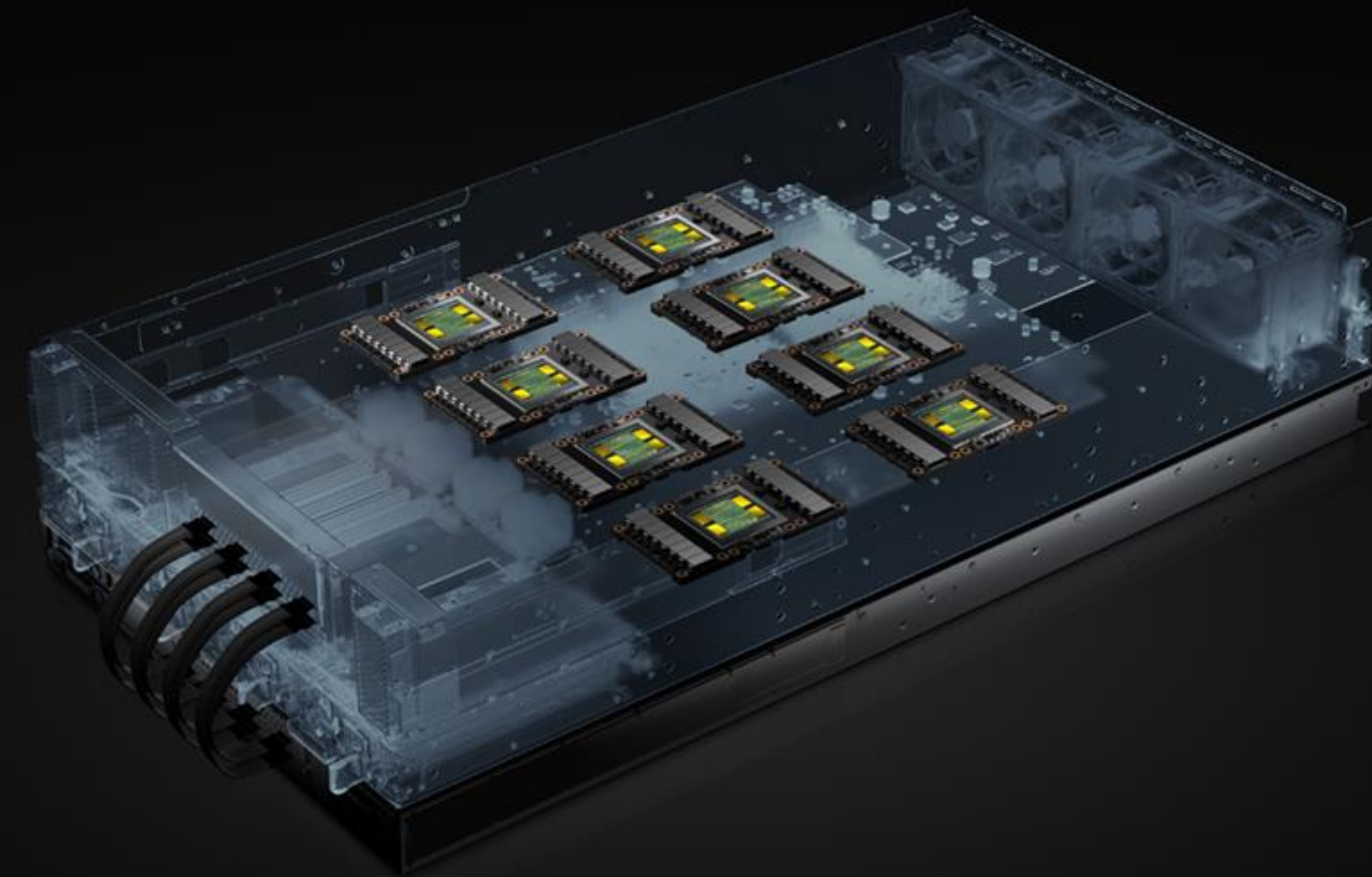
Order today: nvidia.com/DGX-Station



ANNOUNCING HGX-1 WITH TESLA V100

VERSATILE GPU CLOUD COMPUTING

8x Tesla V100 with NVLINK Hybrid Cube
2C:8G | 2C:4G | 1C:2G Configurable
NVIDIA Deep Learning, GRID graphics, CUDA HPC stacks



Microsoft Azure

JASON ZANDER

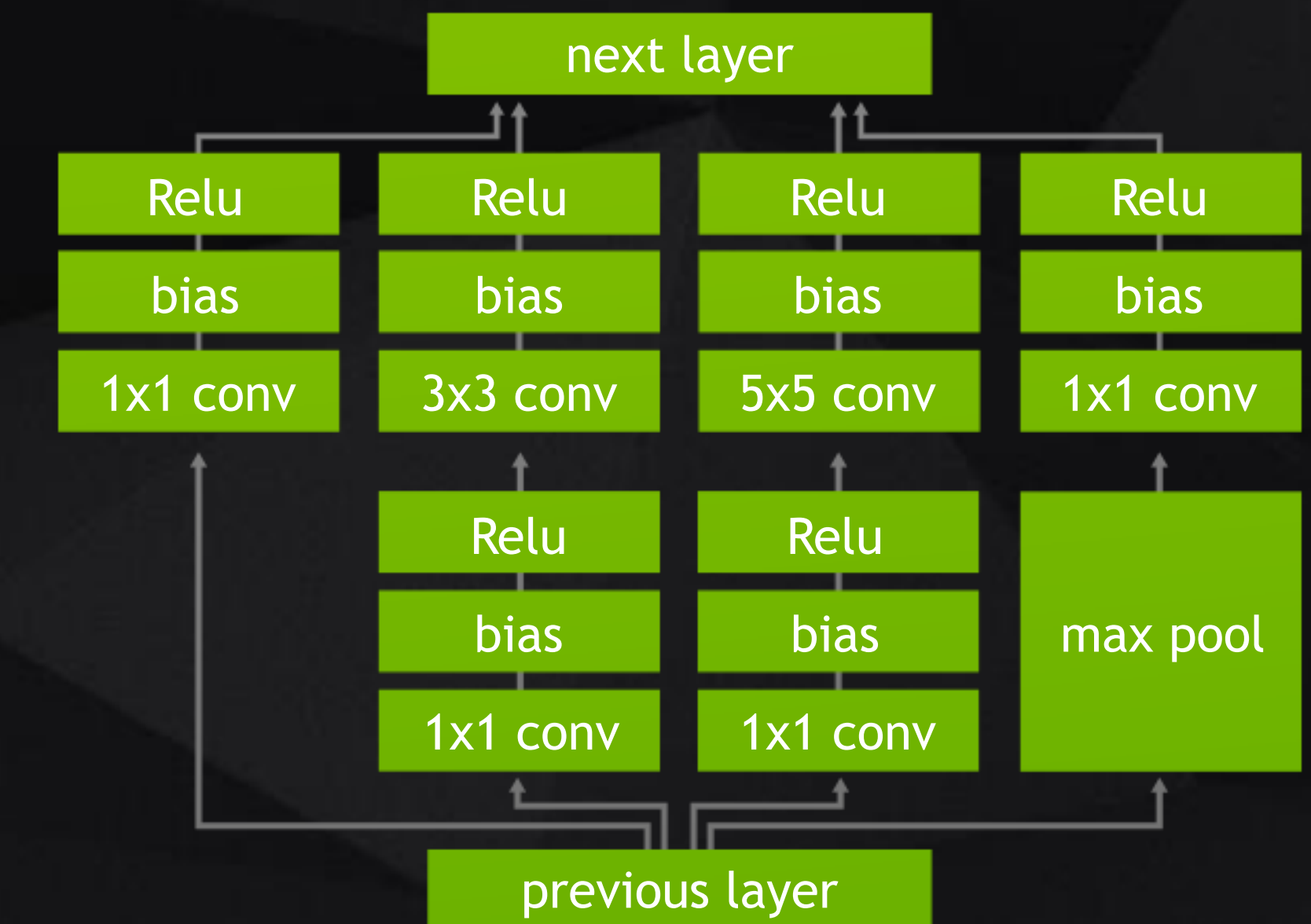
Corporate Vice President
Microsoft Azure, Microsoft



ANNOUNCING TENSORRT FOR TENSORFLOW

COMPILER FOR DEEP LEARNING INFERENCE

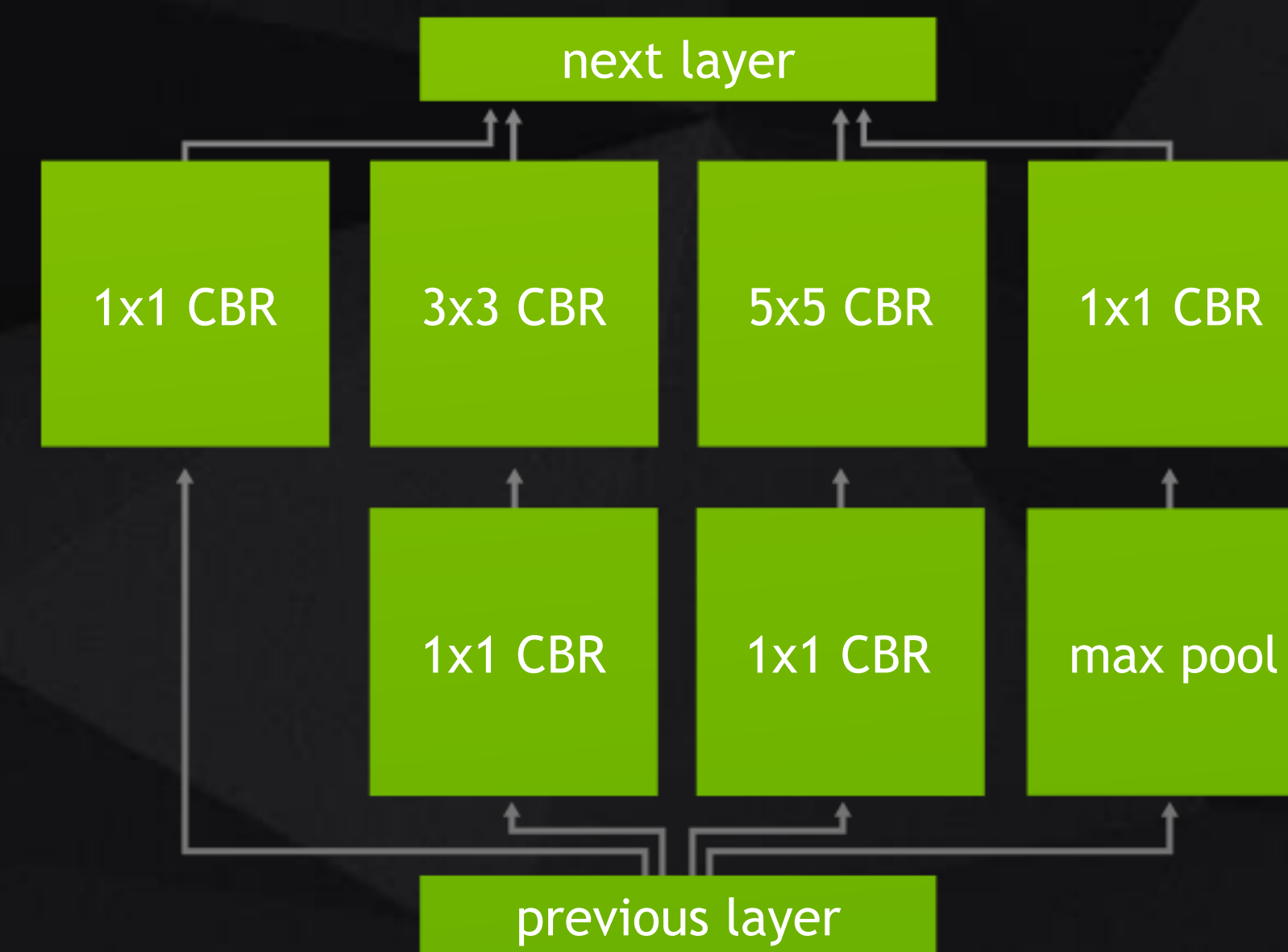
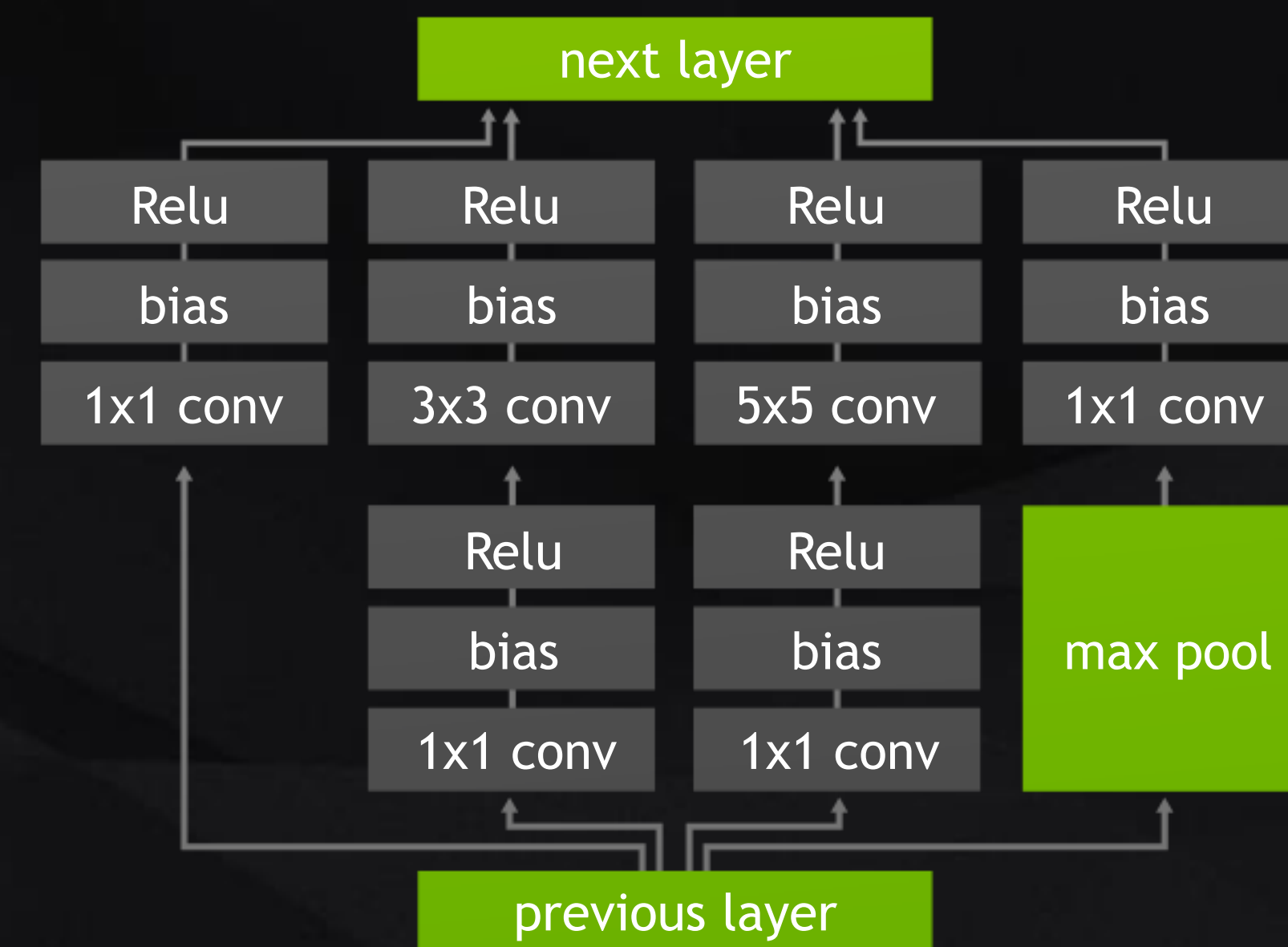
Graph optimizations for vertical and horizontal layer fusion
GPU-specific optimizations
Import models from Caffe and TensorFlow



ANNOUNCING TENSORRT FOR TENSORFLOW

COMPILER FOR DEEP LEARNING INFERENCE

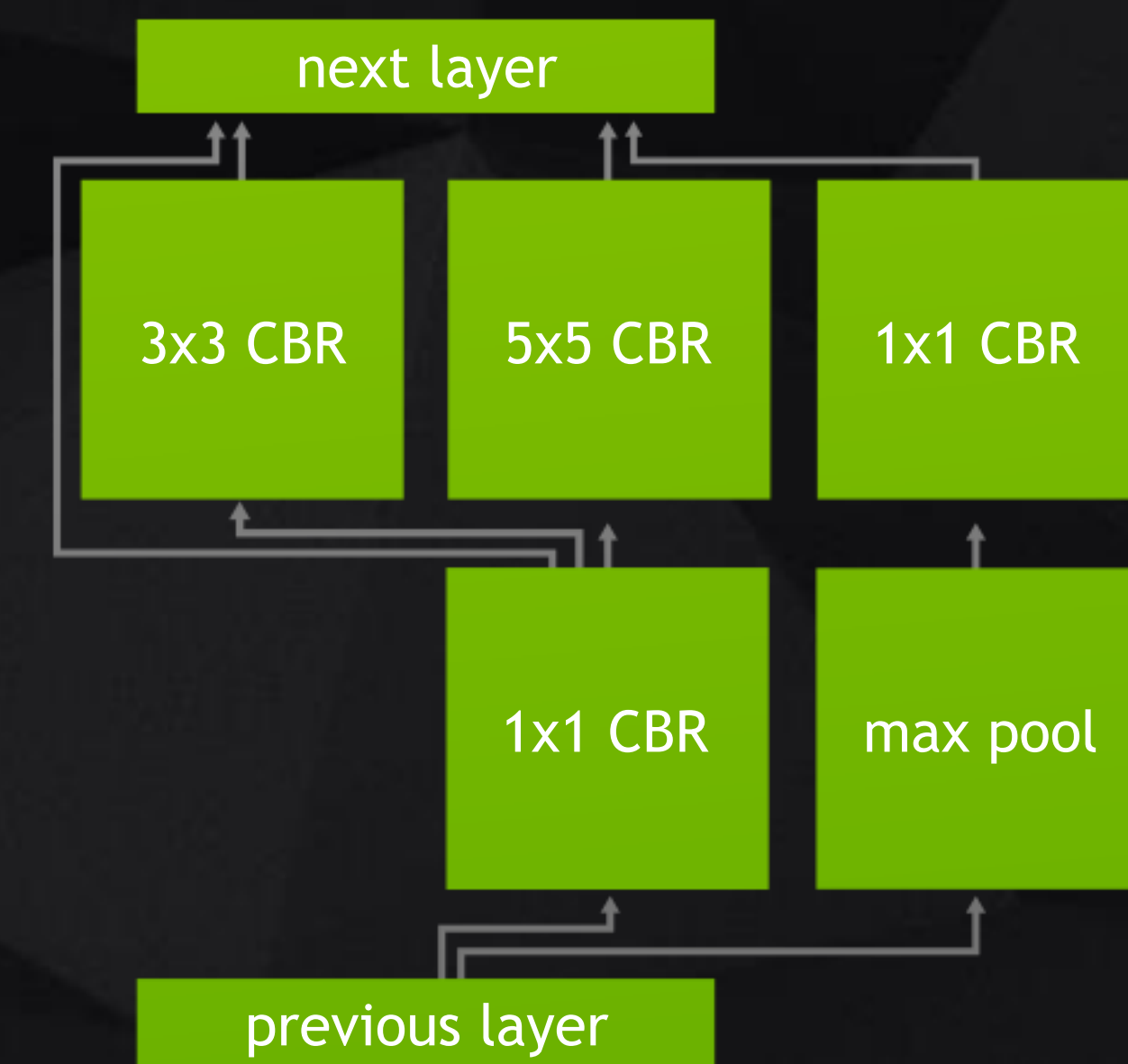
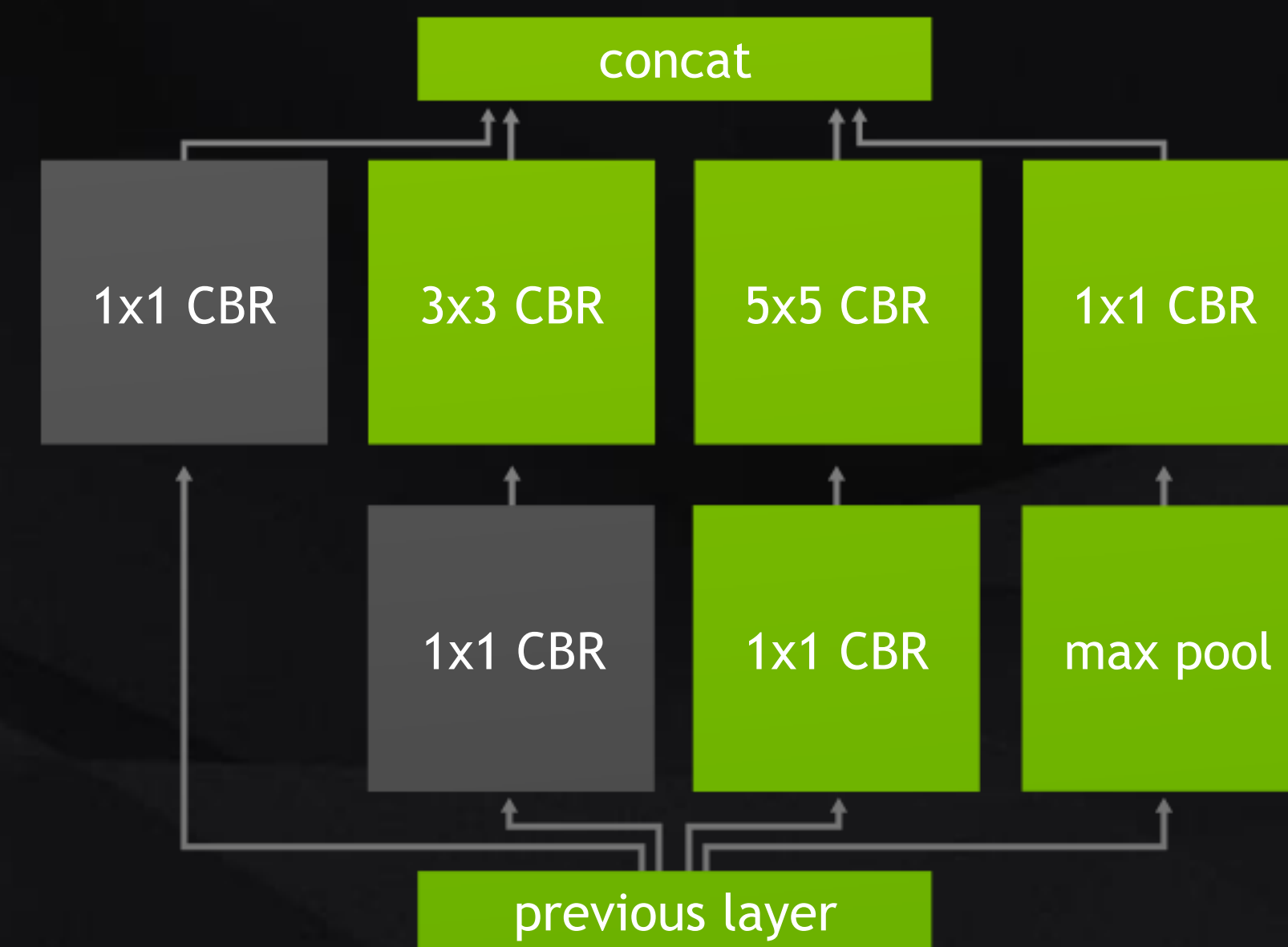
Graph optimizations for vertical and horizontal layer fusion
GPU-specific optimizations
Import models from Caffe and TensorFlow



ANNOUNCING TENSORRT FOR TENSORFLOW

COMPILER FOR DEEP LEARNING INFERENCE

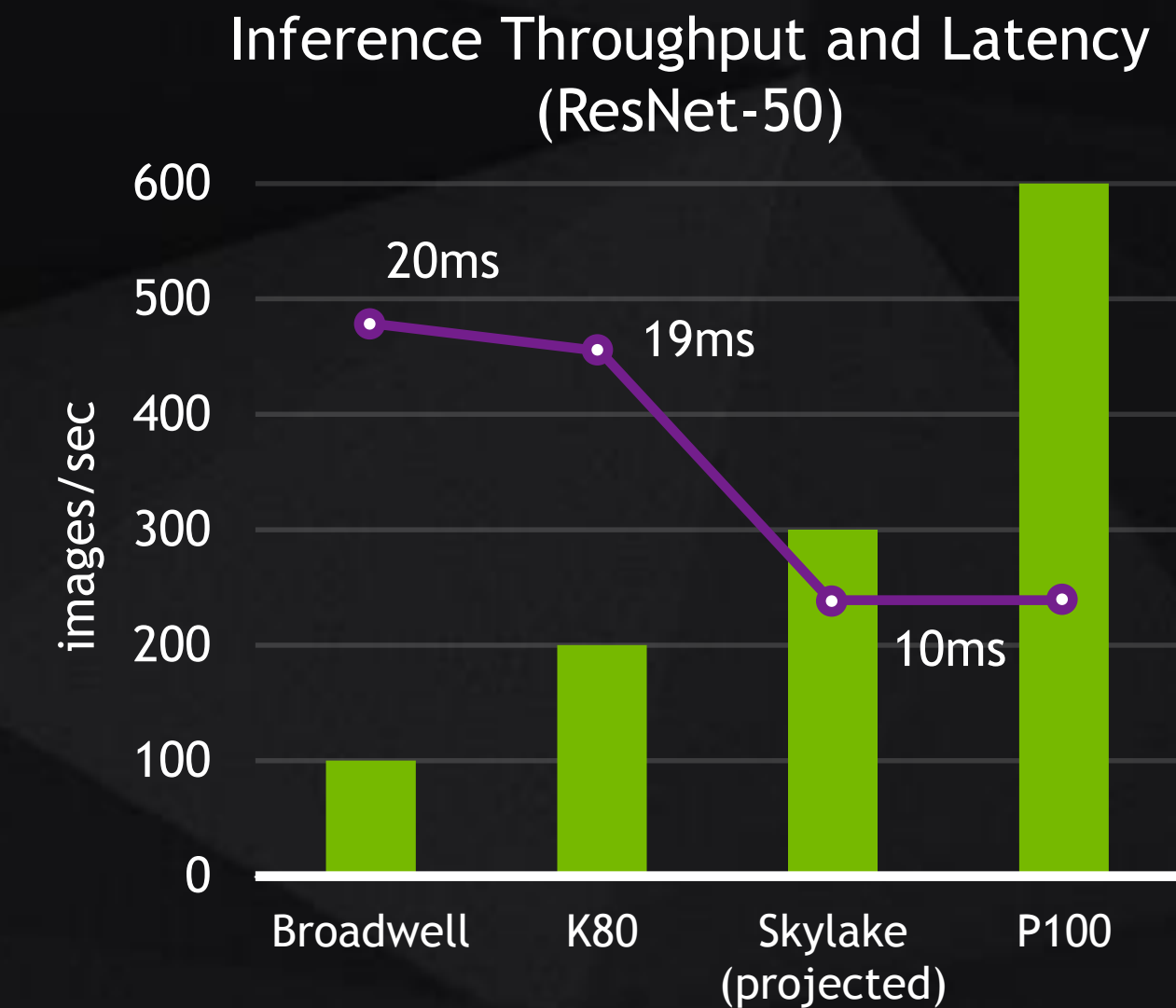
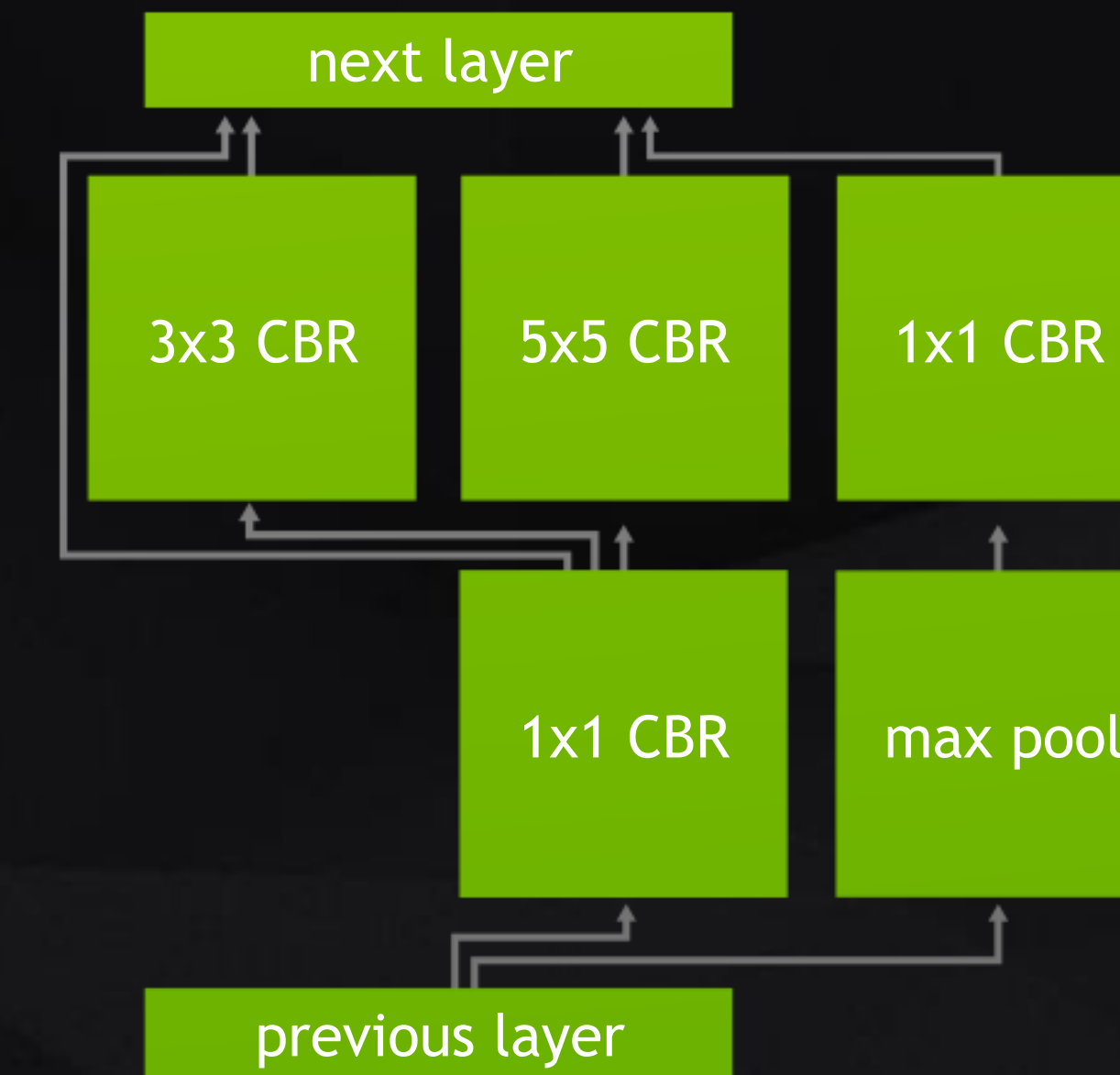
Graph optimizations for vertical and horizontal layer fusion
GPU-specific optimizations
Import models from Caffe and TensorFlow



ANNOUNCING TENSORRT FOR TENSORFLOW

COMPILER FOR DEEP LEARNING INFERENCE

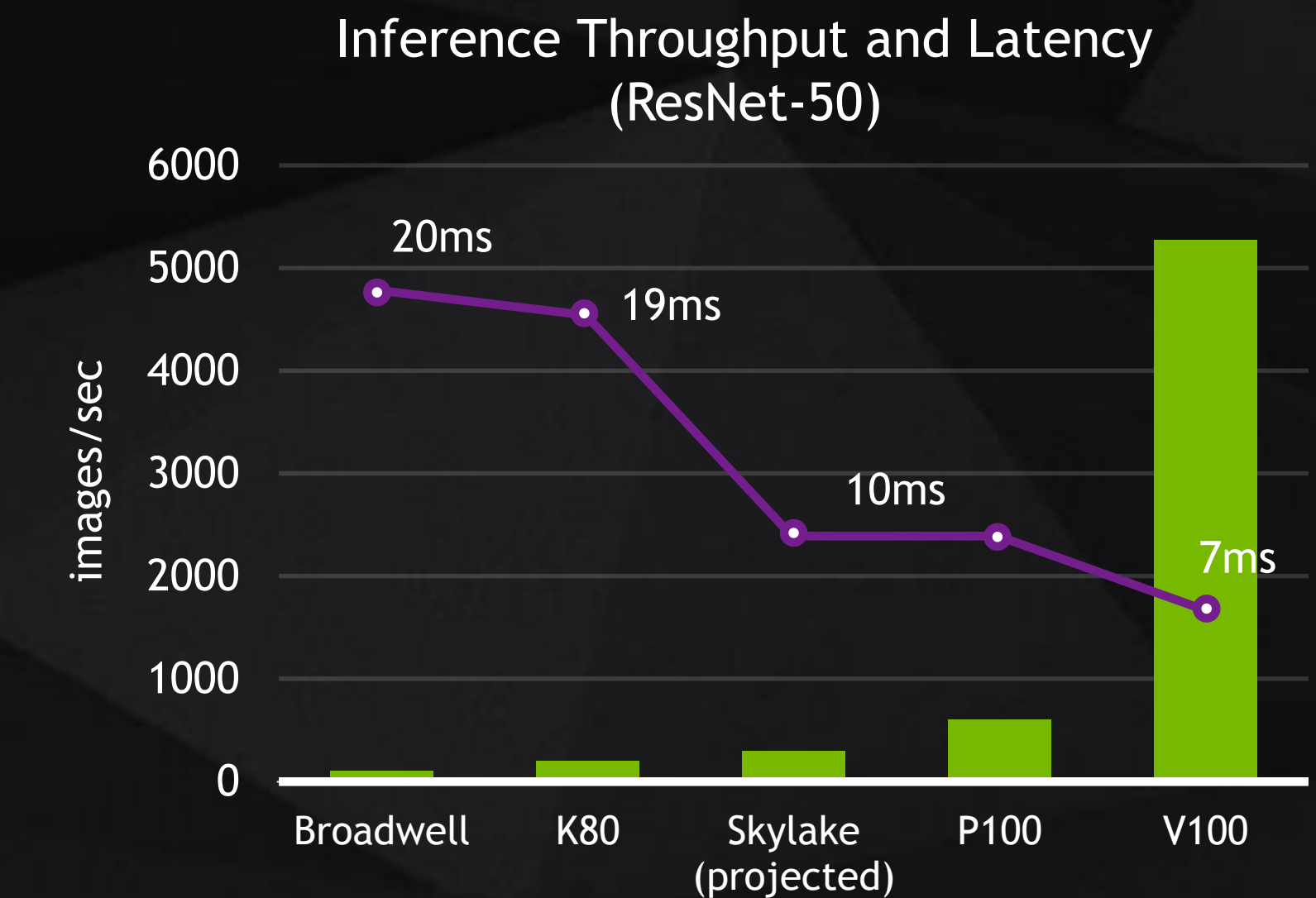
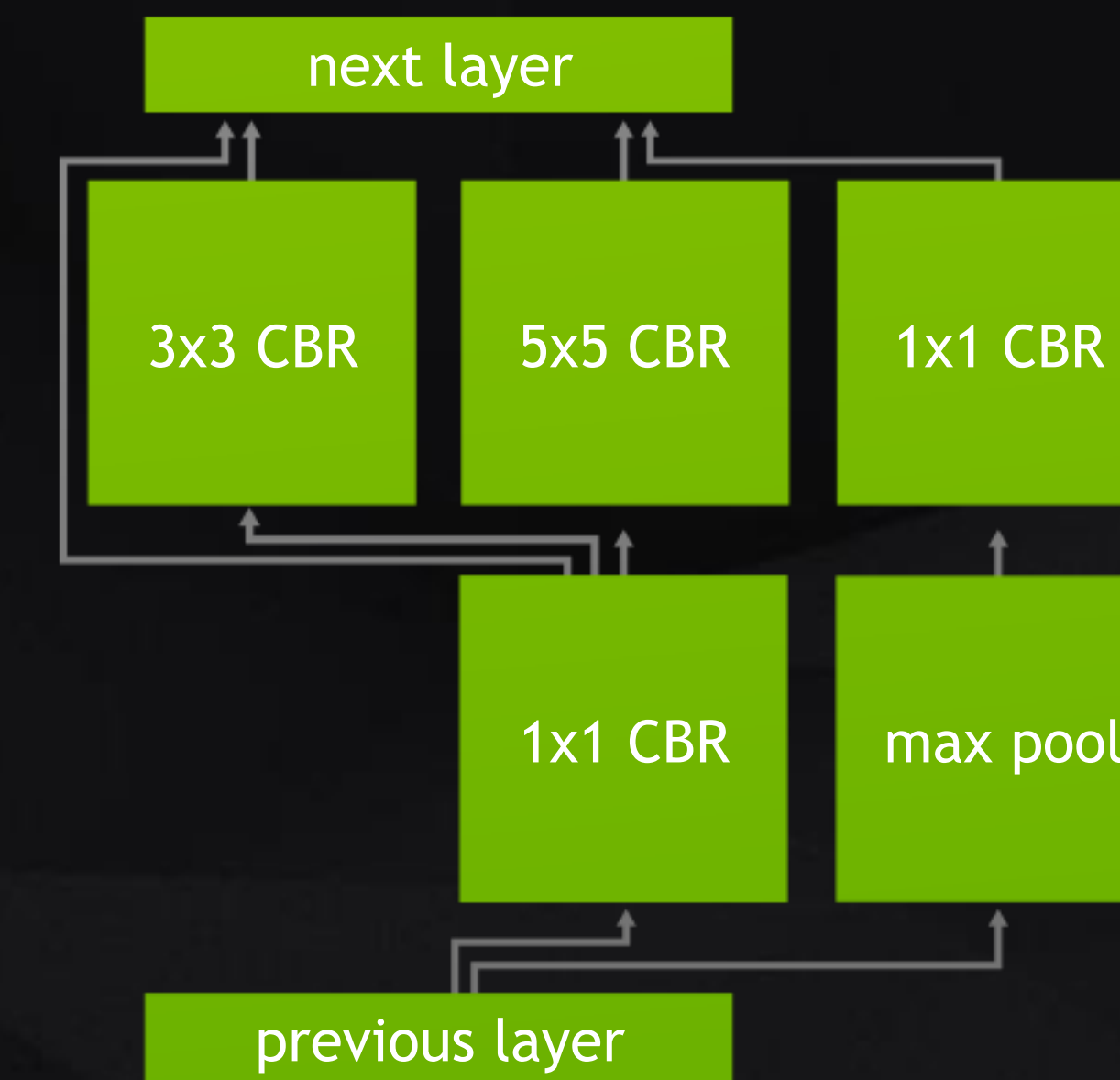
Graph optimizations for vertical and horizontal layer fusion
GPU-specific optimizations
Import models from Caffe and TensorFlow



ANNOUNCING TENSORRT FOR TENSORFLOW

COMPILER FOR DEEP LEARNING INFERENCE

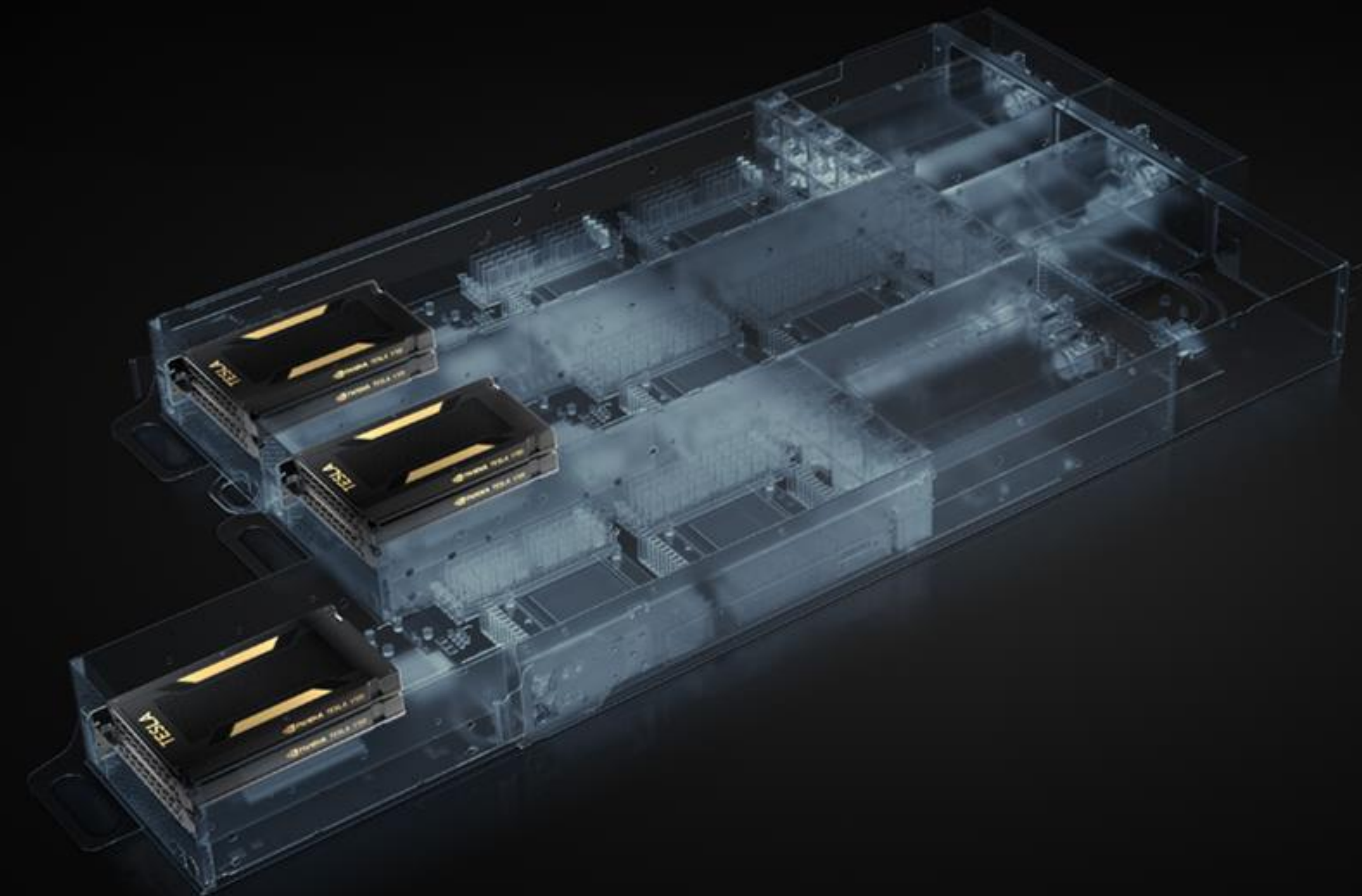
Graph optimizations for vertical and horizontal layer fusion
GPU-specific optimizations
Import models from Caffe and TensorFlow



ANNOUNCING TESLA V100 FOR HYPERSCALE INFERENCE

15-25X INFERENCE SPEED-UP VS SKYLAKE

150W | FHHL PCIe



THE CASE FOR GPU ACCELERATED DATACENTERS



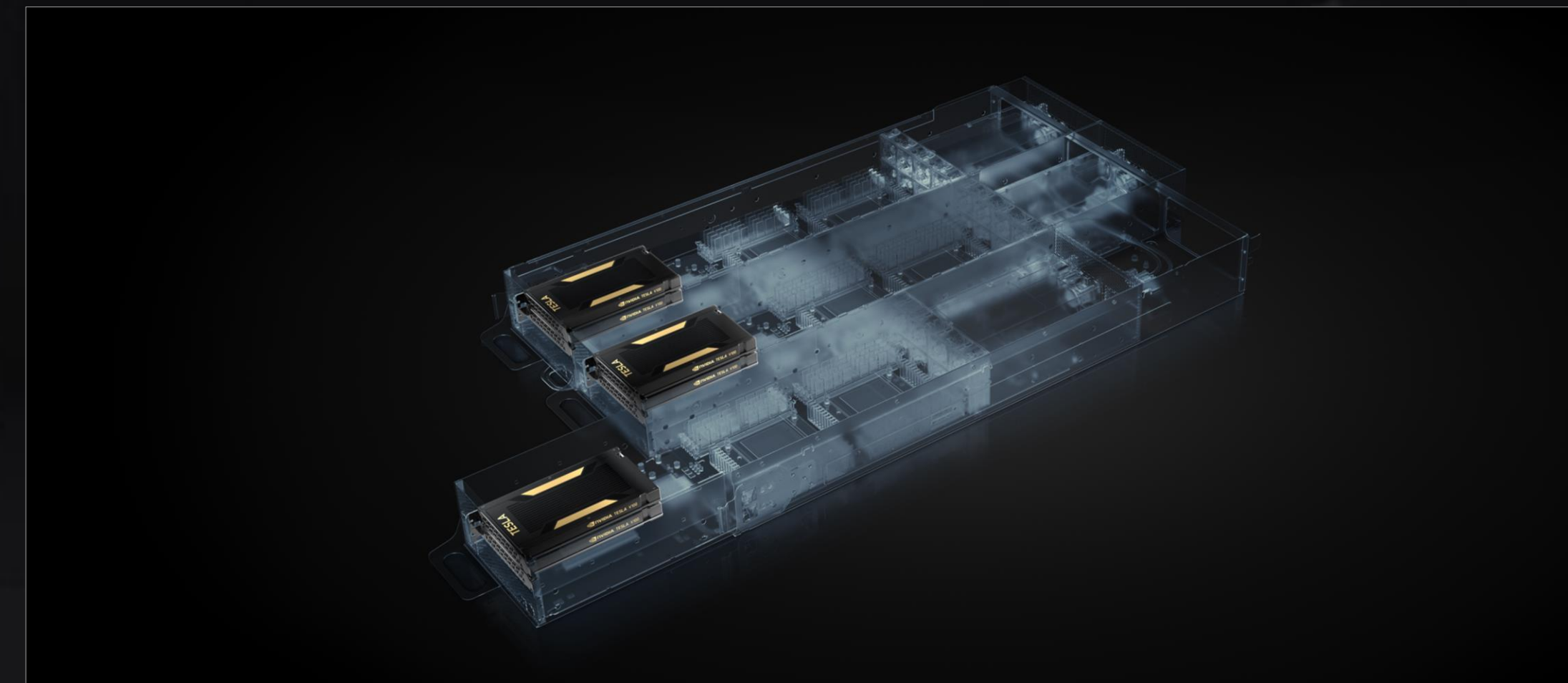
300K inf/s Datacenter

@300 inf/s > 1K CPUs

1K CPUs > 500 Nodes

@\$3K = \$1.5M

@500W = 250KW



500 Nodes CPU Servers



Tesla V100 Reduce 15X



33 Nodes GPU Accelerated Server

NVIDIA DEEP LEARNING STACK



DEEP LEARNING FRAMEWORKS

DEEP LEARNING LIBRARIES

NVIDIA cuDNN, NCCL,
cuBLAS, TensorRT

CUDA DRIVER

OPERATING SYSTEM

GPU

SYSTEM

ANNOUNCING NVIDIA GPU CLOUD

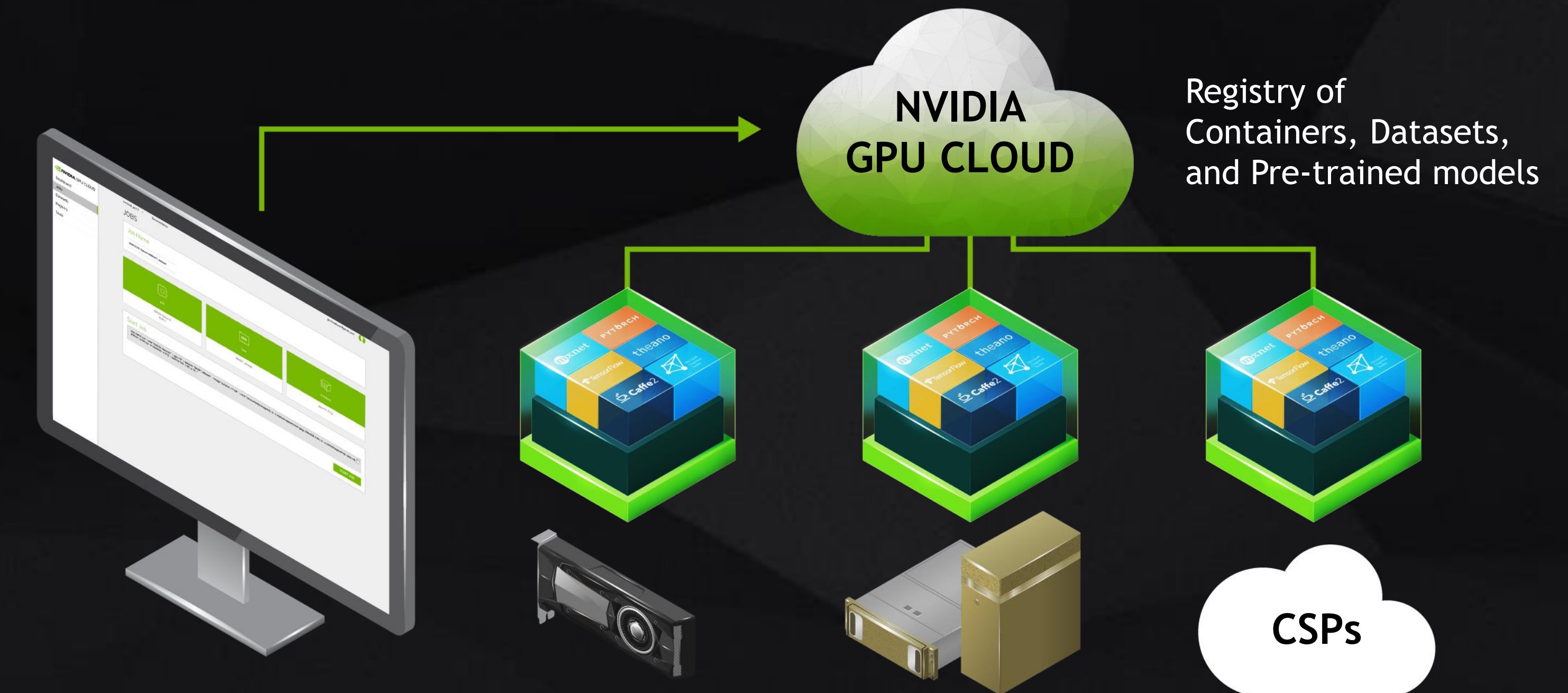
GPU-ACCELERATED CLOUD PLATFORM
OPTIMIZED FOR DEEP LEARNING

Containerized in NVDocker

Optimization across the full stack

Always up-to-date

Fully tested and maintained by NVIDIA



ANNOUNCING NVIDIA GPU CLOUD

GPU-ACCELERATED CLOUD PLATFORM
OPTIMIZED FOR DEEP LEARNING

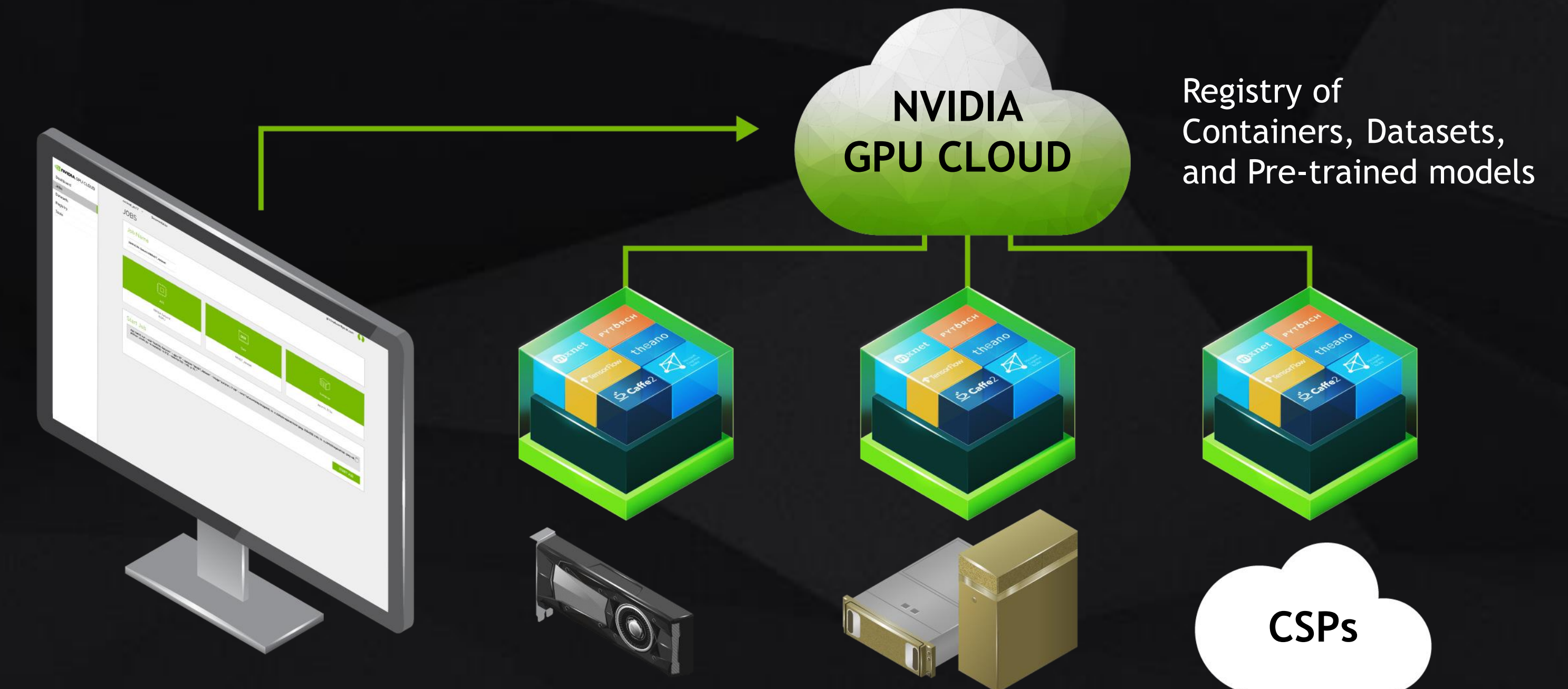
Containerized in NVDocker

Optimization across the full stack

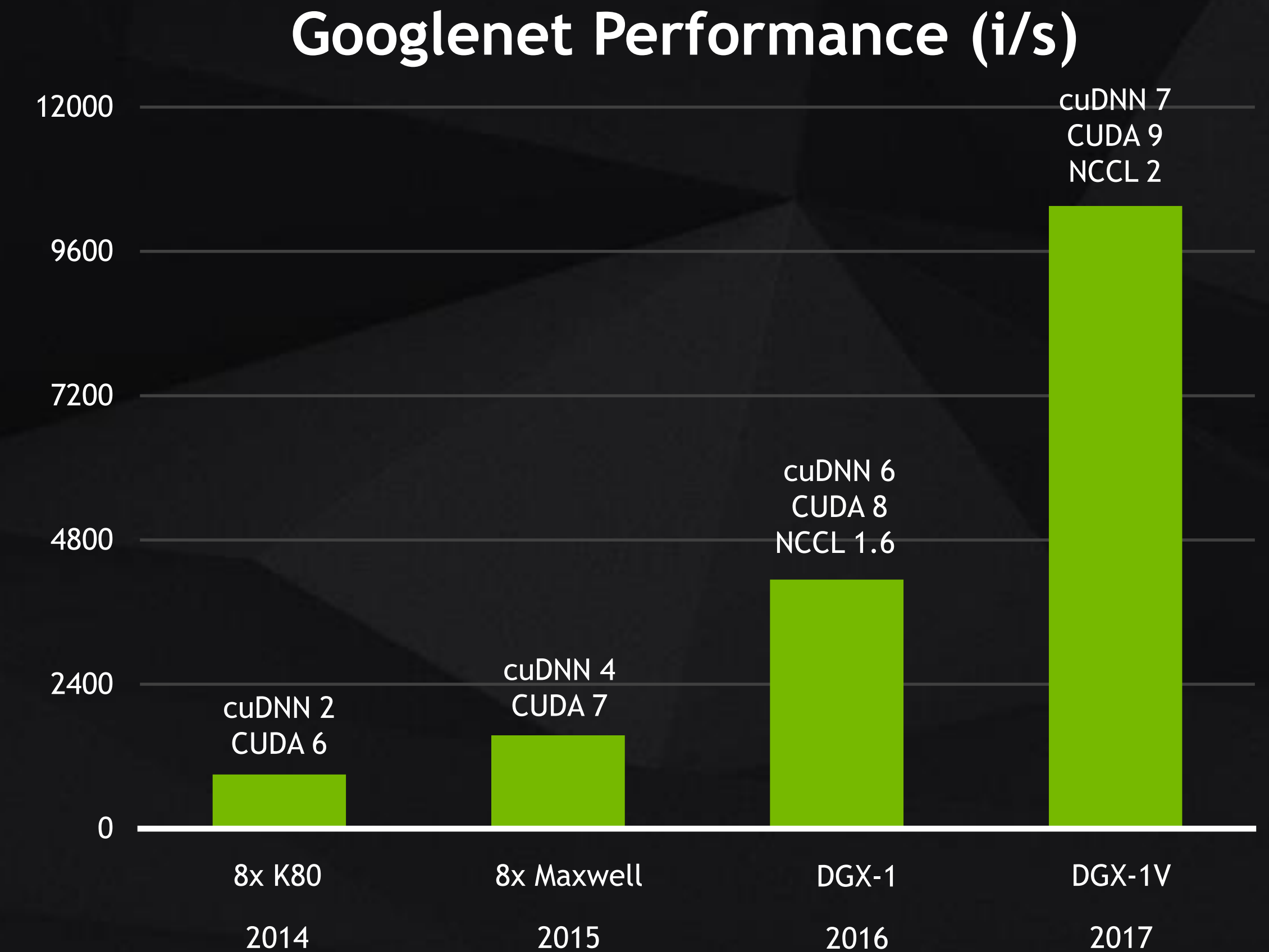
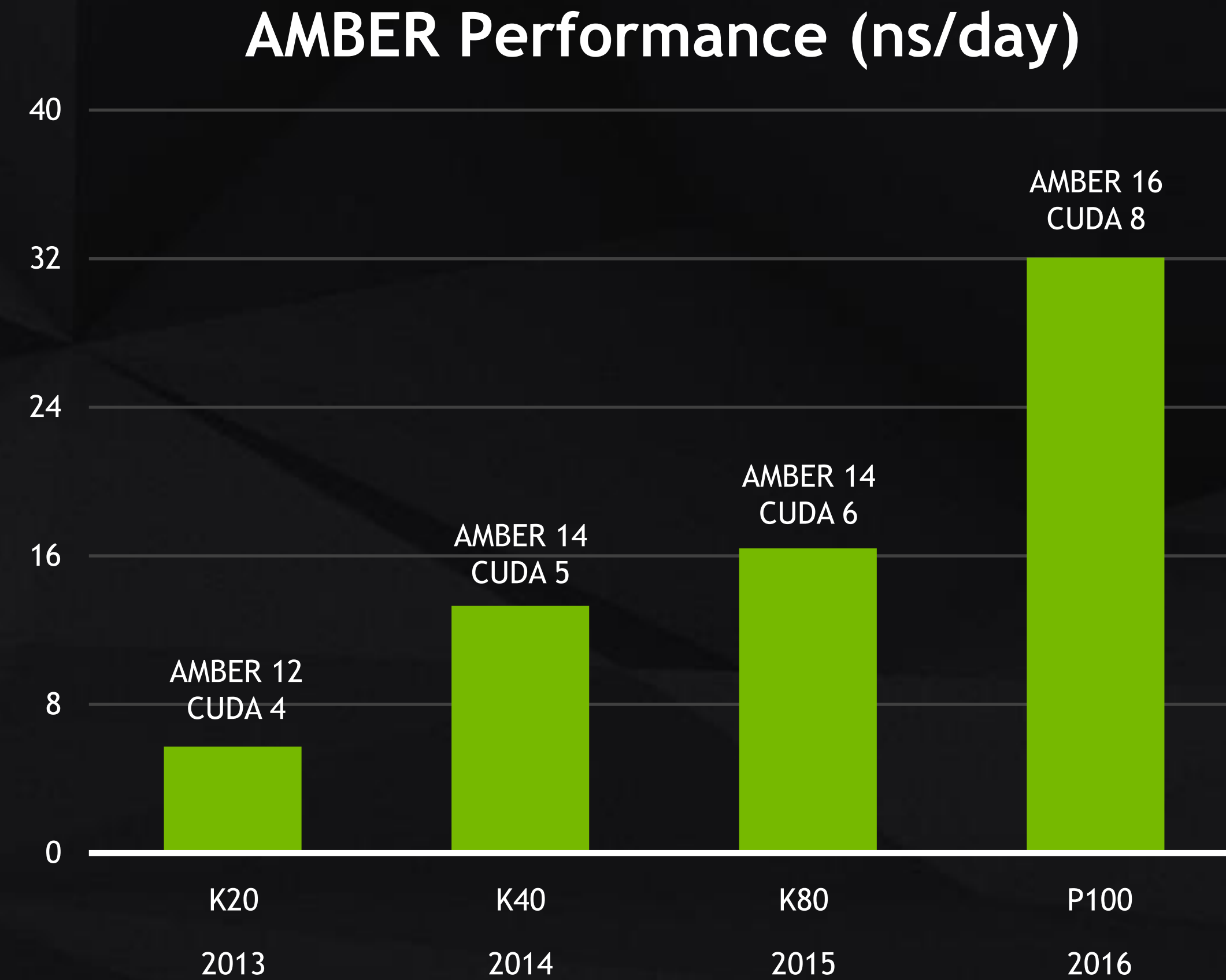
Always up-to-date

Fully tested and maintained by NVIDIA

Beta in July

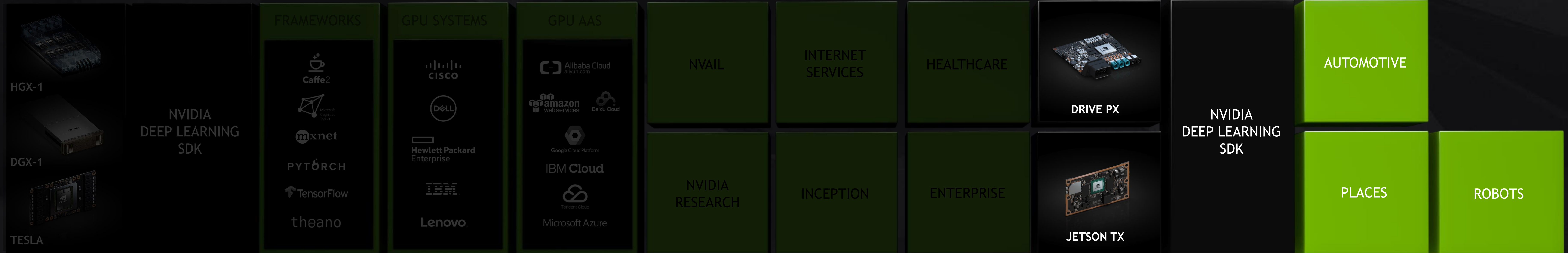


POWER OF GPU COMPUTING



POWERING THE AI REVOLUTION

AI at the Edge



AI REVOLUTIONIZING TRANSPORTATION



280B Miles per Year



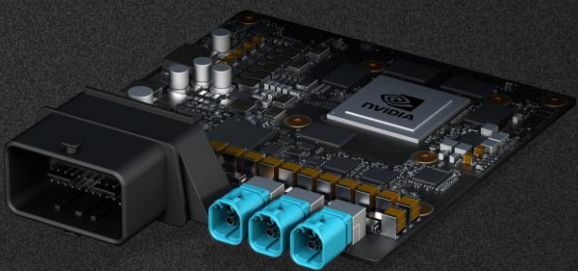
800M Parking Spots for 250M Cars in U.S.



Domino's: 1M Pizzas Delivered per Day

NVIDIA DRIVE — AI CAR PLATFORM

100 TOPS



DRIVE PX Xavier
Level 4/5

10 TOPS

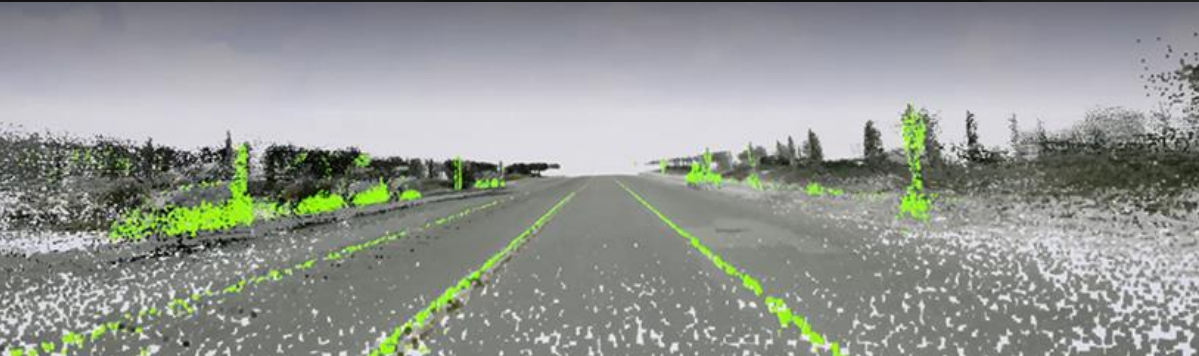


DRIVE PX 2 Parker
Level 2/3

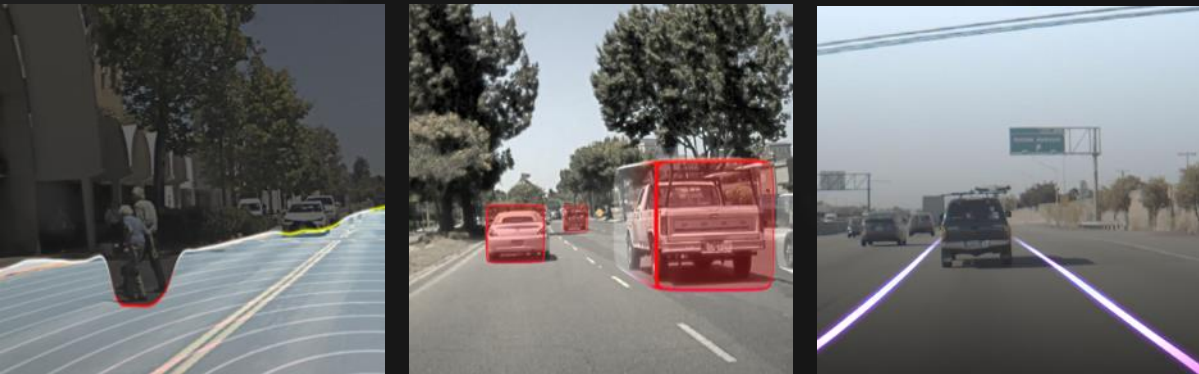
1 TOPS

Localization

Path Planning



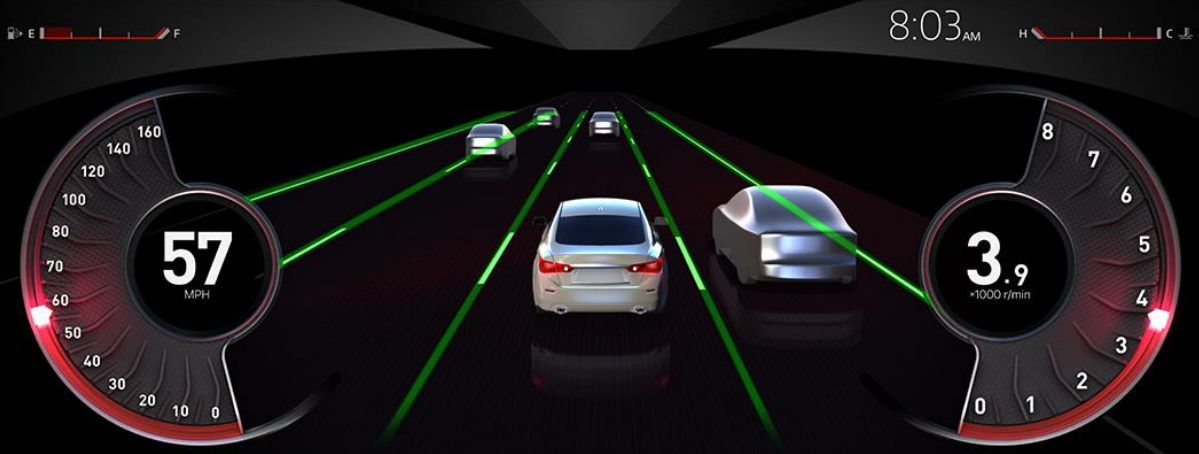
Perception AI



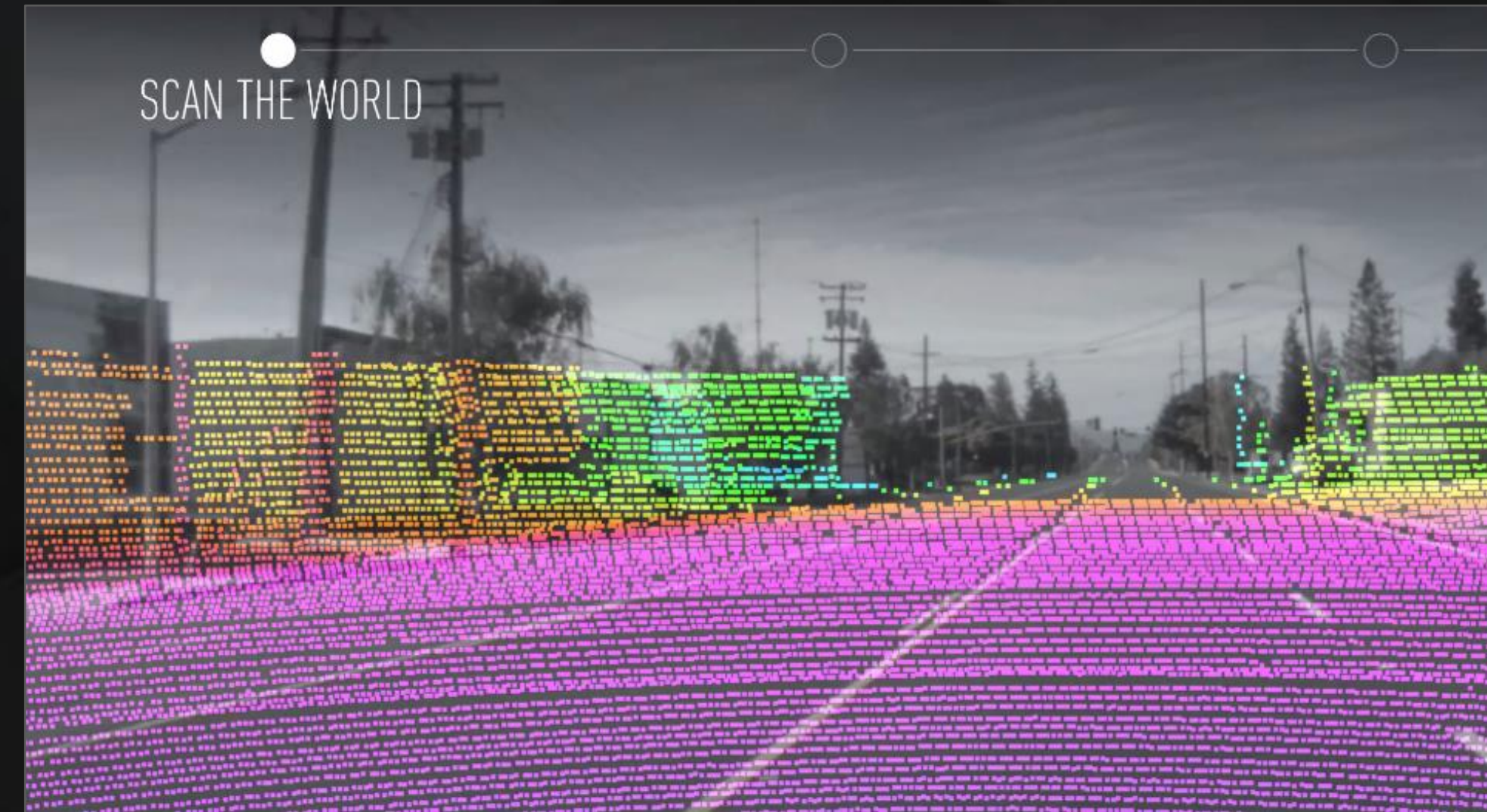
Computer Vision Libraries

CUDA, cuDNN, TensorRT

OS



NVIDIA DRIVE



Mapping-to-Driving



Co-Pilot



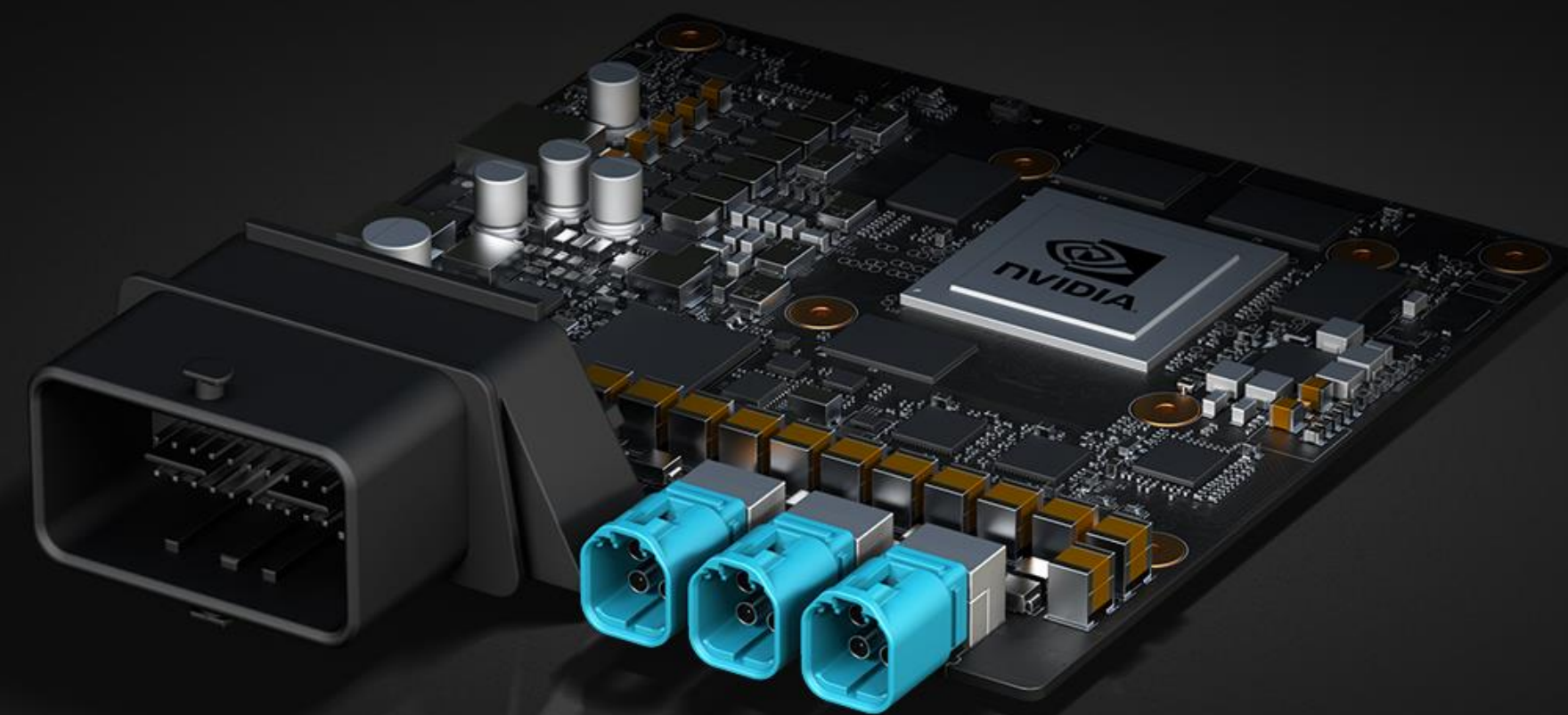
Guardian Angel

**ANNOUNCING
TOYOTA SELECTS
NVIDIA DRIVE PX FOR
AUTONOMOUS VEHICLES**

TOYOTA



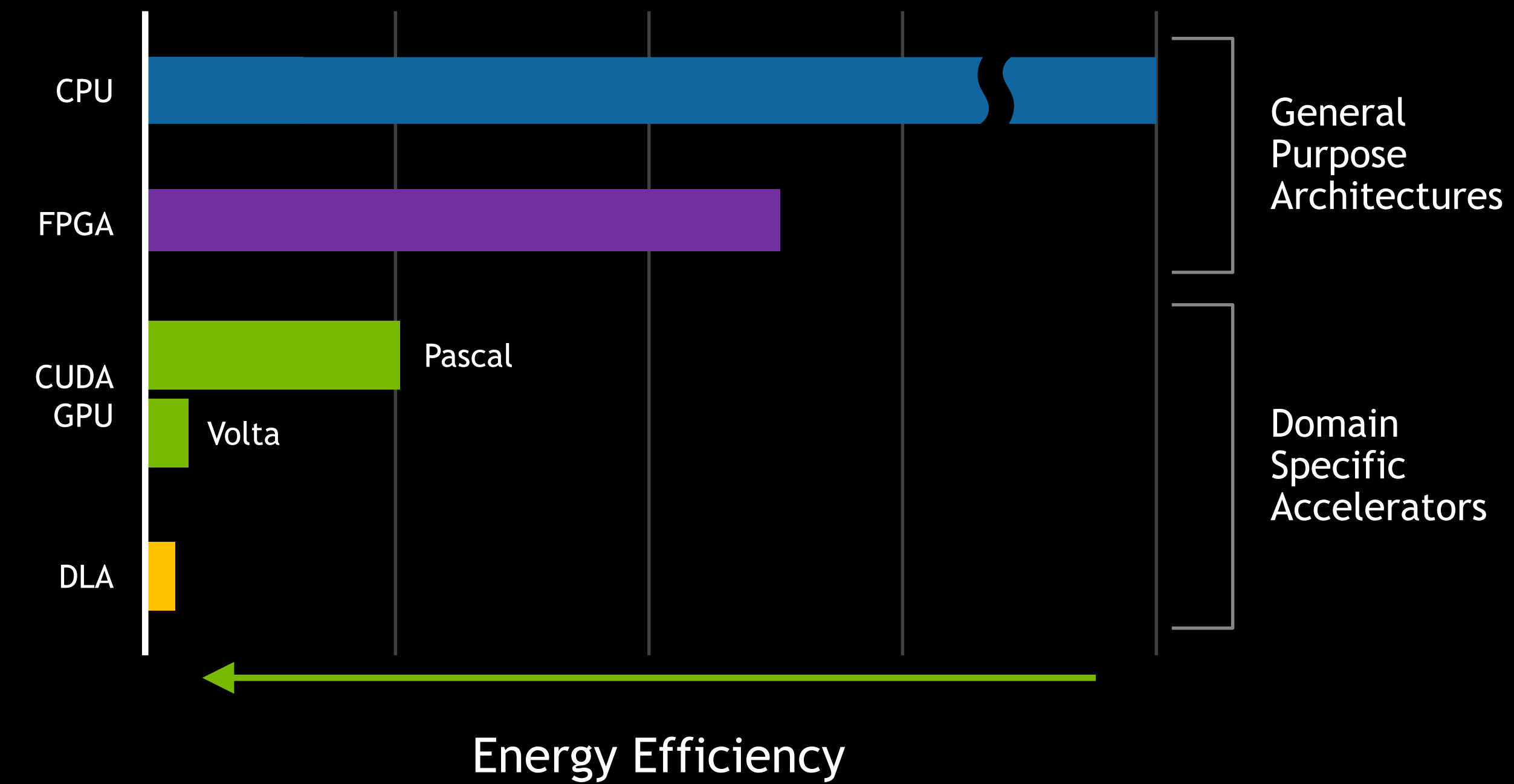
AI PROCESSOR FOR AUTONOMOUS MACHINES



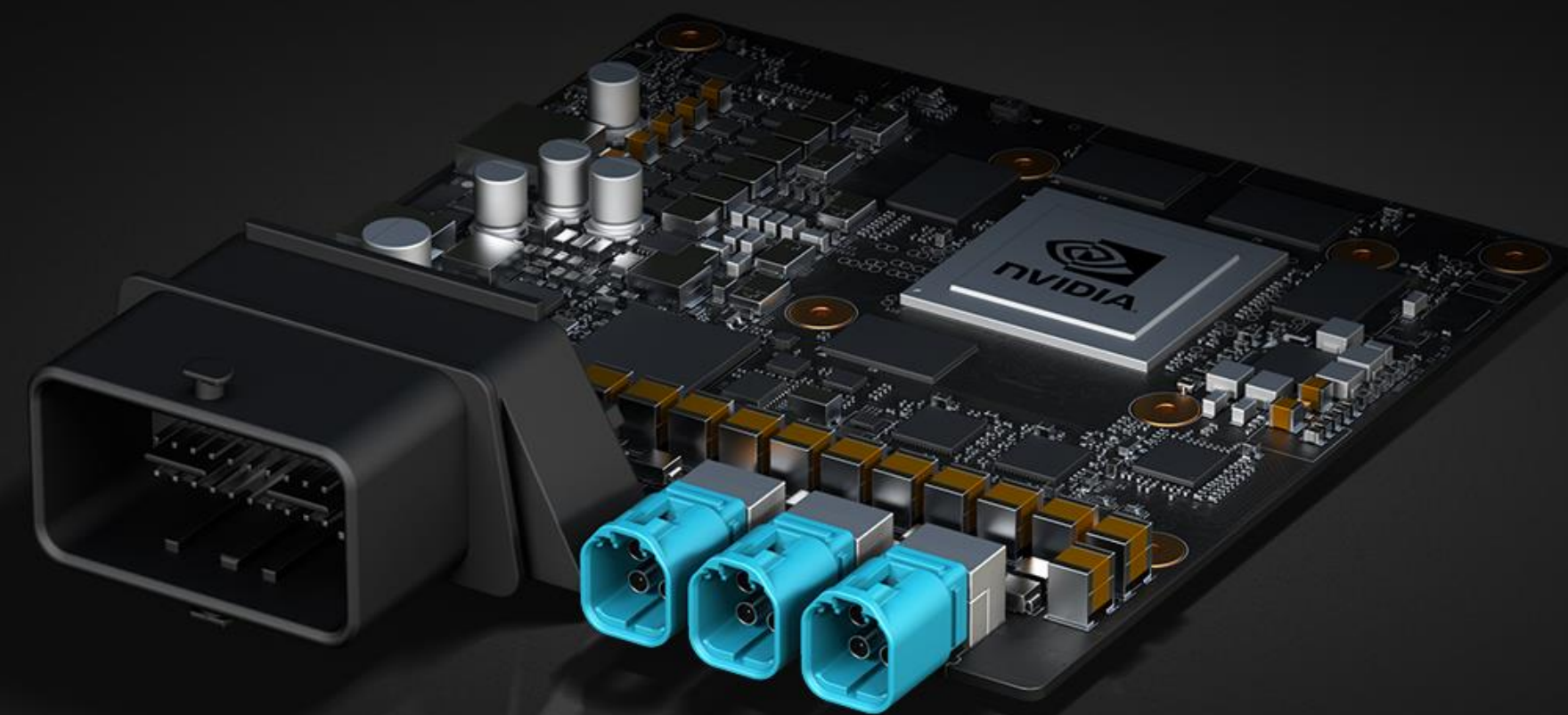
XAVIER

30 TOPS DL | 30W

Custom ARM64 CPU | 512 Core Volta GPU | 10 TOPS DL Accelerator



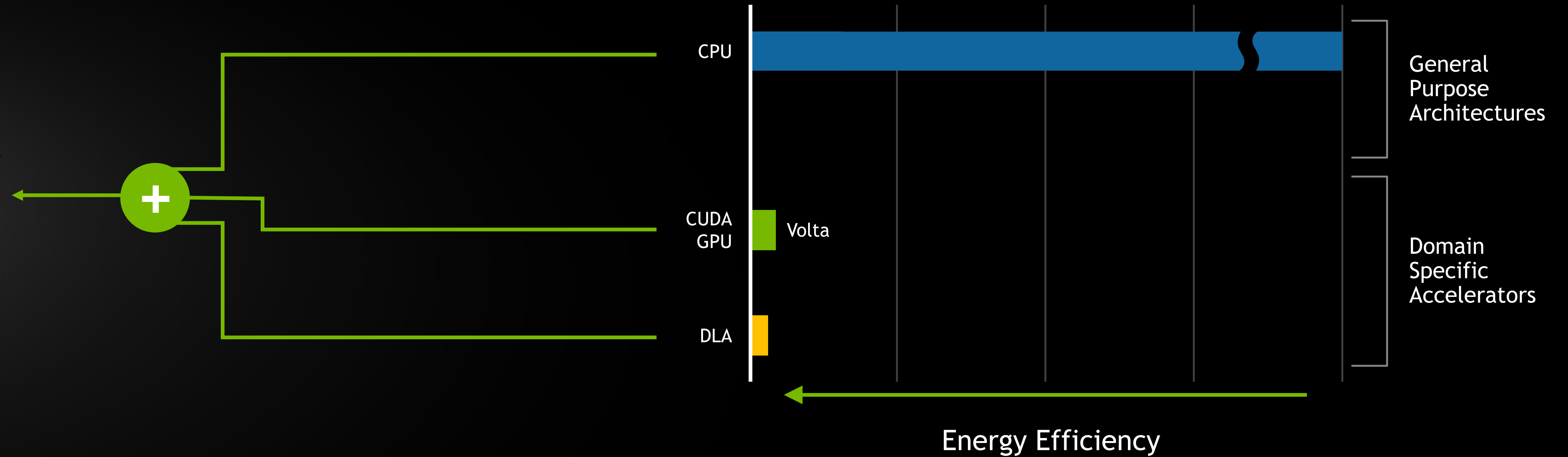
AI PROCESSOR FOR AUTONOMOUS MACHINES



XAVIER

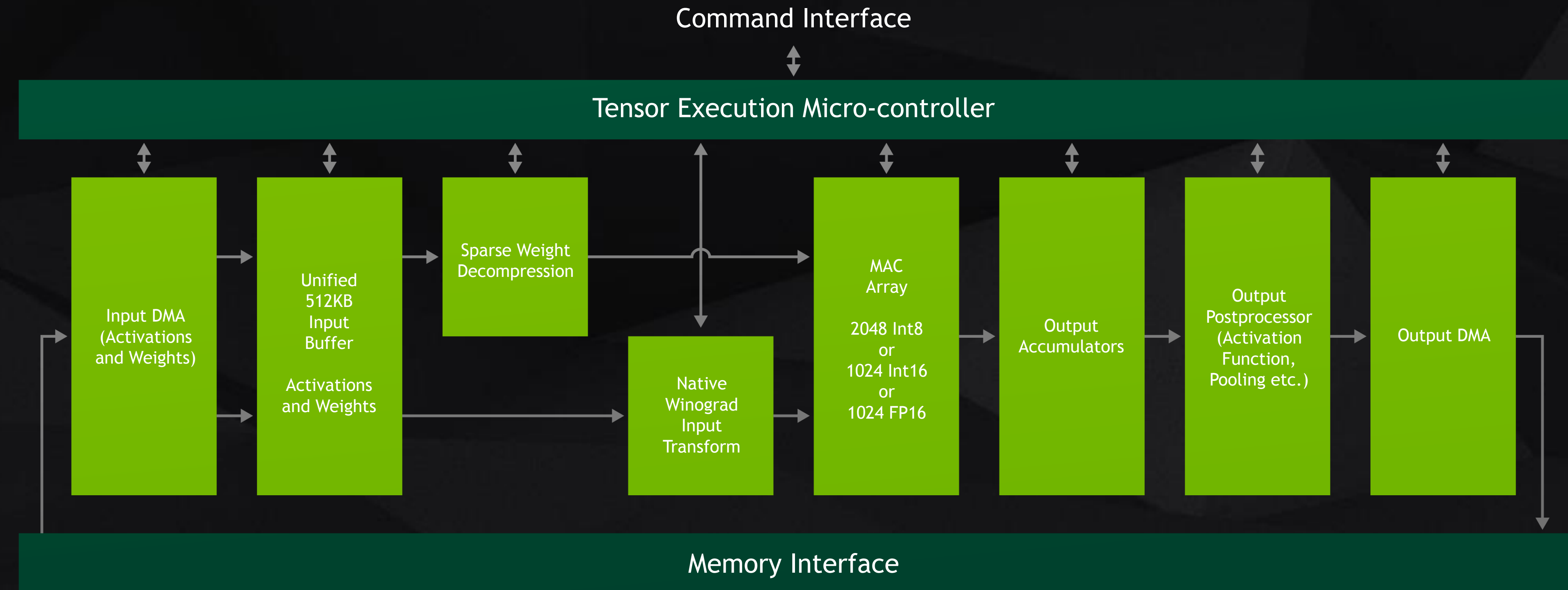
30 TOPS DL | 30W

Custom ARM64 CPU | 512 Core Volta GPU | 10 TOPS DL Accelerator

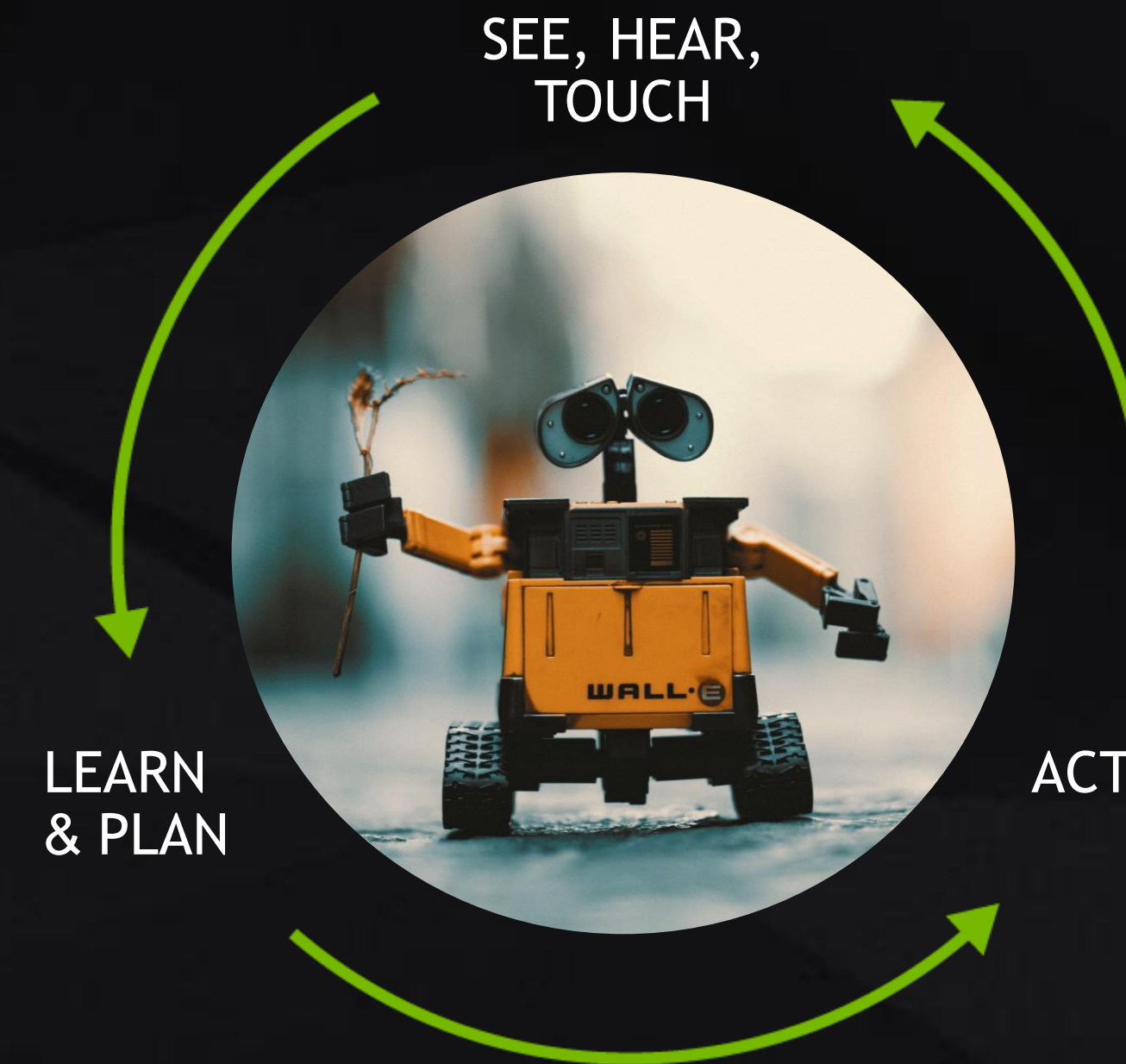


ANNOUNCING XAVIER DLA NOW OPEN SOURCE

Early Access July
General Release September



ROBOTS

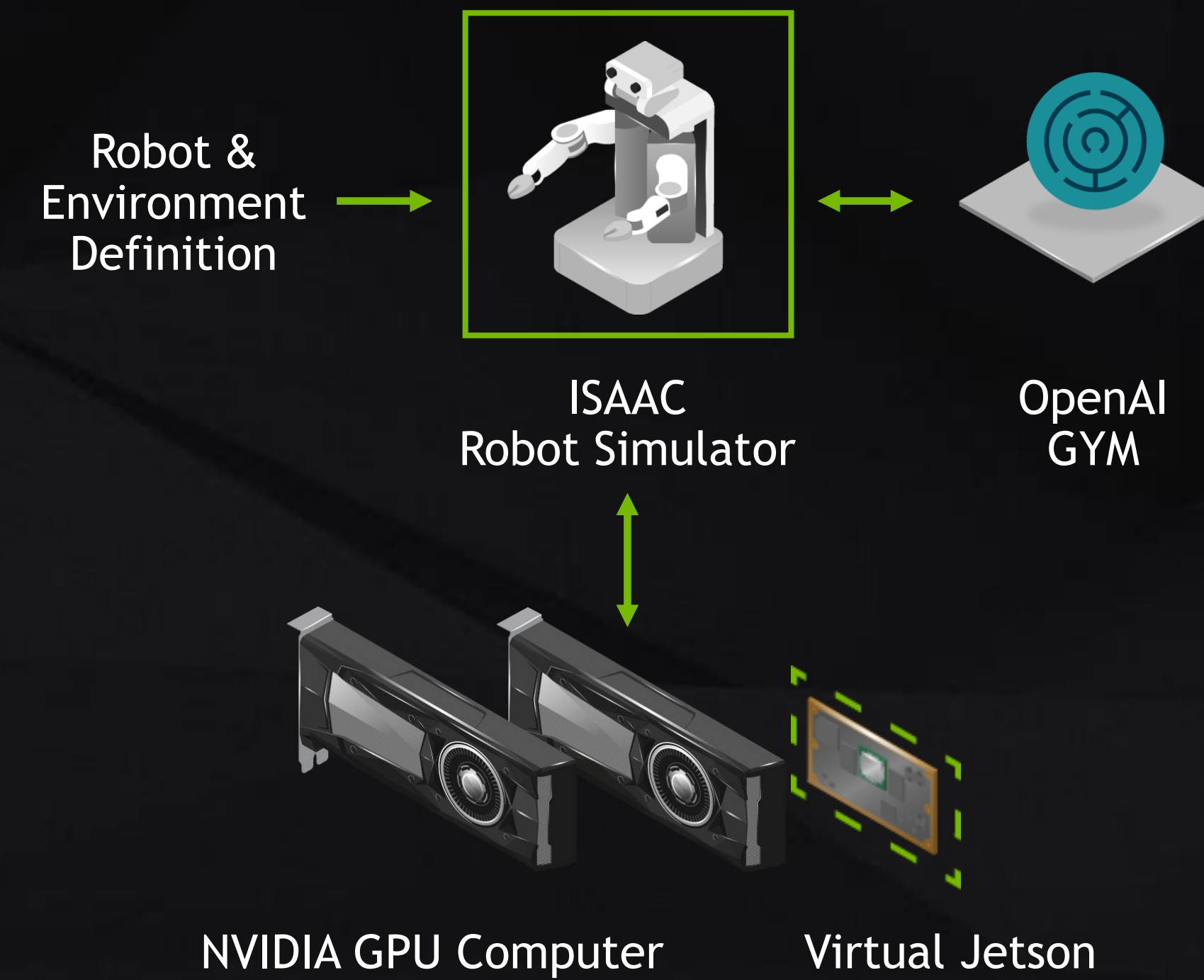


ROBOTS

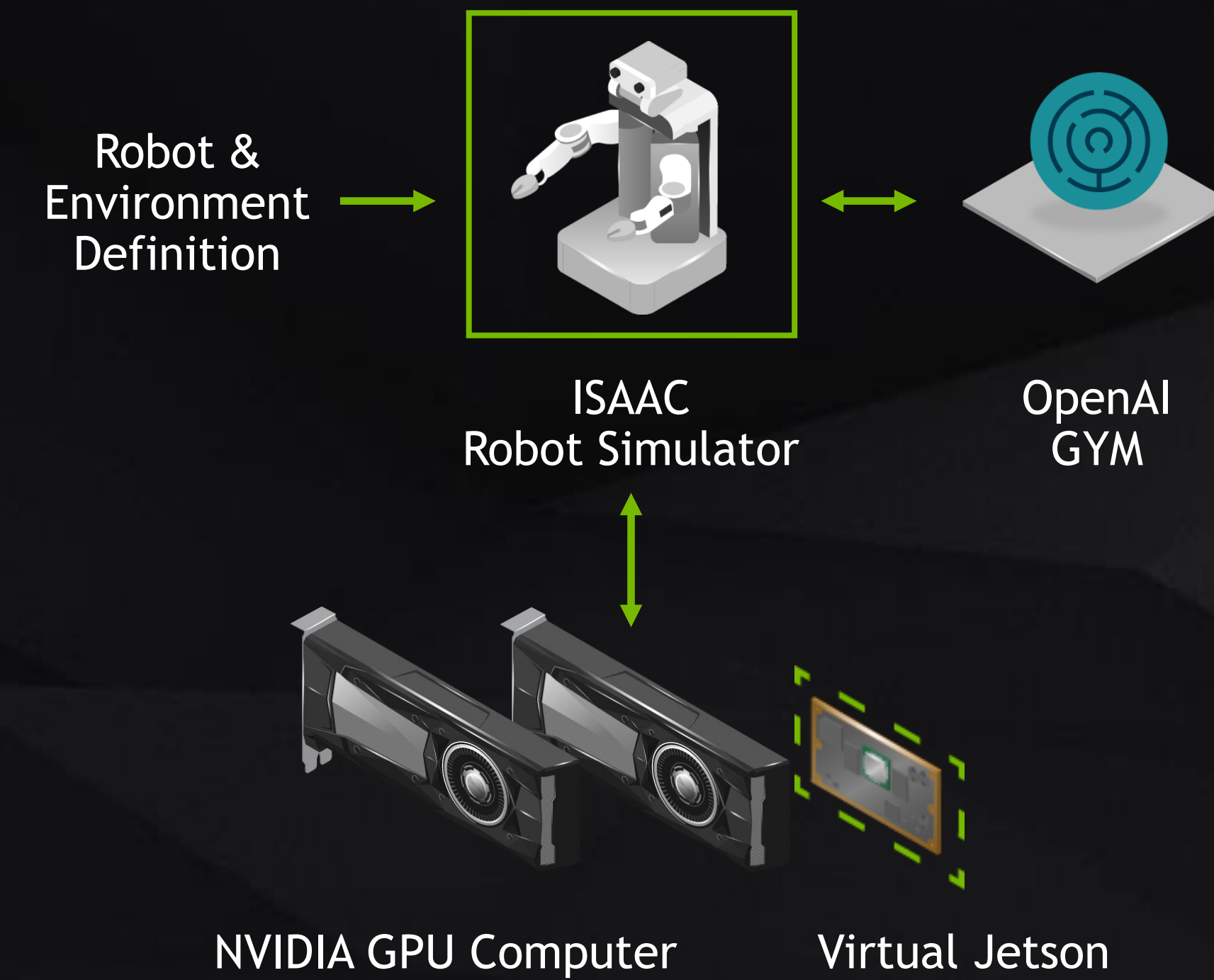
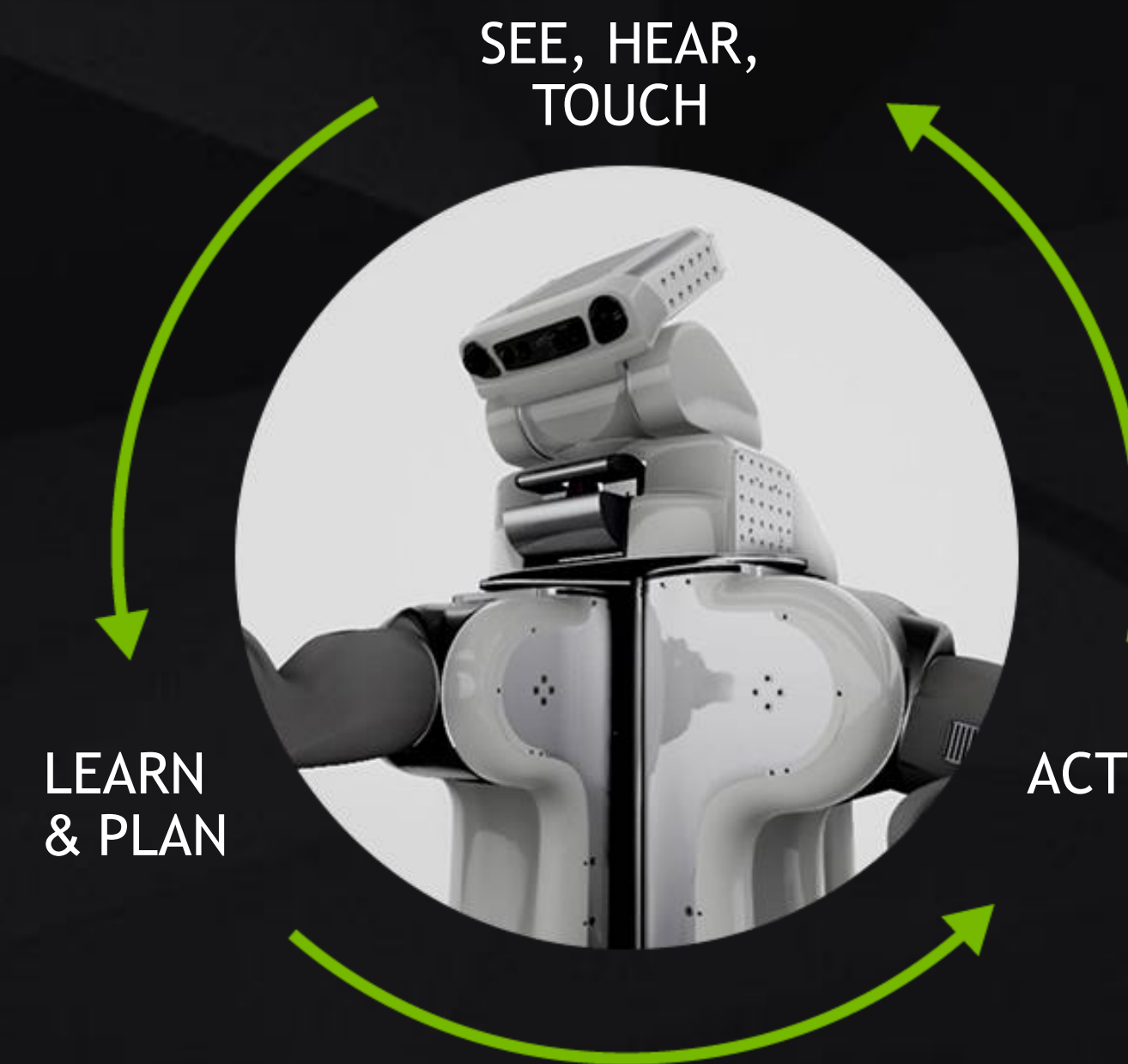


*Credit: Yevgen Chebotar, Karol Hausman,
Marvin Zhang, Sergey Levine*

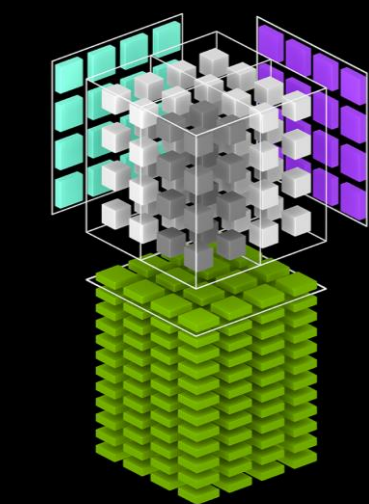
ANNOUNCING ISAAC ROBOT SIMULATOR



ANNOUNCING ISAAC ROBOT SIMULATOR



NVIDIA POWERING THE AI REVOLUTION

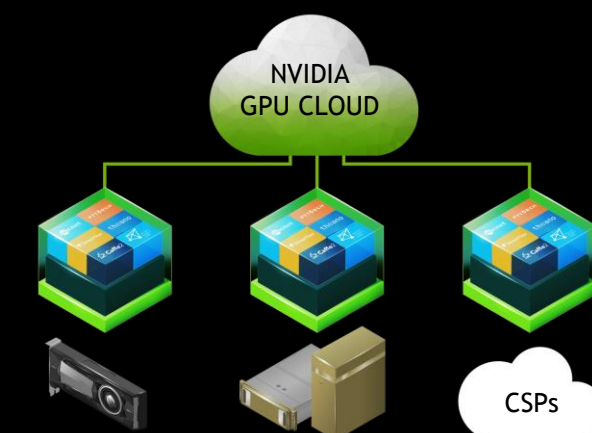


Tensor Core

TensorRT

Alibaba amazon Bai 百度 facebook
Google Microsoft Tencent

NVIDIA GPU in Every Cloud



NVIDIA GPU Cloud



Xavier DLA
Open Source

