# GPU TECHNOLOGY CONFERENCE

## A NEW COMPUTING MODEL

JEN-HSUN HUANG, CO-FOUNDER & CEO | GTC 2016

# LEAPS IN ADOPTION

**2X GTC Attendees**

2012: 2,350
2016: 5,500

Legend: Auto, Internet, Gov't / Labs, Academia, M&E, Finance, Aerospace / Defense, Manufacturing, Oil & Gas, IT / Hardware / Software, Medical

**2X Accelerated Systems, 96% of New Systems on NVIDIA**

TOP 500 The List.

# accelerated systems: Nov 2013, Nov 2014, Nov 2015

**4X CUDA Developers, 10X in Hyperscale + Auto**

2012 → 4x → 2016: 300K

Legend: Academia, Games, Finance, Manufacturing, Internet, Oil & Gas, National Labs, Automotive, Defense, M & E

# 5 THINGS

A Toolbox

VR

A Deep Learning Chip

A Deep Learning Box

A Deep Learning Car

COMPUTE**WORKS**    GAME**WORKS**    VR**WORKS**    DESIGN**WORKS**    DRIVE**WORKS**    JETPACK

PhysX

Hair**Works**

Wave**Works**

Flame**Works**

And other technologies such as:
**Clothing, VXGI, Flex, Destruction**

## NVIDIA GAMEWORKS

Volumetric Lighting

Voxel Accelerated Ambient Occlusion

Hybrid Frustum Traced Shadows

Available Now

CUDA

cuDNN

nvGRAPH

IndeX

And other technologies such as:
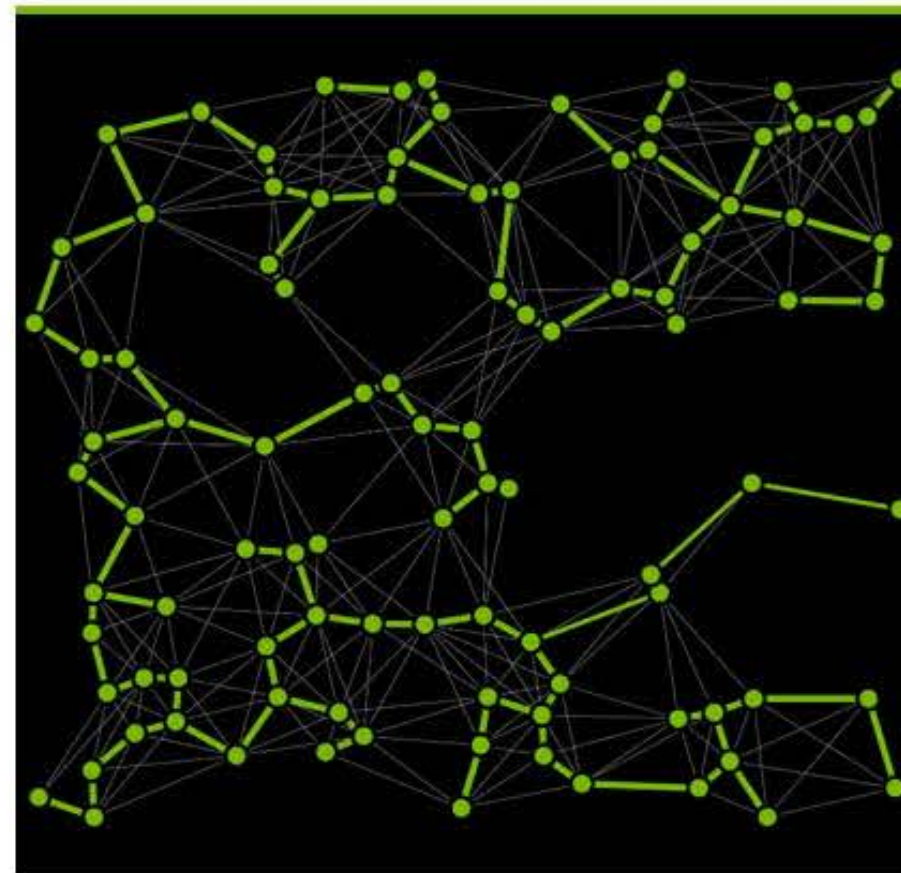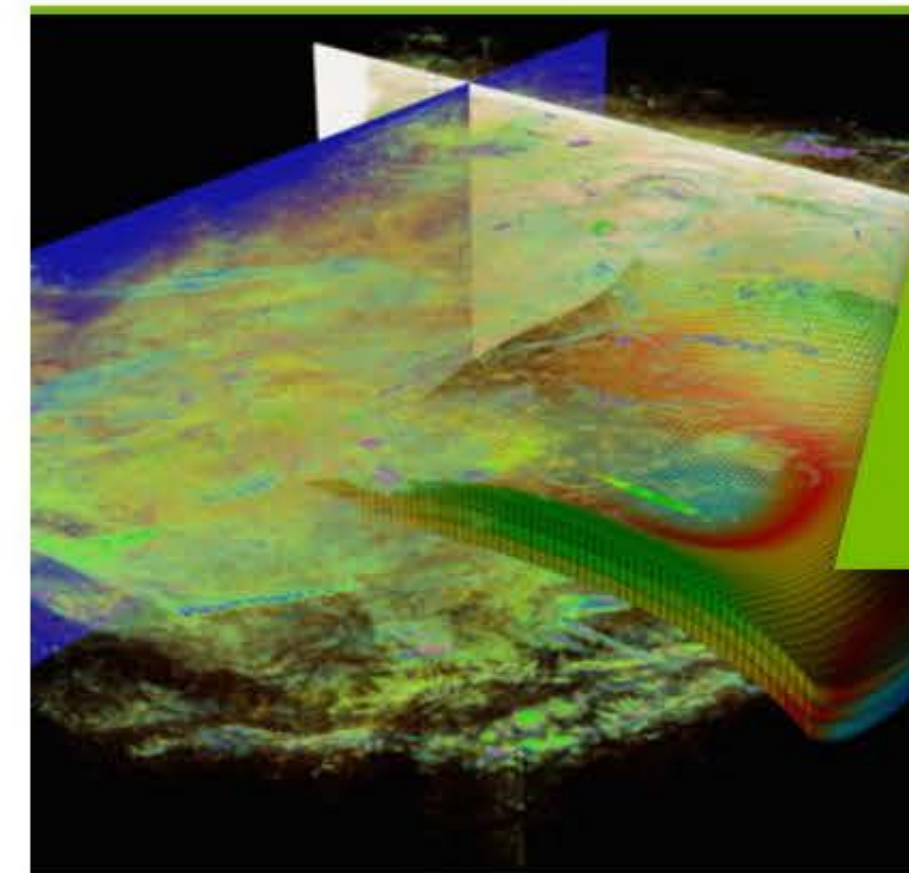**AMGx, cuSOLVER, cuSPARSE, OpenACC, NSIGHT, THRUST**

# NVIDIA COMPUTEWORKS

CUDA 8 — Available June

cuDNN 5 — Available April

nvGRAPH — Available June

IndeX plug-in for ParaView — Available May

COMPUTE**WORKS** | GAME**WORKS** | VR**WORKS** | DESIGN**WORKS** | DRIVE**WORKS** | JET**PACK**

**Deep Learning SDK**

**DIGITS Workflow**

**VisionWorks**

**Jetson Media SDK**

and other technologies such as:
**Linux4Tegra, NSIGHT EE, OpenCV4Tegra, OpenGL, System Trace, Visual Profiler, Vulkan**

# NVIDIA JETPACK

GIE - GPU Inference Engine — Available May

Jetson TX1: 24 images/s/W

# 5 THINGS



**NVIDIA SDK**

**A Deep Learning Chip**

**A Deep Learning Box**

**VR**

**A Deep Learning Car**

# A START OF A NEW PLATFORM

Samsung, Oculus, HTC release headsets

Google announces Jump VR camera platform

New York Times ships Cardboard to subscribers

Microsoft demonstrates Holoportation

VR Startups Raise $1.5B in funding

**EVEREST VR**

MARS 2030

Pre-render light probes surrounding region of interest

Rasterize depth buffer at headset eye positions

Reconstruct image for new viewpoint from depth and multiple probes

IRAY VR
BREAKTHROUGH
PHOTOREAL VR

Available starting in June

**IRAY VR**

# IRAY VR LITE

1. Design in 3ds Max

2. Download Iray for 3ds Max Plug-in

NVIDIA Iray VR Lite

3. Download Android Viewer

4. Get VR HMD

*Available in June*

# 5 THINGS

**NVIDIA SDK**

**A Deep Learning Chip**

**A Deep Learning Box**

**IRAY VR**

**A Deep Learning Car**

# AN AMAZING YEAR IN AI



Microsoft & Google
"Superhuman" Image Recognition



Microsoft
"Super Deep Network"



Berkeley's Brett
One network, everything robotics



Deep Speech 2
One network, 2 languages



A New Computing Model
Hits Pop Culture



AlphaGo
Rivals a World Champion

# A NEW COMPUTING MODEL

**Traditional Computer Vision**
Experts + Time

**Deep Learning Object Detection**
DNN + Data + HPC

**Deep Learning Achieves "Superhuman" Results**

# THE EXPANDING UNIVERSE
# OF MODERN AI

**INDUSTRY LEADERS**

**SIEME**

Ford

TAR

**Alibaba.com**

GE

TES

AstraZeneca

gsk

**START-UPS**

Morpho

Tech
computer vision
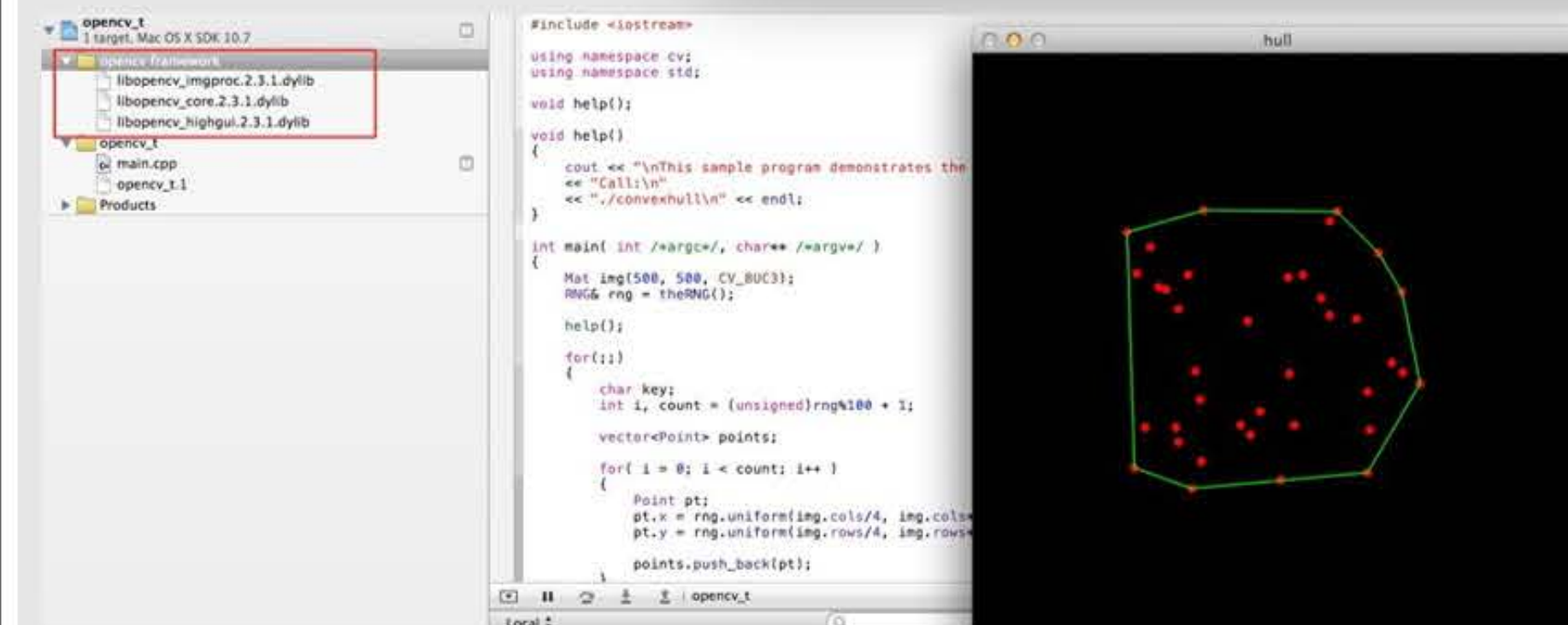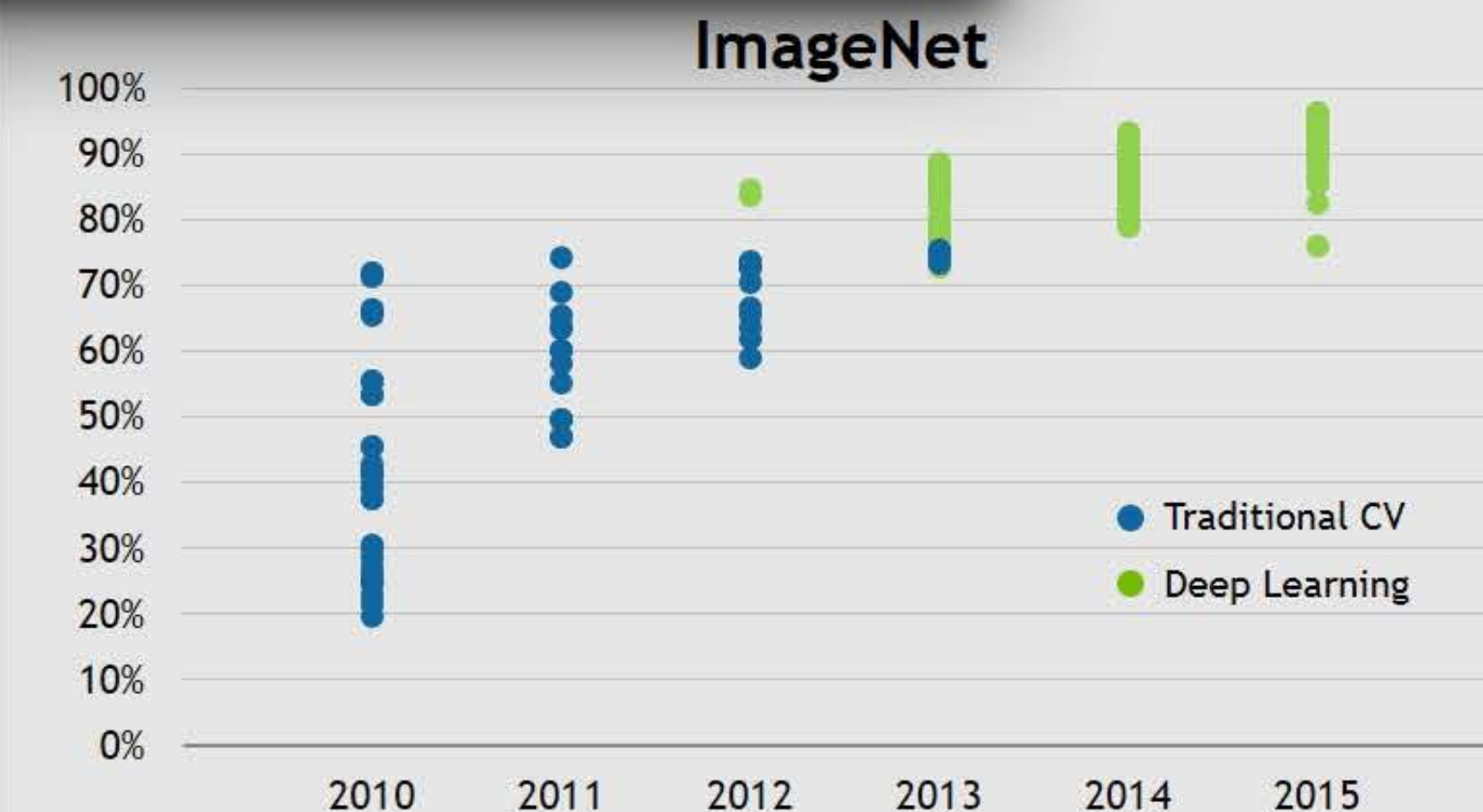
**drive.ai**

Automotive
computer vision

Audi

TOY

**" THE BIG BANG "**

Big Data
GPU
Algorithms

**RESEARCH**

**CORE TECHNOLOGY / FRAMEWORKS**

**AI-as-a-PLATFORM**

Preferred
Networks

Chainer

api.ai

Personal Assistants
conversational interface

MetaMind

eCommerce & Medical
recommendation engines

**Orbital Insight**

Geospatial
predictions from images

Baidu 百度

THE
HOME
DEPOT

Berkeley
UNIVERSITY OF CALIFORNIA

OpenAI

facebook.

torch

Université
de Montréal

theano

amazon
web services

DEEPMIND

Université
de Montréal

nervana

Tech
AI-as-a-service

Bloomberg

MASSACHUSETTS
GENERAL HOSPITAL
MGH

UBE

Carnegie
Mellon
University

Massachusetts
Institute of
Technology

Google

TensorFlow

Berkeley
UNIVERSITY OF CALIFORNIA

Caffe

IBM Watson

**BLUE RIVER**
TECHNOLOGY

Agriculture
crop-yield optimization

deep
genomics

Genomics
genetic interpretation

SADAKO

Waste Management
sorting robots

BMW

Mercedes-Benz

VOLV

NYU

UNIVERSITY OF
OXFORD

Microsoft

CNTK

UNIVERSITY OF
OXFORD

MatConvNet

Google

**clarifai**

Tech
visual recognition platform

charles SCHWAB

MERCK

Walmar

UNIVERSITY OF
TORONTO

NVIDIA.

cuDNN

Microsoft Azure

SocialEyes®

Medical
diabetic retinopathy

CISCO

Pinterest

YAHO

**1,000+ AI START-UPS**
**$5B IN FUNDING**

Source: Venture Scanner

HOW ARE YOU

Education
teaching robots

ebay

Schlumberger

Yand

FANUC
ROBOTICS

yel

# $500B OPPORTUNITY OVER 10 YRS

**Deep Learning Total Revenue by Segment**

Legend:
- Software
- Services
- Hardware

Y-axis: ($ Millions) — $-, $20,000, $40,000, $60,000, $80,000, $100,000, $120,000
X-axis: 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024

**Deep Learning Software Revenue by Industry**

Pie chart segments: Ad Service Technology, Investment, Media, Oil and Gas, Manufacturing, Retail, Other

Cognitive opens new opportunities on top of traditional IT

Opportunity for decision-making support 2025 — ~$2T

Decision Support

Traditional global IT spend 2016 — ~$1.2T

Productivity

Data center systems, Client relationship management, Infrastructure, Process automation, Enterprise resource planning

© 2016 IBM Corporation    Note 1    6

**IBM: "Cognitive business represents a $2T opportunity"**

# NVIDIA GPU FOR HYPERSCALE

## TESLA M40 + TESLA M4

**10X Speed up | 20 images/s/W**

Alibaba/Aliyun   Amazon   Baidu   eBay   Facebook   Flickr   Google   iFLYTEK   iQIYI   JD.com

Orange   Periscope   Pinterest   Qihoo 360   Shazam   Skype   Sogou   Twitter   Yahoo Supermarket   Yandex   Yelp

**Cloud Services Powered by AI**

Facebook AI Research

Soumith Chintala
AI Research Engineer, Facebook

Figure 8: A "turn" vector was created from four averaged samples of faces looking left vs looking right. By adding interpolations along this axis to random samples we were able to reliably transform their pose.

"Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks"

— Soumith Chintala, Facebook AI Research
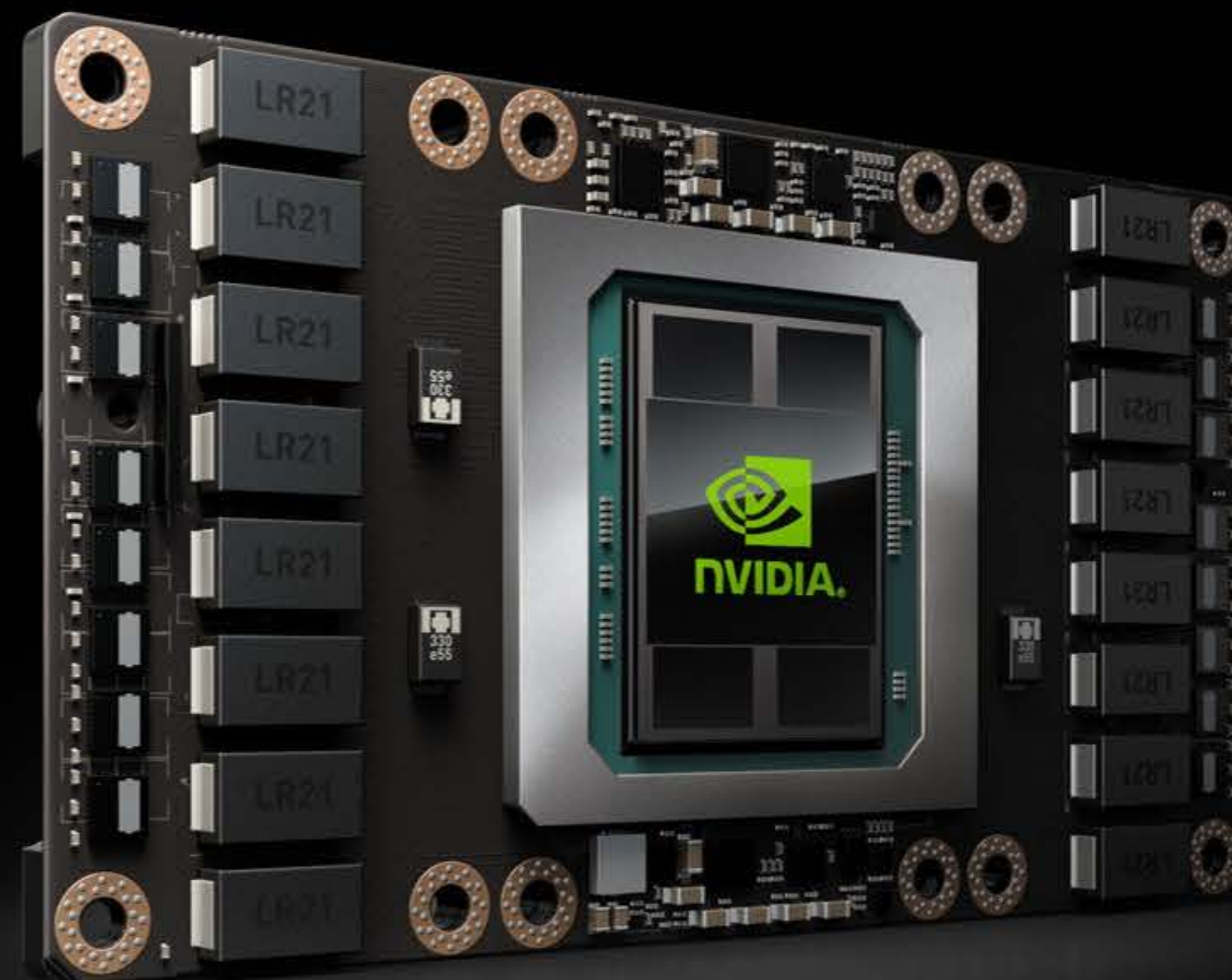Alec Radford & Luke Metz indico Research

Facebook AI Research

UNSUPERVISED LEARNING

# TESLA P100

## THE MOST ADVANCED
## HYPERSCALE DATACENTER GPU EVER BUILT

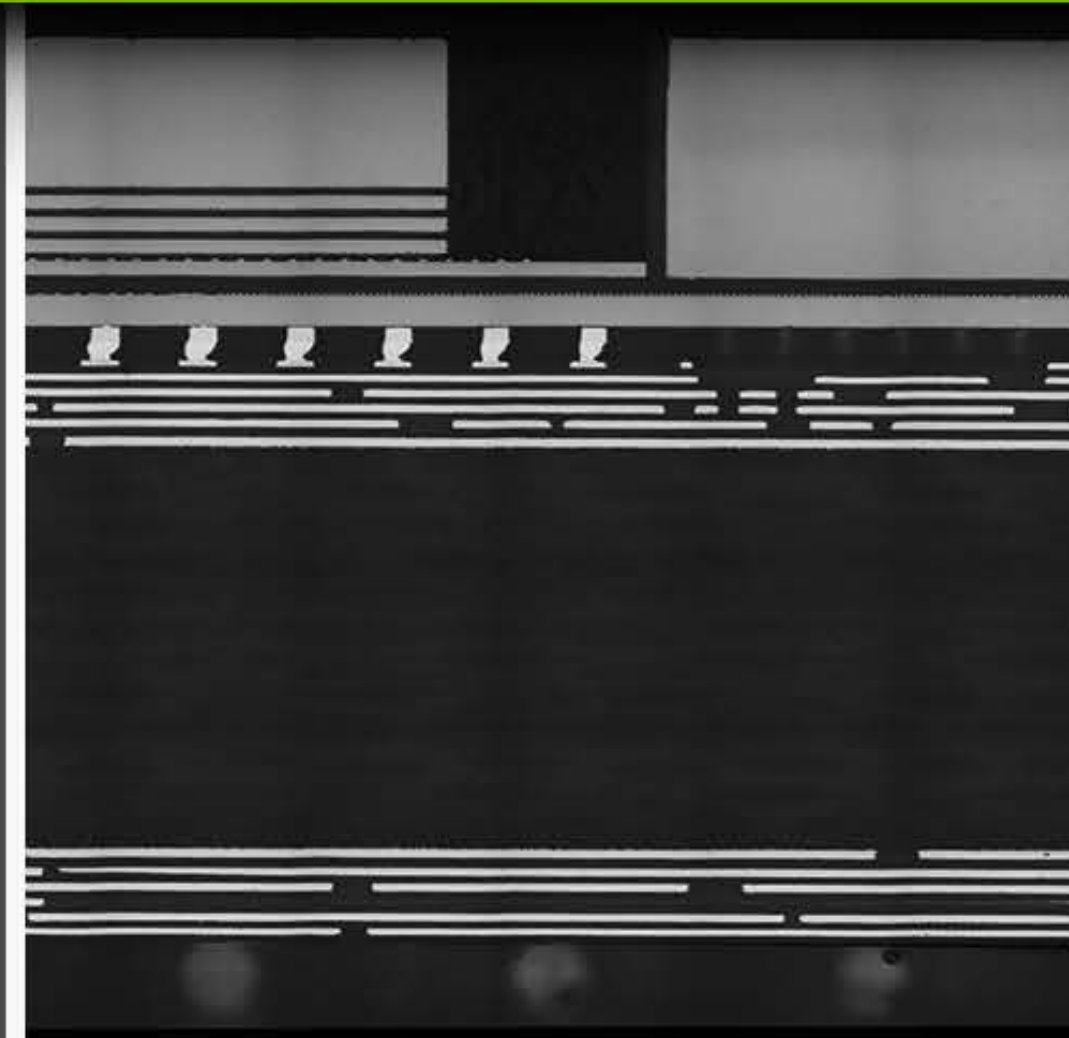150B XTORS | 5.3TF FP64 | 10.6TF FP32 | 21.2TF FP16 | 14MB SM RF | 4MB L2 Cache
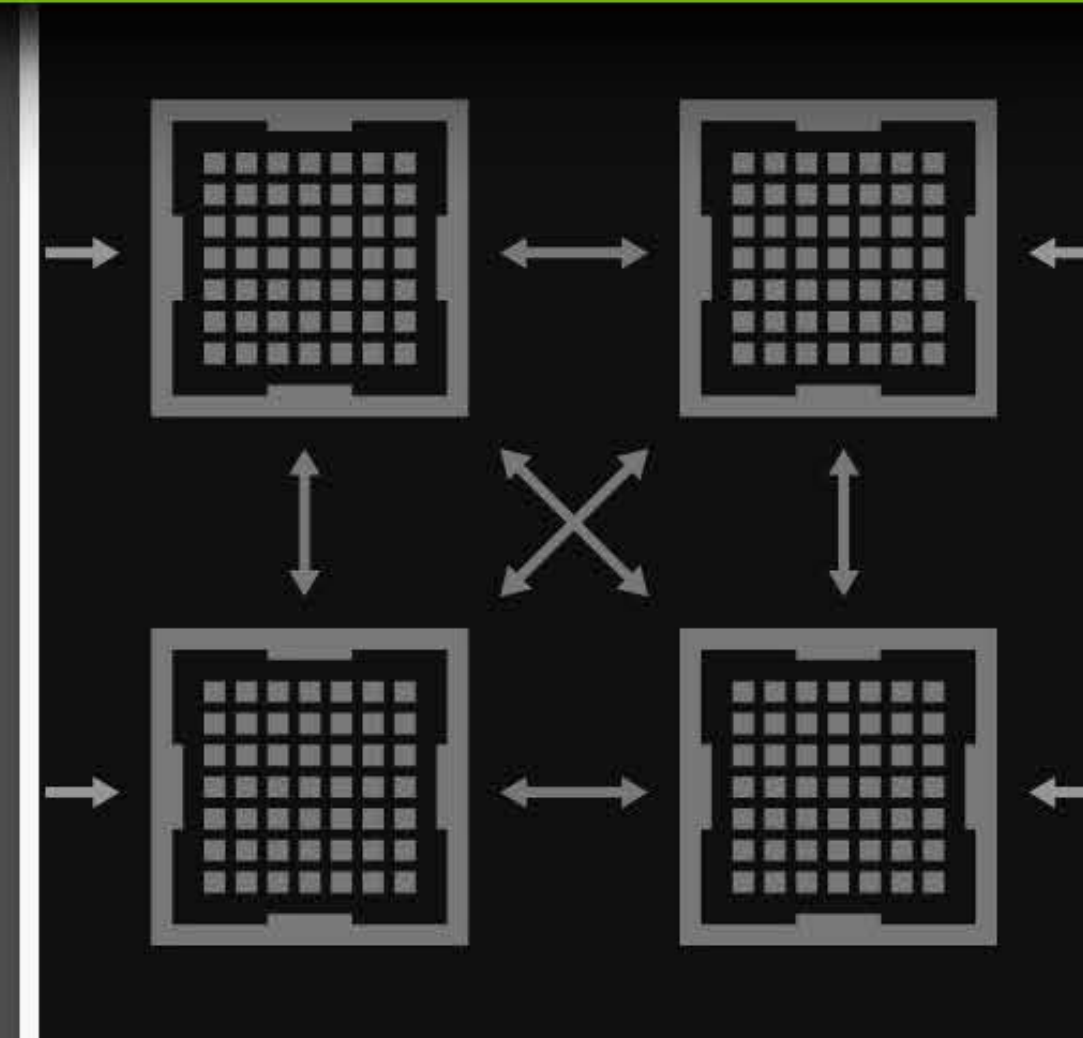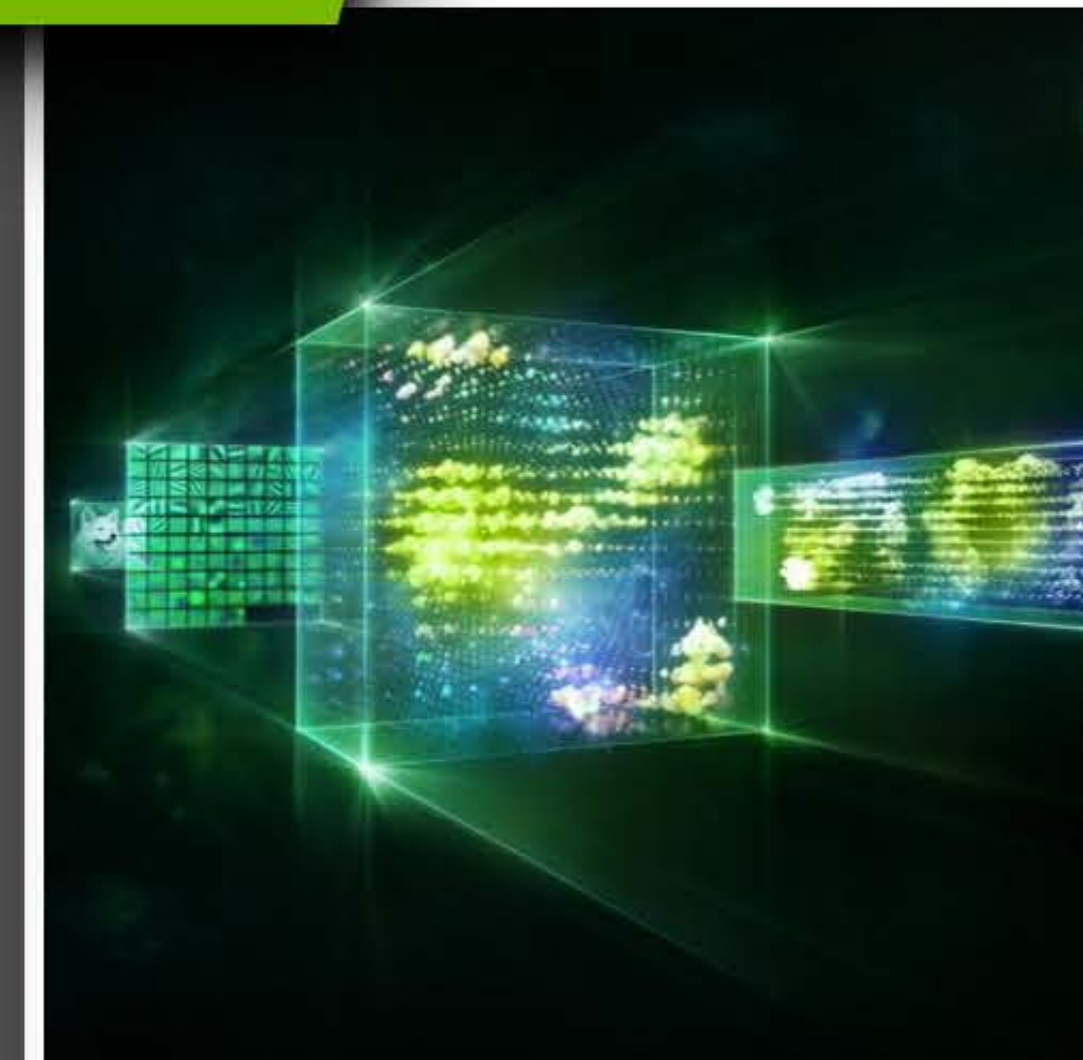
# "FIVE MIRACLES"

Pascal Architecture
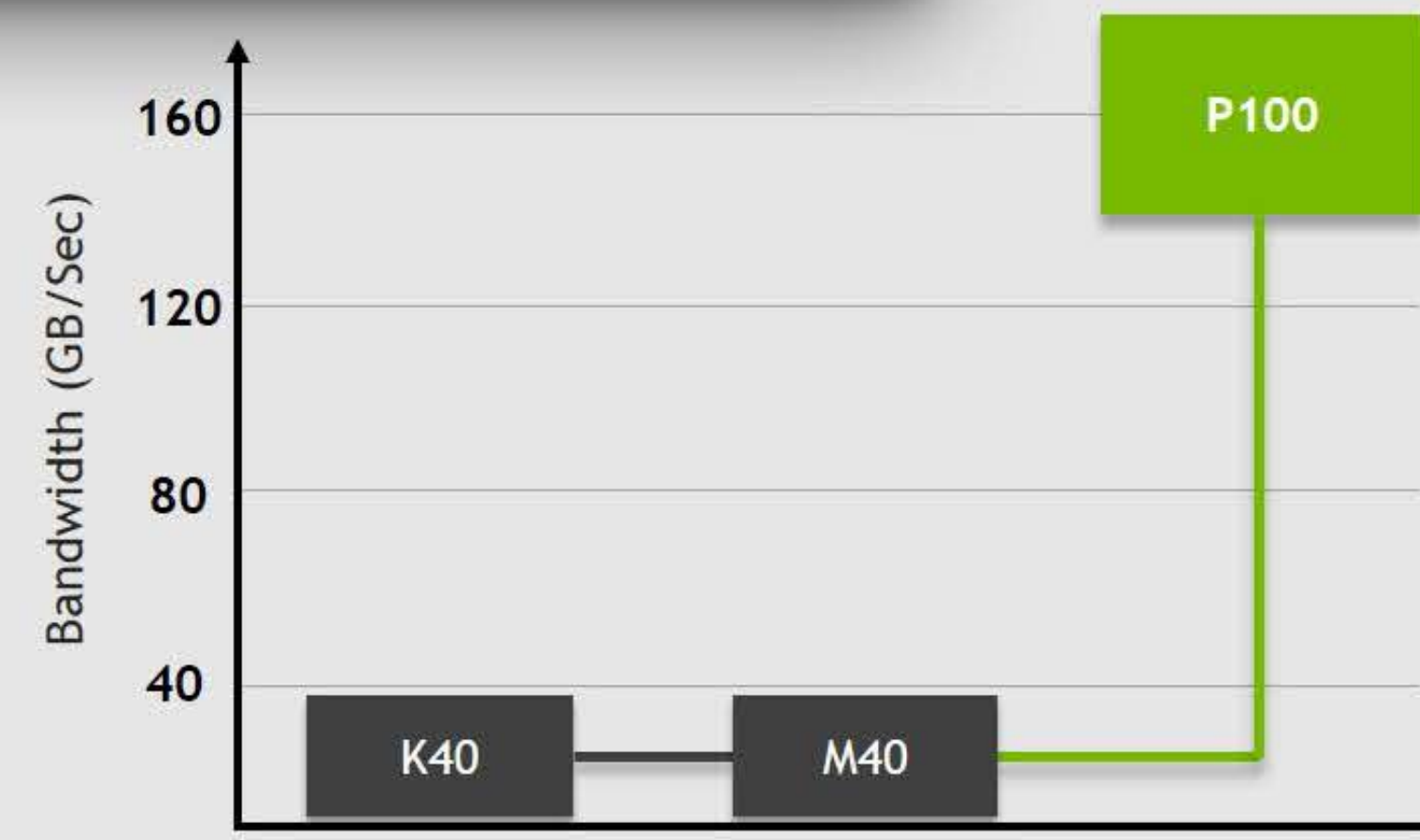
16nm FinFET

CoWoS with HBM2

NVLink

New AI Algorithms

# GIANT LEAPS IN EVERYTHING

**3x Compute**

Teraflops (FP32/FP16)

- P100 (FP16)
- P100 (FP32)
- K40
- M40

**3x GPU Mem BW**

Bandwidth

- P100
- K40
- M40

**5x GPU-GPU BW**

Bandwidth (GB/Sec)

- P100
- K40
- M40

"NVIDIA GPU is accelerating progress in AI. As neural nets become larger and larger, we not only need faster GPUs with larger and faster memory, but also much faster GPU-to-GPU communication, as well as hardware that can take advantage of reduced-precision arithmetic. This is precisely what Pascal delivers."

Yann LeCun, Director of AI Research, Facebook

"AI computers are like space rockets: The bigger the better. Pascal's throughput and interconnect will make the biggest rocket we've seen yet."

Andrew Ng, Chief Scientist, Baidu

"This is a new era of computing. New approaches to the underlying technologies will be required for AI and cognitive. The combination of NVIDIA Pascal GPUs and IBM POWER accelerates Watson's learning of new skills. Together, IBM and NVIDIA will advance the artificial intelligence industry."

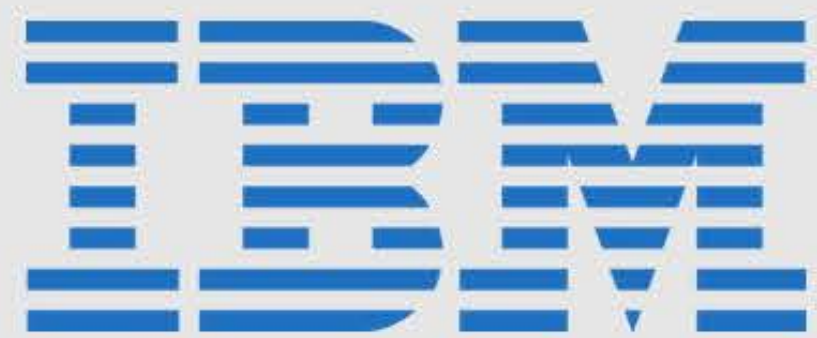Dr. John Kelly III, SVP, Cognitive Solutions & IBM Research

"Microsoft is developing super deep neural networks that are more than 1000 layers. NVIDIA Tesla P100's impressive horsepower will enable Microsoft's CNTK to accelerate AI breakthroughs."

Xuedong Huang, Chief Speech Scientist, Microsoft Research

TESLA P100 SERVERS — COMING IN Q1'17

# 5 THINGS

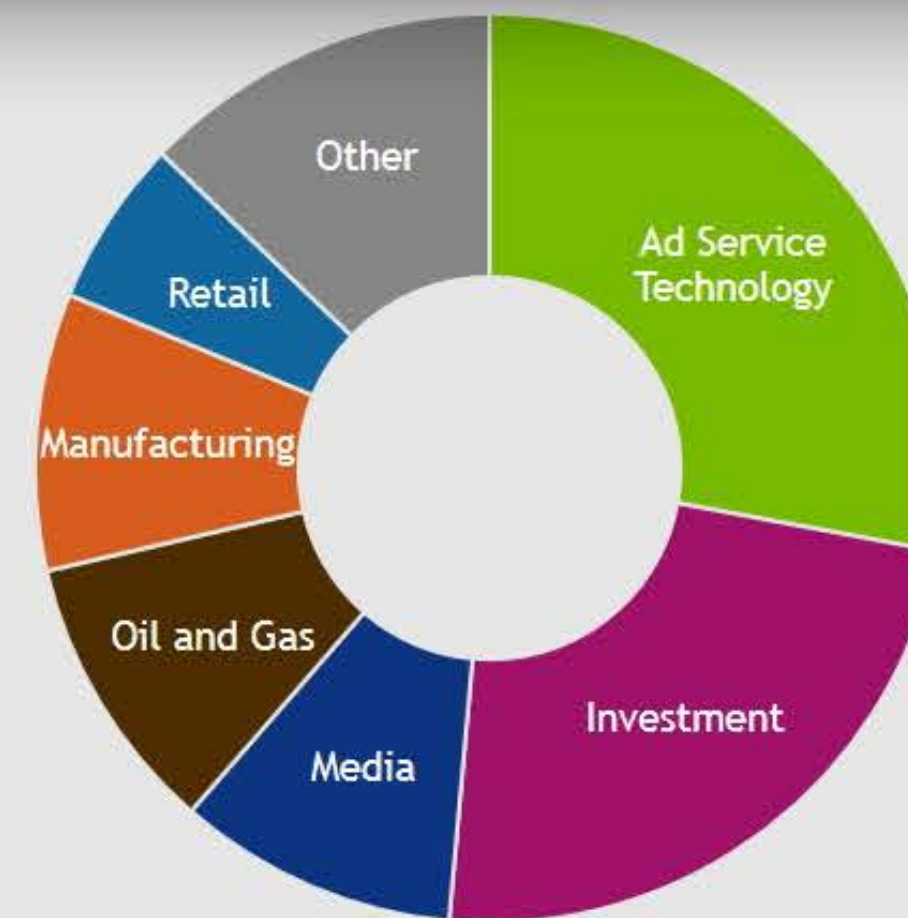NVIDIA SDK

TESLA P100

A Deep Learning Box

IRAY VR

A Deep Learning Car

# GPU-ACCELERATED DL FOR EVERY MARKET



**Deep Learning in the Cloud**

**IBM: "Cognitive business represents a $2T opportunity"**

Pie chart segments: Ad Service Technology, Investment, Media, Oil and Gas, Manufacturing, Retail, Other

**Deep Learning for Enterprise**

SOURCE: "Deep Learning for Enterprise Applications," 4Q 2015, Tractica

# NVIDIA DGX-1

## WORLD'S FIRST DEEP LEARNING SUPERCOMPUTER

Engineered for deep learning  |  170TF FP16  |  8x Tesla P100  |  NVLink hybrid cube mesh  |  Accelerates major AI frameworks

7 TB SSD

8x Tesla P100 16GB

2x Xeon

NVIDIA DGX-1
WORLD'S FIRST DEEP LEARNING
SUPERCOMPUTER | 170 TFLOPS

3U - 3200W

NVLink Hybrid Cube Mesh

Quad IB 100Gbps, Dual 10GbE

|  | DUAL XEON | DGX-1 |
|---|---|---|
| FLOPS (CPU + GPU) | 3 TF | 170 TF |
| AGGREGATE NODE BW | 76 GB/s | 768 GB/s |
| ALEXNET TRAIN TIME | 150 HOURS | 2 HOURS |
| TRAIN IN 2 HOURS | >250 NODES* | 1 NODE |

## "250 SERVERS IN-A-BOX"

Bryan Catanzaro
Senior Researcher, Baidu

"Time series output"

Time series input

**Recurrent Neural Nets**

Model Parallel
GPU0
GPU1

Data Parallel

**Model + Data Parallelism**

**Persistent RNNs: Peak FLOPs at batch of 8**

keep in registers

weights

repeat ~300 times

**Add Model Parallelism over NVLINK**

GPU0
GPU1
GPU2
GPU3

repeat ~300 times

**Compose with Data Parallelism**

Data Parallel

Strong scale to 32X more processors

Rajat Monga
TensorFlow Technical Lead & Manager, Google

# NVIDIA DGX-1

## WORLD'S FIRST
## DEEP LEARNING SUPERCOMPUTER

170TF | "250 servers in-a-box" | nvidia.com/dgx1

**$129,000**

Berkeley UNIVERSITY OF CALIFORNIA

Carnegie Mellon University

香港中文大學
The Chinese University of Hong Kong

MIT Massachusetts Institute of Technology

NYU

Stanford University

Université de Montréal

UNIVERSITY OF OXFORD

USI/SUPSI IDSIA

UNIVERSITY OF TORONTO

Frameworks for Multi-GPU Pascal

Large-scale Deep Learning

Reinforcement Learning

Unsupervised and Transfer Learning

Natural Language Understanding

Autonomous Driving

Medical Applications

# PIONEERS IN AI RESEARCH

# TESLA FAMILY

| M40 + M4 | K80 | | |
|---|---|---|---|
| Hyperscale HPC | Multi-App HPC | Strong-Scale HPC | Researchers / Early Adopters |

# 5 THINGS

**NVIDIA SDK**

**TESLA P100**

**NVIDIA DGX-1**

**IRAY VR**

**A Deep Learning Car**

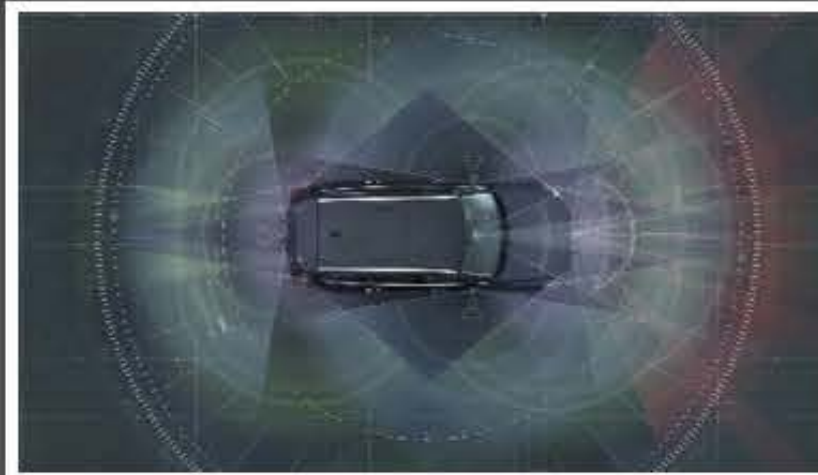# AN AMAZING YEAR FOR SELF-DRIVING CARS



Uber Enters the Race

Toyota Invests $1B in AI Lab

Volvo Drive Me on Public Roads in 2017

Tesla Model 3: 300K pre-orders

NHTSA: Computer Counts as Driver

Audi, BMW, Daimler Buy HERE

Tesla Model S Auto-pilot

Baidu Enters the Race

Honda, Nissan, Toyota Team Up

GM Buys Cruise

# SELF-DRIVING LOOPS

**MAP**

**LOCALIZE**

**SEE**

**DRIVE**

# NVIDIA DRIVE PX
# AI CAR COMPUTER



Caffe
CNTK
KALDI
TensorFlow
theano
torch

**Training on DGX-1**

NVIDIA DGX-1

NVIDIA DRIVE PX

MAPPING

LOCALIZATION

DRIVENET

DAVENET

**Driving with DriveWorks**

World's first DL-powered car computing platform

One scalable architecture — from DNN training to cluster, infotainment, ADAS, autonomous driving, and mapping
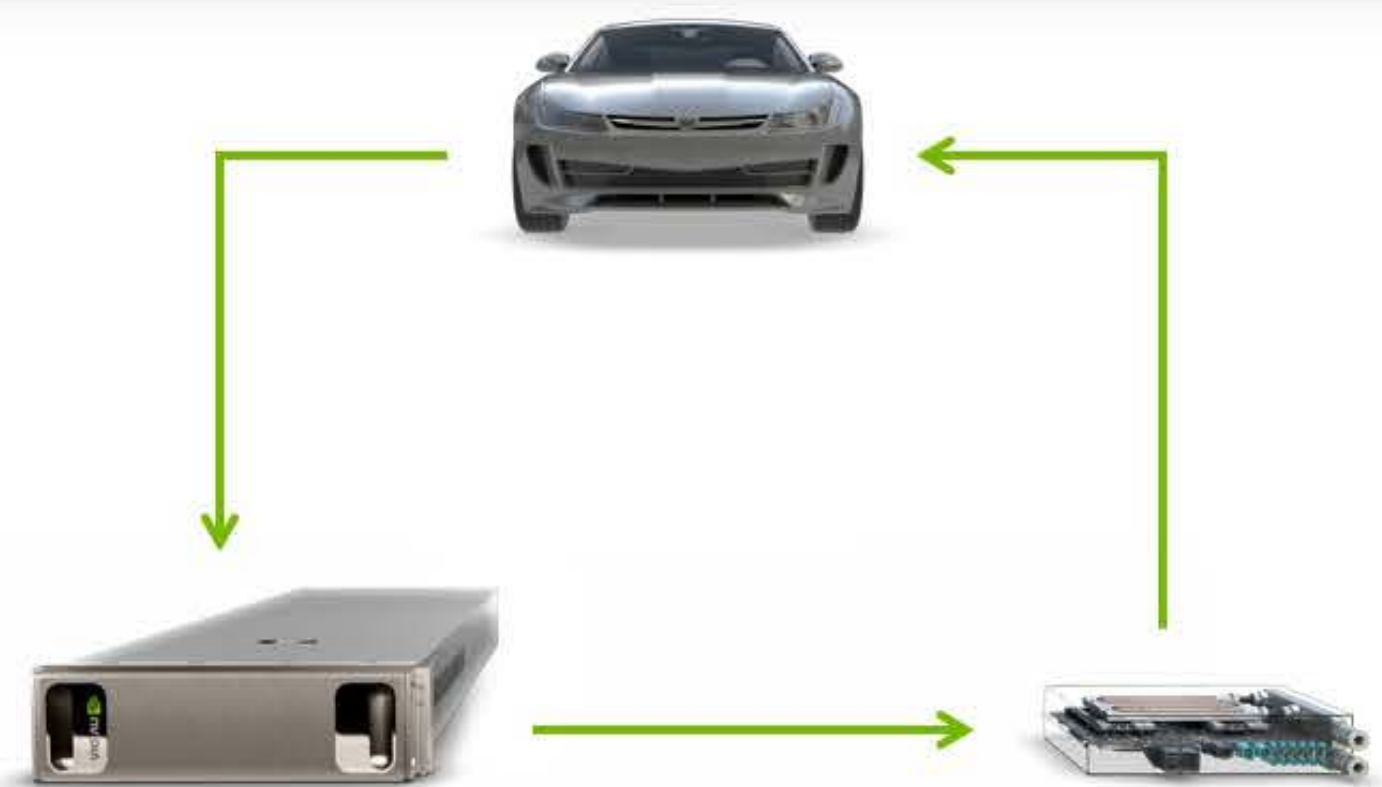
Open platform

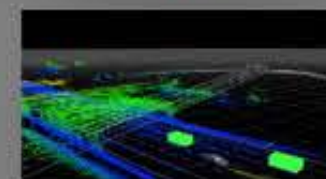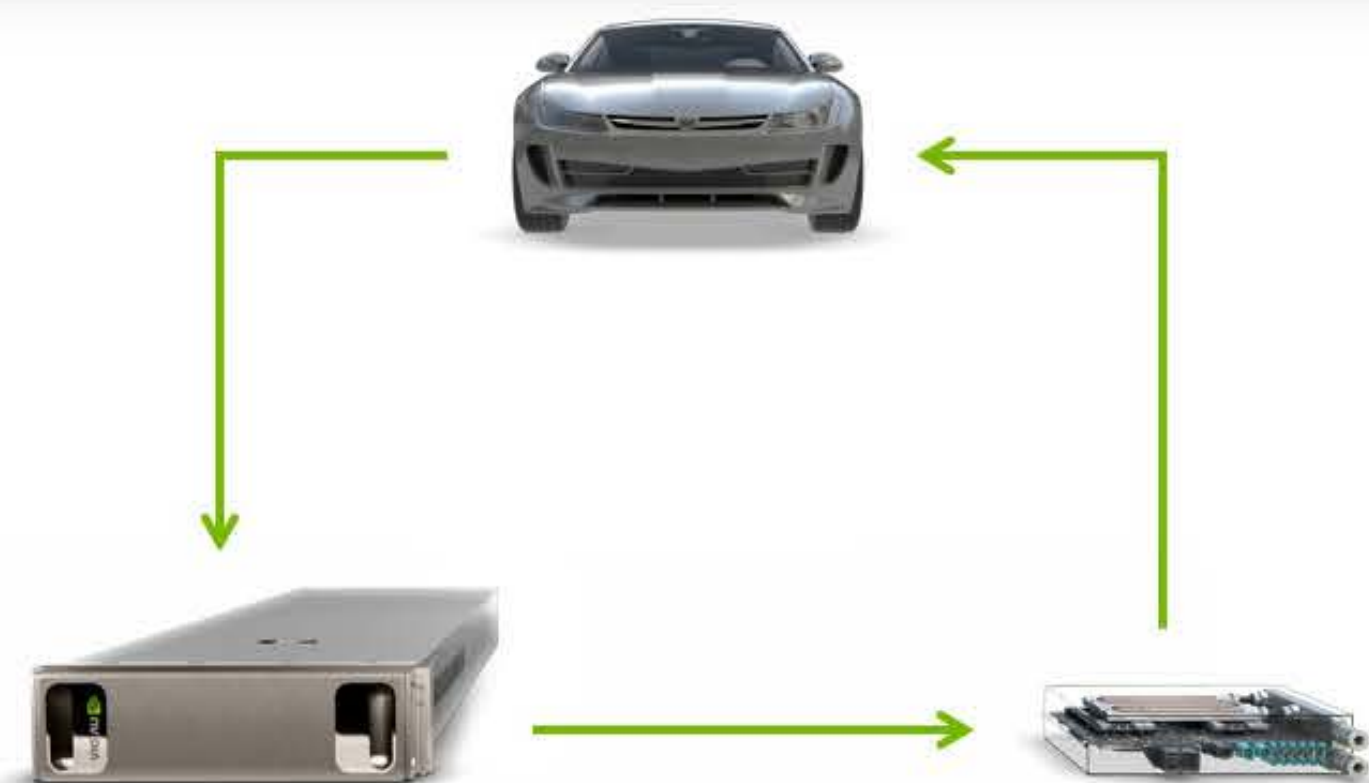# NVIDIA DRIVE PX PERCEPTION

Caffe
CNTK
KALDI
TensorFlow
theano
torch

**Training on DGX-1**

NVIDIA DGX-1          NVIDIA DRIVE PX

MAPPING
LOCALIZATION
DRIVENET
DAVENET

**Driving with DriveWorks**

## NVIDIA DRIVENET
### #1 accuracy score for KITTI car detection

| | Method | Hard | Moderate | Easy | Environment |
|---|---|---|---|---|---|
| 1 | NVDriveNet-H | 83.76 % | 89.81 % | 90.92 % | GPU @ 2.5 Ghz (Python + C/C++) |
| 2 | sensekitti | 79.99 % | 89.72 % | 91.42 % | GPU @ 2.5 Ghz (Python + C/C++) |
| 3 | SDP+RPN | 78.38 % | 88.85 % | 90.14 % | GPU @ 2.5 Ghz (Python + C/C++) |
| 4 | Mono3D | 78.96 % | 88.66 % | 92.33 % | GPU @ 2.5 Ghz (Matlab + C/C++) |
| 5 | 3DOP | 79.10 % | 88.64 % | 93.04 % | GPU @ 2.5 Ghz (Matlab + C/C++) |

NEW
END-TO-END HD MAPPING

Caffe
CNTK
KALDI
TensorFlow
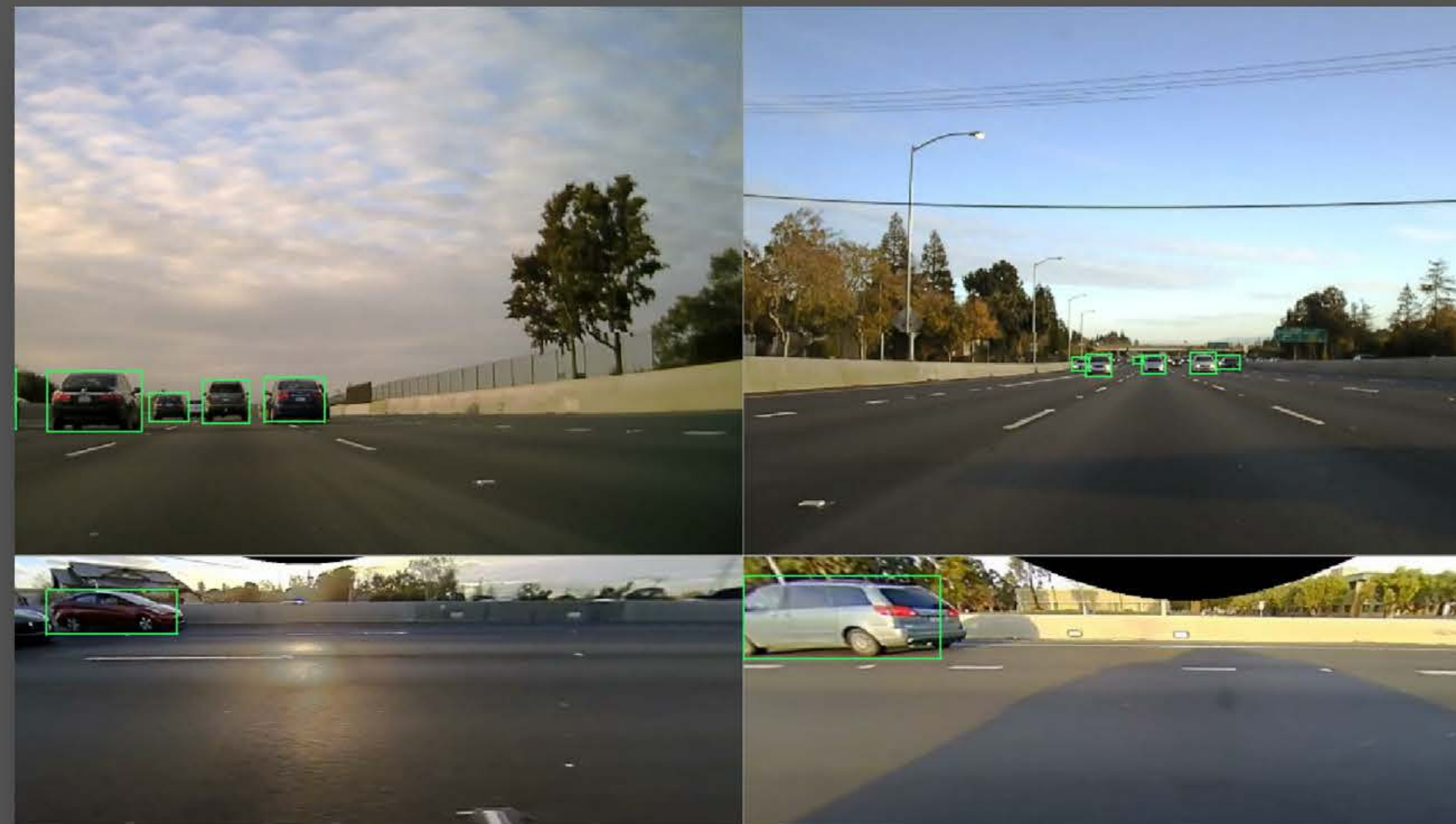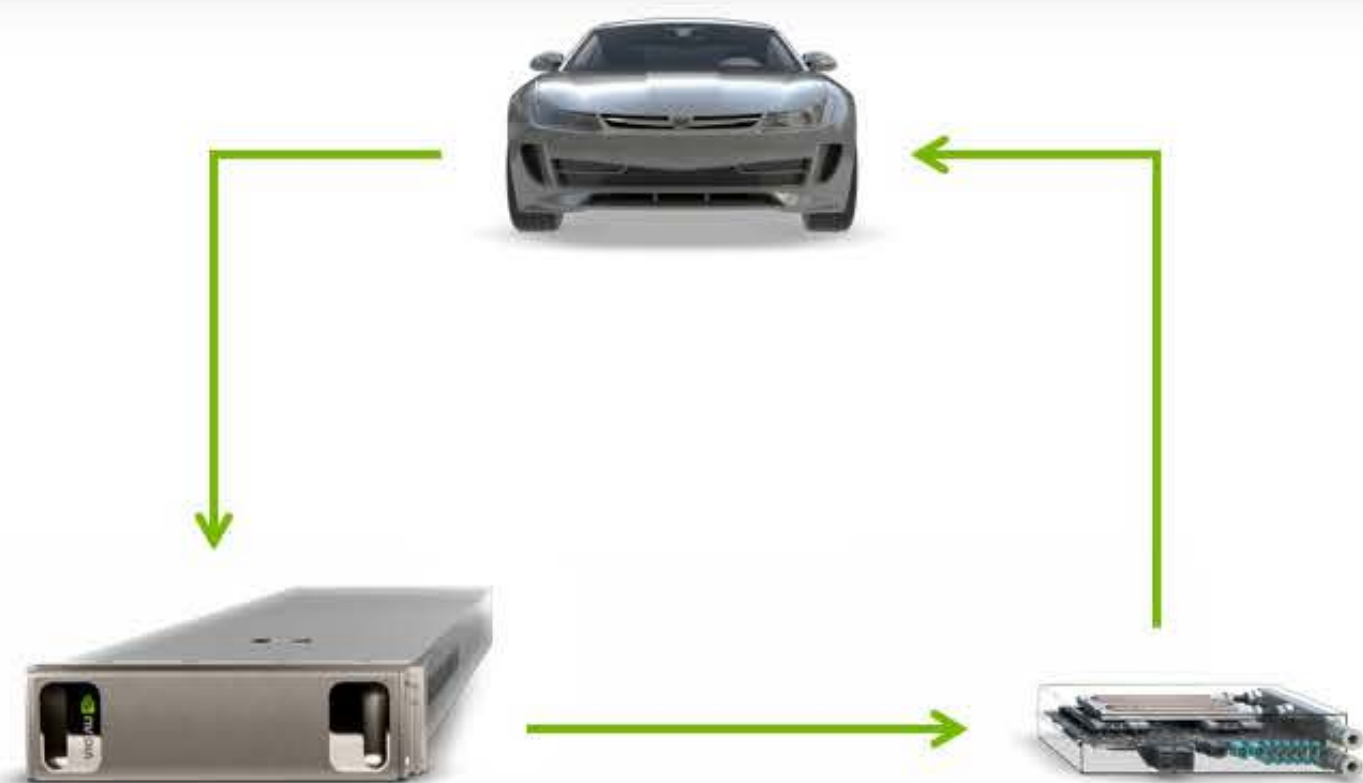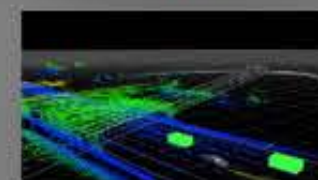theano
torch

Training on DGX-1

NVIDIA DGX-1

NVIDIA DRIVE PX
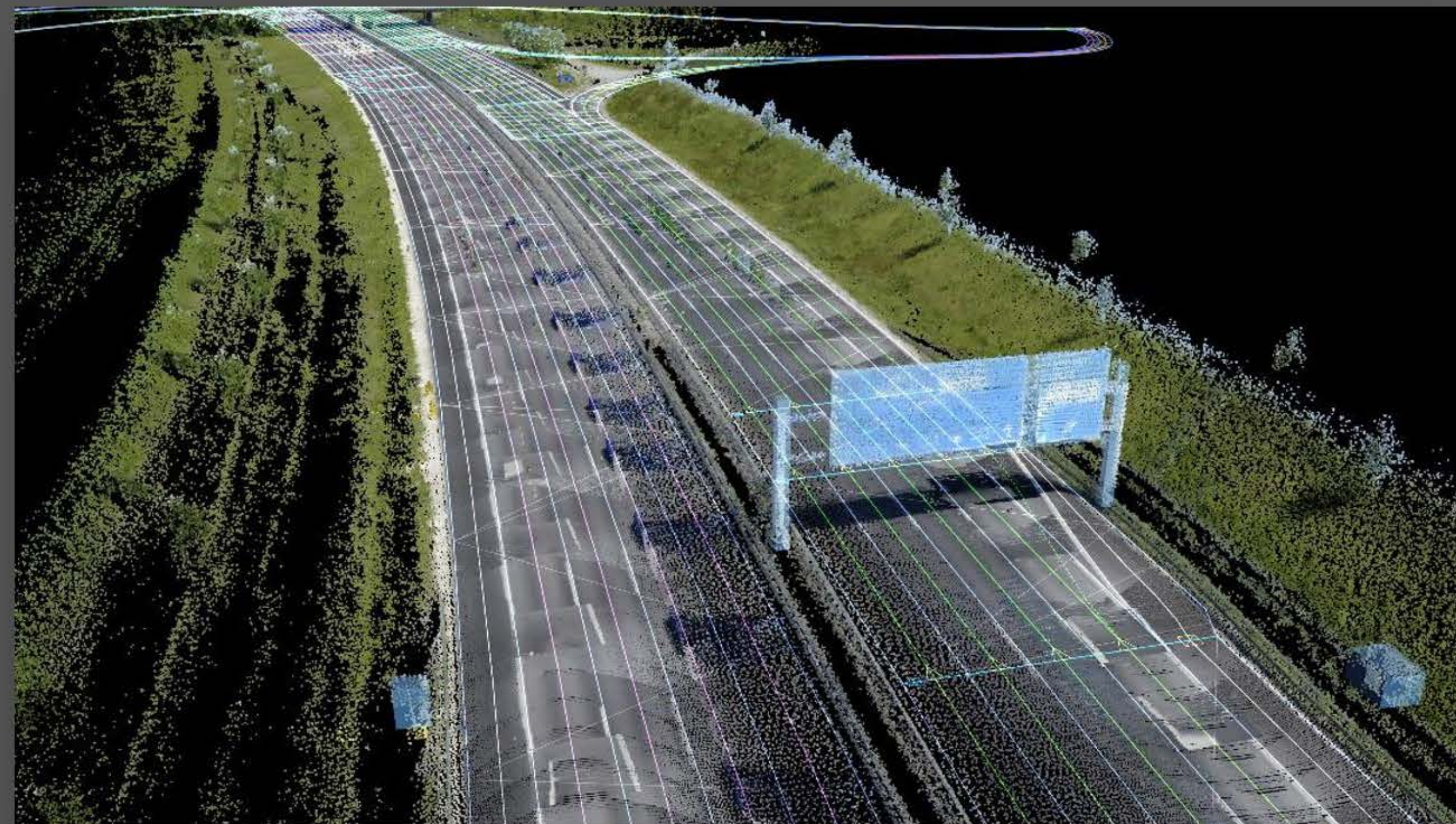
MAPPING

LOCALIZATION

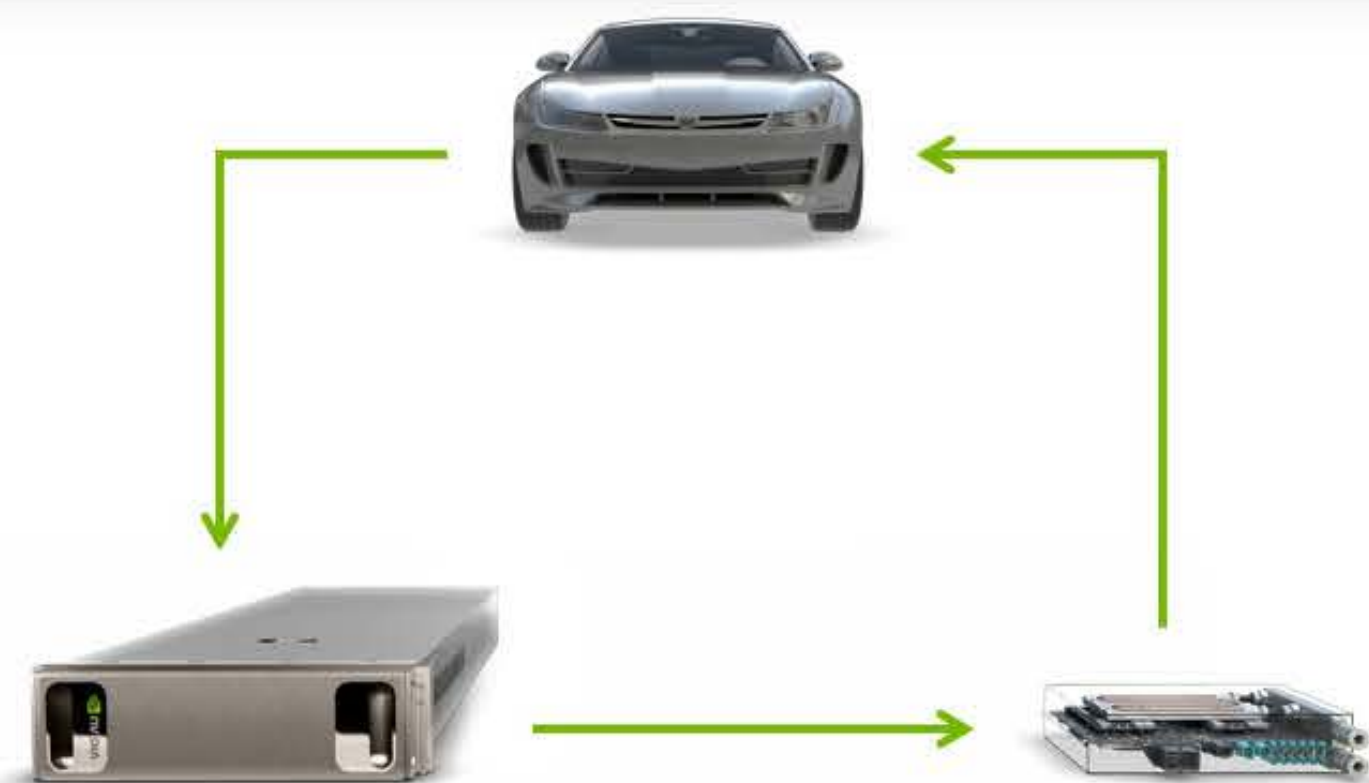DRIVENET

DAVENET

Driving with DriveWorks

BAIDU SELF-DRIVING
CAR COMPUTER

Baidu百度

NEW
END-TO-END HD MAPPING

Caffe
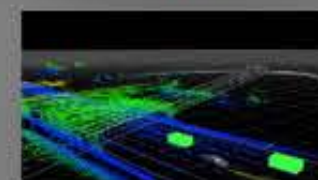CNTK
KALDI
TensorFlow
theano
torch

Training on DGX-1

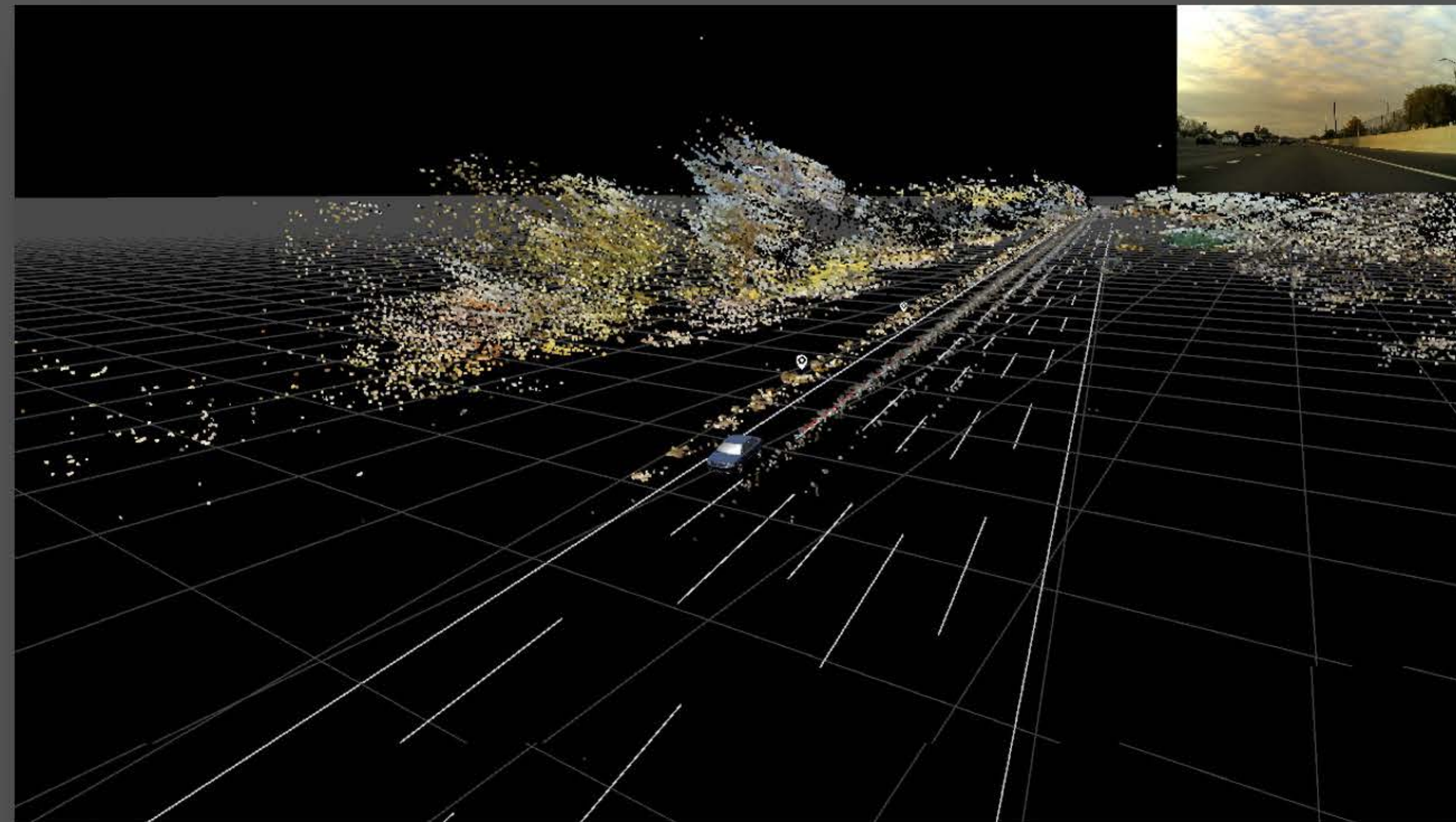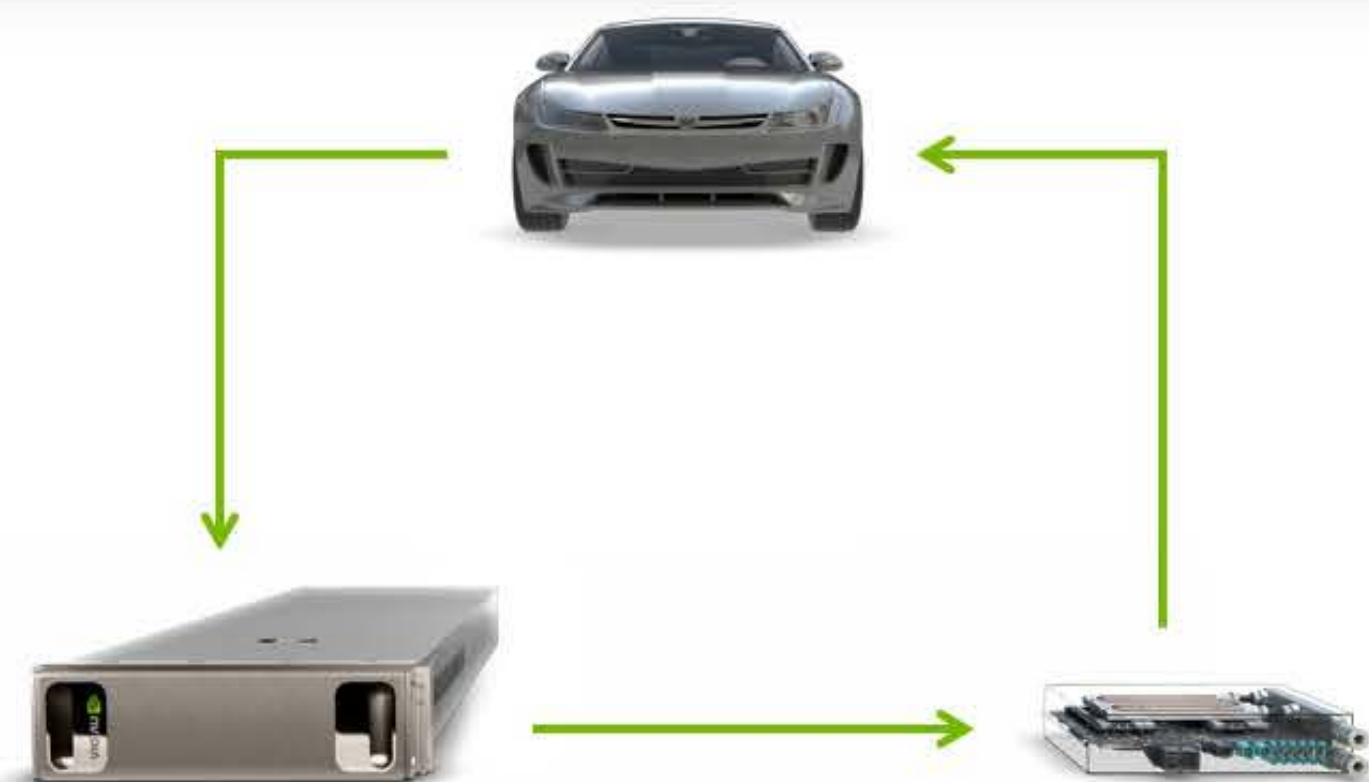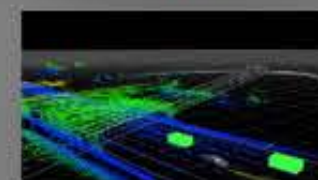NVIDIA DGX-1

NVIDIA DRIVE PX

MAPPING
LOCALIZATION
DRIVENET
DAVENET

Driving with DriveWorks

here

TomTom®

ZENRIN

# WORLD'S FIRST AUTONOMOUS RACE CAR

Designed by Daniel Simon

2,200 lbs

Blazing fast

# WORLD'S FIRST AUTONOMOUS CAR RACE

10 teams, 20 identical cars

DRIVE PX 2: The "brain" of every car

2016/17 Formula E season

**ROBORACE**

# VR, AI, SELF-DRIVING CARS


NVIDIA SDK


TESLA P100


NVIDIA DGX-1


IRAY VR


HD MAPPING, AI DRIVING

GPU TECHNOLOGY CONFERENCE