# GPU in every workload

Mason Wu, Solution Architect
April. 2018

**NVIDIA**

# AGENDA

- NVIDIA Introduction

- Tesla for HPC and AI

- GRID vGPU for Virtualization environment

- Quadro for Professional visualization

AI COMPUTING

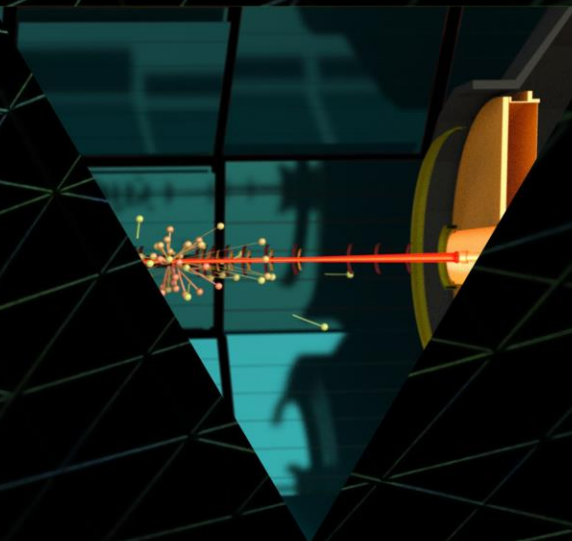GPU COMPUTING

PC GRAPHICS

**NVIDIA**

# A LEARNING MACHINE

GPU sparked the growth of PC gaming, redefined modern computer graphics, and revolutionized parallel computing

GPU deep learning ignited modern AI — the next era of computing

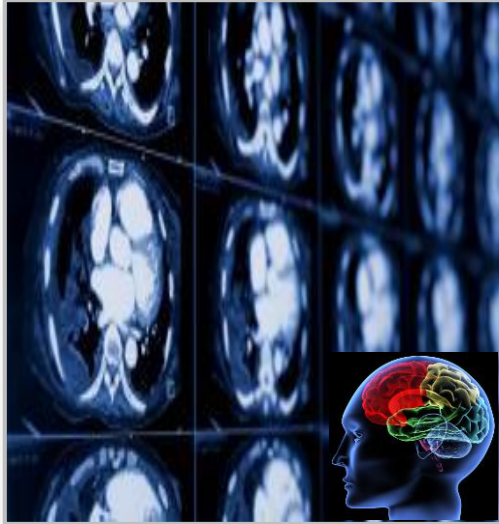1996                    2006                    2016

# TESLA FOR HPC & AI

# DEEP LEARNING IS SWEEPING ACROSS INDUSTRIES

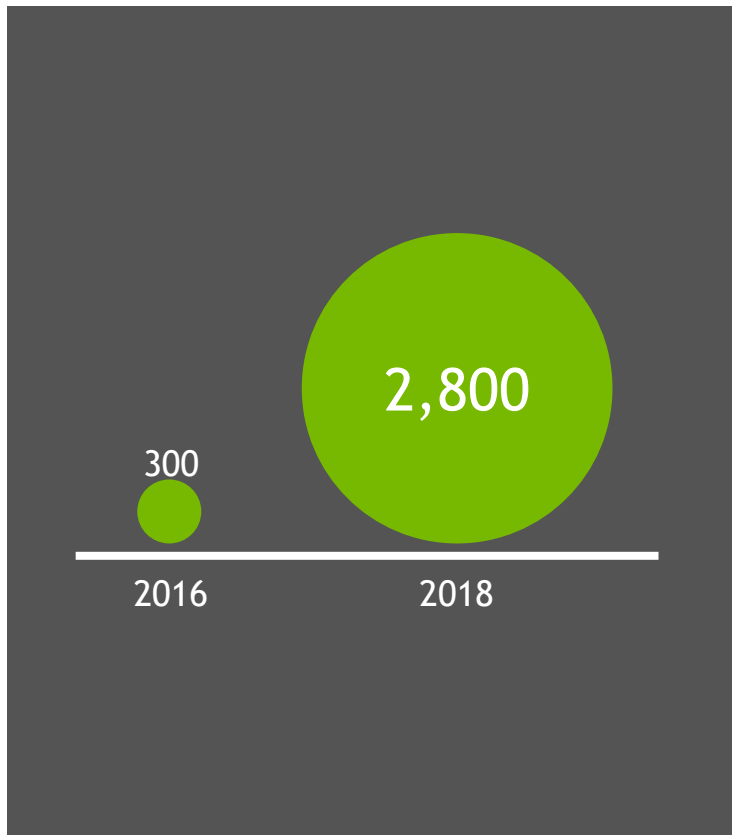**Internet Services**     **Medicine**     **Media & Entertainment**     **Security & Defense**     **Autonomous Machines**
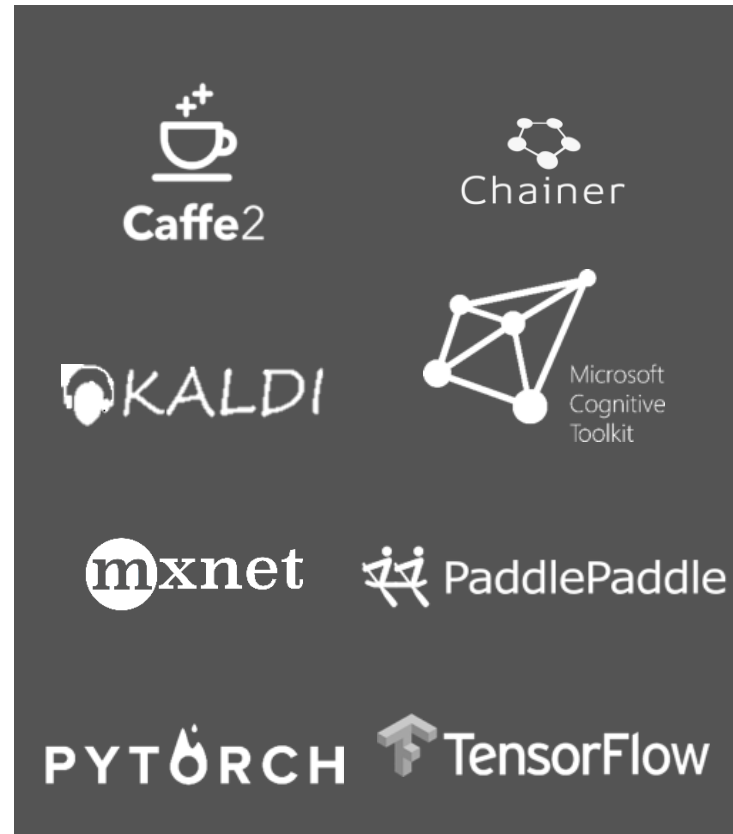
| | | | | |
|---|---|---|---|---|
| ➢ Image/Video classification | ➢ Cancer cell detection | ➢ Video captioning | ➢ Face recognition | ➢ Pedestrian detection |
| ➢ Speech recognition | ➢ Diabetic grading | ➢ Content based search | ➢ Video surveillance | ➢ Lane tracking |
| ➢ Natural language processing | ➢ Drug discovery | ➢ Real time translation | ➢ Cyber security | ➢ Recognize traffic signs |

# MOST ADOPTED PLATFORM FOR ACCELERATING AI



**300** · 2016
**2,800** · 2018

**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

Desktops

**AVAILABLE EVERYWHERE**

# TESLA STACK

## World's Leading Data Center Platform for Accelerating HPC and AI
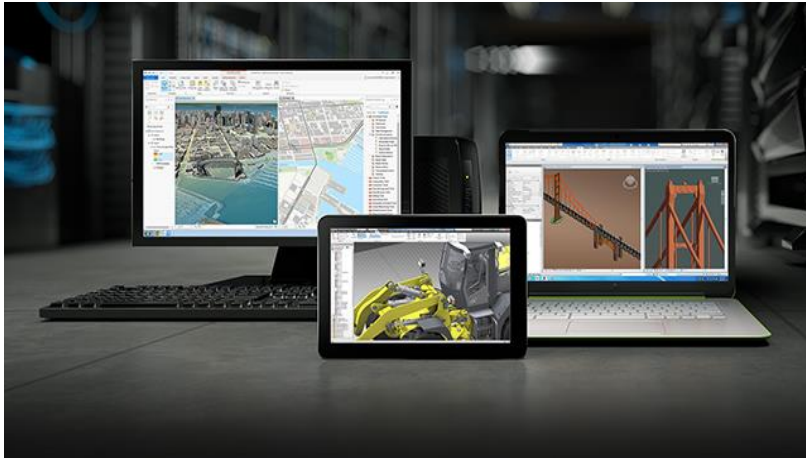
# GRID VCPU FOR VIRTUALIZATION ENVIRONMENT

# **GPU** VIRTUALIZATION FOR ANY WORKLOAD
## NVIDIA delivers GPU virtualization for both graphics and compute workloads
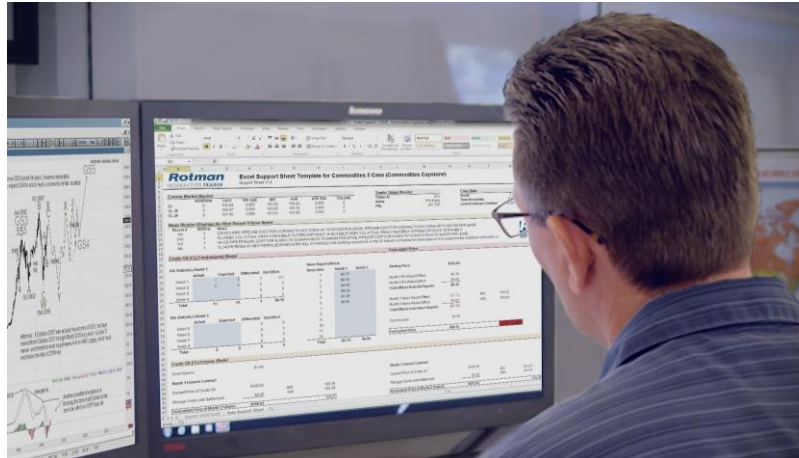
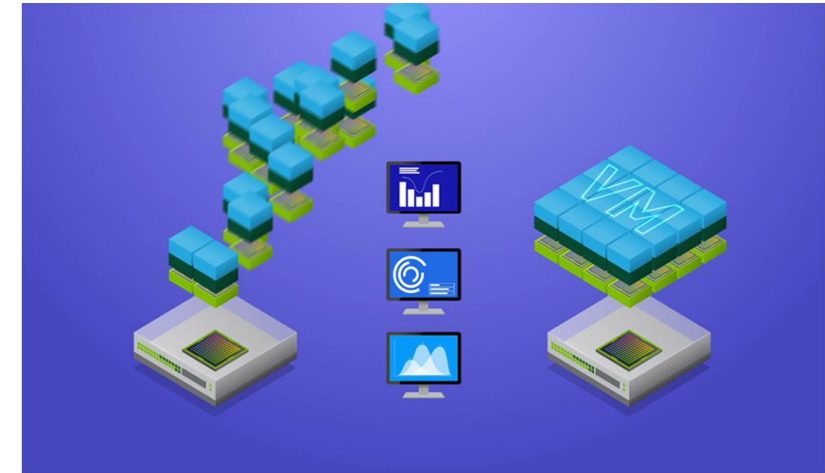# NVIDIA VIRTUAL GPU

## Most Powerful Virtual Workstation

Quadro vDWS on Tesla V100, the most advanced data center GPU

## Best Virtual PC for Knowledge Workers

Expanded GRID vPC delivers more value to more enterprise users

## Agile GPU Accelerated Data Center

Live Migration for the Mission Critical Data Center

# HOW IT WORKS

NVIDIA virtual GPU products deliver a GPU Experience to every Virtual Desktop



CPU Only VDI

Apps and VMs

Hypervisor

Server

With NVIDIA Virtual GPU

Office

Apps and VMs

NVIDIA Graphics Drivers

NVIDIA Virtual GPU

NVIDIA virtualization software

Hypervisor

NVIDIA Tesla GPU
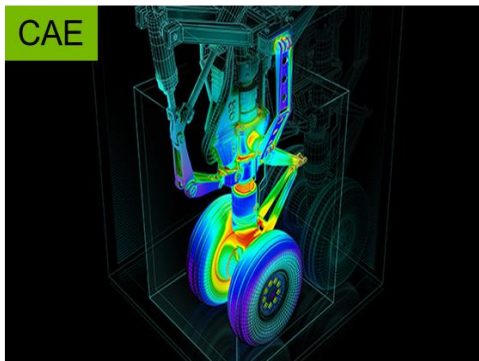
Server

QUADRO FOR PROFESSIONAL VISUALIZATION

# DEMAND FROM THE WORLD'S LARGEST INDUSTRIES

- Desire to work with full-fidelity models driving dramatic increase in data size
- Demand for photorealism throughout the design process for better insights
- Need for immersive environments for faster, higher-quality decision making
- Robust platform for rapid AI research, development & deployment
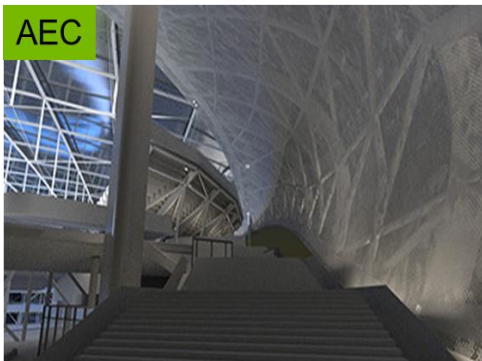
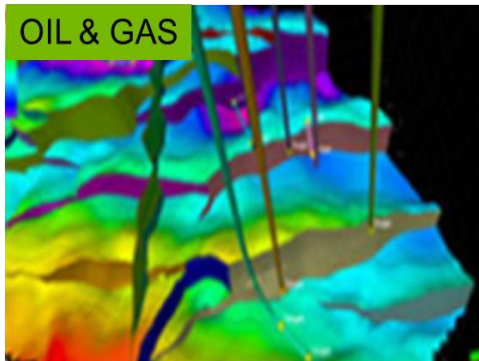MANUFACTURING

CAE

MEDIA & ENTERTAINMENT

AUTOMOTIVE

AEC

OIL & GAS

SCI-VIZ

HEALTHCARE

# NVIDIA QUADRO GV100

- Delivers world's fastest real-time ray tracing on workstations for millions of artists & designers
- Opens new opportunities with accelerated AI for training and inference on workstations
- Accelerates CAE workloads for simulation analysts to bring optimal products to market faster
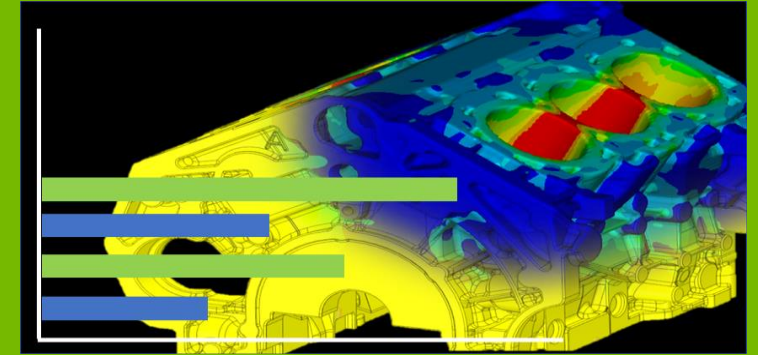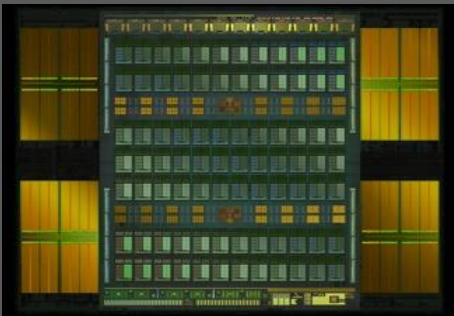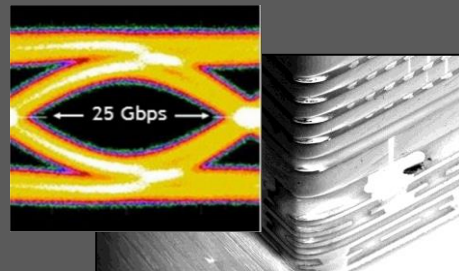
## RAY TRACING

## AI

## SIMULATION

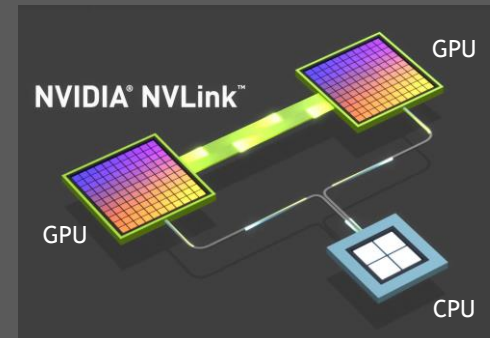# TECHNOLOGIES POWERING THE AI WORKSTATION

| Volta GPU | HBM2 | NVLink | Tensor Cores |
|:---:|:---:|:---:|:---:|
| Efficient, Powerful | Largest GPU Memory | Scalable Performance | Accelerated AI |