

# NVIDIA GRAPHICS HPC AI

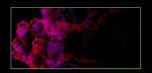








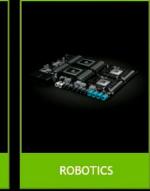




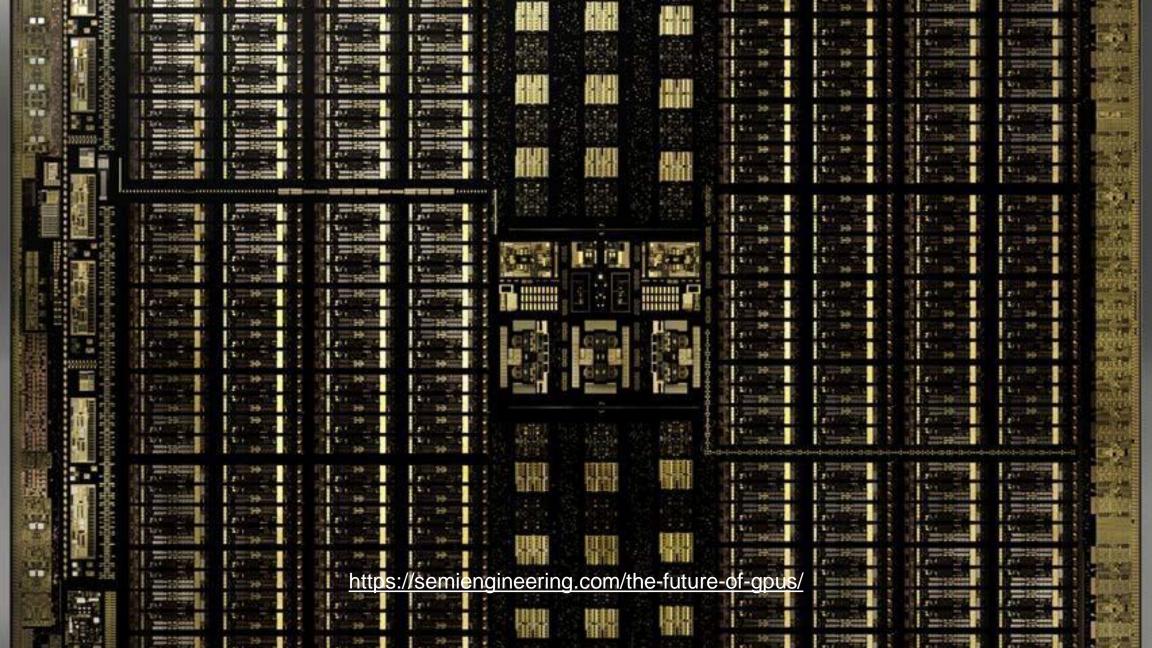














## www.FrontierDevelopmentLab.org





















## SELECTING THE RIGHT GPU SOLUTION

## Unparalleled Data Science Performance and Productivity

#### **ML Enthusiast**



TITAN RTX

**48GB** 

PC solution, easy to acquire, deploy and get started experimenting

#### Machine Learning Developer Data Science Workstations



Quadro Workstation

Enterprise workstation for experienced data scientists



**DGX Station** 

Enterprise ML workgroups, largest memory on a workstation

Data Center Machine Learning
Shared infrastructure for Data Science Teams



DGX-1 / OEM

Enterprise server, proven 8-way configuration, modular approach for scale, multi-node training



DGX-2 / OEM

Largest compute and memory capacity in single node, fastest training solution

	The second second second	
GPII	Memory	
UI U	memory	

Benefit

GPU Fabric 2-way

64GB

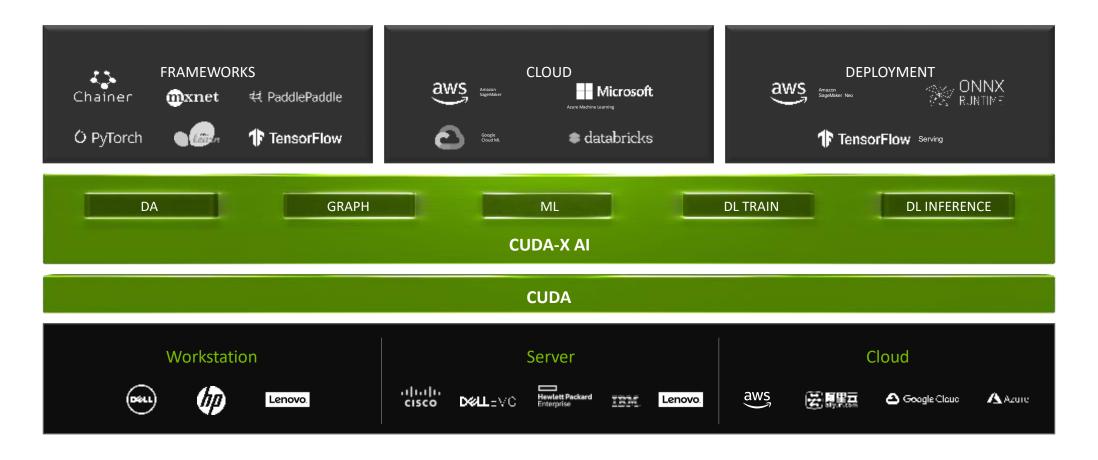
2-way NVLINK 128GB

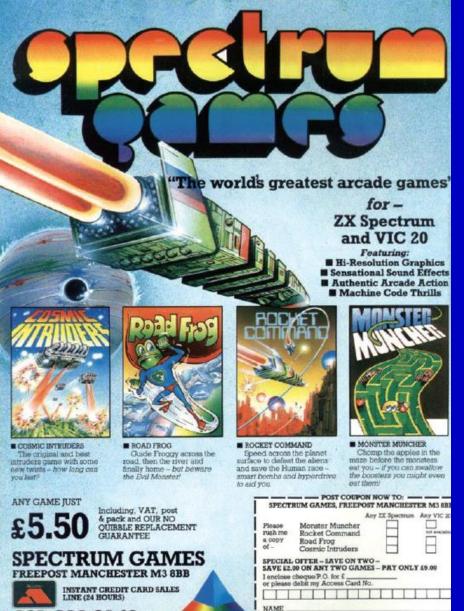
4-way NVLINK 256GB

8-way NVLINK 512GB

16-way NVSWITCH

## **NVIDIA CUDA-X AI ECOSYSTEM**





ADDRESS

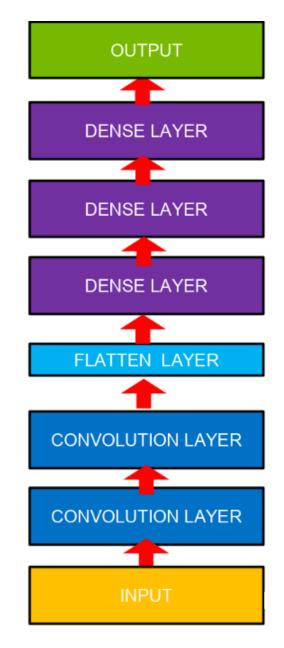
FREEPOST NO STAMP NECESSARY



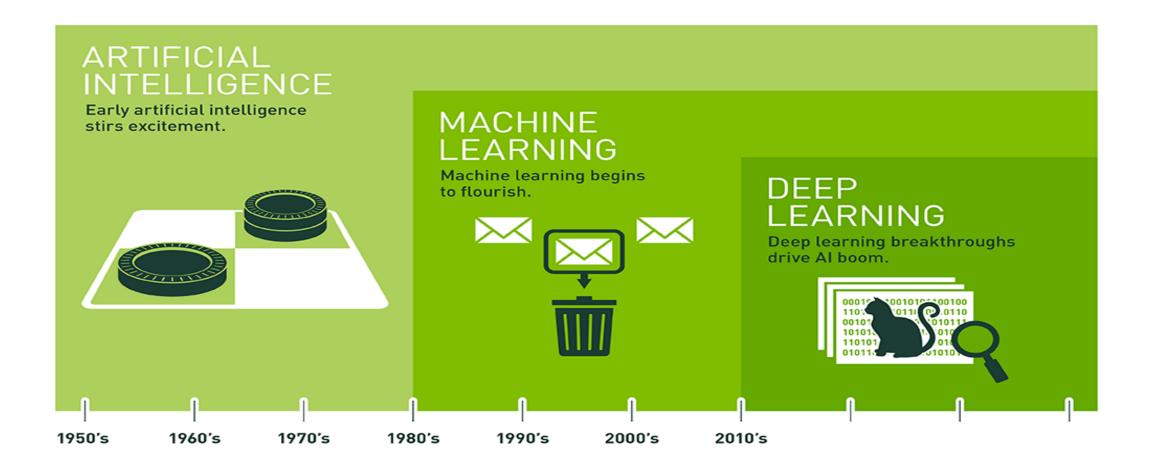


# AI == CODE

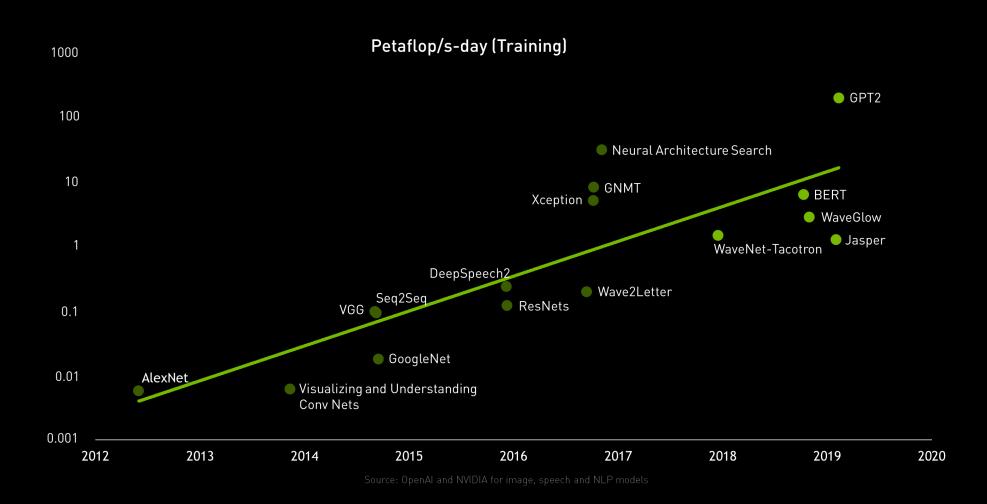
```
class LeNet(nn.Module):
    def forward(self, x):
        out = F.relu(self.conv1(x))
        out = F.relu(self.conv2(out))
        out = out.view(out.size(0), -1)
        out = F.relu(self.fc1(out))
        out = F.relu(self.fc2(out))
        out = self.fc3(out)
        return out
```



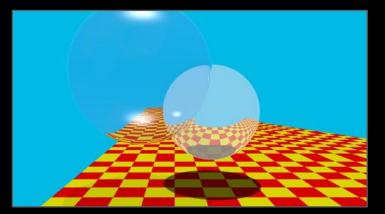
## **Definitions**

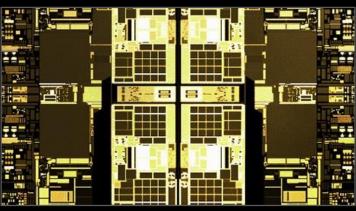


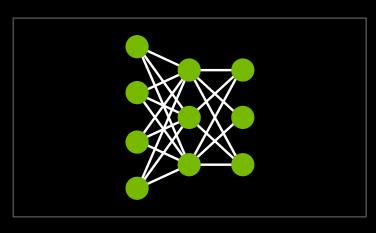
# **NETWORK COMPLEXITY IS EXPLODING**



# **NVIDIA RESEARCH**







RTX

NVSwitch

CuDNN



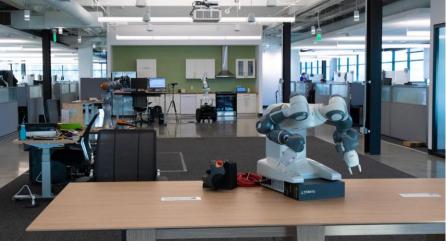




**I**mage

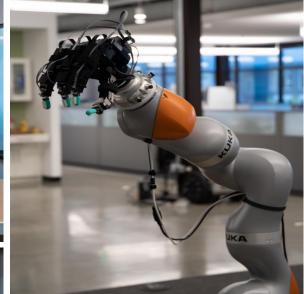
Noise-to-Noise Denoising

Progressive GAN





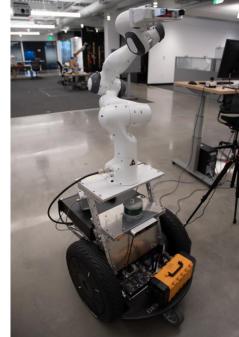






#### NVIDIA ROBOTICS RESEARCH LAB SEATTLE

Drive breakthrough robotics research to enable the next-generation of robots that safely work alongside humans, transforming industries such as manufacturing, logistics, healthcare, and more



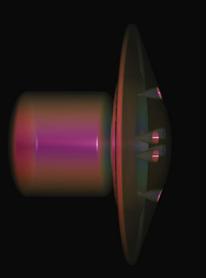
# DexPilot: Depth-Based Teleoperation of Dexterous Robotic Hand-Arm System

Ankur Handa\*<sup>†</sup>, Karl Van Wyk\*<sup>†</sup>, Wei Yang<sup>†</sup>, Jacky Liang<sup>‡</sup>, Yu-Wei Chao<sup>†</sup>, Qian Wan<sup>†</sup>, Stan Birchfield<sup>†</sup>, Nathan Ratliff<sup>†</sup>, Dieter Fox<sup>†</sup>





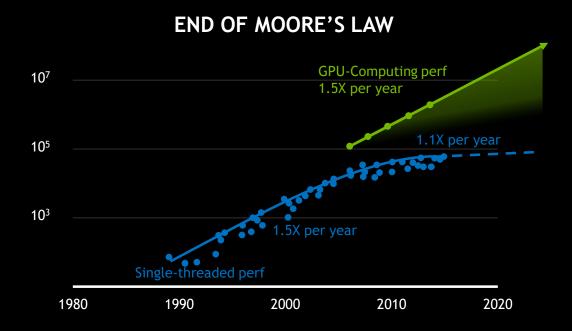




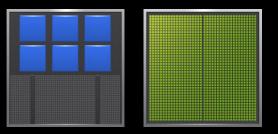


## THE RISE OF GPU COMPUTING

#### Big Data Needs Algorithms and Compute That Scales



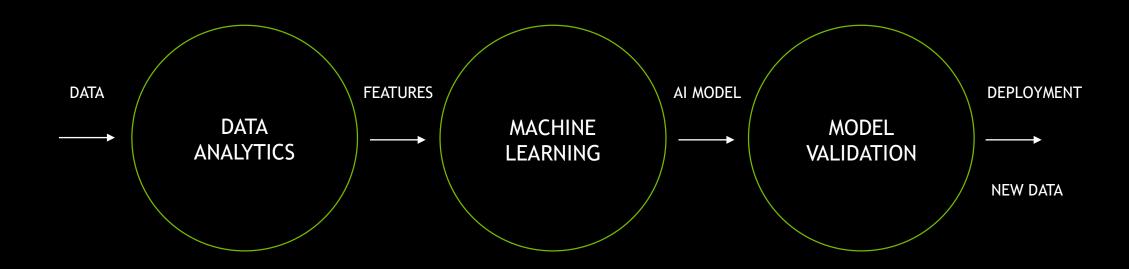
CPU vs. GPU



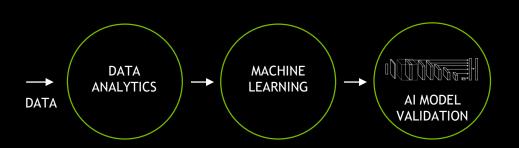
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

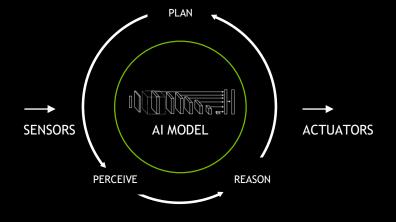


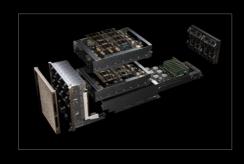
# BUILDING AN AI MODEL

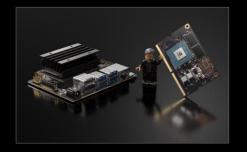


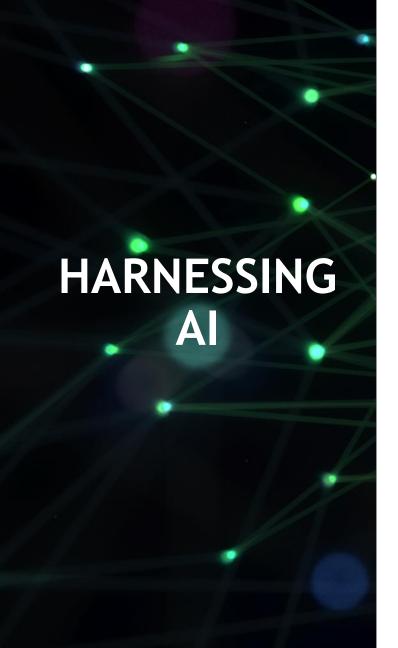
# **BUILDING AN AI PRODUCT**











Step I: Build a data fabric for your organization

Step II: Define your objective

Step III: Hire the right talent

Step IV: Identify key processes to augment with Al

Step V: Create a sandbox lab environment

Step VI: Operationalize successful pilots

Step VII: Scale up for enterprise-wide adoption

Step VIII: Drive cultural change





Designing and building state-of-art Automated Machine Learning software and services for real world problems

#### Background

Spin out from the University of Oxford

Combining leading edge machine learning research with advanced software engineering

Working with leading companies to deliver insight & quantifiable benefits

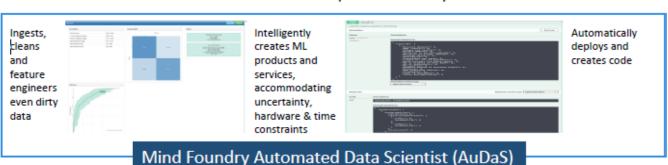
Setting benchmarks for a new wave of machine learning

#### Vision

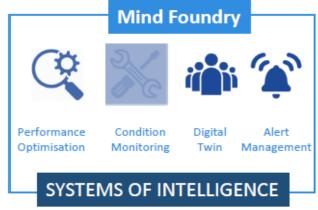
Build services and software that Automate the Data Science lifecycle

Services which understand the entanglement, uncertainty, complexity and power of machine learning

Services that are evergreen and quantum ready





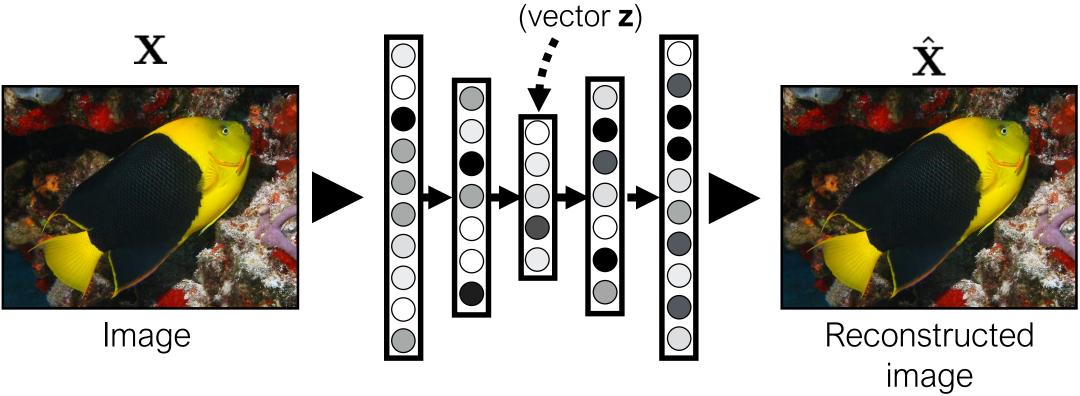




Mind Foundry, www.mindfoundry.ai

# Representation Learning

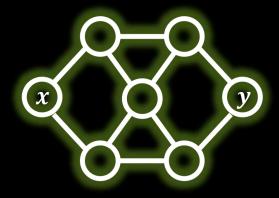
compressed image code



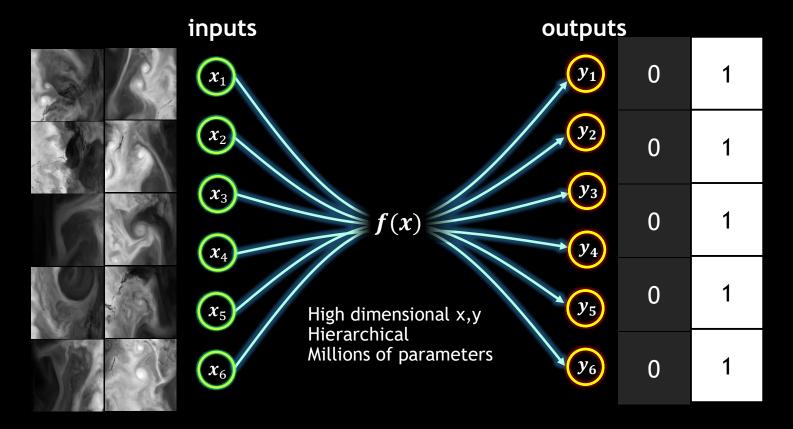


"Autoencoder"

Find f, given x and y



Supervised Deep Learning



# tgr strp

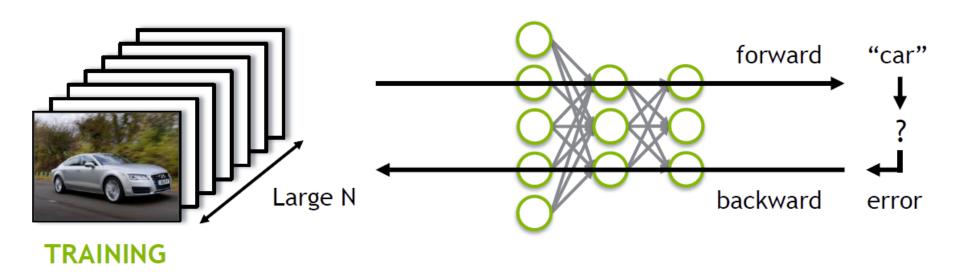
# tiger stripe

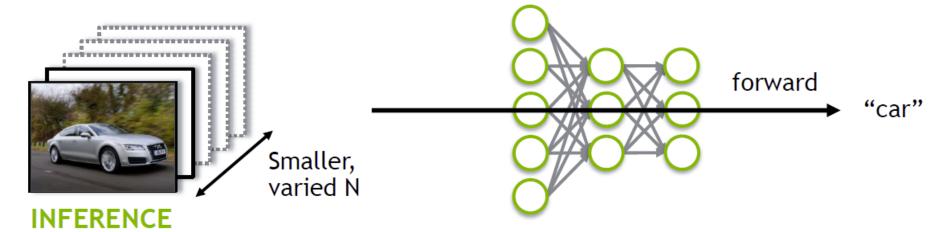
# Black & Orange



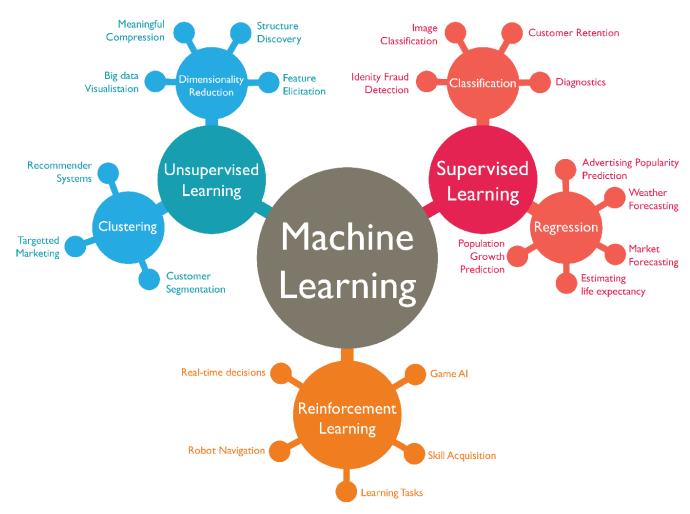


# TRAINING VS INFERENCE

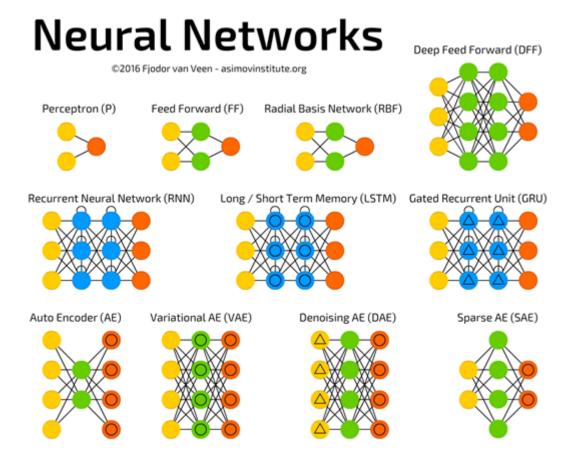




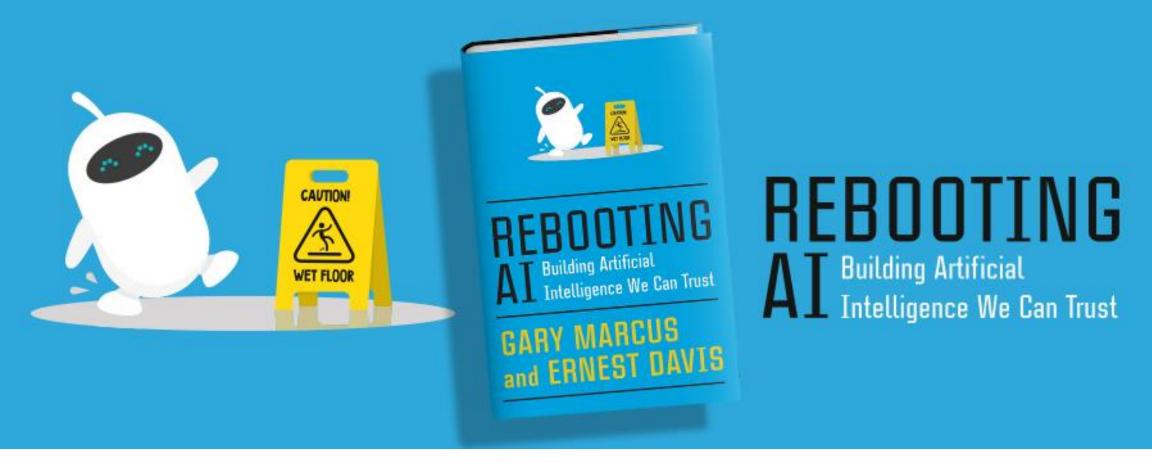
# Types of ML/DL



## **ARCHITECTURES**



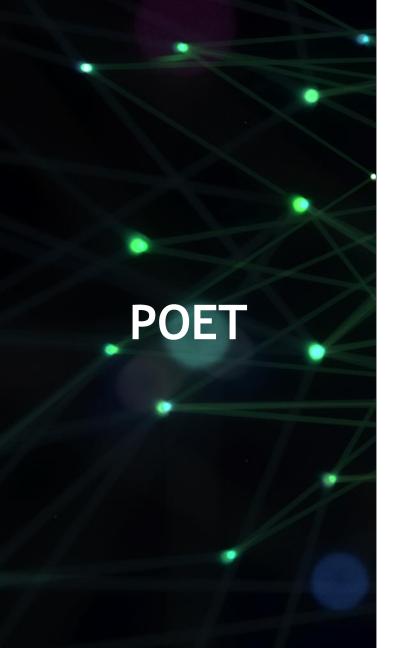
Larger image: http://www.asimovinstitute.org/neural-network-zoo/



#### www.Robust.ai

Gray Marcus, Rodney Brooks, Steven Pinker et al





#### Paired Open-Ended Trailblazer (POET): Endlessly Generating Increasingly Complex and Diverse Learning Environments and Their Solutions

Rui Wang Joel Lehman Jeff Clune\* Kenneth O. Stanley\*

Uber AI Labs
San Francisco, CA 94103
ruiwang, joel.lehman, jeffclune, kstanley@uber.com

\*co-senior authors

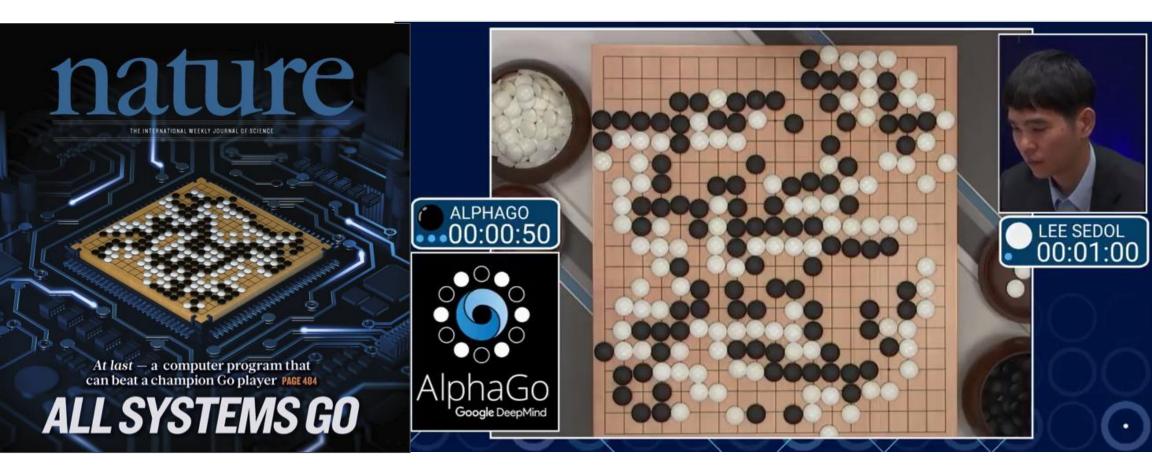
#### Abstract

While the history of machine learning so far encompasses a series of problems posed by researchers and algorithms that learn their solutions, an important question is whether the problems themselves can be generated by the algorithm at the same time as they are being solved. Such a process would in effect build its own diverse and expanding curricula, and the solutions to problems at various stages would become stepping stones towards solving even more challenging problems later in the process. The Paired Open-Ended Trailblazer (POET) algorithm introduced in this paper does just that: it pairs the generation of environmental challenges and the optimization of agents to solve those challenges. It simultaneously explores many different paths through the space of possible problems and solutions and, critically, allows these stepping-stone solutions to transfer between problems if better, catalyzing innovation. The term open-ended signifies the intriguing potential for algorithms like POET to continue to create novel and increasingly complex capabilities without bound. We test POET in a 2-D bipedal-walking obstacle-course domain in which POET can modify the types of challenges and their difficulty. At the same time, a neural network controlling a biped walker is optimized for each environment. The results show that POET produces a diverse range of sophisticated behaviors that solve a wide range of environmental challenges, many of which cannot be solved by direct optimization alone, or even through a direct, single-path curriculum-based control algorithm introduced to highlight the critical role of open-endedness in solving ambitious challenges. The ability to transfer solutions from one environment to another proves essential to unlocking the full potential of the system as a whole, demonstrating the unpredictable nature of fortuitous stepping stones. We hope that POET will inspire a new push towards open-ended discovery across many domains, where algorithms like POET can blaze a trail through their interesting possible manifestations and solutions.



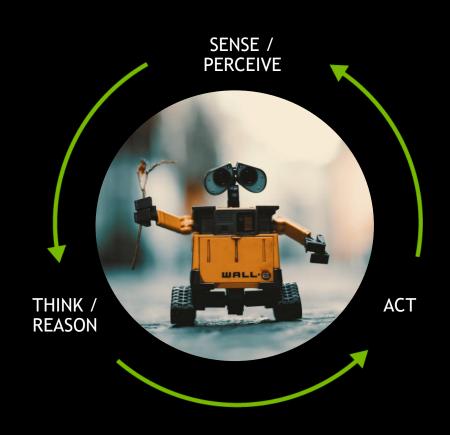


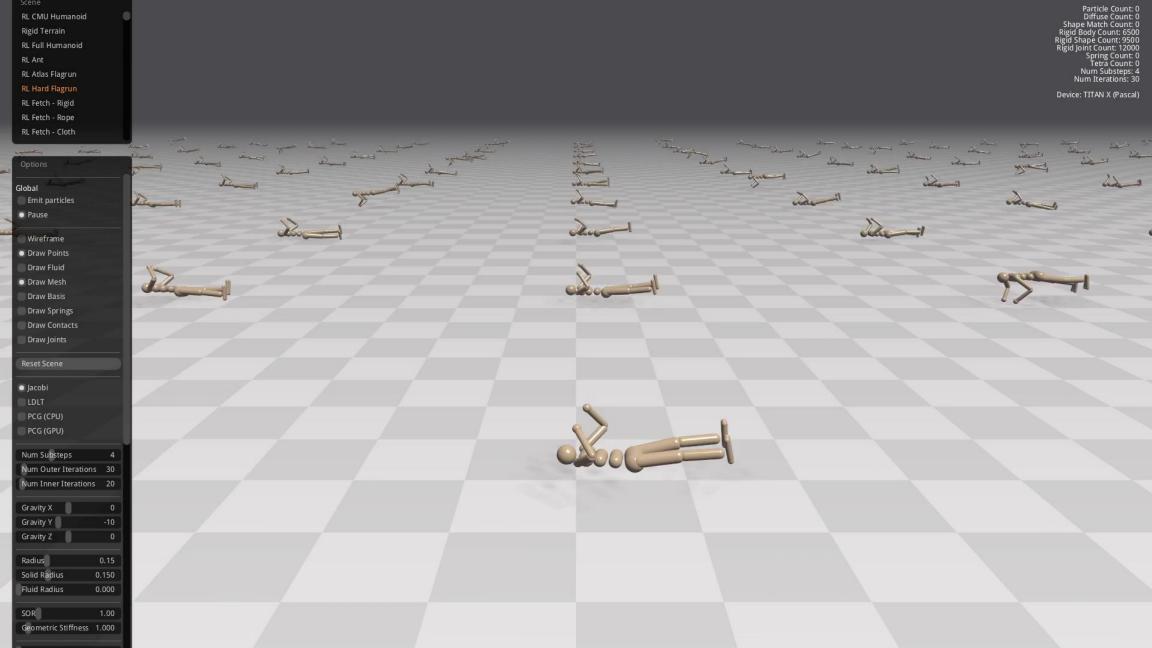
# **DEEPMIND ALPHA\***

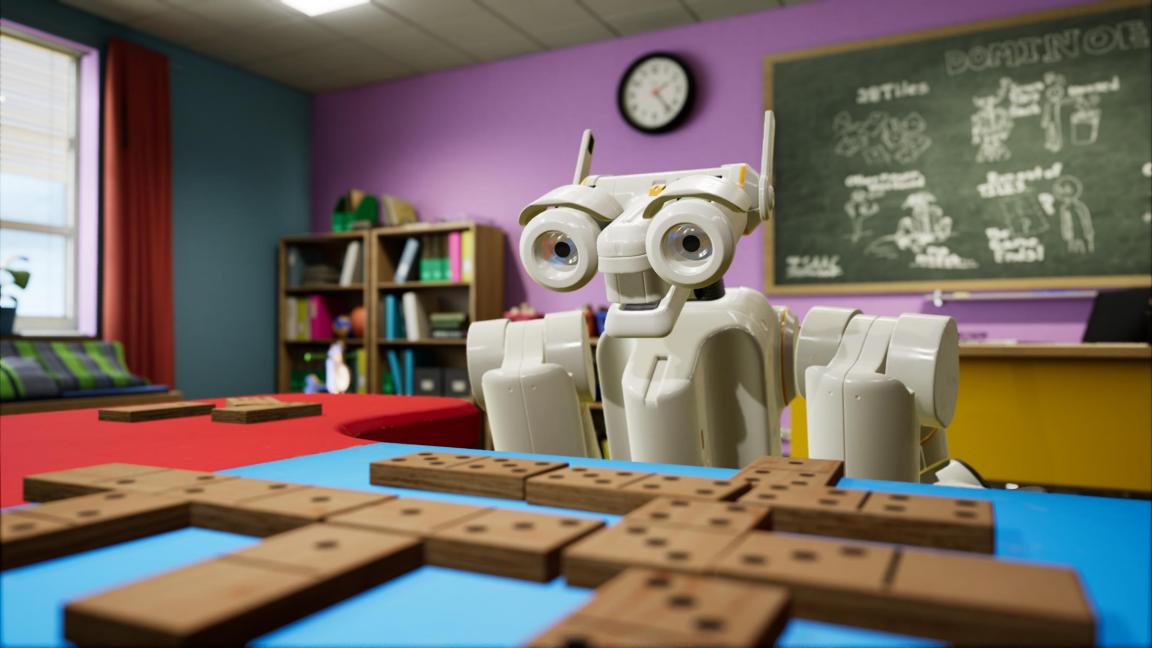


<sup>\* ..</sup>a deeply structured hybrid (Gary Marcus, Jan 2018)

## **ROBOTS**









## ANNOUNCING ISAAC OPEN SDK







Isaac Robot Engine Isaac Sim Isaac Gym

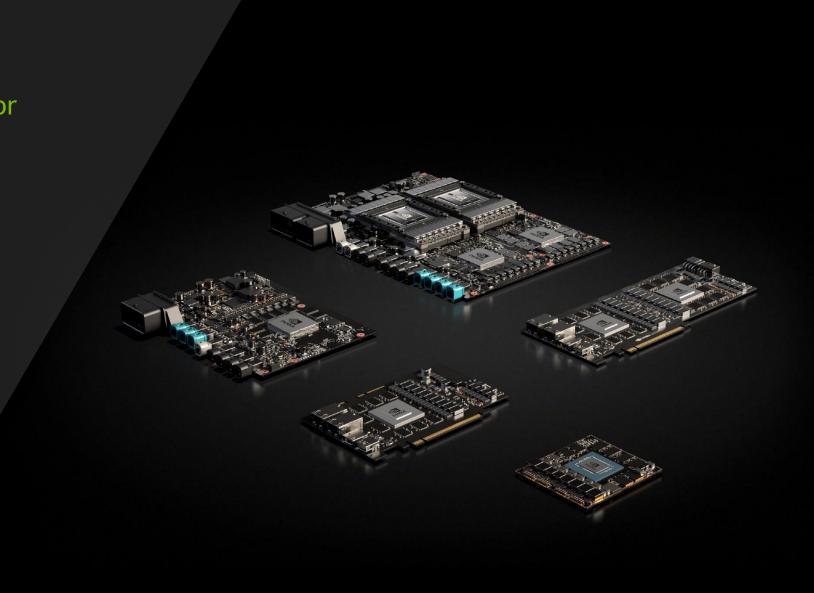
Isaac Robot Engine - Modular robot framework | Isaac Sim - Virtual robotics laboratory
Isaac Gym - Reinforcement learning simulator | Isaac Robot Apps - Kaya, Carter and Link

Available at developer.nvidia.com/isaac-sdk

### **NVIDIA AGX**

Family of Systems for Embedded AI HPC

Self-driving cars Robotics Smart Cities Healthcare

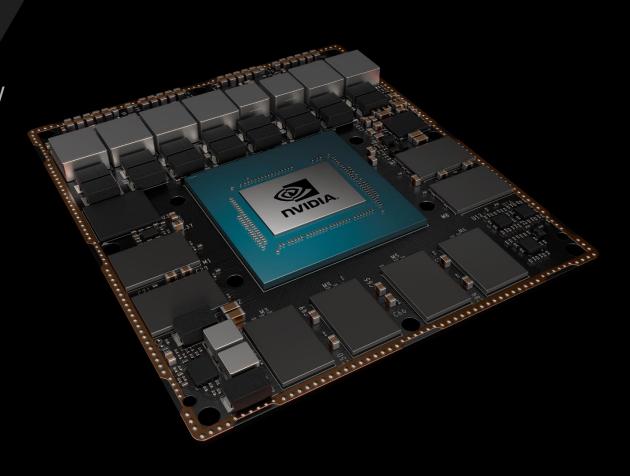


### **JETSON XAVIER**

### for Autonomous Machines

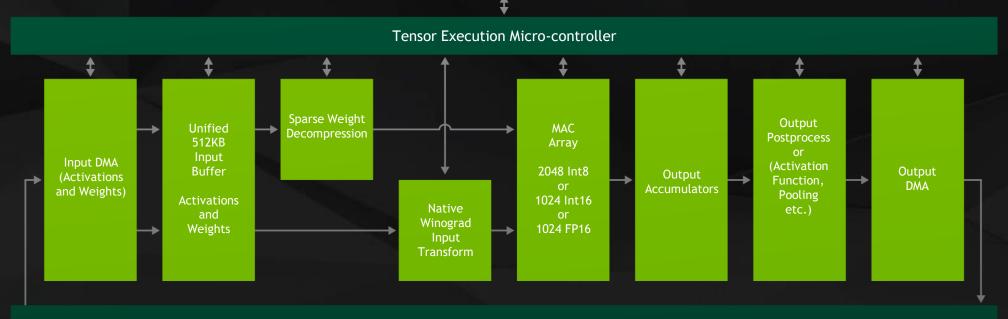
Al Server Performance in 30W • 15W • 10W 512 Volta CUDA Cores • 2x NVDLA 8 core CPU 32 DL TOPS

developr.nvidia.com/jetson-xavier



# XAVIER DLA NOW OPEN SOURCE

Command Interface

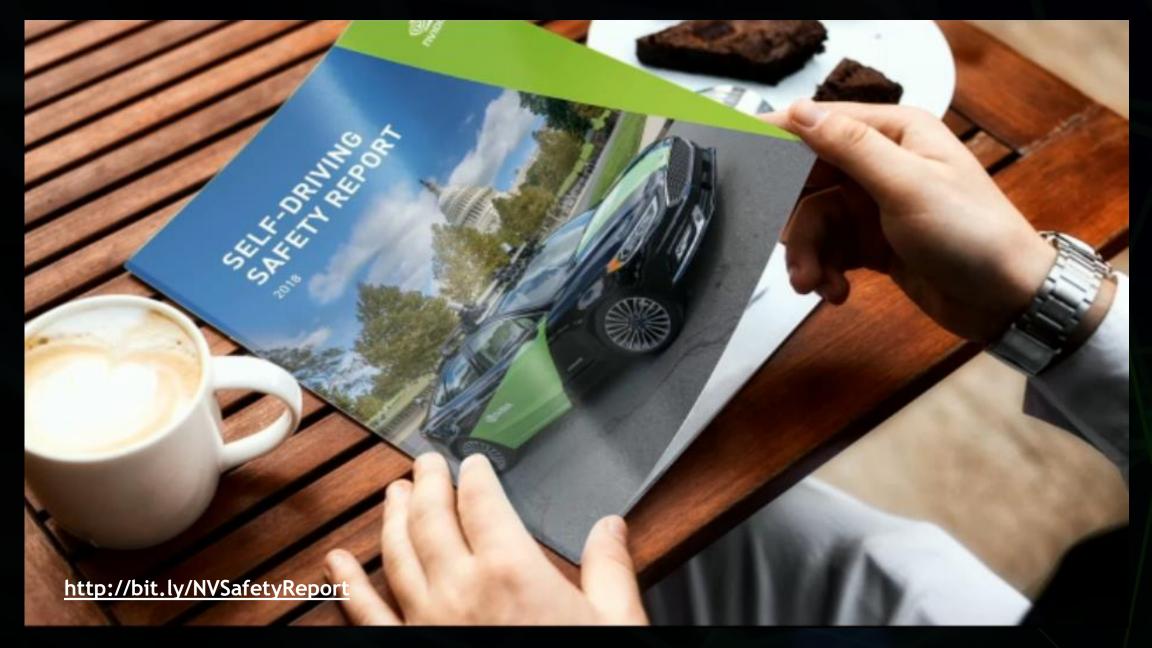


Memory Interface

WWW.NVDLA.ORG







### **DRIVE CONSTELLATION**

Virtual Reality AV Simulator

Runs DRIVE Sim Simulator

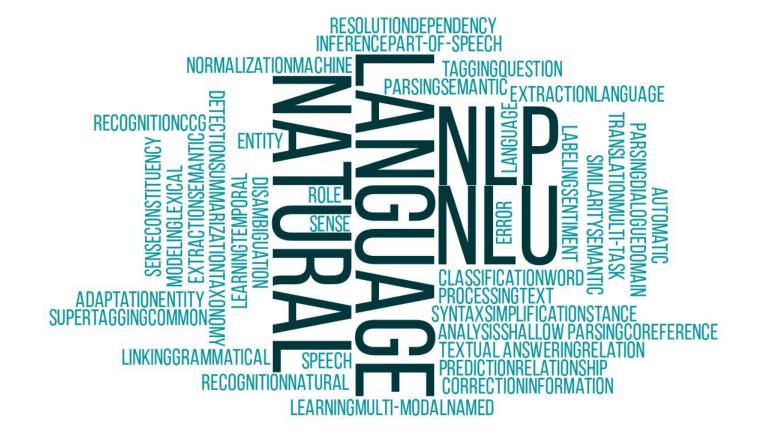
Hardware in the Loop System Level Simulator

Simulate Rare and Difficult Conditions

Scalable Platform | Data Center Solution

Timing Accurate and Bit Accurate







**Detection** 



**Planning** 



Acceleration



**Assimilation** 





**Enhancement Parametrization** 



**Prediction** 



Augmentation



#### Monitor Environmental Change



drought flooding deforestation urbanification melting glaciers sea-level change

Detection



**Planning** 



Acceleration



Assimilation



......



**Enhancement Parametrization** 

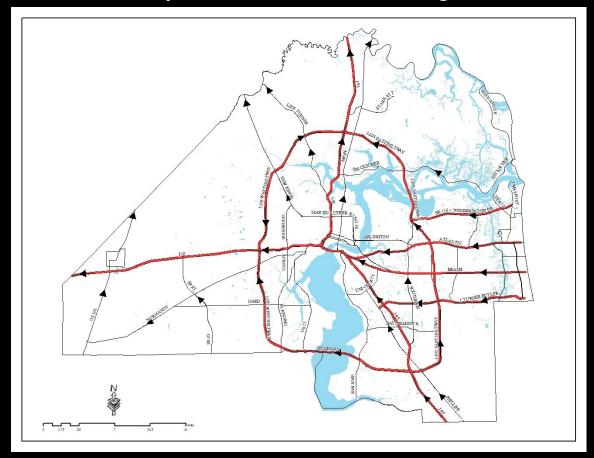


Augmentation





### **Optimize Disaster Planning**



Detection



**Planning** 



Acceleration



**Assimilation** 







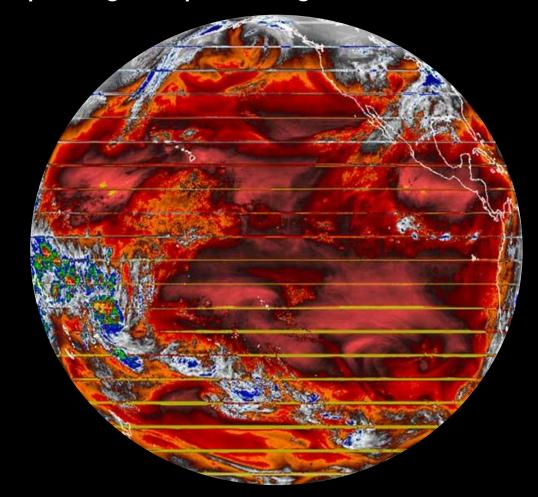
**Prediction** 



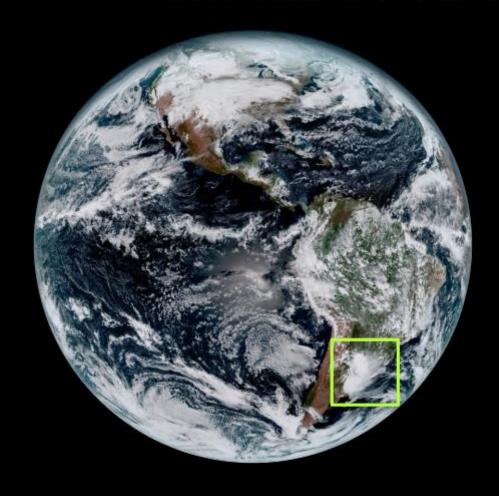
Augmentation



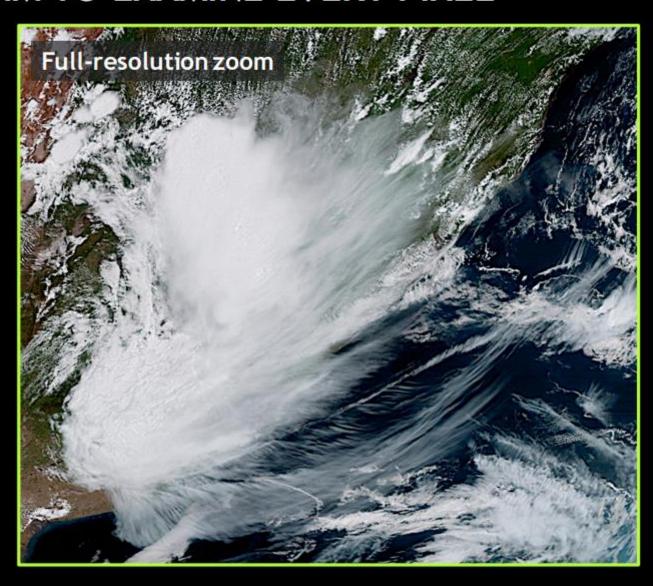
### Use Inpainting to Repair Damaged GOES-17 Observations



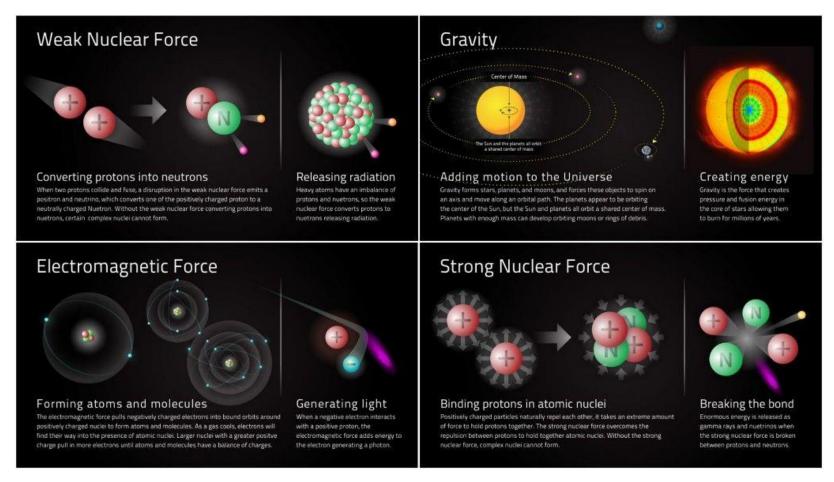
### TRAIN AN ALGORITHM TO EXAMINE EVERY PIXEL



GOES-16: 4k x 4k x 11 channels



### THE FOUR FUNDAMENTAL FORCES OF THE UNIVERSE





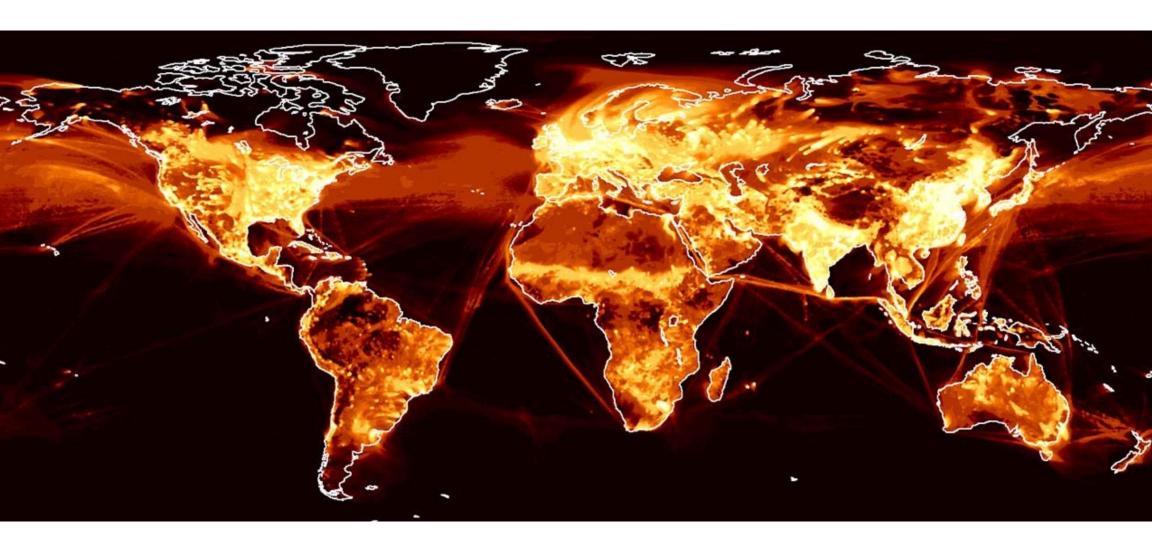
### TWO FUNDAMENTAL NEEDS



Fast filtering, FFTs, correlations, convolutions, resampling, etc to process increasingly larger bandwidths of signals at increasingly fast rates and do increasingly cool stuff we couldn't do before

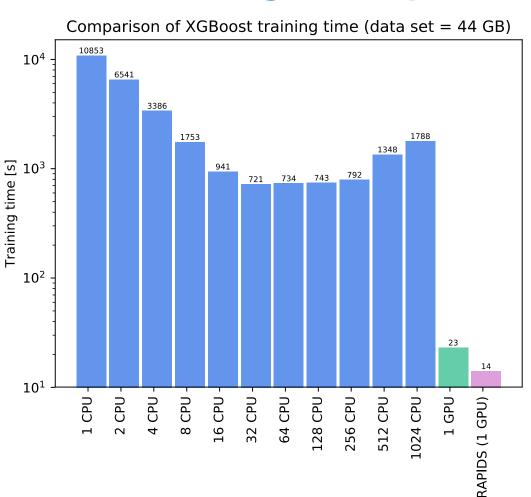


Artificial Intelligence techniques applied to spectrum sensing, signal identification, spectrum collaboration, and anomaly detection





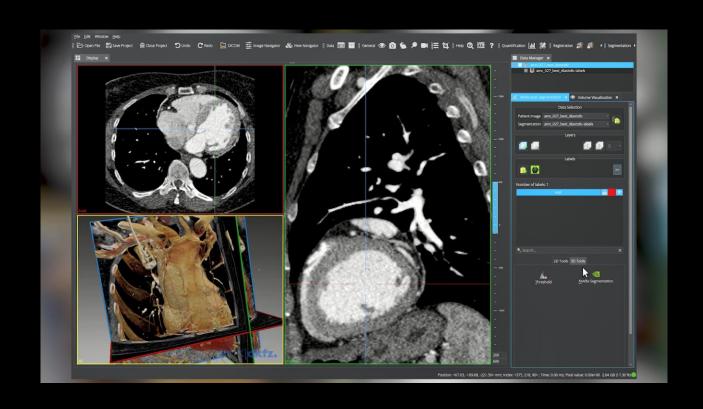
### XGBoost for simulating atmospheric chemistry







# **AI-ASSISTED ANNOTATION**



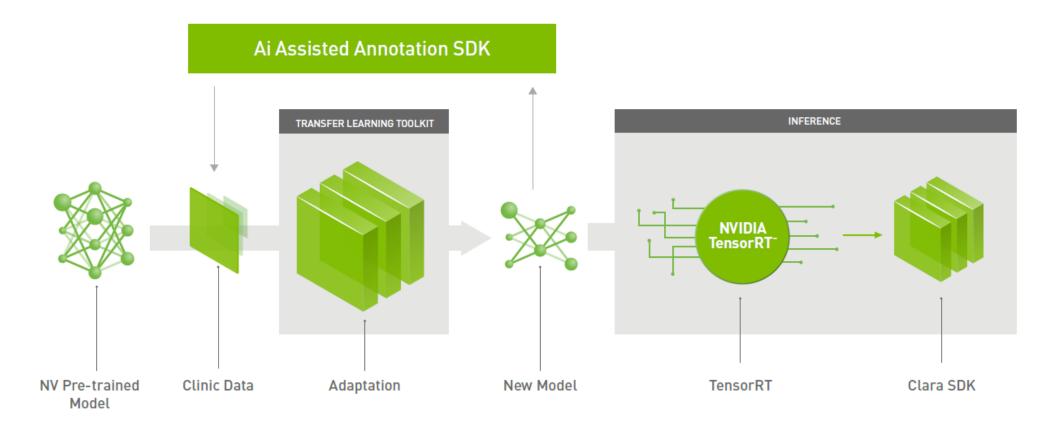
## **CLARA AI TOOLKIT**

### Build, Manage And Deploy AI Applications For Radiology



# End to End NVIDIA Deep Learning Workflow

Pre-Trained models \* Annotation Assistant \* Training & adaptation \* Applications ready to integrate with Clara Platform



### **CONVERGENCE OF HPC AND AI**

### Integrating the Third and Fourth Pillars of Scientific Discovery

#### **HPC**

40+ years of algorithms based on first principles theory

#### Al

New algorithms and models with potential to increase model size and accuracy



Dramatically Improves Accuracy and /or Time-to-Solution at Large Scale



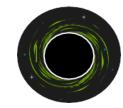
Commercially viable fusion energy



Improve or validate the Standard Model of Physics



Clinically viable precision medicine



Understanding cosmological dark energy and matter

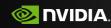


Climate/weather forecasts with ultrahigh fidelity

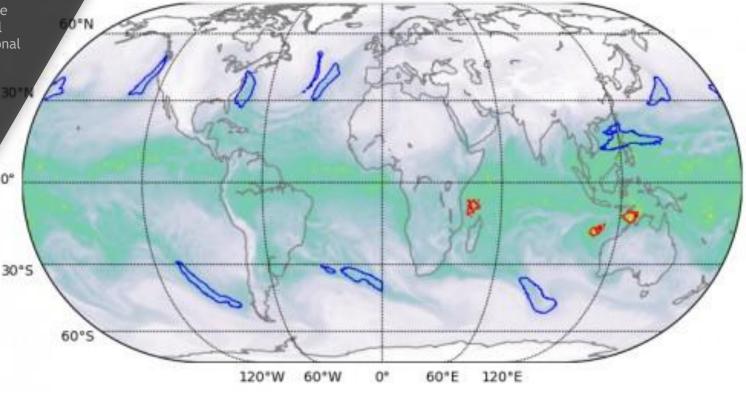


# EXASCALE AI FOR CLIMATE PREDICTION

The ability to accurately predict the path of extreme weather systems can save lives and safeguard global economies. Researchers at Lawrence Berkeley National Laboratory used a climate dataset on the Summit supercomputer with NVIDIA Volta Tensor Core GPUs to train a deep neural network to identify extreme weather patterns from high-resolution climate simulations. They achieved a performance of 1.13 exaflops, the fastest deep learning algorithm reported.







Pictured: high-quality segmentation results produced by deep learning on climate datasets. Image credit: NERSC

### FIVE ROADS TO GPU COMPUTING

#### **GPU Libraries**

Drop-in replacement for existing libraries

cuBLAS, CUDA Math, cuSPARSE, cuRAND, cuSOLVER, nvGRAPH, cuDNN, cuFFT, Thrust

#### **OPEN-ACC**

Comment-based directives in C / C++ / Fortran

Single source code parallelization for multiple architectures

#### **CUDA**

Parallel Programming Model for GPUs in C, C++, Fortran, Python, MATLAB

Specialized Kernels for general purpose GPU

#### **RAPIDS**

GPU Acceleration of Traditional Machine Learning

Accelerate Scikit-Learn style ML algorithms

#### **DEEP LEARNING**

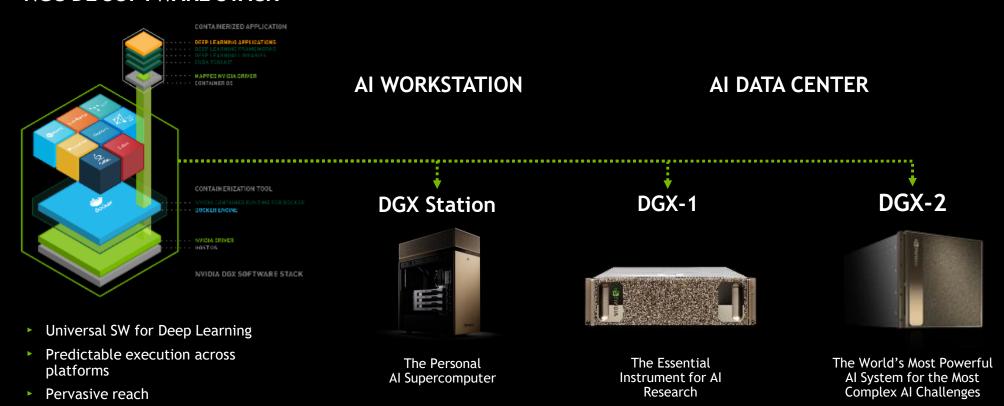
GPU accelerated deep learning frameworks TensorFlow, Pytorch

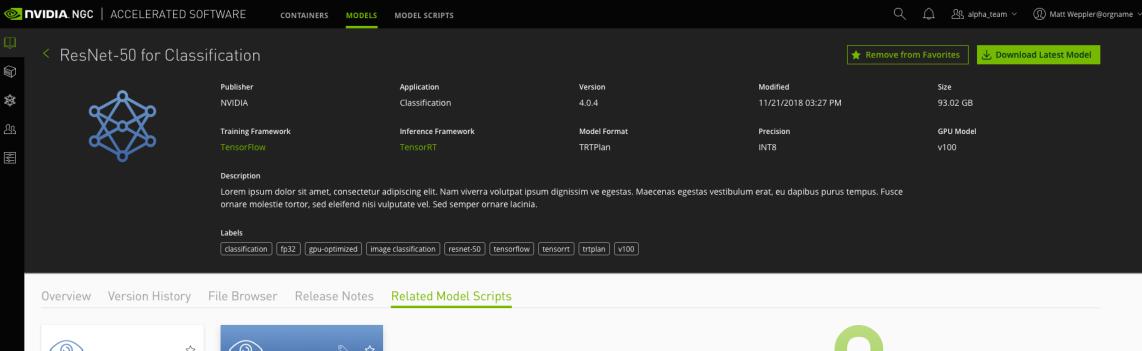
Build GPU-accelerated functions directly from data



### PURPOSE-BUILT AI SUPERCOMPUTERS

#### **NGC DL SOFTWARE STACK**







\$

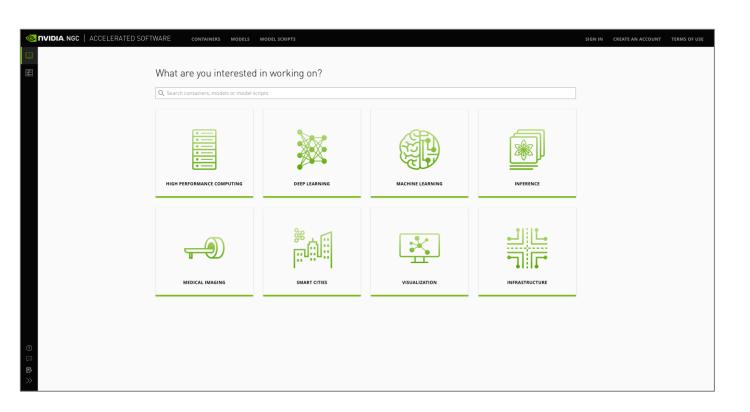






### **GET STARTED WITH NGC**

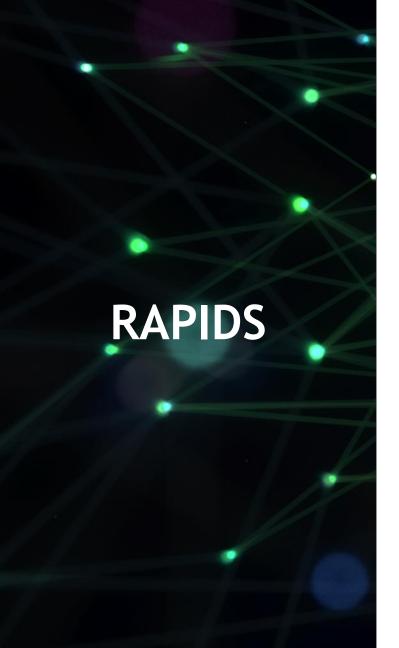
### Explore the NGC Registry for DL, ML & HPC



Deploy containers: ngc.nvidia.com

Learn more about NGC offering: nvidia.com/ngc

Technical information: developer.nvidia.com

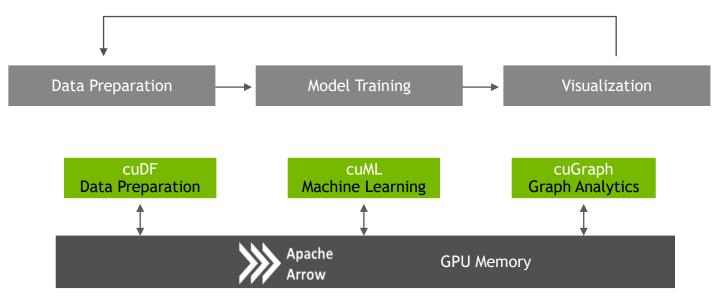


### **RAPIDS**

#### **GPU Accelerated End-to-End Data Science**

RAPIDS is a set of open source libraries for GPU accelerating data preparation and machine learning.

OSS website: rapids.ai



#### **NVIDIA SDK**

The Essential Resource for GPU Developers



#### **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



#### **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



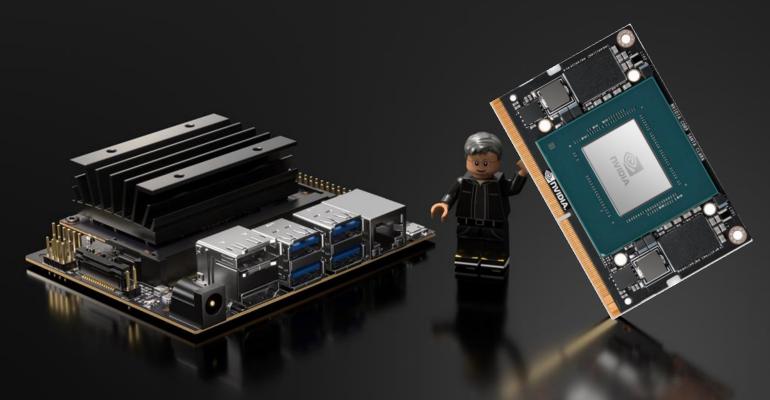
### developer.nvidia.com

# **NVIDIA DGX-2**



# JETSON NANO DEVKIT & XAVIER NX SOM

Up to 21 DL TOPS (15w) | NX: 8 GB Memory | 45x70mm



CUDA-X acceleration stack | High-resolution sensor support | Runs all CUDA-X AI models NX available from nvidia.com and distributors worldwide in March 2020

# DEEPSTREAM ON JETSON NANO



### **JETSON - START NOW**



JETSON DEVELOPER KIT

AGX Xavier Developer Kit \$699
Xavier NX software patch
developer.nvidia.com/
buy-jetson



#### TWO DAYS TO A DEMO

Create your first demo today developer.nvidia.com/ embedded/twodaystoademo



#### **DEEP LEARNING INSTITUTE**

Training • Labs Nanodegrees nvidia.com/DLI



#### GTC

Largest event for GPU developers
gputechconf.com

# NVIDIA DEEP LEARNING INSTITUTE

Online self-paced labs and instructor-led workshops on deep learning and accelerated computing

Take self-paced labs at www.nvidia.co.uk/dlilabs

View upcoming workshops and request a workshop onsite at www.nvidia.co.uk/dli

Educators can join the University Ambassador Program to teach DLI courses on campus and access resources. Learn more at www.nvidia.com/dli

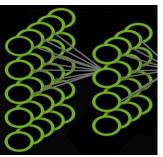








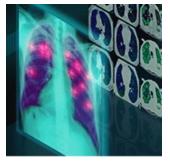




**Fundamentals** 



**Autonomous Vehicles** 



Healthcare



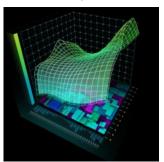
Intelligent Video Analytics



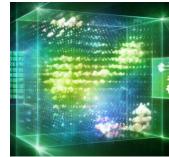
**Robotics** 



Game Development & Digital Content



**Finance** 



**Accelerated Computing** 



Virtual Reality





March 22-26, 2020 | San Jose, CA



#### CONNECT

Connect with experts from NVIDIA, GE Healthcare, NSF Carnegie Mellon, Google, and other leading organizations



**LEARN** 

Gain insight and valuable hands-on training through over 100 sessions



#### **DISCOVER**

See how GPU technologies are creating amazing breakthroughs in important fields such as deep learning



#### **INNOVATE**

Explore disruptive innovations that can transform your work

Join us | Use VIP code NVALOWNDES for 25% off

Don't miss the premier Al conference.

nvidia.com/en-us/gtc/



alowndes@nvidia.com