

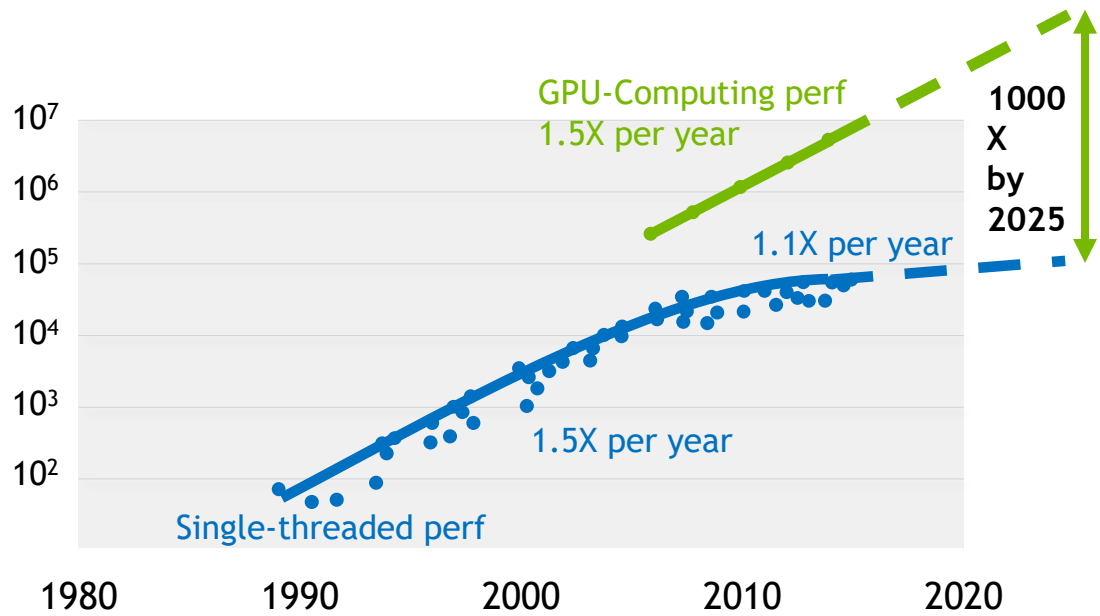


# NVIDIA GPU: WHAT AND WHY

Reynaldo Gomez

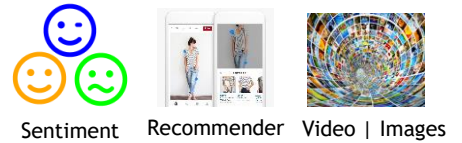
[reynaldog@nvidia.com](mailto:reynaldog@nvidia.com)

# ACCELERATING HIGH-GROWTH CLOUD WORKLOADS

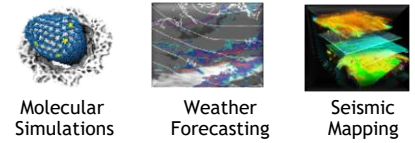


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

Artificial Intelligence



High Performance Computing



Rendering & Graphics



CPU No Longer Keeping Pace With Demand

Growth Workloads Require Acceleration

# WHY PARALLEL PROCESSING ?



Art, Science and GPU's  
Adam Savage & Jamie Hyneman  
Explain Parallel Processing



# NVIDIA DATA CENTER PLATFORM

Single Platform Drives Utilization and Productivity

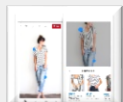
## CUSTOMER USE CASES



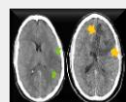
Speech



Translate



Recommender



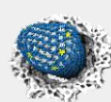
Healthcare



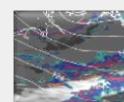
Manufacturing



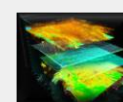
Finance



Molecular Simulations



Weather Forecasting



Seismic Mapping



Creative & Technical



Knowledge Workers

CONSUMER INTERNET & INDUSTRY APPLICATIONS

SCIENTIFIC APPLICATIONS

VIRTUAL GRAPHICS

## APPS & FRAMEWORKS



python



TensorFlow



mxnet



Chainer



ONNX

RAPIDS

PYTORCH

Amber  
NAMD

+600  
Applications



CATIA



Ps



Windows 10

AUTODESK  
3DS MAX

MACHINE LEARNING

cuDF

cuML

cuGRAPH

DEEP LEARNING

cuDNN

CUTLASS

TensorRT

HPC

OpenACC

cuFFT

VIRTUAL GPU

vDWS

vPC

vAPPS

## CUDA-X & NVIDIA SDKs

CUDA & CORE LIBRARIES - cuBLAS | NCCL

## TESLA GPUs & SYSTEMS



TESLA GPU



NVIDIA DGX FAMILY



NVIDIA HGX



DELL



Hewlett Packard  
Enterprise



IBM



inspur

SYSTEM OEM



CLOUD

# AI TRAINING REQUIRES FULL STACK INNOVATION



Workload: Time (in minutes) to train ResNet-50 network.

# TESLA V100 TENSOR CORE GPU

World's Most Powerful  
Data Center GPU

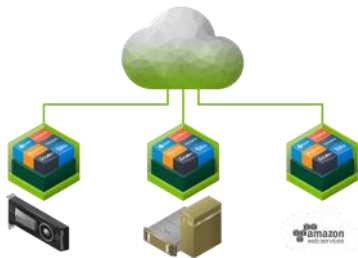
5,120 CUDA cores  
640 NEW Tensor cores  
7.8 FP64 TFLOPS | 15.7 FP32 TFLOPS  
| 125 Tensor TFLOPS  
20MB SM RF | 16MB Cache  
32 GB HBM2 @ 900GB/s |  
300GB/s NVLink



# THE DGX FAMILY OF AI SUPERCOMPUTERS

## CLOUD-SCALE AI

### NVIDIA GPU Cloud



Simple Access to GPU Accelerated Software

## AI WORKSTATION

### DGX Station



with



AI Workstation for Data Science Teams

## AI DATA CENTER

### DGX-1



with



The Essential Instrument for AI Research

### DGX-2



with



The World's Most Powerful AI System for the Most Complex AI Challenges

# CONTAINERS: SIMPLIFYING WORKFLOWS

## WHY CONTAINERS

### Simplifies Deployments

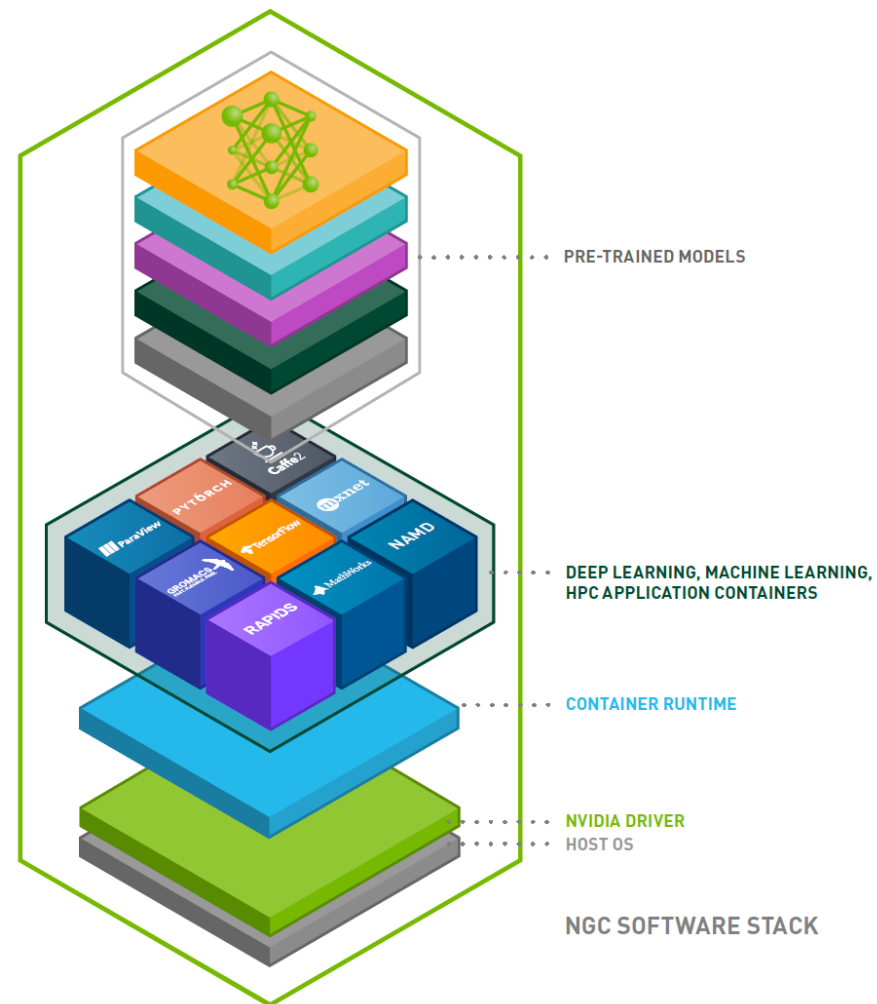
- Eliminates complex, time-consuming builds and installs

### Get started in minutes

- Simply Pull & Run the app

### Portable

- Deploy across various environments, from test to production with minimal changes



# NGC CONTAINERS: ACCELERATING WORKFLOWS

## WHY CONTAINERS

### Simplifies Deployments

- Eliminates complex, time-consuming builds and installs

### Get started in minutes

- Simply Pull & Run the app

### Portable

- Deploy across various environments, from test to production with minimal changes

## WHY NGC CONTAINERS

### Optimized for Performance

- Monthly DL container releases offer latest features and superior performance on NVIDIA GPUs

### Scalable Performance

- Supports multi-GPU & multi-node systems for scale-up & scale-out environments

### Designed for Enterprise & HPC environments

- Supports Docker & Singularity runtimes

### Run Anywhere

- Pascal/Volta/Turing-powered NVIDIA DGX, PCs, workstations, and servers
- From Core to the Edge
- On-Prem to Hybrid to Cloud

# DRAMATICALLY MORE FOR YOUR MONEY

**CPU-Only Cluster**

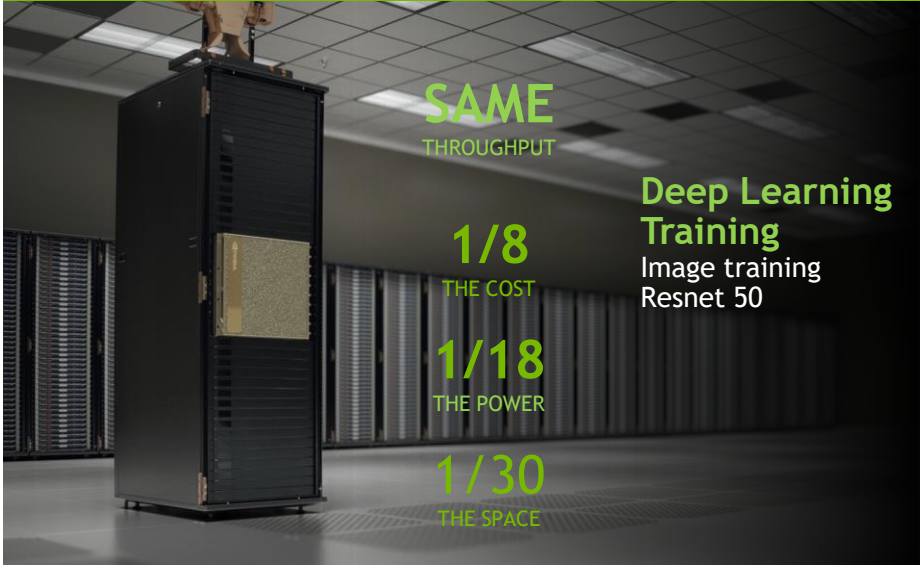


**Deep Learning Training**  
Image training  
Resnet 50

**300 Self-hosted Broadwell CPU Servers**  
**180 KWatts**

=

**GPU-Accelerated**



**SAME**  
THROUGHPUT

**1/8**  
THE COST

**1/18**  
THE POWER

**1/30**  
THE SPACE

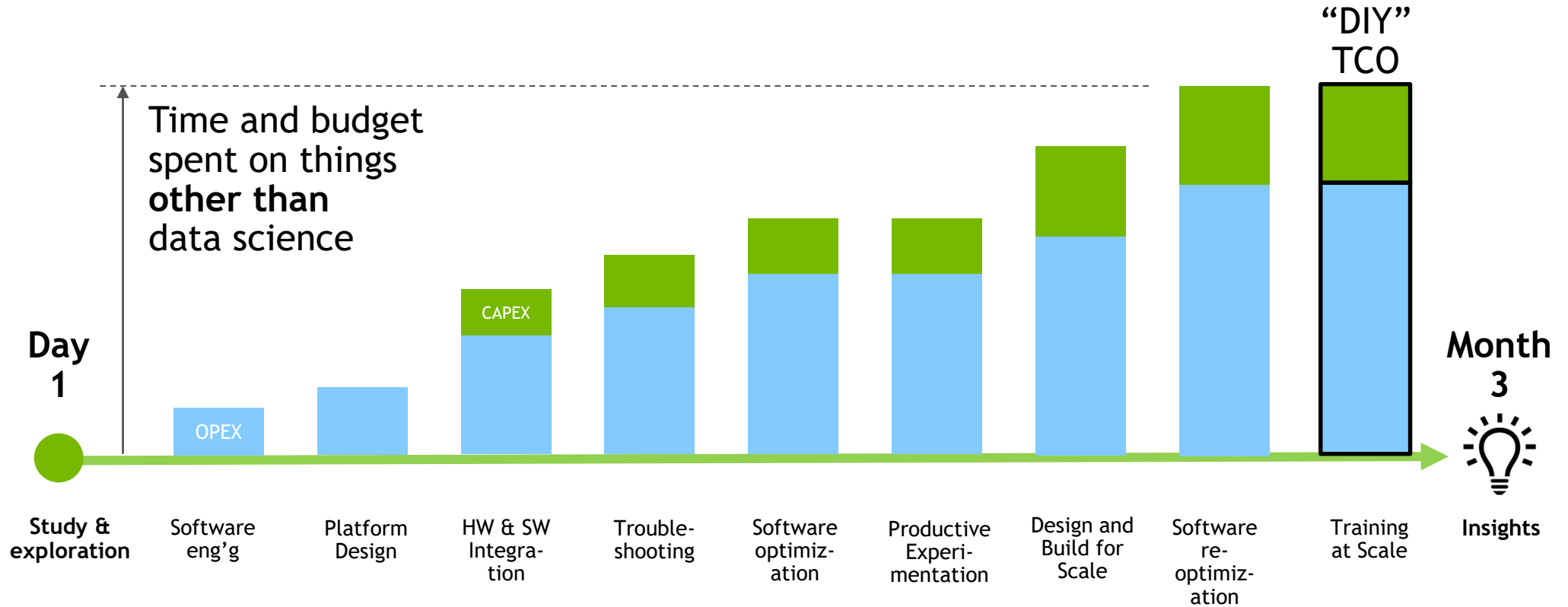
**Deep Learning Training**  
Image training  
Resnet 50

**1 DGX-2**  
**10 KWatts**

# HOW TESLA SAVES MONEY

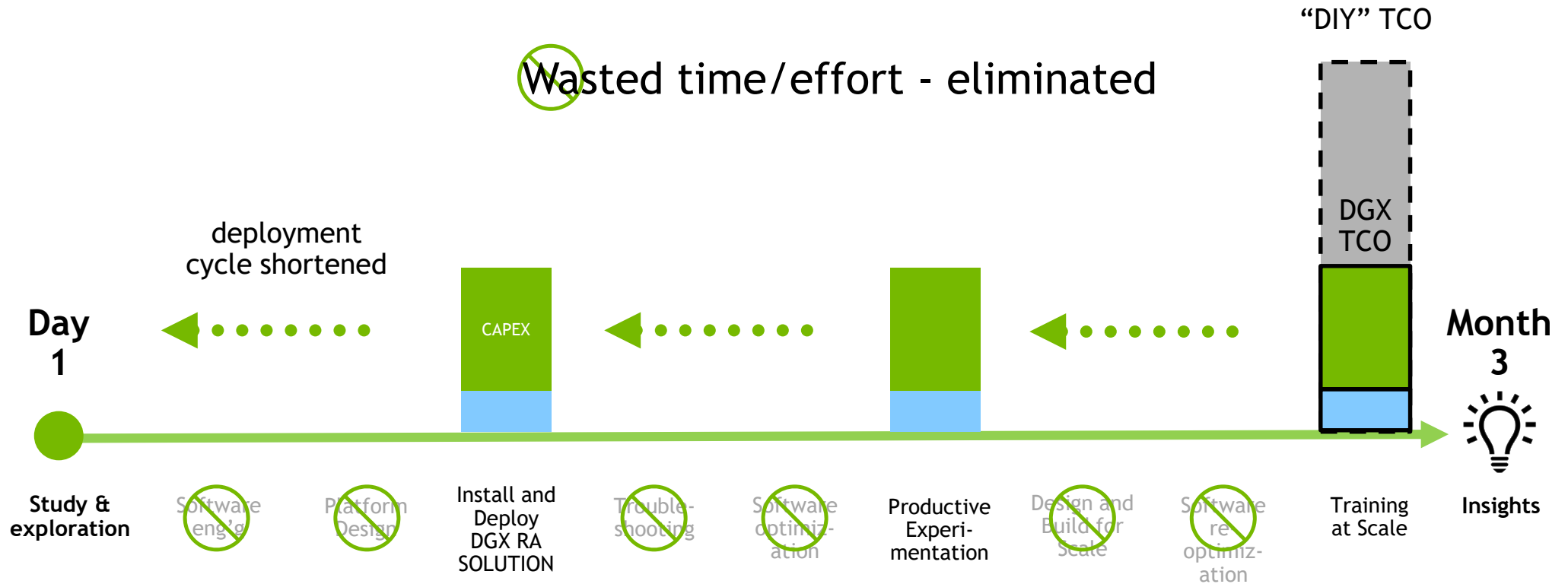
	GPU Node	CPU Node	# CPU Nodes To Match 1 GPU Node	\$ Spend on CPU Nodes	\$ Saved with GPU Node
AMBER	\$~45K	\$~9K	74 CPU Nodes	\$666K	\$621K
GTC			26 CPU Nodes	\$234K	\$189K

# AI SUCCESS DELAYED BY DEPLOYMENT COMPLEXITY



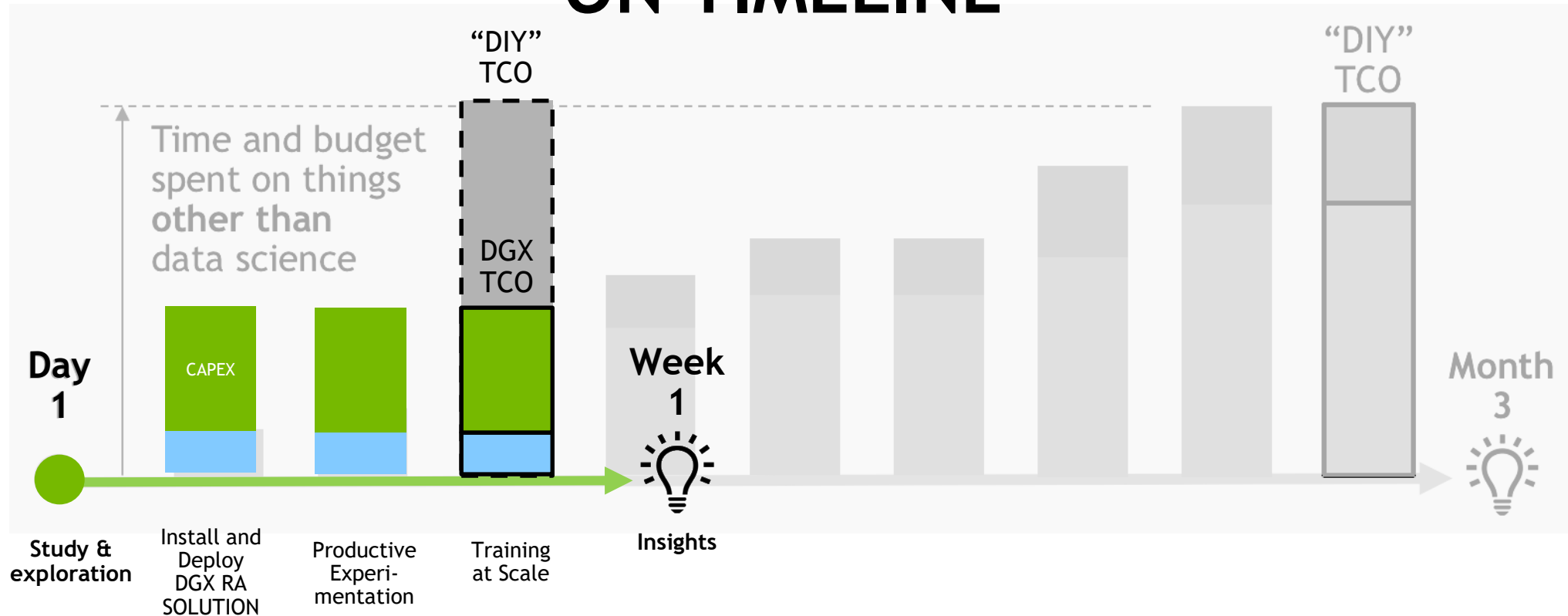
Designing, Building and Supporting an AI Infrastructure - from Scratch

# THE IMPACT OF DGX R/A SOLUTIONS ON TIMELINE



## 2. Deploying an Integrated, Full-Stack AI Solution using a DGX Reference Architecture



# THE IMPACT OF DGX R/A SOLUTIONS ON TIMELINE



## 2. Deploying an Integrated, Full-Stack AI Solution using a DGX Reference Architecture

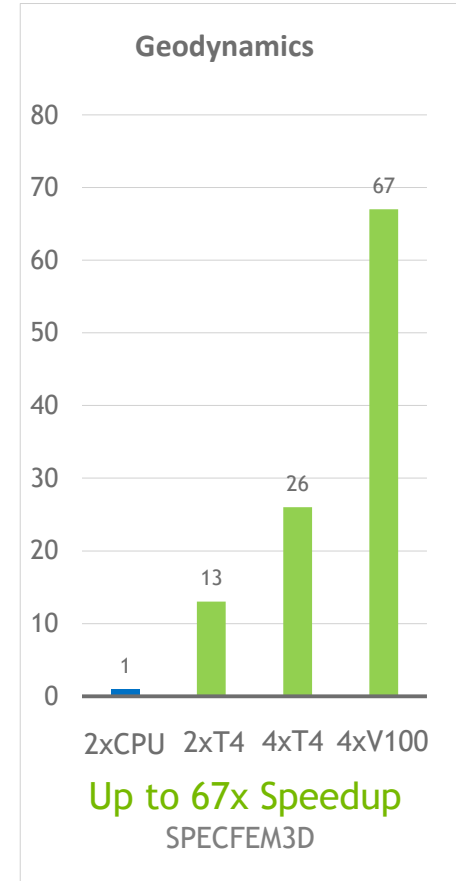
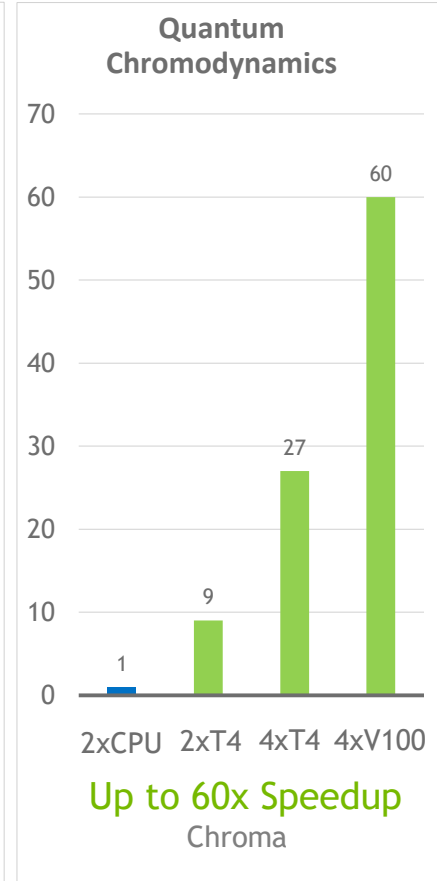
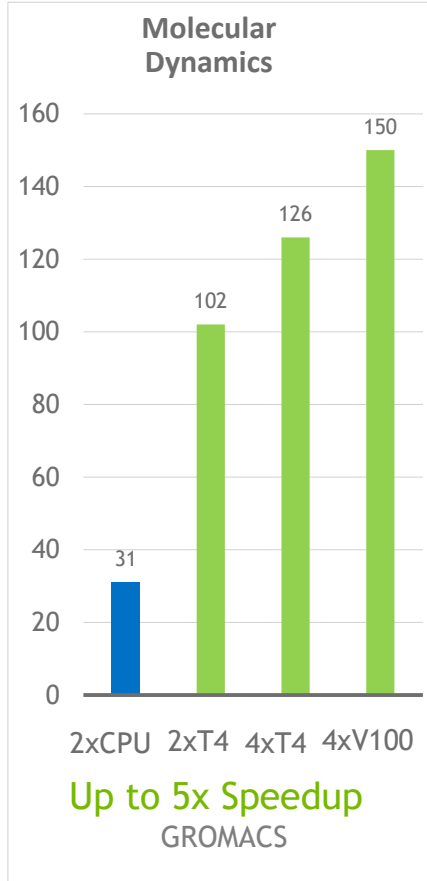
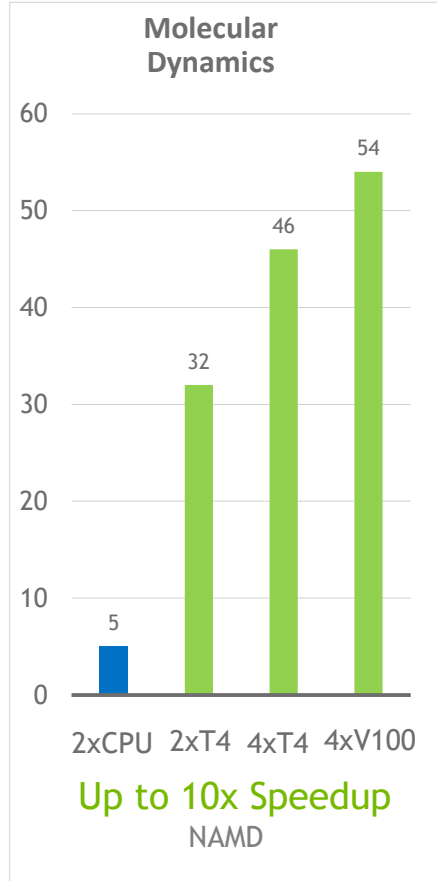
# NVIDIA GPUS ARE ON EVERY CLOUD

Over 30 Offerings Across USA and China

	K520	K80	P40	M60	P4	P100	T4	V100	NGC
 Alibaba Cloud					●	●		●	●
 AWS	●	●		●				●	●
 Baidu Cloud			●		●		●		
 Google Cloud		●			●	●	●	●	●
 IBM Cloud		●		●		●		●	
 Microsoft Azure		●	●	●		●		●	●
 Oracle Cloud						●		●	●
 Tencent Cloud			●		●				

# SUPERCHARGED INSTANCES FOR HPC

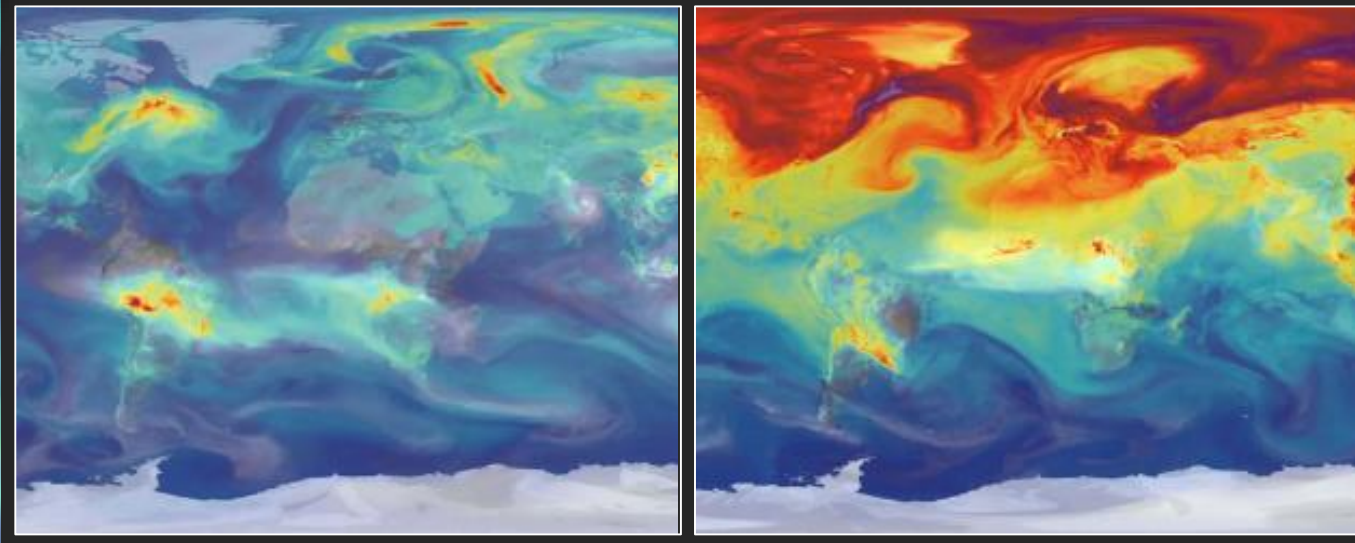
## Universal Platform Drives Utilization and Productivity



# AN AI MONITOR OF EARTH'S VITALS

NASA Ames uses satellite imagery to measure the effects of carbon and greenhouse gas emissions on the planet. They developed DeepSat—a deep learning framework for satellite image classification trained on a GPU-powered supercomputer. The enhanced satellite imagery will help scientists plan how to protect ecosystems and improve crop production.

 NVIDIA.



NASA: Late summer 2016, forest fires in Africa produce plumes of CO<sub>2</sub>  
Left: CO<sub>2</sub> - 10/14/2016 / Right: CO<sub>2</sub> - 12/24/2016

Source: [https://climate.nasa.gov/climate\\_resources/142/](https://climate.nasa.gov/climate_resources/142/)

# EXAMPLE: OILFIELD PREDICTIVE MAINTENANCE

Baker Hughes GE

60K oil wells WW equipped with Electric Submersible Pumps (ESP)

Average Non Producing Time (NPT) due to ESP Failure costs > \$150K per day per well

ML techniques used historically (Rule based, Fuzzy logic, traditional ML), but they don't scale

>50% False Alarms + detect failures too late

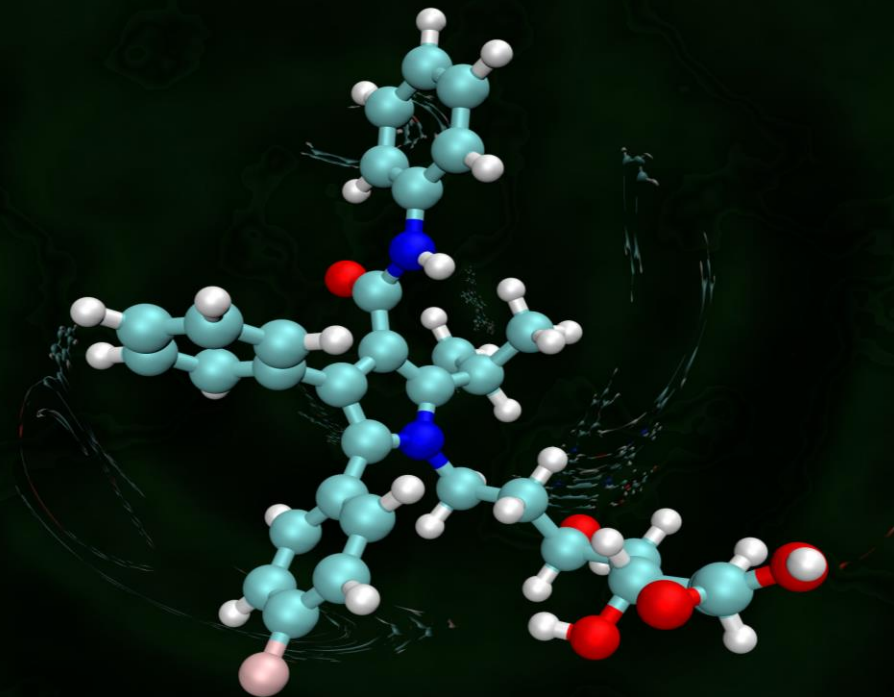
**Deep Learning: 93% detection accuracy with 2 months lead time at 5% False Alarm**

**At 10% DL based anomaly detection yields \$300K per well of lost productivity annually**

# CHASING $10^{60}$ CHEMICAL COMPOUNDS

Identifying molecules with desirable chemical properties is central to many industries. In the chemical space of  $10^{60}$  conceivable compounds, only  $10^8$  have been synthesized. Screening even a small fraction of the remaining compounds with legacy methods would take 100 node-seconds per compound.

Researchers at Dow are using GPU-powered deep learning to deliver completely novel molecular structures with specific properties. The AI produced 3M promising chemical leads in 1 day on an NVIDIA DGX-1.



# HUNTING GHOST PARTICLES WITH DEEP LEARNING

Understanding the properties of Neutrinos is the focus of a world-wide campaign. Observing these 'ghost particles' requires instruments of incredible size and scale. Fermilab's NOvA experiment applies two detectors snapping two million photos/second and analyzing them for neutrino activity.

Fermilab developed deep neural networks trained on NVIDIA GPUs and improved detection rate by 33% – increasing the discovery potential of NOvA and other experiments probing fundamental questions of the universe.



# SMARTER INSPECTION SERVICES

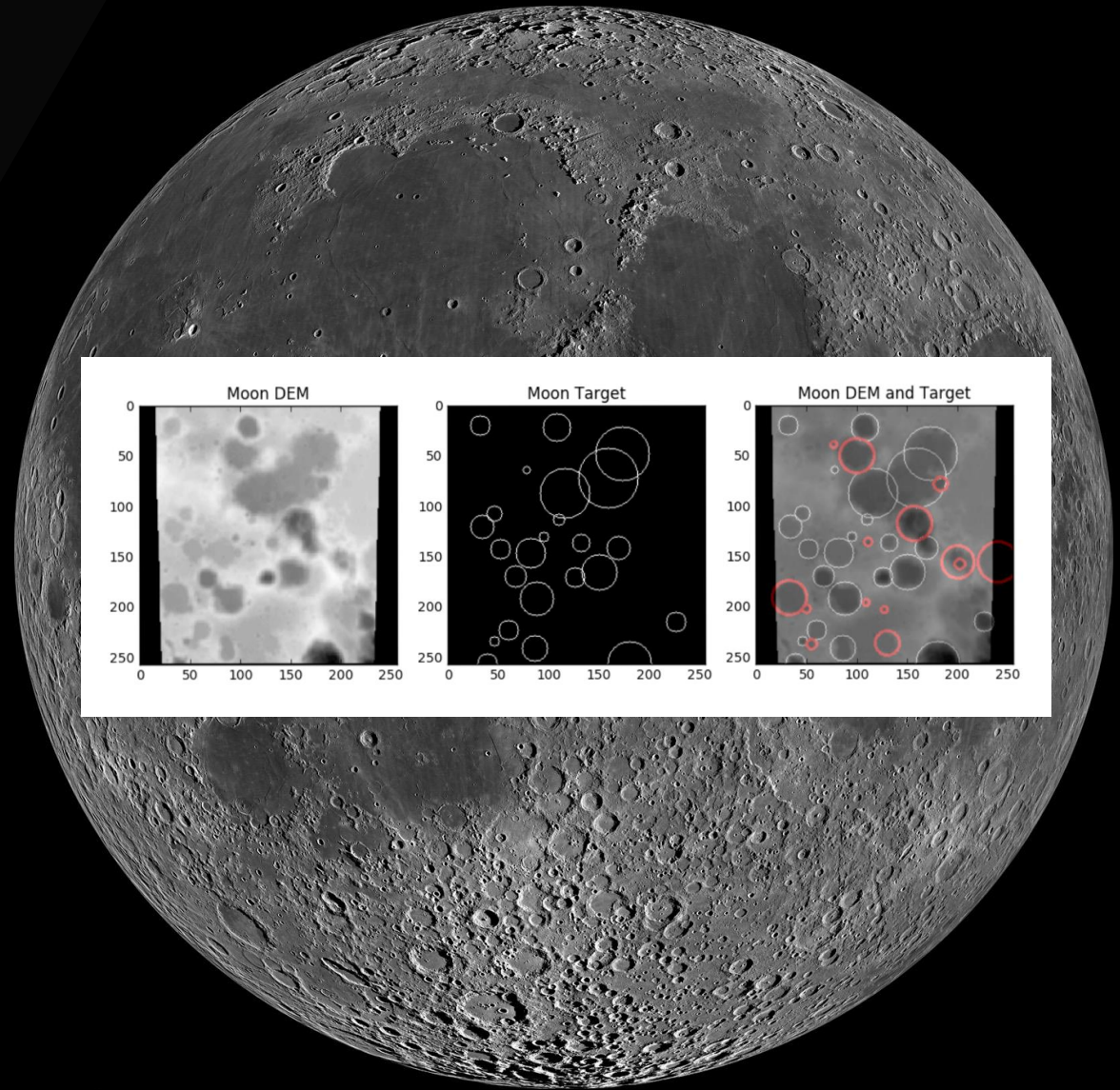
In business, ensuring equipment uptime and meeting safety and regulatory compliance is non-negotiable. Using deep neural networks developed on NVIDIA DGX-1 in the data center that can easily extend to NVIDIA DGX Station in the field, Avitas Systems delivers inspection services using robotic-based autonomous inspection and advanced data analytics. In addition to safeguarding workers, Avitas Systems AI solutions can reduce inspection costs by 25% and reduce maintenance downtime by 15%.



*The robots can handle the heat and use infrared cameras and chemical and other sensing technologies to inspect assets under dangerous conditions and keep production running. Image credit: Avitas Systems*

# MAPPING MOON CRATERS

Studying moon craters provides insight into the history of our solar system. Researchers at U of T and Penn State developed a CNN, powered by NVIDIA Tesla P100 GPUs on the SciNet P8 supercomputer, that automatically detects and classifies characteristics of craters from lunar digital elevation map images. Upon implementation the system identified 6,000 new craters in just a few hours, making it orders of magnitude faster than human counting.

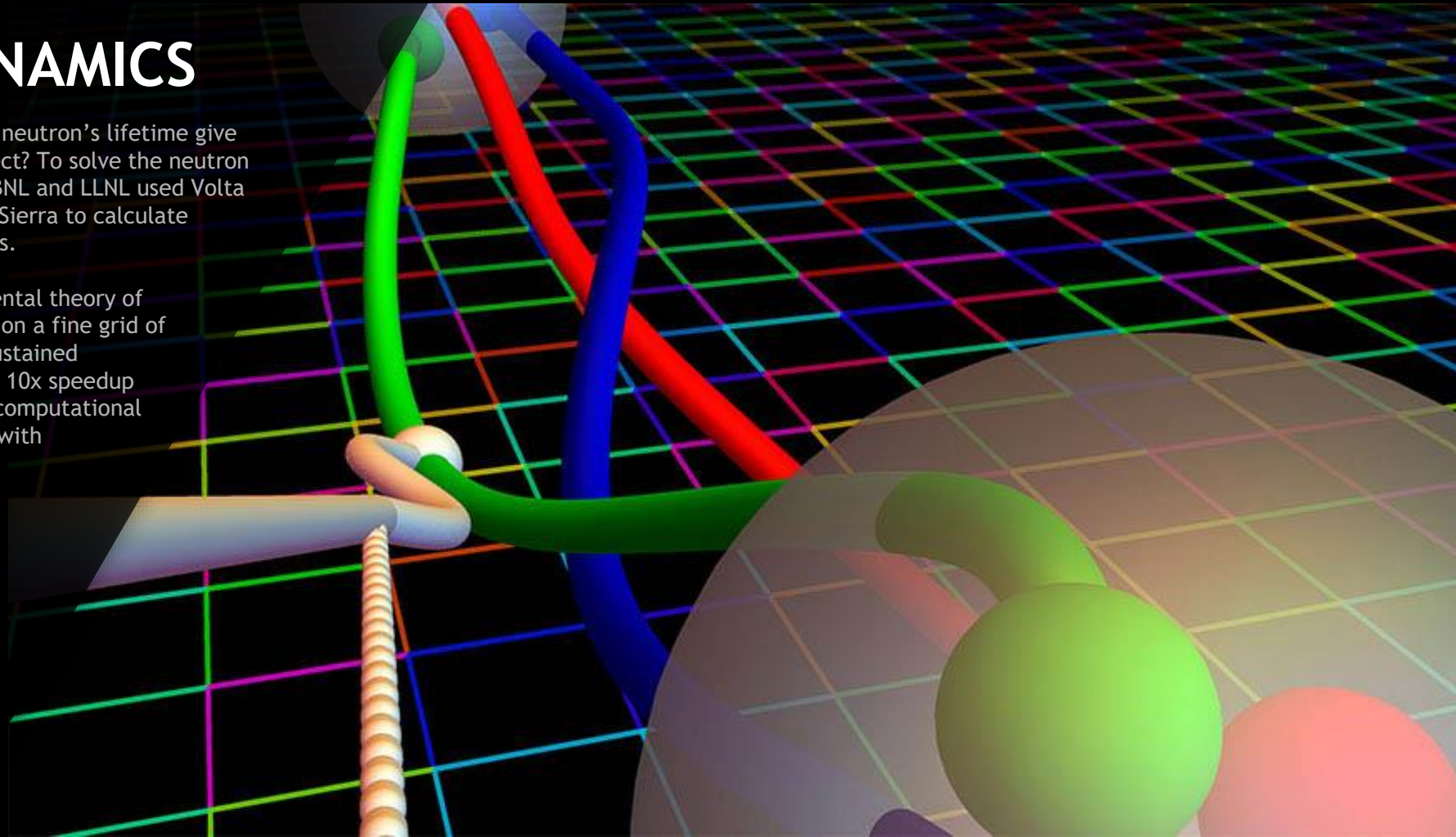


Inset: Sample Moon image (left) and target (center) from the dataset, with the two overlaid (right). Red circles show craters detected by the AI system that are absent from previous datasets.

# QUANTUM CHROMODYNAMICS

The two methods of measuring a neutron's lifetime give different answers. Which is correct? To solve the neutron lifetime puzzle scientists from LBNL and LLNL used Volta GPU-based nodes of Summit and Sierra to calculate the physics of subatomic particles.

The team simulated the fundamental theory of quantum chromodynamics (QCD) on a fine grid of space-time points. Achieving a sustained performance of ~20 petaflops – a 10x speedup over previous-gen systems – the computational breakthrough supplies physicists with the compute power to search for new physics.



Pictured: beta decay, the decay of a neutron ( $n$ ) to a proton ( $p$ ) with the emission of an electron ( $e$ ) and an electron-anti-neutrino ( $\bar{\nu}$ ). In the figure  $g_A$  is depicted as the white node on the red line. The square grid indicates the lattice. Image credit: Evan Berkowitz/Forschungszentrum Jülich/Institut für Kernphysik /Institute for Advanced Simulation

# “SEEING” GRAVITY FOR THE FIRST TIME

In September 2015, 100 years after Einstein predicted them, gravitational waves were observed for the first time. Astronomers at the Laser Interferometer Gravitational-wave Observatory have since used GPU-powered deep learning to process gravitational wave data 100x faster than previous methods, making real-time analysis possible and putting us one step closer to understanding the universe’s oldest secrets.



[Physics Letters B - Deep learning for real-time gravitational wave detection and parameter estimation: Results with advanced LIGO data](#)

*Daniel George, E.A. Huerta*

# OPTIMIZING QUALITY INSPECTIONS WITH AI

Due to the increasingly sophisticated design of its cars and the high-quality standards at Audi, the company inspects all components —doors, engine hoods, fenders, etc.— in its press shop.

In addition to visual inspection by Audi employees, several small deep learning-based cameras trained on NVIDIA GPUs installed directly in the presses detect the finest cracks in sheet metal with the utmost precision in a matter of seconds.



# AN EYE FOR CRISIS MANAGEMENT

Natural disasters are increasingly causing major destruction to life, property and economies. DFKI is using the NVIDIA DGX-2 to evolve DeepEye – which uses satellite images enriched with social media content to identify the effects of natural disasters– into a crisis management solution. With the increased GPU memory and fully connected GPUs based on the NVSwitch architecture, DFKI can build bigger models and process more data to aid rescuers in their decision-making for faster, more efficient dispatching of resources.



Geographical Information

Countries



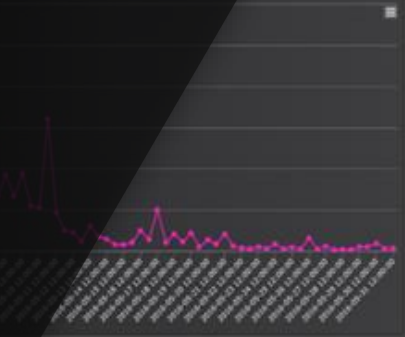
Languages

Countries

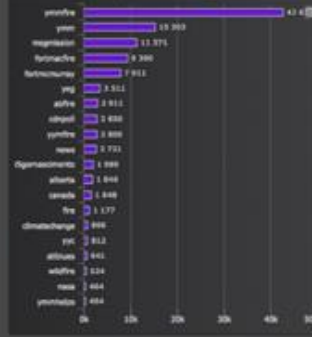
24

195

Word-Timing



Top #Hashtags



Top Tweets

Test

... of Fort McMurray fire incident, rendering services for humanity?

... Murray Fire. Gov't is matching all donations to: [https://www.govt.ca/](#)

... by the fire in Fort McMurray tonight. Stay safe and remember to [https://www.govt.ca/](#)

... Alberta history underway in Fort McMurray [https://www.govt.ca/](#)

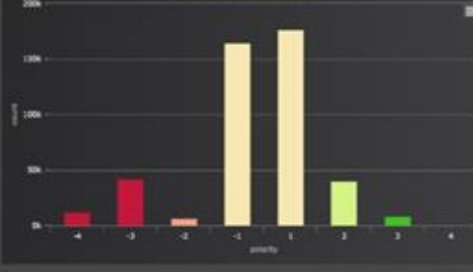
... Please donate to the Red Cross if you can. [https://www.govt.ca/](#)

... in Alberta is bearing down on Canada's Oilfields capital [https://www.govt.ca/](#)

... the road in FMM. Posted by @burg15. LIVE [https://www.govt.ca/](#)

... affected by the Fort McMurray fire - donate on item from this list [https://www.govt.ca/](#)

Sentiment Score



Top Users

802	Columns of smoke continue to rise up from the myriad of fires in Alberta, Canada. Details: <a href="#">https://www.cbc.ca/</a>	98	WMS
807	Authorities call for the evacuation of a city in Alberta, Canada, amid massive wildfire <a href="#">https://www.cbc.ca/</a>	99	Huffington Post
744	Heat Fuels Fire at Fort McMurray <a href="#">https://www.cbc.ca/</a> #NAGA <a href="#">https://www.cbc.ca/</a>	99	CNN
742	Canada's Syrian refugees raise money to help Fort McMurray fire victims. <a href="#">https://www.cbc.ca/</a>	99	The Economist
801	A very useful up-to-the-minute map of Fort McMurray, using space-based images <a href="#">https://www.cbc.ca/</a>	99	Wired
		99	ABC News

Top-Images in Tweets



Image duplicates

#Duplicates	5
#Occurrences of Image URL	6
#Duplicates	5
#Occurrences of Image URL	27
#Duplicates	5
#Occurrences of Image URL	29

Adjective-Noun-Pairs



Top noun-adjective combinations



Top Adjective-Noun-Pairs

Adjective	Noun Pair	Count	Adjective	Noun Pair	Count
bad	storm	350	violent	crime	191
incredible	summit	328	extreme	fire	178
awesome	summit	298	bad	accident	170
crazy	storm	286	amazing	scene	154
colorful	clouds	213	stormy	clouds	142
crazy	fire	211	violent	protest	127
horizontal	leaf	201	successful	business	108
wired	clouds	198	stunning	summit	108
crazy	clouds	195	dirty	window	106

# IMPROVING DEMAND FORECASTS

With >100,000 different products in its 4,700 U.S. stores, the Walmart Labs data science team predicts demand for 500 million item-by-store combinations every week.

By performing forecasting with the open-source RAPIDS data processing and machine learning libraries built on CUDA-X AI on NVIDIA GPUs, Walmart speeds up feature engineering 100x and trains machine learning algorithms 20x faster, resulting in faster delivery of products, real-time reaction to shopper trends, and inventory cost savings at scale.





**THANK YOU**

Reynaldo Gomez  
[reynaldog@nvidia.com](mailto:reynaldog@nvidia.com)