**NVIDIA**

# TRACKING: A PERFECT USE CASE FOR GPUS?

Peter Messmer, 7/1/2019

TrackML Grand Finale, CERN

# THANK YOU FOR THE PARTICIPATION!



kaggle    Search    **Competitions    Datasets    Kernels    Discussion    Courses    •••**

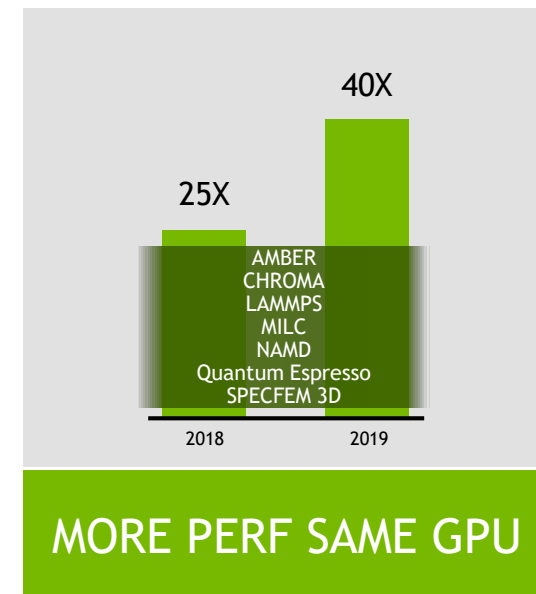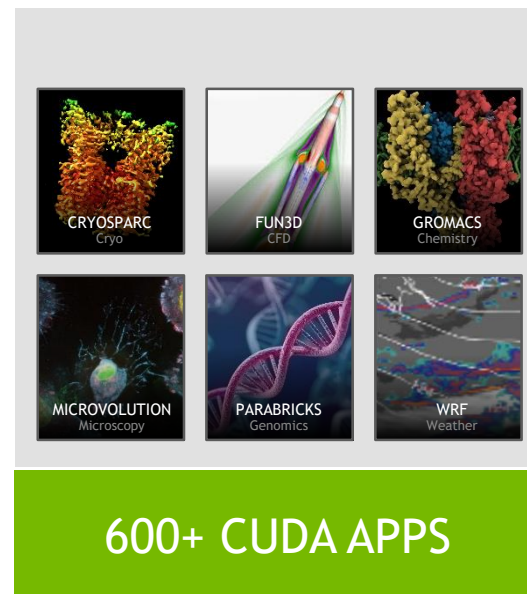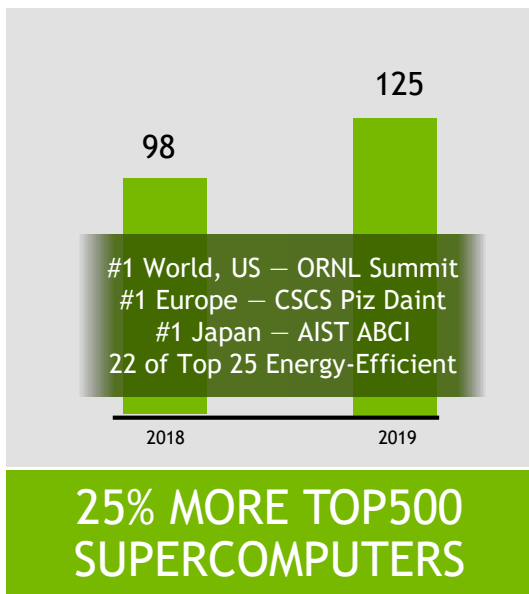Featured Prediction Competition

## TrackML Particle Tracking Challenge
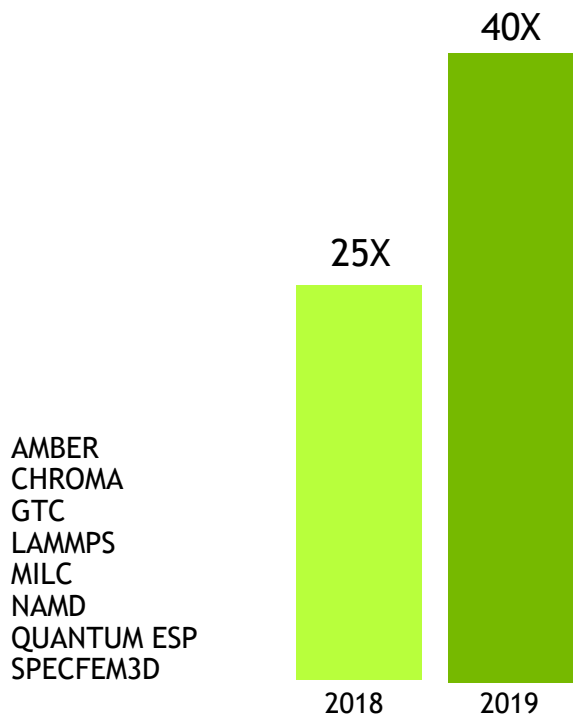High Energy Physics particle tracking in CERN detectors

CERN · 653 teams · a year ago

# A YEAR OF RAPID GROWTH

**125**

**98**

#1 World, US — ORNL Summit
#1 Europe — CSCS Piz Daint
#1 Japan — AIST ABCI
22 of Top 25 Energy-Efficient

2018    2019

## 25% MORE TOP500 SUPERCOMPUTERS

800K

1.2M
DEVELOPERS
+50%

2018    2019

8M

13M
CUDA
DOWNLOADS
+60%

2018    2019

## 50% GROWTH OF NVIDIA DEVELOPERS

CRYOSPARC
Cryo

FUN3D
CFD

GROMACS
Chemistry

MICROVOLUTION
Microscopy

PARABRICKS
Genomics

WRF
Weather

## 600+ CUDA APPS

**40X**

**25X**

AMBER
CHROMA
LAMMPS
MILC
NAMD
Quantum Espresso
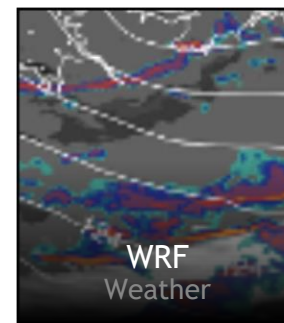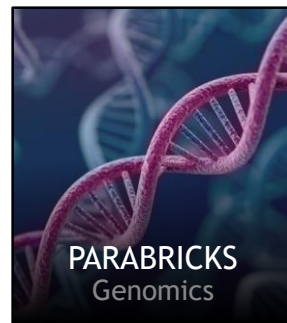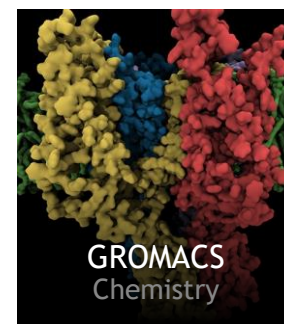SPECFEM 3D
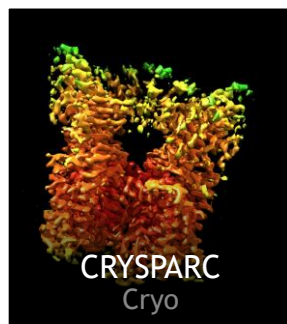
2018    2019

## MORE PERF SAME GPU

# EXPANDING VALUE FOR HPC CUSTOMERS

## Partnering With HPC Development Community

40X

25X

22X

AMBER
CHROMA
GTC
LAMMPS
MILC
NAMD
QUANTUM ESP
SPECFEM3D

CRYSPARC
Cryo

FUN3D
CFD

GROMACS
Chemistry

MICROVOLUTION
Microscopy

PARABRICKS
Genomics

WRF
Weather

| | |
|---|---|
| CRYOSPARC | 24x |
| FUN3D | 24x |
| GROMACS | 7x |
| MICROVOLUTION | 48x |
| PARABRICKS | 22x |
| WRF | 8x |

2018    2019

2019

MORE PERFORMANCE WITH SAME GPU

ADDING NEW AND IMPROVED TOP APPLICATIONS

*CPU Server: Dual Xeon Gold 6140@2.30GHz, GPU Servers: same CPU server w/ 4 NVIDIA V100 PCIe or SXM2 GPUs*

# ANNOUNCING CUDA TO ARM

## ENERGY-EFFICIENT SUPERCOMPUTING

NVIDIA GPU Accelerated Computing Platform On ARM

Optimized CUDA-X HPC & AI Software Stack

CUDA, Development Tools and Compilers

Available End of 2019

# INTERSECTION OF HPC & AI TRANSFORMING SCIENCE

## HPC

- > Algorithms based on first principles theory
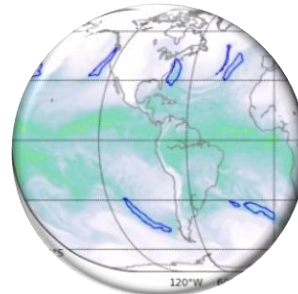- > Proven models for accurate results

## AI

- > Neural networks that learn patterns from large data sets
- > Improve predictive accuracy and faster response time

**SPEEDING PATH TO FUSION ENERGY**

PRINCETON UNIVERSITY

**90% Prediction Accuracy Publish in Nature April 2019**
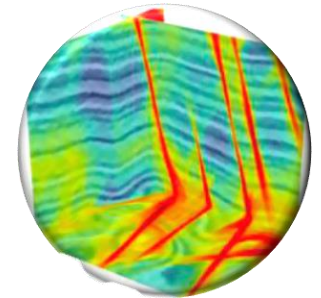
**EXASCALE WEATHER MODELING**

BERKELEY LAB

**Tensor Cores Achieved 1.13 EF 2018 Gordon Bell Winner**

**INDENTIFYING CHEMICAL COMPOUNDS**

DOW
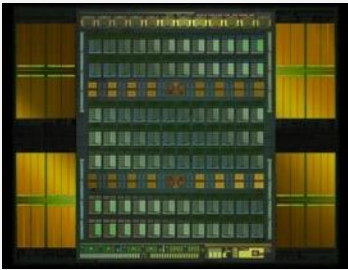
**Orders Of Magnitude Speedup 3M New Compounds In 1 Day**

**O&G FAULT INTERPRETATION**

BUREAU OF ECONOMIC GEOLOGY

**Time-to-solution Reduced From Weeks To 2 Hours**

# INTRODUCING TESLA V100

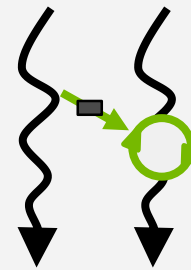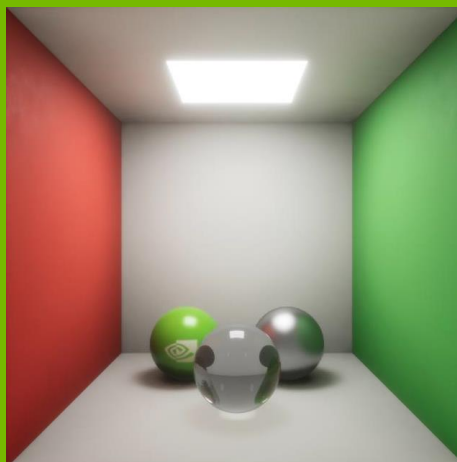| Volta Architecture | Improved NVLink & HBM2 | Volta MPS | Improved SIMT Model | Tensor Core |
|---|---|---|---|---|
| Most Productive GPU | Efficient Bandwidth | Inference Utilization | New Algorithms | 120 Programmable TFLOPS Deep Learning |

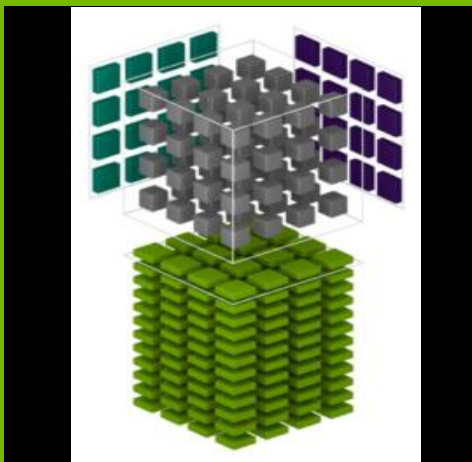## The Fastest and Most Productive GPU for Deep Learning and HPC

# NVIDIA TURING: GRAPHICS REINVENTED

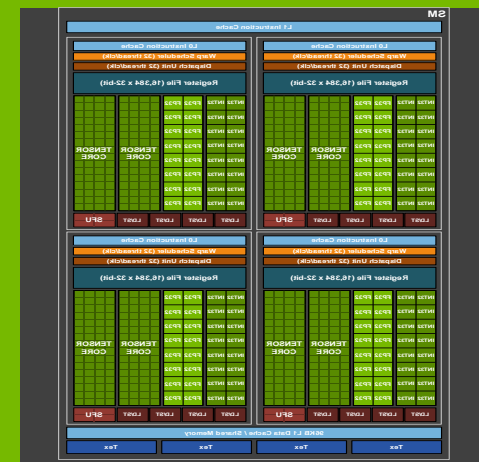## Built to Revolutionize the Work of Creative Professionals

RT Cores

Tensor Cores

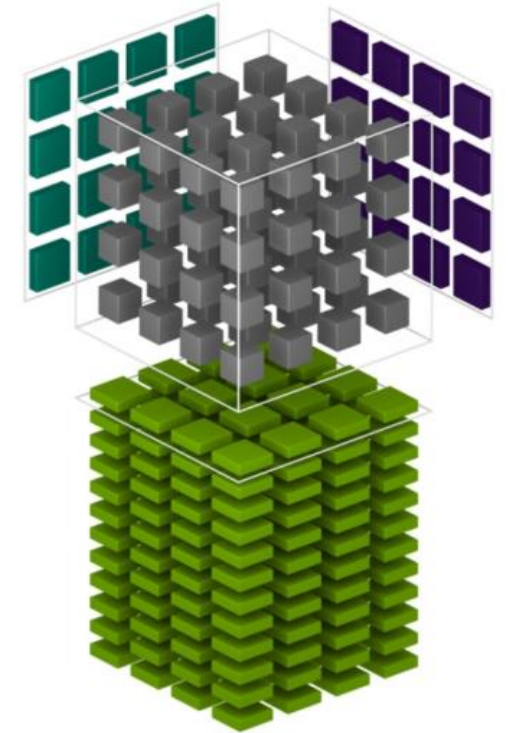CUDA Cores

Streaming Multiprocessors

# TENSOR CORES IN TURING

## New in Volta, Upgraded in Turing

| GPU | SMs | Total | Peak Half FLOPS | PEAK INT8 OPS | PEAK INT4 OPS | PEAK B1 OPS |
|---|---|---|---|---|---|---|
| V100 | 80 | 640 | | N.A. | N.A. | N.A. |
| Quadro RTX 6000/8000 | 72 | 576 | 130.5 TFLOPS* | 260 TOPS* | 521 TOPS* | 2087 TOPS* |

Matrix Multiplication Pipeline,
half precision inputs → half / float accumulator
8bit/4bit INT inputs → 32bit INT accumulator
1bit Binary inputs → 32 bit INT accumulator (XOR + POPC)

Used in CUBLAS, CUDNN, CUTLASS
Exposed in CUDA 10 (4bit INT and 1bit binary are experimental)

* Using 1.77GHz Boost Clock
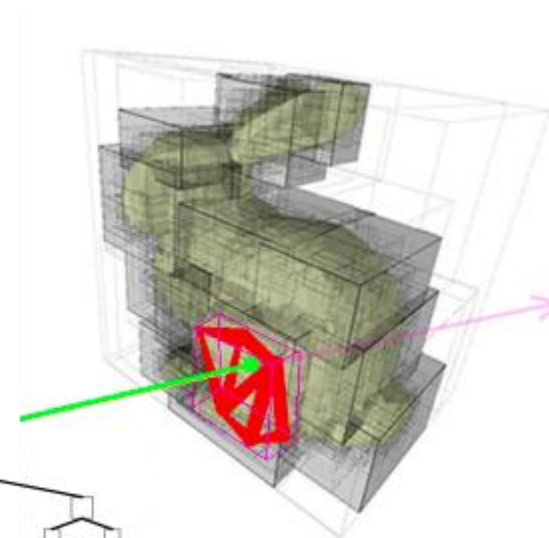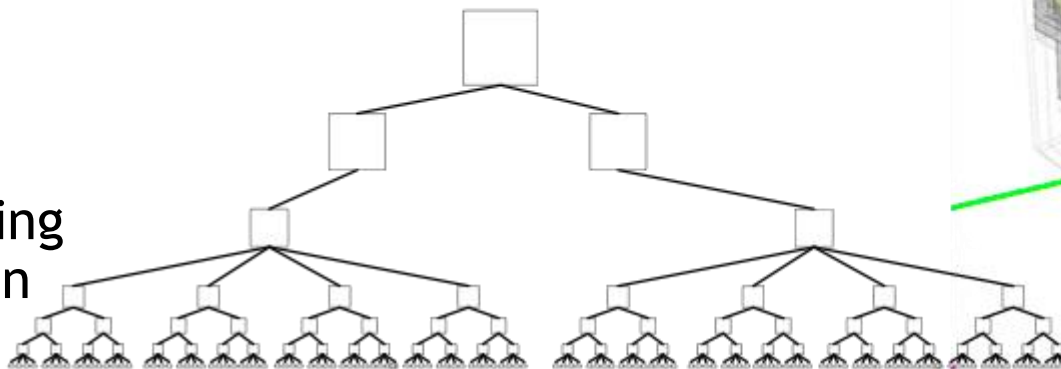
nVIDIA

# TURING RT CORES

## Hardware Accelerated Ray Tracing

RT Cores perform

- Ray-BVH Traversal
- Instancing:  1 Level
- Ray-Triangle Intersection

Return to SM for

- Multi-level Instancing
- Custom Intersection
- Shading

Programming via OptiX RT framework
Low overhead interop with CUDA

# HOW TO START WITH GPUS

| **1** Applications | | |
|---|---|---|
| **2** Libraries | **3** Compiler Directives | **4** Programming Languages |
| Easy to use | Easy to Start | Most Performance |
| Most Performance | Portable Code | Most Flexibility |
| | **OpenACC** | **CUDA** |

1. Review available GPU-accelerated applications

2. Check for GPU-Accelerated applications and libraries

3. Add OpenACC Directives for quick acceleration results and portability

4. Dive into CUDA for highest performance and flexibility

⊛ nVIDIA.

# SINGLE CODE FOR MULTIPLE PLATFORMS
## OpenACC - Performance Portable Programming Model for HPC
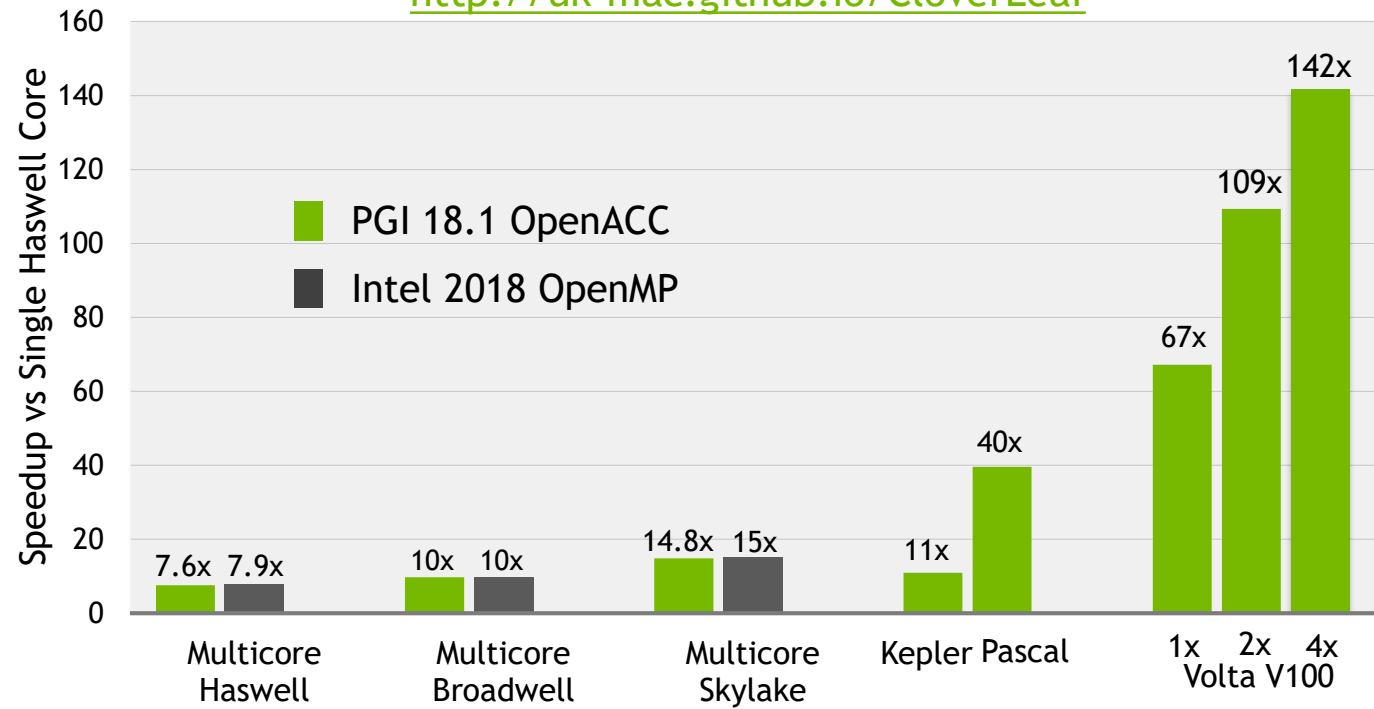
OpenPOWER

Sunway

x86 CPU

x86 Xeon Phi

NVIDIA GPU

AMD GPU

PEZY-SC

AWE Hydrodynamics CloverLeaf mini-App, bm32 data set
http://uk-mac.github.io/CloverLeaf



Speedup vs Single Haswell Core

- PGI 18.1 OpenACC
- Intel 2018 OpenMP

Multicore Haswell: 7.6x, 7.9x
Multicore Broadwell: 10x, 10x
Multicore Skylake: 14.8x, 15x
Kepler Pascal: 11x, 40x
Volta V100: 1x 67x, 2x 109x, 4x 142x

NVIDIA.

# NVIDIA CUDA-X UPDATES

## Software To Deliver Acceleration For HPC & AI Apps; 500+ New Updates

| Machine Learning & Deep Learning | Computational Physics & Chemistry | Computational Fluid Dynamics | Life Sciences & Bioinformatics | Structural Mechanics | Weather & Climate | Geoscience, Seismology & Imaging | Numerical Analytics | Electronic Design Automation |
|---|---|---|---|---|---|---|---|---|



**600+ Apps**

| Linear Algebra | Parallel Algorithms | Signal Processing | Deep Learning | Machine Learning | Visualization |
|---|---|---|---|---|---|

**CUDA-X HPC & AI**

**40+ GPU Acceleration Libraries**

**CUDA**

| Desktop Development | Data Center | Supercomputers | GPU-Accelerated Cloud |
|---|---|---|---|

# RAPIDS — OPEN GPU DATA SCIENCE

Software Stack

# CUML — OPEN GPU DATA SCIENCE
## Broad range of GPU accelerated algorithms



https://rapids.ai/
https://github.com/rapidsai/cuml

# GPU-ACCELERATED XGBOOST

## Unleashing the Power of NVIDIA GPUs for Users of XGBoost

**Faster Time To Insight**
XGBoost training on GPUs is significantly faster than CPUs, completely transforming the timescales of machine learning workflows.

**Better Predictions, Sooner**
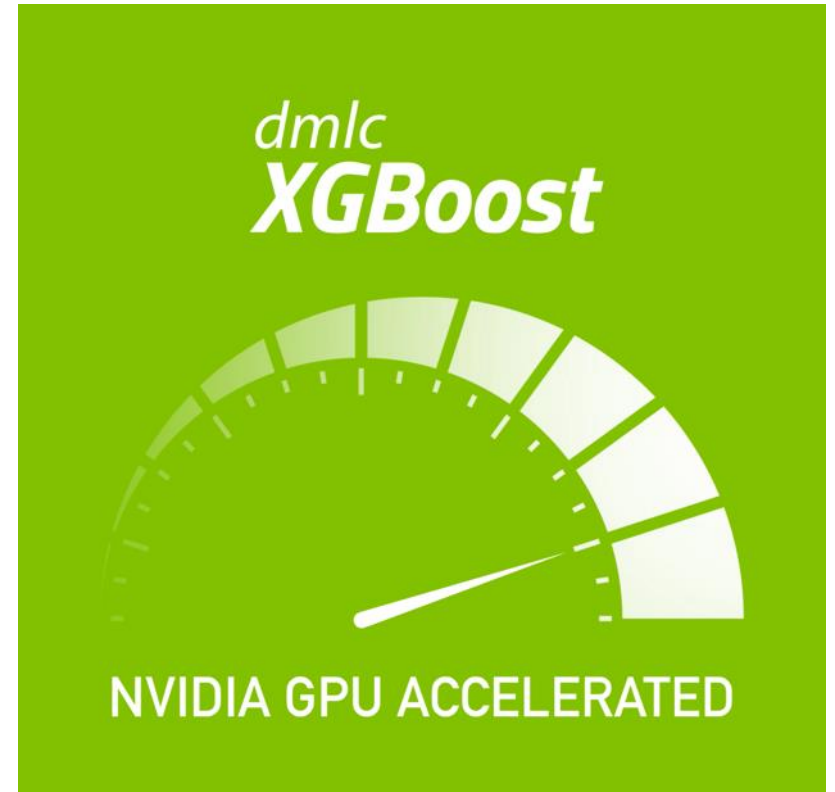Work with larger datasets and perform more model iterations without spending valuable time waiting.

**Lower Costs**
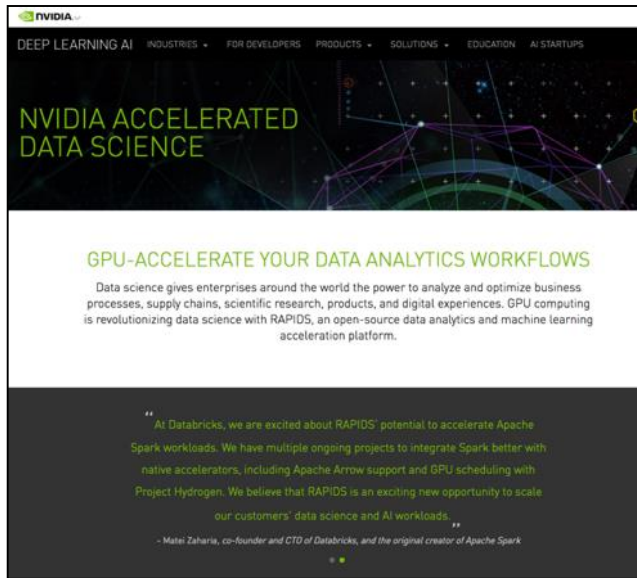Reduce infrastructure investment and save money with improved business forecasting.

**Easy to Use**
Works seamlessly with the RAPIDS open source data processing and machine learning libraries and ecosystem for end-to-end GPU-accelerated workflows.



dmlc
**XGBoost**

**NVIDIA GPU ACCELERATED**

# FOR MORE INFORMATION



nvidia.com/datascience

rapids.ai

nvidia.com/en-us/technologies/cuda-x/

# GPUS ARE SPOT ON FOR TRACKING WITH ML!

Convergence of scientific computing and ML/AI is happening now

GPUs: Accelerator for both HPC & AI/ML

Pick the right programming model: CUDA, OpenACC, libraries, frameworks

Frameworks increasingly Python driven

DNN frameworks fully supported on GPU

RAPIDS: broad framework for machine learning applications


What's your story with GPUs?

# TESLA PRODUCTS DECODER

| | P100 (SXM2) | P100 (PCIE) | P40 | P4 | T4 | V100 (PCIE) | V100 (SXM2) | V100 (FHHL) |
|---|---|---|---|---|---|---|---|---|
| GPU CHIP | GP100 | GP100 | GP102 | GP104 | TU104 | GV100 | GV100 | GV100 |
| PEAK FP64 (TFLOPs) | 5.3 | 4.7 | NA | NA | NA | 7 | 7.8 | 6.5 |
| PEAK FP32 (TFLOPs) | 10.6 | 9.3 | 12 | 5.5 | 8.1 | 14 | 15.7 | 13 |
| PEAK FP16 (TFLOPs) | 21.2 | 18.7 | NA | NA | 65 | 112 | 125 | 105 |
| PEAK TOPs | NA | NA | 47 | 22 | 260 | NA | NA | NA |
| Memory Size | 16 GB HBM2 | 16/12 GB HBM2 | 24 GB GDDR5 | 8 GB GDDR5 | 16 GB GDDR6 | 32 GB HBM2 | 32 GB HBM2 | 16GB HBM2 |
| Memory BW | 732 GB/s | 732/549 GB/s | 346 GB/s | 192 GB/s | 320GB/s | 900 GB/s | 900 GB/s | 900 GB/s |
| Interconnect | NVLINK + PCIe Gen3 | PCIe Gen3 | PCIe Gen3 | PCIe Gen3 | PCIe Gen3 | PCIe Gen3 | NVLINK + PCIe Gen3 | PCIe Gen3 |
| ECC | Internal + HBM2 | Internal + HBM2 | GDDR5 | GDDR5 | GDDR6 | Internal + HBM2 | Internal + HBM2 | Internal + HBM2 |
| Form Factor | SXM2 | PCIE Dual Slot | PCIE Dual Slot | PCIE LP | PCIE LP | PCIE Dual Slot | SXM2 | PCIE Single Slot Full Height Half Length |
| Power | 300 W | 250 W | 250 W | 50-75 W | 70 W | 250W | 300W | 150W |