**NVIDIA.**

# ACCELERATED COMPUTING WITH NVIDIA GPUS

Jesse Tetreault, Solutions Architect
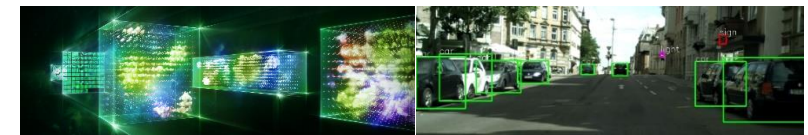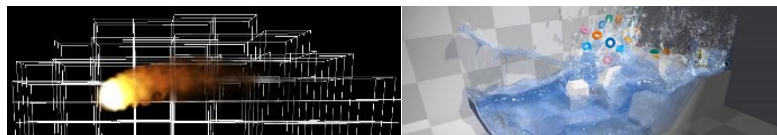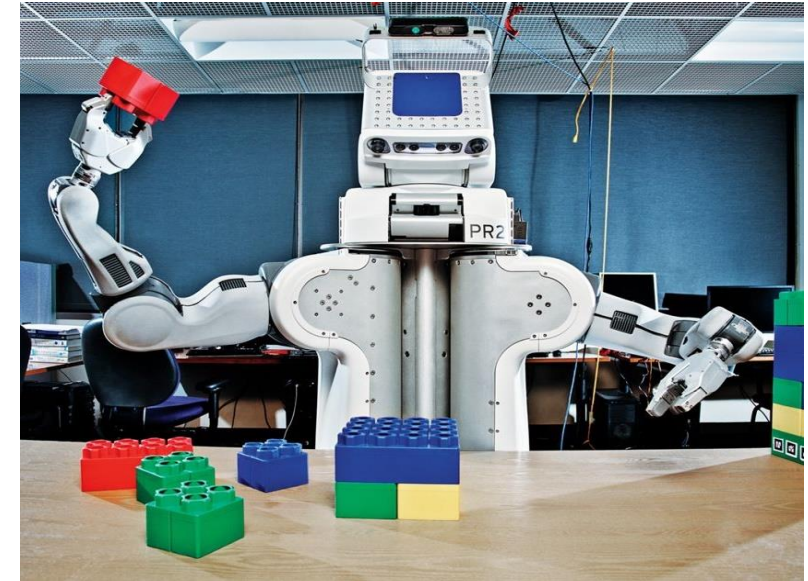
October 2019

# ACCELERATED COMPUTING
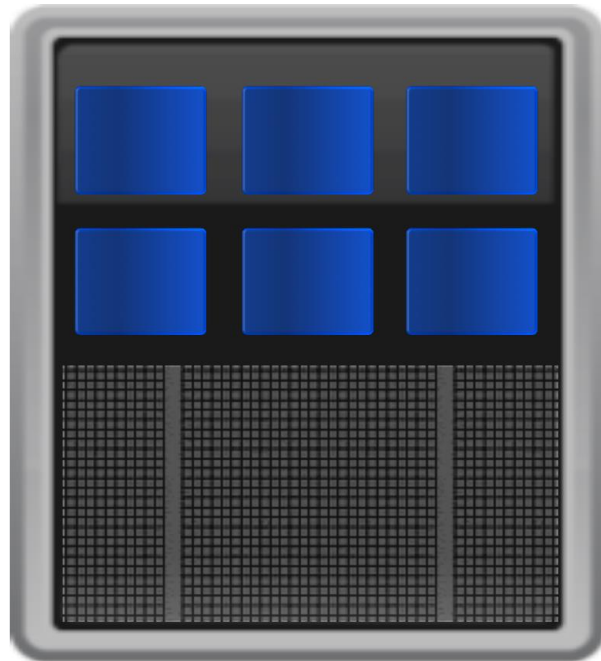
# NVIDIA
# "THE AI COMPUTING COMPANY"
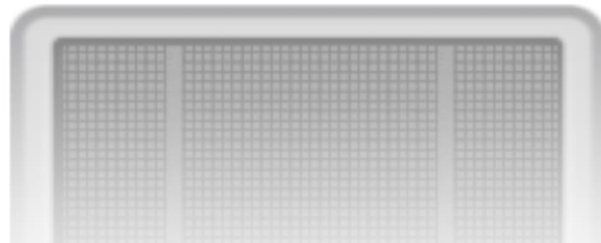


GPU Computing

Computer Graphics

Artificial Intelligence

# Add GPUs: Accelerate Science Applications

## CPU

## GPU

# HOW GPU ACCELERATION WORKS

**Application Code**

Compute-Intensive Functions

**5% of Code**

Rest of Sequential
CPU Code

**GPU**

**CPU**

+

# HOW TO START WITH GPUS

| 1 Applications | | |
|---|---|---|
| **2** Libraries | **3** Compiler Directives | **4** Programming Languages |
| Easy to use<br><br>Most Performance | Easy to Start<br><br>Portable Code | Most Performance<br><br>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 CUDA-X LIBRARIES

## 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 |
|---|---|---|---|

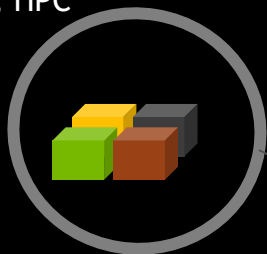# NVIDIA DEEP LEARNING SOFTWARE STACK

# NGC: GPU-OPTIMIZED SOFTWARE HUB

## Ready-to-run GPU Optimized Software, Anywhere

**50+ Containers**
DL, ML, HPC

**15+ Model Training Scripts**
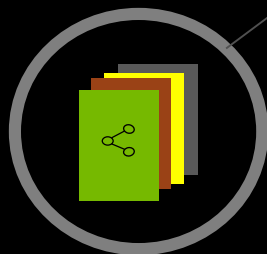NLP, Image Classification, Object Detection & more

**NGC**

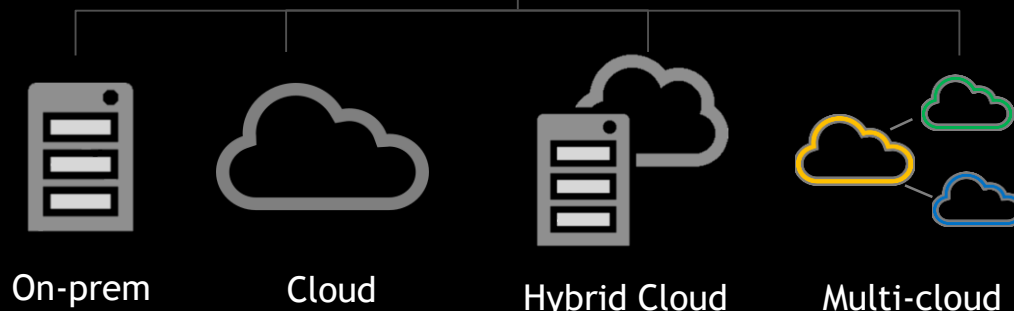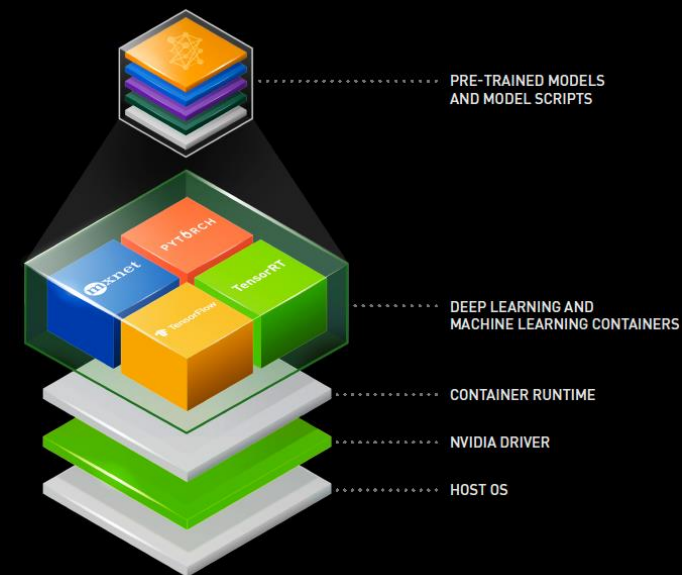**60 Pre-trained Models**
NLP, Image Classification, Object Detection & more

**Industry Workflows**
Medical Imaging, Intelligent Video Analytics

PRE-TRAINED MODELS AND MODEL SCRIPTS

DEEP LEARNING AND MACHINE LEARNING CONTAINERS

CONTAINER RUNTIME

NVIDIA DRIVER

HOST OS

On-prem          Cloud          Hybrid Cloud          Multi-cloud

# GPU-ACCELERATED DATA SCIENCE PLATFORMS

## Unparalleled Performance and Productivity



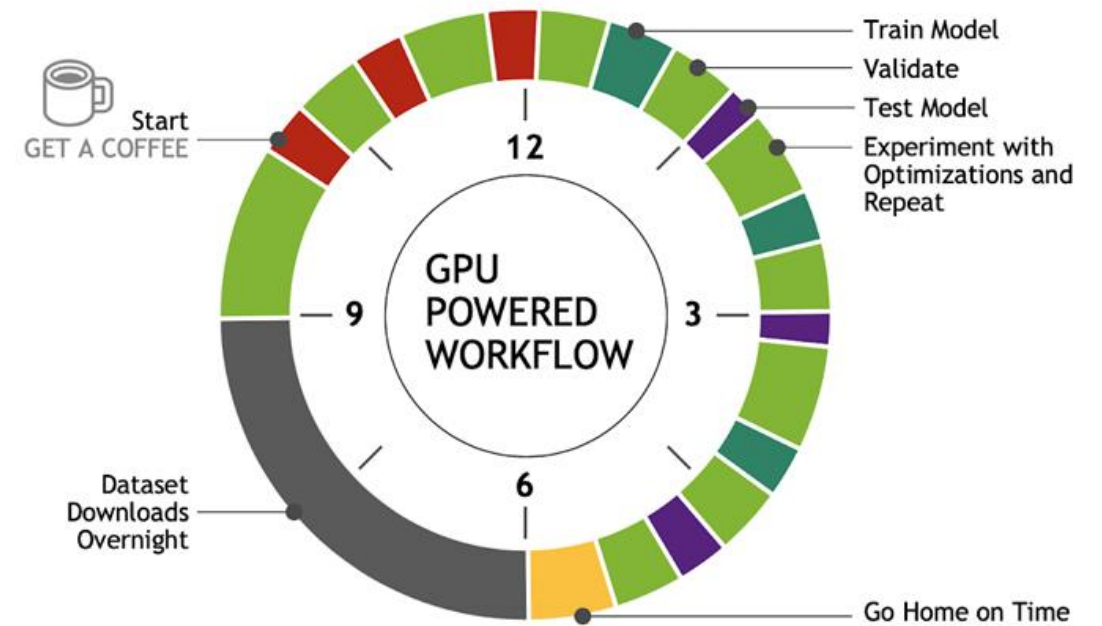| | **ML in the Cloud**<br>All the top CSPs | **ML Enthusiast**<br>High-end PCs | | **Enterprise Desktop**<br>Individual Workstations | **Enterprise Data Center**<br>Shared Infrastructure for Data Science Teams | | |
|---|---|---|---|---|---|---|---|
| | | | | | | **Max Flexibility** | **Max Performance** |
| | **NVIDIA GPUs in the Cloud** | **GeForce** | **TITAN RTX** | **NVIDIA-Powered Data Science Workstations** | **T4 Enterprise Servers** | **DGX Station, DGX-1 / HGX-1** | **DGX-2 / HGX-2** |
| **Benefit** | Ease of getting started, low/no barrier to entry, elasticity of resources | Enthusiast PC solution, easy to acquire, low cost, great performance | The ultimate PC GPU for data scientists. Easy to acquire, deploy and get started experimenting. | Enterprise workstation for experienced data scientists | Standard GPU-accelerated data center infrastructures with the world's leading servers | Enterprise server, proven 4 or 8-way configuration, modular approach for scale-up, fastest multi-GPU & multi-node training | Largest compute and memory capacity in a single node, fastest training solution |
| **Typical GPU Memory (system dependent)** | **varies depending on offering** | **22GB** | **48GB** | **96GB** | **64 GB** (4 x 16 GB) | **128GB-256GB** | **512GB** |
| **GPU Fabric** | varies depending on offering | 2-way NVLink | 2-way NVLink | 2-way NVLink | PCIe 3.0 | 4- and 8-way NVLink | 16-way NVSwitch |

# ACCELERATING DATA SCIENCE IN HEALTHCARE

# DAY IN THE LIFE OF A DATA SCIENTIST

# CHALLENGES IN DATA SCIENCE

# RAPIDS IN DATA SCIENCE



**Wrangle Data**

**Data Preparation**

**Train**

**Deploy**

Imaging

Genomic

Medical Records

Claims

Wearables

ETL

Data Lake

**cuDF**

Dataframe Manipulation
Feature Engineering

Train

Evaluate

**cuML**

Cross Validation
Hyperparameter Tuning

Inference

Performance in these two domains is typically a pain point for Data Scientists

# cuML Algorithms

| cuML | Single-GPU | Multi-GPU | Multi-Node-Multi-GPU |
|---|:---:|:---:|:---:|
| Gradient Boosted Decision Trees (GBDT) | ✓ | ✓ | ✓ |
| GLM | ✓ | ✓ | |
| Logistic Regression | ✓ | | |
| Random Forest | ✓ | ✓ | ✓ |
| K-Means | ✓ | ✓ | ✓ |
| K-NN | ✓ | ✓ | |
| DBSCAN | ✓ | | |
| UMAP | ✓ | | |
| Holt-Winters | ✓ | | |
| Kalman Filter | ✓ | | |
| t-SNE | ✓ | | |
| Principal Components | ✓ | | |
| Singular Value Decomposition | ✓ | ✓ | |

# Data Processing Evolution
## Faster data access, less data movement

Hadoop Processing, Reading from disk

| HDFS Read | Query | HDFS Write | HDFS Read | ETL | HDFS Write | HDFS Read | ML Train |

Spark In-Memory Processing

| HDFS Read | Query | ETL | ML Train |

**25-100x Improvement**
Less code
Language flexible
Primarily In-Memory

RAPIDS

| Arrow Read | Query | ETL | ML Train |

**50-100x Improvement**
Same code
Language flexible
Primarily on GPU

NVIDIA.

# Disk → Memory → GPUs



Scalable, but slow due to repeated reads & writes to disk

Faster, by keeping data always in host memory instead of on disk
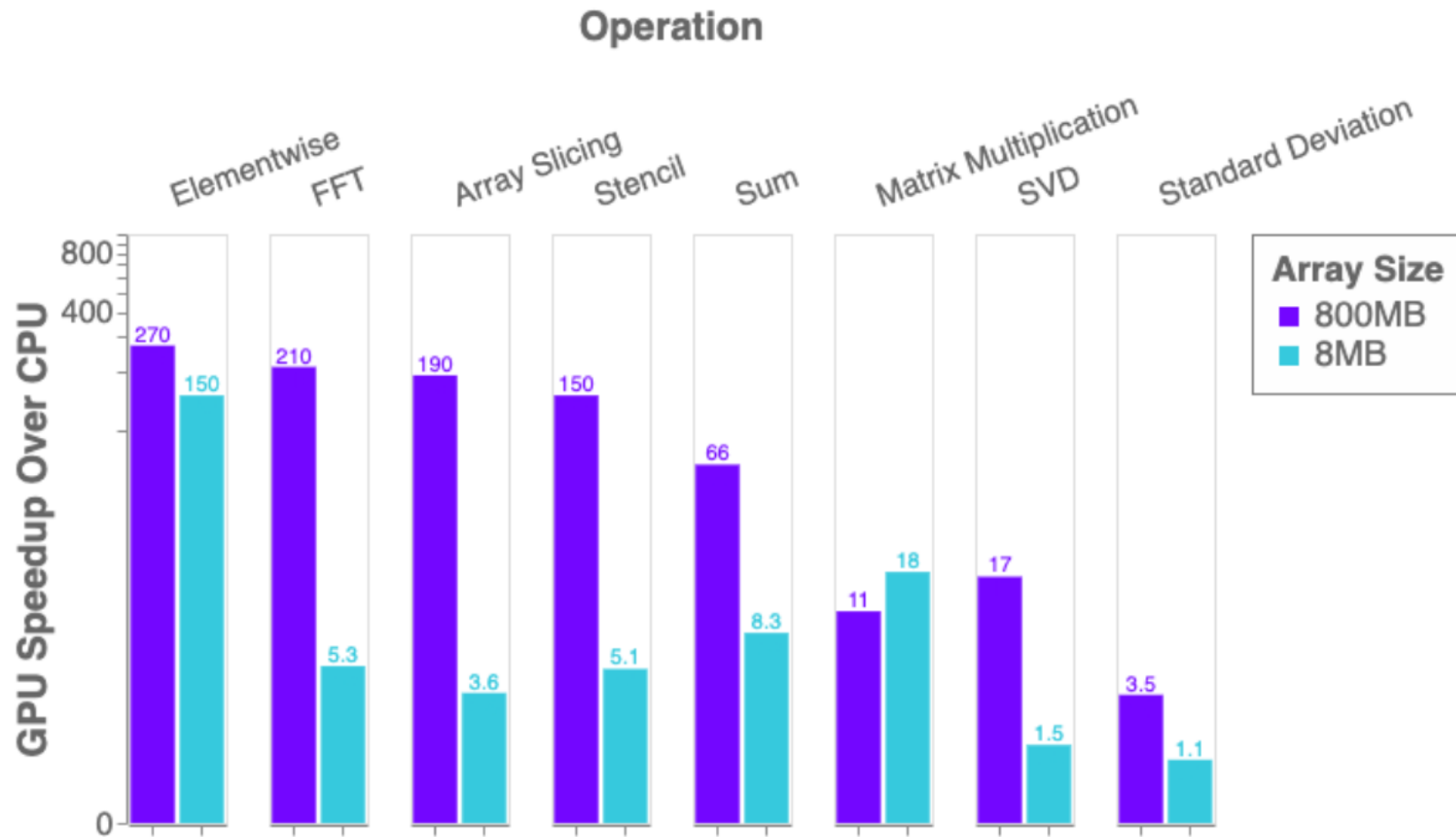
Performance limited by CPUs

Keeps data in GPU memory instead of CPU memory
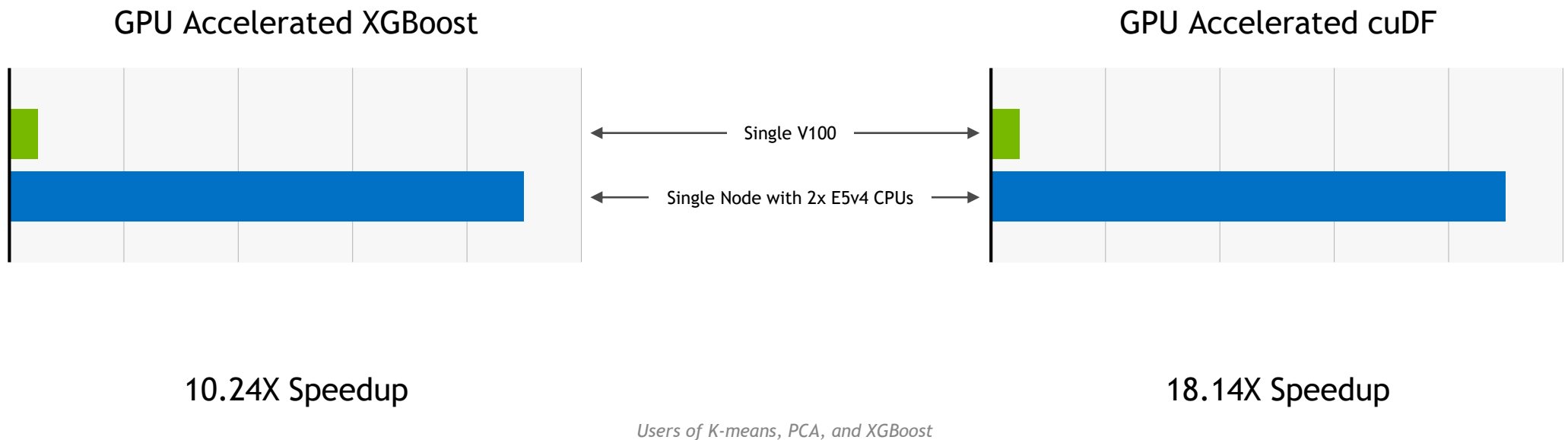
Computations are GPU accelerated

NVIDIA.

# Real Outcomes using Accelerated Machine Learning

# cuPy Acceleration

# TRANSFORM GENETICS WITH RAPIDS
## Personalize Immunotherapy for Cancer Patients

GPU Accelerated XGBoost

GPU Accelerated cuDF

Single V100

Single Node with 2x E5v4 CPUs

10.24X Speedup
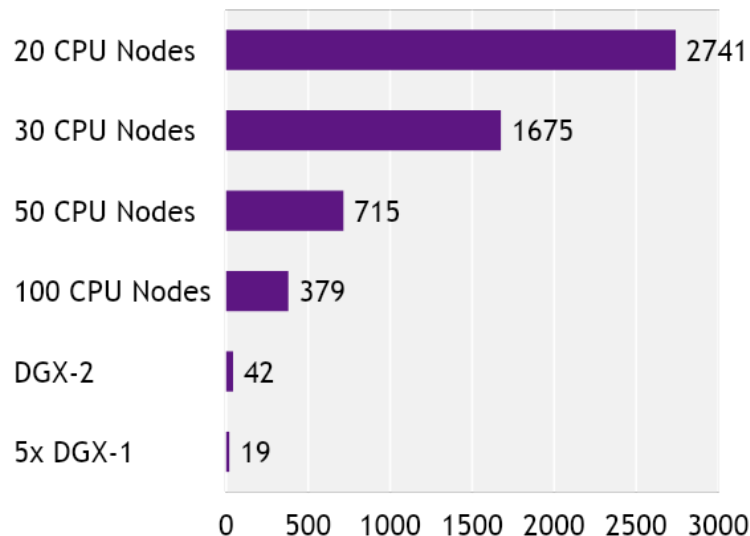
18.14X Speedup

*Users of K-means, PCA, and XGBoost*

*"We see close to 20x speedup using XGBoost on DGX-1. This helps us significantly improve our personalized immunotherapy and expand our analysis to millions of peptide candidates."*
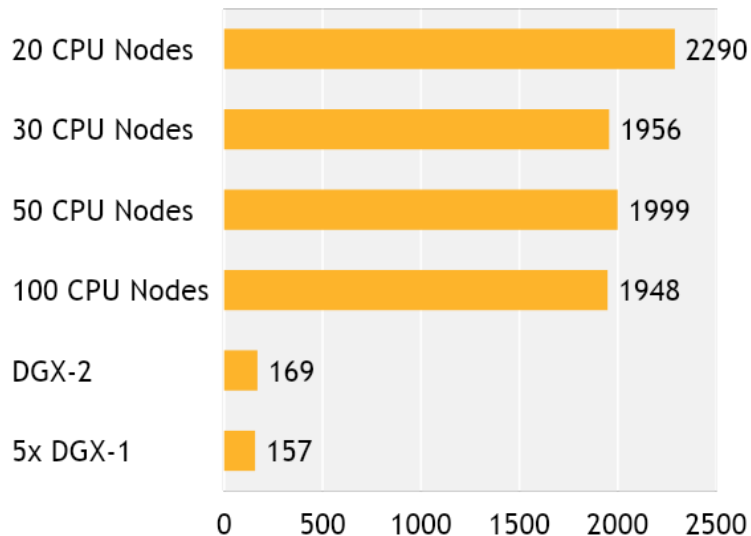
Yong Hou, Duty Director of BGI Research

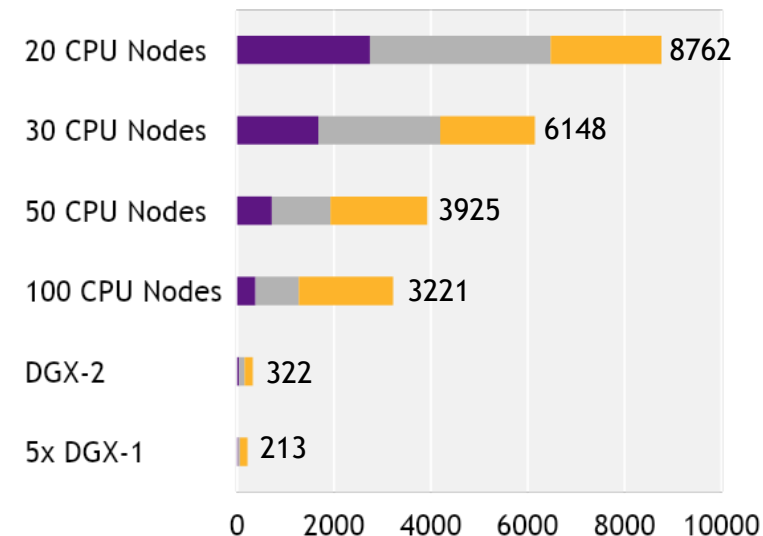华大基因
BGI

# Faster Speeds, Real-World Benefits

### cuIO/cuDF – Load and Data Preparation

| | Time (s) |
|---|---|
| 20 CPU Nodes | 2741 |
| 30 CPU Nodes | 1675 |
| 50 CPU Nodes | 715 |
| 100 CPU Nodes | 379 |
| DGX-2 | 42 |
| 5x DGX-1 | 19 |

### cuML - XGBoost

| | Time (s) |
|---|---|
| 20 CPU Nodes | 2290 |
| 30 CPU Nodes | 1956 |
| 50 CPU Nodes | 1999 |
| 100 CPU Nodes | 1948 |
| DGX-2 | 169 |
| 5x DGX-1 | 157 |

### End-to-End

| | Time (s) |
|---|---|
| 20 CPU Nodes | 8762 |
| 30 CPU Nodes | 6148 |
| 50 CPU Nodes | 3925 |
| 100 CPU Nodes | 3221 |
| DGX-2 | 322 |
| 5x DGX-1 | 213 |

**Time in seconds (shorter is better)**

■ cuIO/cuDF (Load and Data Prep)  ■ Data Conversion  ■ XGBoost

**Benchmark**

200GB CSV dataset; Data prep includes joins, variable transformations

**CPU Cluster Configuration**

CPU nodes (61 GiB memory, 8 vCPUs, 64-bit platform), Apache Spark

**DGX Cluster Configuration**

5x DGX-1 on InfiniBand network

NVIDIA.