



# RAPIDS PITCH DECK

# 6 QUESTIONS FACING EVERY AI ENTERPRISE

## Top Challenges for AI, Big Data, and Enterprise Transformation



### DATA DELUGE

Is your data doubling each year?



### PROLONGED TRAINING TIME

Is ML training prohibitively long, delaying time-to-predictions?



### COMPLEX WORKLOADS

Is Spark workloads creating relentless infrastructure sprawl?



### DELAYED INTELLIGENCE

Are you an intelligent enterprise needing real time predictive analytics?



### TEDIOUS DATA PREP

Do you have oceans of data, that take lifetimes to wrangle?



### SHRINKING BUDGET

Is your CAPEX budget shrinking amidst escalating infrastructure demand?

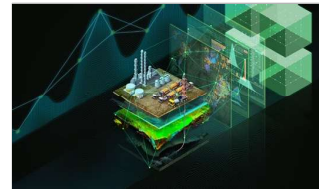
# GPU-ACCELERATED DATA SCIENCE

## Use Cases in Every Industry



### CONSUMER INTERNET

- Ad Personalization
- Click Through Rate Optimization
- Churn Reduction



### OIL & GAS

- Sensor Data Tag Mapping
- Anomaly Detection
- Robust Fault Prediction



### FINANCIAL SERVICES

- Claim fraud
- Customer service chatbots/routing
- Risk evaluation



### MANUFACTURING

- Remaining Useful Life Estimation
- Failure Prediction
- Demand Forecasting



### HEALTHCARE

- Improve Clinical Care
- Drive Operational Efficiency
- Speed Up Drug Discovery



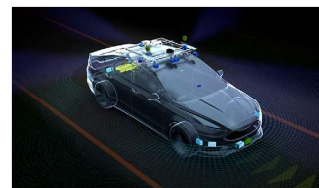
### TELCO

- Detect Network/Security Anomalies
- Forecasting Network Performance
- Network Resource Optimization (SON)



### RETAIL

- Supply Chain & Inventory Management
- Price Management / Markdown Optimization
- Promotion Prioritization And Ad Targeting



### AUTOMOTIVE

- Personalization & Intelligent Customer Interactions
- Connected Vehicle Predictive Maintenance
- Forecasting, Demand, & Capacity Planning

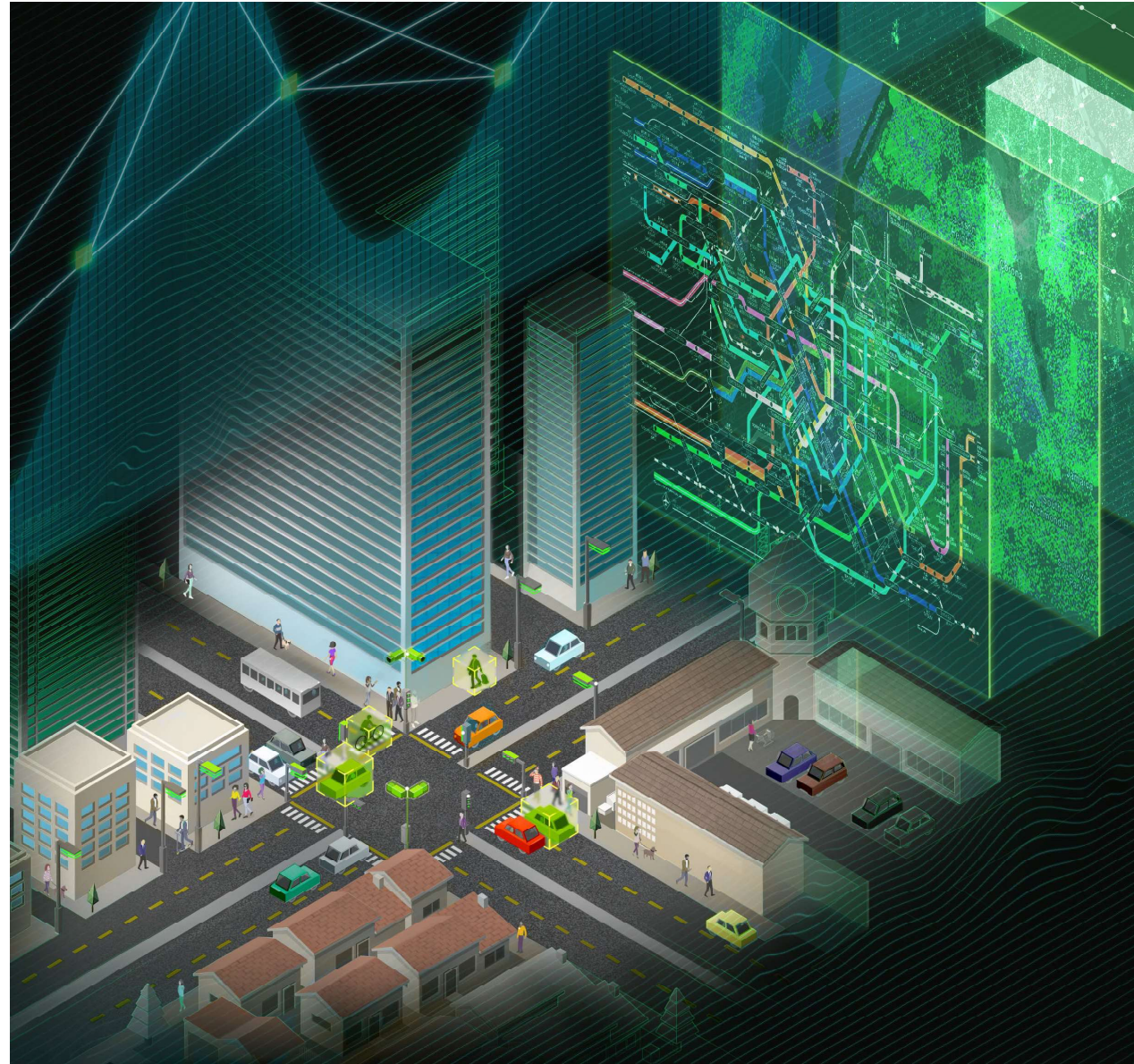
# DATA SCIENCE IN CONSUMER INTERNET

Ad Personalization

Click Through Rate Optimization

Customer Life Time Value (LTV)  
Prediction

Churn Prevention



# MACHINE LEARNING CHALLENGES

30+

Hours to  
Build GBDT  
(Gradient Boosted Tree Regression)

SLOW PROCESSES

Days

Data Transformation

Weeks

Feature Engineering

Months

Scoring Pipelines

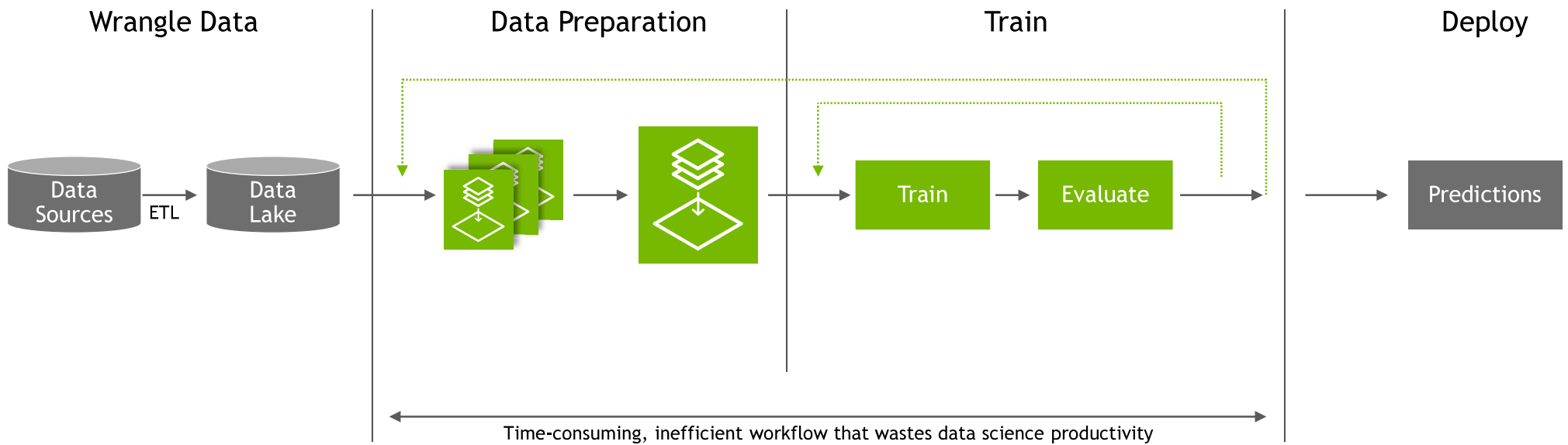
MODEL COMPLEXITY

\$3M+

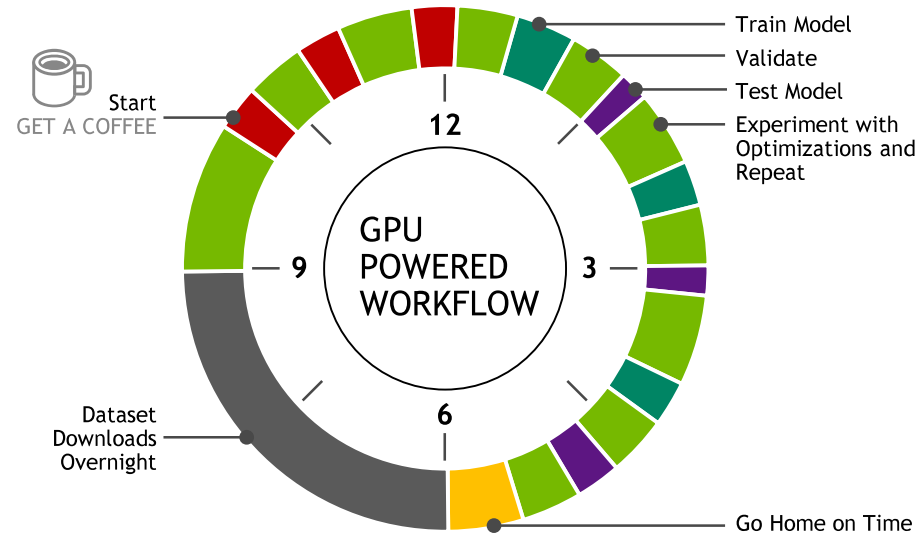
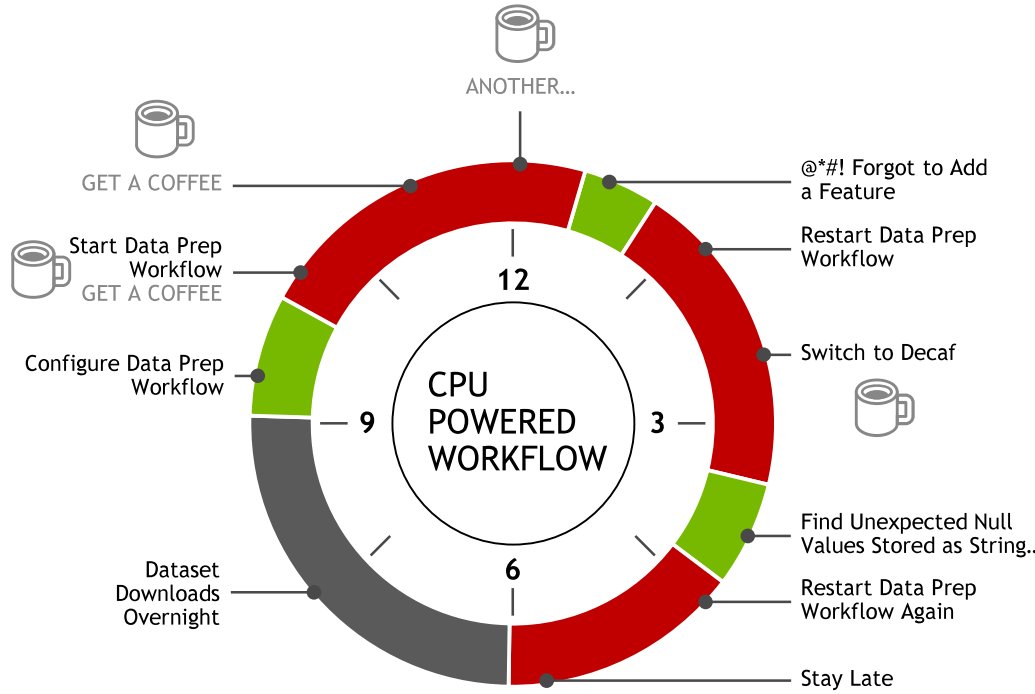
More Servers and Infrastructure  
Yielding Diminishing Returns

ESCALATING TCO

# ML WORKFLOW STIFLES INNOVATION



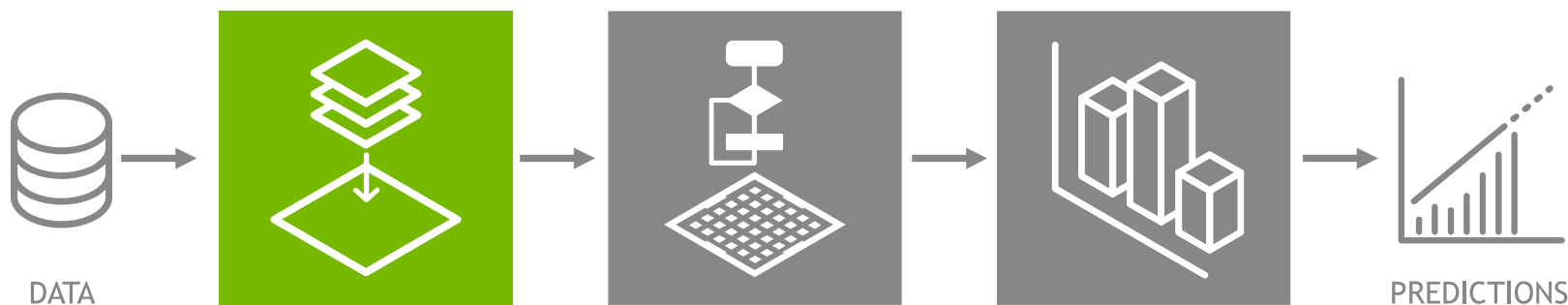
# DAY IN THE LIFE OF A DATA SCIENTIST



- Dataset Collection
- Analysis
- Data Prep
- Train
- Inference

# DATA SCIENCE WORKFLOW WITH RAPIDS

Open Source, End-to-end GPU-accelerated Workflow Built On CUDA



## DATA PREPARATION

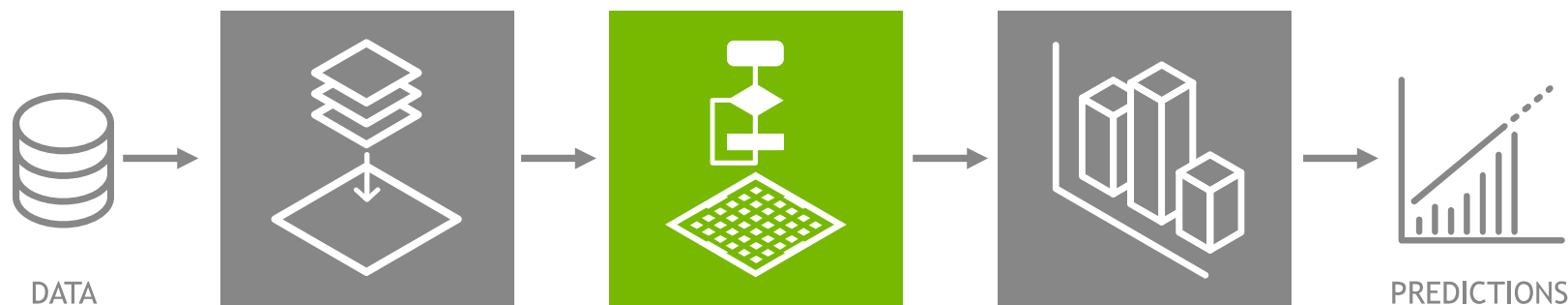
GPUs accelerated compute for in-memory data preparation

Simplified implementation using familiar data science tools

Python drop-in Pandas replacement built on CUDA C++. GPU-accelerated Spark (in development)

# DATA SCIENCE WORKFLOW WITH RAPIDS

Open Source, End-to-end GPU-accelerated Workflow Built On CUDA



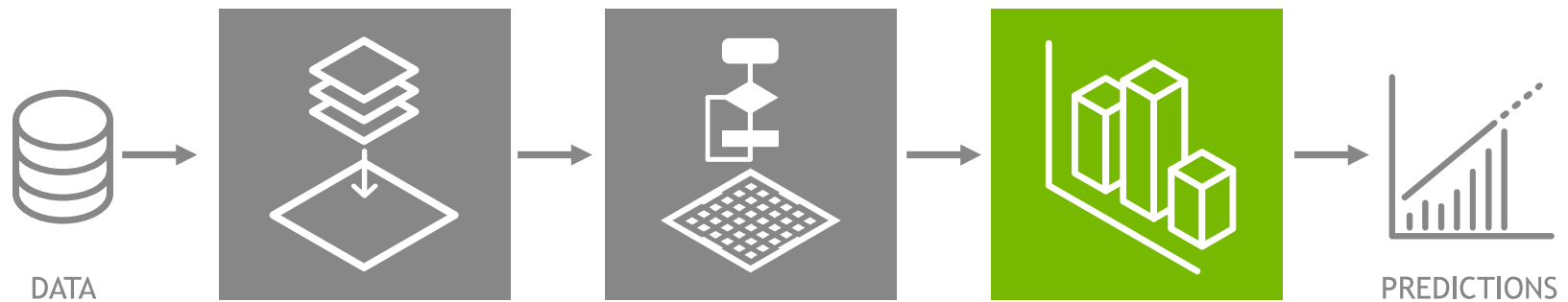
## MODEL TRAINING

GPU-acceleration of today's most popular ML algorithms

XGBoost, PCA, Kalman, K-means, k-NN, DBScan, tSVD ...

# DATA SCIENCE WORKFLOW WITH RAPIDS

Open Source, End-to-end GPU-accelerated Workflow Built On CUDA



## VISUALIZATION

Effortless exploration of datasets, billions of records in milliseconds

Dynamic interaction with data = faster ML model development

Data visualization ecosystem (Graphistry & OmniSci), integrated with RAPIDS

# TRADITIONAL DATA SCIENCE CLUSTER

Workload Profile:

Fannie Mae Mortgage Data:

- 192GB data set
- 16 years, 68 quarters
- 34.7 Million single family mortgage loans
- 1.85 Billion performance records
- XGBoost training set: 50 features

300 Servers | \$3M | 180 kW



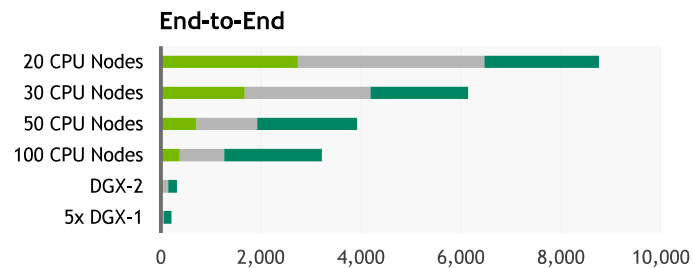
# GPU-ACCELERATED MACHINE LEARNING CLUSTER

DGX-2 and RAPIDS for  
Predictive Analytics

1 DGX-2 | 10 kW

1/8 the Cost | 1/15 the Space

1/18 the Power



# RAPIDS: DELIVERING DATA SCIENCE VALUE



Maximized Productivity

Oak Ridge National Labs
<b>215x</b> Speedup Using RAPIDS with XGBoost



Top Model Accuracy

Global Retail Giant
<b>\$1B</b> Saving with 4% Error Rate Reduction

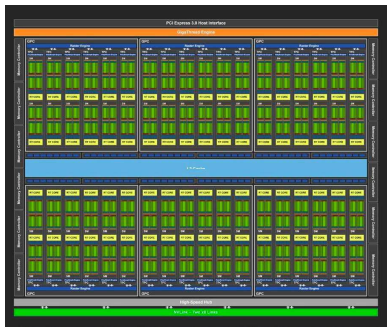


Lowest TCO

Streaming Media Company
<b>\$1.5M</b> Infrastructure Cost Saving

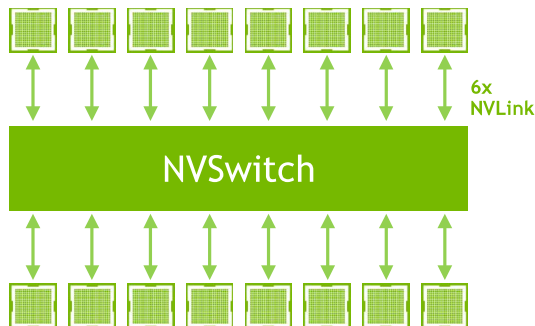
# PILLARS OF RAPIDS PERFORMANCE

## CUDA Architecture



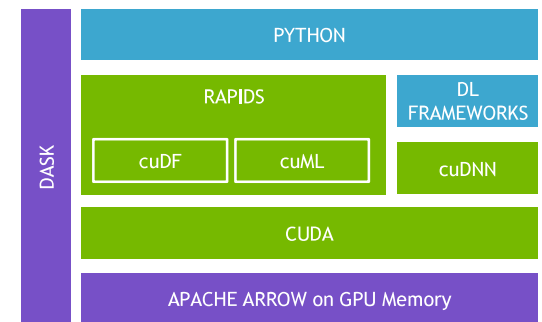
Massively Parallel Processing

## NVLink/NVSwitch



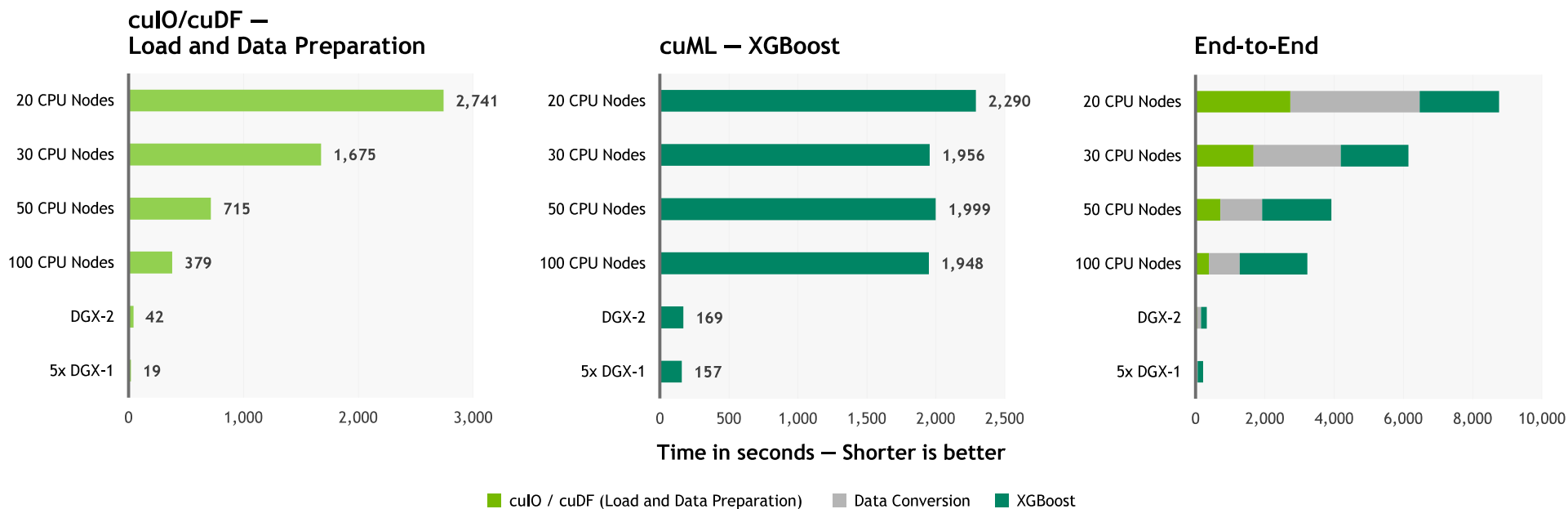
High Speed Connecting between GPUs for Distributed Algorithms

## Integrated Software



Fully Integrated Software and Hardware for Instant Productivity

# FASTER SPEEDS, REAL WORLD BENEFITS



## Benchmark

200GB CSV dataset; Data preparation includes joins, variable transformations.

## CPU Cluster Configuration

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

## DGX Cluster Configuration

5x DGX-1 on InfiniBand network

# WIDESPREAD SUPPORT FOR RAPIDS

## Open Source Community



## Enterprise Data Science Platforms



## Startups



## Deep Learning Integration



# RAPIDS

## GPU Servers



## Storage Partners



\* Spark and Hadoop support coming soon