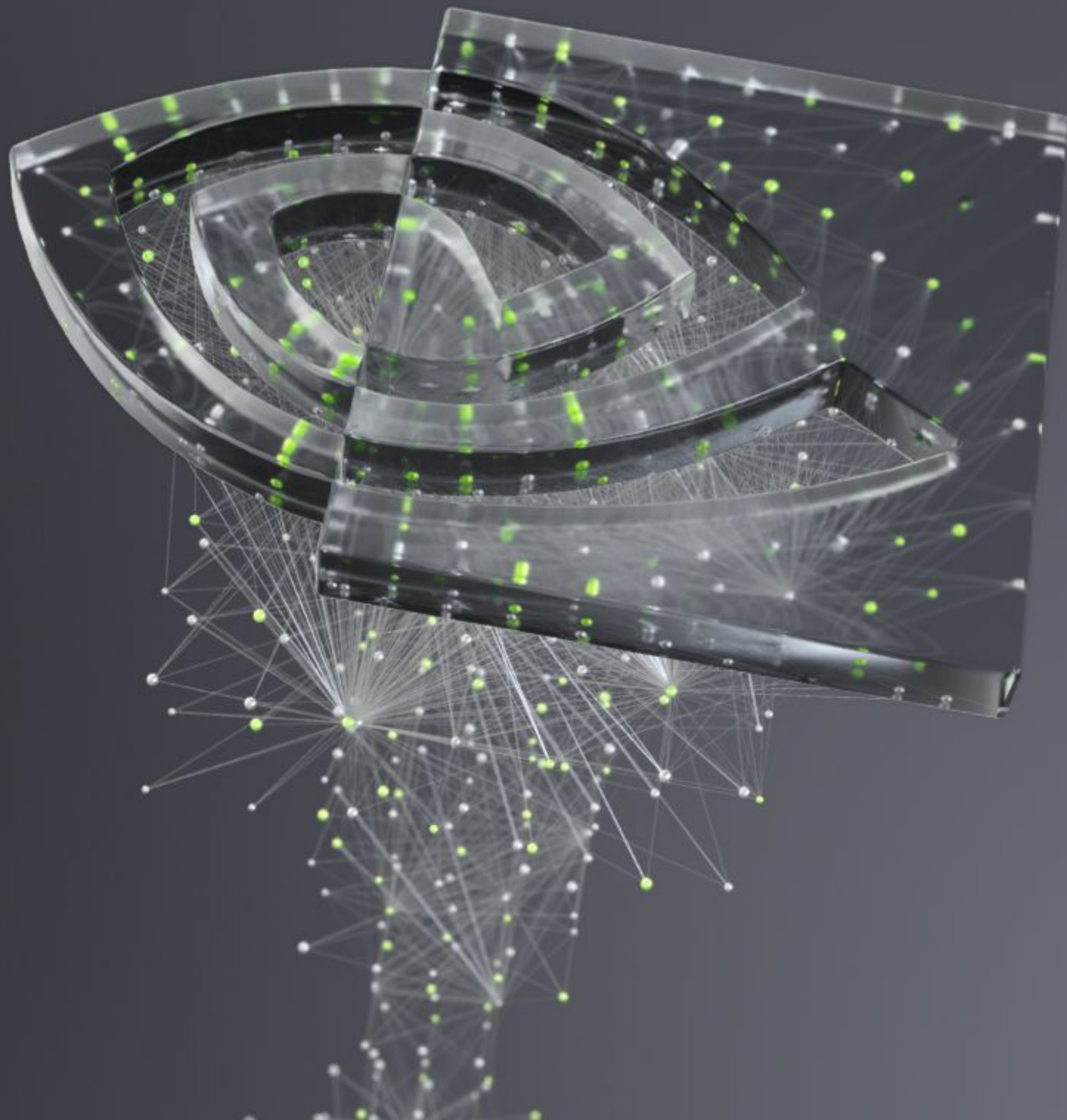# Intro to RAPIDS:

PyHEP 2021

Benjamin Zaitlen

RAPIDS Foundations

# The Evolution of Data Processing

## Faster Data Access, Less Data Movement

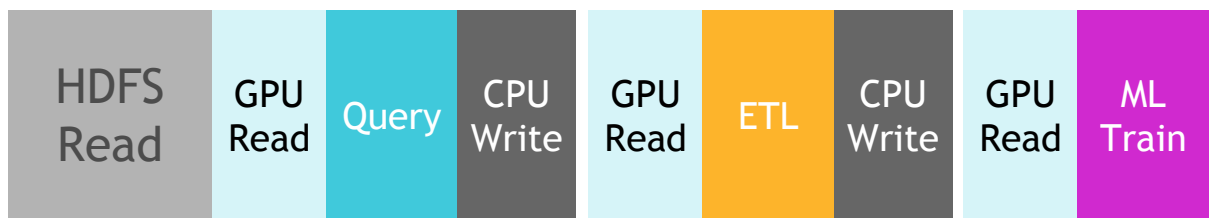**Hadoop Processing, Reading from Disk**

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

**CPU-Based Spark In-Memory Processing**

| HDFS Read | Query | ETL | ML Train |

**25-100x Improvement**
Less Code
Language Flexible
Primarily In-Memory

**Traditional GPU Processing**

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

**5-10x Improvement**
More Code
Language Rigid
Substantially on GPU
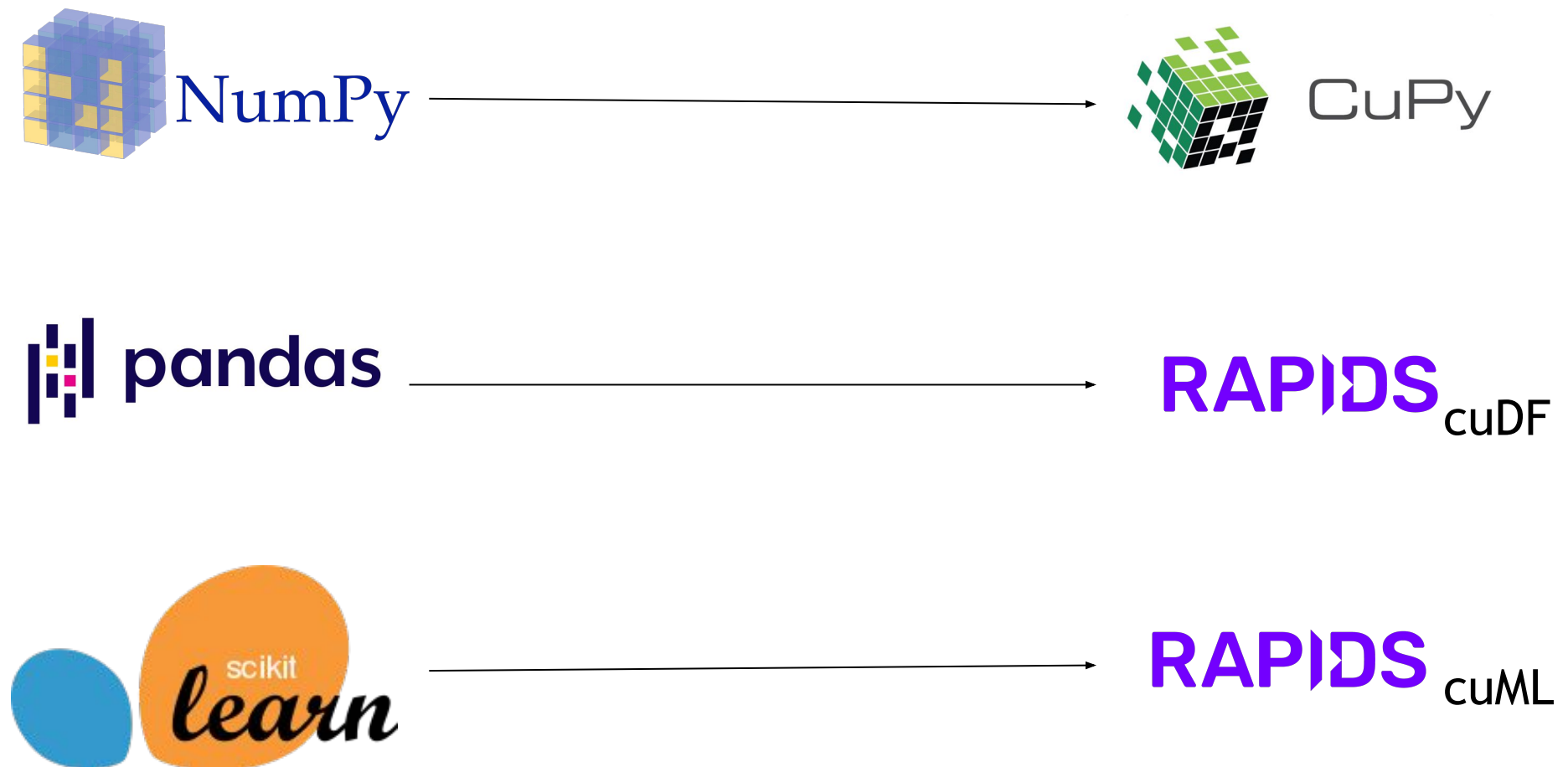
*RAPIDS*

| Arrow Read | Query | ETL | ML Train |

*50-100x Improvement*
Same Code
Language Flexible
Primarily on GPU

**?**

## Why Use GPUs

*GPUs are built for intensive parallel processing. As datasets continue to grow, data scientists are limited by the sequential nature of CPU compute. GPUs provide the power and parallelism necessary for today's data science.*
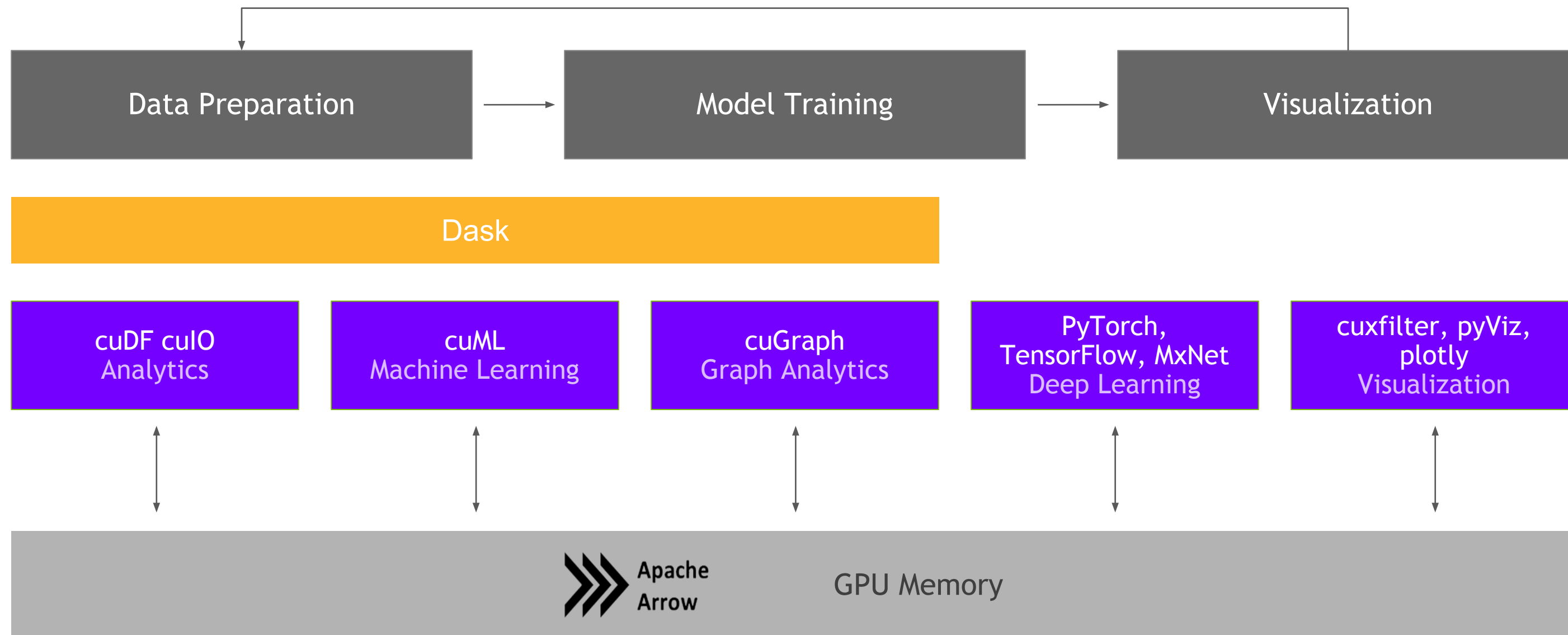
3 NVIDIA.
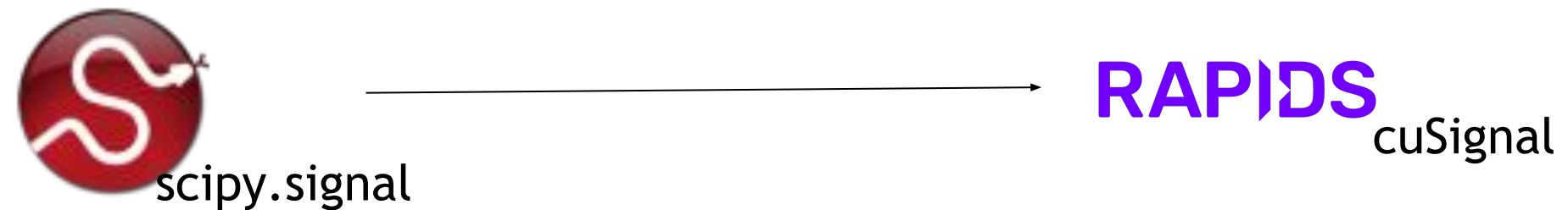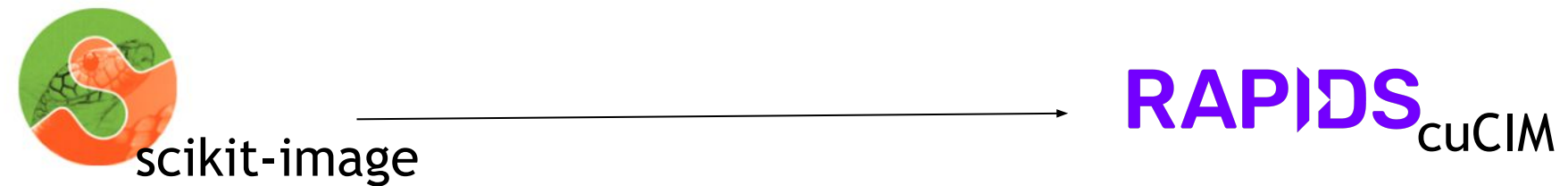
# RAPIDS
## GPU Accelerated Data Science



NumPy → CuPy

pandas → **RAPIDS**cuDF

scikit learn → **RAPIDS**cuML

NVIDIA.

# RAPIDS

## End-to-End GPU Accelerated Data Science

| Data Preparation | Model Training | Visualization |
|---|---|---|

**Dask**

| cuDF cuIO<br>Analytics | cuML<br>Machine Learning | cuGraph<br>Graph Analytics | PyTorch,<br>TensorFlow, MxNet<br>Deep Learning | cuxfilter, pyViz,<br>plotly<br>Visualization |
|---|---|---|---|---|

Apache Arrow    GPU Memory

# RAPIDS
## GPU Accelerated Data Science



NumPy → CuPy

pandas → RAPIDS cuDF

scikit learn → RAPIDS cuML

scikit-image → RAPIDS cuCIM

scipy.signal → RAPIDS cuSignal

6

# Growth and Adoption

## Growing community engagement

**RAPIDS Monthly Downloads**



116k

@RAPIDSai

Followers: 10,345

github.com/rapidsai

Stars: 8K+     Contributors: 150+

# RAPIDS

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

| Data Preparation | Model Training | Visualization |
|---|---|---|

**Dask**

| cuDF cuIO<br>Analytics | cuML<br>Machine Learning | cuGraph<br>Graph Analytics | PyTorch,<br>TensorFlow, MxNet<br>Deep Learning | cuxfilter, pyViz,<br>plotly<br>Visualization |
|---|---|---|---|---|

Apache Arrow    GPU Memory

# What is cuDF?

## Expandable platform for GPU data science

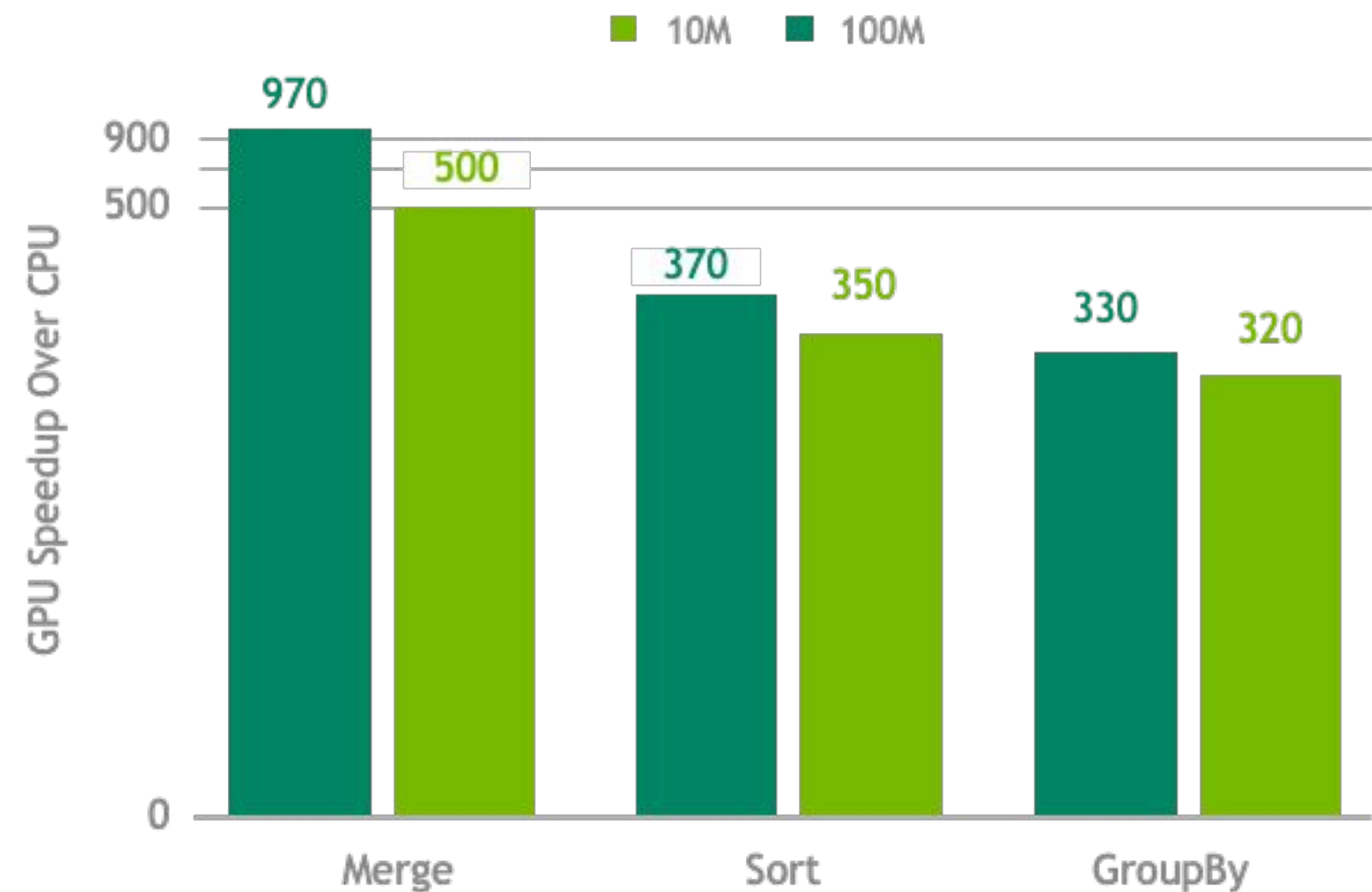| | |
|---|---|
| **Python** | |
| **Cython** | → Dask cuDF<br>cuDF<br>Pandas |
| **libcudf C++** | |
| **CUDA Libraries** | → Thrust<br>Cub<br>Jitify |
| **CUDA** | |

- Familiar pandas-like Python API
- Table (dataframe) and column types and algorithms
- High-performance C++ layer provides GPU-optimized CUDA kernels, data types, operations, and primitives
- CUDA/C++ is top level supported and used by many for integrating RAPIDS

# ACCELERATED PRE-PROCESSING

## A FAMILIAR EXPERIENCE FOR DATA ENGINEERS

RAPIDS provides a GPU DataFrame library with a pandas-like API while providing *significant* performance improvements.
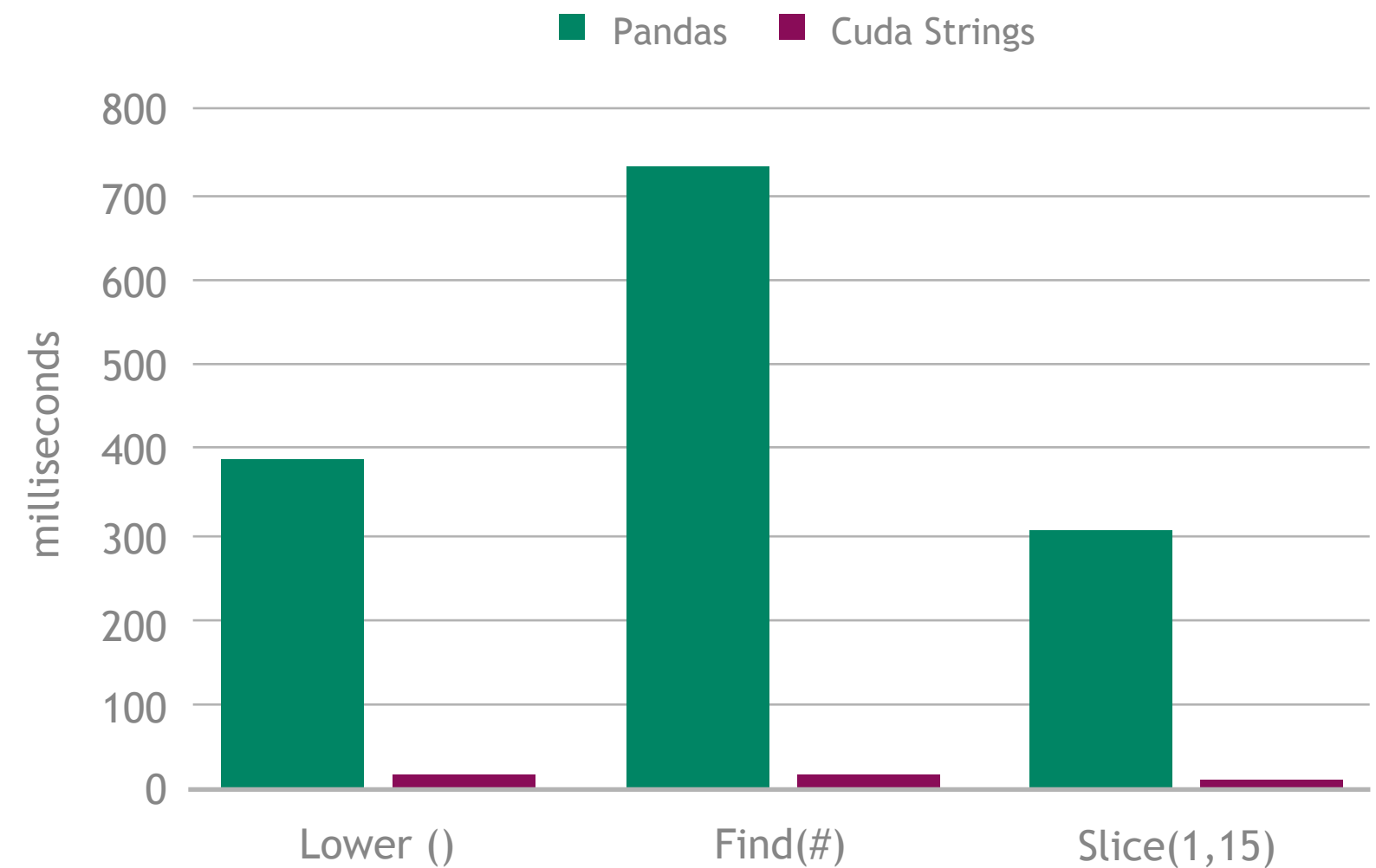


*Single GPU Speed-Ups vs pandas*

GPU: NVIDIA Tesla V100 32GB on DGX-1
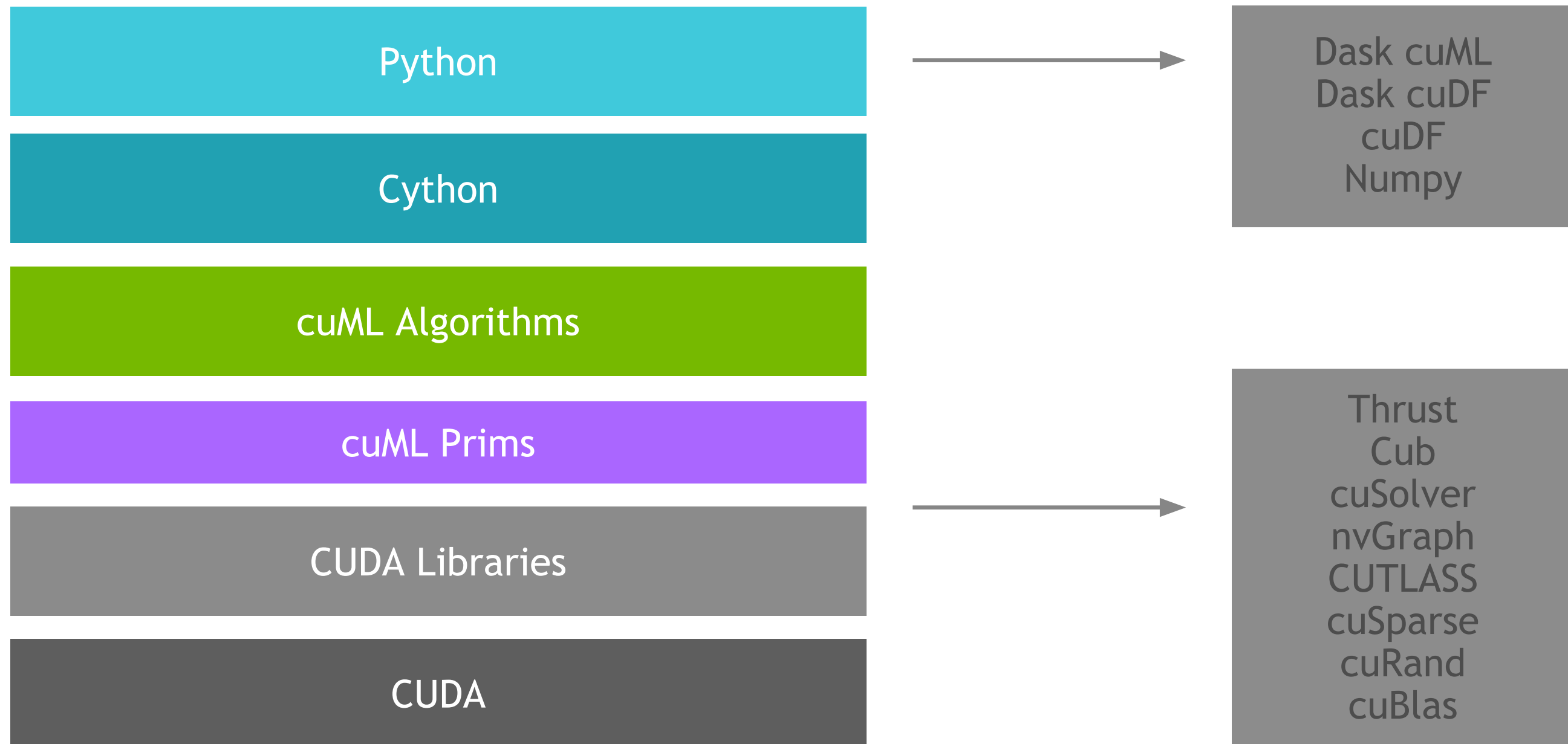CPU: Intel(R) Xeon(R) CPU E5-2698 v4 @ 2.20GHz

# Comprehensive String Support

## Backbone of ETL: Strings

- Regular Expressions

- Element-wise operations
  - Split, Find, Extract, Cat, Typecasting, etc...

- String GroupBys, Joins, Sorting, etc.

- Categorical columns fully on GPU

- NLP Preprocessors
  - Tokenizers, Normalizers, Edit Distance, Porter Stemmer, etc.

# ML Technology Stack

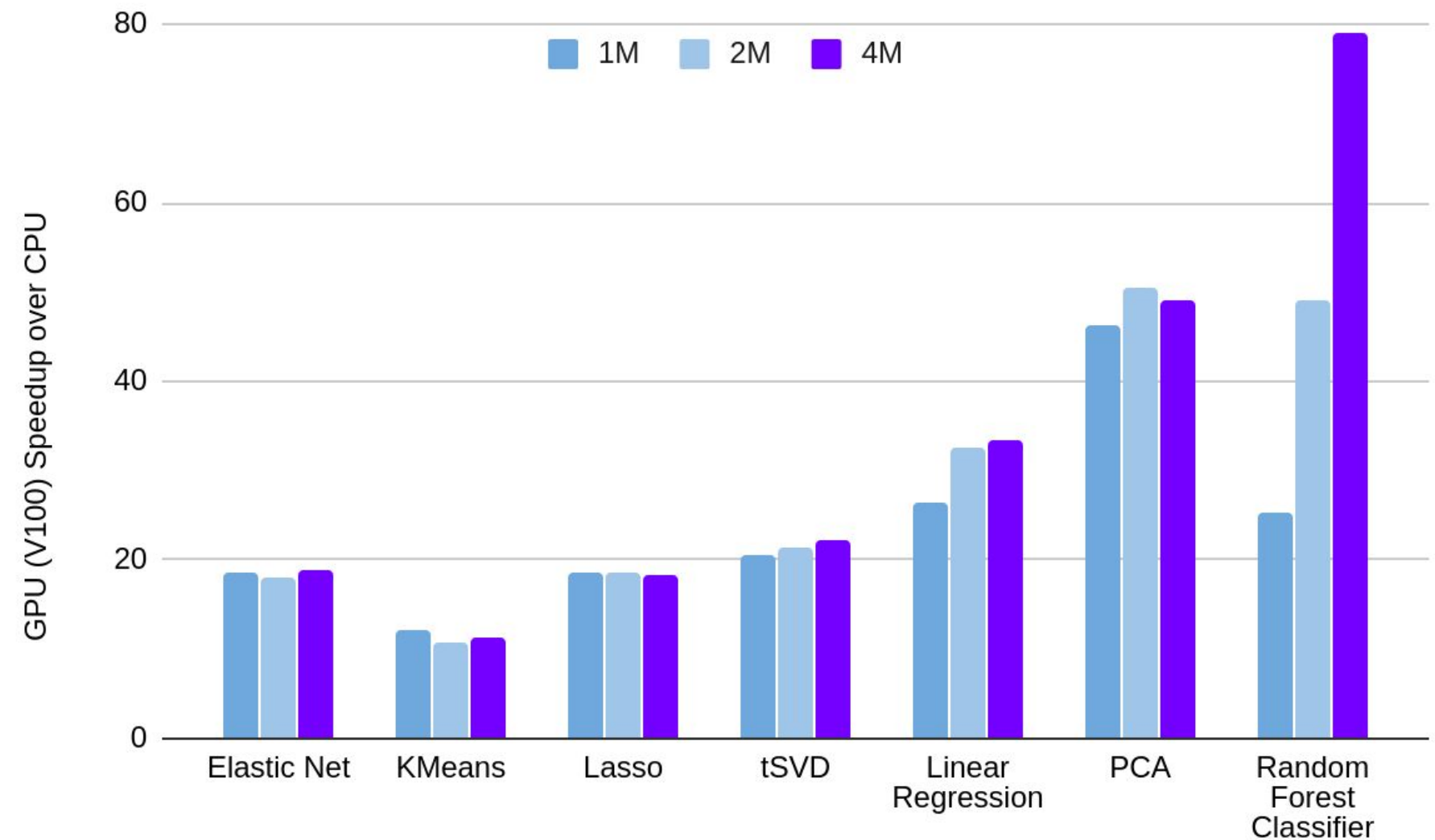| | | |
|---|---|---|
| Python | → | Dask cuML<br>Dask cuDF<br>cuDF<br>Numpy |
| Cython | | |
| cuML Algorithms | | |
| cuML Prims | | Thrust<br>Cub<br>cuSolver<br>nvGraph<br>CUTLASS<br>cuSparse<br>cuRand<br>cuBlas |
| CUDA Libraries | → | |
| CUDA | | |

# ACCELERATED MACHINE LEARNING
## GPU-POWER WITH THE FEEL OF SCIKIT-LEARN

RAPIDS provides a GPU ML library with a scikit-learn API while providing *significant* performance improvements.

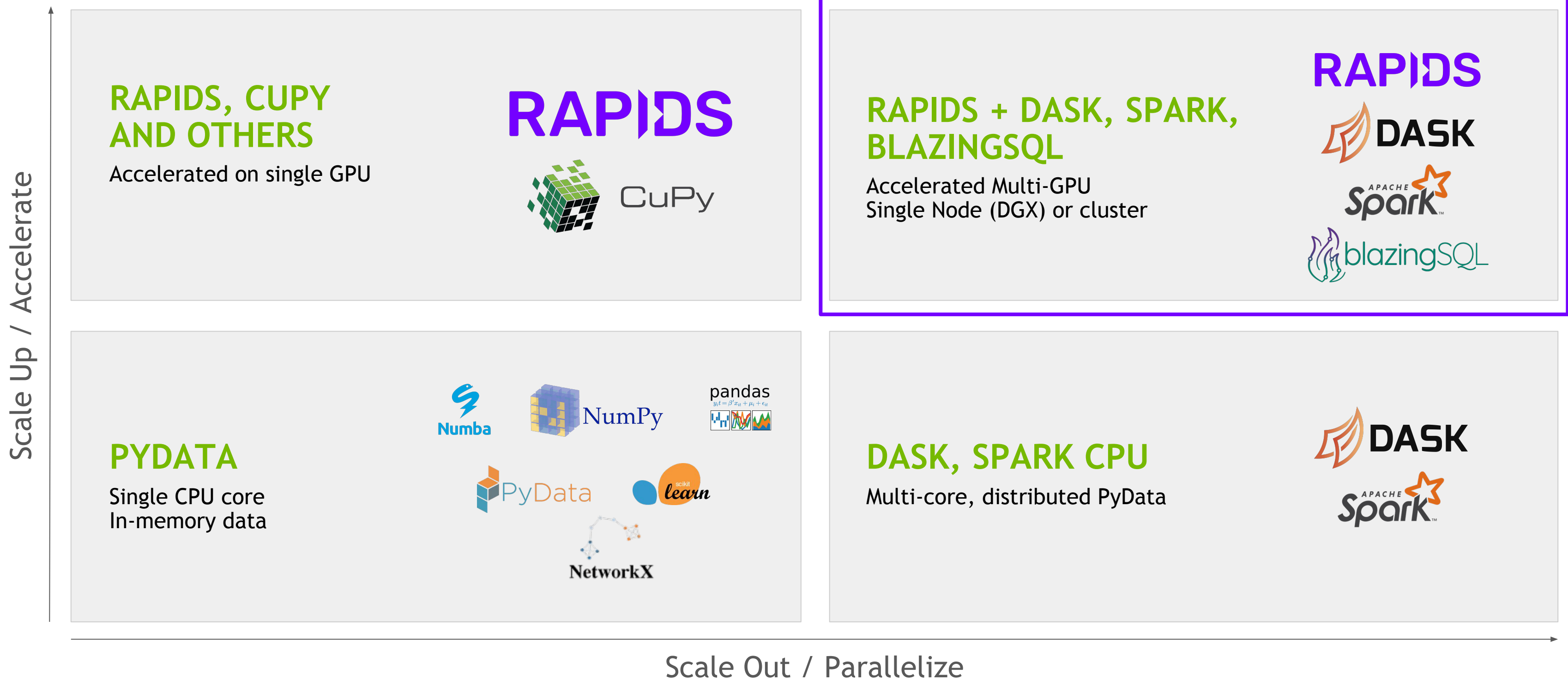**26 GPU-Accelerated Algorithms & Growing**



1x V100 vs. 2x 20 Core CPUs (DGX-1, RAPIDS 0.15)

NVIDIA.

# Scale Out with RAPIDS

## Multiple options to scale, from multi-GPU to a whole cluster

**Scale Up / Accelerate** (vertical axis)

### RAPIDS, CUPY AND OTHERS
Accelerated on single GPU

**RAPIDS** · CuPy

### RAPIDS + DASK, SPARK, BLAZINGSQL
Accelerated Multi-GPU
Single Node (DGX) or cluster

**RAPIDS** · DASK · Apache Spark · blazingSQL

### PYDATA
Single CPU core
In-memory data

Numba · NumPy · pandas · PyData · scikit learn · NetworkX

### DASK, SPARK CPU
Multi-core, distributed PyData

DASK · Apache Spark

**Scale Out / Parallelize** (horizontal axis)

NVIDIA.

Python library for parallel computing

Scales Numpy, Pandas, and Scikit-Learn
Accelerates custom systems

Easy for beginners,
Secure and trusted for institutions

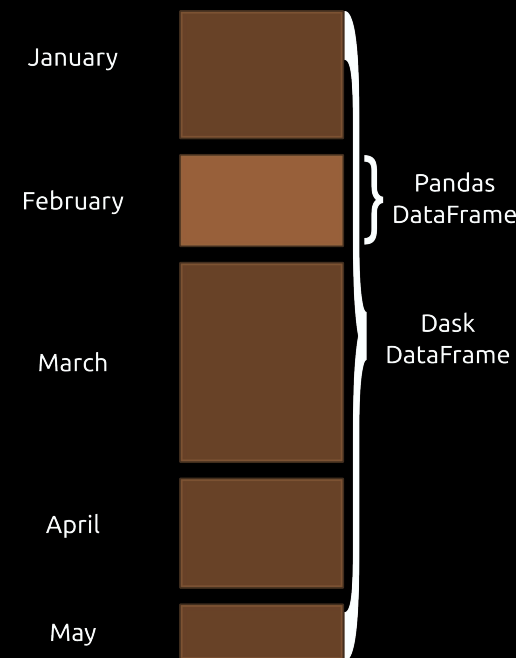# Dask accelerates the existing Python ecosystem

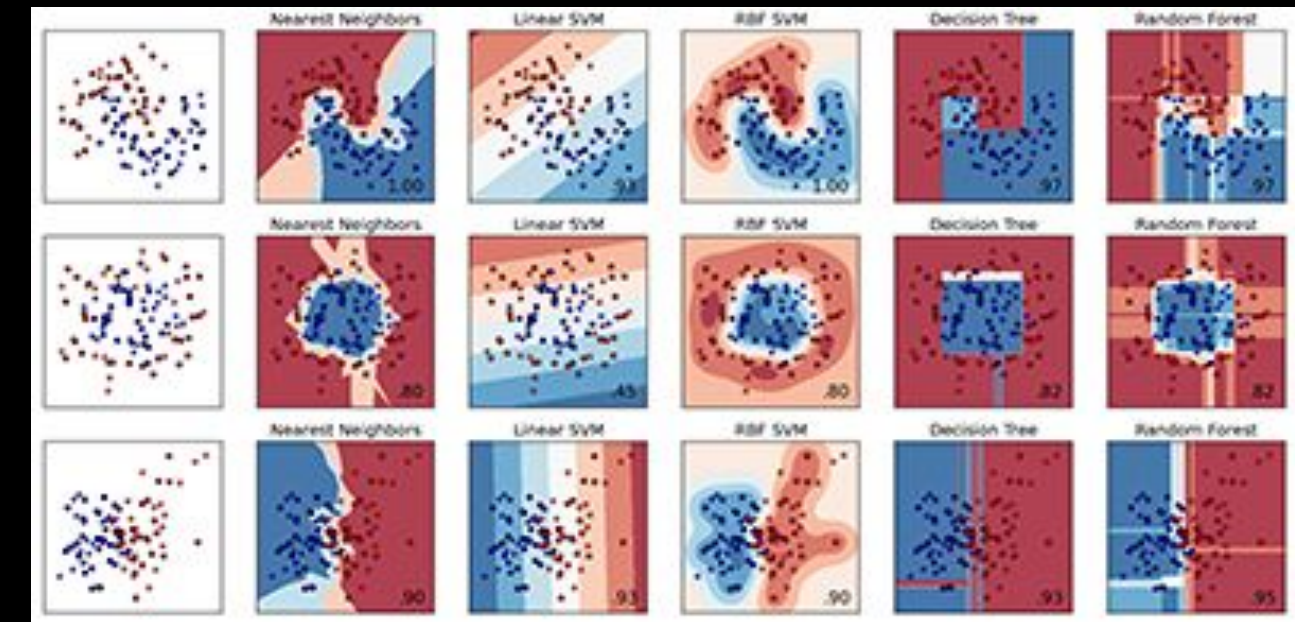Built alongside the current community

## Numpy



```
import numpy as np

x = np.ones((1000, 1000))

x + x.T - x.mean(axis=0)
```

## Pandas



```
import pandas as pd

df = pd.read_csv("file.csv")

df.groupby("x").y.mean()
```

## Scikit-Learn



```
from scikit_learn.linear_model \
        import LogisticRegression

lr = LogisticRegression()

lr.fit(data, labels)
```
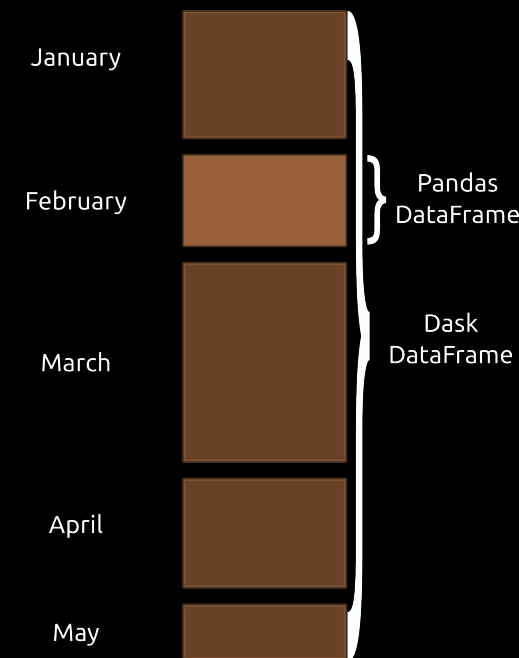
DASK

# Dask accelerates the existing Python ecosystem

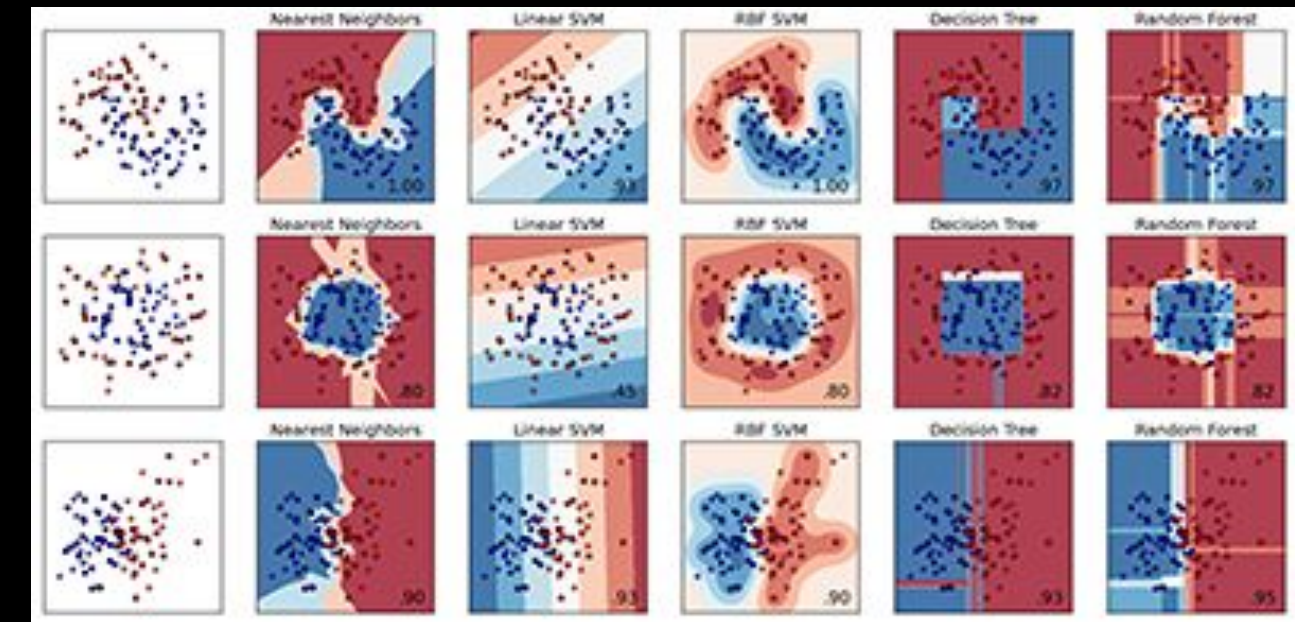Built alongside the current community

## Numpy



## Pandas



## Scikit-Learn



```
import dask.array as da

x = da.ones((10000, 10000))
x + x.T - x.mean(axis=0)
```

```
import dask.dataframe as dd

df = dd.read_csv("s3://*.csv")
df.groupby("x").y.mean()
```

```
from dask_ml.linear_model \
     import LogisticRegression

lr = LogisticRegression()
lr.fit(data, labels)
```

# Parallelize existing complex code

Dask scales existing codebases with modest changes

```
def f(data, model) -> pd.DataFrame:
    …


def g(data, model) -> pd.DataFrame:
    …


results = []

for x in A:
  for y in B:
    if x < y:
        results.append(f(x, y))
    else:
        results.append(g(x, y))
```

Many codebases have opportunities for parallelism

But the problem doesn't look like a big array or big dataframe

DASK

# Parallelize existing complex code

Dask scales existing codebases with modest changes

```
@dask.delayed
def f(data, model) -> pd.DataFrame:
    …
@dask.delayed
def g(data, model) -> pd.DataFrame:
    …


results = []

for x in A:
  for y in B:
    if x < y:
        results.append(f(x, y))
    else:
        results.append(g(x, y))


results = dask.compute(results)
```

Dask Delayed adds parallelism without changing existing logic.

Dask lazily traverses your code to build a recipe for future execution.
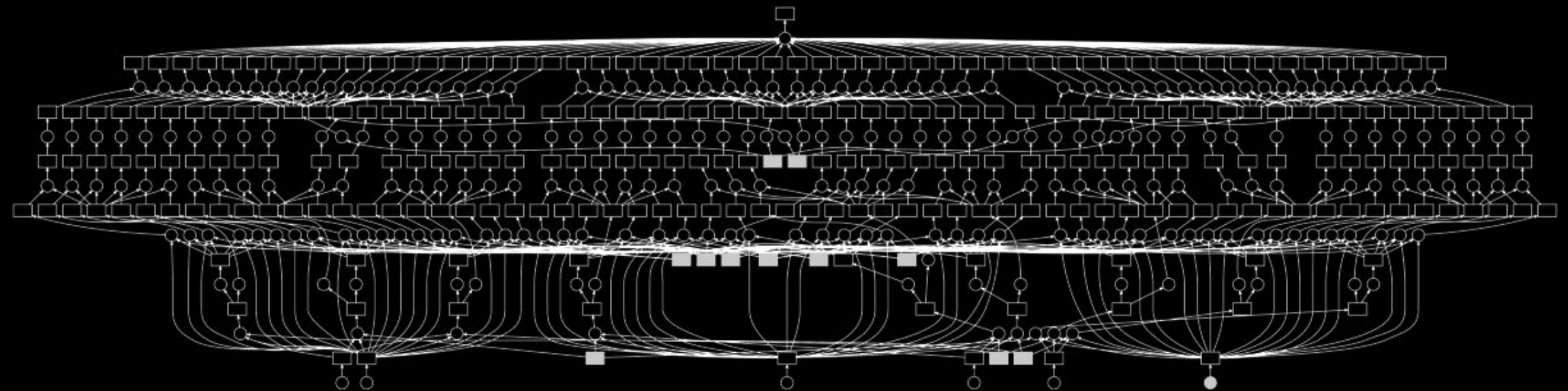
DASK

# Parallelize existing complex code

Dask scales existing codebases with modest changes

```python
@dask.delayed
def f(data, model) -> pd.DataFrame:
    …

@dask.delayed
def g(data, model) -> pd.DataFrame:
    …


results = []

for x in A:
  for y in B:
    if x < y:
        results.append(f(x, y))
    else:
        results.append(g(x, y))

results = dask.compute(results)
```

Your code creates a task graph for future execution.
Each node is one Python function.
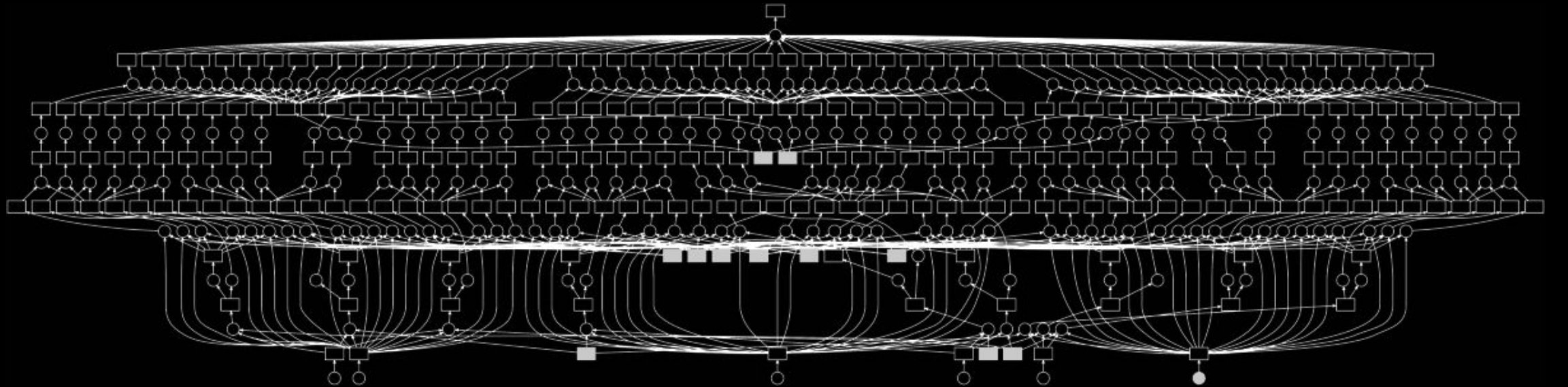
# Parallelize existing complex code

Dask scales existing codebases with modest changes

```
@dask.delayed
def f(data, model) -> pd.DataFrame:
    …

@dask.delayed
def g(data, model) -> pd.DataFrame:
    …


results = []


for x in A:
  for y in B:
    if x < y:
        results.append(f(x, y))
    else:
        results.append(g(x, y))


results = dask.compute(results)
```



Dask then executes that graph on parallel hardware

**DASK**

# Dask deploys on all major resource managers

Cloud, HPC, or Yarn, it's all the same to Dask

| **Cloud** | **HPC** | **Hadoop/Spark** |
|:---:|:---:|:---:|

```
cluster = KubeCluster()

cluster = ECSCluster()



df = dd.read_parquet(...)
```
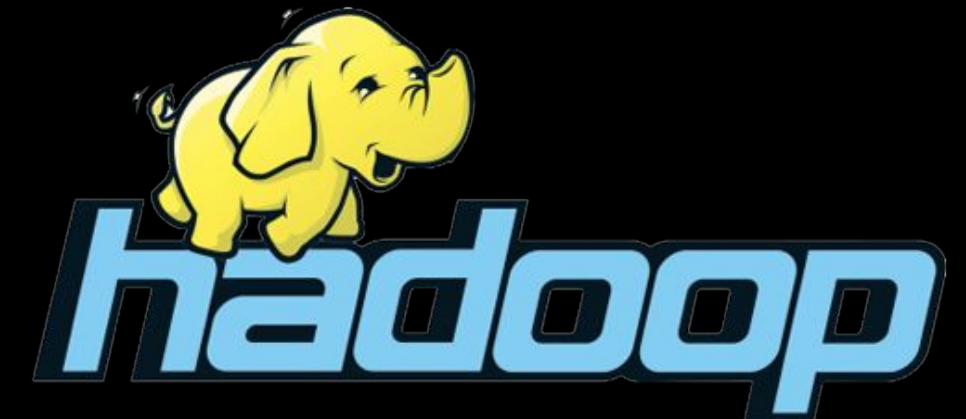
```
cluster = PBSCluster()

cluster = LSFCluster()

cluster = SLURMCluster()

…

df = dd.read_parquet(...)
```

```
cluster = YarnCluster()




df = dd.read_parquet(...)
```
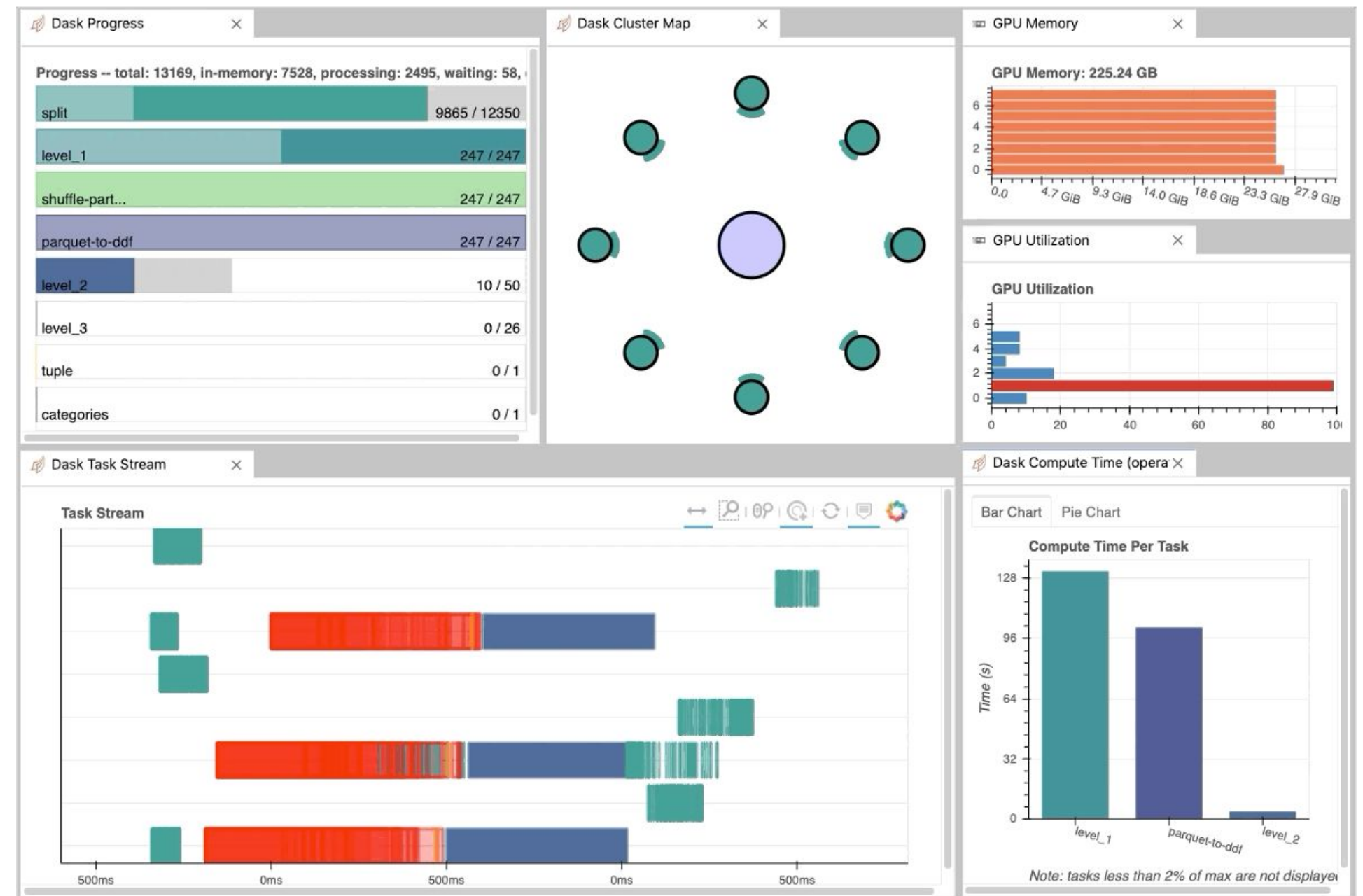
# Dask + RAPIDS

## PyData-native scalable analytics

- **Deployable**: Kubernetes, Yarn, SLURM

- **PyData native**: Easy migration, built on top of NumPy, Pandas, Scikit-learn

- **Easy scalability**: Easy to install; scales to thousands of nodes

- **Popular**: most Common parallelism framework in PyData and SciPy community

# RAPIDS Dev Environment

## JupyterLab + Friends



- JupyterLab
- Dask Extension
- NVDashboard Extension

# Faster Speeds, Real World Benefits

## Faster Data Access, Less Data Movement

### cuIO/cuDF –
### Load and Data Preparation

| | Time (sec) |
|---|---|
| 20 CPU Nodes | 2,74 |
| 30 CPU Nodes | 1,675 |
| 50 CPU Nodes | 715 |
| 100 CPU | 379 |
| 16x A100 | 30 |

### XGBoost Machine Learning

| | Time (sec) |
|---|---|
| 20 CPU | 2,290 |
| 30 CPU | 1,956 |
| 50 CPU | 1,999 |
| 100 CPU | 1,948 |
| 16x A100 | 73 |

### End-to-End

| | Time (secs) |
|---|---|
| 20 CPU | 8,763 |
| 30 CPU | 6,147 |
| 50 CPU | 3,926 |
| 100 CPU | 3,221 |
| 16x A100 | 116 |

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

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

| Benchmark | CPU Cluster Configuration | A100 Cluster Configuration | RAPIDS Version |
|---|---|---|---|
| 200GB CSV dataset; Data prep includes joins, variable transformations | CPU nodes (61 GiB memory, 8 vCPUs, 64-bit platform), Apache Spark | 16 A100 GPUs (40GB each) | RAPIDS 0.19 |

# RAPIDS/Dask End-to-End Performance

## Reducing Data Science Processes from Hours to Seconds

RAPIDS delivers massive speed-ups across the end-to-end data science lifecycle. Conducting benchmarks in a commercial cloud environment, we're able to get incredible performance running a common ML model training pipeline.

Between loading and cleansing data, engineering features, and training a classifier using a 200GB CSV dataset, a RAPIDS-based pipeline completed these operations in *just over two minutes*. The same process takes two and half hours on a similar CPU-configuration.

### RAPIDS End-to-End Workflow Runtimes

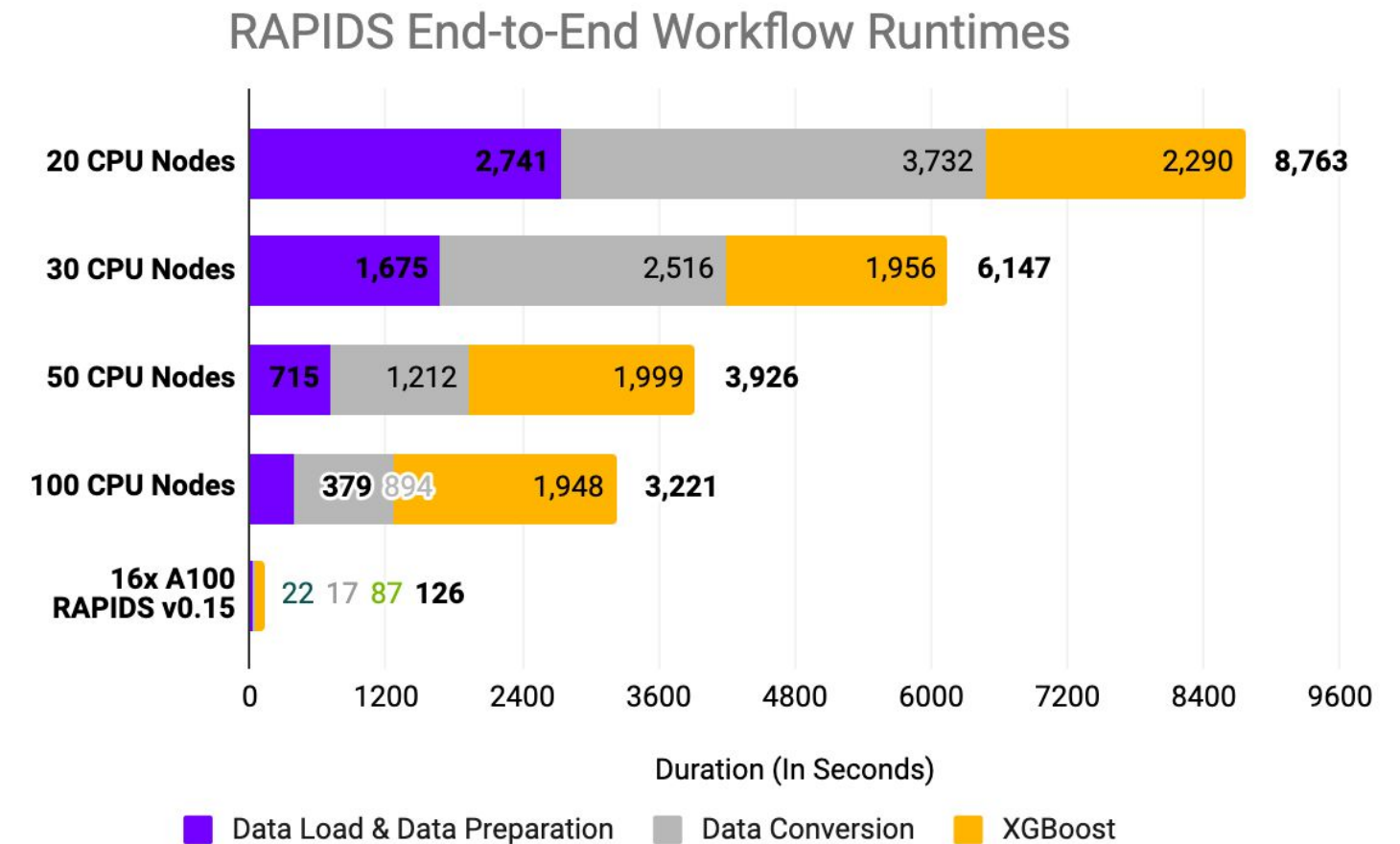| | Data Load & Data Preparation | Data Conversion | XGBoost | Total |
|---|---|---|---|---|
| 20 CPU Nodes | 2,741 | 3,732 | 2,290 | 8,763 |
| 30 CPU Nodes | 1,675 | 2,516 | 1,956 | 6,147 |
| 50 CPU Nodes | 715 | 1,212 | 1,999 | 3,926 |
| 100 CPU Nodes | | 379 | 894 1,948 | 3,221 |
| 16x A100 RAPIDS v0.15 | 22 | 17 | 87 | 126 |

Duration (In Seconds)

■ Data Load & Data Preparation  ■ Data Conversion  ■ XGBoost

## 16
### A100s Provide More Power than 100 CPU Nodes
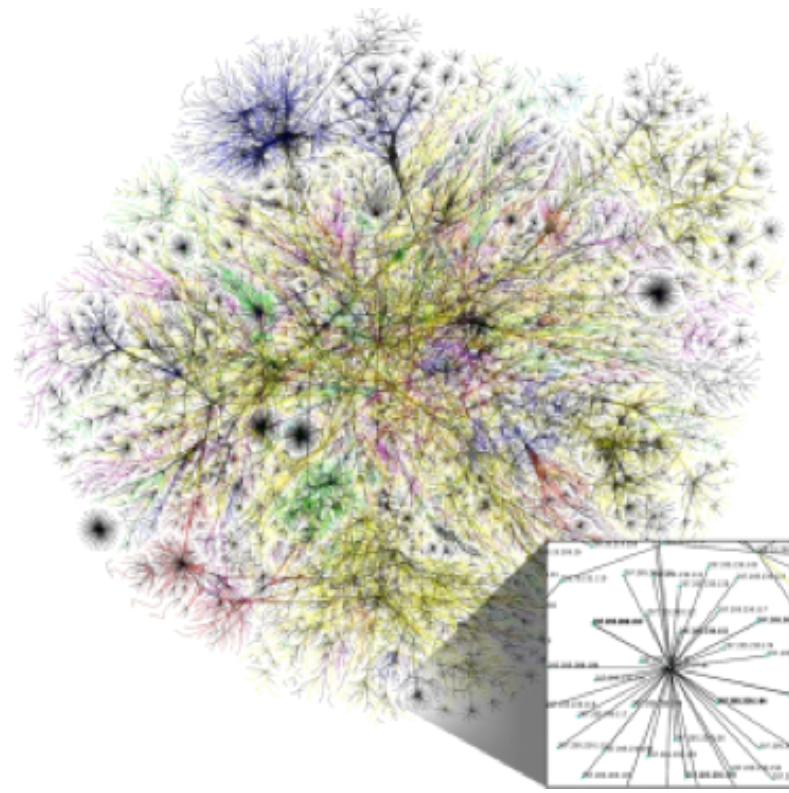
## 70x
### Faster Performance than Similar CPU Configuration

## 20x
### More Cost-Effective than Similar CPU Configuration

*CPU approximate to n1-highmem-8 (8 vCPUs, 52GB memory) on Google Cloud Platform. TCO calculations-based on Cloud instance costs.
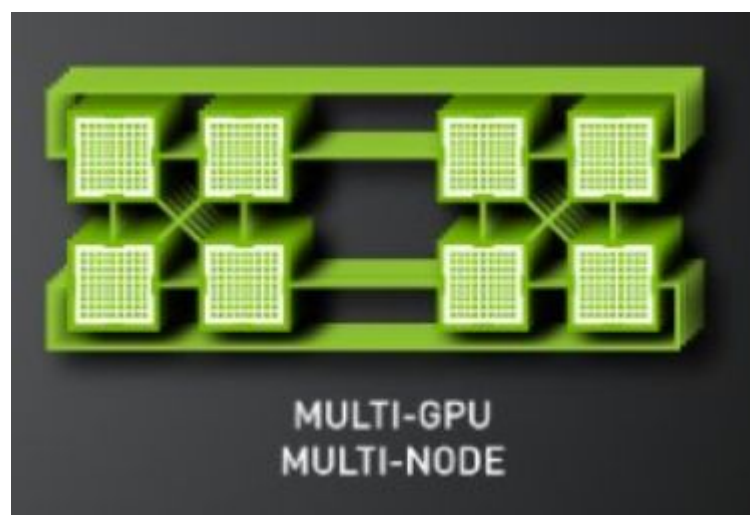
# cuGraph - Algorithms

## GPU-accelerated NetworkX

| | | Spectral Clustering |
|---|---|---|
| | **Community** | - Balanced Cut and Modularity Maximization |
| | | Louvain (Multi-GPU) and Leiden |
| | | Ensemble Clustering for Graphs |
| | | KCore and KCore Number |
| | | Triangle Counting |
| | | K-Truss |

**Traveling Salesman** | **Routing**

| **Components** | Weakly Connected Components |
|---|---|
| | Strongly Connected Components |

Minimum Spanning Tree
Maximum Spanning Tree | **Tree**

| **Link Analysis** | Page Rank (Multi-GPU) |
|---|---|
| | Personal Page Rank (Multi-GPU) |
| | HITS |

Graph Classes
Subgraph Extraction
**Egonet** | **Structure**

| **Link Prediction** | Jaccard |
|---|---|
| | Weighted Jaccard |
| | Overlap Coefficient |

Force Atlas 2
Hungarian Algorithm | **Other**

| **Traversal** | Single Source Shortest Path (SSSP) (Multi-GPU) |
|---|---|
| | Breadth First Search (BFS) (Multi-GPU) |

Renumbering
Auto-Renumbering
NetworkX converters | **Utilities**

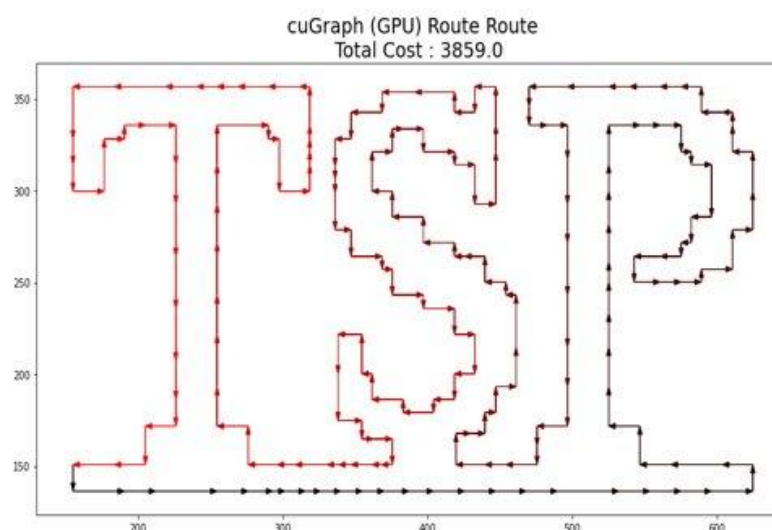| **Centrality** | Katz (Multi-GPU) |
|---|---|
| | Betweenness Centrality (Vertex and Edge) |

# Scaling and Expanding Graph Analytics



**Multi-node, Multi-GPU Scaling**

New graph primitives will underpin all algorithms

PageRank performance up to 180x faster than CPU
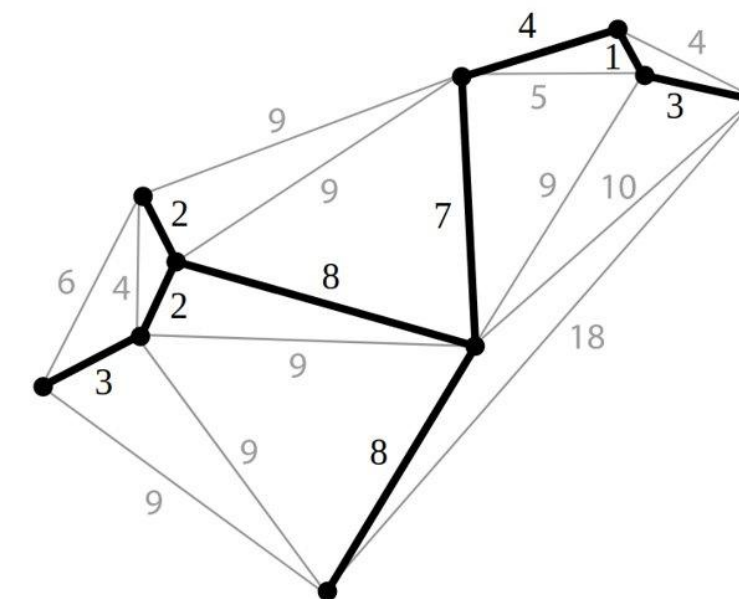
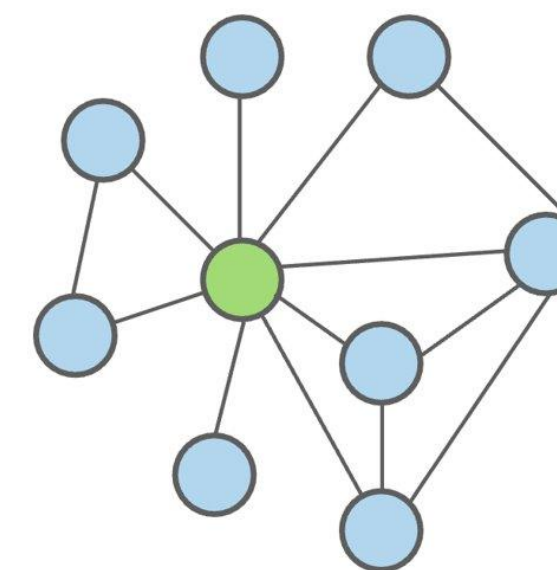New 2d partitioning methods for large graphs



**Traveling Salesperson Solver**

Up to 32x faster than CPU alternatives



**Improved NetworkX API Compatibility**



**Minimum Spanning Tree**



**EGONet**

# Visualization and NodeJS



## Plot.ly Dash

A Python visualization framework able to use RAPIDS libraries for notebooks and hosted dashboard applications.

Ideal for purpose built analytics applications, but also useful for notebook workflows.

Integrated RAPIDS backend for large datasets.

## RAPIDS cuXFilter

A Python notebook based crossfilter dashboard library, using cuDF. Incorporates many chart libraries such as Datashader, HvPlot, Holoviews, Bokeh, and Deck.gl.

Easy integration with RAPIDS notebook based workflows.

## RAPIDS Node.js (early alpha)

Experimental Node.js Javascript bindings for RAPIDS and related GPU libraries. Usable for both visualization and general-purpose compute on Node.js platforms.

# Build End-to-End Data Science Applications

## Leverage RAPIDS Core Libraries to Build Custom Solutions

| | Description | Similar To | Problem Domain | Maturity | Performance | Example User | API Docs |
|---|---|---|---|---|---|---|---|
| **cuDF** | Dataframes & ETL | pandas | Data Preparation | | | Walmart | Read the Docs |
| **Apache Spark 3.0 Plugin** | ETL | Apache Spark | Data Preparation | | | CLOUDERA | Read the Docs |
| **BlazingSQL** | ANSI SQL | SQL | Data Preparation | | | OAK RIDGE National Laboratory | Read the Docs |
| **cuML** | Machine Learning | scikit-learn | Model Training | | | Capital One | Read the Docs |
| **cuGraph** | Graph Analytics | NetworkX | Model Training | | | VISA | Read the Docs |
| **XGBoost** | GBMs | XGBoost | Model training | | | Scotiabank | Read the Docs |
| **RAPIDSViz** | Large-Scale Visualization | Bokeh, DataShader, HoloViews | Visualization | | | plotly | Read the Docs |

# Use RAPIDS-Enabled Tools & Frameworks

## High-Performance Solutions for a Wide Variety of Domains

| | Description | Similar To | Problem Domain | Maturity | Performance | Example User | API Docs |
|---|---|---|---|---|---|---|---|
| **CLX** | Cyber log parsing & analytics | N/A | Cybersecurity | | | BEST BUY | Read the Docs |
| **cuCIM** | Image processing & analytics | scikit-image | Image Processing | | | QuanSight | Read the Docs |
| **cuSignal** | Signal processing & analytics | N/A | Signal Processing | | | LOCKHEED MARTIN | Read the Docs |
| **cuSpatial** | Spatial processing & analytics | N/A | Spatial Data Processing | | | tsmc | Read the Docs |
| **cuStreamz** | Stream processing & analytics | Streamz & Kafka | Stream Processing | | | NVIDIA | Read the Docs |
| **Node-RAPIDS** | Server-side JavaScript | Node.js | Web Development | | | Technical Preview | Read the Docs |
| **NVTabular** | Feature engineering and data loading | N/A | Recommender Systems | | | Postmates | Read the Docs |

31

# cuDF

## Updates + Improvements



### Optimizations

- Complex Hash Aggregations
- Character Parallel String Algorithms
- Parquet GPU Direct Storage Support

### New Features

- List, Struct, Dictionary, And Decimal types and operations
- Expanded GroupBy Operations
- Improved API and Developer Docs

### Build Infrastructure

- CUDA 11.2 Support And CUDA Enhanced Compatibility
- CMake Refactored for easier source builds

## Upcoming Improvements

- Abstract Syntax Tree Evaluation
- ORC GPU Direct Storage
- Reduce Python Overheads

- Improved CUDA Stream Support
- Time Series Support
- Conditional Joins

- Upgrade to C++17
- CUDA 11.4

# The Rapidly Growing RAPIDS Ecosystem

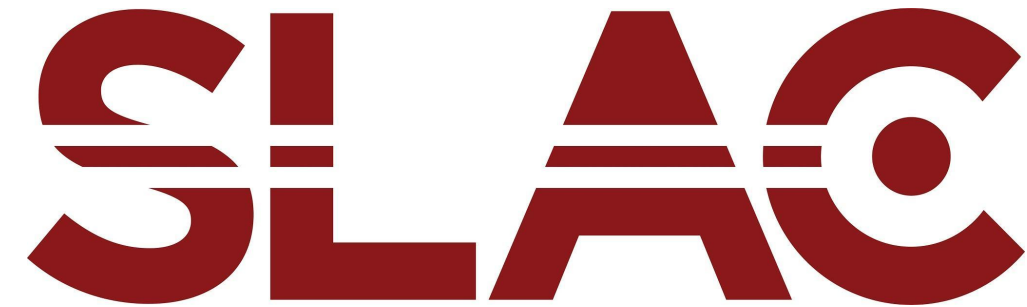## Supported, Used, & Extended by a Wide Variety of Partners

**CONTRIBUTORS**

ANACONDA · Capital One · CuPy · Chainer · Deepwave Digital · GUNROCK · NVIDIA · QuanSight · Walmart

**ADOPTERS**

anyscale · BEST BUY · BMW GROUP · Booz | Allen | Hamilton · Capital One · CLOUDERA · databricks · Deepwave Digital · DOMINO

graphistry · H2O.ai · IBM · iguazio · Informatica · kinetica · Inria · MAPR · NASA · omni·sci

Preferred Networks · PyTorch · RAY · SaturnCloud · Spotify · splunk> · Uber · URSA LABS · VW · Walmart Global Tech

**OPEN SOURCE**

APACHE ARROW · blazingSQL · CuPy · DASK · HoloViz · node JS · nuclio · Numba · scikit learn · dmlc XGBoost

NVIDIA

# Dask/RAPIDS in HPC

Supported, Used, and Extended in Research and Academia

Dask in HPC (Recording)

Dask in HEP (Recording)

# Deploy RAPIDS Everywhere
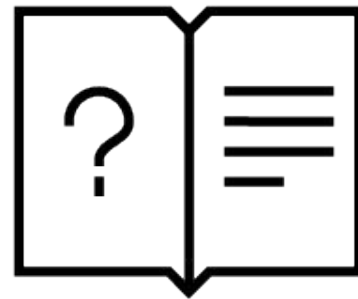## Focused on Robust Functionality, Deployment, and User Experience

Integration with major cloud providers | Both containers and cloud specific machine instances
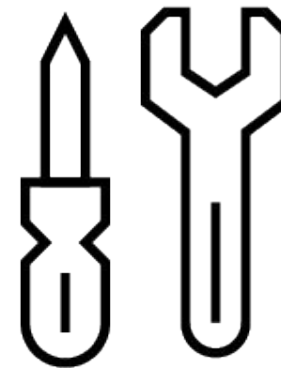Support for Enterprise and HPC Orchestration Layers

# How to Get Started with RAPIDS
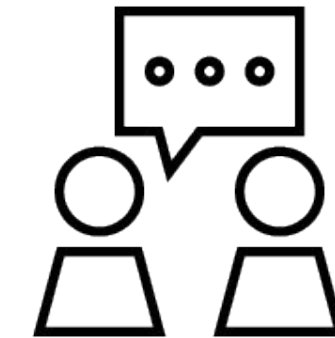
## A Variety of Ways to Get Up & Running

### More about RAPIDS

- Learn more at RAPIDS.ai
- Read the API docs
- Check out the RAPIDS blog
- Read the NVIDIA DevBlog

### Self-Start Resources

- Get started with RAPIDS
- Deploy on the Cloud today
- Start with Google Colab
- Look at the cheat sheets

### Discussion & Support

- Check the RAPIDS GitHub
- Use the NVIDIA Forums
- Reach out on Slack
- Talk to NVIDIA Services

*Keep in touch with us on the Dask Slack workspace: link here*

NVIDIA.

Let's get started!