

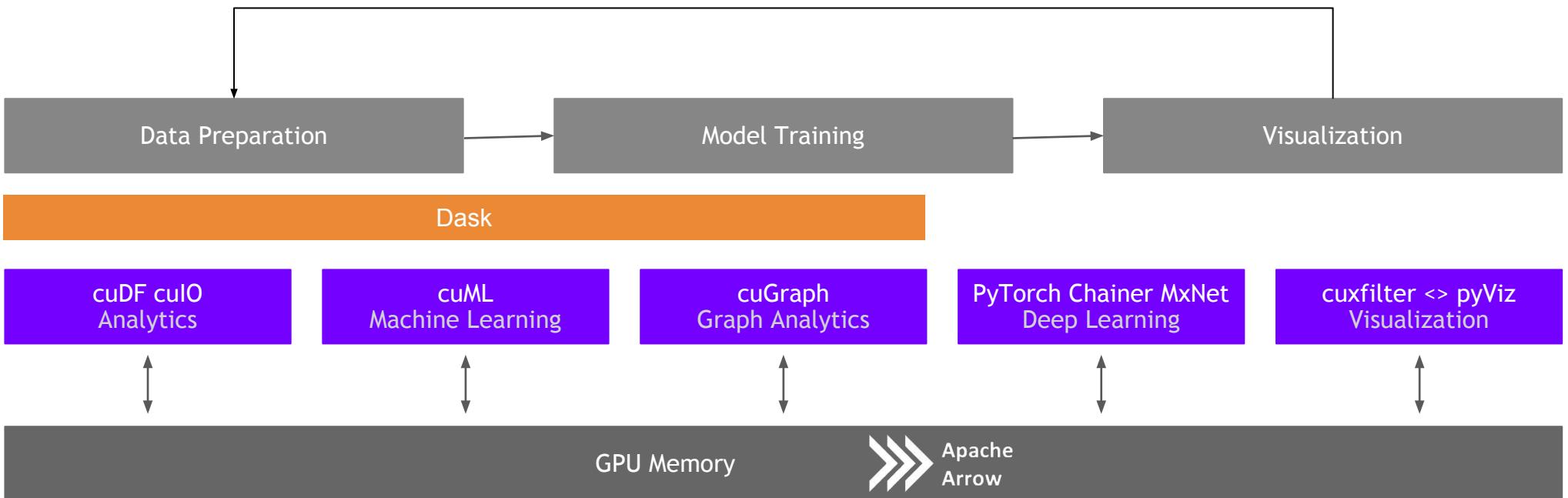
RAPIDS

The Platform Inside and Out
Release 0.13

Joshua Patterson - Director, RAPIDS Engineering

RAPIDS

End-to-End Accelerated GPU Data Science



Data Processing Evolution

Faster data access, less data movement

Hadoop Processing, Reading from disk

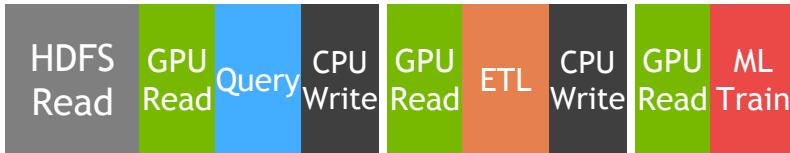


Spark In-Memory Processing



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

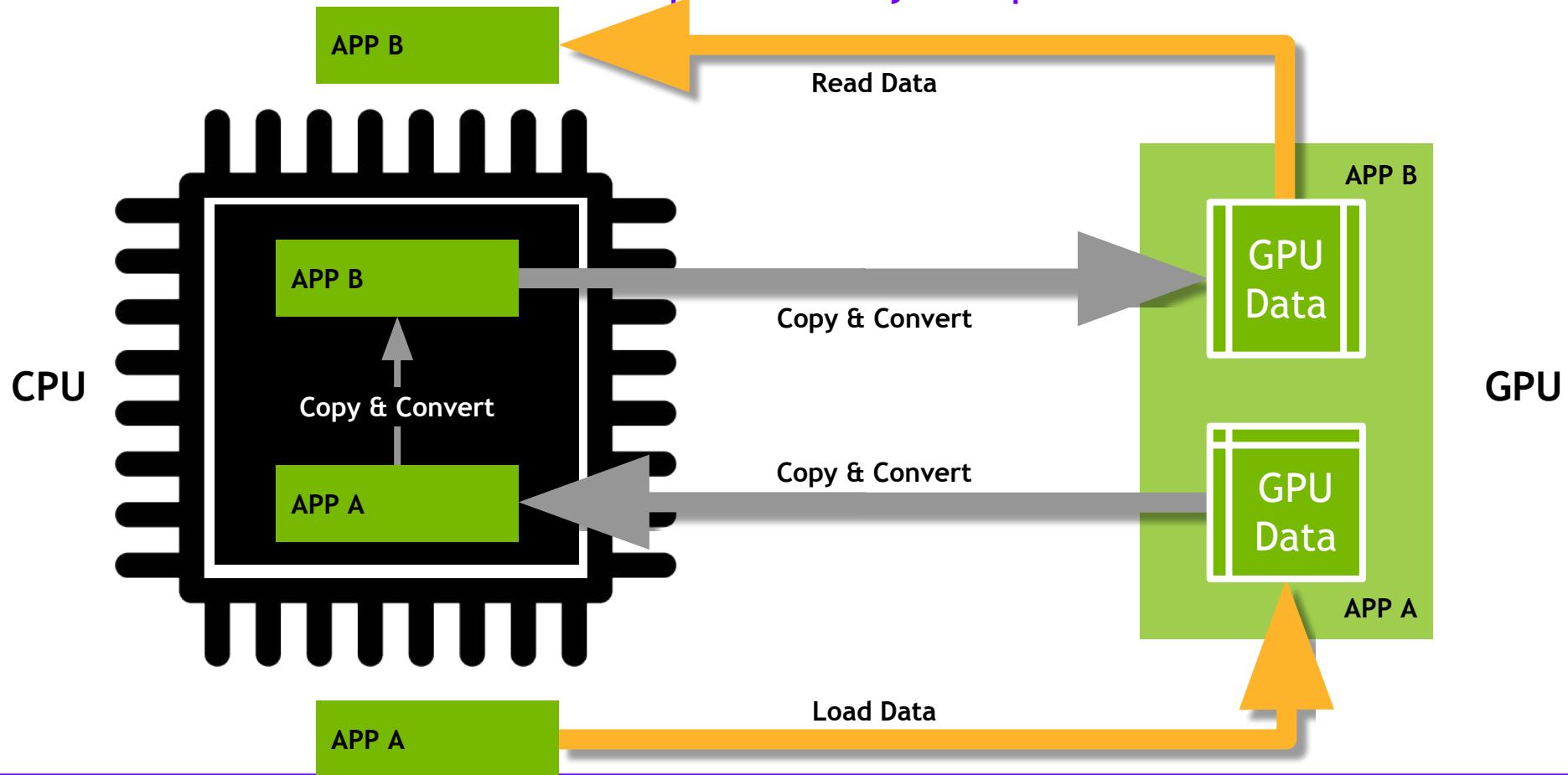
Traditional GPU Processing



5-10x Improvement
More code
Language rigid
Substantially on GPU

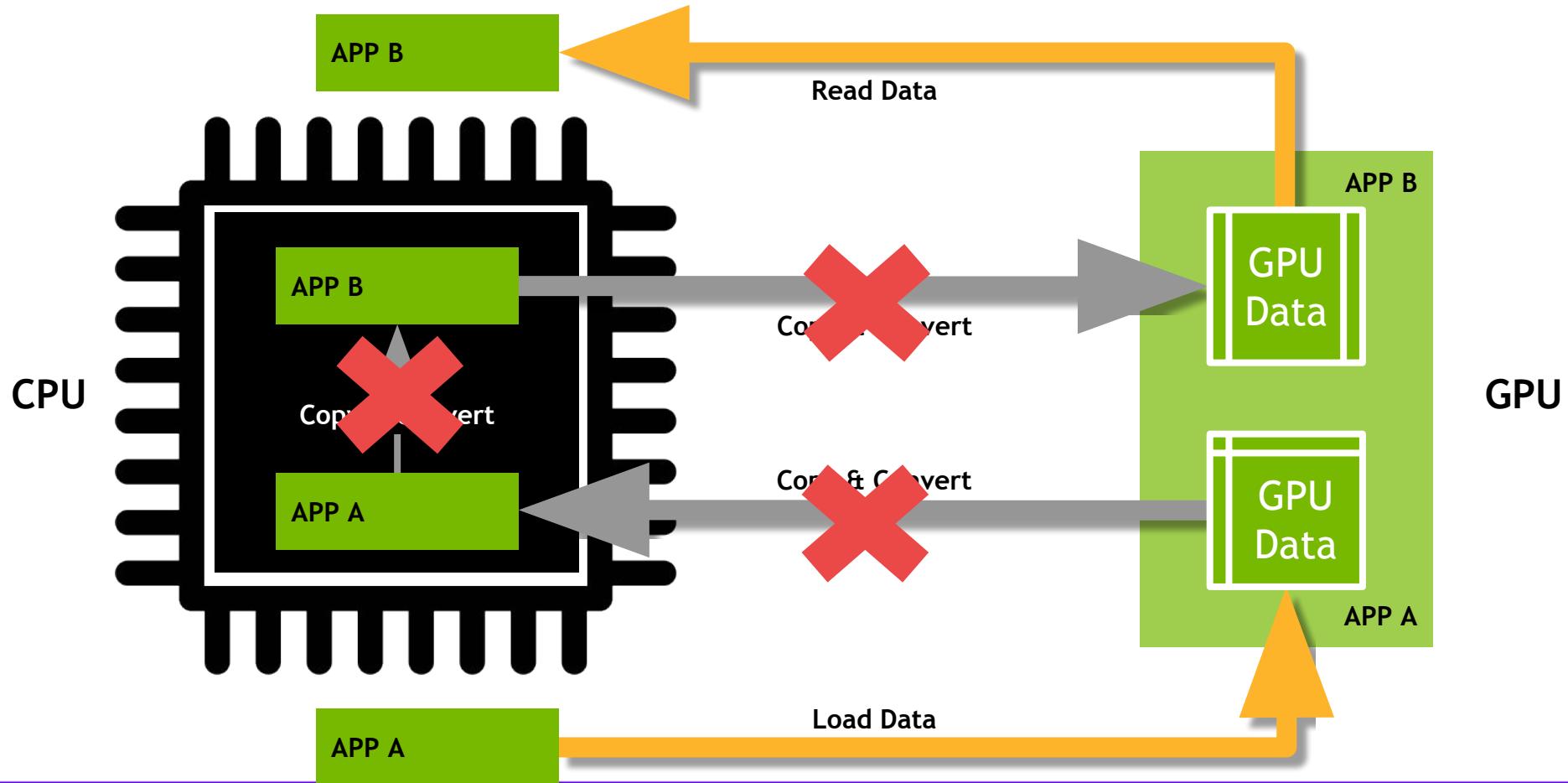
Data Movement and Transformation

The bane of productivity and performance

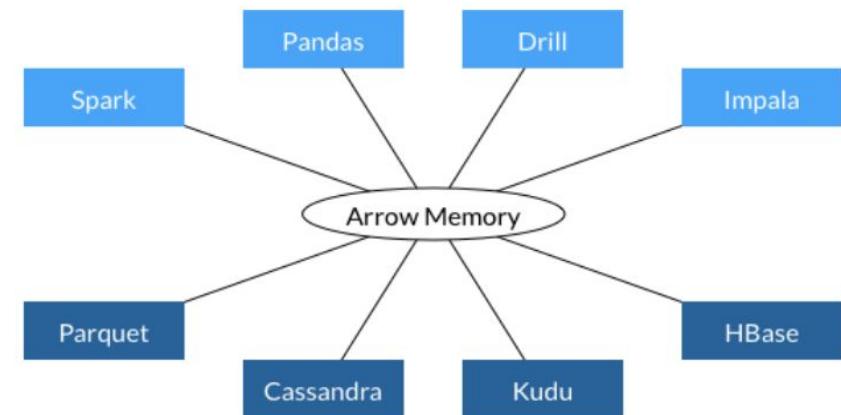
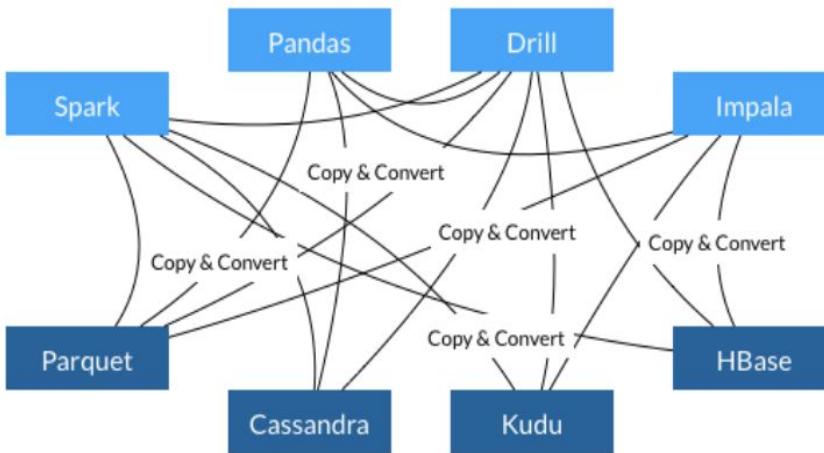


Data Movement and Transformation

What if we could keep data on the GPU?



Learning from Apache Arrow ➤➤➤



- Each system has its own internal memory format
- 70-80% computation wasted on serialization and deserialization
- Similar functionality implemented in multiple projects

- All systems utilize the same memory format
- No overhead for cross-system communication
- Projects can share functionality (eg, Parquet-to-Arrow reader)

From Apache Arrow Home Page - <https://arrow.apache.org/>

Data Processing Evolution

Faster data access, less data movement

Hadoop Processing, Reading from disk

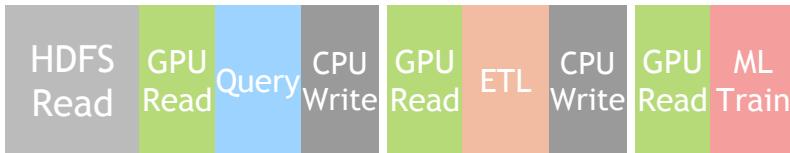


Spark In-Memory Processing



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

Traditional GPU Processing



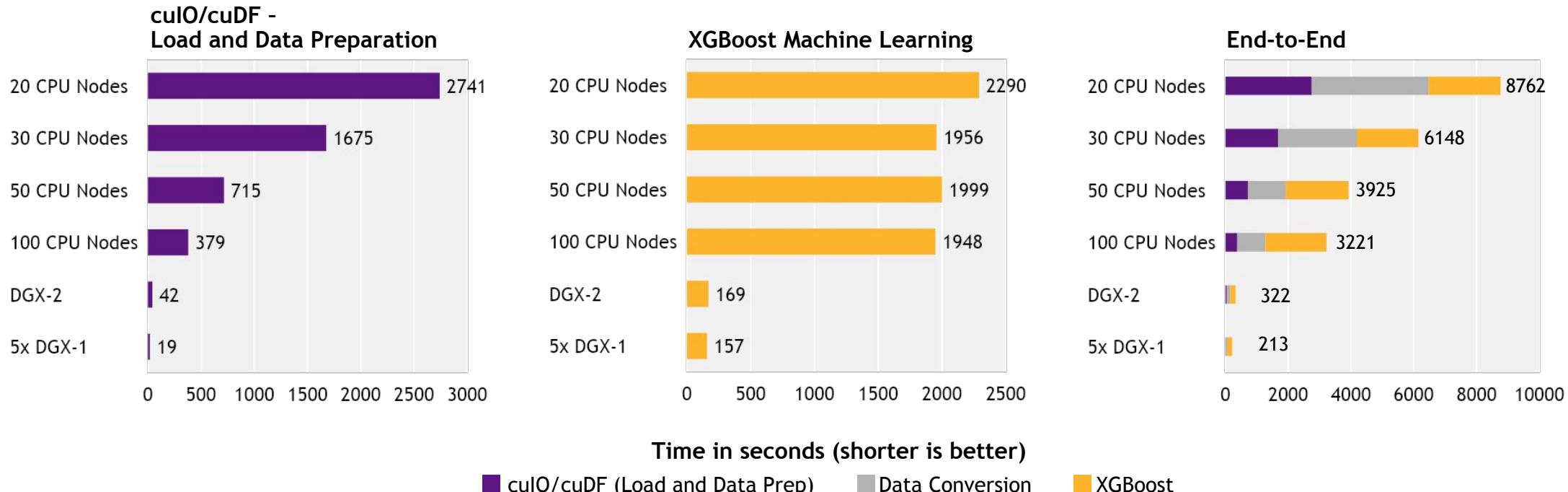
5-10x Improvement
More code
Language rigid
Substantially on GPU

RAPIDS



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

Faster Speeds, Real-World Benefits



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

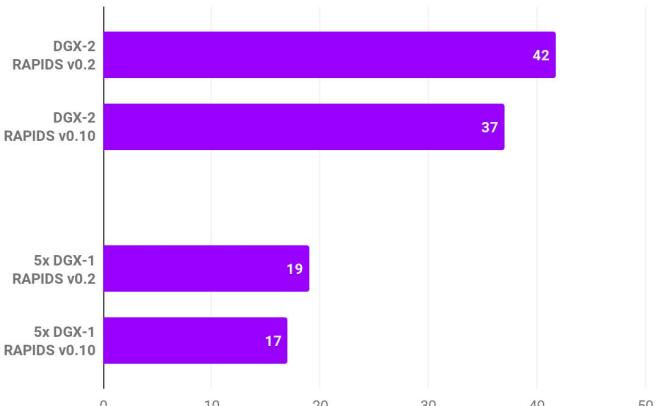
DGX Cluster Configuration

5x DGX-1 on InfiniBand network

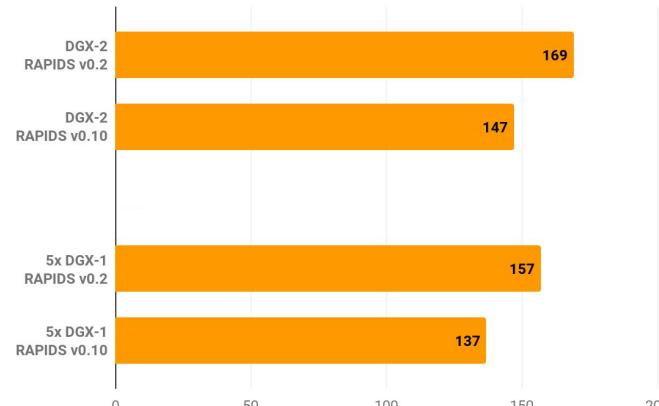
Faster Speeds, Real-World Benefits

Improving Over Time

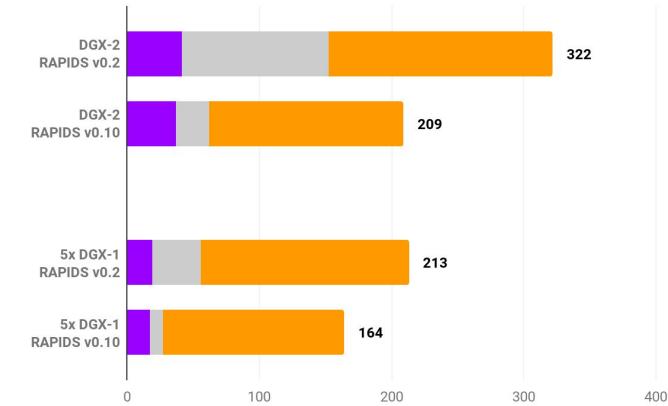
**cuIO/cuDF -
Load and Data Preparation**



XGBoost Machine Learning



End-to-End



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

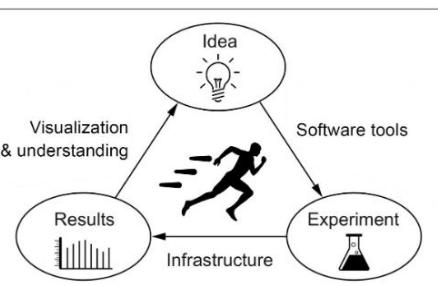
Speed, Ease of Use, and Iteration

The Way to Win at Data Science

François Chollet · @fchollet · Following

Winners are those who went through *more iterations* of the "loop of progress" -- going from an idea, to its implementation, to actionable results. So the winning teams are simply those able to run through this loop *faster*.

And this is where Keras gives you an edge.



12:31 PM - 3 April 2019

50 Retweets 158 Likes

5 50 158

François Chollet · @fchollet · Apr 3

We often talk about how following UX best practices for API design makes Keras more accessible and easier to use, and how this helps beginners. But those who stand to benefit most from good UX *aren't* the beginners. It's actually the very best practitioners in the world.

François Chollet · @fchollet · Apr 3

Because good UX reduces the overhead (development overhead & cognitive overhead) to setting up new experiments. It means you will be able to iterate faster. You will be able to try more ideas.

François Chollet · @fchollet · Apr 3

And ultimately, that's how you win competitions or get papers published.

François Chollet · @fchollet · Apr 3

So I don't think it's mere personal preference if Kaggle champions are overwhelmingly using Keras.

François Chollet · @fchollet · Apr 3

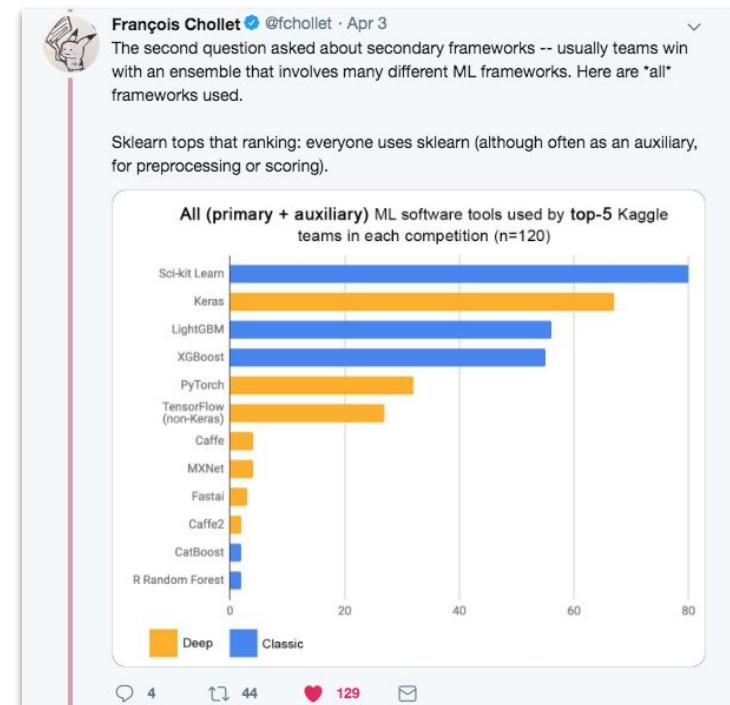
Using Keras means you're more likely to win, and inversely, those who practice the sort of fast experimentation strategy that sets them up to win are more likely to prefer Keras.

Joshua Patterson · @datametrician · Apr 3

Replying to @fchollet

This is the fundamental belief that drives @RAPIDSai. @nvidia #GPU infrastructure is fast, people need to iterate quickly, people want a known #python interface. Combine them and you're off to the races!

4 2 11



kaggle

RAPIDS 0.13 Release Summary

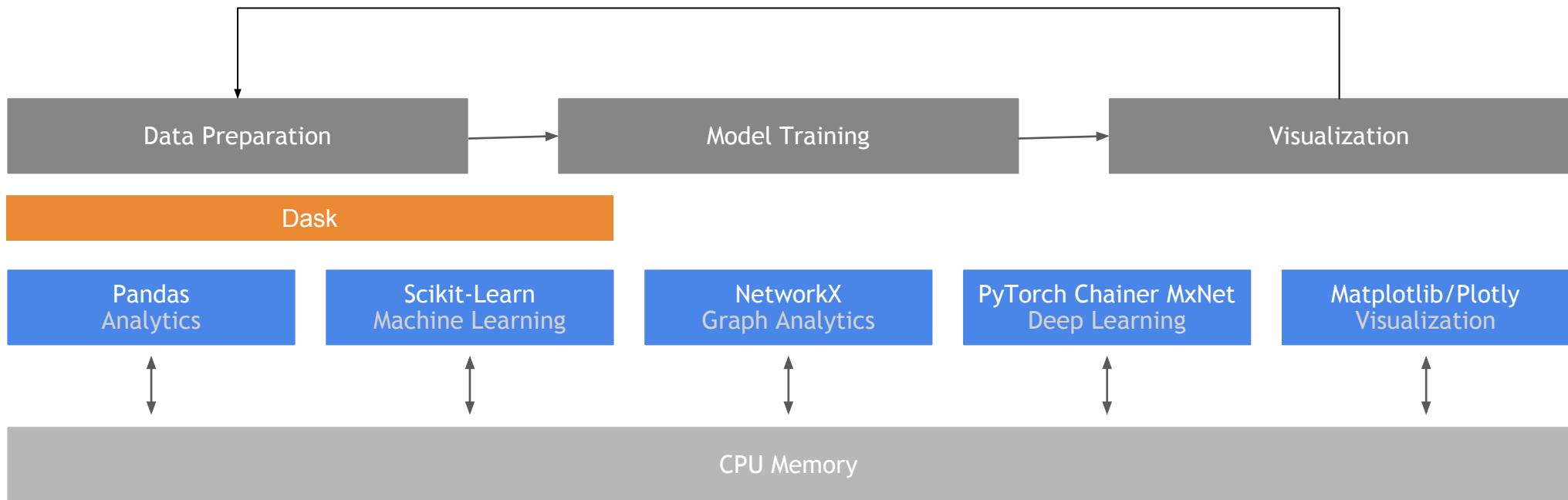
What's New in Release 0.13?

- **cuDF** Dataframes library python code base now ported to use the libcudf++ API, adds expanded groupby aggregations, join methods, concatenate optimizations, distributed multi column sorting and multi column hash repartitioning
- **cuML** machine learning library adds new Dask API to XGBoost 1.0, new configurable types for estimators, and multi-node, multi-GPU support for linear models (PCA, tSVD, OLS, Ridge)
- **cuGraph** graph analytics library adds betweenness centrality, K-Truss community detection, and code refactoring
- **UCX-Py** focused on resolving Multi-Node Multi-GPU InfiniBand tests
- **BlazingSQL** adds ROUND(), CASE with Strings, and for distributed queries AVG() support
- **cuSignal** adds conda install support, polyphase resampler, and new acoustics module
- **cuSpatial** adds batched cubic spline interpolation for trajectories and major documentation improvements

RAPIDS Core

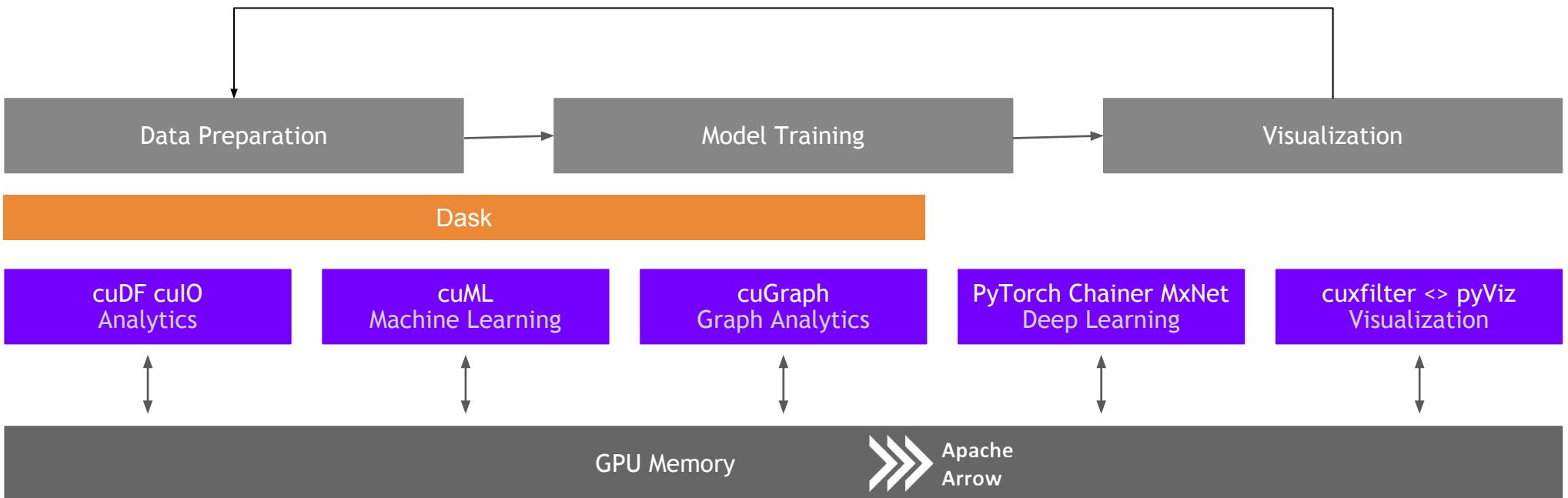
Open Source Data Science Ecosystem

Familiar Python APIs



RAPIDS

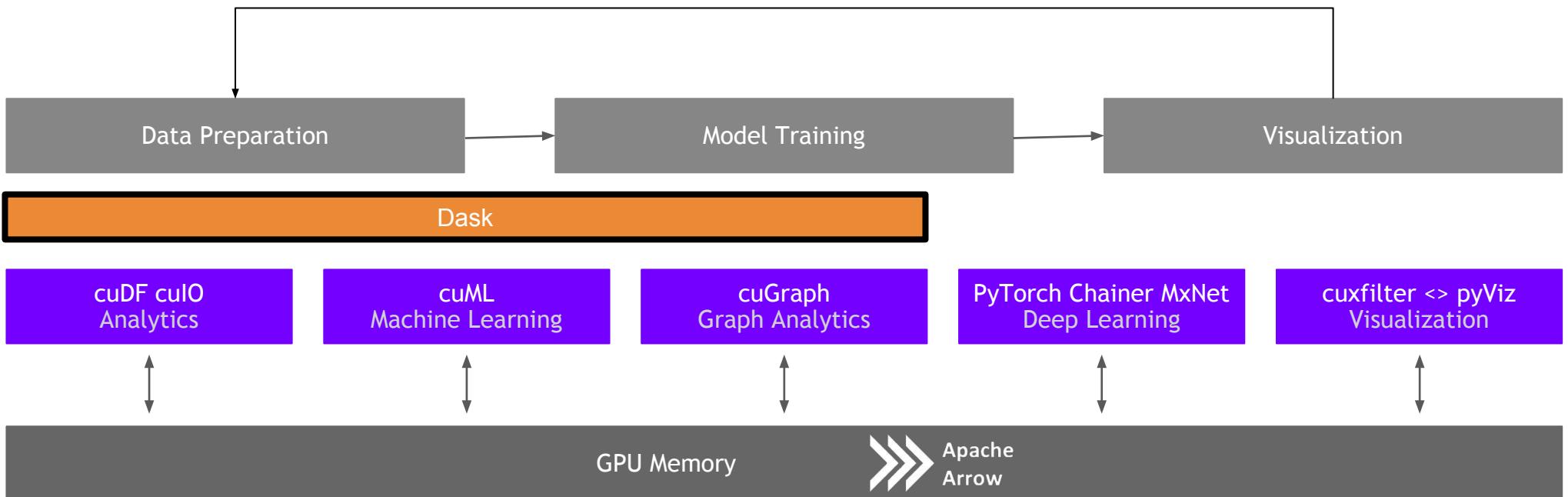
End-to-End Accelerated GPU Data Science



Dask

RAPIDS

Scaling RAPIDS with Dask



Why Dask?

PyData Native

- **Easy Migration:** Built on top of NumPy, Pandas, Scikit-Learn, etc.
- **Easy Training:** With the same APIs
- **Trusted:** With the same developer community

Deployable

- HPC: SLURM, PBS, LSF, SGE
- Cloud: Kubernetes
- Hadoop/Spark: Yarn



Easy Scalability

- Easy to install and use on a laptop
- Scales out to thousand-node clusters

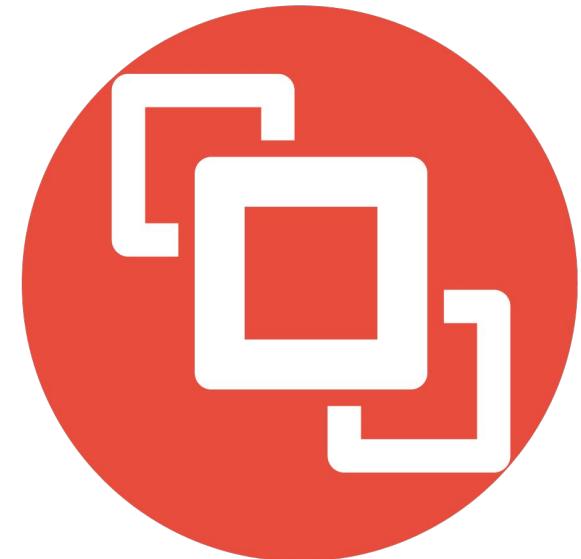
Popular

- Most common parallelism framework today in the PyData and SciPy community

Why OpenUCX?

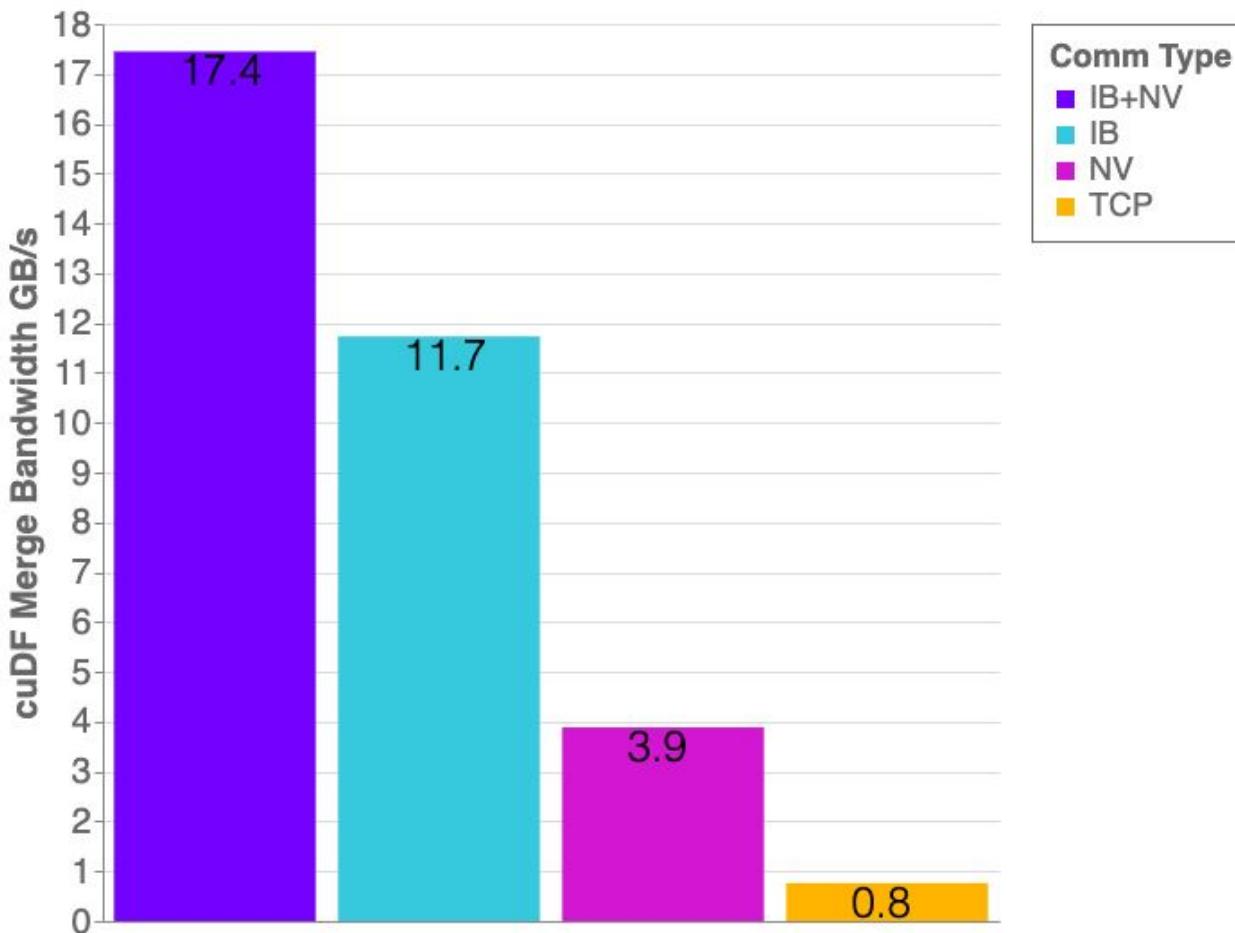
Bringing hardware accelerated communications to Dask

- TCP sockets are slow!
- UCX provides uniform access to transports (TCP, InfiniBand, shared memory, NVLink)
- Alpha Python bindings for UCX (ucx-py)
- Will provide best communication performance, to Dask based on available hardware on nodes/cluster



```
conda install -c conda-forge -c rapidsai \
    cudatoolkit=<CUDA version> ucx-proc=*gpu ucx ucx-py
```

Benchmarks: Distributed cuDF Random Merge



cuDF v0.13, UCX-PY 0.13

Running on NVIDIA DGX-1 (8GPUs):

GPU: NVIDIA Tesla V100 32GB
CPU: Intel(R) Xeon(R) CPU 8168
@ 2.70GHz

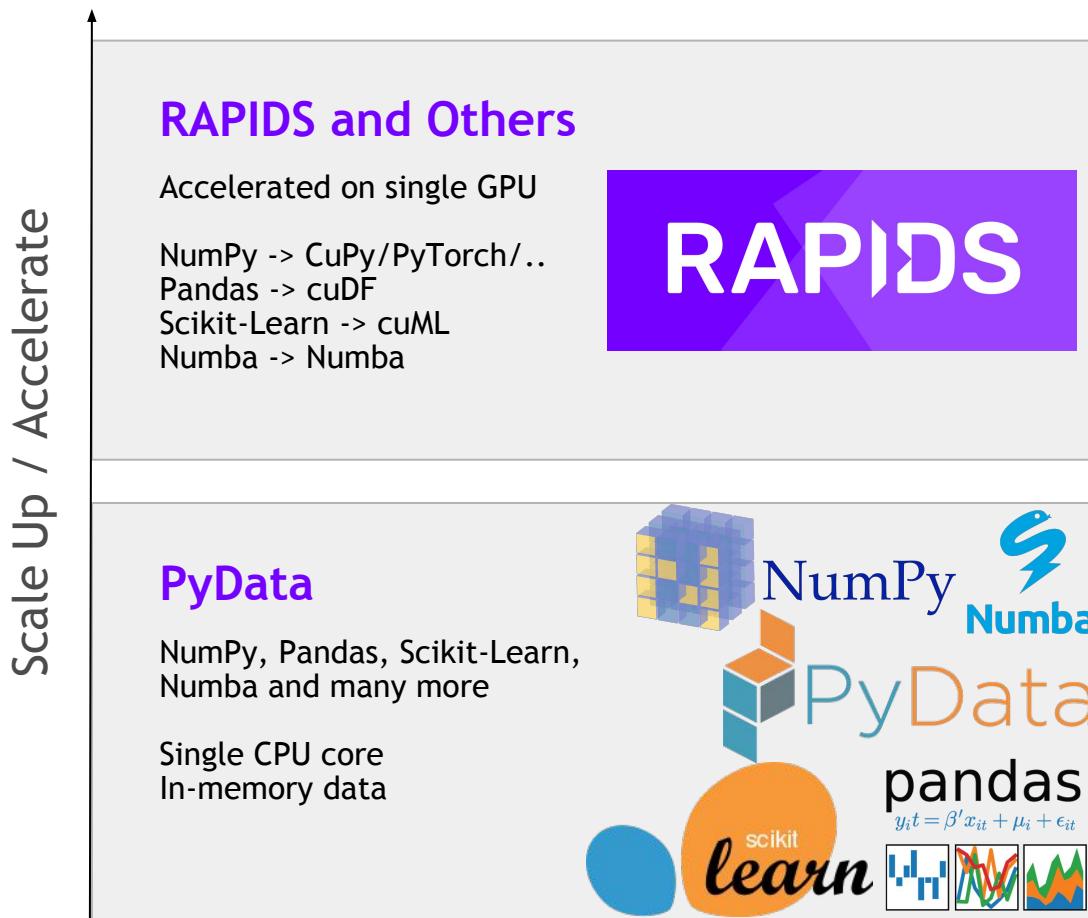
Benchmark Setup:

DataFrames: Left/Right 1x int64 column key column, 1x int64 value columns

Merge: inner

30% of matching data balanced across each partition

Scale up with RAPIDS



Scale out with RAPIDS + Dask with OpenUCX

Scale Up / Accelerate

RAPIDS and Others

Accelerated on single GPU

NumPy -> CuPy/PyTorch/..
Pandas -> cuDF
Scikit-Learn -> cuML
Numba -> Numba



RAPIDS + Dask with OpenUCX

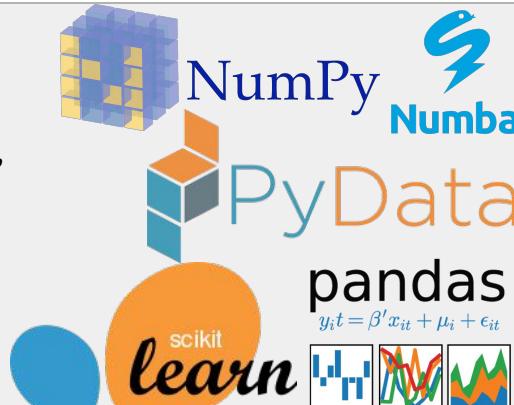
Multi-GPU
On single Node (DGX)
Or across a cluster



PyData

NumPy, Pandas, Scikit-Learn,
Numba and many more

Single CPU core
In-memory data



Dask

Multi-core and Distributed PyData

NumPy -> Dask Array
Pandas -> Dask DataFrame
Scikit-Learn -> Dask-ML
... -> Dask Futures

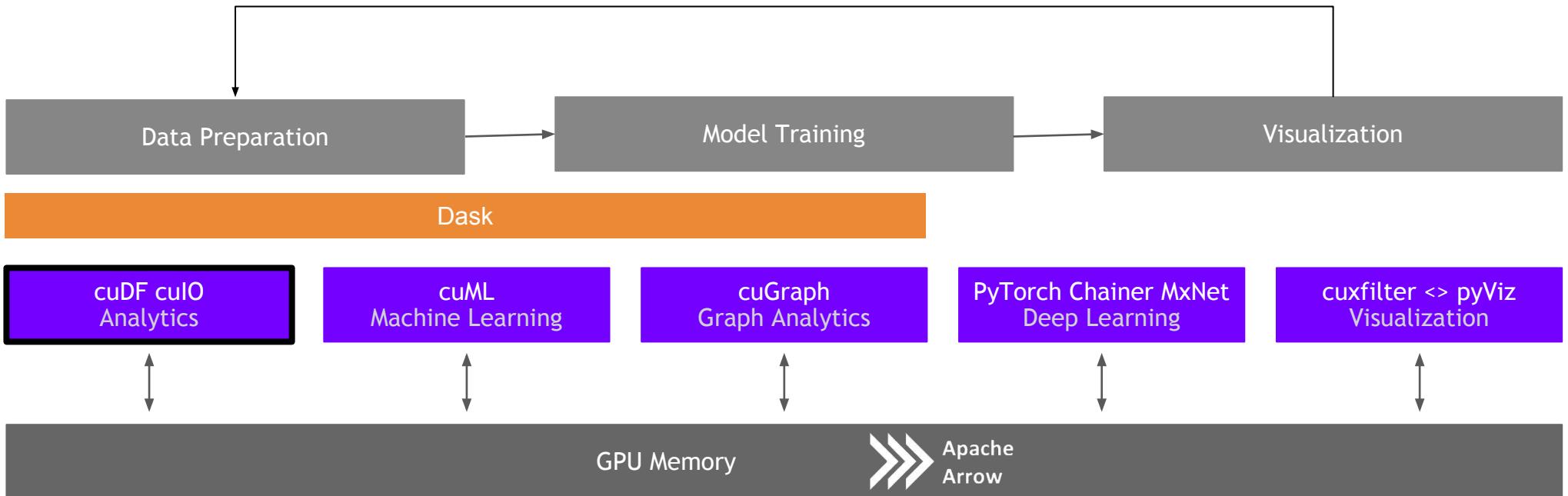


Scale out / Parallelize

cuDF

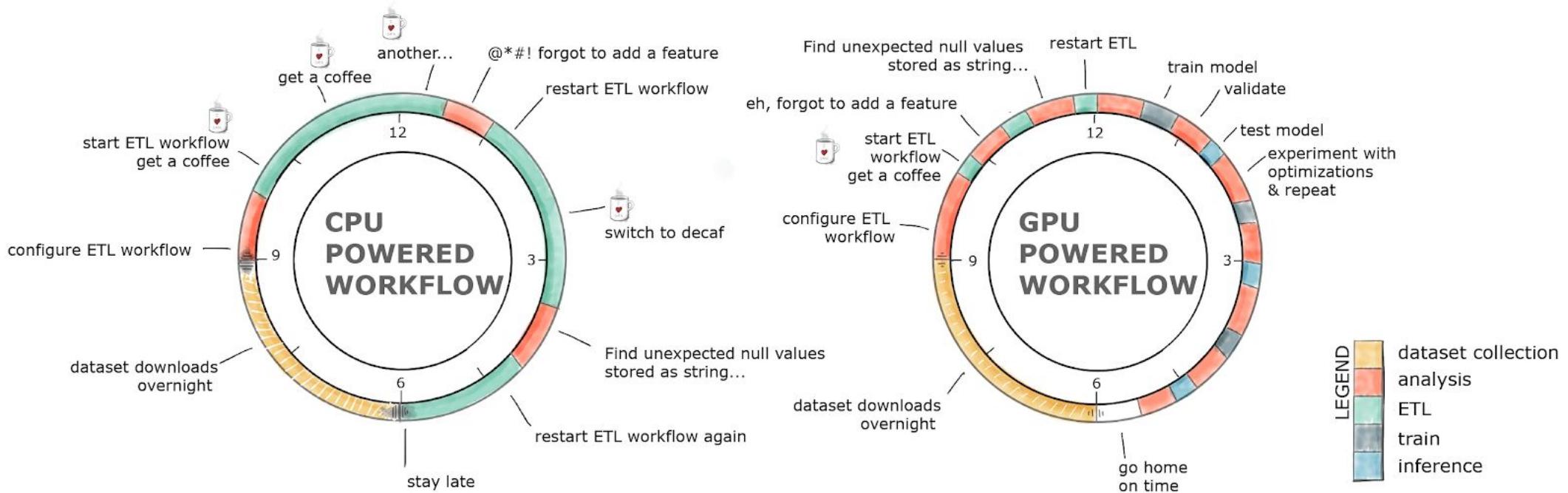
RAPIDS

GPU Accelerated data wrangling and feature engineering

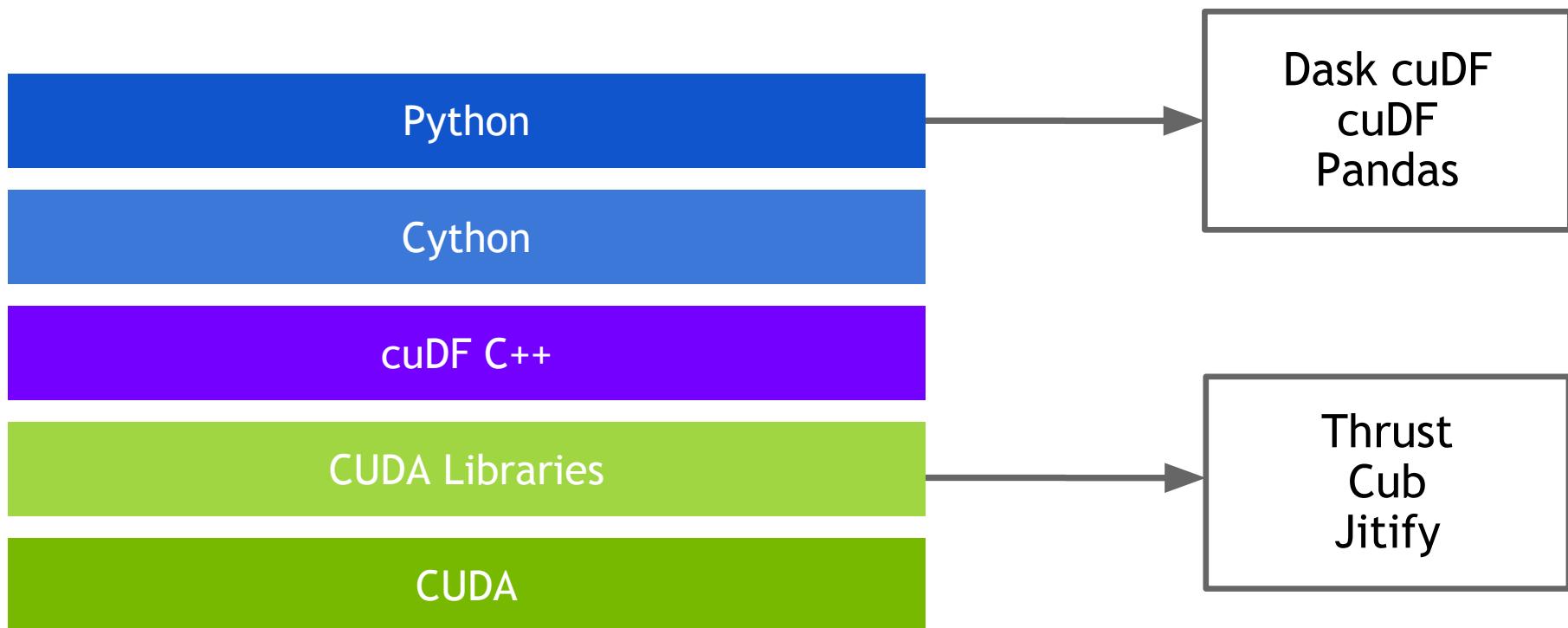


GPU-Accelerated ETL

The average data scientist spends 90+% of their time in ETL as opposed to training models



ETL Technology Stack



ETL: the Backbone of Data Science

libcuDF is...

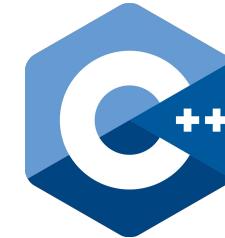
CUDA C++ Library

- Table (dataframe) and column types and algorithms
- CUDA kernels for sorting, join, groupby, reductions, partitioning, elementwise operations, etc.
- Optimized GPU implementations for strings, timestamps, numeric types (more coming)
- Primitives for scalable distributed ETL

```
std::unique_ptr




```



ETL: the Backbone of Data Science

cuDF is...

```
In [2]: #Read in the data. Notice how it decompresses as it reads the data into memory.  
gdf = cudf.read_csv('/rapids/Data/black-friday.zip')
```

```
In [3]: #Taking a look at the data. We use "to_pandas()" to get the pretty printing.  
gdf.head().to_pandas()
```

Out[3]:

	User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Cat
0	1000001	P00069042	F	0-17	10	A	2	0	3
1	1000001	P00248942	F	0-17	10	A	2	0	1
2	1000001	P00087842	F	0-17	10	A	2	0	12
3	1000001	P00085442	F	0-17	10	A	2	0	12
4	1000002	P00285442	M	55+	16	C	4+	0	8

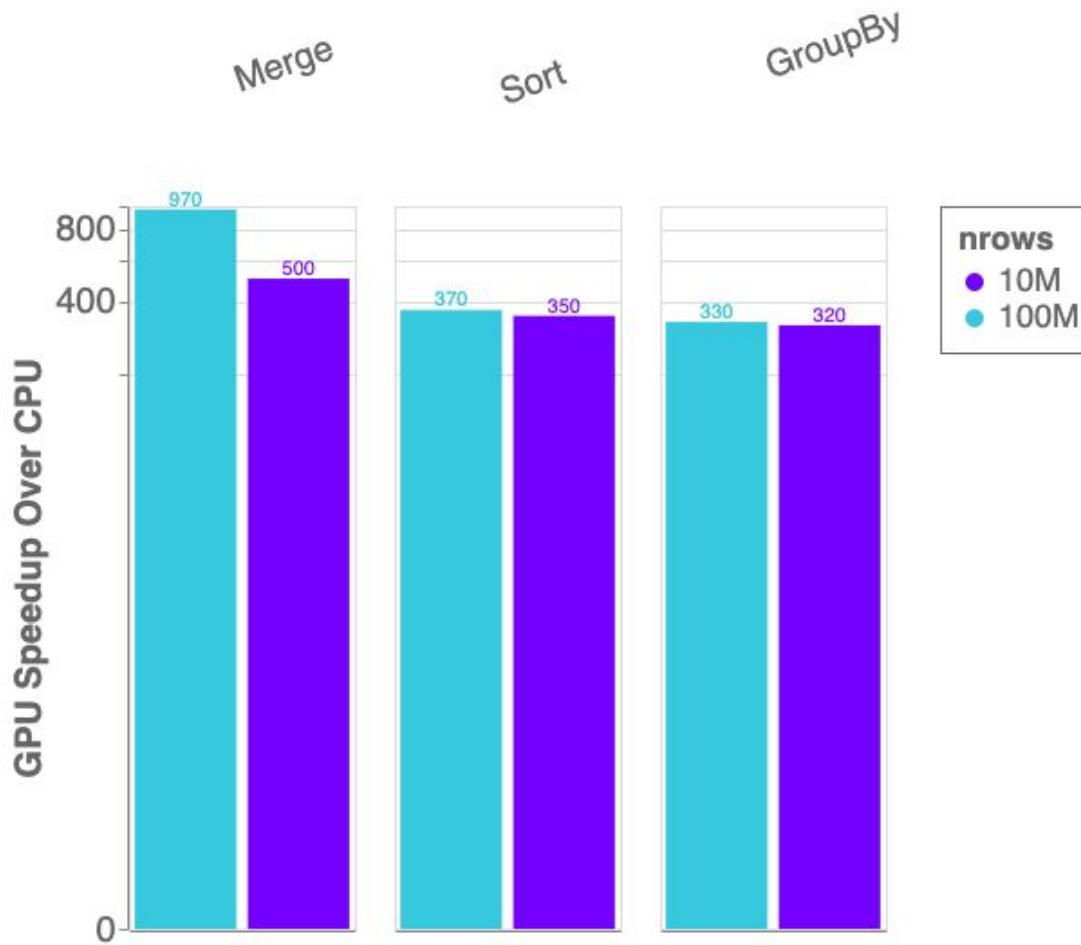
```
In [6]: #grabbing the first character of the years in city string to get rid of plus sign, and converting  
#to int  
gdf['city_years'] = gdf.Stay_In_Current_City_Years.str.get(0).stoi()
```

```
In [7]: #Here we can see how we can control what the value of our dummies with the replace method and turn  
#strings to ints  
gdf['City_Category'] = gdf.City_Category.str.replace('A', '1')  
gdf['City_Category'] = gdf.City_Category.str.replace('B', '2')  
gdf['City_Category'] = gdf.City_Category.str.replace('C', '3')  
gdf['City_Category'] = gdf['City_Category'].str.stoi()
```

Python Library

- A Python library for manipulating GPU DataFrames following the Pandas API
- Python interface to CUDA C++ library with additional functionality
- Creating GPU DataFrames from Numpy arrays, Pandas DataFrames, and PyArrow Tables
- JIT compilation of User-Defined Functions (UDFs) using Numba

Benchmarks: single-GPU Speedup vs. Pandas



cuDF v0.13, Pandas 0.25.3

Running on NVIDIA DGX-1:

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

Benchmark Setup:

RMM Pool Allocator Enabled

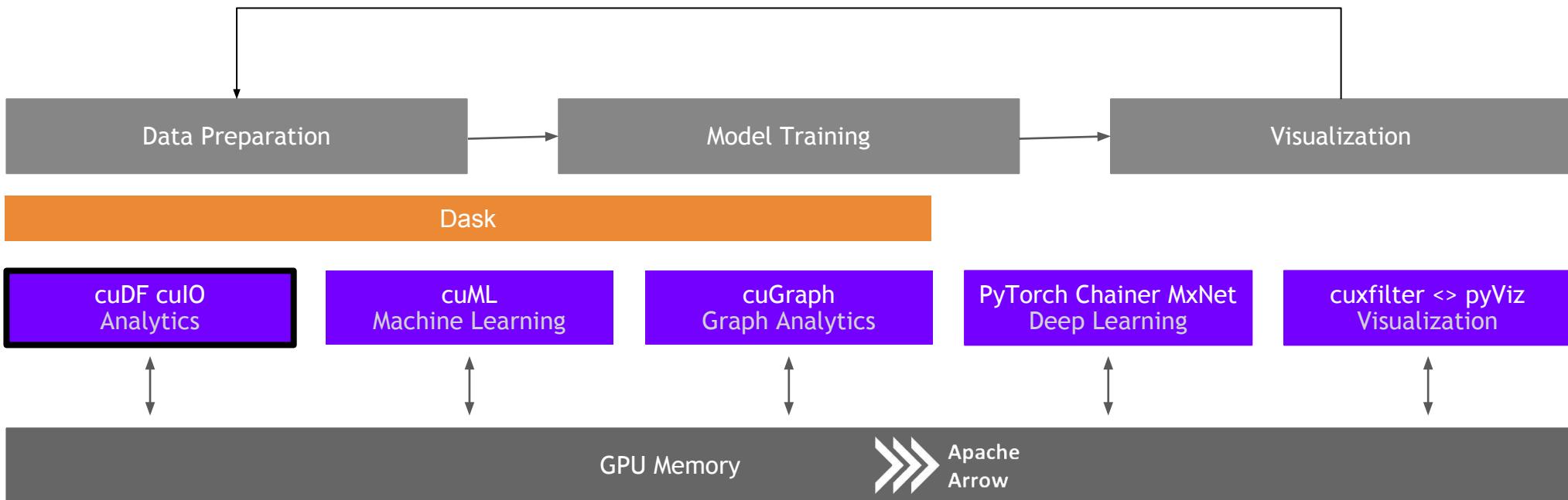
DataFrames: 2x int32 columns key columns,
3x int32 value columns

Merge: inner

GroupBy: count, sum, min, max calculated
for each value column

ETL: the Backbone of Data Science

cuDF is not the end of the story



ETL: the Backbone of Data Science

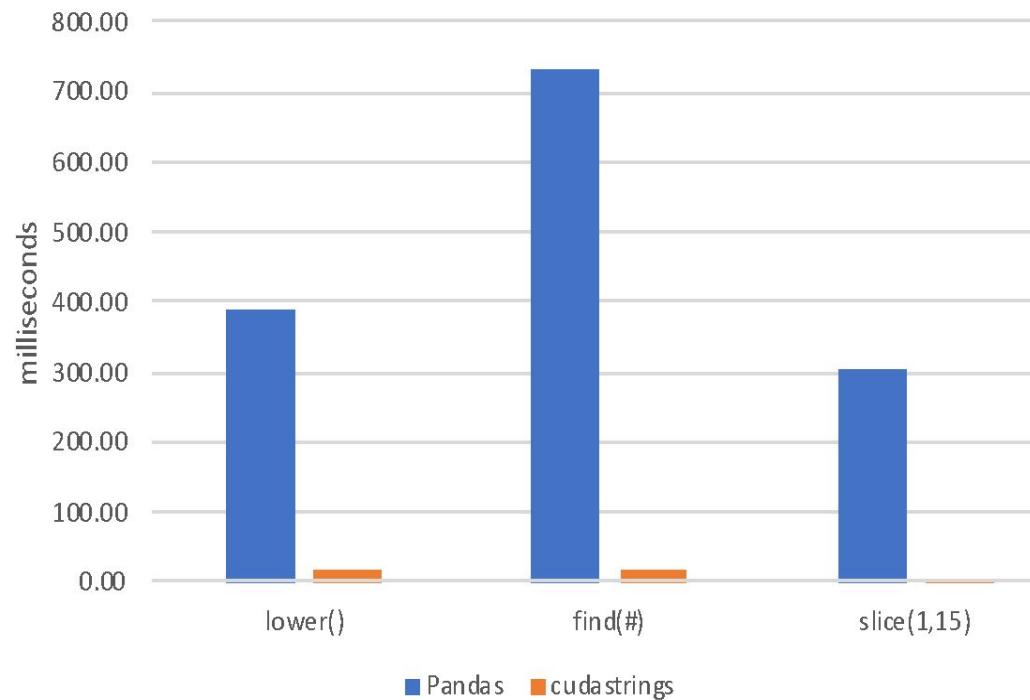
String Support

Current v0.13 String Support

- Regular Expressions
- Element-wise operations
 - Split, Find, Extract, Cat, Typecasting, etc...
- String GroupBys, Joins, Sorting, etc.
- Categorical columns fully on GPU
- Native String type in libcudf C++

Future v0.14+ String Support

- Further performance optimization
- JIT-compiled String UDFs



Extraction is the Cornerstone

cuDF I/O for Faster Data Loading

- Follow Pandas APIs and provide >10x speedup
- CSV Reader - v0.2, CSV Writer v0.8
- Parquet Reader - v0.7, Parquet Writer v0.12
- ORC Reader - v0.7, ORC Writer v0.10
- JSON Reader - v0.8
- Avro Reader - v0.9
- GPU Direct Storage integration in progress for bypassing PCIe bottlenecks!
- Key is GPU-accelerating both parsing and decompression wherever possible

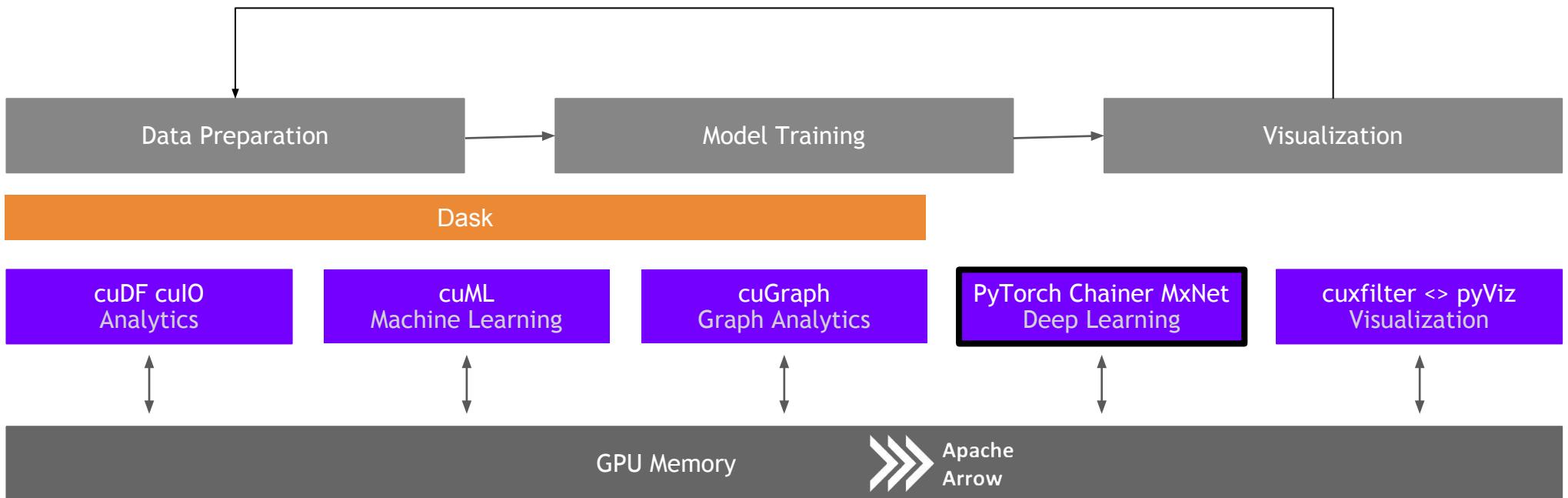
```
1]: import pandas, cudf
2]: %time len(pandas.read_csv('data/nyc/yellow_tripdata_2015-01.csv'))
CPU times: user 25.9 s, sys: 3.26 s, total: 29.2 s
Wall time: 29.2 s
2]: 12748986
3]: %time len(cudf.read_csv('data/nyc/yellow_tripdata_2015-01.csv'))
CPU times: user 1.59 s, sys: 372 ms, total: 1.96 s
Wall time: 2.12 s
3]: 12748986
4]: !du -hs data/nyc/yellow_tripdata_2015-01.csv
1.9G    data/nyc/yellow_tripdata_2015-01.csv
```

Source: Apache Crail blog: [SQL Performance: Part 1 - Input File Formats](#)

ETL is not just DataFrames!

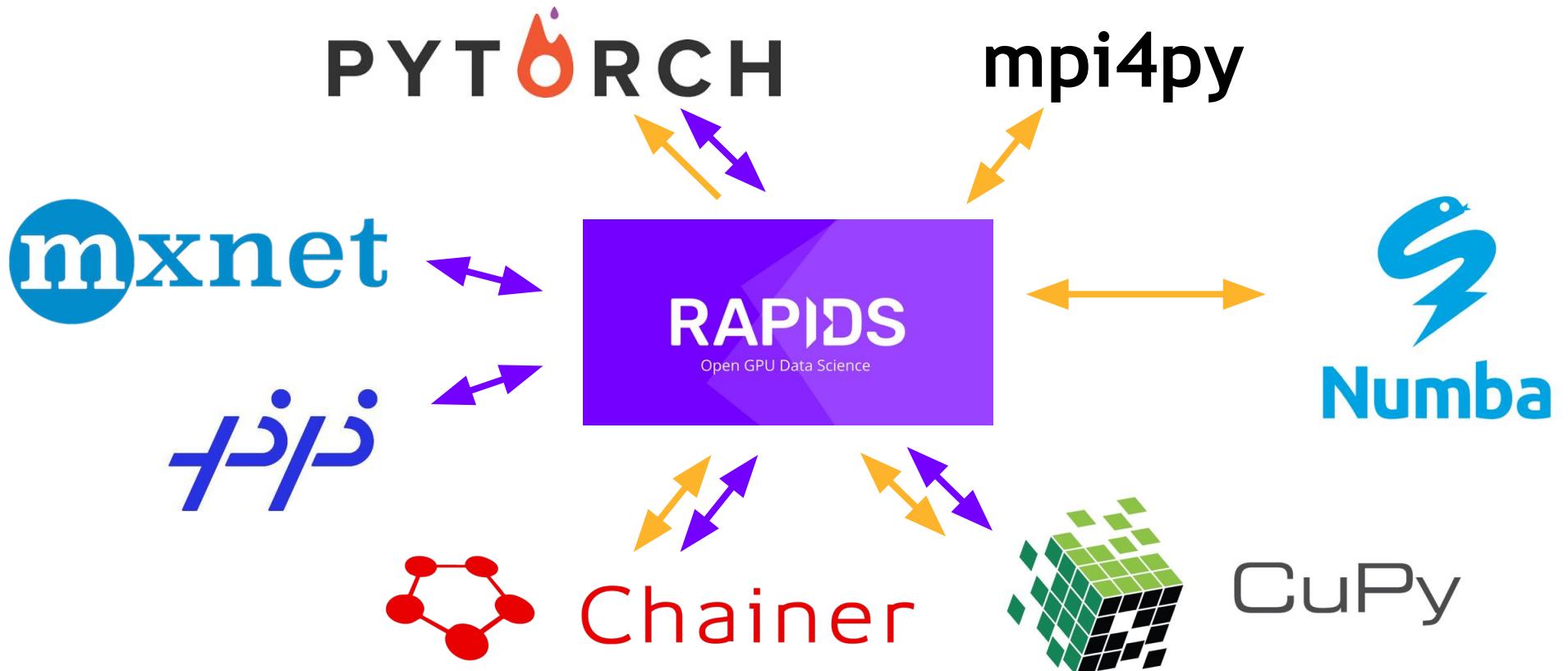
RAPIDS

Building bridges into the array ecosystem



Interoperability for the Win

DLPack and `__cuda_array_interface__`



Interoperability for the Win

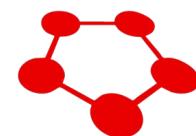
DLPack and __cuda_array_interface__

P Y T  R C H

m p i 4 p y

 m x n e t





Chainer



 Numba

CuPy

ETL: Arrays and DataFrames

Dask and CUDA Python arrays



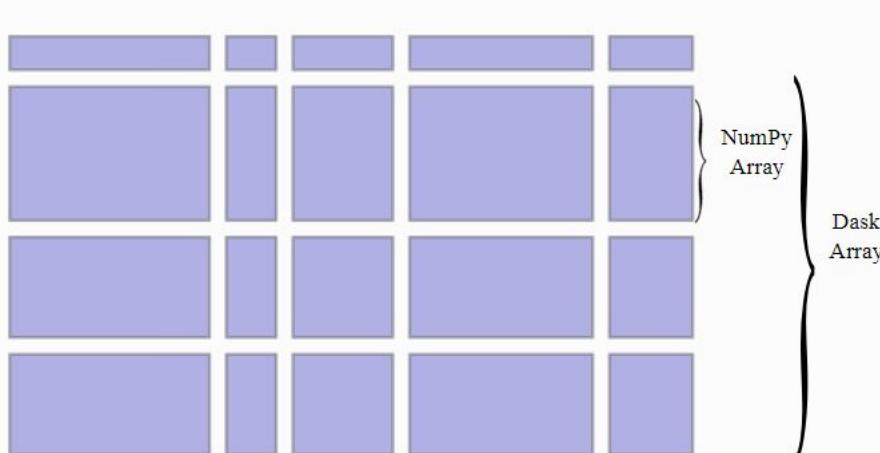
Chainer



CuPy

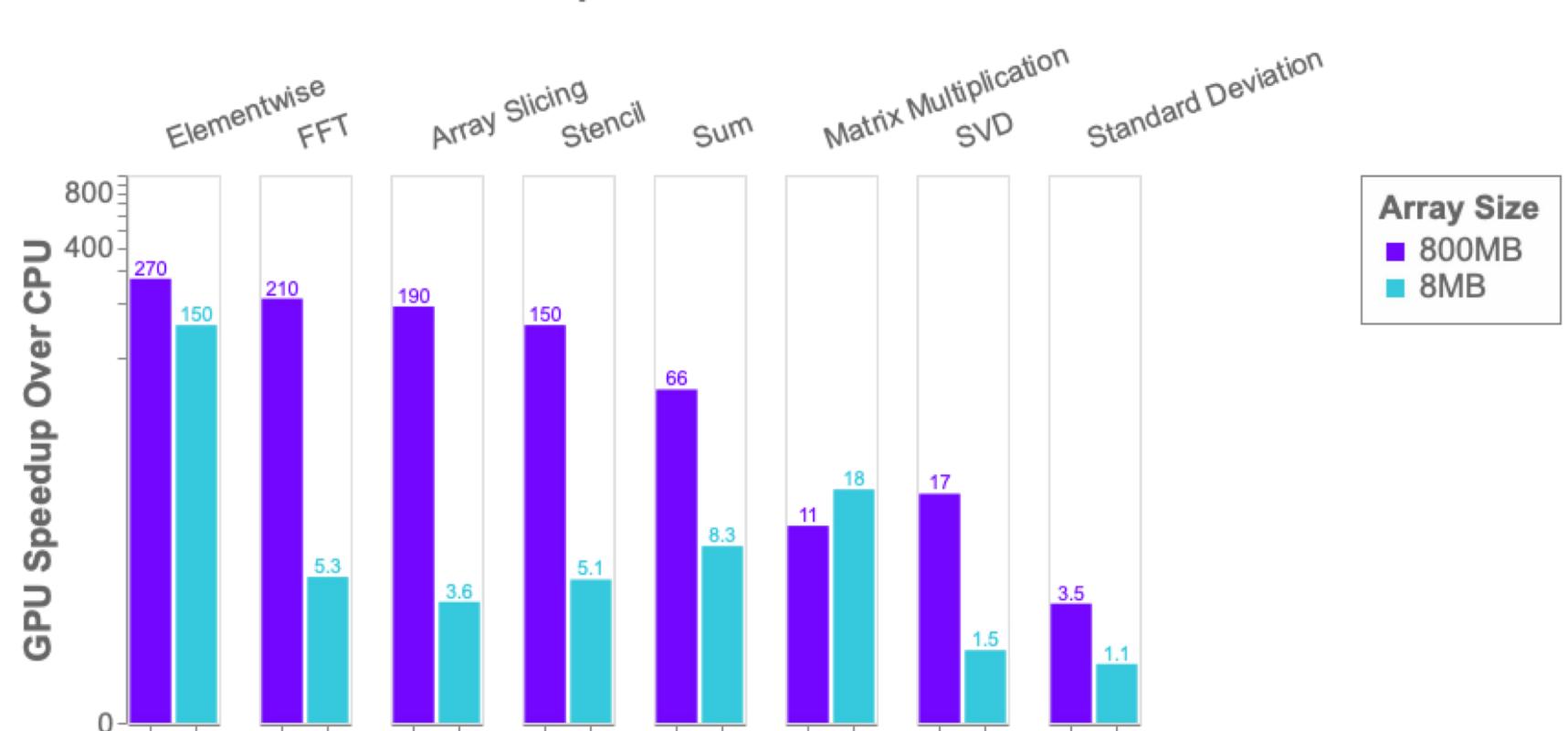


Numba



- Scales NumPy to distributed clusters
- Used in climate science, imaging, HPC analysis up to 100TB size
- Now seamlessly accelerated with GPUs

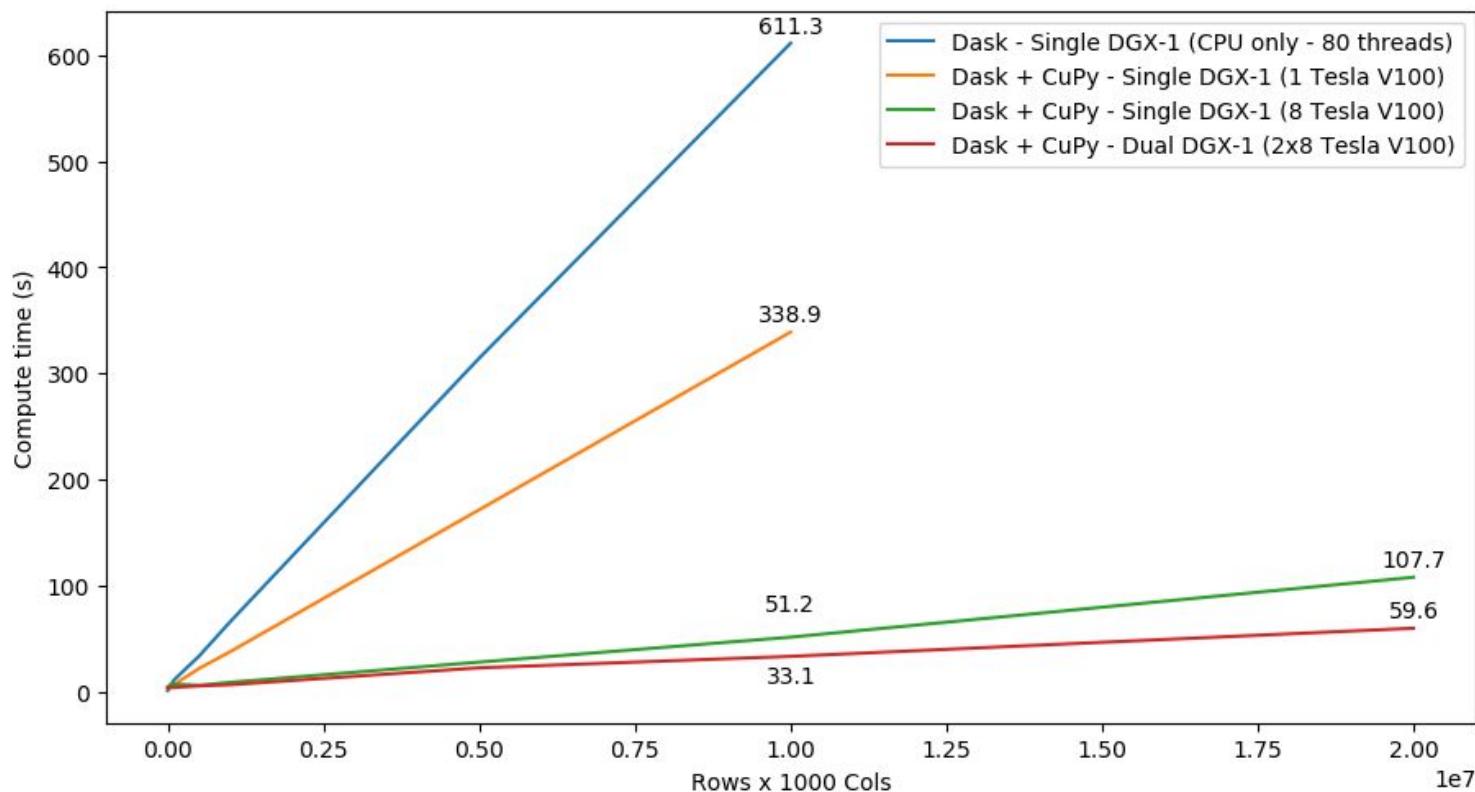
Benchmark: single-GPU CuPy vs NumPy



More details: <https://blog.dask.org/2019/06/27/single-gpu-cupy-benchmarks>

SVD Benchmark

Dask and CuPy Doing Complex Workflows



Petabyte Scale Analytics with Dask and CuPy

Architecture	Time
Single CPU Core	2hr 39min
Forty CPU Cores	11min 30s
One GPU	1min 37s
Eight GPUs	19s

<https://blog.dask.org/2019/01/03/dask-array-gpus-first-steps>



3.2 PETABYTES IN LESS THAN 1 HOUR

Distributed GPU array | parallel reduction | using 76x GPUs

Array size	Wall Time (data creation + compute)
3.2 PB (20M x 20M doubles)	54 min 51 s

Cluster configuration: 20x GCP instances, each instance has:
CPU: 1 VM socket (Intel Xeon CPU @ 2.30GHz), 2-core, 2 threads/core, 132GB mem, GbE ethernet, 950 GB disk
GPU: 4x NVIDIA Tesla P100-16GB-PCIe (total GPU DRAM across nodes 1.22 TB)
Software: Ubuntu 18.04, RAPIDS 0.5.1, Dask=1.1.1, Dask-Distributed=1.1.1, CuPY=5.2.0, CUDA 10.0.130

ETL: Arrays and DataFrames

More Dask Awesomeness from RAPIDS

RAPIDS-Dask and cuDF NYCTaxi Screencast

```
[3]: import dask_cudf
# df = dask_cudf.read_csv('gs://anaconda-public-data/nyc-taxi/csv//2015/yellow_*.csv')
df = dask_cudf.read_csv('data/nyc/yellow_tripdata_2015-*.csv')
df = df.persist()

Time full-pass computation

Most of the time here is spent reading data from disk and parsing it.

[4]: %time df.passenger_count.sum().compute()
CPU times: user 100 ms, sys: 28 ms, total: 128 ms
Wall time: 137 ms
[4]: 245566747

[5]: %time df.groupby('passenger_count').trip_distance.mean().compute()
[5]: ...to_pandas()

How well do people tip?

[6]: df2 = df[['trip_pickup_datetime', 'trip_distance', 'tip_amount', 'fare_amount']]
df2 = df2.query('tip_amount > 0')
df2['hour'] = df2.trip_pickup_datetime.dt.hour.astype('int32')

df2['tip_fraction'] = df2.tip_amount / df2.fare_amount
hour = df2.groupby('hour').tip_fraction.mean().compute().to_pandas()

[7]: %matplotlib inline
hour.plot(figsize=(10, 6), title="Tip Fraction by Hour")

Dask Progress
Progress - total: 1032, in-memory: 93, processing: 1, errred: 0
assign 182 / 182
series-group... 91 / 91
getitem 182 / 182
series-group... 14 / 14
read_csv 91 / 91
series-group... 14 / 14
getitem-read...
query 91 / 91
series-group... 1 / 1
astype-dt-ho...
91 / 91
truediv-session 91 / 91
series-group... 91 / 91
```

3:13 / 4:10

RAPIDS-Dask and CuPy SVD Screencast

```
[2]: import dask
import dask.array
import numpy
import cupy

rs = dask.array.random(1000000, 1000, chunks=(1000, 1000))
x = rs.random(1000000, 1000, chunks=(1000, 1000))
x = x.persist()

Create Random Dataset

[3]: import dask.array.linalg
u, s, v = dask.array.linalg.svd(x)

[4]: dask.visualize(u, s, v)

[4]: u, s, v = dask.persist(u, s, v)

Inspect output

[1]: print(u[10, :10].compute())
print(s[10].compute())
print(v[10, :10].compute())

[1]: fro.distributed.utils import format_bytes

[1]: print(format_bytes(u nbytes))
print(format_bytes(s nbytes))
print(format_bytes(v nbytes))

Dask Progress
Progress - total: 756, in-memory: 205, processing: 12, errred: 0
getitem 424 / 434
dot 100 / 100
series-group... 100 / 100
nr 100 / 100
stack-ds 11 / 11
getitem-in... 1 / 1
```

3:12 / 3:42

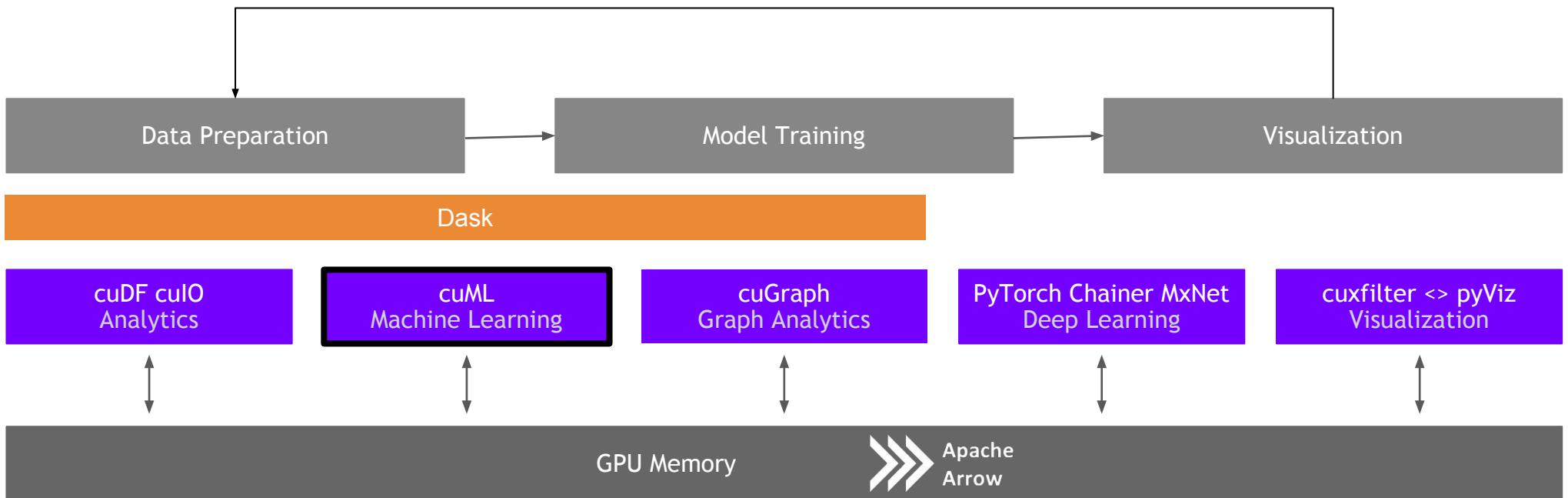
<https://youtu.be/gV0cykgsTPM>

https://youtu.be/R5CiXti_MWo

cuML

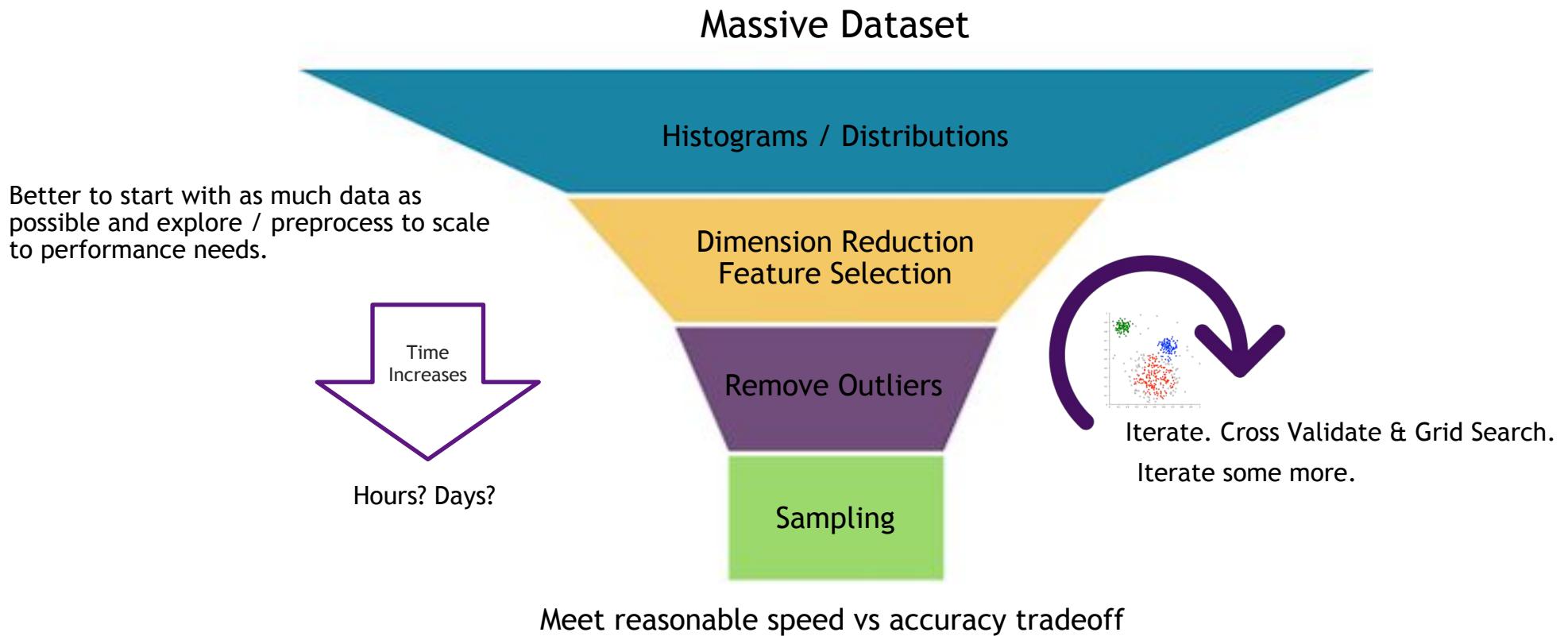
Machine Learning

More models more problems

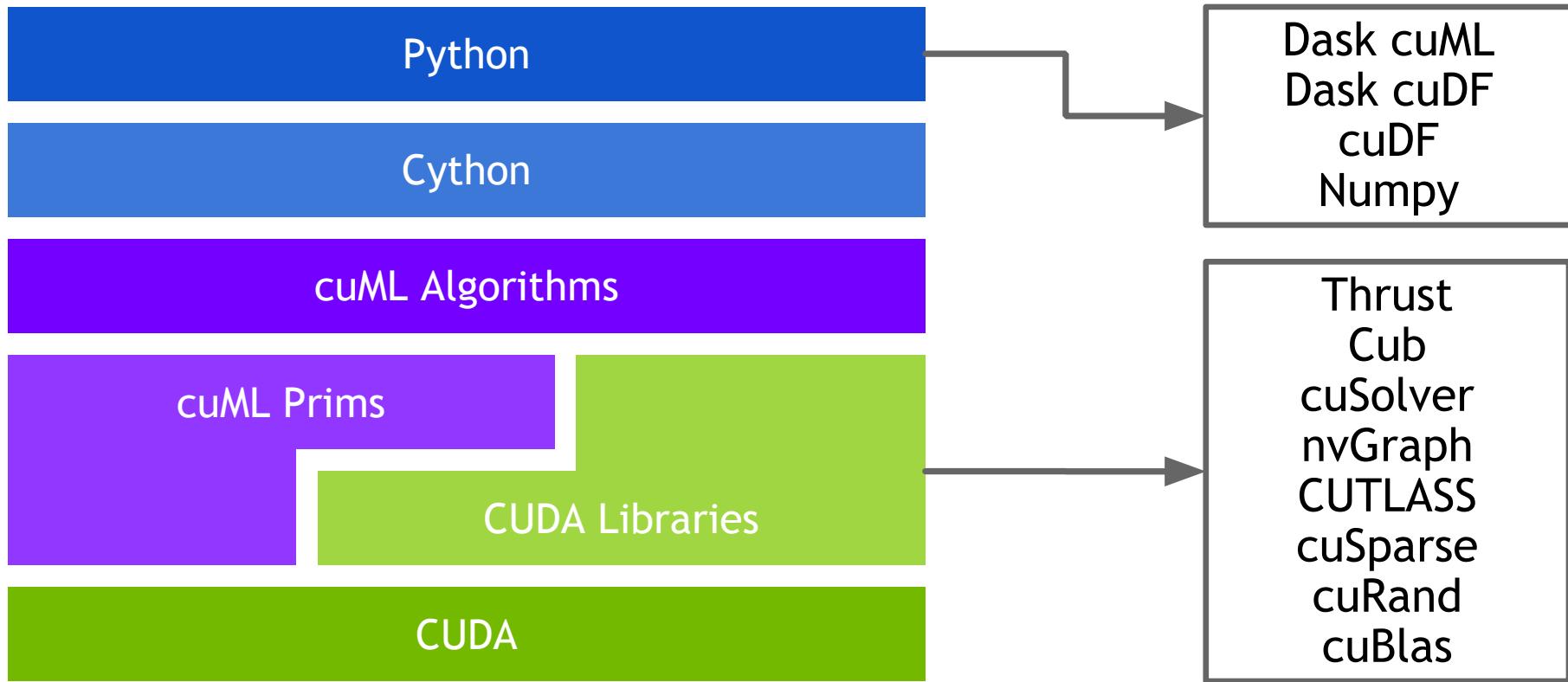


Problem

Data sizes continue to grow

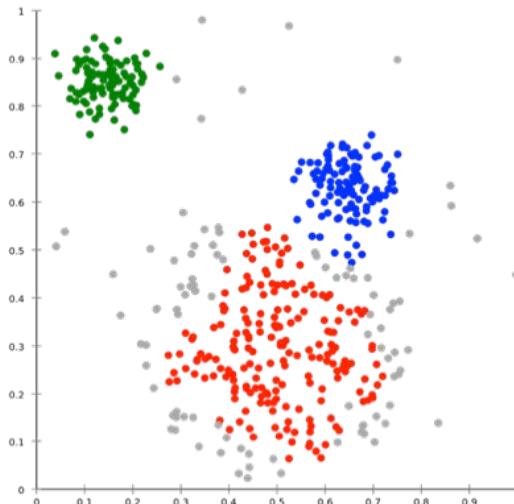


ML Technology Stack



Algorithms

GPU-accelerated Scikit-Learn



Cross Validation

Hyper-parameter Tuning

More to come!

Classification / Regression

Inference

Clustering

Decomposition & Dimensionality Reduction

Time Series

Decision Trees / Random Forests
Linear Regression
Logistic Regression
K-Nearest Neighbors
Support Vector Machines

Random forest / GBDT inference

K-Means
DBSCAN
Spectral Clustering

Principal Components
Singular Value Decomposition
UMAP
Spectral Embedding
T-SNE

Holt-Winters
Seasonal ARIMA

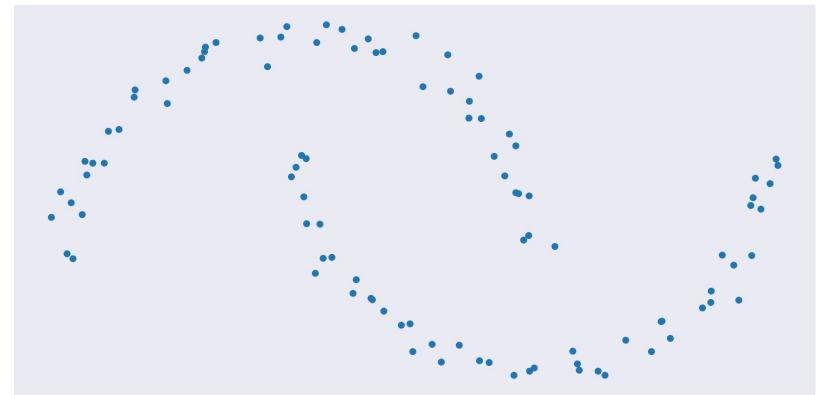
Key:

- Preexisting
- NEW or enhanced for 0.13

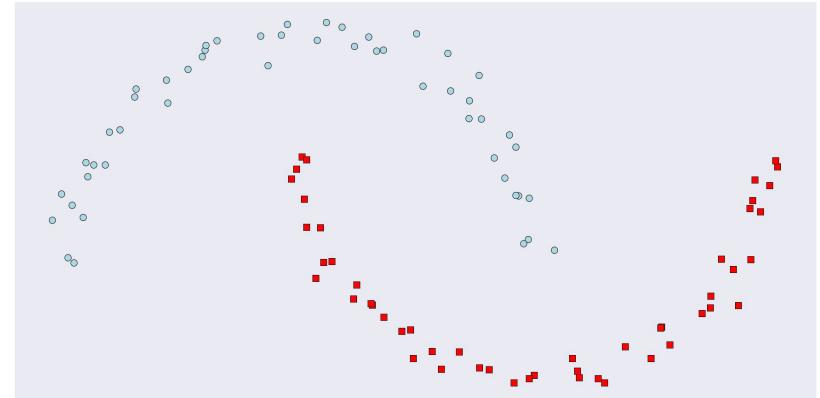
RAPIDS matches common Python APIs

CPU-Based Clustering

```
from sklearn.datasets import make_moons  
import pandas  
  
X, y = make_moons(n_samples=int(1e2),  
                   noise=0.05, random_state=0)  
  
X = pandas.DataFrame({'fea%d'%i: X[:, i]  
                      for i in range(X.shape[1])})
```



```
from sklearn.cluster import DBSCAN  
dbSCAN = DBSCAN(eps = 0.3, min_samples = 5)  
  
dbSCAN.fit(X)  
  
y_hat = dbSCAN.predict(X)
```



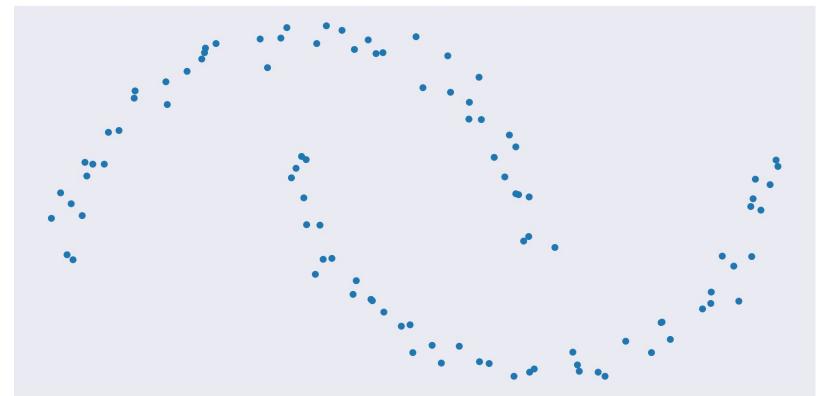
RAPIDS matches common Python APIs

GPU-Accelerated Clustering

```
from sklearn.datasets import make_moons
import cudf

X, y = make_moons(n_samples=int(1e2),
                   noise=0.05, random_state=0)

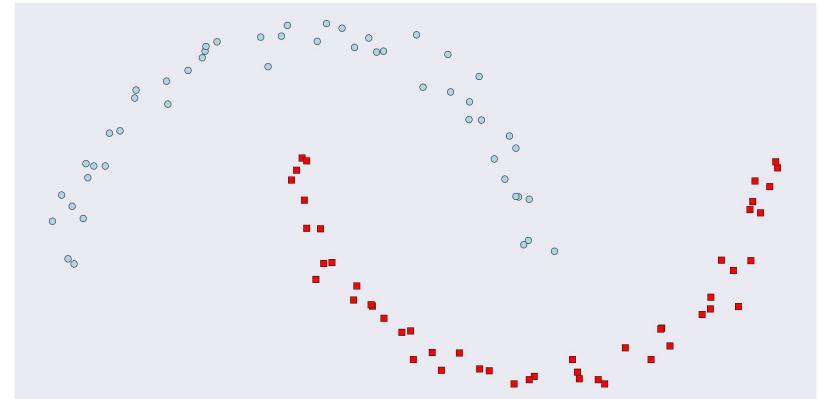
X = cudf.DataFrame({'fea%d'%i: X[:, i]
                    for i in range(X.shape[1])})
```



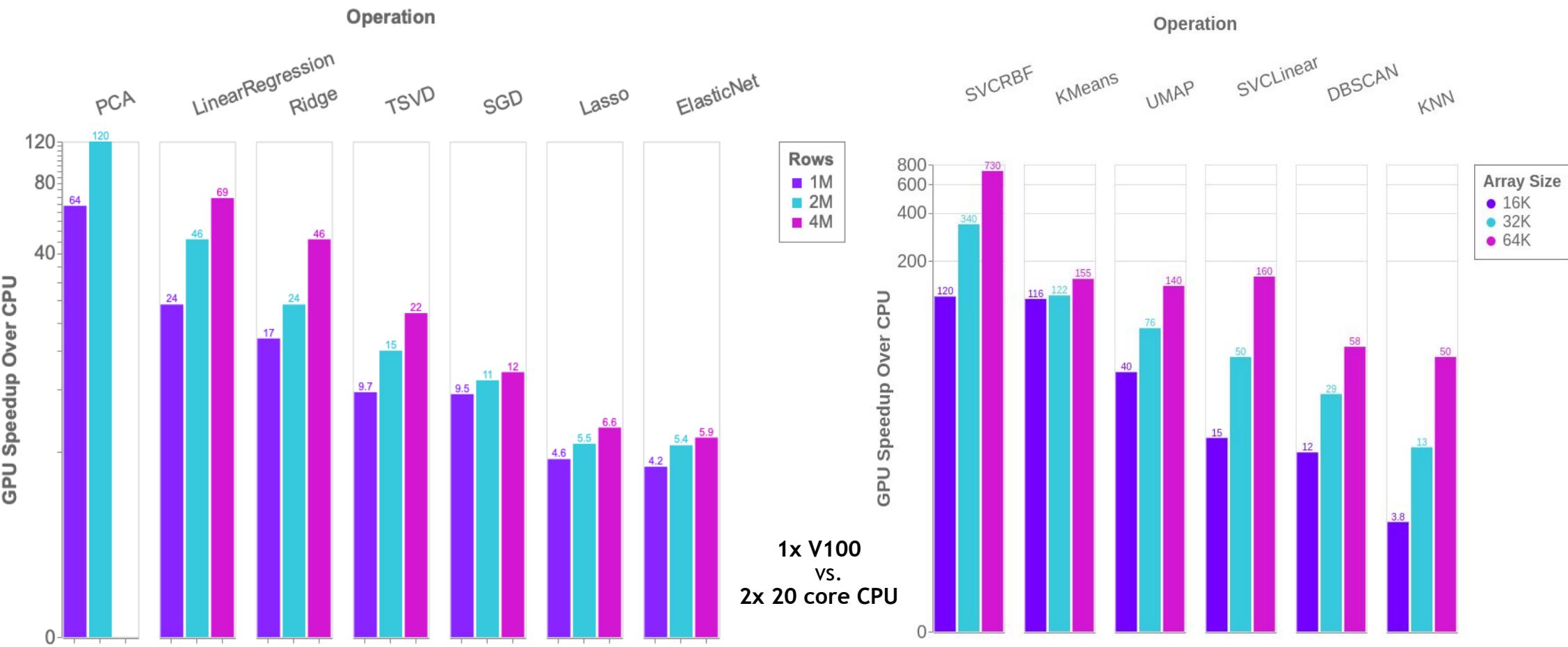
```
from cuml import DBSCAN
dbscan = DBSCAN(eps = 0.3, min_samples = 5)

dbscan.fit(X)

y_hat = dbscan.predict(X)
```



Benchmarks: single-GPU cuML vs scikit-learn

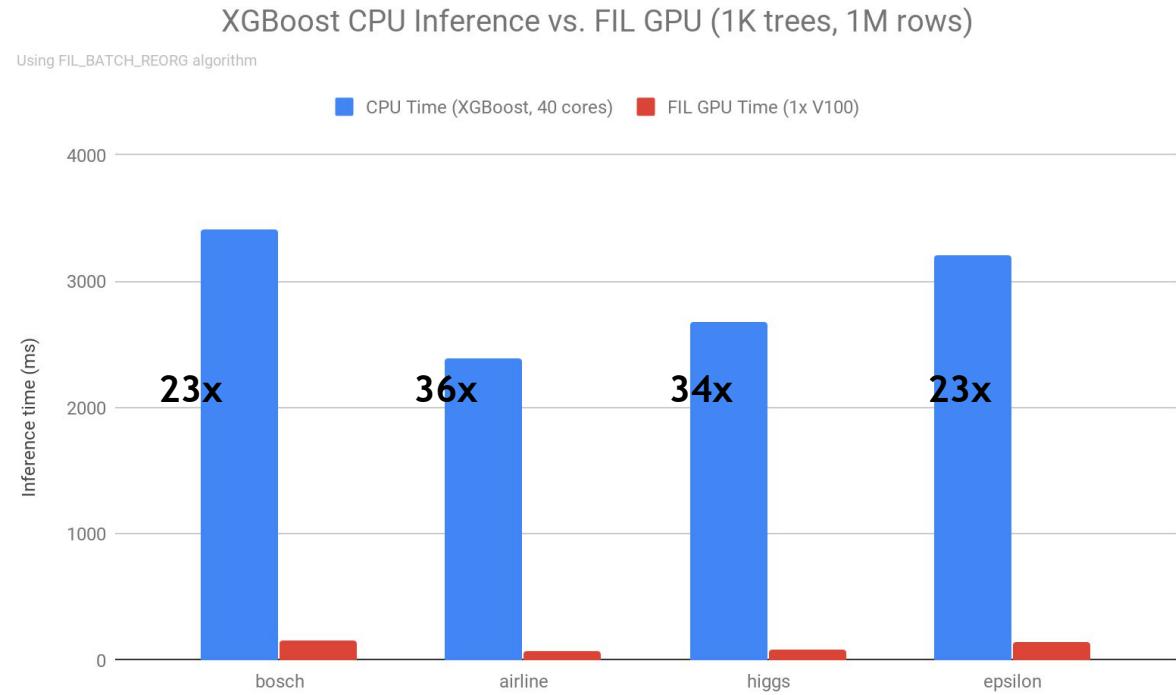


Forest Inference

Taking models from training to production

cuML's **Forest Inference Library** accelerates prediction (inference) for random forests and boosted decision trees:

- Works with existing saved models (XGBoost, LightGBM, scikit-learn RF cuML RF soon)
- Lightweight Python API
- Single V100 GPU can infer up to 34x faster than XGBoost dual-CPU node
- Over 100 million forest inferences per sec (with 1000 trees) on a DGX-1 for large (sparse) or dense models

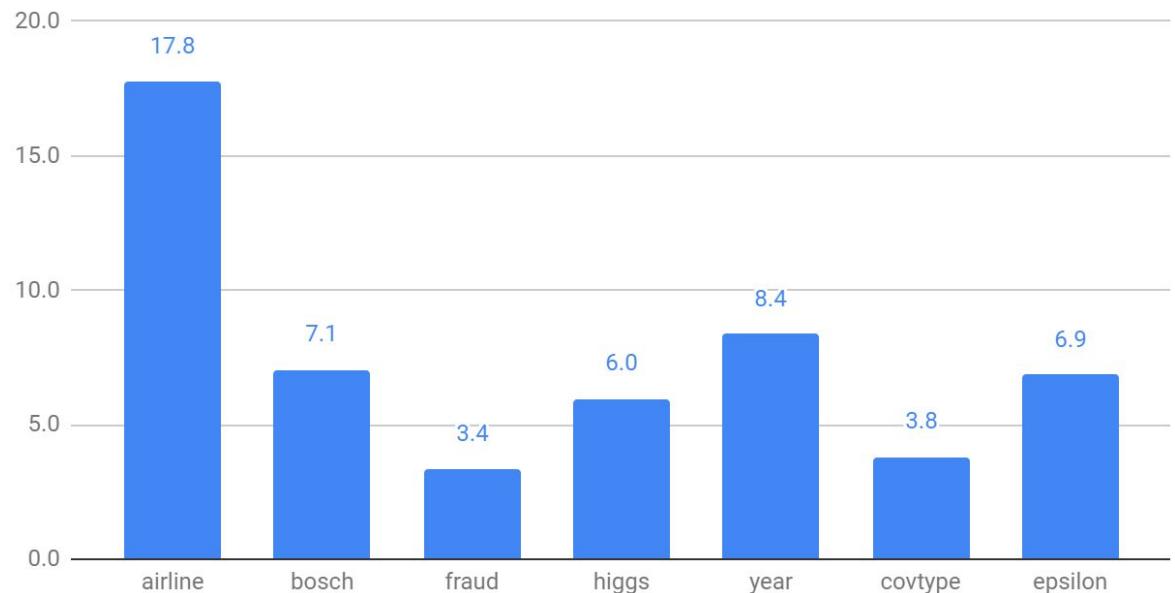


RAPIDS

+ *dmlc*
XGBoost

- RAPIDS works closely with the XGBoost community to accelerate GBDTs on GPU
- The default rapids conda metapackage includes XGBoost
- XGBoost can seamlessly load data from cuDF dataframes and cuPy arrays
- Dask allows XGBoost to scale to arbitrary numbers of GPUs
- With the *gpu_hist* tree method, a single GPU can outpace 10s to 100s of CPUs
- Version 1.0 of XGBoost launched with the RAPIDS 0.13 stack.

XGBoost GPU Speedup - Single V100 vs Dual 20-core Xeon E5-2698



Road to 1.0

March 2020 - RAPIDS 0.13

cuML	Single-GPU	Multi-GPU	Multi-Node-Multi-GPU
Gradient Boosted Decision Trees (GBDT)			
Linear Regression			
Logistic Regression			
Random Forest			
K-Means			
K-NN			
DBSCAN			
UMAP			
Holt-Winters			
ARIMA			
t-SNE			
Principal Components			
Singular Value Decomposition			
SVM			

Road to 1.0

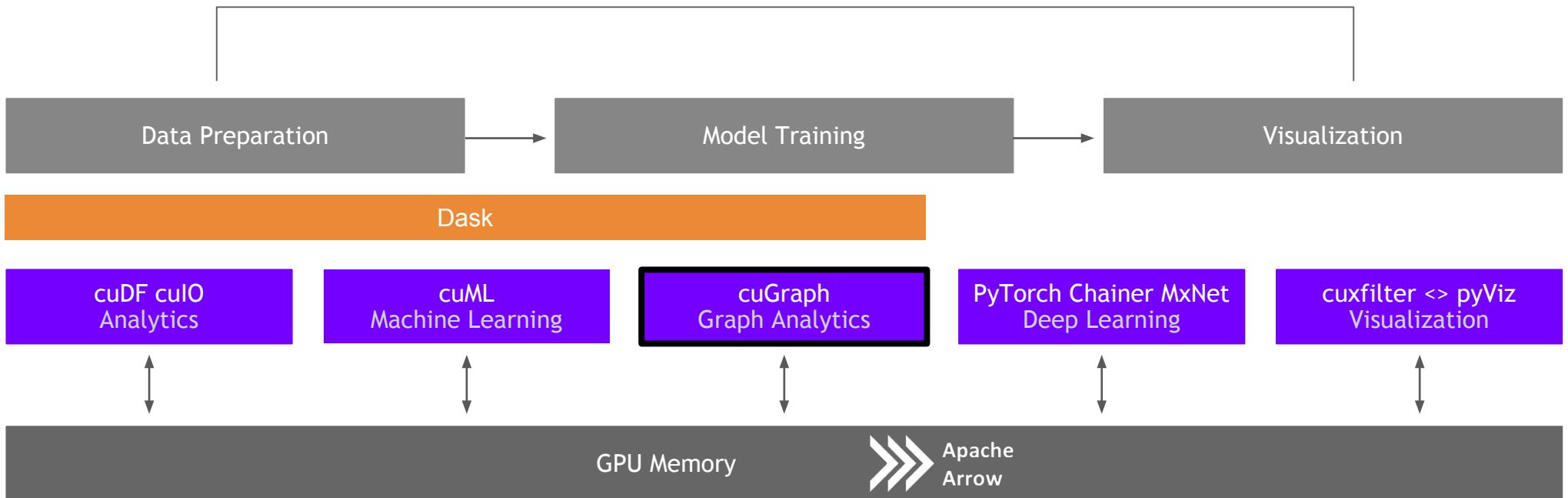
2020 - RAPIDS 1.0

cuML	Single-GPU	Multi-Node-Multi-GPU
Gradient Boosted Decision Trees (GBDT)		
Linear Regression (regularized)		
Logistic Regression		
Random Forest		
K-Means		
K-NN		
DBSCAN		
UMAP		
Holt-Winters		
ARIMA		
t-SNE		
Principal Components		
Singular Value Decomposition		
SVM		

cuGraph

Graph Analytics

More connections more insights



Goals and Benefits of cuGraph

Focus on Features and User Experience

Breakthrough Performance

- Up to 500 million edges on a single 32GB GPU
- Multi-GPU support for scaling into the billions of edges

Multiple APIs

- Python: Familiar NetworkX-like API
- C/C++: lower-level granular control for application developers

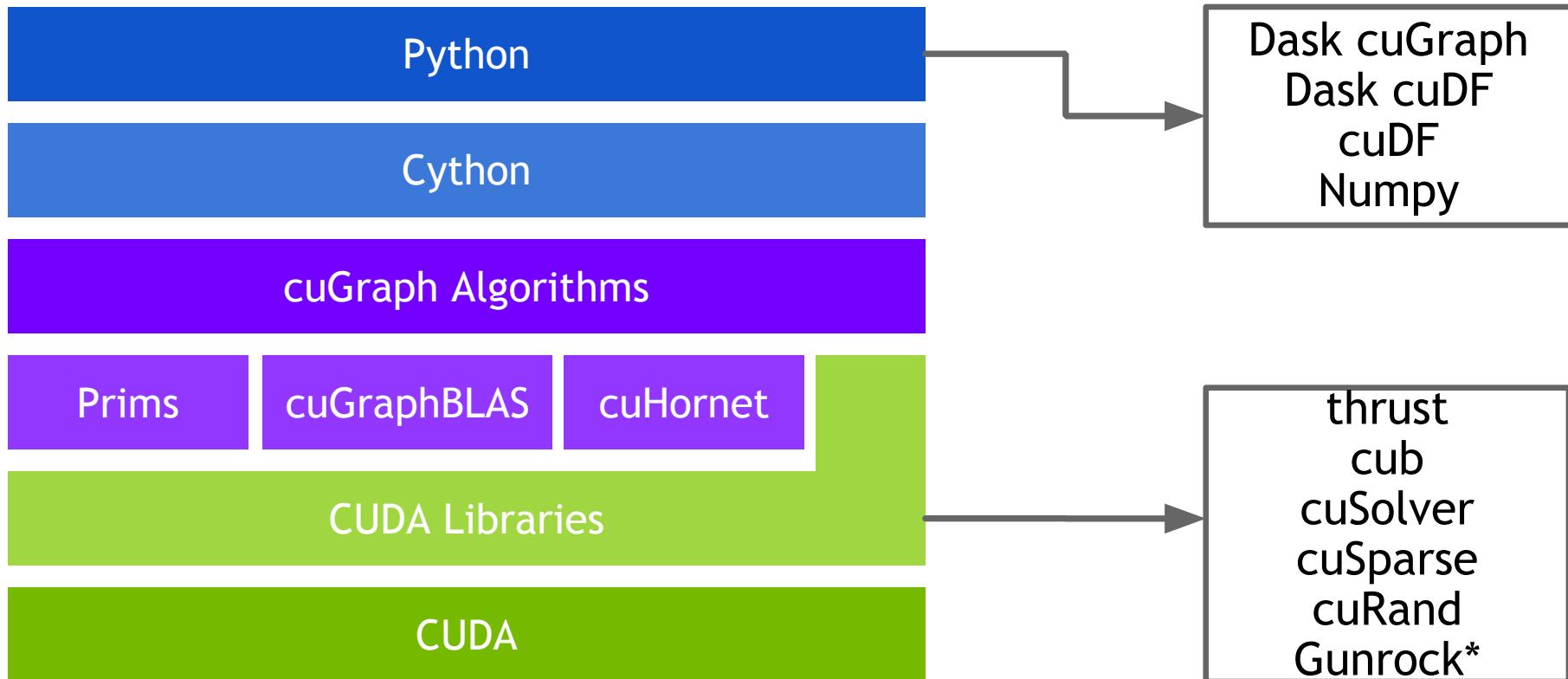
Seamless Integration with cuDF and cuML

- Property Graph support via DataFrames

Growing Functionality

- Extensive collection of algorithm, primitive, and utility functions

Graph Technology Stack

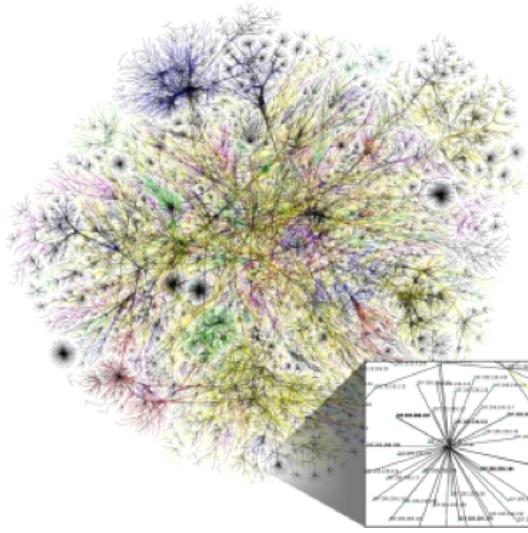


nvGRAPH has been Opened Sourced and integrated into cuGraph. A legacy version is available in a RAPIDS GitHub repo

* Gunrock is from UC Davis

Algorithms

GPU-accelerated NetworkX



Graph Classes

Structure

Renumbering
Auto-renumbering

Multi-GPU

Utilities

More to come!

Community

Components

Link Analysis

Link Prediction

Traversal

Centrality

Spectral Clustering
Balanced-Cut
Modularity Maximization
Louvain
Ensemble Clustering for Graphs
Subgraph Extraction
KCore and KCore Number
Triangle Counting
K-Truss

Weakly Connected Components
Strongly Connected Components

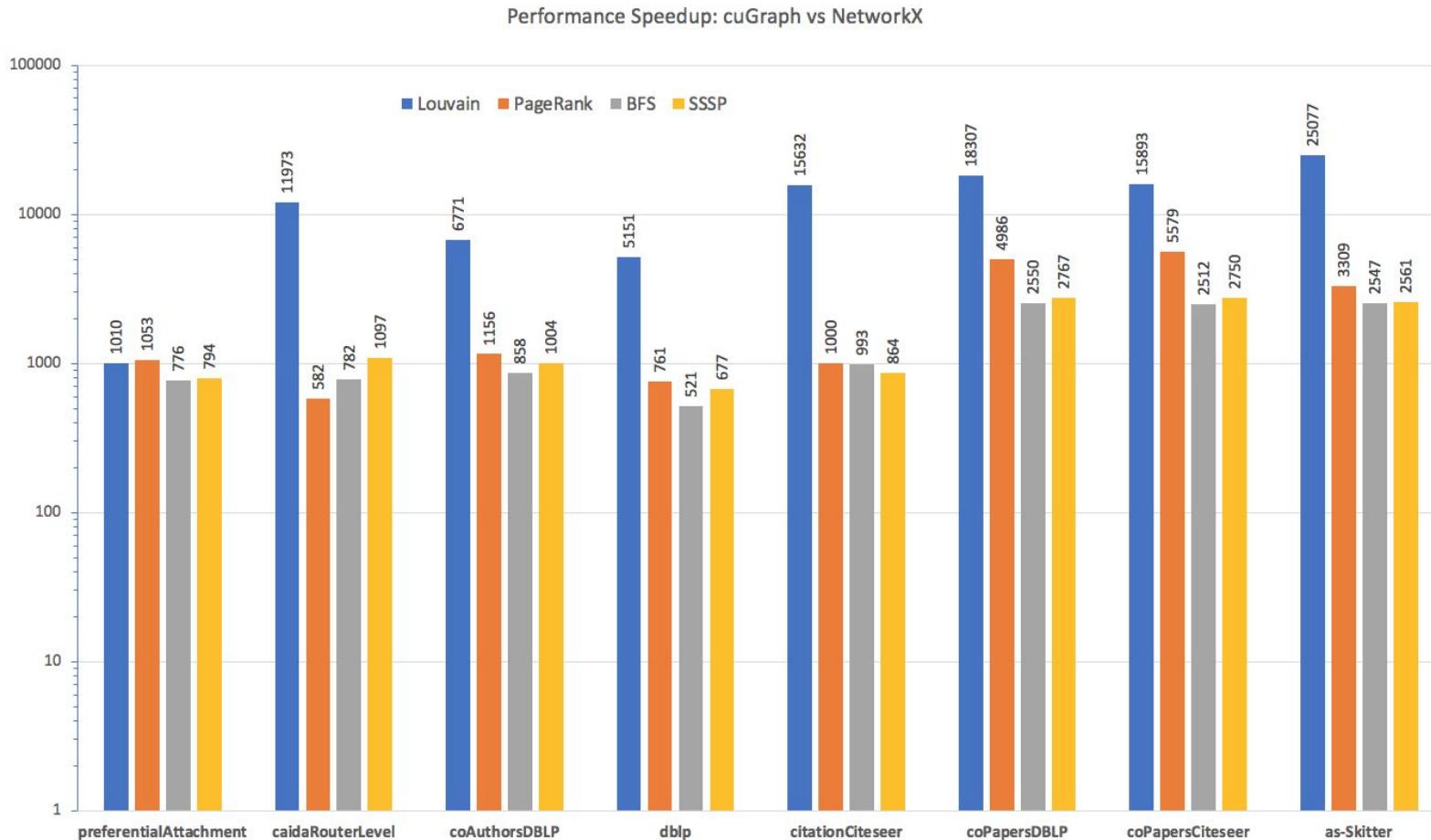
Page Rank (Multi-GPU)
Personal Page Rank

Jaccard
Weighted Jaccard
Overlap Coefficient

Single Source Shortest Path (SSSP)
Breadth First Search (BFS)

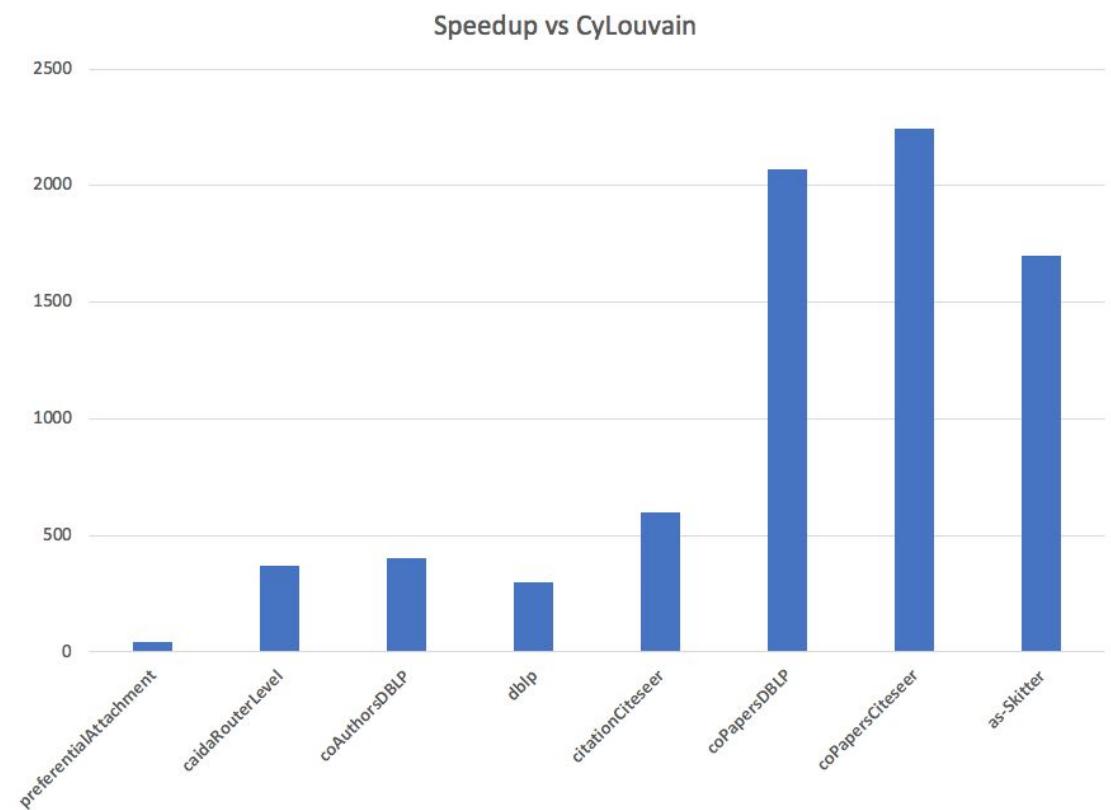
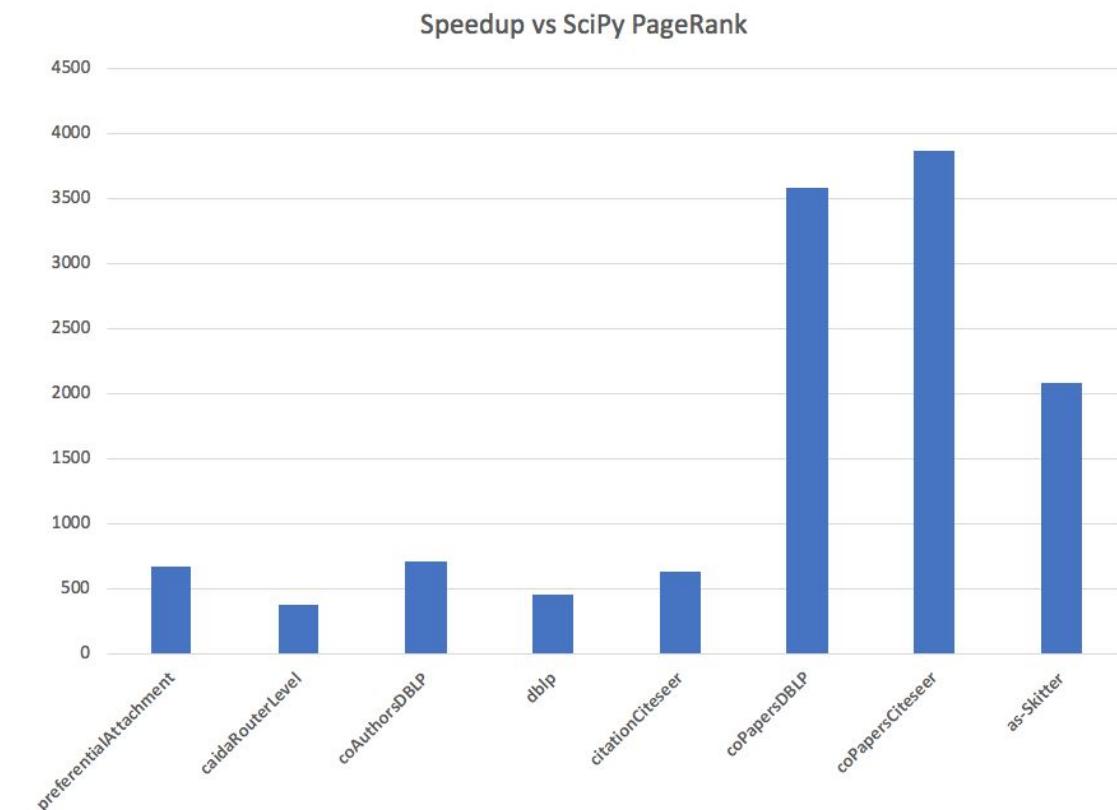
Katz
Betweenness Centrality

Benchmarks: single-GPU cuGraph vs NetworkX



Dataset	Nodes	Edges
preferentialAttachment	100,000	999,970
caidaRouterLevel	192,244	1,218,132
coAuthorsDBLP	299,067	299,067
dblp-2010	326,186	1,615,400
citationCiteseer	268,495	2,313,294
coPapersDBLP	540,486	30,491,458
coPapersCiteseer	434,102	32,073,440
as-Skitter	1,696,415	22,190,596

Benchmarks: cyLouvain and SciPy PageRank



Multi-GPU PageRank Performance

PageRank portion of the HiBench benchmark suite

HiBench Scale	Vertices	Edges	CSV File (GB)	# of V100 GPUs	# of CPU Threads	PageRank for 3 Iterations (secs)
Huge	5,000,000	198,000,000	3	1		1.1
BigData	50,000,000	1,980,000,000	34	3		5.1
BigData x2	100,000,000	4,000,000,000	69	6		9.0
BigData x4	200,000,000	8,000,000,000	146	12		18.2
BigData x8	400,000,000	16,000,000,000	300	16		31.8
BigData x8	400,000,000	16,000,000,000	300	800*		5760*

*BigData x8, 100x 8-vCPU nodes, Apache Spark GraphX \Rightarrow 96 mins!

Road to 1.0

March 2020 - RAPIDS 0.13

cuGraph	Single-GPU	Multi-GPU	Multi-Node-Multi-GPU
PageRank			
Personal Page Rank			
Katz			
Betweenness Centrality			
Spectral Clustering			
Louvain			
Ensemble Clustering for Graphs			
K-Core			
K-Truss			
Triangle Counting			
Connected Components (Weak and Strong)			
Jaccard			
Overlap Coefficient			
Single Source Shortest Path (SSSP)			
Breadth First Search (BFS)			

Road to 1.0

March 2020 - RAPIDS 1.0

cuGraph	Single-GPU	Multi-Node-Multi-GPU
PageRank		
Personal Page Rank		
Katz		
Betweenness Centrality		
Spectral Clustering		
Louvain		
Ensemble Clustering for Graphs		
K-Core		
K-Truss		
Triangle Counting		
Connected Components (Weak and Strong)		
Jaccard		
Overlap Coefficient		
Single Source Shortest Path (SSSP)		
Breadth First Search (BFS)		

cuSpatial

cuSpatial

Breakthrough Performance & Ease of Use

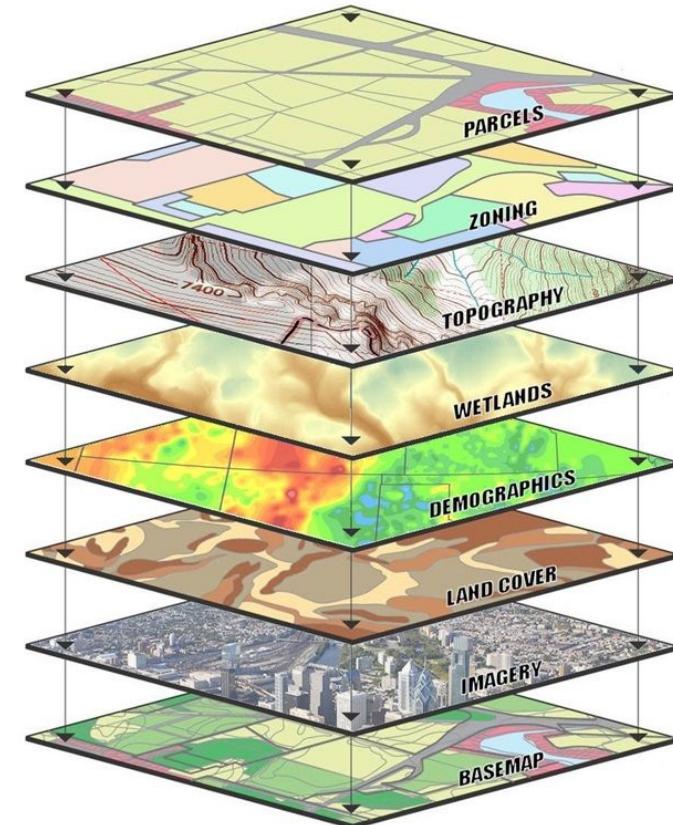
- Up to 1000x faster than CPU spatial libraries
- Python and C++ APIs for maximum usability and integration

Growing Functionality

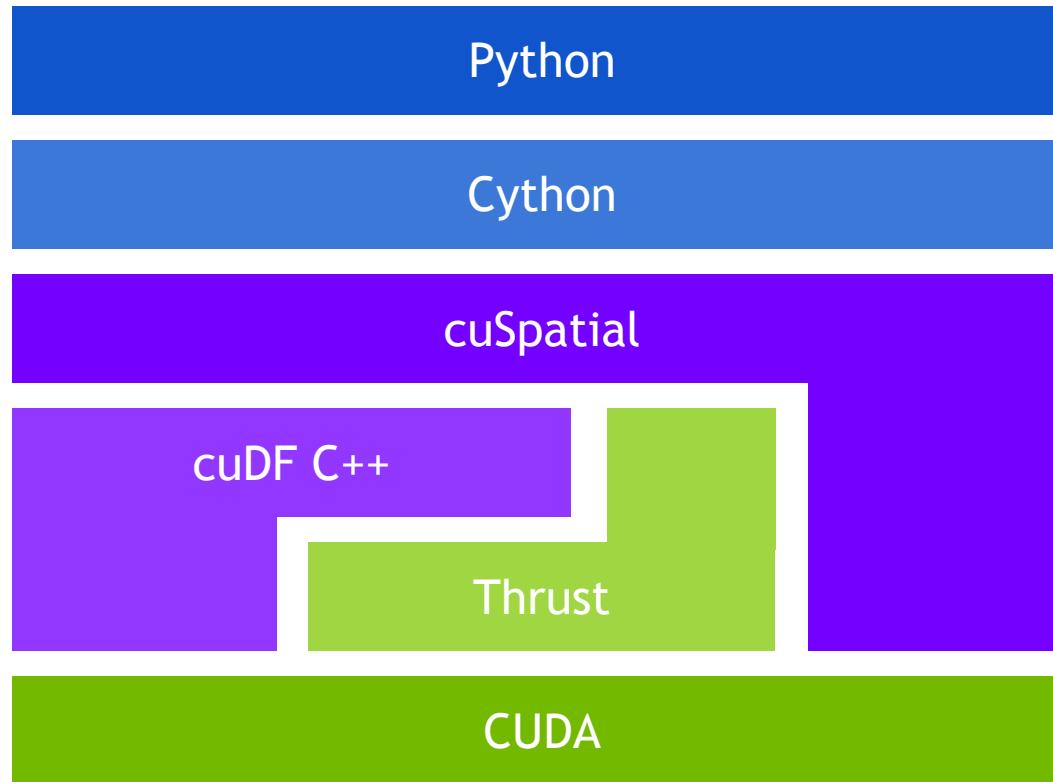
- Extensive collection of algorithm, primitive, and utility functions for spatial analytics

Seamless Integration into RAPIDS

- cuDF for data loading, cuGraph for routing optimization, and cuML for clustering are just a few examples



cuSpatial Technology Stack



cuSpatial

0.13 and Beyond

Layer	0.13 Functionality	Functionality Roadmap (2020)
High-level Analytics	C++ Library w. Python bindings enabling distance, speed, trajectory similarity, trajectory clustering	C++ Library w. Python bindings for additional spatio-temporal trajectory clustering, acceleration, dwell-time, salient locations, trajectory anomaly detection, origin destination, etc.
Graph layer	cuGraph	Map matching, Djikstra algorithm, Routing
Query layer	Spatial Window	Nearest Neighbor, KNN, Spatiotemporal range search and joins
Index layer		Grid, Quad Tree, R-Tree, Geohash, Voronoi Tessellation
Geo-operations	Point in polygon (PIP), Haversine distance, Hausdorff distance, lat-lon to xy transformation, spline interpolation	Line intersecting polygon, Other distance functions, Polygon intersection, union
Geo-representation	Shape primitives, points, polylines, polygons	Additional shape primitives

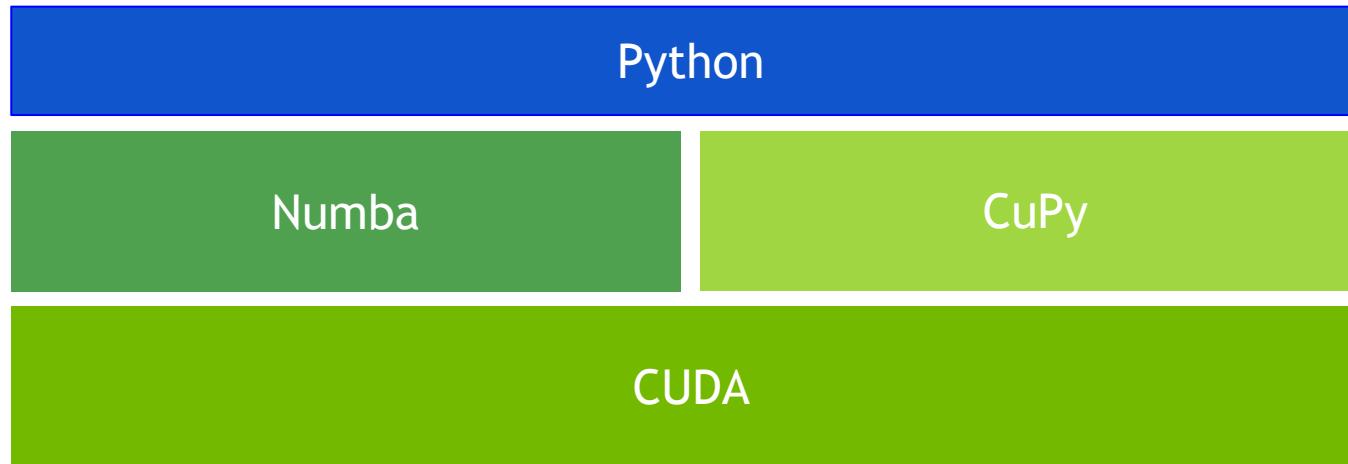
cuSpatial

Performance at a Glance

cuSpatial Operation	Input data	cuSpatial Runtime	Reference Runtime	Speedup
Point-in-Polygon Test	1.3+ million vehicle point locations and 27 Region of Interests	1.11 ms (C++) 1.50 ms (Python) [Nvidia Titan V]	334 ms (C++, optimized serial) 130468.2 ms (python Shapely API, serial) [Intel i7-7800X]	301X (C++) 86,978X (Python)
Haversine Distance Computation	13+ million Monthly NYC taxi trip pickup and drop-off locations	7.61 ms (Python) [Nvidia T4]	416.9 ms (Numba) [Nvidia T4]	54.7X (Python)
Hausdorff Distance Computation (for clustering)	52,800 trajectories with 1.3+ million points	13.5s [Quadro V100]	19227.5s (Python SciPy API, serial) [Intel i7-6700K]	1,400X (Python)

cuSignal

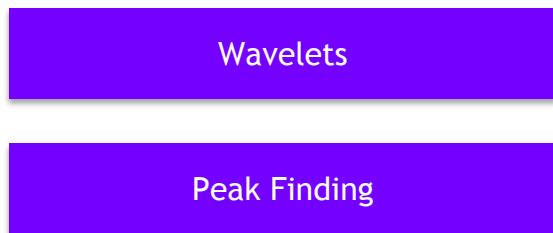
cuSignal Technology Stack



Unlike other RAPIDS libraries, cuSignal is purely developed in Python with custom CUDA Kernels written with Numba and CuPy (notice no Cython layer).

cuSignal - Selected Algorithms

GPU-accelerated SciPy Signal



More to come!



- Convolve/Correlate
- FFT Convolve
- Convolve/Correlate 2D
- Resampling - Polyphase, Upfirdn, Resample
- Hilbert/Hilbert 2D
- Wiener
- Firwin
- Chirp
- Square
- Gaussian Pulse
- Kaiser
- Blackman
- Hamming
- Hanning
- Periodogram
- Welch
- Spectrogram

Speed of Light Performance - V100

timeit (7 runs) rather than *time*. Benchmarked with ~1e8 sample signals on a DGX Station

Method	Scipy Signal (ms)	cuSignal (ms)	Speedup (xN)
fftconvolve	28400	92.2	308.0
correlate	16800	48.4	347.1
resample	14700	51.1	287.7
resample_poly	3110	13.7	227.0
welch	4620	53.7	86.0
spectrogram	2520	28	90.0
cwt	46700	277	168.6

Learn more about cuSignal functionality and performance by browsing the [notebooks](#)

Efficient Memory Handling

Seamless Data Handoff from cuSignal to PyTorch >=1.4

Leveraging the `__cuda_array_interface__` for Speed of Light End-to-End Performance

```
from numba import cuda
import cupy as cp
import torch
from cusignal import resample_poly

# Create CuPy Array on GPU
gpu_arr = cp.random.randn(100_000_000, dtype=cp.float32)

# Polyphase Resample
resamp_arr = resample_poly(gpu_arr, up=2, down=3, window=('kaiser', 0.5))

# Zero Copy to PyTorch
torch_arr = torch.as_tensor(resamp_arr, device='cuda')

# Confirm Pointers
print('Resample Array: ', resamp_arr.__cuda_array_interface__['data'])
print('Torch Array: ', torch_arr.__cuda_array_interface__['data'])

Resample Array: (140516096213080, False)
Torch Array: (140516096213080, False)
```

Enabling Online Signal Processing with Zero-Copy Memory

CPU <-> GPU Direct Memory Access with Numba's Mapped Array

```
import numpy as np
import cupy as cp
from numba import cuda
import cusignal

# Create CPU/GPU Shared Memory, similar to numpy.zeros()
N = 2**18
shared_arr = cusignal.get_shared_mem(N, dtype=np.complex64)
print('CPU Pointer: ', shared_arr.__array_interface__['data'])
print('GPU Pointer: ', shared_arr.__cuda_array_interface__['data'])

CPU Pointer: (139895179837440, False)
GPU Pointer: (139895179837440, False)
```

```
# Load Shared Array with Numpy Array
shared_arr = np.random.randn(N) + 1j*np.random.randn(N)
```

```
%timeit
# Perform CPU FFT
cpu_fft = np.fft.fft(shared_arr)

8 ms ± 60.2 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)
```

```
%timeit
# Perform GPU FFT
gpu_fft = cp.fft.fft(cp.asarray(shared_arr))
cp.cuda.Device(0).synchronize()
```

866 µs ± 38 µs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

Community

Ecosystem Partners

CONTRIBUTORS



ADOPTERS

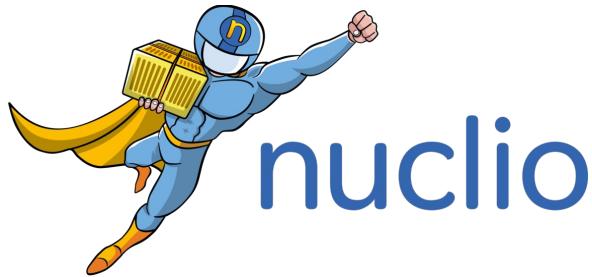


OPEN SOURCE



Building on top of RAPIDS

A bigger, better, stronger ecosystem for all



High-Performance
Serverless event and
data processing that
utilizes RAPIDS for GPU
Acceleration

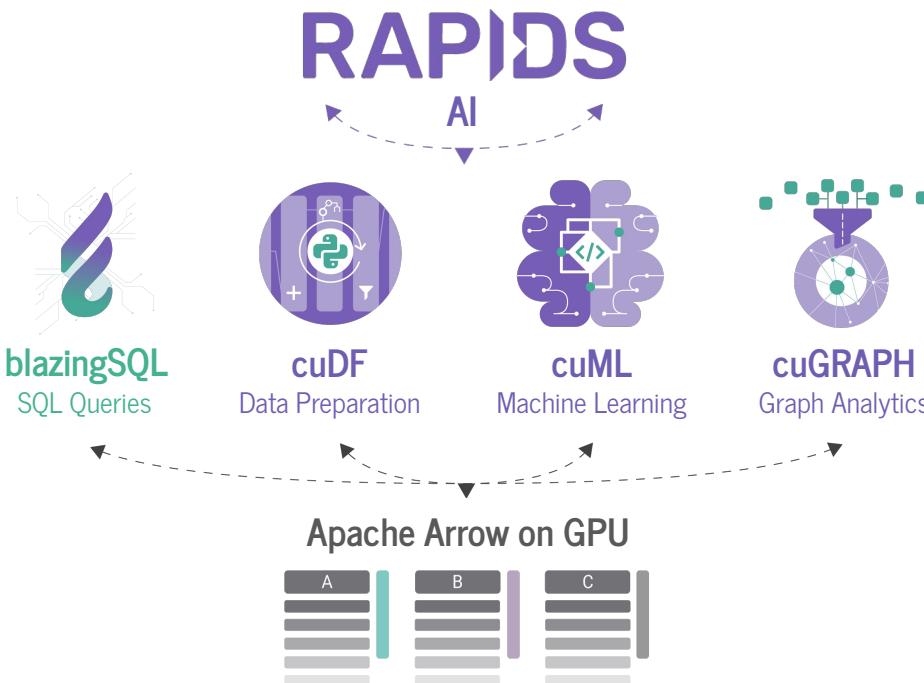


GPU accelerated SQL
engine built on top of
RAPIDS

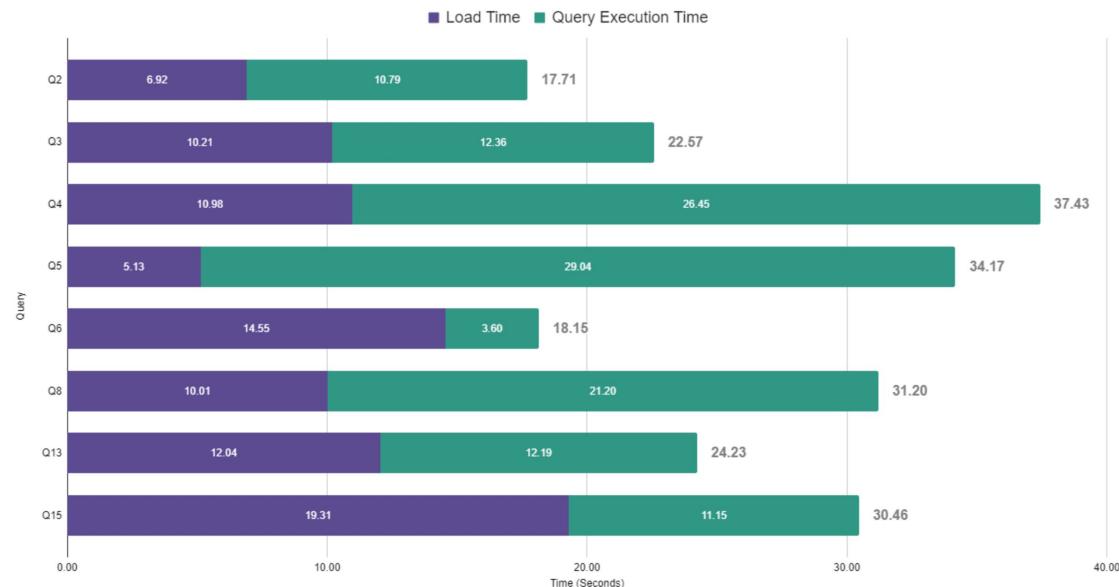
Streamz

Distributed stream
processing using
RAPIDS and Dask

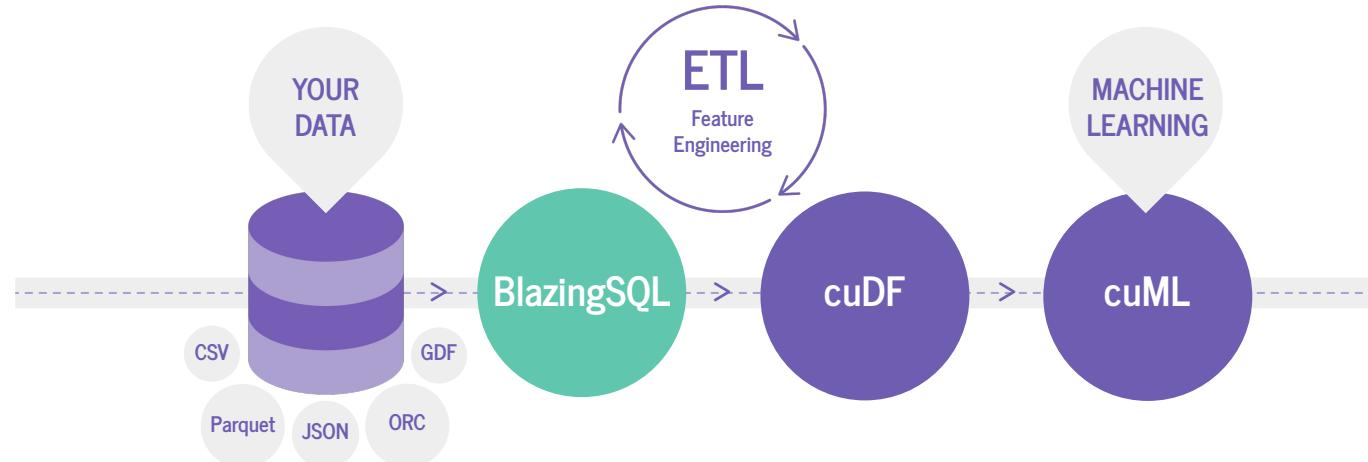
BlazingSQL



TPC-H SF1000 Query Times - GCS Storage



BlazingSQL



```
from blazingsql import BlazingContext
import cudf

bc = BlazingContext()

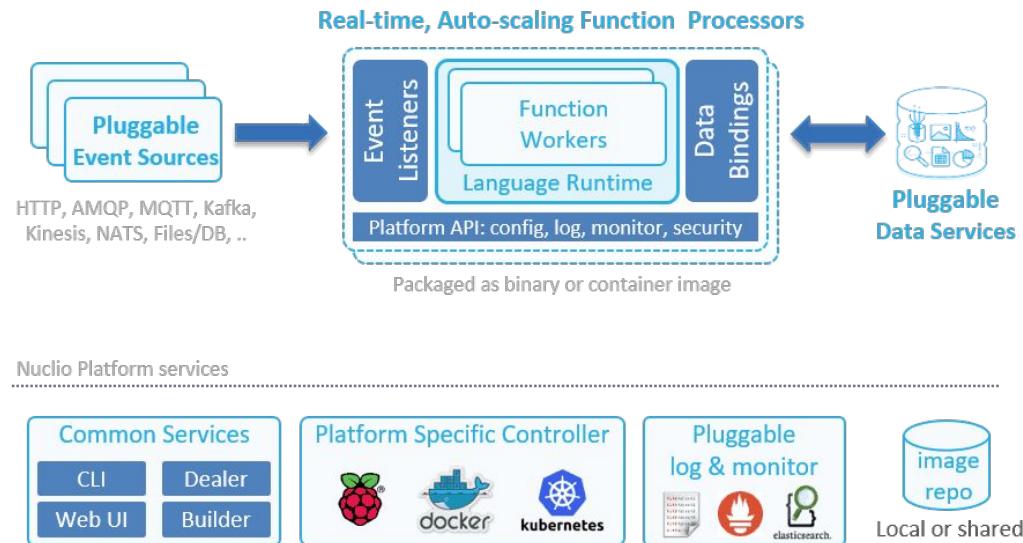
bc.s('bsql', bucket_name='bsql', access_key_id='<access_key>', secret_key='<secret_key>')

bc.create_table('orders', s3://bsql/orders/')

gdf = bc.sql('select * from orders').get()
```

RAPIDS + Nuclio

Serverless meets GPUs



<https://towardsdatascience.com/python-pandas-at-extreme-performance-912912b1047c>

Deploy RAPIDS Everywhere

Focused on robust functionality, deployment, and user experience



Google Cloud



Kubeflow



Azure



Azure Machine Learning



Cloud Dataproc



Amazon SageMaker



Alibaba Cloud

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

5 Steps to getting started with RAPIDS

1. Install RAPIDS on using [Docker](#), [Conda](#), or [Colab](#)
2. Explore our [walk through videos](#), [blog content](#), our [github](#), the [tutorial notebooks](#), and our [examples workflows](#),
3. Build your own data science workflows.
4. Join our community conversations on [Slack](#), [Google](#), and [Twitter](#)
5. Contribute back. Don't forget to ask and answer questions on [Stack Overflow](#)

Easy Installation

Interactive Installation Guide

RAPIDS RELEASE SELECTOR

RAPIDS is available as conda packages, docker images, and from source builds. Use the tool below to select your preferred method, packages, and environment to install RAPIDS. Certain combinations may not be possible and are dimmed automatically. Be sure you've met the required [prerequisites above](#) and see the [details below](#).

METHOD	Conda 📺	Docker + Examples 🎨	Docker + Dev Env 🏠	Source ⚙️			
RELEASE	Stable (0.13)		Nightly (0.14a)				
PACKAGES	All Packages	cuDF	cuML	cuGraph	cuSignal	cuSpatial	cuxfilter
LINUX	Ubuntu 16.04 🐧	Ubuntu 18.04 🐧	CentOS 7 🐧	RHEL 7 🐧			
PYTHON	Python 3.6	Python 3.7					
CUDA	CUDA 10.0	CUDA 10.1.2	CUDA 10.2				

i **NOTE:** Ubuntu 16.04/18.04 & CentOS 7 use the same `conda install` commands.

```
COMMAND    conda install -c rapidsai -c nvidia -c conda-forge \
                  -c defaults rapids=0.13 python=3.6
```

COPY COMMAND ↗

DETAILS BELOW



Explore: RAPIDS Github

<https://github.com/rapidsai>

The screenshot shows the GitHub profile for the RAPIDS organization. The profile header features a purple background with the word "RAPIDS" in white. Below the header, the text "Open GPU Data Science" and a link "http://rapids.ai" are displayed. The navigation bar includes links for "Repositories 92", "Packages", "People 135", "Teams 138", and "Projects 6". The main content area is titled "Pinned repositories" and lists six repositories:

- cudf**: cuDF - GPU DataFrame Library. Cuda, 2.5k stars, 336 forks.
- cuml**: cuML - RAPIDS Machine Learning Library. C++ 1.1k stars, 169 forks.
- cugraph**: cuGraph - RAPIDS Graph Analytics Library. Cuda, 331 stars, 64 forks.
- cusignal**: cuSignal. Jupyter Notebook, 229 stars, 23 forks.
- cuspatial**: CUDA-accelerated GIS and spatiotemporal algorithms. Python, 90 stars, 21 forks.
- notebooks**: RAPIDS Sample Notebooks. Jupyter Notebook, 319 stars, 144 forks.

Explore: RAPIDS Docs

Improved and easier to use!

The screenshot shows a web browser displaying the RAPIDS Docs website at docs.rapids.ai/api/cudf/stable/10min.html. The page title is "10 Minutes to cuDF and Dask-cuDF". The left sidebar contains a navigation menu with sections like "Home", "cudf", "stable (0.13)", "Search docs", "CONTENTS:", "API Reference", and a detailed "10 Minutes to cuDF and Dask-cuDF" section. This section includes links to "What are these Libraries?", "When to use cuDF and Dask-cuDF", "Persisting Data", "Wait", "Multi-GPU with Dask-cuDF", "10 Minutes to Dask-XGBoost", "10 Minutes to cuDF and CuPy", "Overview of User Defined Functions with cuDF", and "Developer Documentation". The main content area starts with a brief introduction: "Modeled after 10 Minutes to Pandas, this is a short introduction to cuDF and Dask-cuDF, geared mainly for new users." It then discusses "What are these Libraries?", explaining cuDF as a GPU DataFrame library and Dask as a flexible parallel computing library. It also covers the integration of Dask-cuDF, mentioning its ability to process DataFrame partitions on multiple GPUs. A section titled "When to use cuDF and Dask-cuDF" provides guidance based on workflow requirements. At the bottom, there is a code block demonstrating the import statements for cuDF and Dask-cuDF:

```
[1]: import os  
import numpy as np  
import pandas as pd  
import cudf  
import dask_cudf  
  
np.random.seed(12)  
  
#### Portions of this were borrowed and adapted from the  
#### cuDF cheatsheet, existing cuDF documentation,  
#### and 10 Minutes to Pandas.
```

<https://docs.rapids.ai>

Explore: RAPIDS Code and Blogs

Check out our code and how we use it

README.md

RAPIDS cuDF - GPU DataFrames

build running

NOTE: For the latest stable README.md ensure you are on the master branch.

Built based on the Apache Arrow columnar memory format, cuDF is a GPU DataFrame library for loading, joining, aggregating, filtering, and otherwise manipulating data.

cuDF provides a pandas-like API that will be familiar to data engineers & data scientists, so they can use it to easily accelerate their workflows without going into the details of CUDA programming.

For example, the following snippet downloads a CSV, then uses the GPU to parse it into rows and columns and run calculations:

```
import cudf, io, requests
from io import StringIO

url="https://github.com/plotly/datasets/raw/master/tips.csv"
content = requests.get(url).content.decode('utf-8')

tips_df = cudf.read_csv(StringIO(content))
tips_df['tip_percentage'] = tips_df['tip']/tips_df['total_bill']*100

# display average tip by dining party size
print(tips_df.groupby('size').tip_percentage.mean())

Output:
size
```



RAPIDS Release 0.8: Same Community New Freedoms

Making more friends and building more bridges to more ecosystems. It's now easier than ever to get started with RAPIDS.



Josh Patterson

Jul 19 · 7 min read



gQuant—GPU Accelerated examples for Quantitative Analyst Tasks

A simple trading strategy backtest for 5000 stocks using GPUs and getting 20X speedup



Yi Dong

Jul 16 · 6 min read ★



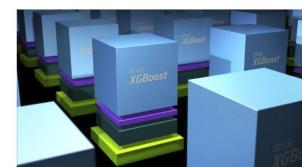
Financial data modeling with RAPIDS.

See how RAPIDS was used to place 17th in the Banco Santander Kaggle Competition



Jiwei Liu

Jul 3 · 5 min read

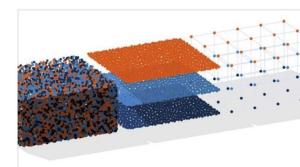


NVIDIA GPUs and Apache Spark, One Step Closer

RAPIDS XGBoost4j-Spark Package Now Available

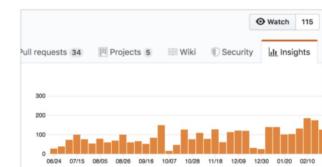


Karthikeyan Rajendran



When Less is More: A brief story about XGBoost feature engineering

A glimpse into how a Data Scientist makes decisions about featuring engineering an XGBoost machine



Nightly News: CI produces latest packages

Release code early and often. Stay current on latest features with our nightly conda and container releases.

<https://github.com/rapidsai>

<https://medium.com/rapids-ai>

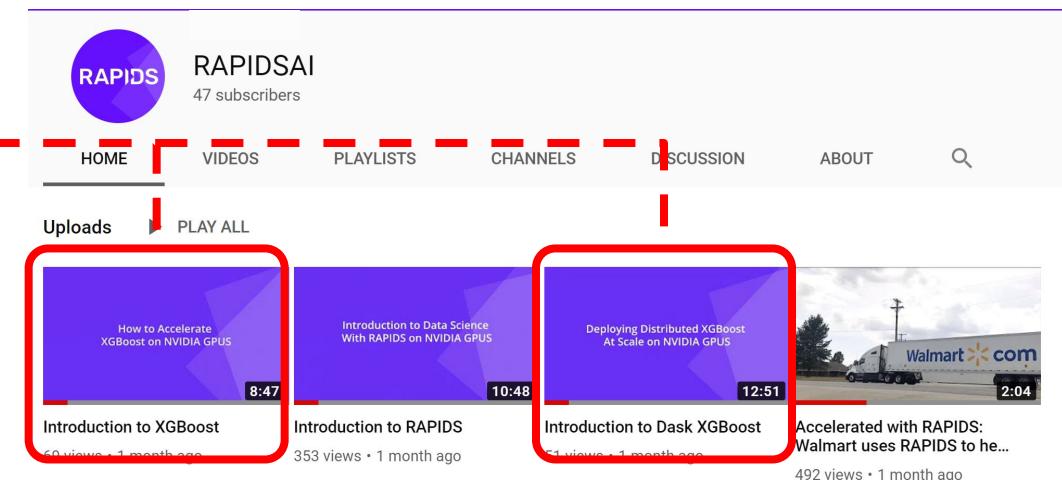
Explore: Notebooks Contrib

Tutorials, examples, and various E2E demos available, with Youtube explanations, code walkthroughs and use cases

intro_tutorials	05_Introduction_to_Dask_cuDF	This notebook shows how to work with cuDF DataFrames distributed across multiple GPUs using Dask.
intro_tutorials	06_Introduction_to_Supervised_Learning	This notebook shows how to do GPU accelerated Supervised Learning in RAPIDS.
intro_tutorials	07_Introduction_to_XGBoost	This notebook shows how to work with GPU accelerated XGBoost in RAPIDS.
intro_tutorials	08_Introduction_to_Dask_XGBoost	This notebook shows how to work with Dask XGBoost in RAPIDS.
intro_tutorials	09_Introduction_to_Dimensionality_Reduction	This notebook shows how to do GPU accelerated Dimensionality Reduction in RAPIDS.
intro_tutorials	10_Introduction_to_Clustering	This notebook shows how to do GPU accelerated Clustering in RAPIDS.

Intermediate Notebooks:

Folder	Notebook Title	Description
examples	DBSCAN_Demo_FULL	This notebook shows how to use DBSCAN algorithm and its GPU accelerated implementation present in RAPIDS.
examples	Dask_with_cuDF_and_XGBoost	In this notebook we show how to quickly setup Dask and train an XGBoost model using cuDF.



Join the Conversation



[Google Groups](#)



[Docker Hub](#)



[Slack Channel](#)



[Stack Overflow](#)

Contribute Back

Issues, feature requests, PRs, Blogs, Tutorials, Videos, QA...bring your best!

The GitHub interface shows four repositories:

- cuml**: cuML - RAPIDS Machine Learning Library. Last updated 9 minutes ago.
- cuDF**: cuDF - GPU DataFrame Library. Last updated 31 minutes ago.
- notebooks-contrib**: RAPIDS Community Notebooks. Last updated 40 minutes ago.
- cugraph**: cuGraph. Last updated 1 day ago.

TECH BLOG Walmart Labs

How GPU Computing literally saved me at work?

Python+GPU = Power, 2 Days to 20 seconds

Abhishek Mungoli [Follow](#) May 9 · 9 min read

John Murray [@MurrayData](#) [Follow](#)

Comparison CPU vs GPU @rapidsai to project 100 million x,y points to lat/lon to 0.01mm accuracy. CPU 1 core c 65 mins, multicore c 13 mins, GPU #RAPIDS 2 seconds. I optimised the code since previous run. Dell T7910 Xeon E5-2640V4x2/NVIDIA Titan Xp cc @NvidiaAI @marc_stampfli

```
john@plato:~/Source/Python/misc$ python crs_test.py
Generating Data
CPU Iterative
4005.0377202 seconds
CPU mapped
3957.19386101 seconds
CPU multiprocessing
788.550751209 seconds
GPU Rapids
2.103230476 seconds
```

Getting Started with cuDF (RAPIDS)

Darren Ramsook [Follow](#) Jun 9 · 3 min read

Getting Started

RAPIDS Docs

Improved and easier to use!

The screenshot shows a web browser window displaying the RAPIDS Docs website at docs.rapids.ai/api/cudf/stable/10min.html. The page title is "10 Minutes to cuDF and Dask-cuDF". The left sidebar contains a navigation menu with sections like Home, cudf, stable (0.13), and a search bar. The main content area starts with a brief introduction: "Modeled after 10 Minutes to Pandas, this is a short introduction to cuDF and Dask-cuDF, geared mainly for new users." It then asks, "What are these Libraries?", explaining cuDF and Dask. It continues with "When to use cuDF and Dask-cuDF", mentioning Multi-GPU with Dask-cuDF, 10 Minutes to Dask-XGBoost, 10 Minutes to cuDF and CuPy, Overview of User Defined Functions with cuDF, and Developer Documentation. Below this, there's a code block demonstrating how to import libraries and set a random seed:

```
[1]: import os  
import numpy as np  
import pandas as pd  
import cudf  
import dask_cudf  
  
np.random.seed(12)  
  
#### Portions of this were borrowed and adapted from the  
#### cuDF cheatsheet, existing cuDF documentation,  
#### and 10 Minutes to Pandas.
```

<https://docs.rapids.ai>

Get Started

Easier than ever to get started with cuDF

The screenshot shows a web browser displaying the cuDF documentation at docs.rapids.ai/api/cudf/stable/10min.html. The page title is "10 Minutes to cuDF and Dask-cuDF". The left sidebar contains a navigation menu with sections like Home, cudf, stable (0.13), and a search bar. The main content area starts with a brief introduction: "Modeled after 10 Minutes to Pandas, this is a short introduction to cuDF and Dask-cuDF, geared mainly for new users." It then asks, "What are these Libraries?", explaining cuDF and Dask. It continues with "When to use cuDF and Dask-cuDF", noting that cuDF is suitable for single GPU or memory-limited workflows, while Dask is for distributed computing. A code block shows Python imports for os, numpy, pandas, cudf, and dask_cudf, followed by a note about borrowing from cuDF's cheatsheet.

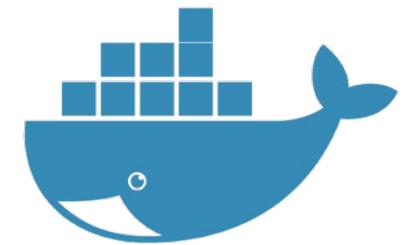
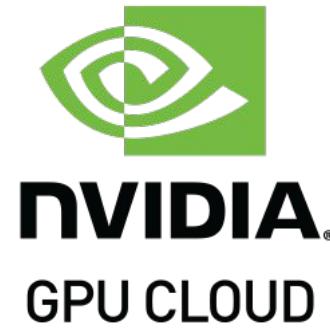
```
[1]: import os
import numpy as np
import pandas as pd
import cudf
import dask_cudf
np.random.seed(12)

#### Portions of this were borrowed and adapted from the
#### cuDF cheatsheet, existing cuDF documentation,
#### and 10 Minutes to Pandas.
```

<https://docs.rapids.ai>

RAPIDS

How do I get the software?



- <https://github.com/rapidsai>
- <https://anaconda.org/rapidsai/>
- <https://ngc.nvidia.com/registry/nvidia-rapidsai-rapidsai>
- <https://hub.docker.com/r/rapidsai/rapidsai/>

Join the Movement

Everyone can help!



APACHE ARROW

<https://arrow.apache.org/>

@ApacheArrow



RAPIDS

<https://rapids.ai>

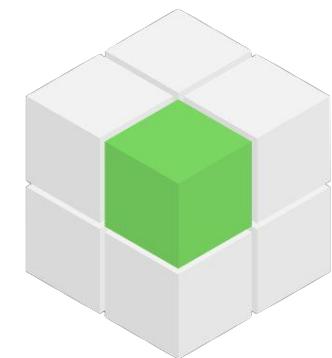
@RAPIDSAI



Dask

<https://dask.org>

@Dask_dev



GPU Open Analytics Initiative

<http://gpuopenanalytics.com/>

@GPUOAI

Integrations, feedback, documentation support, pull requests, new issues, or code donations welcomed!

THANK YOU

Joshua Patterson

joshuap@nvidia.com

@datametrician



RAPIDS