

RAPIDS

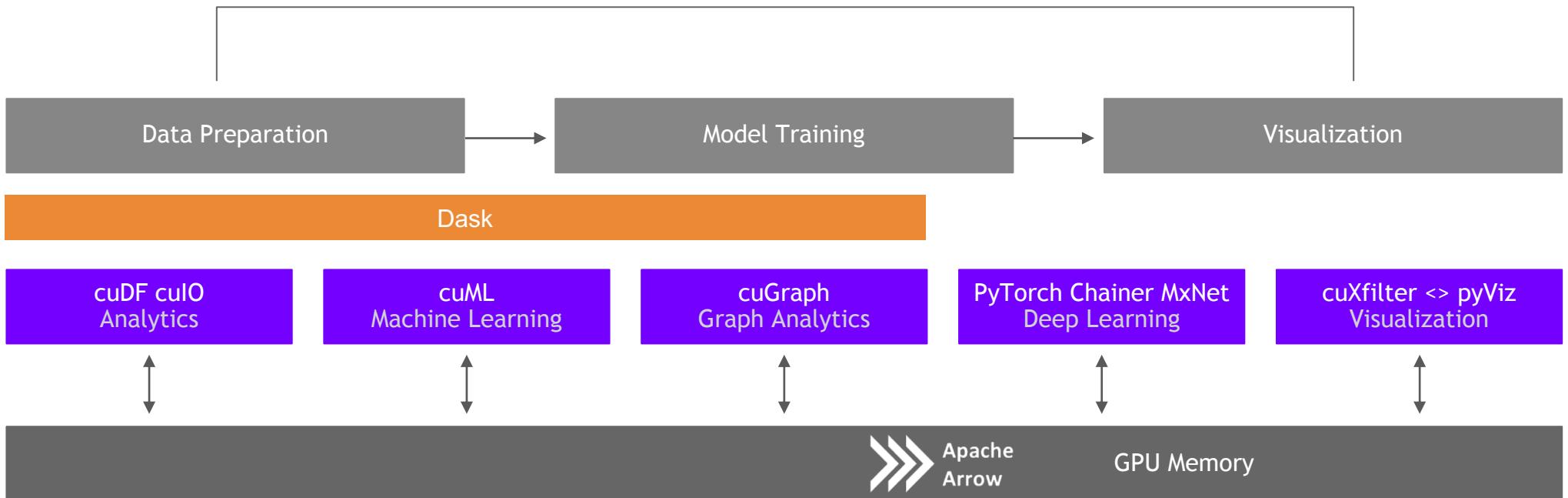
Accelerating Network Analysis
(and all your data science pipelines)

Part 2 - Core and Examples

Brad Rees and Corey Nolet

RAPIDS

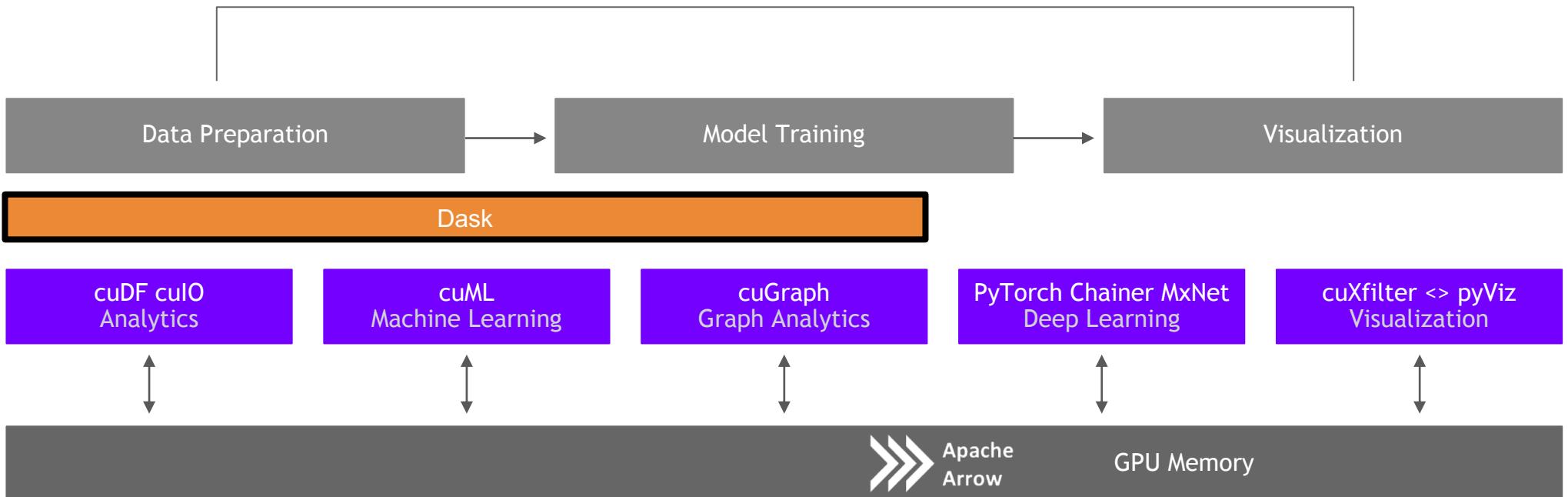
End-to-End Accelerated GPU Data Science



Dask

RAPIDS

Scaling RAPIDS with Dask





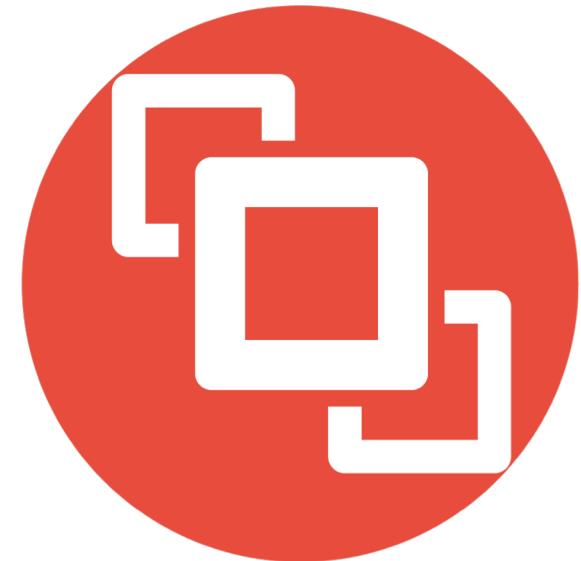
Why Dask?

- **PyData Native**
 - Built on top of NumPy, Pandas, Scikit-Learn, etc. (easy to migrate)
 - With the same APIs (easy to train)
 - With the same developer community (well trusted)
- **Scales**
 - Easy to install and use on a laptop
 - Scales out to thousand-node clusters
- **Popular**
 - Most common parallelism framework today at PyData and SciPy conferences
- **Deployable**
 - HPC: SLURM, PBS, LSF, SGE
 - Cloud: Kubernetes
 - Hadoop/Spark: Yarn

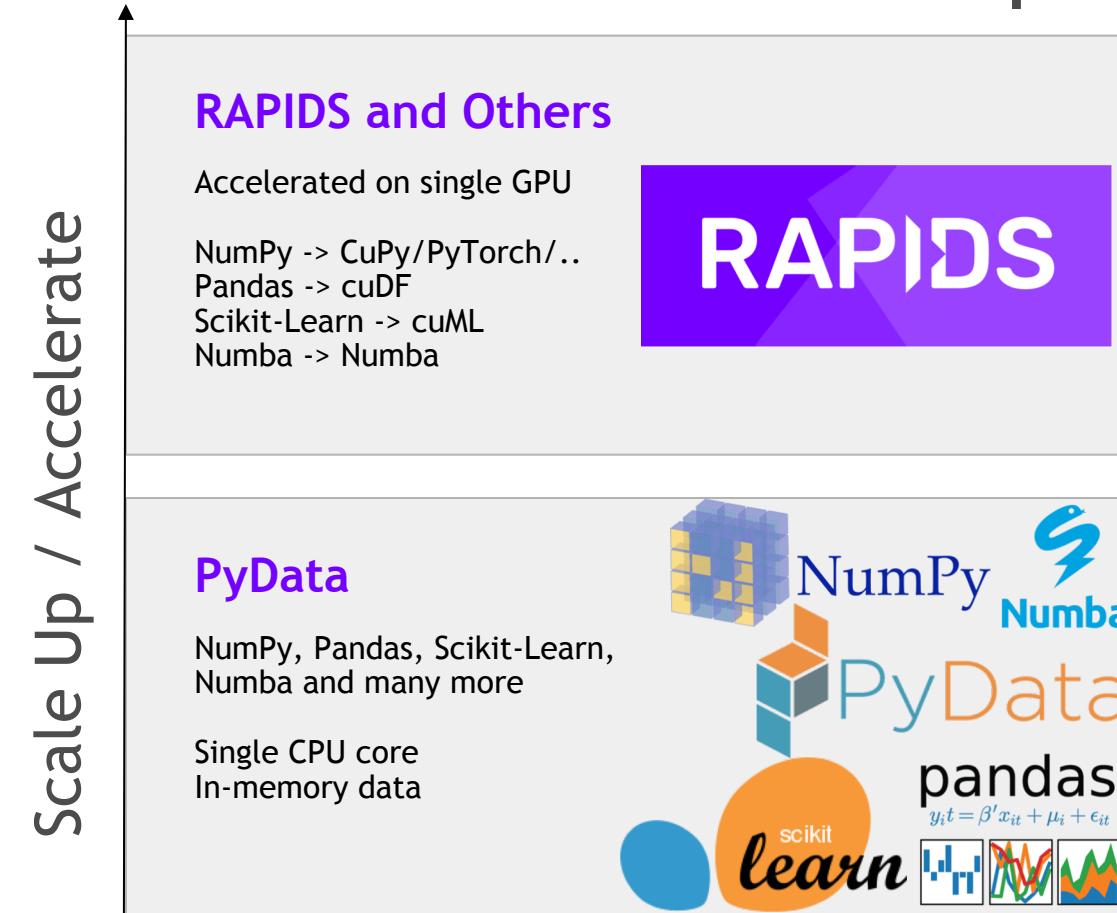
Why OpenUCX?

Bringing hardware accelerated communications to Dask

- TCP sockets are slow!
- UCX provides uniform access to transports (TCP, InfiniBand, shared memory, NVLink)
- Python bindings for UCX (ucx-py) in the works
<https://github.com/rapidsai/ucx-py>
- Will provide best communication performance, to Dask based on available hardware on nodes/cluster



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

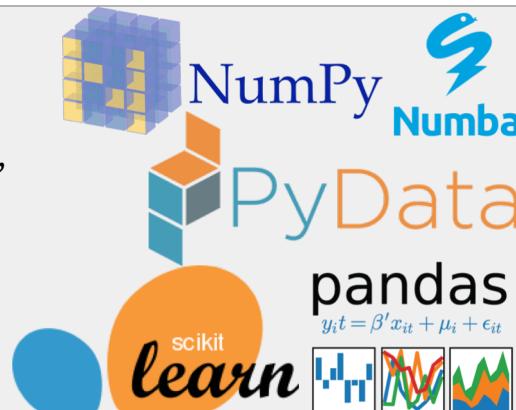
Multi-GPU
On single Node (DGX)
Or across a cluster



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

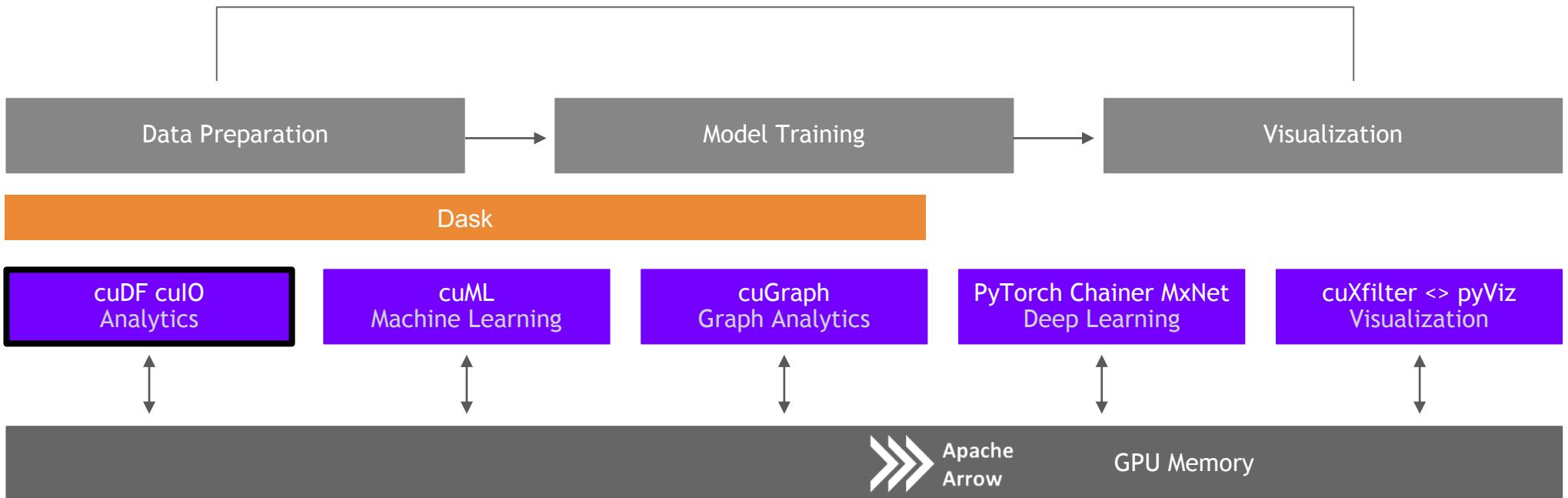


Scale out / Parallelize

cuDF

RAPIDS

GPU Accelerated data wrangling and feature engineering



ETL - the Backbone of Data Science

libcuDF is...

CUDA C++ Library

- Low level library containing function implementations and C/C++ API
- Importing/exporting Apache Arrow in GPU memory using CUDA IPC
- CUDA kernels to perform element-wise math operations on GPU DataFrame columns
- CUDA sort, join, groupby, reduction, etc. operations on GPU DataFrames

```
void some_function( cudf::column const* input,  
                    cudf::column * output,  
                    args...)  
{  
    // Do something with input  
    // Produce output  
}
```



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()
```

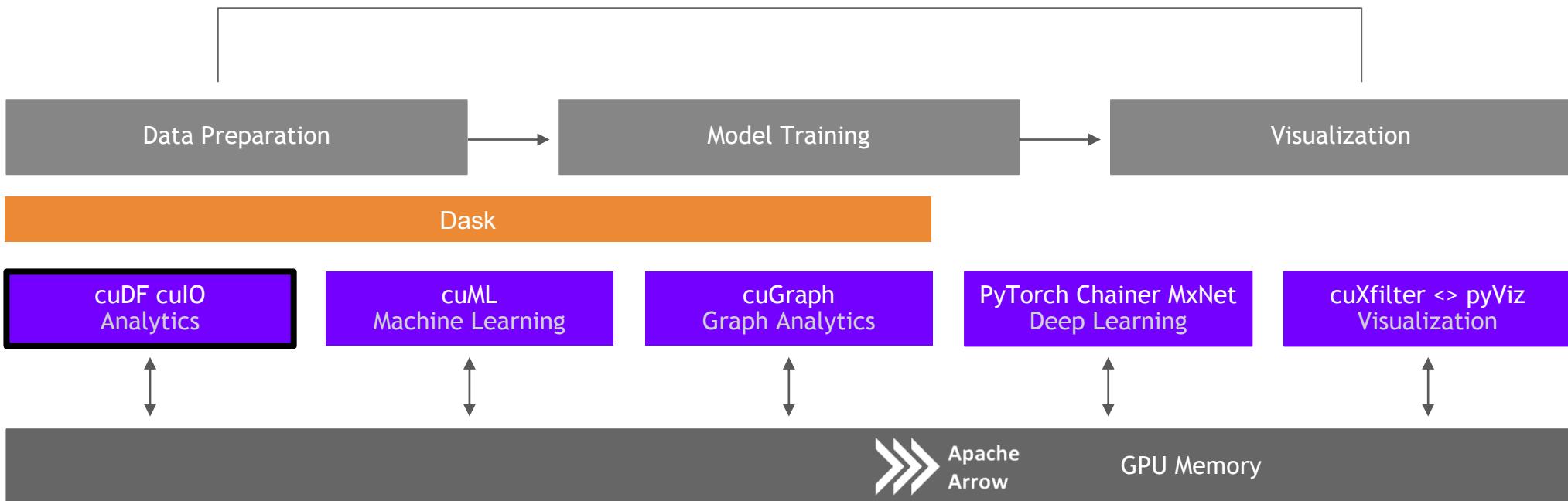
```
In [7]: #Here we can see how we can control what the value of our dummies with the replace method and turn  
#strings to ints  
gdf['City_Category'] = gdf.City_Category.str.replace('A', '1')  
gdf['City_Category'] = gdf.City_Category.str.replace('B', '2')  
gdf['City_Category'] = gdf.City_Category.str.replace('C',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

ETL - the Backbone of Data Science

cuDF is not the end of the story



ETL - the Backbone of Data Science

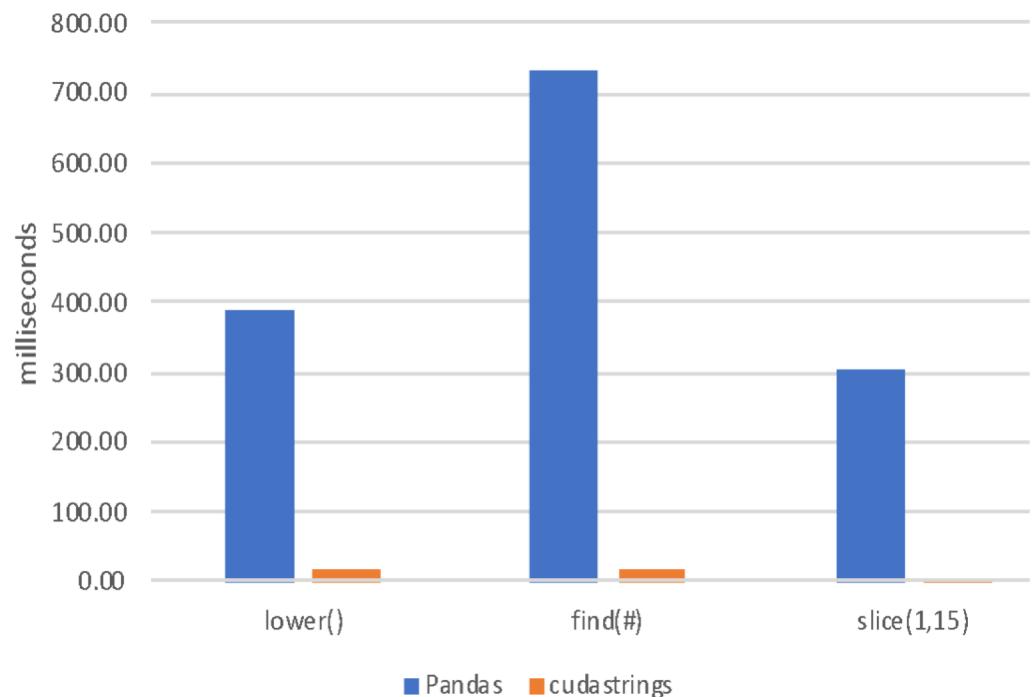
String Support

Now v0.8 String Support:

- Regular Expressions
- Element-wise operations
 - Split, Find, Extract, Cat, Typecasting, etc...
- String GroupBys, Joins

Future v0.9+ String Support:

- Combining cuStrings into libcudf
- Extensive performance optimization
- More Pandas String API compatibility
- Improved Categorical column support



Extraction is the Cornerstone of ETL

culO is born

- Follows the APIs of Pandas and provide >10x speedup
- CSV Reader - v0.2, CSV Writer v0.8
- Parquet Reader - v0.7
- ORC Reader - v0.7
- JSON Reader - v0.8
- Avro Reader - v0.9
- HDF5 Reader - v0.10
- Key is GPU-accelerating both parsing and decompression wherever possible

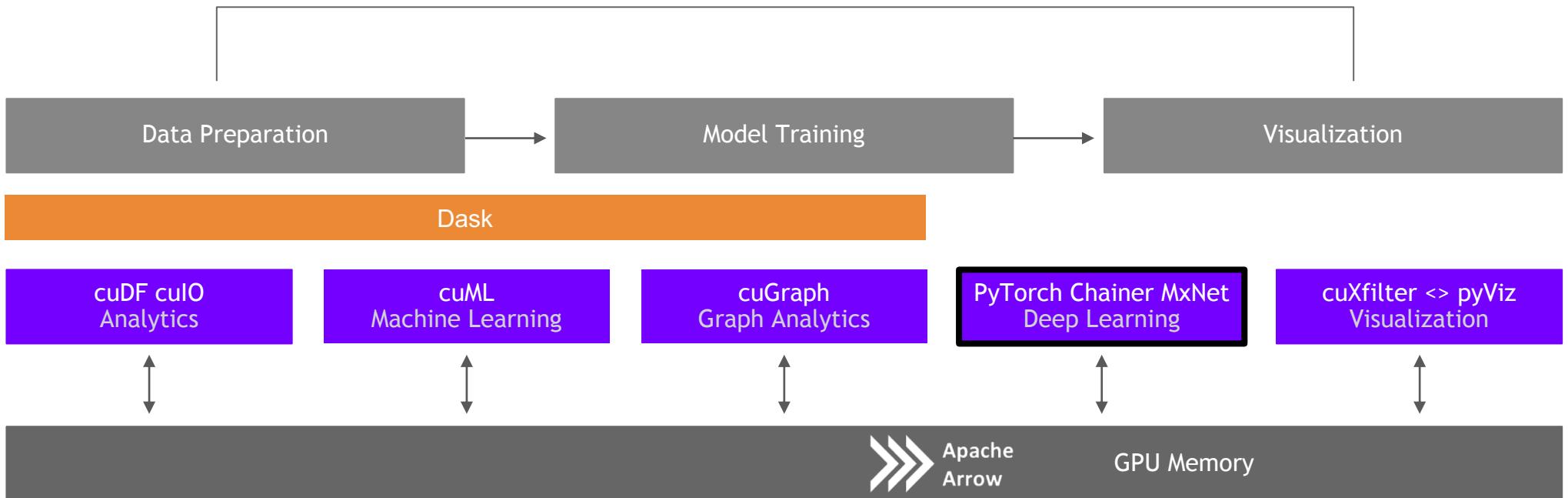
```
1]: import pandas, cudf
2]: %time len(pandas.read_csv('data/nyc/yellow_tripdata_2015-01.csv'))
CPU times: user 25.9 s, sys: 3.26 s, total: 29.2 s
Wall time: 29.2 s
2]: 12748986
3]: %time len(cudf.read_csv('data/nyc/yellow_tripdata_2015-01.csv'))
CPU times: user 1.59 s, sys: 372 ms, total: 1.96 s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

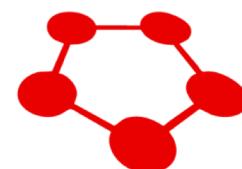


Interoperability for the Win

DLPack and __cuda_array_interface__



mpi4py

The mxnet logo is the word "mxnet" in a large, white, sans-serif font, enclosed within a solid blue circular background.

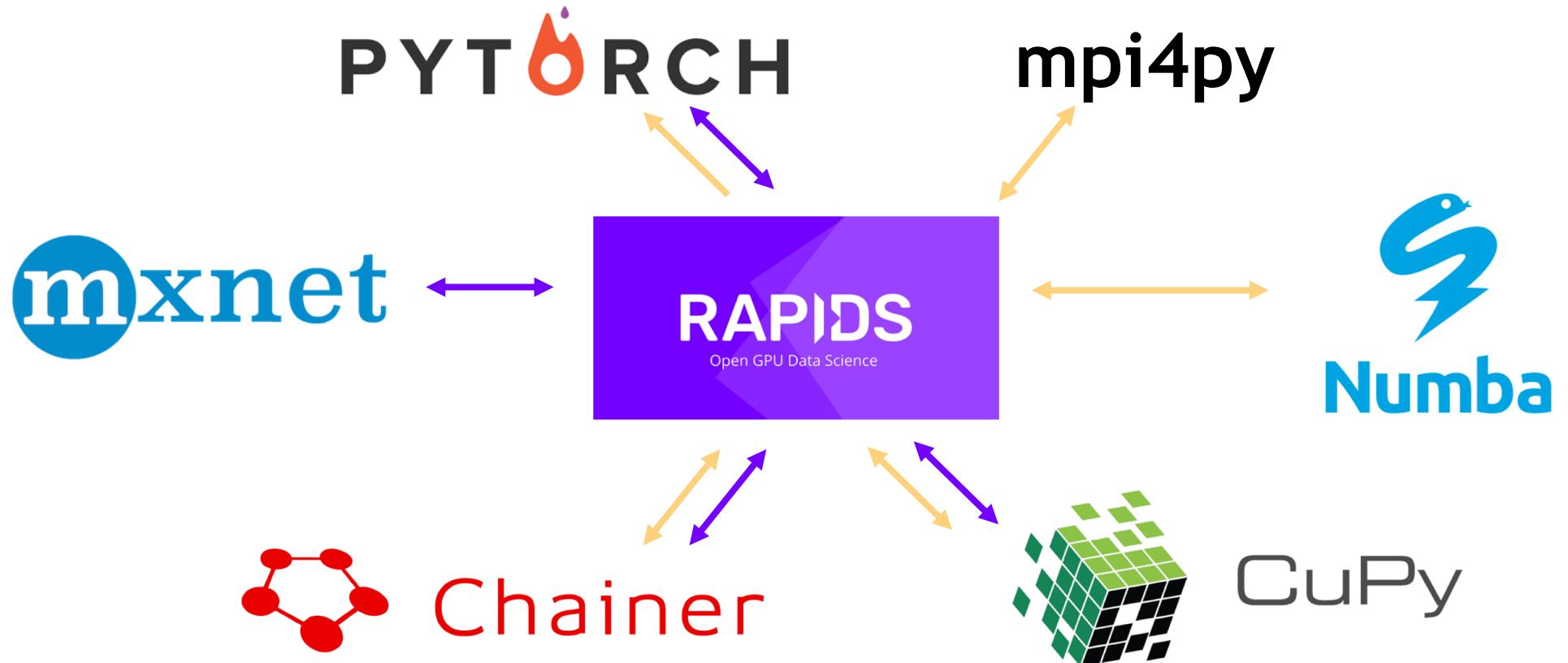
Chainer



CuPy

Interoperability for the Win

DLPack and `__cuda_array_interface__`



ETL - Arrays and DataFrames

Dask and CUDA Python arrays



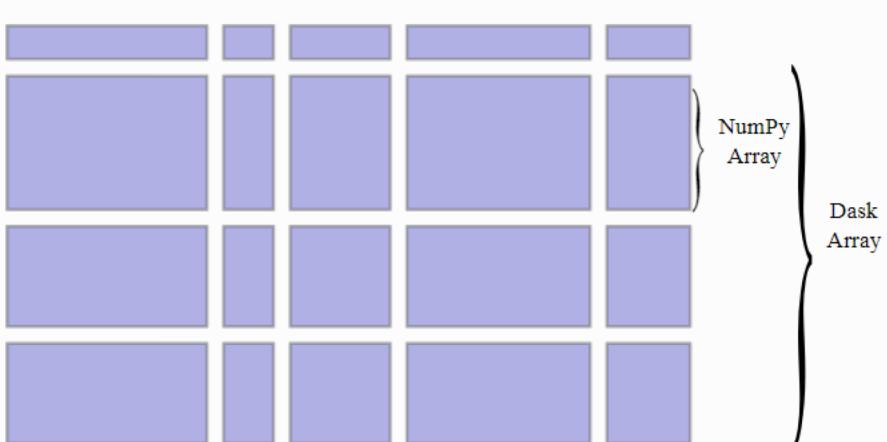
Chainer



CuPy

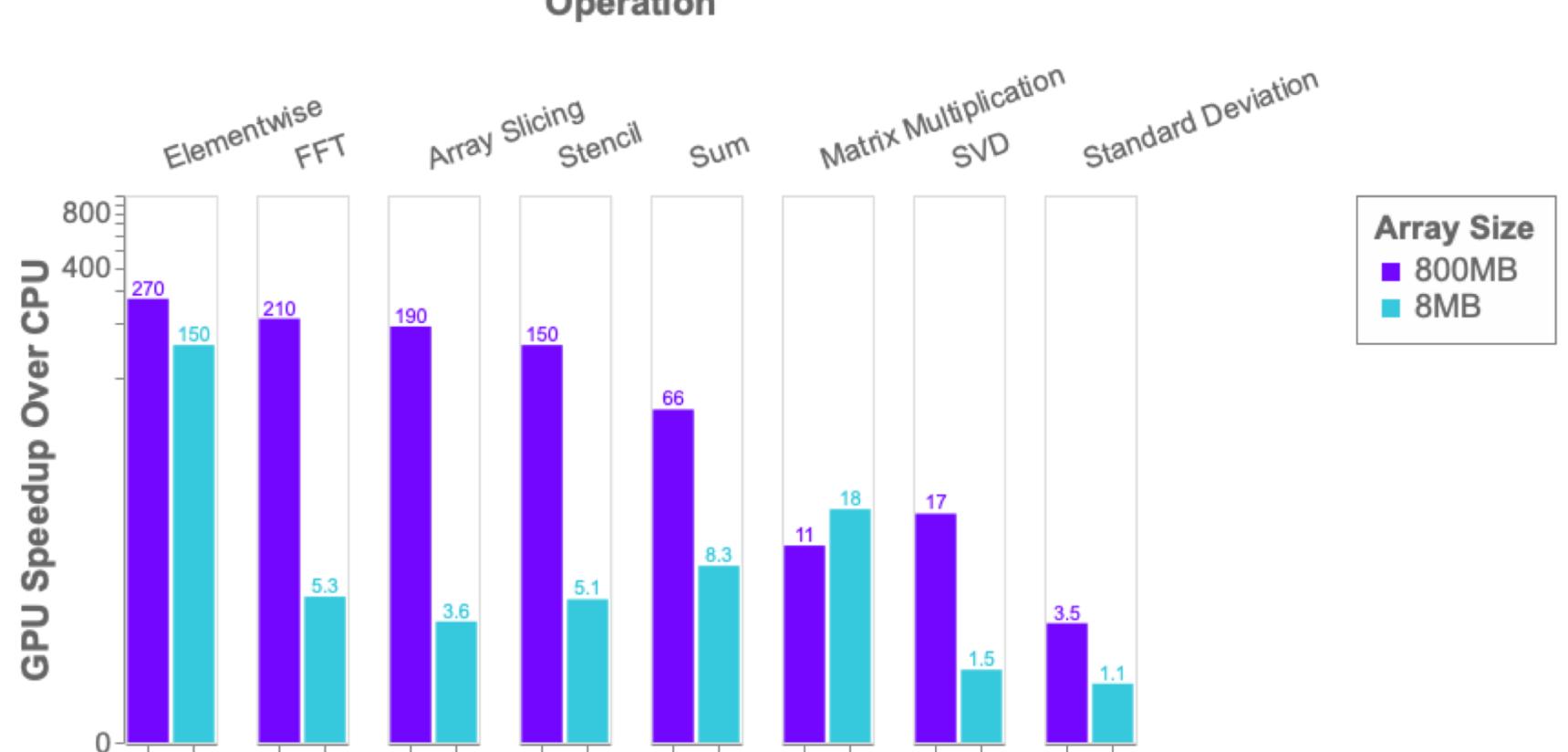


Numba



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

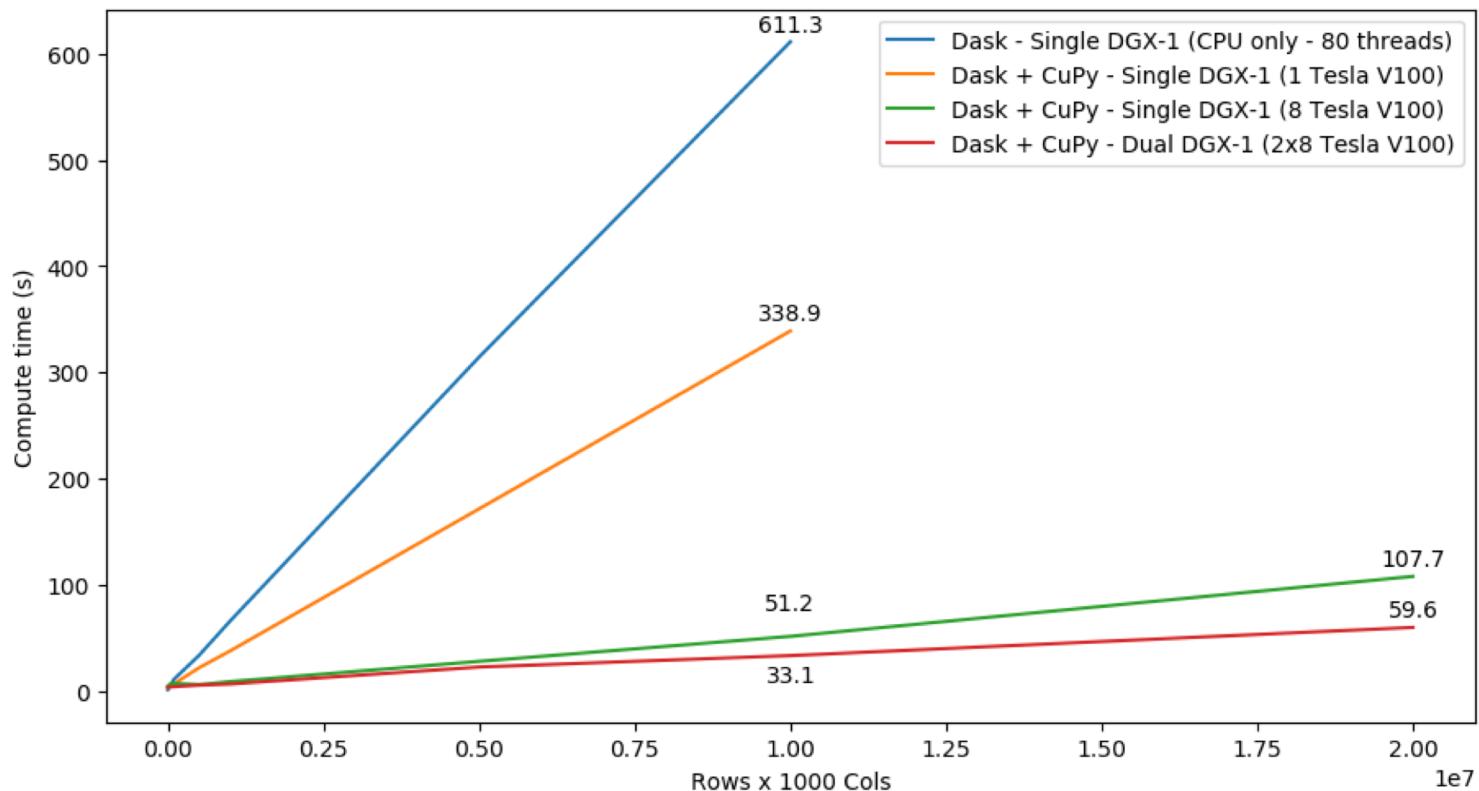
Benchmark: single-GPU CuPy vs NumPy



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

SVD Benchmark

Dask and CuPy Doing Complex Workflows



Also...Achievement Unlocked:

Petabyte Scale Data Analytics with Dask and CuPy

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



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

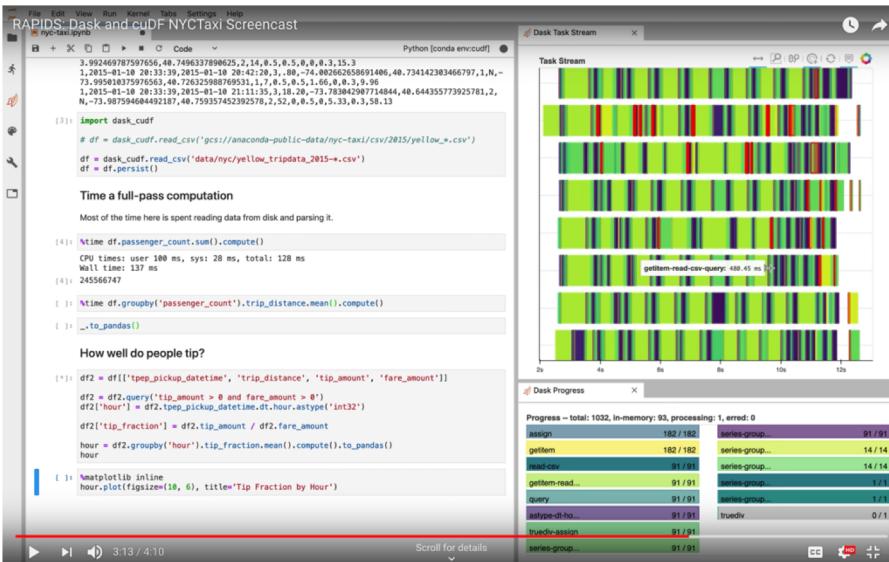
GPU: 4x NVIDIA Tesla P100-16GB-PCIe (total GPU DRAM across nodes 1.22 TB)

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

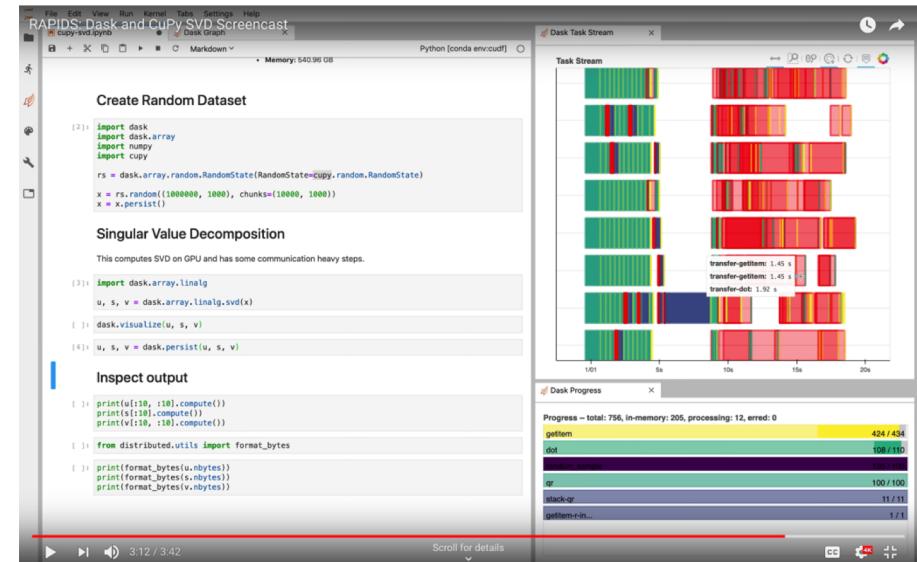
<https://blog.dask.org/2019/01/03/dask-array-gpus-first-steps>

ETL - Arrays and DataFrames

More Dask Awesomeness from RAPIDS



<https://youtu.be/gV0cykgsTPM>

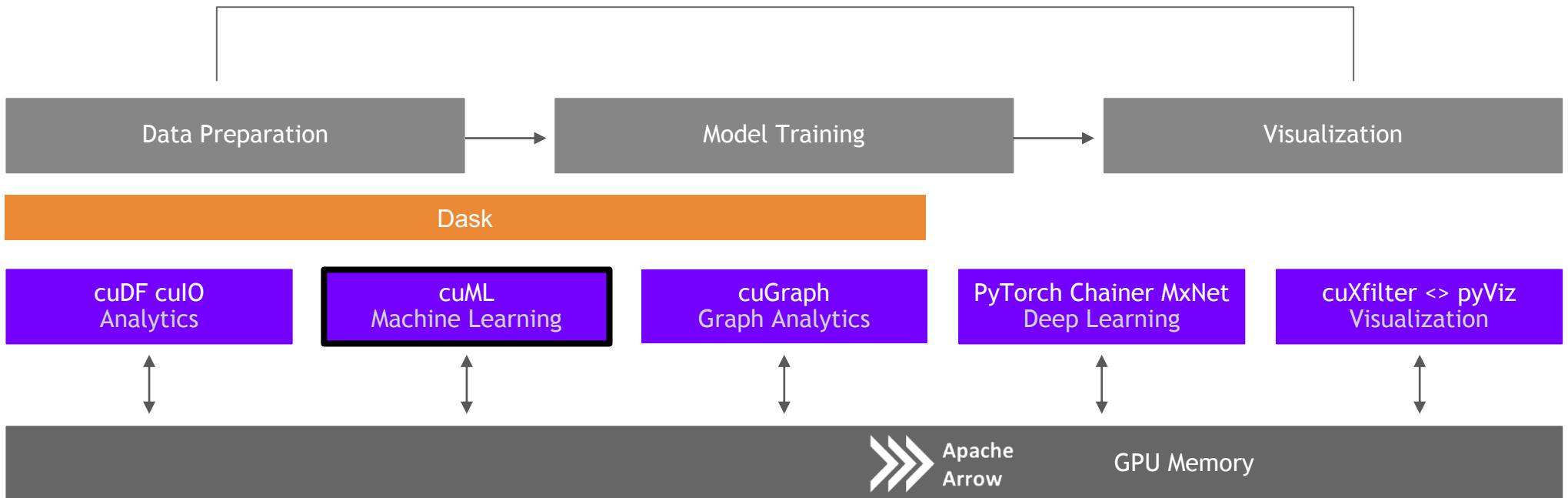


https://youtu.be/R5CiXti_MWo

cuML

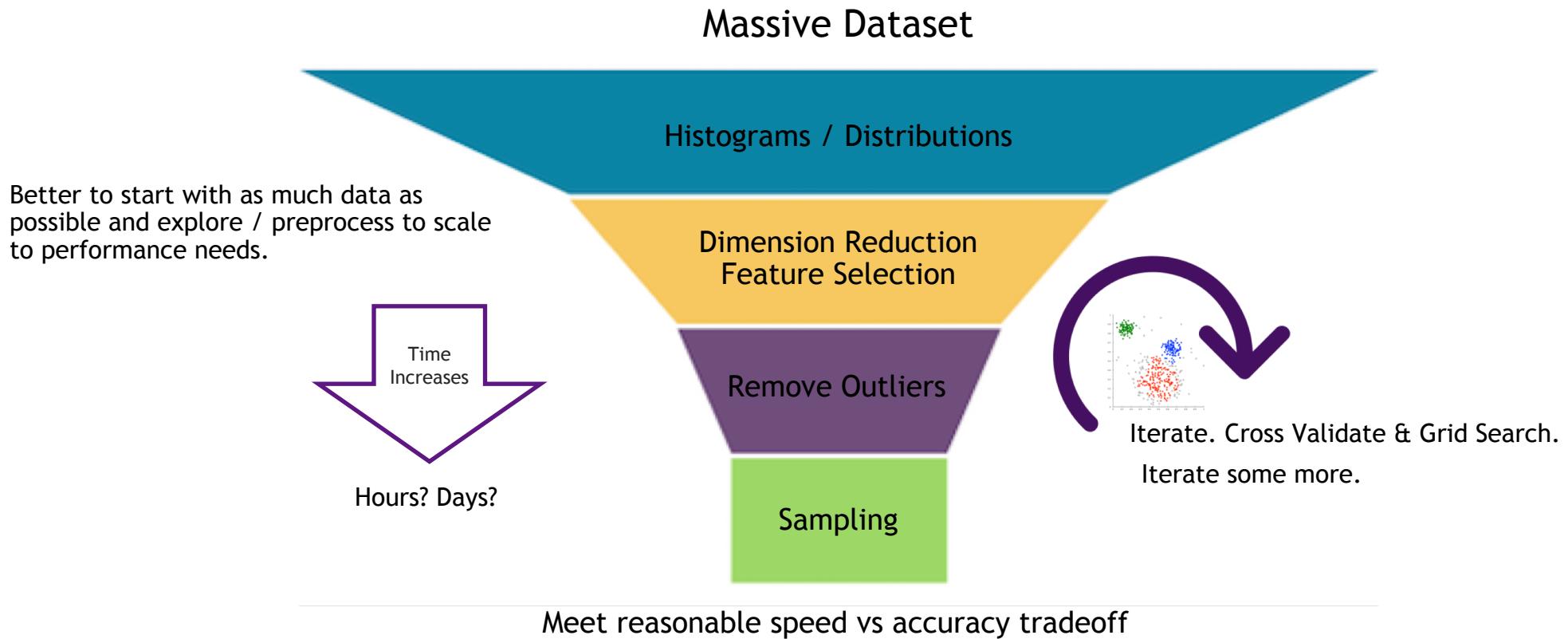
Machine Learning

More models more problems

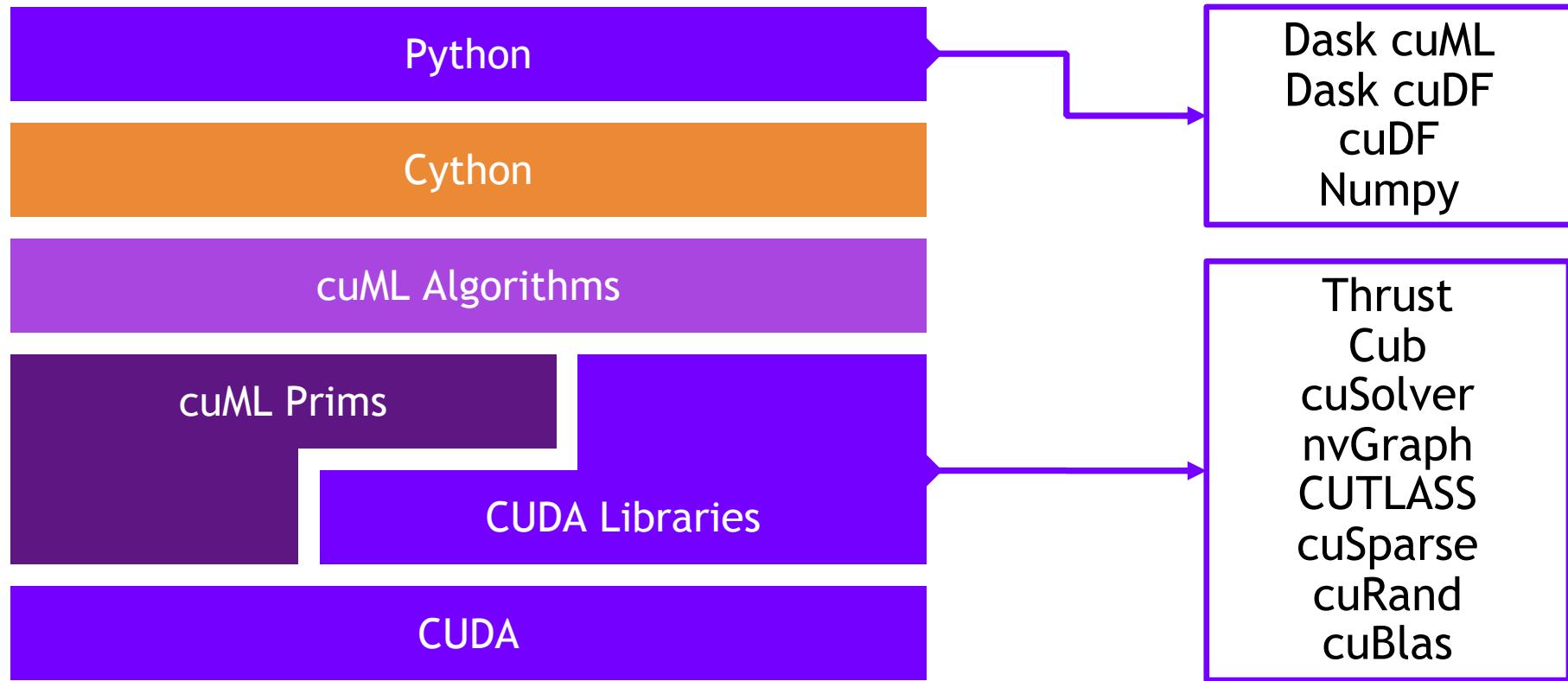


Problem

Data sizes continue to grow

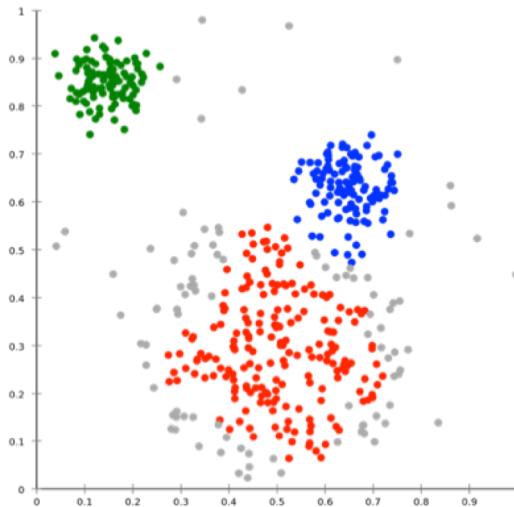


ML Technology Stack



Algorithms

GPU-accelerated Scikit-Learn



Cross Validation

Hyper-parameter Tuning

More to come!

Classification / Regression

Statistical Inference

Clustering

Decomposition & Dimensionality Reduction

Time Series Forecasting

Recommendations

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

Kalman Filtering
Bayesian Inference
Gaussian Mixture Models
Hidden Markov Models

K-Means
DBSCAN
Spectral Clustering

Principal Components
Singular Value Decomposition
UMAP
Spectral Embedding

ARIMA
Holt-Winters

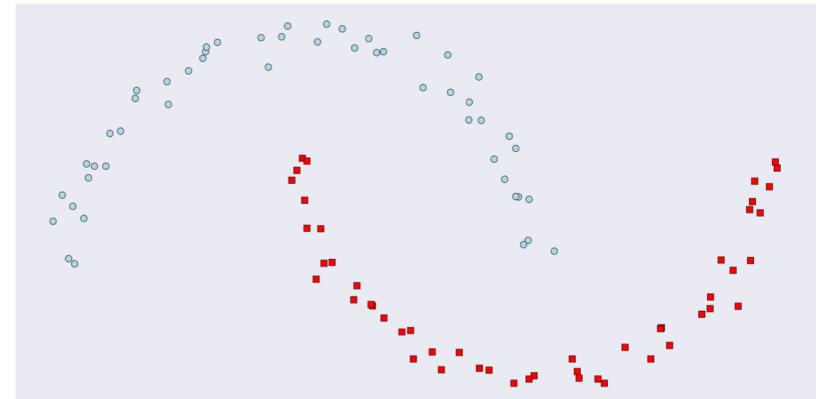
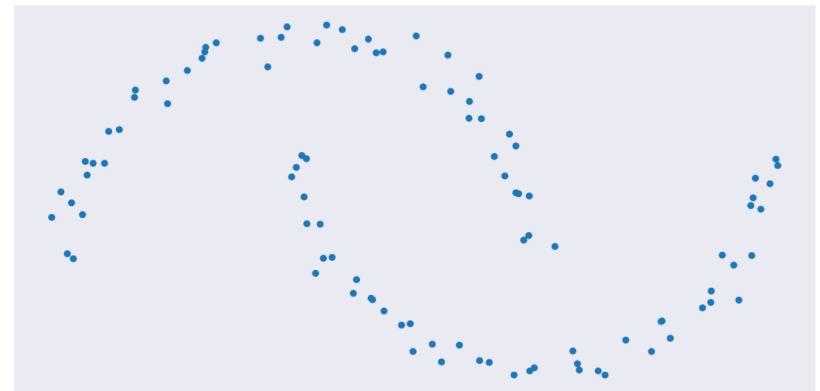
Implicit Matrix Factorization

Clustering

Code Example

```
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)
```



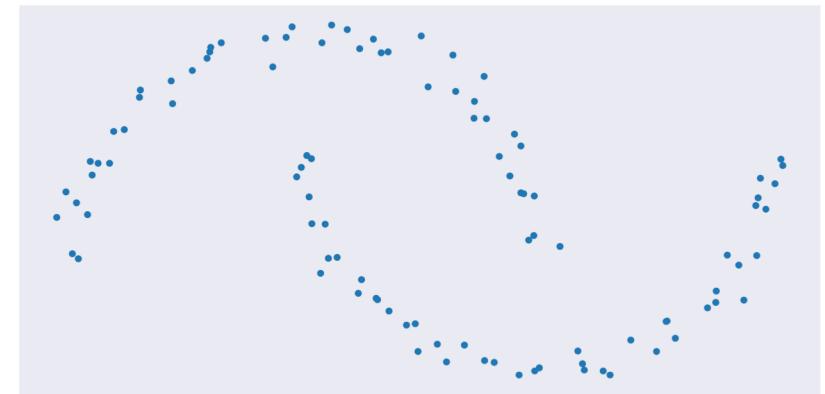
GPU-Accelerated Clustering

Code Example

```
from sklearn.datasets import make_moons
import cudf

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

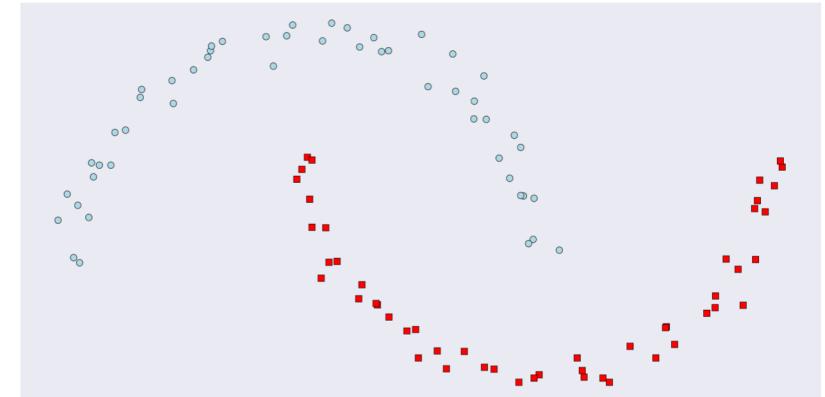
X = cudf.DataFrame({'fea%d'%i: X[:, i]
                    for i in range(X.shape[1])})
```



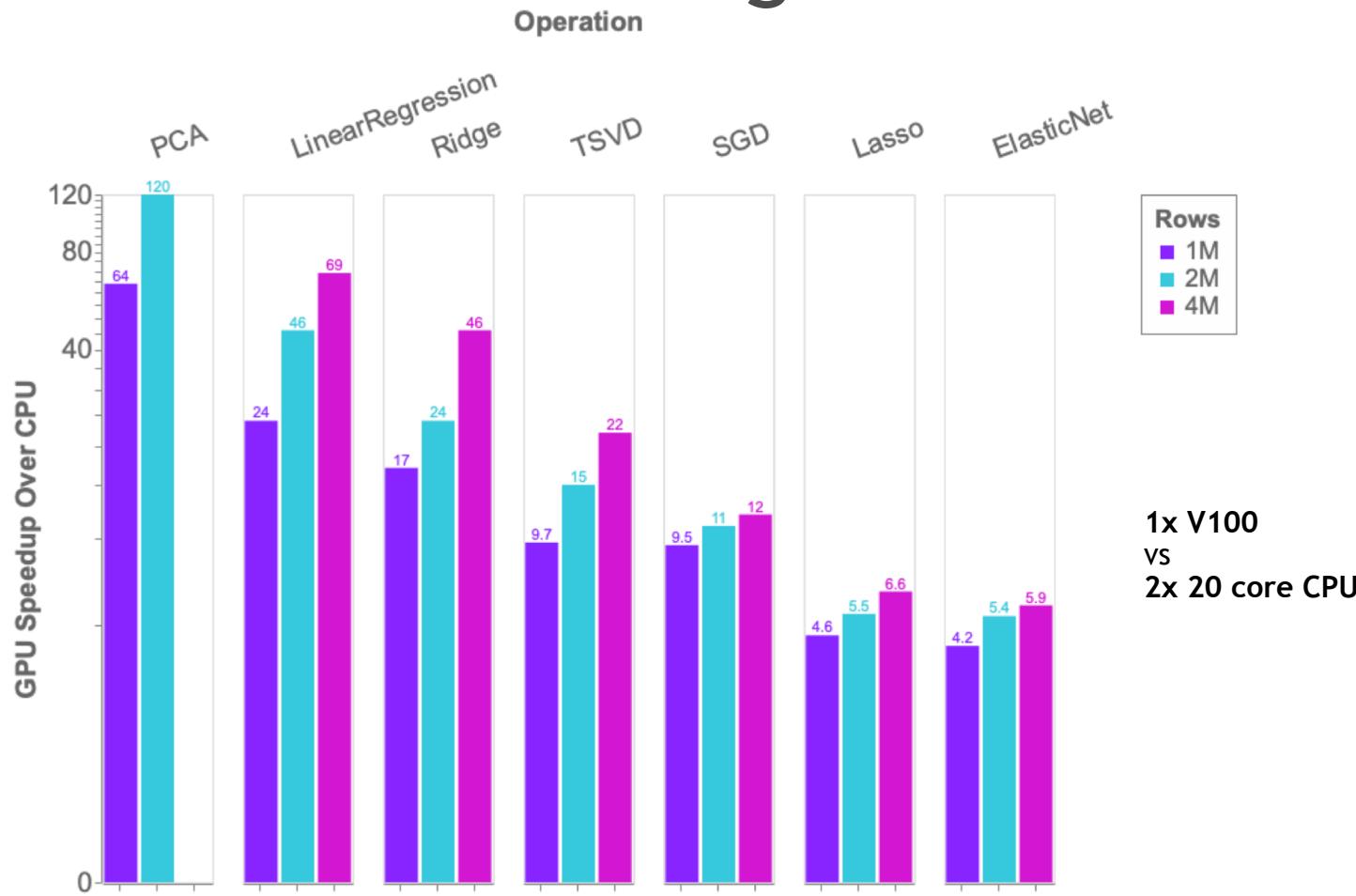
```
from cuml import DBSCAN
dbSCAN = DBSCAN(eps = 0.3, min_samples = 5)

dbSCAN.fit(X)

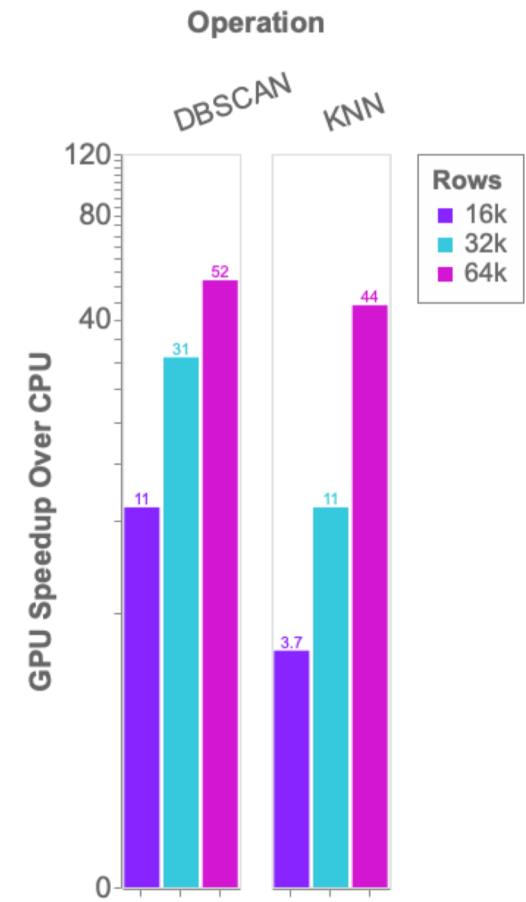
y_hat = dbSCAN.predict(X)
```



Benchmarks: single-GPU cuML vs scikit-learn



1x V100
vs
2x 20 core CPU



Road to 1.0

October 2018 - RAPIDS 0.1

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			
ARIMA & Holts-Winters			
Kalman Filter			
t-SNE			
Principal Components			
Singular Value Decomposition			

Road to 1.0

June 2019 - RAPIDS 0.8

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			
ARIMA & Holts-Winters			
Kalman Filter			
t-SNE			
Principal Components			
Singular Value Decomposition			

Road to 1.0

August 2019 - RAPIDS 0.9

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			
ARIMA & Holts-Winters			
Kalman Filter			
t-SNE			
Principal Components			
Singular Value Decomposition			

Road to 1.0

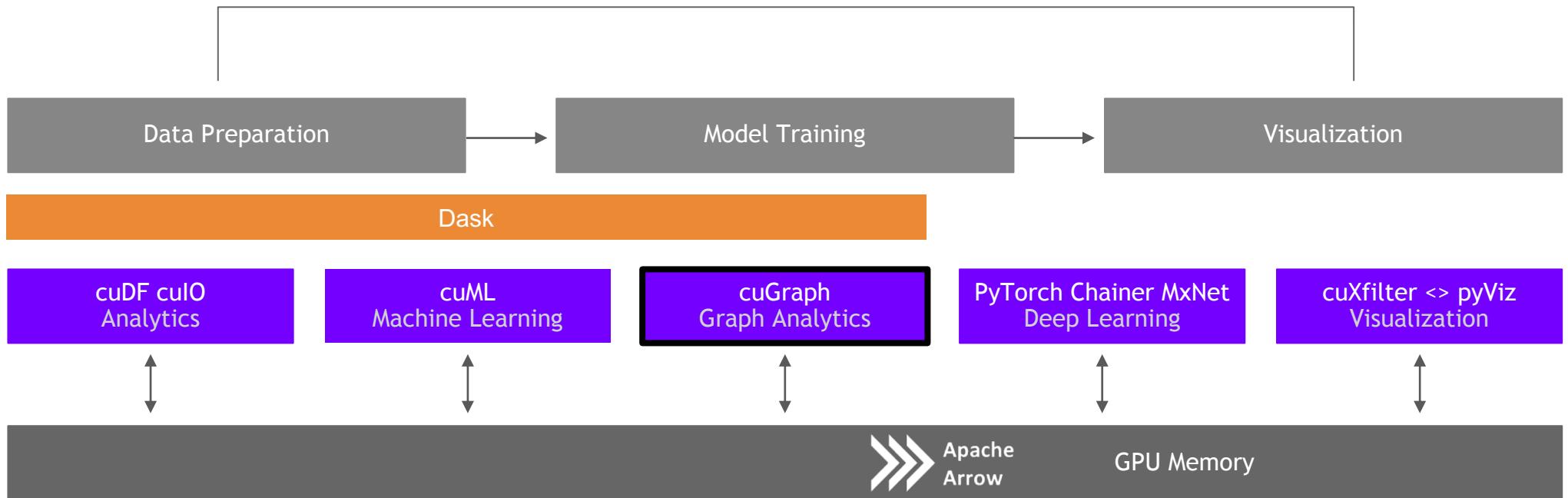
January 2020 - RAPIDS 0.12?

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			
ARIMA & Holts-Winters			
Kalman Filter			
t-SNE			
Principal Components			
Singular Value Decomposition			

cuGraph

Graph Analytics

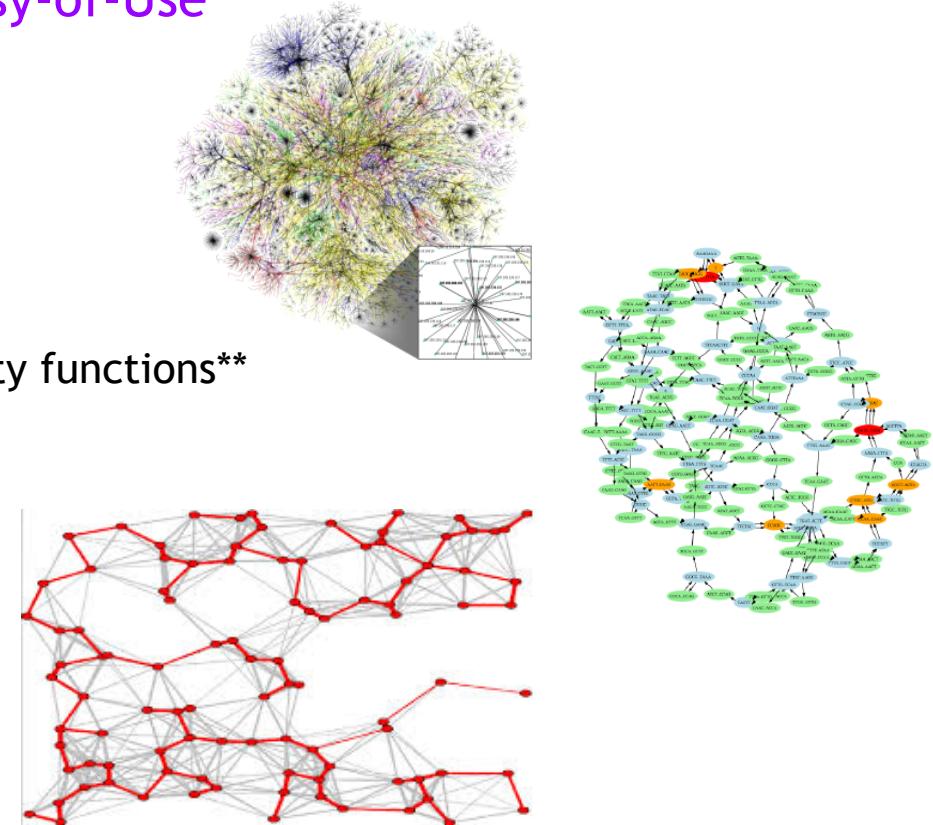
More connections more insights



GOALS AND BENEFITS OF CUGRAPH

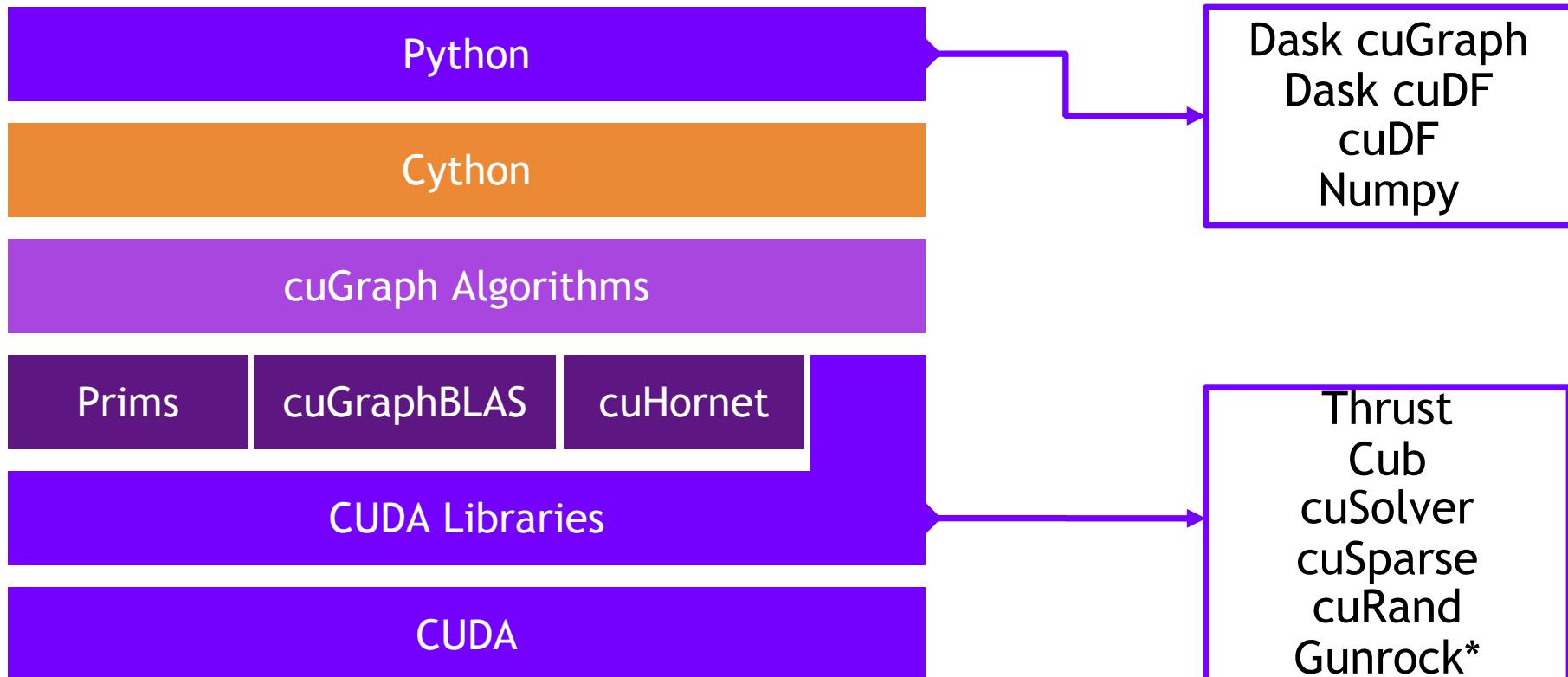
Focus on Features an Easy-of-Use

- Seamless integration with cuDF and cuML
 - Python APIs accepts and returns cuDF DataFrames
- Features
 - Extensive collection of algorithm, primitive, and utility functions**
 - Python API:
 - Multiple APIs: NetworkX, Pregel**, Frontier**
 - Graph Query Language**
 - C/C++
 - Full featured C++ API
 - Ease-of-use and lower-level granular control
- Breakthrough Performance



** On Roadmap

Graph Technology Stack

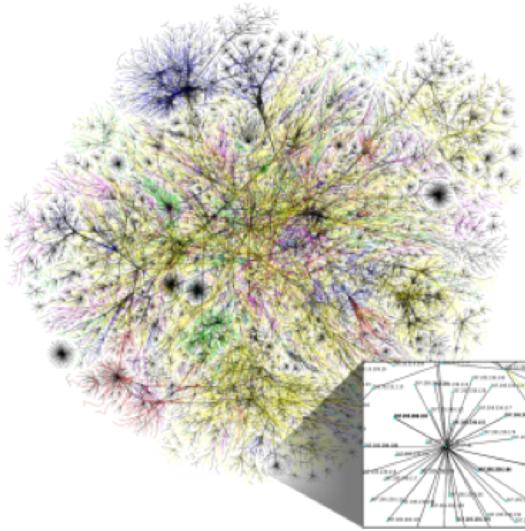


nvGRAPH has been Opened Sourced and integrated into cuGraph. A legacy version is available in a RAPIDS GitHub repo

* Gunrock is from UC Davis

Algorithms

GPU-accelerated NetworkX



Query Language

Multi-GPU

Utilities

More to come!

Community

Components

Link Analysis

Link Prediction

Traversal

Structure

Spectral Clustering
Balanced-Cut
Modularity Maximization
Louvain
Subgraph Extraction
Triangle Counting

Weakly Connected Components
Strongly Connected Components

Page Rank
Personal Page Rank

Jaccard
Weighted Jaccard
Overlap Coefficient

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

COO-to-CSR
Transpose
Renumbering

Louvain Single Run

```
G = cugraph.Graph()  
G.add_edge_list(gdf["src_0"], gdf["dst_0"], gdf["data"])  
df, mod = cugraph.nvLouvain(G)
```

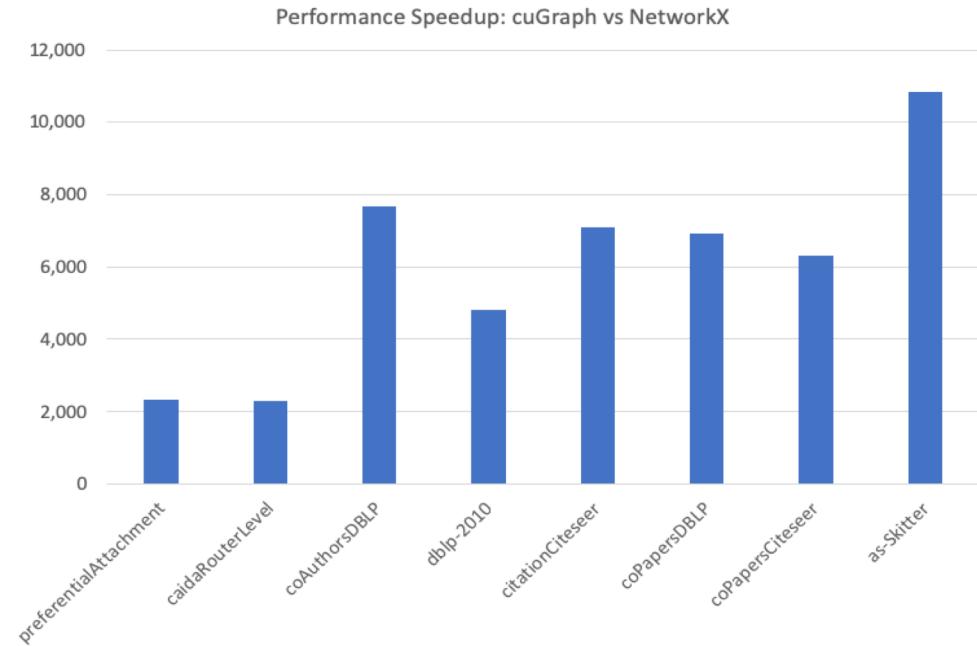
Louvain returns:

cudf.DataFrame with two names columns:

louvain["vertex"]: The vertex id.

louvain["partition"]: The assigned partition.

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

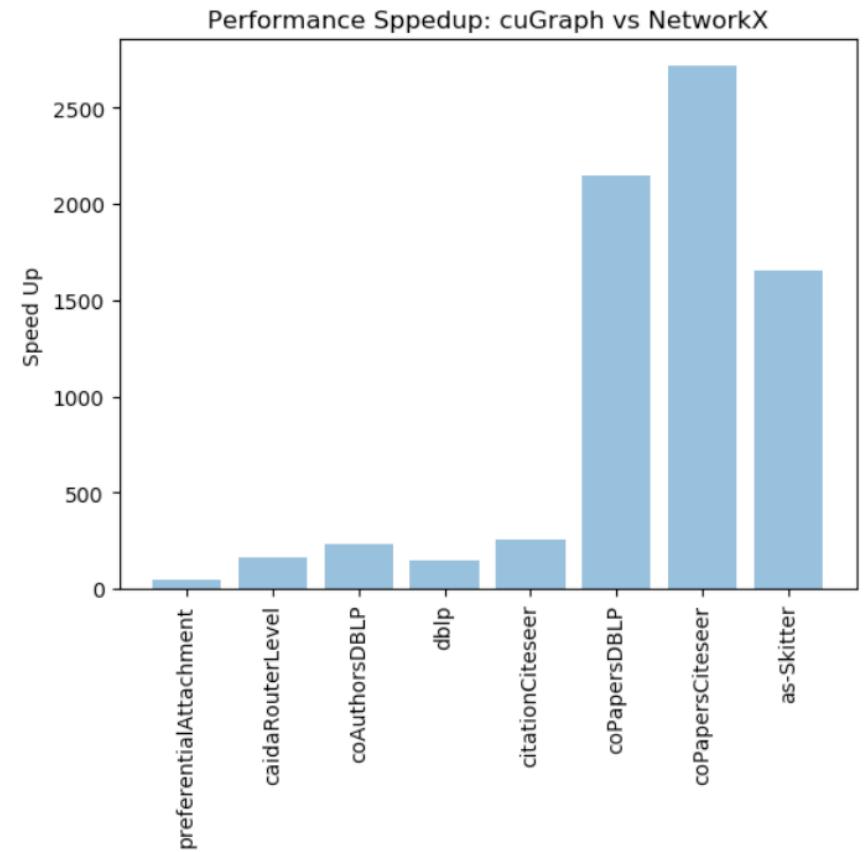


PageRank Speedup

cuGraph PageRank vs NetworkX PageRank

```
G = cugraph.Graph()  
G.add_edge_list(gdf['src'], gdf['dst'], None)  
  
df = cugraph.pagerank(G, alpha, max_iter, tol)
```

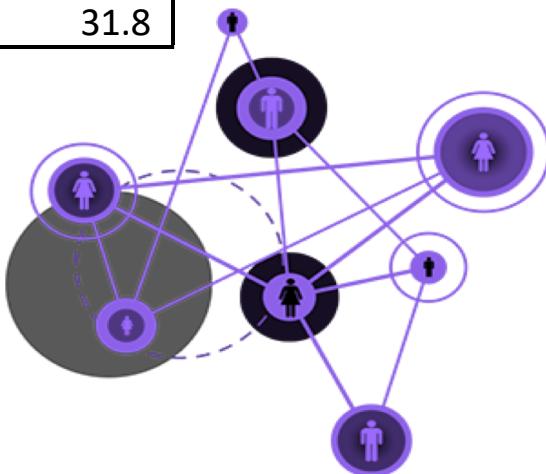
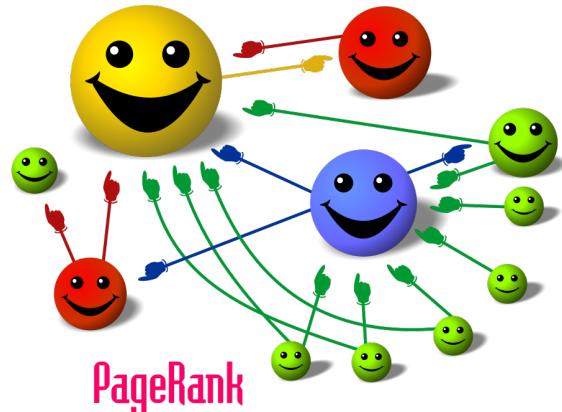
File Name	Num of Vertices	Num of Edges
preferentialAttachment	100,000	999,970
caidaRouterLevel	192,244	1,218,132
coAuthorsDBLP	299,067	1,955,352
dblp-2010	326,186	1,615,400
citationCiteseer	268,495	2,313,294
coPapersDBLP	540,486	30,491,458
coPapersCiteseer	434,102	32,073,440
as-Skitter	1,696,415	22,190,596



Multi-GPU PageRank Performance

HiBench Websearch benchmark

HiBench Scale	Vertices	Edges	CSV File (GB)	# GPUs	PageRank (3 Iterations)
huge	5,000,000	198,000,000	3	1	1.1
bigdata	50,000,000	1,980,000,000	34	3	5.1
bigdata x 2	100,000,000	4,000,000,000	69	6	9.0
bigdata x 4	200,000,000	8,000,000,000	146	12	18.2
bigdata x 8	400,000,000	16,000,000,000	300	16	31.8



Road to 1.0

June 2019 - RAPIDS 0.8

cuGraph	Single-GPU	Multi-GPU	Multi-Node-Multi-GPU
Jaccard and Weighted Jaccard			
Page Rank			
Personal Page Rank			
SSSP			
BFS			
Triangle Counting			
Subgraph Extraction			
Katz Centrality			
Betweenness Centrality			
Connected Components			
Louvain			
Spectral Clustering			
InfoMap			
K-Cores			

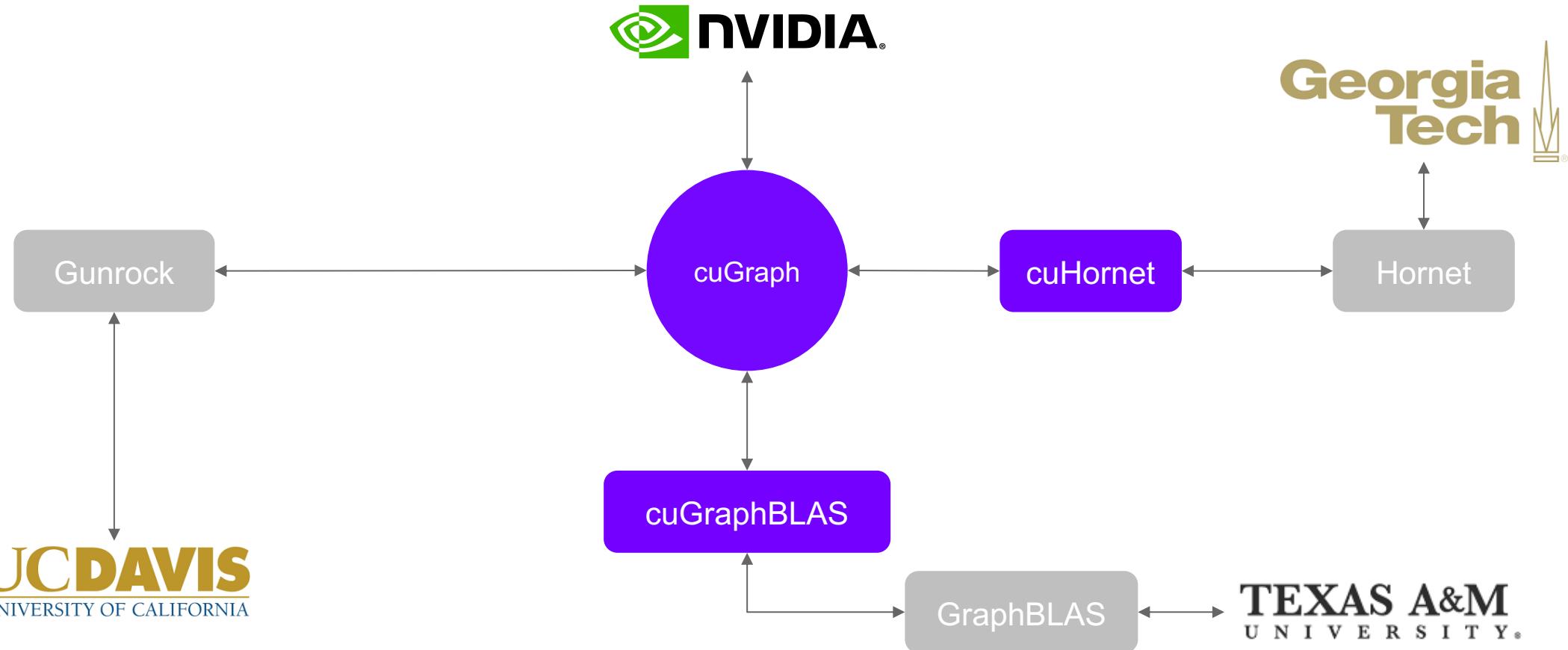
Road to 1.0

August 2019 - RAPIDS 0.9

cuGraph	Single-GPU	Multi-GPU	Multi-Node-Multi-GPU
Jaccard and Weighted Jaccard			
Page Rank			
Personal Page Rank			
SSSP			
BFS			
Triangle Counting			
Subgraph Extraction			
Katz Centrality			
Betweenness Centrality			
Connected Components			
Louvain			
Spectral Clustering			
InfoMap			
K-Cores			

Bringing in leading researchers

Leveraging the great work of others



Road to 1.0

January 2020 - RAPIDS 0.12?

cuGraph	Single-GPU	Multi-GPU	Multi-Node-Multi-GPU
Jaccard and Weighted Jaccard			
Page Rank			
Personal Page Rank			
SSSP			
BFS			
Triangle Counting			
Subgraph Extraction			
Katz Centrality			
Betweenness Centrality			
Connected Components			
Louvain			
Spectral Clustering			
InfoMap			
K-Cores			

Example Notebooks

- Cyber
- Music
- Bicliques / Frequent Item Set
- Weighted Link Prediction
- DASK Example



How to Get Started

RAPIDS Docs

New, improved, and easier to use

The screenshot shows a web browser window for the cuDF - Stable Docs - RAPIDS Docs API Reference. The URL is <https://docs.rapids.ai/api/cudf/stable/>. The page title is "cuDF - Stable Docs - RAPIDS Docs". The top navigation bar includes links for "APIS", "cuDF", "cuML", "cuGraph", "nvStrings", "LIBS", "libcudf", and "RMM". Below this, there are tabs for "STABLE" (selected), "NIGHTLY", and "LEGACY". A dropdown menu for "APIS" is open. On the left, a sidebar titled "RAPIDS Docs" has a "RETURN TO API DOCS" link. It contains a search bar and a "CONTENTS:" section with a tree view of API Reference, including "DataFrame", "Series", "Groupby", "IO", "GpuArrowReader", "10 Minutes to cuDF", "Multi-GPU with Dask-cuDF", and "Developer Documentation". The main content area is titled "API Reference" and "DataFrame". It shows the class definition `class cudf.dataframe.DataFrame(name_series=None, index=None)`, a description ("A GPU Dataframe object."), and examples. The first example shows how to build a DataFrame with `__setitem__`:

```
>>> import cudf
>>> df = cudf.DataFrame()
>>> df['key'] = [0, 1, 2, 3, 4]
>>> df['val'] = [float(i + 10) for i in range(5)] # insert column
>>> print(df)
   key    val
0    0  10.0
1    1  11.0
2    2  12.0
3    3  13.0
4    4  14.0
```

The second example shows how to build a DataFrame with an initializer:

```
>>> import cudf
>>> import numpy as np
>>> from datetime import datetime, timedelta
>>> ids = np.arange(5)
```

<https://docs.rapids.ai>

RAPIDS Docs

Easier than ever to get started with cuDF

The screenshot shows a web browser window titled "cuDF - Stable Docs - RAPIDS". The URL is <https://docs.rapids.ai/api/cudf/stable/>. The page content is the "10 Minutes to cuDF" guide. The sidebar on the left lists various topics under "10 Minutes to cuDF", including Object Creation, Viewing Data, Selection, Getting, Selection by Label, Selection by Position, Boolean Indexing, Setting, Missing Data, Operations, Stats, Applymap, Histogramming, String Methods, Merge, Concat, Join, Append, Grouping, Reshaping, Time Series, Categoricals, Plotting, and Converting Data Representation. The main content area starts with a heading "10 Minutes to cuDF" and a sub-section "Object Creation". It includes code snippets and explanatory text.

10 Minutes to cuDF

Modeled after 10 Minutes to Pandas, this is a short introduction to cuDF, geared mainly for new users.

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

### Portions of this were borrowed from the
### cuDF cheatsheet, existing cuDF documentation,
### and 10 Minutes to Pandas.
### Created November, 2018.
```

Object Creation

Creating a `Series`.

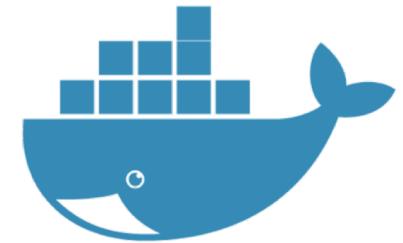
```
[2]: s = cudf.Series([1,2,3,None,4])
print(s)
```

0	1
1	2
2	3
3	
4	4

Creating a `DataFrame` by specifying values for each column.

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

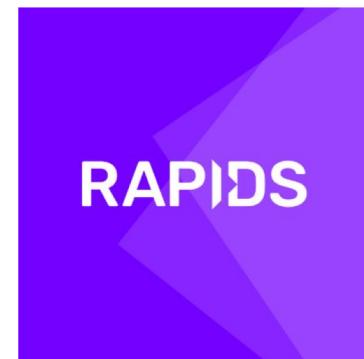
Everyone can help!



APACHE ARROW

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

@ApacheArrow



RAPIDS

<https://rapids.ai>

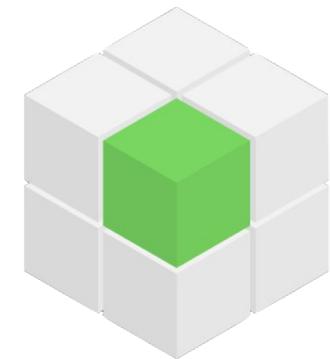
@RAPIDSAI



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

Brad Rees

brees@nvidia.com

@BradReesWork

Corey Nolet

cnolet@nvidia.com

@cjnolet



RAPIDS