

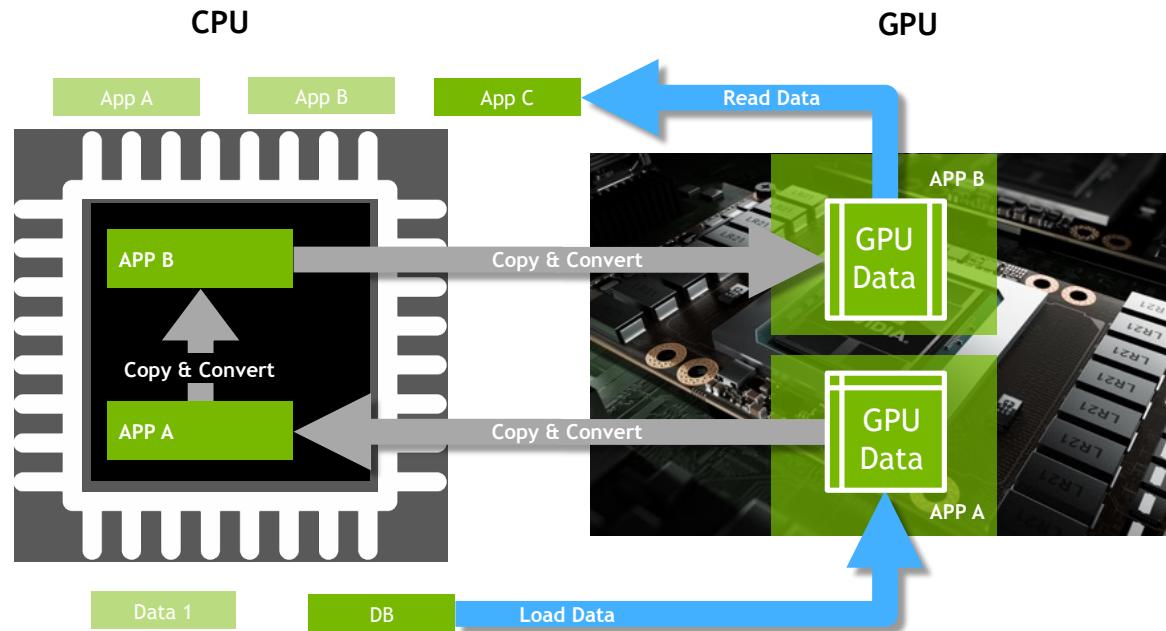


WHY RAPIDS

IMPETUS

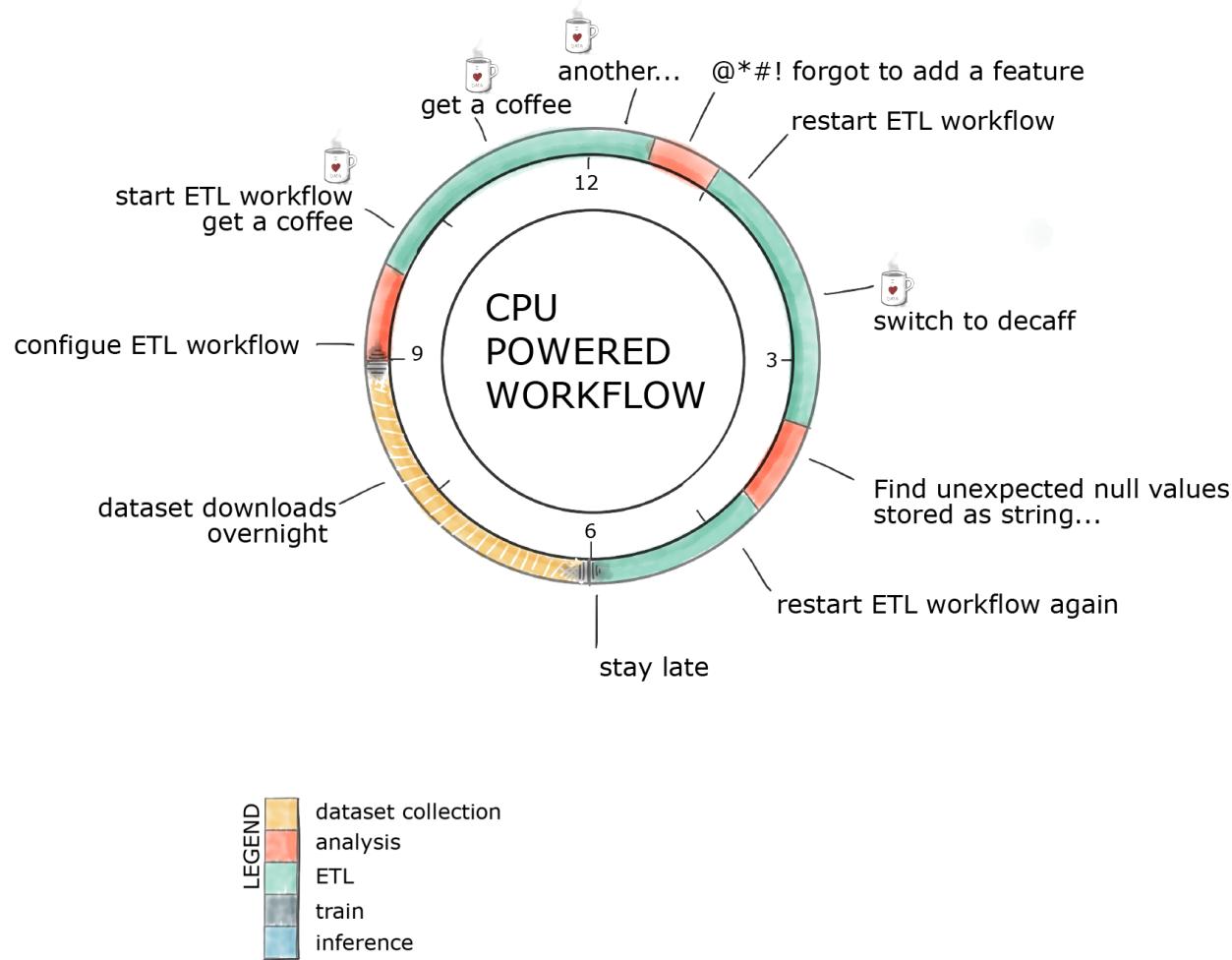
No one was *addressing* the big issue

- Alternative tools and frameworks
 - Too many data formats (data marshalling)
 - Too much data movement (in and out of GPU to perform marshalling)
 - ❖ Too much focus on individual algorithms and not end-to-end performance



DAY IN THE LIFE

Or Why I want to be a Data Scientist

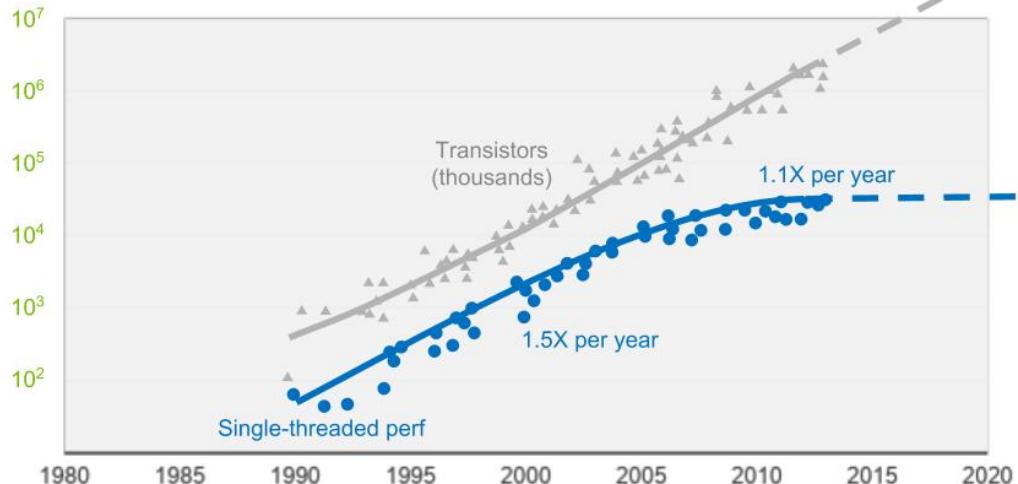


The average data scientist spends 90+% of their time in ETL as opposed to training models



SEEMS LIKE A LOT OF UNPRODUCTIVE TIME?

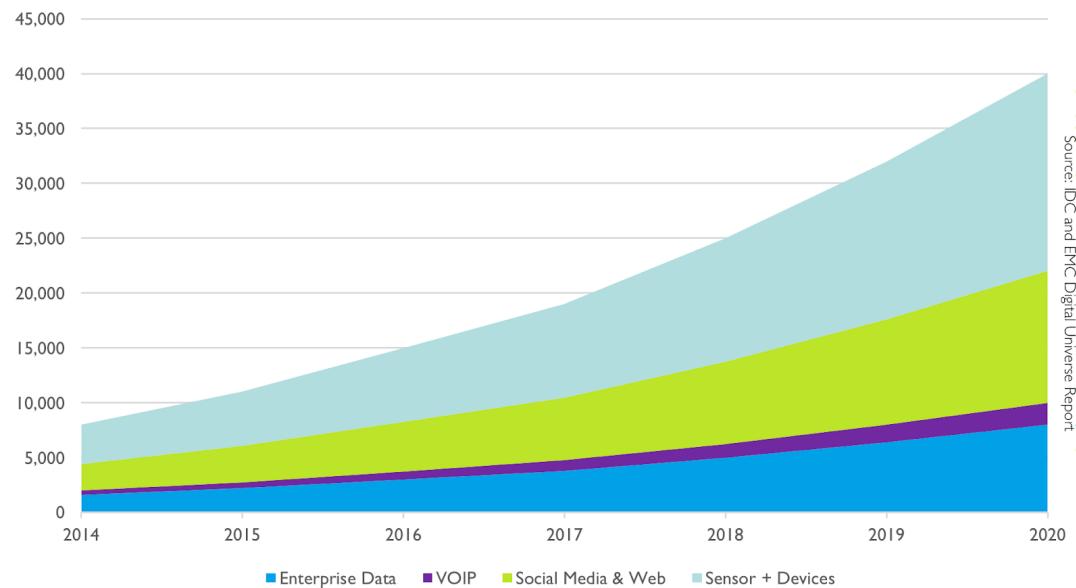
Post-Moore's law
CPU Performance Has Plateaued



Moore's law is no longer a predictor of capacity in CPU market growth

- Distributing CPUs exacerbates the problem

Data sizes continue to grow
Data Growth and Source in Exabytes



WE NEED MORE COMPUTE!

Basic workloads are bottlenecked by the CPU

- In a simple benchmark consisting of aggregating data, the **CPU is the bottleneck**
- This is after the data is parsed and cached into memory which is another common bottleneck
- The CPU bottleneck is even worse in more complex workloads!

SELECT cab_type, count(*) FROM trips_orc GROUP BY cab_type;

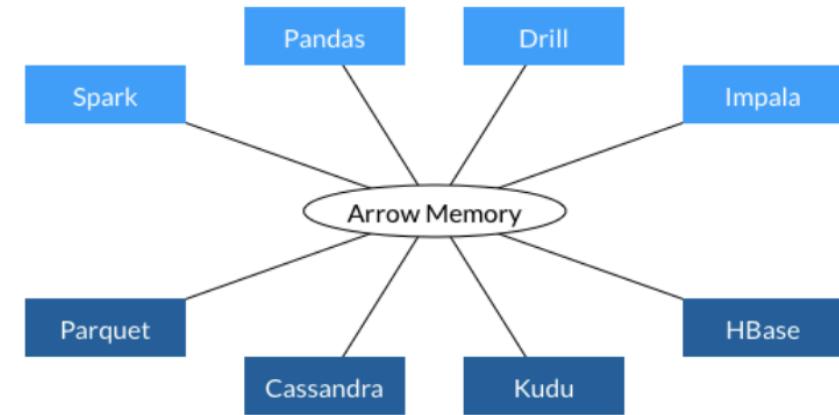
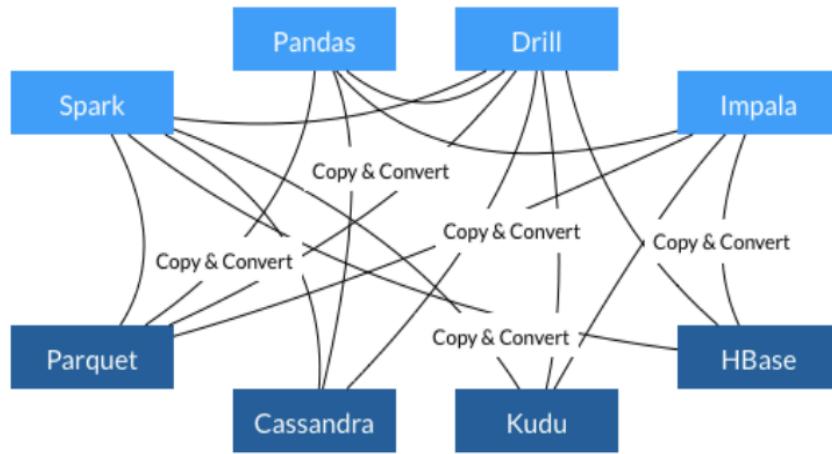
```
top - 08:54:14 up 1:50, 4 users, load average: 0.20, 1.64, 6.43
Tasks: 360 total, 2 running, 358 sleeping, 0 stopped, 0 zombie
%Cpu0 : 94.7 us, 1.7 sy, 0.0 ni, 3.3 id, 0.0 wa, 0.0 hi, 0.3 si, 0.0
%Cpu1 : 95.0 us, 1.7 sy, 0.0 ni, 3.4 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu2 : 98.3 us, 0.3 sy, 0.0 ni, 1.3 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu3 : 87.3 us, 4.3 sy, 0.0 ni, 8.4 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu4 : 95.0 us, 1.3 sy, 0.0 ni, 3.7 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu5 : 98.3 us, 0.0 sy, 0.0 ni, 1.7 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu6 : 96.7 us, 1.3 sy, 0.0 ni, 2.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu7 : 92.7 us, 1.0 sy, 0.0 ni, 5.6 id, 0.3 wa, 0.0 hi, 0.3 si, 0.0
%Cpu8 : 93.7 us, 1.3 sy, 0.0 ni, 5.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu9 : 92.3 us, 0.7 sy, 0.0 ni, 7.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu10 : 97.3 us, 0.7 sy, 0.0 ni, 2.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu11 : 97.3 us, 0.7 sy, 0.0 ni, 2.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu12 : 92.0 us, 3.0 sy, 0.0 ni, 5.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu13 : 94.9 us, 1.0 sy, 0.0 ni, 4.0 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu14 : 88.3 us, 3.0 sy, 0.0 ni, 8.7 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu15 : 92.6 us, 2.3 sy, 0.0 ni, 4.7 id, 0.0 wa, 0.0 hi, 0.3 si, 0.0
%Cpu16 : 94.7 us, 2.3 sy, 0.0 ni, 2.6 id, 0.0 wa, 0.0 hi, 0.3 si, 0.0
%Cpu17 : 93.0 us, 0.7 sy, 0.0 ni, 6.0 id, 0.0 wa, 0.0 hi, 0.3 si, 0.0
%Cpu18 : 93.0 us, 3.7 sy, 0.0 ni, 3.3 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
%Cpu19 : 91.2 us, 0.7 sy, 0.0 ni, 8.1 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0
```

HOW CAN WE DO BETTER?

- Focus on the full Data Science workflow
 - Data Loading
 - Data Transformation
 - Data Analytics
- Python
 - Provide close to a drop-in replacement for existing tools
- Performance - Leverage GPUs



LEARNING FROM APACHE ARROW ➤

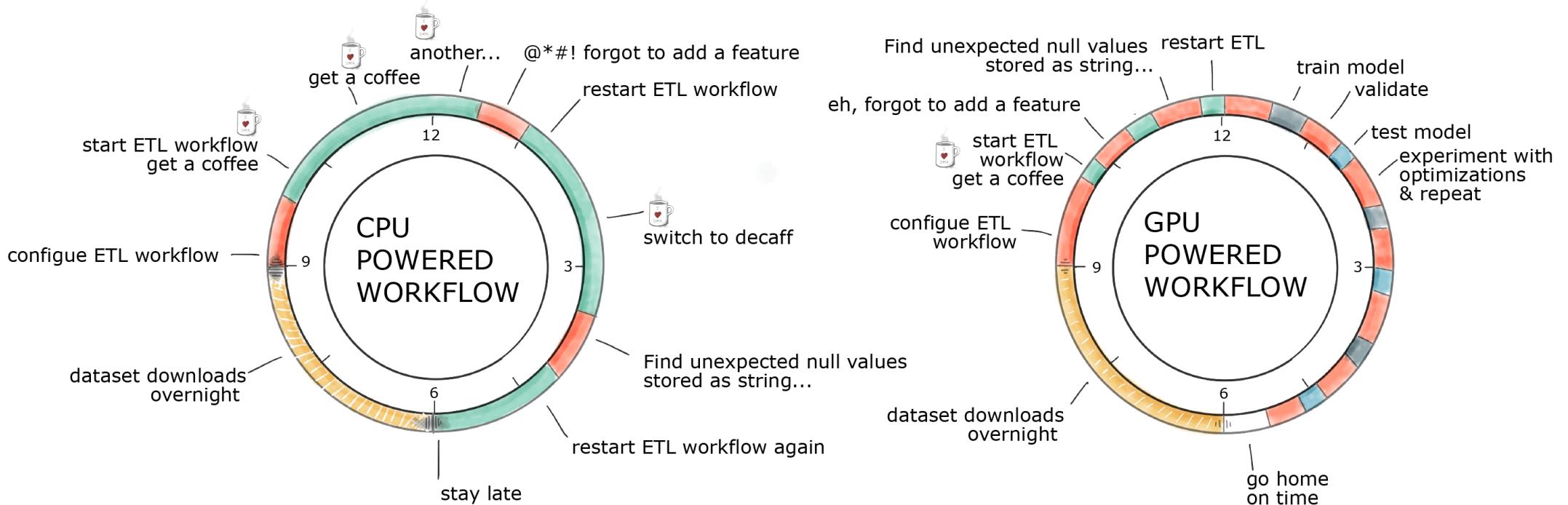


- Each system has its own internal memory format
- 70-80% computation wasted on serialization and deserialization
- Similar functionality implemented in multiple projects

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

DAY IN THE LIFE

Data Scientist Using RAPIDS

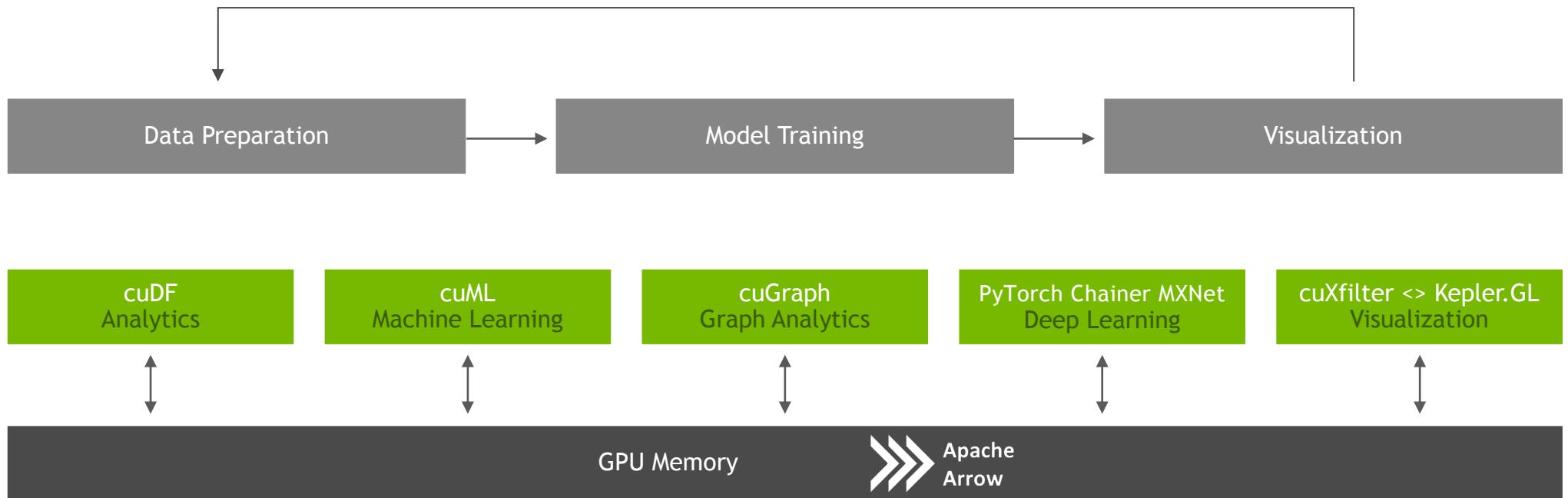


LEGEND

Yellow	dataset collection
Red	analysis
Teal	ETL
Blue	train
Cyan	inference



RAPIDS OPEN SOURCE SOFTWARE



We will touch on many of these components

“More data requires better approaches!”

~ Xavier Amatriain (CTO of Curai)

CUDF

Sorry about the upcoming wall of words

DATA PROCESSING ENGINE

cuDF

The World's First Accelerated Data Frame Library

Over a decade of NVIDIA Engineering has made this possible
Ease of use of pandas, scalability of Spark, in Python the most used data science language
Performs up to 50x faster than Spark on a single V100*

Key Features

- Complete support for all standard data frame routines found in Pandas, Spark, R, or SAS
- Single, Double, INT, Float, String, Dictionary, Date Time support and more
- Support for multiple GPUs and concurrent kernels
- Supports CUDA streams for concurrent operations
- Python Bindings, Numba
- Device API that can be called with your own CUDA kernels
- I/O & Accelerated parsing of standard formats: CSV, Parquet, ORC, JSON, XML
- JDBC, ODBC, and HDFS Connectors
- Direct integrations with cuArray, cuGraph, & cuML (no copies needed)

CUDF

GPU DataFrame library

	Area Abbreviation	Area Code	Area	Item Code	Item	Element Code	Element	Unit	latitude	longitude	...	Y2004	Y2005	Y2006	Y2007	Y2008
0	AF	2	Afghanistan	2511	Wheat and products	5142	Food	1000 tonnes	33.94	67.71	...	3249.0	3486.0	3704.0	4164.0	4252.0
1	AF	2	Afghanistan	2805	Rice (Milled Equivalent)	5142	Food	1000 tonnes	33.94	67.71	...	419.0	445.0	546.0	455.0	490.0
2	AF	2	Afghanistan	2513	Barley and products	5521	Feed	1000 tonnes	33.94	67.71	...	58.0	236.0	262.0	263.0	230.0
3	AF	2	Afghanistan	2513	Barley and products	5142	Food	1000 tonnes	33.94	67.71	...	185.0	43.0	44.0	48.0	62.0
4	AF	2	Afghanistan	2514	Maize and products	5521	Feed	1000 tonnes	33.94	67.71	...	120.0	208.0	233.0	249.0	247.0
5	AF	2	Afghanistan	2514	Maize and products	5142	Food	1000 tonnes	33.94	67.71	...	231.0	67.0	82.0	67.0	69.0
6	AF	2	Afghanistan	2517	Millet and products	5142	Food	1000 tonnes	33.94	67.71	...	15.0	21.0	11.0	19.0	21.0
7	AF	2	Afghanistan	2520	Cereals, Other	5142	Food	1000 tonnes	33.94	67.71	...	2.0	1.0	1.0	0.0	0.0
8	AF	2	Afghanistan	2531	Potatoes and products	5142	Food	1000 tonnes	33.94	67.71	...	276.0	294.0	294.0	260.0	242.0
9	AF	2	Afghanistan	2536	Sugar cane	5521	Feed	1000 tonnes	33.94	67.71	...	50.0	29.0	61.0	65.0	54.0
10	AF	2	Afghanistan	2537	Sugar beet	5521	Feed	1000 tonnes	33.94	67.71	...	0.0	0.0	0.0	0.0	0.0

- Apache Arrow data format
- Pandas-like API
- Unary and Binary Operations
- Joins / Merges
- Group Bys
- Filters
- User-Defined Functions (UDFs)
- Accelerated file readers
- Etc.

Methods

<code>add_column</code> (name, data[, forceindex])	Add a column
<code>apply_chunks</code> (func, incols, outcols[, ...])	Transform user-specified chunks using the user-provided function.
<code>apply_rows</code> (func, incols, outcols, kwargs[, ...])	Apply a row-wise user defined function.
<code>as_gpu_matrix</code> ([columns, order])	Convert to a matrix in device memory.
<code>as_matrix</code> ([columns])	Convert to a matrix in host memory.
<code>assign</code> (**kwargs)	Assign columns to DataFrame from keyword arguments.
<code>copy</code> ([deep])	Returns a copy of this dataframe
<code>drop</code> (labels)	Drop column(s)
<code>drop_column</code> (name)	Drop a column by name
<code>fillna</code> (value[, method, axis, inplace, limit])	Fill null values with <code>value</code> .
<code>from_arrow</code> (table)	Convert from a PyArrow Table.
<code>from_gpu_matrix</code> (data[, index, columns, ...])	Convert from a numba gpu ndarray.
<code>from_pandas</code> (dataframe[, nan_as_null])	Convert from a Pandas DataFrame.
<code>from_records</code> (data[, index, columns, nan_as_null])	Convert from a numpy recarray or structured array.
<code>groupby</code> ([by, sort, as_index, method, level])	Groupby
<code>hash_columns</code> ([columns])	Hash the given <code>columns</code> and return a new Series
<code>head</code> ([n])	Returns the first n rows as a new DataFrame
<code>iteritems</code> ()	Iterate over column names and series pairs
<code>join</code> (other[, on, how, lsuffix, rsuffix, ...])	Join columns with other DataFrame on index or on a key column.
<code>label_encoding</code> (column, prefix, cats[, ...])	Encode labels in a column with label encoding.
<code>merge</code> (right[, on, how, left_on, right_on, ...])	Merge GPU DataFrame objects by performing a database-style join operat
<code>nlargest</code> (n, columns[, keep])	Get the rows of the DataFrame sorted by the n largest value of <code>columns</code>
<code>nsmallest</code> (n, columns[, keep])	Get the rows of the DataFrame sorted by the n smallest value of <code>columns</code>
<code>one_hot_encoding</code> (column, prefix, cats[, ...])	Expand a column with one-hot-encoding.
<code>partition_by_hash</code> (columns, nparts)	Partition the dataframe by the hashed value of data in <code>columns</code> .
<code>quantile</code> ([q, interpolation, columns, exact])	Return values at the given quantile.
<code>query</code> (expr)	Query with a boolean expression using Numba to compile a GPU kernel.
<code>rename</code> ([mapper, columns, copy, inplace])	Alter column labels.
<code>replace</code> (to_replace, value)	Replace values given in <code>to_replace</code> with <code>value</code> .
<code>select_dtypes</code> ([include])	Return a subset of the DataFrame's columns based on the column dtypes.
<code>set_index</code> (index)	Return a new DataFrame with a new index
<code>sort_index</code> ([ascending])	Sort by the index
<code>sort_values</code> (by[, ascending, na_position])	Sort by the values row-wise.
<code>tail</code> ([n])	Returns the last n rows as a new DataFrame
<code>to_arrow</code> ([preserve_index])	Convert to a PyArrow Table.
<code>to_dlpack</code> ()	Converts a cuDF object into a DLPack tensor.
<code>to_feather</code> (path, *args, **kwargs)	Write a DataFrame to the feather format.
<code>to_gpu_matrix</code> ()	Convert to a numba gpu ndarray
<code>to_hdf</code> (path_or_buf, key, *args, **kwargs)	Write the contained data to an HDF5 file using HDFStore.
<code>to_json</code> ([path_or_buf])	Convert the cuDF object to a JSON string.
<code>to_pandas</code> ()	Convert to a Pandas DataFrame.
<code>to_parquet</code> (path, *args, **kwargs)	Write a DataFrame to the parquet format.
<code>to_records</code> ([index])	Convert to a numpy recarray
<code>to_string</code> ([nrows, ncols])	Convert to string
<code>transpose</code> ()	Transpose index and columns.

Methods

<code>append (arbitrary)</code>	Append values from another <code>Series</code> or array-like object.	<code>one_hot_encoding (cats[, dtype])</code>	Perform one-hot-encoding
<code>applymap (udf[, out_dtype])</code>	Apply a elementwise function to transform the values in the Column.	<code>product ([axis, skipna])</code>	Compute the product of the series
<code>argsort ([ascending, na_position])</code>	Returns a Series of int64 index that will sort the series.	<code>quantile ([q[, interpolation, exact, quant_index]])</code>	Return values at the given quantile.
<code>as_mask ()</code>	Convert booleans to bitmask	<code>rename ([index, copy])</code>	Alter Series name.
<code>astype (dtype)</code>	Convert to the given <code>dtype</code> .	<code>replace (to_replace, value)</code>	Replace values given in <code>to_replace</code> with <code>value</code> .
<code>ceil ()</code>	Rounds each value upward to the smallest integral value not less than the original.	<code>reset_index ([drop])</code>	Reset index to RangeIndex
<code>count ([axis, skipna])</code>	The number of non-null values	<code>reverse ()</code>	Reverse the Series
<code>digitize (bins[, right])</code>	Return the indices of the bins to which each value in series belongs.	<code>scale ()</code>	Scale values to [0, 1] in float64
<code>factorize ([na_sentinel])</code>	Encode the input values as integer labels	<code>set_index (index)</code>	Returns a new Series with a different index.
<code>fillna (value[, method, axis, inplace, limit])</code>	Fill null values with <code>value</code> .	<code>set_mask (mask[, null_count])</code>	Create new Series by setting a mask array.
<code>find_first_value (value)</code>	Returns offset of first value that matches	<code>sort_index ([ascending])</code>	Sort by the index.
<code>find_last_value (value)</code>	Returns offset of last value that matches	<code>sort_values ([ascending, na_position])</code>	Sort by the values.
<code>floor ()</code>	Rounds each value downward to the largest integral value not greater than the original.	<code>std ([ddof, axis, skipna])</code>	Compute the standard deviation of the series
<code>from_categorical (categorical[, codes])</code>	Creates from a pandas.Categorical	<code>sum ([axis, skipna])</code>	Compute the sum of the series
<code>from_masked_array (data, mask[, null_count])</code>	Create a Series with null-mask.	<code>tail ([n])</code>	Returns the last n rows as a new Series
<code>hash_encode (stop[, use_name])</code>	Encode column values as ints in [0, stop) using hash function.	<code>take (indices[, ignore_index])</code>	Return Series by taking values from the corresponding <code>indices</code> .
<code>hash_values ()</code>	Compute the hash of values in this column.	<code>to_array ([fillna])</code>	Get a dense numpy array for the data.
<code>label_encoding (cats[, dtype, na_sentinel])</code>	Perform label encoding	<code>to_dlpark ()</code>	Converts a cuDF object into a DLPack tensor.
<code>masked_assign (value, mask)</code>	Assign a scalar value to a series using a boolean mask <code>df[df < 0] = 0</code>	<code>to_frame ([name])</code>	Convert Series into a DataFrame
<code>max ([axis, skipna])</code>	Compute the max of the series	<code>to_gpu_array ([fillna])</code>	Get a dense numba device array for the data.
<code>mean ([axis, skipna])</code>	Compute the mean of the series	<code>to_hdf (path_or_buf, key, *args, **kwargs)</code>	Write the contained data to an HDF5 file using HDFStore.
<code>mean_var ([ddof])</code>	Compute mean and variance at the same time.	<code>to_json ([path_or_buf])</code>	Convert the cuDF object to a JSON string.
<code>min ([axis, skipna])</code>	Compute the min of the series	<code>to_string ([nrows])</code>	Convert to string
<code>nlargest ([n, keep])</code>	Returns a new Series of the <code>n</code> largest element.	<code>unique ([method, sort])</code>	Returns unique values of this Series.
<code>nsmallest ([n, keep])</code>	Returns a new Series of the <code>n</code> smallest element.	<code>value_counts ([method, sort])</code>	Returns unique values of this Series.
<code>nunique ([method, dropna])</code>	Returns the number of unique values of the Series: approximate version, and	<code>values_to_string ([nrows])</code>	Returns a list of string for each element.
		<code>var ([ddof, axis, skipna])</code>	Compute the variance of the series

CUDF

Today

CUDA

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

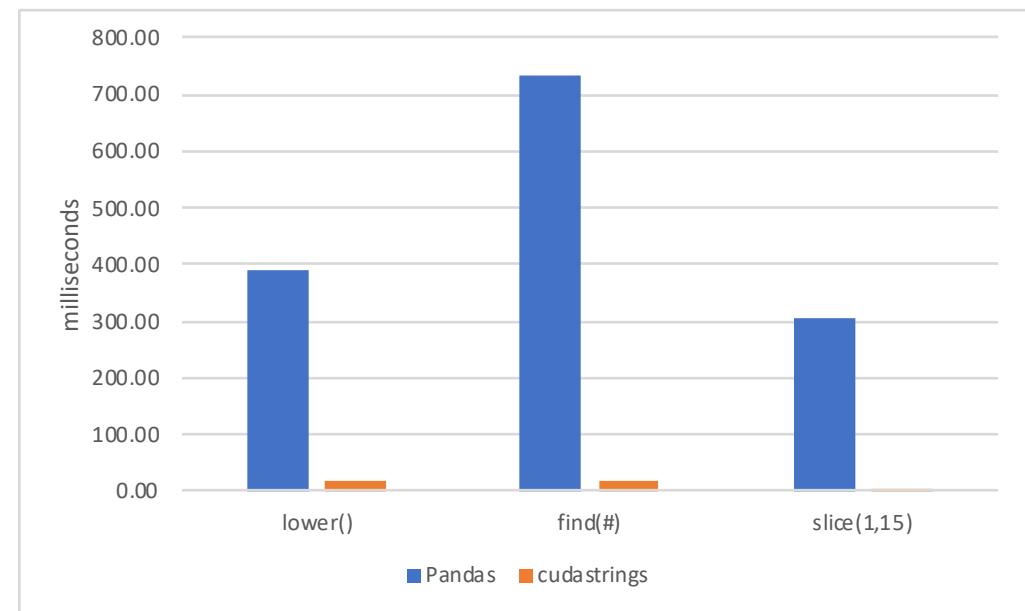
With Python Bindings

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

STRING SUPPORT

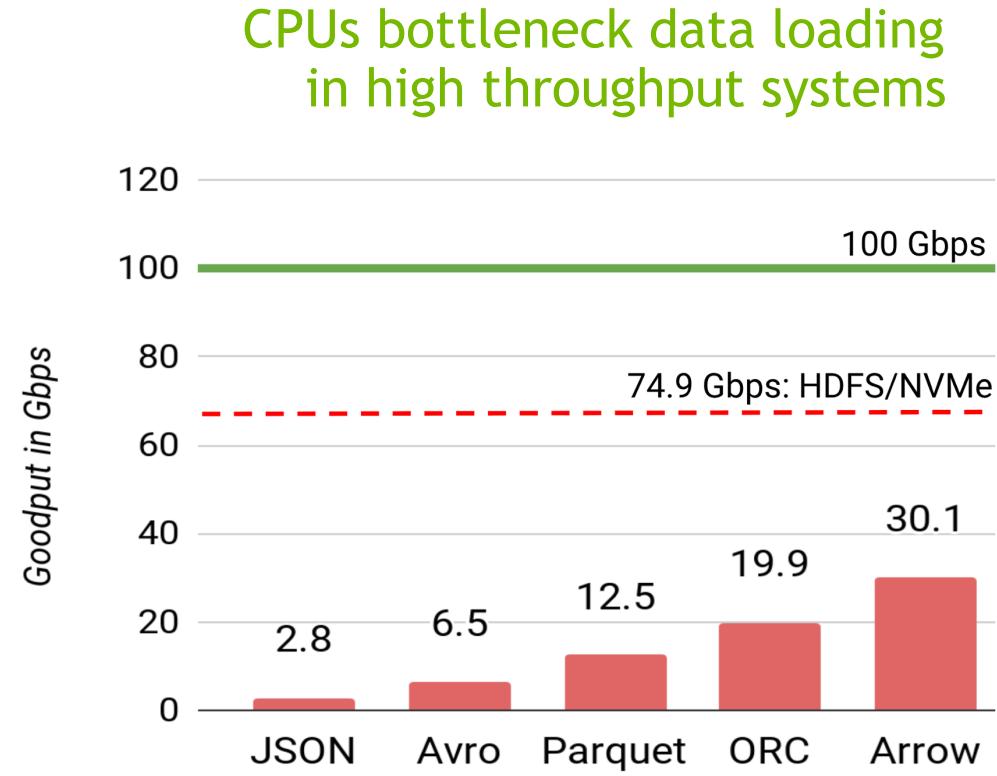
GPU-Accelerated string functions with a Pandas-like API

- API and functionality is following Pandas:
<https://pandas.pydata.org/pandas-docs/stable/api.html#string-handling>
- Handles ingesting and exporting typical Python objects (Pandas Series, Numpy arrays, PyArrow arrays, Python lists, etc.)
- Initial performance results:
 - `lower()`: ~22x speedup
 - `find()`: ~40x speedup
 - `slice()`: ~100x speedup



ACCELERATED DATA LOADING

- CSV Reader
 - CSV >10x speed improvement over pandas
 - Parquet
 - ORC Reader
 - JSON Reader
 - Feather Reader
 - HDF5 Reader
 - Arrow Reader
- Matches Pandas-API
- Decompression of the data ***will be*** GPU-accelerated as well!



Also performs writing

ACCELERATED DATA LOADING

- CSV Reader
 - CSV
 - Parquet
 - ORC Reader
 - JSON Reader
 - Feather Reader
 - HDF5 Reader
 - Arrow Reader

>10x speed improvement over pandas

```
import pandas as pd
import cudf
...
df = pd.read_csv('foo.csv', names=['index', 'A'], dtype=['date', 'float64'])
gdf = cudf.read_csv('foo.csv', names=['index', 'A'], dtype=['date', 'float64'])
...
----- Or just -----
gdf = cudf.read_csv('foo.csv')
```

CPUs bottleneck data loading
in high throughput systems



- Matches Pandas-API
- Decompression of the data **will be** GPU-accelerated as well!

Also performs writing

PYTHON CUDA ARRAY INTERFACE

Interoperability for Python GPU Array Libraries

- The CUDA array interface is a standard format that describes a GPU array to allow sharing GPU arrays between different libraries without needing to copy or convert data
- Native ingest and export of `__cuda_array_interface__` compatible objects via Numba device arrays in cuDF
- Numba, CuPy, and PyTorch are the first libraries to adopt the interface:
 - https://numba.pydata.org/numba-doc/dev/cuda/cuda_array_interface.html
 - <https://github.com/cupy/cupy/releases/tag/v5.0.0b4>
 - <https://github.com/pytorch/pytorch/pull/11984>



Numba



CuPy

DLPACK

Interoperability with Deep Learning Libraries

- DLPack is an open-source memory tensor structure designed to allow sharing tensors between deep learning frameworks
- Currently supported by PyTorch, MXNet, and Chainer / CuPy
- cuDF supports ingesting and exporting column-major DLPack tensors
 - If you're interested in row-major tensor support please let us know!

PYTORCH

mxnet

Chainer

“Details are confusing. It is only by selection, by elimination, by emphasis,
that we get to the real meaning of things.”

~ Georgia O'Keefe

CUML

cuML API

GPU-accelerated machine learning at every layer

Python

Scikit-learn-like interface for data scientists
utilizing cuDF & Numpy

Algorithms

CUDA C++ API to leverage accelerated machine
learning algorithms.

Primitives

Reusable building blocks for
composing machine learning algorithms.

PRIMITIVES

GPU-accelerated math optimized for feature matrices

Linear Algebra

- Element-wise operations
- Matrix multiply
- Norms
- Eigen Decomposition
- SVD/RSVD
- Transpose
- QR Decomposition

Probability & Statistics

Random Sampling

Distances

Scores & Metrics

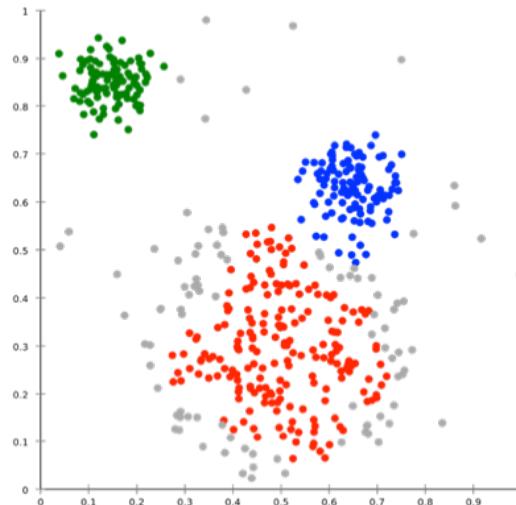
Objective Functions

Sparse Matrices

More to come!

ALGORITHMS

GPU-accelerated Scikit-Learn



Classification / Regression

Decision Trees / Random Forests
Linear Regression
Logistic Regression
K-Nearest Neighbors
Kalman Filtering
Bayesian Inference
Gaussian Mixture Models
Hidden Markov Models

Statistical Inference

K-Means
DBSCAN
Spectral Clustering
Principal Components
Singular Value Decomposition
UMAP
Spectral Embedding

Clustering

Decomposition & Dimensionality Reduction

ARIMA
Holt-Winters

Cross Validation

Timeseries Forecasting

Hyperparameter Tuning

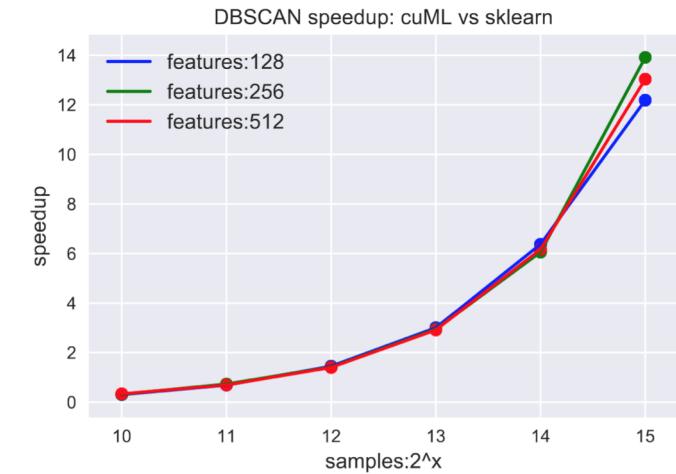
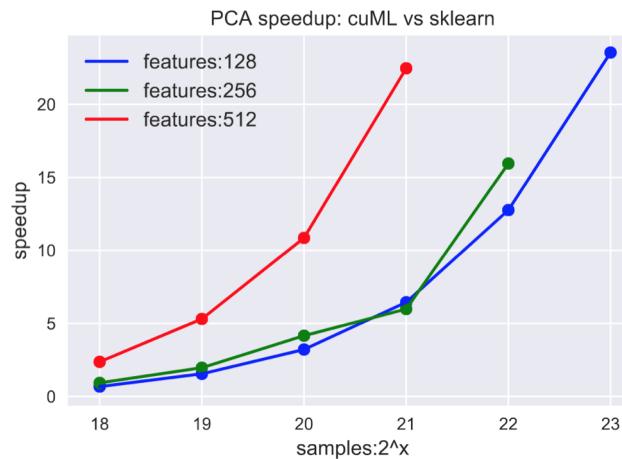
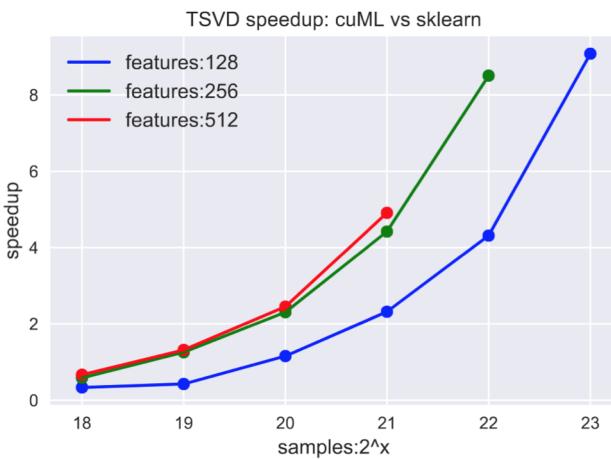
Recommendations

Implicit Matrix Factorization

More to come!

CUML

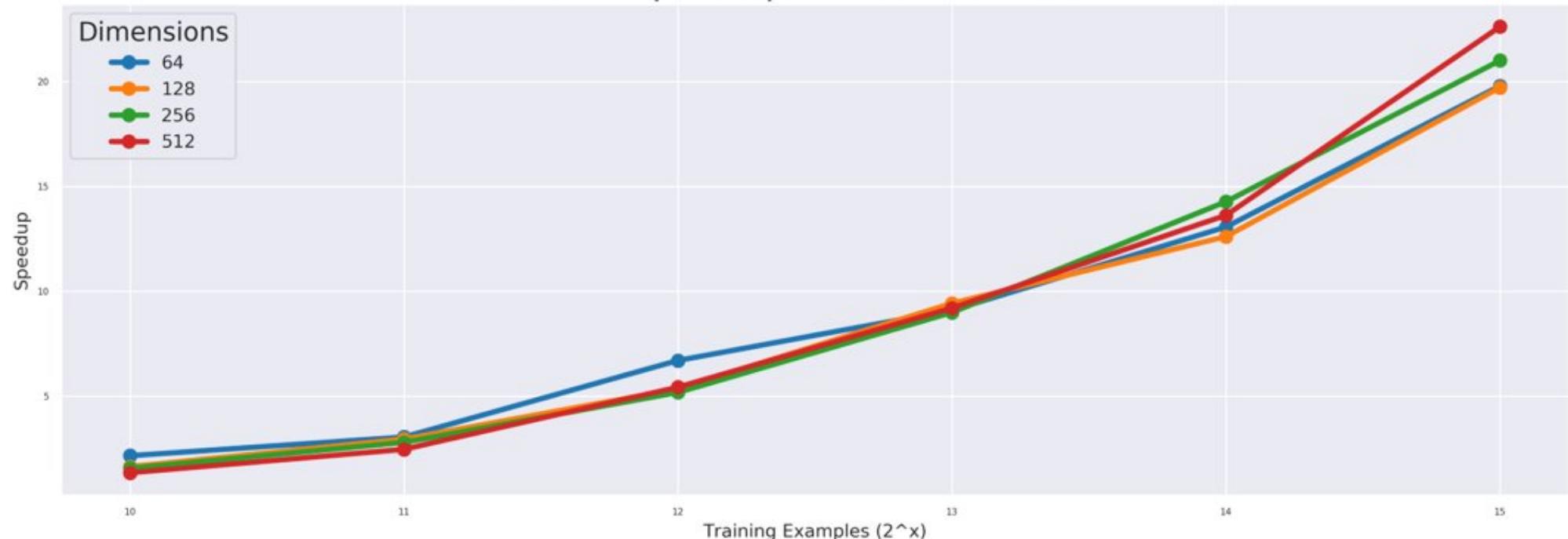
Speedup comparison to scikit-learn



UMAP

Speedup comparison to scikit-learn

UMAP Speedup: cuML vs SKLearn



SINGLE NODE MULTI-GPU

Linear Regression

- Reduction: 40mins -> 1min
- Size: 225gb
- System: DGX2

tSVD

- Reduction: 1.6hrs-> 1.5min
- Size: 220gb
- System: DGX2

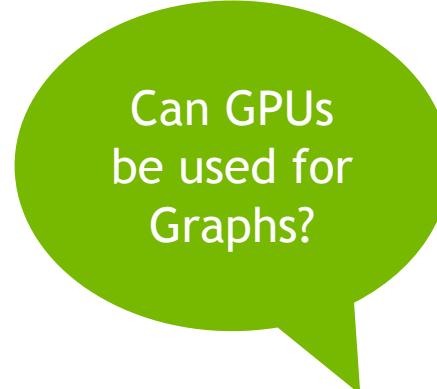
Nearest Neighbors

- Reduction: 4+hrs-> 30sec
- Size: 128gb
- System: DGX1



CUGRAPH

Because we know everything is a graph



Can GPUs
be used for
Graphs?



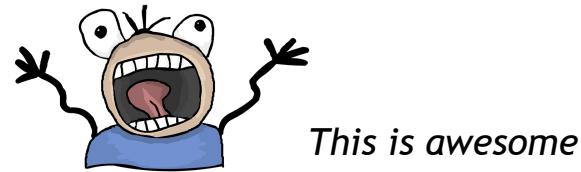
GRAPHS ON GPU

- Yes!
 - Larger amount of GPU memory supports bigger graphs
 - CUDA supports faster computation
 - BLAS and BSP/GAS model for algorithms
 - Higher memory bandwidth supports faster pointer chasing

GOALS AND BENEFITS OF CUGRAPH

Focus on Features an Easy-of-Use

- Seamless integration with cuDF and cuML
 - Inputs and outputs are cuDF Dataframes
- Features
 - targeting the most commonly used algorithms from NetworkX, and Spark, as the initial set
 - Overlap Coefficient, for example, is not in NetworkX
- NetworkX-like API*
- Breakthrough Performance
 - 10 to 10000x faster (*your mileage may vary*)



KDD will be a week before 0.9 is released

WHAT'S IN CUGRAPH

Current Single GPU Algorithms - as of release 0.8

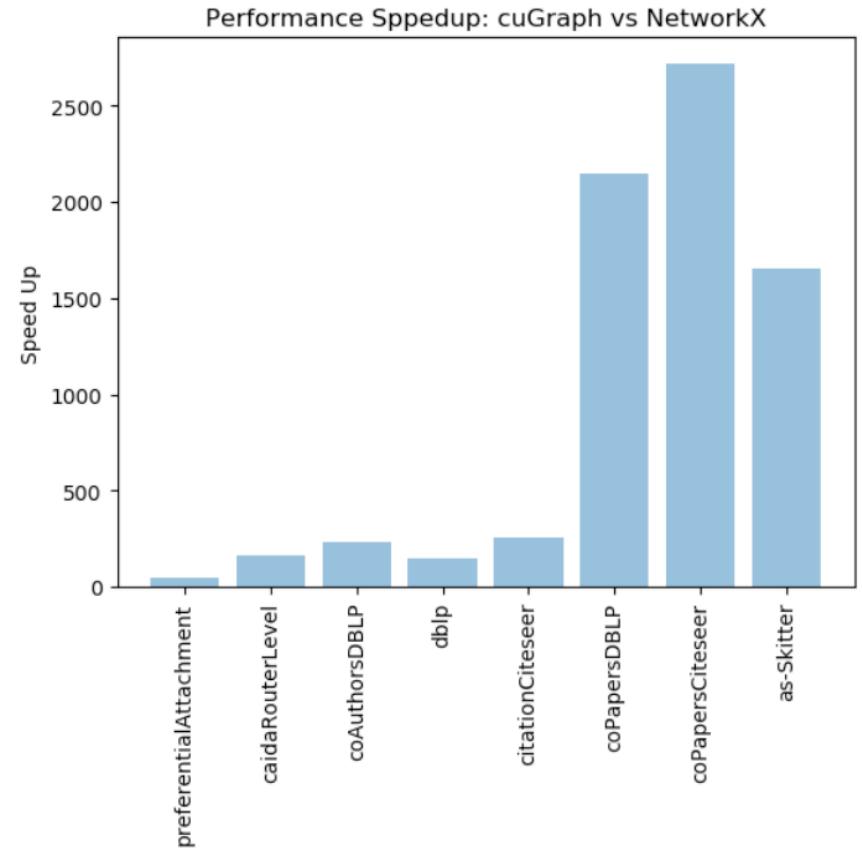
- Centrality
 - PageRank
 - MG being released in 0.9 (mid-Aug)
- Similarity
 - Jaccard
 - Weighted Jaccard
 - Overlap Coefficient
- Transveral
 - Single Source Shortest Path (SSSP)
 - Breadth First Search (BFS)
- Triangle Counting (TC)
- Subgraph Extraction
- Clustering
 - Spectral Clustering
 - Balanced-Cut
 - Modularity Maximization
 - Louvain
- Renumbering
- Graph functions
 - Size, Order, Degree
 - Get two hop neighbors

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



PageRank PERFORMANCE

Looking to the future - Multi-GPU processing

Single and Dual GPU on Commodity Workstation (PCIe)

Nodes	Edges	Single	Dual
1,048,576	16,777,216	0.019	0.020
2,097,152	33,554,432	0.047	0.035
4,194,304	67,108,864	0.114	0.066
8,388,608	134,217,728	0.302	0.162
16,777,216	268,435,456	0.771	0.353
33,554,432	536,870,912	1.747	0.821
67,108,864	1,073,741,824		1.880

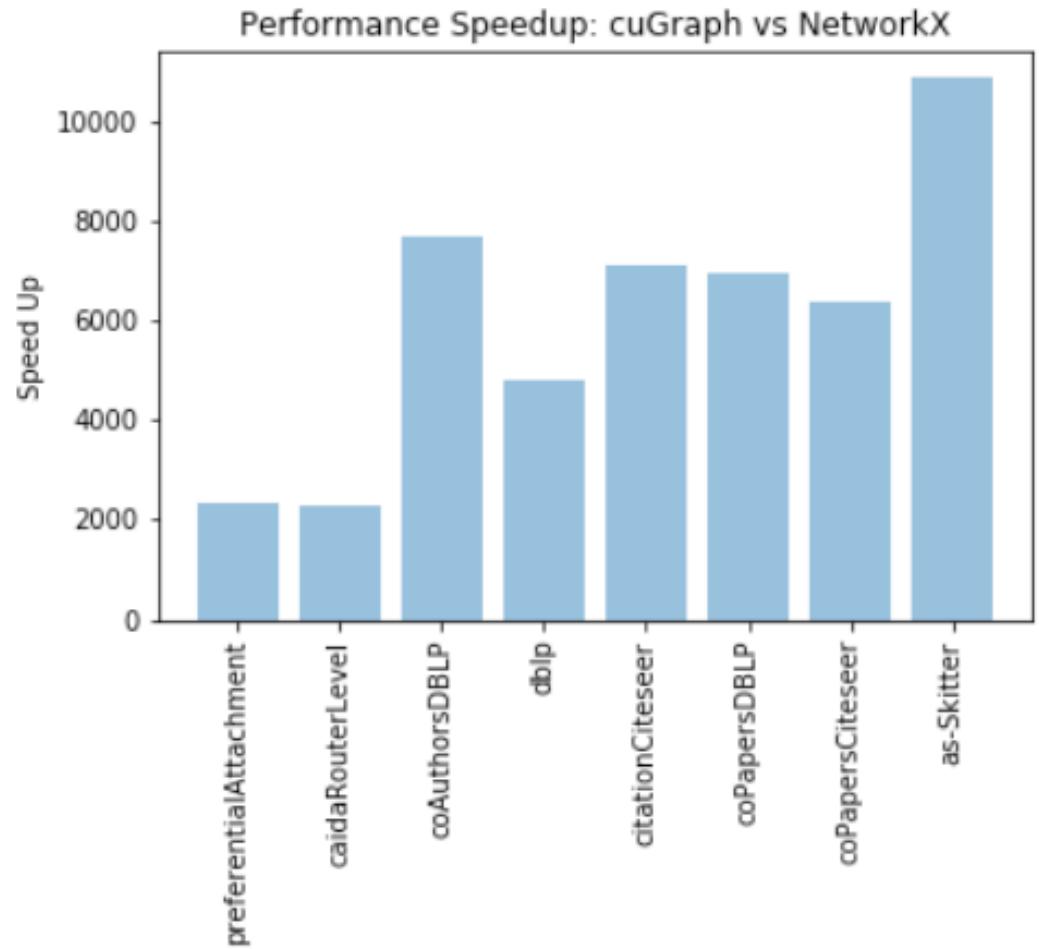
DGX-2 (SMX)

GPUS	RMAT	Nodes	Edges	Runtime
1	25	33,554,432	536,870,912	1.4052
2	26	67,108,864	1,073,741,824	1.3891
4	27	134,217,728	2,122,307,214	1.3891
8	28	268,435,456	4,294,967,296	1.4103
16	29	536,870,912	8,589,934,592	1.4689

LOUVAIN SINGLE RUN

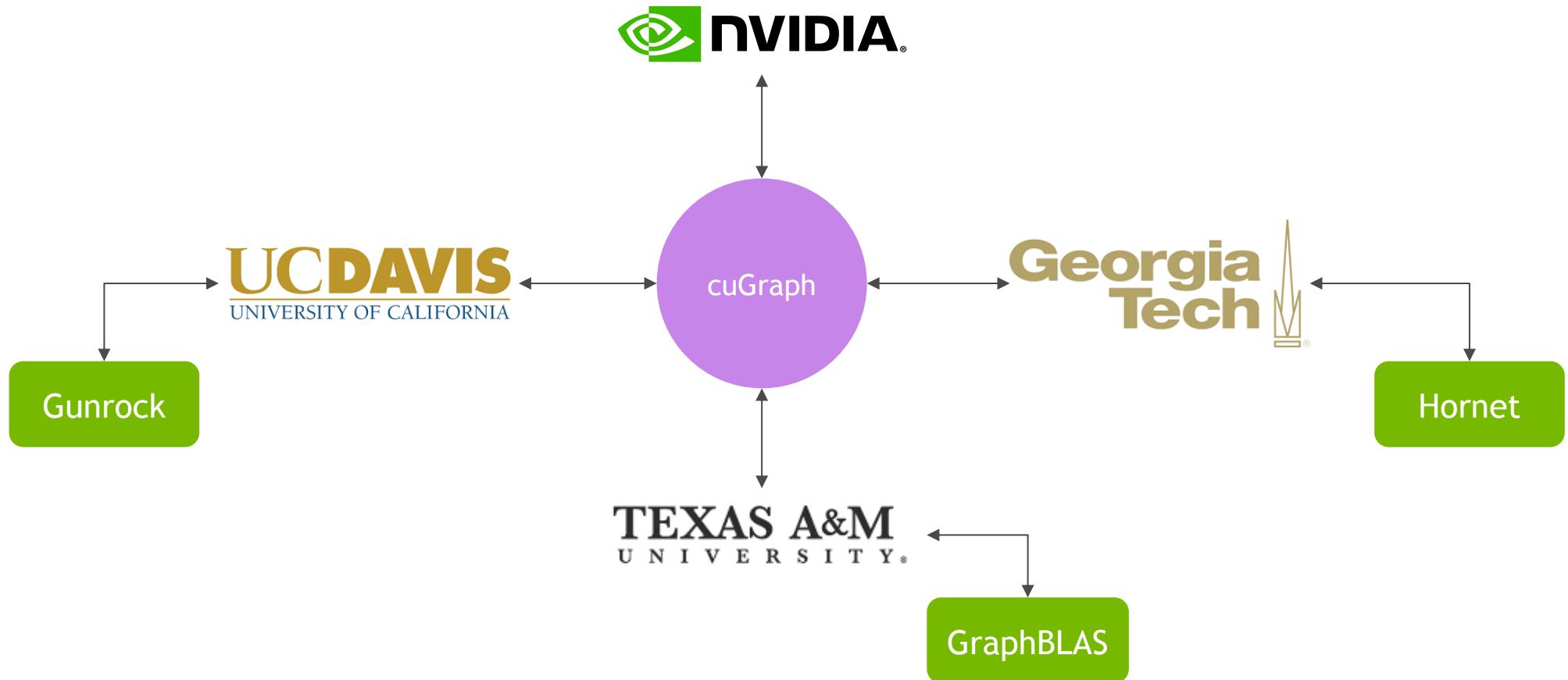
```
G = cugraph.Graph()  
G.add_edge_list(gdf["src_0"], gdf["dst_0"], gdf["data"])  
df, mod = cugraph.nvLouvain(G)
```

louvain: cudf.DataFrame with two names columns:
louvain["vertex"]: The vertex id.
louvain["partition"]: The assigned partition.



BRINGING IN LEADING RESEARCHERS

Leveraging the great work of others





Scalable Python and now RAPIDS

DASK

DASK

What is Dask and why does RAPIDS use it for scaling out?

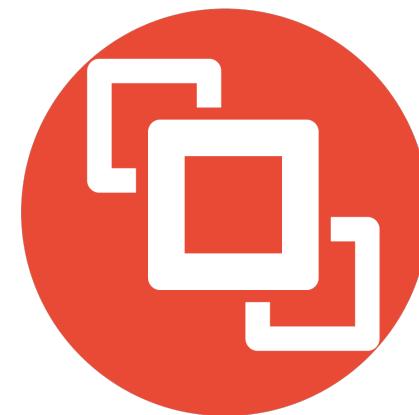
- Dask is a distributed computation scheduler built to scale Python workloads from laptops to supercomputer clusters.
- Extremely modular with scheduling, compute, data transfer, and out-of-core handling all being disjointed allowing us to plug in our own implementations.
- Can easily run multiple Dask workers per node to allow for an easier development model of one worker per GPU regardless of single node or multi node environment.



DASK

Scale up and out with cuDF

- Use cuDF primitives underneath in map-reduce style operations with the same high level API
- Instead of using typical Dask data movement of pickling objects and sending via TCP sockets, take advantage of hardware advancements using a communications framework called OpenUCX:
 - For intranode data movement, utilize NVLink and PCIe peer-to-peer communications
 - For internode data movement, utilize GPU RDMA over Infiniband and RoCE



<http://www.openucx.org/>



[https://github.com/rapidsai/
dask-cudf](https://github.com/rapidsai/dask-cudf)



LET'S START BUILD ANALYTICS USING RAPIDS