

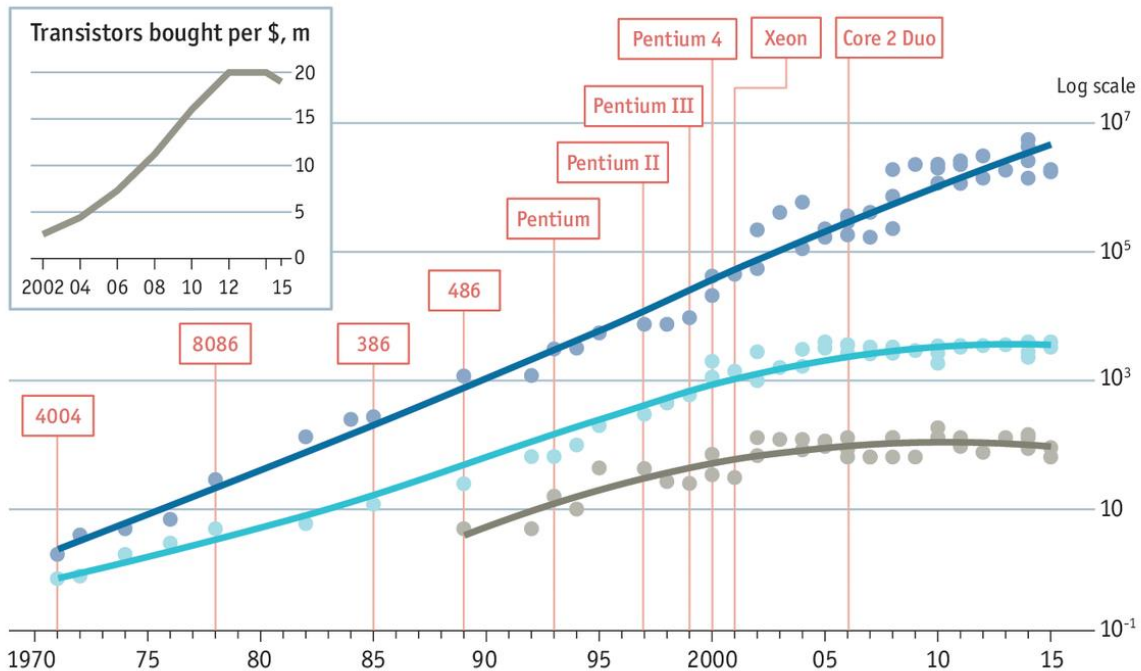
RAPIDS

GPU Accelerated Data Analytics in Python

Mads R. B. Kristensen, NVIDIA

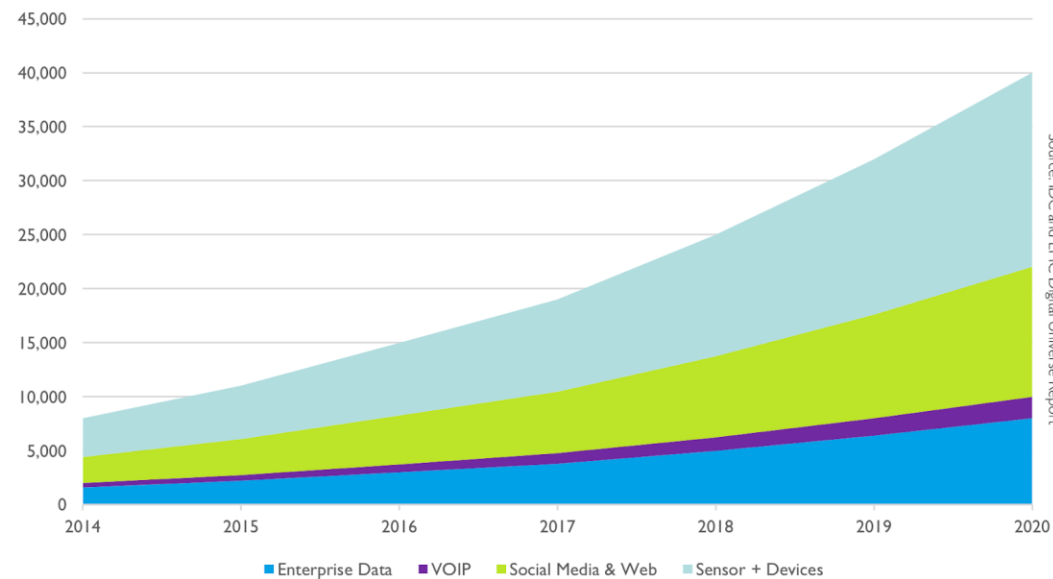
Stuttering

● Transistors per chip, '000 ● Clock speed (max), MHz ● Thermal design power*, w □ Chip introduction dates, selected



Sources: Intel; Bob Colwell; Linley Group; International Business Strategies; *The Economist* *Maximum safe power consumption
Economist.com

Data Growth and Source in Exabytes



Source: IDC and EMC Digital Universe Report

Scale up and out with RAPIDS and Dask

Scale Up / Accelerate

RAPIDS and Others

Accelerated on single GPU

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



Dask + RAPIDS

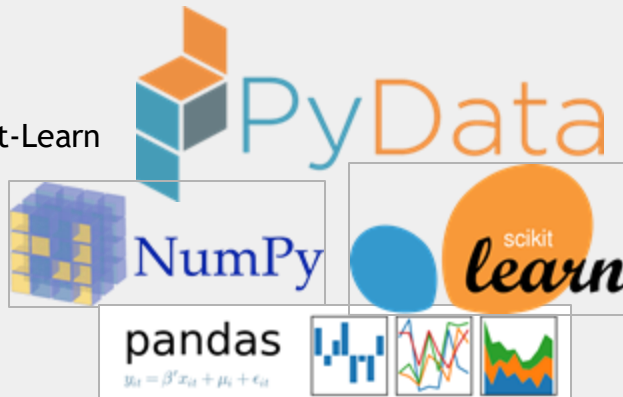
Multi-GPU
On single Node (DGX)
Or across a cluster



PyData

NumPy, Pandas, Scikit-Learn
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

Scale up and out with RAPIDS and Dask

Scale Up / Accelerate

RAPIDS and Others

Accelerated on single GPU

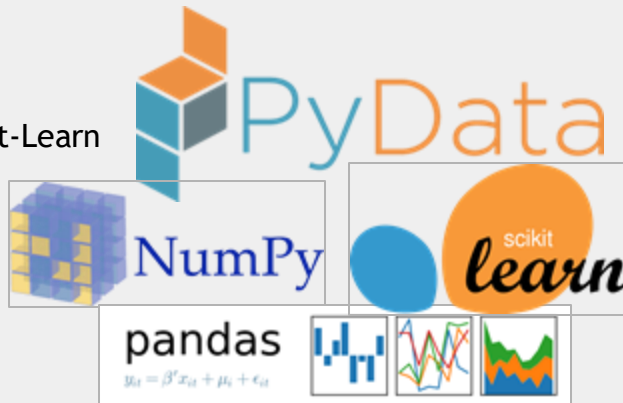
NumPy -> CuPy/PyTorch/..
Pandas -> cuDF
Scikit-Learn -> cuML
Numba -> Numba



PyData

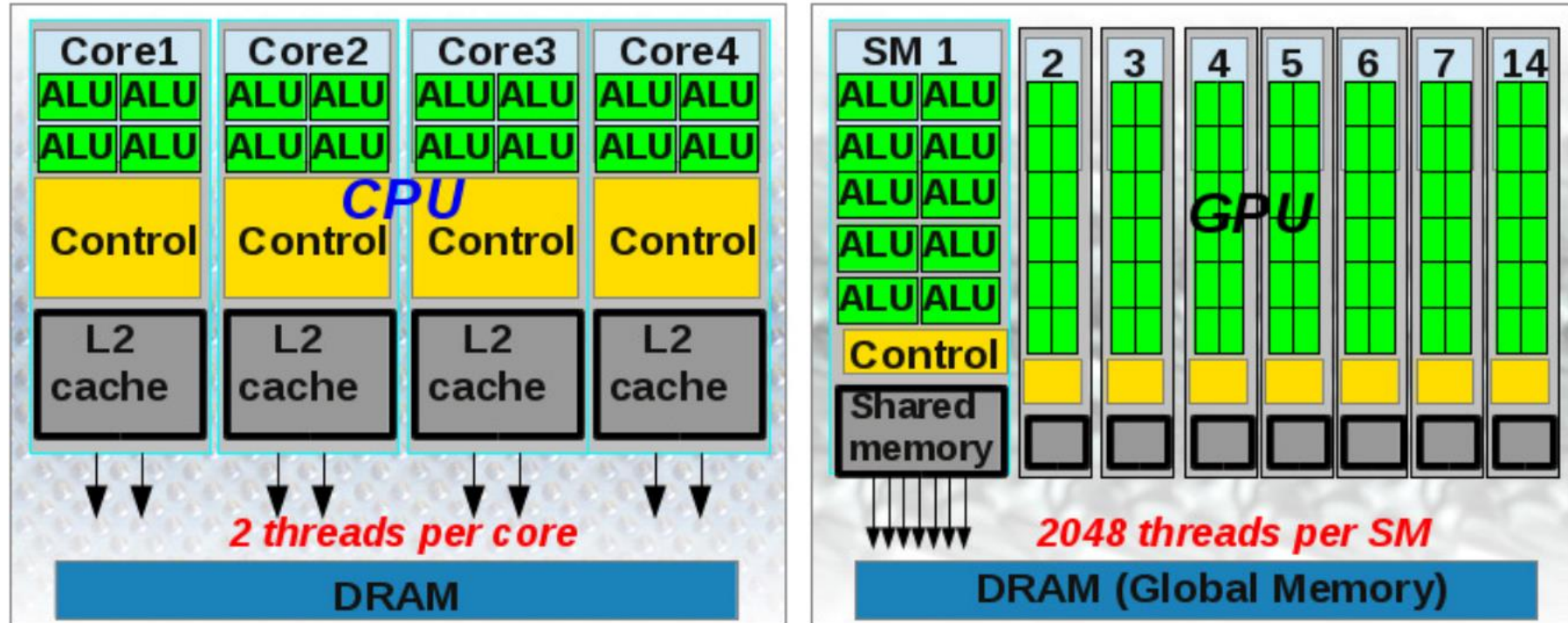
NumPy, Pandas, Scikit-Learn
and many more

Single CPU core
In-memory data



Scale out / Parallelize

CPU vs GPU



Data Processing Evolution

Faster data access, less data movement

Hadoop Processing, Reading from disk



Spark In-Memory Processing

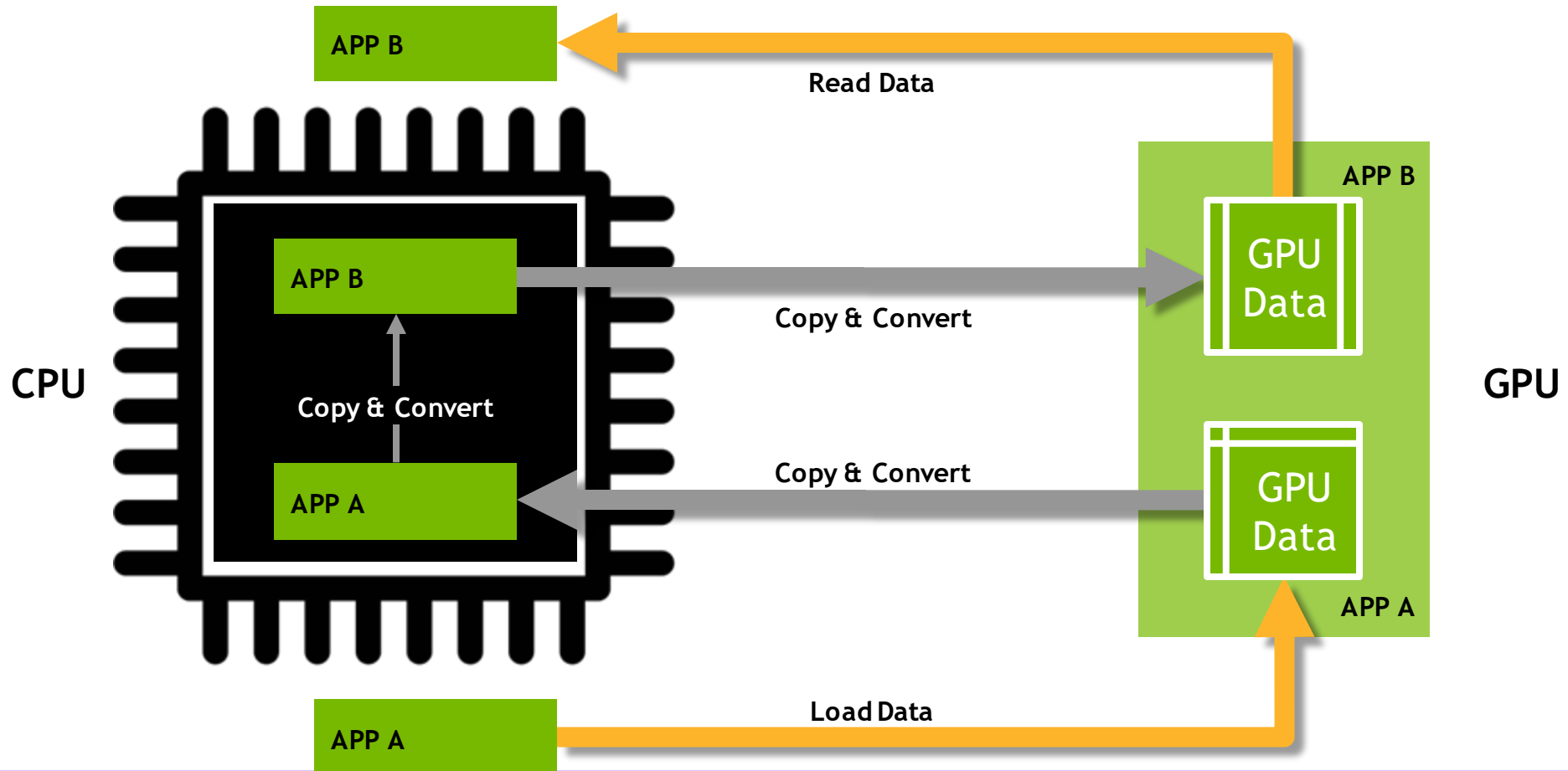


Traditional GPU Processing



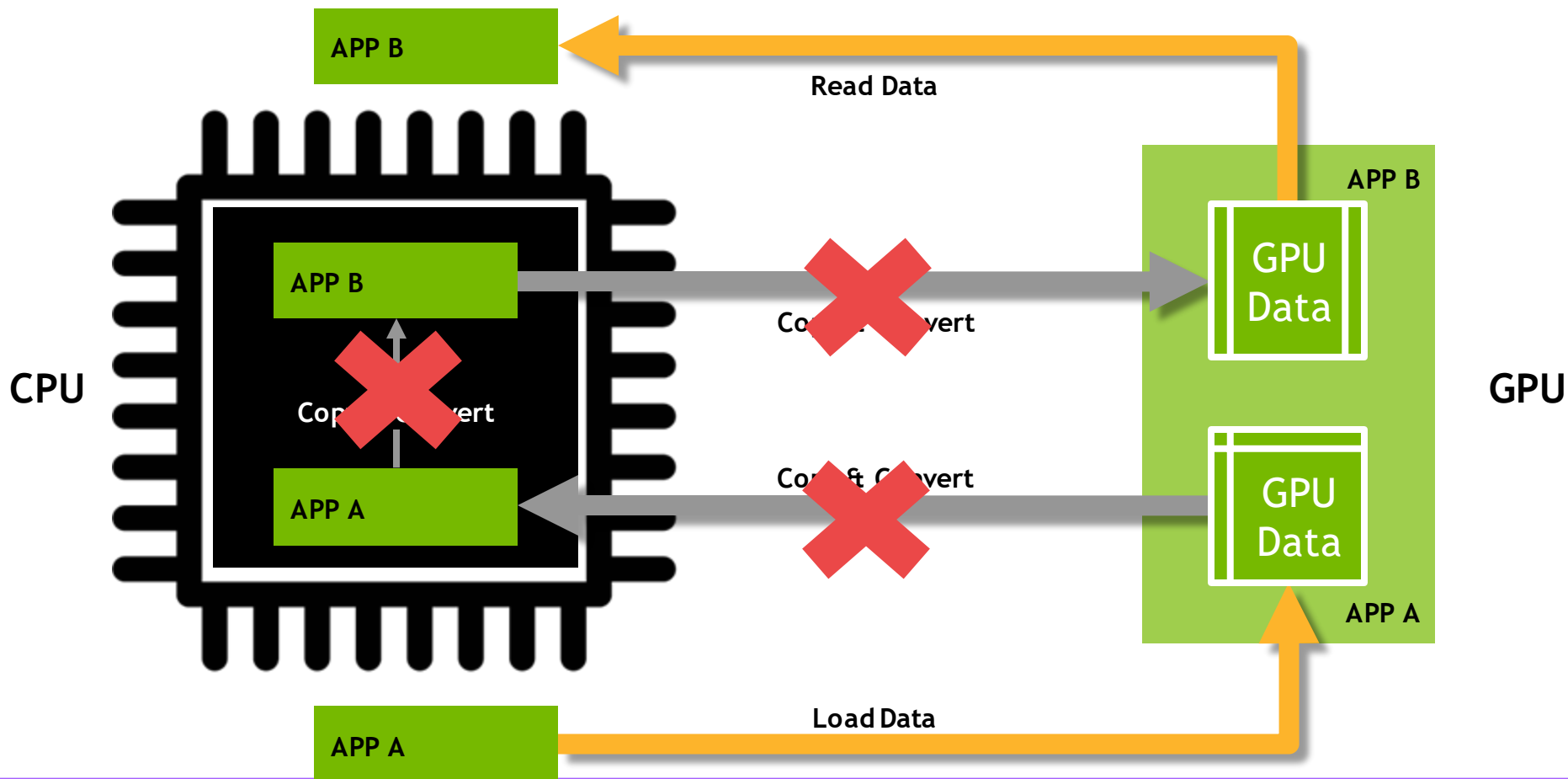
Data Movement and Transformation

What if we could keep data on the GPU?



Data Movement and Transformation

What if we could keep data on the GPU?



Data Processing Evolution

Faster data access, less data movement

Hadoop Processing, Reading from disk



Spark In-Memory Processing



Traditional GPU Processing



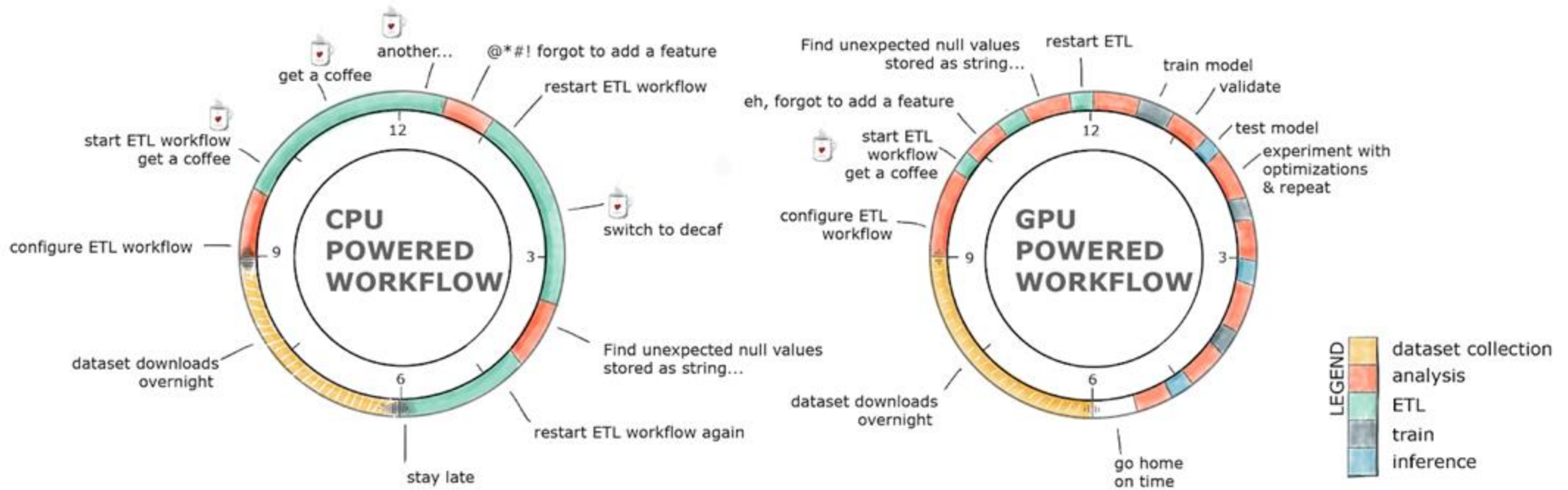
RAPIDS



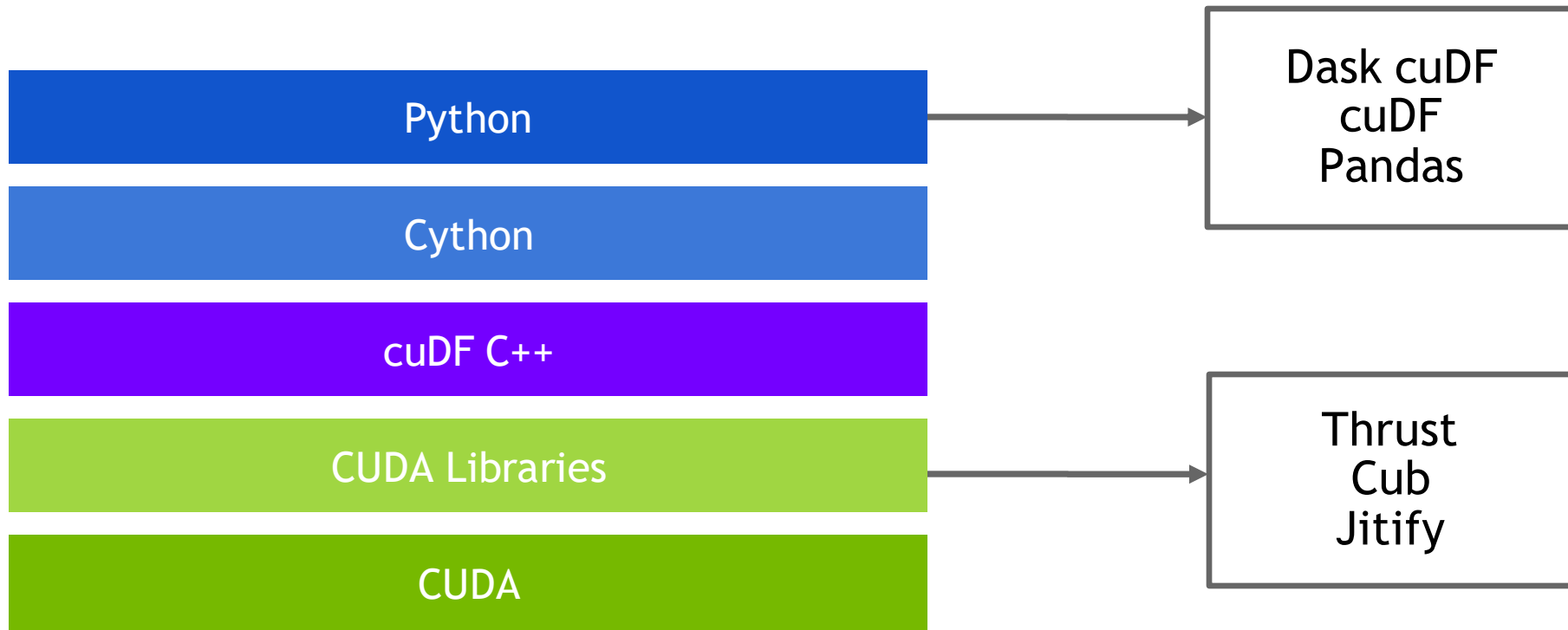
CuDF - Pandas on GPU

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<table>
gather(table_view const& input,
      column_view const& gather_map, ...)
{
    // return a new table containing
    // rows from input indexed by
    // gather_map
}
```



ETL: the Backbone of Data Science

cuDF is...

Python Library

```
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]:
```

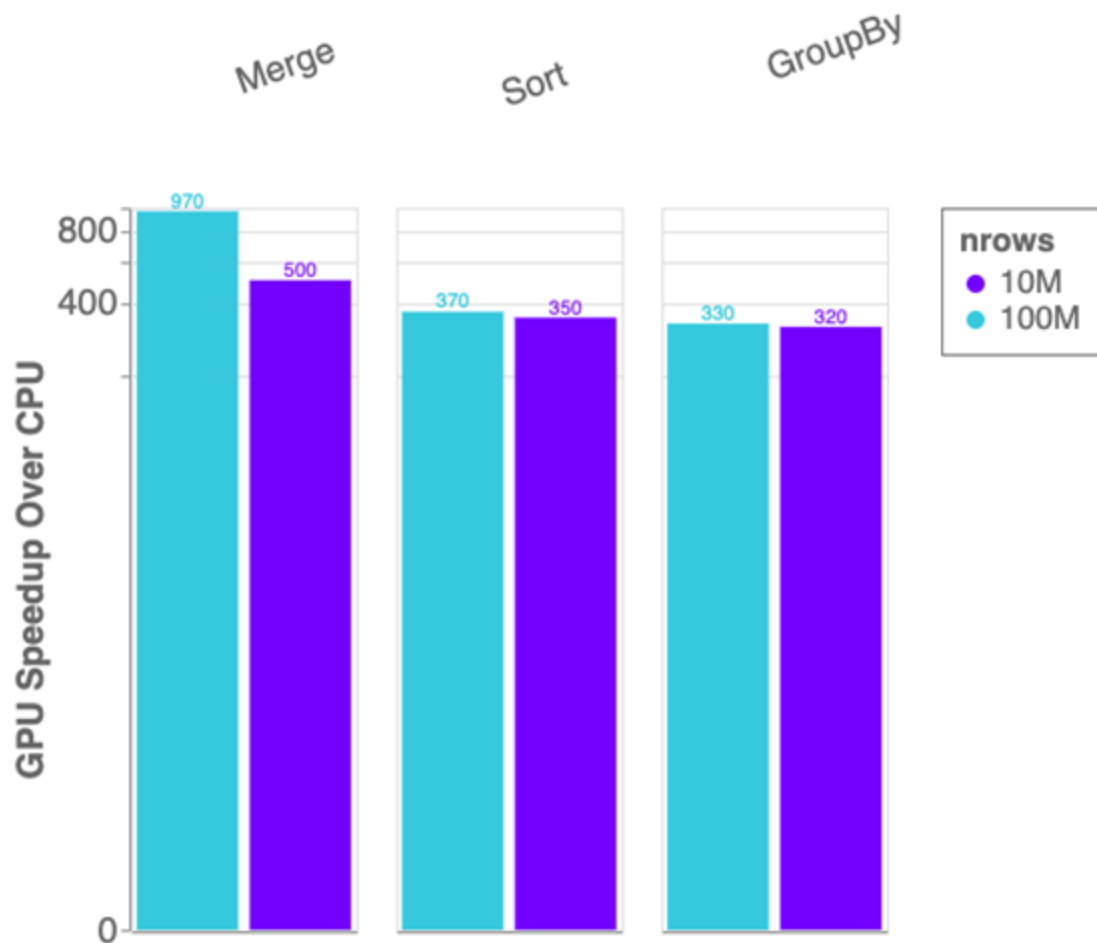
	User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Cat
0	1000001	P00069042	F	0-17	10	A	2	0	3
1	1000001	P00248942	F	0-17	10	A	2	0	1
2	1000001	P00087842	F	0-17	10	A	2	0	12
3	1000001	P00085442	F	0-17	10	A	2	0	12
4	1000002	P00285442	M	55+	16	C	4+	0	8

```
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()
```

- 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

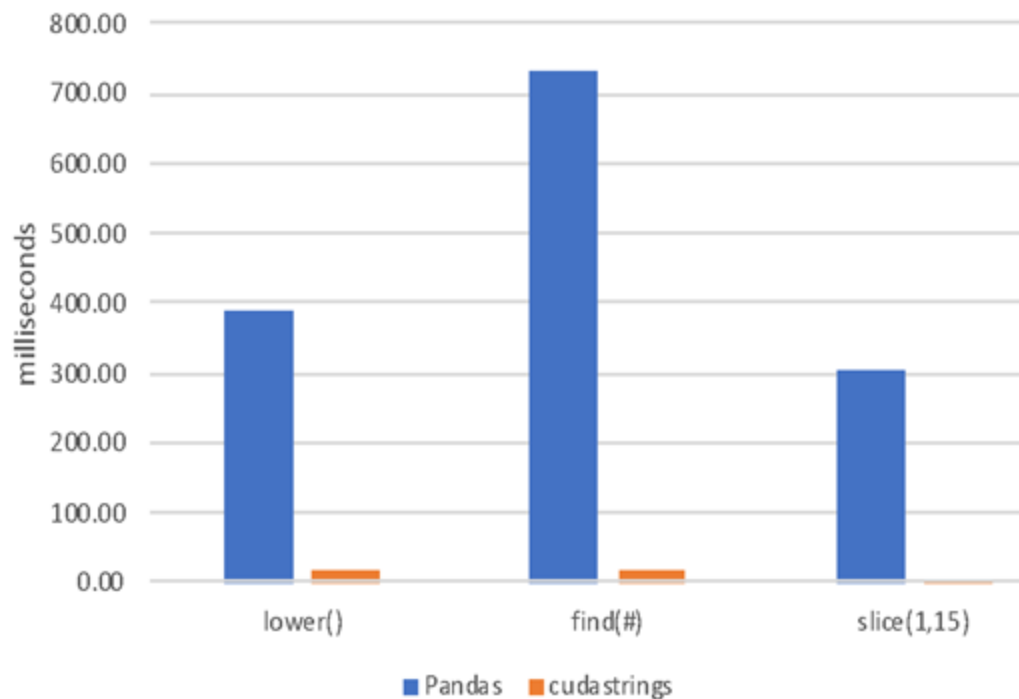
GroupBy: count, sum, min, max calculated
for each value column

ETL: the Backbone of Data Science

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++

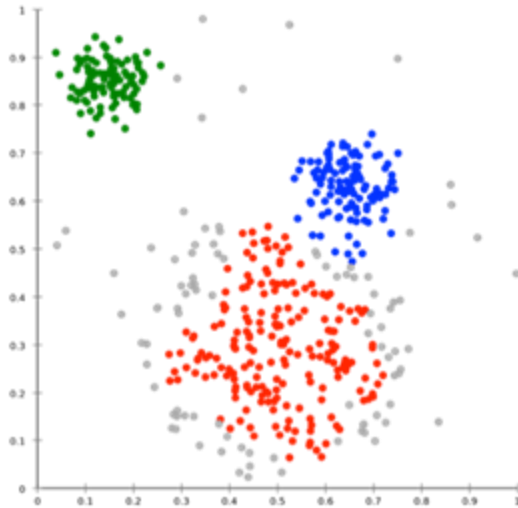
- Further performance optimization
- JIT-compiled String UDFs



CuML - Scikit-Learn on GPU

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

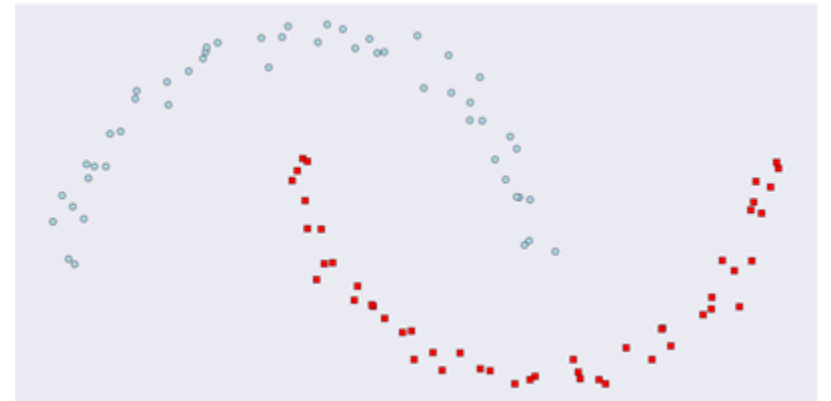
X, y = make_moons(n_samples=int(1e2),
                  noise=0.05, random_state=0)

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

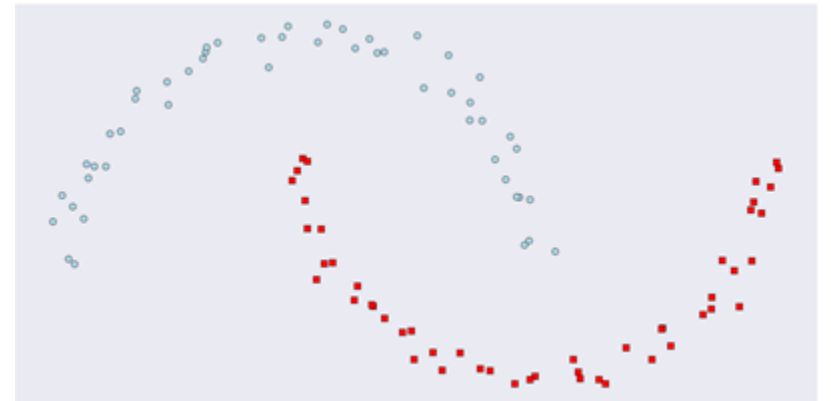
X, y = make_moons(n_samples=int(1e2),
                  noise=0.05, random_state=0)

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)

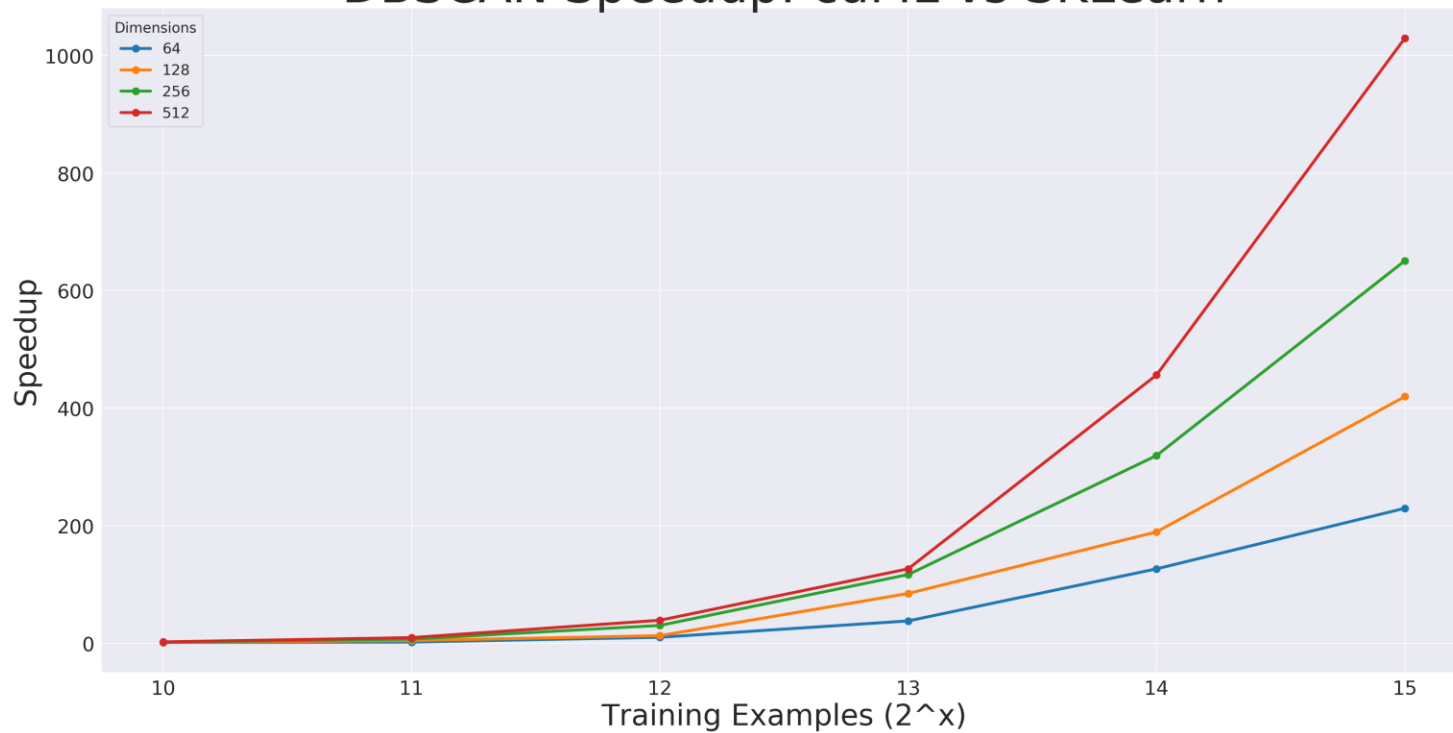
y_hat = dbscan.predict(X)
```



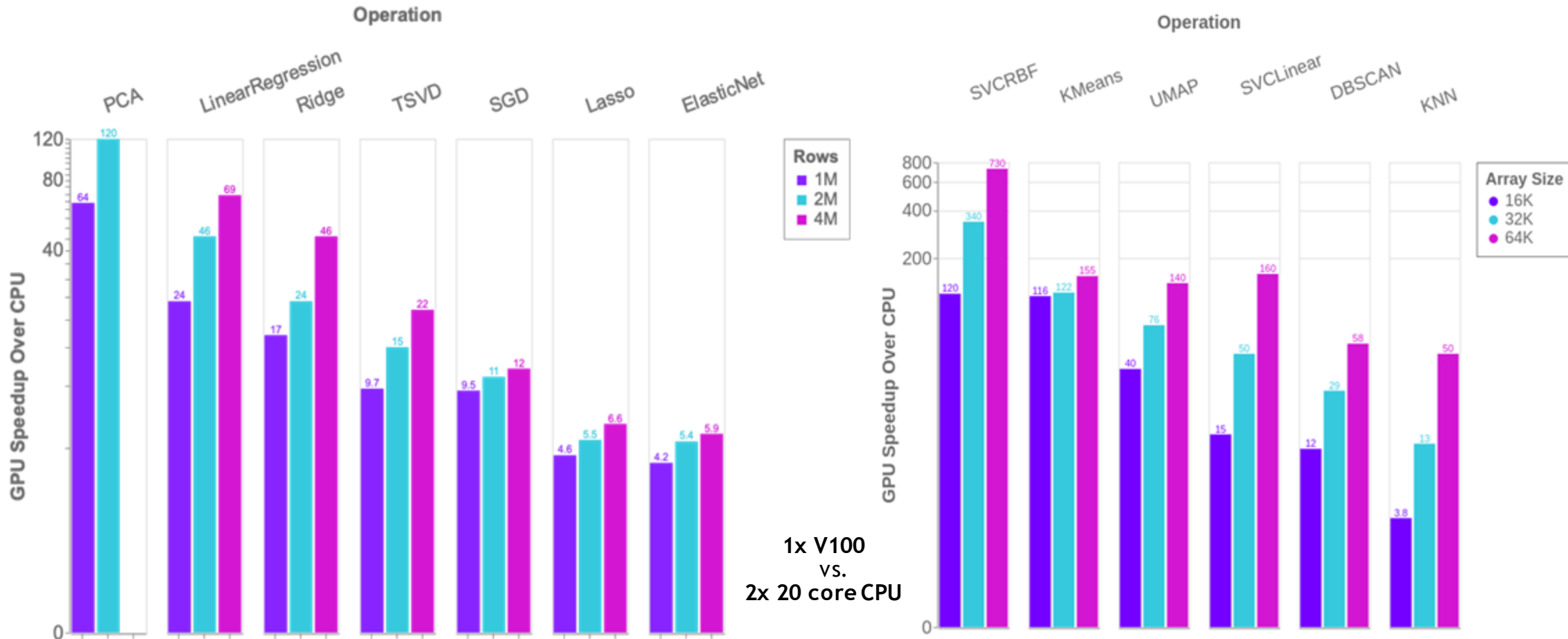
CLUSTERING

Benchmark

DBSCAN Speedup: cuML vs SKLearn



Benchmarks: single-GPU cuML vs scikit-learn



Scale up and out with RAPIDS and Dask

Scale Up / Accelerate

RAPIDS and Others

Accelerated on single GPU

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



Dask + RAPIDS

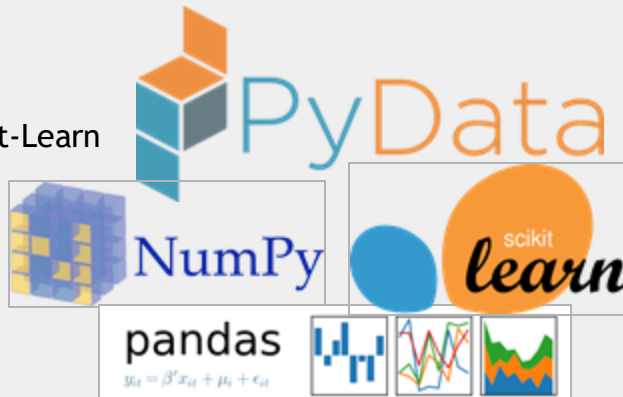
Multi-GPU
On single Node (DGX)
Or across a cluster



PyData

NumPy, Pandas, Scikit-Learn
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

Dask

Distributing Python Libraries


Scale up and out with RAPIDS and Dask

Scale Up / Accelerate

PyData

NumPy, Pandas, Scikit-Learn
and many more

Single CPU core
In-memory data



pydata

NumPy

scikit learn


pandas

$y_i = \beta^T x_i + \mu_i + \epsilon_i$

Dask

Multi-core and Distributed PyData

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



DASK

Scale out / Parallelize

Dask Parallelizes

Natively



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

Parallel NumPy

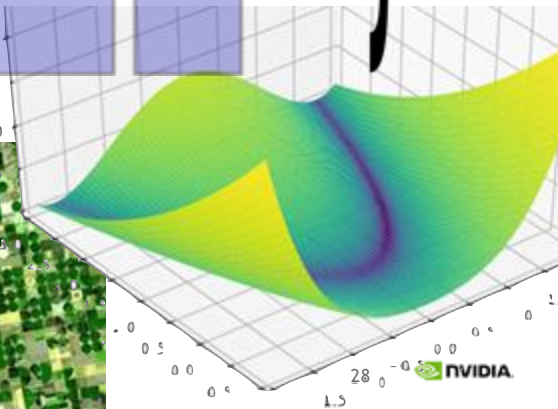
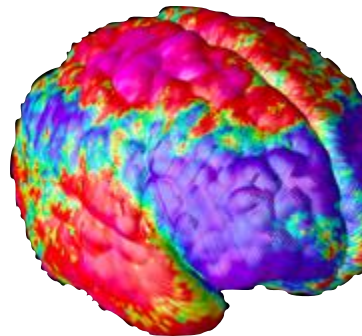
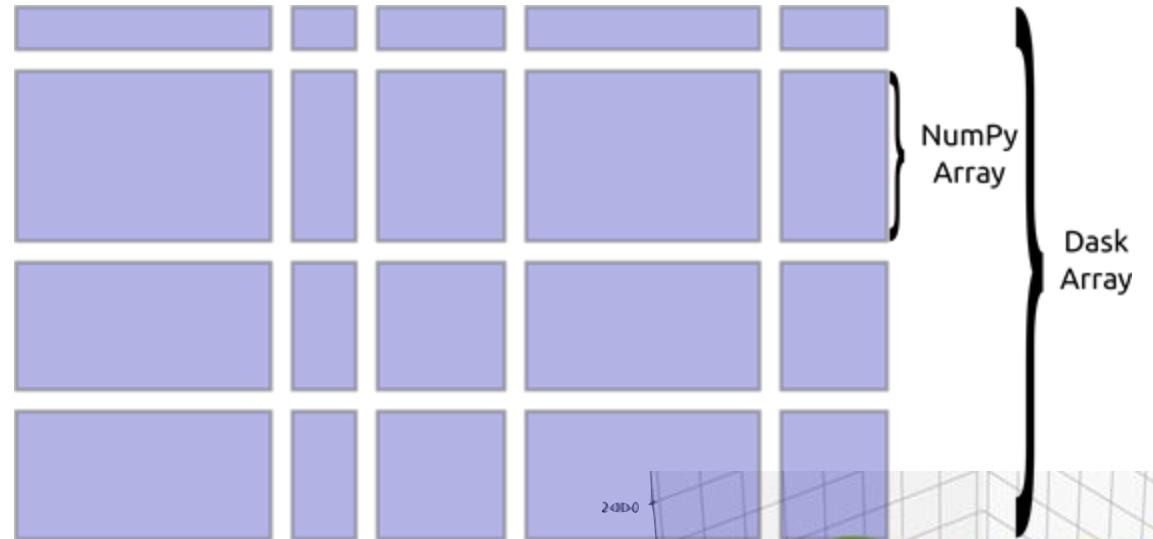
For imaging, simulation analysis, machine learning

- Same API as NumPy

```
import dask.array as da
x = da.from_hdf5(...)
x + x.T - x.mean(axis=0)
```

- One Dask Array is built from many NumPy arrays

Either lazily fetched from disk
Or distributed throughout a cluster



Parallel Pandas

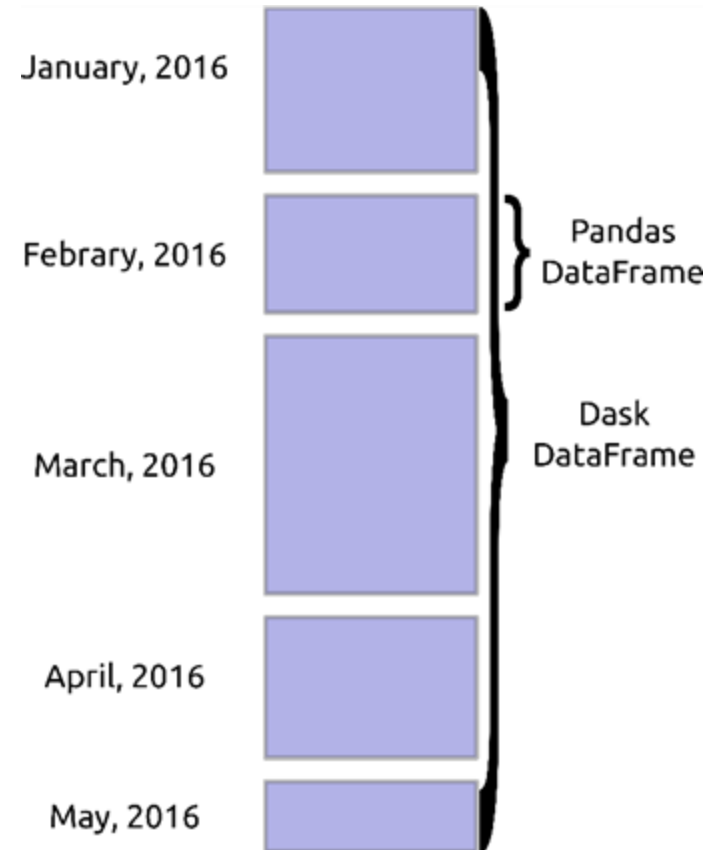
For ETL, time series, data munging

- Same API as Pandas

```
import dask.dataframe as dd
df = dd.read_csv(...)
df.groupby('name').balance.max()
```

- One Dask DataFrame is built from many Pandas DataFrames

Either lazily fetched from disk
Or distributed throughout a cluster

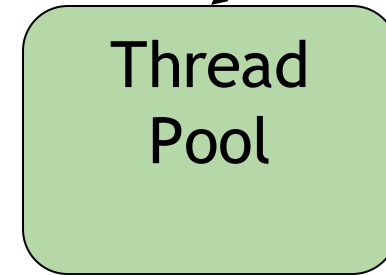


Parallel Scikit-Learn

For Hyper-Parameter Optimization, Random Forests, ...

- Same API

```
estimator = RandomForest()  
estimator.fit(data, labels)
```



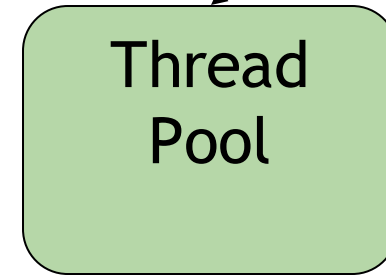
Parallel Scikit-Learn

For Hyper-Parameter Optimization, Random Forests, ...

- Same API

```
from scikit_learn.externals import joblib
with joblib.parallel_backend('dask'):
    estimator = RandomForest()
    estimator.fit(data, labels)
```

- Same exact code, just wrap in a "with" block
- Replaces default threaded execution with Dask
Allowing scaling onto clusters
- Available in most Scikit-Learn algorithms where joblib is used



Parallel Python

For custom systems, ML algorithms, workflow engines

- Parallelize existing codebases

```
results = {}  
  
for x in X:  
    for y in Y:  
        if x < y:  
            result = f(x, y)  
        else:  
            result = g(x, y)  
        results.append(result)
```

Parallel Python

For custom systems, ML algorithms, workflow engines

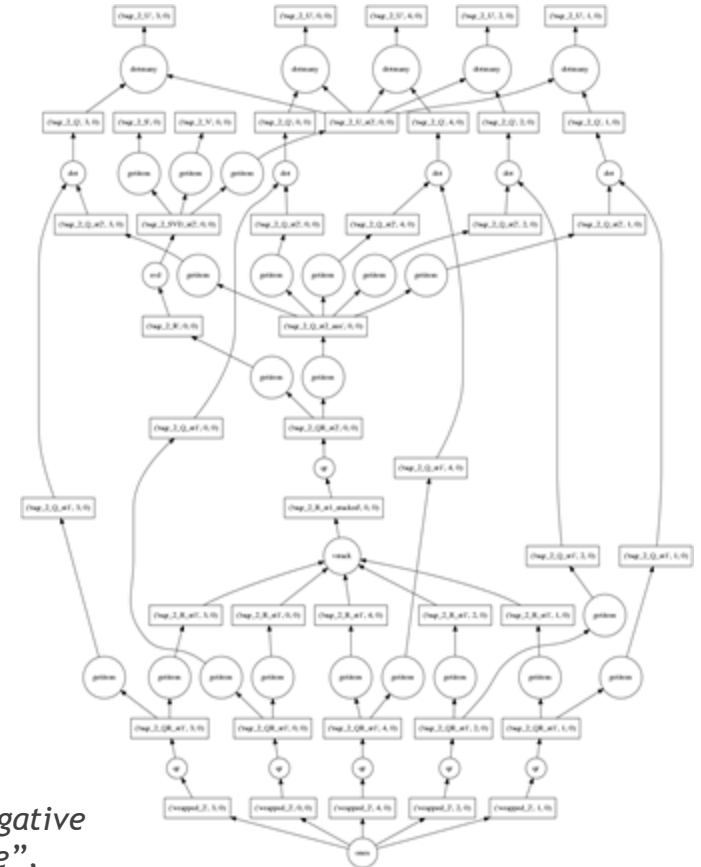
- Parallelize existing codebases

```
f = dask.delayed(f)
g = dask.delayed(g)

results = {}

for x in X:
    for y in Y:
        if x < y:
            result = f(x, y)
        else:
            result = g(x, y)
        results.append(result)

result = dask.compute(results)
```



M Tepper, G Sapiro “Compressed nonnegative matrix factorization is fast and accurate”, IEEE Transactions on Signal Processing, 2016

Dask Connects Python users to Hardware



User

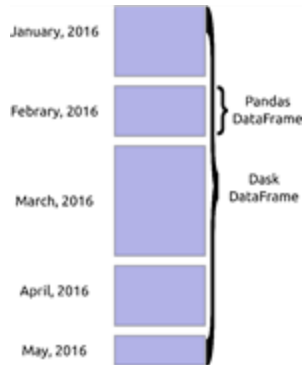


Execute on distributed hardware

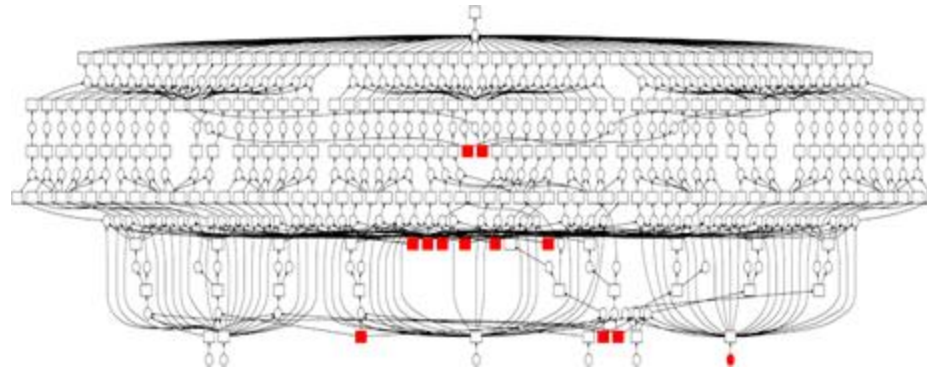
Dask Connects Python users to Hardware



User



Writes high level code
(NumPy/Pandas/Scikit-Learn)



Turns into a task graph



Execute on distributed
hardware

Scale up and out with RAPIDS and Dask

Scale Up / Accelerate

RAPIDS and Others

Accelerated on single GPU

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



Dask + RAPIDS

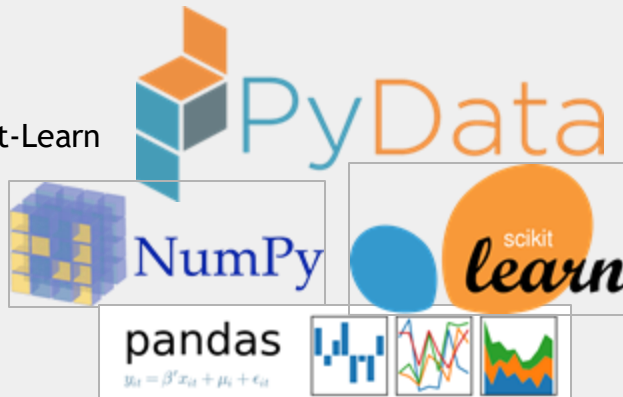
Multi-GPU
On single Node (DGX)
Or across a cluster



PyData

NumPy, Pandas, Scikit-Learn
and many more

Single CPU core
In-memory data



Dask

Multi-core and Distributed PyData

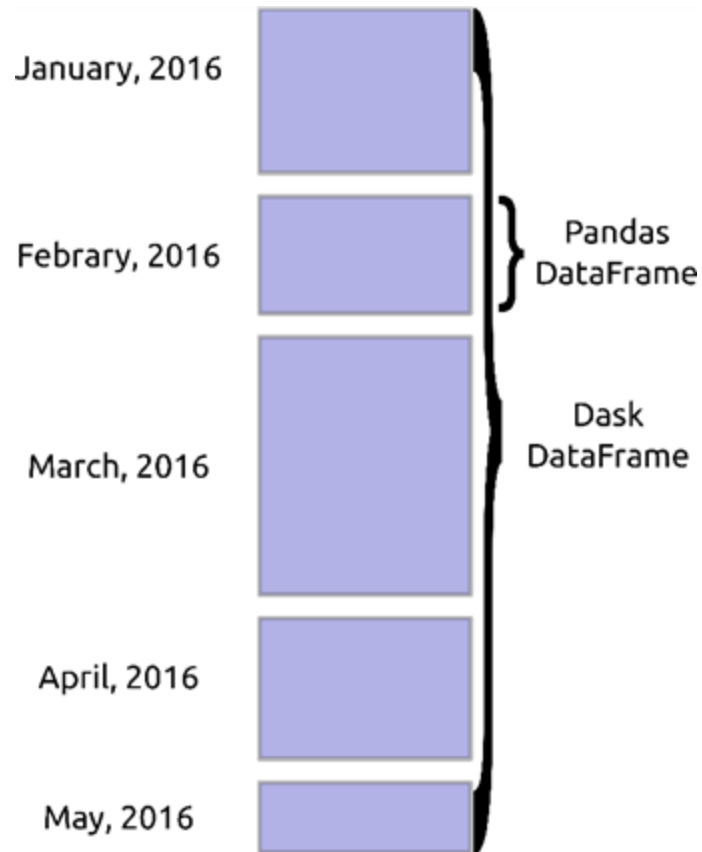
NumPy -> Dask Array
Pandas -> Dask DataFrame
Scikit-Learn -> Dask-ML
... -> Dask Futures



Scale out / Parallelize

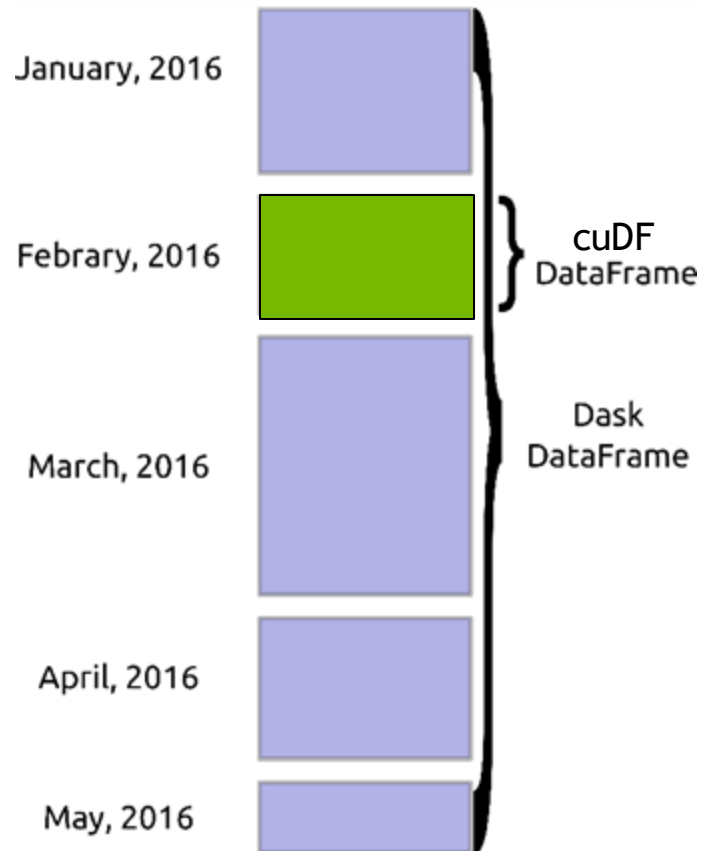
Combine Dask with cuDF

Many GPU DataFrames form a distributed DataFrame

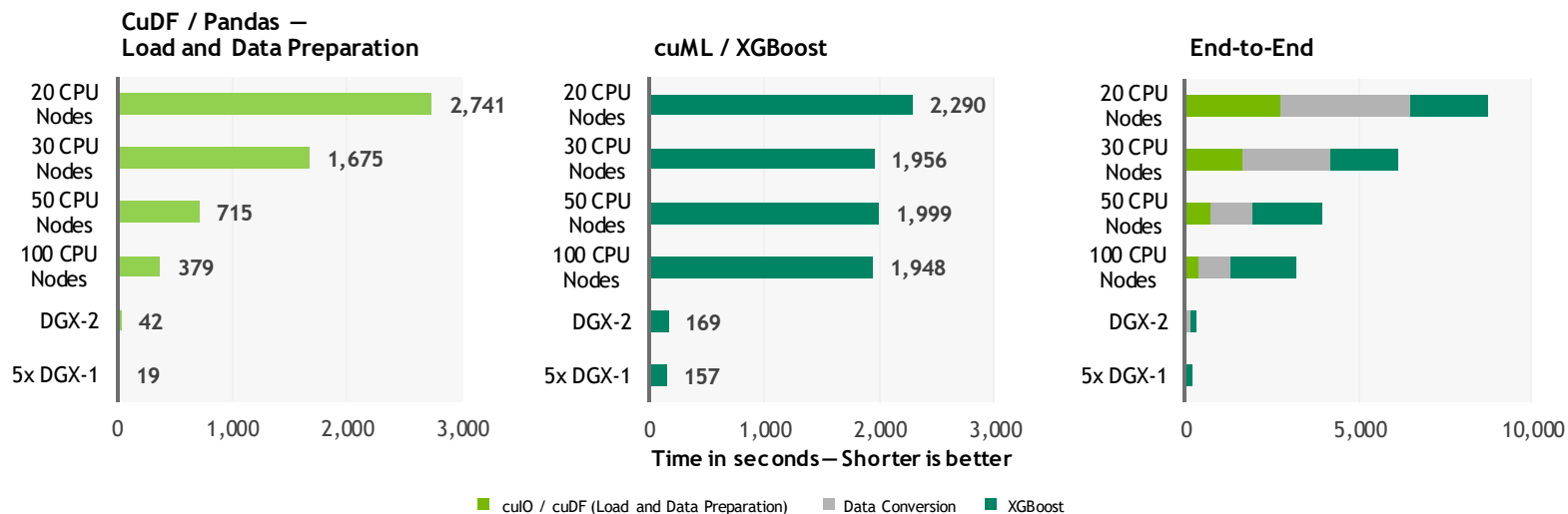


Combine Dask with cuDF

Many GPU DataFrames form a distributed DataFrame



END-TO-END BENCHMARKS



Benchmark

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

CPU Cluster Configuration

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

DGX Cluster Configuration

5x DGX-1 on InfiniBand network

Getting Started

Explore: RAPIDS Github

<https://github.com/rapidsai>



RAPIDS

Open GPU Data Science

<http://rapids.ai>

Repositories 92

Packages

People 135

Teams 138

Projects 6

Pinned repositories

cuda

cuDF - GPU DataFrame Library

Cuda 2.5k 336

cuml

cuML - RAPIDS Machine Learning Library

C++ 1.1k 169

cugraph

cuGraph - RAPIDS Graph Analytics Library

Cuda 331 64

cusignal

cuSignal

Jupyter Notebook 229 23

cuspatial

CUDA-accelerated GIS and spatiotemporal algorithms

Python 90 21

notebooks

RAPIDS Sample Notebooks

Jupyter Notebook 319 144

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](#).

	Preferred		Advanced				
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			

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

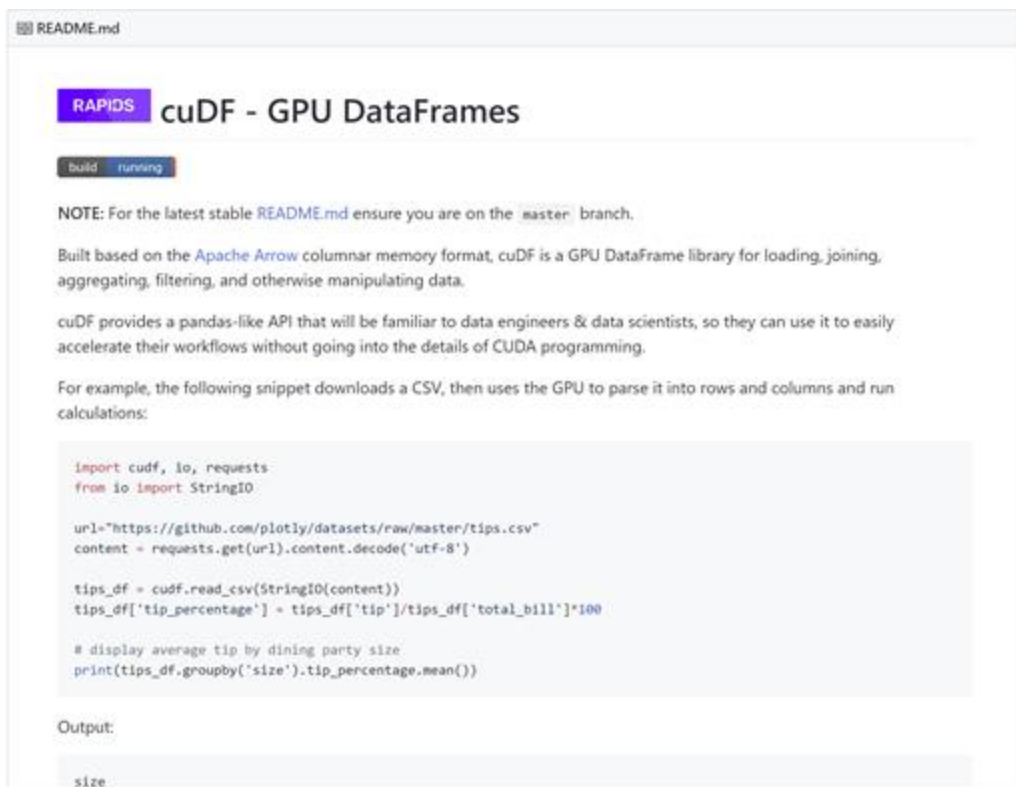
```
conda install -c rapidsai -c nvidia -c conda-forge \
-c defaults rapids=0.13 python=3.6
```

COPY COMMAND

DETAILS BELOW

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
```

<https://github.com/rapidsai>



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

Karthikayan 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://medium.com/rapids-ai>