



DEEP
LEARNING
INSTITUTE



DLI Accelerated Data Science Teaching Kit

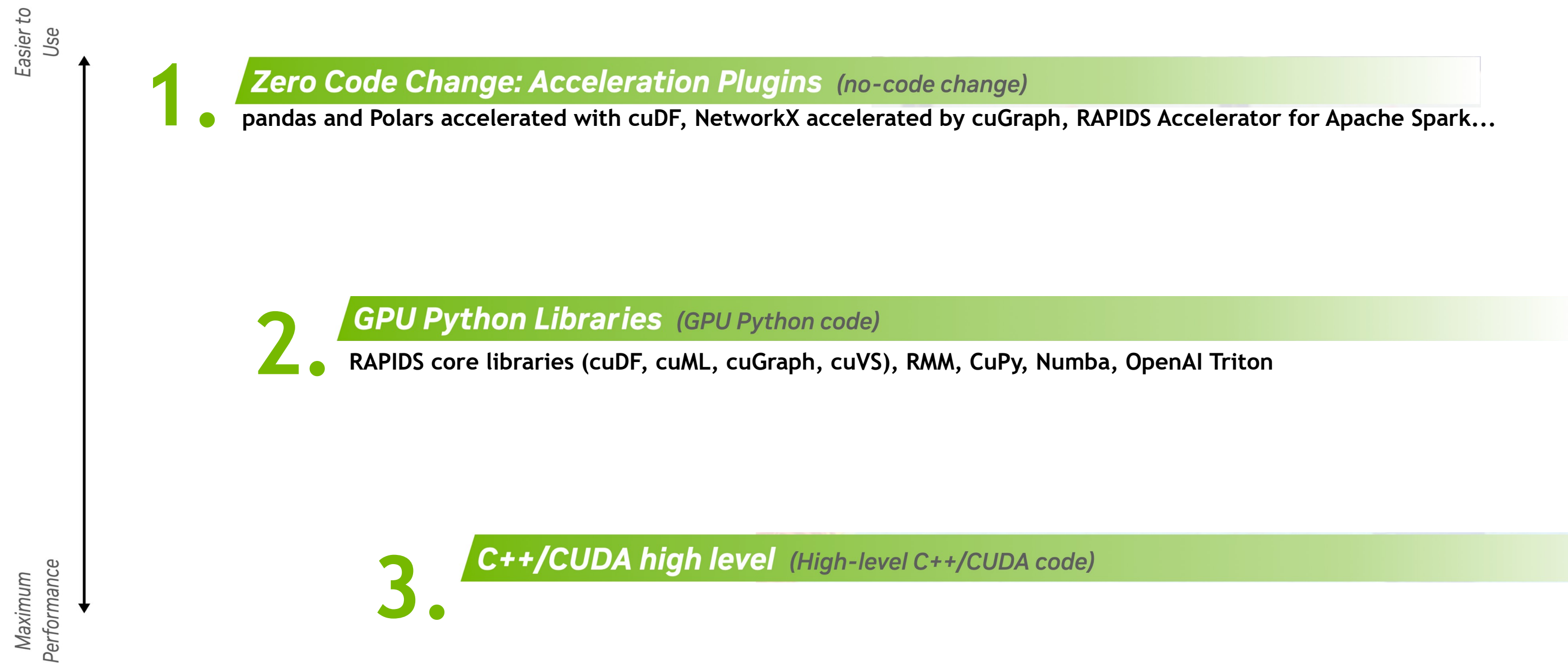
Lecture 21.2 - Refactoring Workloads



The Accelerated Data Science Teaching Kit is licensed by NVIDIA, Georgia Institute of Technology, and Prairie View A&M University under the [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/).

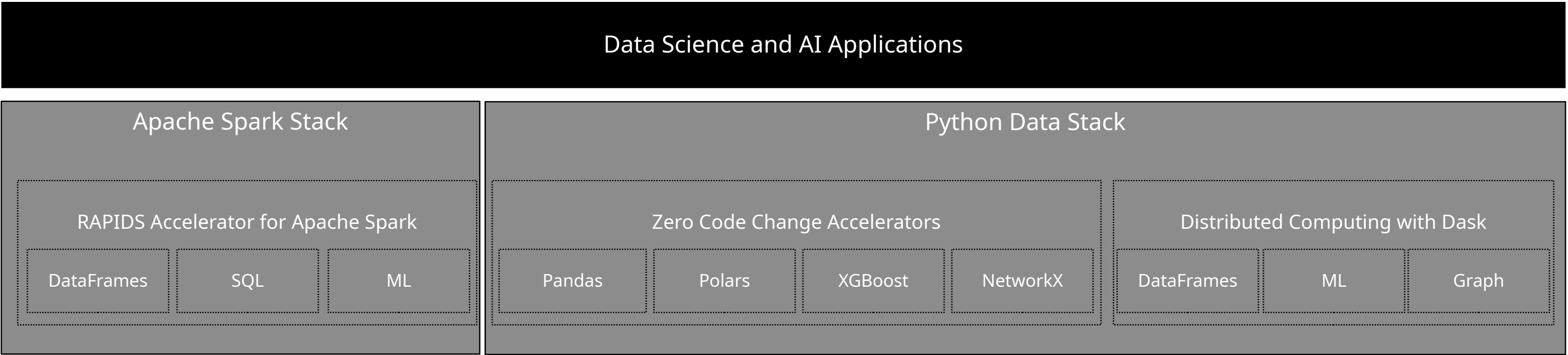
As we covered before...

There are 3 levels to access this acceleration

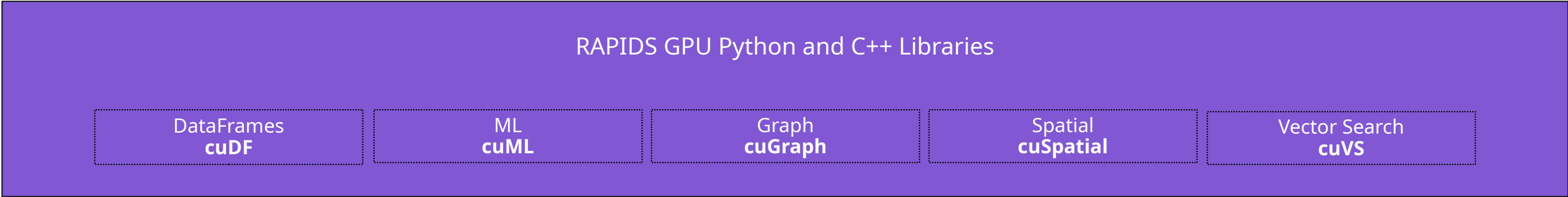


The GPU Accelerated Data Science Stack

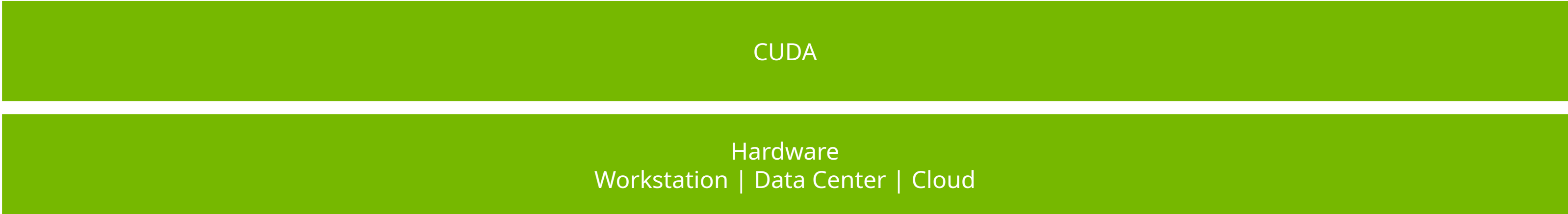
1.



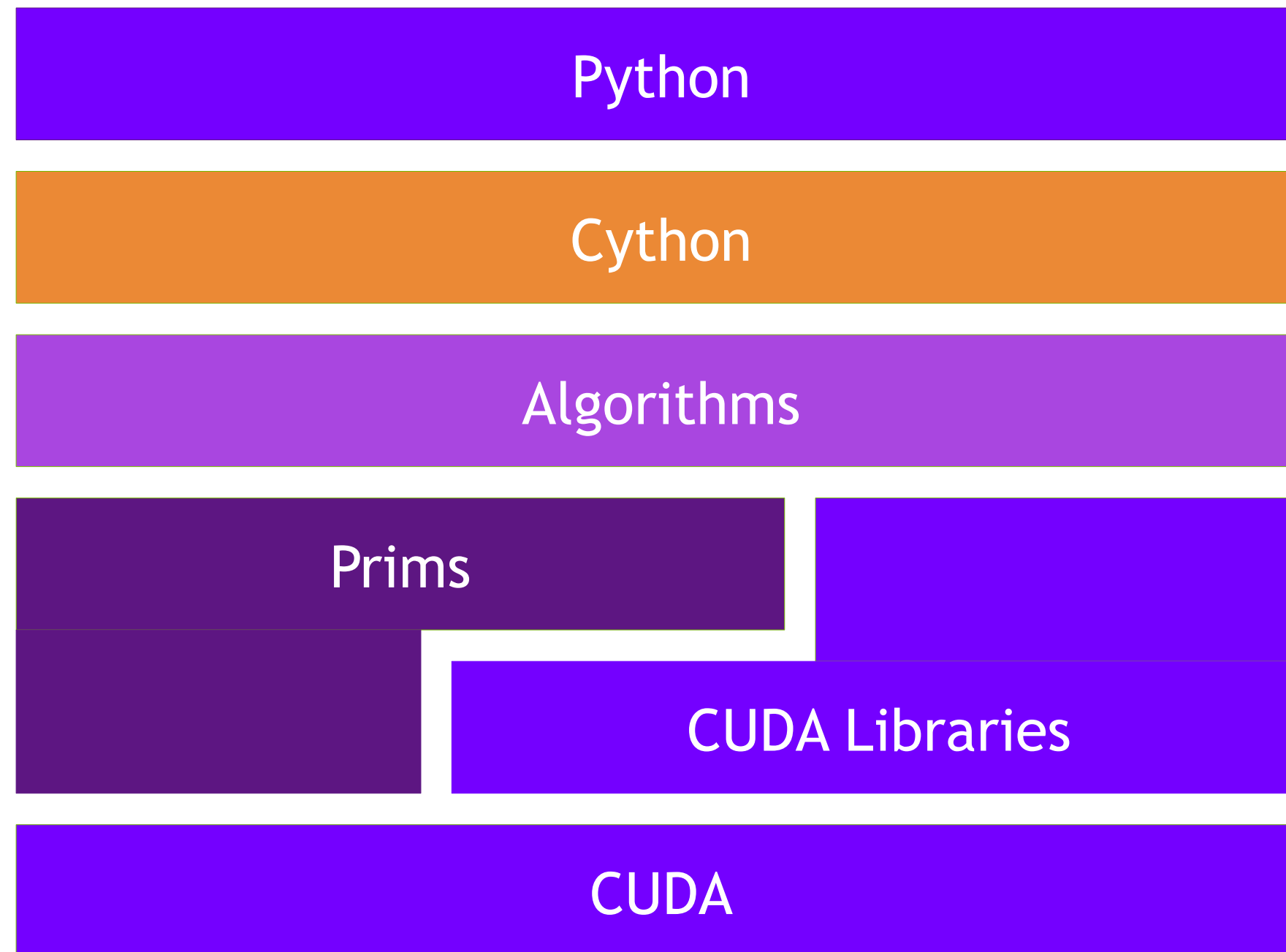
2.



3.

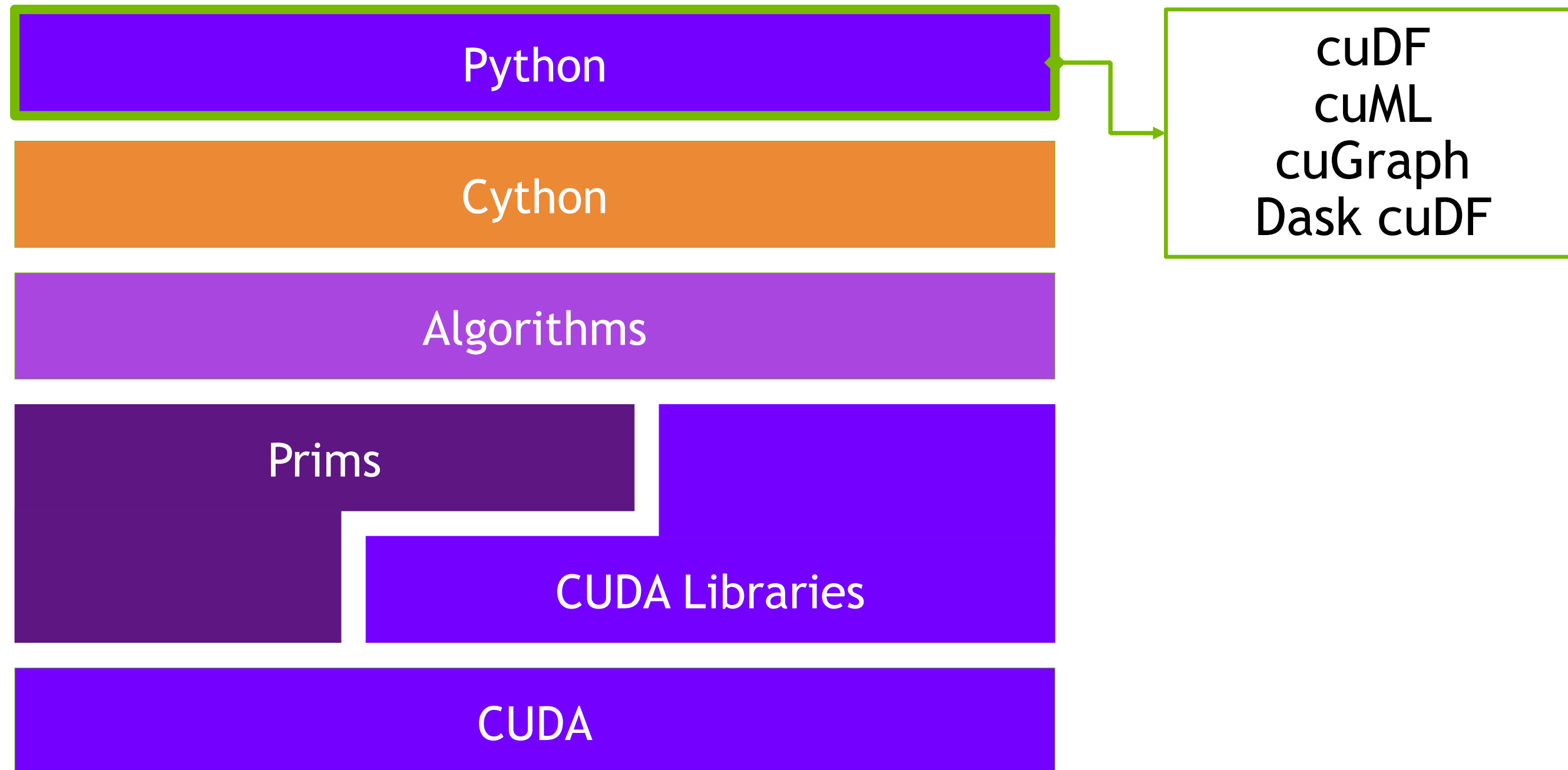


The RAPIDS Library Generalized Technology Stack



2. GPU Python Libraries

Overview of the RAPIDS GPU Library Python Stack



Refactoring

CPU to GPU Data Science

- Large amounts of existing code in PyData (Numpy, pandas, scikit-learn, etc.)
- RAPIDS uses Pandas-like API
- Very easy and straightforward
- Simple changes in a few lines of code
- Replace import statements

<code>import pandas as pd</code>	→	<code>import cudf</code>
<code>import numpy as np</code>	→	<code>import cupy as cp</code>

- Use the new imports in place of previous libraries

Example 1

Pandas to cuDF

- Use the cudf df like pandas df
 - Examples: sort_values, concat, merge, unique, std, iloc, groupby

```
import pandas as pd
df = pd.read_csv('df.csv')
df1 = pd.read_csv('df1.csv')
pd.concat([df, df1])
df.fillna(0)
df.head(10)
```



```
import cudf
df = cudf.read_csv('df.csv')
df1 = cudf.read_csv('df1.csv')
cudf.concat([df, df1])
df.fillna(0)
df.head(10)
```

Same output, but faster!

Example 2

Numpy to cuPY

- Use the cupy array like numpy array
 - Examples: randint, arrange, zeros, shape, max, flatten, sort

```
import numpy as np
choices = range(6)
```

```
probs = np.random.rand(6)
s = sum(probs)
probs = [e / s for e in probs]
selected = np.random.choice(choices, 10000, p=probs)

print(selected.shape)
```



```
import cupy as cp
choices = range(6)
```

```
probs = cp.random.rand(6)
s = sum(probs)
probs = [e / s for e in probs]
selected = cp.random.choice(choices, 10000, p=probs)

print(selected.shape)
```

Same output, but faster!

Example 3

Scikit learn to cuML

- cuML has similar capabilities as sklearn
 - Examples: train_test_split, SVC, KMeans, LinearRegression, LabelBinarizer, NearestNeighbors

```
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test =
train_test_split(X, y, random_state =0)
```

```
model = LinearRegression()
```

```
model.fit(X_train, y)
y_pred = model.predict(X_test)
```

```
import cuml.LinearRegression
from cuml.preprocessing.model_selection import
train_test_split
```

```
X_train, X_test, y_train, y_test =
train_test_split(X, y, random_state =0)
```

```
model = cuml.LinearRegression()
```

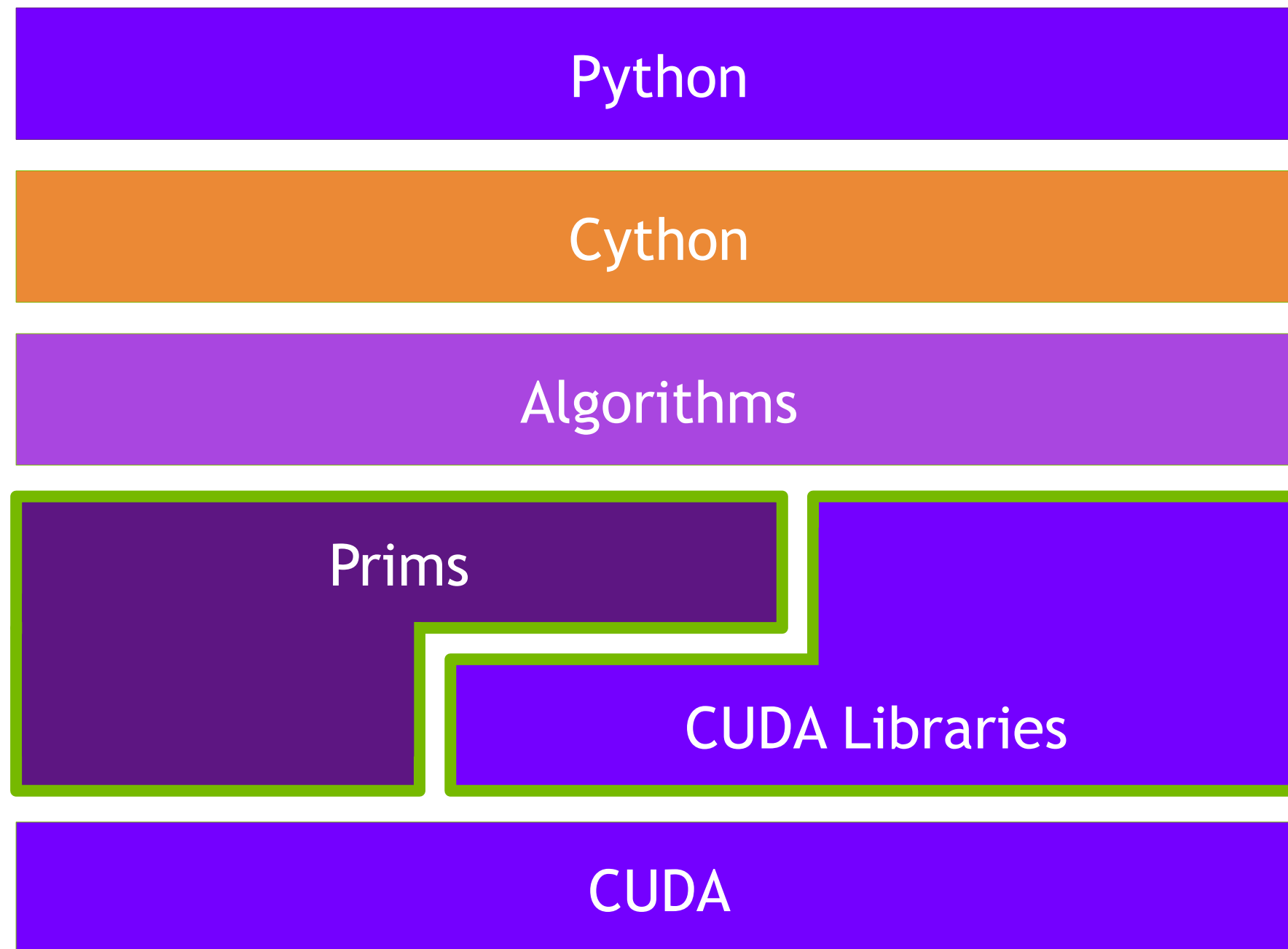
```
model.fit(X_train, y)
y_pred = model.predict(X_test)
```

Same output, but faster!

3. C++/CUDA High Level

Overview of C++/CUDA High Level

Building Data Science Applications and Expanding the RAPIDS Python Libraries



Target Users:

- Application Developer
- C++ Analytics Library

libcuDF Overview

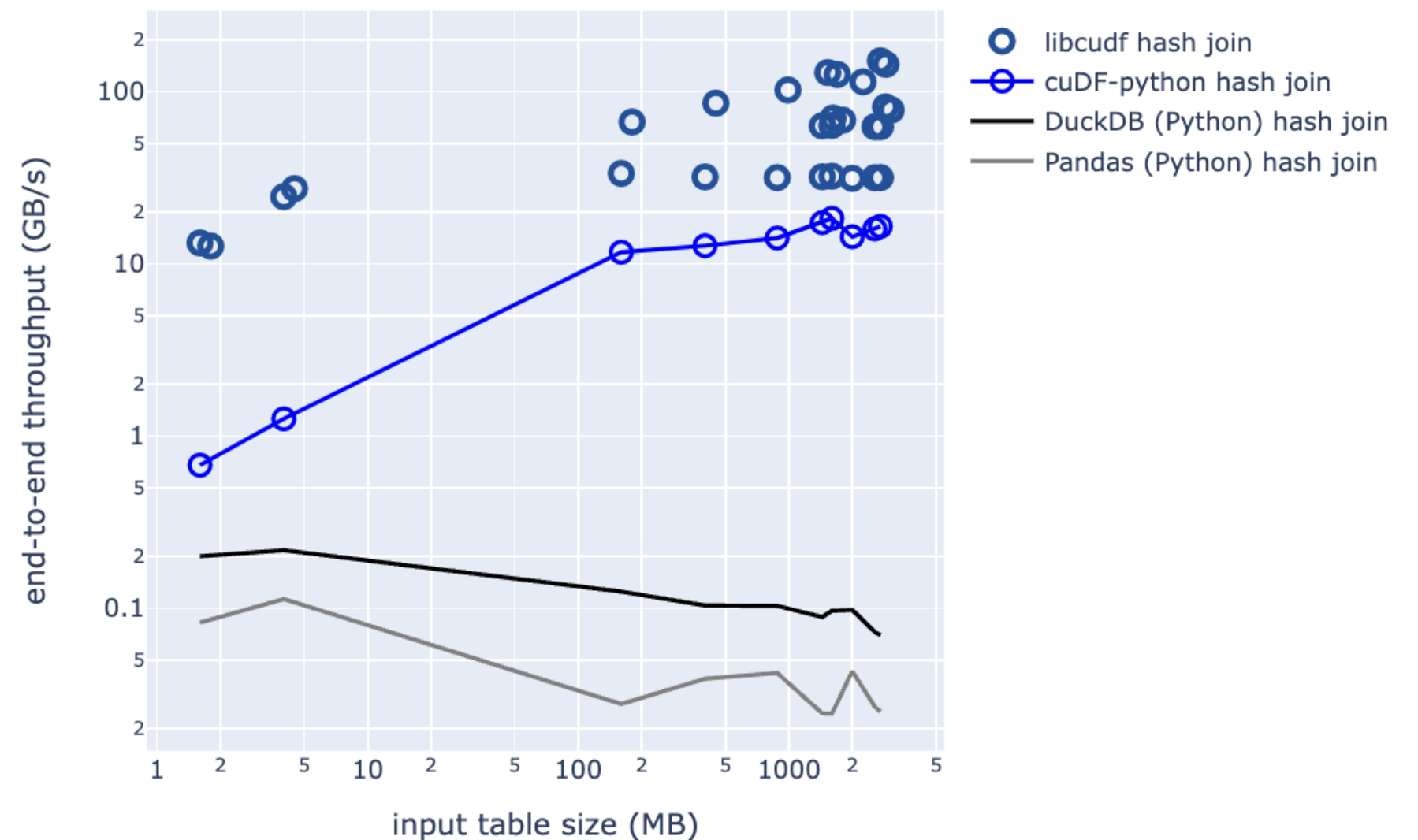
The engine powering GPU-accelerated Apache Spark, Dask, and high-performance data analytics

libcudf is the CUDA/C++ framework for tabular data analysis

- Data ingestion and parsing, joins, aggregations, filters, window functions, regular expressions, nested types, and more
- Built on the Apache Arrow memory specification
- Consistent **C++17** RAII-based APIs

Fastest library for joins, aggregations, sorting, and more

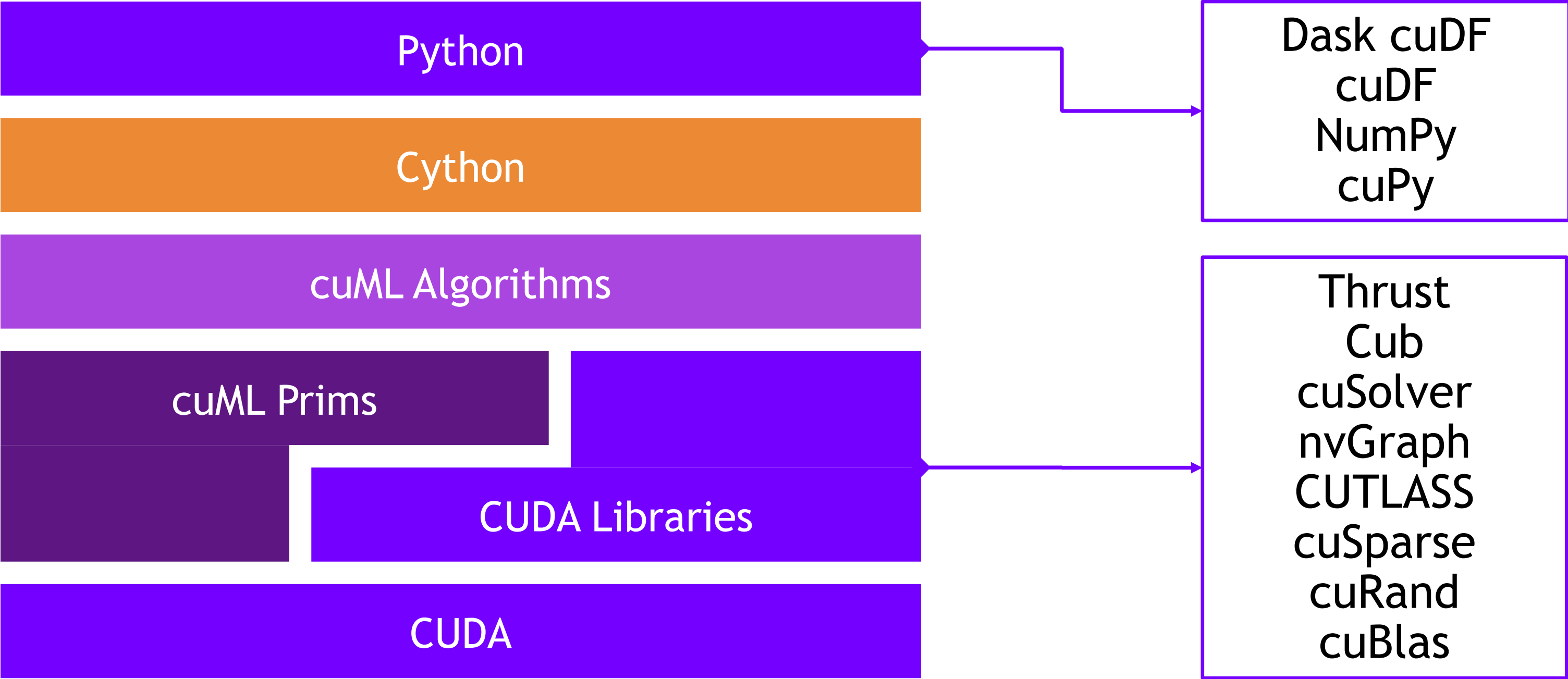
- Up to 10x faster than cuDF
- Traditional and conditional joins
- Nested-type sorting and aggregations



Explore: <https://github.com/rapidsai/cudf/tree/branch-24.12/cpp/examples>

cuML Technology Stack

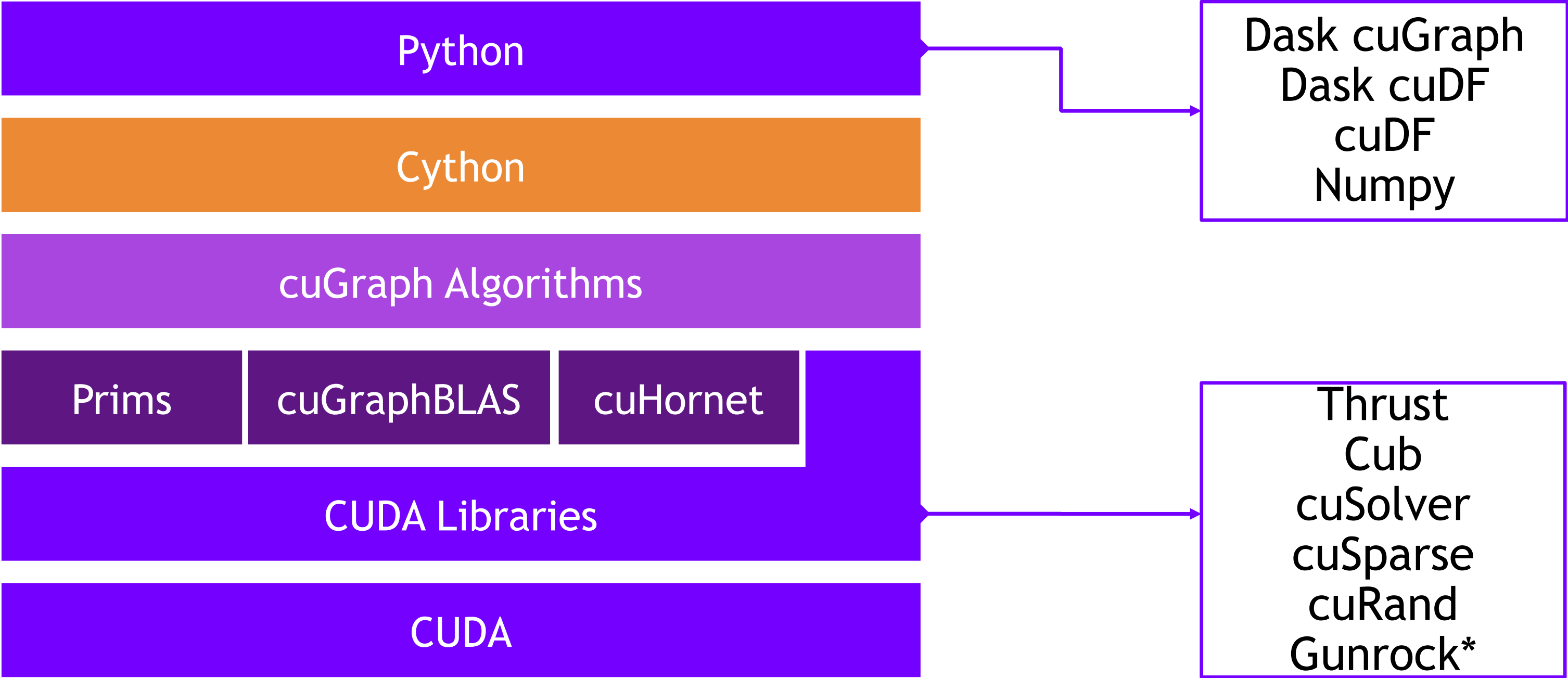
<https://github.com/rapidsai/cuml/tree/branch-24.12/cpp>



Explore: <https://github.com/rapidsai/cuml/tree/branch-24.12/cpp/examples>

cuGraph Technology Stack

<https://github.com/rapidsai/cugraph/tree/branch-24.12/cpp>



Explore: <https://github.com/rapidsai/cugraph/tree/branch-24.12/cpp/examples>



DEEP
LEARNING
INSTITUTE



PRAIRIE VIEW
A&M UNIVERSITY

DLI Accelerated Data Science Teaching Kit

Thank You