# RAPIDS

## The Platform Inside and Out
## Release 0.18

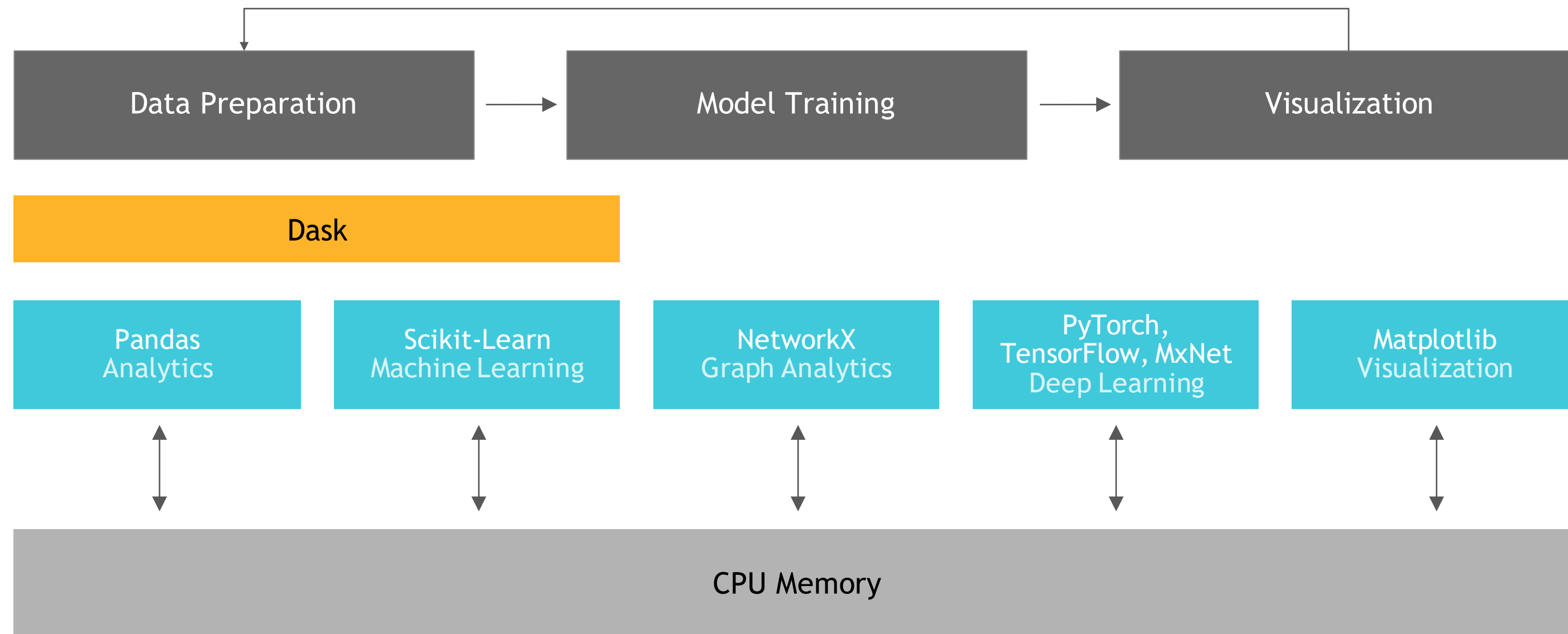# Open Source Data Science Ecosystem
## Familiar Python APIs

# RAPIDS
## End-to-End GPU Accelerated Data Science

| Data Preparation | Model Training | Visualization |
| --- | --- | --- |

| Dask | | |
| --- | --- | --- |

| cuDF cuIO Analytics | cuML Machine Learning | cuGraph Graph Analytics | PyTorch, TensorFlow, MxNet Deep Learning | cuxfilter, pyViz, plotly Visualization |
| --- | --- | --- | --- | --- |

GPU Memory

Apache Arrow

# Data Processing Evolution
## Faster Data Access, Less Data Movement

**Hadoop Processing, Reading from Disk**

| HDFS Read | Query | HDFS Write | HDFS Read | ETL | HDFS Write | HDFS Read | ML Train |
|---|---|---|---|---|---|---|---|

**Spark In-Memory Processing**

| HDFS Read | Query | ETL | ML Train |
|---|---|---|---|

**25-100x Improvement**
Less Code
Language Flexible
Primarily In-Memory

**Traditional GPU Processing**

| HDFS Read | GPU Read | Query | CPU Write | GPU Read | ETL | CPU Write | GPU Read | ML Train |
|---|---|---|---|---|---|---|---|---|

**5-10x Improvement**
More Code
Language Rigid
Substantially on GPU

# Data Processing Evolution
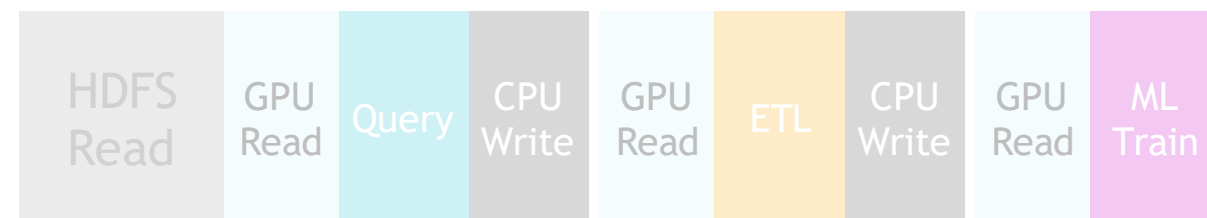## Faster Data Access, Less Data Movement

**Hadoop Processing, Reading from Disk**

| HDFS Read | Query | HDFS Write | HDFS Read | ETL | HDFS Write | HDFS Read | ML Train |

**Spark In-Memory Processing**

| HDFS Read | Query | ETL | ML Train |

**25-100x Improvement**
Less Code
Language Flexible
Primarily In-Memory

**Traditional GPU Processing**

| HDFS Read | GPU Read | Query | CPU Write | GPU Read | ETL | CPU Write | GPU Read | ML Train |

**5-10x Improvement**
More Code
Language Rigid
Substantially on GPU

**RAPIDS**

| Arrow Read | Query | ETL | ML Train |

**50-100x Improvement**
Same Code
Language Flexible
Primarily on GPU

# Faster Speeds, Real World Benefits
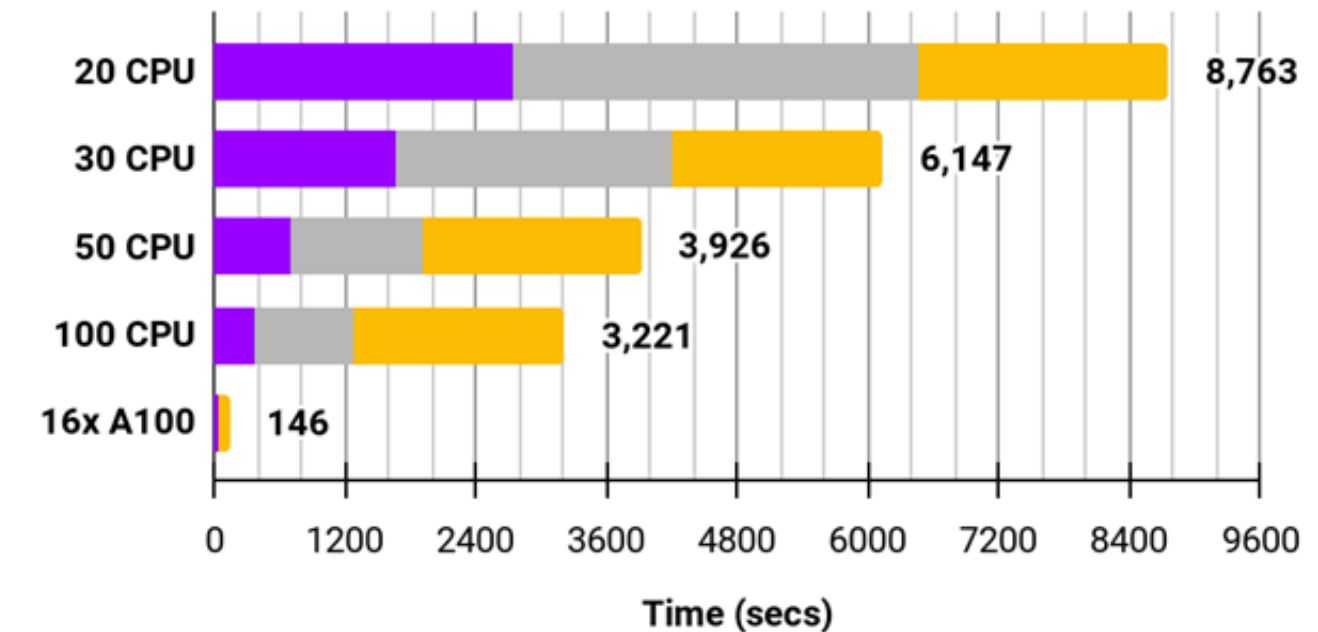## Faster Data Access, Less Data Movement

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

| | Time (sec) |
|---|---|
| 20 CPU Nodes | 2,74 |
| 30 CPU Nodes | 1,675 |
| 50 CPU Nodes | 715 |
| 100 CPU | 379 |
| 16x A100 | 32 |

**XGBoost Machine Learning**

| | Time (sec) |
|---|---|
| 20 CPU | 2,290 |
| 30 CPU | 1,956 |
| 50 CPU | 1,999 |
| 100 CPU | 1,948 |
| 16x A100 | 99 |

**End-to-End**

| | Time (secs) |
|---|---|
| 20 CPU | 8,763 |
| 30 CPU | 6,147 |
| 50 CPU | 3,926 |
| 100 CPU | 3,221 |
| 16x A100 | 146 |

**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**
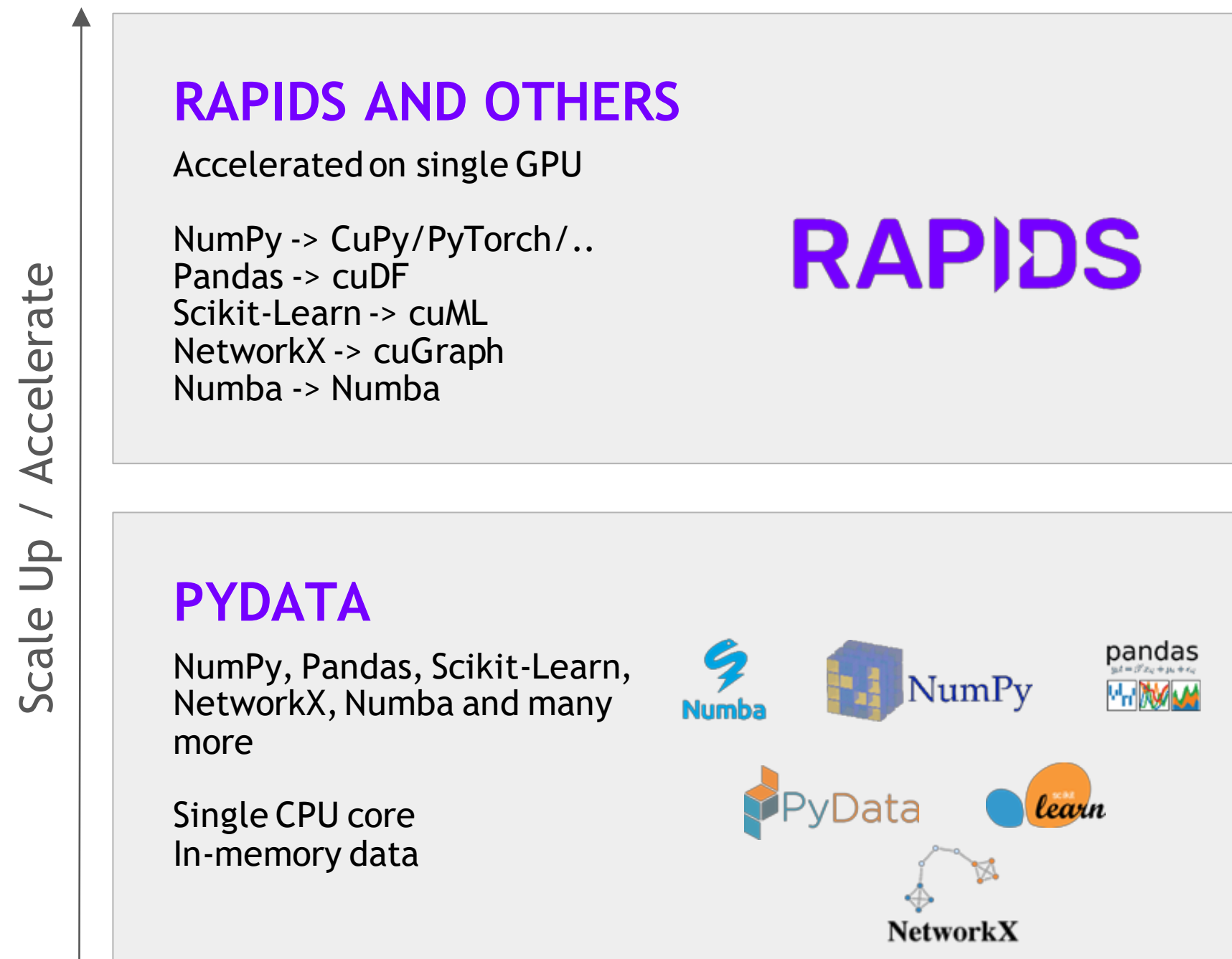CPU nodes (61 GiB memory, 8 vCPUs, 64-bit platform), Apache Spark
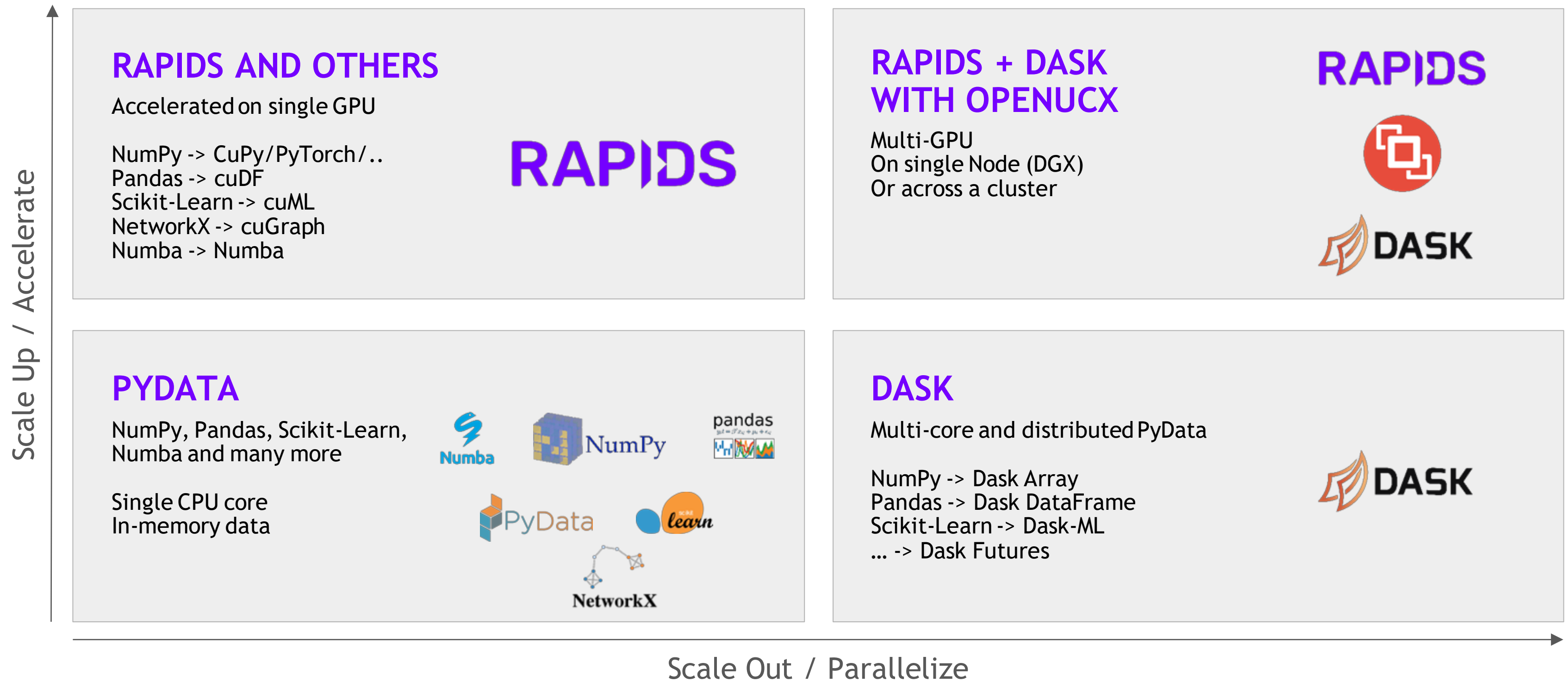
**A100 Cluster Configuration**
16 A100 GPUs (40GB each)
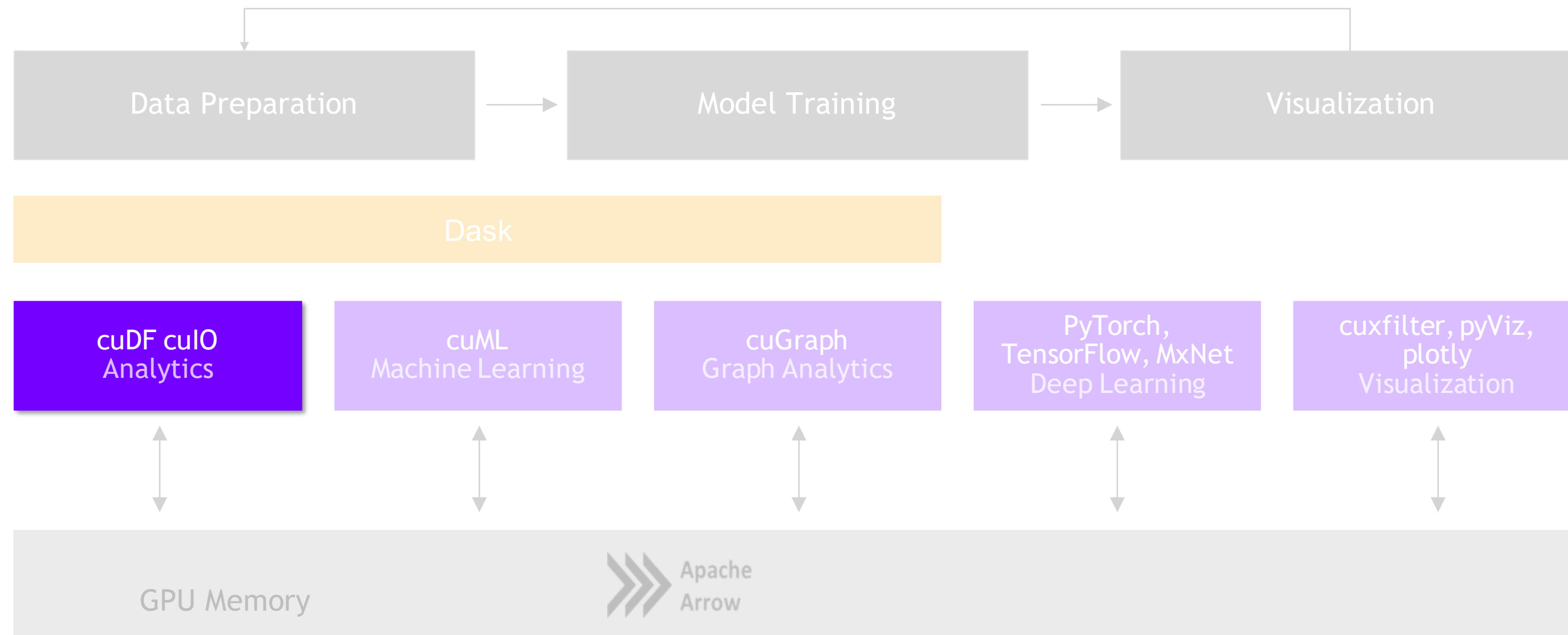
**RAPIDS Version**
RAPIDS 0.17

# Scale Up with RAPIDS



Scale Up / Accelerate

## RAPIDS AND OTHERS

Accelerated on single GPU

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

## PYDATA

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

Single CPU core
In-memory data

# Scale Out with RAPIDS + Dask with OpenUCX

**Scale Up / Accelerate** (vertical axis label)

## RAPIDS AND OTHERS

Accelerated on single GPU

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

**RAPIDS**

## RAPIDS + DASK WITH OPENUCX

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

**RAPIDS**

**DASK**

## PYDATA

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

Single CPU core
In-memory data

Numba NumPy pandas

PyData learn

NetworkX

## DASK

Multi-core and distributed PyData

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

**DASK**

**Scale Out / Parallelize**

cuDF

# RAPIDS

## GPU Accelerated Data Wrangling and Feature Engineering

| Data Preparation | Model Training | Visualization |
|---|---|---|

| Dask |
|---|

| cuDF cuIO<br>Analytics | cuML<br>Machine Learning | cuGraph<br>Graph Analytics | PyTorch,<br>TensorFlow, MxNet<br>Deep Learning | cuxfilter, pyViz,<br>plotly<br>Visualization |
|---|---|---|---|---|

GPU Memory

Apache
Arrow

# ETL - the Backbone of Data Science
## cuDF is...

### 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

```
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_Ca |
|---|---------|------------|--------|-----|-----------|---------------|---------------------------|----------------|------------|
| 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()
```

# 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



■ 10M   ■ 100M

GPU Speedup Over CPU

| | Merge | Sort | GroupBy |
|---|---|---|---|
| 100M | 970 | 370 | 330 |
| 10M | 500 | 350 | 320 |

# cuML

# Machine Learning
## More Models More Problems



| Data Preparation | Model Training | Visualization |

**Dask**

| cuDF cuIO<br>Analytics | cuML<br>Machine Learning | cuGraph<br>Graph Analytics | PyTorch,<br>TensorFlow, MxNet<br>Deep Learning | cuxfilter, pyViz,<br>plotly<br>Visualization |

GPU Memory    Apache Arrow

# 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)


y_hat = dbscan.fit_predict(X)
```

# RAPIDS Matches Common Python APIs
## GPU-accelerated Clustering

```python
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])})
```

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


y_hat = dbscan.fit_predict(X)
```

# Algorithms
## GPU-accelerated Scikit-Learn

**Classification / Regression**

Decision Trees / **Random Forests**
Linear/Lasso/Ridge/**LARS**/ElasticNet Regression
Logistic Regression
K-Nearest Neighbors (**exact or approximate**)
Support Vector Machine Classification and Regression
Naive Bayes

**Inference**

Random Forest / GBDT Inference (FIL)

**Preprocessing**

Text vectorization (TF-IDF / Count)
Target Encoding
Cross-validation / splitting

**Clustering
Decomposition &
Dimensionality Reduction**

K-Means
**DBSCAN**
Spectral Clustering
Principal Components (including iPCA)
Singular Value Decomposition
UMAP
Spectral Embedding
T-SNE

**Cross Validation**

**Hyper-parameter Tuning**

**Time Series**

Holt-Winters
Seasonal ARIMA / Auto ARIMA

More to come!

# Benchmarks: Single-GPU cuML vs Scikit-learn



1x V100 vs. 2x 20 Core CPUs (DGX-1, RAPIDS 0.15)

# Getting Started

# Easy Installation
## Interactive Installation Guide



https://rapids.ai/start.html

# RAPIDS Docs
## Up to Date and Easy to Use



*https://docs.rapids.ai*