# RAPIDS: Accelerated Data Science and Data Processing
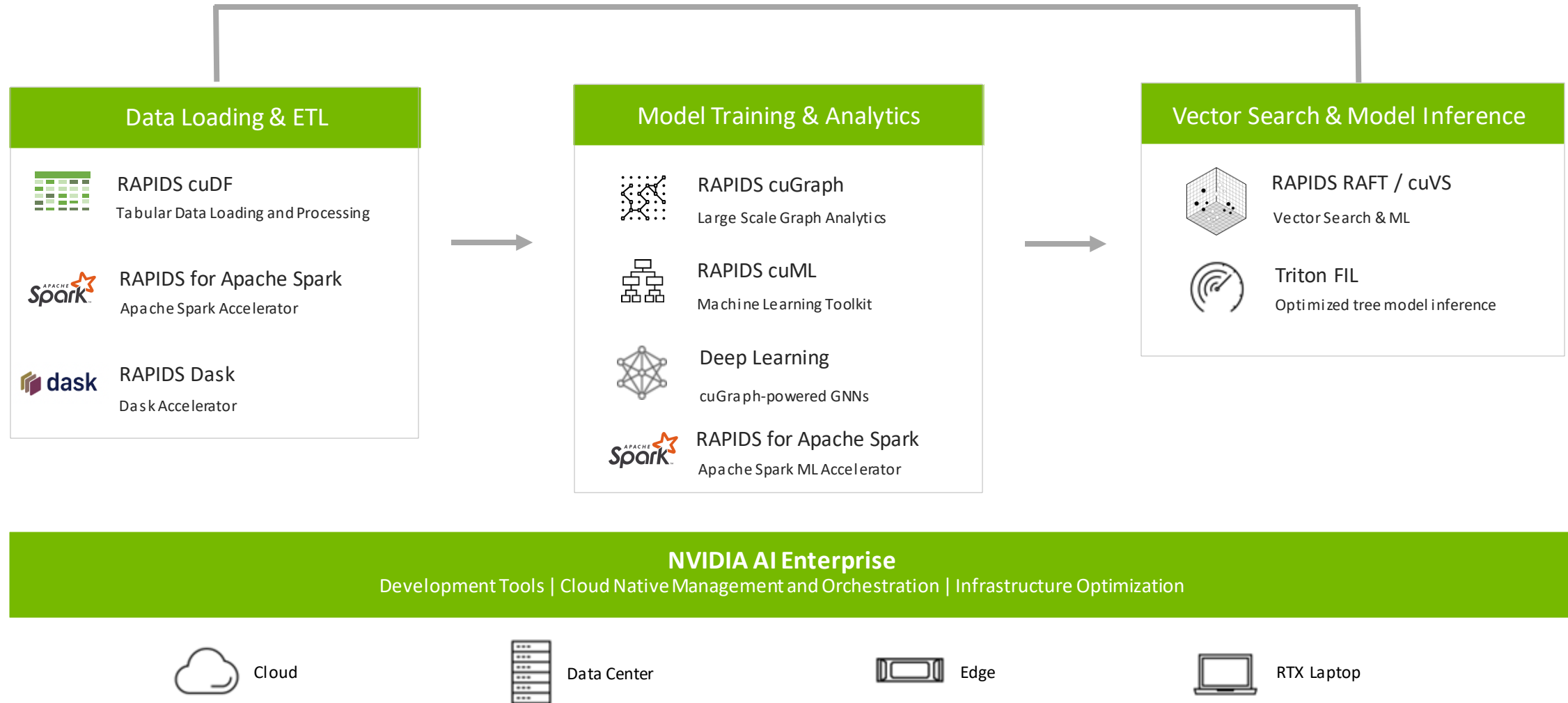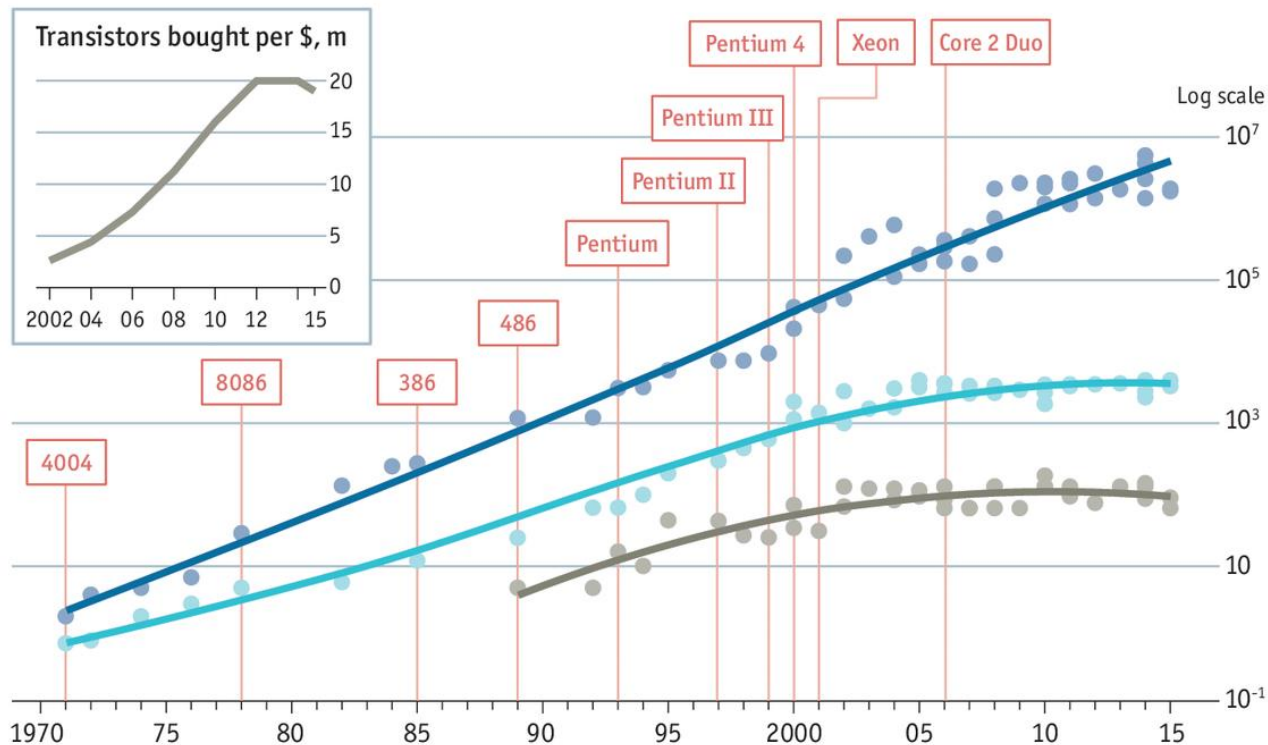
Mads R. B. Kristensen, NVIDIA

# RAPIDS Accelerates Data Science End-to-End
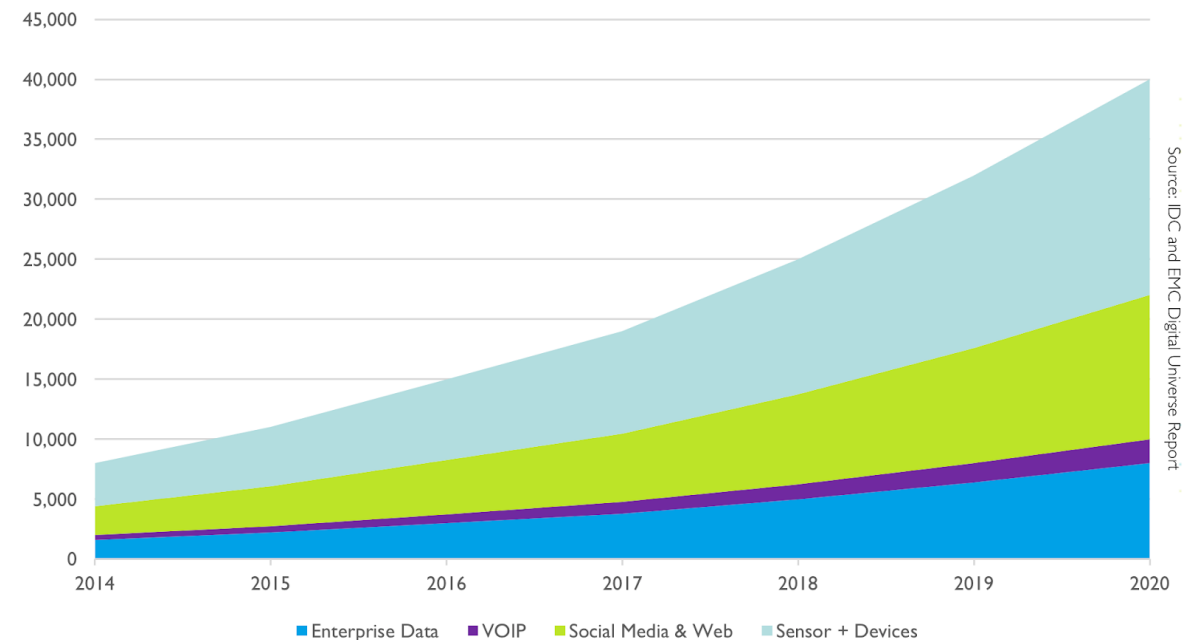
## Data Loading & ETL

**RAPIDS cuDF**
Tabular Data Loading and Processing

**RAPIDS for Apache Spark**
Apache Spark Accelerator

**RAPIDS Dask**
Dask Accelerator

## Model Training & Analytics

**RAPIDS cuGraph**
Large Scale Graph Analytics

**RAPIDS cuML**
Machine Learning Toolkit

**Deep Learning**
cuGraph-powered GNNs

**RAPIDS for Apache Spark**
Apache Spark ML Accelerator

## Vector Search & Model Inference

**RAPIDS RAFT / cuVS**
Vector Search & ML

**Triton FIL**
Optimized tree model inference

## NVIDIA AI Enterprise
Development Tools | Cloud Native Management and Orchestration | Infrastructure Optimization

Cloud          Data Center          Edge          RTX Laptop

# Stuttering

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



Transistors bought per $, m

Pentium 4   Xeon   Core 2 Duo

Pentium III

Pentium II

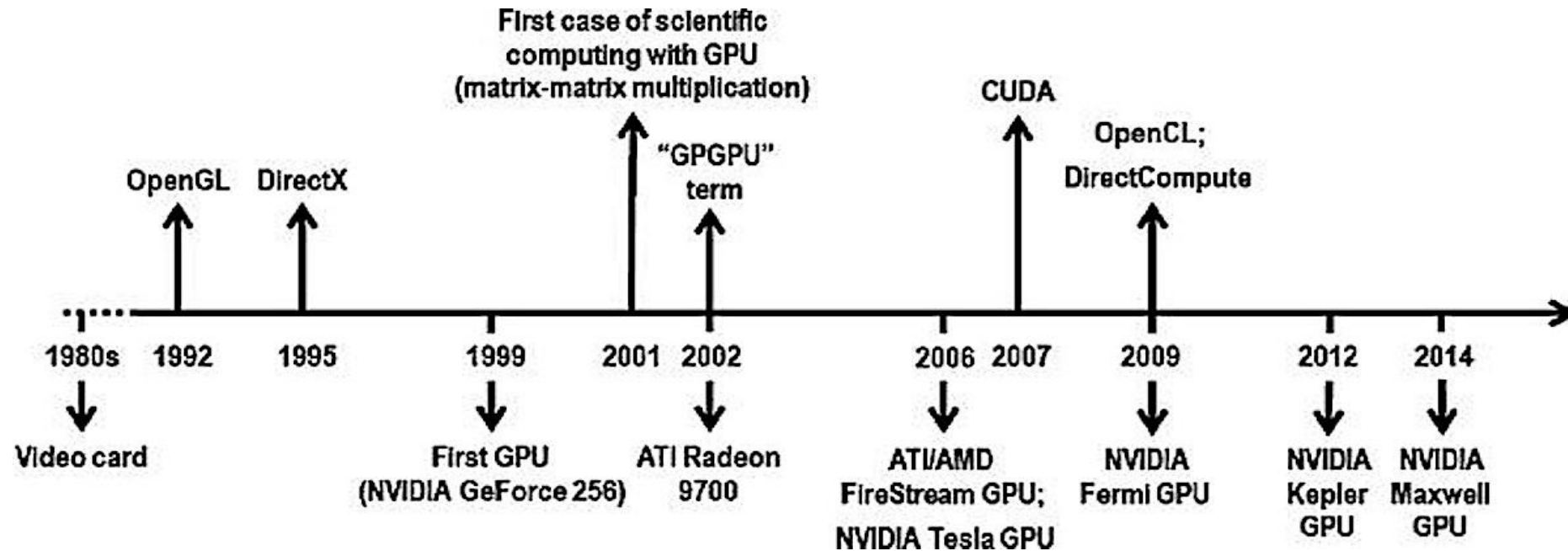Pentium

486

386

8086

4004

Log scale

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

Economist.com



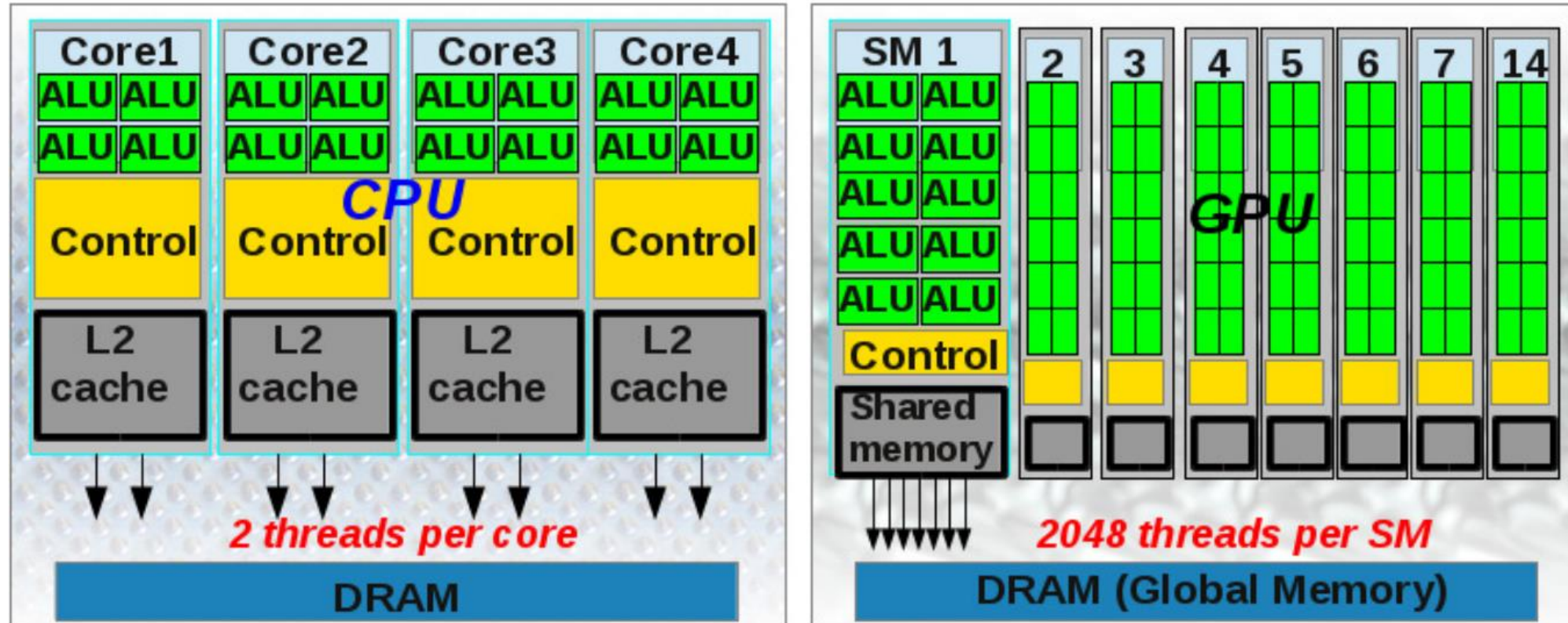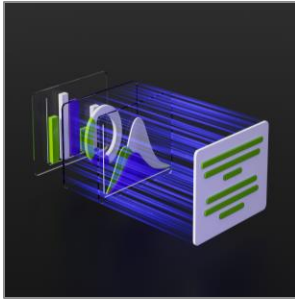## Data Growth and Source in Exabytes

■ Enterprise Data   ■ VOIP   ■ Social Media & Web   ■ Sensor + Devices

Source: IDC and EMC Digital Universe Report

3

# History of the GPU

# CPU vs GPU

# Modern Enterprise Applications Need Accelerated Computing

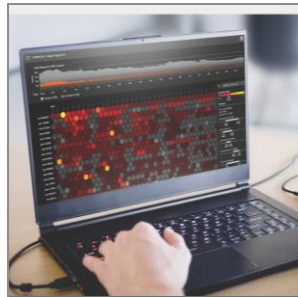Internet scale data | Massive models | Real-time performance
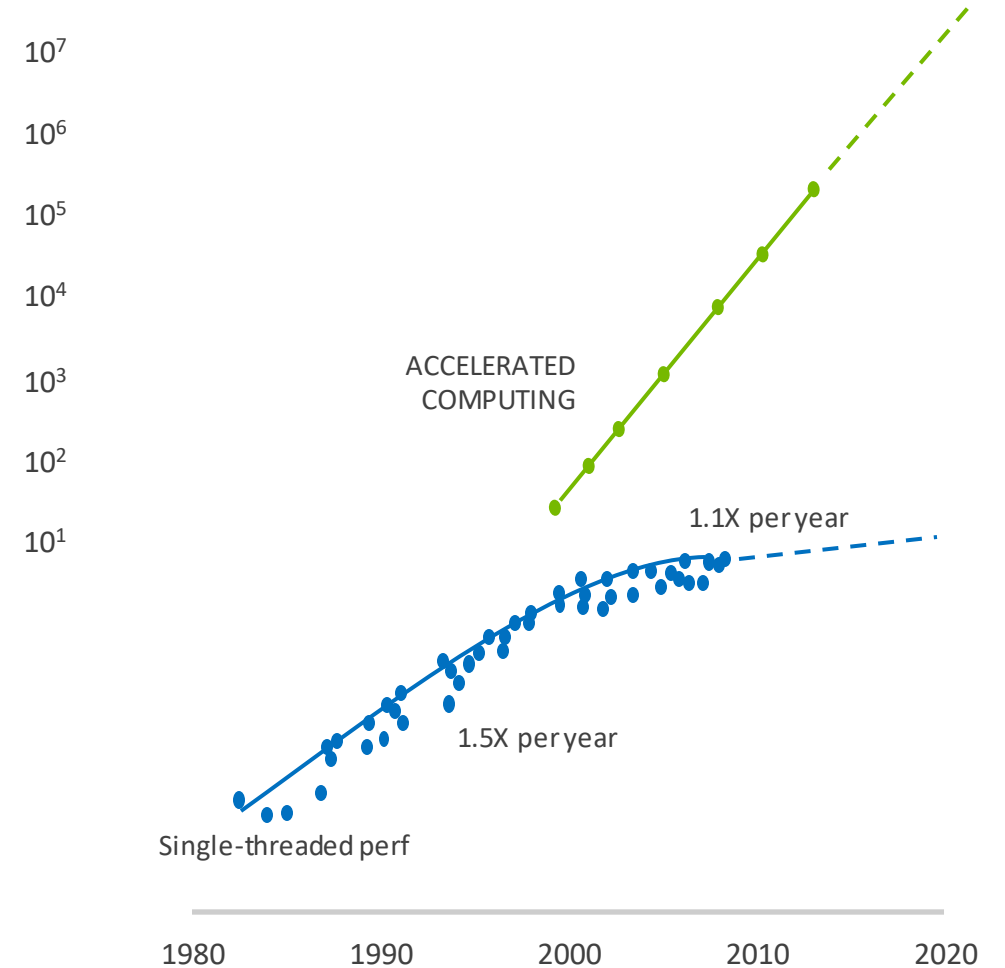


Recommenders


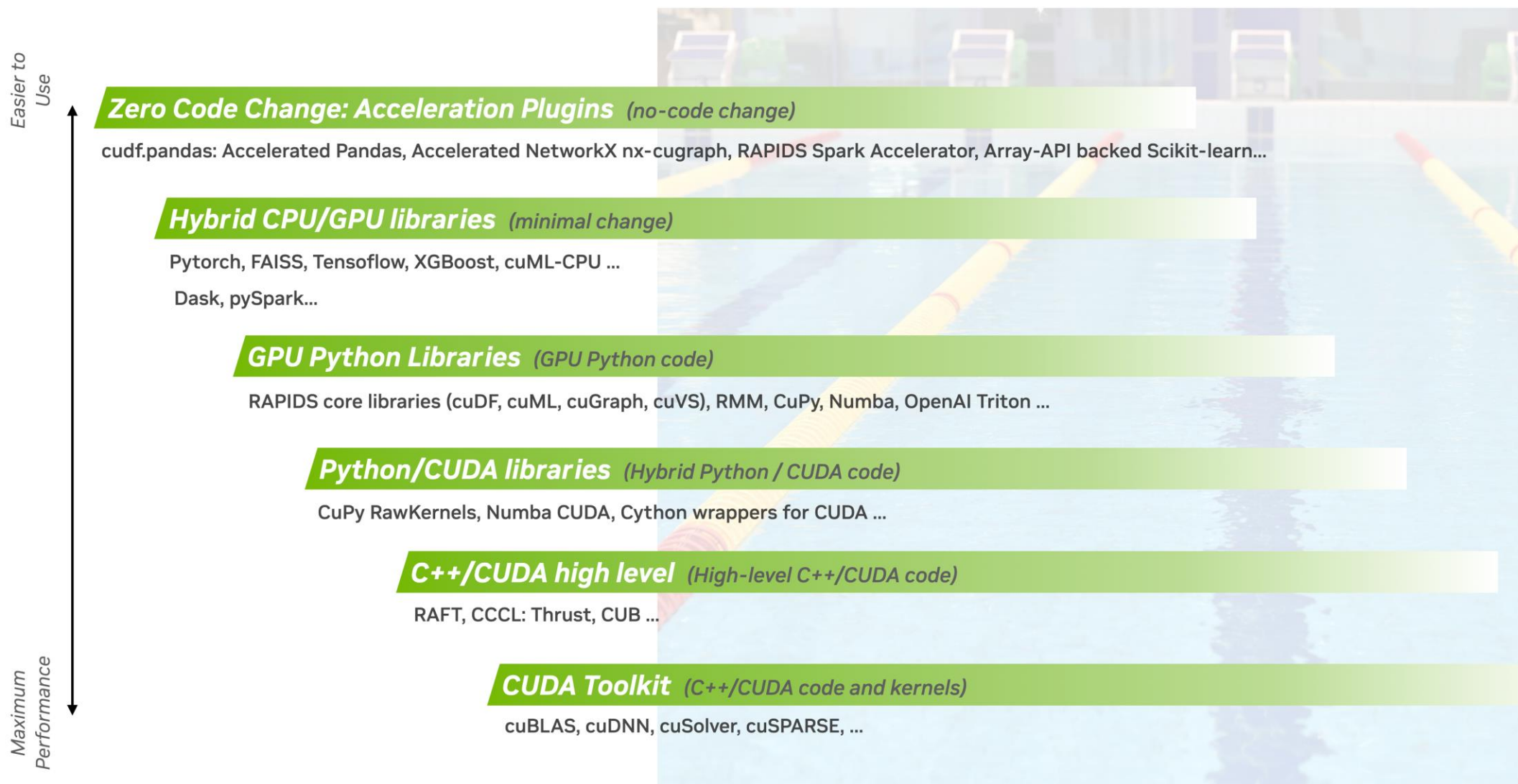
LLMs



Forecasting



Fraud Detection



Genomic Analysis



Cybersecurity



ACCELERATED COMPUTING

1.1X per year

1.5X per year

Single-threaded perf

$10^7$

$10^6$

$10^5$

$10^4$

$10^3$

$10^2$

$10^1$

1980    1990    2000    2010    2020

# Accelerated Computing Swim Lanes

RAPIDS makes accelerated computing more seamless while enabling specialization for maximum performance

*Easier to Use*

**Zero Code Change: Acceleration Plugins** *(no-code change)*

cudf.pandas: Accelerated Pandas, Accelerated NetworkX nx-cugraph, RAPIDS Spark Accelerator, Array-API backed Scikit-learn...

**Hybrid CPU/GPU libraries** *(minimal change)*

Pytorch, FAISS, Tensoflow, XGBoost, cuML-CPU ...

Dask, pySpark...

**GPU Python Libraries** *(GPU Python code)*

RAPIDS core libraries (cuDF, cuML, cuGraph, cuVS), RMM, CuPy, Numba, OpenAI Triton ...

**Python/CUDA libraries** *(Hybrid Python / CUDA code)*

CuPy RawKernels, Numba CUDA, Cython wrappers for CUDA ...

**C++/CUDA high level** *(High-level C++/CUDA code)*

RAFT, CCCL: Thrust, CUB ...

**CUDA Toolkit** *(C++/CUDA code and kernels)*

cuBLAS, cuDNN, cuSolver, cuSPARSE, ...

*Maximum Performance*

NVIDIA.

# RAPIDS ETL
Extract, transform, and load

# Pandas

## Python's Preeminent DataFrame Library

**9.5M**
Pandas users

**65%**
Users love
using pandas
(HuggingFace
most loved at
72%)

**135M+**
Monthly
downloads

**25%**
Y/Y downloads
growth



DataFrame

row

column

```python
import pandas as pd

df = (
    pd.read_csv('data.csv')
    .melt(id_vars=['id', 'name'])
    .rename(columns={
        'variable':'var',
        'value':'val'})
    .query('val >= 200')
    .sort_values('val', ascending=False)
)
```


NVIDIA.

# cuDF - GPU DataFrames

```python
%%time
import pandas as pd

df = pd.read_csv("rows.csv",
                 usecols=[
                     "Registration State",
                     "Violation Description",
                     "Vehicle Body Type"
                 ])
(df[["Registration State", "Violation Description"]]
 .value_counts()   # count of offences per state per offence
 .groupby("Registration State")   # grouped by state
 .head(1)   # offence with largest count
 .sort_index()   # sort by state name
 .reset_index()
 .head(5)
)
```

```python
%%time
import cudf

df = cudf.read_csv("rows.csv",
                   usecols=[
                       "Registration State",
                       "Violation Description",
                       "Vehicle Body Type"
                   ])
(df[["Registration State", "Violation Description"]]
 .value_counts()   # count of offences per state per offence
 .groupby("Registration State")   # grouped by state
 .head(1)   # offence with largest count
 .sort_index()   # sort by state name
 .reset_index()
 .head(5)
)
```
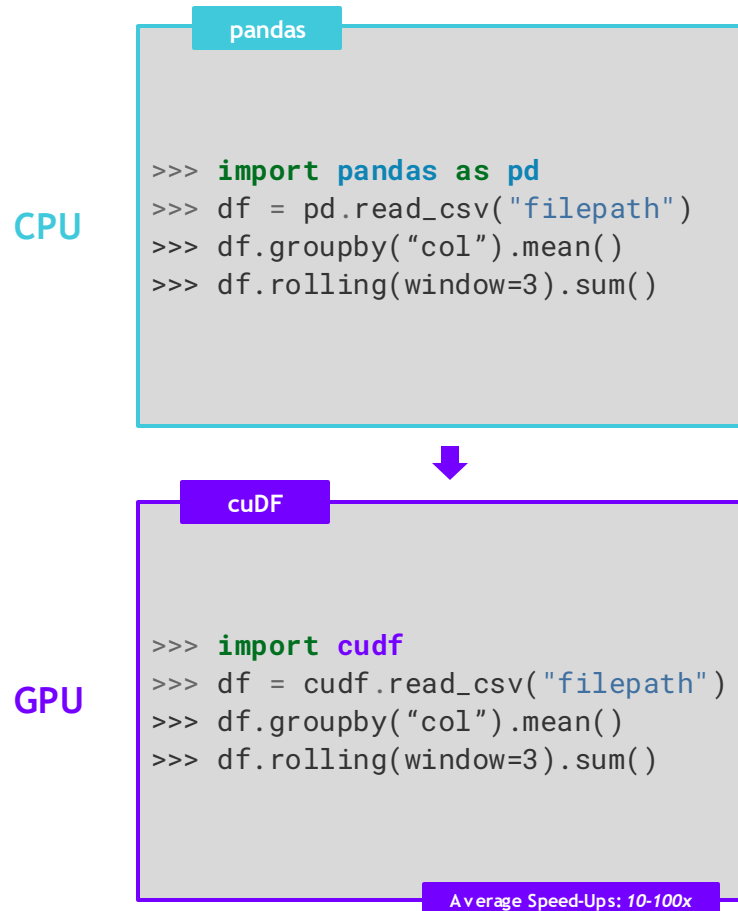
# cuDF Problems



*cuDF coverage of the Pandas API (green=implemented, gray=not implemented)*

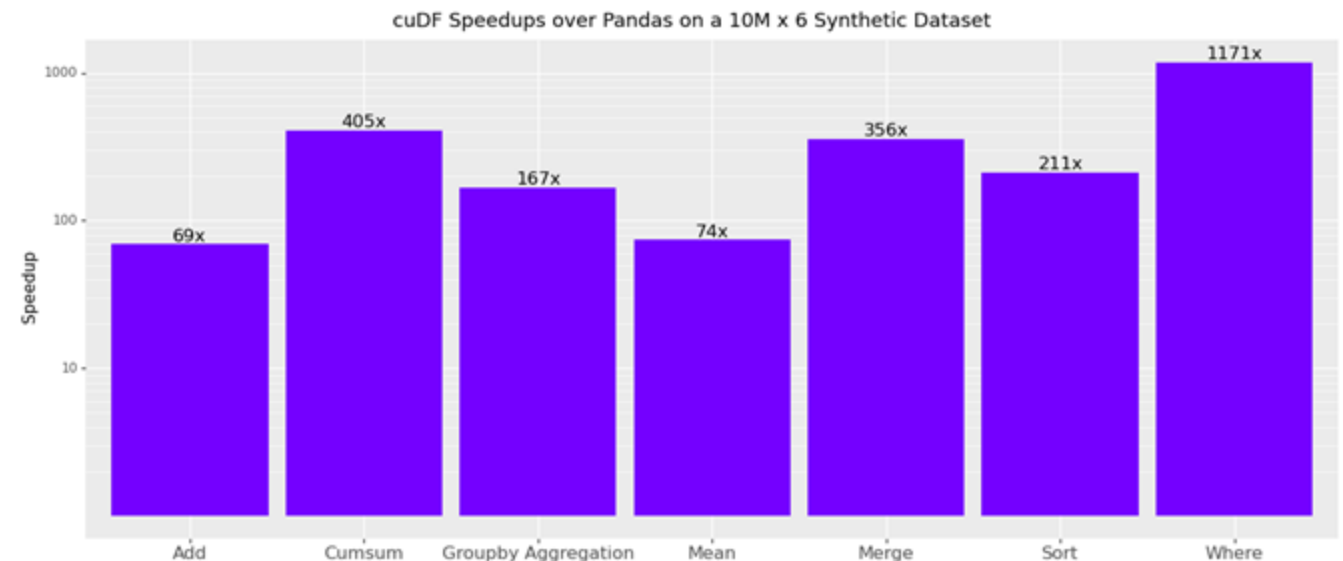*"We just don't have time to rewrite our code in a new paradigm."*

# cuDF

A GPU DataFrame library in Python with a pandas-like API built into the PyData ecosystem

## Pandas-like API on the GPU

**pandas**

**CPU**

```
>>> import pandas as pd
>>> df = pd.read_csv("filepath")
>>> df.groupby("col").mean()
>>> df.rolling(window=3).sum()
```

**cuDF**

**GPU**

```
>>> import cudf
>>> df = cudf.read_csv("filepath")
>>> df.groupby("col").mean()
>>> df.rolling(window=3).sum()
```

*Average Speed-Ups: 10-100x*

## Best-in-Class Performance ([Benchmark](#))



cuDF Speedups over Pandas on a 10M x 6 Synthetic Dataset

| Groupby | Strings and Regex | UDFs | Nested Types | Time Series |
|---|---|---|---|---|

| Indexing | Missing Data | CuPy Interoperability | Rolling Windows |
|---|---|---|---|

[10 Minutes to cuDF](#)

NVIDIA A100 vs. AMD EPYC 7642 48-Core Processor
cuDF Python vs. Pandas

# Accelerated pandas

- Requires *no changes* to existing pandas code. Just
  - `%load_ext cudf.pandas`
  - `$ python –m cudf.pandas <script.py>`
- 100% of the pandas API
- Accelerates workflows up to 150x using the GPU
- Compatible with code that uses third-party libraries
- Falls back to using pandas on the CPU for unsupported functions and methods

```
[ ]:
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

data = pd.read_parquet("data.parquet")
subset = data.index.indexer_between_time("09:30", "16:00")
data = data.iloc[subset]
results = data.groupby(pd.Grouper(freq="1D")).mean()

sns.lineplot(results)
plt.xticks(rotation=30)
```

Load cudf.pandas

Module spoofing the pandas import

Your Code

Third-party Library Using Pandas

Call Pandas Function/Method

Use Equivalent cuDF Function/Method

If operation unsupported

Copy to CPU

Fallback to pandas

Copy to GPU

Result

NVIDIA

# Pandas Accelerator Mode for cuDF (cudf.pandas)

```python
%%time
import pandas as pd

df = pd.read_csv("rows.csv",
                 usecols=[
                     "Registration State",
                     "Violation Description",
                     "Vehicle Body Type"
                 ])
(df[["Registration State", "Violation Description"]]
 .value_counts()  # count of offences per state per offence
 .groupby("Registration State")  # grouped by state
 .head(1)  # offence with largest count
 .sort_index()  # sort by state name
 .reset_index()
 .head(5)
)
```

```python
%%time
%load_ext cudf.pandas
import pandas as pd

df = pd.read_csv("rows.csv",
                 usecols=[
                     "Registration State",
                     "Violation Description",
                     "Vehicle Body Type"
                 ])
(df[["Registration State", "Violation Description"]]
 .value_counts()  # count of offences per state per offence
 .groupby("Registration State")  # grouped by state
 .head(1)  # offence with largest count
 .sort_index()  # sort by state name
 .reset_index()
 .head(5)
)
```
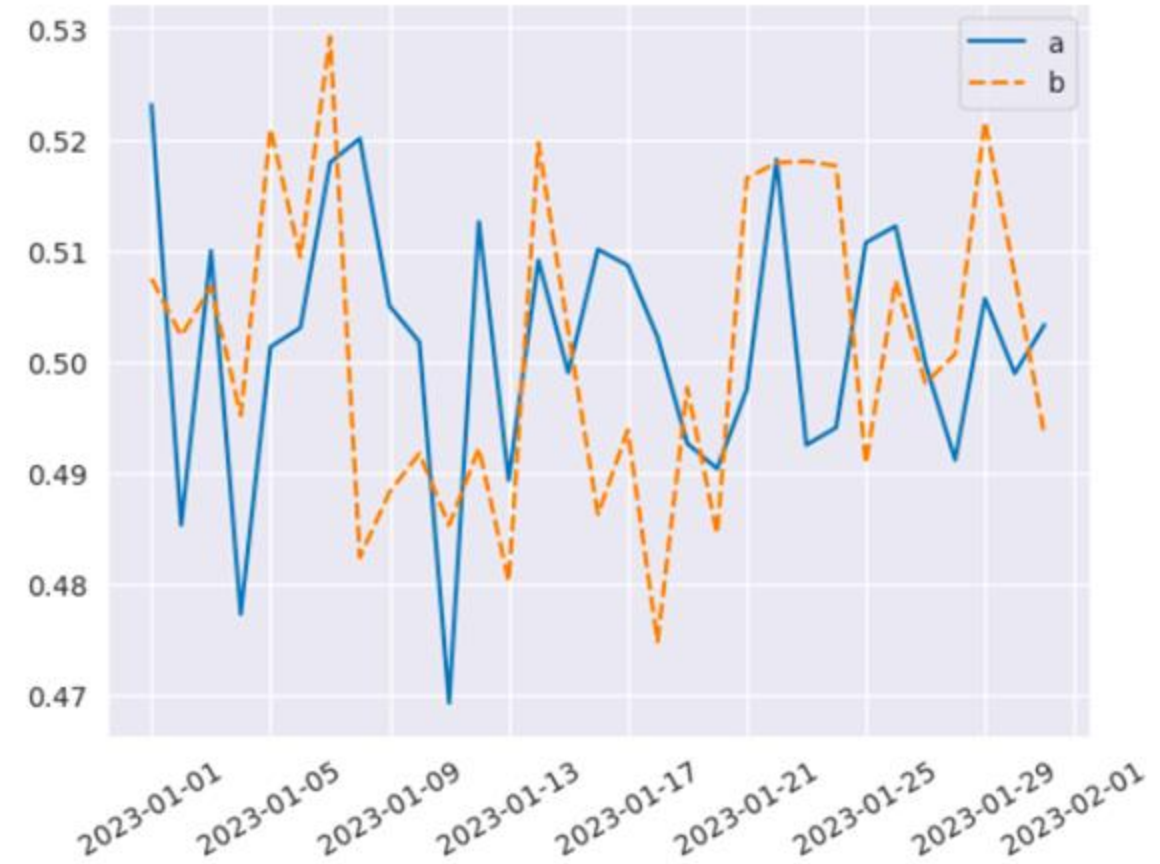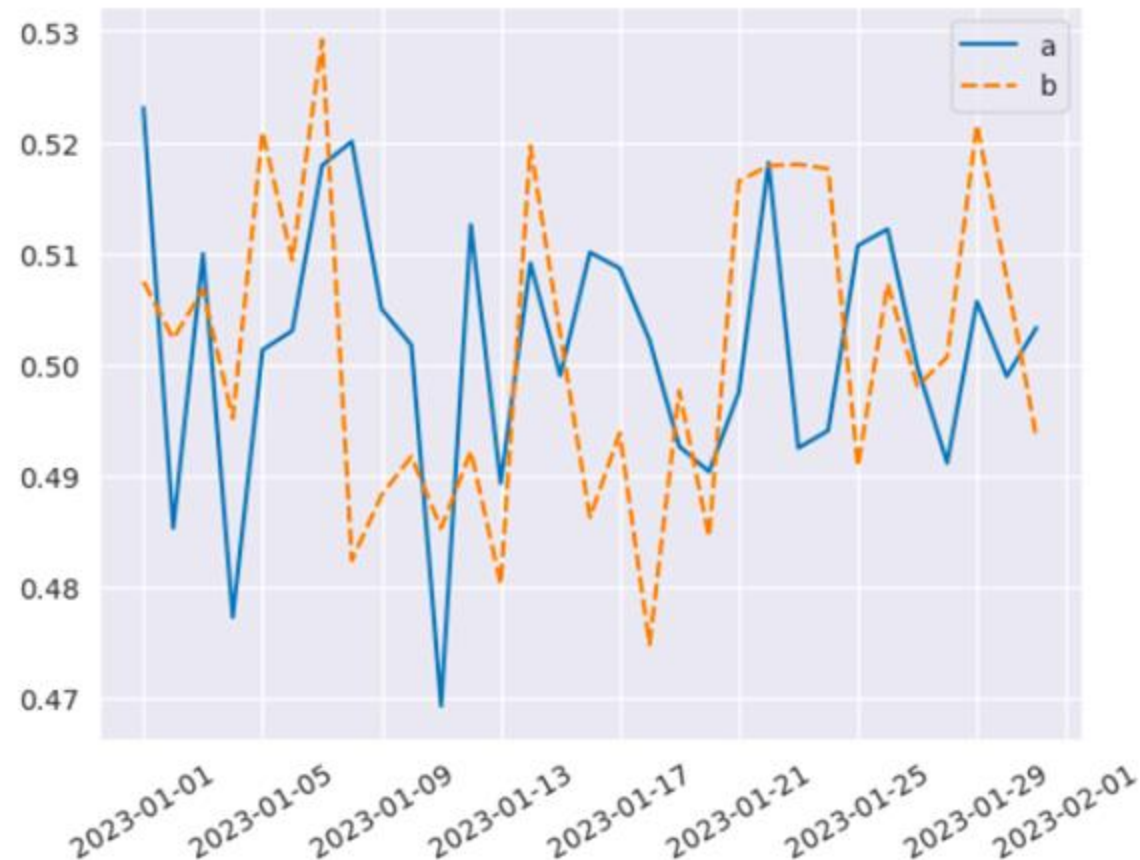
# cudf.pandas in Action

## A brief example

```python
%%load_ext cudf.pandas

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

rng = pd.date_range("2023-01-01", "2023-02-01",  freq="1T")
df = pd.DataFrame({
    "a": np.random.rand(len(rng)),
    "b": np.random.rand(len(rng))
    },
    index=rng
)

df = df.iloc[rng.indexer_between_time("09:30", "16:00")]
results = df.groupby(pd.Grouper(freq="1D")).mean()

_ = sns.lineplot(results)
_ = plt.xticks(rotation=30)
```
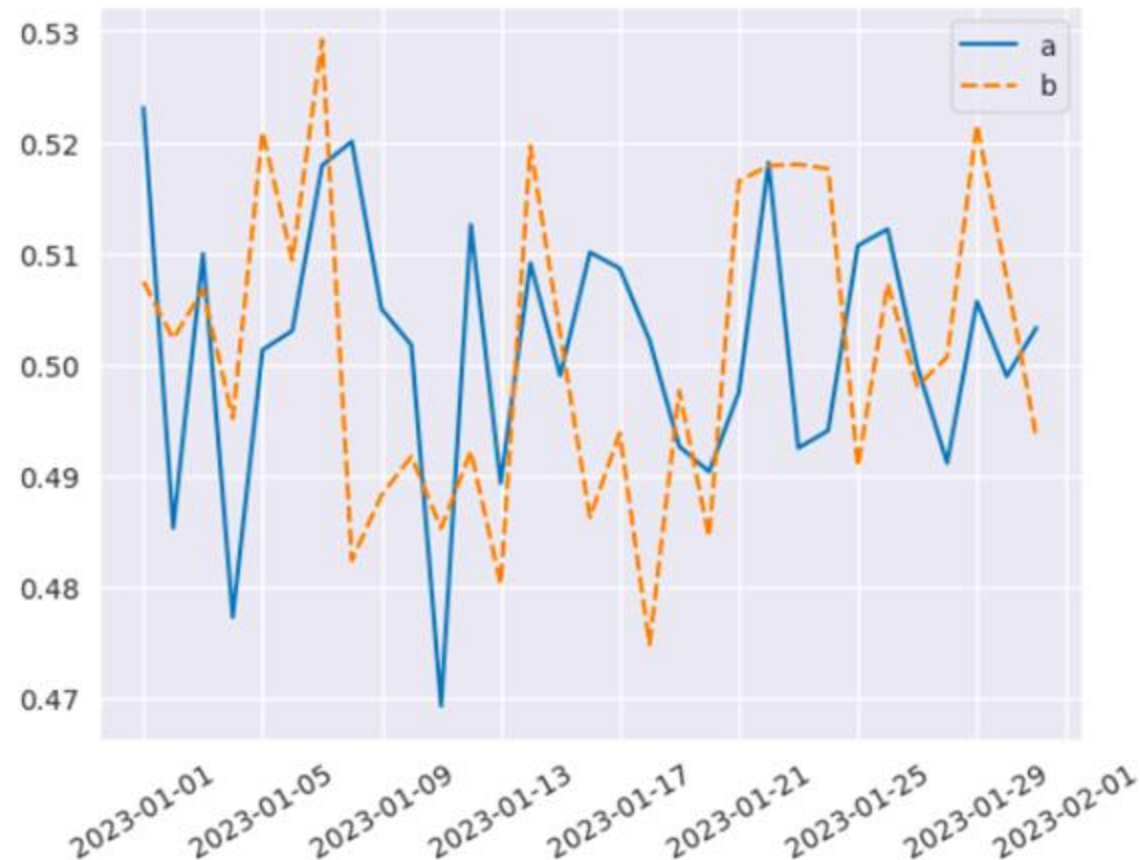
# cudf.pandas in Action

## A brief example

```
%%load_ext cudf.pandas

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

rng = pd.date_range("2023-01-01", "2023-02-01",  freq="1T")
df = pd.DataFrame({
    "a": np.random.rand(len(rng)),
    "b": np.random.rand(len(rng))
    },
    index=rng
)

df = df.iloc[rng.indexer_between_time("09:30", "16:00")]
results = df.groupby(pd.Grouper(freq="1D")).mean()

_ = sns.lineplot(results)
_ = plt.xticks(rotation=30)
```

This part runs entirely on the **GPU**
cuDF supports all these operations

NVIDIA
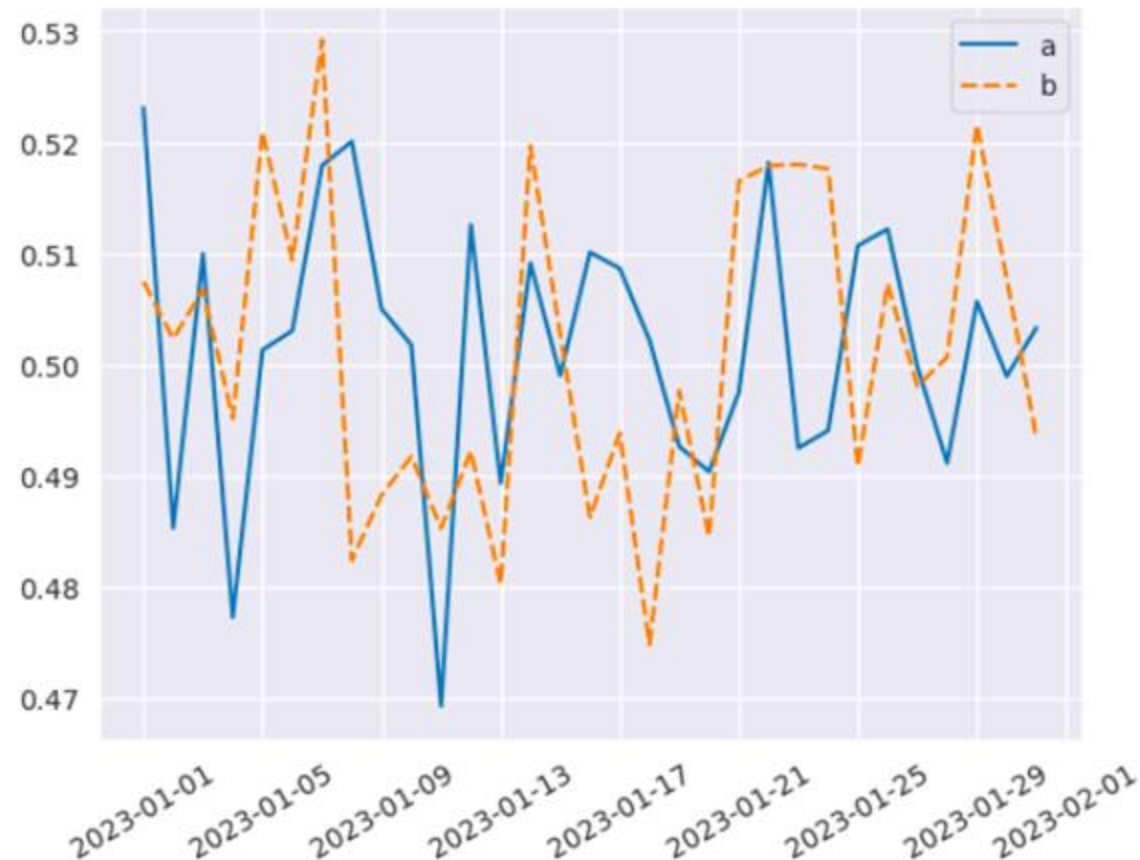
# `cudf.pandas` in Action

## A brief example

```python
%%load_ext cudf.pandas

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

rng = pd.date_range("2023-01-01", "2023-02-01",  freq="1T")
df = pd.DataFrame({
    "a": np.random.rand(len(rng)),
    "b": np.random.rand(len(rng))
    },
    index=rng
)

df = df.iloc[rng.indexer_between_time("09:30", "16:00")]
results = df.groupby(pd.Grouper(freq="1D")).mean()

_ = sns.lineplot(results)
_ = plt.xticks(rotation=30)
```

`indexer_between_time` isn't supported on
the **GPU** — so it runs on the **CPU**



NVIDIA.

# `cudf.pandas` in Action

## A brief example
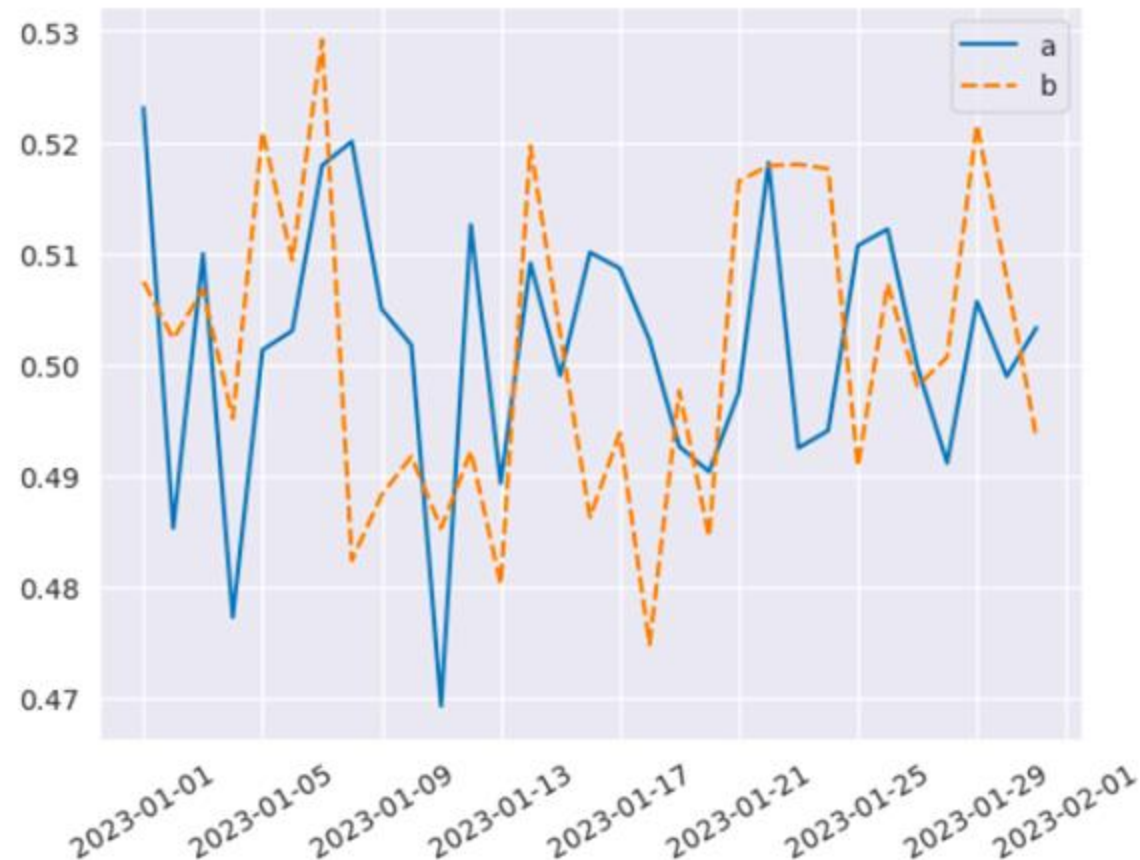
```python
%%load_ext cudf.pandas

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

rng = pd.date_range("2023-01-01", "2023-02-01",  freq="1T")
df = pd.DataFrame({
    "a": np.random.rand(len(rng)),
    "b": np.random.rand(len(rng))
    },
    index=rng
)

df = df.iloc[rng.indexer_between_time("09:30", "16:00")]
results = df.groupby(pd.Grouper(freq="1D")).mean()

_ = sns.lineplot(results)
_ = plt.xticks(rotation=30)
```



But this part happens on the **GPU**. The result of
`indexer_between_time` is copied back from CPU to GPU
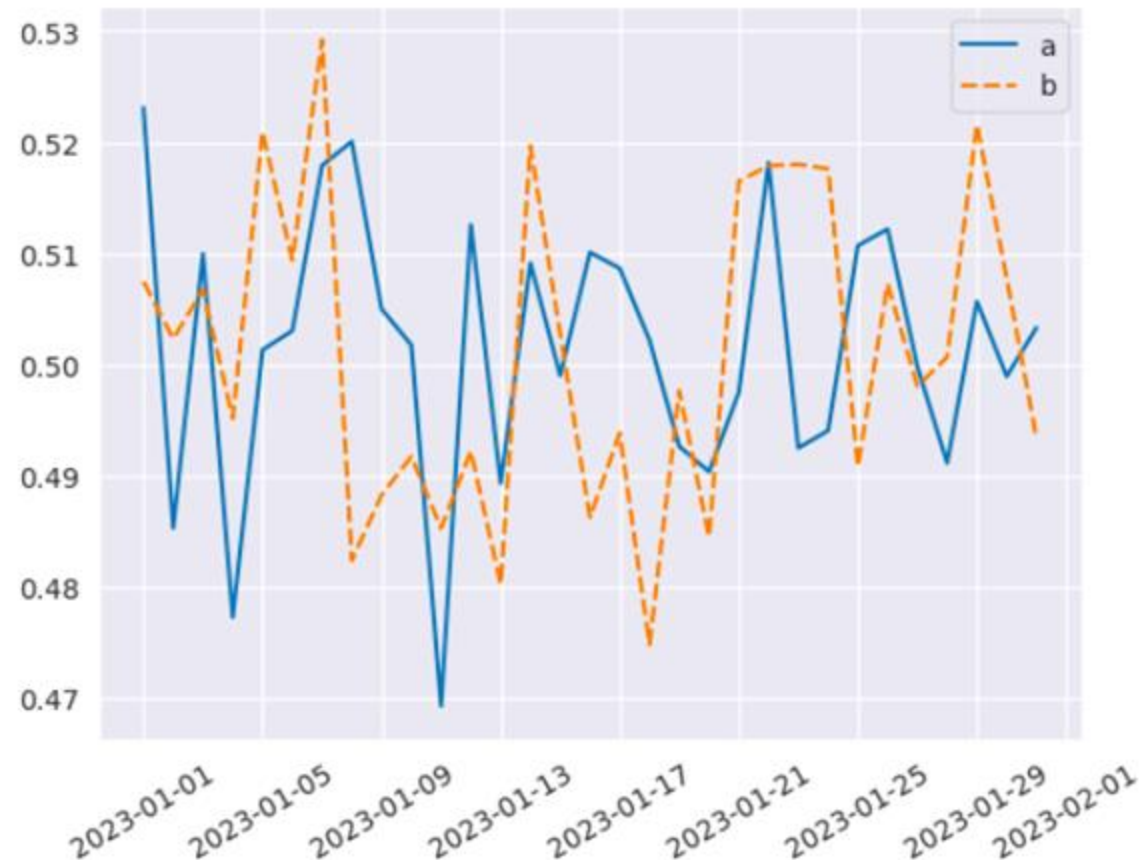
# cudf.pandas in Action

## A brief example

```
%%load_ext cudf.pandas

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

rng = pd.date_range("2023-01-01", "2023-02-01",  freq="1T")
df = pd.DataFrame({
    "a": np.random.rand(len(rng)),
    "b": np.random.rand(len(rng))
    },
    index=rng
)

df = df.iloc[rng.indexer_between_time("09:30", "16:00")]
results = df.groupby(pd.Grouper(freq="1D")).mean()

_ = sns.lineplot(results)
_ = plt.xticks(rotation=30)
```

This part runs entirely on the **GPU**

**NVIDIA**

# cudf.pandas in Action

## A brief example

```python
%%load_ext cudf.pandas

import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns

rng = pd.date_range("2023-01-01", "2023-02-01",  freq="1T")
df = pd.DataFrame({
    "a": np.random.rand(len(rng)),
    "b": np.random.rand(len(rng))
    },
    index=rng
)

df = df.iloc[rng.indexer_between_time("09:30", "16:00")]
results = df.groupby(pd.Grouper(freq="1D")).mean()

_ = sns.lineplot(results)
_ = plt.xticks(rotation=30)
```



We can seamlessly interoperate with third-party libraries like Seaborn

NVIDIA.

# cudf.pandas summary

Provides **all** of the Pandas API

Uses the GPU (via cuDF) for operations supported by cuDF

Uses the CPU (via Pandas) for operations not supported by cuDF
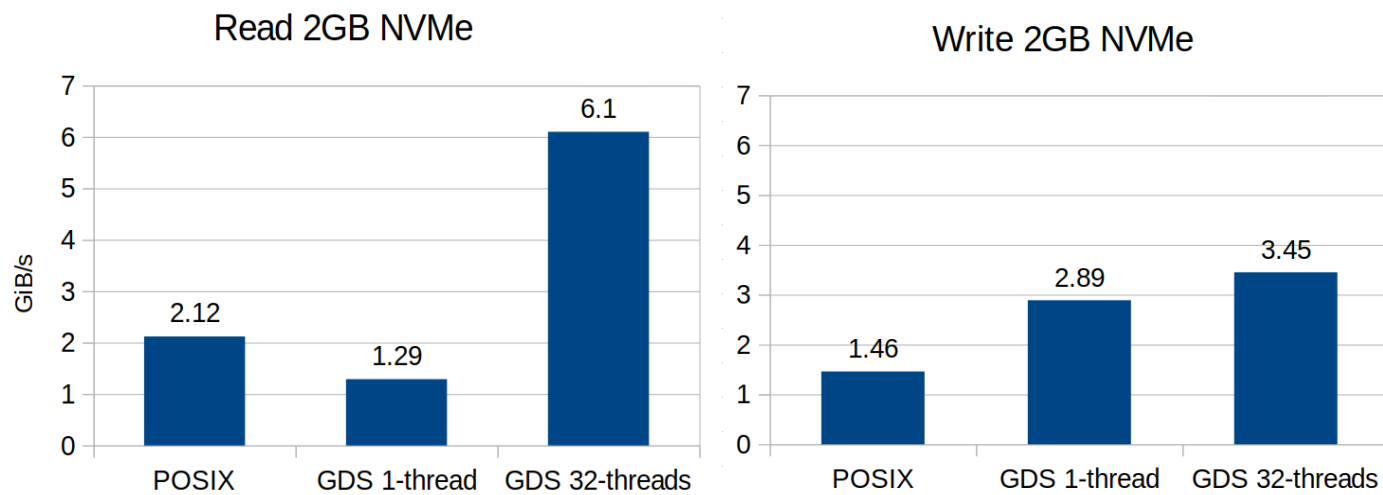
Data movement is completely hidden from the user

Zero code change: accelerates Pandas "in-place"

# RAPIDS KvikIO

**KvikIO** is a C++ and Python frontend for cuFile that provide features such as an object-oriented API, exception handling, RAII semantic, multithreading IO, fallback mode, and a Zarr backend.

Using KvikIO should feel natural to C++ and Python developers.

Comparing KvikIO's Zarr backend versus manually copying between GPU and host memory before accessing the Zarr array using POSIX

### Read 2GB NVMe

### Write 2GB NVMe

KvikIO: https://github.com/rapidsai/kvikio

```cpp
1   #include <cuda_runtime.h>
2   #include <kvikio/file_handle.hpp>
3   using namespace std;
4
5   int main() {
6     void *a = nullptr;
7     cudaMalloc(&a, 80);
8     // Read file into `a` in parallel using 16 threads
9     kvikio::default_thread_pool::reset(16);
10    {
11      kvikio::FileHandle f("/nvme/input.raw", "r");
12      future<size_t> fut = f.pread(a, sizeof(a), 0);
13      size_t read = fut.get(); // Blocking
14      // Note, `f` closes automatically on destruction.
15    }
16  }
```

```python
1   # Write CuPy array to disk
2   import cupy
3   import kvikio
4   a = cupy.arange(10)
5   with kvikio.CuFile("/nvme/input.raw", "w") as f:
6       f.write(a)
7
8   # Write same CuPy array to a Zarr store
9   import zarr
10  from kvikio.zarr import GDSStore
11  z = zarr.array(a,
12      compressor=None,
13      store=GDSStore("/nvme/store"),
14      meta_array=cupy.empty(()),
15  )
16  # We can not access the Zarr array `z` as a
17  # regular CuPy array.
18  b = z[:]  # Read from disk to GPU seamlessly
```

# RAPIDS ML and Graph Analytics

# cuML

Accelerated machine learning with a scikit-learn API
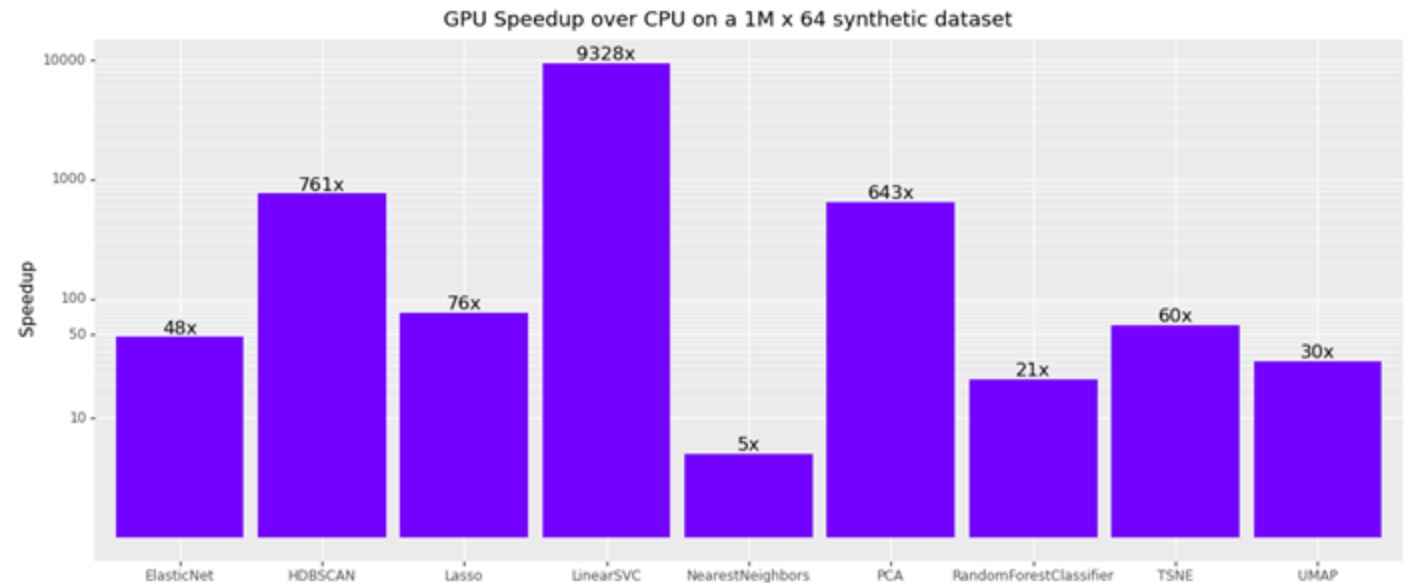
## 50+ GPU-Accelerated Algorithms

**Scikit-learn**

**CPU**

```
>>> from sklearn.ensemble import
RandomForestClassifier
>>> clf = RandomForestClassifier()
>>> clf.fit(x, y)
```

**cuML**

**GPU**

```
>>> from cuml.ensemble import
RandomForestClassifier
>>> clf = RandomForestClassifier()
>>> clf.fit(x, y)
```

GPU Speedup over CPU on a 1M x 64 synthetic dataset



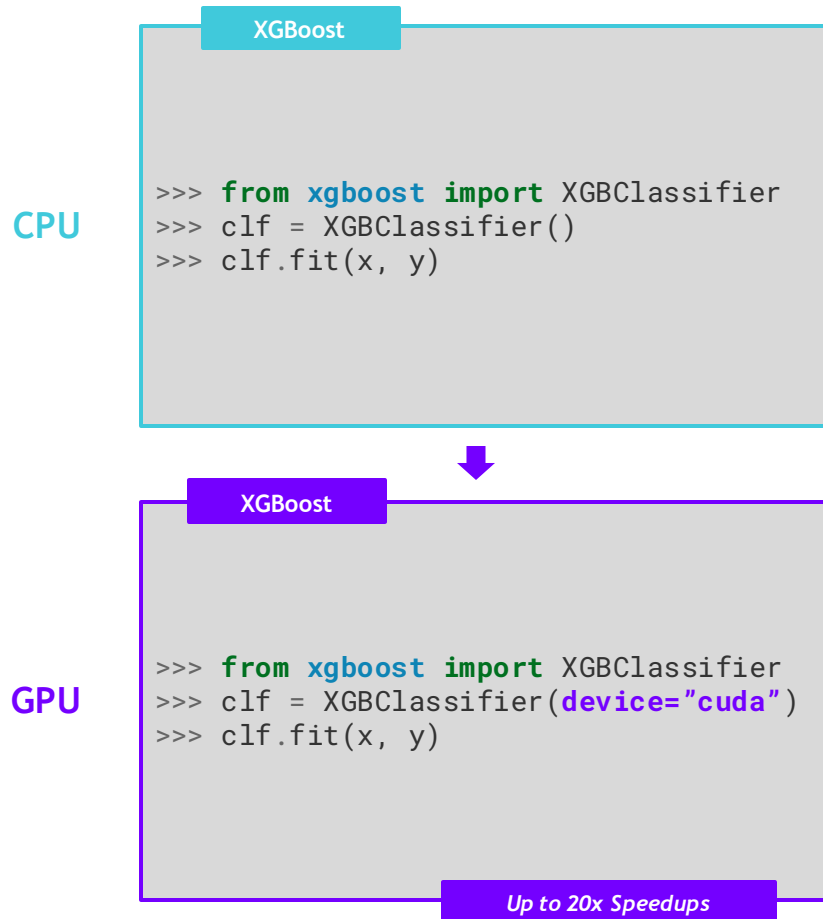| Time Series | Classification | Regression | Clustering | Preprocessing |
|---|---|---|---|---|
| Cross Validation | Tree Models | Dimensionality Reduction | Explainability | |

A100 GPU vs. AMD EPYC 7642 (96 logical cores)
cuML 23.04, scikit-learn 1.2.2, umap-learn 0.5.3

NVIDIA.

# Accelerated XGBoost

*"XGBoost is All You Need" – Bojan Tunguz, 4x Kaggle Grandmaster*

**CPU**

XGBoost

```
>>> from xgboost import XGBClassifier
>>> clf = XGBClassifier()
>>> clf.fit(x, y)
```

XGBoost

```
>>> from xgboost import XGBClassifier
>>> clf = XGBClassifier(device="cuda")
>>> clf.fit(x, y)
```

**GPU**

*Up to 20x Speedups*

- One line of code change to unlock up to 20x speedups with GPUs

- Scalable to the world's largest datasets with Dask and PySpark

- Built-in SHAP support for model explainability

- Deployable with Triton for lighting-fast inference in production
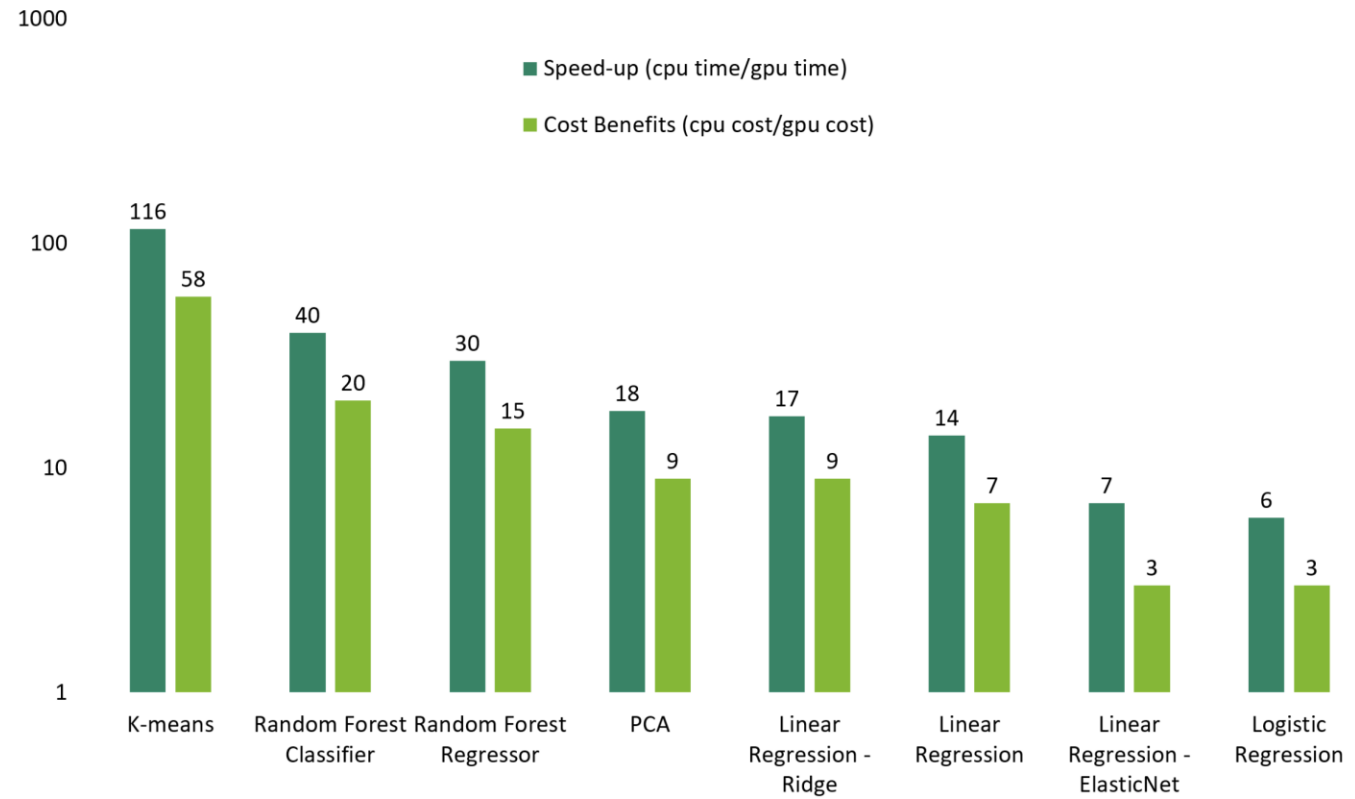
- RAPIDS helps maintain the XGBoost project

**RAPIDS** ♥ **dmlc XGBoost**

Scale up and out with RAPIDS and Dask

## RAPIDS and Others

Accelerated on single GPU

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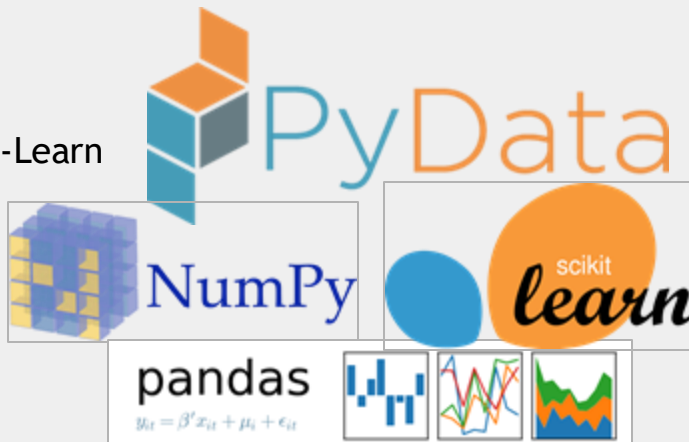Pandas -> cuDF
Scikit-Learn -> cuML
Numba -> Numba

**RAPIDS**

## Dask + RAPIDS

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

**RAPIDS**

**DASK**

## PyData

NumPy, Pandas, Scikit-Learn
and many more

Single CPU core
In-memory data

**PyData**

NumPy

scikit learn

pandas

$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$

## Dask

Multi-core and Distributed PyData

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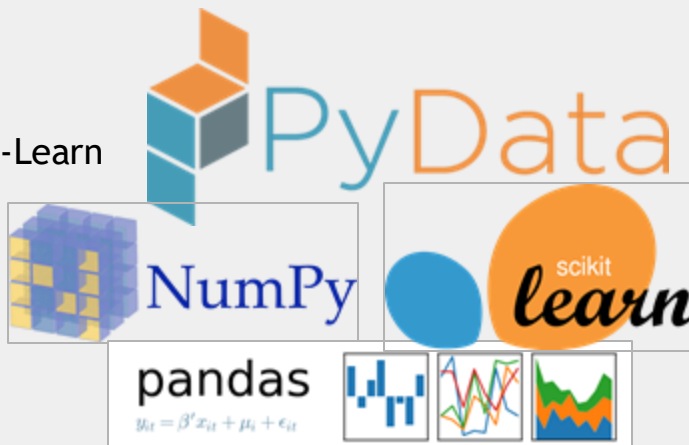... -> Dask Futures

**DASK**

**Scale Up / Accelerate**

**Scale out / Parallelize**

Scale up and out with RAPIDS and Dask



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

**Scale out / Parallelize** (horizontal axis)

**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
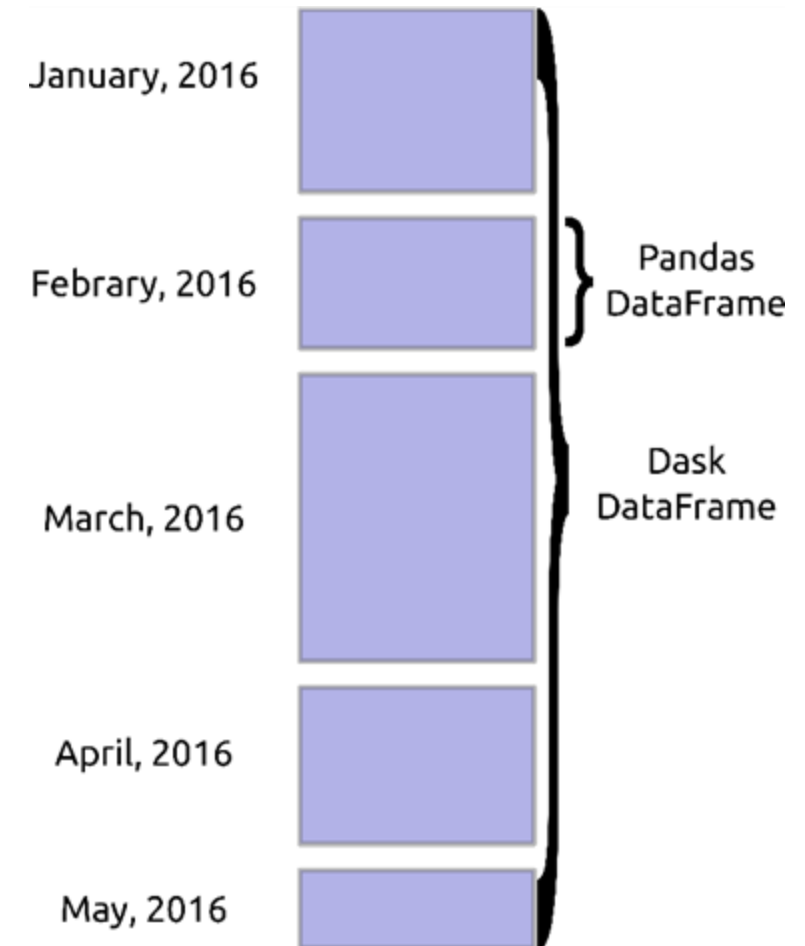Scikit-Learn -> Dask-ML
... -> Dask Futures

● **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**

January, 2016

Febrary, 2016

} Pandas
DataFrame

March, 2016

Dask
DataFrame

April, 2016

May, 2016

# Parallel Python
## For custom systems, ML algorithms, workflow engines
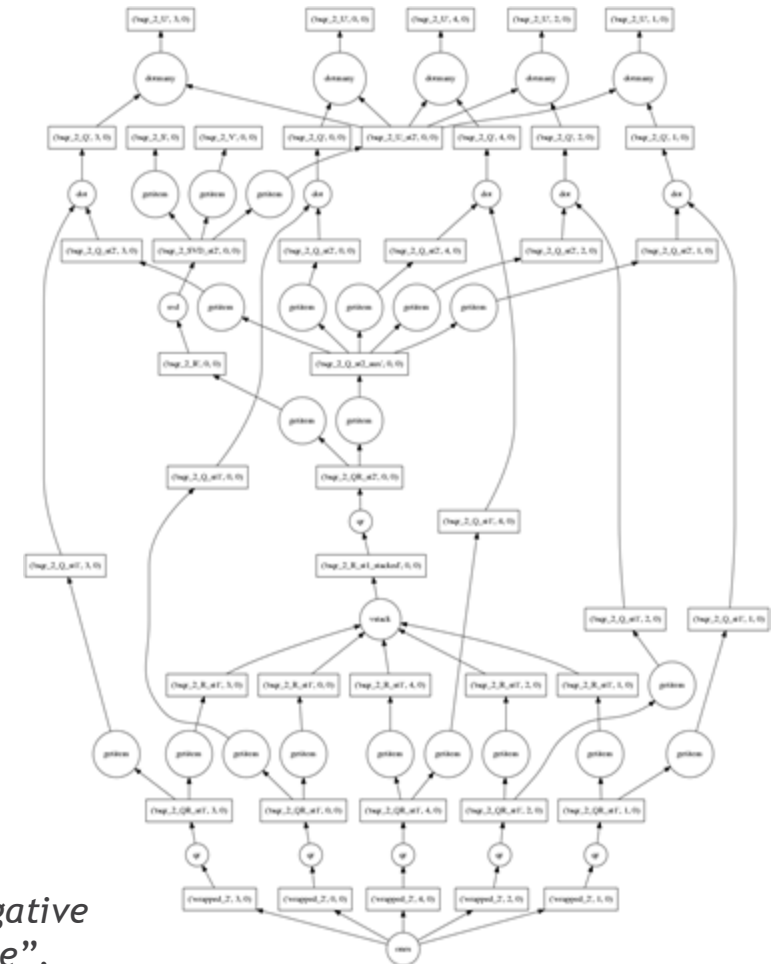
- ## Parallelize existing codebases

```
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)
```
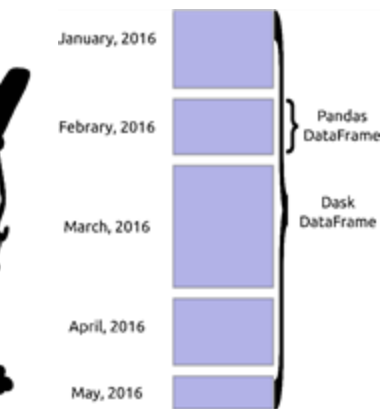
M Tepper, G Sapiro *"Compressed nonnegative matrix factorization is fast and accurate"*, IEEE Transactions on Signal Processing, 2016
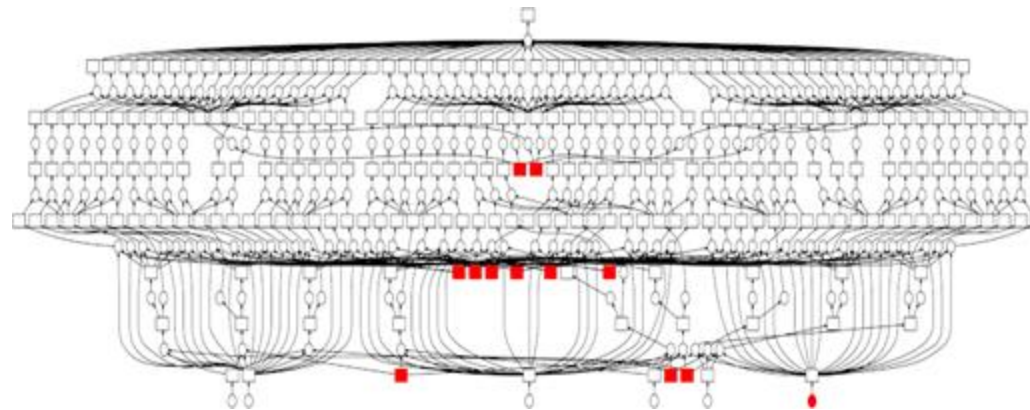
# Dask Connects Python users to Hardware



| User | 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** (vertical axis)

**Scale out / Parallelize** (horizontal axis)

**RAPIDS and Others**

Accelerated on single GPU

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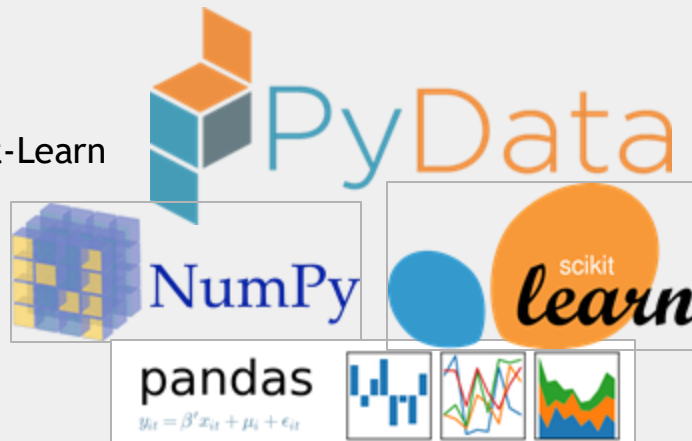Numba -> Numba

**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

# Accelerated Dask

Just set "cudf" as the backend and use Dask-CUDA Workers

- Configurable Backend and GPU-Aware Workers

- Memory Spilling (GPU->CPU->Disk)

- Optimized Memory Management

- Accelerated RDMA and Networking (UCX)

```python
import dask
from dask_cuda import LocalCUDACluster
from dask.distributed import Client
import dask.dataframe as dd

dask.config.set({"dataframe.backend": "cudf"})

cluster = LocalCUDACluster(...)
client = Client(cluster)
```

```python
from dask_cuda import LocalCUDACluster
cluster = LocalCUDACluster(...)
cluster
```

### LocalCUDACluster
382454ff

**Dashboard:** http://127.0.0.1:8787/status          **Workers:** 1

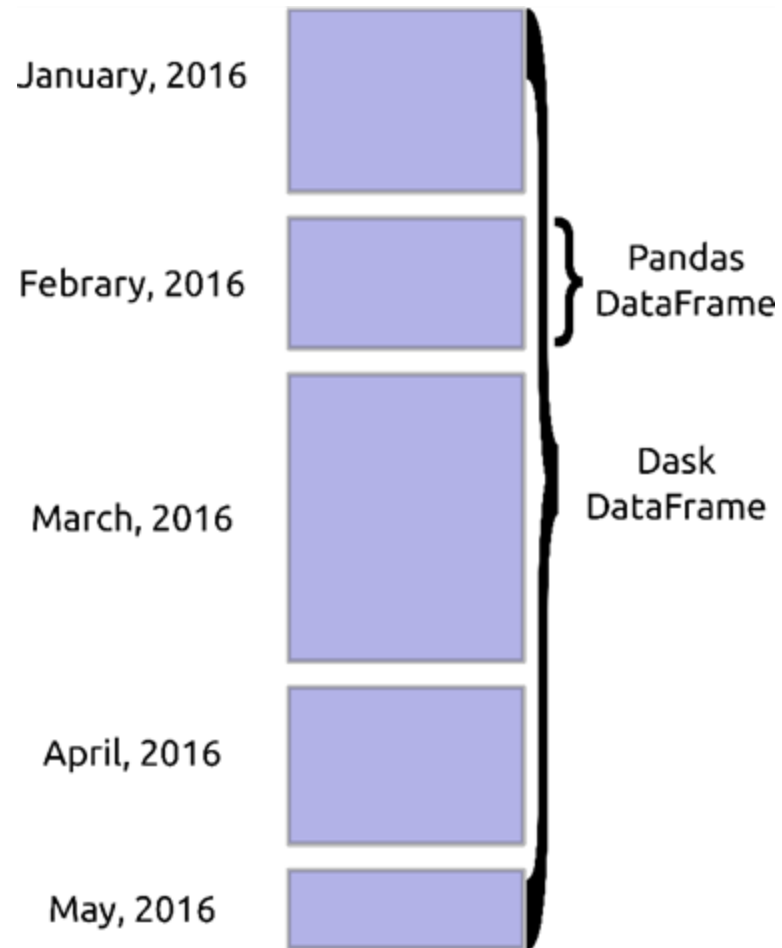**Total threads:** 1                                 **Total memory:** 0.98 TiB
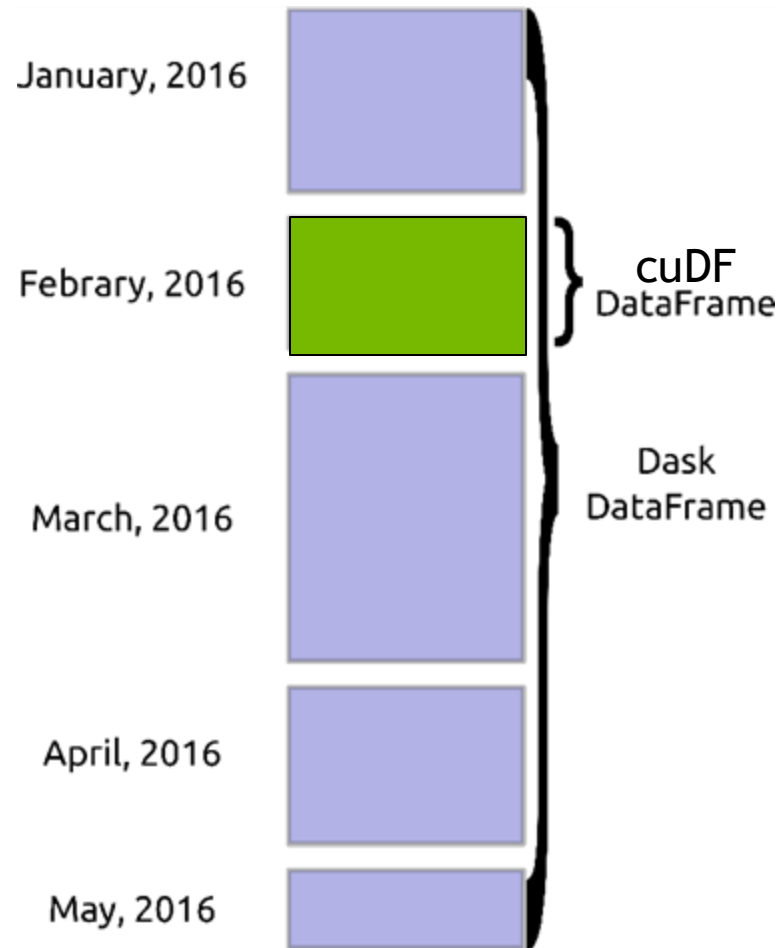
**Status:** running                                  **Using processes:** True
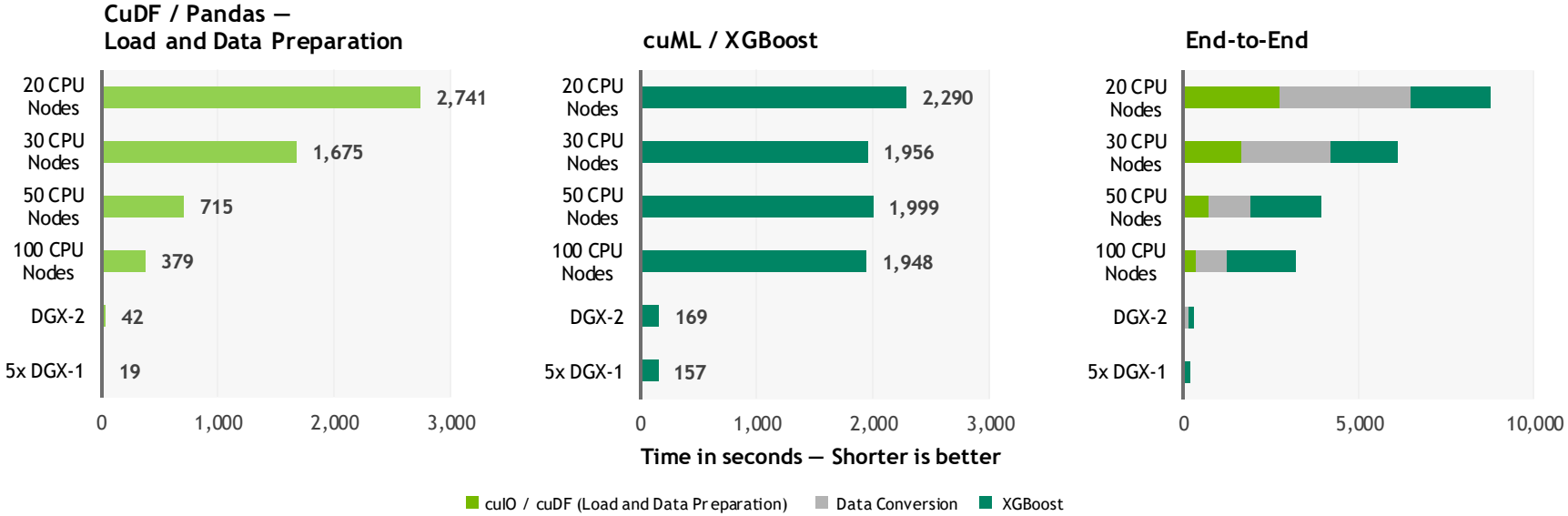
▶ Scheduler Info

## Combine Dask with cuDF
## Many GPU DataFrames form a distributed DataFrame

# Combine Dask with cuDF
## Many GPU DataFrames form a distributed DataFrame



January, 2016

Febrary, 2016 — cuDF DataFrame

March, 2016 — Dask DataFrame

April, 2016

May, 2016

# End-to-End Benchmarks

### CuDF / Pandas —
### Load and Data Preparation

| | |
|---|---|
| 20 CPU Nodes | 2,741 |
| 30 CPU Nodes | 1,675 |
| 50 CPU Nodes | 715 |
| 100 CPU Nodes | 379 |
| DGX-2 | 42 |
| 5x DGX-1 | 19 |

0    1,000    2,000    3,000

### cuML / XGBoost

| | |
|---|---|
| 20 CPU Nodes | 2,290 |
| 30 CPU Nodes | 1,956 |
| 50 CPU Nodes | 1,999 |
| 100 CPU Nodes | 1,948 |
| DGX-2 | 169 |
| 5x DGX-1 | 157 |

0    1,000    2,000    3,000

**Time in seconds — Shorter is better**

### End-to-End

| | |
|---|---|
| 20 CPU Nodes | |
| 30 CPU Nodes | |
| 50 CPU Nodes | |
| 100 CPU Nodes | |
| DGX-2 | |
| 5x DGX-1 | |

0    5,000    10,000

■ cuIO / cuDF (Load and Data Preparation)    ■ Data Conversion    ■ XGBoost

---

**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

# Accelerated Apache Spark

Zero code change acceleration for Spark DataFrames and SQL



CPU Spark

```
spark.sql("""
select
    order
    count(*) as order_count
from
    orders"""
)
```
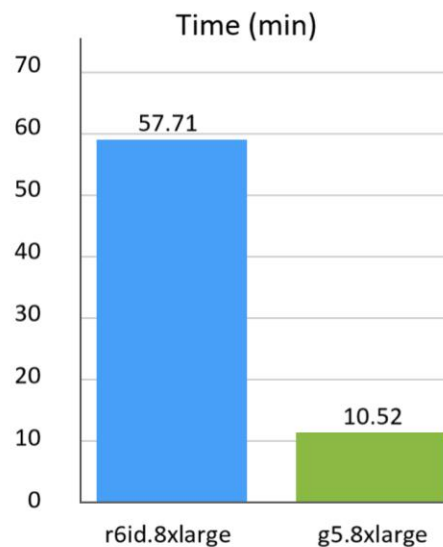
GPU Spark

```
spark.conf.set("spark.plugins",
"com.nvidia.spark.SQLPlugin")

spark.sql("""
select
    order
    count(*) as order_count
from
    orders"""
)
```

Average Speed-Ups: >5x

- RAPIDS operates as a software plugin to the popular Apache Spark platform
- Automatically accelerates supported operations (with CPU fallback if needed)
- Requires no code changes
- Works with Spark standalone, YARN clusters, Kubernetes clusters
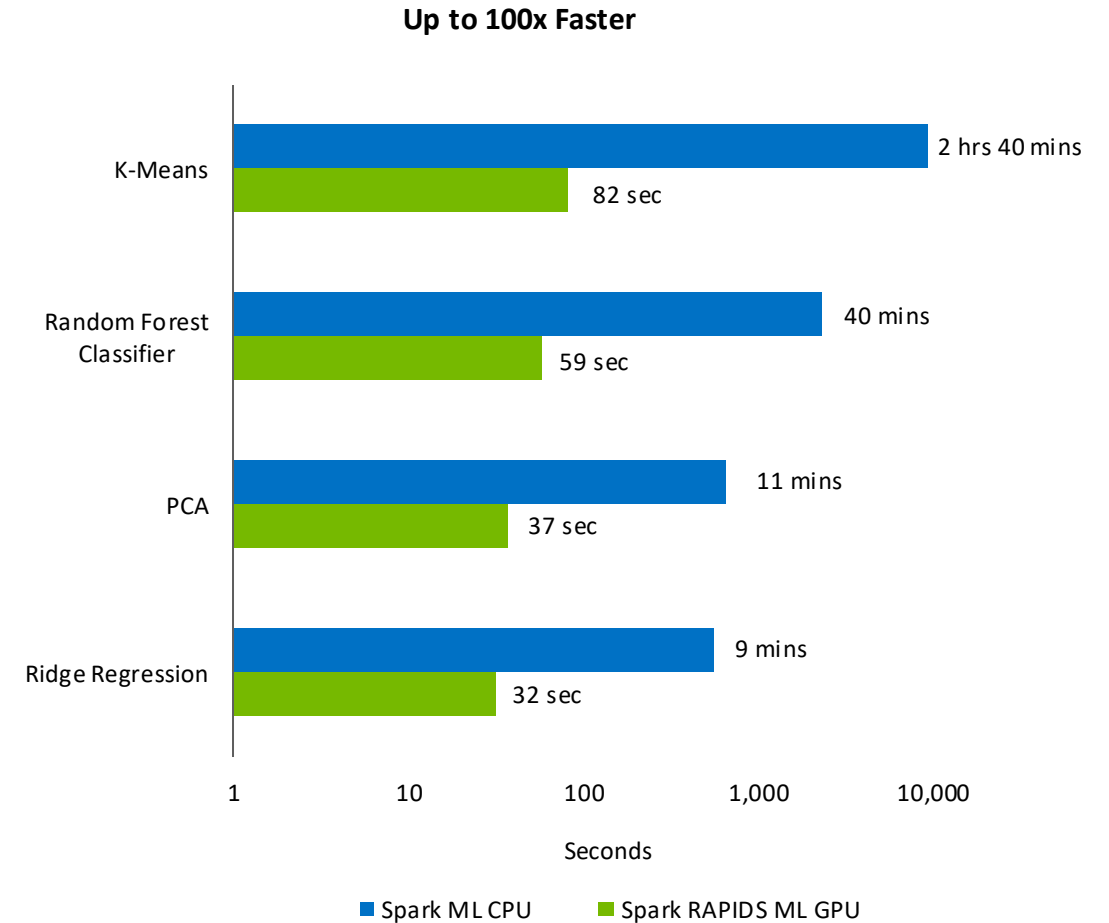
## NVIDIA Decision Support Benchmark 3TB (Public Cloud)



Time (min)

57.71    10.52

r6id.8xlarge    g5.8xlarge

**5.5x faster**

Cost

$19.06    $3.76

r6id.8xlarge    g5.8xlarge

**80% cost savings**

Apache Spark 3.4.1, RAPIDS Spark release 24.04
See GTC session S62257 for details

*Hybrid CPU/GPU libraries*

NVIDIA.

# Accelerated Apache Spark ML

Bringing GPU-accelerated machine learning to every Apache Spark user

# Accelerated Apache Spark ML

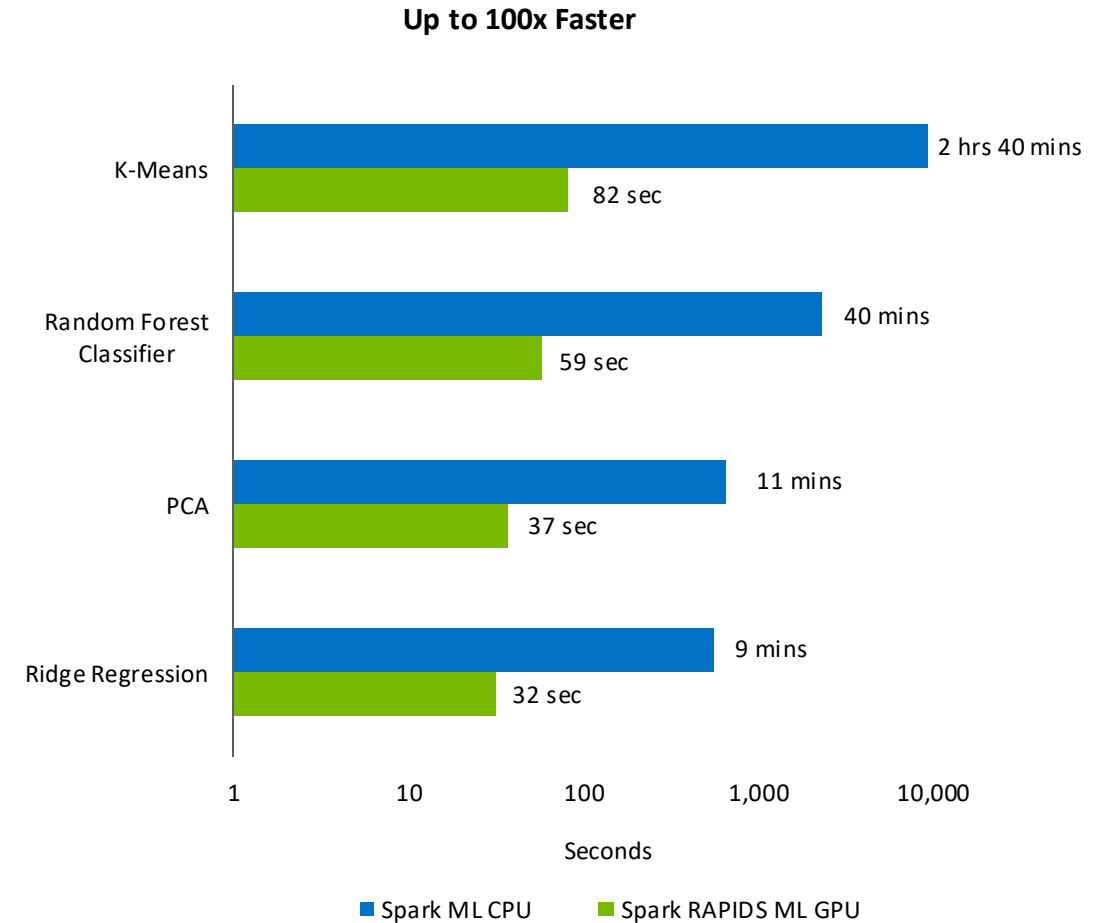## Bringing GPU-accelerated machine learning to every Apache Spark user

```python
from spark_rapids_ml.clustering import Kmeans

kmeans_estm = KMeans()\
.setK(100)\
.setFeaturesCol("features")\
.setMaxIter(30)

kmeans_model = kmeans_estm.fit(pyspark_data_frame)

kmeans_model.write().save("saved-model")

transformed = kmeans_model.transform(pyspark_data_frame)
```

**Up to 100x Faster**

| Category | Spark ML CPU | Spark RAPIDS ML GPU |
|---|---|---|
| K-Means | 2 hrs 40 mins | 82 sec |
| Random Forest Classifier | 40 mins | 59 sec |
| PCA | 11 mins | 37 sec |
| Ridge Regression | 9 mins | 32 sec |

Seconds (1, 10, 100, 1,000, 10,000)

■ Spark ML CPU    ■ Spark RAPIDS ML GPU

CPU Cluster: Intel Xeon Platinum 8000 series, 32GB RAM
GPU Cluster: NVIDIA A10 24GB, Dataset: 12GB Synthetic Dataset

Hybrid CPU/GPU libraries

NVIDIA

# Getting Started and Learning More

# Deploying RAPIDS

Documentation to get you and up and running RAPIDS anywhere

### Local Machine
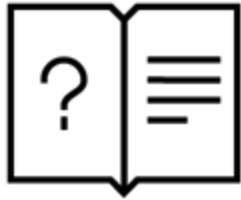Use RAPIDS on your local workstation or server.

docker   conda   pip   WSL2

### Cloud
Use RAPIDS on the cloud.

Amazon Web Services
Google Cloud Platform
Microsoft Azure   IBM Cloud

### HPC
Use RAPIDS on high performance computers and supercomputers.

SLURM

### Platforms
Use RAPIDS on compute platforms.

Kubernetes   Kubeflow   Coiled
Databricks

### Tools
There are many tools to deploy RAPIDS.

containers   dask-kubernetes
dask-operator   dask-helm-chart
dask-gateway

### Cloud ML Examples
See our example notebooks repo with opinionated deployments of RAPIDS to boost machine learning workflows.

xgboost   optuna   mlflow
ray tune

### Guides
Detailed guides on how to deploy and optimize RAPIDS.

Microsoft Azure   Infiniband   MIG

## RAPIDS Deployment Documentation

NVIDIA.
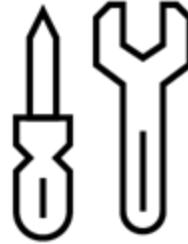
# How to Get Started with RAPIDS
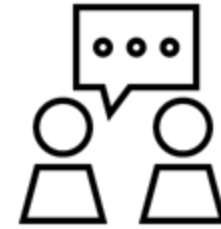
## A Variety of Ways to Get Up & Running

### More about RAPIDS

- Learn more at RAPIDS.ai
- Read the API docs
- Check out the RAPIDS blog
- Read the NVIDIA DevBlog

### Self-Start Resources

- Get started with RAPIDS
- Deploy on the Cloud today
- Start with Google Colab
- Look at the cheat sheets

### Discussion & Support

- Check the RAPIDS GitHub
- Use the NVIDIA Forums
- Reach out on Slack
- Talk to NVIDIA Services

## Get Engaged

@RAPIDSai

https://github.com/rapidsai

https://rapids.ai/slack-invite/

https://rapids.ai