



RAPIDS: Accelerated Data Science and Model Inference

Roman Enzmann, Solutions Architect at NVIDIA



Agenda

- Why RAPIDS in Data Science
- How is speedup achieved?
- Success Stories
- Data Loading & ETL (cuDF and Spark RAPIDS)
- Machine Learning (cuML)
- Graph Analytics (cuGraph)
- Triton FIL
- Fraud Detection With RAPIDS
- How to Get Started

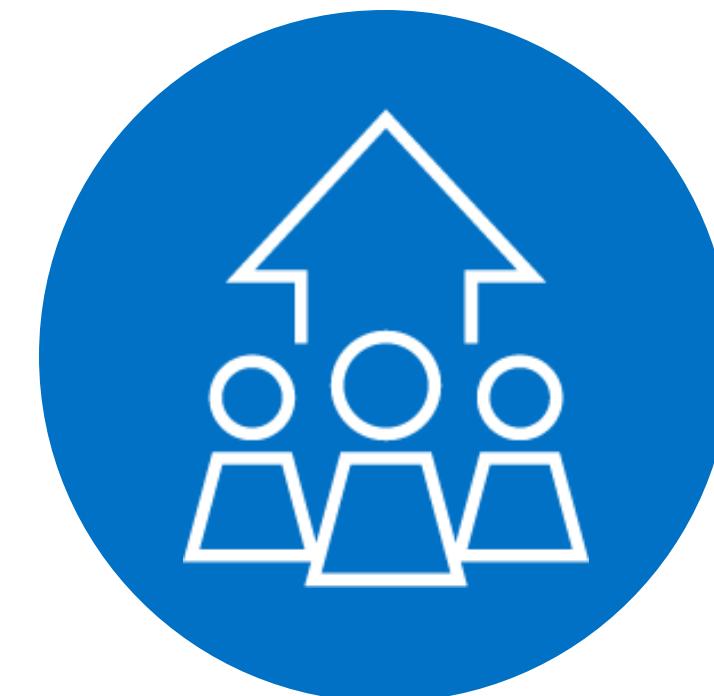
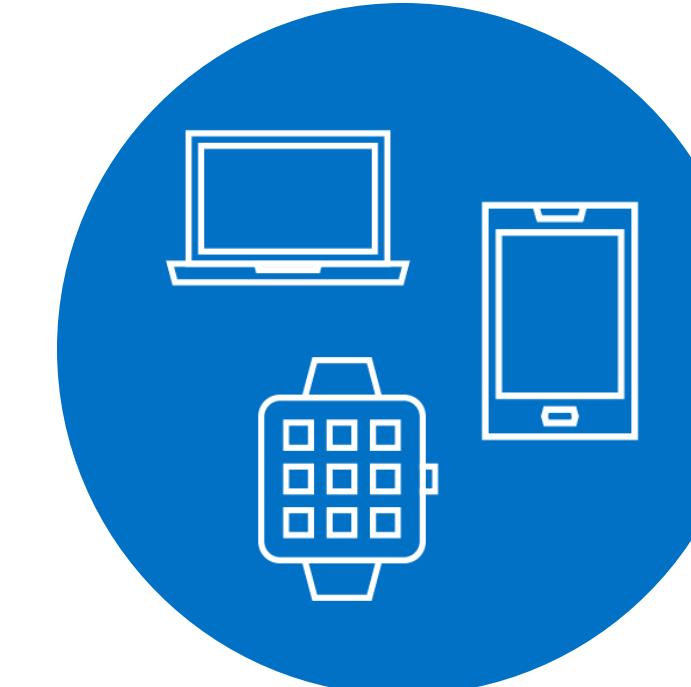
Data Center Traffic Continues to Expand

Cyberattacks Increased by 72% from 2021 to 2023¹

Explosion of Data

16B

Connected Things (IoT)²



5B+

Internet Users³



291ZB

of data generated
by 2027⁴

Consequences of Cyberattacks

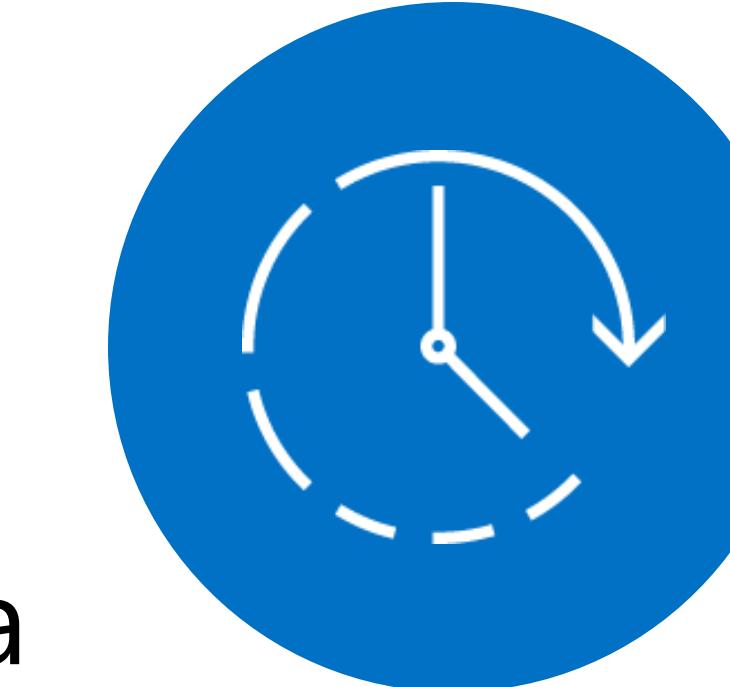
\$4.45M

Average cost of a breach⁵



277

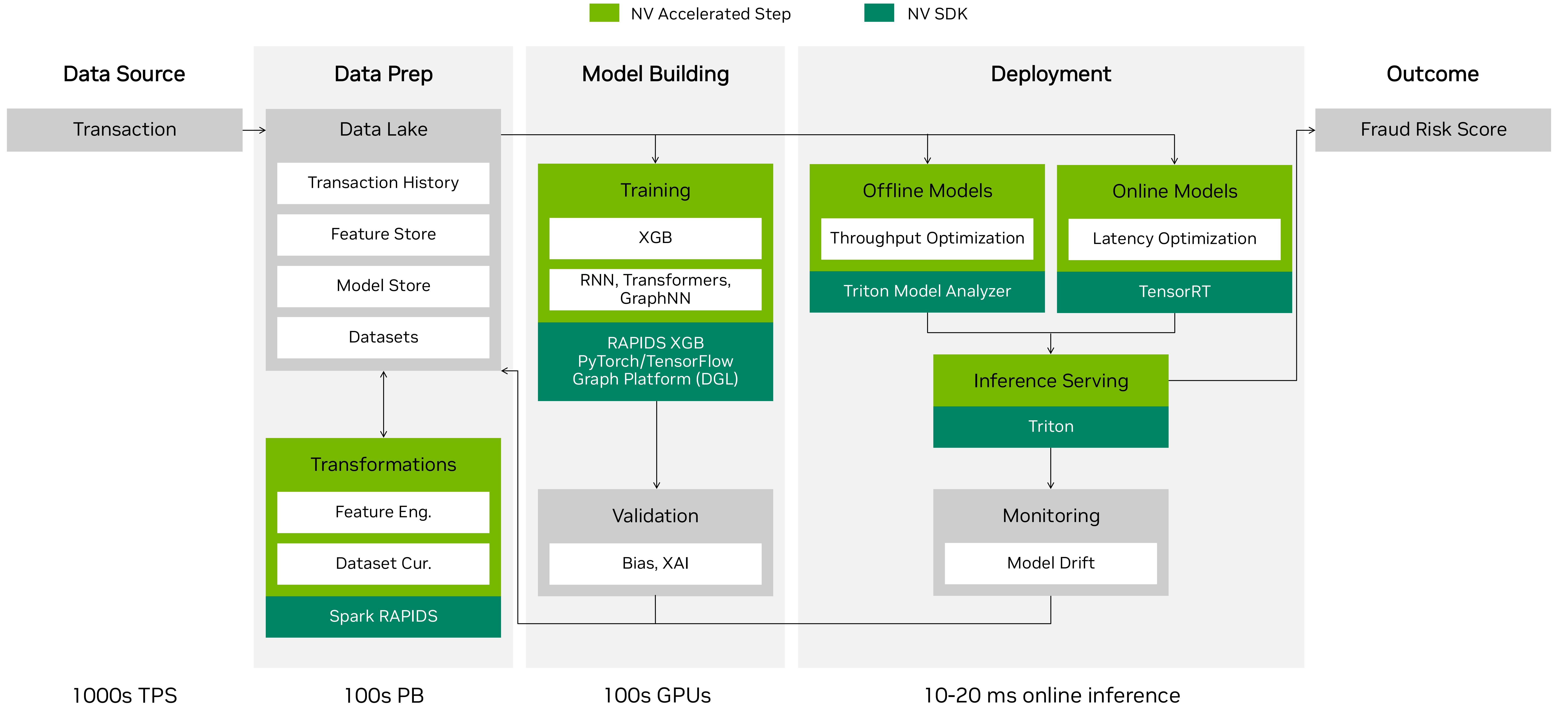
days on average to
identify and contain a
breach⁵



343M

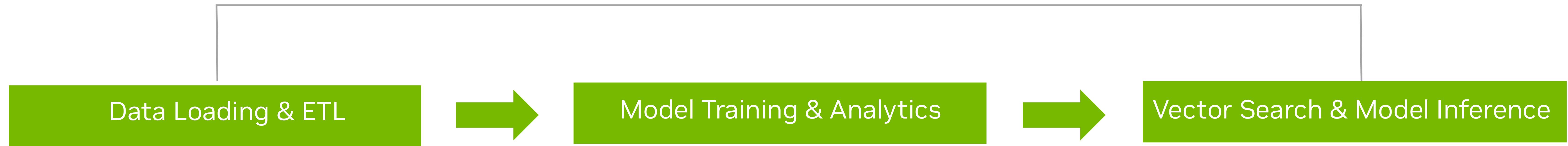
Victims of 2,365
cyberattacks in 2023⁶

Transaction Fraud Detection Workflow Example



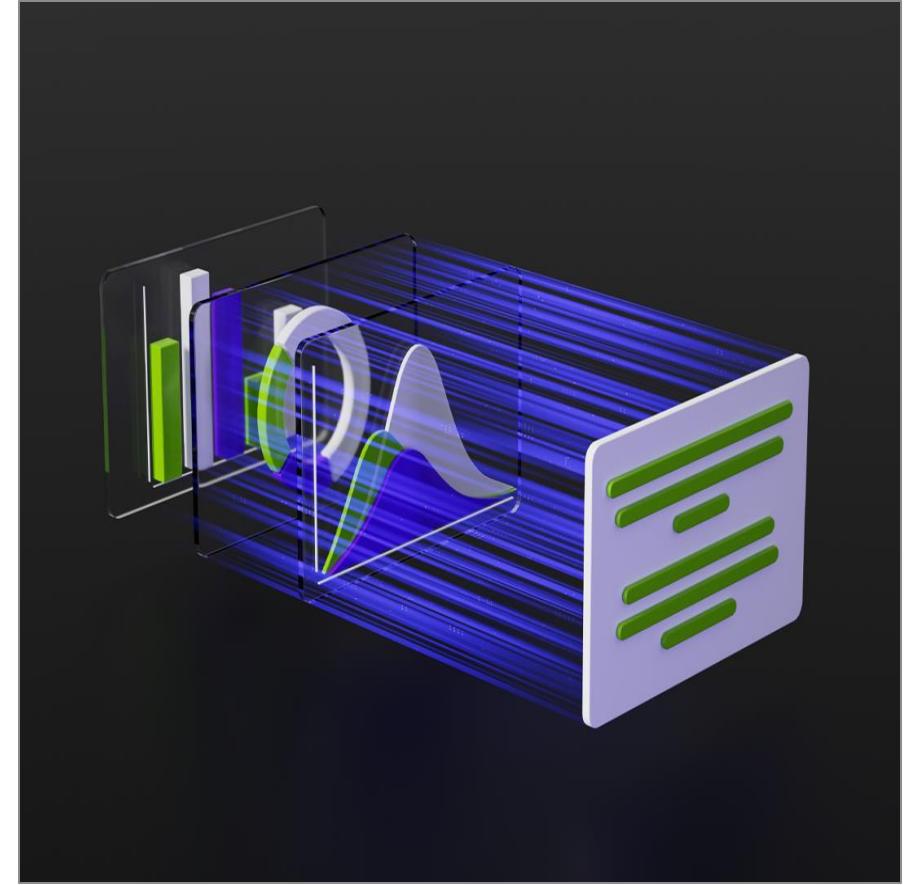
Modelling Process

Simple three step process at face value

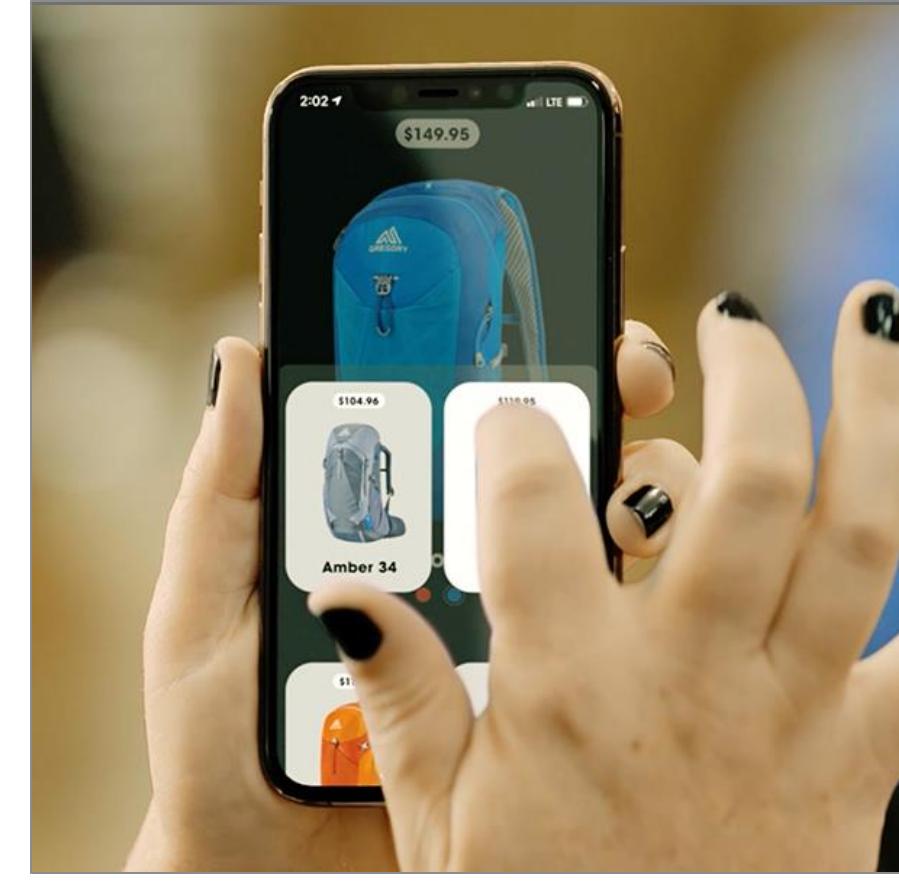


Modern Enterprise Applications Need Accelerated Computing

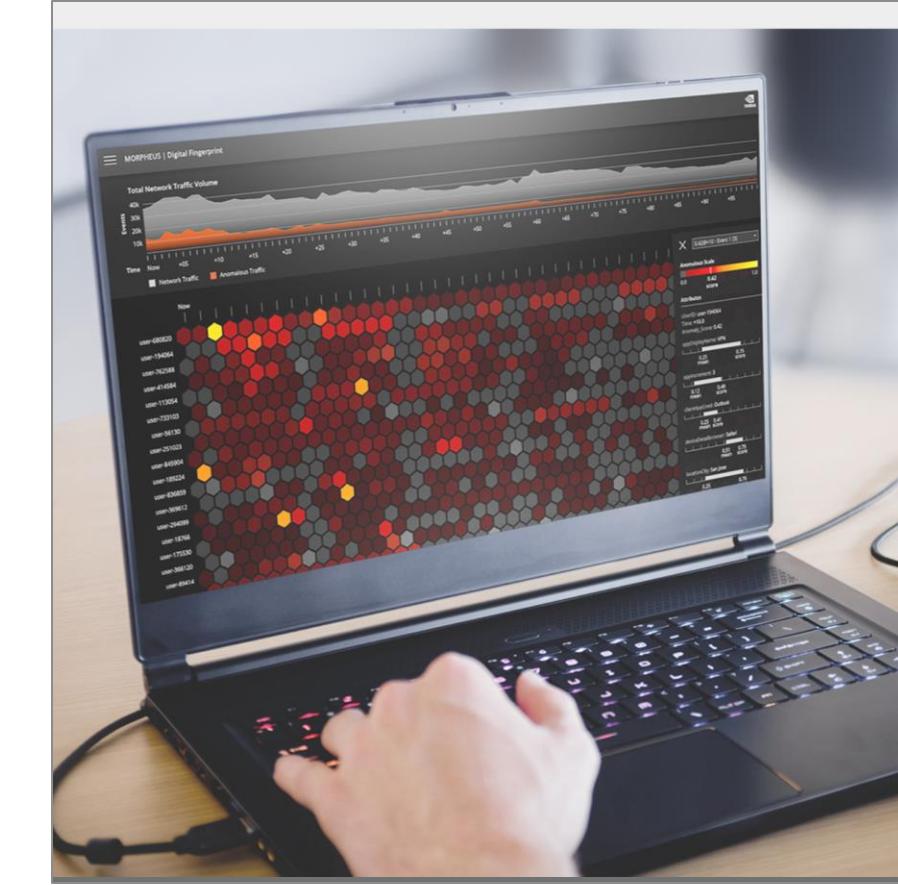
Internet scale data | Massive models | Real-time performance



LLMs



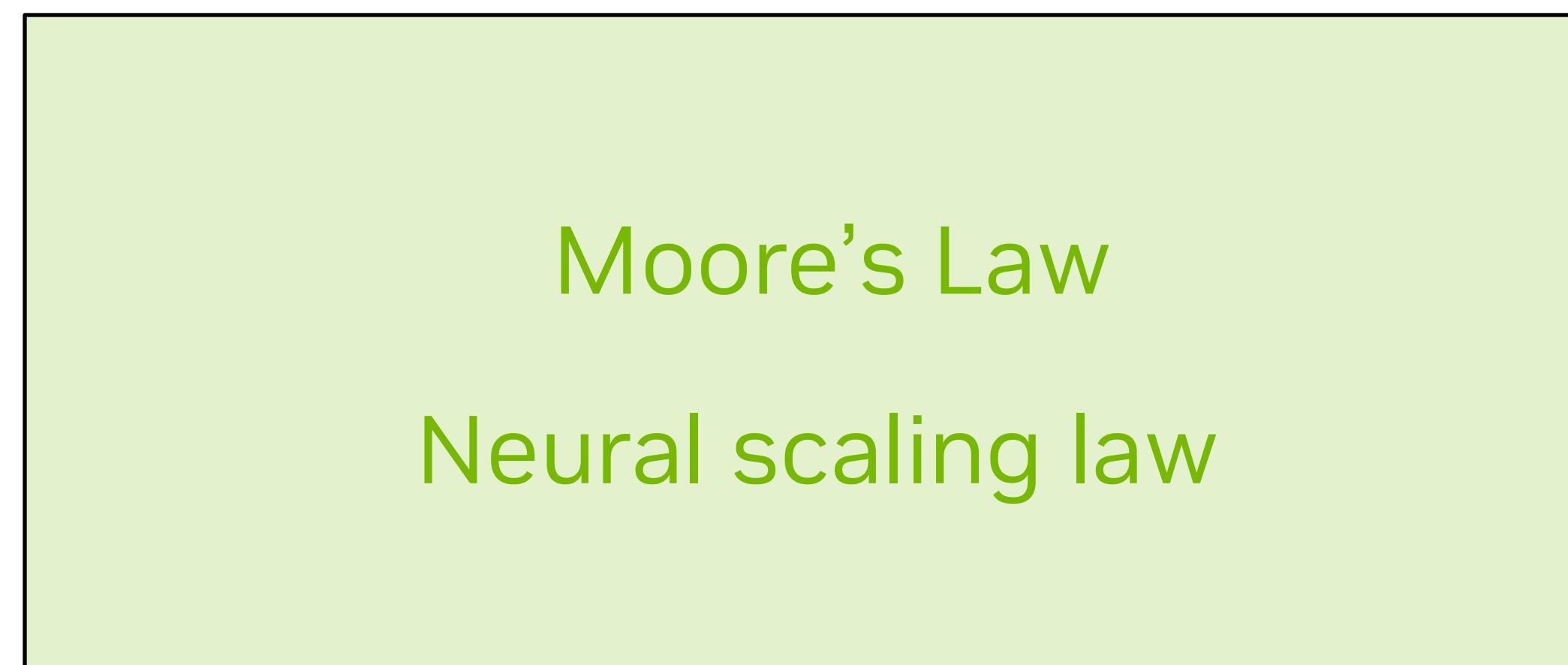
Recommenders



Cybersecurity

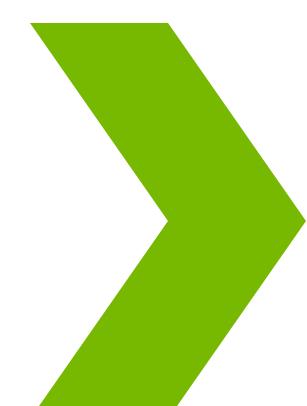


Fraud detection

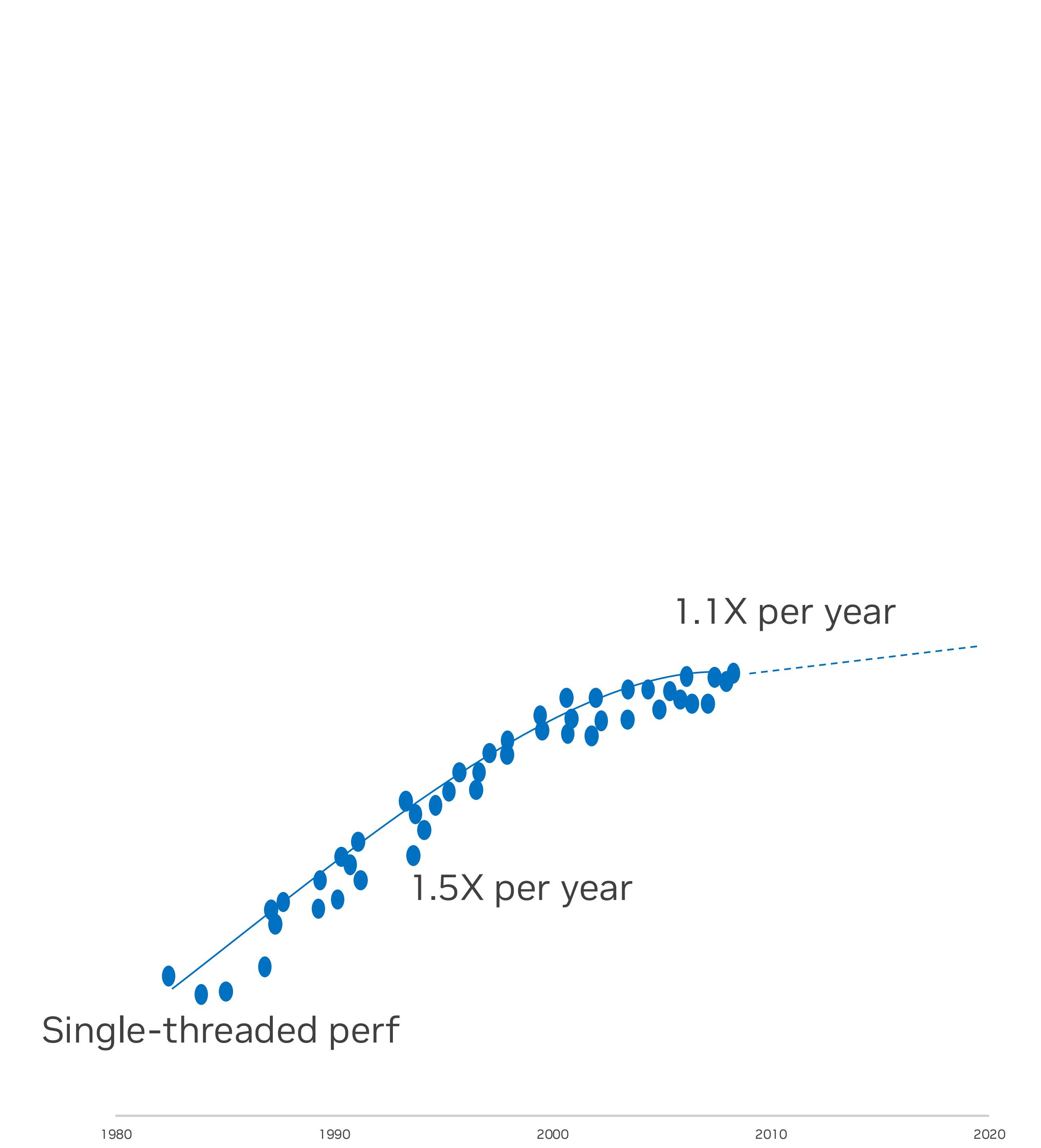


Drivers:

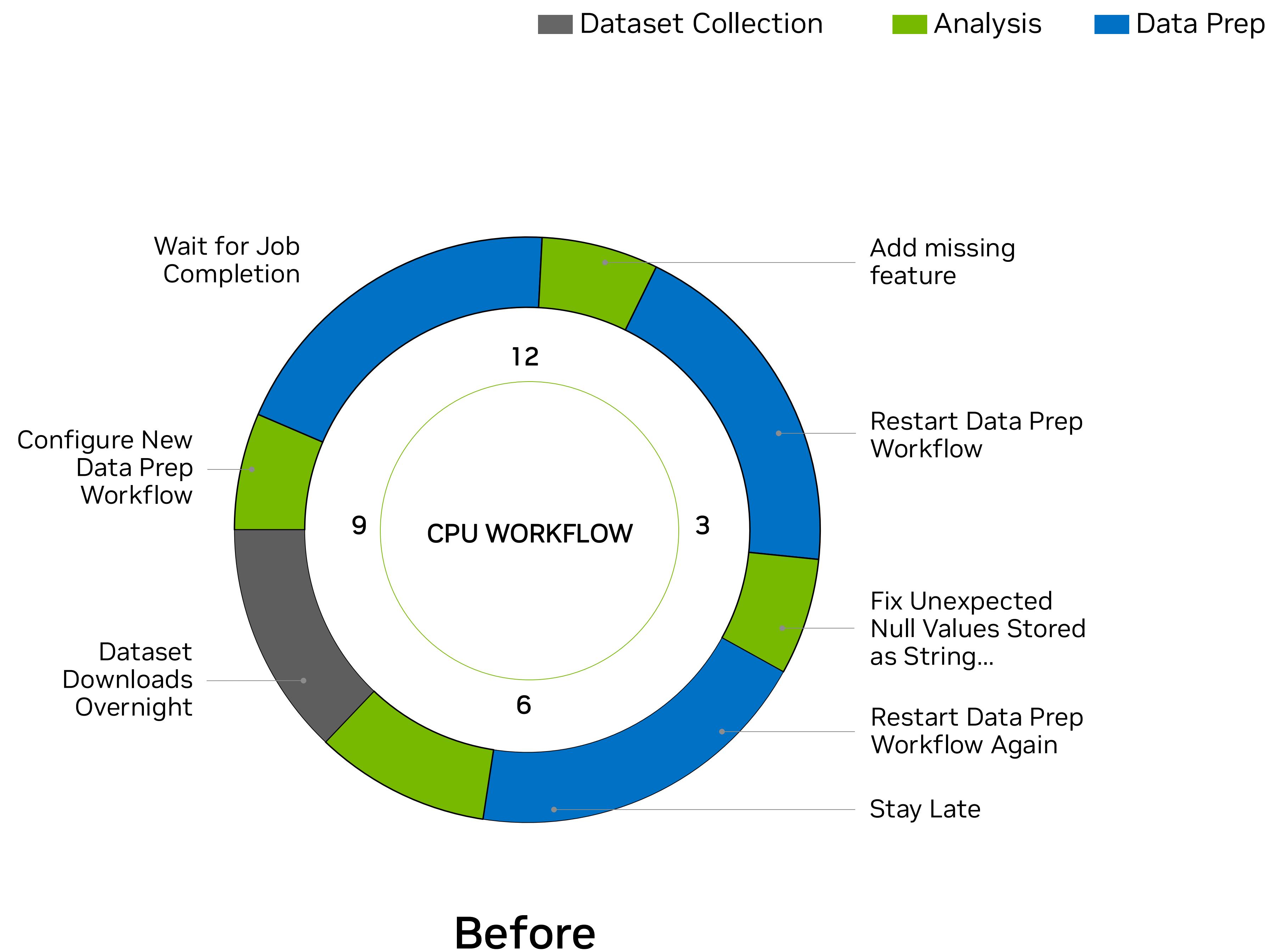
1. Data availability
2. Model complexity
3. Transaction speed



Accelerated Performance

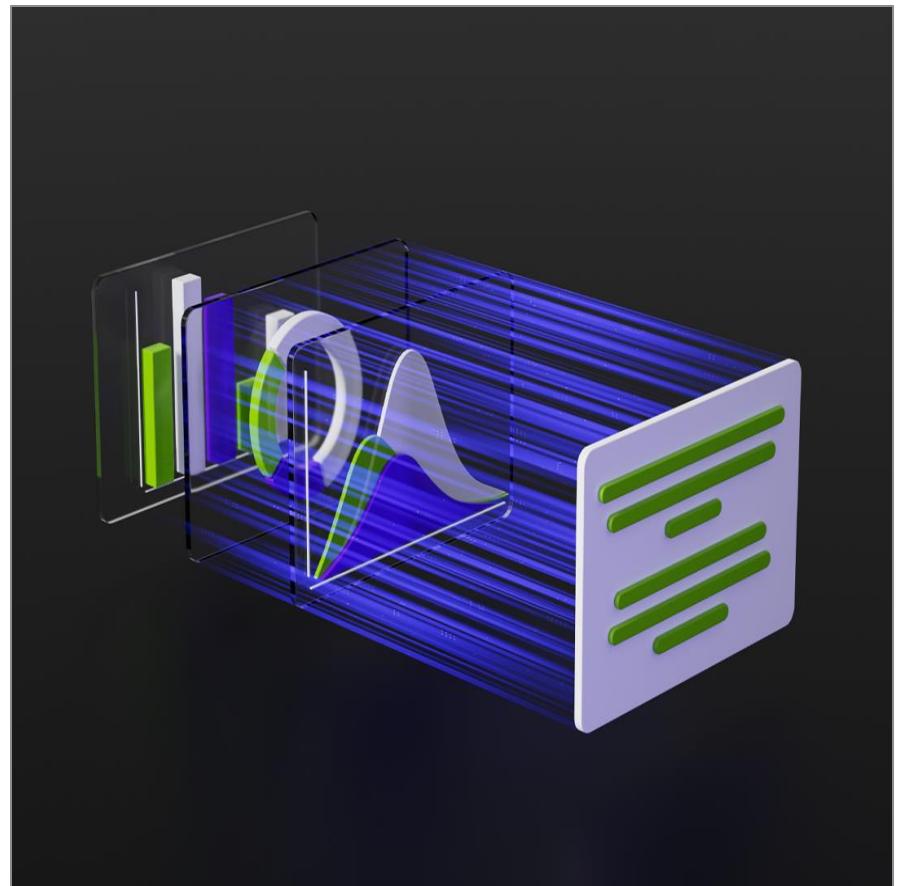


Accelerated Computing Empowers Modern Data Science Teams

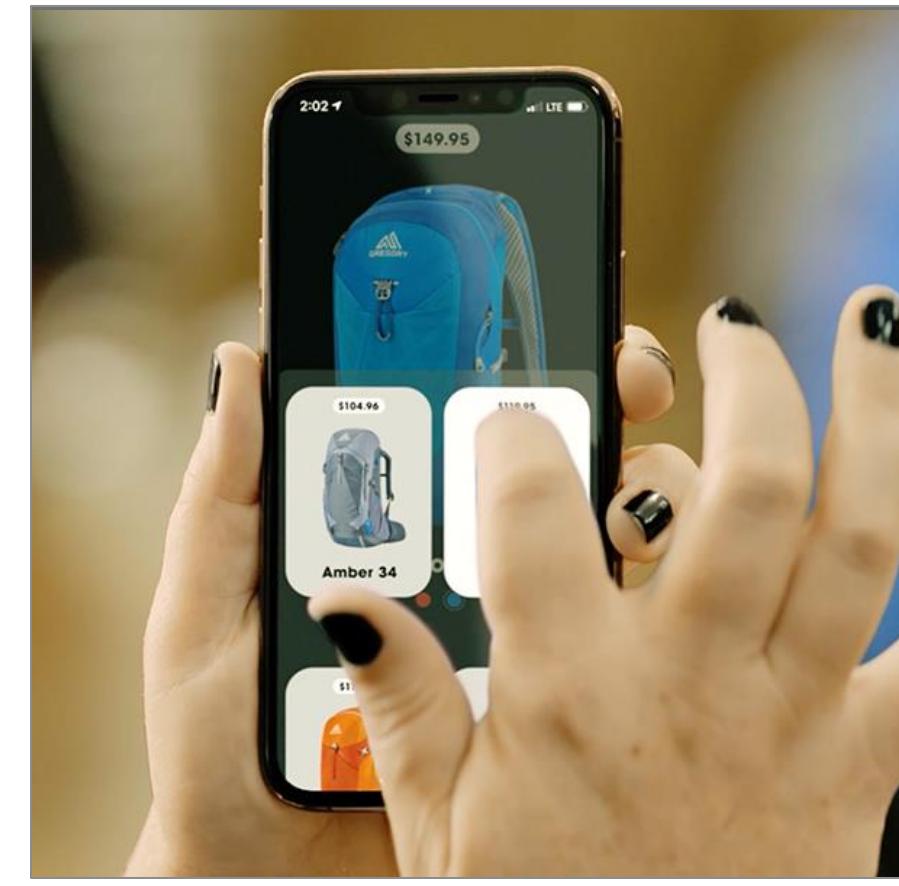


Modern Enterprise Applications Need Accelerated Computing

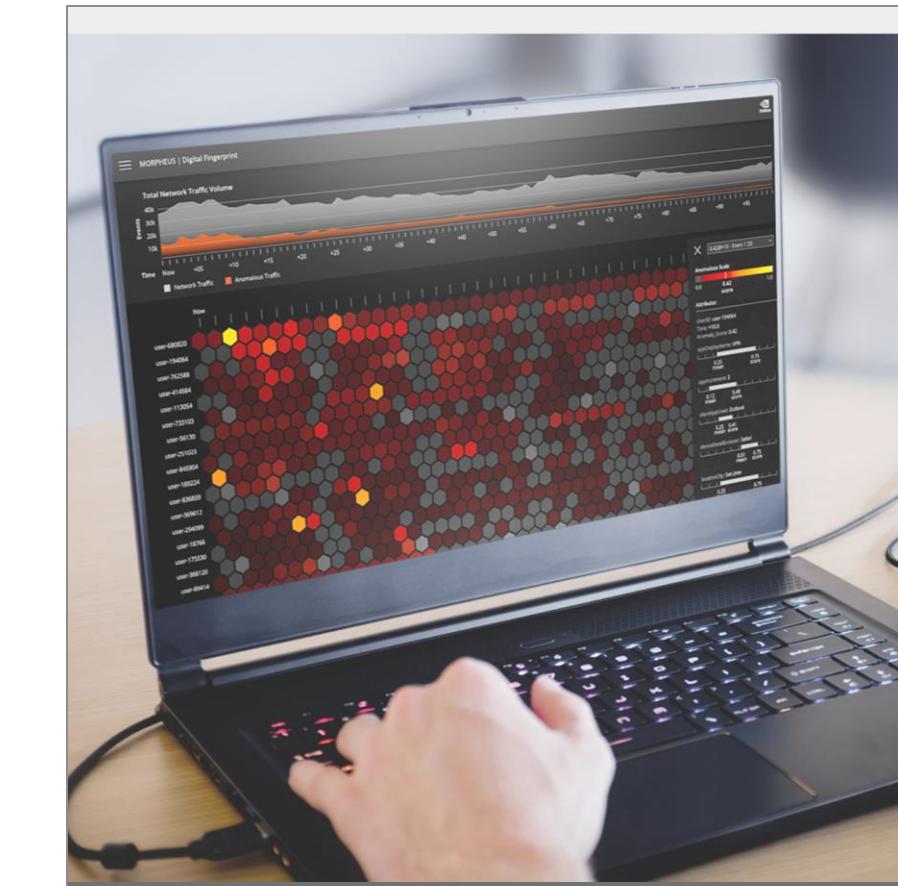
Internet scale data | Massive models | Real-time performance



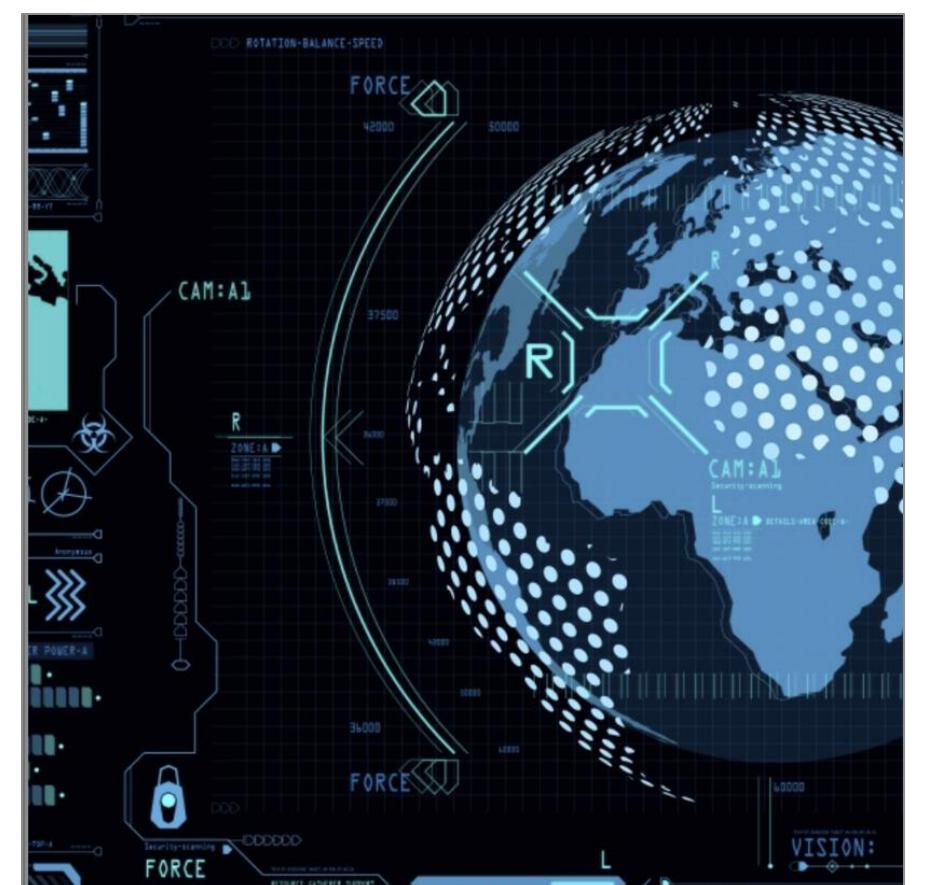
LLMs



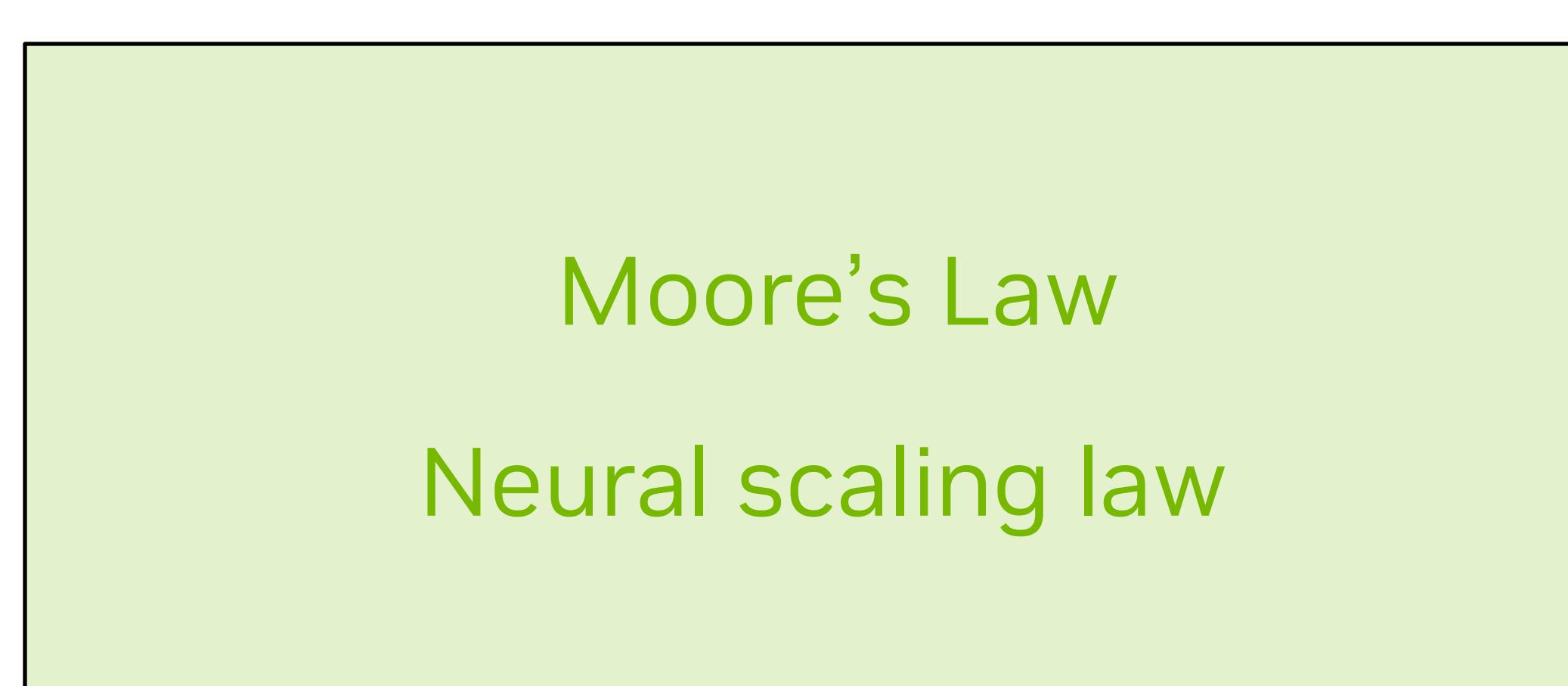
Recommenders



Cybersecurity

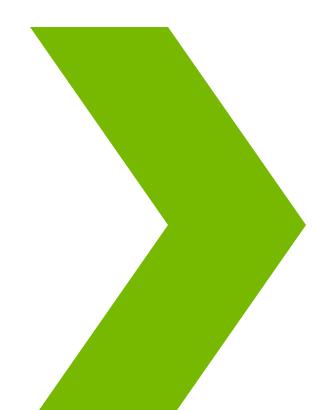


Fraud detection

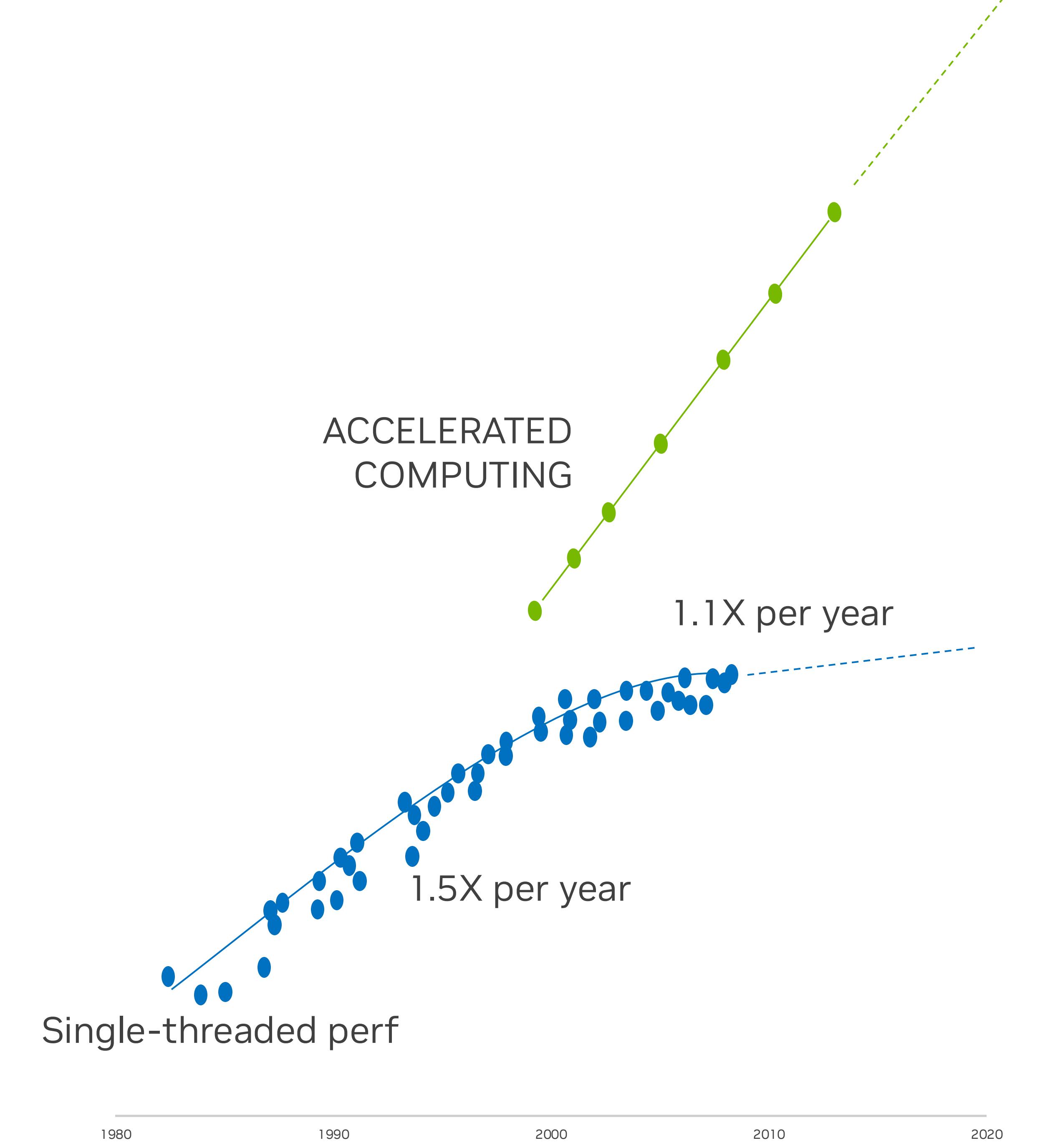


Drivers:

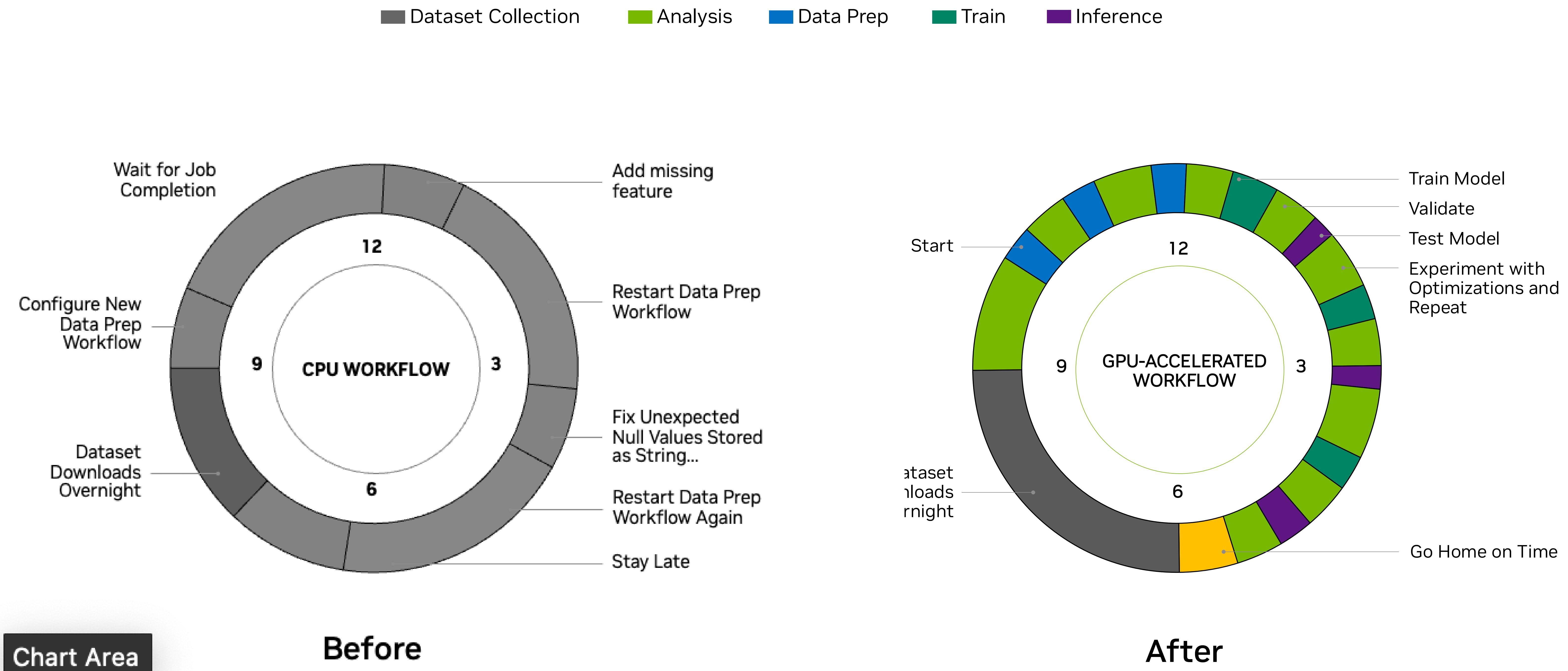
1. Data availability
2. Model complexity
3. Transaction speed



Accelerated Performance



Accelerated Computing Empowers Modern Data Science Teams



How RAPIDS Fits into Data Science

NVIDIA RAPIDS

Improve Existing Data Processing Workflows



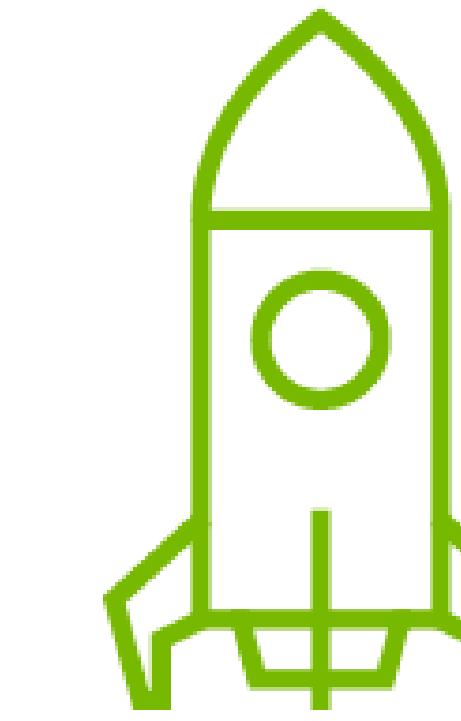
Faster Results

Accelerate data replication
and AI pipelines



Lower Costs

Save on infrastructure
Reduce power consumption and
carbon footprint



Quick Time-to-Value

No code changes required

RAPIDS Accelerates Data Science End-to-End

Data Loading & ETL

Model Training & Analytics

Vector Search & Model Inference

RAPIDS



NVIDIA AI Enterprise

Development Tools | Cloud Native Management and Orchestration | Infrastructure Optimization



Cloud



Data Center

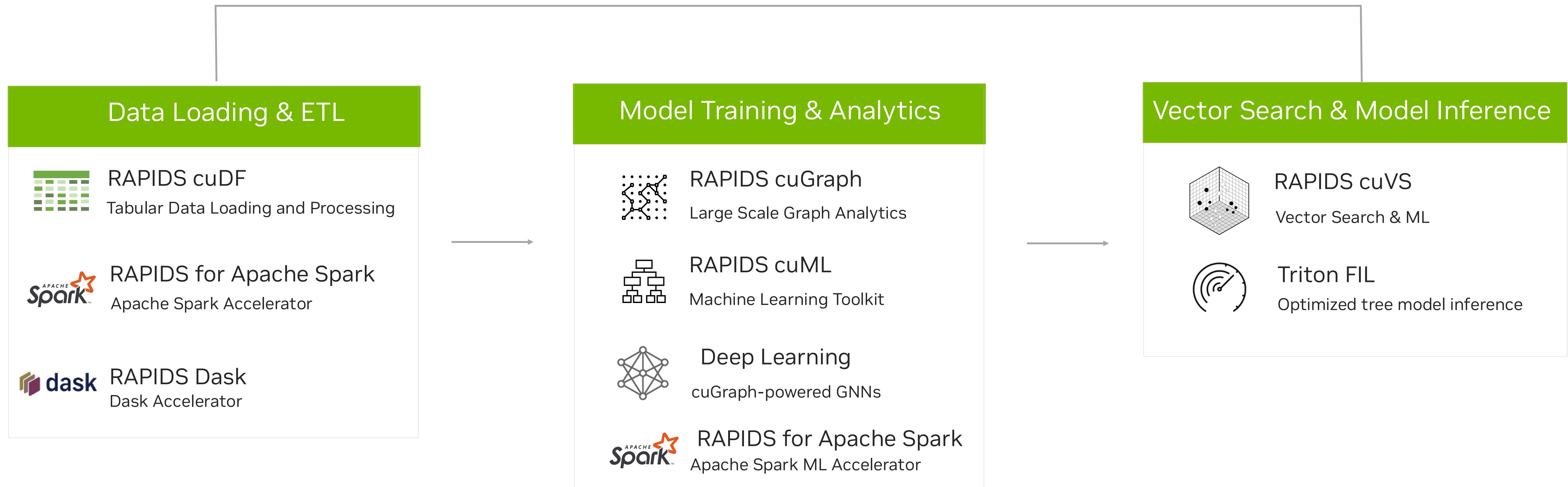


Edge



RTX Laptop

RAPIDS Accelerates Data Science End-to-End



NVIDIA AI Enterprise

Development Tools | Cloud Native Management and Orchestration | Infrastructure Optimization



Cloud



Data Center



Edge

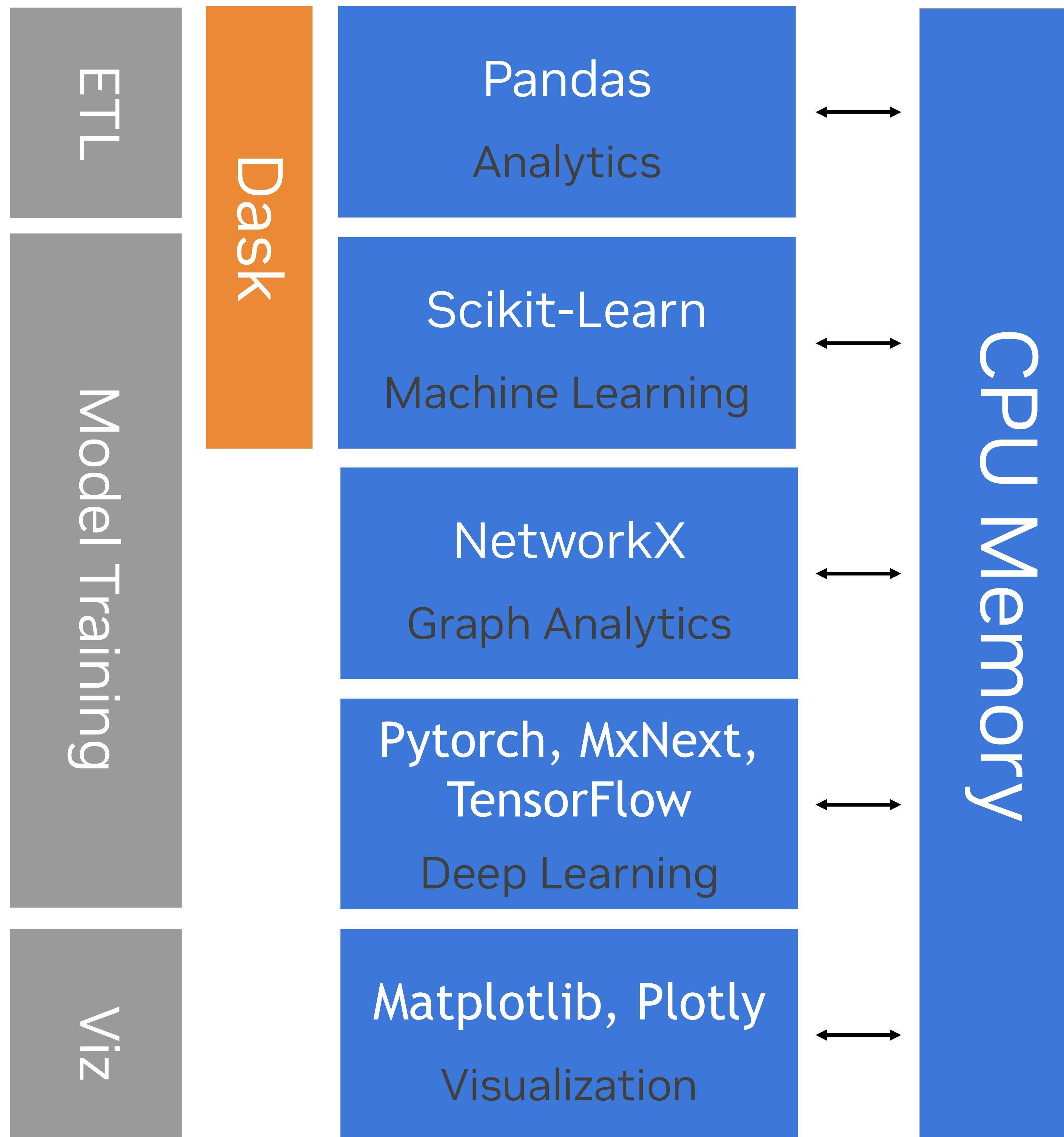


RTX Laptop

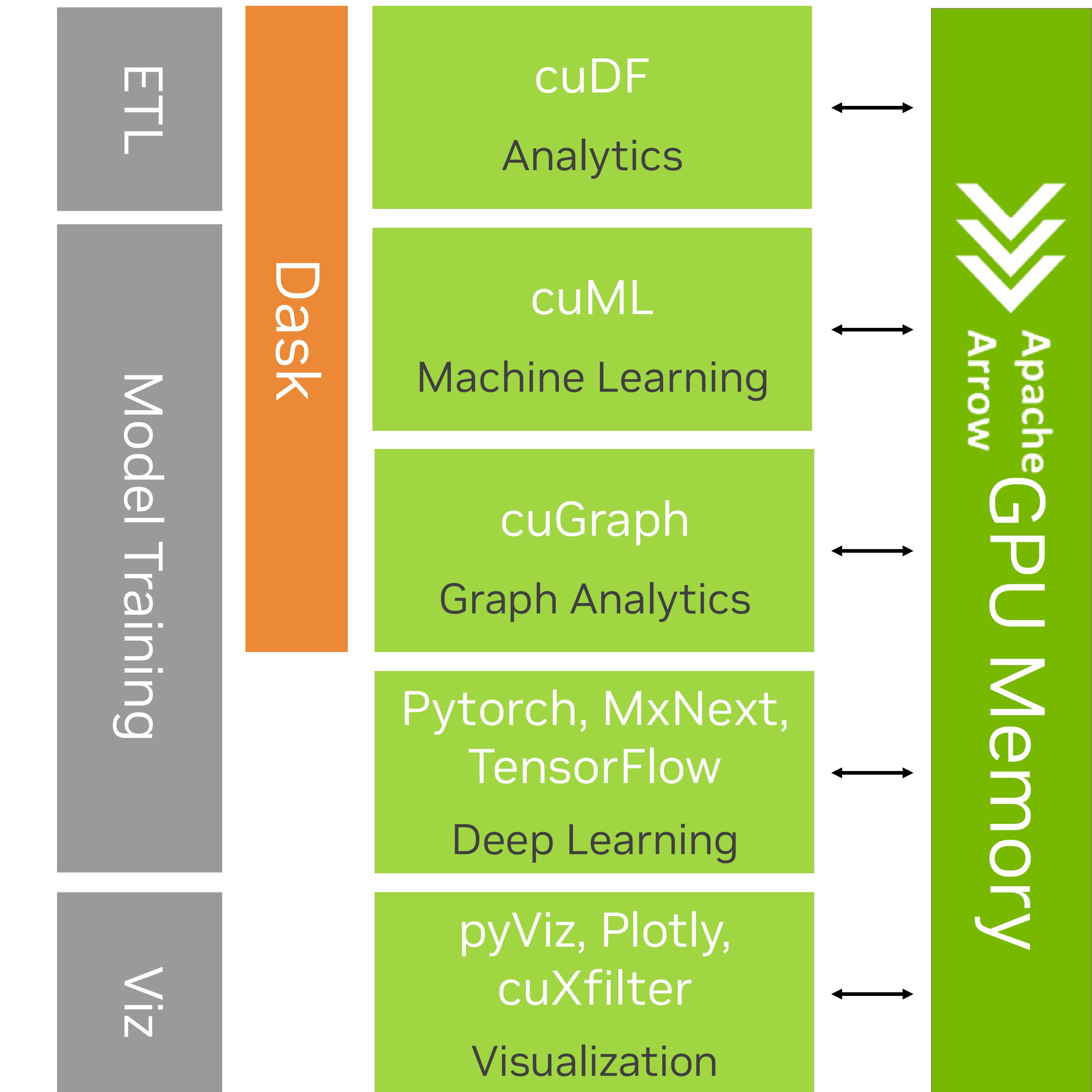
Side by Side Comparison

Same Familiar Python APIs – No/Low Code Change

CPU Approach

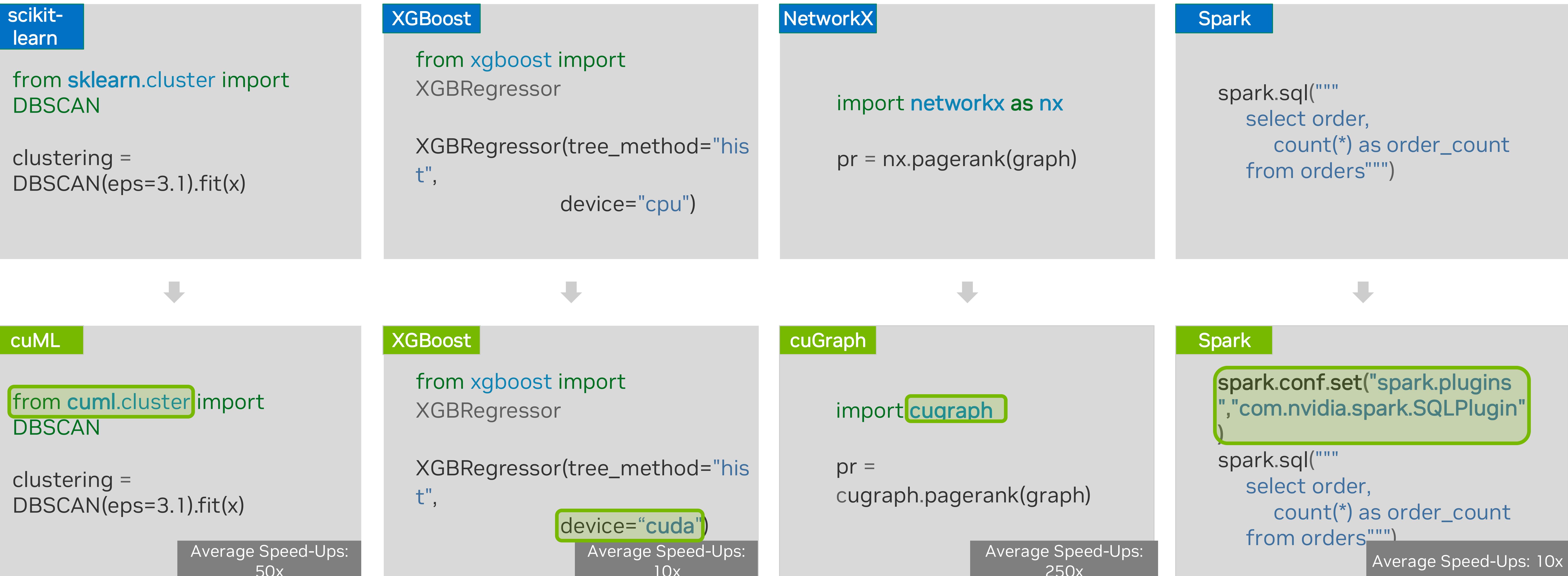


GPU-Accelerated



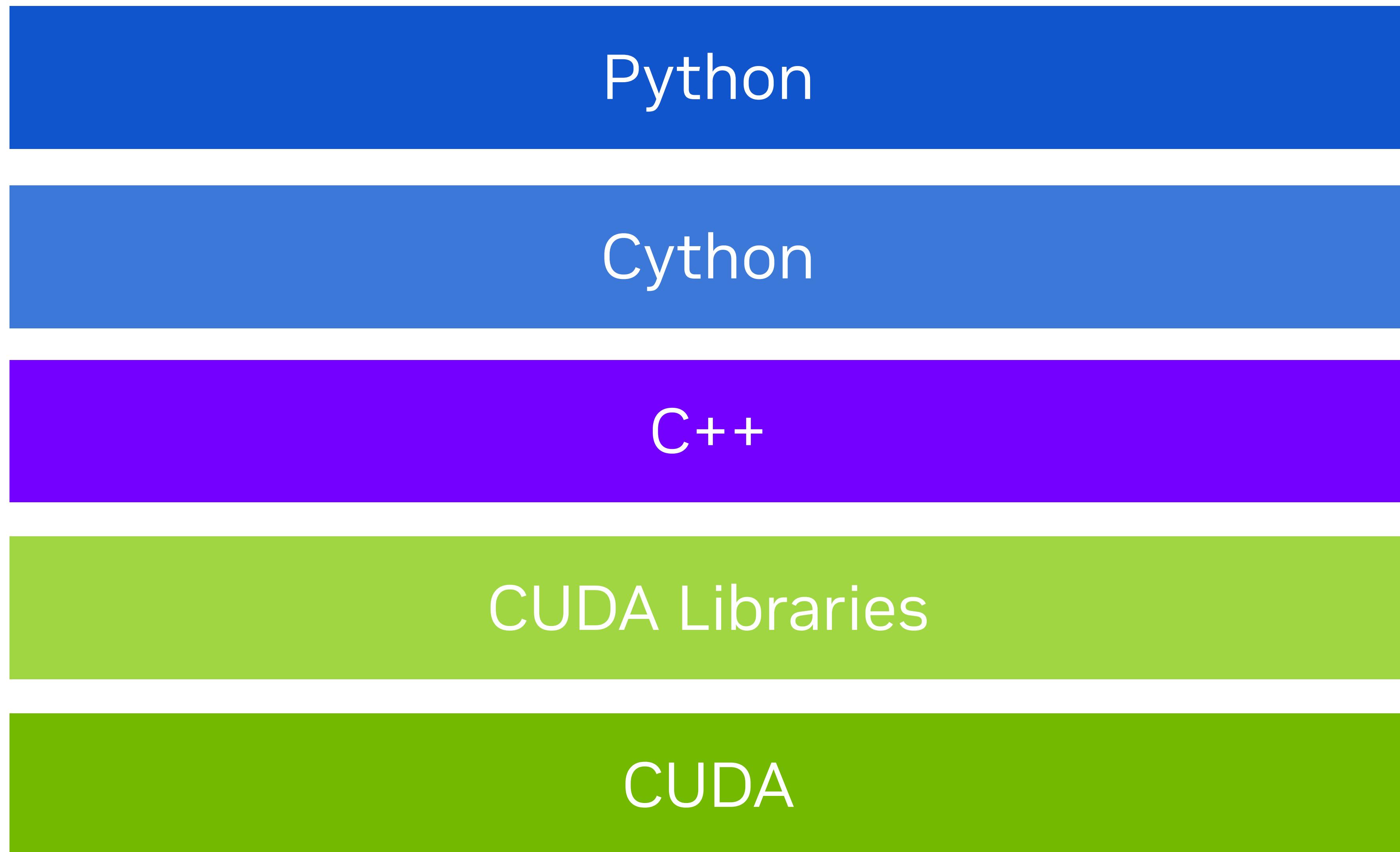
Minor Code Changes for Major Benefits

Abstracting Accelerated Compute through Familiar Interfaces



Technology Stack

RAPIDS



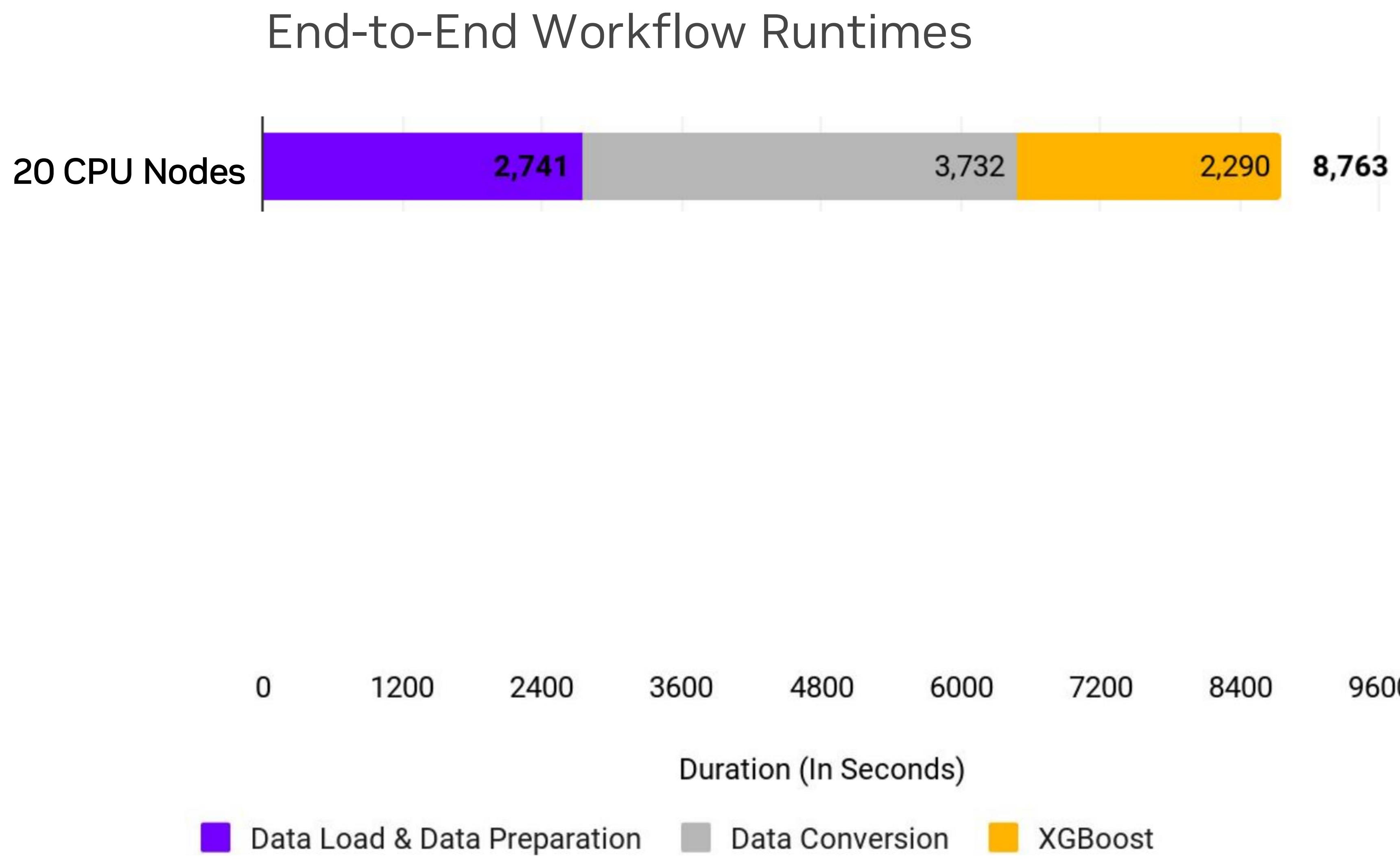
Accelerated Computing Swim Lanes

Making accelerated computing more seamless while enabling deep customization for maximum performance



Lightning-fast End-to-end Performance

Reducing Data Science Processes from Hours to Seconds



16 GPUs Provide More Power than 100 CPU Nodes

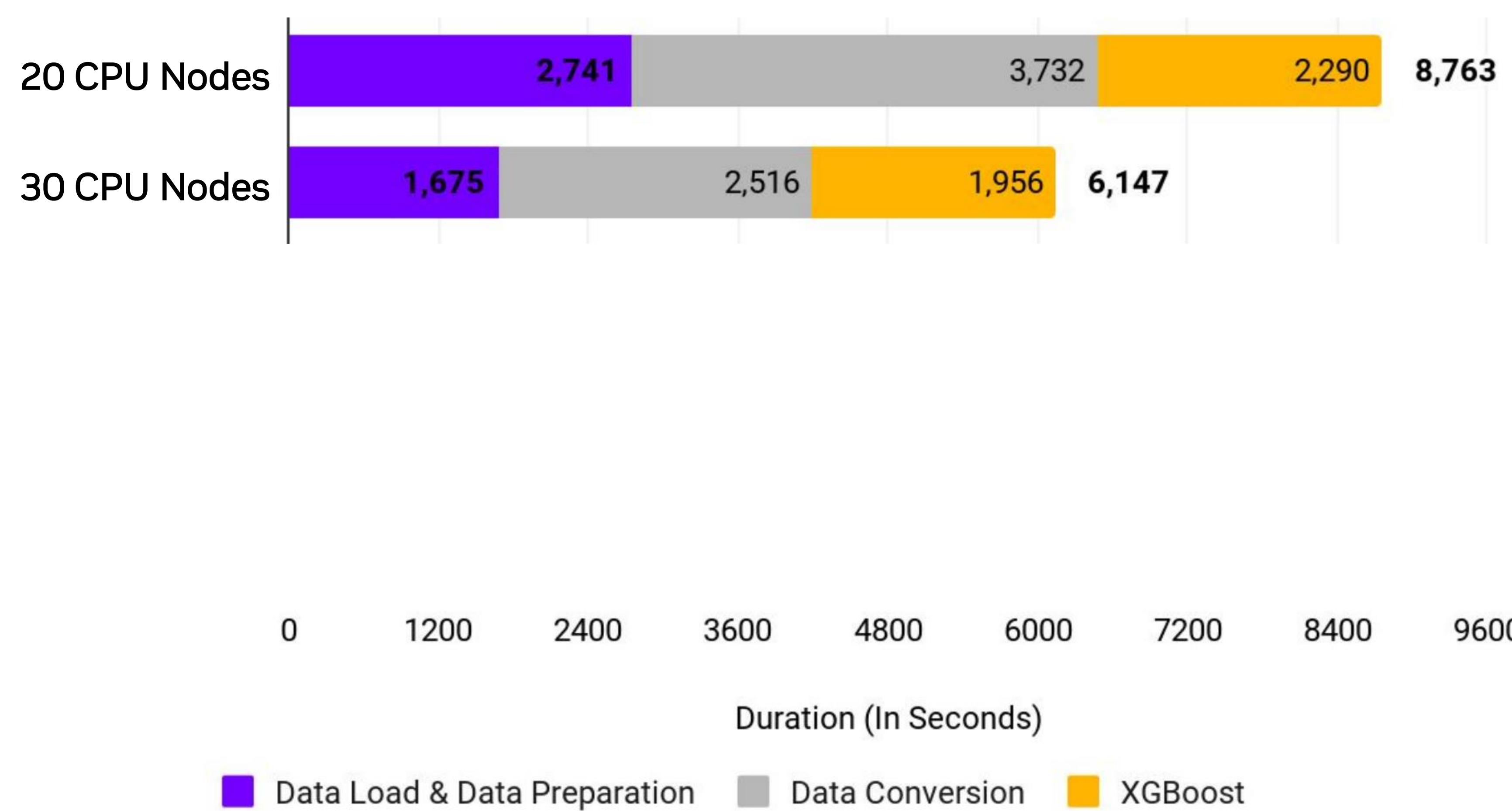
70x Faster Performance than Similar CPU Configuration

20x More Cost-Effective than Similar CPU Configuration

Lightning-fast End-to-end Performance

Reducing Data Science Processes from Hours to Seconds

End-to-End Workflow Runtimes



16 GPUs Provide More Power than 100 CPU Nodes

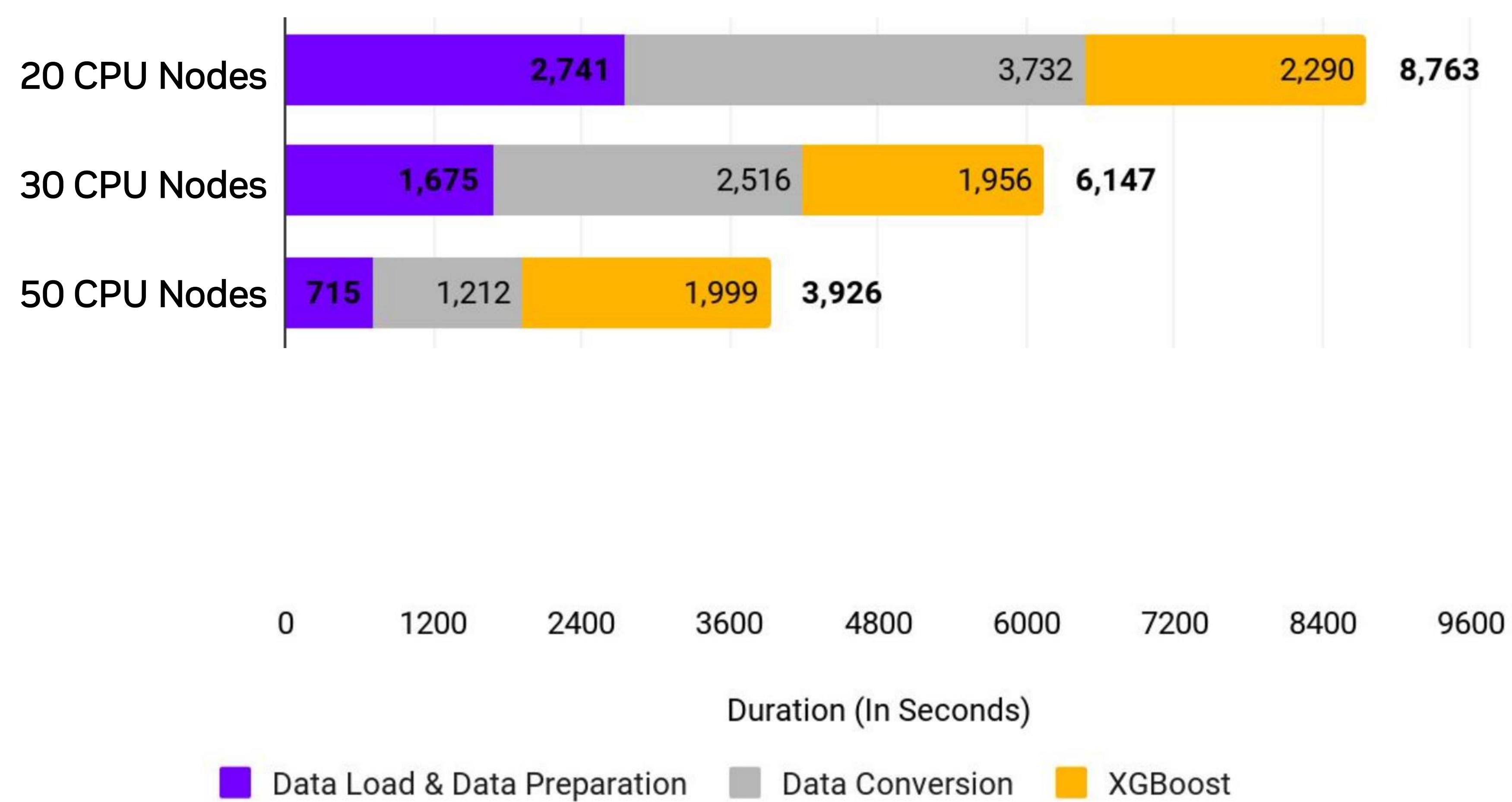
70x Faster Performance than Similar CPU Configuration

20x More Cost-Effective than Similar CPU Configuration

Lightning-fast End-to-end Performance

Reducing Data Science Processes from Hours to Seconds

End-to-End Workflow Runtimes



16 GPUs Provide More Power than 100 CPU Nodes

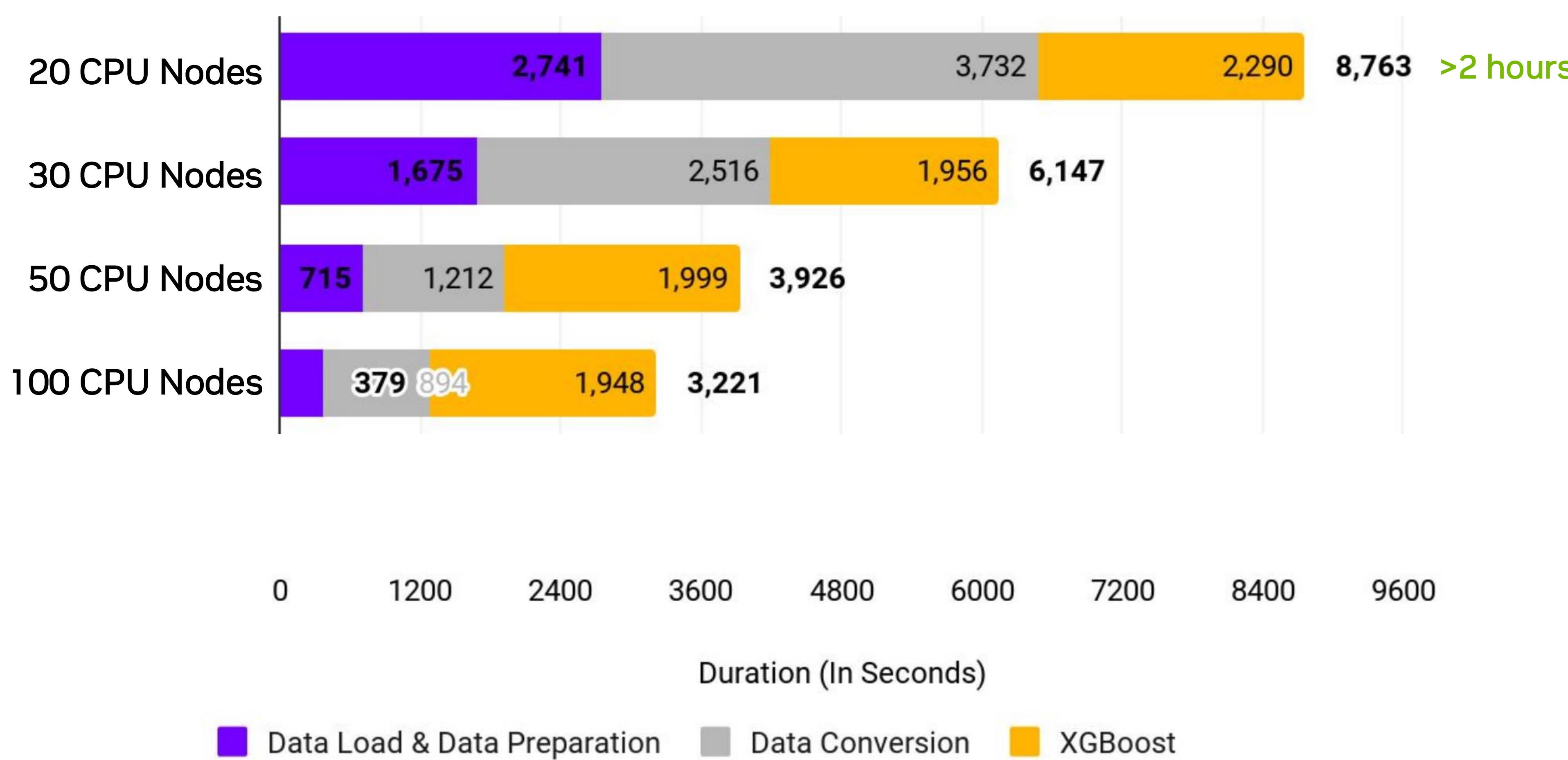
70x Faster Performance than Similar CPU Configuration

20x More Cost-Effective than Similar CPU Configuration

Lightning-fast End-to-end Performance

Reducing Data Science Processes from Hours to Seconds

End-to-End Workflow Runtimes



16 GPUs Provide More Power than 100 CPU Nodes

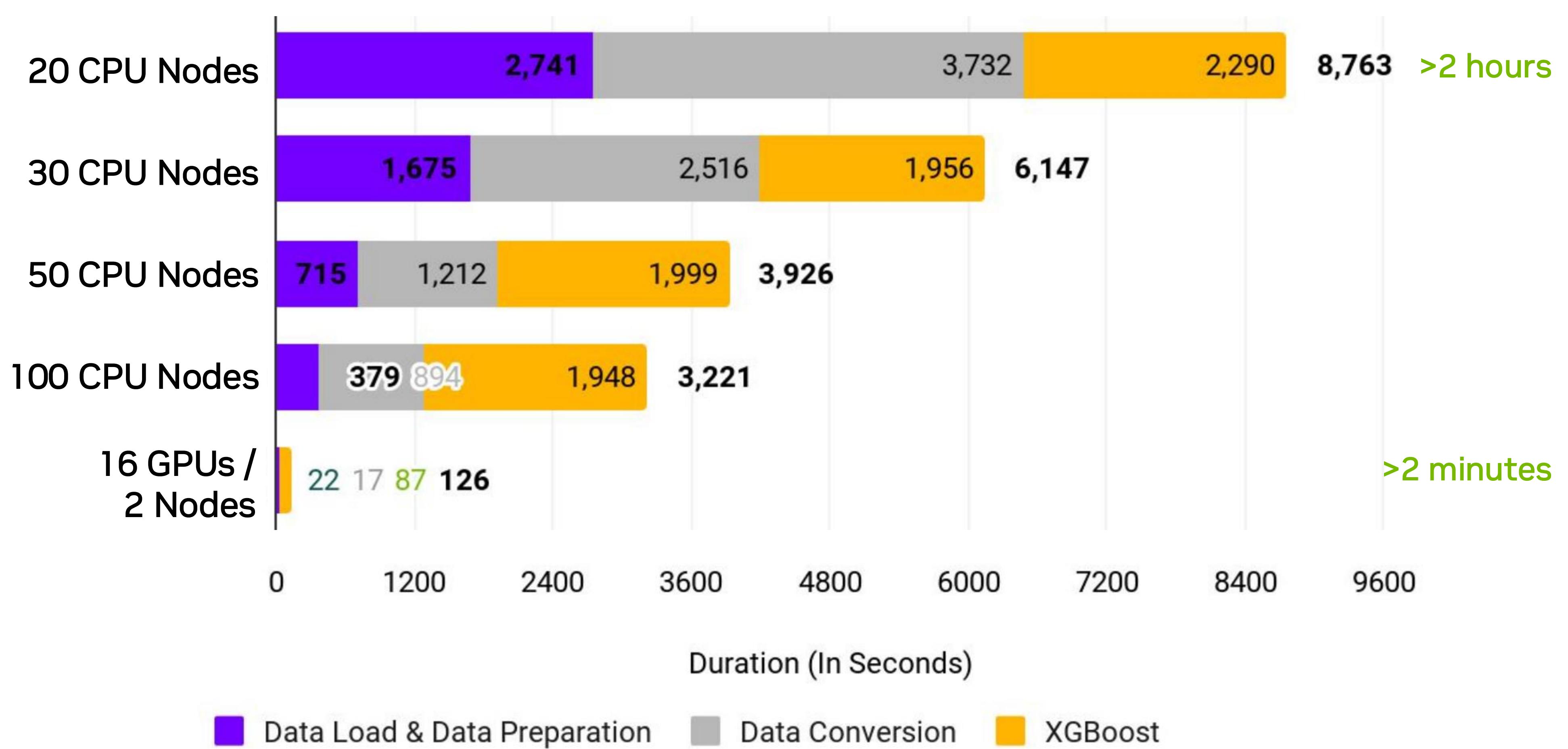
70x Faster Performance than Similar CPU Configuration

20x More Cost-Effective than Similar CPU Configuration

Lightning-fast End-to-end Performance

Reducing Data Science Processes from Hours to Seconds

End-to-End Workflow Runtimes



16 GPUs Provide More Power than 100 CPU Nodes

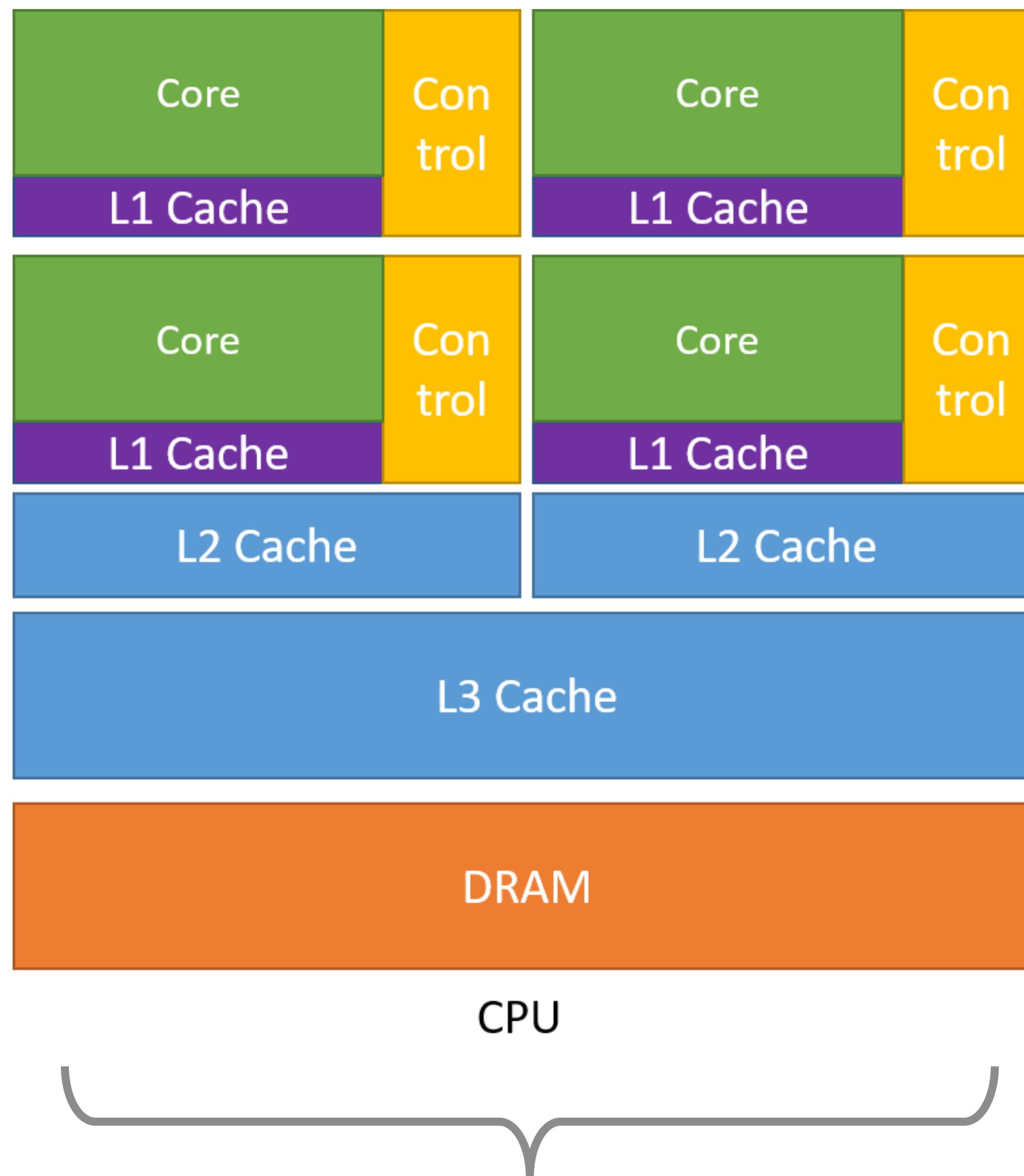
70x Faster Performance than Similar CPU Configuration

20x More Cost-Effective than Similar CPU Configuration

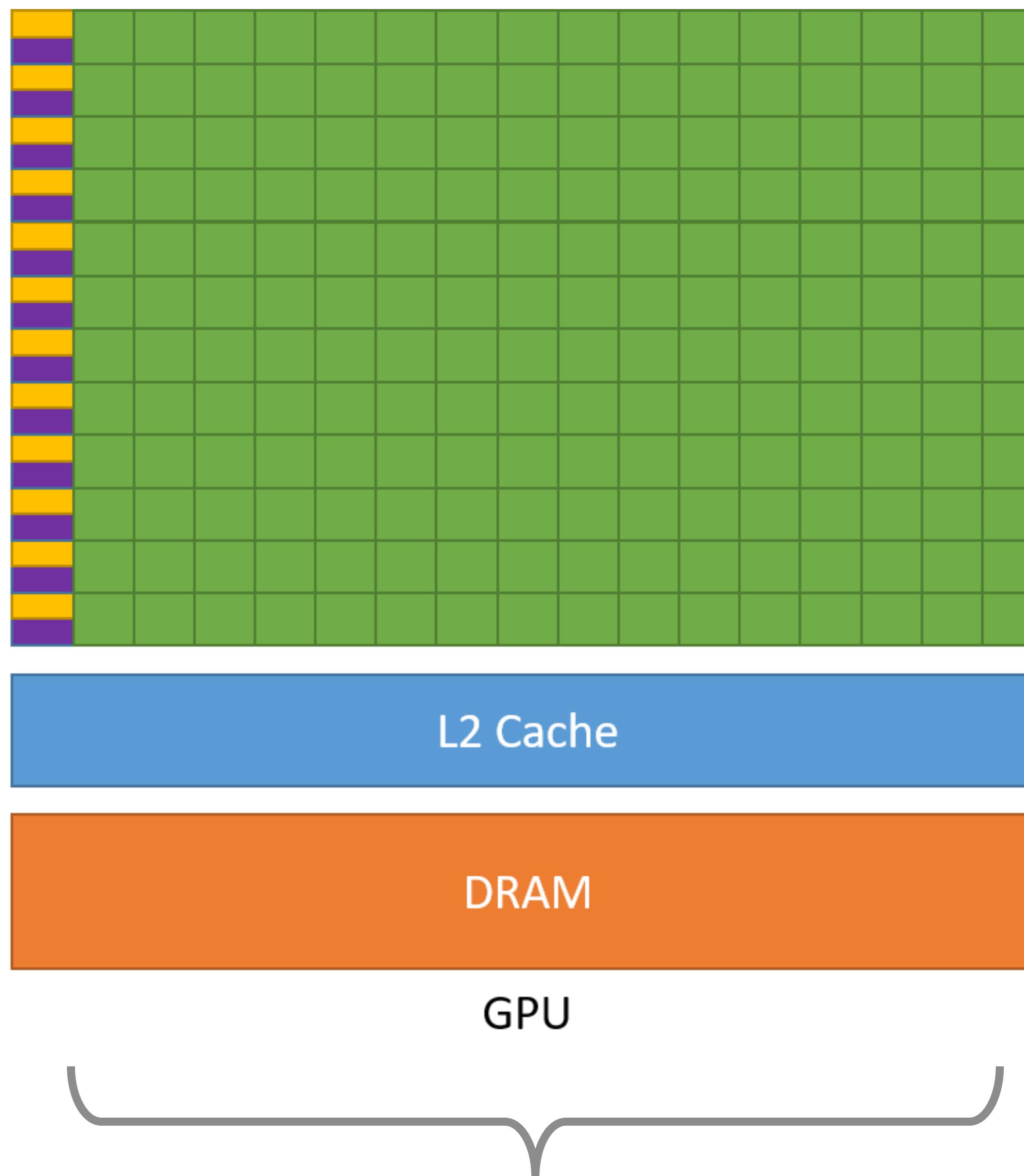
Speedups and Architecture

GPU vs CPU

RAPIDS



Sequential operations (tens
of threads)



Parallel operations
(thousands of threads)

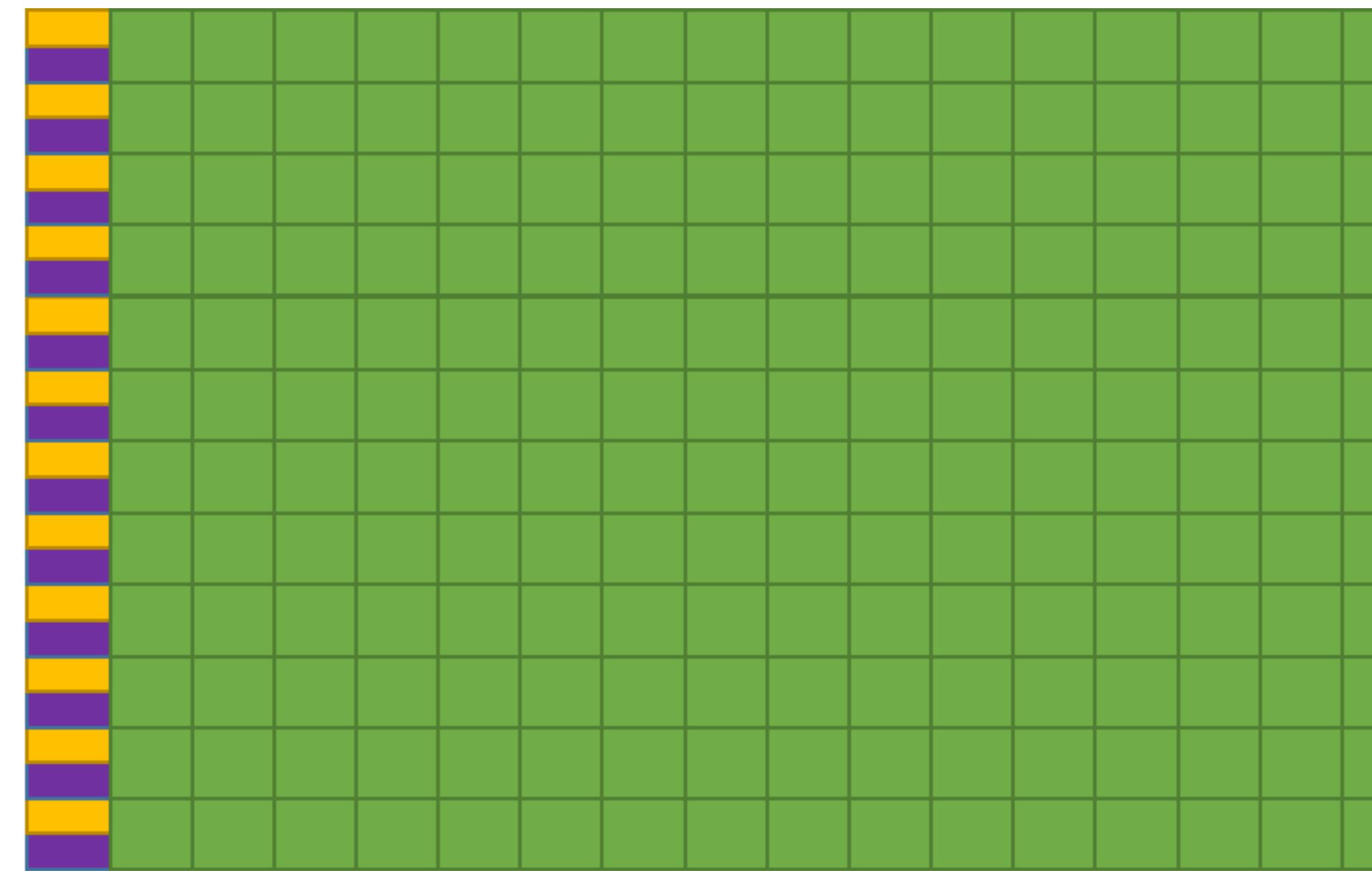


CUDA

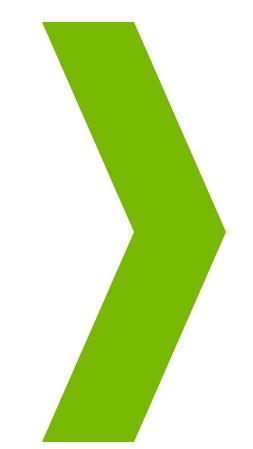
GPU Computing Applications

GPU vs CPU

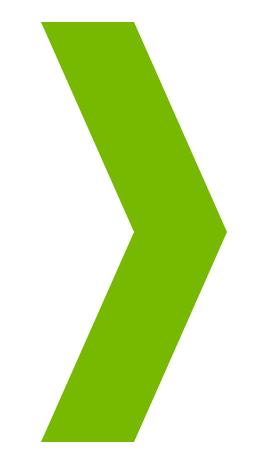
RAPIDS



GPU



Multiple streaming processors (SPs) (32 or 48 CUDA cores) make up 1 streaming multiprocessor (SM)



At time 0 a single warp of 32 threads is executed



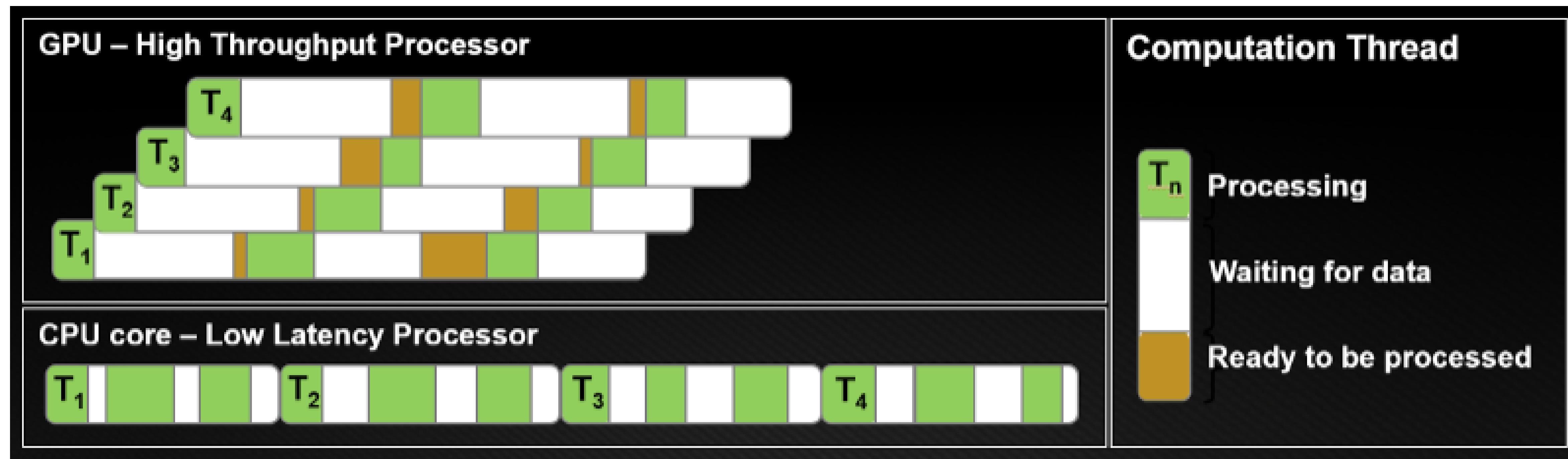
All threads in a warp execute the same instruction



Each thread performs instruction on its own piece of data

GPU vs CPU

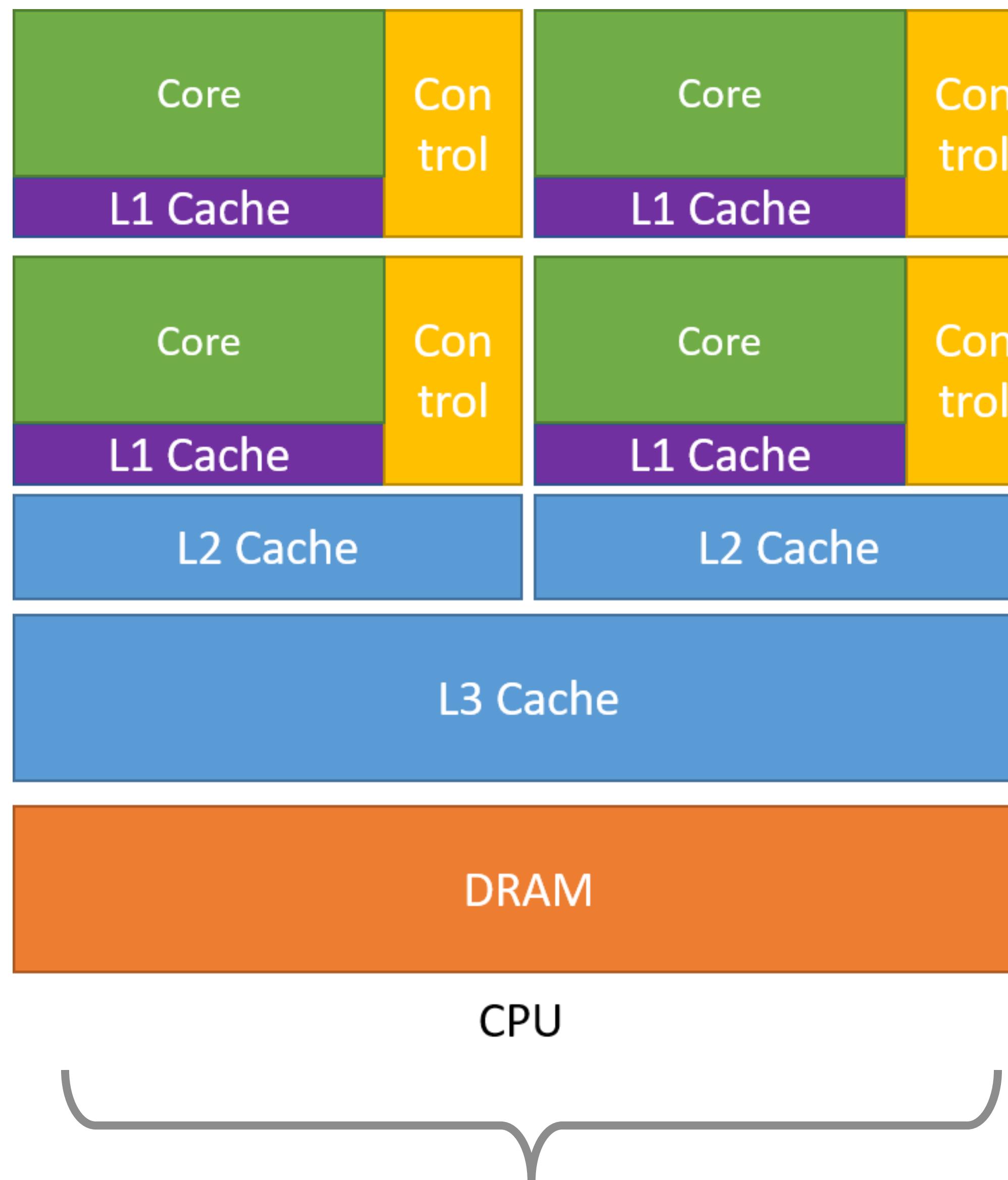
RAPIDS



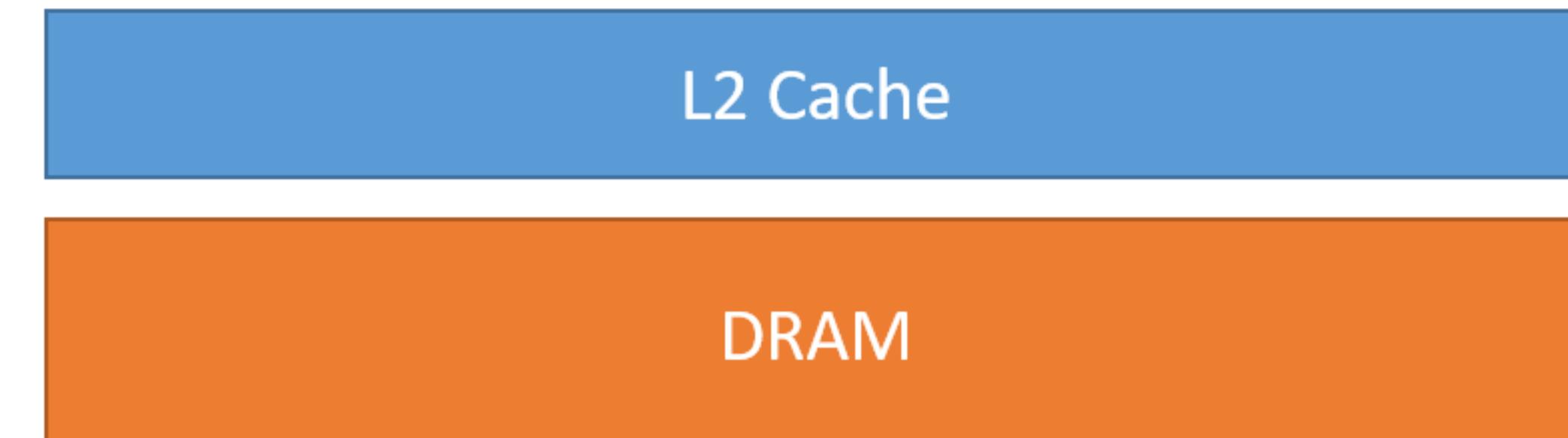
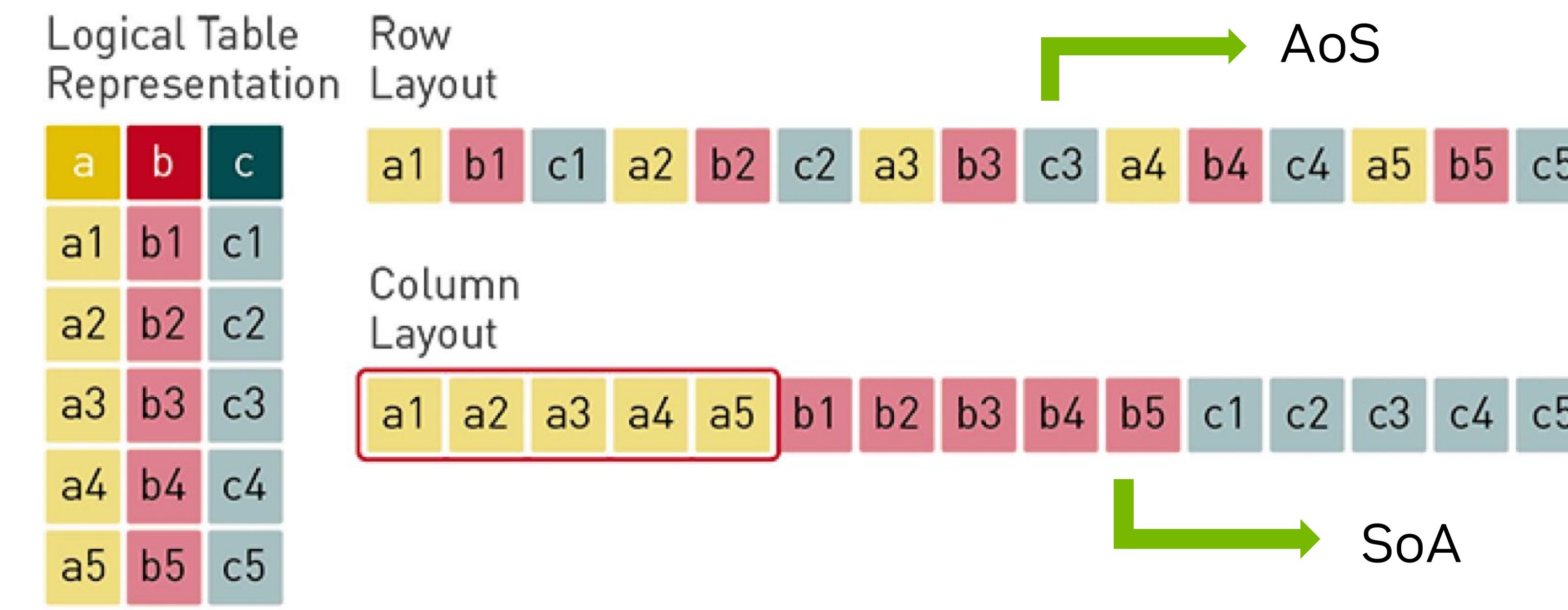
Keep GPUs busy

GPU vs CPU

RAPIDS



Sequential operations (tens
of threads)



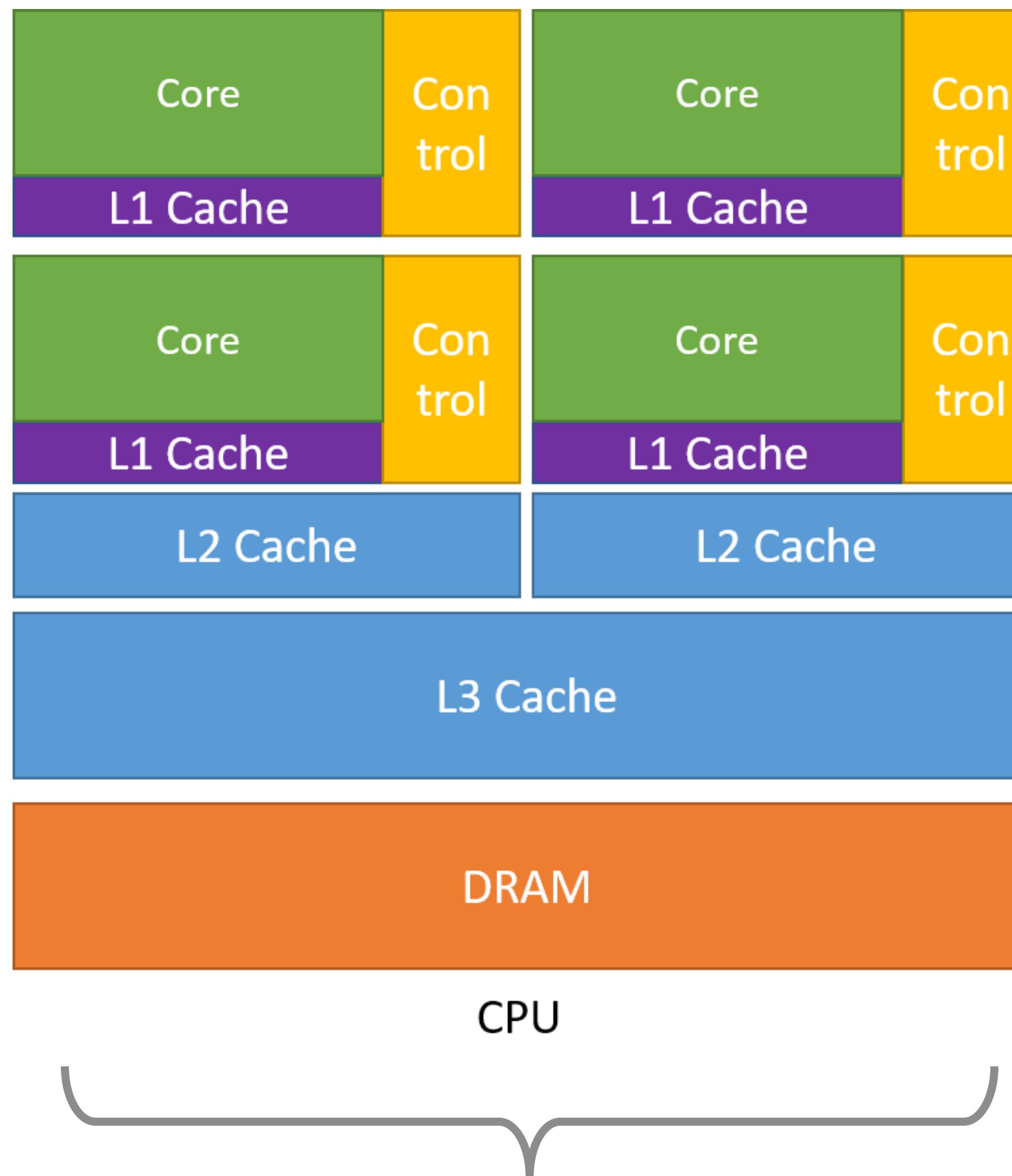
Parallel operations
(thousands of threads)

GPU Computing Applications

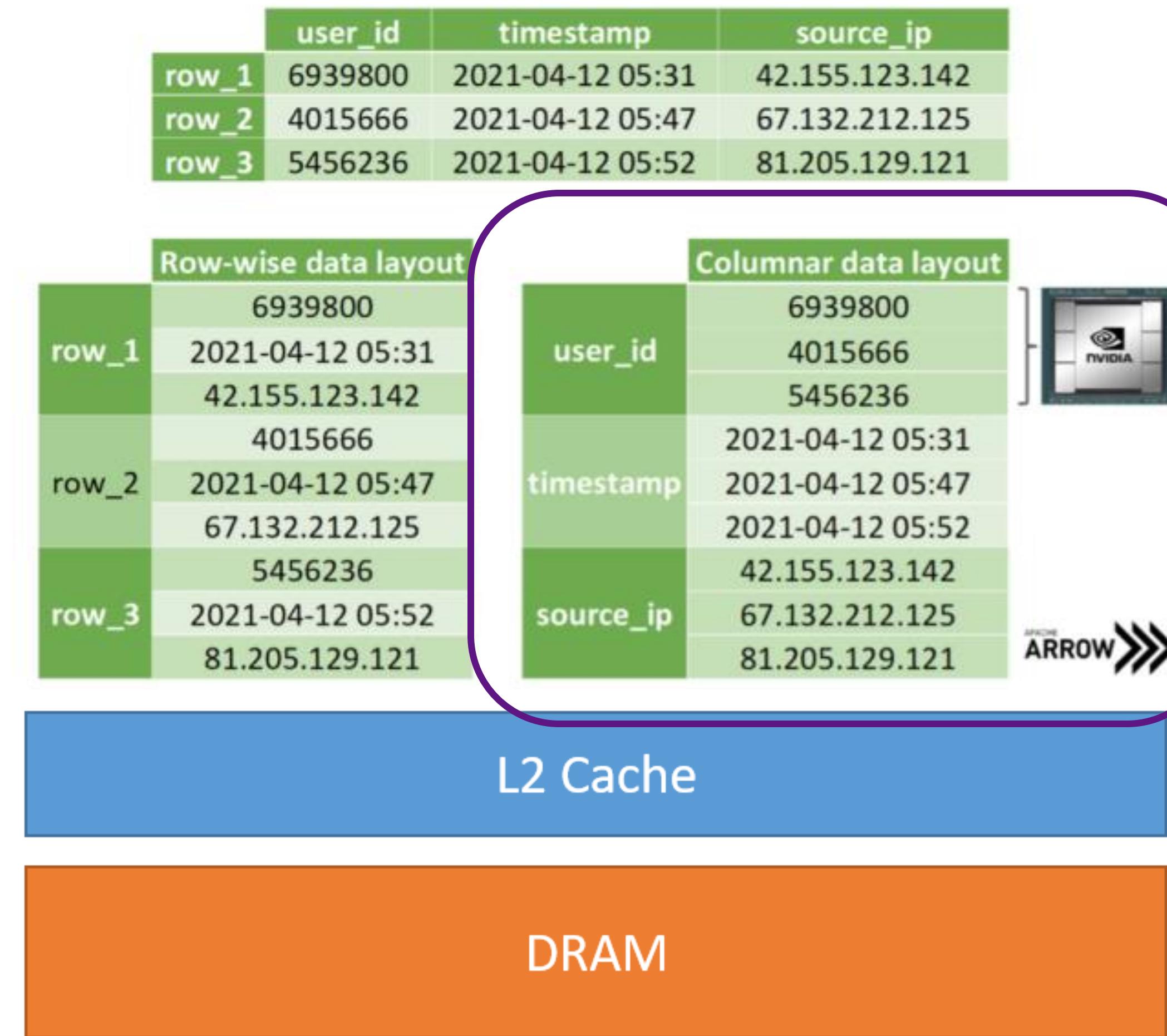
CUDA

GPU vs CPU

RAPIDS



Sequential operations (tens
of threads)



Parallel operations
(thousands of threads)



CUDA

GPU Computing Applications

Apache Arrow

Cross-language Development Platform for In-Memory Analytics



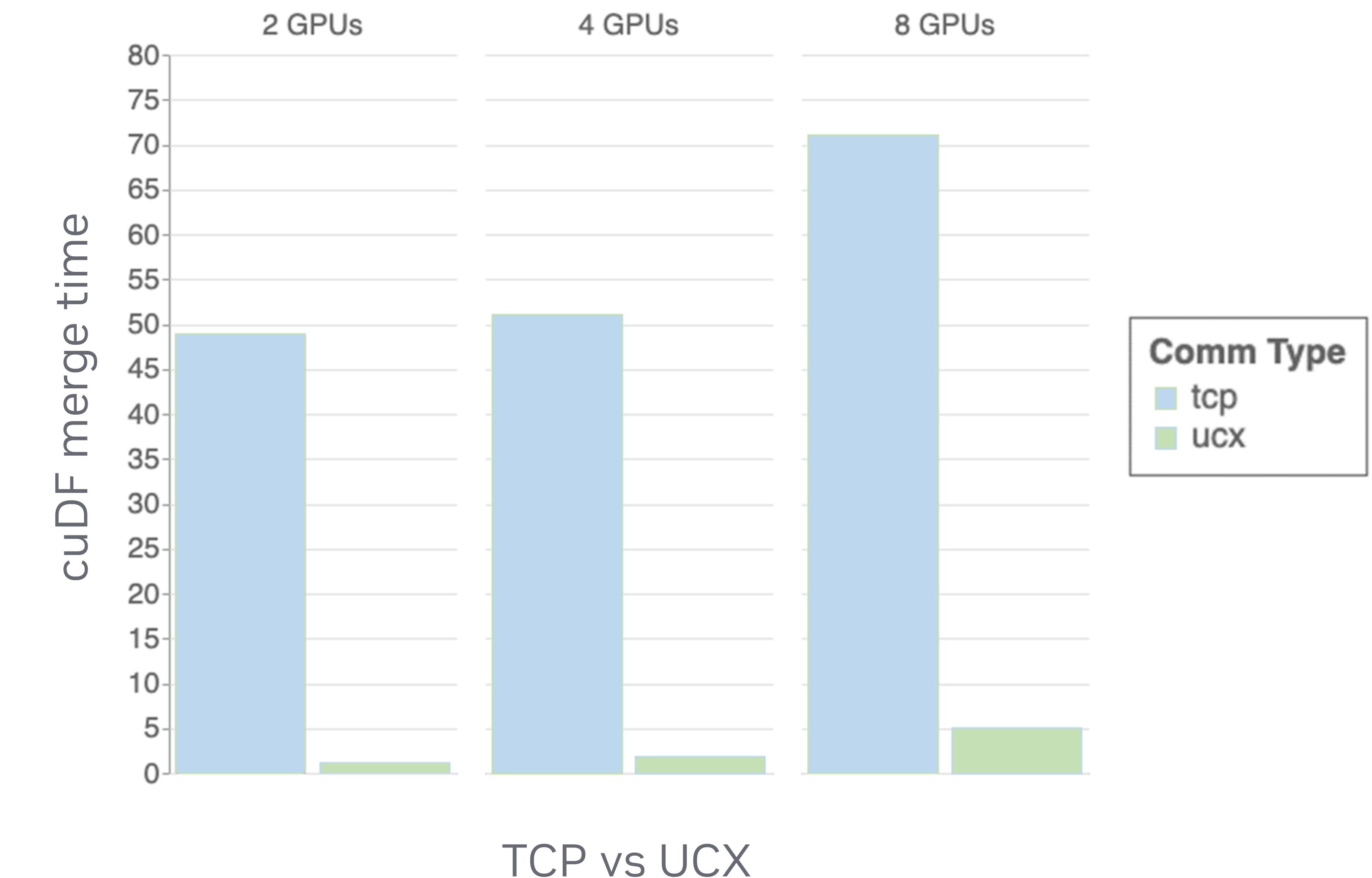
- Each system has its own internal memory format
- 70-80% computation wasted on serialization and deserialization
- Similar functionality implemented in multiple projects
- All systems utilize the same memory format
- No overhead for cross-system communication
- Projects can share functionality (e.g., Parquet-to-Arrow reader)

Unified Communication X (UCX)



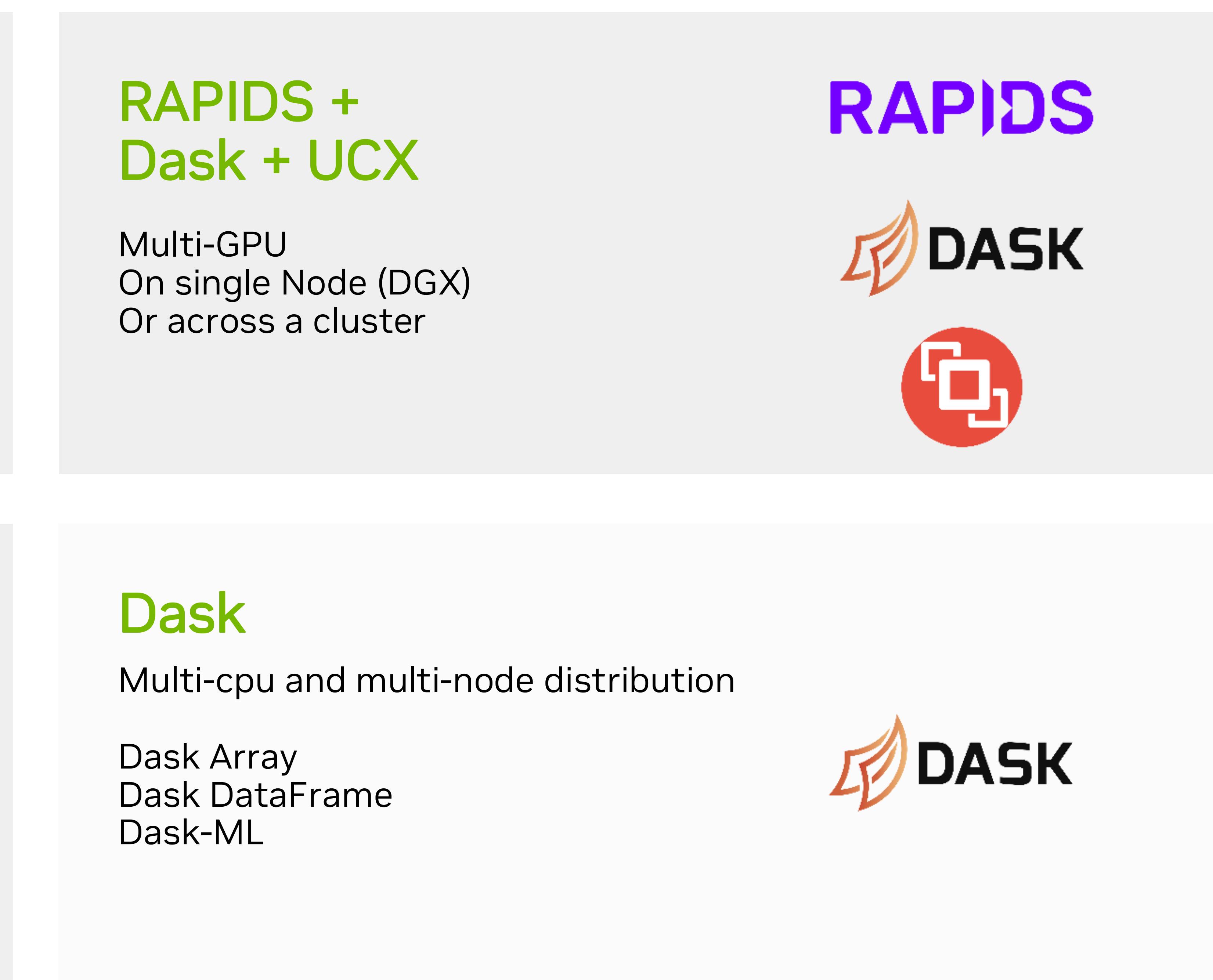
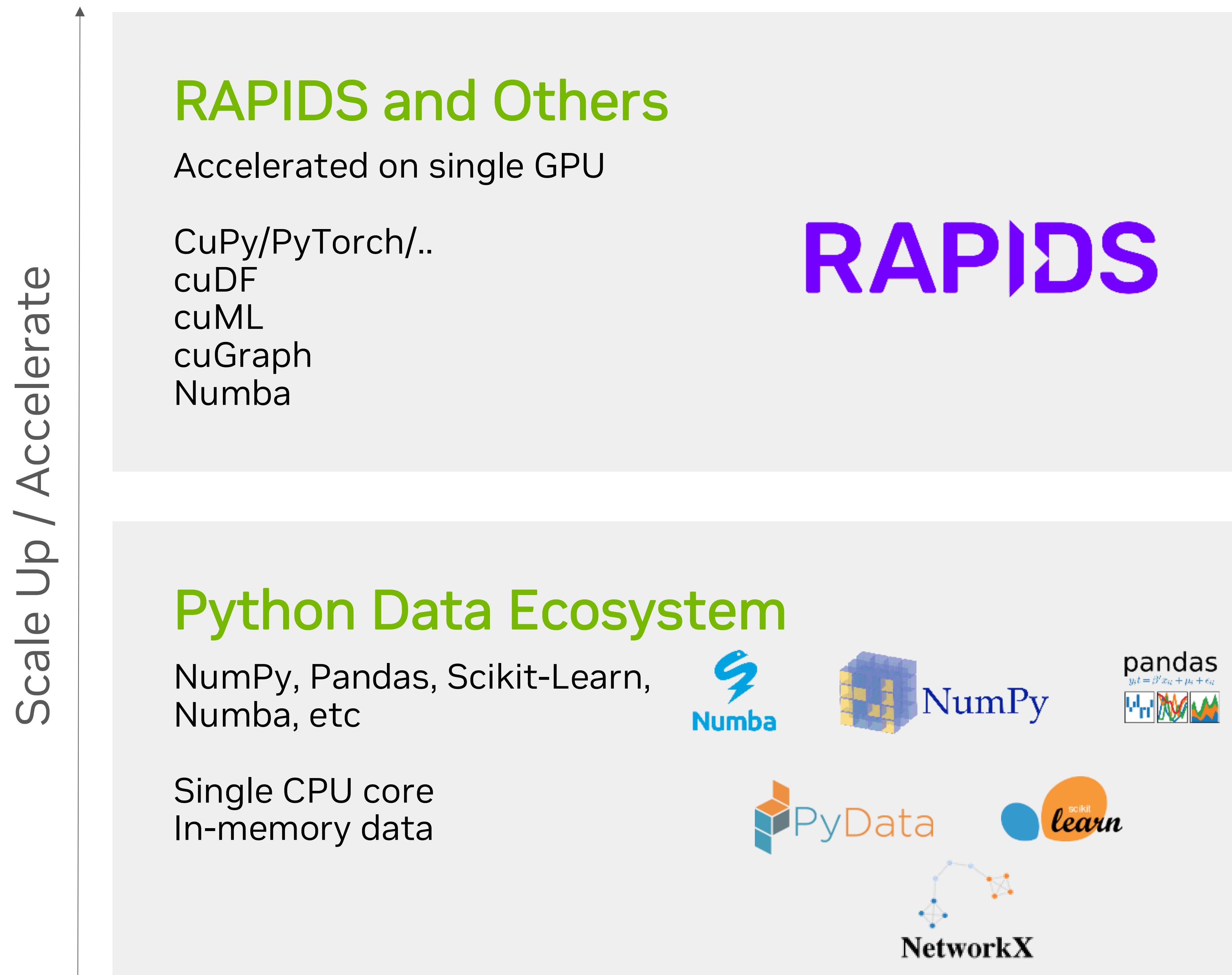
Bringing Hardware Accelerated Communications to Dask

- Dask-CUDA supports integration with UCX
 - Only two-line change necessary
- UCX provides uniform access to transports:
 - NVLink (direct GPU-to-GPU) and Infiniband (direct GPU-Fiber)
 - TCP, Shared memory
- Zero-copy GPU memory transfers over RDMA



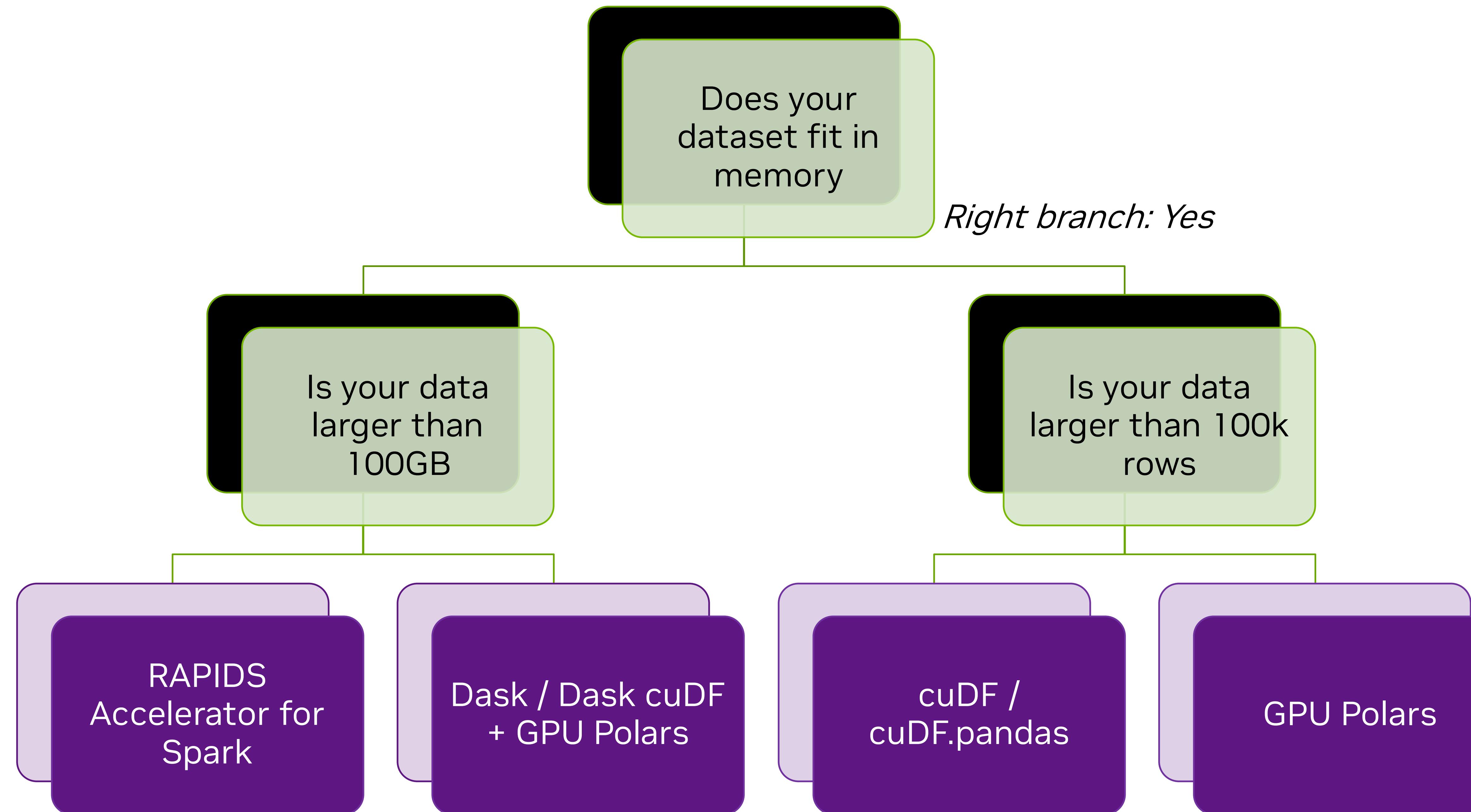
Scaling Up and Out

RAPIDS + Dask + UCX



How to Choose The Correct API

Different APIs are specialized for different workloads



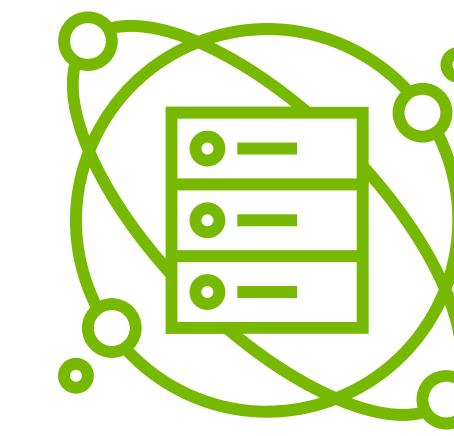
Success Stories

Customer Spotlight: PayPal

Accelerating computing with NVIDIA drives tangible business benefits



With hundreds of thousands of Apache Spark jobs running hourly, CPU based compute is both expensive and time-consuming



By focusing on operations and queries that were computation-bound, PayPal was able to utilize GPU architecture to accelerate operations



NVIDIA technology improved business outcomes by decreasing run time, requiring fewer machines, and increasing machine utilization

Customer Pain Point

Running Spark jobs in the cloud is expensive and time-consuming

NVIDIA Solution

Migration to GPUs allows PayPal to leverage compute

Tangible Business Benefit

Spark RAPIDS and GPUs drove ~70% cost savings



Customer Spotlight: bunq

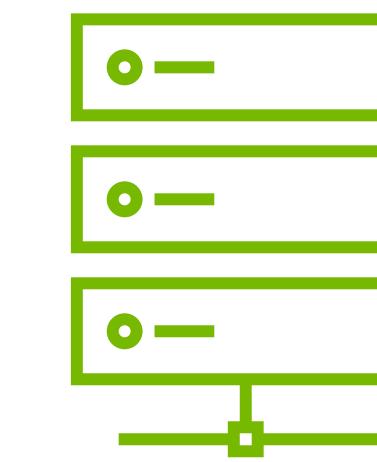
Digital bank debunks financial fraud with generative AI



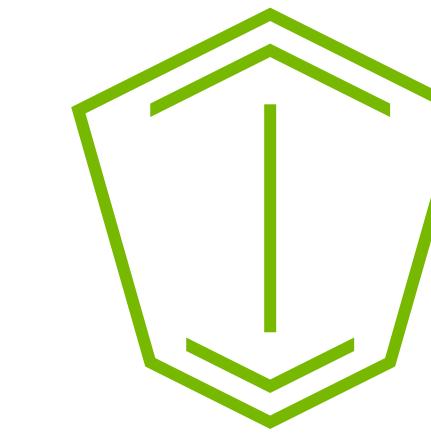
With more than 12 million customers and 8 billion euros' worth of deposits made to date, bunq, an online bank, has become one of the largest neobanks in the EU

Customer Pain Point

Financial fraud and money laundering increasingly prevalent



Using supervised and unsupervised learning, bunq's AI-powered transaction-monitoring system is completely automated and easily scalable



Beyond model training, bunq utilizes NVIDIA solutions for Finn, bunq's personal AI-assistant, automated marketing efforts, and handling customer service tickets

bunq

NVIDIA®

Realized Business Benefits

Data processing accelerated 5x;
fraud detection model trained 100x faster

NVIDIA RAPIDS Acceleration Libraries

Data Preparation

cuDF

Less Waiting. More Building.

“ Data preparation and exploration takes 60 to 80 percent of the analytical pipeline in a typical machine learning or deep learning project ”

Accelerated data processing

Optimized libraries for various data sizes

100k of rows



More than 100k of rows



More than 100GB of data



Accelerated pandas

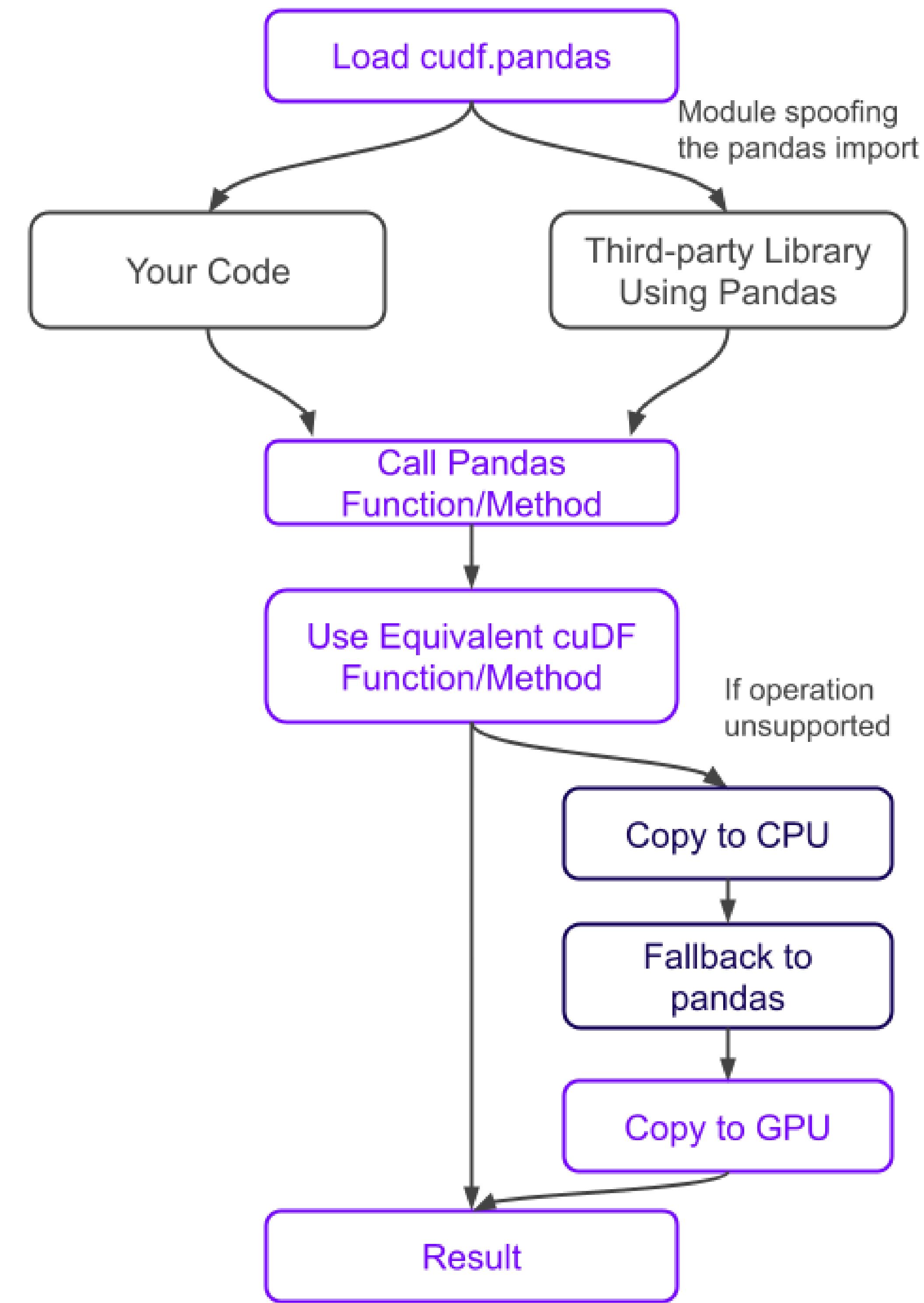
cudf.pandas: the zero code change GPU accelerator for pandas built on cuDF

- 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

```
[ ]: [REDACTED] 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)
```

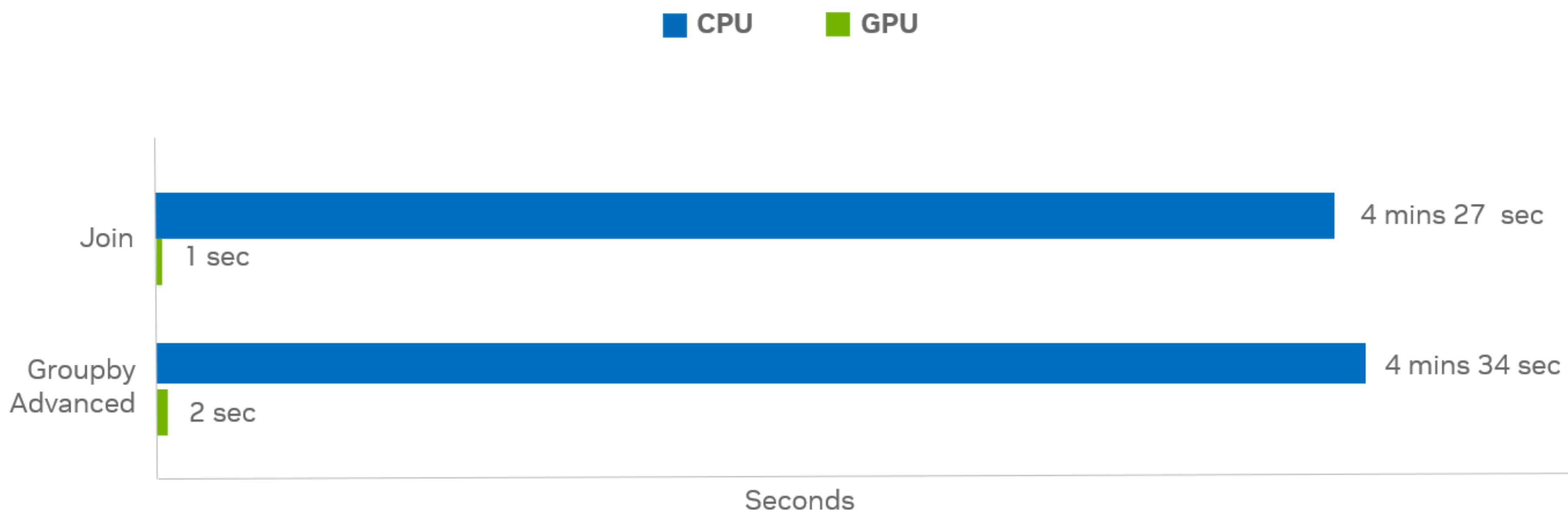


Accelerated pandas

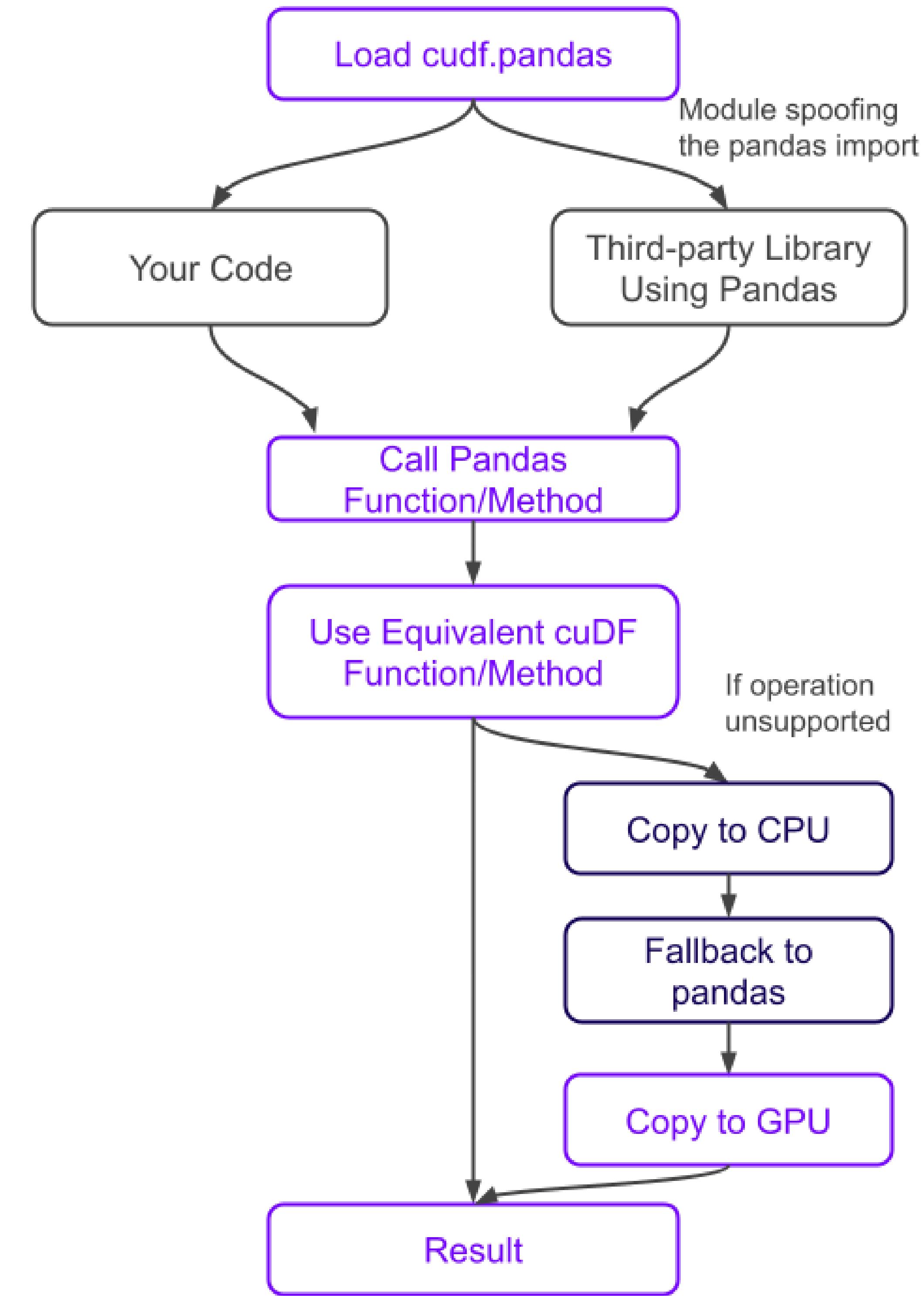
cudf.pandas: the zero code change GPU accelerator for pandas built on cuDF

- 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

150x Faster pandas



Standard DuckDB Data Benchmark (5 GB) performance comparison
between cuDF.pandas and traditional pandas v2.2



cuDF pandas accelerator mode (cudf.pandas)

Bringing the speed of cuDF to every Pandas user



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

weekday_names = {
    0: "Monday",
    1: "Tuesday",
    2: "Wednesday",
    3: "Thursday",
    4: "Friday",
    5: "Saturday",
    6: "Sunday",
}

df = pd.read_parquet("/tmp/nyc_parking_violations_2022.parquet")
df["Issue Date"] = df["Issue Date"].astype("datetime64[ms]")
df["issue_weekday"] = df["Issue Date"].dt.weekday.map(weekday_names)
df.groupby(["issue_weekday"])["Summons Number"].count().sort_values()

CPU times: user 601 ms, sys: 141 ms, total: 742 ms
Wall time: 720 ms

issue_weekday
Sunday      462992
Saturday    1108385
Monday      2488563
Wednesday   2760088
Tuesday     2809949
Friday       2891679
Thursday    2913951
Name: Summons Number, dtype: int64
```

- Constantly ensured full Pandas compatibility:
 - Unit-testing with Pandas test suite
 - Integration test for key PyData libraries

cuDF pandas accelerator mode (cudf.pandas)

Bringing the speed of cuDF to every Pandas user

```
%%cudf.pandas.profile
import pandas as pd

weekday_names = {
    0: "Monday",
    1: "Tuesday",
    2: "Wednesday",
    3: "Thursday",
    4: "Friday",
    5: "Saturday",
    6: "Sunday",
}

df = pd.read_parquet("/tmp/nyc_parking_violations_2022.parquet")
df["Issue Date"] = df["Issue Date"].astype("datetime64[ms]")
df["issue_weekday"] = df["Issue Date"].dt.weekday.map(weekday_names)
df.groupby(["issue_weekday"])["Summons Number"].count().sort_values()

Total time elapsed: 1.854 seconds
14 GPU function calls in 1.200 seconds
2 CPU function calls in 0.027 seconds
```

Function	GPU ncalls	GPU cumtime	GPU percall	CPU ncalls	CPU cumtime	CPU percall
read_parquet	1	0.841	0.841	0	0.000	0.000
DataFrame.__getitem__	2	0.002	0.001	0	0.000	0.000
NDFrame.astype	1	0.043	0.043	0	0.000	0.000
DataFrame.__setitem__	2	0.006	0.003	0	0.000	0.000
Series	1	0.000	0.000	0	0.000	0.000
Series.map	1	0.253	0.253	0	0.000	0.000
DataFrame.groupby	1	0.001	0.001	0	0.000	0.000
DataFrameGroupBy.__getitem__	1	0.002	0.002	0	0.000	0.000
GroupBy.count	1	0.028	0.028	0	0.000	0.000
Series.sort_values	1	0.005	0.005	0	0.000	0.000
Series.__repr__	1	0.017	0.017	0	0.000	0.000
Series.to_frame	1	0.003	0.003	0	0.000	0.000
DataFrame._repr_html_	0	0.000	0.000	1	0.020	0.020
NDFrame._repr_latex_	0	0.000	0.000	1	0.007	0.007

Not all pandas operations ran on the GPU. The following functions required CPU fallback:

- DataFrame._repr_html_
- NDFrame._repr_latex_

To request GPU support for any of these functions, please file a Github issue here:
<https://github.com/rapidsai/cudf/issues/new/choose>

- Profiler built in
- Accessible in command line and notebook
- Possible to see how much time each function spends divided on CPU and GPU

Accelerated polars

Zero code change accelerator for Polars

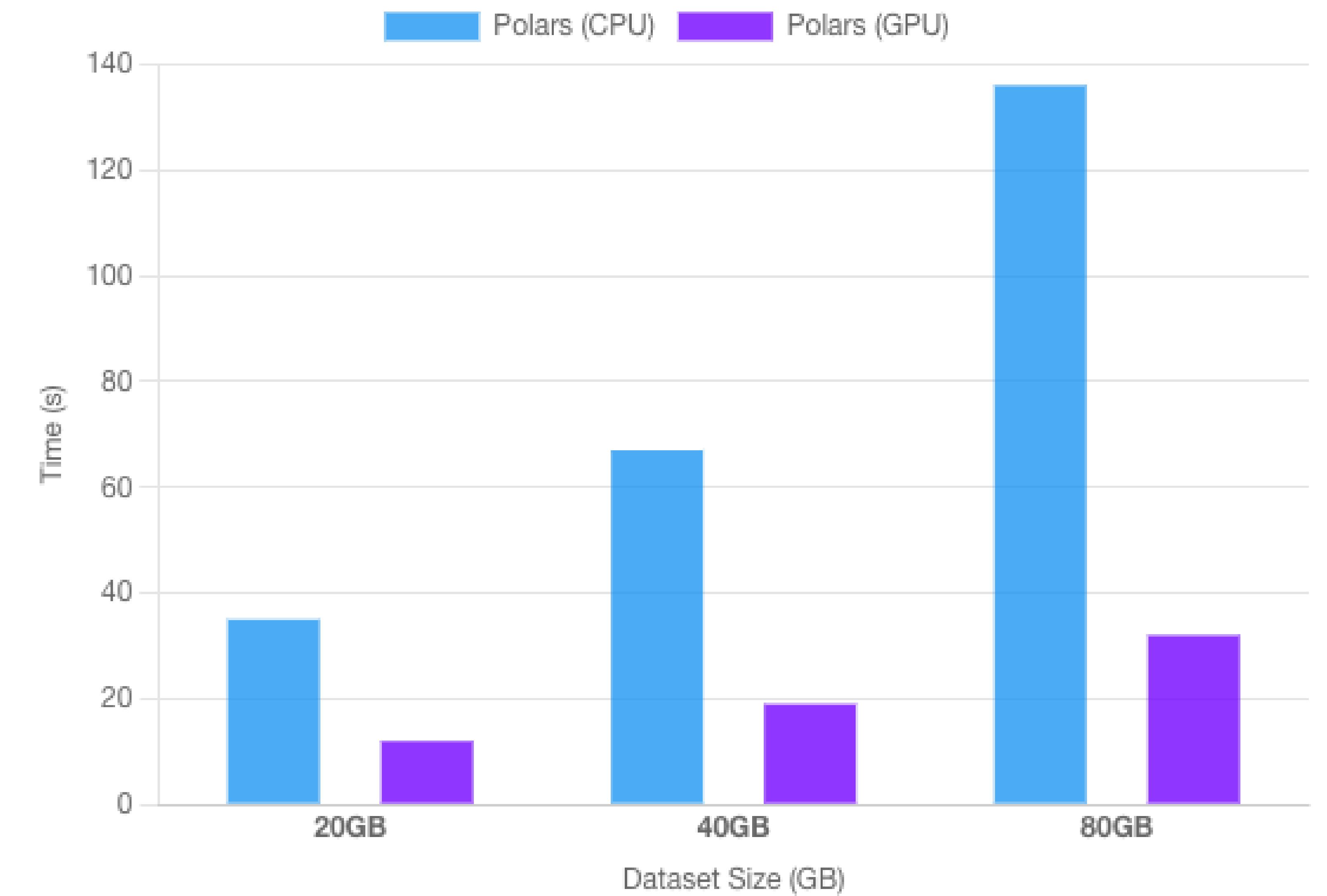
- Zero code change to existing Polars code
- Polars attempts to execute operations on GPU first and falls back to CPU if necessary
- Inspection of optimized query plan to attempt to build execution plan. If fails → fallback to CPU

```
import polars as pl

l df = pl.LazyFrame({"a": [1.242, 1.535]})

print(
    l df.select(
        pl.col("a").round(1)
    ).collect(engine="gpu")
)
```

Query Processing time across a range of dataset sizes



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)

```
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)
```

```
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

Data Preparation

NVIDIA Spark RAPIDS

How Spark RAPIDS Works

Key technologies for GPU acceleration

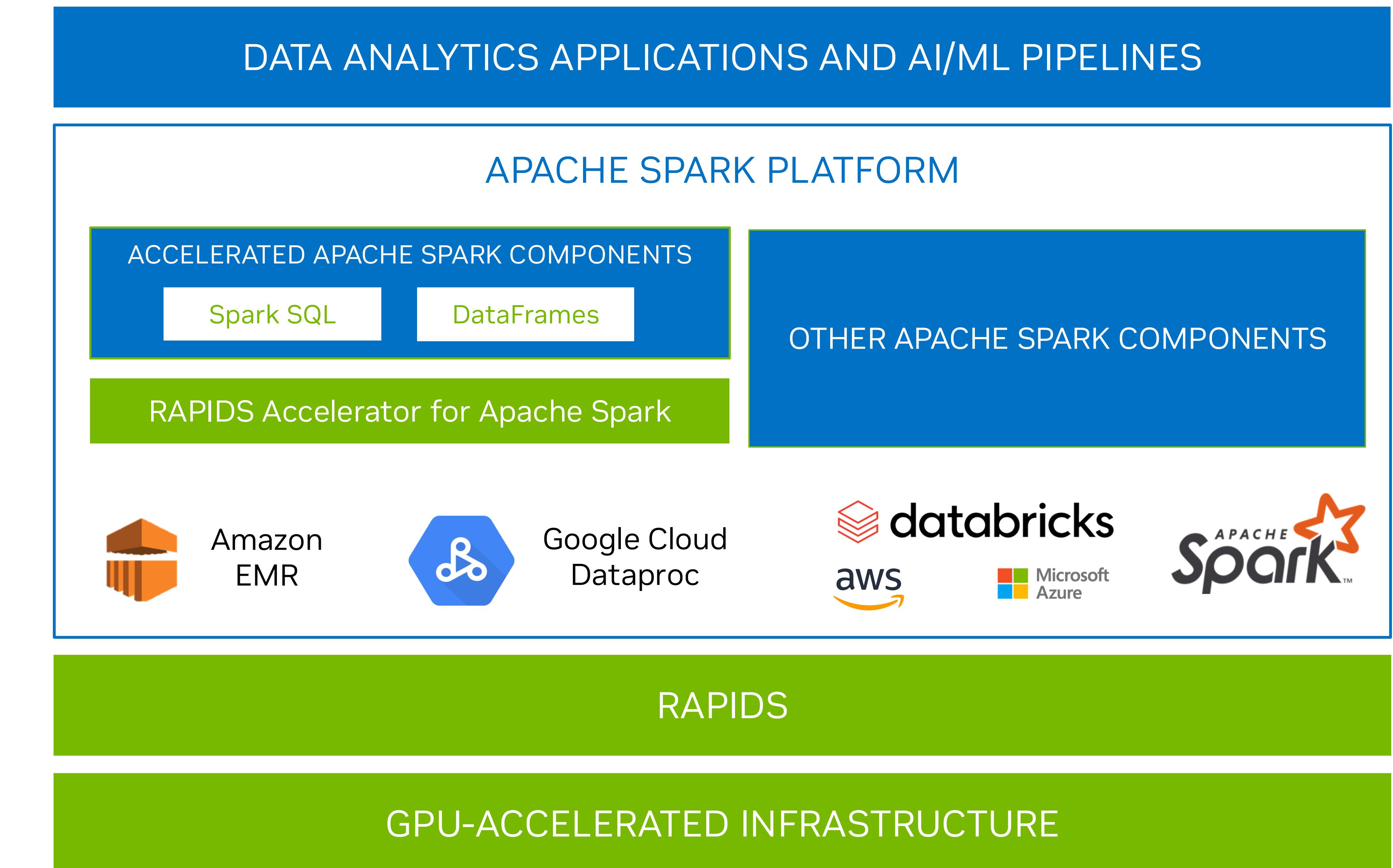
NVIDIA RAPIDS Accelerator

- Operates as a **software plugin** to popular Apache Spark platforms
 - Automatically accelerates supported operations
 - Requires no code changes
- Operations currently accelerated
 - Spark SQL
 - DataFrame
- Works with Spark standalone, YARN clusters, Kubernetes clusters

Key Spark 3 Innovations

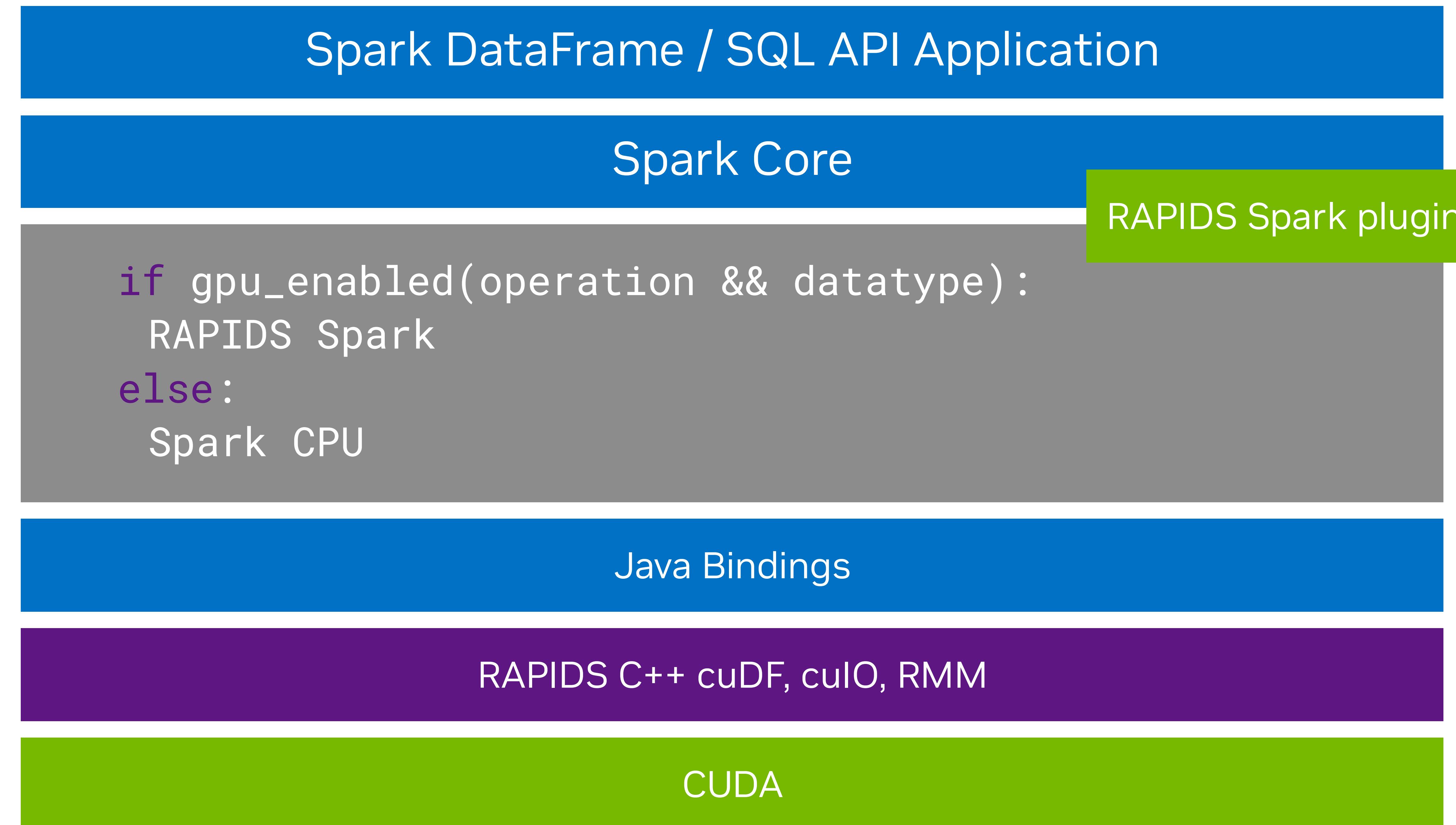
Columnar processing support in the Catalyst query optimizer – allows efficient GPU acceleration

GPU-aware scheduling of executors with a specified number of GPUs and how many GPUs for each task



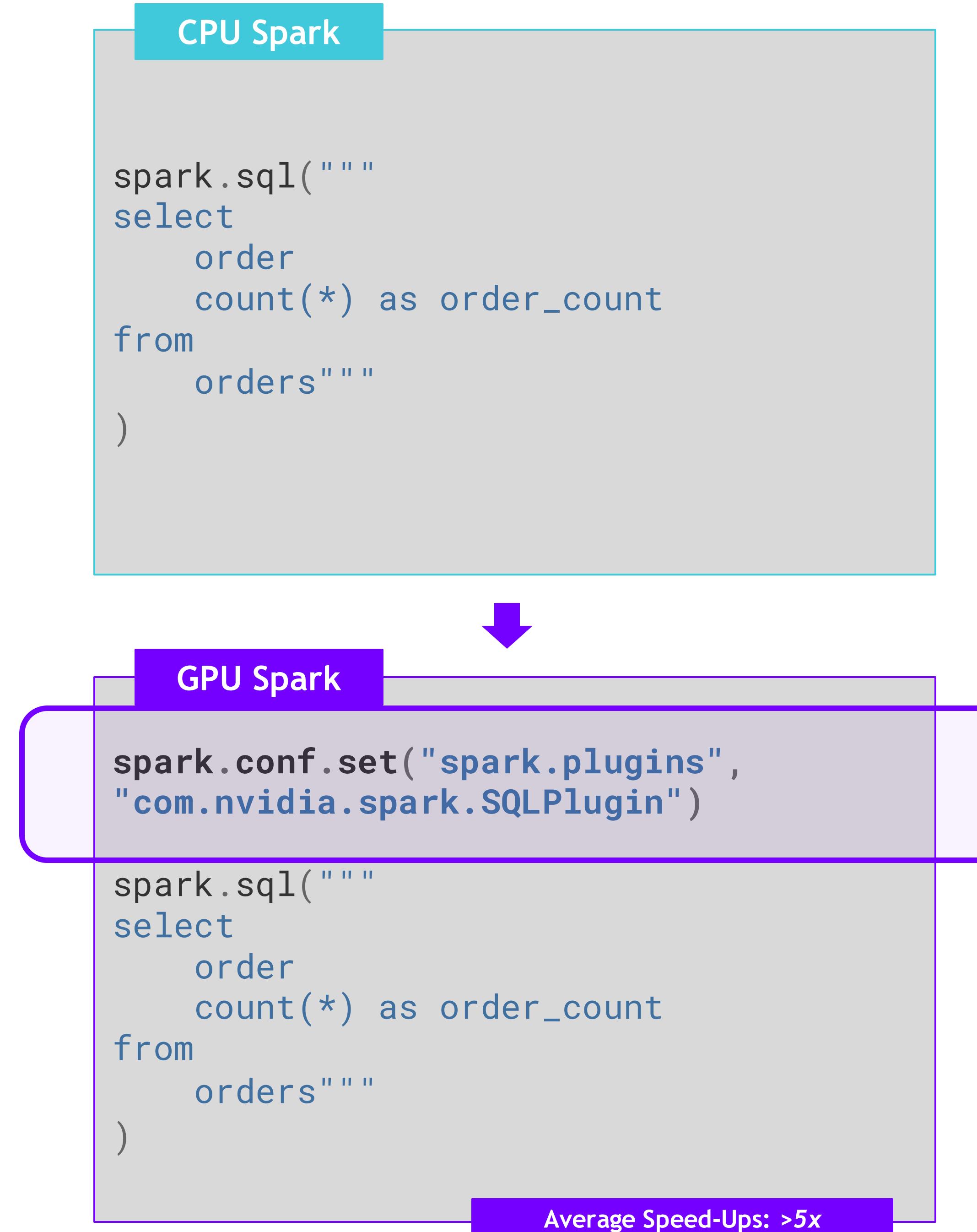
RAPIDS Accelerator for Apache Spark

Spark Plugin for GPU Acceleration



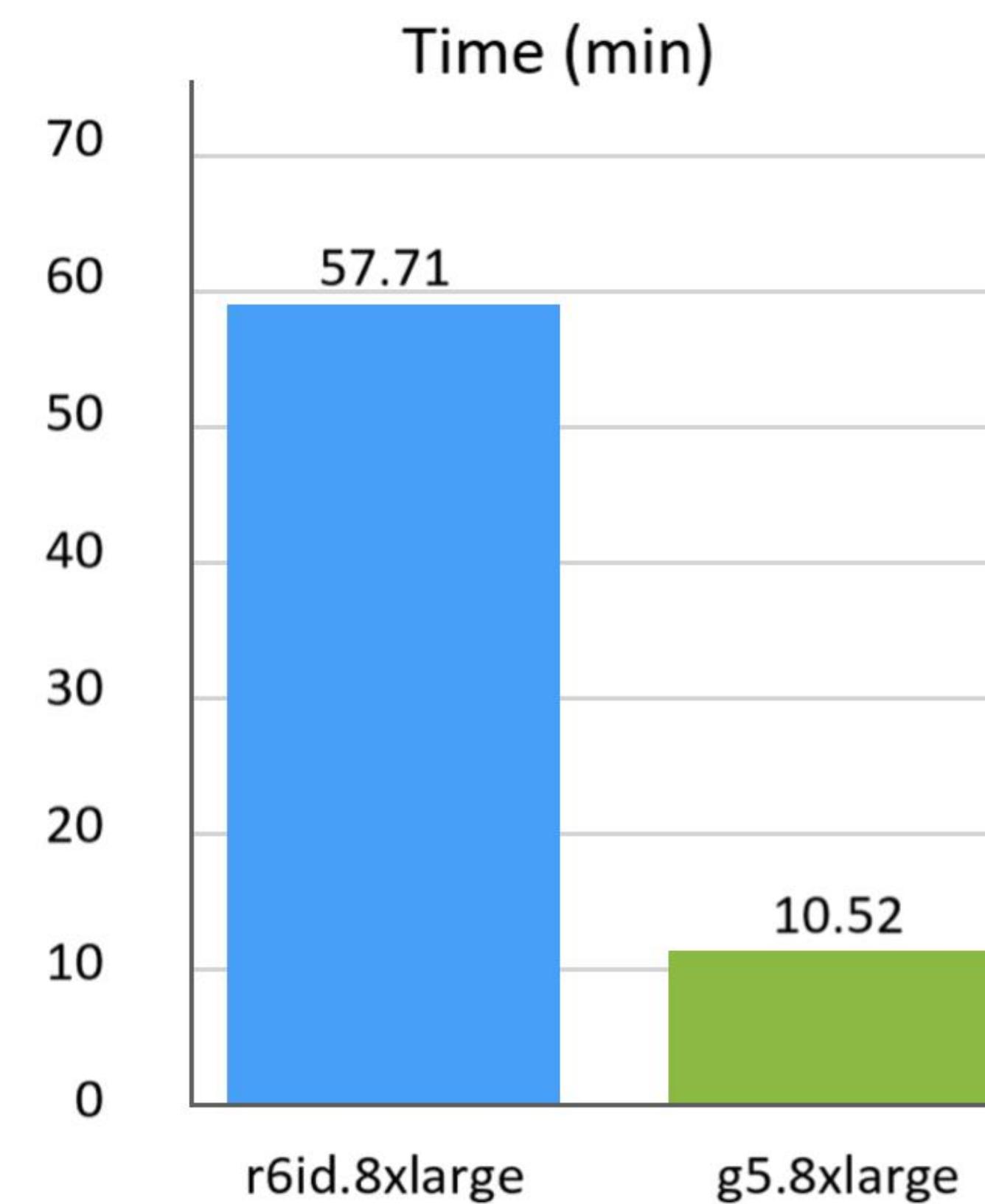
Accelerated Apache Spark

Zero code change acceleration for Spark DataFrames and SQL



- RAPIDS operates as a software plugin
- Automatically accelerates supported operations (with CPU fallback if needed)
- Works with Spark standalone, YARN clusters, Kubernetes clusters
- Compatible with PySpark, SparkR, Java, Scala and other DataFrame-based APIs
- How to run
 - Add jar to classpath and set spark.plugins config

NVIDIA Decision Support Benchmark 3TB (Public Cloud)



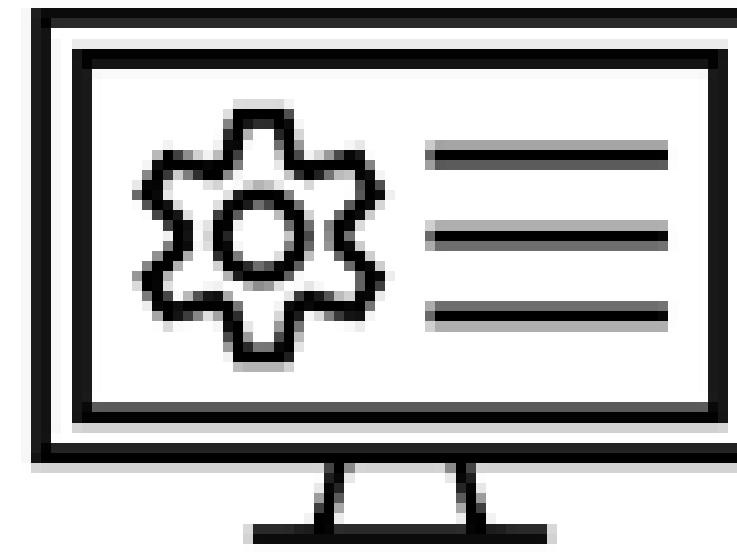
5.5x faster



80% cost savings

Simple Steps to Savings

Guided by Accelerated Spark Analysis Tool



Qualification

Identify and estimate the cost savings and acceleration potential of your Spark workloads based on an analysis of the log files.

Bootstrap

Get recommended configuration parameters for GPU acceleration on your Spark cluster. Use them for initial run.

Profiling / Tuning

See optimized configuration per application based on the initial (bootstrap) job run.

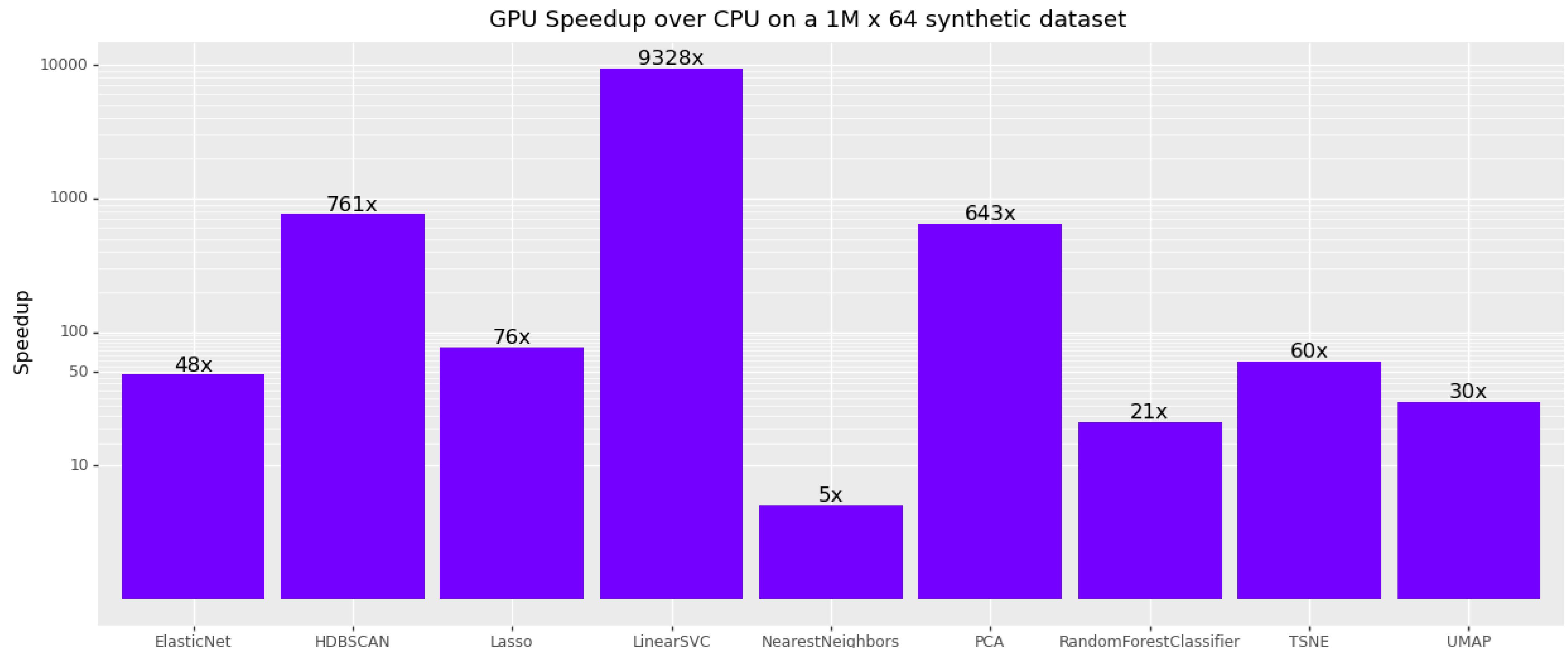
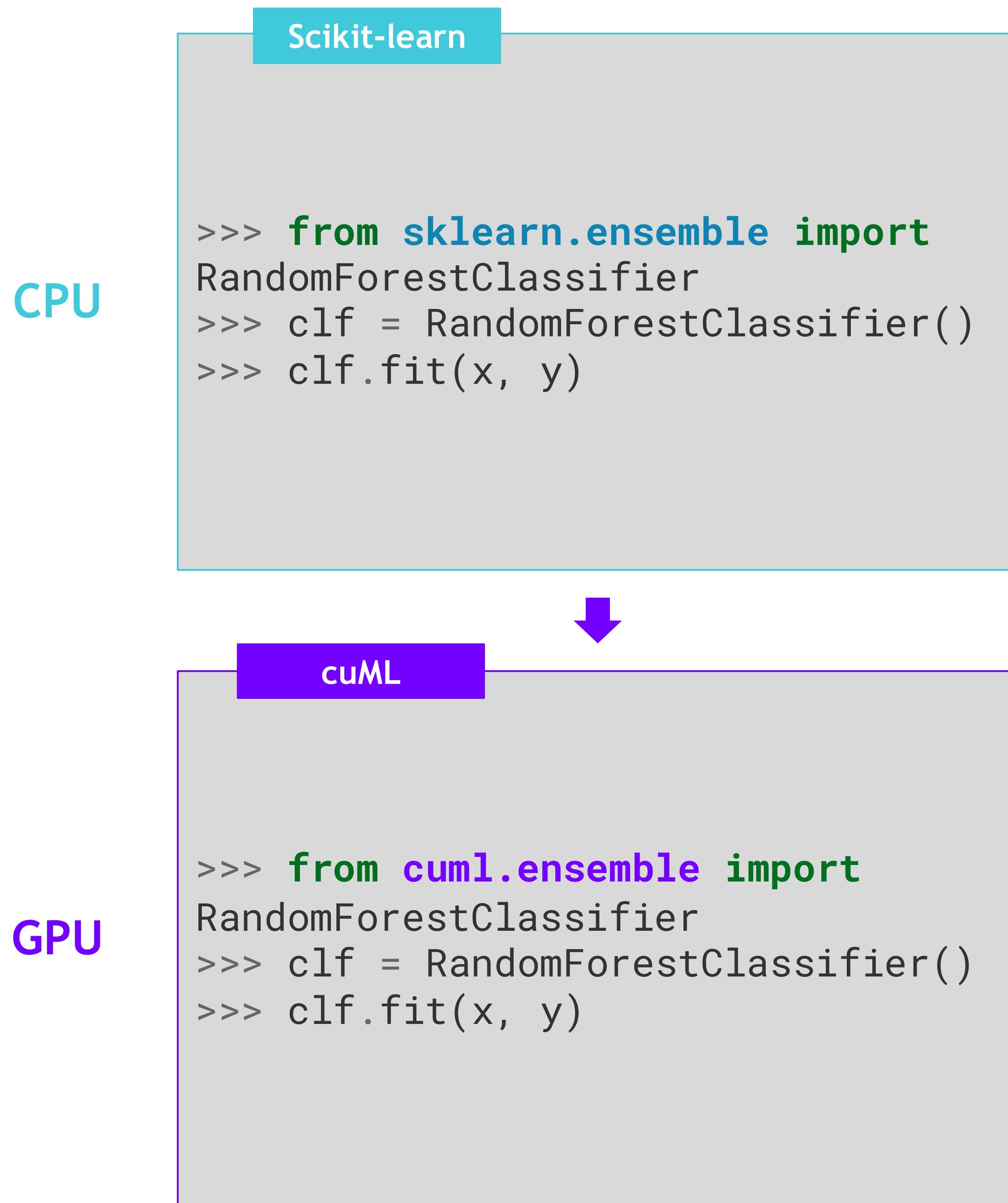
Machine Learning

cuML

cuML

Accelerated machine learning with a scikit-learn API

50+ GPU-Accelerated Algorithms

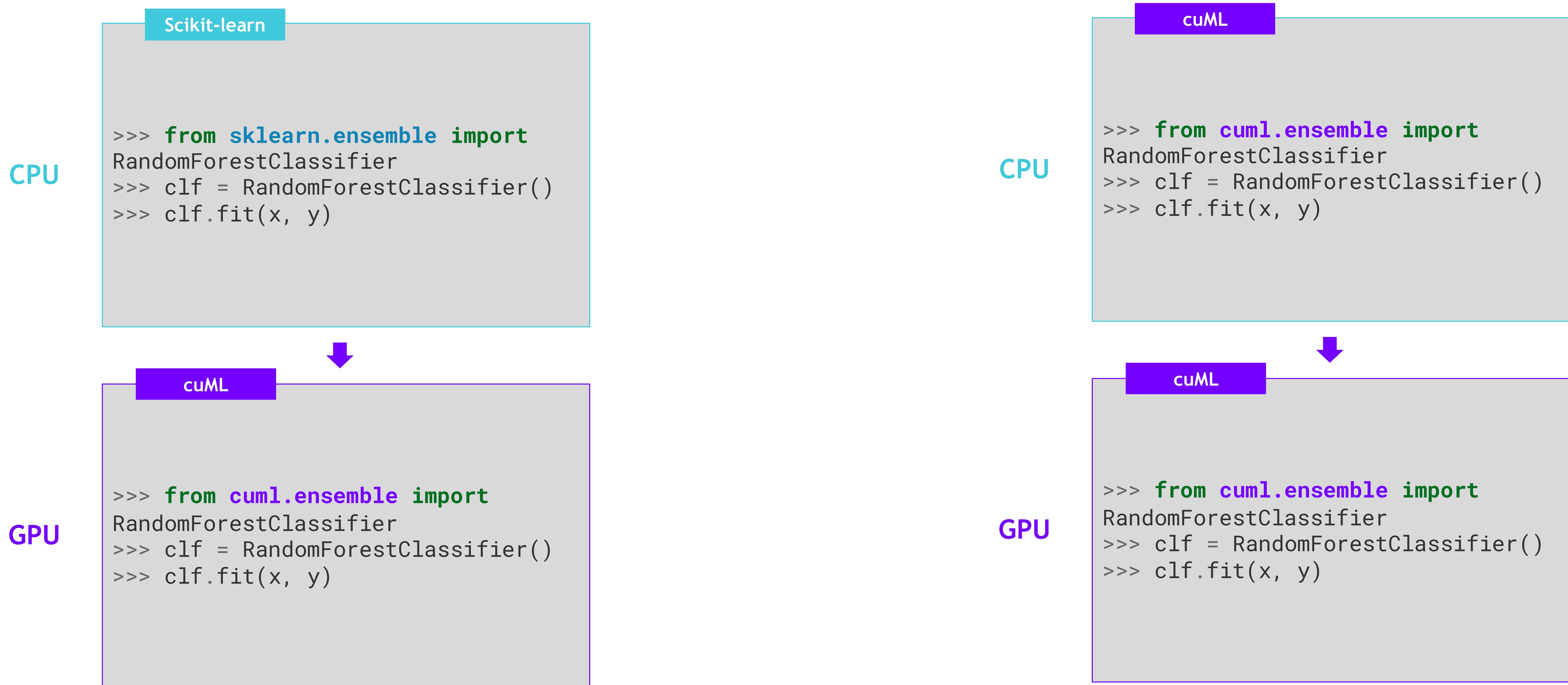


A100 GPU vs. AMD EPYC 7642 (96 logical cores)
cuML 23.04, scikit-learn 1.2.2, umap-learn 0.5.3

cuML-CPU

CPU version of cuML

Conda install -c rapidsai -c nvidia -c conda-forge cuml-cpu=23.10



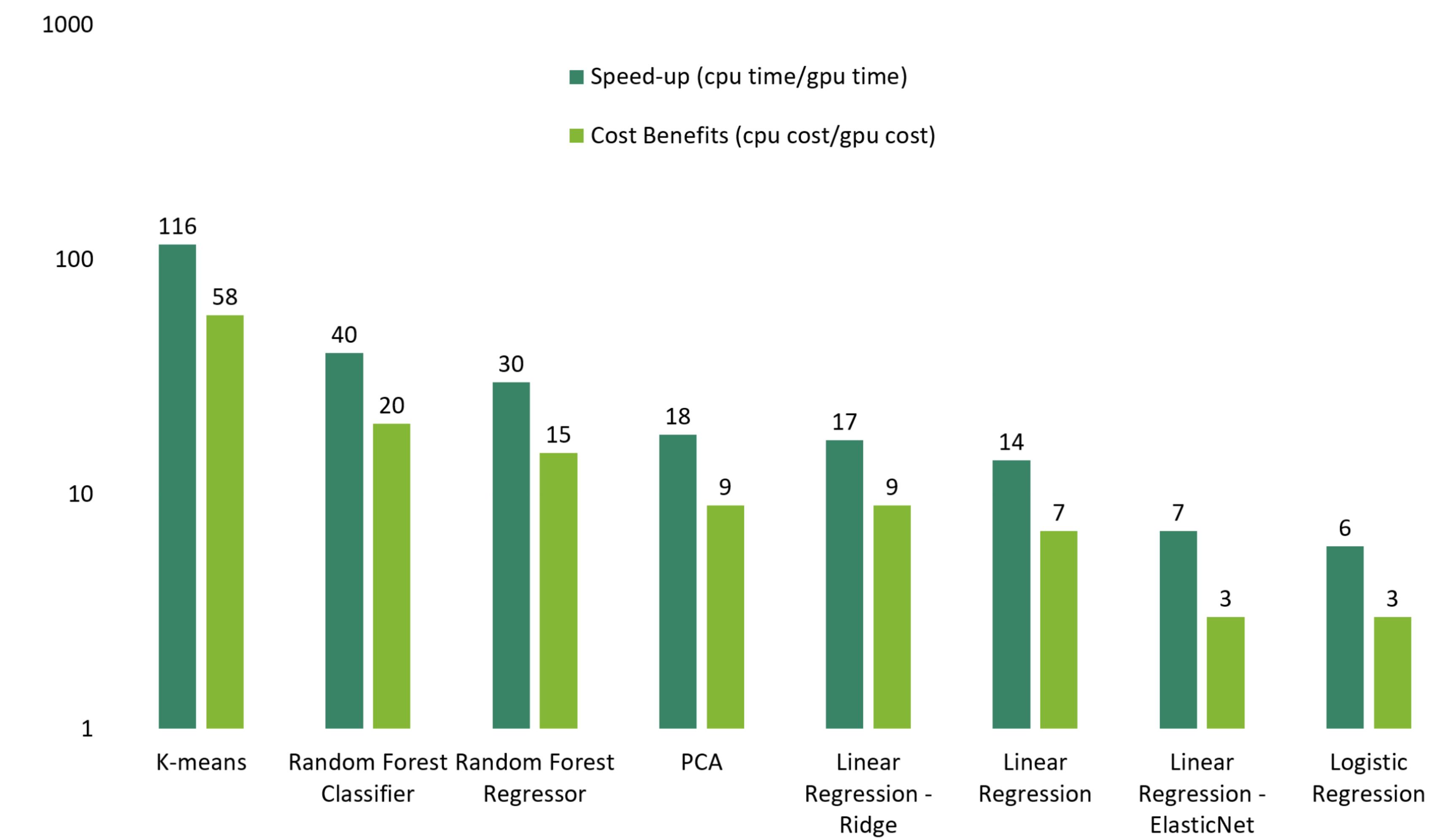
Accelerated Spark Machine Learning

SparkML built on tops of RAPIDS cuML allows no-change speedups for Spark

- Compatible with pyspark.ml DataFrame APIs
- Built on top of RAPIDS
- Requires no application code change
- Package import change:

```
● ● ●  
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)
```

Available Algorithms



Accelerated XGBoost

“XGBoost is All You Need” – Bojan Tunguz, 4x Kaggle Grandmaster

CPU

XGBoost

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

GPU

XGBoost

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

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

XGBoost 2.0 improvements:

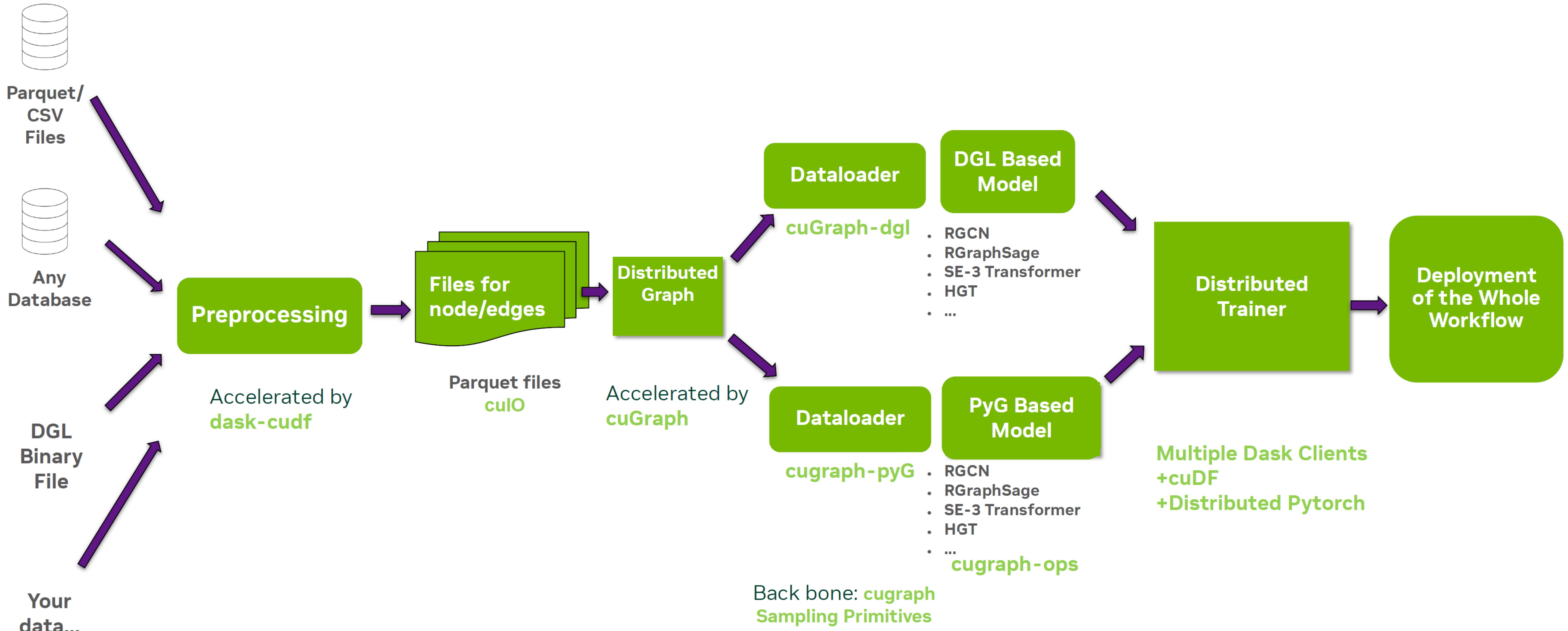
- Column-based split for federated learning
- Multi-target trees with vector-leaf outputs
- `approx` tree method is now GPU accelerated
- Improved external memory support

Graph Analytics

cuGraph

cuGraph + GNN Frameworks

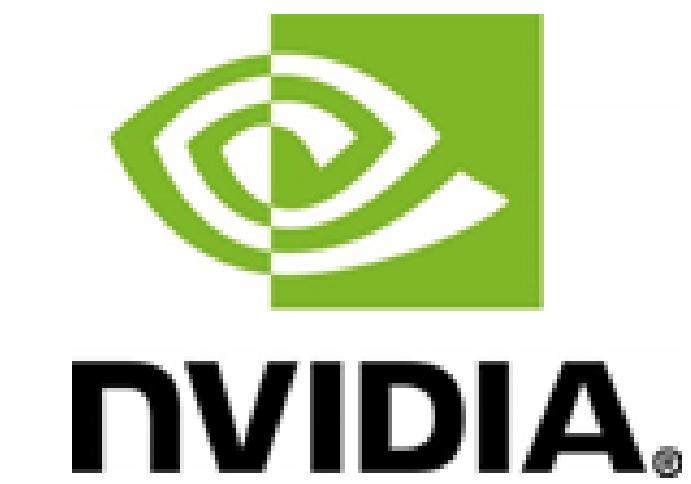
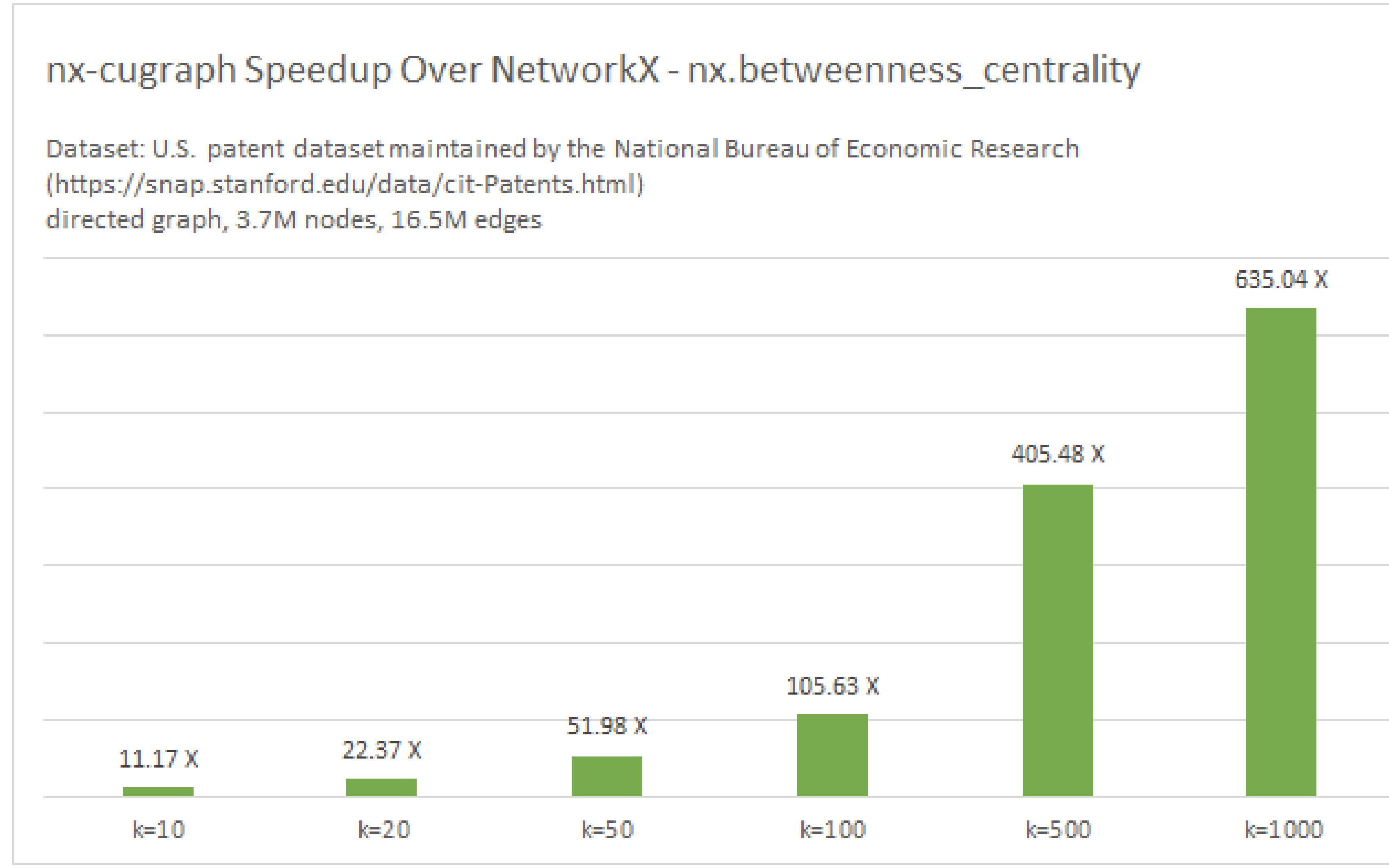
End-to-End Workflows



Accelerated NetworkX

Zero code change acceleration for NetworkX, powered by cuGraph

- Zero code change GPU-acceleration for NetworkX code
- Accelerates algorithms up to 600x, based on algorithm and graph size
- Support for 60 popular graph algorithms and growing
- Fallback to using CPU NetworkX for any unsupported algorithms



```
import pandas as pd
import networkx as nx

url = "https://data.rapids.ai/cugraph/datasets/cit-Patents.csv"
df = pd.read_csv(url, sep=" ", names=["src", "dst"], dtype="int32")
G = nx.from_pandas_edgelist(df, source="src", target="dst")

%time result = nx.betweenness_centrality(G, k=10)
```

```
user@machine:/# ipython bc_demo.ipynb
CPU times: user 7min 38s, sys: 5.6 s, total: 7min 44s
Wall time: 7min 44s

user@machine:/# NETWORKX_BACKEND_PRIORITY=cugraph ipython bc_demo.ipynb
CPU times: user 18.4 s, sys: 1.44 s, total: 19.9 s
Wall time: 20 s
```

pip install nx-cugraph-cu12 --extra-index-url <https://pypi.nvidia.com>
conda install -c rapidsai -c conda-forge -c nvidia nx-cugraph

Accelerated Inference

Triton Inference Server



- Part of NVIDIA AI platform
- Available with NVAIE
- Open-source software
- Standardizes and unifies AI model deployment and execution

Concurrent execution

- Enables high GPU utilization and throughput

Dynamic batching

- Increases throughput under latency constraints

Model analyzer

- Finds best configuration for models to maximize target function (possible targets: latency, throughput, memory footprint)

Multi-GPU Multi-node inference

- Inter-layer parallelism
- Intra-layer parallelism

RAPIDS Forest Inference Library backend

- Backend support for tree-based ML workloads

Amazon SageMaker integration

- Native integration of NVIDIA Triton Inference Server in Amazon SageMaker

RAPIDS FOREST INFERENCE LIBRARY (FIL)

Inference on Tree-Based ML Models

Native support

- XGBoost
- LightGBM
- Support with Treelite serialization
 - GradientBoostinClassifier / -Regressor
 - IsolationForest
 - RandomForestRegressor
 - ExtraTreesClassifier / -Regressor
 - Any tree model when exported to Treelite's checkpoint serialization format

Supported serialization formats

- XGBoost JSON (< 1.7)
- XGBoost Binary
- LightGBM Text
- Treelite binary checkpoint

By deploying a model on GPU more accurate models can be used in production without going over latency budget

Optimized Inference for Tree-based Models

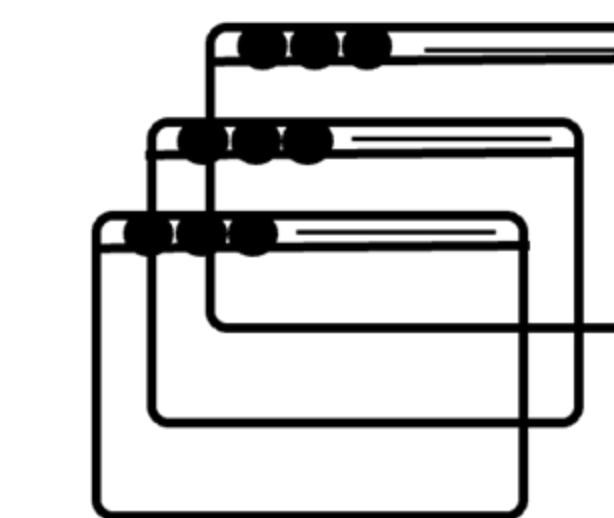
Forest Inference Library (FIL) Backend for both CPUs and GPUs



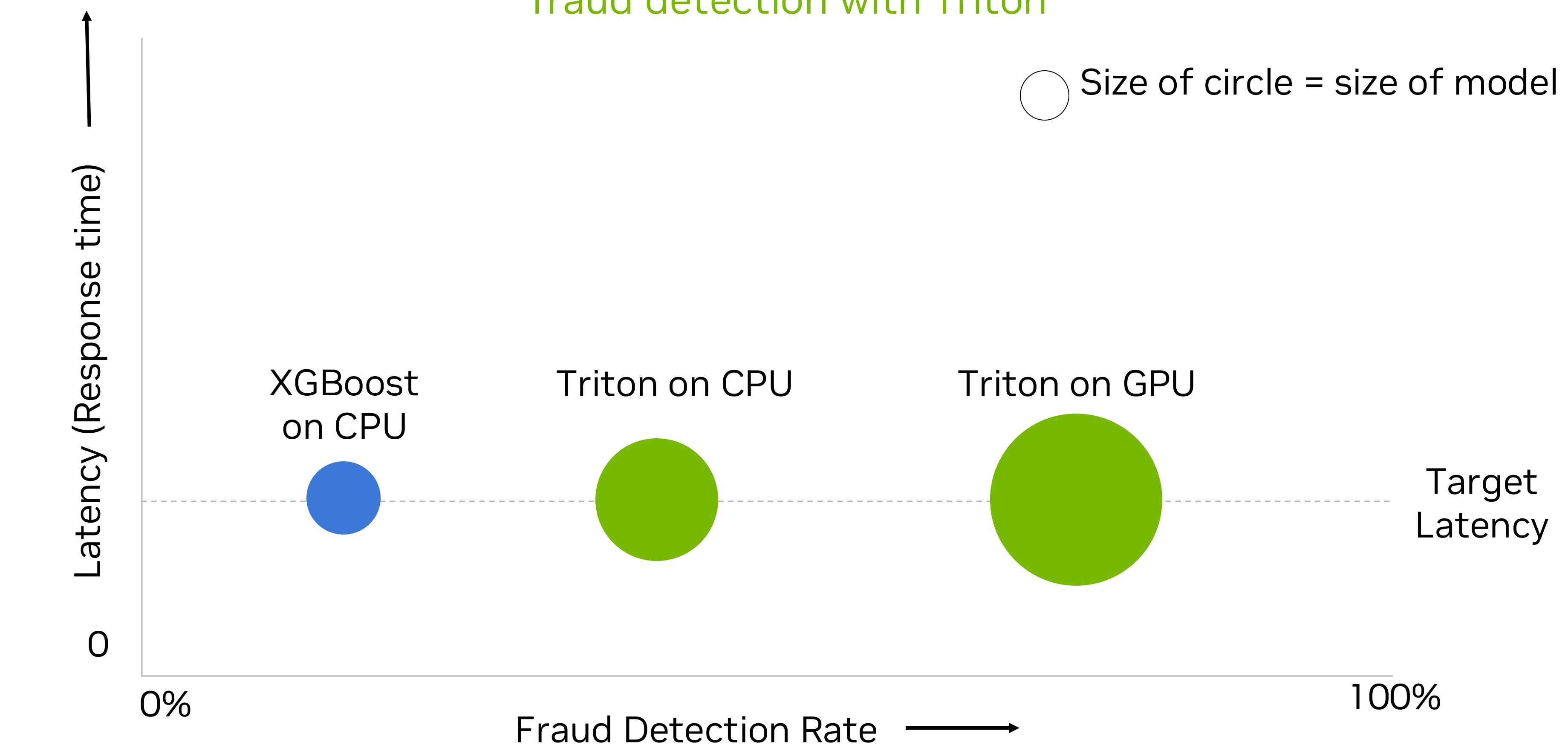
Real time inference
for models of all
sizes



Model explainability
with SHAP values



Run multiple tree-based
models concurrently on
the same CPU/GPU



dmlc
XGBoost

scikit
learn
RandomForest

LightGBM

RAPIDS
cuML RandomForest

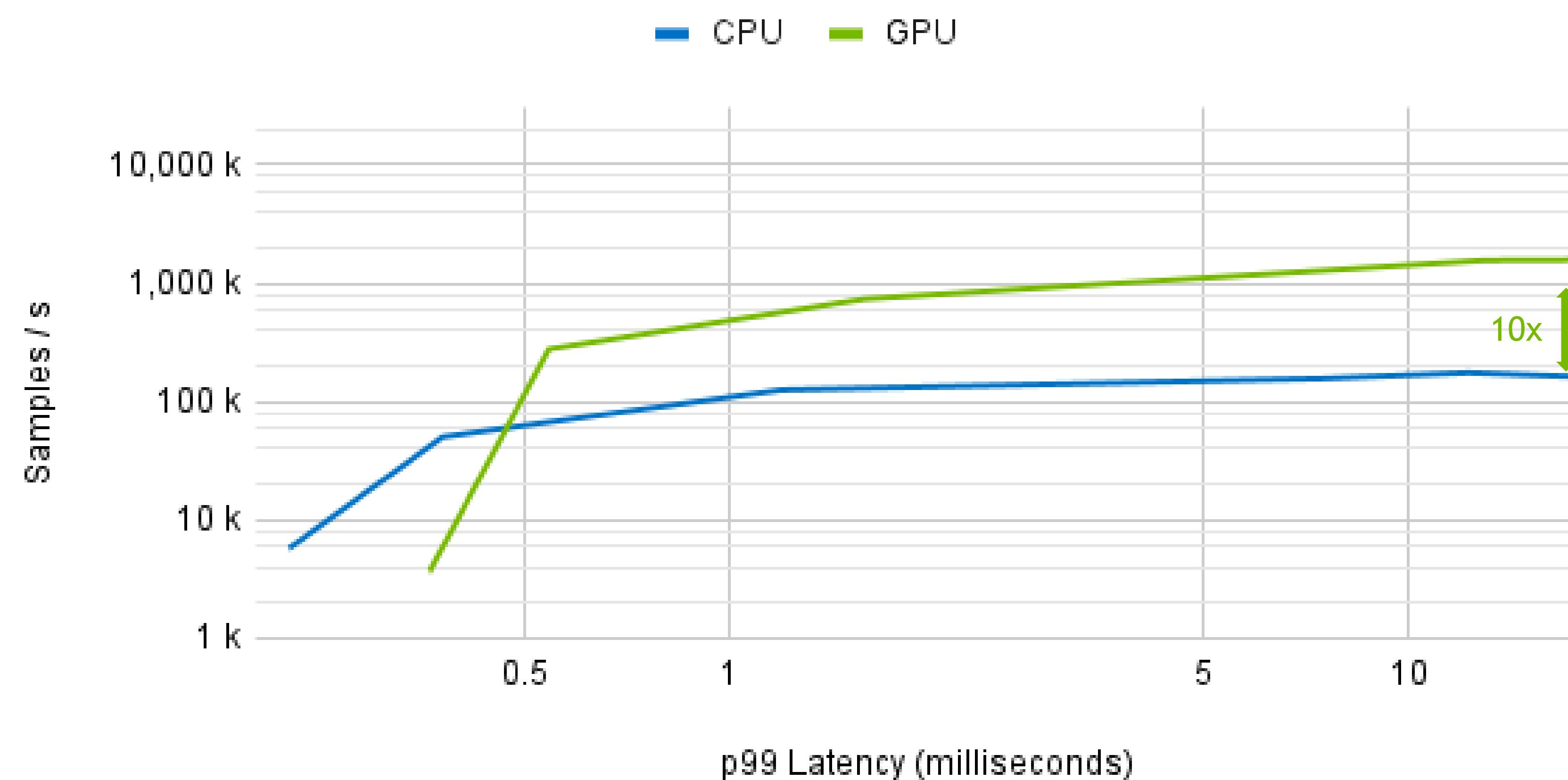
[Blog and Notebook](#)
[NVIDIA Launchpad lab](#)

Optimizing Throughput and Latency with Triton and FIL

Get maximum throughput for your latency budget with Triton

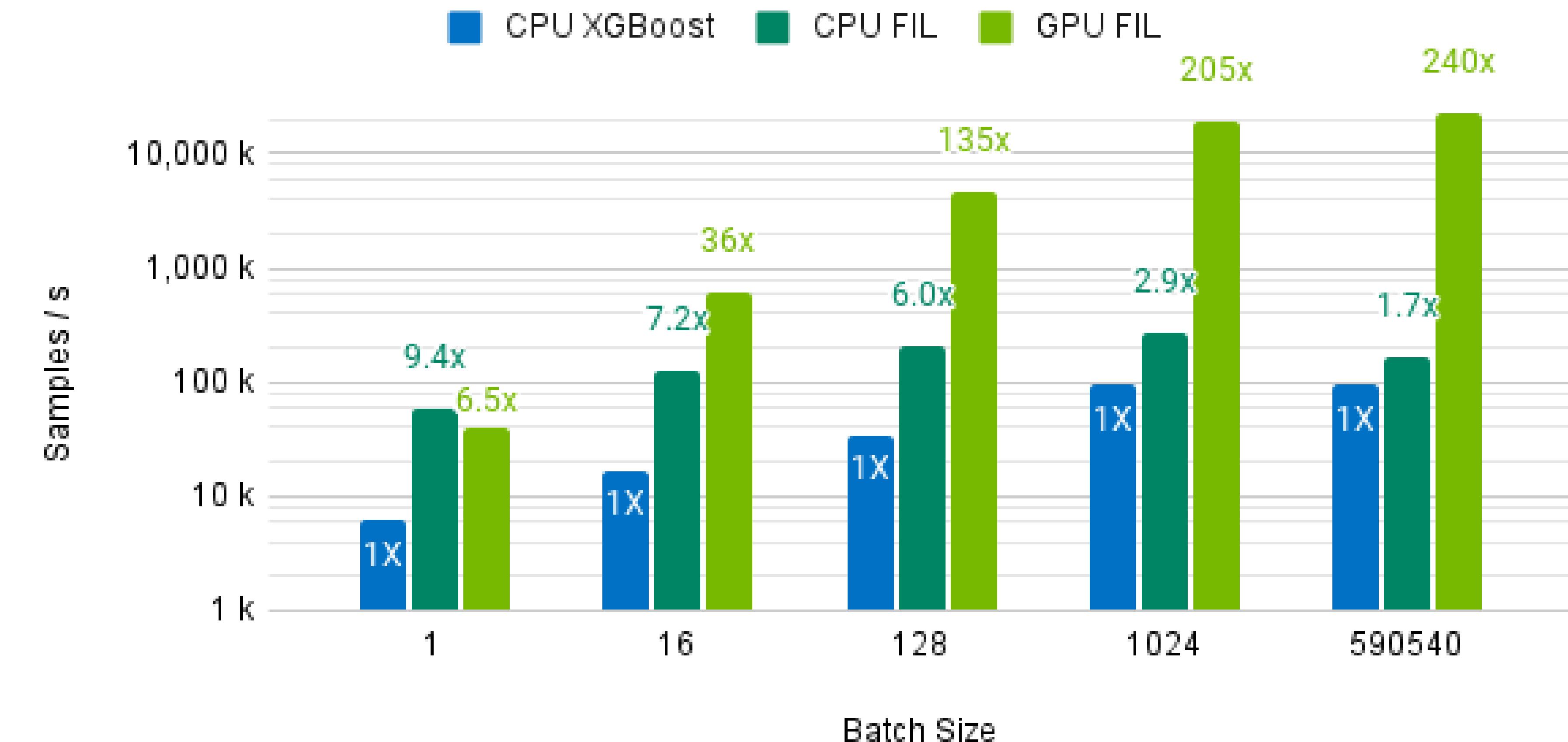
Throughput vs p99 Latency With Triton

Log-Log Scale



Throughput vs Batch Size for XGBoost/FIL Only

Log Scale



Throughput/latency curve for serving a ~200K node XGBoost model over GRPC with Triton

Performance of FIL library relative to CPU XGBoost

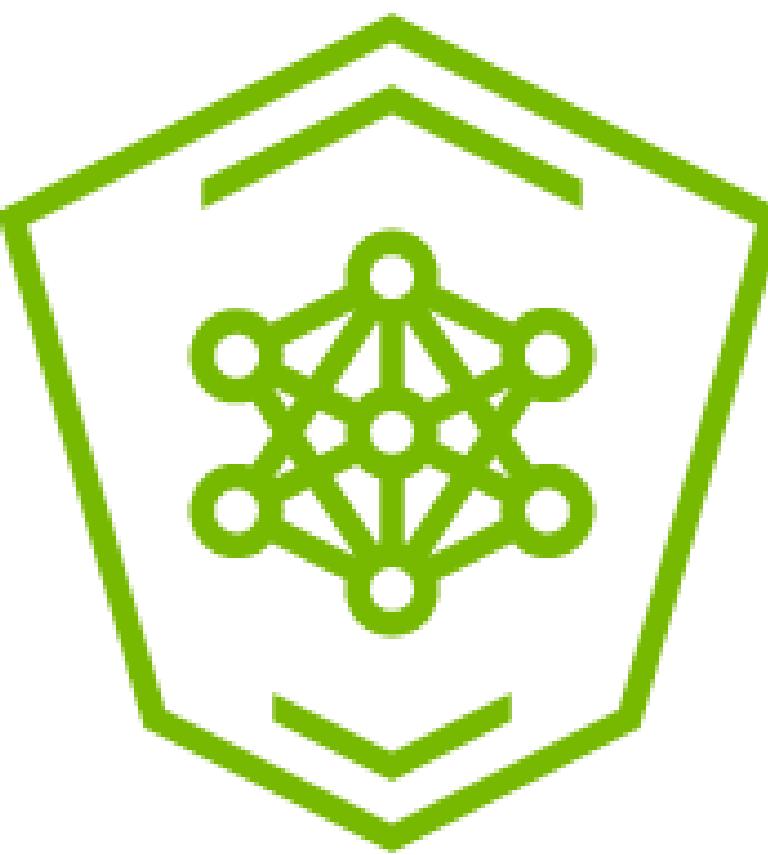
GPU benchmarks run on a GCP a2 instance with single A100; CPU benchmarks run on a GCP n2-standard-16
Max model depth: 10, Features: 393, Trees: 1000
Model trained on IEEE Kaggle Fraud Detection dataset



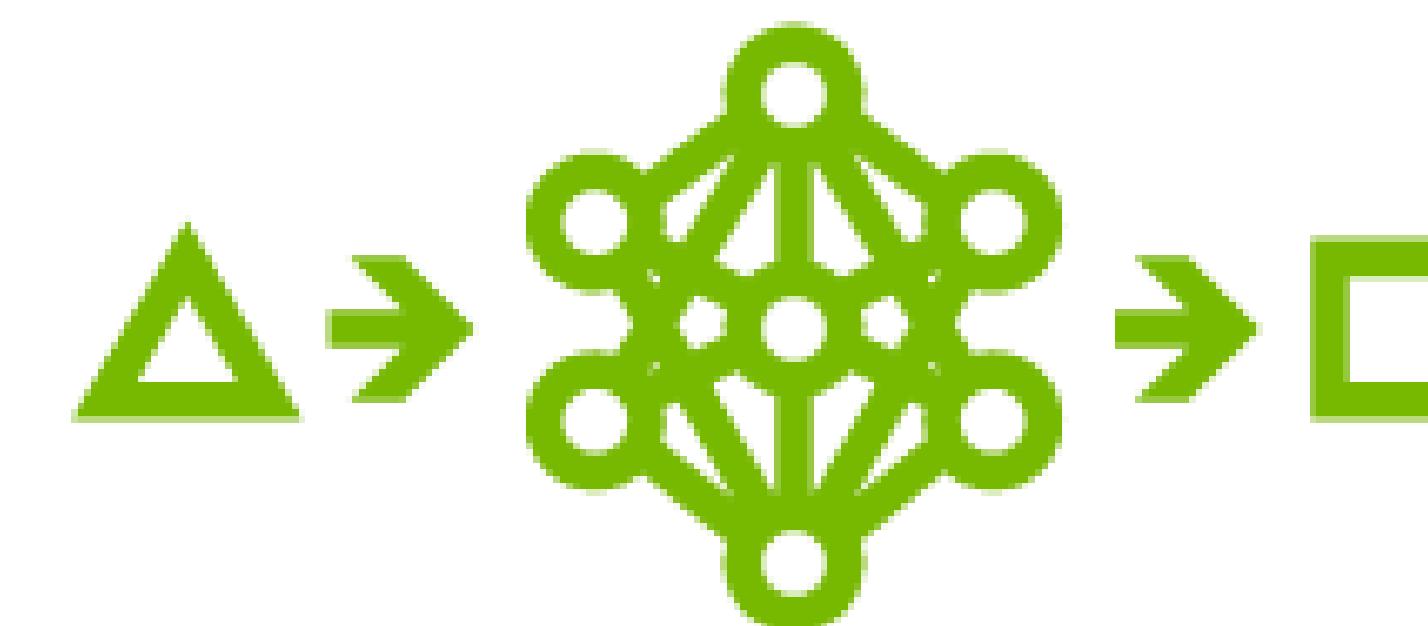
NVIDIA Software Stack

Overcome Challenges in Cybersecurity

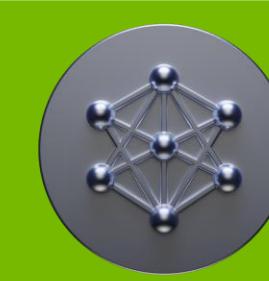
Generative AI unlocks new possibilities in modern cyber defense



NVIDIA Morpheus



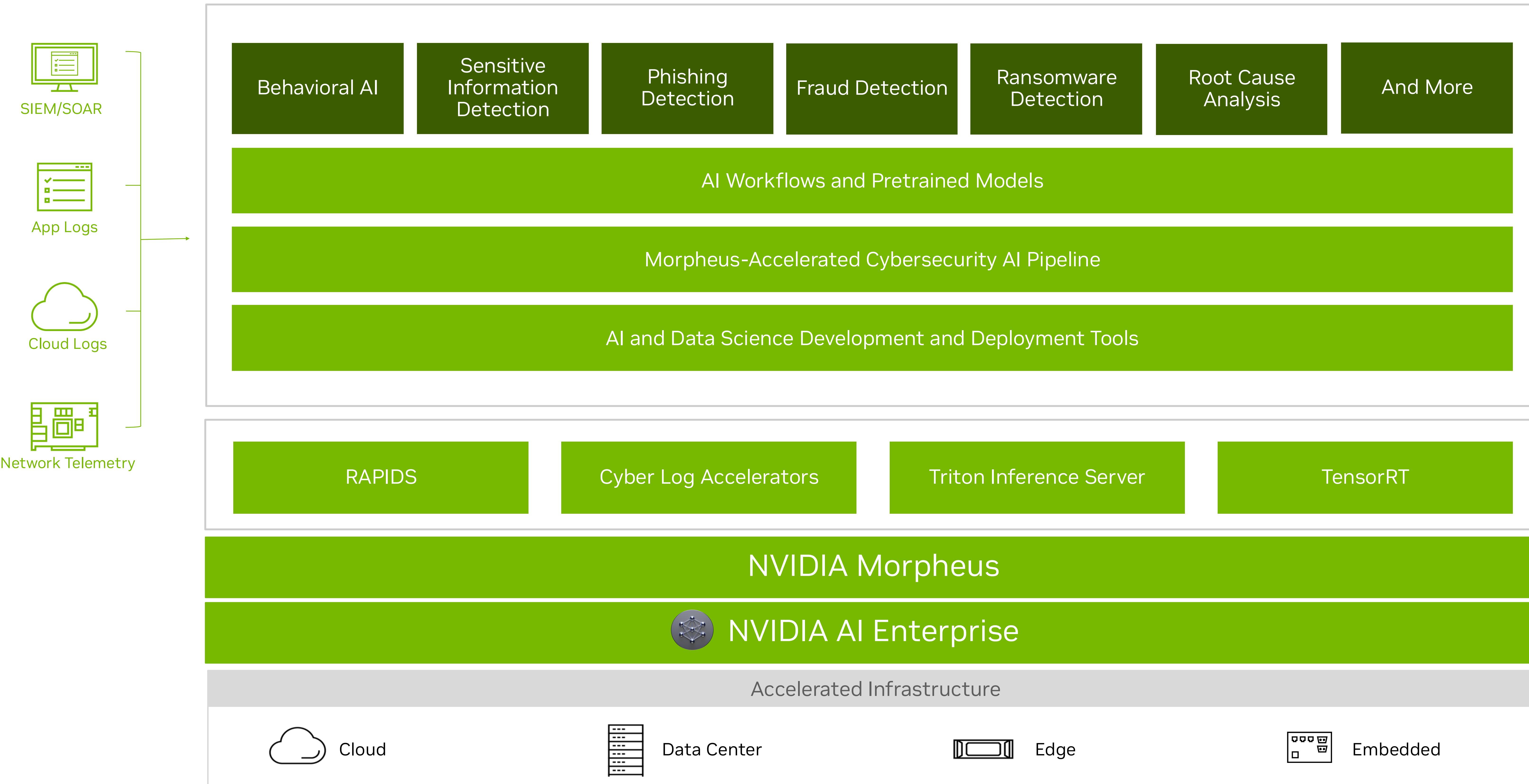
NVIDIA NeMo



NVIDIA AI Enterprise

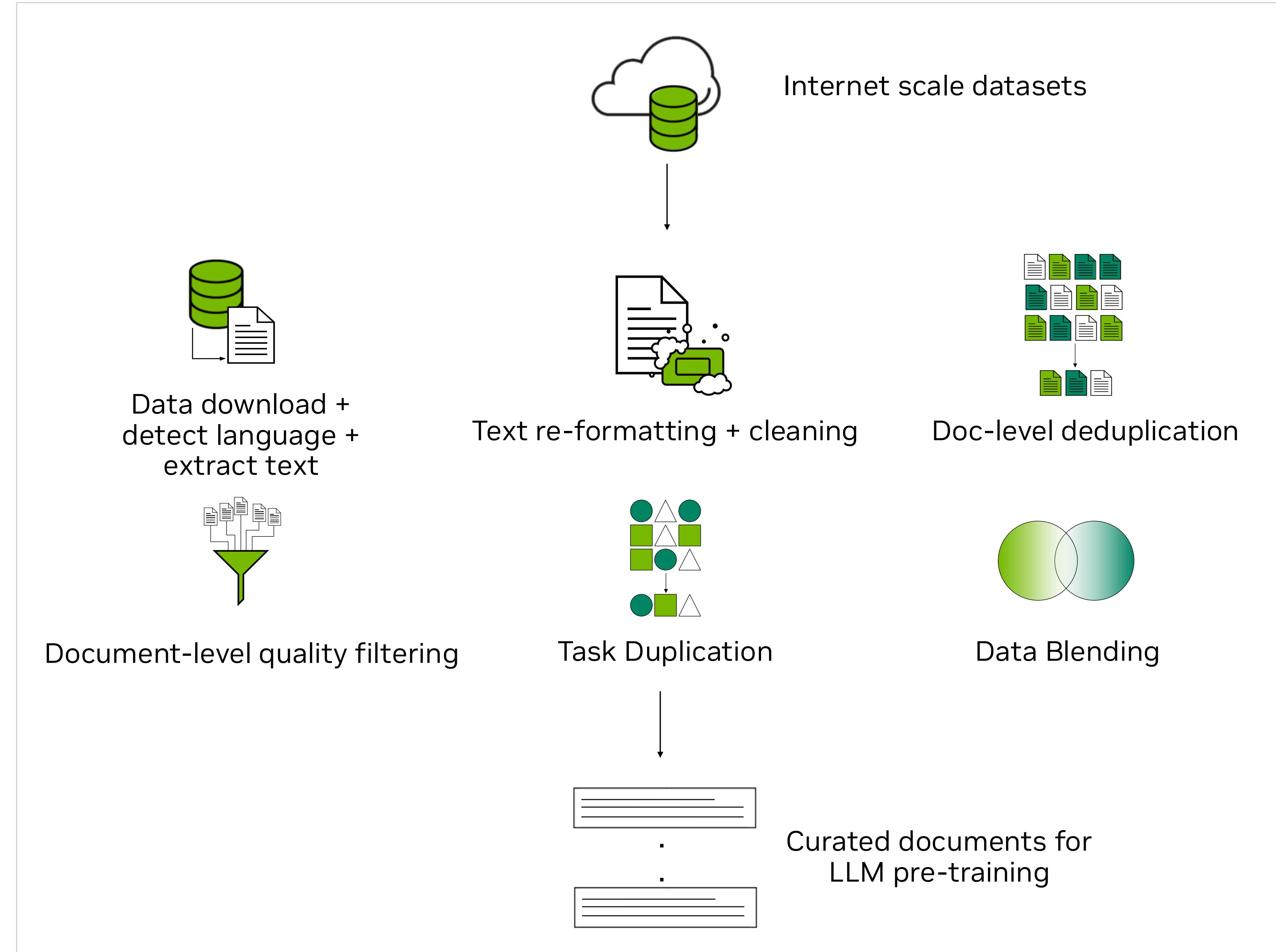
NVIDIA Morpheus Cybersecurity AI Framework

Automated, real-time threat detection at scale



NeMo Curator

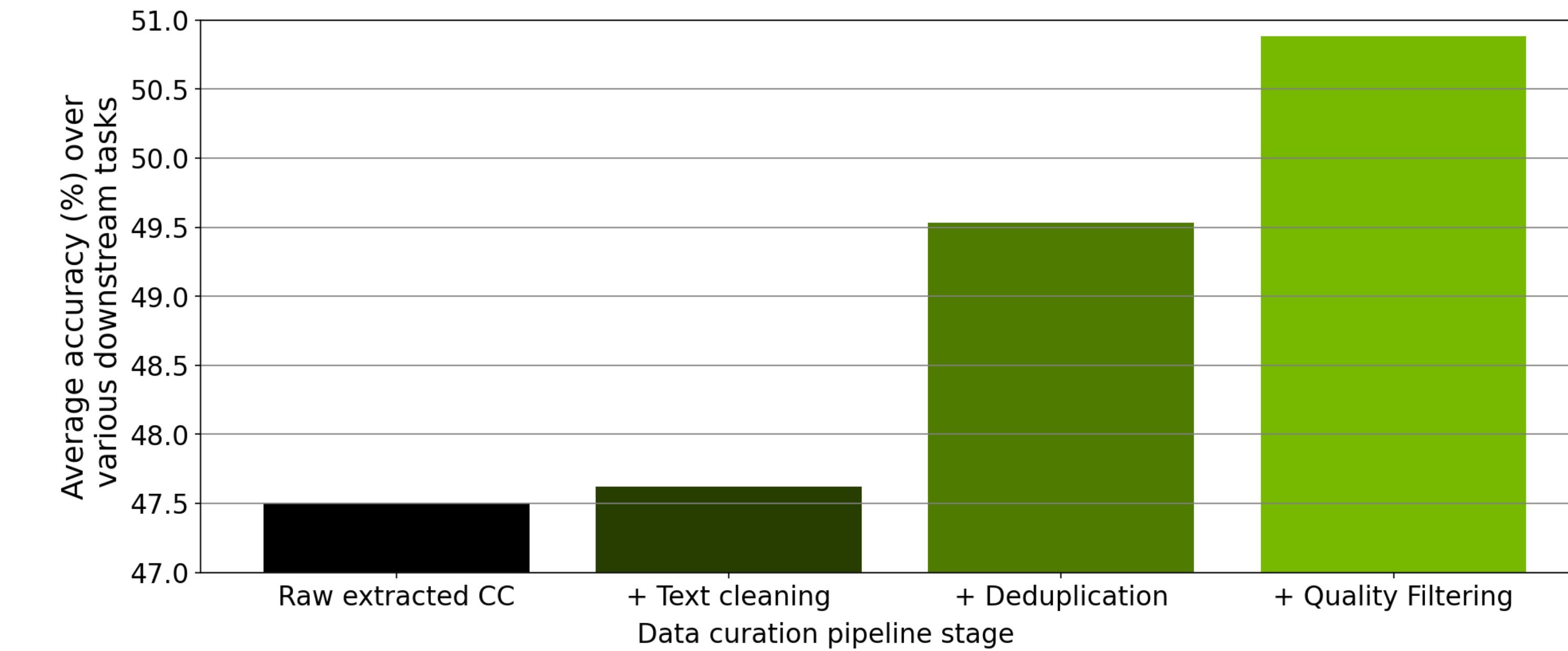
Python library of scalable data-mining modules for curating language data for training LLMs



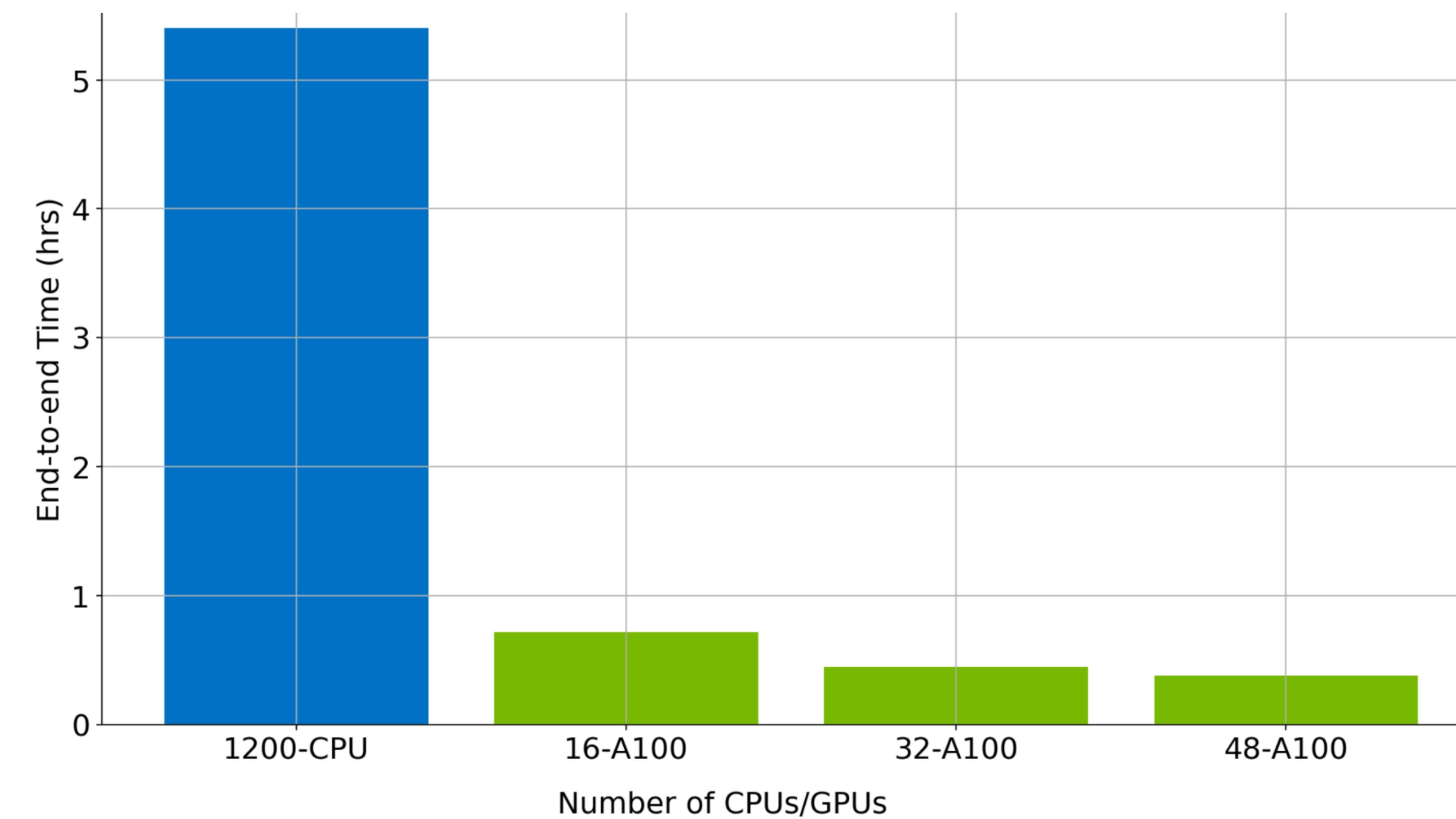
<https://github.com/NVIDIA/NeMo-Curator>

<https://catalog.ngc.nvidia.com/orgs/nvidia/containers/nemo>

Downstream task performance

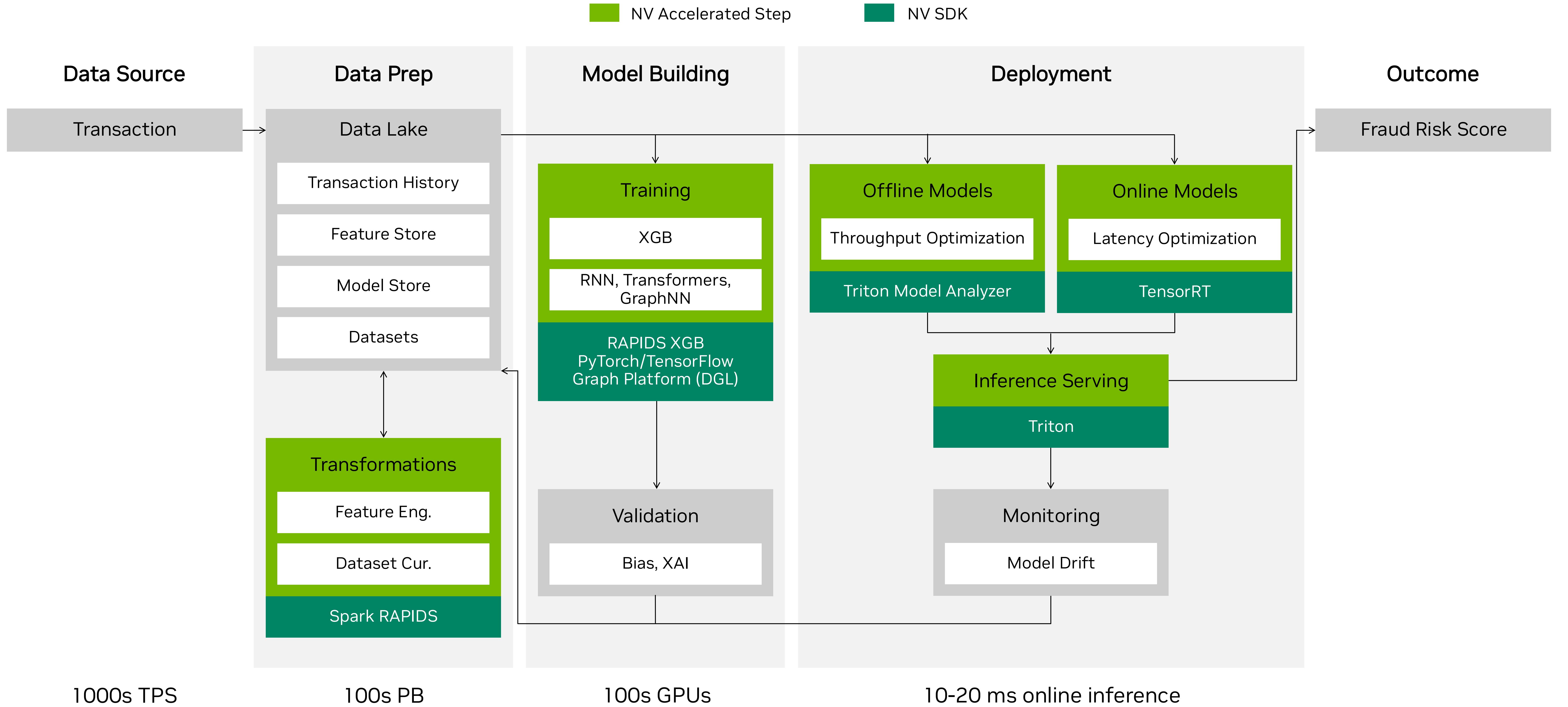


Fuzzy Deduplication of 500 GB of Text Data



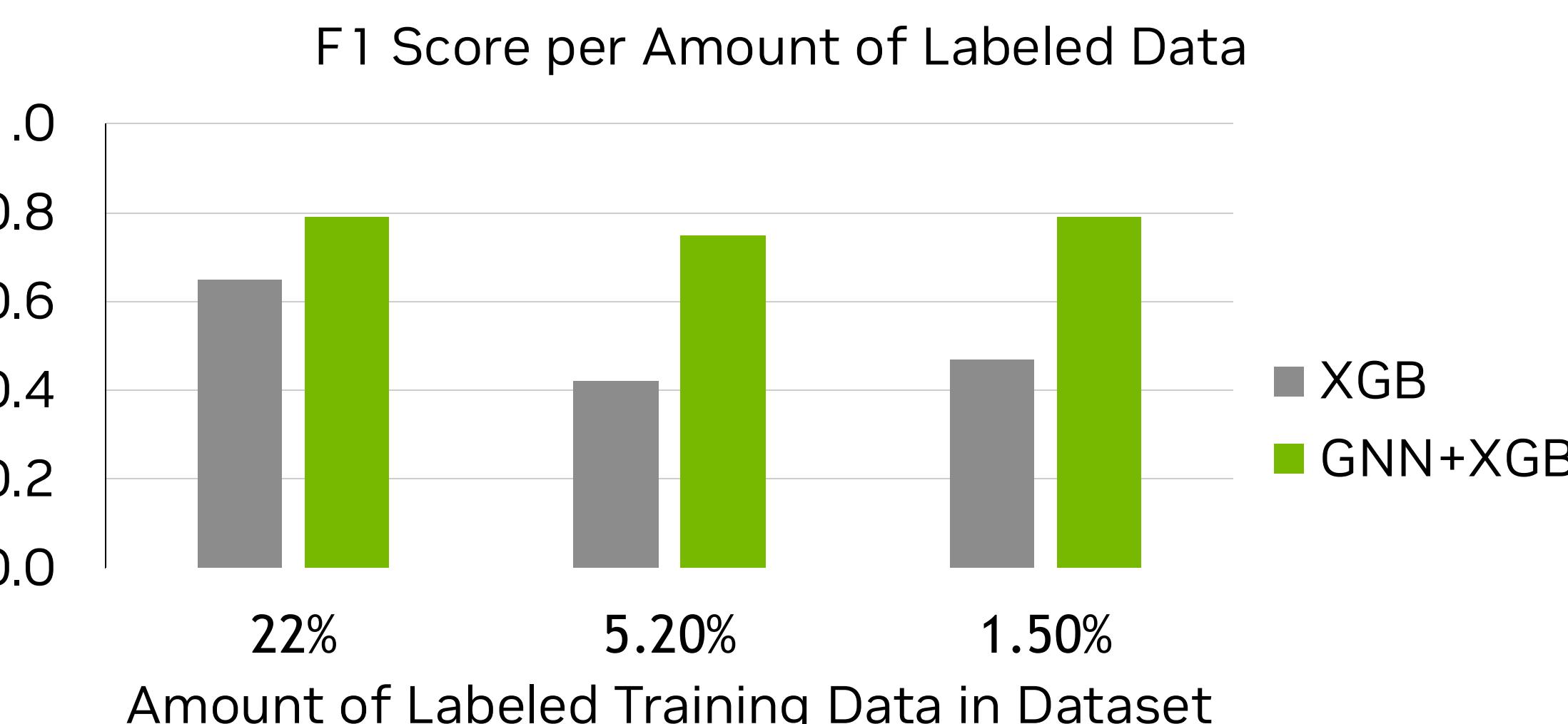
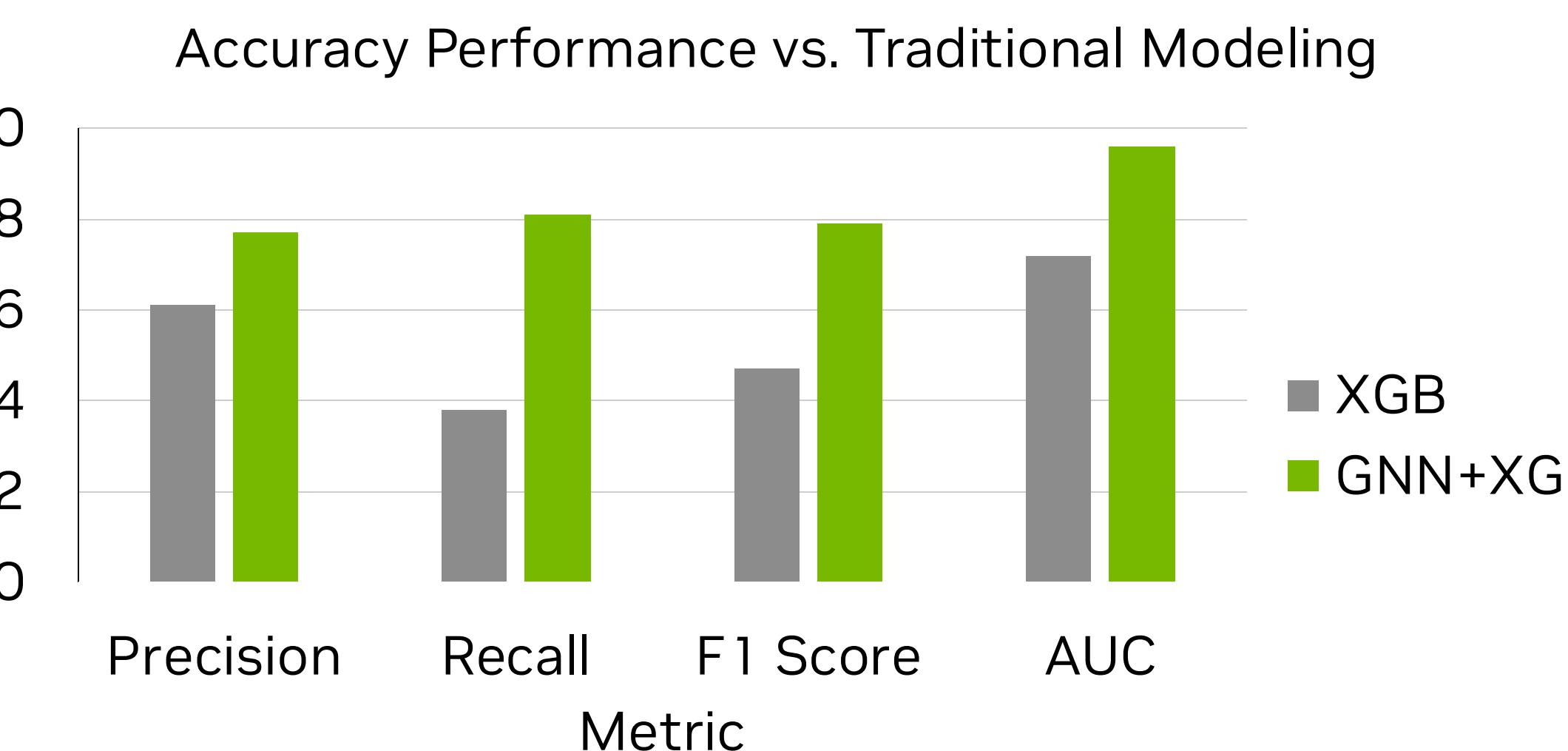
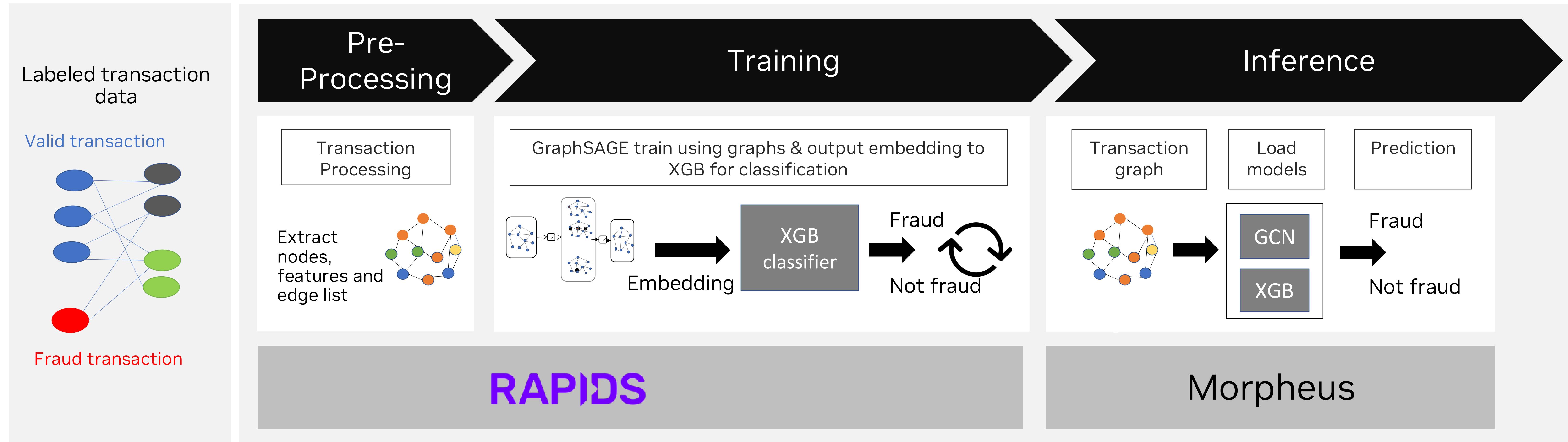
Accelerating Fraud Detection

Transaction Fraud Detection Workflow Example



FRAUD DETECTION

Morpheus Use Case Using GNNs and XGBoost for Rapid Detection of Fraudulent Transactions



How to Get Started

RAPIDS GitHub

Get started with the RAPIDS GitHub repository

<https://github.com/rapidsai>

The screenshot shows the GitHub repository page for 'rapidsai'. At the top, there's a navigation bar with links for Overview, Repositories (126), Projects (16), Packages, and People (49). Below the navigation is a purple header with the RAPIDS logo and the text 'Open GPU Data Science' and 'Verified'. It shows 1.4k followers and links to the RAPIDS website and Twitter. A 'Pinned' section displays six repositories: cuDF (cuDF - GPU DataFrame Library), rmm (RAPIDS Memory Manager), cuML (cuML - RAPIDS Machine Learning Library), cuGraph (cuGraph - RAPIDS Graph Analytics Library), cuspatial (CUDA-accelerated GIS and spatiotemporal algorithms), and raft (RAFT contains fundamental widely-used algorithms and primitives for machine learning and information retrieval). Each pinned repository has a small description, language (C++ or Cuda), star count, and commit count. At the bottom, there's a search bar for 'Find a repository...' and filters for Type, Language, and Sort.

<https://rapids.ai/>

The screenshot shows the RAPIDS website homepage. At the top, there's a purple header with the RAPIDS logo and the text 'GPU Accelerated Data Science'. Below the header is a 'QUICK START' button. The main content area features a large purple banner with the text 'DATA SCIENCE AT THE SPEED OF THOUGHT'. Below the banner, there are three main sections: 'WHAT IS RAPIDS', 'WHY USE RAPIDS', and 'OPEN SOURCE ECOSYSTEM'. Each section has a brief description and a 'Jump to About Section' link. There are also four purple callout boxes at the top of the main content area: 'NEW: CUDA NOW PRE-INSTALLED IN GOOGLE COLAB', 'NEW: CUVS VECTOR SEARCH', 'CATCH UP WITH GTC'24 ON-DEMAND', and 'RAPIDS 24.06 RELEASED'. The footer of the page includes the NVIDIA logo.

Google Colab

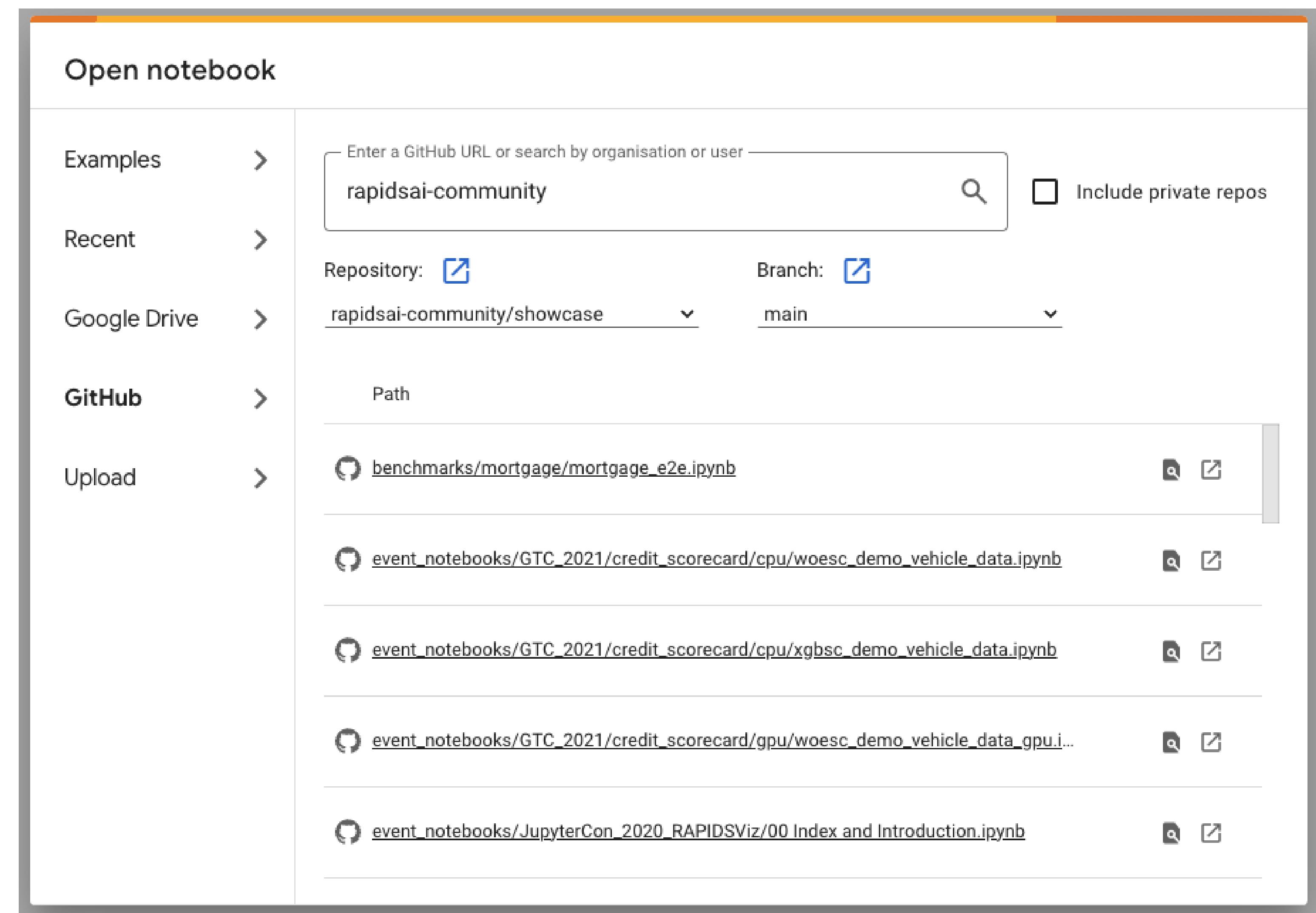
Deploy wherever you run

https://colab.research.google.com/github/rapidsai-community/showcase/blob/main/getting_started_tutorials/cudf_pandas_colab_demo.ipynb

The screenshot shows a Google Colab notebook titled "cudf_pandas_colab_demo.ipynb". The left sidebar contains a table of contents with several sections: "10 Minutes to RAPIDS cuDF's pandas accelerator mode (cudf.pandas)", "Verify your setup", "Download the data", "Analysis using Standard Pandas", "Using cudf.pandas", "Understanding Performance", "Behind the scenes: What's going on here?", "Using third-party libraries with cudf.pandas", and "Visualizing which states have more pickup trucks relative to other vehicles?". The main content area displays the first section, "10 Minutes to RAPIDS cuDF's pandas accelerator mode (cudf.pandas)". It includes a warning message about verifying setup, a code cell for running nvidia-smi, and a terminal output showing GPU information. The terminal output is as follows:

```
[ ] !nvidia-smi # this should display information about available GPUs
Wed May 8 17:01:36 2024
+-----+
| NVIDIA-SMI 535.104.05      Driver Version: 535.104.05    CUDA Version: 12.2 |
| GPU  Name        Persistence-M  Bus-Id      Disp.A  Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap | Memory-Usage | GPU-Util  Compute M. |
| |          %          /   |           /   |          /   |          /   |
| 0  Tesla T4           Off  00000000:00:04.0 Off   0% Default N/A |
| N/A  46C   P8       9W / 70W |    0MiB / 15360MiB |      0%     Default |
+-----+
| Processes:                               |
+-----+
```

<https://colab.research.google.com/github/rapidsai-community/showcase/blob/>



NVIDIA Launchpad

Deploy wherever you run

<https://www.nvidia.com/en-us/launchpad/>

NVIDIA LaunchPad Collections ▾

Introduction Choose a Lab How to Join Features Get Started

Request Access to Hands-On Labs

Ready for a hands-on experience with NVIDIA software solutions? Pick a lab to begin your journey.

Filters: Industry Technologies Categories Products Use Cases **rapids** X

11 result(s) found



Accelerating Apache Spark with Zero Code Changes

Best for: AI practitioner
Included products: NVIDIA AI Enterprise, NVIDIA RAPIDS, NVIDIA-Certified Systems
Included technologies: Apache Spark, NVIDIA RAPIDS

[Take This Lab >](#)



Build AI-Based Cybersecurity Solutions

Best for: AI practitioner
Included products: NVIDIA AI Enterprise, NVIDIA Morpheus, NVIDIA-Certified Systems
Included technologies: Cyber Log Accelerators, NVIDIA RAPIDS, NVIDIA TensorRT, NVIDIA Triton Inference Server

[Take This Lab >](#)



Data Processing, Tokenization, and Sentiment Analysis

Best for: Data engineer, Data scientist
Included products: NVIDIA AI Enterprise, NVIDIA-Certified Systems
Included technologies: NVIDIA RAPIDS, NVIDIA Triton Inference Server, Pandas, PyTorch

[Take This Lab >](#)



Deploy Fraud Detection XGBoost Model with NVIDIA Triton

Best for: AI practitioner
Included products: NVIDIA AI Enterprise, NVIDIA RAPIDS, NVIDIA-Certified Systems
Included technologies: NVIDIA RAPIDS, NVIDIA Triton Inference Server

[Take This Lab >](#)

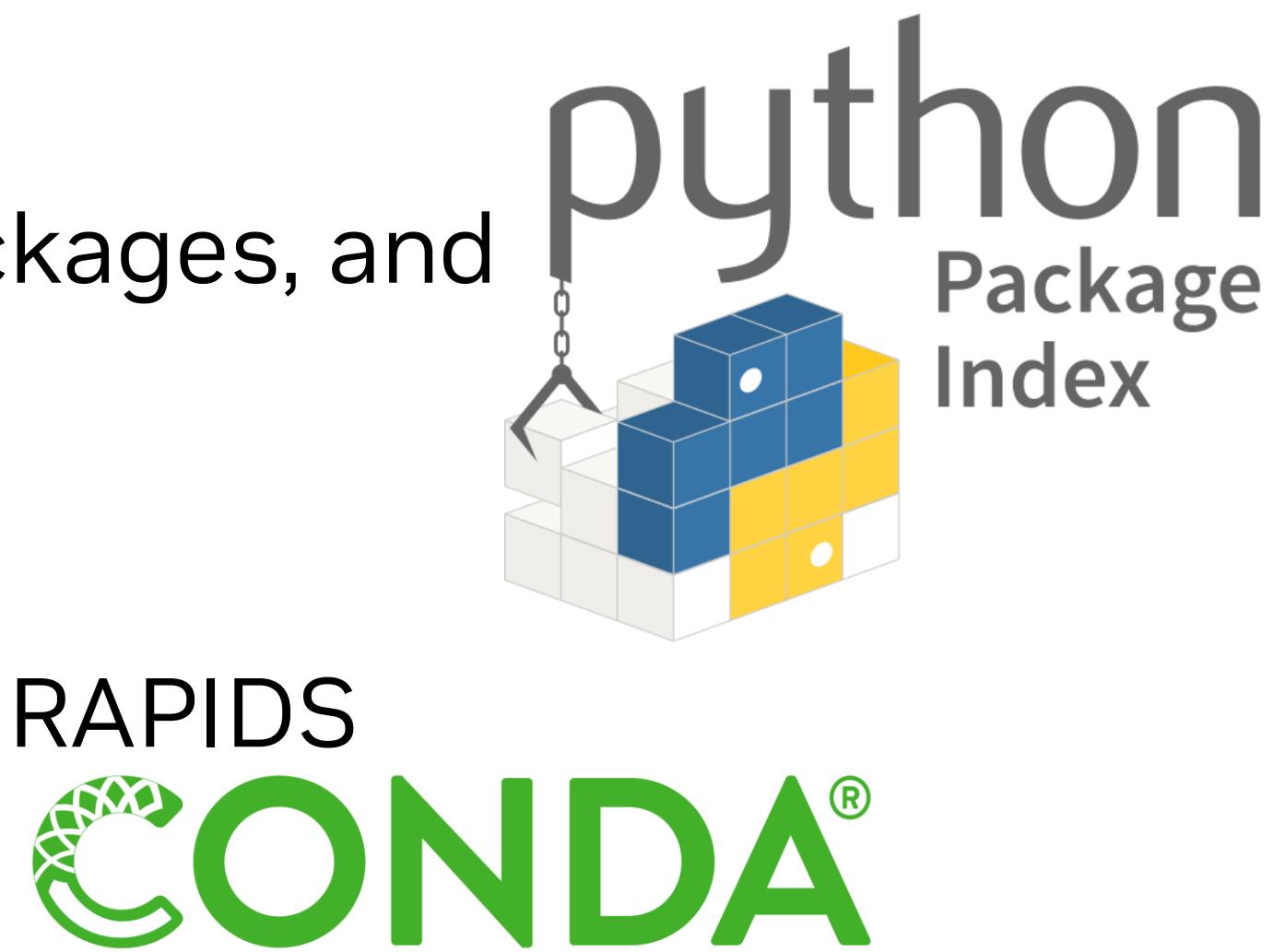


Smooth Installation and Packaging

Meeting a wide array of needs with easy to install and use distribution methods

Standard binary distributions

- One line installation methods on rapids.ai
 - *Conda*:
 - With conda-forge CUDA packages, and rapidsai RAPIDS packages
 - *Pip* installation
 - Custom index to support all RAPIDS ecosystem.



Containers

- Thin CUDA containers
- *RAPIDS* and other *application-level* containers
- *Devcontainers* for developers wanting to contribute
- *NVIDIA AI Enterprise*



Source code on GitHub

<https://github.com/rapidsai>

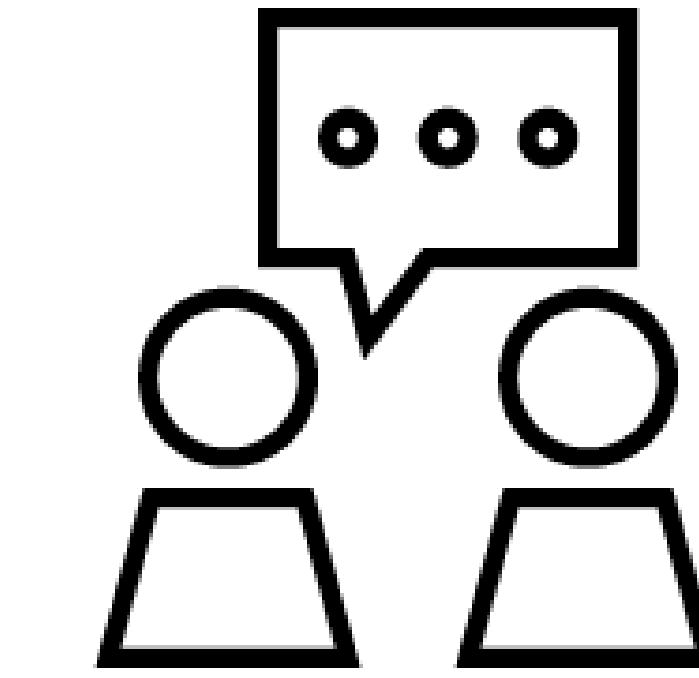
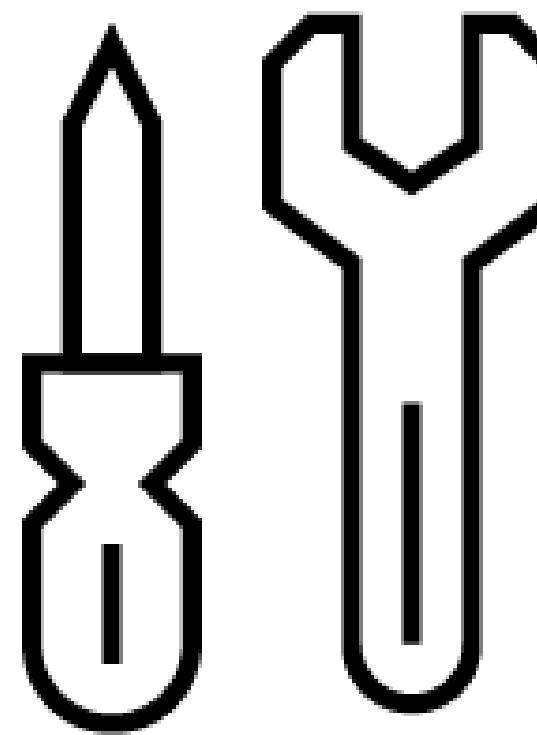
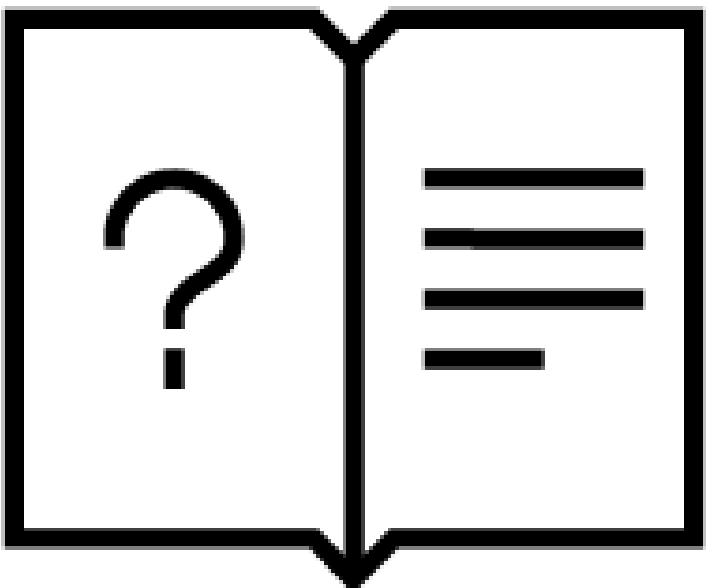


Cloud Deployment

In the cloud

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/>

RAPIDS

<https://rapids.ai>

*Thank
you!*

renzmann@nvidia.com