# Accelerate ETL and Machine Learning in Apache Spark

Erik Ordentlich, Sameer Raheja | GTC | March 19, 2024

# RAPIDS Spark

- Data Processing Challenges
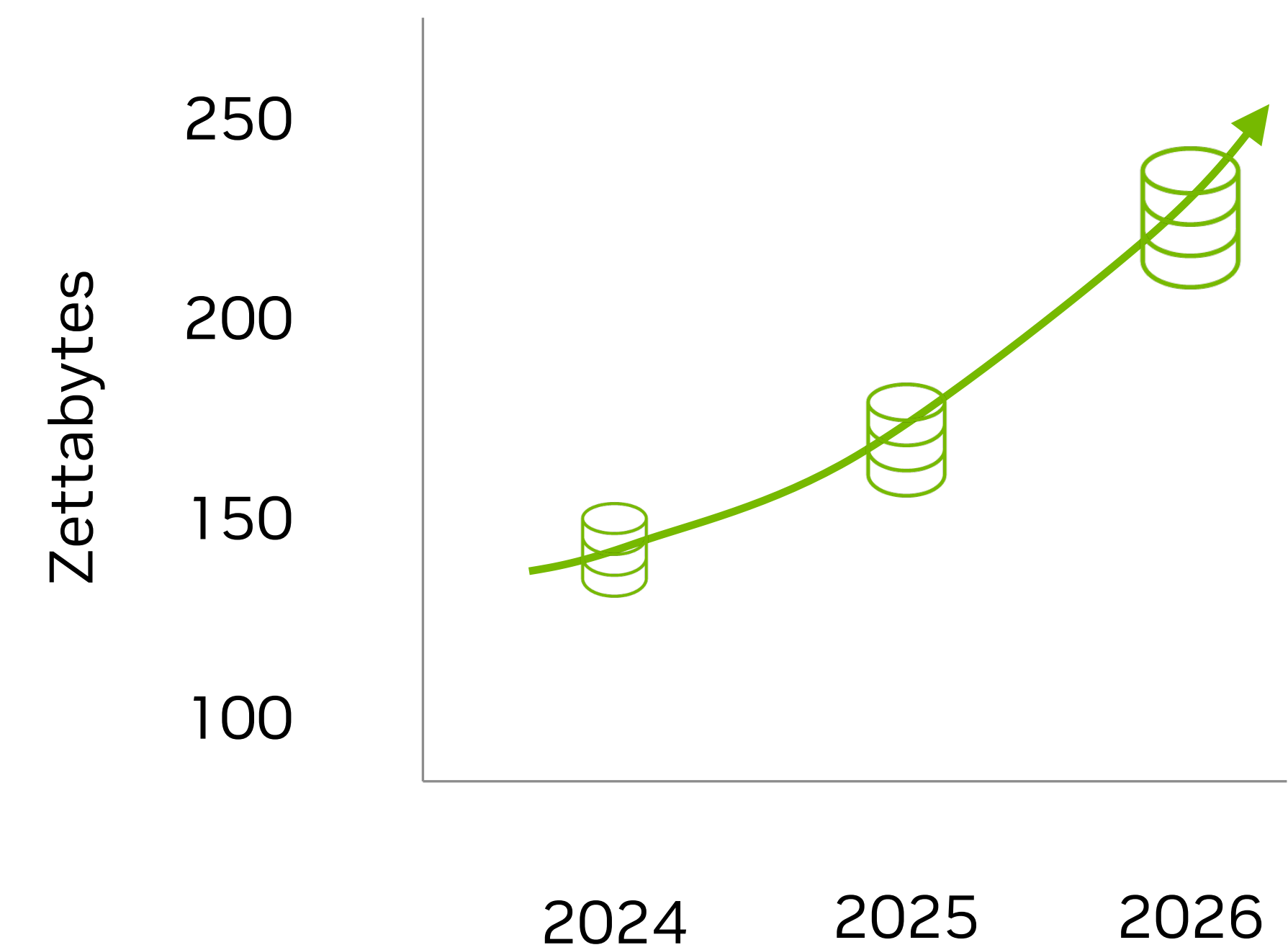
- RAPIDS Accelerator for Apache Spark Data Processing

- RAPIDS Accelerator for Apache Spark ML

- Additional Information

**NVIDIA.**

# Data Processing Compute Challenges



Zettabytes

250
200
150
100

2024   2025   2026

## Data Growth

221 ZB of data by 2026

## Scaling Compute

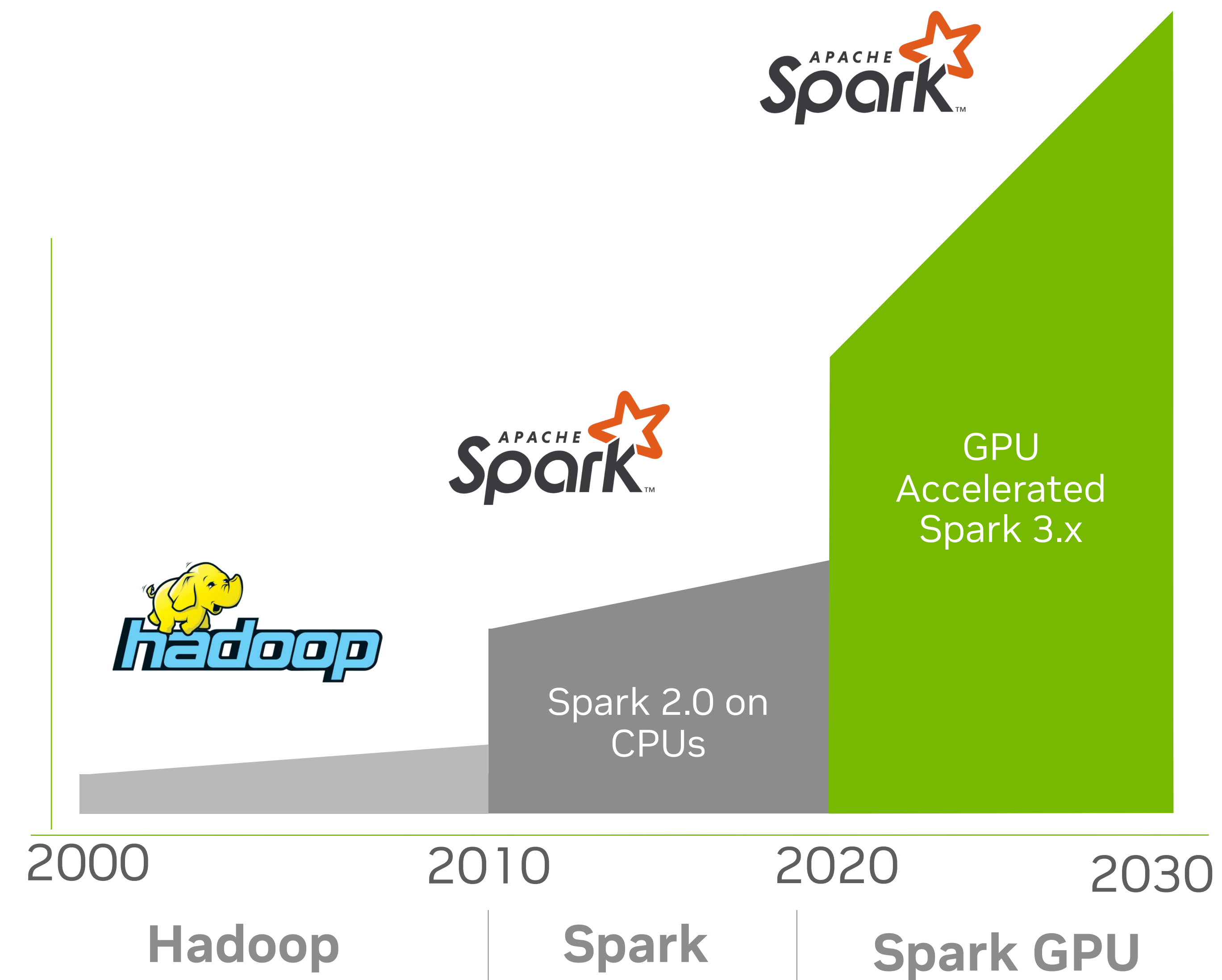Moore's law has slowed

## Power Consumption

Data centers account for 2% of the total US electricity use

# Scaling ETL Processing With Apache Spark with GPUs

RAPIDS Accelerator for Apache Spark



Growth in Requirement for Data Processing

GPU Accelerated Spark 3.x

Spark 2.0 on CPUs

2000        2010        2020        2030

**Hadoop**        **Spark**        **Spark GPU**

4   NVIDIA.

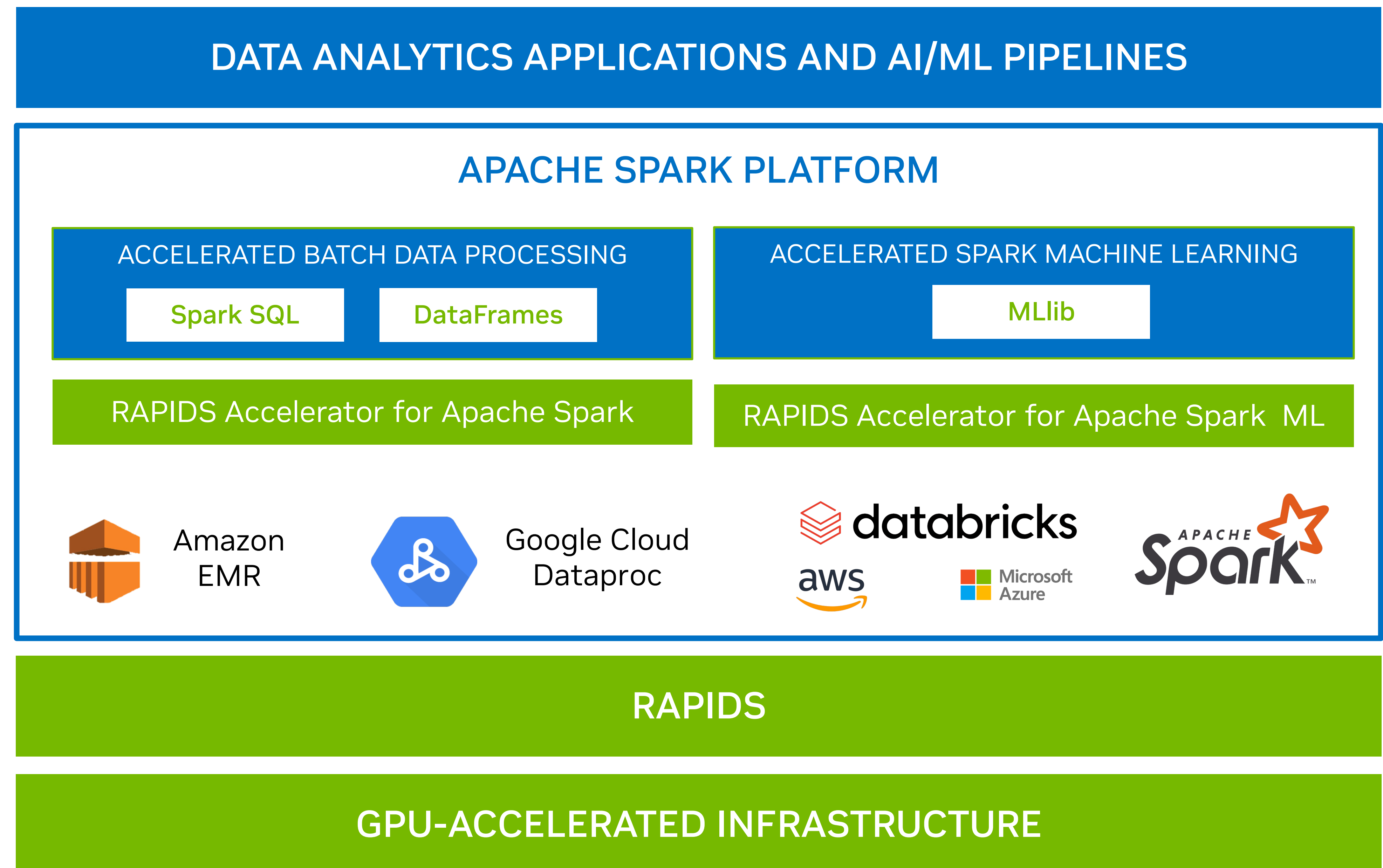# NVIDIA RAPIDS Accelerator

## Key technologies for GPU acceleration

## How it works

- Operates as a software plugin to popular Apache Spark platforms
  - Automatically accelerates supported operations
  - Requires no code changes
- Operations 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*

---

**DATA ANALYTICS APPLICATIONS AND AI/ML PIPELINES**

**APACHE SPARK PLATFORM**

| ACCELERATED BATCH DATA PROCESSING | ACCELERATED SPARK MACHINE LEARNING |
|---|---|
| Spark SQL    DataFrames | MLlib |

RAPIDS Accelerator for Apache Spark | RAPIDS Accelerator for Apache Spark ML

Amazon EMR    Google Cloud Dataproc    databricks    aws    Microsoft Azure    Spark

**RAPIDS**

**GPU-ACCELERATED INFRASTRUCTURE**

# RAPIDS Spark

- Data Processing Challenges

- RAPIDS Spark Data Processing

- RAPIDS Spark ML

- Additional Information

# No Query Changes

- Add jar to classpath and set spark.plugins config

- Same SQL and DataFrame code

- Compatible with PySpark, SparkR, Java, Scala and other DataFrame-based APIs

- Seamless fallback to CPU for unsupported operations

```
spark.sql( """

    SELECT
        o_order_priority
        count(*) as order_count
    FROM
        orders
    WHERE
        o_orderdate >= DATE '1993-07-01'
        AND o_orderdate < DATE '1993-07-01' +
interval '3' month
        AND EXISTS (
            SELECT
                *
            FROM lineitem
            WHERE
                l_orderkey = o_orderkey
                AND l_commitdate < l_receiptdate
            )
    GROUP BY
        o_orderpriority ORDER BY o_orderpriority

""" ).show()
```

# NVIDIA Decision Support Benchmark

NVIDIA Decision Support (NDS) is our adaptation of the TPC-DS benchmark often used by Spark customers and providers

NDS consists of the same 100+ SQL queries as the industry standard benchmark but has modified parts for execution scripts.

The NDS benchmark is derived from the TPC-DS benchmark and as such is not comparable to published TPC-DS results, as the NDS results do not comply with the TPC-DS Specification

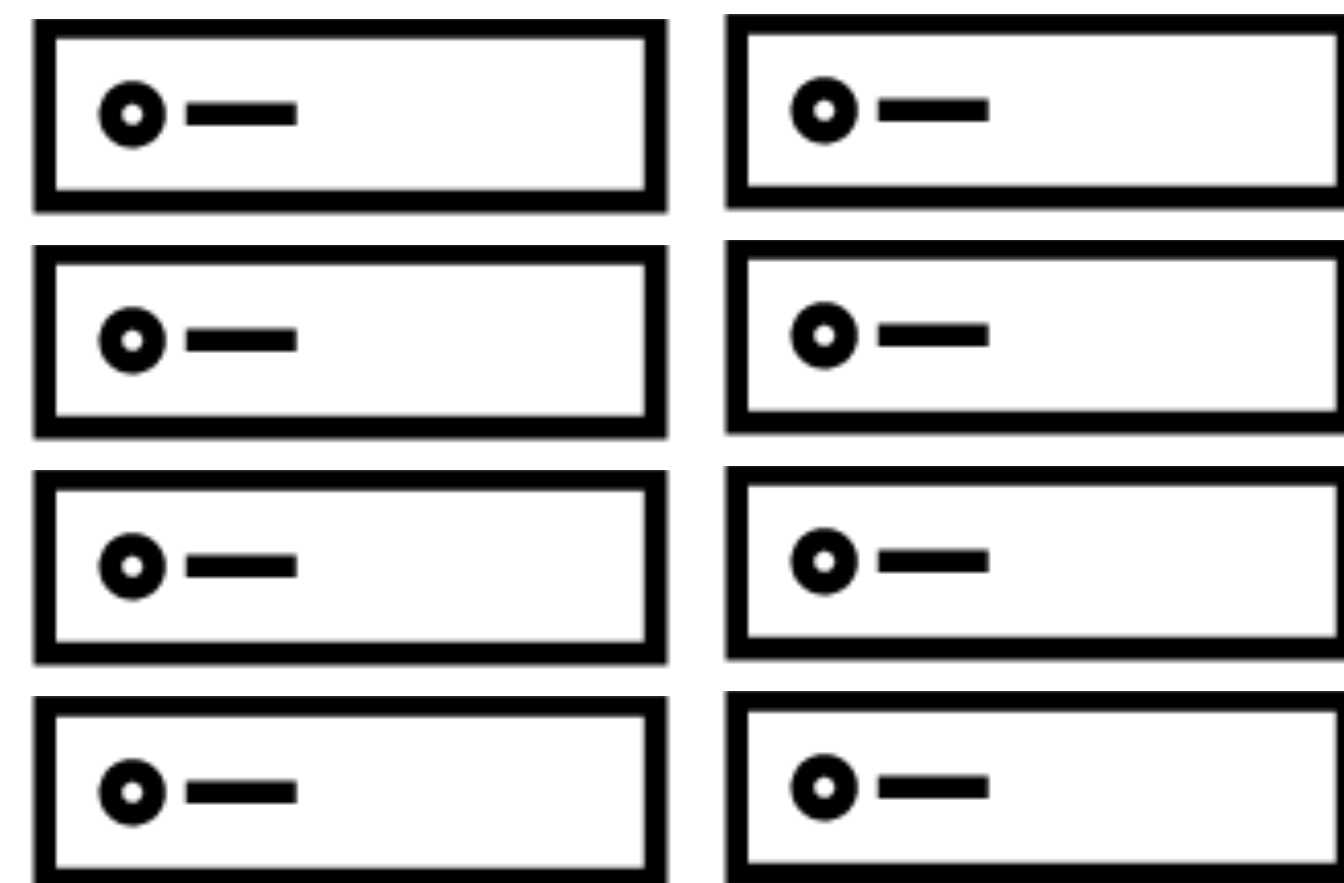https://github.com/nvidia/spark-rapids-benchmarks

NVIDIA

# AWS EC2 cluster

Parquet data, scale factor 3k, stored on S3

8 x r6id.8xlarge

8 x g5.8xlarge

S3

| | CPU Cores | CPU Mem (GB) | Network BW (Gbps) | Storage | GPU | On Demand $ Cost / Hr |
|---|---|---|---|---|---|---|
| r6id.8xlarge | 32 | 256 | 12.5 | 1900GB local SSD | | $2.419 |
| g5.8xlarge | 32 | 128 | 25 | 900GB local SSD | A10 | $2.448 |

# AWS EC2
Configurations

| | CPU | GPU | Config type |
|---|---|---|---|
| spark.driver.memory | 16G | 16G | Resource |
| spark.executor.cores | 16 | 16 | |
| spark.executor.instances | 16 | 8 | |
| spark.executor.memory | 64G | 64G | |
| spark.rapids.filecache.enabled | | true | |
| spark.executor.resource.gpu.amount | | 1 | |
| spark.task.resource.gpu.amount | | 0.0625 | |
| spark.scheduler.minRegisteredResourcesRatio | 1.0 | 1.0 | Scheduling |
| spark.locality.wait | 0s | 0s | |
| spark.sql.files.maxPartitionBytes | 128mb (default) | 2GB | Shuffle |
| spark.shuffle.manager | | com.nvidia.spark.rapids.spark341.RapidsShuffleManager | |
| spark.rapids.shuffle.multiThreaded.{reader\|writer}.threads | | 32 | |
| spark.rapids.sql.multiThreadedRead.numThreads | | 100 | |
| spark.plugins | | com.nvidia.spark.SQLPlugin | GPU |
| spark.rapids.memory.host.spillStorageSize | | 16G | |
| spark.rapids.memory.pinnedPool.size | | 8G | |
| spark.rapids.sql.concurrentGpuTasks | | 3 | |

NVIDIA.

# NVIDIA Decision Support Benchmark 3TB, AWS EC2

## Apache Spark 3.4.1, RAPIDS Spark release 24.04



Time (min)

70
60    57.7
50
40
30
20
10              10.3
0
      r6id.8xlarge    g5.8xlarge

5.5x faster

Cost

$25.00
$20.00    $19.06
$15.00
$10.00
$5.00              $3.76
$0.00
      r6id.8xlarge    g5.8xlarge

80% cost savings

# Grace Hopper (GH200) 16 Node Cluster

## Parquet data stored on HDFS

16 x Grace Hopper Nodes

|  | CPU Cores | CPU Mem (GB) | Network BW (Gbps) | Storage | GPU | Retail Price / node |
|---|---|---|---|---|---|---|
| Quanta GH200 | 72 | 512 | 100 | 4 x 3.8TB local SSD | H100 | $45763 |

# Grace Hopper

## Configurations

| | GPU | Config type |
|---|---|---|
| spark.driver.memory | 50G | Resource |
| spark.driver.maxResultSize | 2G | |
| spark.executor.cores | 16 | |
| spark.executor.memory | 16G | |
| spark.rapids.filecache.enabled | true | |
| spark.executor.resource.gpu.amount | 1 | |
| spark.task.resource.gpu.amount | 0.0625 | |
| | | |
| spark.locality.wait | 0s | Scheduling |
| | | |
| spark.sql.files.maxPartitionBytes | 2GB | Shuffle |
| spark.shuffle.manager | com.nvidia.spark.rapids.spark341.RapidsShuffleManager | |
| spark.rapids.shuffle.multiThreaded.{reader\|writer}.threads | 32 | |
| | | |
| spark.plugins | com.nvidia.spark.SQLPlugin | GPU |
| spark.rapids.memory.host.spillStorageSize | 32G | |
| spark.rapids.memory.pinnedPool.size | 8G | |
| spark.rapids.sql.concurrentGpuTasks | 4 | |

# Grace Hopper (GH200) 16 Node Cluster

RAPIDS Spark 24.04



Time (min)

| | SF3K | SF10K | SF30K | SF100K |
|---|---|---|---|---|
| | 3.4 | 5.4 | 14.6 | 37.6 |

Data Growth Rate

Time Growth Rate

| | SF3K | SF10K | SF30K | SF100K |
|---|---|---|---|---|
| Data Growth Rate | 1.0 | 3.3 | 10.0 | 33.3 |
| Time Growth Rate | | 1.6 | 4.3 | 11.0 |

# RAPIDS Accelerator for Apache Spark

## Spark Plugin for GPU Acceleration

Spark DataFrame / SQL API Application

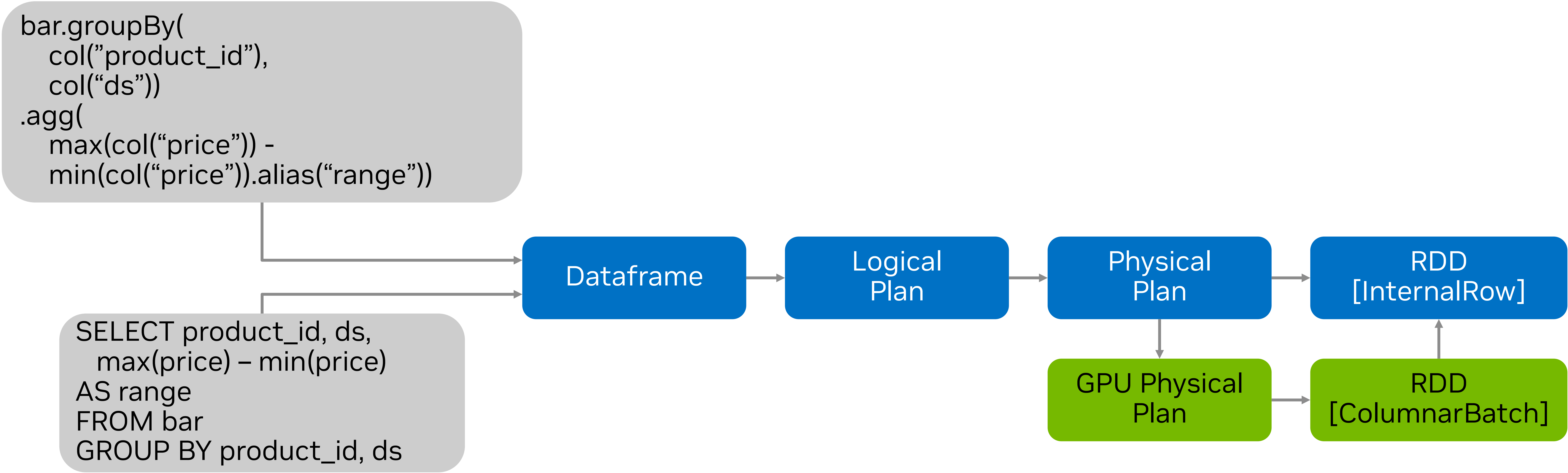Spark Core

RAPIDS Spark plugin

```
if gpu_enabled(operation && datatype):
  RAPIDS Spark
else:
  Spark CPU
```

Java Bindings

RAPIDS C++ cuDF, cuIO, RMM

CUDA

# Spark SQL & DataFrame Query Execution

bar.groupBy(
    col("product_id"),
    col("ds"))
.agg(
    max(col("price")) -
    min(col("price")).alias("range"))

SELECT product_id, ds,
    max(price) – min(price)
AS range
FROM bar
GROUP BY product_id, ds

Dataframe → Logical Plan → Physical Plan → RDD [InternalRow]

GPU Physical Plan → RDD [ColumnarBatch]

16

# RAPIDS Spark Qualification Tool

## Predicting the benefit of Spark + GPUs

**Spark CPU cluster**

Spark 2.x or Spark 3.x

→

**Event logs**
App info
Execs
Expressions

→

**Qualification tool**
Analyze logs with exec-level speedup factors

→

**Output**
App recommendations
Estimated speedup
Estimated cost savings
Exec and stage details

# spark-rapids-user-tools 24.2.1

```
pip install spark-rapids-user-tools
```

Released: Mar 14, 2024

A simple wrapper process around cloud service providers to run tools for the RAPIDS Accelerator for Apache Spark.

## Navigation

- ☰ Project description
- ⟲ Release history
- ⬇ Download files

## Statistics

View statistics for this project via Libraries.io ↗, or by using our public dataset on Google BigQuery ↗

## Meta

**License:** Apache Software License

**Author:** NVIDIA Corporation ✉

**Requires:** Python >=3.8

## Project description

## spark-rapids-user-tools

User tools to help with the adoption, installation, execution, and tuning of RAPIDS Accelerator for Apache Spark.

The wrapper improves end-user experience within the following dimensions:

1. **Qualification**: Educate the CPU customer on the cost savings and acceleration potential of RAPIDS Accelerator for Apache Spark. The output shows a list of apps recommended for RAPIDS Accelerator for Apache Spark with estimated savings and speed-up.
2. **Bootstrap**: Provide optimized RAPIDS Accelerator for Apache Spark configs based on GPU cluster shape. The output shows updated Spark config settings on driver node.
3. **Tuning**: Tune RAPIDS Accelerator for Apache Spark configs based on initial job run leveraging Spark event logs. The output shows recommended per-app RAPIDS Accelerator for Apache Spark config settings.
4. **Diagnostics**: Run diagnostic functions to validate the Dataproc with RAPIDS Accelerator for Apache Spark environment to make sure the cluster is healthy and ready for Spark jobs.

## Getting started

NVIDIA.

# Apache Spark Ecosystem

## Supported Distributions

| Open-source | Cloud | | | On-prem |
|:---:|:---:|:---:|:---:|:---:|



| Apache Spark 3+ Community Release | Amazon EMR | Databricks | Google Cloud Dataproc | Cloudera CDP |

# Improvements Over the Last Year

- Reliability
  - Spill framework to reduce OOM issues to minimize OOM or GPU specific config changes
  - OOM retry framework for automatic OOM handling in memory-intensive operators

- Performance
  - Dynamic repartitioning in large/skewed hash joins
  - File caching
  - Improved I/O and larger chunk handling for Parquet

- Usability
  - Tooling support on Azure and AWS Databricks, Google Dataproc and AWS EMR

- Scaling to 100s of TB and beyond

- ARM support

- JSON handling improvements

- Support for Delta Lake

# Roadmap

- Blackwell GPU will have hardware decompression for snappy and zstd

- Qualification tool with an ML model to improve prediction

- Support for Apache Hudi and Apache Iceberg

- Performance improvements reading from cloud object stores

- Scaling improvements for GPU hardware

NVIDIA.

# RAPIDS Spark

- Data Processing Challenges

---

- RAPIDS Spark Data Processing

- RAPIDS Spark ML

- Additional Information

---

**NVIDIA**

# NVIDIA RAPIDS Accelerator

## Key technologies for GPU acceleration



**DATA ANALYTICS APPLICATIONS AND AI/ML PIPELINES**

**APACHE SPARK PLATFORM**

ACCELERATED BATCH DATA PROCESSING

Spark SQL | DataFrames

ACCELERATED SPARK MACHINE LEARNING

MLlib

RAPIDS Accelerator for Apache Spark

RAPIDS Accelerator for Apache Spark ML

Amazon EMR | Google Cloud Dataproc | databricks | aws | Microsoft Azure | Spark

NVIDIA / spark-rapids-ml

**RAPIDS**

**GPU-ACCELERATED INFRASTRUCTURE**

# RAPIDS Spark ML

## Motivation

```
spark.sql("SELECT * FROM range(10) where id > 7")

df.join(df2, 'name').select(df.weight, df2.height)

pyspark.ml.clustering.KMeans().fit(df)
```
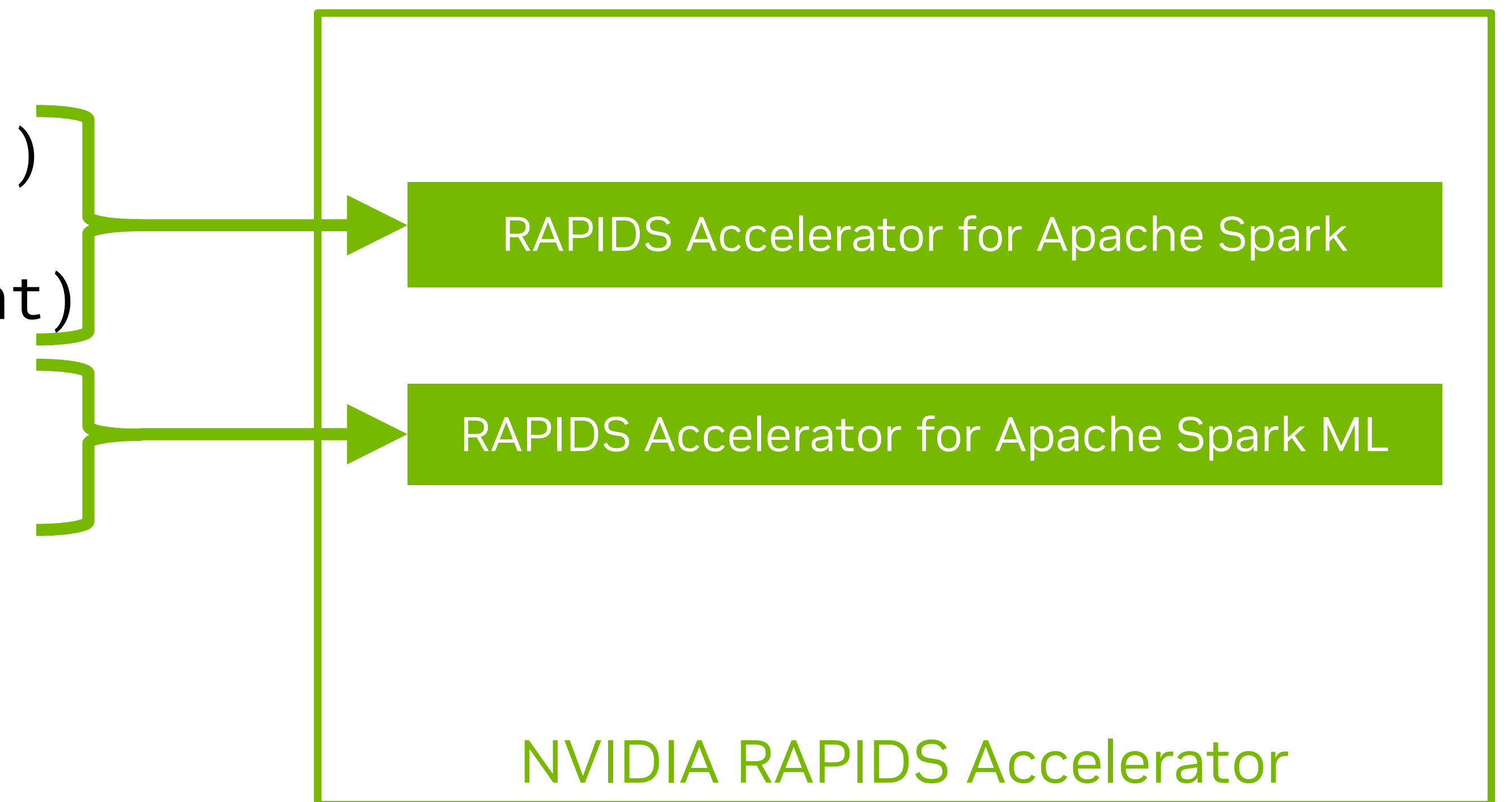
RAPIDS Accelerator for Apache Spark

RAPIDS Accelerator for Apache Spark ML

NVIDIA RAPIDS Accelerator

# Package import change

- Compatible with pyspark.ml DataFrame APIs
- Requires no application code change
- Package import change

```
from pyspark.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)
```

# Package import change

- Compatible with pyspark.ml DataFrame APIs
- Requires no application code change
- Package import change for acceleration

```
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)
```

# Distributed cuML integration

**PySpark MLlib API**

PCA KMeans …

**One-Task-Per-GPU scheduling on Spark DataFrame**

**cuML MNMG classes / Raft NCCL communication**

GPU GPU GPU

Map PySpark MLlib API calls to cuML for supported Algos
Use PySpark APIs (task-per-gpu scheduling, repartition, mapInPandas, barrier, and broadcast) to setup the cuML Multi-Node Multi-GPU cluster

Process data on GPUs using NCCL for communication

NVIDIA.

# RAPIDS Spark ML

## Supported Algorithms

| In Spark MLlib only | In Spark MLlib and cuML | In cuML only |
|---|---|---|
| CrossValidator | K-Means | Exact k-NN |
| | Linear Regression | UMAP |
| | Logistic Regression | |
| | PCA | |
| | Random Forest Classifier | |
| | Random Forest Regressor | |

# Example: MLlib-like API for GPU exact k-NN

```
>>> from spark_rapids_ml.knn import NearestNeighbors
>>> topk = 2
>>> gpu_knn = NearestNeighbors().setInputCol("features").setIdCol("id").setK(topk)
>>> gpu_model = gpu_knn.fit(data_df)
>>> (_, _, knn_df) = gpu_model.kneighbors(query_df)
>>> knnjoin_df = gpu_model.exactNearestNeighborsJoin(query_df, distCol="EuclideanDistance")
>>> knnjoin_df.show()
+--------------+--------------+-----------------+
|       item_df|      query_df|EuclideanDistance|
+--------------+--------------+-----------------+
|{1, [2.0, 2.0]}|{3, [1.0, 1.0]}|        1.4142135|
|{0, [1.0, 1.0]}|{3, [1.0, 1.0]}|              0.0|
|{2, [3.0, 3.0]}|{4, [3.0, 3.0]}|              0.0|
|{1, [2.0, 2.0]}|{4, [3.0, 3.0]}|        1.4142135|
+--------------+--------------+-----------------+
```
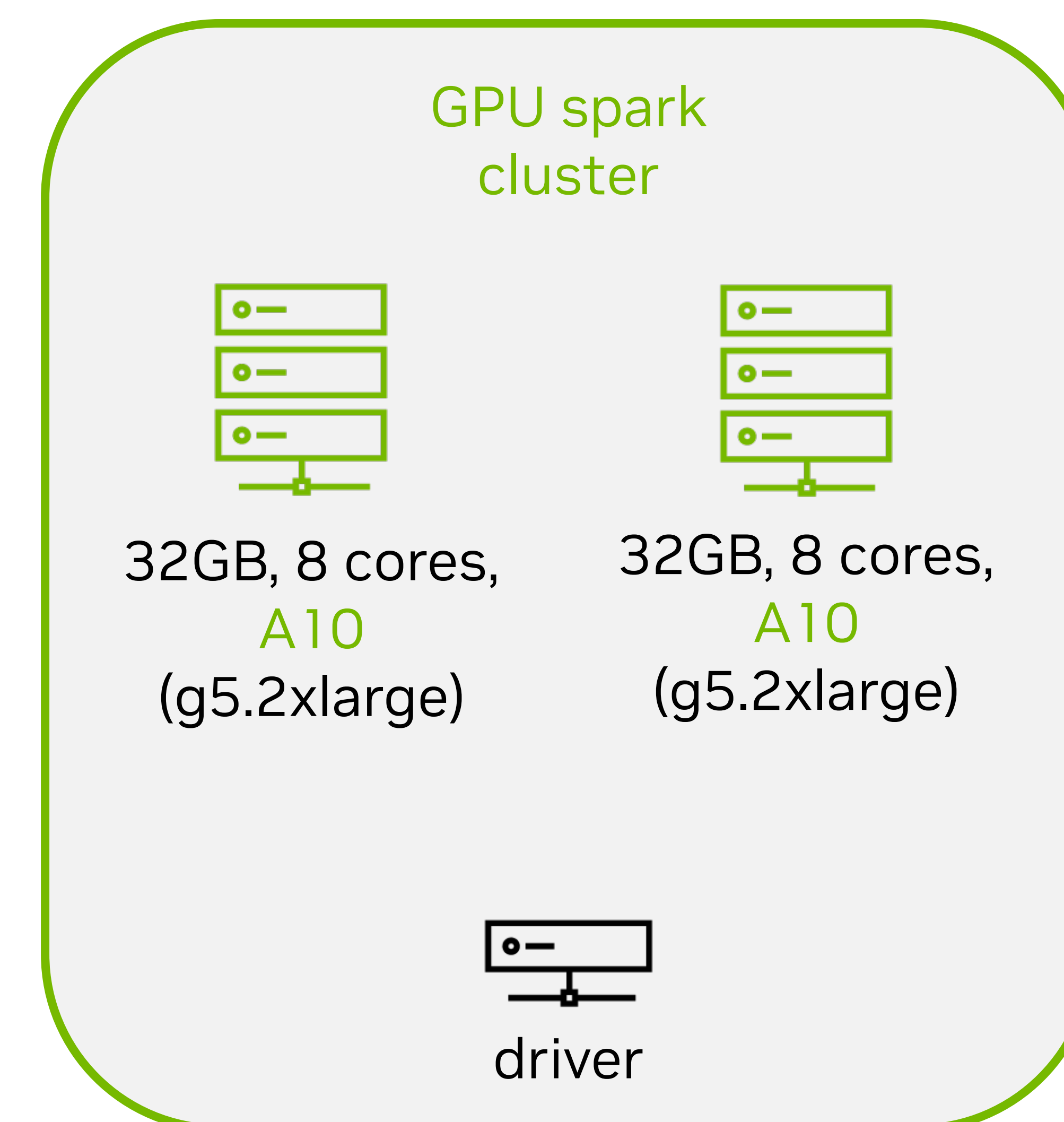
# Microbenchmarking

**CPU spark cluster**

32GB, 8 cores
(m5.2xlarge)

32GB, 8 cores
(m5.2xlarge)

driver

- **Environment**
  - Databricks AWS hosted Spark

- **Workload & Data**
  - estimator.fit(data_df) [i.e. training]
  - data_df read from Parquet format in AWS S3
  - Compute intensive synthetic workloads:
    - 1 million rows
    - 3000 dimensional vectors
  - Data available in S3 public bucket.

**GPU spark cluster**

32GB, 8 cores,
A10
(g5.2xlarge)

32GB, 8 cores,
A10
(g5.2xlarge)

driver

- Instructions and scripts to reproduce: https://github.com/NVIDIA/spark-rapids-ml/tree/main/python/benchmark#databricks

- [Repo also has instructions for GCP Dataproc and AWS EMR]

# Training/fit time: 6x-100x faster

■ Spark ML CPU (sec)   ■ Spark-Rapids-ML GPU (sec)

Seconds

| Model | Spark ML CPU (sec) | Spark-Rapids-ML GPU (sec) |
|---|---|---|
| K-means | 9526 | 82 |
| Random Forest Classifier | 2364 | 59 |
| Random Forest Regressor | 1572 | 52 |
| PCA | 661 | 37 |
| Linear Regression - Ridge | 558 | 32 |
| Linear Regression | 594 | 41 |
| Linear Regression - ElasticNet | 551 | 79 |
| Logistic Regression | 406 | 69 |

◎ NVIDIA

# Cost Benefits and Speedups



■ Speed-up (cpu time/gpu time)

■ Cost Benefits (cpu cost/gpu cost)

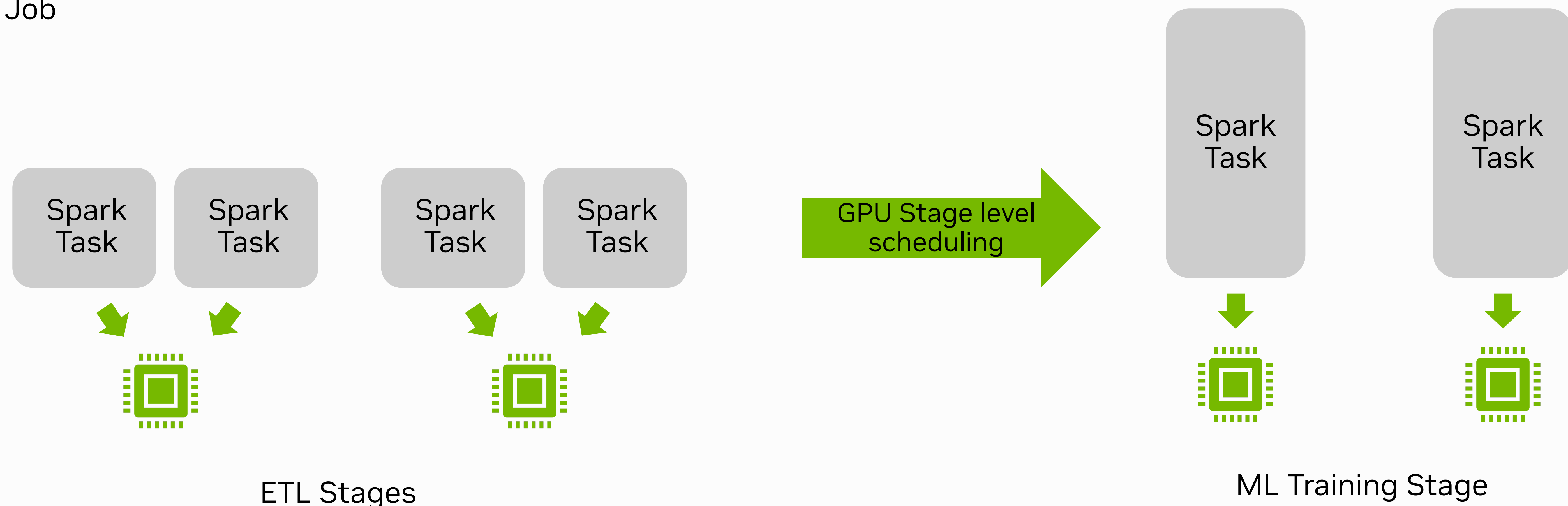| Category | Speed-up | Cost Benefits |
|---|---|---|
| K-means | 116 | 58 |
| Random Forest Classifier | 40 | 20 |
| Random Forest Regressor | 30 | 15 |
| PCA | 18 | 9 |
| Linear Regression - Ridge | 17 | 9 |
| Linear Regression | 14 | 7 |
| Linear Regression - ElasticNet | 7 | 3 |
| Logistic Regression | 6 | 3 |

# End-to-End Acceleration
## Stage Level Scheduling

- ML training stage runs all tasks at the same time with one Task for each GPU
  - Required by cuML/NCCL layer
- ETL can benefit from multiple Tasks per GPU
- Stage level scheduling:
  - Allows different tasks per GPU on a per stage basis within the same Spark Job.
- GPU aware stage level scheduling

Single Spark Job

| Spark Task | Spark Task | | Spark Task | Spark Task |

GPU Stage level scheduling →

Spark Task       Spark Task

ETL Stages

ML Training Stage
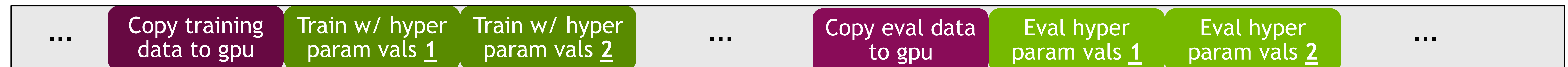
# End-to-End Acceleration
## Accelerated CrossValidator

- PySpark API compatibility allows all accelerated Algos to leverage PySpark's built in CrossValidator for hyper parameter tuning
- We can do better:

### PySpark MLlib CrossValidator

| ... | Copy training data to gpu | Train w/ hyper param vals **1** | Copy eval data to gpu | Eval hyper param vals **1** | Copy training data to gpu | Train w/ hyper param vals **2** | Copy eval data to gpu | Eval hyper param vals **2** |
|---|---|---|---|---|---|---|---|---|

### RAPIDS Spark ML CrossValidator

| ... | Copy training data to gpu | Train w/ hyper param vals **1** | Train w/ hyper param vals **2** | ... | Copy eval data to gpu | Eval hyper param vals **1** | Eval hyper param vals **2** | ... |
|---|---|---|---|---|---|---|---|---|

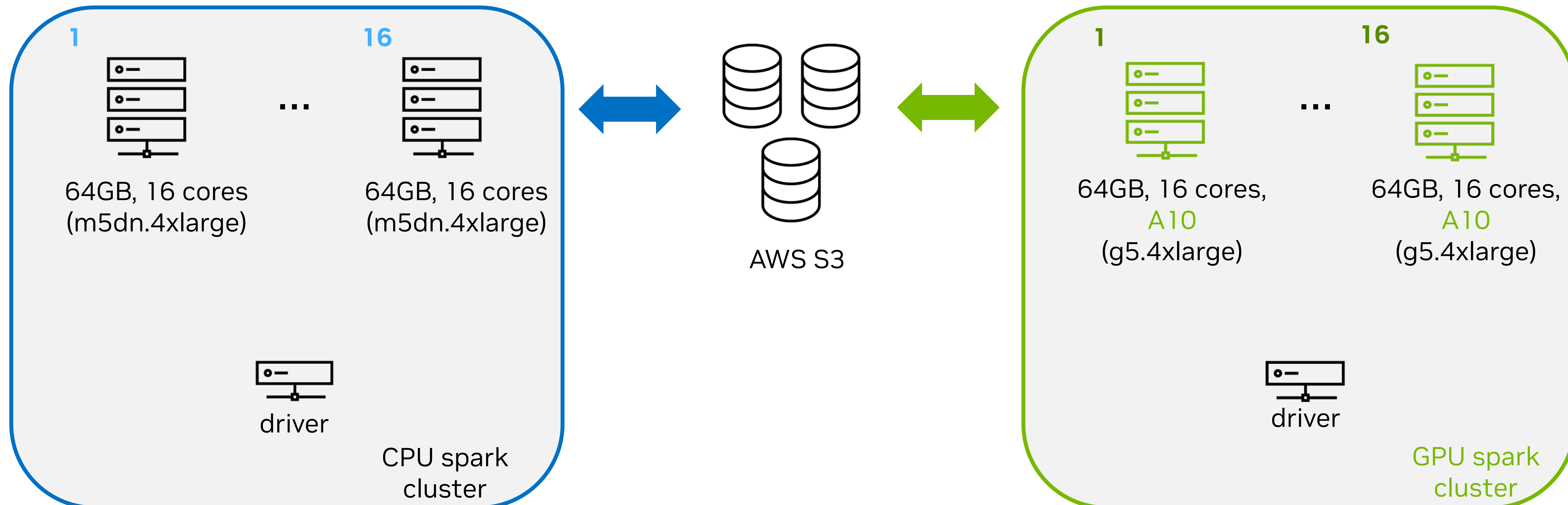NVIDIA.

# Fannie Mae Mortgage Benchmark
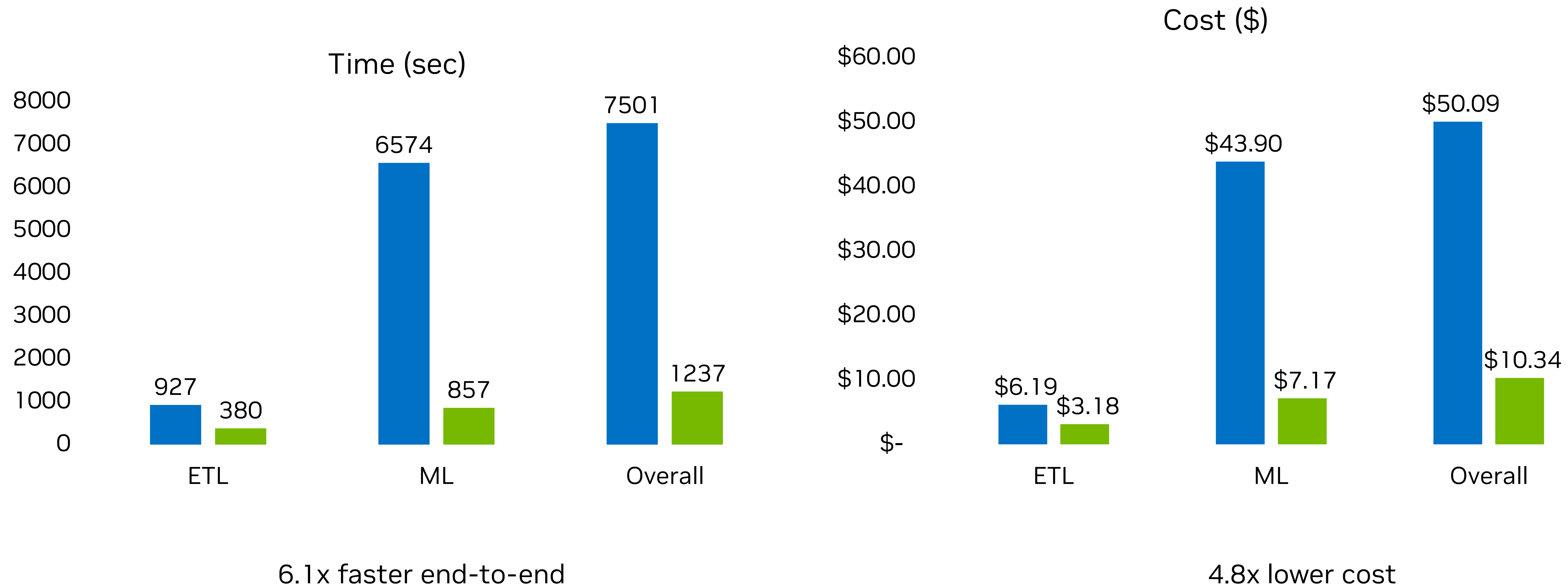## ETL + ML

- End-to-end ETL – ML workload:
  - ETL:
    - Starts with compressed Fannie Mae Single-Family Loan Performance Data ~ 800GB csv dataset converted to 26.8 GB compressed Parquet.
    - Feature engineering to 2.6 billion records by 27 feature dataset with loan delinquency as label. (as in this example: https://github.com/NVIDIA/spark-rapids-examples/.../MortgageETL.ipynb).
  - ML
    - Logistic regression with 3 fold cross validation wrt to log loss over 8 algo parameter choices.
  - Environment:  Databricks AWS 16 node clusters:



**1** ... **16**

64GB, 16 cores
(m5dn.4xlarge)          64GB, 16 cores
(m5dn.4xlarge)

driver

CPU spark cluster

AWS S3

**1** ... **16**

64GB, 16 cores,
A10
(g5.4xlarge)          64GB, 16 cores,
A10
(g5.4xlarge)

driver

GPU spark cluster

35

# Fannie Mae Mortgage Benchmark
## ETL + ML

Time (sec)

Cost ($)

8000
7000
6000
5000
4000
3000
2000
1000
0

$60.00

$50.00

$40.00

$30.00

$20.00

$10.00

$-

927    380    6574    857    7501    1237

ETL    ML    Overall

$6.19    $3.18    $43.90    $7.17    $50.09    $10.34

ETL    ML    Overall

6.1x faster end-to-end

4.8x lower cost

Accuracy:  GPU – CPU ave CV score < 0.004%

# Roadmap

- Better out-of-core

- Blackwell and Grace-Hopper optimizations

- Spark APIs for more algorithms from cuML:
  - Batch approximate nearest neighbor vector search
  - DBSCAN clustering

- GPU optimized Pipelines

- Additional Spark MLlib algos

# Broader GPU Accelerated ML/DL on Spark Ecosystem

🤗 **Hugging Face**

Introducing
Llama 2

*dmlc*
**XGBoost**

**TensorFlow**

**PyTorch**

**APACHE Spark** ™

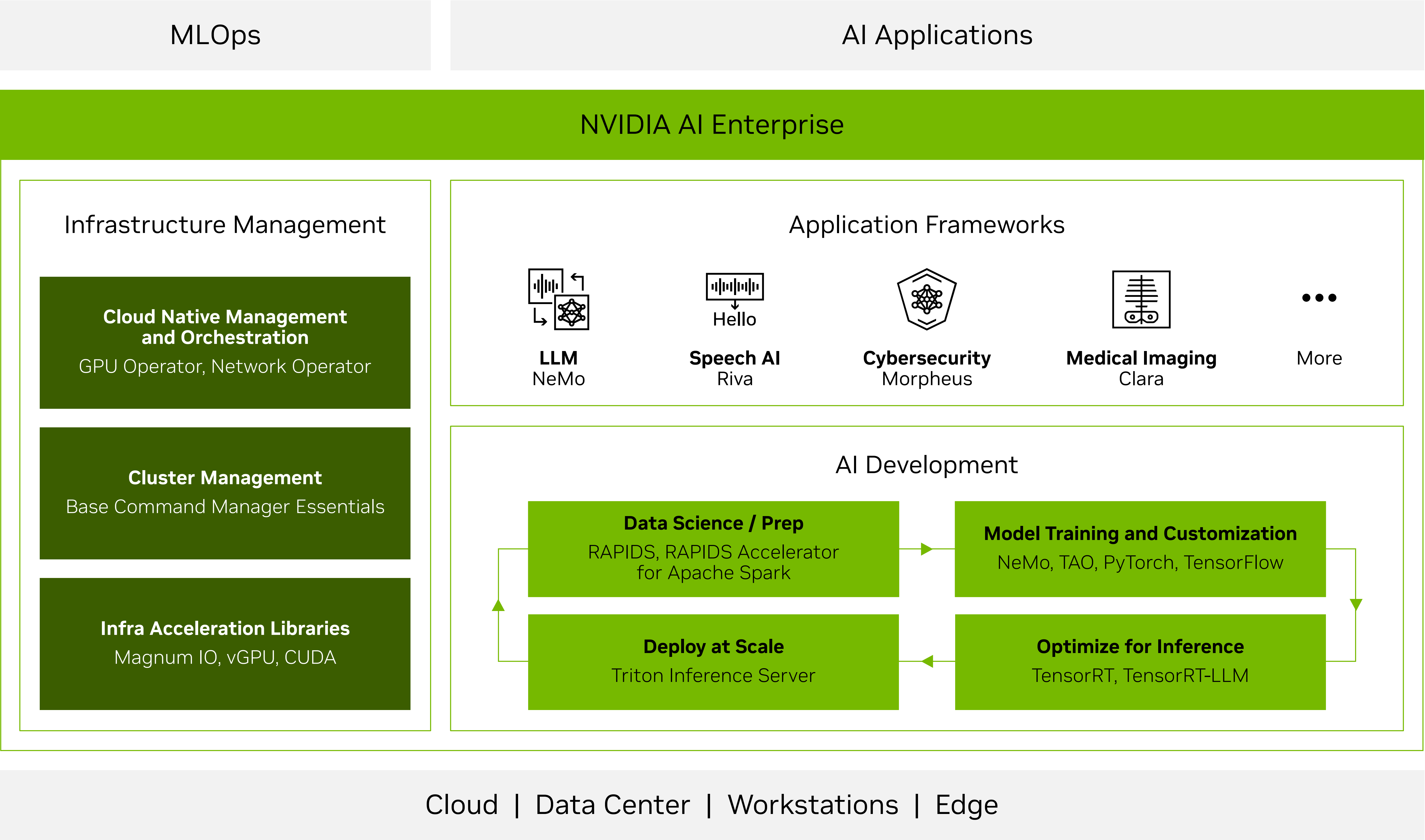| NVIDIA RAPIDS | NVIDIA TRITON INFERENCE SERVER | • • • | NVIDIA TENSOR RT |

**NVIDIA CUDA and GPUS**

38

# RAPIDS Spark

- Data Processing Challenges

- RAPIDS Spark Data Processing

- RAPIDS Spark ML

- Additional Information

# NVIDIA AI Enterprise
## End to end software platform for AI and Data Science

| MLOps | AI Applications |
|---|---|

### NVIDIA AI Enterprise

#### Infrastructure Management

**Cloud Native Management and Orchestration**
GPU Operator, Network Operator

**Cluster Management**
Base Command Manager Essentials

**Infra Acceleration Libraries**
Magnum IO, vGPU, CUDA

#### Application Frameworks

**LLM**
NeMo

**Speech AI**
Riva

Hello

**Cybersecurity**
Morpheus

**Medical Imaging**
Clara

••• More

#### AI Development

**Data Science / Prep**
RAPIDS, RAPIDS Accelerator for Apache Spark

**Model Training and Customization**
NeMo, TAO, PyTorch, TensorFlow

**Deploy at Scale**
Triton Inference Server

**Optimize for Inference**
TensorRT, TensorRT-LLM

Cloud | Data Center | Workstations | Edge

NVIDIA

# Learn More!

PayPal: How PayPal Reduces Cloud Costs by up to 70% with Spark RAPIDS [S62506]
   Wed 8:30-8:55

Baidu: From SQL to Chat: How to Revolutionize Enterprise Data Analysis with NVIDIA [S61622]
   Wed 9:00-9:25

north.io: How AI and Accelerated Computing are Revolutionizing Oceanographic Data Processing [S61391]
   Wed 3:00-3:25

Taboola: RAPIDS Accelerator for Apache Spark Propels Data Center Efficiency and Cost Savings [S62130]
   Thu 10:00-10:25

**Sessions**

Reduce Apache Spark MLlib Costs with NVIDIA GPUs [CWE62407]
   Wed 4:00-4:50

Cost Savings and Speedup with the RAPIDS Accelerator for Apache Spark [CWE62404]
   Thu 9:00-9:50

**Connect with Experts**

Accelerating Data Analytics on GPUs with the RAPIDS Accelerator for Apache Spark [DLIT61196]
   Wed 8:00-9:40

**Tutorial**

# For More Information

spark-rapids-support@nvidia.com

https://docs.nvidia.com/spark-rapids

https://nvidia.github.io/spark-rapids

https://github.com/NVIDIA/spark-rapids-ml

https://nvidia.github.io/spark-rapids-ml