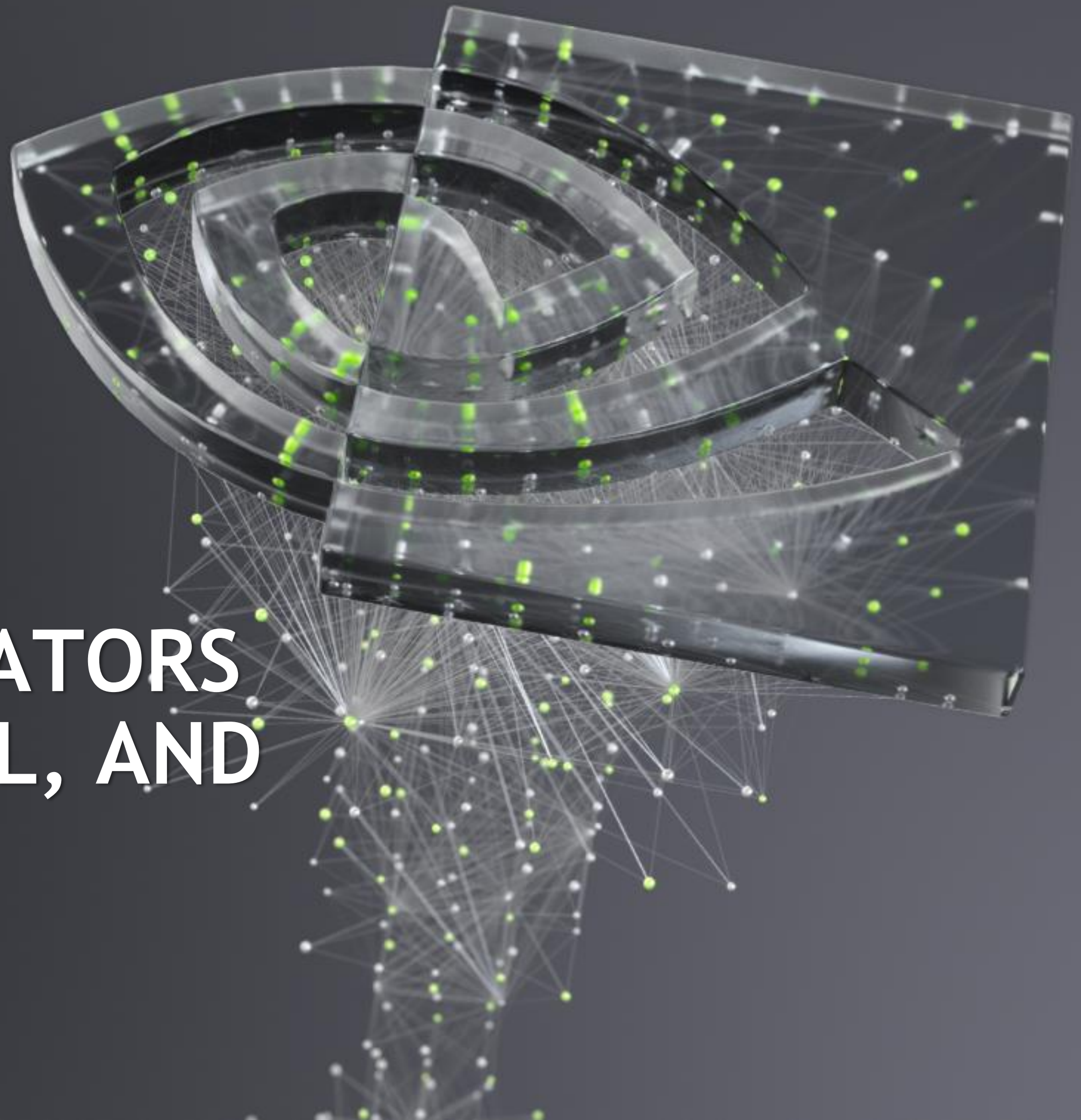




UTILIZING ACCELERATORS TO SPEEDUP ETL, ML, AND DL APPLICATIONS

Jason Lowe and Robert Evans, 05/19/2020





AGENDA

Accelerated ETL

Accelerated SQL/Dataframe

Accelerated Shuffle

What's Next

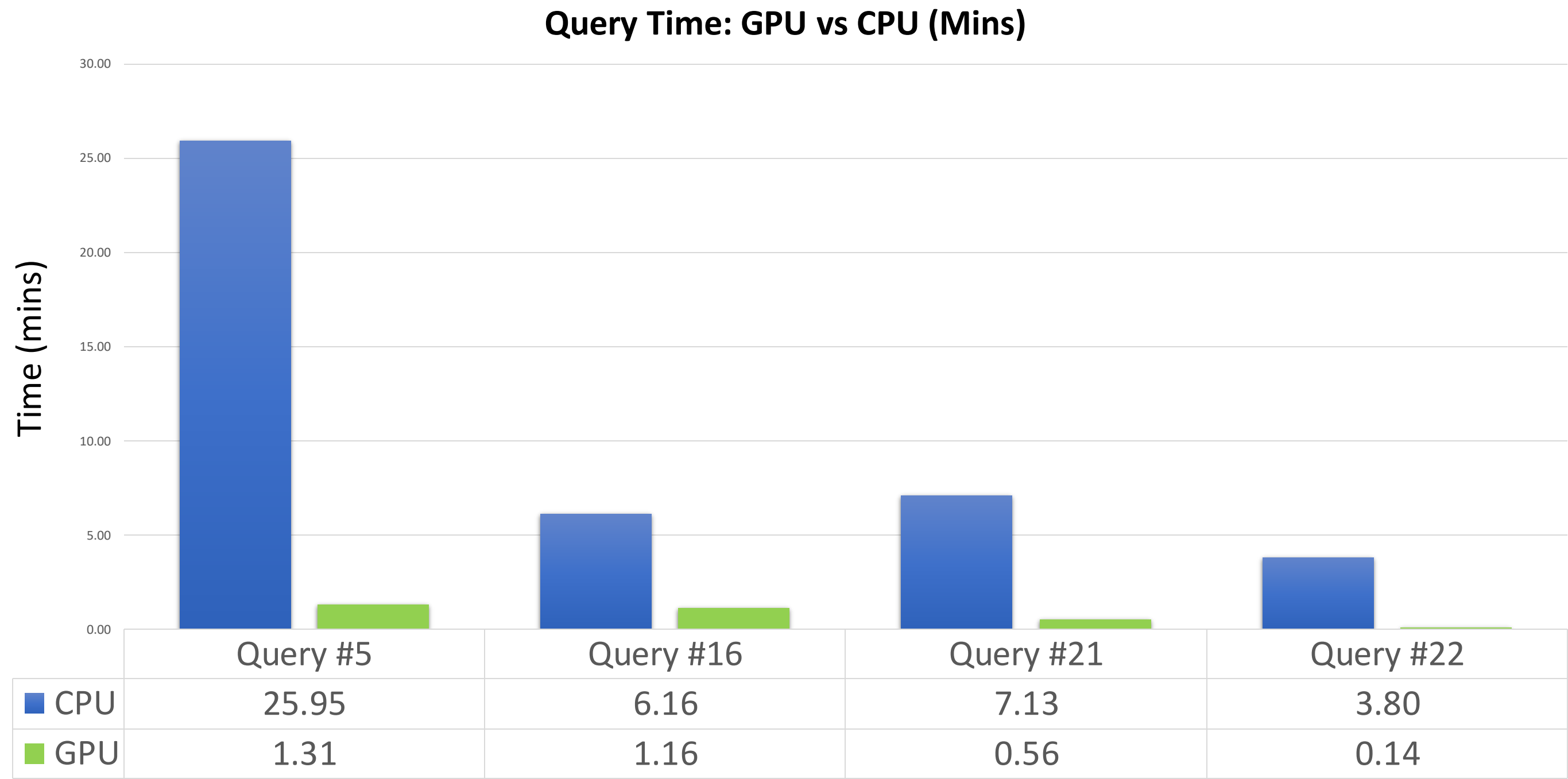
ACCELERATED ETL?

Can a GPU make an elephant fast?



YES

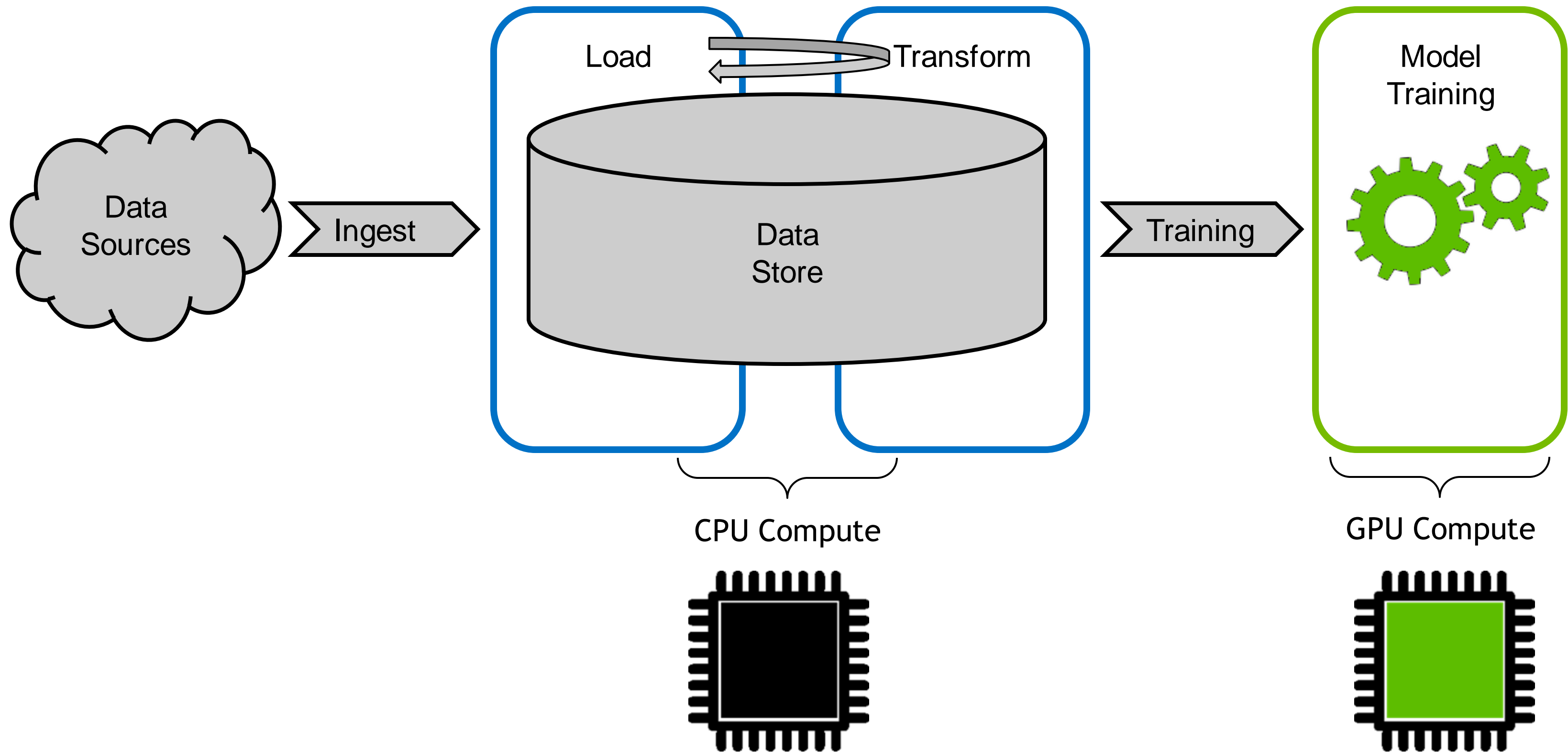
TPCx–BB Like Benchmark Results (10TB Dataset, Two Nodes DGX-2 Cluster)*



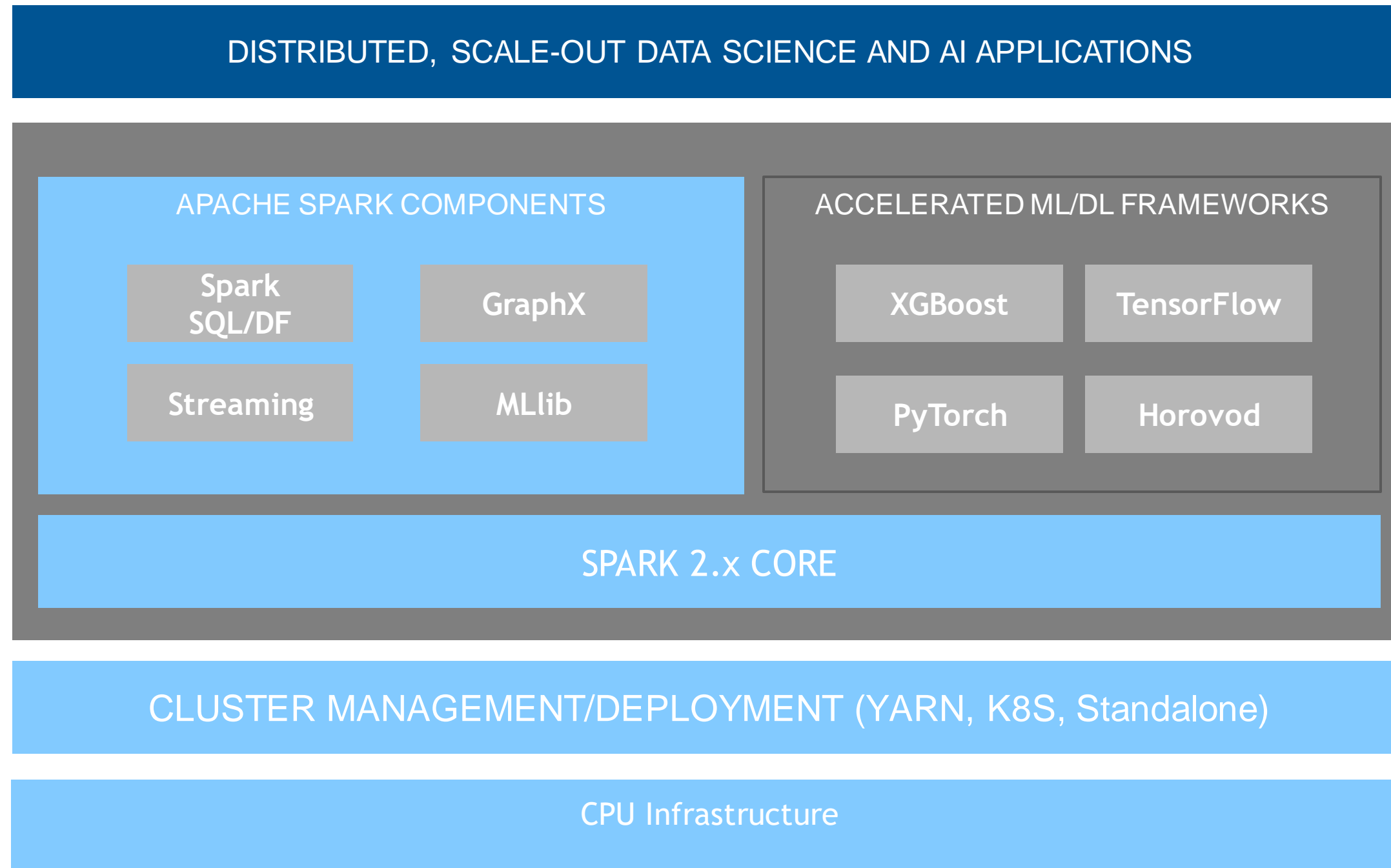
Environment: Two DGX-2 (96 CPU Cores, 1.5TB Host memory, 16 V100 GPUs, 512 GB GPU Memory)

* Not official or complete TPCx-BB runs (ETL power only).

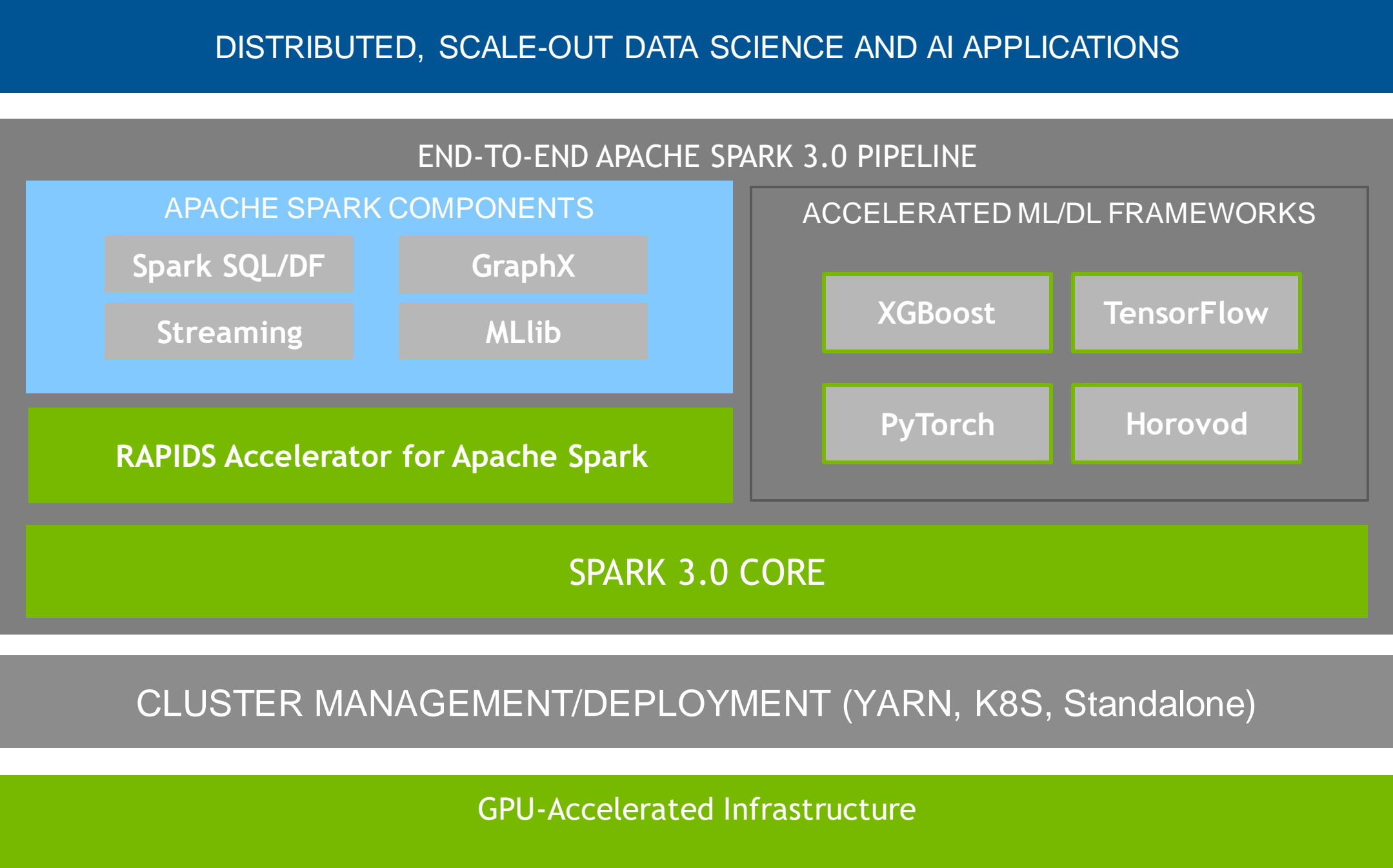
MODERN ML/DL WORKFLOW



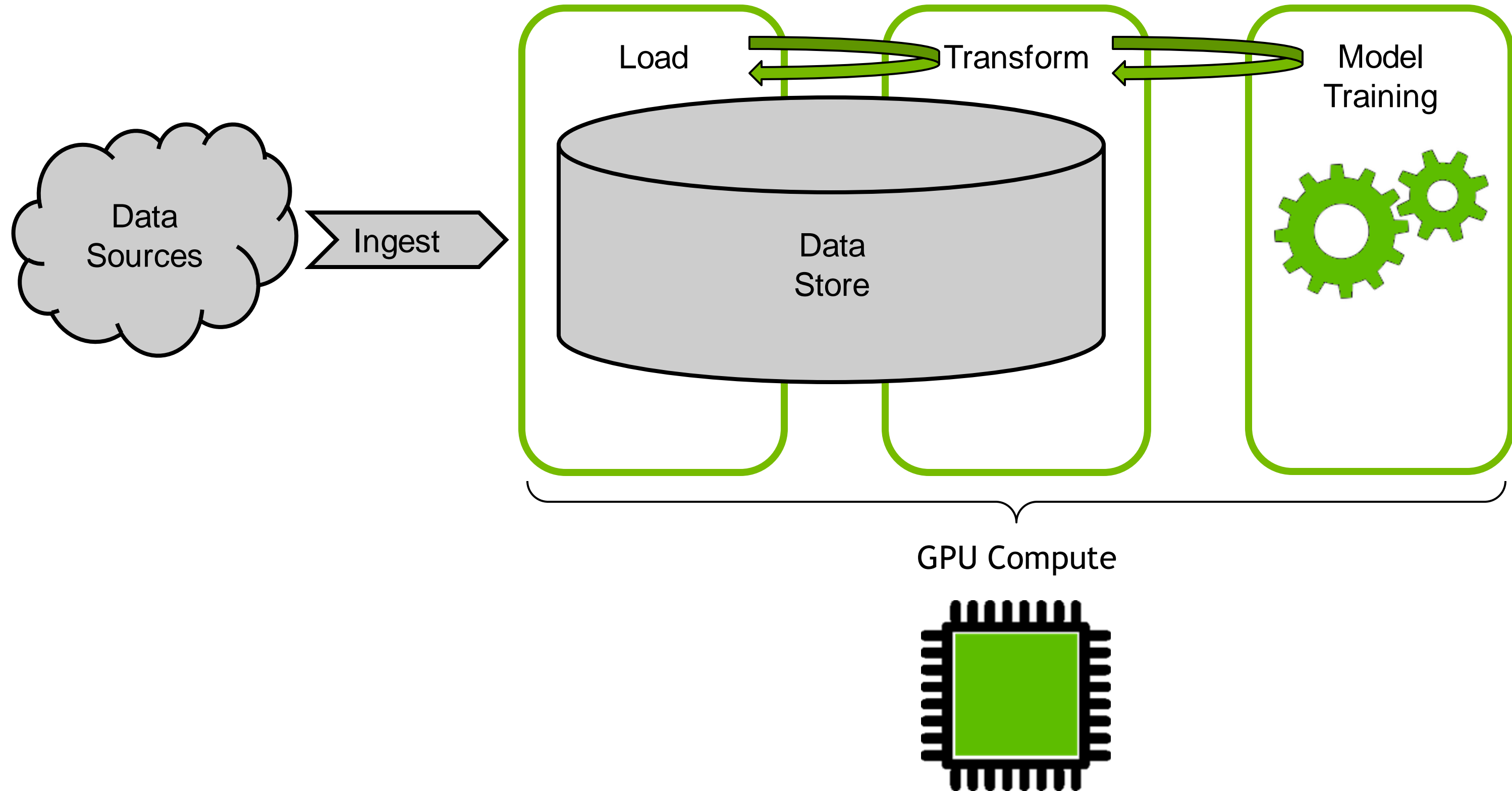
APACHE SPARK 2.X



SPARK 3.X IS A UNIFIED AI PLATFORM



ETL + ML/DL WORKFLOW



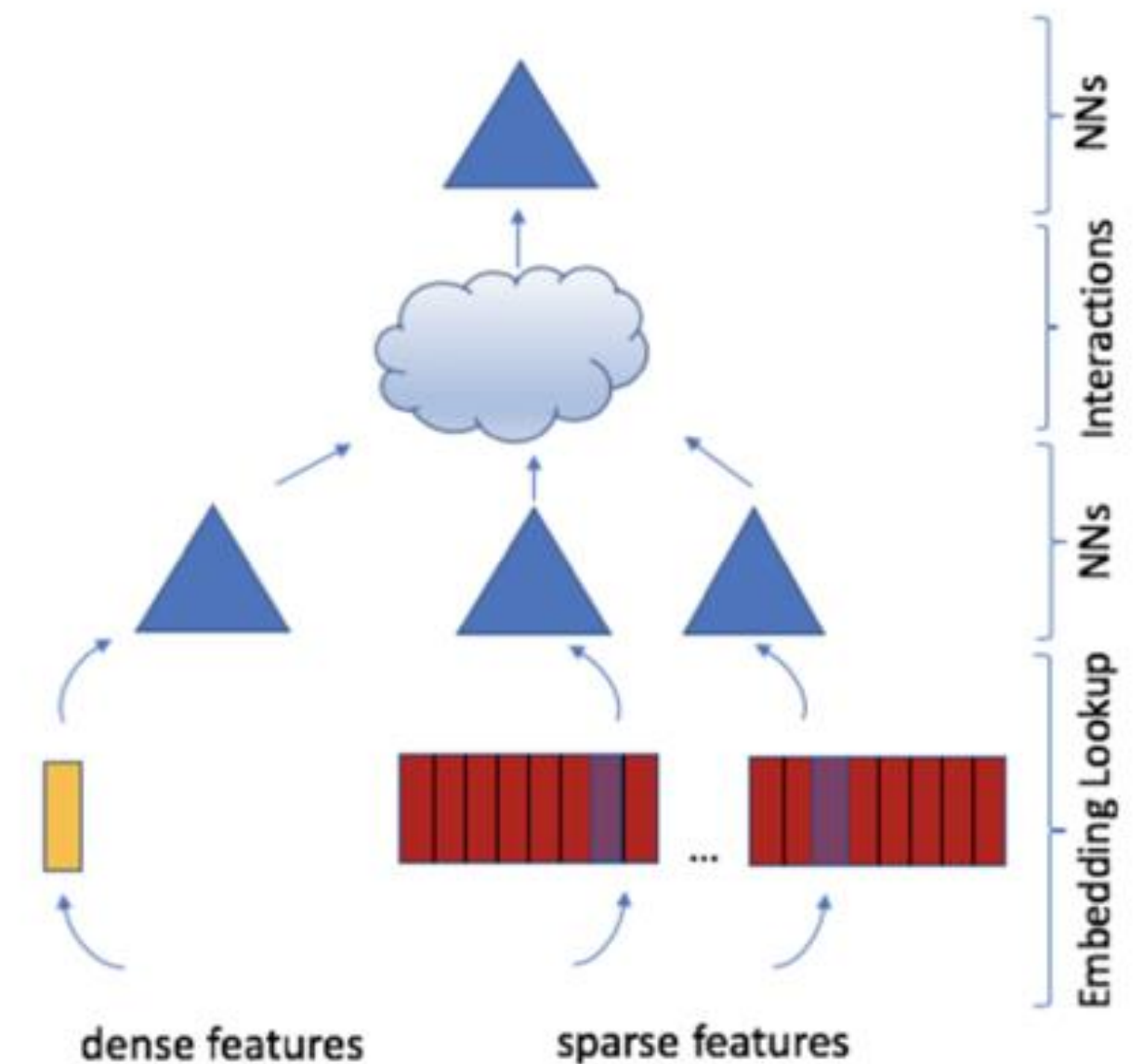
DEEP LEARNING RECOMMENDATION MACHINES

Example use case: Criteo dataset

Anonymized 7-day clickstream dataset (1 TB)

Convert high cardinality string categorical data to contiguous integer ids

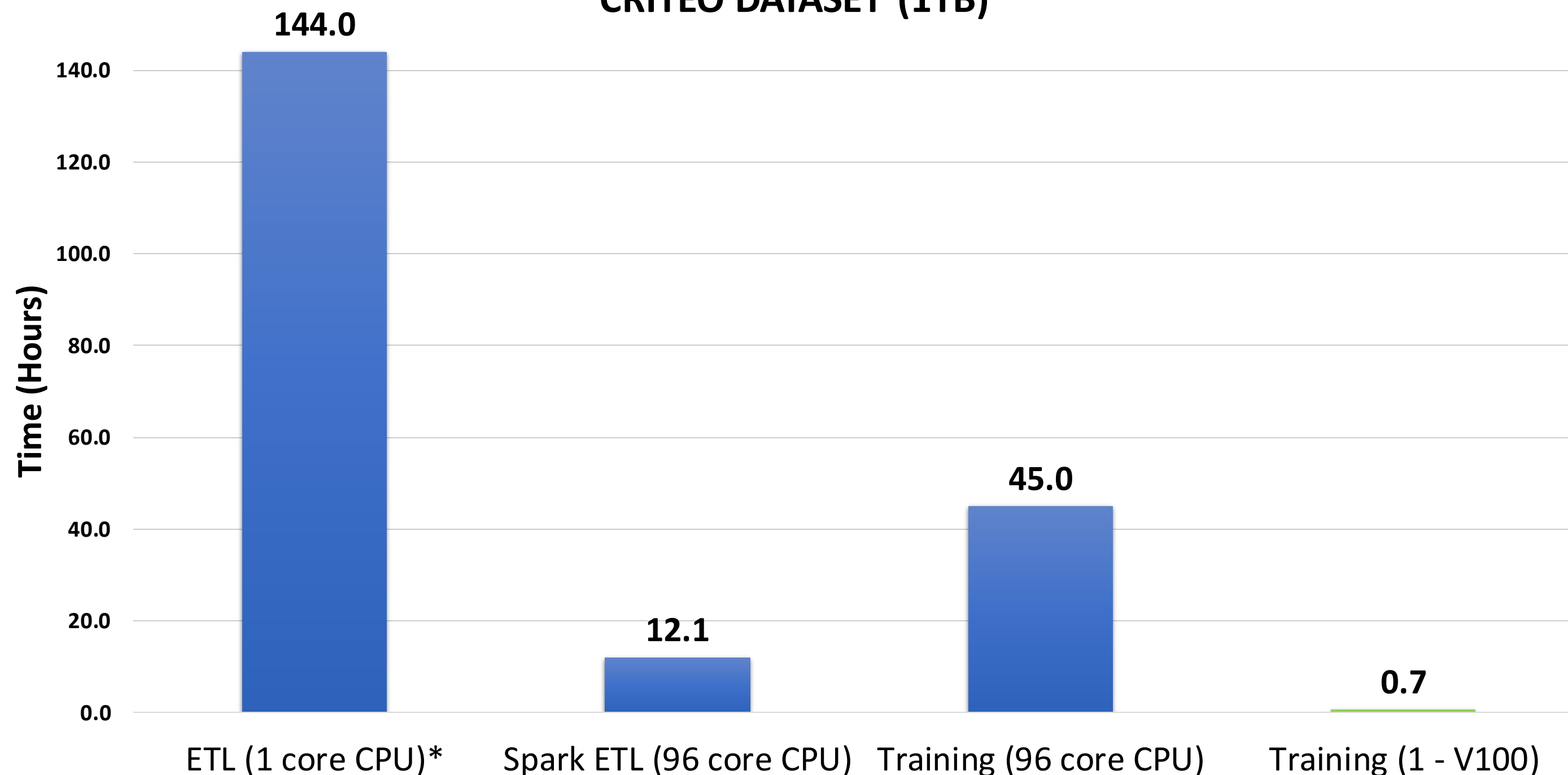
DLRM github repo has scripts for this out of the box



<https://medium.com/analytics-vidhya/deep-learning-recommendation-machines-dlrm-4fec2a5e7ef8>

DLRM ON CRITEO DATASET (PAST)

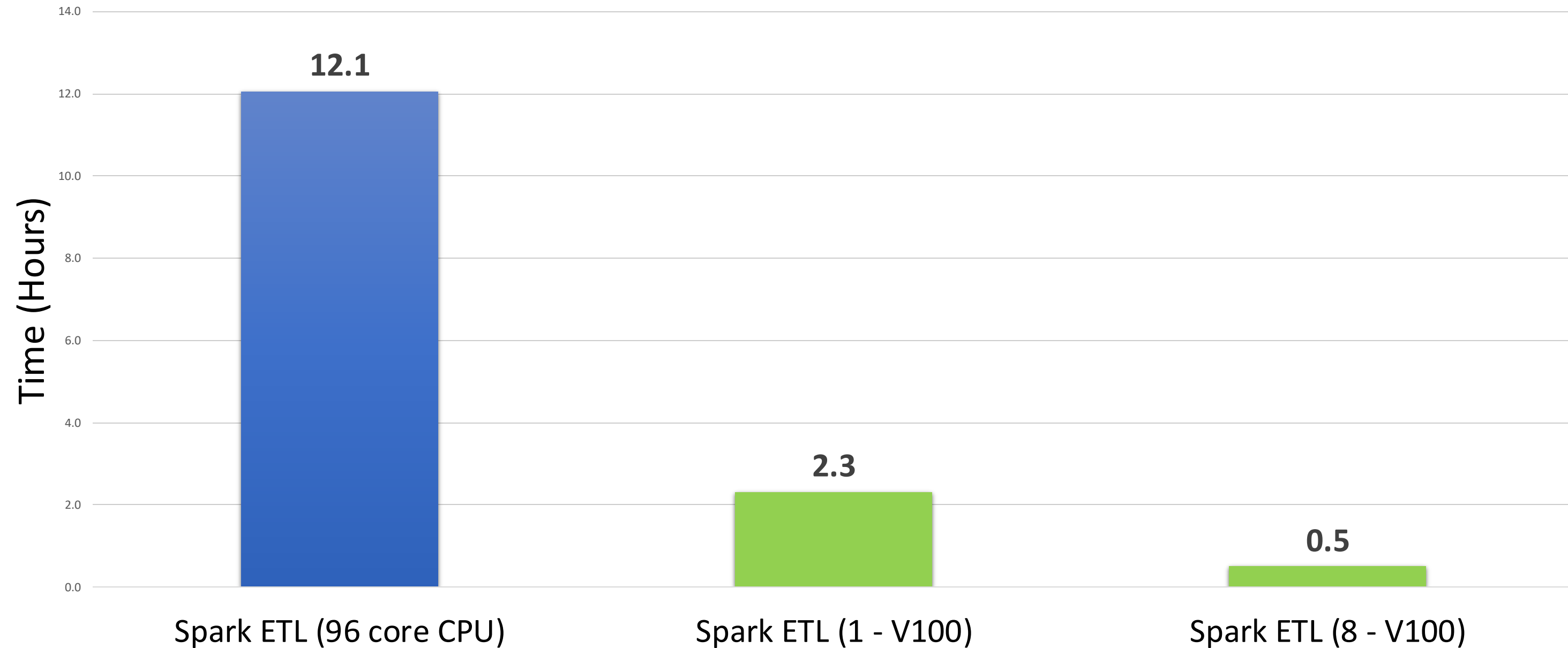
ETL & Training Run Time for CPU & GPU
CRITEO DATASET (1TB)



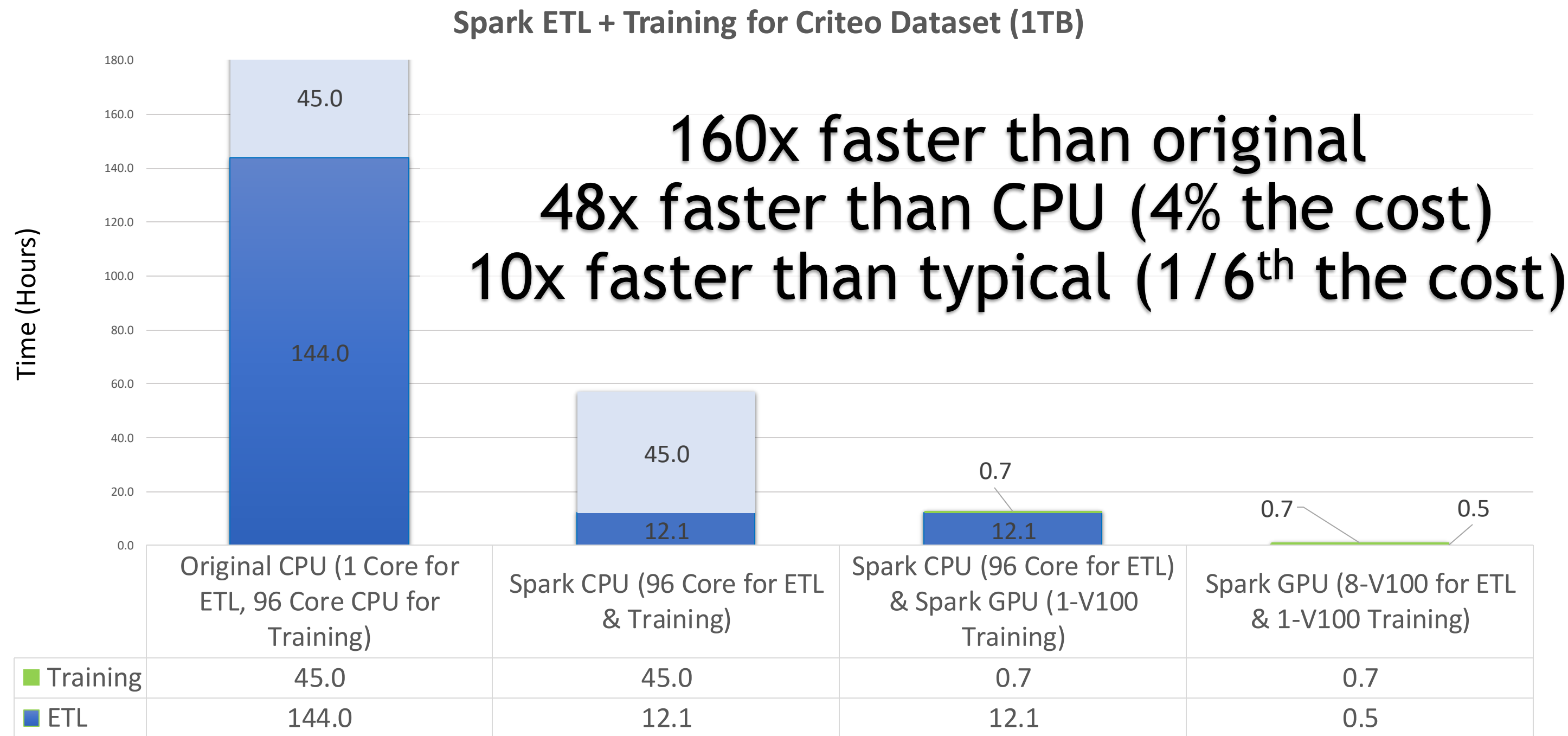
* Extrapolated couldn't convince anyone to wait that long


DLRM ETL ON CRITEO DATASET (PRESENT)

Spark ETL for CRITEO DATASET (1TB)



DLRM END-TO-END ON CRITEO DATASET (PRESENT)





“The more you buy, the more you
save.”

— Jensen Huang, GTC 2018

RAPIDS ACCELERATOR FOR APACHE SPARK (PLUGIN)

DISTRIBUTED SCALE-OUT SPARK APPLICATIONS

APACHE SPARK CORE

Spark SQL API

DataFrame API

Spark Shuffle

```
if gpu_enabled(operation, data_type)
    call-out to RAPIDS
else
    execute standard Spark operation
```

RAPIDS Accelerator
for Spark

- Custom Implementation of Spark Shuffle
- Optimized to use RDMA and GPU-to-GPU direct communication

JNI bindings
Mapping From Java/Scala to C++

JNI bindings
Mapping From Java/Scala to C++

RAPIDS C++ Libraries

UCX Libraries

CUDA



SQL/DATAFRAME PLUGIN

No Code Changes
(none)

Same SQL and Dataframe code.

```
spark.conf.set("spark.rapids.sql.enabled", "true")
```

```
start = time.time()
spark.sql("""
select
    o_orderpriority,
    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()
time.time() - start
```

WHAT WE SUPPORT

and growing...

!	^	concat	double	input_file_block_length	locate	nanvl	rand*	sinh	ucase	TimeSub for time ranges
%	abs	cos	e	input_file_block_start	log	negative	regex_replace*	smallint	upper	startswith
&	acos	cosh	exp	input_file_name	log10	not	replace	spark_partition_id	when	endswith
*	and	cot	expm1		log1p	now	rint	sqrt	window	contains
+	asin	count	first	int	log2	nullif	rollup	string	year	limit
-	atan	cube	first_value	isnan	lower	nvl	row_number	substr		order by
/	avg	current_date	float	isnotnull	max	nvl2	second	substring	~	group by
<	bigint	current_timestamp	floor	isnull	mean	or	shiftright	sum	CSV Reading*	filter
<=	boolean	date	from_unixtime	last	min	pi	shiftright	tan	Orc Reading	union
<=>	cast	datediff	hour	last_value	minute	posexplode*	shiftrightunsigned	tanh	Orc Writing	repartition
=	cbrt	day	if	lcase	mod	position	sign	timestamp	Parquet Reading	equi-joins
==	ceil	dayofmonth	ifnull	like	monotonically_increasing_id	pow	signum	tinyint	Parquet Writing	select
>	ceiling	degrees	in	ln	month	power		trim		
>=	coalesce		initcap			radians	sin		ANSI casts	

IS THIS A SILVER BULLET?

NO

Small amounts of data

Few hundred MB per partition for GPU

Highly cache coherent processing

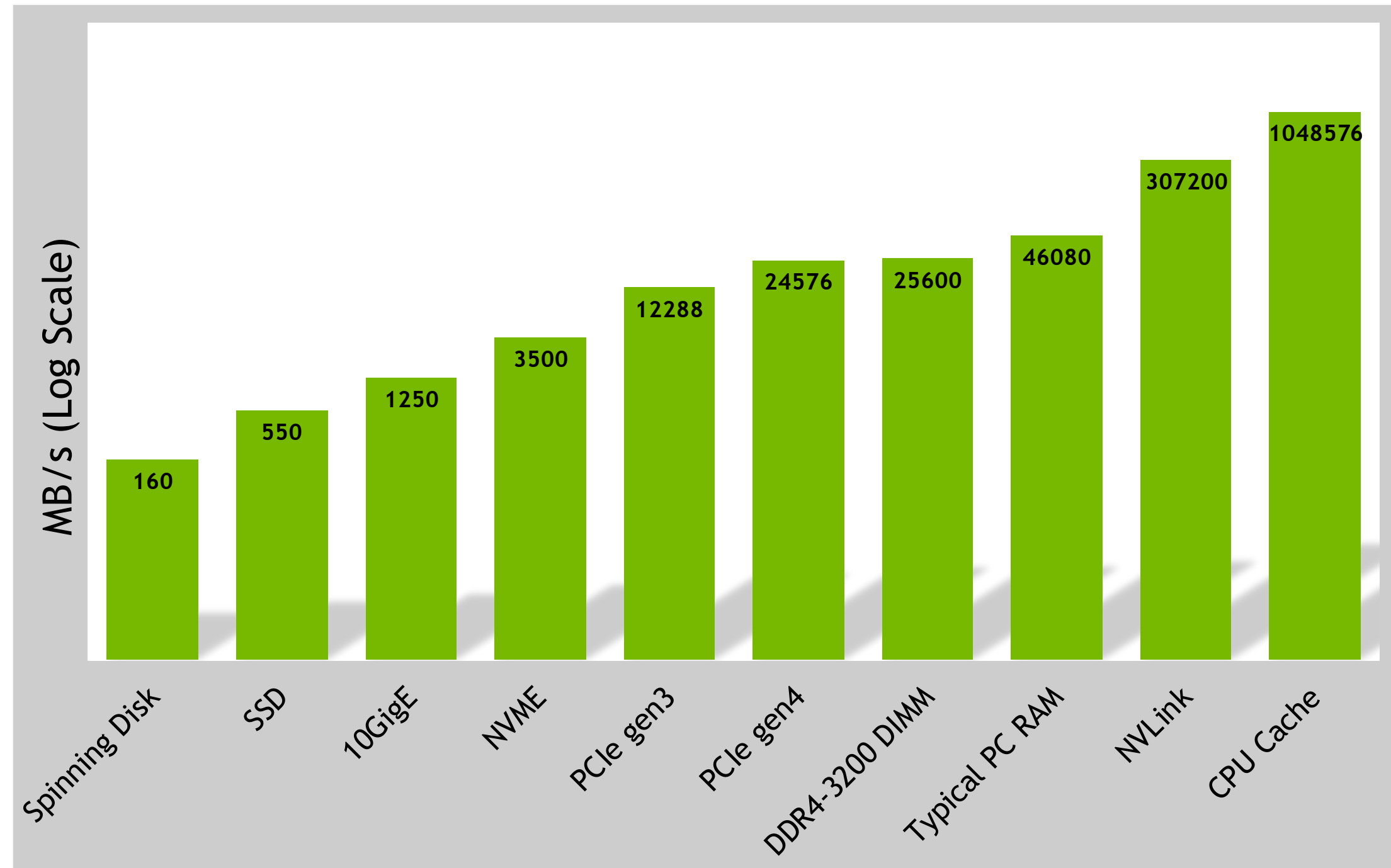
Data Movement

Slow I/O (networking, disks, etc.)

Going back and forth to the CPU (UDFs)

Shuffle

Limited GPU Memory



BUT IT CAN BE AMAZING

What the SQL plugin excels at

High cardinality joins

High cardinality aggregates

High cardinality sort

Window operations (especially on large windows)

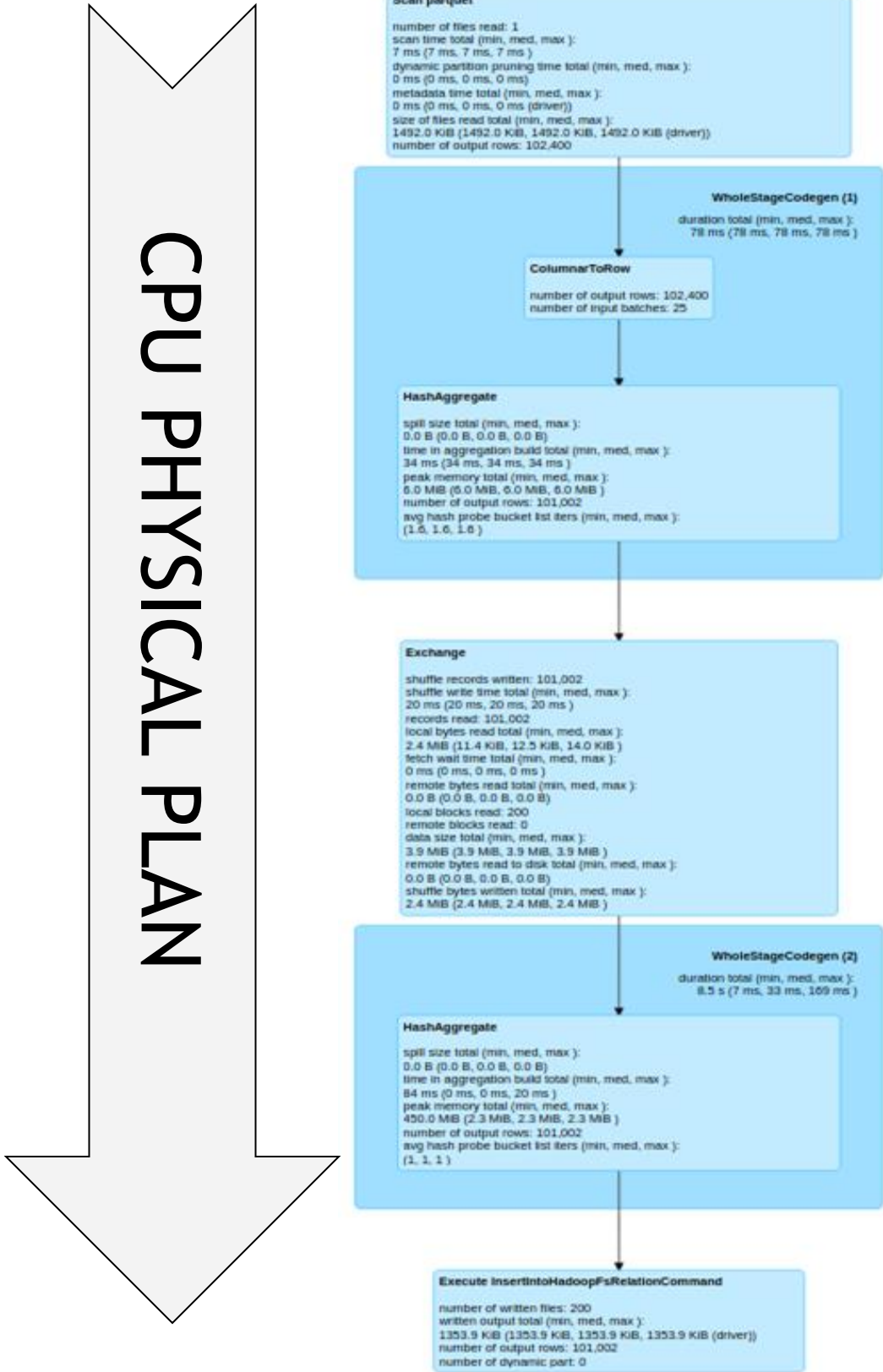
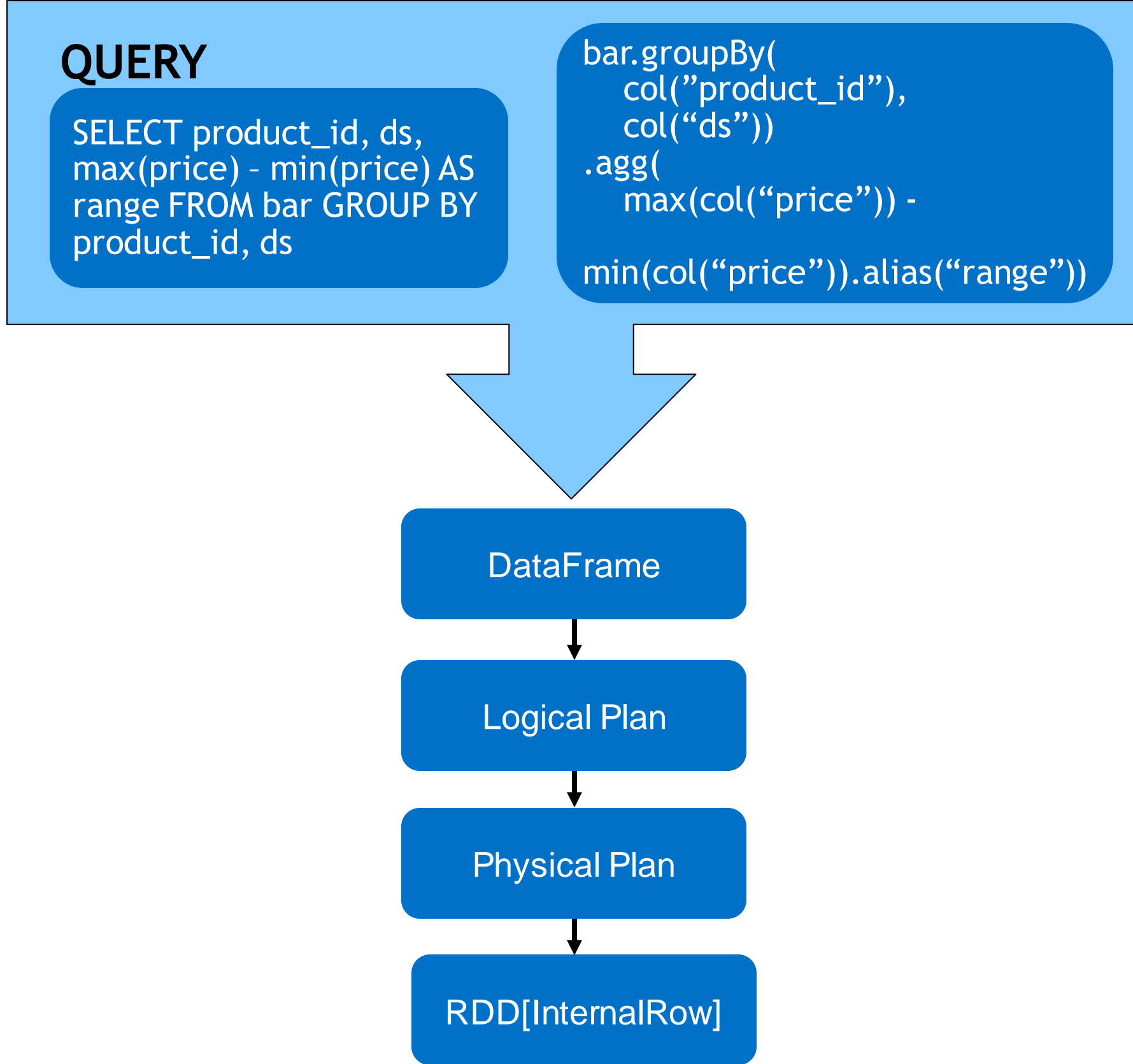
Complicated processing

Transcoding (Writing Parquet and ORC is hard, reading CSV is hard)

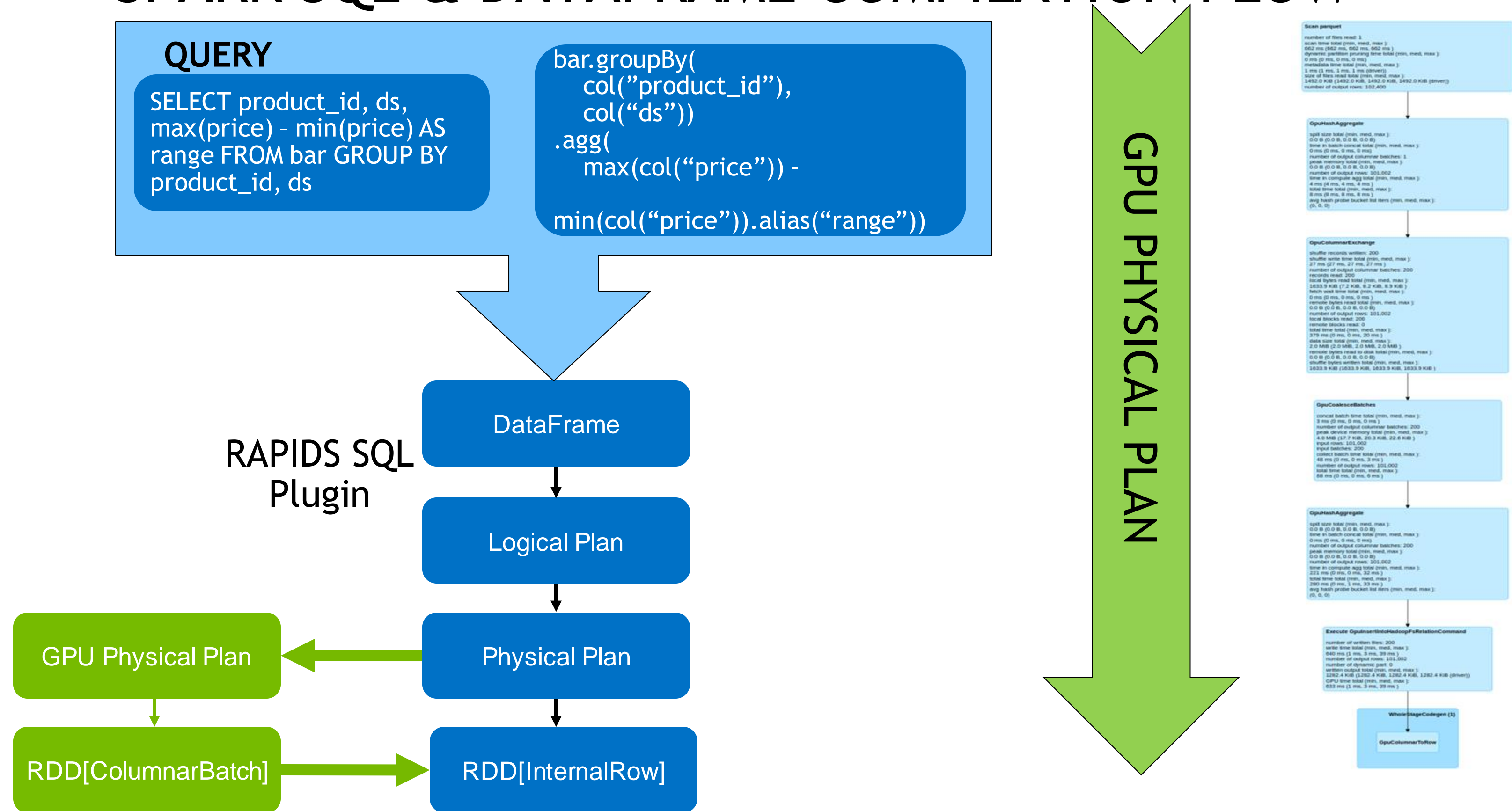


HOW DOES IT WORK

SPARK SQL & DATAFRAME COMPILATION FLOW

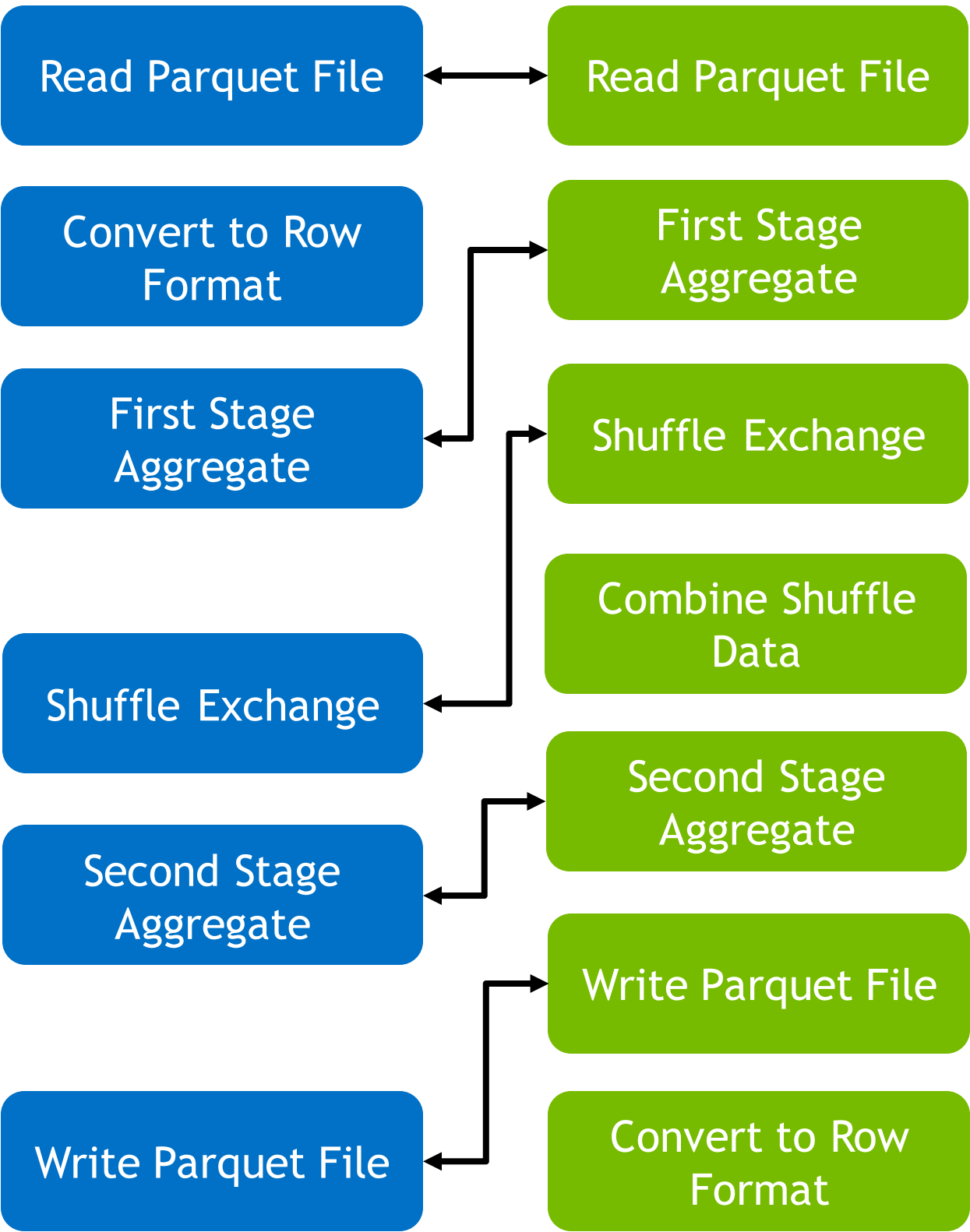
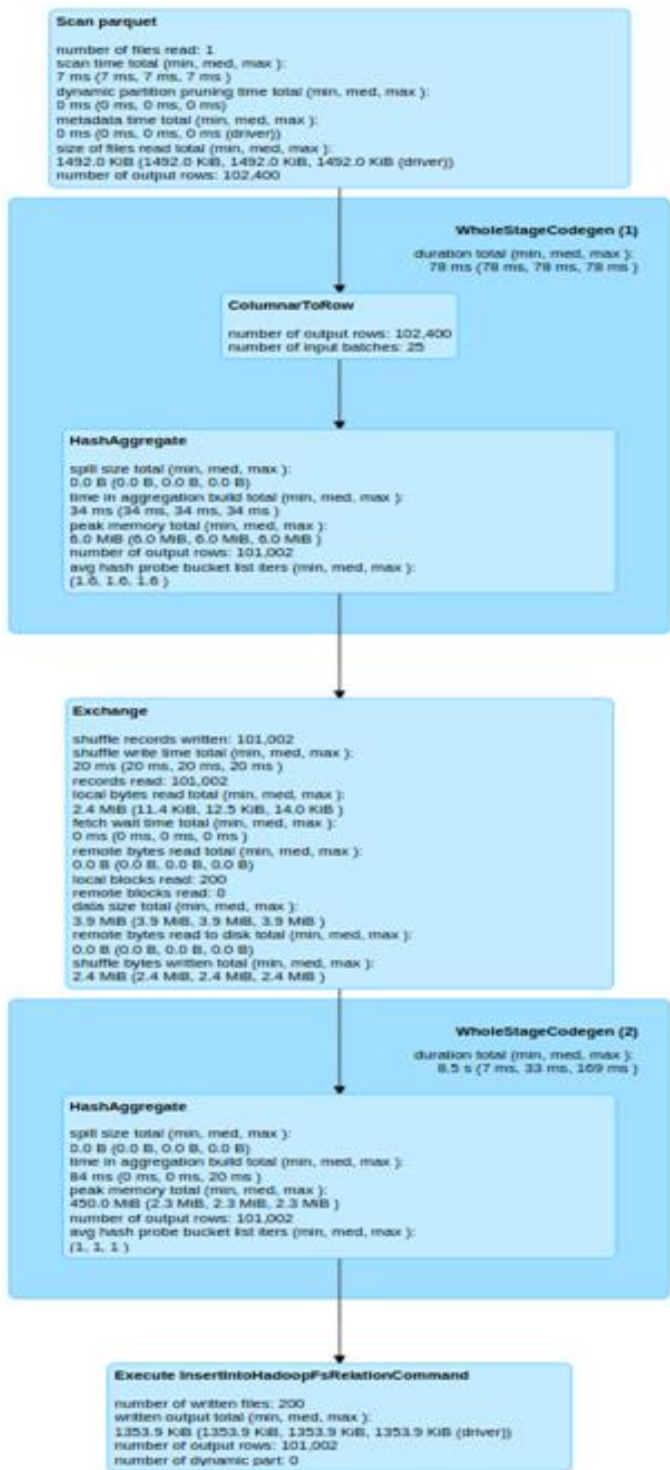


SPARK SQL & DATAFRAME COMPILATION FLOW



SPARK SQL & DATAFRAME COMPILATION FLOW

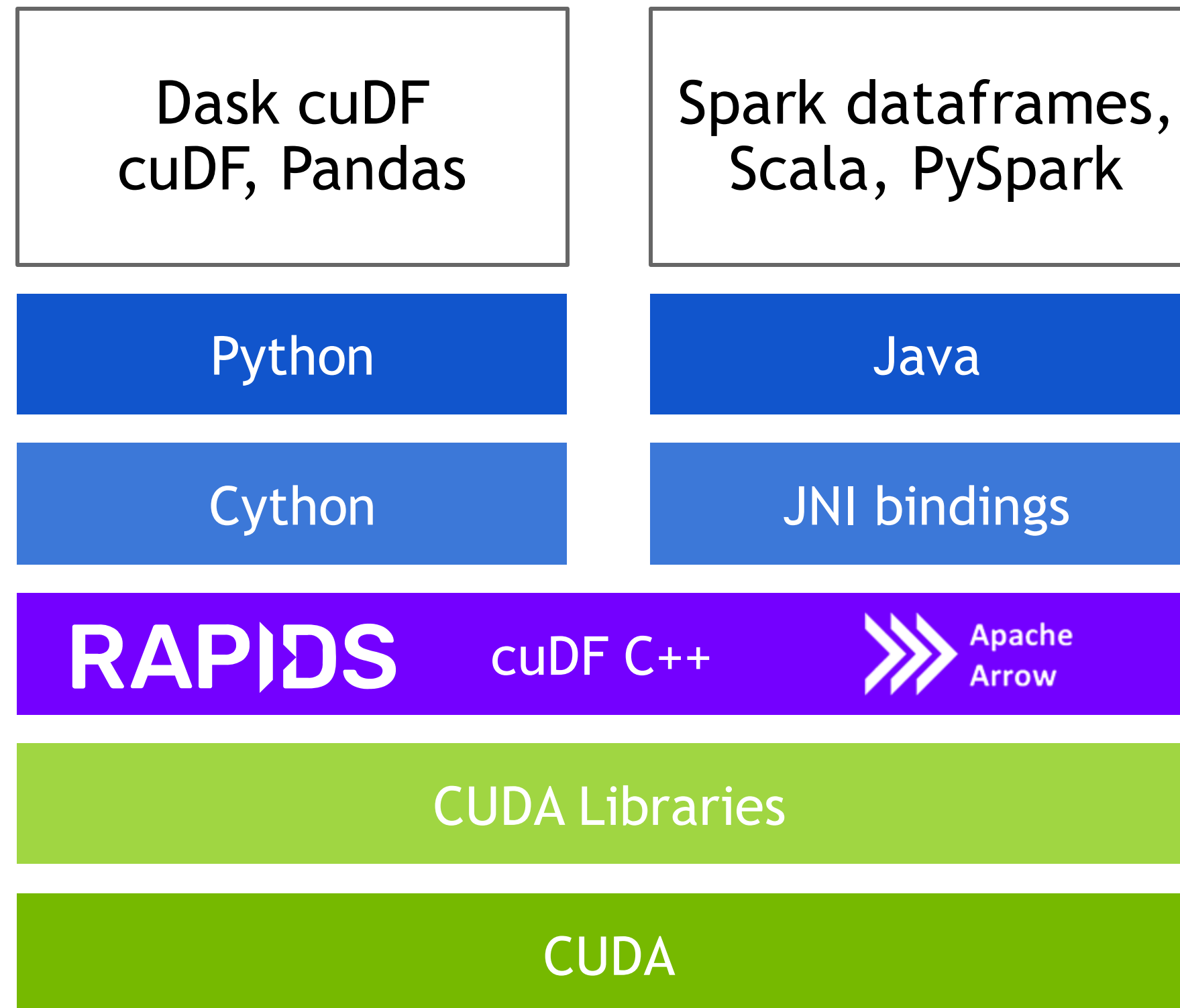
CPU PHYSICAL PLAN



GPU PHYSICAL PLAN



ETL TECHNOLOGY STACK

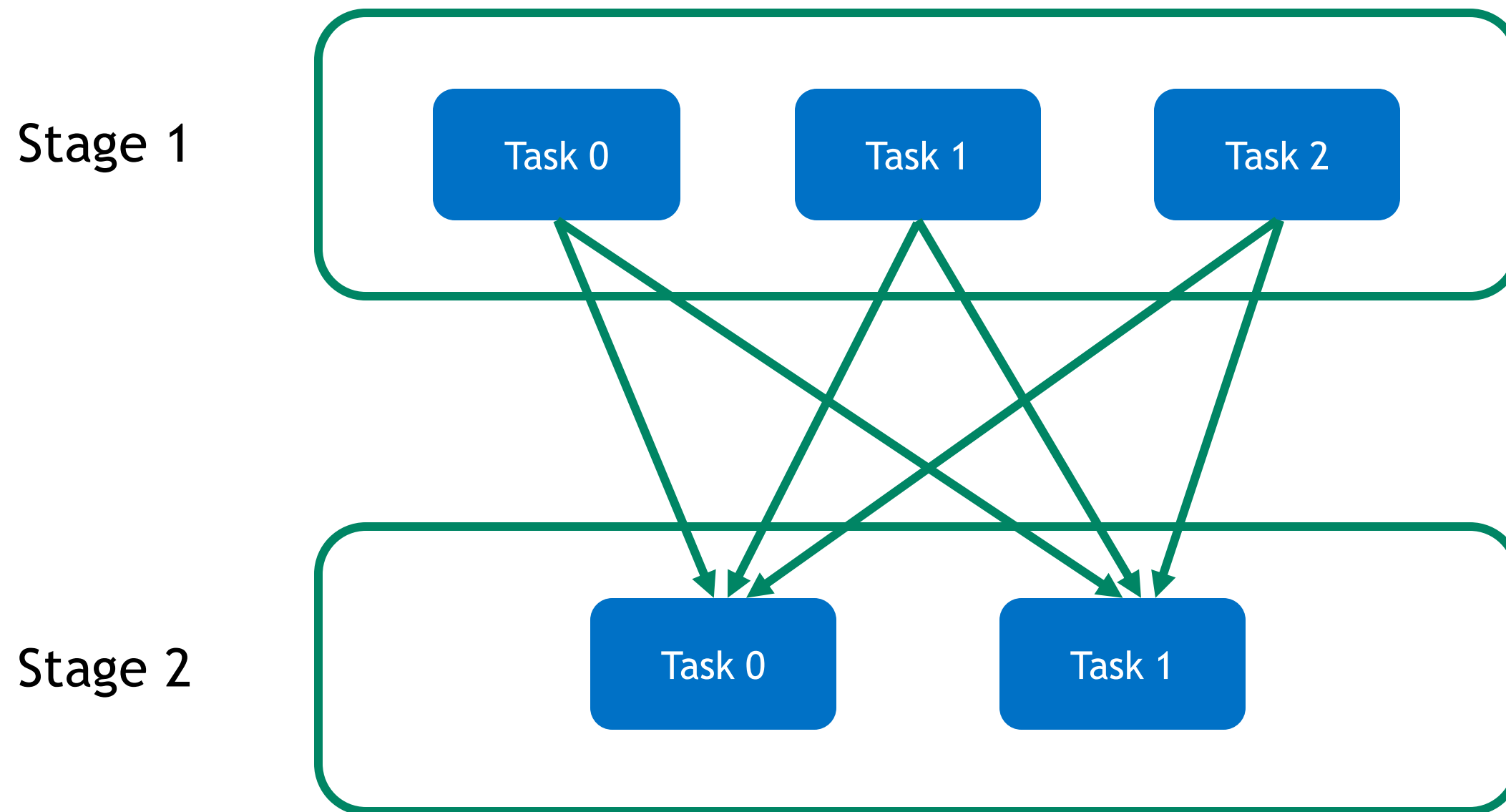




ACCELERATED SHUFFLE

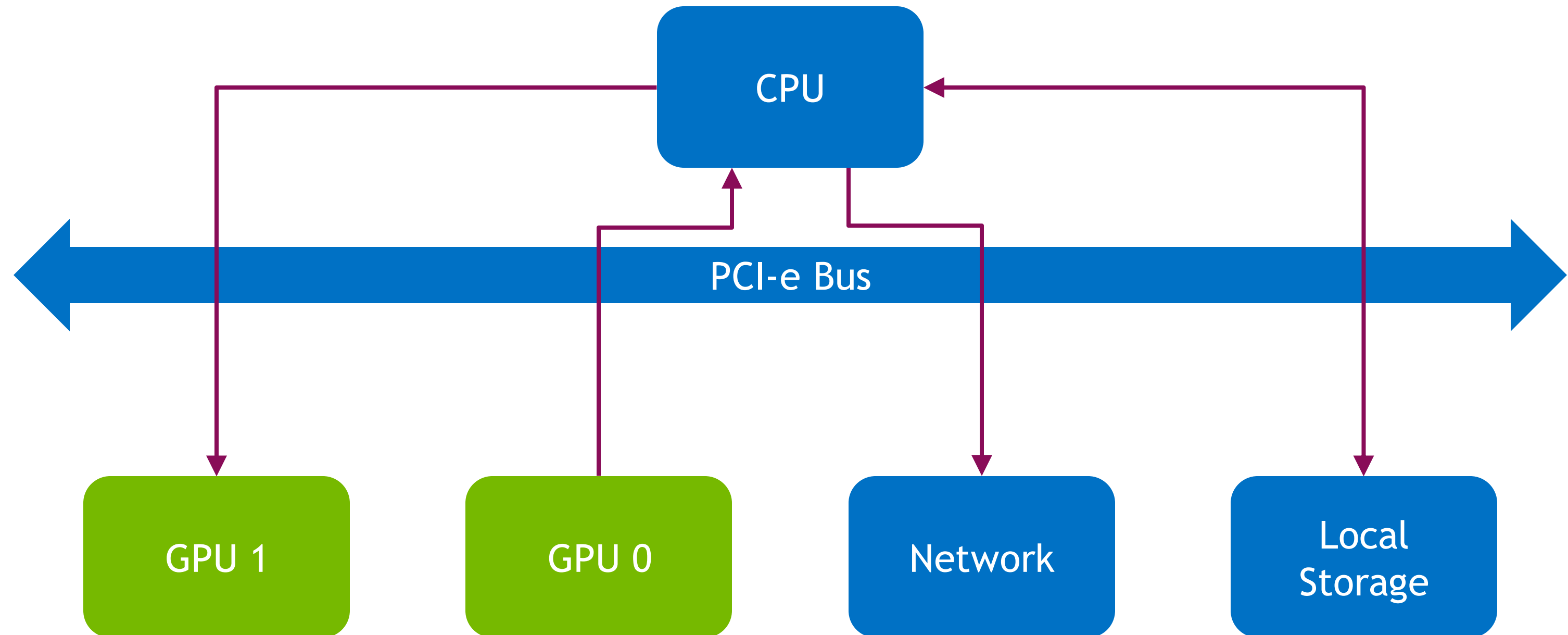
SPARK SHUFFLE

Data Exchange Between Stages



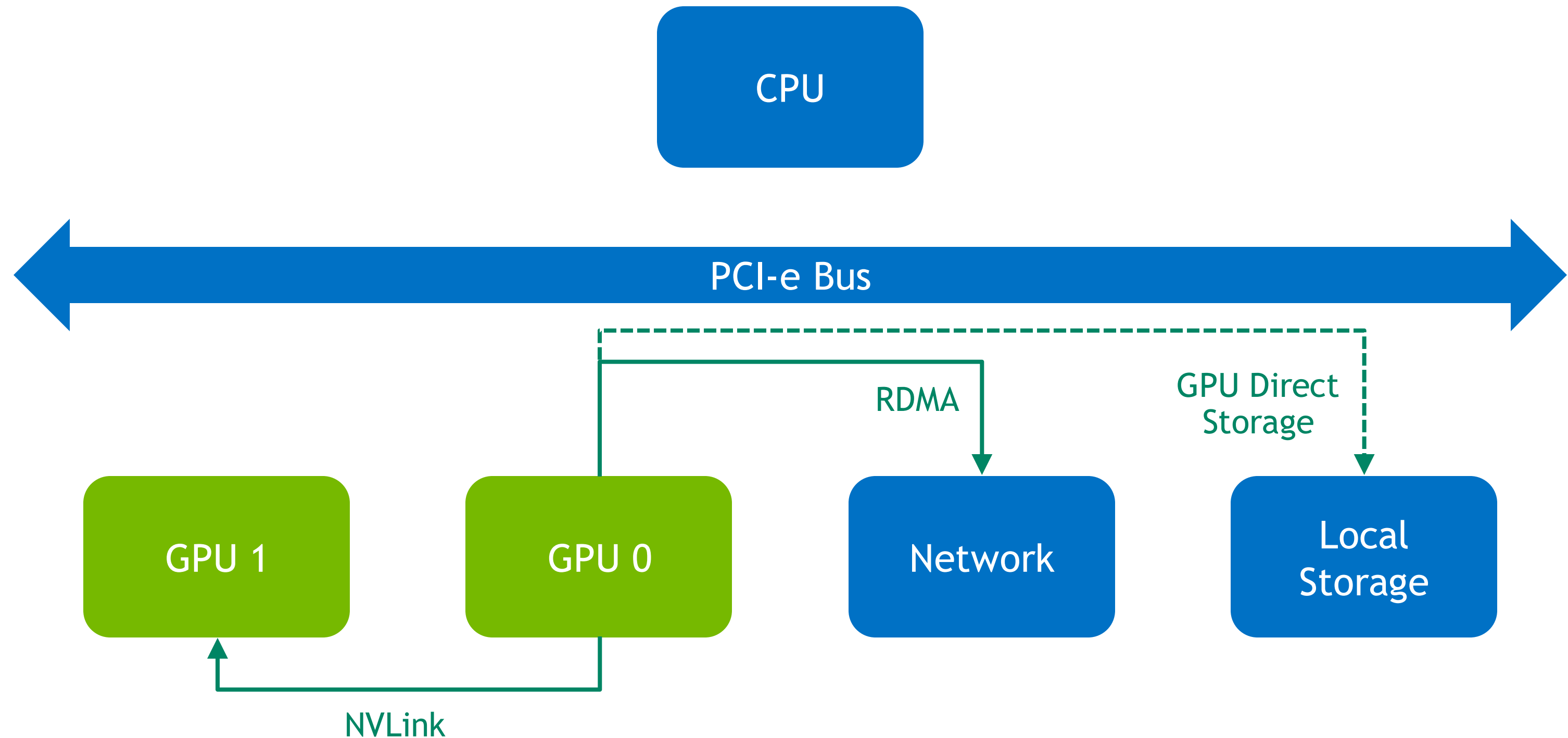
SPARK SHUFFLE

CPU-Centric Data Movement



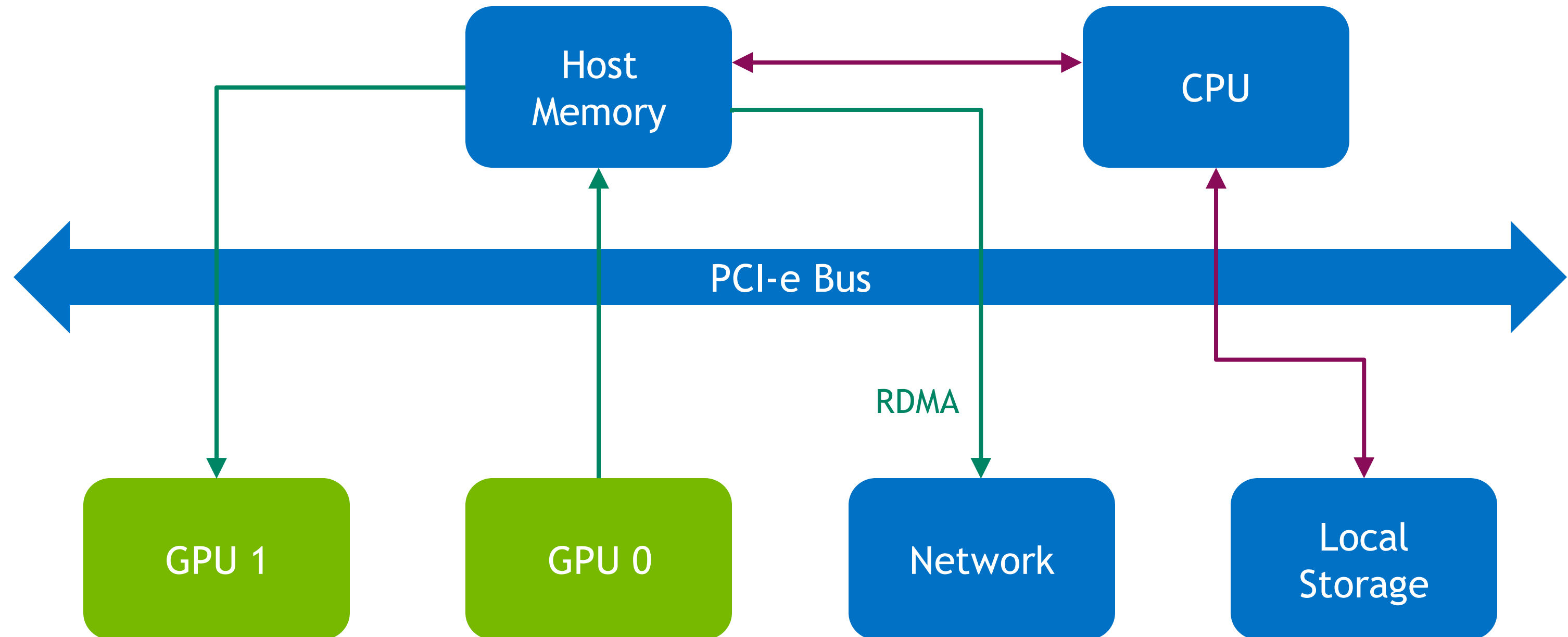
ACCELERATED SPARK SHUFFLE

GPU-Centric Data Movement



ACCELERATED SPARK SHUFFLE

Shuffling Spilled Data



UCX LIBRARY

Unified Communication X

Abstracts communication transports

Selects best available route(s) between endpoints

TCP, RDMA, Shared Memory, GPU

Zero-copy GPU memory transfers over RDMA

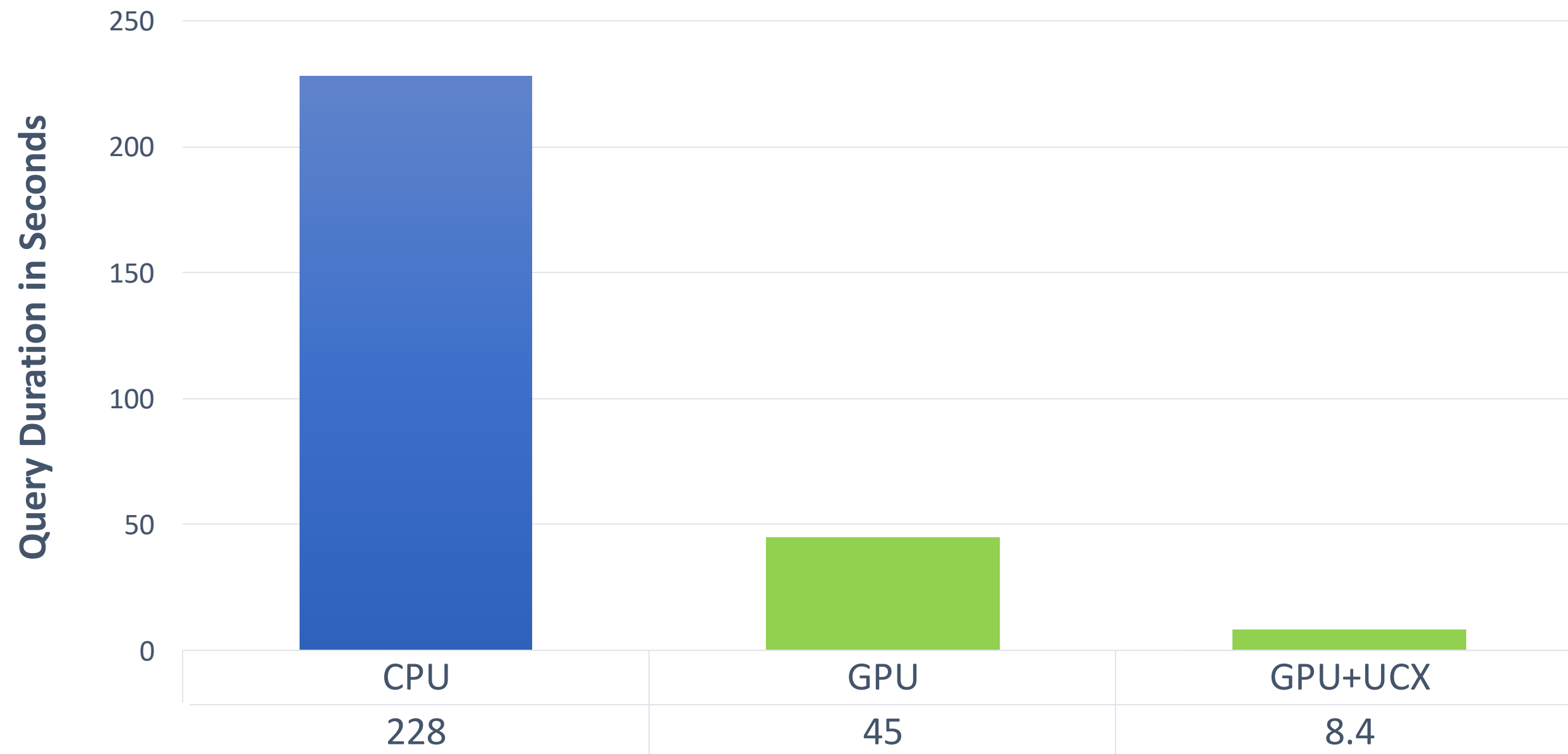
RDMA requires network support (IB or RoCE)

<http://openucx.org>



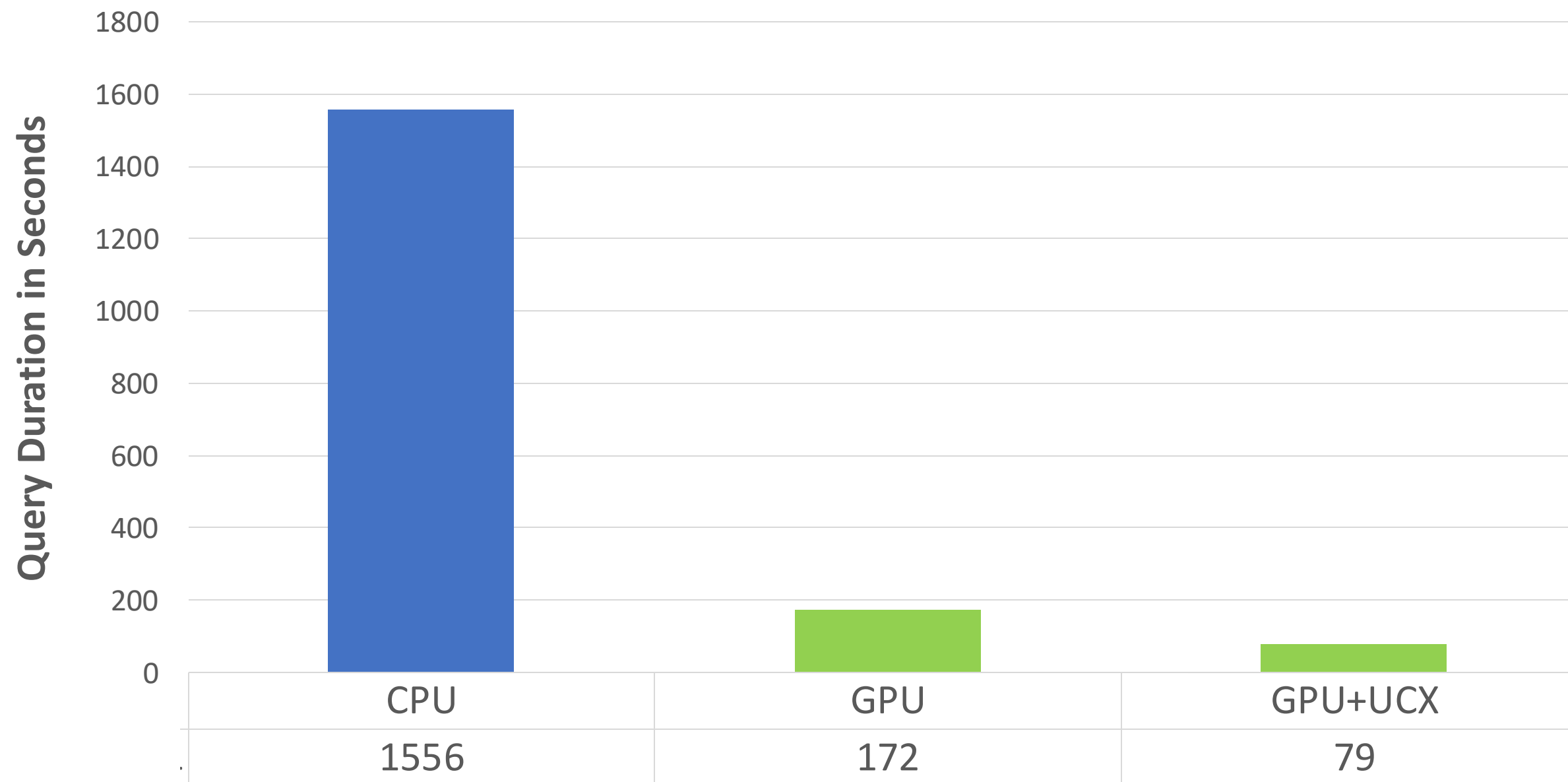
ACCELERATED SHUFFLE RESULTS

Inventory Pricing Query



ACCELERATED SHUFFLE RESULTS

ETL for Logistical Regression Model





WHAT'S NEXT?

WHAT'S NEXT

Open Source/Spark 3.0 Release

Nested types Arrays, Structs, and Maps

Decimal type

More operators



GPU Direct Storage

Time zone support for timestamps
(only UTC for now)

Higher order functions

UDFs

WHERE TO GET MORE INFO

Learn more about the RAPIDS Accelerator for Apache Spark

Visit: [NVIDIA.com/Spark](https://nvidia.com/spark)

Please use the “contact us” to get in touch with NVIDIA’s Spark team

Listen to how Adobe Email Marketing Intelligent Services leverages the RAPIDS Accelerator & Spark 3.0 on Databricks



Upcoming Spark+AI Summit Sessions on GPU support for Apache Spark 3.0:

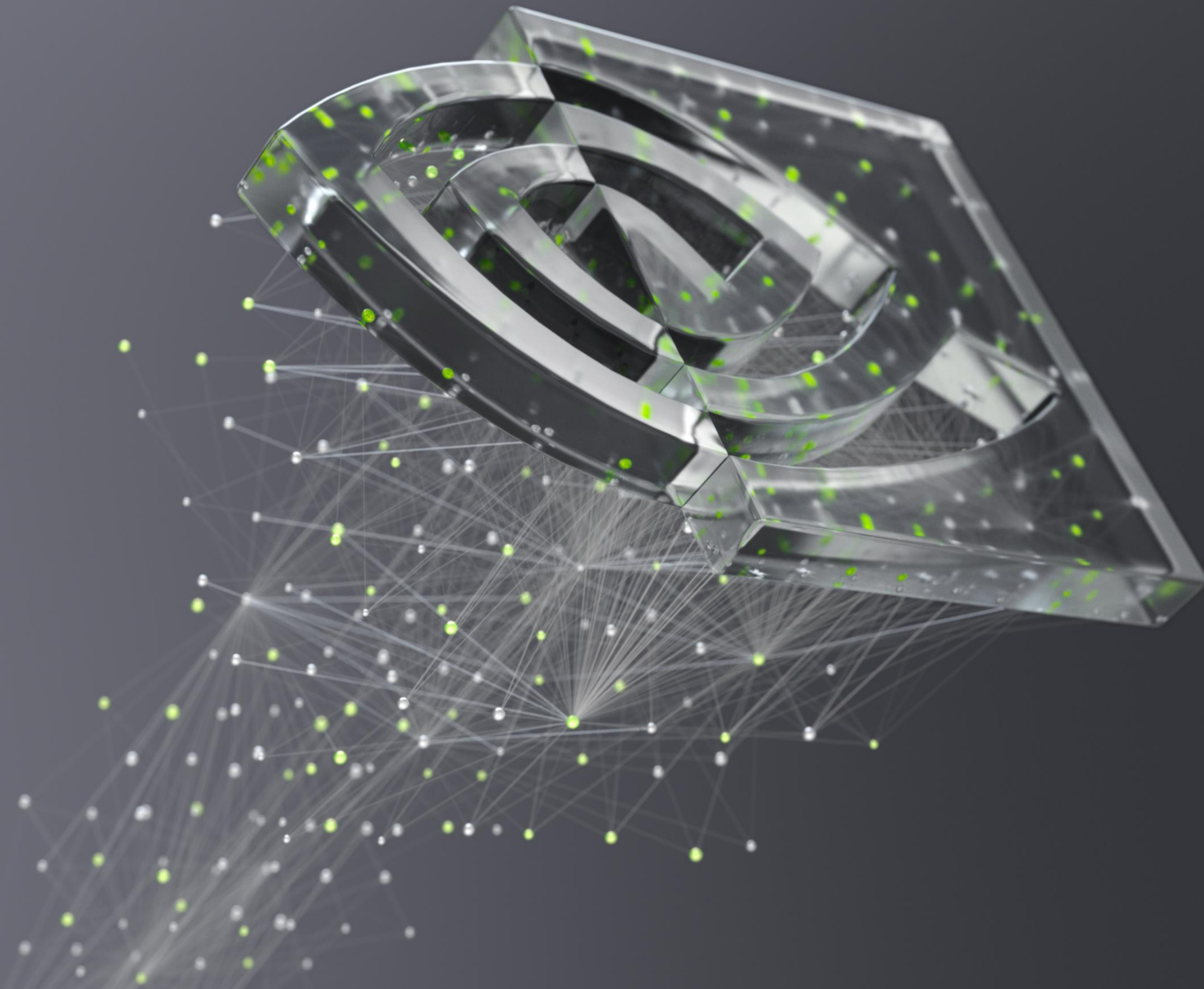
[Deep Dive into GPU Support in Apache Spark 3.x](#)

[Scalable Acceleration of XGBoost Training on Apache Spark GPU Clusters](#)

Preview of Spark 3.0 GPU Features: [NVIDIA.com/Spark-Book](https://nvidia.com/spark-book)

QUESTIONS







BACKUP SLIDES

FAQS

Q: What are the minimum requirements?

A: The RAPIDS accelerator requires:

Apache Spark 3.0

RAPIDS cudf 0.14

CUDA 10.1 or later

NVIDIA GPU with Pascal architecture or later

Ubuntu 16.04+ or CentOS 7+

FAQS

Q: Do all cluster nodes require GPUs?

A: All Spark executors running with the RAPIDS accelerator require their own GPU.

The Spark driver process does not require a node with a GPU.

Q: Can I run more than one executor per GPU?

A: No, there must be a one-to-one mapping between Spark executors and GPUs.

You can run more than one concurrent task per executor.

FAQS

Q: Will the RAPIDS accelerator work in the cloud?

A: Yes, if the VM environment meets the minimum requirements.

Q: Will the RAPIDS accelerator be available for Apache Spark 2.x?

A: No. The columnar processing APIs added in Apache Spark 3.0 are required.

Q: How can I tell if an operation is being accelerated?

A: Accelerated operations appear in the query explanation and SQL UI.

RAPIDS ACCELERATOR CONFIGURATION

`spark.rapids.sql.enabled` is the master enable

`spark.rapids.sql.explain` enables logging of operations not accelerated

`spark.rapids.sql.concurrentGpuTasks` controls concurrent task count per GPU

SPARK ACCELERATOR-AWARE SCHEDULING

Tracking JIRA: SPARK-24615

Request executor and driver resources (GPU, FPGA, etc.)

Resource discovery

Specify task resources

API to determine assigned resources

YARN, Kubernetes, and Standalone

SPARK ACCELERATOR-AWARE SCHEDULING

Sample Command-Line

```
./bin/spark-shell --master yarn --executor-cores 2 \  
--conf spark.driver.resource.gpu.amount=1 \  
--conf spark.driver.resource.gpu.discoveryScript=/opt/spark/getGpuResources.sh \  
--conf spark.executor.resource.gpu.amount=2 \  
--conf spark.executor.resource.gpu.discoveryScript=./getGpuResources.sh \  
--conf spark.task.resource.gpu.amount=1 \  
--files examples/src/main/scripts/getGpusResources.sh
```

SPARK STAGE LEVEL SCHEDULING

Tracking JIRA: SPARK-27495

Specify task resource requirements per RDD operation

Dynamically allocates containers to meet resource requirements

Schedules tasks on appropriate containers