DLI Accelerated Data Science Teaching Kit

# Lecture 11.4 – RAPIDS and Spark

# What is RAPIDS Accelerator for Spark?

- Accelerate Spark distributed computing framework using GPUs, via
  - RAPIDS cuDF library
  - Accelerated shuffle based on UCX (GPU-to-GPU communication)
- Existing Spark applications run with no code change
  - Launch Spark with RAPIDS accelerator (plugin jar)
  - Enable configuration setting

```
spark.conf.set('spark.rapids.sql.enabled','true')
```
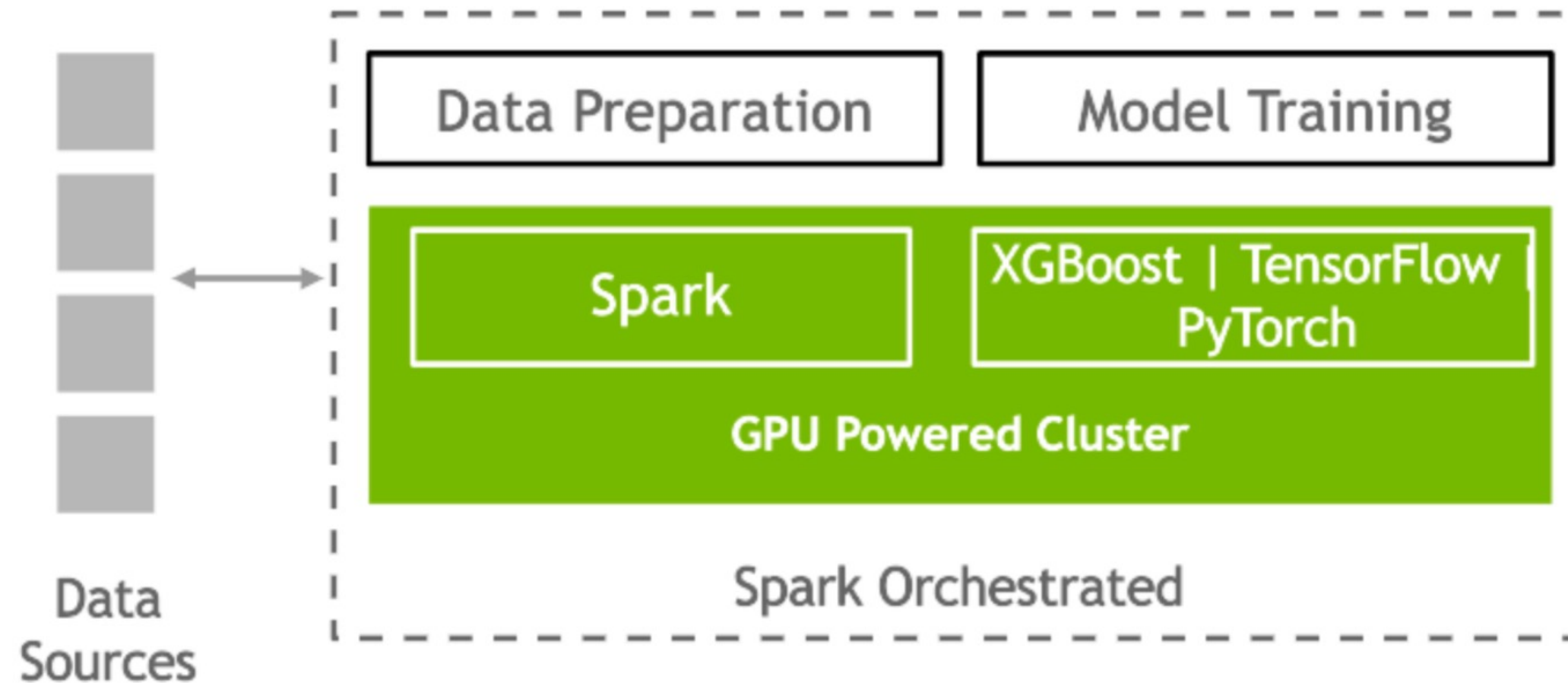
https://nvidia.github.io/spark-rapids/

# Spark 3.0 Offers Unified AI framework for ETL + ML/DL
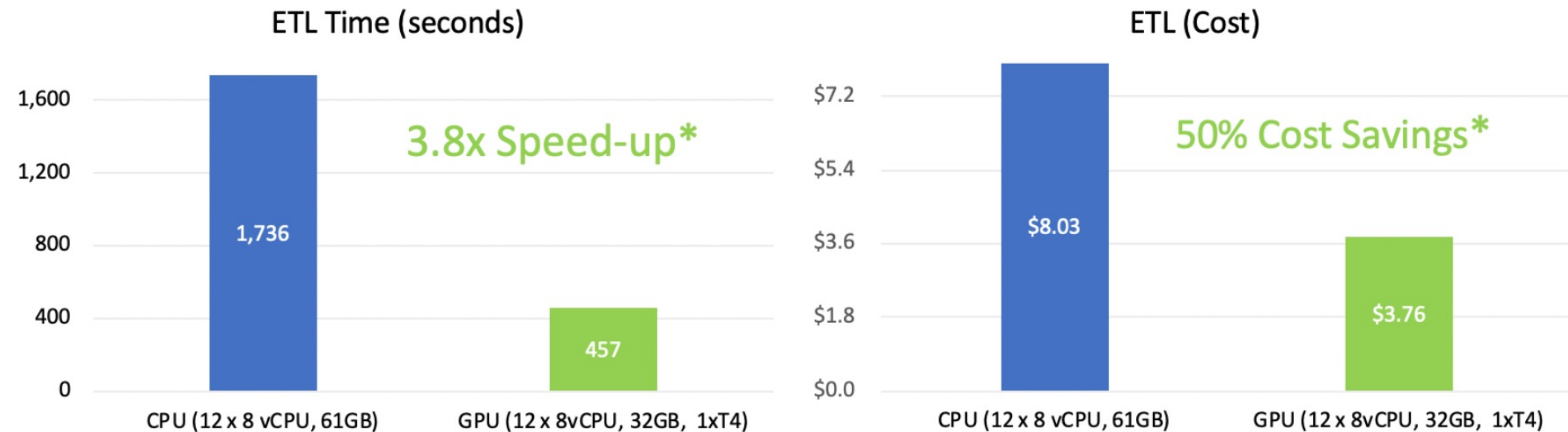
Single pipeline from data input to model training

# Accelerating Spark with RAPIDS

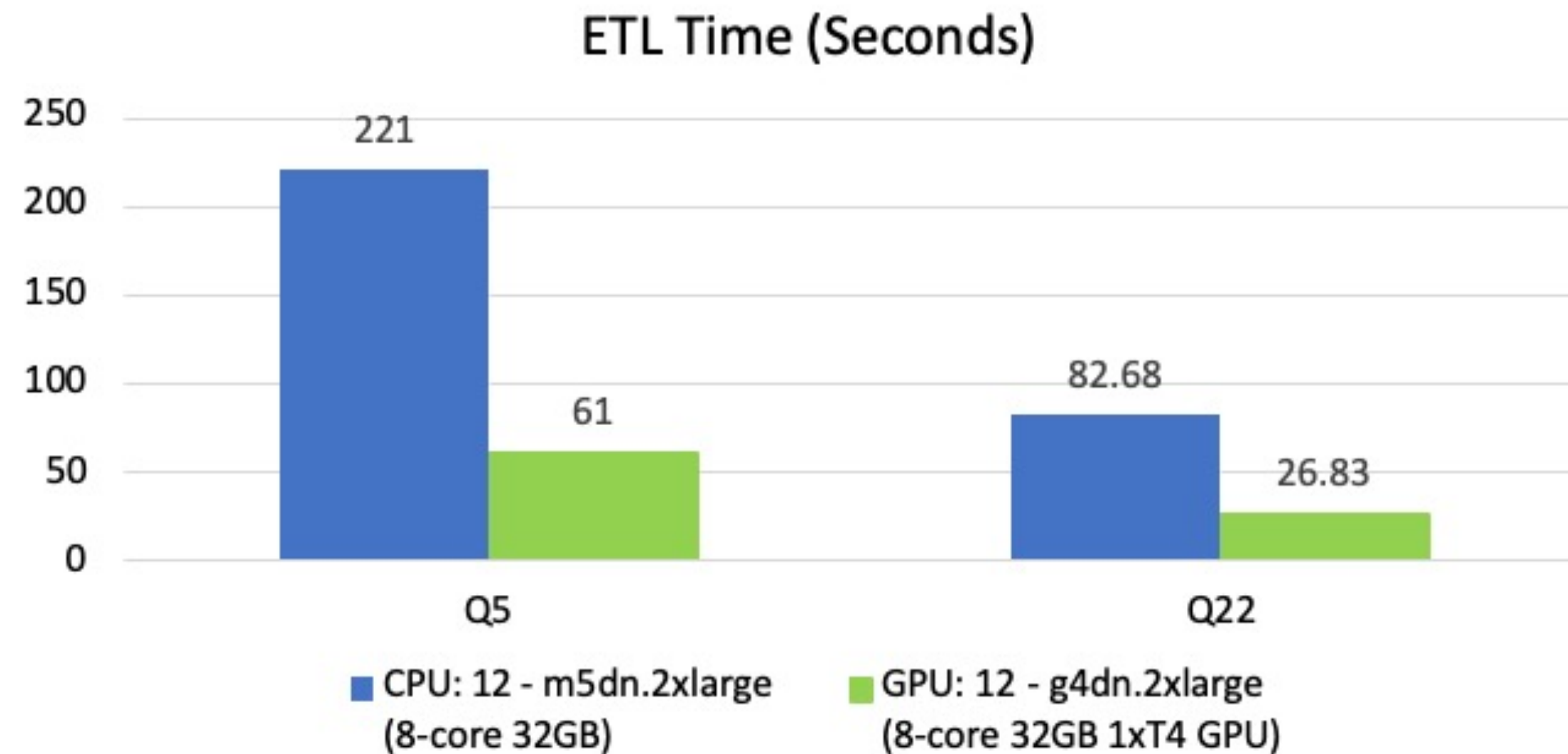Higher speed, and lower costs



ETL Time (seconds)

3.8x Speed-up*

1,736 — CPU (12 x 8 vCPU, 61GB)
457 — GPU (12 x 8vCPU, 32GB, 1xT4)

ETL (Cost)

50% Cost Savings*

$8.03 — CPU (12 x 8 vCPU, 61GB)
$3.76 — GPU (12 x 8vCPU, 32GB, 1xT4)

*ETL for FannieMae Mortgage Dataset (~200GB) as shown in our demo. Costs based on Cloud T4 GPU instance market price & V100 GPU price on Databricks Standard edition

# Accelerating Spark with RAPIDS on AWS

Higher speed, and lower costs

ETL Time (Seconds)

~3.5x Speed-up
~40% Cost Savings

- CPU: 12 - m5dn.2xlarge (8-core 32GB)
- GPU: 12 - g4dn.2xlarge (8-core 32GB 1xT4 GPU)

Based on TPCx-BB like Queries #5 & #22 with 1TB scale factor input

# GPU Scheduling Example: Starting Code

```
./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
```

# GPU Scheduling Example: Discovery Script

```bash
#!/bin/bash  #
# Outputs a JSON formatted string that is expected by the
# spark.{driver/executor}.resource.gpu.discoveryScript config.  #
# Example output: {"name": "gpu", "addresses":["0","1","2","3","4","5","6","7"]}

ADDRS=$(nvidia-smi --query-gpu=index --format=csv,noheader \
        | sed -e :a -e N -e'$!ba' -e 's/\n/","/g')
echo {\"name\": \"gpu\", \"addresses\":[\"$ADDRS\"]}
```

# GPU Scheduling Example: Assignments API

```scala
// Task API
val context = TaskContext.get()
val resources = context.resources()
val assignedGpuAddrs = resources("gpu").addresses
// Pass assignedGpuAddrs into TensorFlow or other AI code


// Driver API
scala> sc.resources("gpu").addresses
Array[String] = Array(0)
```

# GPU Scheduling Example: Schedule API

# Sparks + RAPIDS Resources

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

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

- https://www.nvidia.com/en-us/deep-learning-ai/solutions/data-science/apache-spark-3/

- https://ngc.nvidia.com

DLI Accelerated Data Science Teaching Kit

# Questions?