



## Explore GPU accelerated Data Preparation and Machine Learning with Nvidia RAPIDS

Dustin VanStee - Data Scientist  
Loic Fura - Power System Specialist

*IBM Worldwide Client Experience Centers*

**2019** IBM Systems Technical University  
Oct 7-11 | Las Vegas, NV



# IBM Systems Worldwide Client Experience Centers

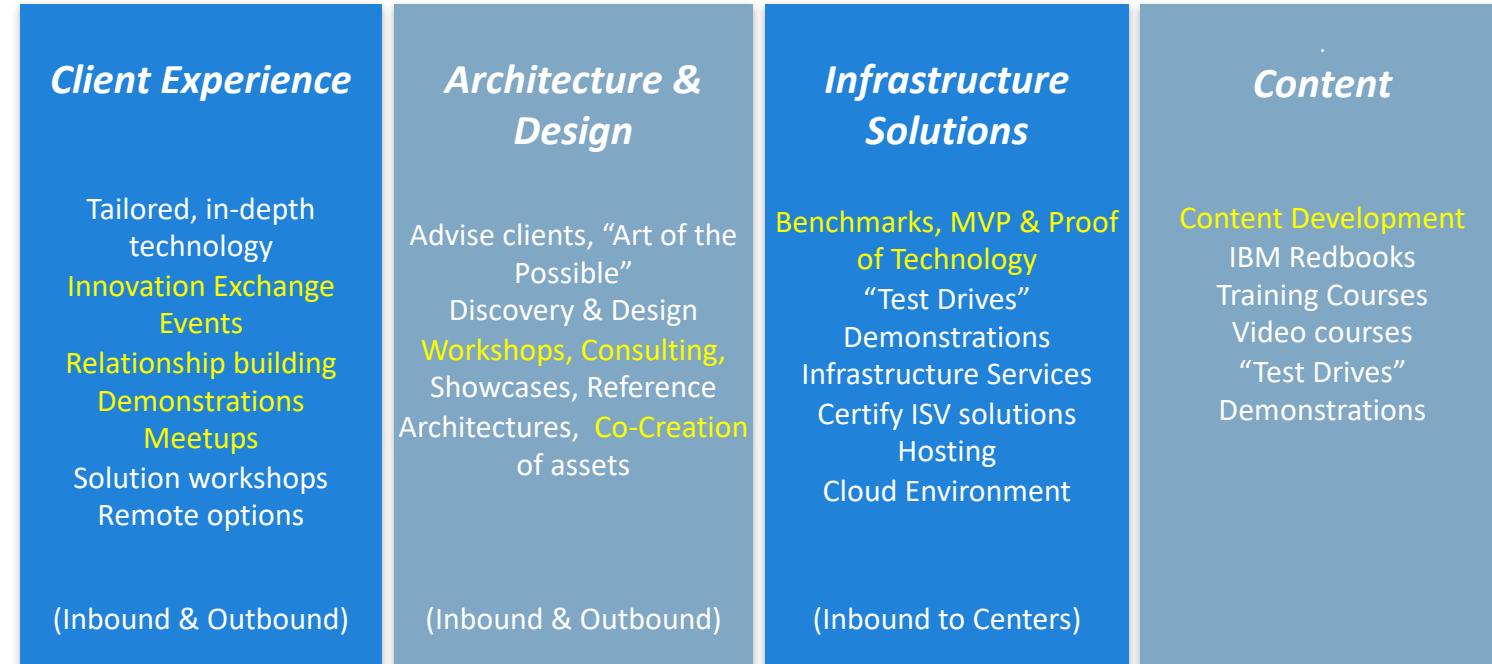


## **IBM Systems Worldwide Client Experience Centers**

maximize IBM Systems competitive advantage in the Cloud and Cognitive era by providing access to world class **technical experts** and **infrastructure services** to assist Clients with the transformation of their IT implementations..

### **9 Worldwide Locations (\* also Infrastructure Hubs):**

Austin TX , \*Poughkeepsie NY, Rochester MN,  
Tucson AZ, \*Beijing CHINA, Boeblingen GERMANY,  
Guadalajara MEXICO,\*Montpellier FRANCE, Tokyo  
JAPAN



**NEW:** *Co-Creation Lab; CEC Cloud; RedHat Center of Competency*

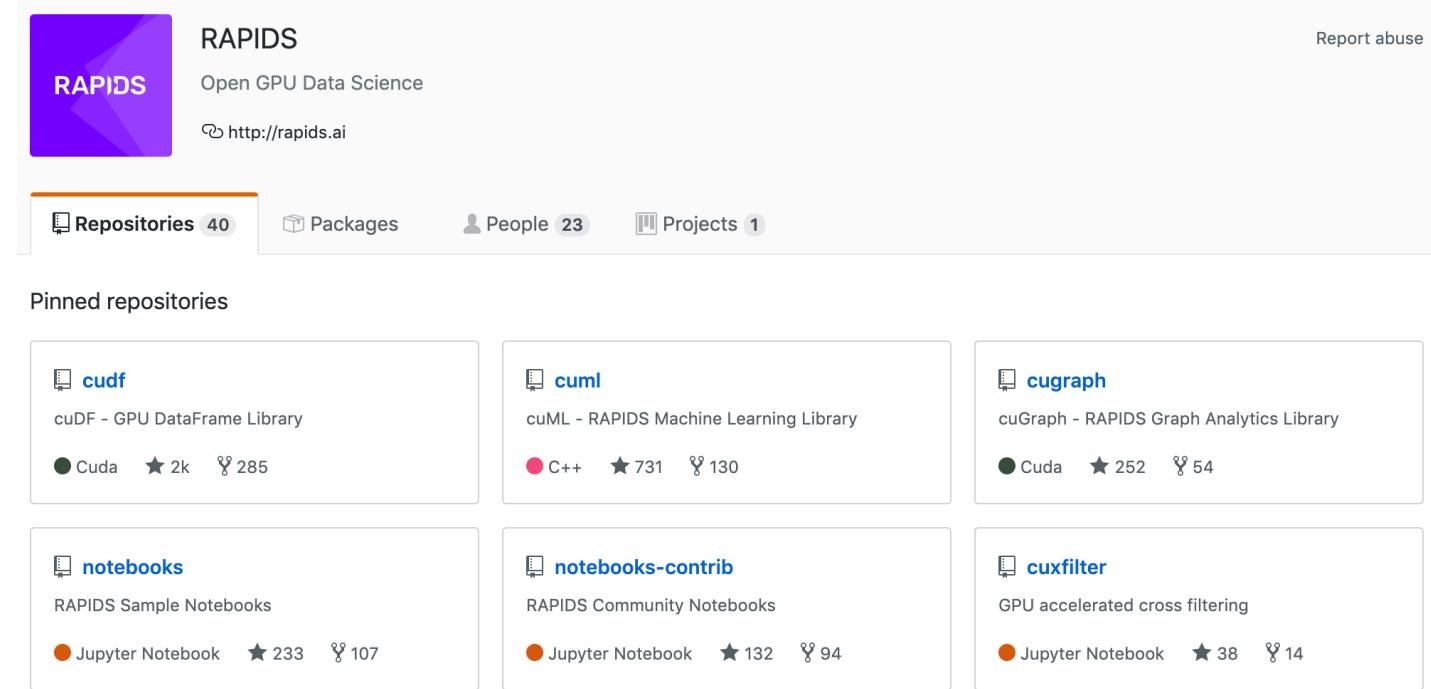
**For further information, please contact the Centers via email at:**  
[ccenter@us.ibm.com](mailto:ccenter@us.ibm.com)

# Session Objectives

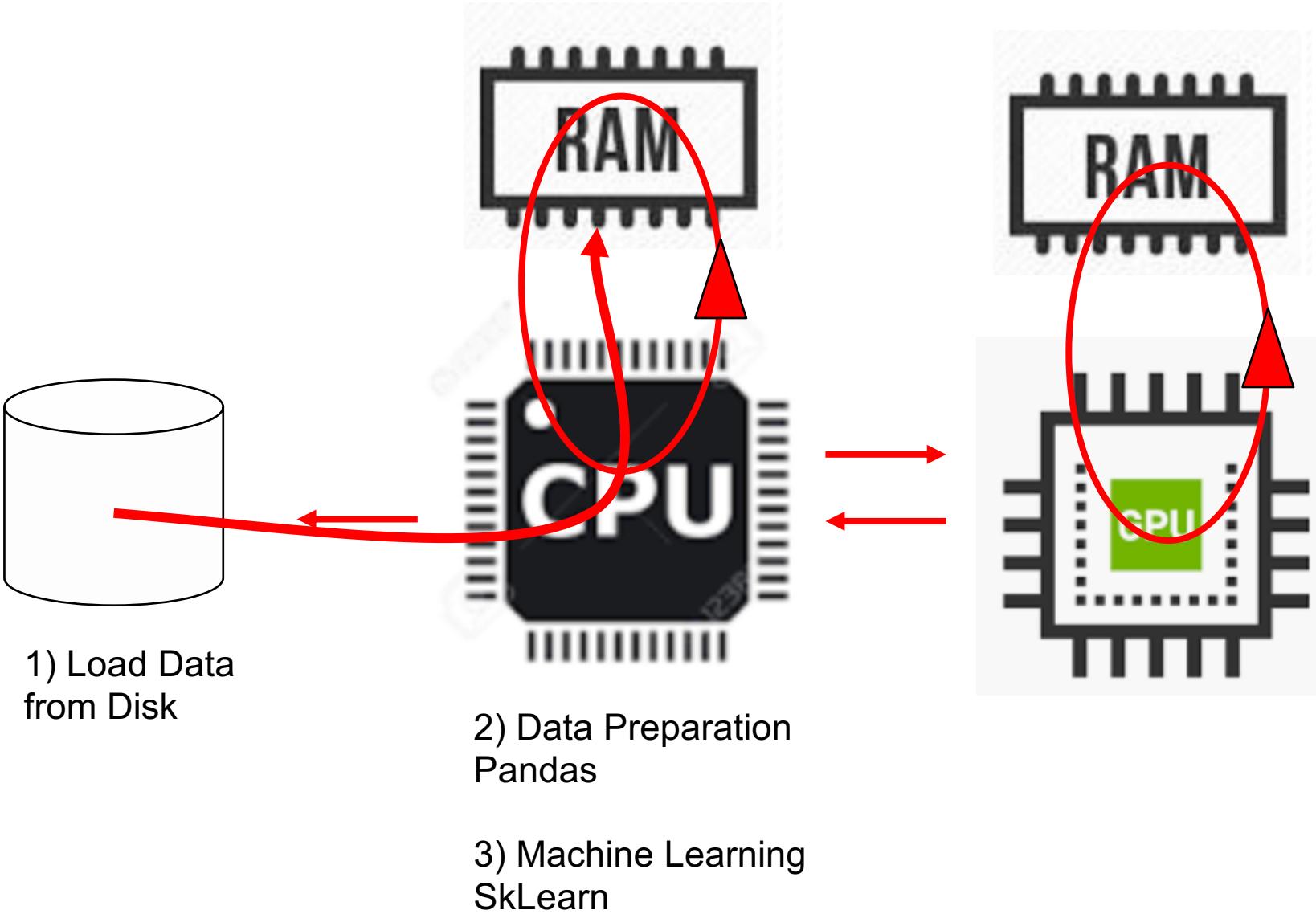
- What is RAPIDS
- Motivation for RAPIDS
- RAPIDS architecture
- How can I get RAPIDS on AC922 ?
- Quick Lab Demo to get started

# What is RAPIDS

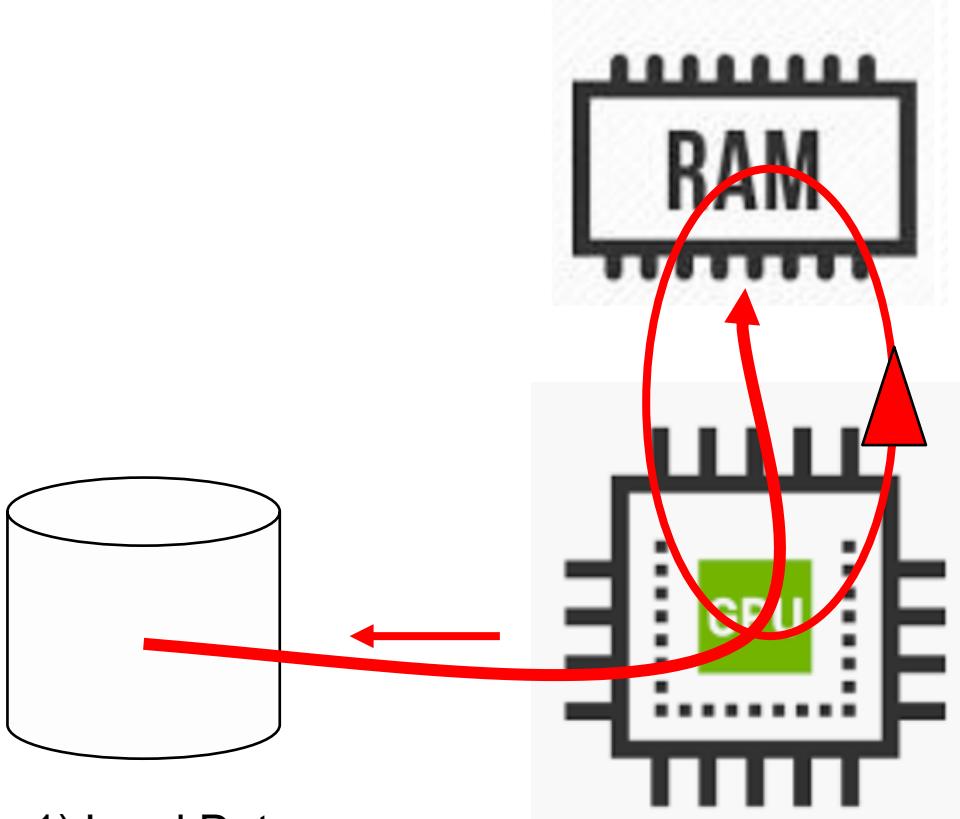
- RAPIDS is a software package from Nvidia
- Accelerates **Data Preparation** and **Machine Learning**
- Eliminates a lot of CPU to GPU transfers
- Meant to be used in place of Pandas and Sklearn python libraries
- Provides 10x-50x performance boost for many tasks
- Scales to multiple GPUs
- Designed to integrate with DL frameworks too



# Typical CPU Based Workflow



# RAPIDS Based Workflow



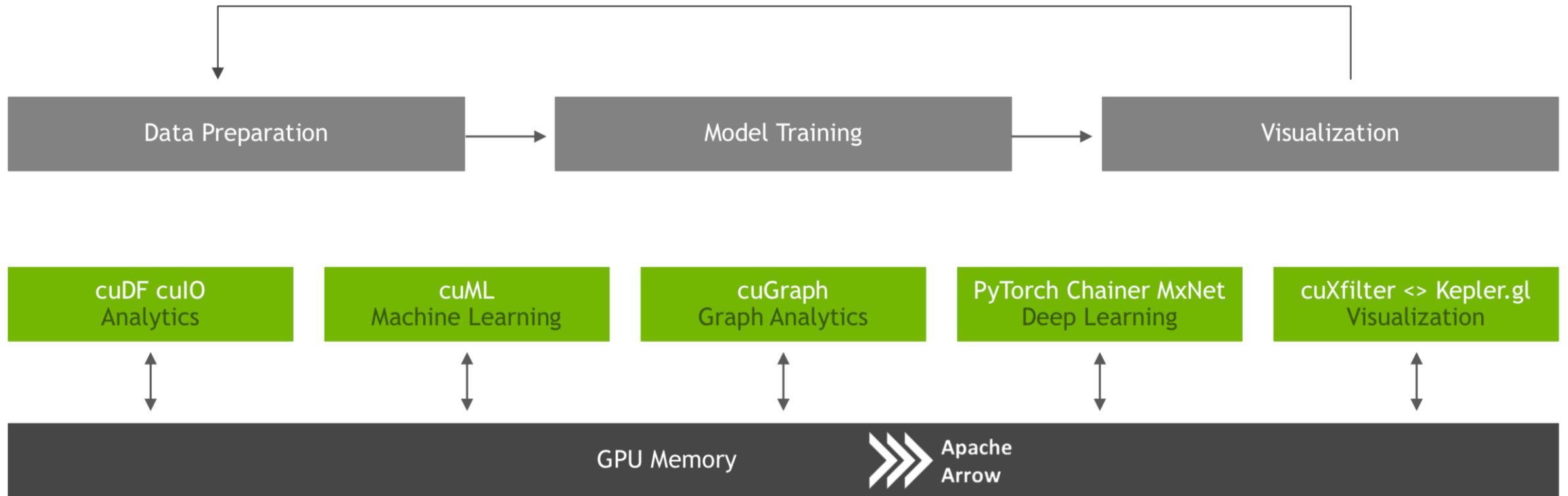
1) Load Data  
from Disk

2) Data Preparation  
Pandas

3) Machine Learning  
SkLearn

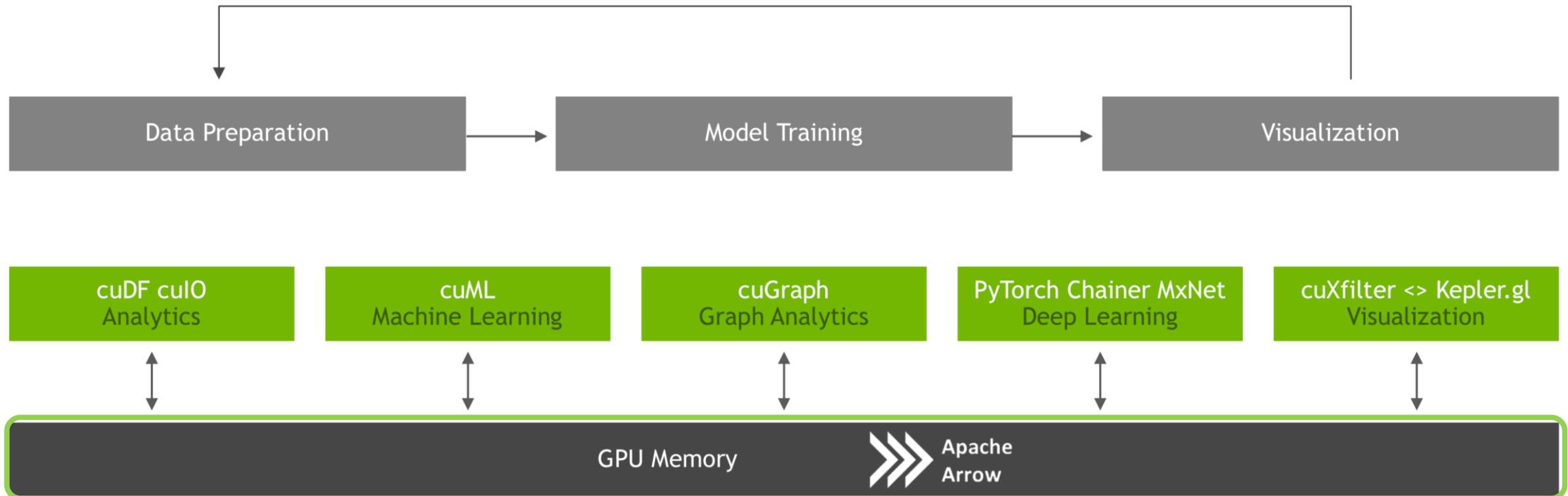
# RAPIDS

## End to End Accelerate GPU Data Science

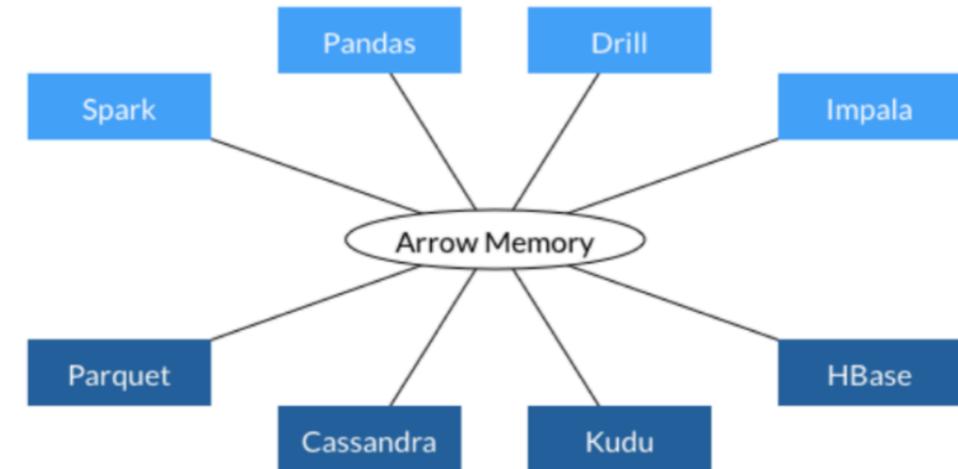
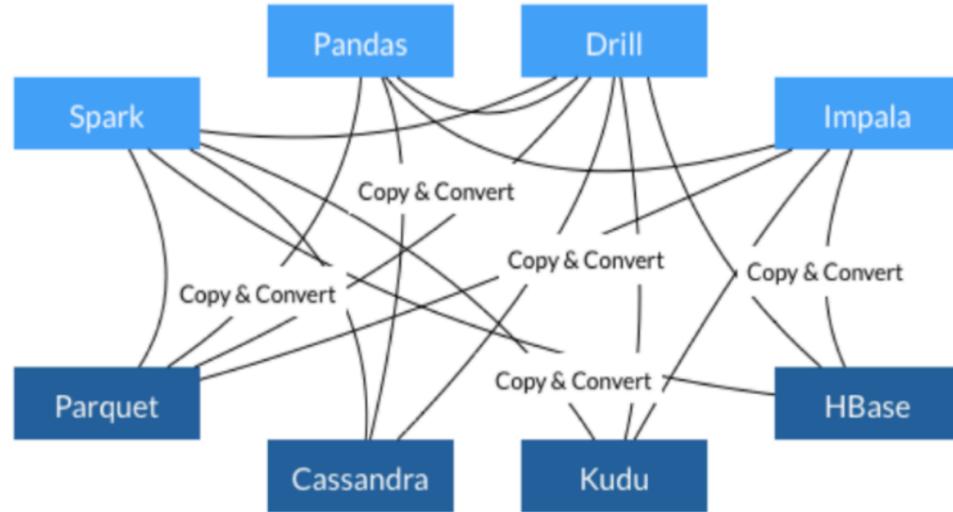


# RAPIDS

## End to End Accelerate GPU Data Science



# LEARNING FROM APACHE ARROW ➤➤➤



- 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 (eg, Parquet-to-Arrow reader)

From Apache Arrow Home Page - <https://arrow.apache.org/>

# Data Preparation – Speeding up the 80%

## cuDF is comprised of

### Python

- Python API for data prep
- Mirror Pandas API
- Creates Apache Arrow data frames
- JIT completion for UDFs via Numba



### CUDA

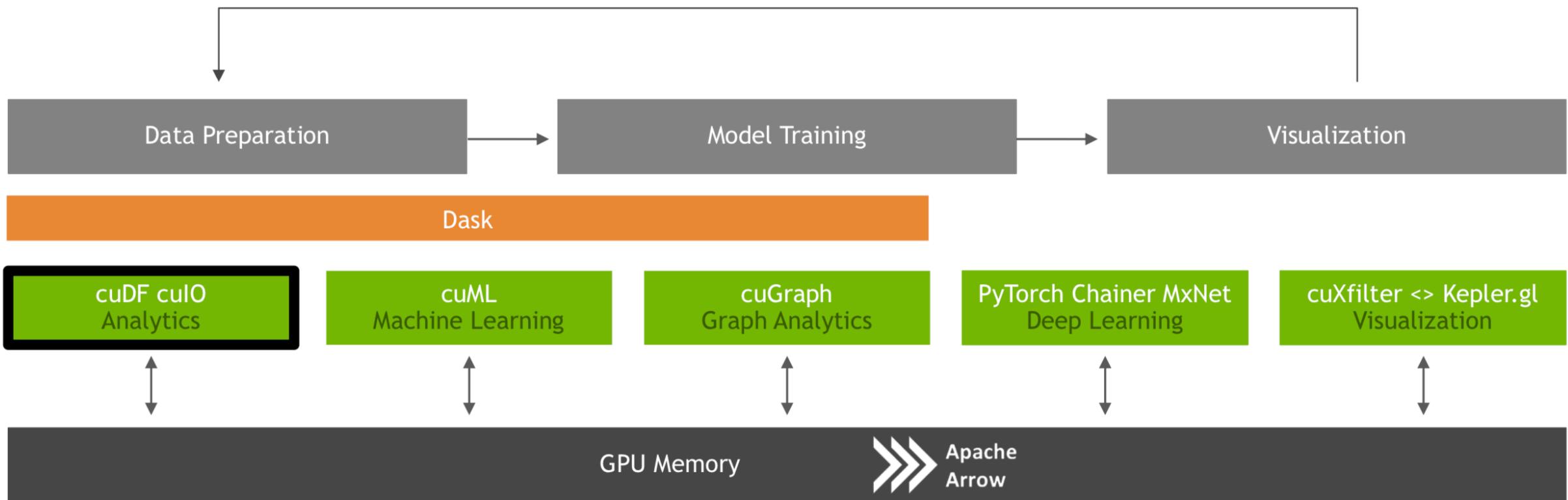
- Low level implementations for data prep on GPU
- Implements functions for data movement using Apache Arrow
- CUDA kernels to perform element wise opps
- CUDA sort/join/groupby/reduce

The screenshot shows the RAPIDS Docs API Reference for the cuDF DataFrame class. The top navigation bar includes tabs for RAPIDS Docs, API DOCS, STABLE, NIGHTLY, and LEGACY, along with links for cuDF, cuML, cuGraph, nvStrings, libcudf, libnvstrings, and RMM. The main content area is titled "API Reference" and "DataFrame". It defines the DataFrame class as "A GPU Dataframe object." and provides examples for building dataframes with \_\_setitem\_\_. A code block shows a Python session where a DataFrame is created with columns 'key' and 'val', and then printed to show the resulting data:

```
>>> import cudf
>>> df = cudf.DataFrame()
>>> df['key'] = [0, 1, 2, 3, 4]
>>> df['val'] = [float(i + 10) for i in range(5)] # insert column
>>> print(df)
   key  val
0    0  10.0
1    1  11.0
2    2  12.0
3    3  13.0
4    4  14.0
```

# ETL - THE BACKBONE OF DATA SCIENCE

cuDF is not the end of the story



## Scale out RAPIDS with DASK!

The logo consists of a stylized orange flame icon followed by the word "DASK" in a large, bold, white sans-serif font.

Dask natively scales Python

Dask provides advanced parallelism for analytics, enabling performance at scale for the tools you love

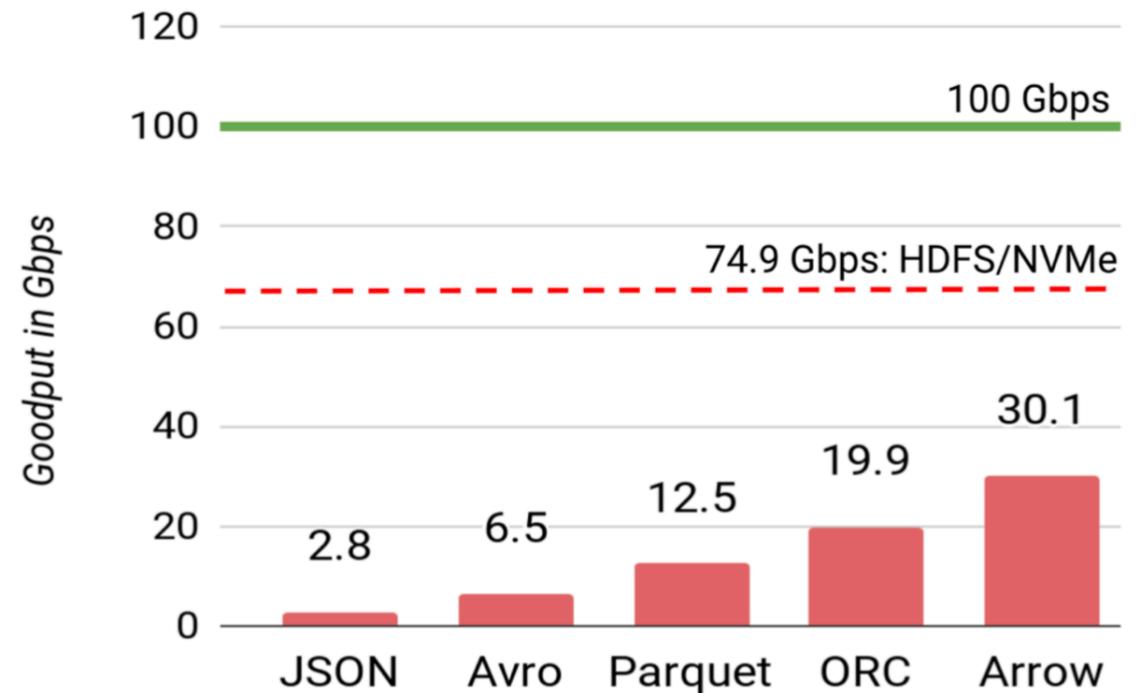
# Scale out RAPIDS with DASK!



- **PyData Native**
  - Built on top of NumPy, Pandas Scikit-Learn, ... (easy to migrate)
  - With the same APIs (easy to train)
  - With the same developer community (well trusted)
- **Scales**
  - Easy to install and use on a laptop
  - Scales out to thousand-node clusters
- **Popular**
  - Most common parallelism framework today at PyData and SciPy conferences
- **Deployable**
  - HPC: SLURM, PBS, LSF, SGE
  - Cloud: Kubernetes
  - Hadoop/Spark: Yarn

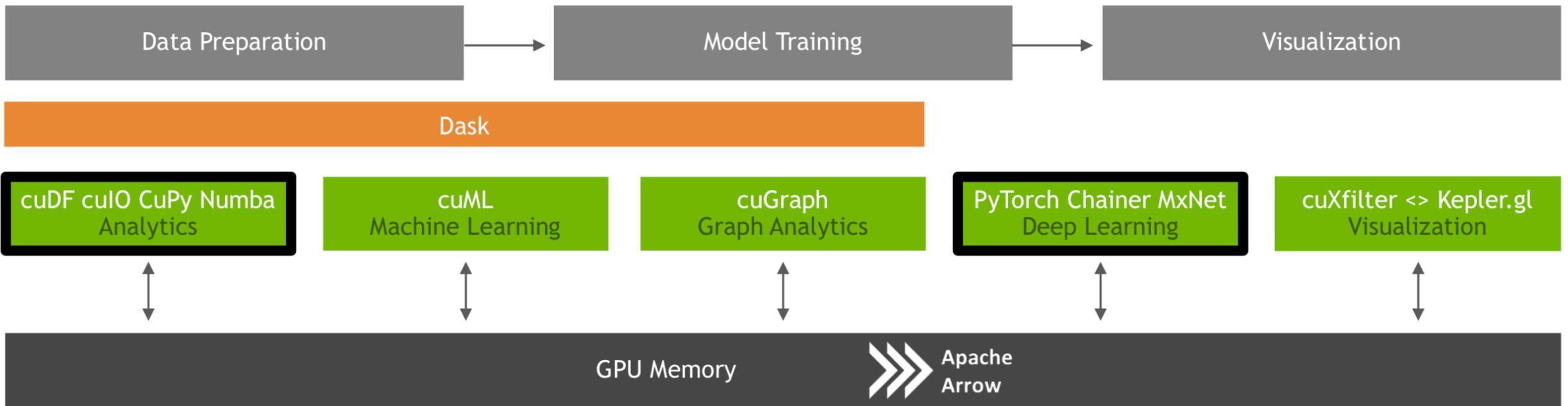
# Data Loading is FAST

- CSV Reader
  - Follows API of `pandas.read_csv`
  - Current implementation is >10x speed improvement over pandas
- Parquet Reader - v0.7
  - Work in progress: Will follow API of `pandas.read_parquet`
- ORC Reader - v0.7
  - Work in progress: Will have similar API of Parquet reader
- Additionally looking towards GPU-accelerating decompression for common compression schemes



Source: Apache Crail blog: [SQL Performance: Part 1 - Input File Formats](#)

# RAPIDS and Deep Learning Frameworks



## Using RAPIDS with PyTorch

Using the GPU for ETL and preprocessing of deep learning workflows



Even Oldridge [Follow](#)  
May 21 · 8 min read

<https://medium.com/rapids-ai/using-rapids-with-pytorch-e602da018285>

# RAPIDS on Power – How can I use it ?

Auto-AI software: PowerAI Vision

Watson  
Machine Learning

Watson ML Accelerator

Watson ML CE

Runtime Environment  
Train, Deploy, Manage Models

XGBoost

Spark

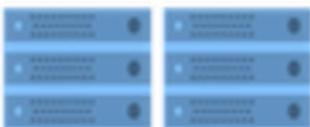
TensorFlow

Keras

PyTorch

Chainer

SnapML



Accelerated AC922  
Power9 Servers



Storage  
(Spectrum Scale ESS)

```
## Install PowerAI ! AKA WML-CE
```

```
conda config --prepend channels \
https://public.dhe.ibm.com/ibmdl/export/pub/software/server/ibm-ai/conda/
```

```
conda create --name wmlce_env python=3.6
conda activate wmlce_env
```

```
conda install powerai
```

```
## Python Program
```

```
import cudf
```

```
import cuml
```

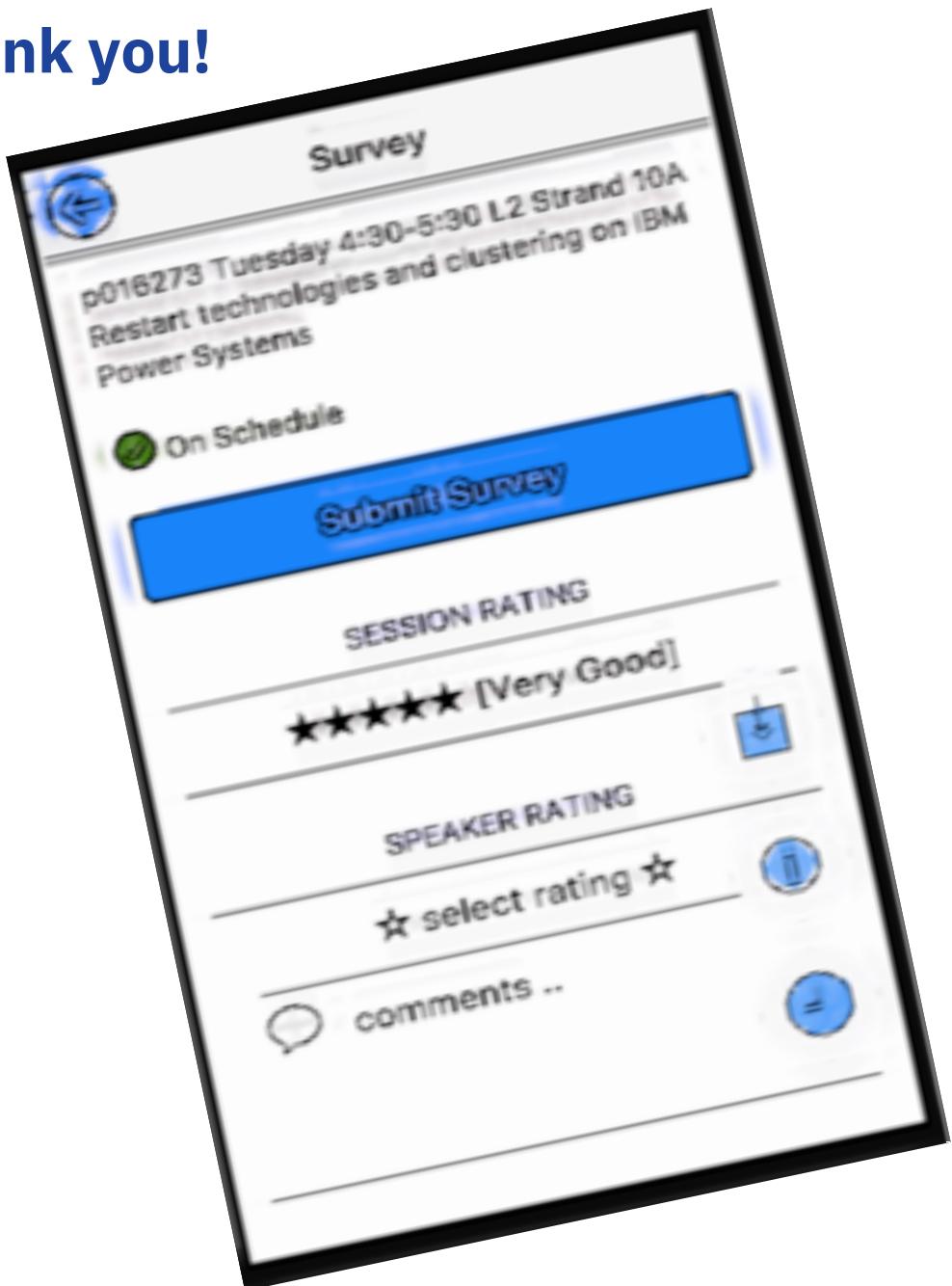
...

...

[https://www.ibm.com/support/knowledgecenter/SS5SF7\\_1.6.1/navigation/wmlce\\_install.htm](https://www.ibm.com/support/knowledgecenter/SS5SF7_1.6.1/navigation/wmlce_install.htm)

## **Jupyter Lab with Rapids**

**Thank you!**



Dustin VanStee  
Data Scientist

[vanstee@us.ibm.com](mailto:vanstee@us.ibm.com)

**Please complete the Session Evaluation!**

# IBM AI IMMERSION EXPERIENCE



## 2 Day Workshop

IMMERSIVE HANDS-ON SESSIONS USING ML/DL ON IBM GPU POWER SERVERS FOR DATA ENGINEERS, DATA SCIENTISTS AND PRACTITIONERS

### Sample Agenda

#### Day 1

- State of AI and General Use Cases
- Machine Learning Overview
- Lab: hands-on with ML using H2O DAI
- Auto AI Demo
- Deep Learning – Applications in Computer Vision
- Lab: Object Detection with PowerAI Vision
- Building an AI Platform (Infrastructure)

Attendees will gain an understanding of the most recent breakthroughs in deep learning and how to apply these techniques using open source deep learning frameworks and IBM AI software.

#### Industry Focused

Financial  
Comms/CSI  
Retail  
Healthcare  
Distribution  
Industrial  
Research  
Utilities

#### Day 2

- Deep Learning with NLP Overview
- Lab: Aspect Based Sentiment using FastAI & PyTorch
- Deep Learning with Computer Vision
- Convolutional Neural Networks
- GAN Overview
- Lab: MNIST Data Generation Using GAN

**Customize the topics & agenda.**

For a full menu of topics go to  
<https://ibm.box.com/v/IBMAIBootcampTopics>

SIGN UP BY CONTACTING  
COGNITIVE SYSTEMS SOLUTIONS CENTER  
[CSSC@US.IBM.COM](mailto:CSSC@US.IBM.COM)

IBM Systems Worldwide  
Client Experience Centers

#IBMSysCEC

