

# Databricks ML Workshop

Samir Gupta, Solutions Architect



# Workshop Agenda

Overview of Databricks for Data Science

The Databricks ML Runtime

Koalas - Pandas interface to Spark Dataframes

R on Databricks

DBconnect - use your favourite IDE on Databricks

Mlflow - train, experiment, track and deploy models

Sample notebooks

Q&A

# Hardest Part of AI isn't AI, it's Data

*“Hidden Technical Debt in Machine Learning Systems,” Google NIPS 2015*

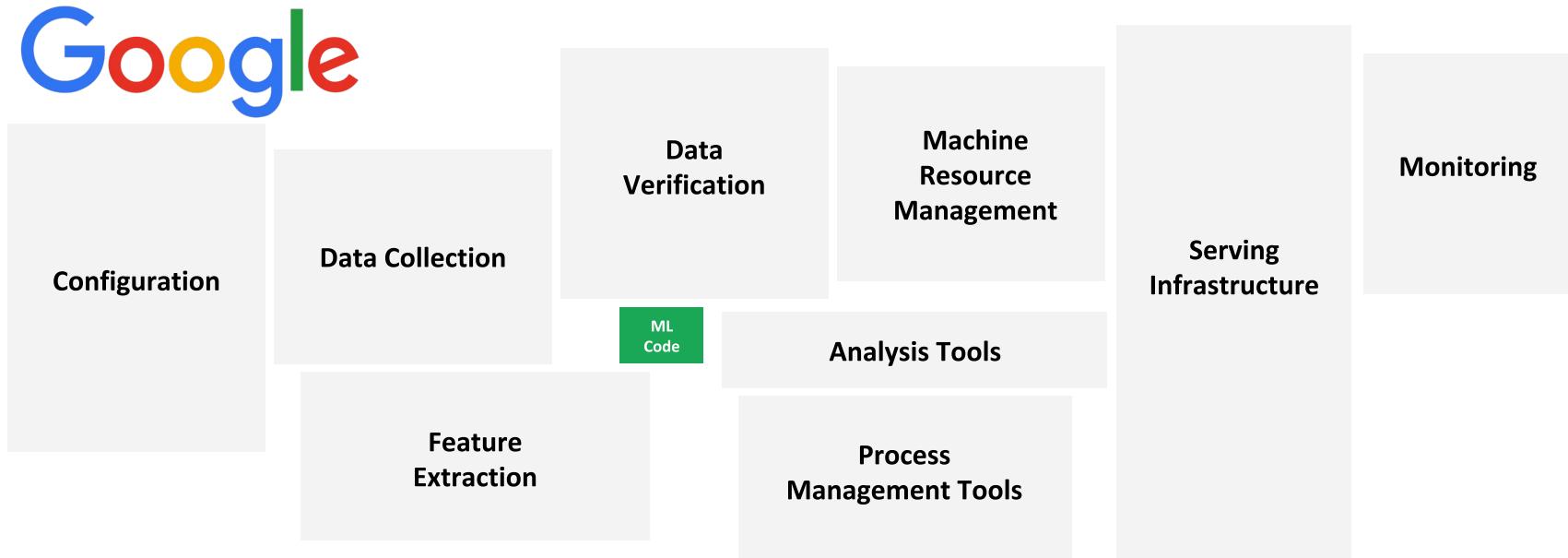
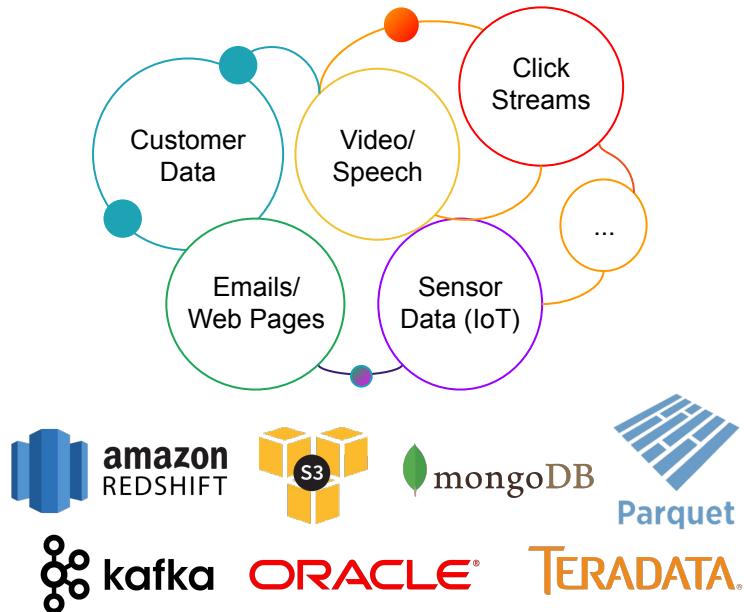


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small green box in the middle. The required surrounding infrastructure is vast and complex.

# And...Data & AI Technologies are in Silos

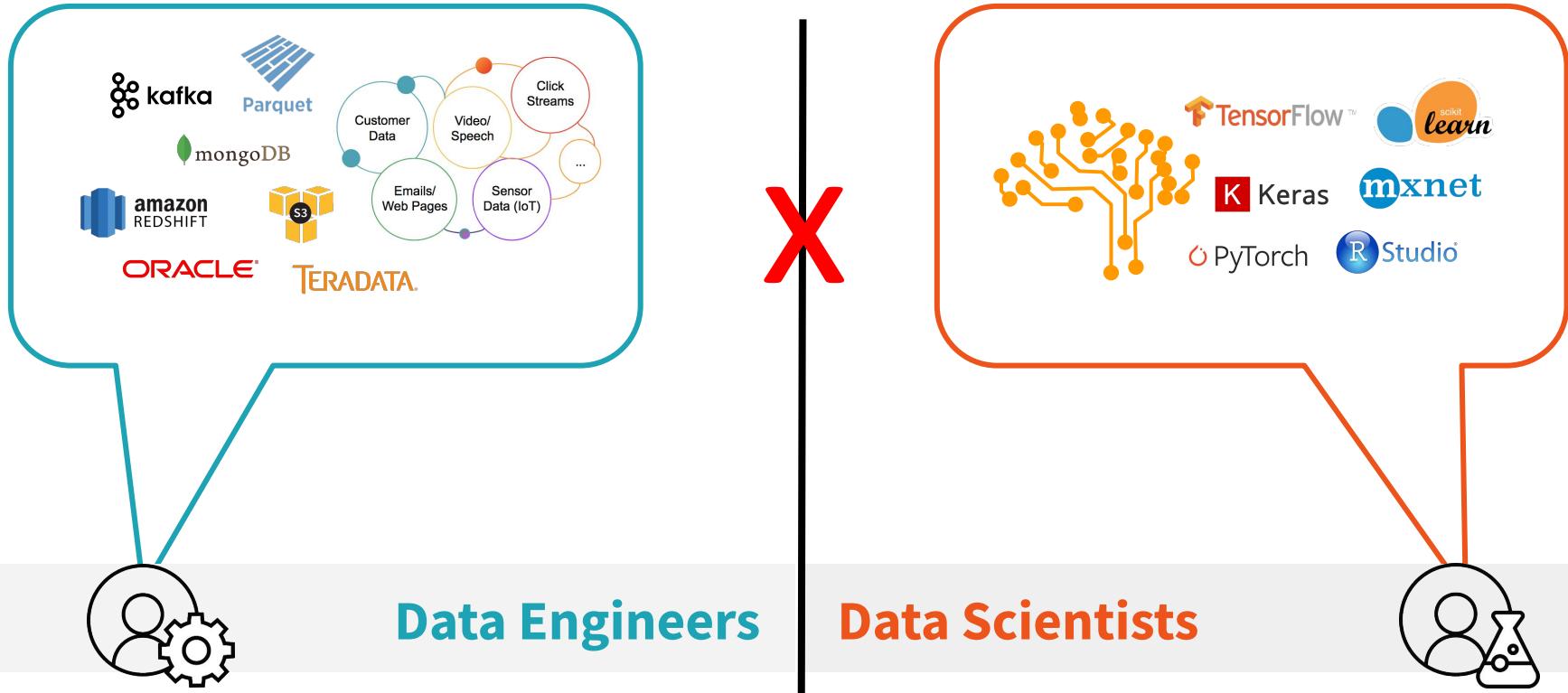


Great for Data, but not AI



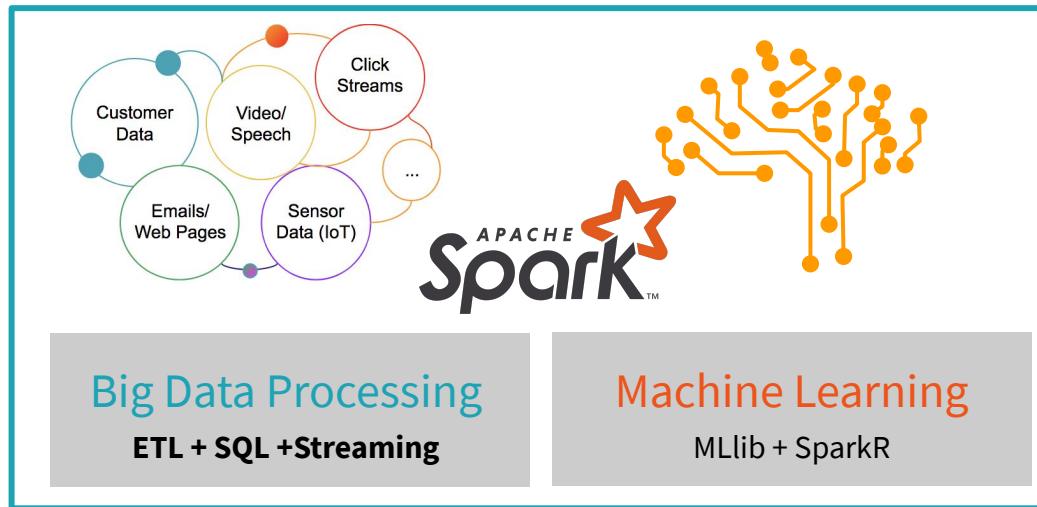
Great for AI, but not for data

# Data & AI People are in Silos

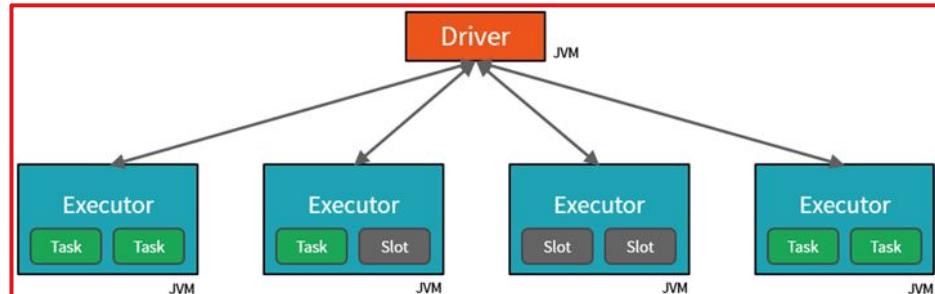


# Apache Spark...

1<sup>st</sup> Unified Analytics Engine that uniquely combines data & AI technologies...



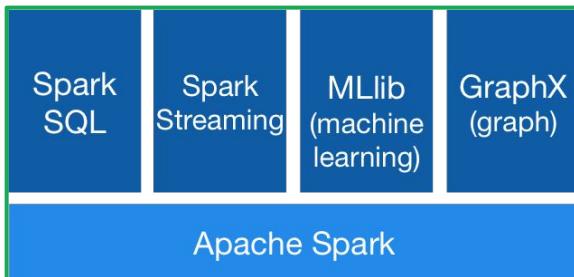
# What is Apache Spark?



**Fast**, distributed processing on **large volumes** of structured and unstructured data

```
df = spark.read.json("logs.json")
df.where("age > 21")
    .select("name.first").show()
```

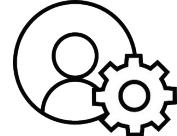
**Rich** and easy to use API programmable in Python, Scala, R, Java, SQL



**Flexible** processing for batch, streaming, and data science/ML workloads

# Apache Spark...

Was only the first step



Data Engineers



Data Scientists

Spark Expertise, Data Pipelines, Infrastructure Mgmt...

Unified Analytics Engine



# Databricks Unified Analytics Platform

## AZURE DATA SOURCES

Blob Storage

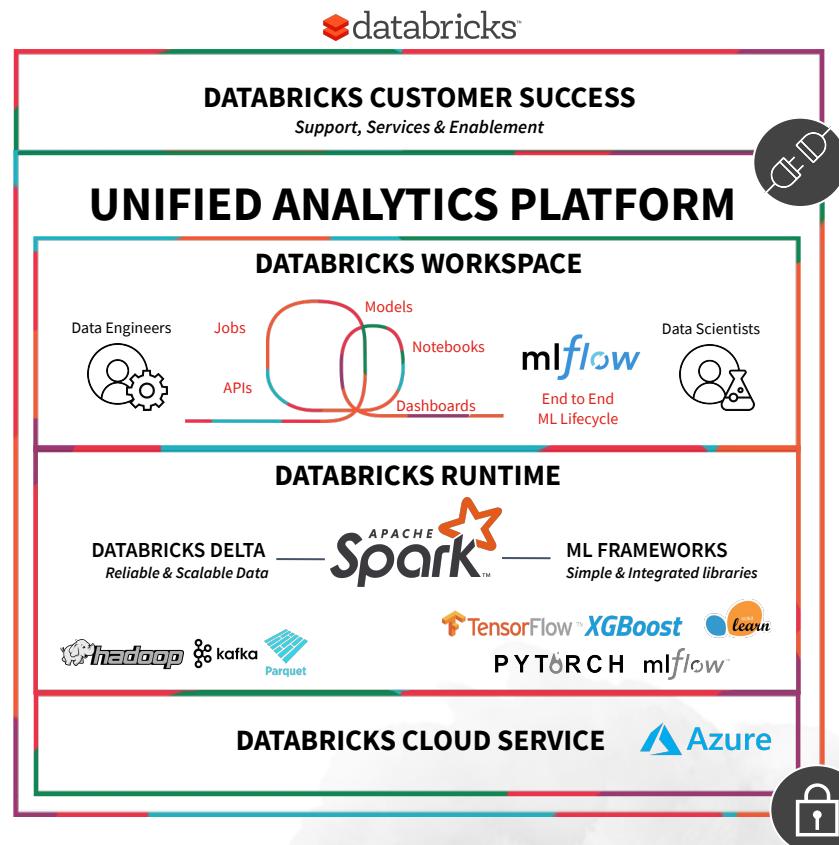
Data Lake Store

SQL Data Warehouse

Cosmos DB

Event Hub

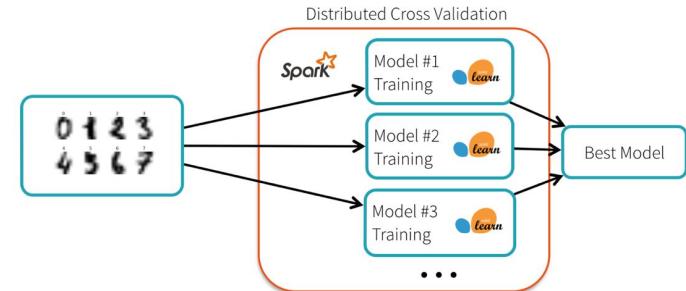
IoT Hub



# Why Distributed Machine Learning?

1. Speed - too slow to train a model on a single machine
2. Resources - data or models are too big for a single machine

**spark-sklearn** - train multiple sklearn models on Spark

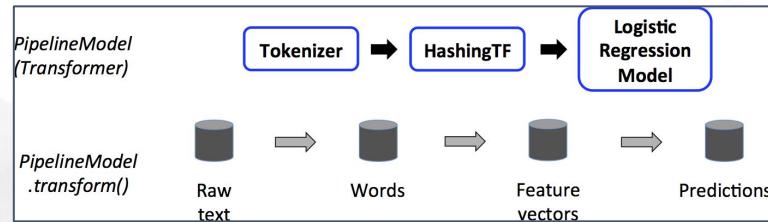


**SparkML** - Native Python and Scala API for ML on Spark

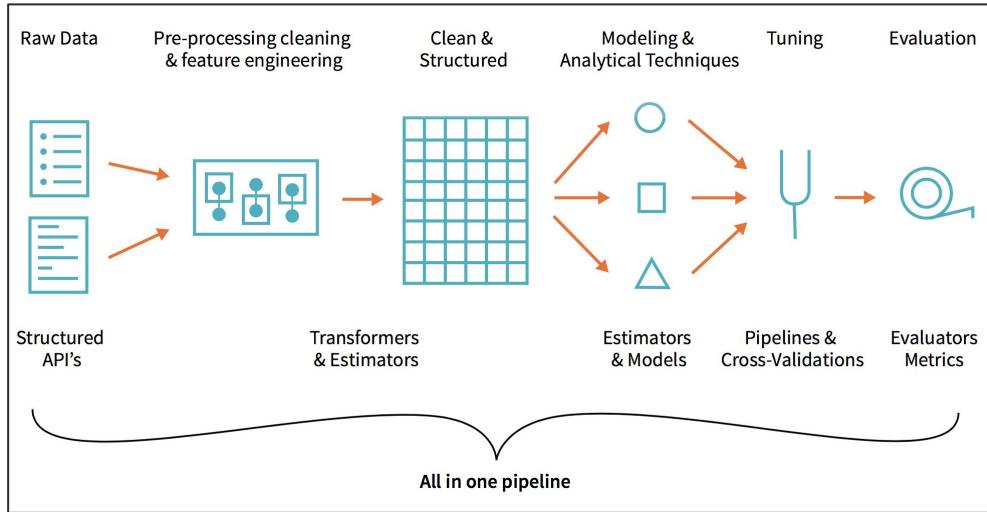
Featuring engineering, regression, classification, clustering, and more written for Spark

Uses transformers (.transform()) and estimators (.fit())

Combine many transformers and estimators into a Pipeline



# Productionizing Spark Models



Latency Requirement	What to Use
10-100ms	Low latency Sparkless Highly Available Prediction Server
100ms-1m	Spark Structured Streaming
1m-1hr+	Spark Batch Jobs

Later: MLflow + AzureML - End-to-end model lifecycle

# Koalas: Pandas on Spark

Pandas is the *de-facto* Python standard (single-node) for DataFrame data processing

Spark is the de facto standard for big data processing

Koalas let's data scientists:

- Be immediately productive in Spark - with no learning curve if already familiar with pandas
- Have a single codebase that works both with pandas (tests, smaller datasets) and with Spark (large, distributed datasets).

```
from databricks import koalas

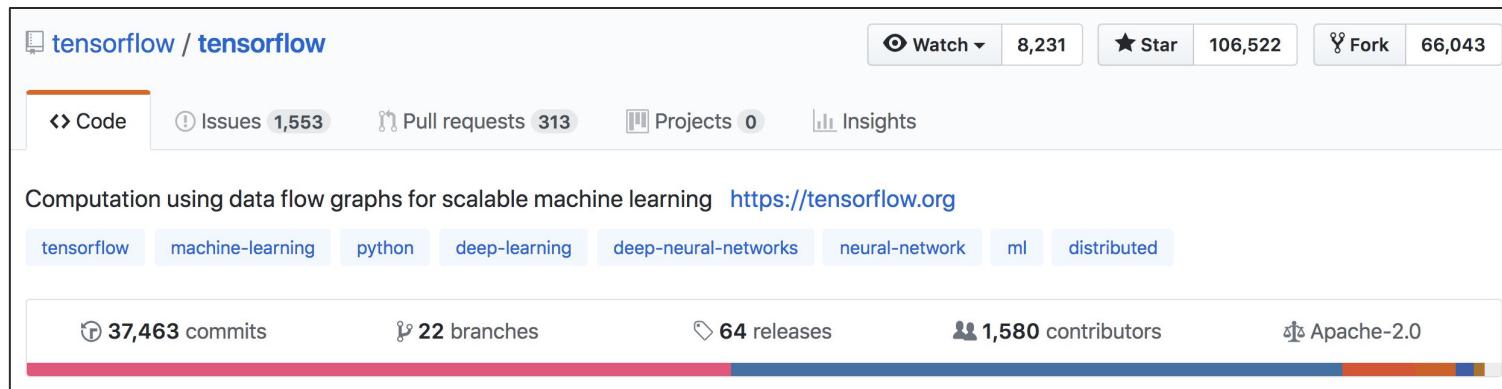
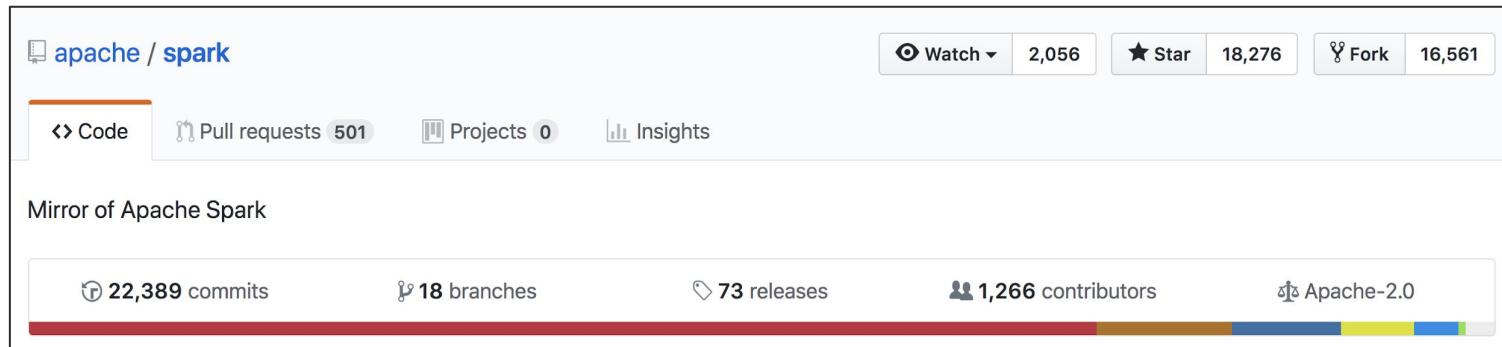
data = koalas.read_csv("/databricks-datasets/sai-summit-2019-sf/fire-calls.csv", header=0)
data = data[["CallType", "Neighborhood", "NumAlarms", "OrigPriority", "UnitType", "Delay"]]

# Rename the columns
data.columns = ['a', 'b', 'c', 'd', 'e', 'f']

# Do some operations in place:
data['g'] = data.c * data.f

display(data)
```

# We love Spark, but what about Deep Learning?



# Databricks ML Runtime

Ready to use clusters with built-in ML Frameworks

including TensorFlow, Keras, Horovod, and [more](#)



XGBoost



P Y T O R C H

Horovod Estimator

for **simplified distributed training** on TensorFlow with Horovod using Apache Spark on Databricks

GPU support

Azure GPU (NC) instances pre-loaded with Tesla/CUDA/cuDNN libraries!



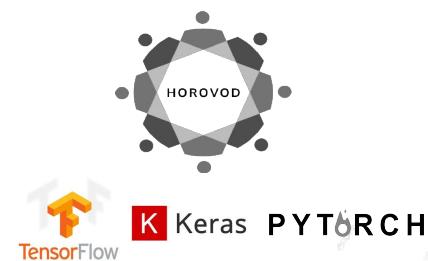
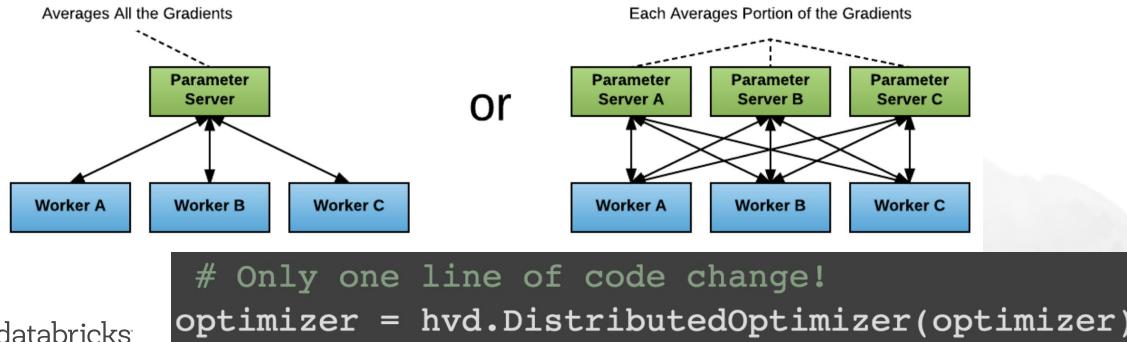
# Options for Deep Learning on Databricks

Option 1: Train on a single node

Option 2: Train on each node of the cluster

- Ex. Distributed hyperparameter tuning & search

Option 3: Distributed Training (ie. HorovodEstimator, Dist-Keras)



# R on Databricks

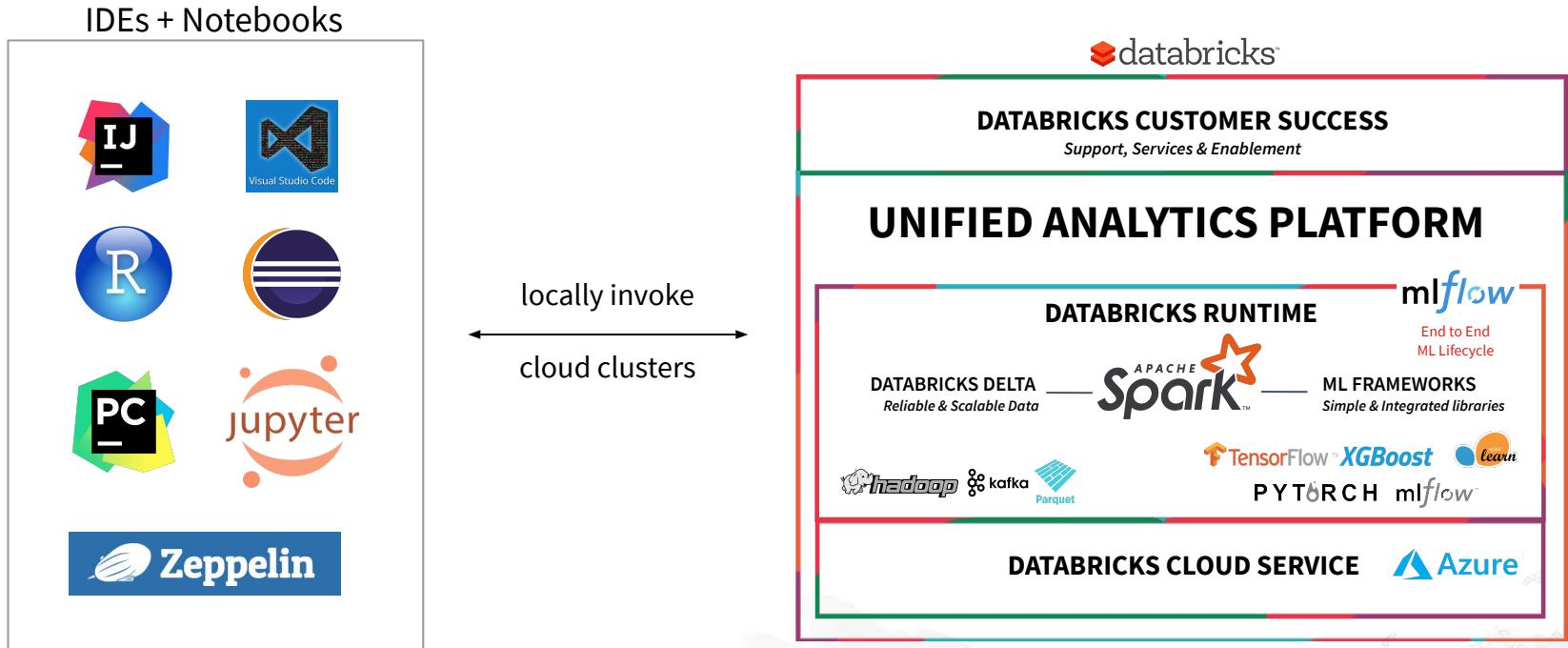
## Databricks RStudio Integration

### SparkR - recommended

- Created by Spark committers
- Optimized performance on Databricks
- Up to 100x faster!
- dapply() to apply in parallel!

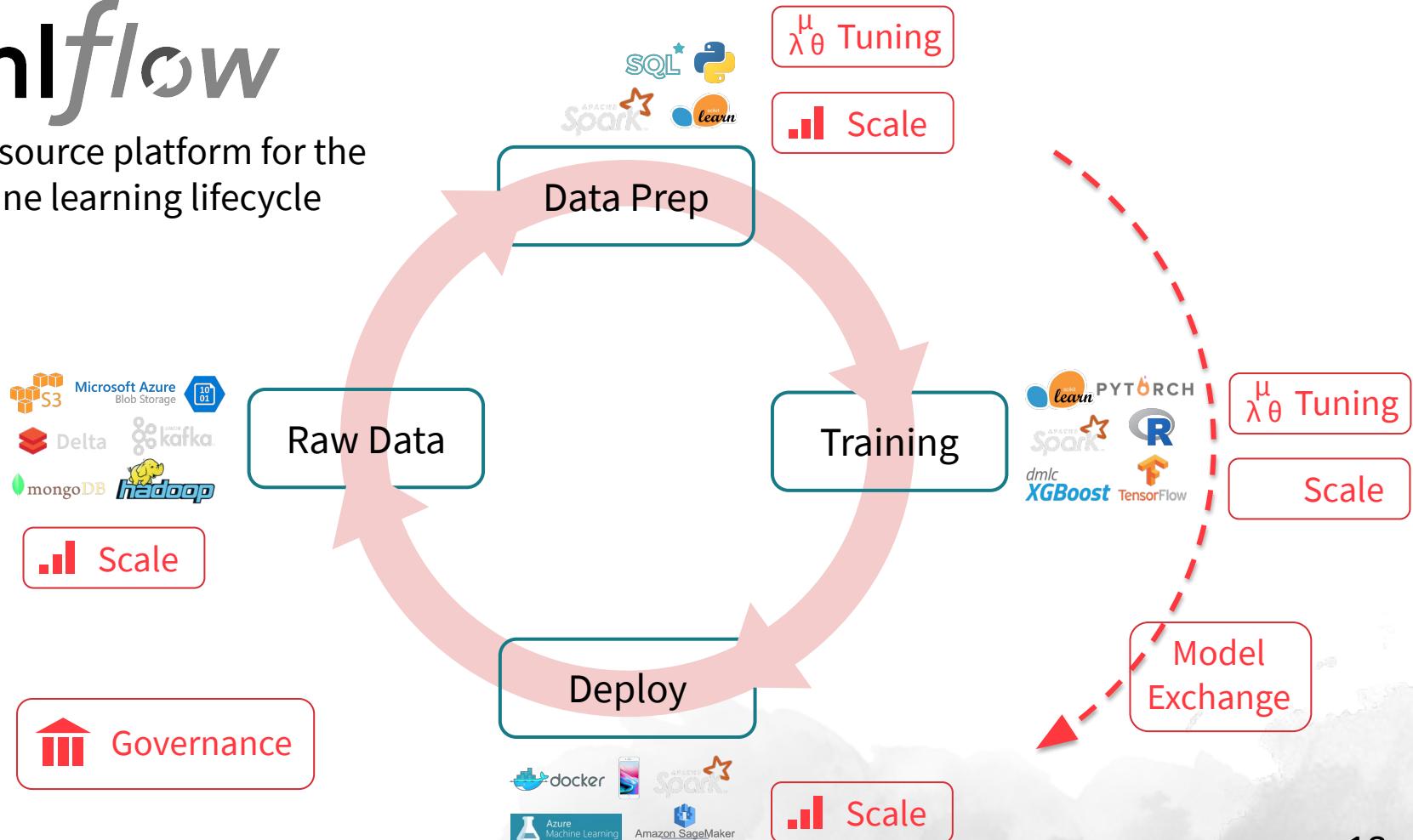
```
# Convert waiting time from hours to seconds.
# Note that we can apply UDF to DataFrame.
schema <- structType(structField("eruptions", "double"), structField("waiting", "double"),
                      structField("waiting_secs", "double"))
df1 <- dapply(df, function(x) { x <- cbind(x, x$waiting * 60) }, schema)
head(collect(df1))
## eruptions waiting waiting_secs
##1 3.600    79     4740
##2 1.800    54     3240
##3 3.333    74     4440
##4 2.283    62     3720
##5 4.533    85     5100
##6 2.883    55     3300
```

# DBconnect: Databricks for Developers



# mlflow

An open source platform for the machine learning lifecycle



# MLflow Components

## *mlflow* Tracking

Record and query experiments: code, data, config, results

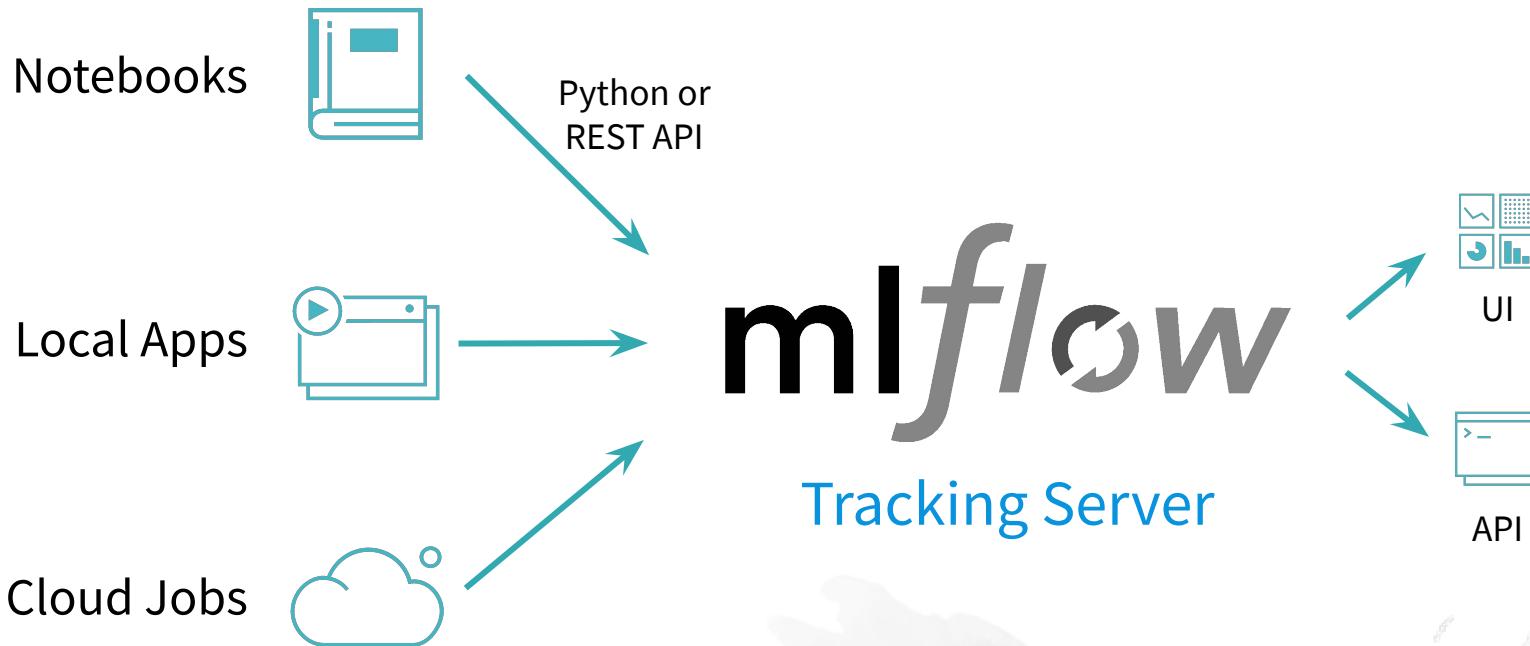
## *mlflow* Projects

Packaging format for reproducible runs on any platform

## *mlflow* Models

General model format that supports diverse deployment tools

# MLflow Tracking



# Key Concepts in Tracking

**Parameters:** key-value inputs to your code

**Metrics:** numeric values (can update over time)

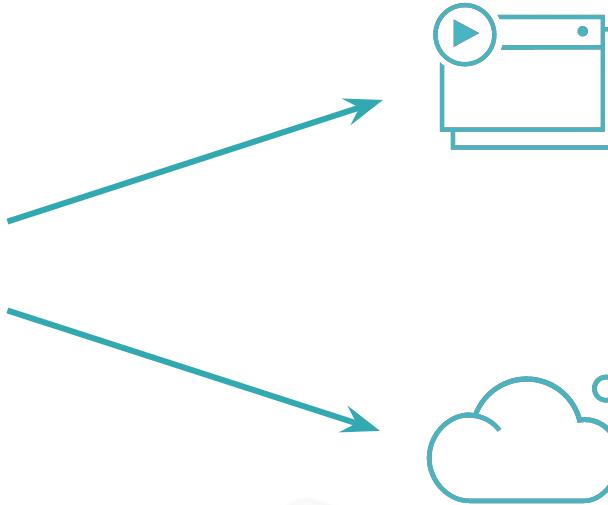
**Artifacts:** arbitrary files, including models

**Source:** what code ran?

The screenshot shows the mlflow UI interface. At the top, there's a navigation bar with the mlflow logo, GitHub, and Docs links. Below it, a sidebar titled 'Experiments' lists 'Default' and 'Something'. The main area is titled 'Default' and shows 'Experiment ID: 0' and 'Artifact Location: /Users/matei/mlflow/miruns/0'. It includes search and filter fields for 'Search Runs' (metrics.rmse < 1 and params.model = "tree") and 'Filter Params' (alpha, lr) and 'Filter Metrics' (rmse, r2). A button for 'Clear' is also present. Below this, a message says '4 matching runs' with buttons for 'Compare Selected' and 'Download CSV'. A table follows, with columns: Date, User, Source, Version, Parameters, and Metrics. The data in the table is as follows:

Date	User	Source	Version	Parameters	Metrics
2018-06-28 17:09:49	matei	matei_test.py	7cff8e	(n/a)	loss 2.123
2018-06-28 17:09:06	matei	matei_test.py	7cff8e	(n/a)	loss 4.543
2018-06-28 17:09:05	matei	matei_test.py	7cff8e	(n/a)	loss 4.543
2018-06-25 13:08:12	matei	matei_test.py	53ccdc	(n/a)	loss 4.543

# MLflow Projects



Local Execution

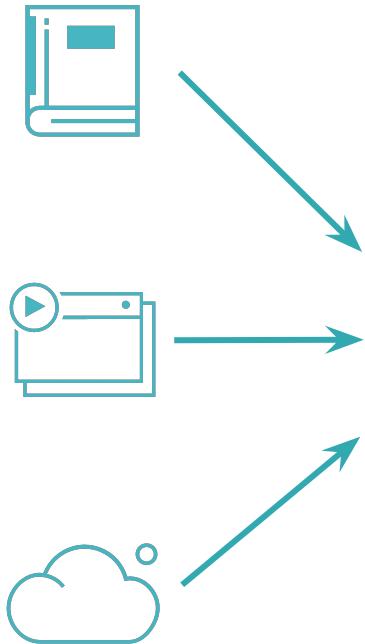
Remote Execution

# Example MLflow Project

```
my_project/
  └── MLproject
      conda_env: conda.yaml
      entry_points:
        main:
          parameters:
            training_data: path
            lambda: {type: float, default: 0.1}
          command: python main.py {training_data} {lambda}
  └── conda.yaml
  └── main.py
  └── model.py
  ...
  
```

```
$ mlflow run git://<my_project>
mlflow.run("git://<my_project>", ...)
```

# MLflow Models



**mlflow**

Model Format

Flavor 1



Flavor 2



Simple model flavors  
usable by many tools



Inference Code



Batch & Stream Scoring



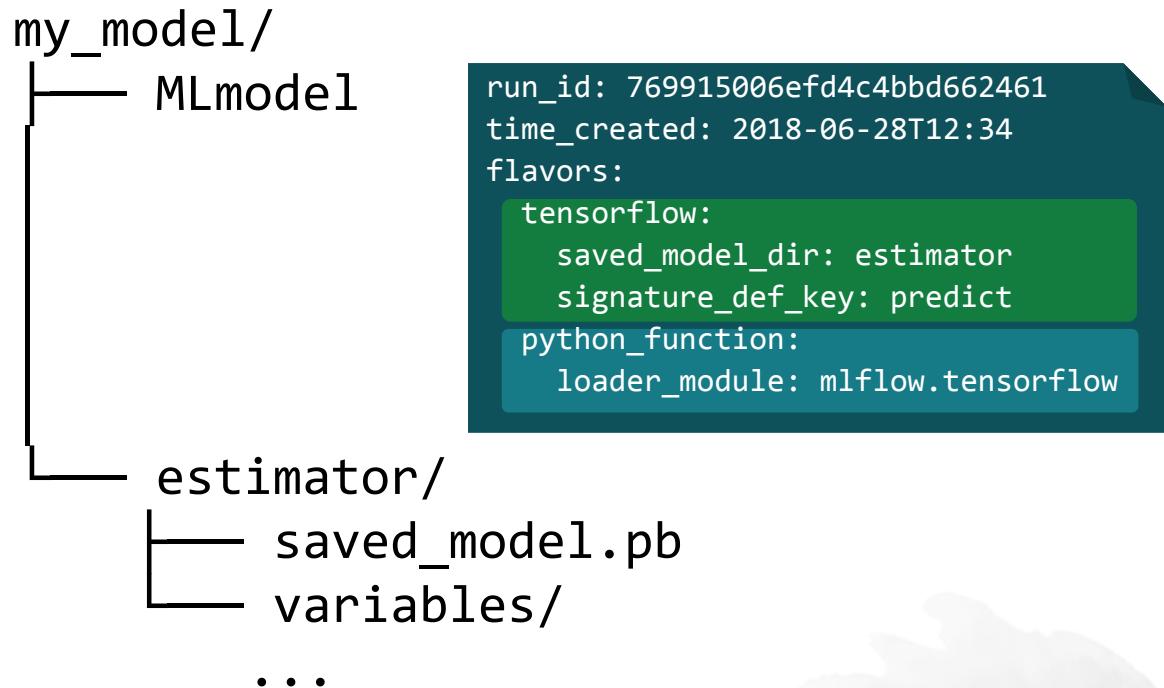
Azure  
Machine Learning



kubernetes

Cloud Serving Tools

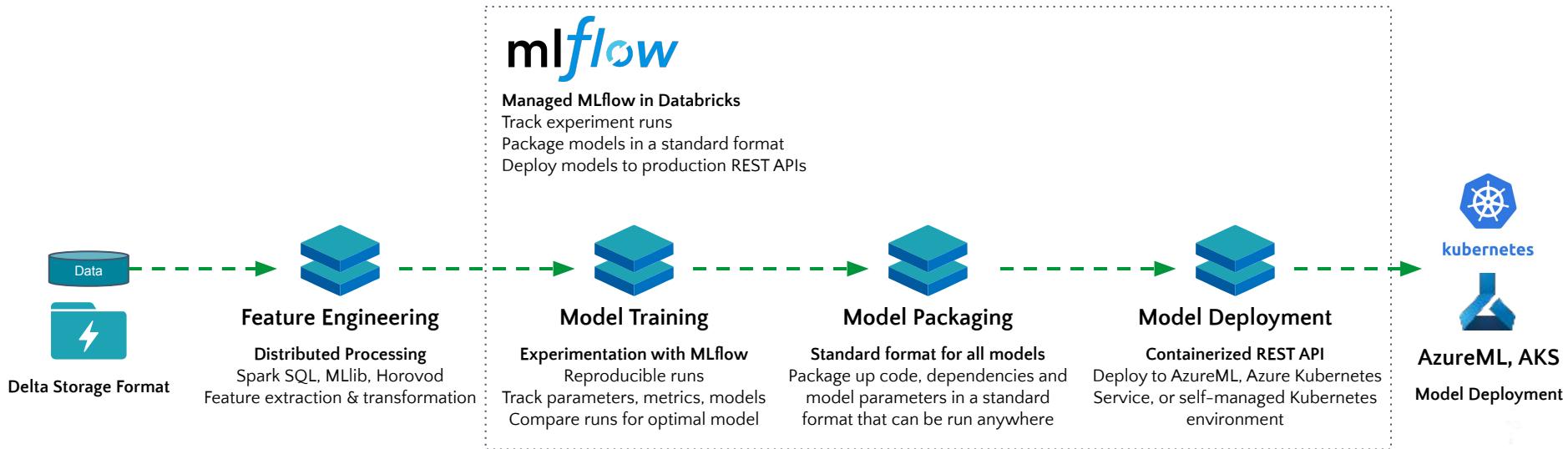
# Example MLflow Model



] Usable by tools that understand  
TensorFlow model format

] Usable by any tool that can run  
Python (Docker, Spark, etc!)

# Machine Learning



# Resources

## Documentation:

[Spark MLlib Guide](#) - docs for feature extraction, ML algorithms, and pipelines in Spark

[Databricks ML Guide](#) - best practices for using MLlib on Databricks

[MLflow Guide](#) - tracking experiments, saving and deploying models on Databricks

[Deep Learning on Databricks](#) - guide to using distributed DL frameworks

[Koalas](#) - Pandas interface to Spark

## Notebooks:

[ML Workshop](#) - GitHub repository with examples on data analysis, ML training, and experiments

More example notebooks - email [samir@databricks.com](mailto:samir@databricks.com)

# Demo and Sample Notebooks