



Clemens Mewald,  
Director of Product Management,  
Databricks

Joel Thomas,  
Senior Solution Architect,  
Databricks

## WEBINAR

Productionizing Machine  
Learning: From Deployment  
to Drift Detection



# Logistics

- We can't hear you...
- Recording will be available...
- Slides will be available...
- Code samples and notebooks will be available...
- Submit your questions...
- Bookmark **[databricks.com/blog](https://databricks.com/blog)**



## VISION

Accelerate innovation by unifying data science, engineering and business

---

## SOLUTION

Unified Analytics Platform

---

## WHO WE ARE

- Original creators of  **APACHE Spark**™,  **DELTA LAKE**, and  **mlflow**™
- 2,000+ global companies use our platform across big data & machine learning lifecycle



## Databricks Workspace

*Collaborative Notebooks, Production Jobs*

## Databricks Runtime



## Databricks Cloud Service



# About Today's Presenters



## **Clemens Mewald, Director of Product Management at Databricks**

- MSc in Computer Science from UAS Wiener Neustadt, Austria
- MBA from MIT Sloan
- Previously Product Lead on the Google Brain Team for TensorFlow, TFX
- Product Management Lead for Data Science and ML at Databricks



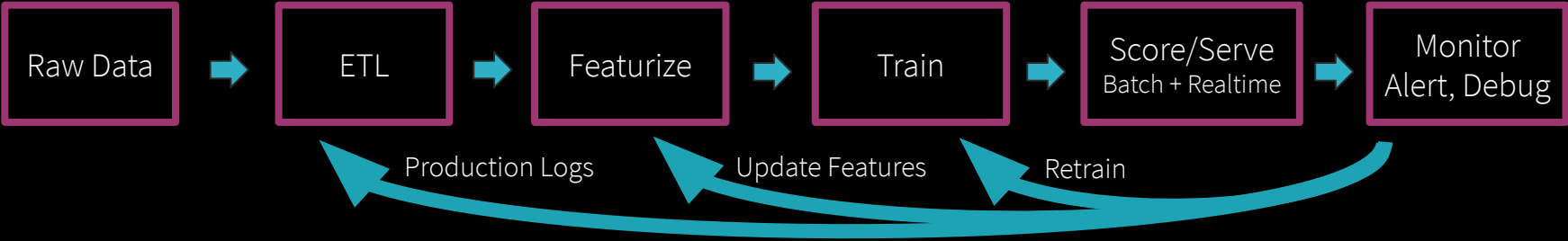
## **Joel Thomas, Senior Solutions Architect at Databricks**

- MS in Electrical Engineering from University of Houston
- Previously Principal Data Scientist at Micron
- Architect and implement data & AI solutions for customers

# Outline

- / ML Development Challenges
- / How MLflow, Delta Lake, + Databricks tackles these
- / Real Life End to End Data + ML Lifecycle
- / Demo of Model Drift Detection on Databricks
- / Q&As

# ML Lifecycle and Challenges



**Zoo of Ecosystem Frameworks**

A collection of logos for various ML ecosystem frameworks, including kafka, S3, Delta, Microsoft Azure Blob Storage, Hadoop, mongoDB, Apache Spark, pandas, SQL, Python, PYTORCH, learn, TensorFlow, Azure Machine Learning, Amazon SageMaker, docker, and Apache Spark.

Tuning

Deploy

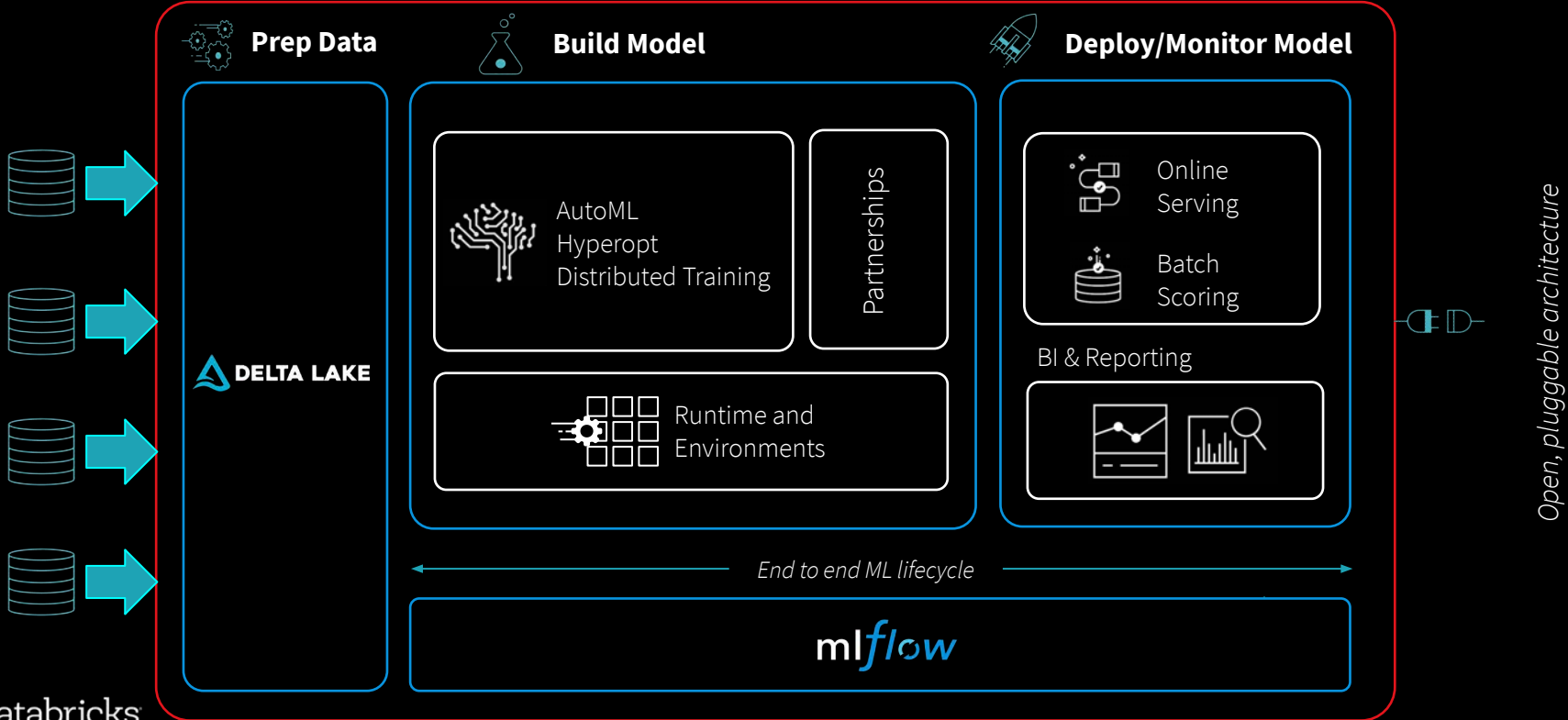
Model Mgmt

Collaboration                      Scale                      Governance

- Feature Repository
- Experiment Tracking
- AutoML, Hyper-p. search
- Remote Cloud Execution
- Project Mgmt (scale teams)
- Model Exchange
- A/B Testing
- CI/CD/Jenkins push to prod
- Orchestration (Airflow, Jobs)
- Lifecycle mgmt.
- Data Drift
- Model Drift

# Unifying the End-to-end Data & ML Lifecycle

## Overview

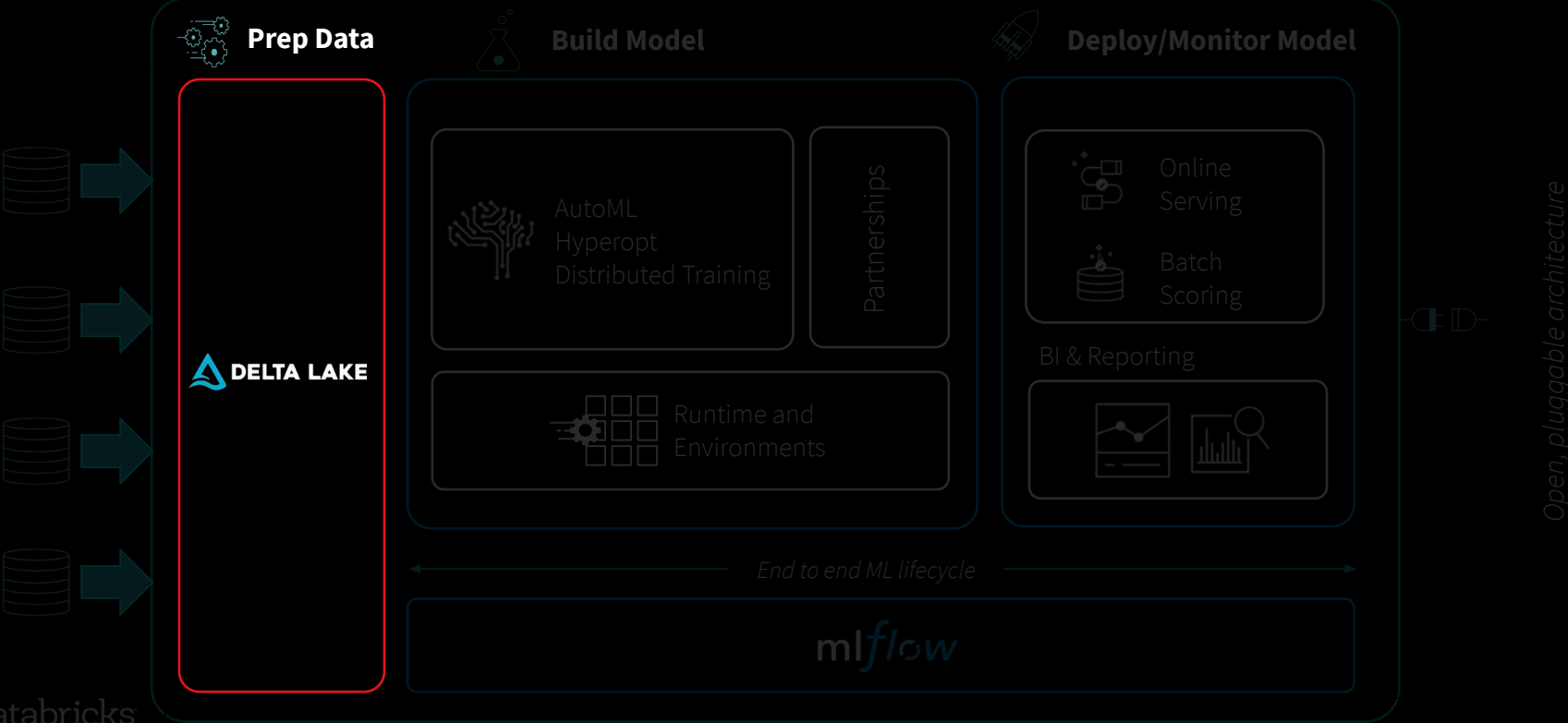


Open, pluggable architecture

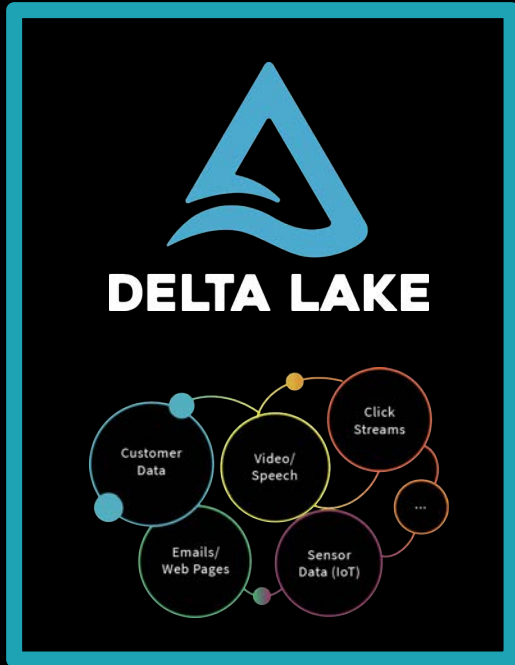


# Unifying the End-to-end Data & ML Lifecycle

## Overview



# Reliability & Performance With Open Standards



## Reliability - High Quality Data

- Transactions and schema enforcement
- Data quality validation and expectation

## Performance - Fast Queries at Scale

- Optimizations - data skipping, caching & indexing
- Compaction optimizes file layout

## Open Standard - Easier Adoption

- Open source, open format; compatible with Spark API's
- Data stays in your data lake under your control



[databricks.com/delta](https://databricks.com/delta)



[delta.io](https://delta.io)

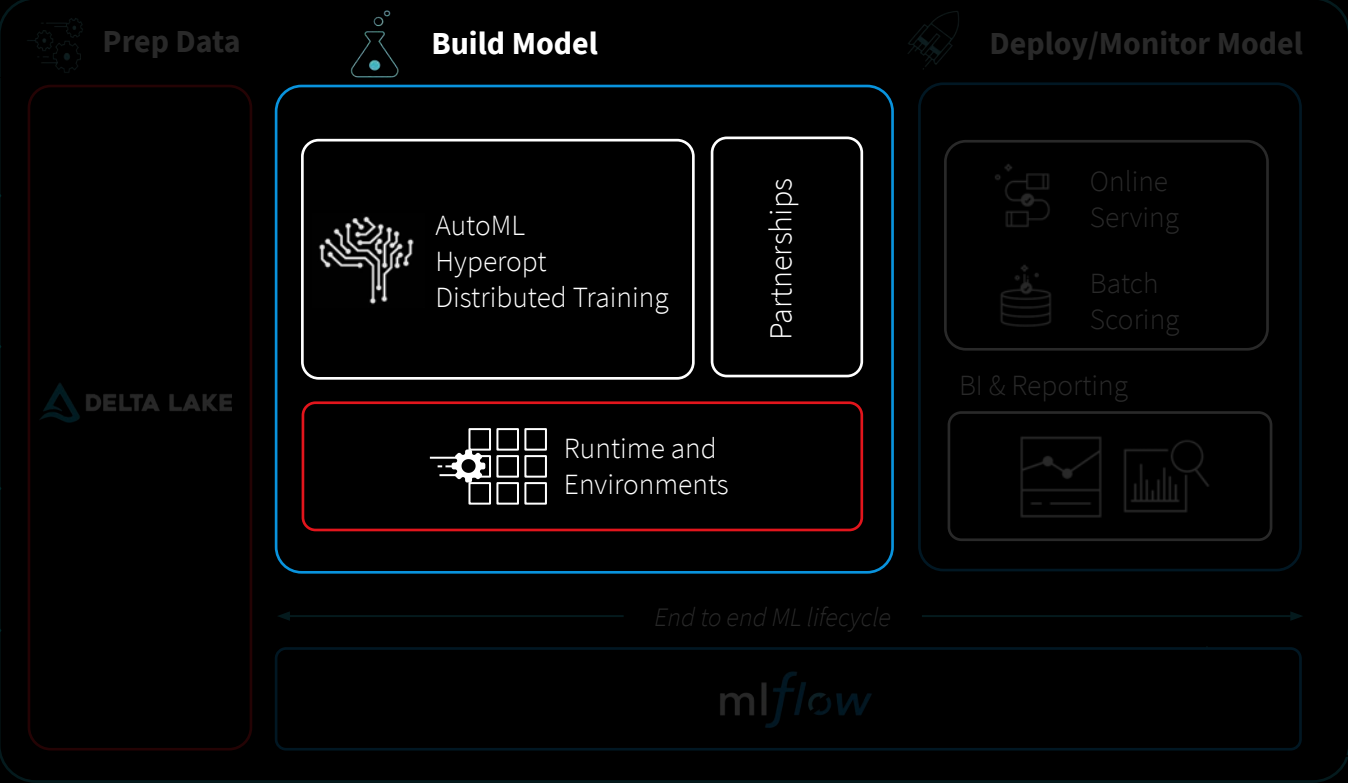


[github.com/delta-io](https://github.com/delta-io)



# Unifying the End-to-end Data & ML Lifecycle

## Overview



Open, pluggable architecture

# ML environment and Model Training

 databricks | Runtime for ML

Packages and optimizes most common ML Frameworks



Built-in Optimization for Distributed Deep Learning



Distribute and Scale any Single-Machine ML Code to 1,000's of machines.

Built-In AutoML and Experiment Tracking



**mlflow**<sup>™</sup>

AutoML and Tracking /  
Visualizations with MLflow

Customized Environments using Conda



requirements.txt  
conda.yaml

Customization



Pre-configured  
Environment



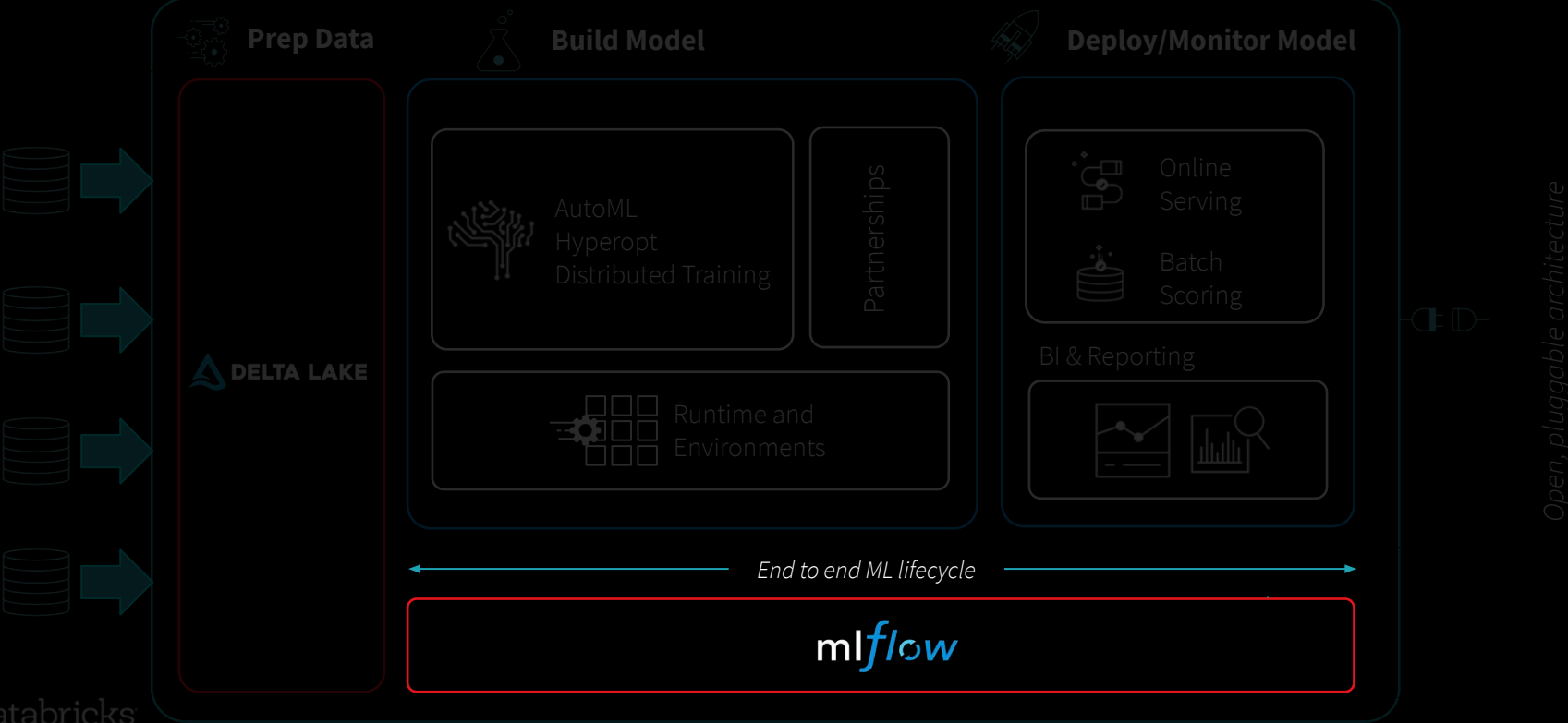
Machine  
Learning



Conda-  
Based

# Unifying the End-to-end Data & ML Lifecycle

## Overview



# An open source framework for the complete ML Lifecycle



Simplified experiment tracking, reproducibility, and deployment,  
from experimentation to production

## 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



[databricks.com/mlflow](https://databricks.com/mlflow)



[mlflow.org](https://mlflow.org)



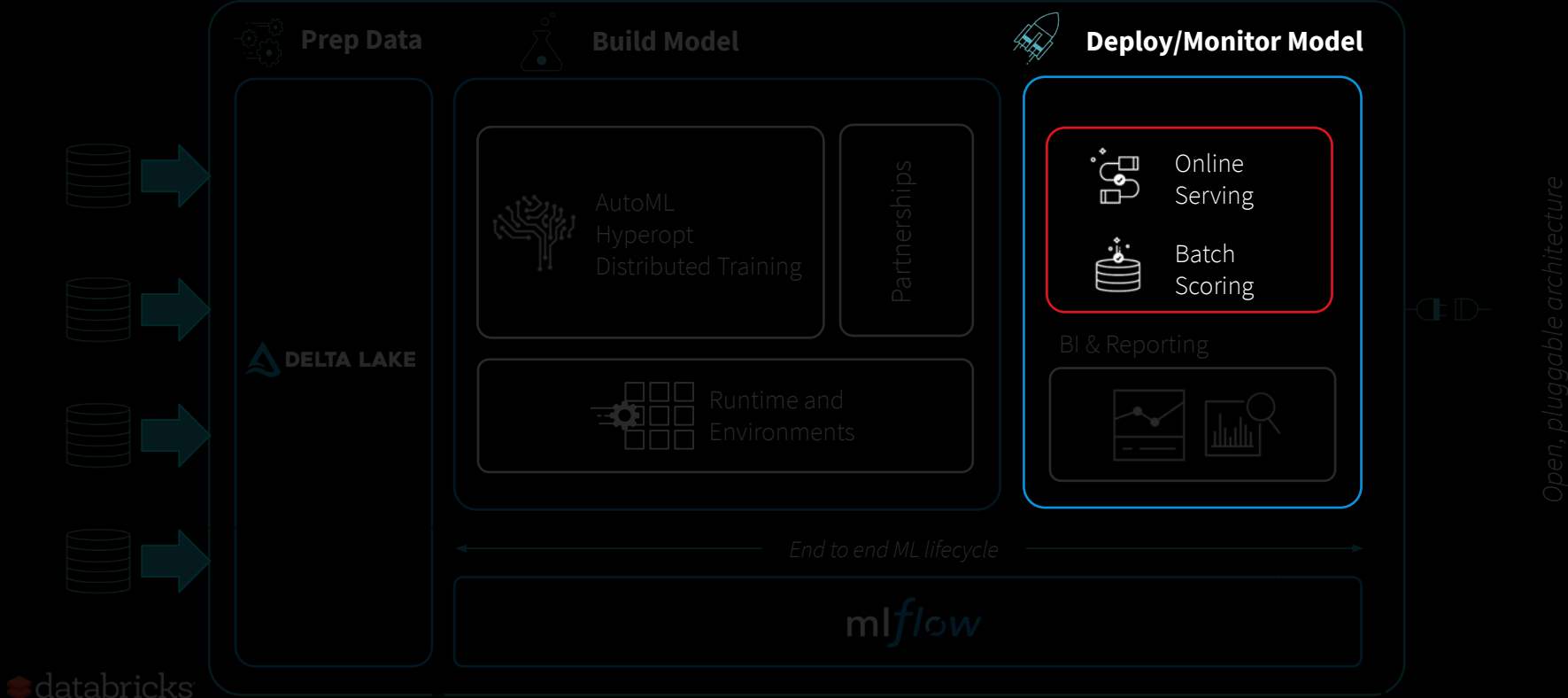
[github.com/mlflow](https://github.com/mlflow)



[twitter.com/MLflow](https://twitter.com/MLflow)

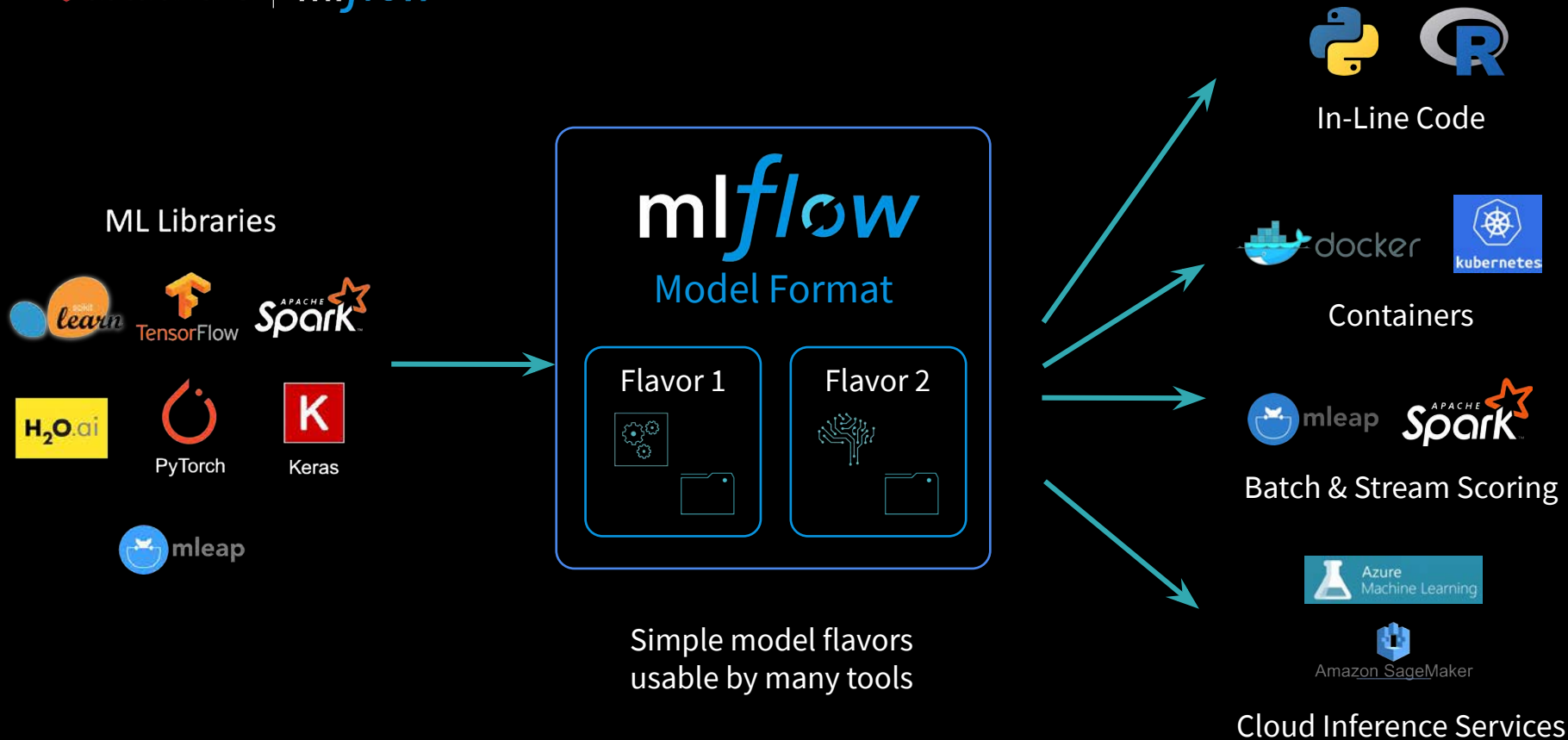
# Unifying the End-to-end Data & ML Lifecycle

## Overview



# Model Deployment and Monitoring

 | 





# What happens after Deployment?



Heraclitus (paraphrased):

“Change is the only constant in life”

# Types of change to worry about in Machine Learning

Concept Drift



- Statistical properties of target variable change (i.e. what you are trying to predict changes)
- E.g. fraud detection

# Types of change to worry about in Machine Learning

Concept Drift



- Statistical properties of target variable change (i.e. what you are trying to predict changes)
- E.g. fraud detection

Data Drift



- Statistical properties of input variables change
- E.g. seasonality, personal preferences, trends change

# Types of change to worry about in Machine Learning

## Concept Drift



- Statistical properties of target variable change (i.e. what you are trying to predict changes)
- E.g. fraud detection

## Data Drift



- Statistical properties of input variables change
- E.g. seasonality, personal preferences, trends change

## Upstream Data Changes



- Encoding of a feature changes (e.g. switch from Fahrenheit to Celsius)
- Features are no longer being generated (leads to missing values)

# Ways to detect and protect against changes

## Monitor

### Training Data

- ▶ Schema & distribution of incoming data
- ▶ Distribution of labels

### Requests & Predictions

- ▶ Schema & distribution of requests
- ▶ Distribution of predictions
- ▶ Quality of predictions

# Ways to detect and protect against changes

## Monitor

### Training Data

- ▶ Schema & distribution of incoming data
- ▶ Distribution of labels

### Requests & Predictions

- ▶ Schema & distribution of requests
- ▶ Distribution of predictions
- ▶ Quality of predictions

## Intervene

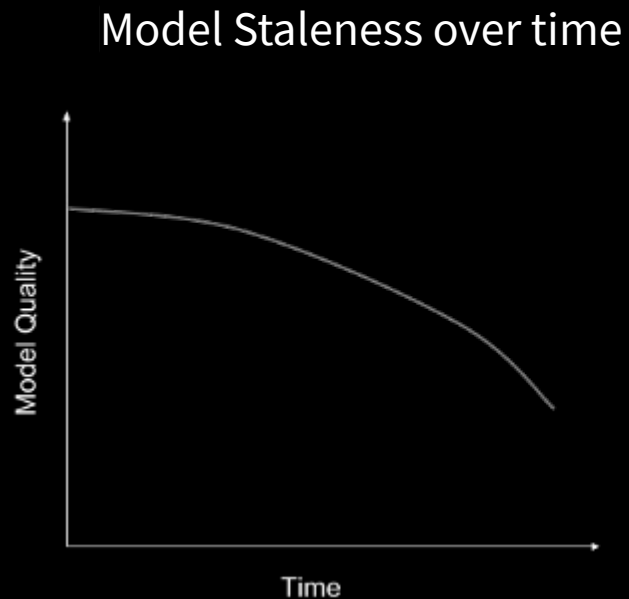
### Correct Data Problems

- ▶ Detect and correct upstream data problems early
- ▶ Clean up erroneous labels

### Retrain / Calibrate Model

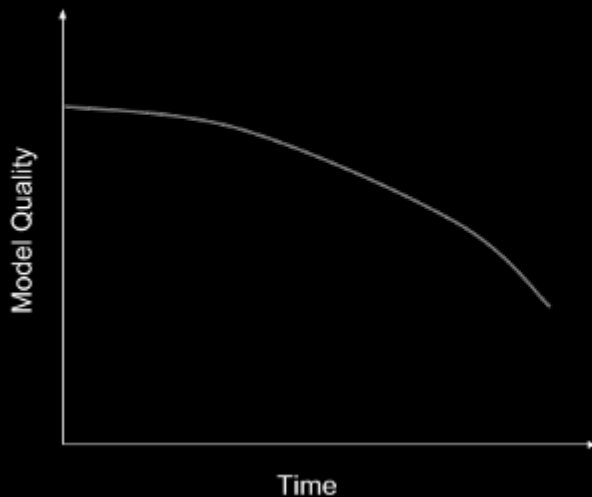
- ▶ In cases where there is concept or data drift, keep models fresh by retraining them

# Your Machine Learning Models become stale

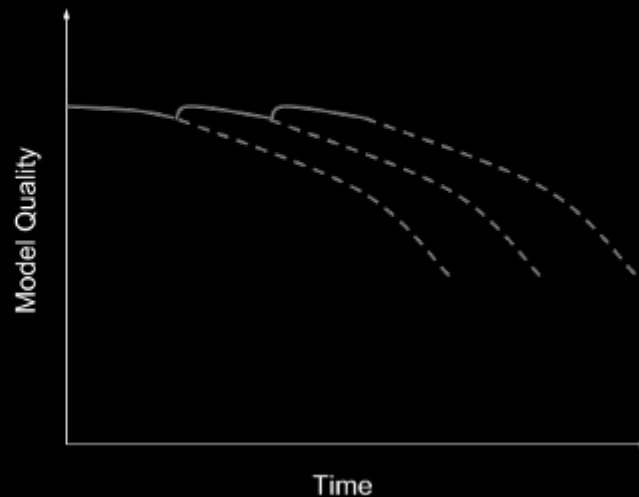


# Your Machine Learning Models become stale

Model Staleness over time



Refreshing models over time





# Real-Life Examples of Drift



IoT

- Sensor readings can change over time (that's why sensors need to be recalibrated)
- Sensor readings may change with seasonality



Retail

- Inventory changes over time
- Customer preferences change over time
- Regional changes and trends

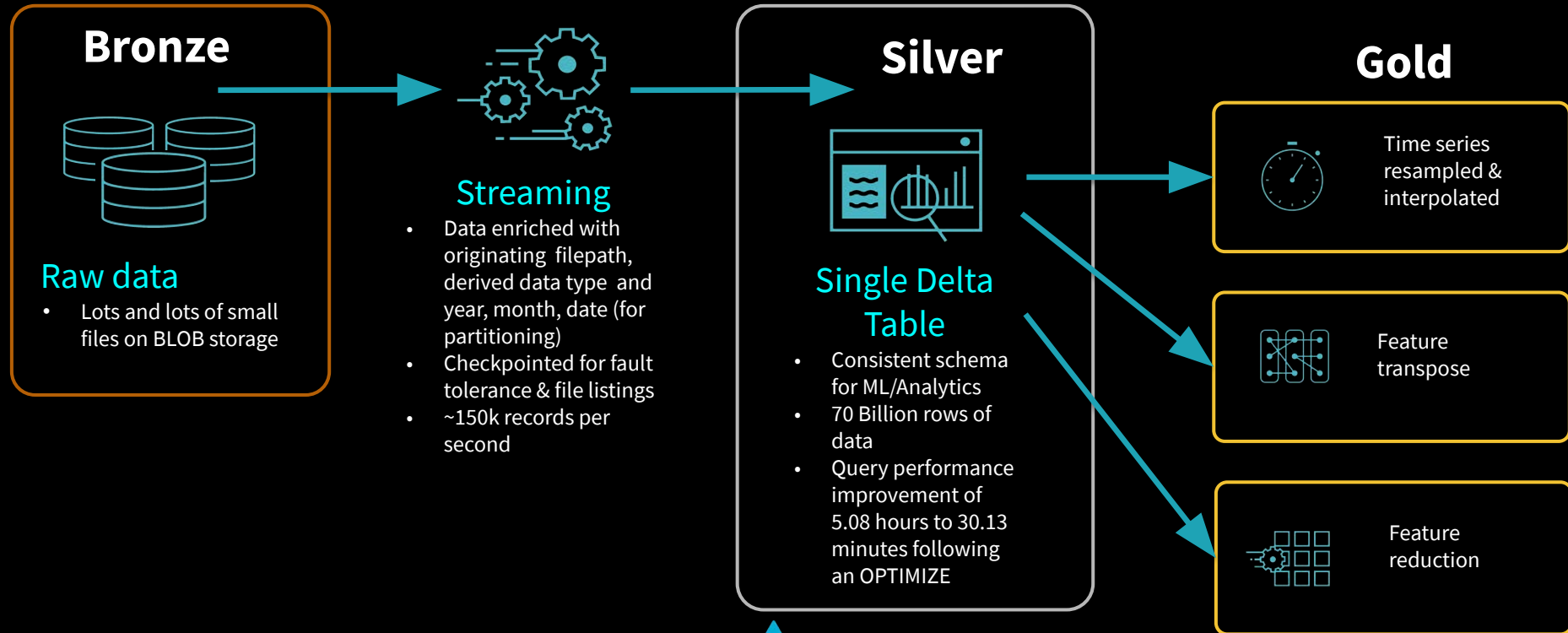


Media &  
Entertainment

- Content changes every second
- Today's patterns may not apply tomorrow

# From Data Prep to Production: Real Life IoT Example

## Process and Prepare Data at Scale



# From Data Prep to Production: Real Life IoT Example

## Train, Deploy, and Monitor at Scale

### Training data

160k sensors / 60 billion rows

Gold feature tables



### Trained models

100k

Train an sklearn model per sensor via a Pandas UDF



Model Training



Log params, metrics & model to MLflow

### Predictions

9 million in 1 hour on a 4 node cluster

Apply an sklearn model per sensor via a Pandas UDF



Model Inference



Select & load model from MLflow

Compare predicted value to rolling average via a Pandas UDF

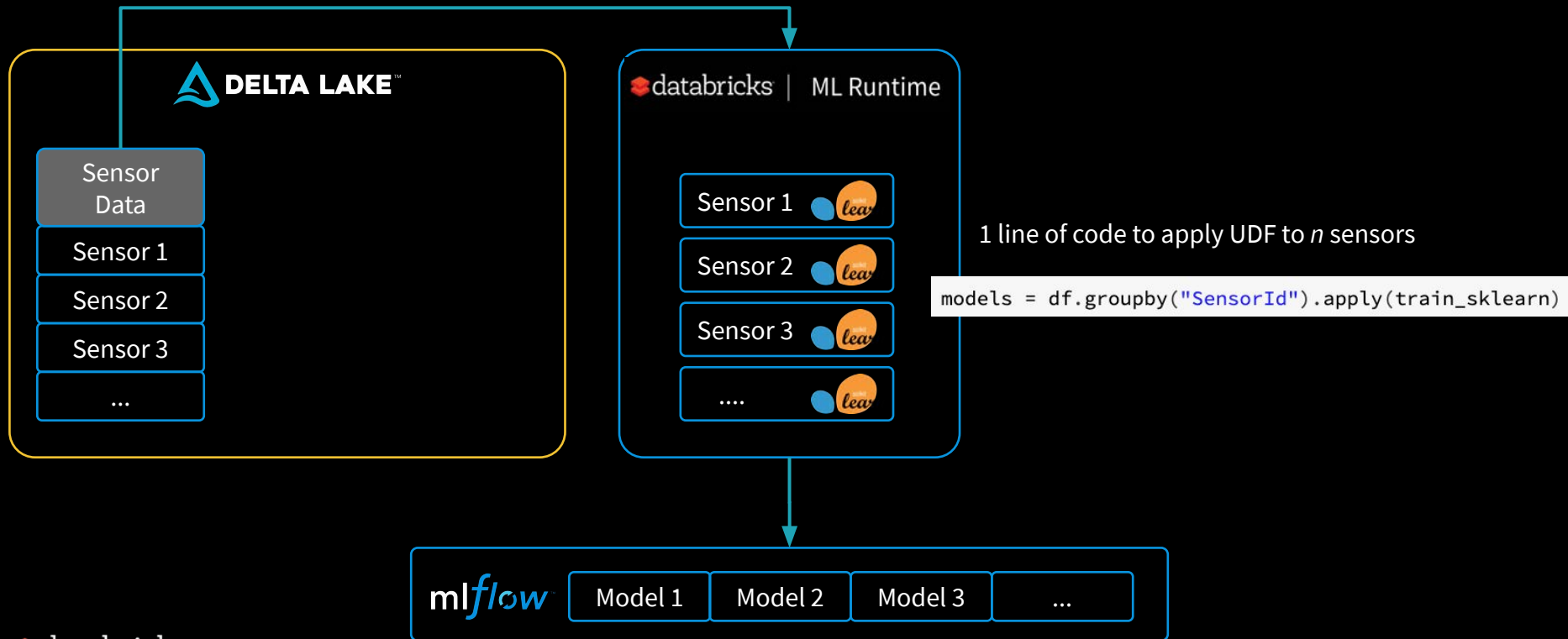


Model Monitoring



# From Data Prep to Production: Real Life IoT Example

## 1) Training at Scale



# From Data Prep to Production: Real Life IoT Example

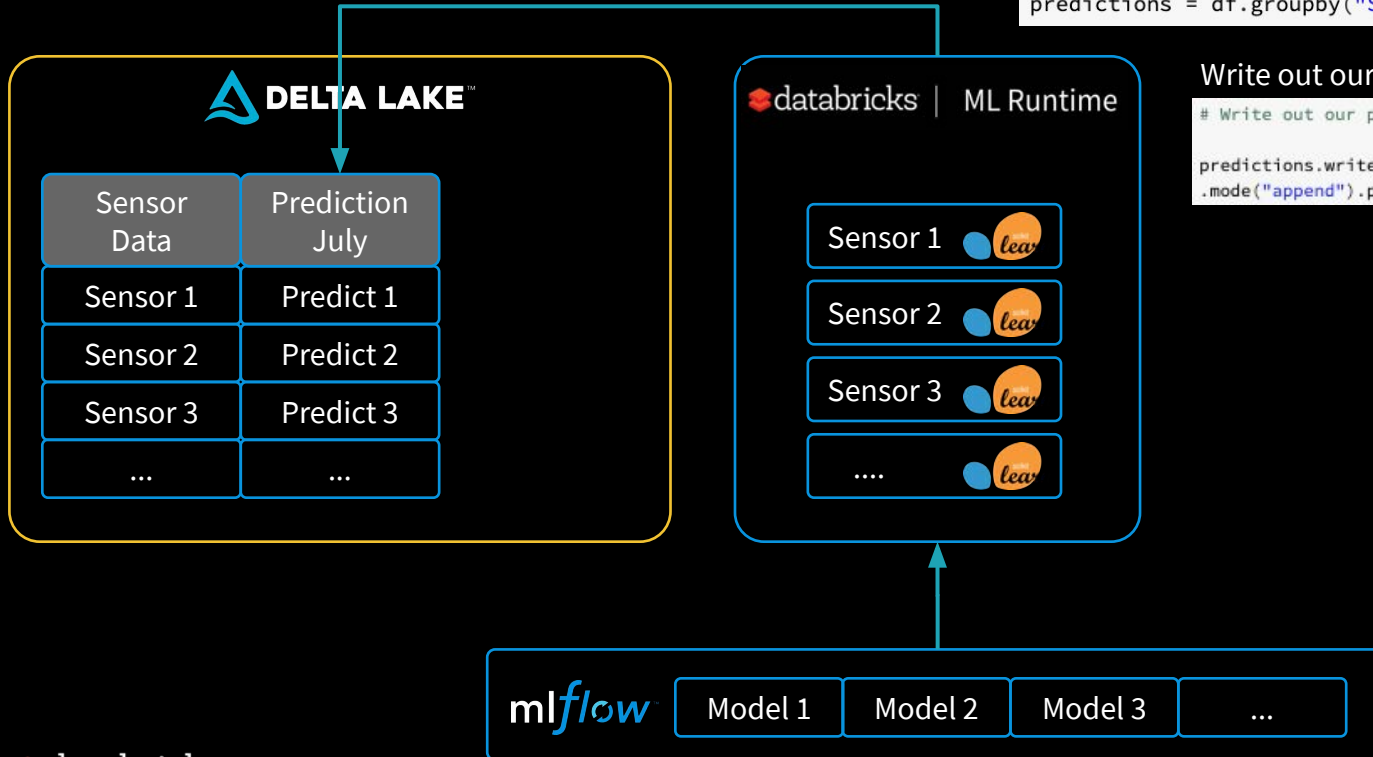
## 2) Deploying at Scale

1 line of code to apply UDF to  $n$  new sensor readings

```
predictions = df.groupby("SensorId").apply(apply_sklearn_model)
```

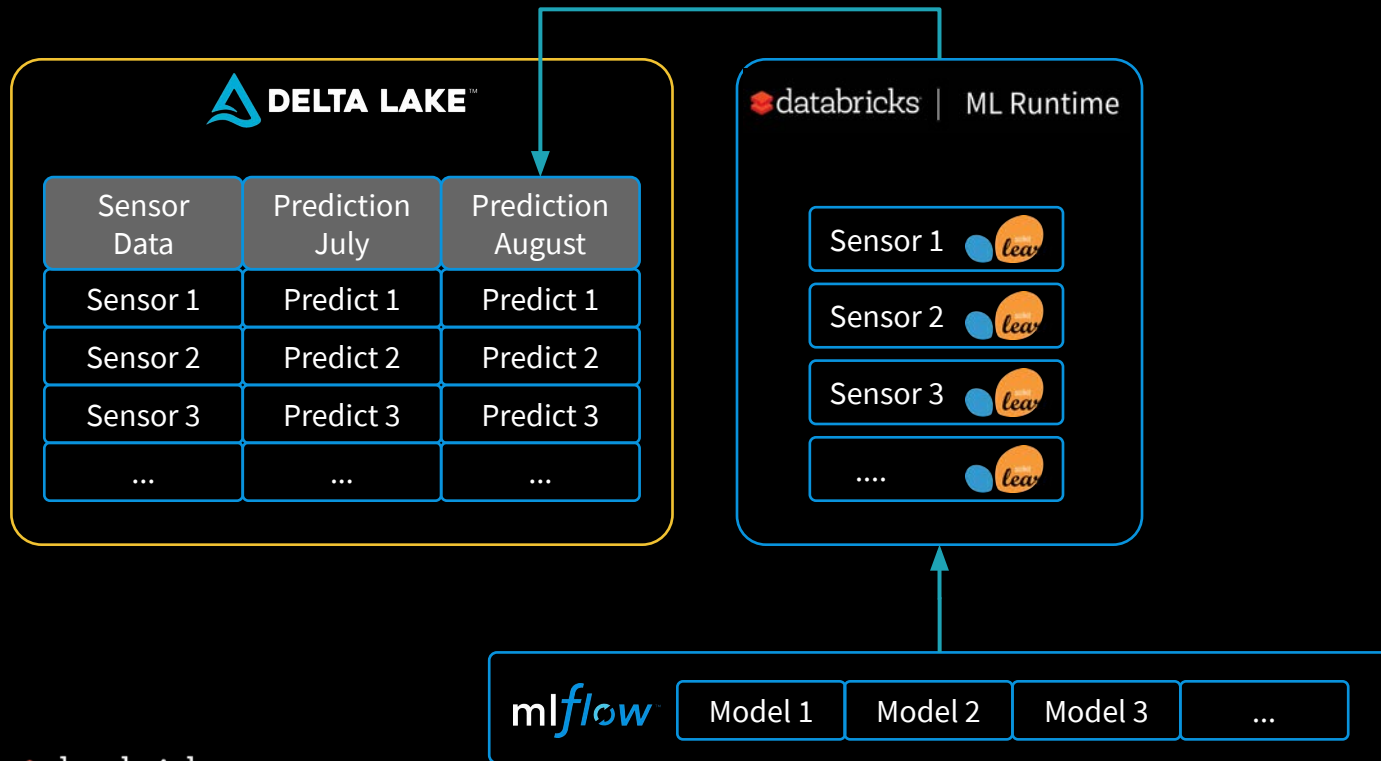
Write out our predictions to a Delta table

```
# Write out our predictions to a table
predictions.write.format("delta") \
    .mode("append").partitionBy("SensorId").save(predictions_path)
```



# From Data Prep to Production: Real Life IoT Example

## 2) Deploying at Scale

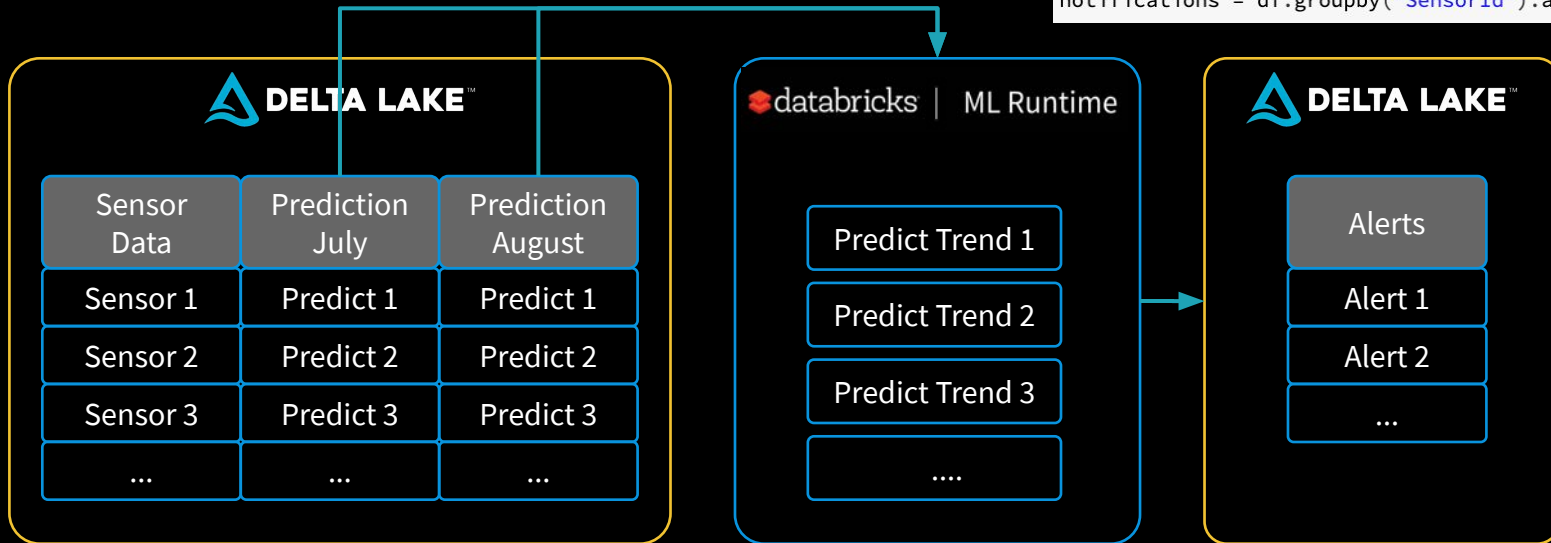


# From Data Prep to Production: Real Life IoT Example

## 3) Monitoring at Scale

1 line of code to apply UDF to  $n$  predictions

```
notifications = df.groupby("SensorId").apply(monitor_predictions)
```



Write out our alerts to a Delta table

```
# Write out our notifications to a table

notifications.write.format("delta") \
  .mode("append").partitionBy("SensorId").save(notifications_path)
```

# Demo



# Who wants a perfect cake?



**Parameters like temperature and duration have an impact on quality of cake**

# Glassware Manufacturing Dataset

Data streaming from IoT sensors

- Temperature
- Pressure
- Duration

These parameters in ideal combination gives good quality product.

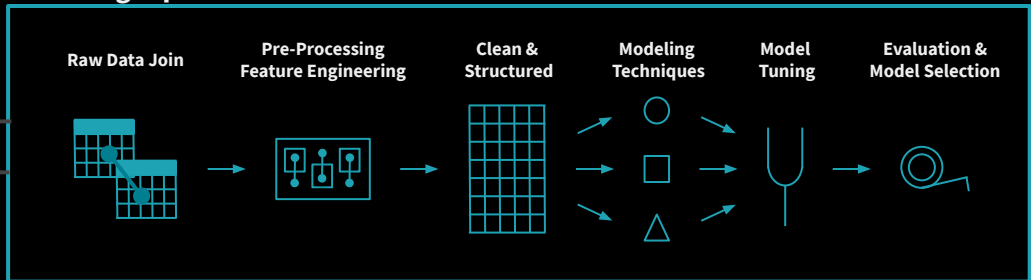


**Goal:** Predict product quality inline to be used for additional manual inspection

Synthesized dataset to showcase model drift



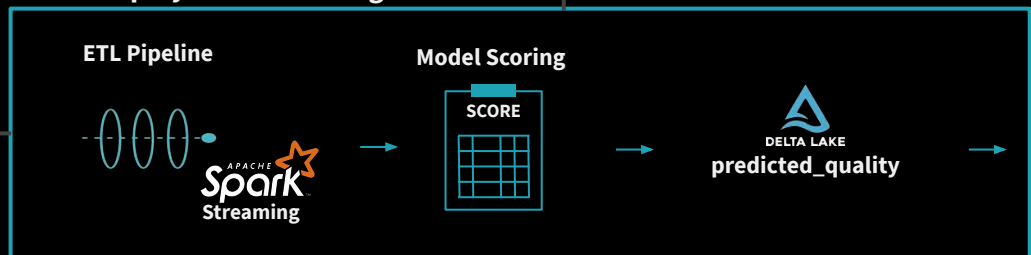
## Modeling Pipeline



Model Registry  
Best Models, Metrics,  
& Artifacts

mlflow

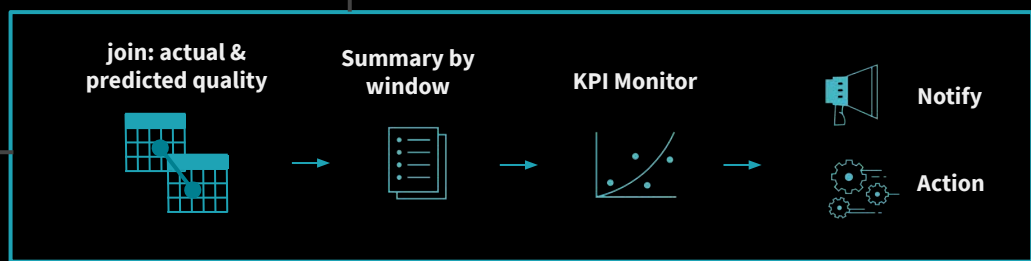
## Model Deployment & Scoring



Downstream Applications



## Model Drift Detection



APACHE  
**Spark**  
Streaming



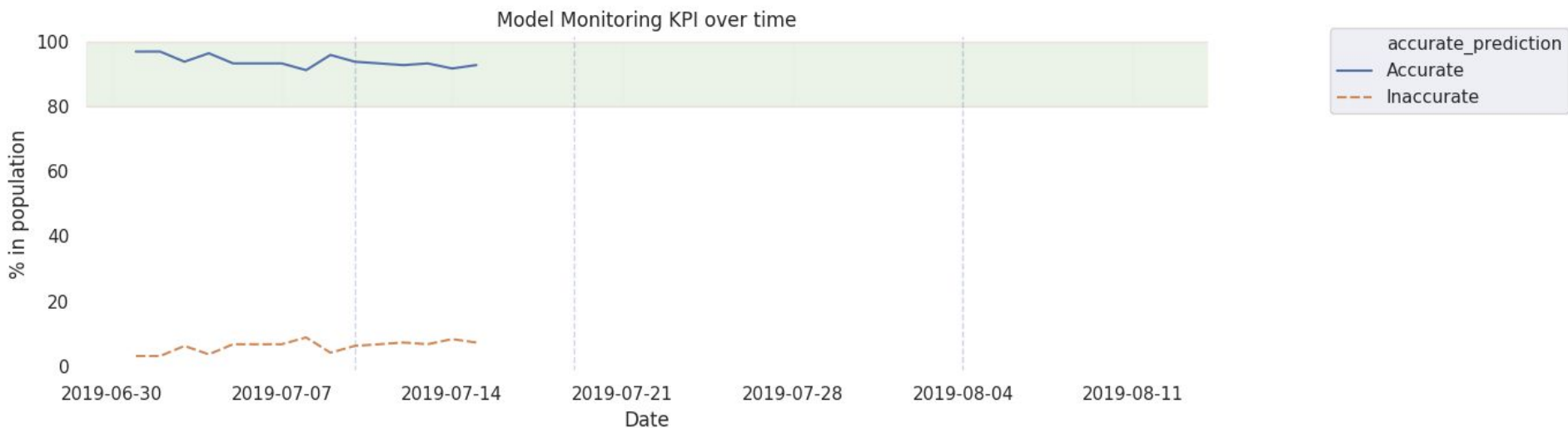
DELTA LAKE  
sensor\_reading

DELTA LAKE  
product\_quality

# Deployment to Drift Detection - a Typical Workflow

1. To understand the data, we start with EDA (Exploratory Data Analysis)
2. Using historical data, we explore various modeling methods, tune its hyperparameters, and identify our best model
3. All the experiment runs are tracked using MLflow and we tag the best model for production use
4. While scoring in a streaming pipeline, production model is accessed from MLflow
5. Model is stable for first 'x' days

# KPI to Monitor Model Quality

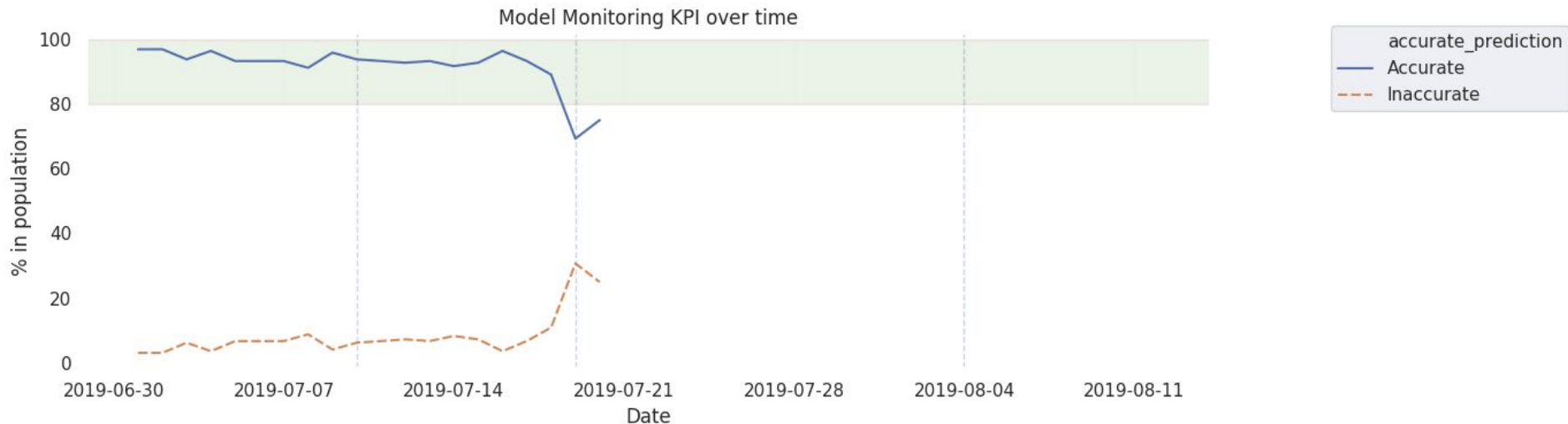


## Model Drift KPIs:

- KPIs and its margin depends on the model and business problem
- Sometimes more than 1 KPI maybe needed at times to capture behavior changes

# Deployment to Drift Detection - a Typical Workflow

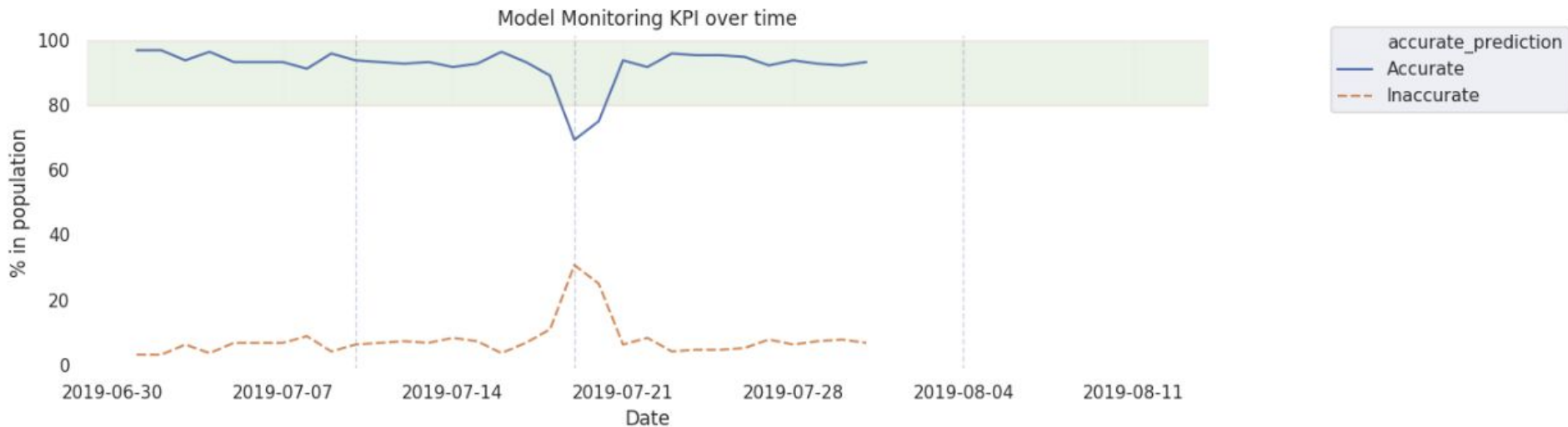
- After 'y' days, we see model drift occur, as identified by tracking KPIs
- This triggers re-training process



# Deployment to Drift Detection - a Typical Workflow

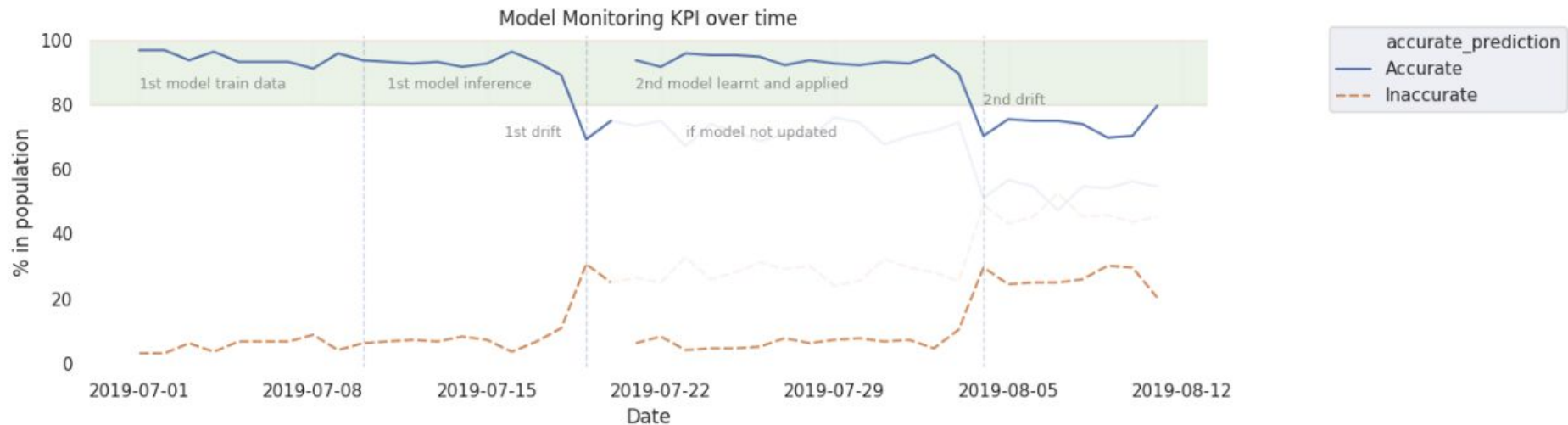
- Once again, we explore various modeling methods, tune its hyperparameters, and identify our new best model
- The new model is tagged as current production model in MLflow
- We once again observe that KPIs are back within acceptable range
- Over time, based on business demands, it may be needed to update KPIs and its acceptable limits

# Deployment to Drift Detection - a Typical Workflow





# Summary - Deployment to Drift Detection



# Model Drift Demo: Source Code

<https://github.com/joelcthomas/modeldrift>

# Summary

- / The only constant is change. Concept drift, data drift, and upstream data problems all affect your ML pipeline.
- / Monitoring data and ML model performance is critical.
- / In many cases, your ML model needs to adapt to changes.
- / Databricks provides an end-to-end platform to deploy, monitor, and automatically retrain your ML models.

# Coming soon

## **mlflow** Model Registry

Tagging, sharing and versioning models in the MLflow tracking server

- Register any model in MLflow Model format
- Deploy to serving systems
- Add metadata (e.g. tags/notes) and track model creators and users



- 15-17 October 2019 | Amsterdam
- New product announcements around Delta Lake, MLflow, etc.
- Sign up at [databricks.com/sparkaisummit](https://databricks.com/sparkaisummit)



# Q&A

**Get Started** - [databricks.com/try](https://databricks.com/try)

**Contact us** - [databricks.com/contact](https://databricks.com/contact)