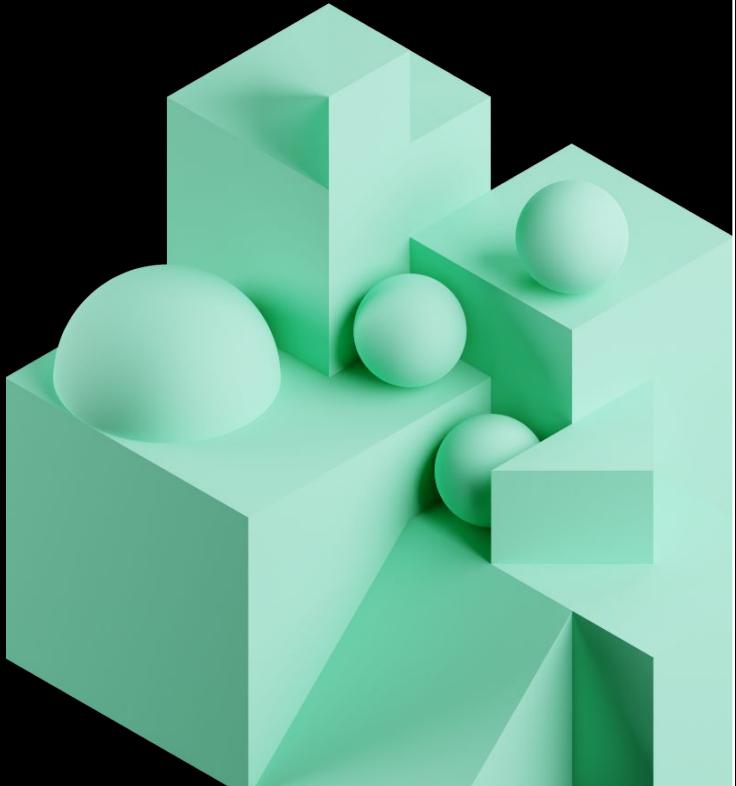




# Machine Learning in Production

---

Databricks  
2023

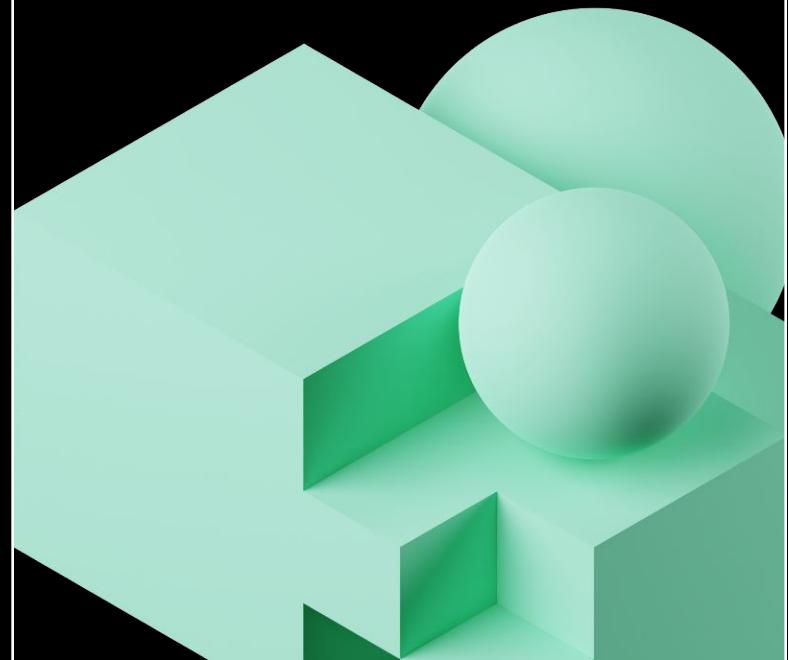


# Agenda

MORNING			AFTERNOON		
	Demo	Lab		Demo	Lab
<b>Introduction</b>			<b>O3. Deployment Paradigms</b>		
<b>O1. Experimentation</b>			<b>DevOps vs. ModelOps</b>		
<b>Delta Lake</b>			<b>Deployment Patterns</b>	✓	
<b>Feature Store</b>	✓		<b>Deployment Paradigms</b>	✓	✓
<b>MLflow Experiment Tracking</b>	✓	✓	<b>O4. Production</b>		
<b>O2. Model Management</b>			<b>CI/CD</b>		
<b>MLflow Model Management</b>	✓	✓	<b>Monitoring</b>	✓	✓
<b>MLflow Model Registry</b>	✓		<b>Model Rollout Strategies</b>		

# LET'S GET STARTED

# Course Introduction



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

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

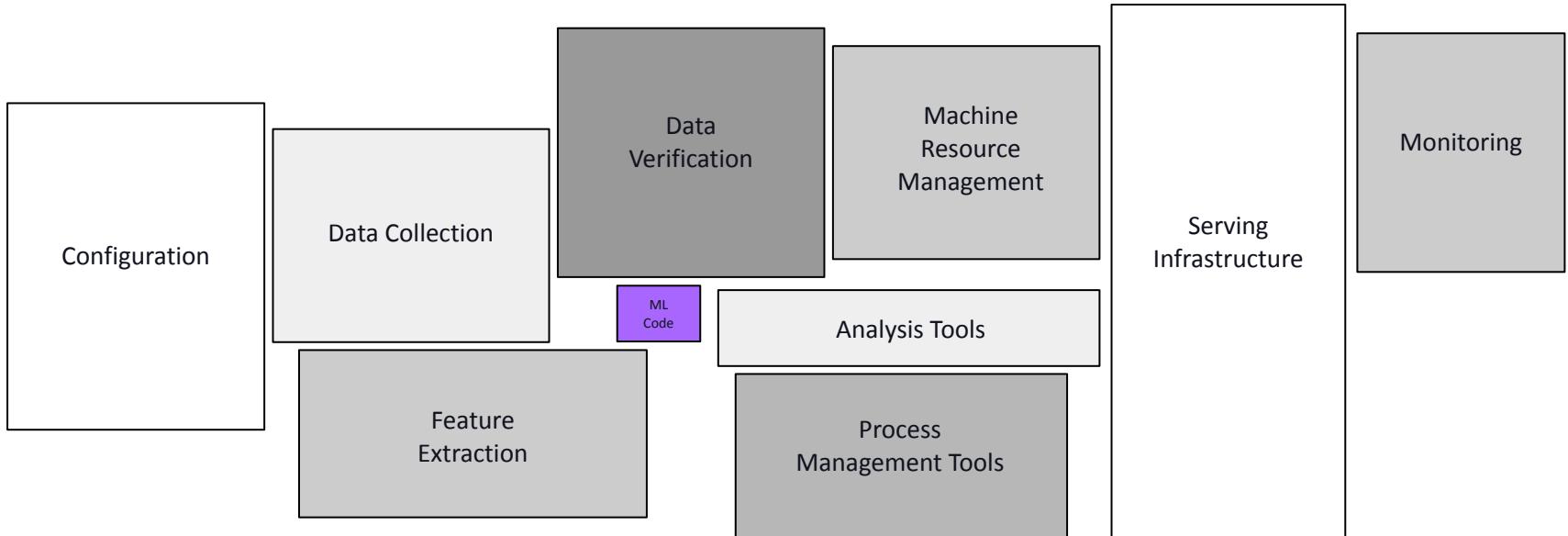


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small green box in the middle. The required surrounding infrastructure is vast and complex.

# What is MLOps and why should we care about it?

# What is MLOps?

MLOps is a set of **processes** and **automation** for managing **code**, **data**, and **models** to improve performance, stability and long-term efficiency of ML systems

**MLOps** = **DevOps** + **DataOps** + **ModelOps**



# Why Should We Care About MLOps?

MLOps helps you **reduce risk**

- **Technical risk** – poorly performing models, fragile infrastructure
- **Compliance risk** – violating regulatory or corporate policies

MLOps improves **long-term efficiency** through **automation**

- Catch errors before they hit production
- Avoid slow, manual processes

# People and Process



Business Stakeholder

→ Responsible for business value of the ML solution



Data Engineer

→ Builds data pipelines



Data Scientist

→ Translates business problem; trains, tunes model



ML Engineer

→ Deploys ML model to production



Data Governance Officer

→ Responsible for data governance and compliance



## databricks

Lakehouse Platform

DevOps



DataOps



ModelOps



# Data-Centric ML Platform

# Guiding Principles



Always keep your business goals in mind



Take a data-centric approach to machine learning



Implement MLOps in a modular fashion



Process should guide automation

# Data & Model Management and Reproducibility

Managing the machine learning lifecycle means:

- **Reproducibility of data** -- data management
- **Reproducibility of code** -- version control, source tracking
- **Reproducibility of models** -- model management
- **Automated integration** with production systems -- CI/CD, testing



# MLflow

An open-source platform for managing end-to-end ML lifecycle



- Open-source platform for machine learning lifecycle
- Operationalizing machine learning
- Developed by Databricks
- Pre-installed on the Databricks Runtime for ML

# MLflow Components

The four components of MLflow



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



Packaging format for reproducible runs on any platform



General model format that supports diverse deployment tools



Centralized and collaborative model lifecycle management

APIs: CLI, Python, R, Java, REST

# Data is never static!



**Model Centric → Data Centric Machine Learning**

# The Full ML Lifecycle

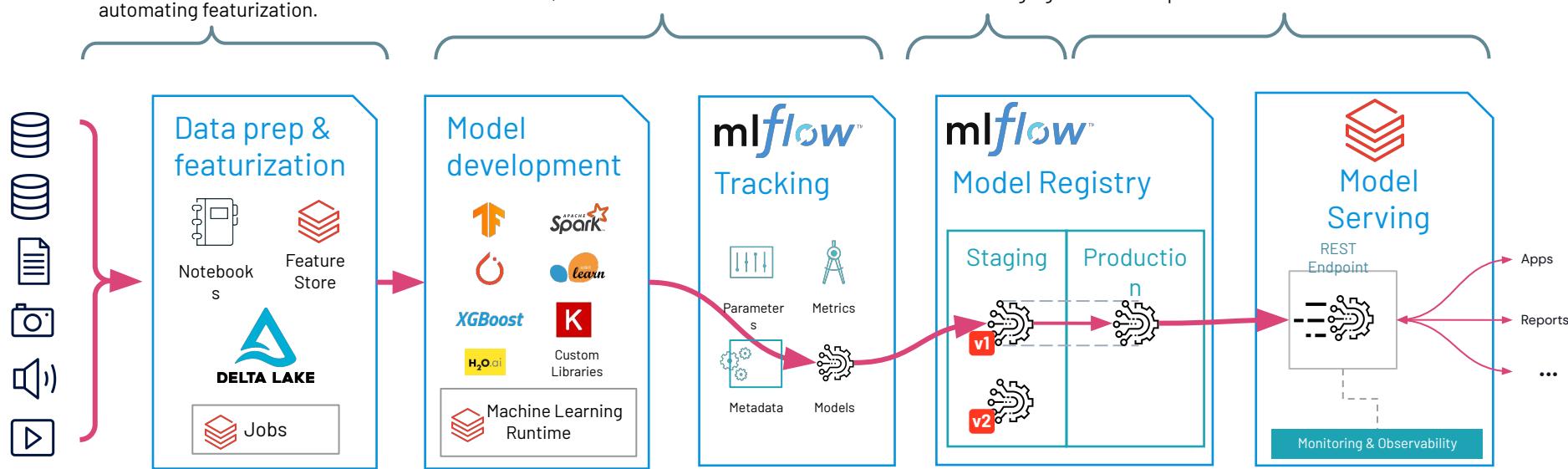
Build, train, test, deploy and monitor machine learning models with ease

**Data Scientists** build features.  
**Data Engineers** provide infra for automating featurization.

**Data Scientists** build models and log them to MLflow, which records environment info.

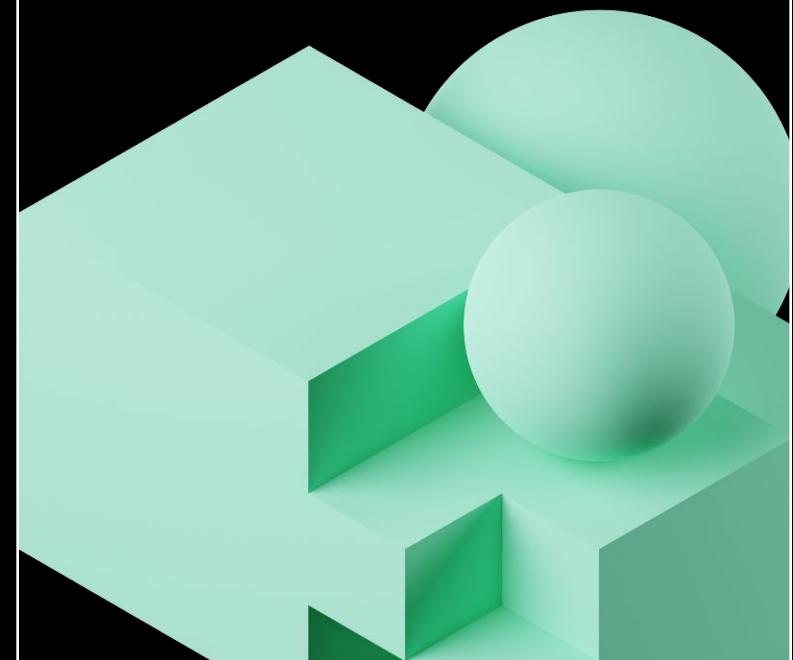
**Data Scientists** move models to Staging.

**ML Engineers** serve models as REST endpoints.

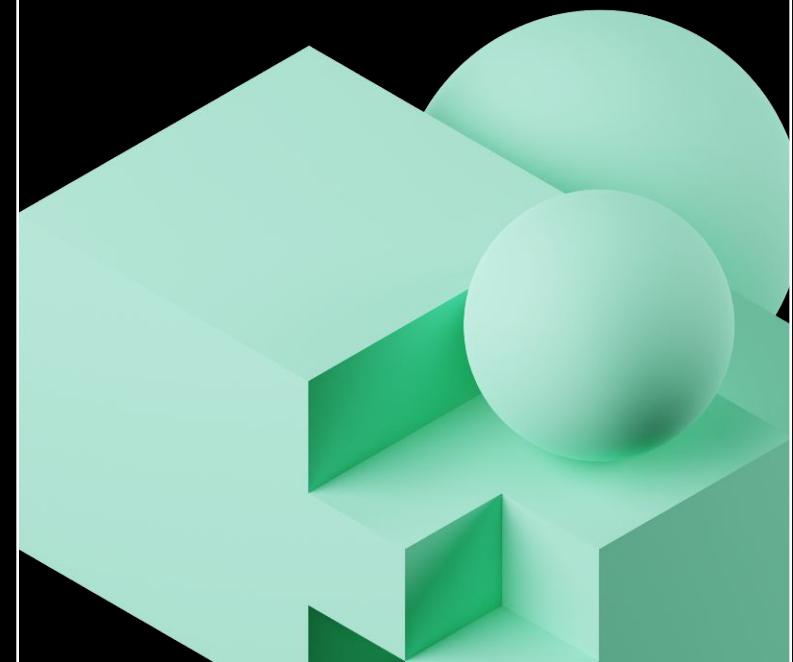


# WORKSPACE SETUP

# 01. Experimentation

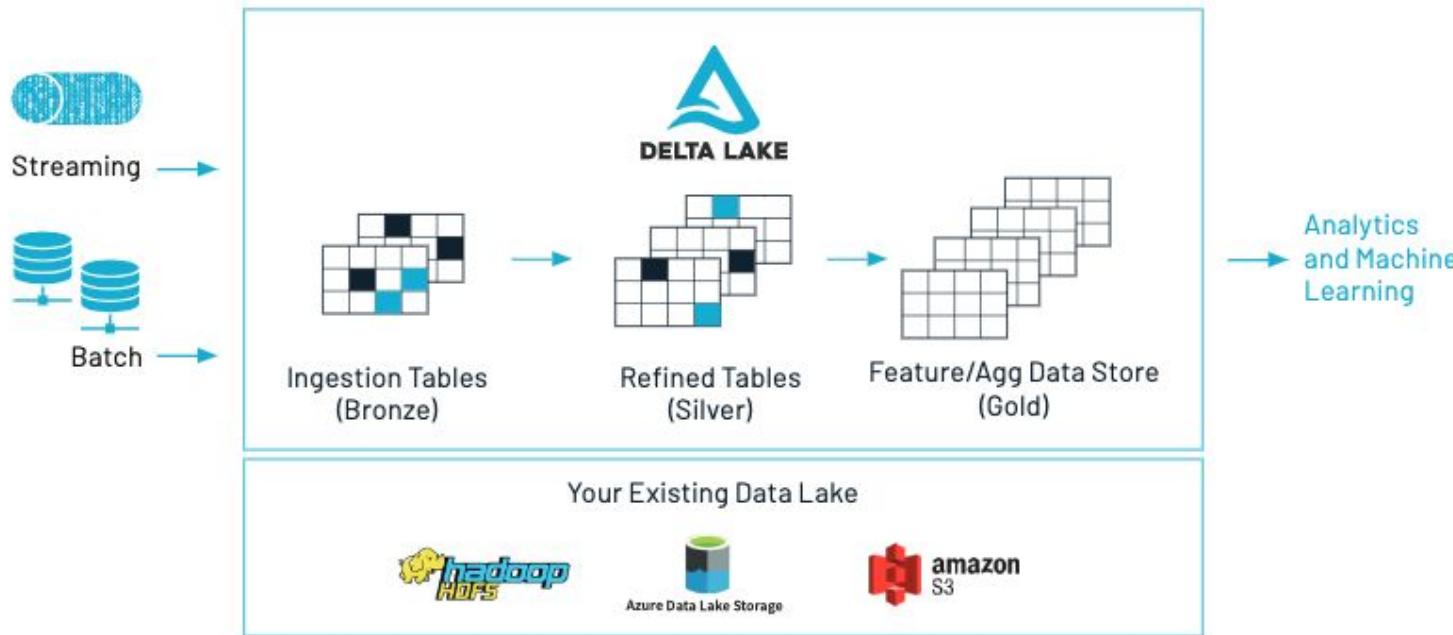


# Delta Lake



# Delta Lake Overview

An open-source storage framework for building a modern Lakehouse



# Delta Lake Overview

## Key Features

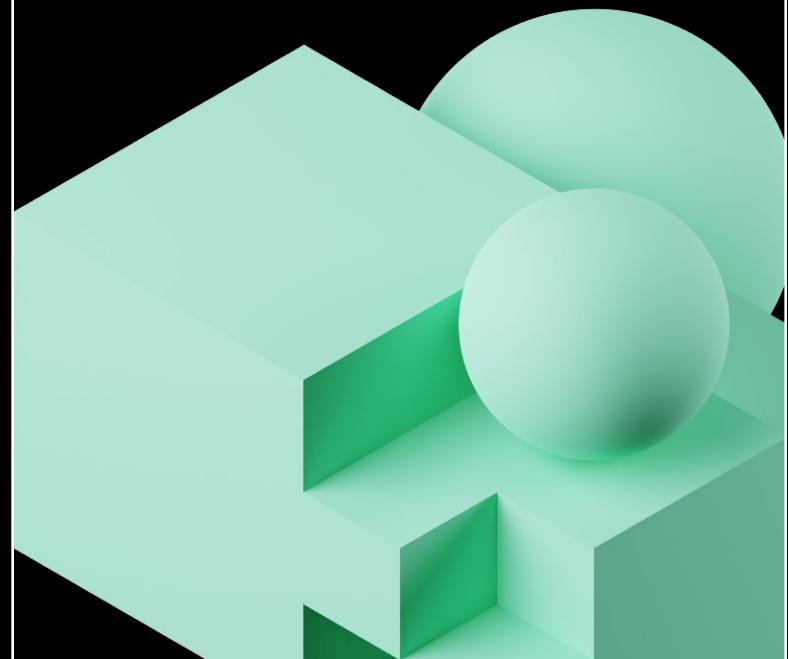
- ACID transactions → reliability
- Time travel (data versioning) → reliability, error recovery, synchronization
- Schema enforcement and evolution → reliability
- Audit history → data governance and compliance
- Parquet format + transaction log
- Compatible with Apache Spark API

# Delta Lake Overview

## Optimization features

- Compaction
- Data Skipping
- Z-ordering
- FAQs on “OPTIMIZE”

# Feature Store



Production Machine Learning  
= Production Software + Data

# What is a Feature?

An example feature in a recommendation system

## Raw data

User profile
Zip code
Payment methods
...
User history
Last 6 months of purchases
Web session
User ID
Shopping cart
Item view in last hour
...

## Types of Features

### Feature Augmentation

e.g. Weather in region

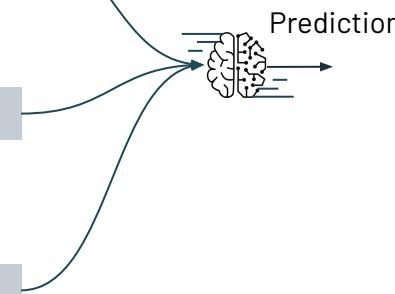
### Pre-computed Features

e.g. Purchases last 7, 14, 21 days

### Transformations of real-time data

e.g. Category Encoding

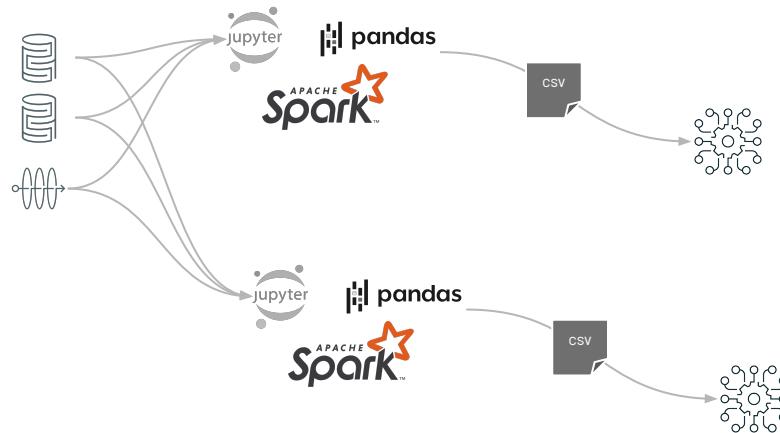
## ML Model



## Outcome

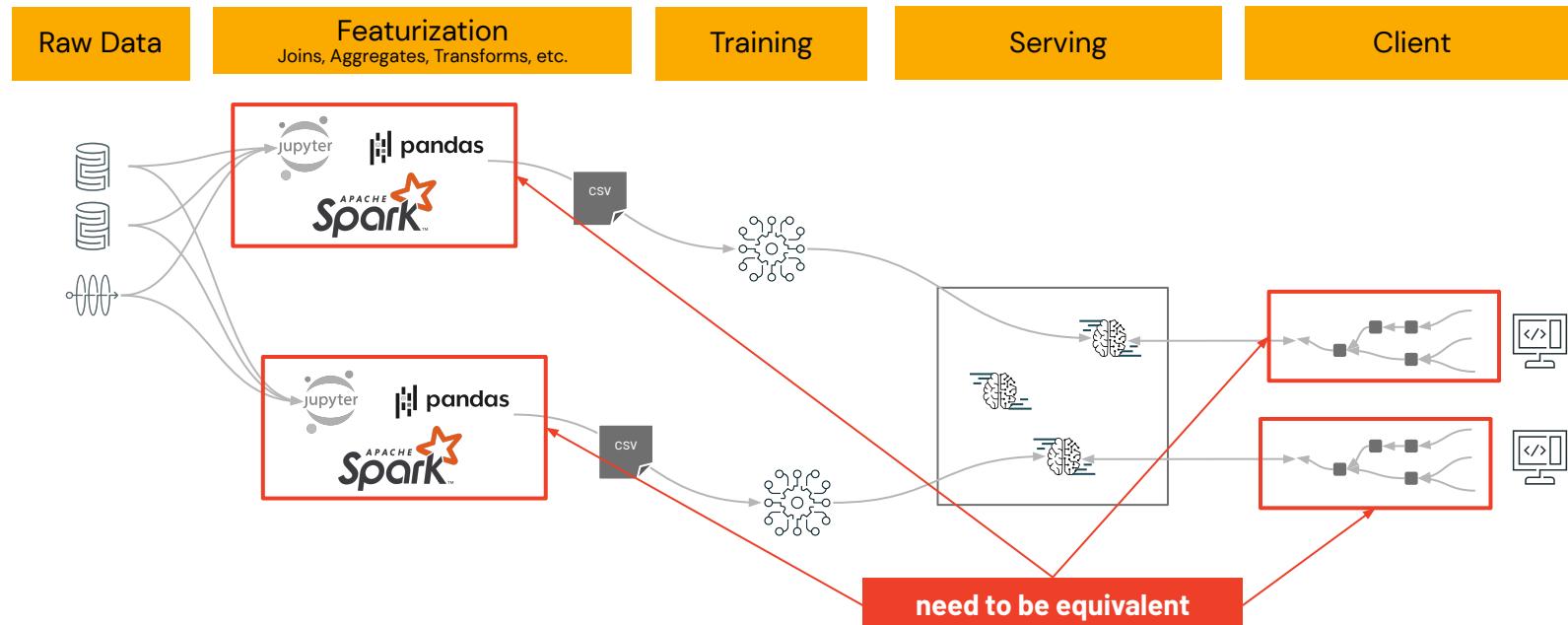
Item	P(purchase user)
Power drill	0.58
Lightbulb	0.13
Wrench	0.12
Screwdriver	0.01

# Data Challenges in ML



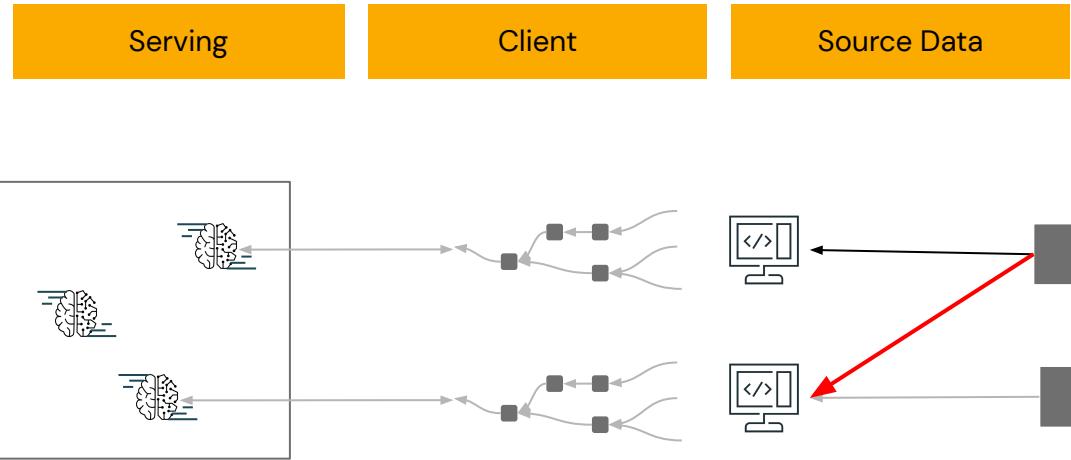
Challenge 1: Data silos

# Data Challenges in ML



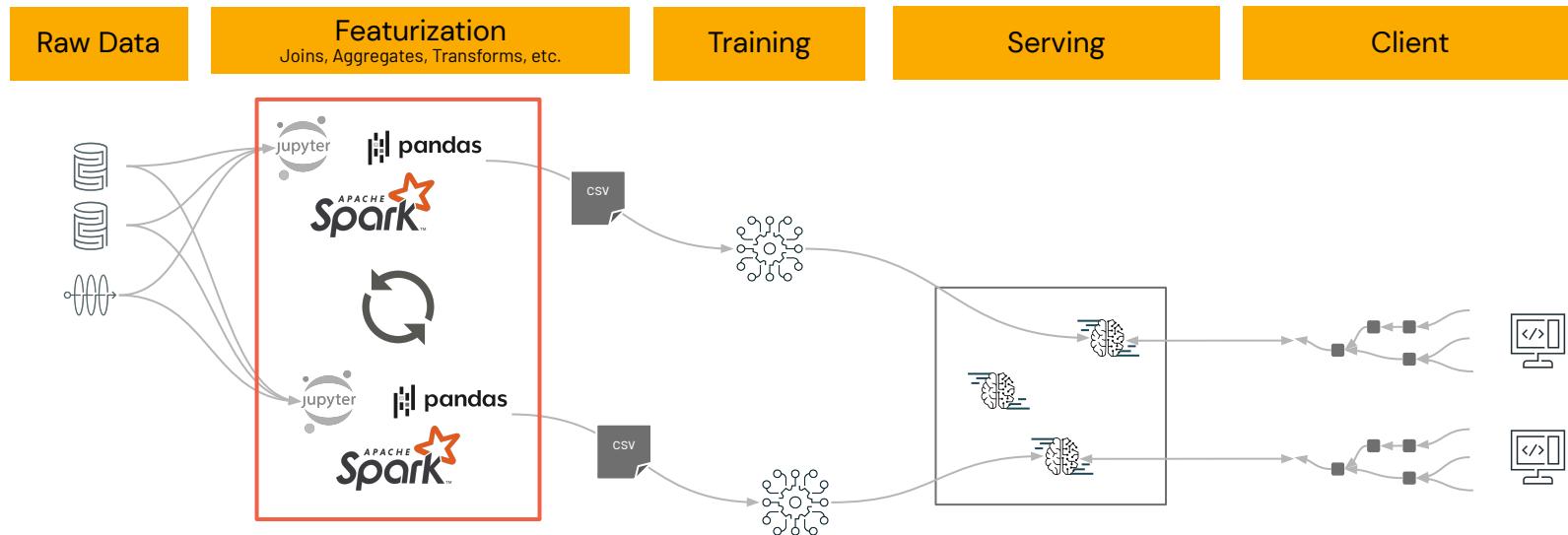
Challenge 2: Online/offline skew

# Data Challenges in ML



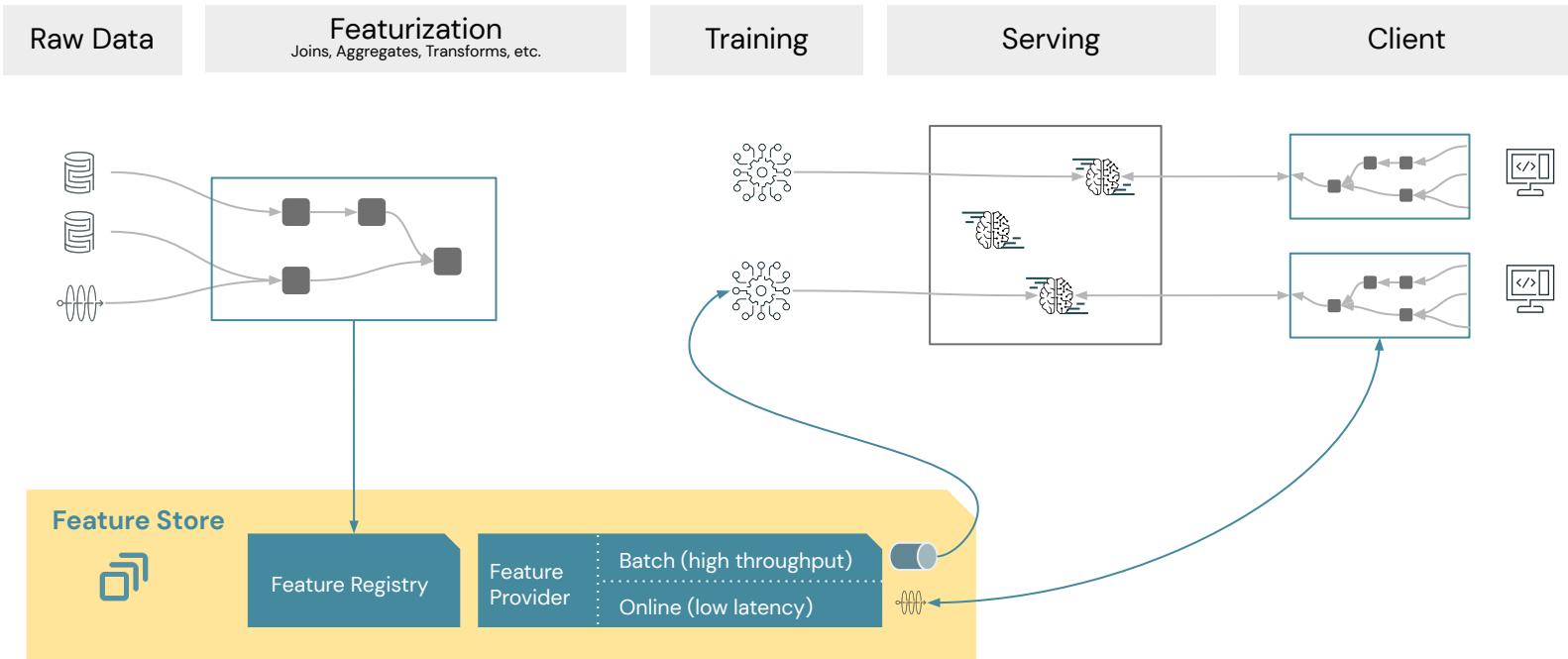
Challenge 3: Client Configuration

# Data Challenges in ML



Challenge 4: Data Freshness

# Solving the Feature Store Problem



# Feature Store

The first Feature Store codesigned with a Data and MLOps Platform



## Feature Registry

- Discoverability and Reusability
- Versioning
- Upstream and downstream Lineage

Co-designed with



- Open format
- Built-in data versioning and governance
- Native access through PySpark, SQL, etc.

## Feature Provider

- Batch and online access to Features
- Feature lookup packaged with Models
- Simplified deployment process

Co-designed with [mlflow](#)

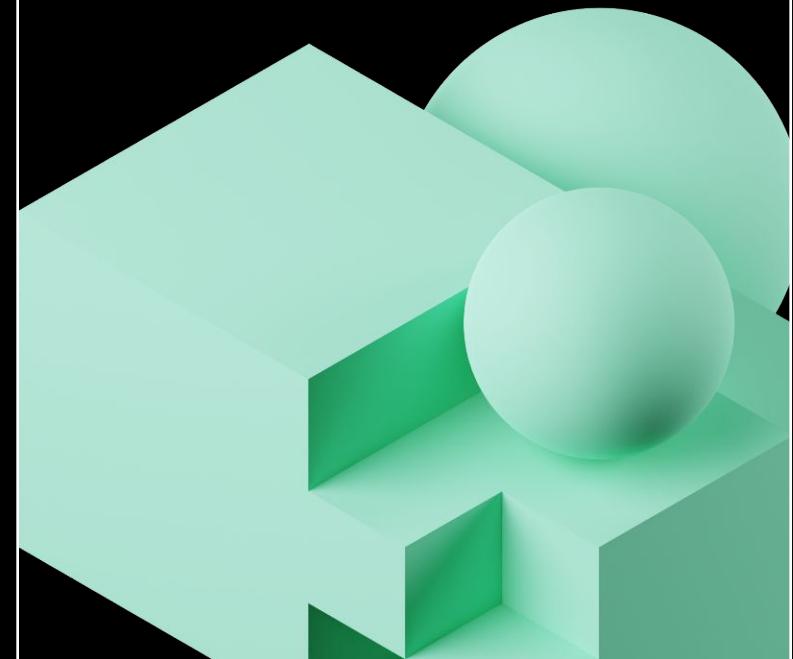
- Open model format that supports all ML frameworks
- Feature version and lookup logic hermetically logged with Model

# DEMO: FEATURE STORE

Notebook: 01-Experimentation → 01-Feature-Store



# MLflow Tracking



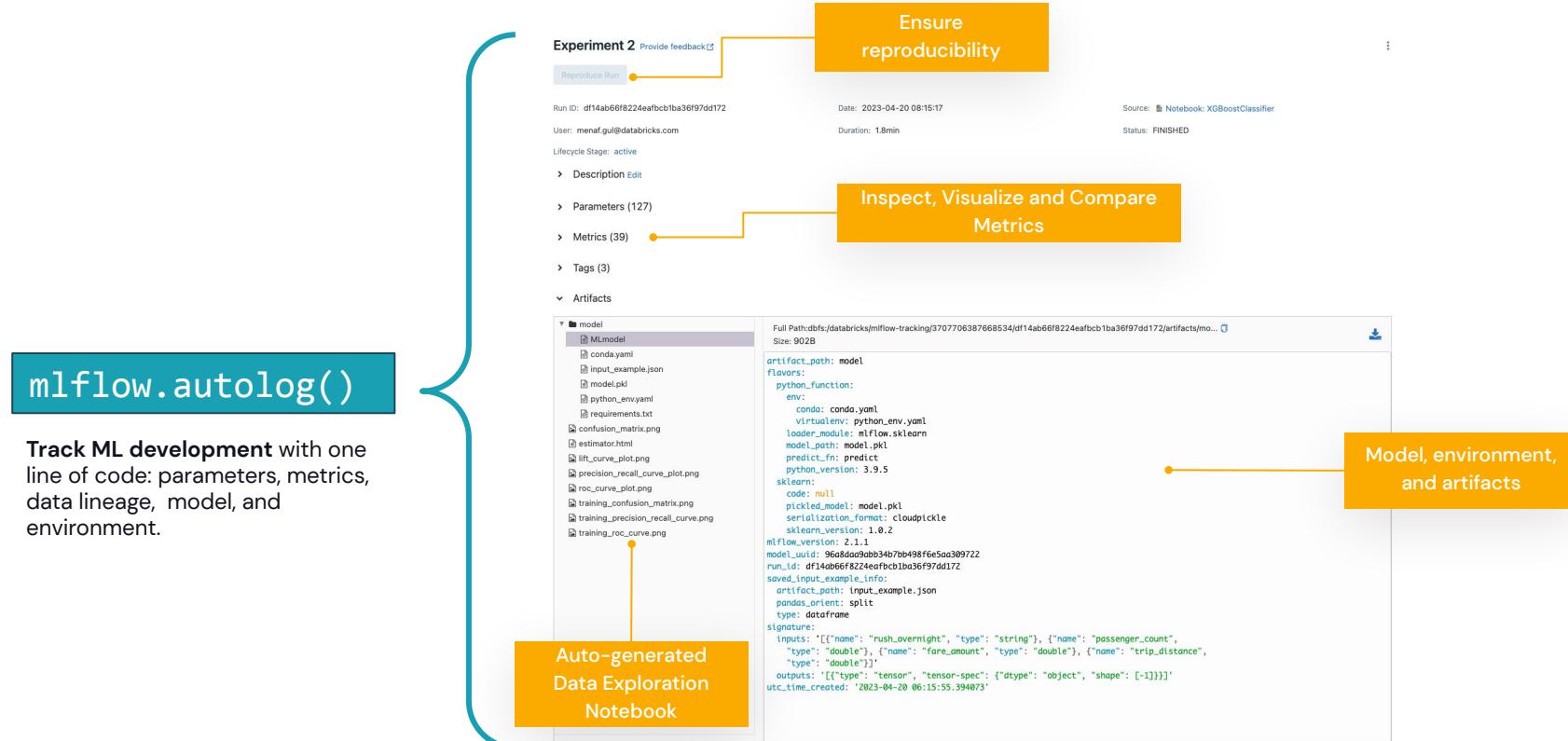
# Experiment Tracking

Record runs, and keep track of models parameters, results, code, and data from each experiment and in one place.

## Provides:

- Pre-configured MLflow tracking server
- Databricks Workspace & Notebooks UI integration
- S3, Azure Blob Storage, Google Cloud for artifacts storage
- Experiments management via role based Access Control Lists (ACLs)

# MLflow Tracking and Auto-logging



# Workflow Overview

- Write code to convert raw data into features and create a Spark DataFrame
- Write the DataFrame as a feature table in Feature Store
- Create a training set based on features from feature tables
- Train model and log it with MLflow
- Perform batch inference on new data. The model automatically retrieves the features it needs from Feature Store.
- (Optional) Publish the features to an online store for batch or real-time inference

# DEMO: **EXPERIMENT TRACKING**

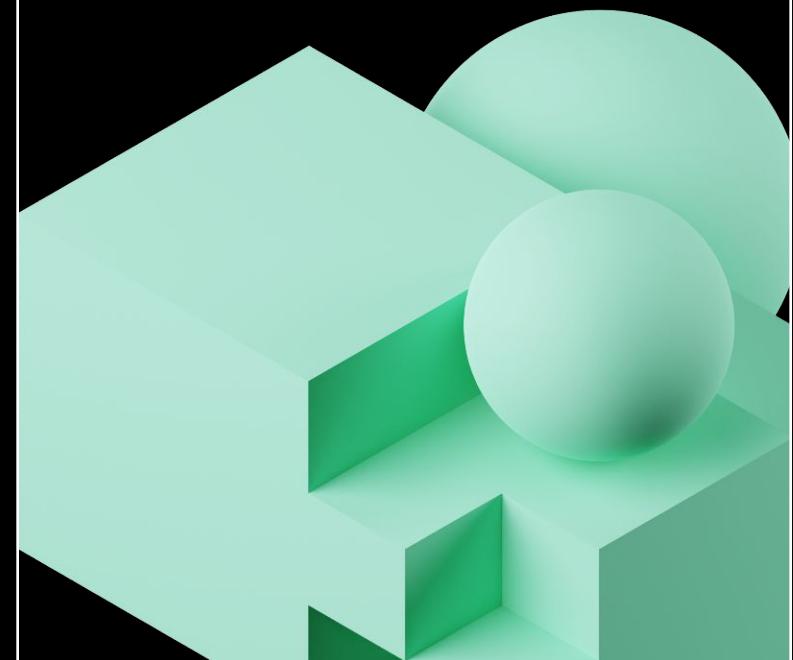
Notebook: 01-Experimentation → 02-Experiment-Tracking



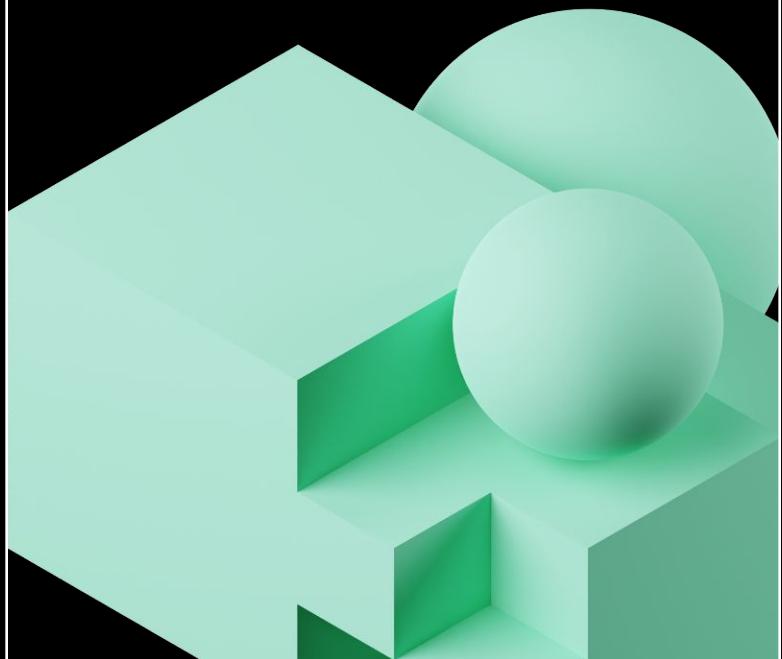
# LAB: EXPERIMENT TRACKING

Notebook: 01-Experimentation → Labs → 02-Experiment-Tracking-Lab

# 02. Model Management



# MLflow Models

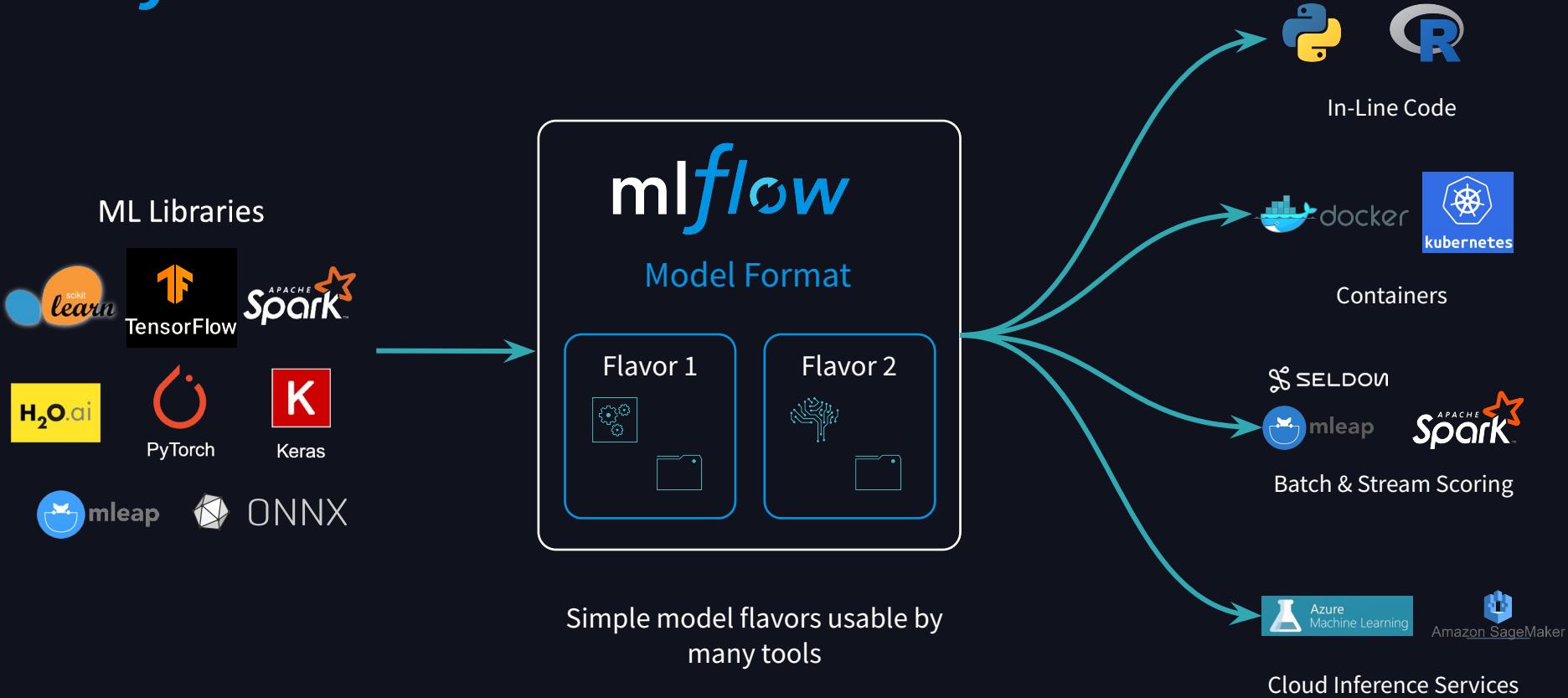


# Data Transformation

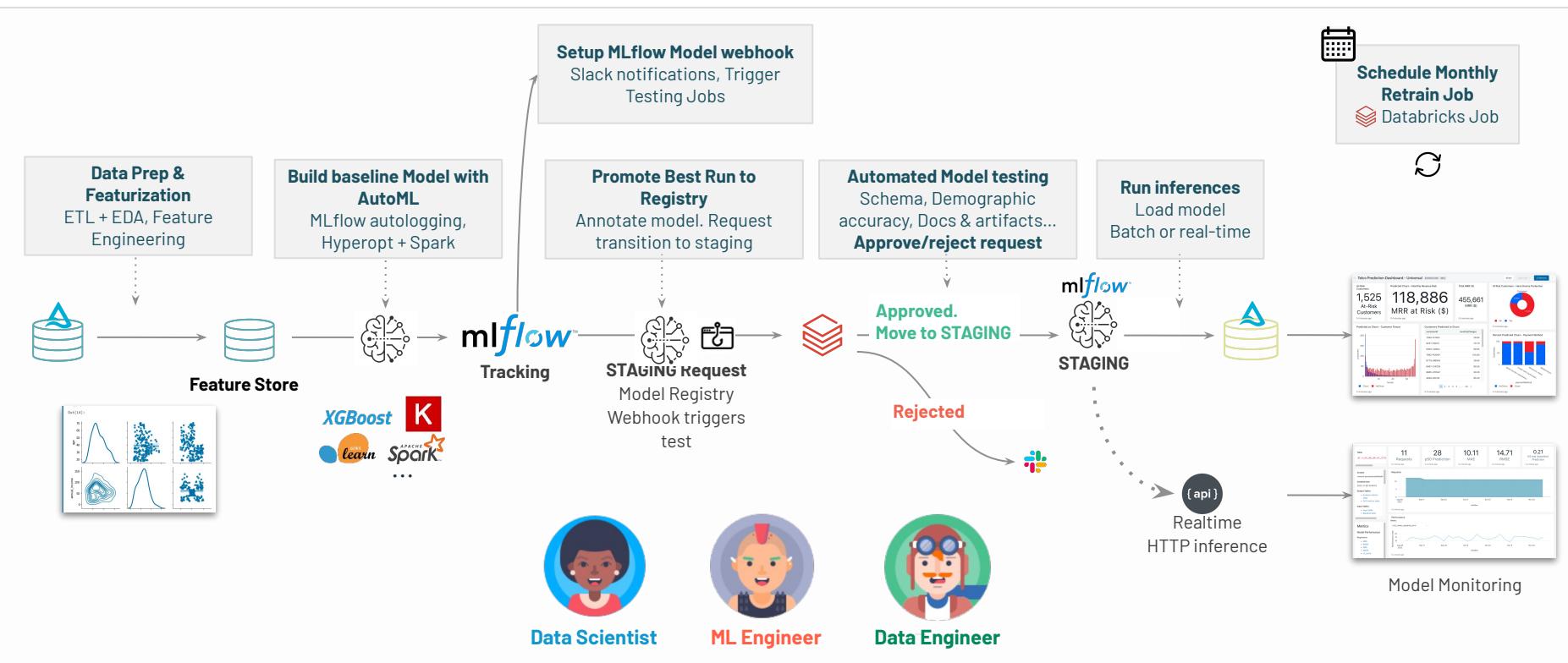
Within the model or prior to model training? It depends!

	Pros	Cons
Prior to Model Training	<ul style="list-style-type: none"><li>- Compute once and reuse</li><li>- Leverages full dataset</li></ul>	<ul style="list-style-type: none"><li>- Increases production footprint</li><li>- Slower to iterate</li></ul>
Within the Model	<ul style="list-style-type: none"><li>- Easier to iterate</li><li>- Smaller production footprint</li></ul>	<ul style="list-style-type: none"><li>- Model latency depends upon transformation overheads</li><li>- Data visibility</li></ul>

# mlflow Models



# End to End Model Management



# Model Management with Model Registry

Use one central place to collaboratively manage ML models



**One Collaborative Hub:** The Model Registry provides a central hub for making models discoverable, improving **collaboration** and **knowledge sharing** across the organization.



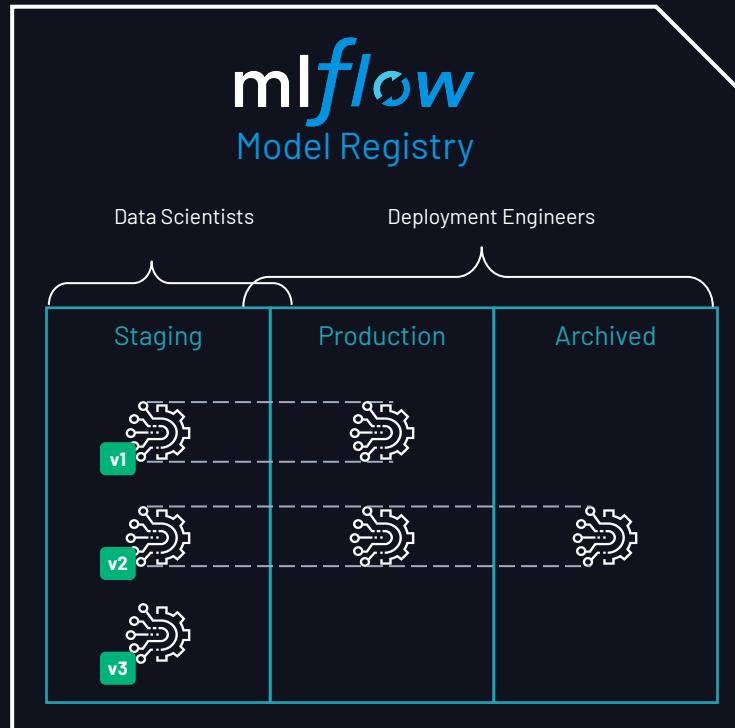
**Manage the entire Model Lifecycle (MLOps):** The Model Registry provides lifecycle management for models from experimentation to deployment, improving **reliability** and robustness of the model deployment process.



**Visibility and Governance:** The Model Registry provides full visibility into the deployment stage of all models, who requested and approved changes, allowing for full governance and auditability.

# mlflow Model Registry

VISION: Centralized and collaborative model lifecycle management



# DEMO: MODEL MANAGEMENT

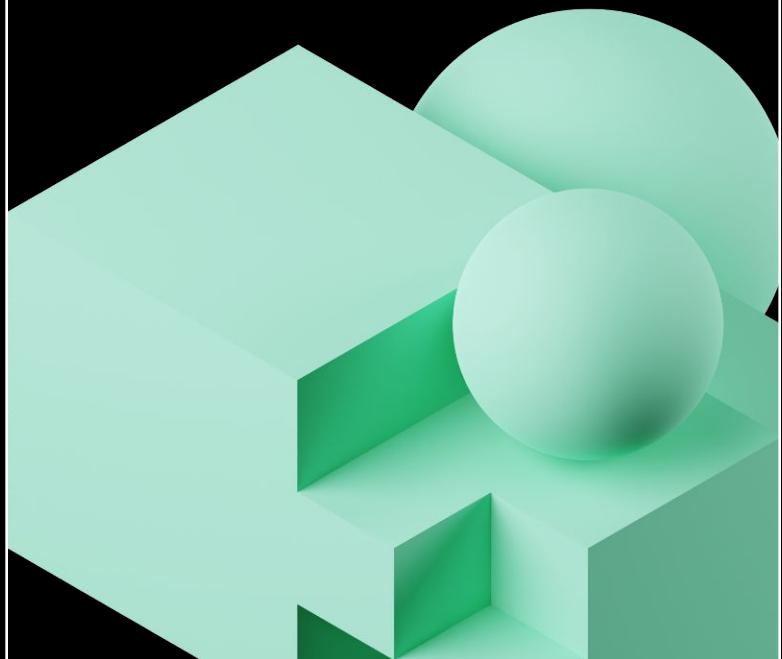
Notebook: 02-Model-Management → 01-Model-Management



# LAB: MODEL MANAGEMENT

Notebook: O2-Model-Management → Labs → 01-Model-Management-Lab

# Model Registry



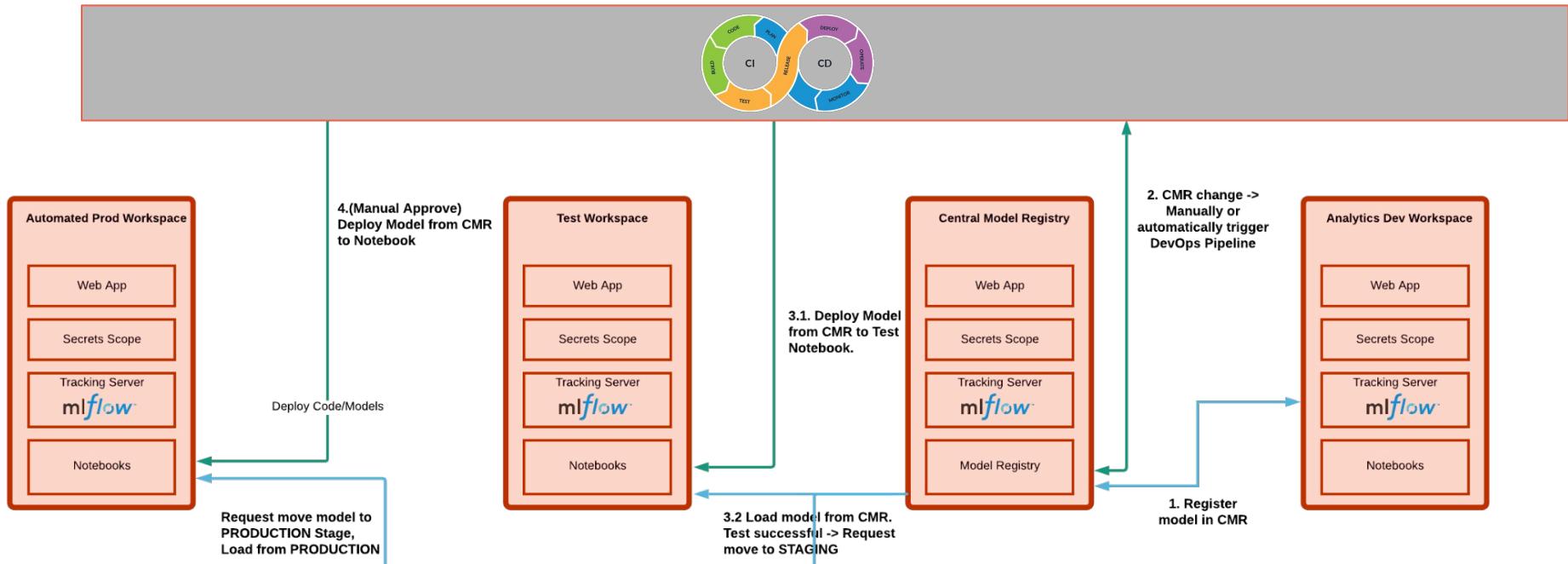
# Model Registry Workflow Options

## Central Model Registry – Model Access Management

- **Central Model Registry:** Dedicated workspace(s) to store models and experiments
  - Refer to this [documentation](#)
- **Per-workspace Model Registries:** Each workspace has its own model registry, model transition is made via CI/CD
- **Model Registry - *Unity Catalog* Integration:** Centrally govern model access permissions, share models across workspaces, model lineage

# Central Model Registry

Dedicated workspace(s) to store models and experiments



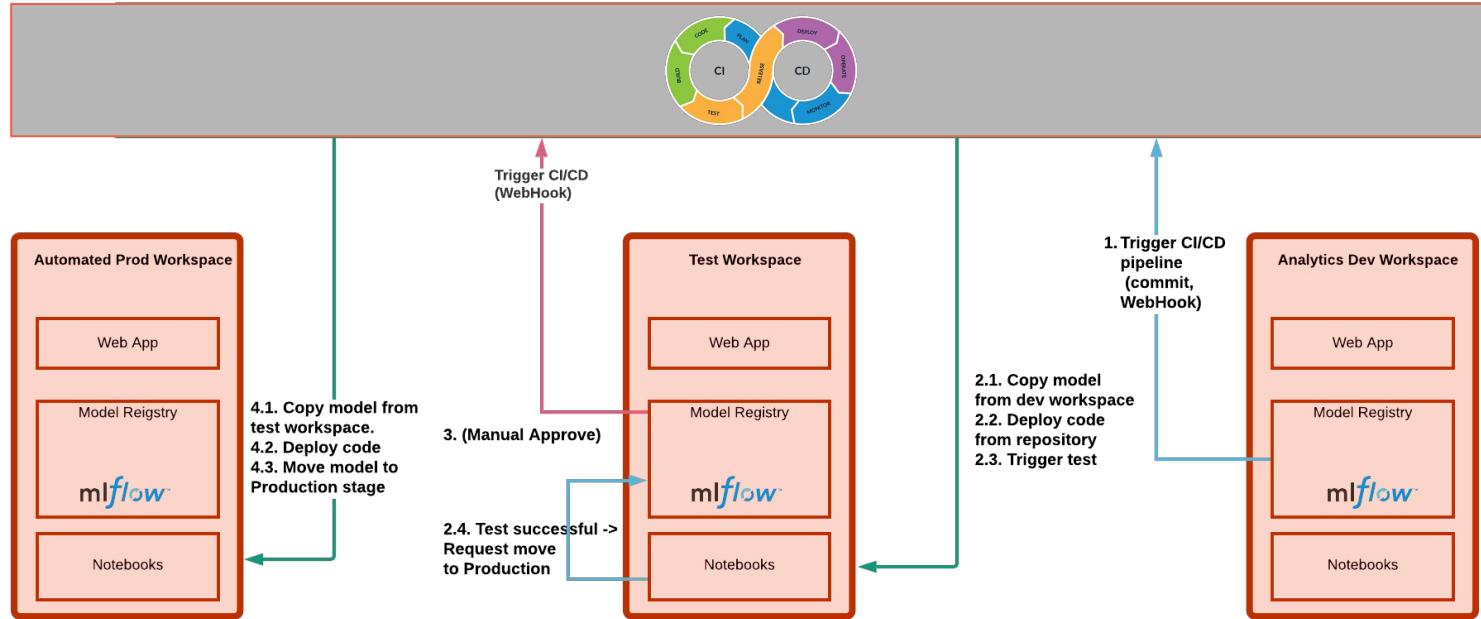
# Central Model Registry

Dedicated workspace(s) to store models and experiments

- Pros:
  - Visibility into current state of all models
  - Central management of models
  - Cross cloud
- Cons:
  - Setup is a bit complicated:
    - We need PAT & other configuration in a secret scope
    - We can have [max 100 secret scopes](#) per workspace
    - Using shared PATs isn't good from compliance point of view
  - Could be cluttered, especially if we'll use it for tracking of experiments
  - Control of permissions could be complex if we have many teams

# Per-workspace Model Registries

Each workspace has its own model registry, model transitions via CI/CD



With Model Registry – Unity Catalog Integration, you can manage model access permissions in single location and share models across workspaces.



# Per-workspace Model Registries

Each workspace has its own model registry, model transitions via CI/CD

- Pros:
  - Narrowed access to staging/production workspace
  - PAT is required only for a system account (service principals, for example)
- Cons:
  - Limited observability – you need to visit staging/production workspaces to get information about model state
  - Approvals needs to be done in relevant workspaces
  - There are no built-in tools for transitioning of models/experiments between workspaces, although there is some tooling
  - Requires implementation as part of CI/CD pipeline

# DEMO: MODEL REGISTRY

Notebook: O2-Model-Management → O2-Model-Registry



# DEMO: WEB HOOKS & TESTING

Notebook: 02-Model-Management → 03a-Webhooks-and-Testing

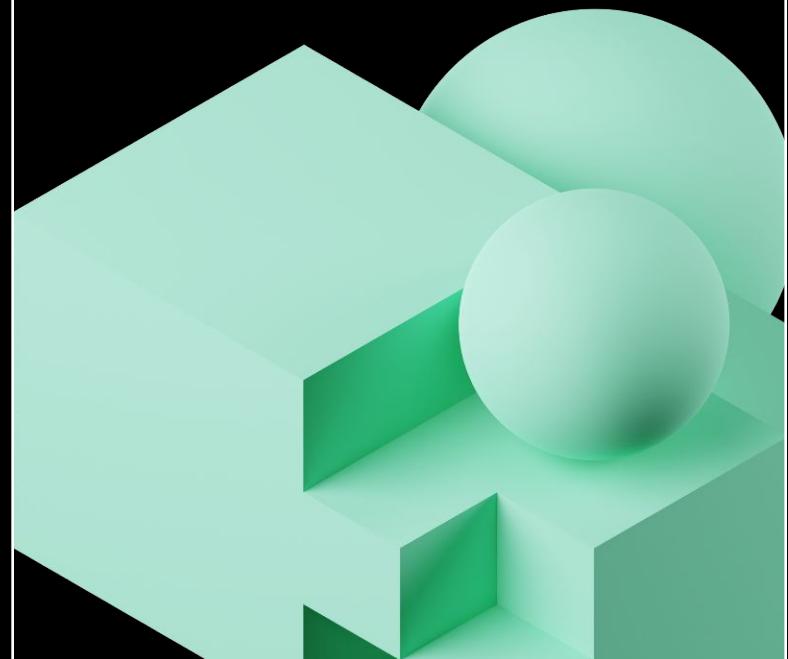


# **DEMO:** **AUTOMATED TESTING**

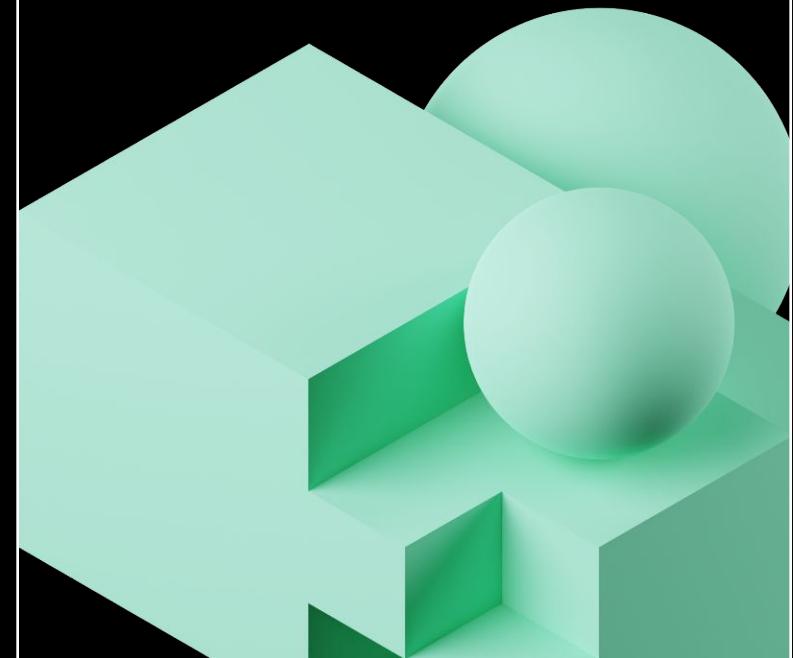
Notebook: 02-Model-Management → 03b-Webhooks-Job-Demo



# O3. Deployment Paradigms



# DevOps vs. ModelOps



# What is ML Deployment?

- Data Science != ML Engineering
- Data science is scientific
  - Business problems → data problems
  - Model mathematically
  - Optimize performance
- ML engineers are concerned with
  - Reliability
  - Scalability
  - Maintainability
  - SLAs
  - ...

# DevOps vs. ModelOps

## DevOps

**Software development + IT operations**

- Manages deployments
- CI/CD of features, patches, updates, rollbacks

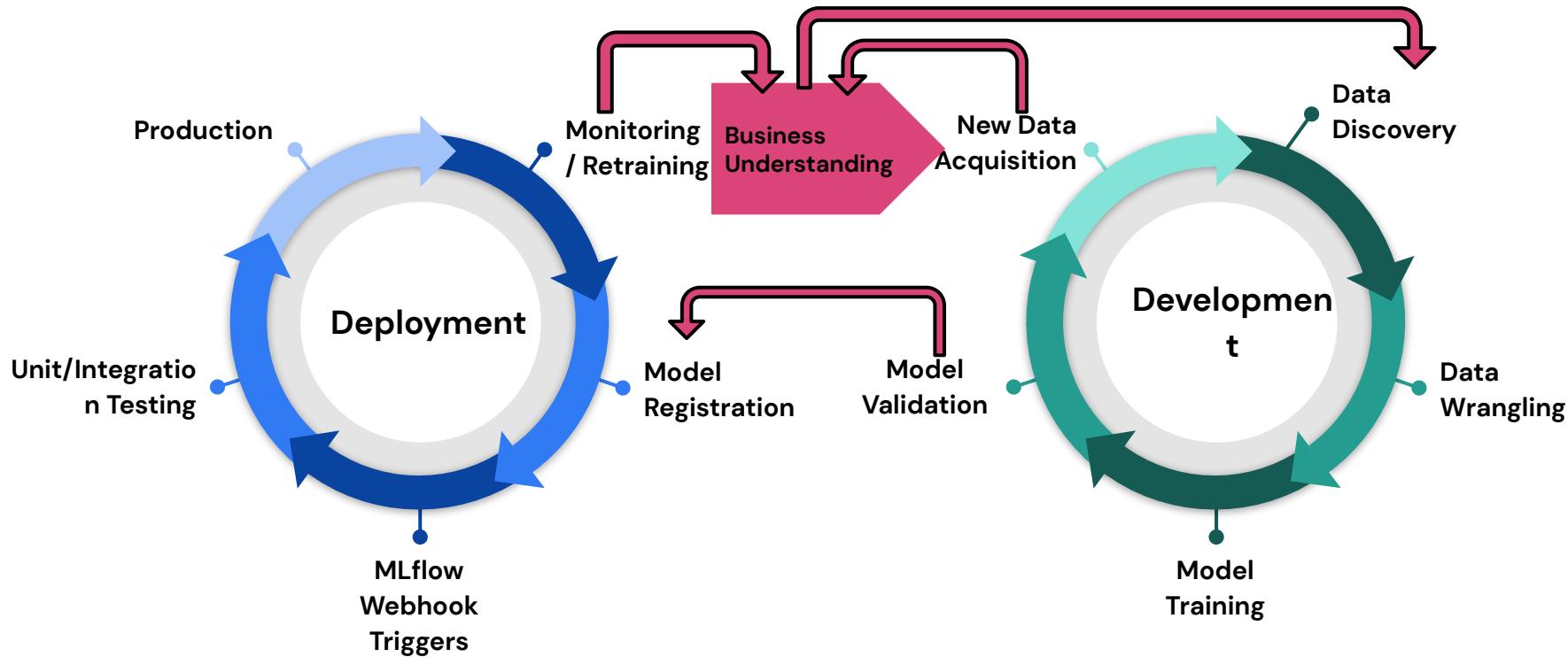
## ModelOps

**Data modeling + deployment operations**

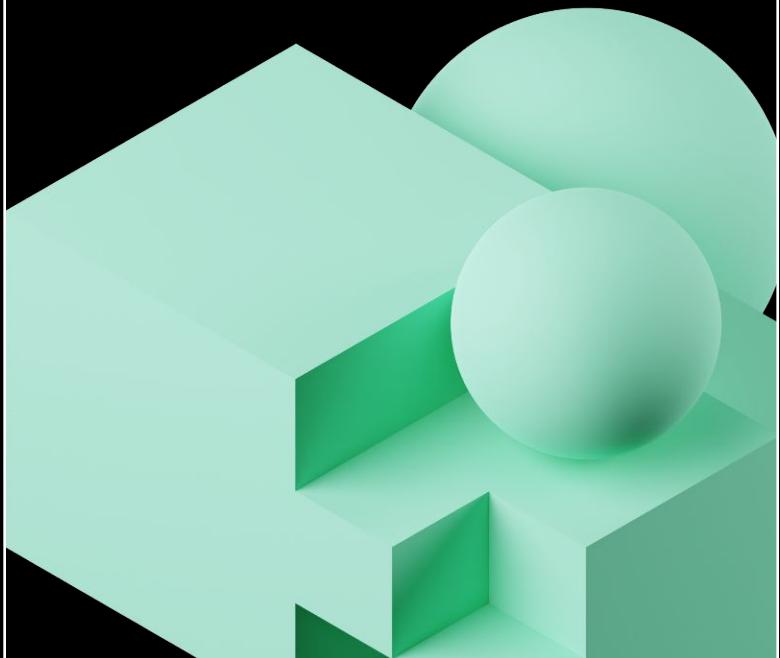
- Artifact management (Continuous Training)
- Model performance monitoring (Continuous Monitoring)
- Data management
- Use of containers and managed services



# Closed Loop Systems



# Deployment Patterns



# Semantics of Dev, Staging, and Production

ML Workflow Assets:



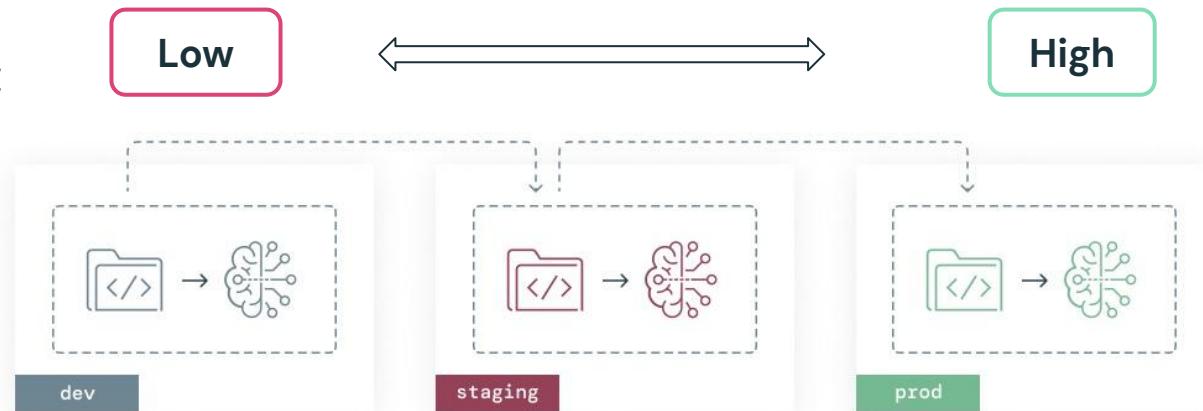
Assets need to be:



# Dev vs. Staging vs. Production

Level of trust, quality, testing and openness

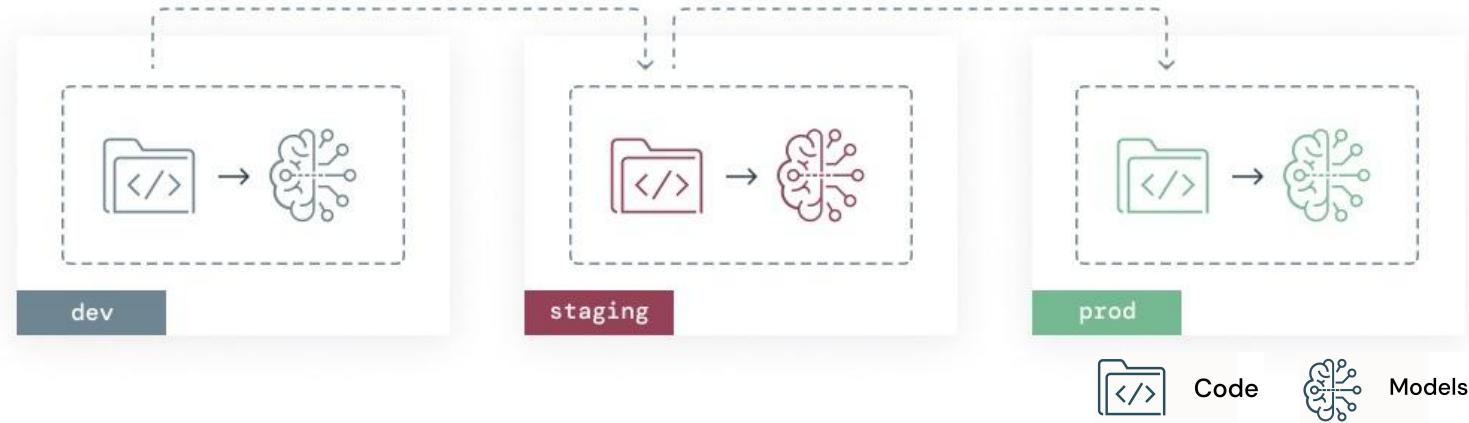
Level of trust,  
quality and testing:



Openness of  
access:



# Model vs. Code Lifecycles



**Model and code lifecycles often operate asynchronously**

- Weekly fraud detection model
  - Model update; no code change
- Computer vision model, large language model
  - No model update; code change

→ mlflow™

# Dev vs. Staging vs. Production

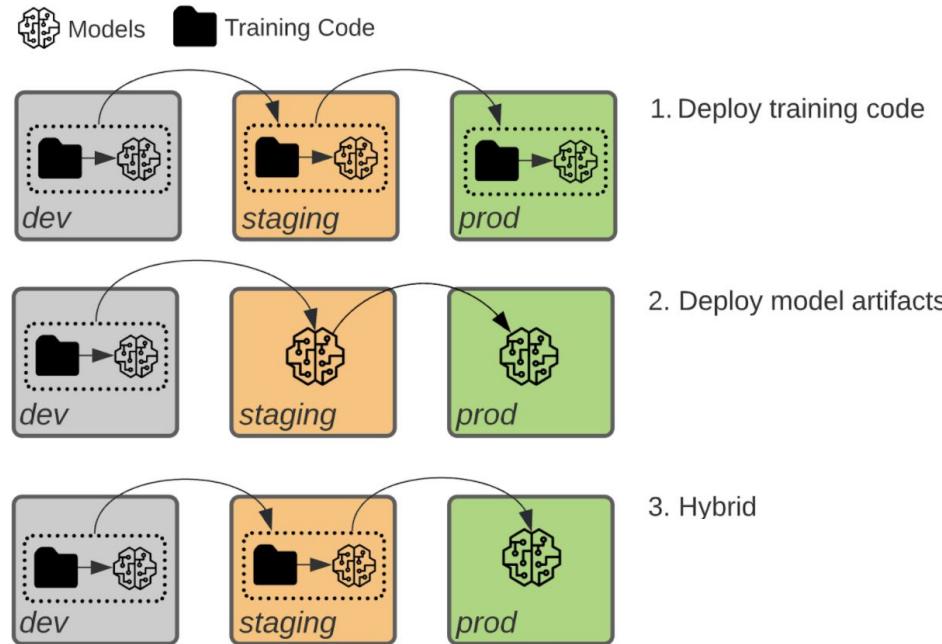
## Managing Assets

ASSET	SEMANTICS	SEPARATED BY
Execution environments	Labeled according to where development, testing and connections with production systems happen	Cloud provider and Databricks Workspace access controls
Models	Labeled according to model lifecycle phase	MLflow access controls or cloud storage permissions
Data	Labeled according to its origin in dev, staging or prod execution environments	Table access controls or cloud storage permissions
Code	Labeled according to software development lifecycle phase	Git repository branches



# Three Model Deployment Patterns

More information on [The Big Book of MLOps](#) published by Databricks



# What Moves Towards Production?

## Pattern 1: Deploy models



Deploy models	
<b>Process</b>	<ol style="list-style-type: none"><li>1. Models are trained in dev and promoted to staging</li><li>2. Models are tested in staging</li><li>3. Models are promoted to prod</li></ol>
<b>Trade-offs</b>	<ul style="list-style-type: none"><li>↑ Simplicity</li><li>↑ DS familiarity</li><li>↑ Computational cost</li> <li>↓ Automation</li><li>↓ Scalability</li><li>↓ Reproducibility</li></ul>

# What Moves Towards Production?

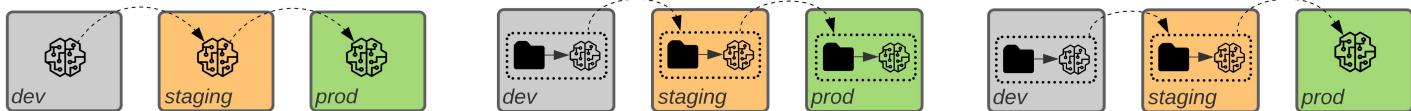
## Pattern 1: Deploy code



	Deploy models	Deploy code
Process	<ol style="list-style-type: none"><li>Models are trained in dev and promoted to staging</li><li>Models are tested in staging</li><li>Models are promoted to prod</li></ol>	<ol style="list-style-type: none"><li>Code is promoted from dev to staging</li><li>Models are retrained and tested in staging</li><li>Code is promoted to prod</li><li>Models are retrained in prod</li></ol>
Trade-offs	<ul style="list-style-type: none"><li>↑ Simplicity</li><li>↑ DS familiarity</li><li>↑ Computational cost</li><li>↓ Automation</li><li>↓ Scalability</li><li>↓ Reproducibility</li></ul>	<ul style="list-style-type: none"><li>↑ Scalability</li><li>↑ Automation</li><li>↑ Reproducibility</li><li>↑ ML Eng familiarity</li><li>↓ Simplicity</li><li>↓ DS familiarity</li><li>↓ Computational cost</li></ul>

# What Moves Towards Production?

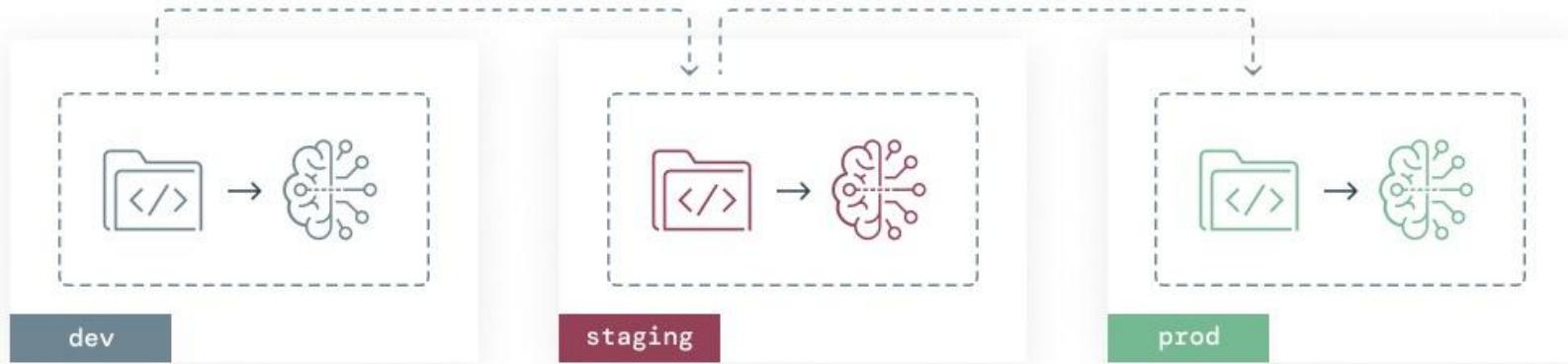
## Pattern 1: Hybrid



	Deploy models	Deploy code	Hybrid
Process	<ol style="list-style-type: none"><li>Models are trained in dev and promoted to staging</li><li>Models are tested in staging</li><li>Models are promoted to prod</li></ol>	<ol style="list-style-type: none"><li>Code is promoted from dev to staging</li><li>Models are retrained and tested in staging</li><li>Code is promoted to prod</li><li>Models are retrained in prod</li></ol>	<ol style="list-style-type: none"><li>Code is promoted from dev to staging</li><li>Models are retrained and tested in staging</li><li>Models are promoted to prod</li></ol>
Trade-offs	<p>↑ Simplicity ↑ DS familiarity ↑ Computational cost</p> <p>↓ Automation ↓ Scalability ↓ Reproducibility</p>	<p>↑ Scalability ↑ Automation ↑ Reproducibility ↑ ML Eng familiarity</p> <p>↓ Simplicity ↓ DS familiarity ↓ Computational cost</p>	<p>↑ Scalability ↑ Automation ↑ Reproducibility ↑ ML Eng familiarity ↑ Computational cost</p> <p>↓ Simplicity ↓ DS familiarity</p>



# Recommended Deploy Code Pattern



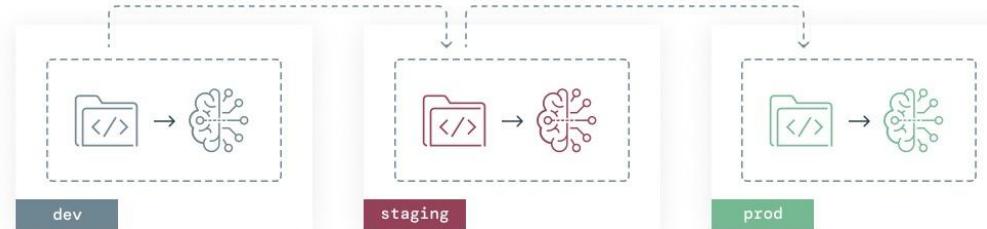
Develop training code  
Develop ancillary code  
→ Promote code

✓ Test model training code  
on subset of data  
✓ Test ancillary code  
→ Promote code

✓ Train model on production  
data  
✓ Test model  
→ Deploy model  
→ Deploy ancillary code

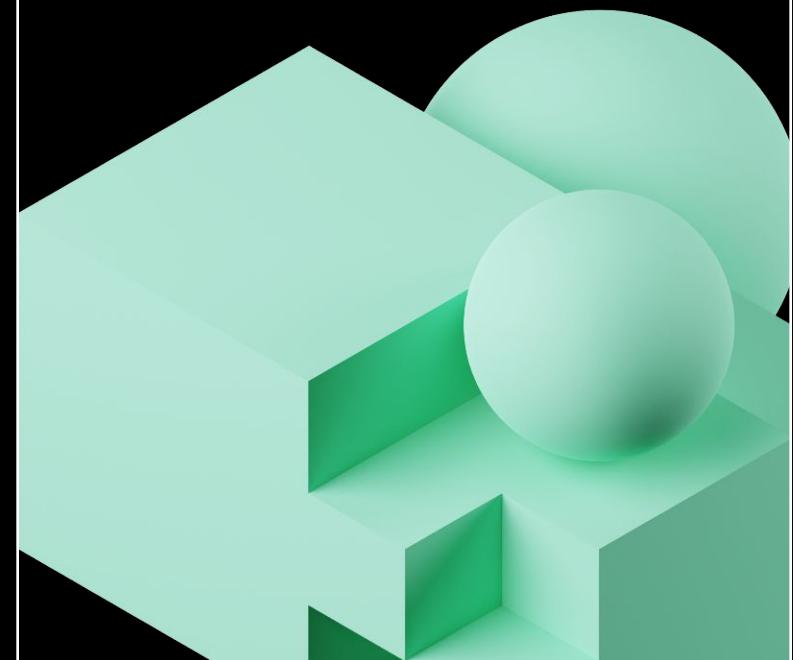


# Deeper Dive into Benefits of “Deploy Code”



<b>Automation</b>	↑ Supports automated retraining in locked-down env.
<b>Data access control</b>	↑ Only prod env needs read access to prod training data.
<b>Reproducible models</b>	↑ Eng control over training env, which helps to simplify reproducibility.
<b>Support for large projects</b>	↑ This pattern forces the DS team to use modular code and iterative testing, which helps with coordination and development in larger projects.
<b>Data science familiarity</b>	↓ DS team must learn to write & hand off modular code to Eng.
<b>Eng setup &amp; maintenance</b>	↓ Requires CI/CD infra for unit and integration tests, even for one-off models.

# Deployment Paradigms



# The Four Deployment Paradigms

## 1. Batch

- 80–90% of deployments
- Leverages databases and object storage
- Fast retrieval of stored predictions

## 2. Streaming (continuous)

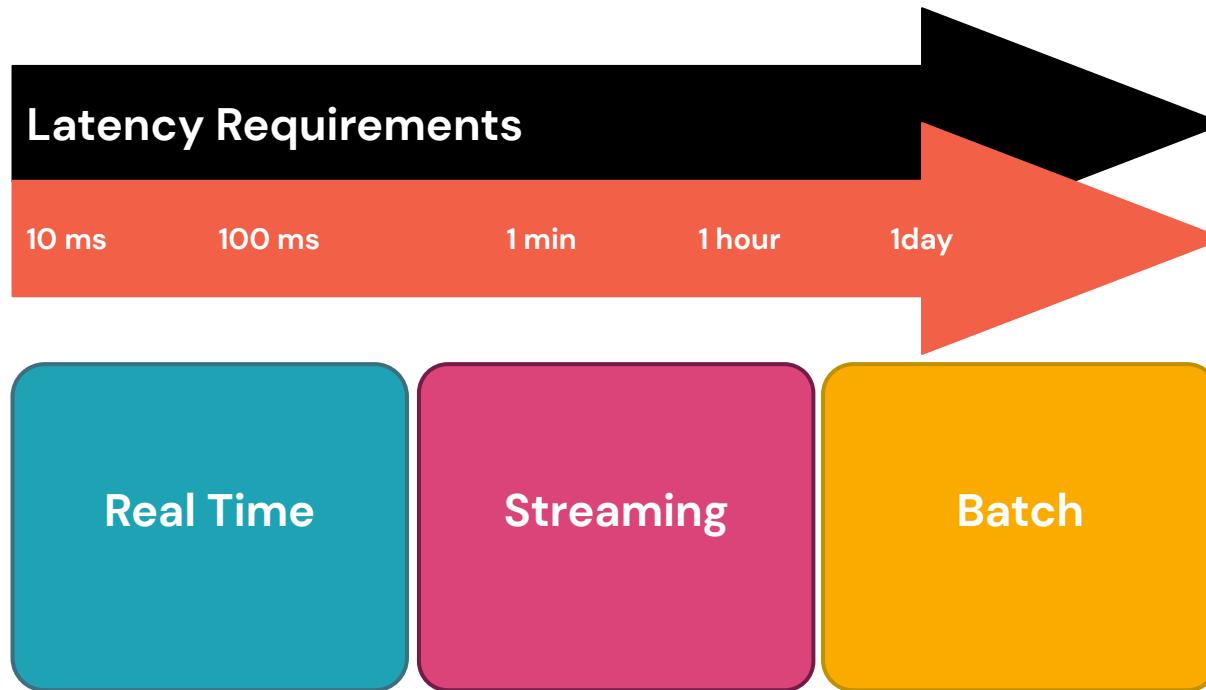
- 10–15% of deployments
- Moderately fast scoring on new data

## 3. Real Time

- 5–10% of deployments
- Usually using REST (Azure ML, SageMaker, containers)

## 4. On-device (edge)

# Latency Requirements (roughly)



# mlflow Model Deployment Options



In-Line Code



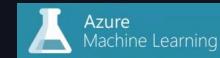
Containers



Batch & Stream  
Scoring



OSS Inference  
Solutions



Cloud  
Inference  
Services

# Databricks Model Serving

## Enable Model Serving with 1 Click

Registered Models > wine-quality-serving-demo

wine-quality-serving-demo



Permissions

Use model for inference

Details Serving Preview

Enable serverless model endpoints behind a REST API interface. This will launch endpoints for all active staging and production versions of this model. You can still use [Classic model serving](#). Click [this link](#) to enable it.

Enable Serverless Model Endpoints

Use model for inference

One click model deployment directly from the Model Registry



# Online/Offline Model Serving with Databricks

Deploy any ML model at large scale AND low latency

## Large-scale batch scoring

Set up model inference [Preview](#) [Provide feedback](#) X

Select one of real-time inference, streaming via Delta Live Tables, or batch.

Real-time Streaming (Delta Live Tables) **Batch inference**

Generates a notebook in your home folder that you can edit.

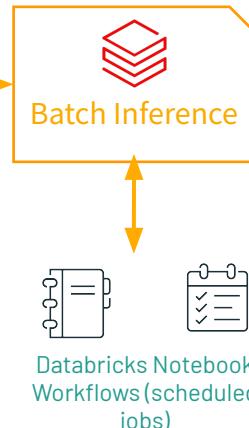
\* Model version  
Model version

\* Input table  
Input table  Browse

\* Output table location  
/FileStore/batch-inference/ Model\_Serving\_Demo\_Enhanced\_Mod

The default output path on DBFS is accessible to everyone in this Workspace. Modify the notebook to disable writing data to DBFS.

Cancel



## Low-latency online serving

Set up model inference [Preview](#) [Provide feedback](#) X

Select one of real-time inference, streaming via Delta Live Tables, or batch.

Real-time Streaming (Delta Live Tables) **Batch inference**

Serve this model behind an endpoint. To configure multiple served models behind this endpoint, update this endpoint's configuration via the endpoint details page. [Learn more](#)

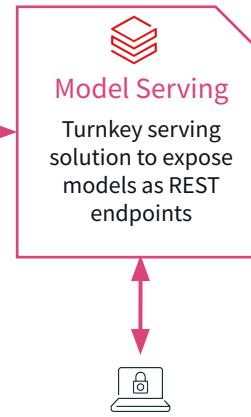
\* Model  
Model\_Serving\_Demo\_Enhanced\_Model\_menaf\_gul\_8vg5\_da\_coms

\* Model version  
Model version

\* Endpoint name  
Endpoint name

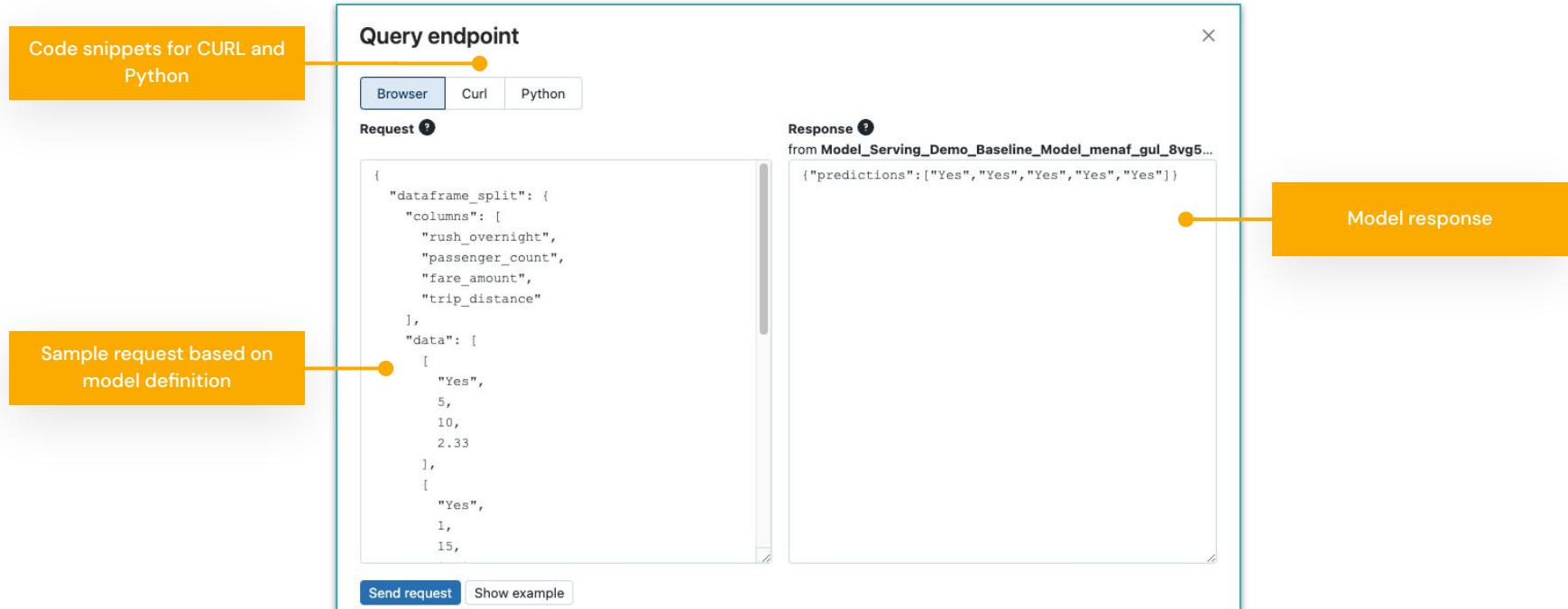
Compute  
Small 0-4 concurrent requests (0-4 BU)  Scale to zero

Cancel



# Validate Model Serving with Input Examples

In-browser testing with auto-populated example



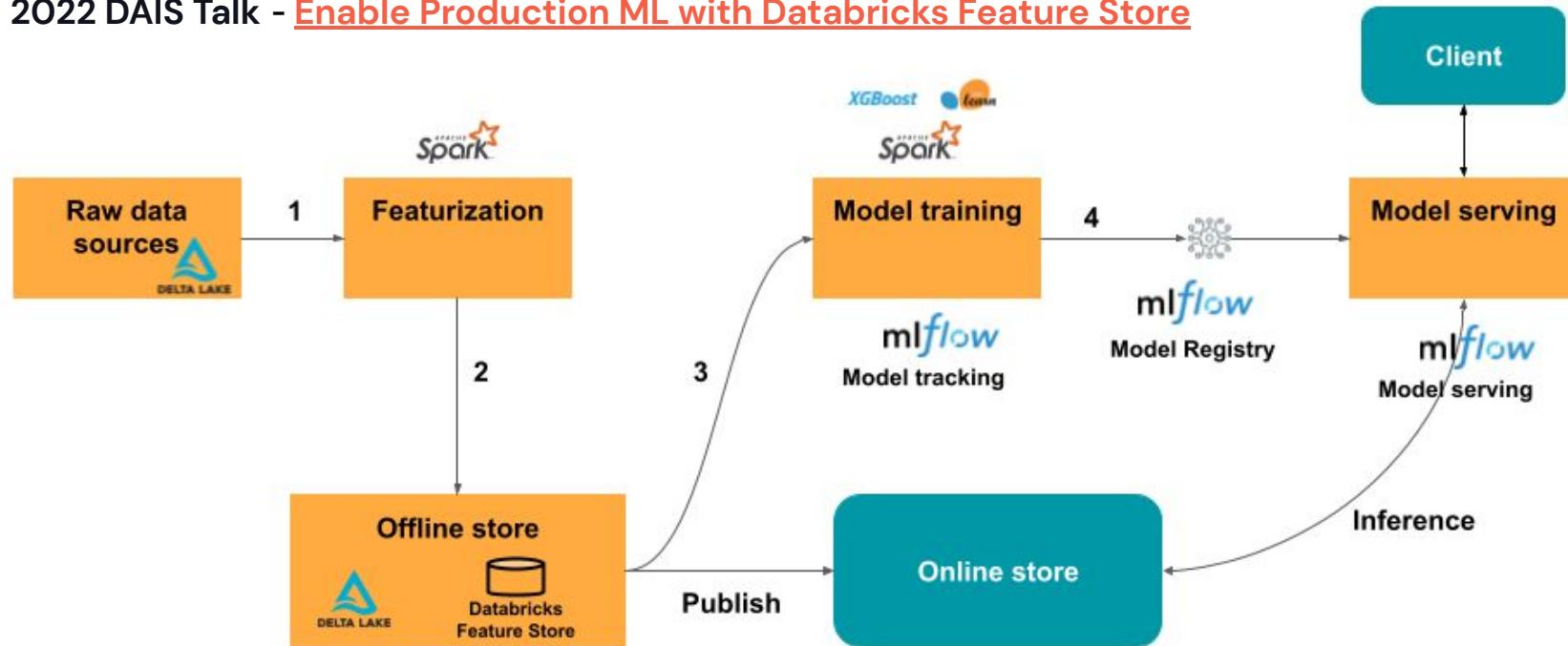
# Serving Models on Databricks

Serverless Model Endpoint – Preview in AWS & Azure

- Production-grade serving
- Low latency, high availability and scalable
- HTTP endpoints accessible with Databricks authentication
- Scoring validation from the UI (same UI across clouds)
- 1-click enable serving
- Autoscaling clusters
- Endpoint latency & QPS monitoring

# Serve Feature Store Models Online

2022 DAIS Talk - [Enable Production ML with Databricks Feature Store](#)



# DEMO: BATCH INFERENCE

Notebook: 03-Deployment-Paradigms → 01-Batch



# DEMO: REAL-TIME INFERENCE

Notebook: O3-Deployment-Paradigms → O2-Real-Time



# LAB: BATCH INFERENCE

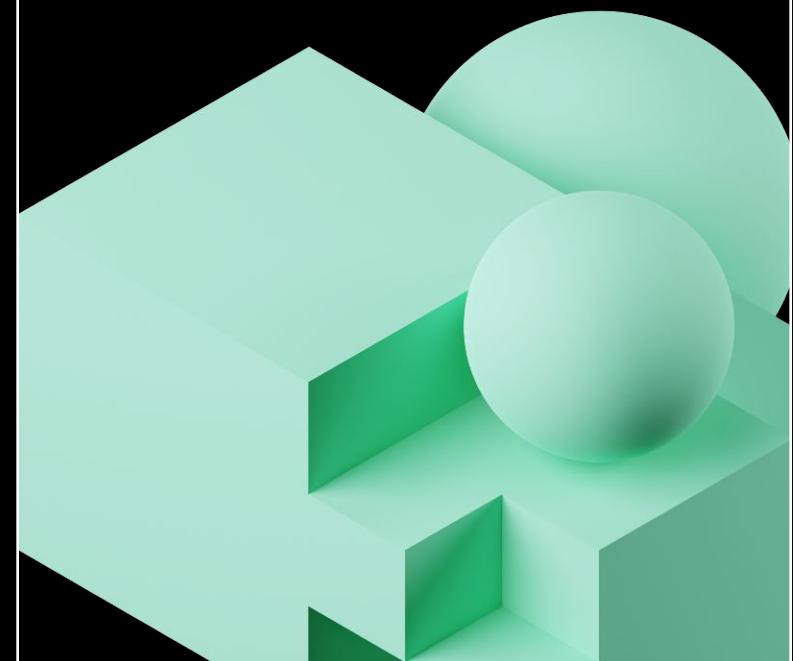
Notebook: O3-Deployment-Paradigms → Labs → O1-Batch-Lab



# LAB: REAL-TIME INFERENCE

Notebook: 03-Deployment-Paradigms → Labs → 02-Real-Time-Lab

# 04. Production



# Production Requirements

## Core Requirements

- Model registry
- Data and model drift
- Interpretability
- Reproducibility
- Security
- Environment management

## Core +

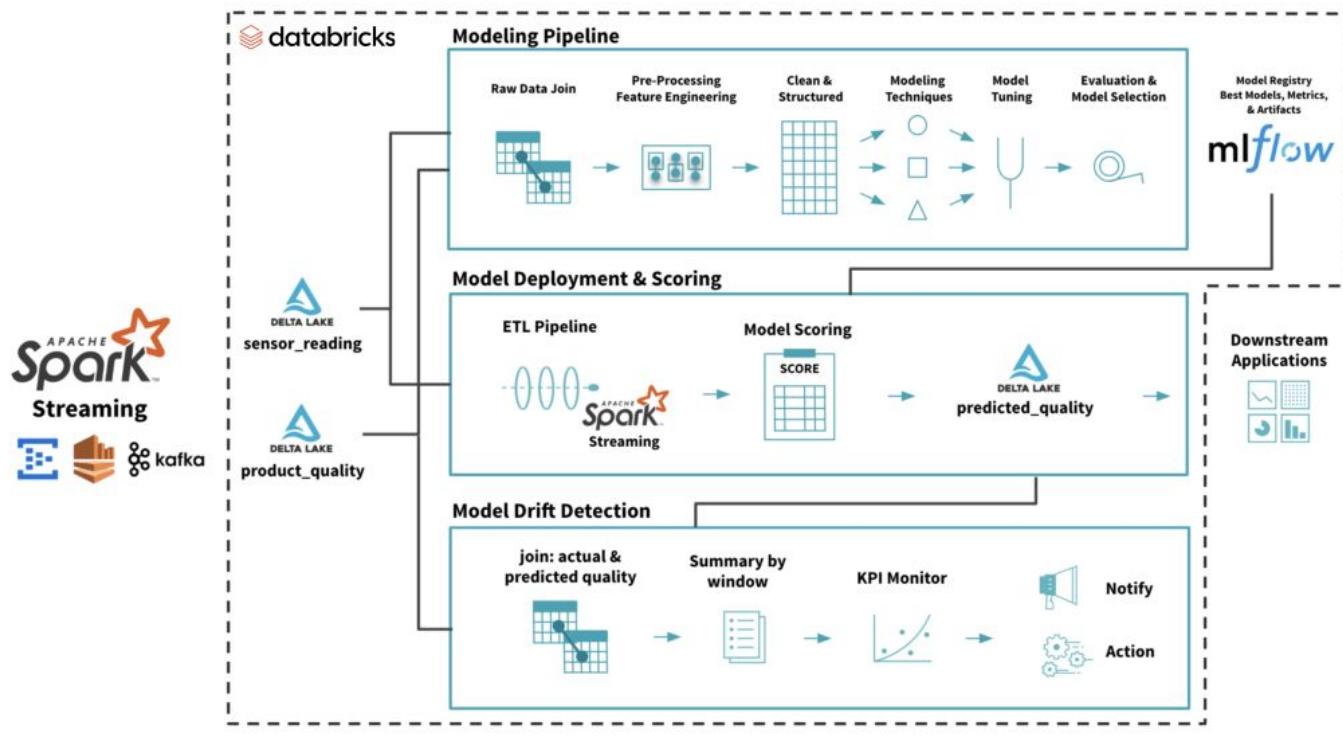
- ML pipeline with featurization logic
- CI/CD pipeline for automation
- Monitoring and alerting
- Testing framework
- Version control

## Specialized

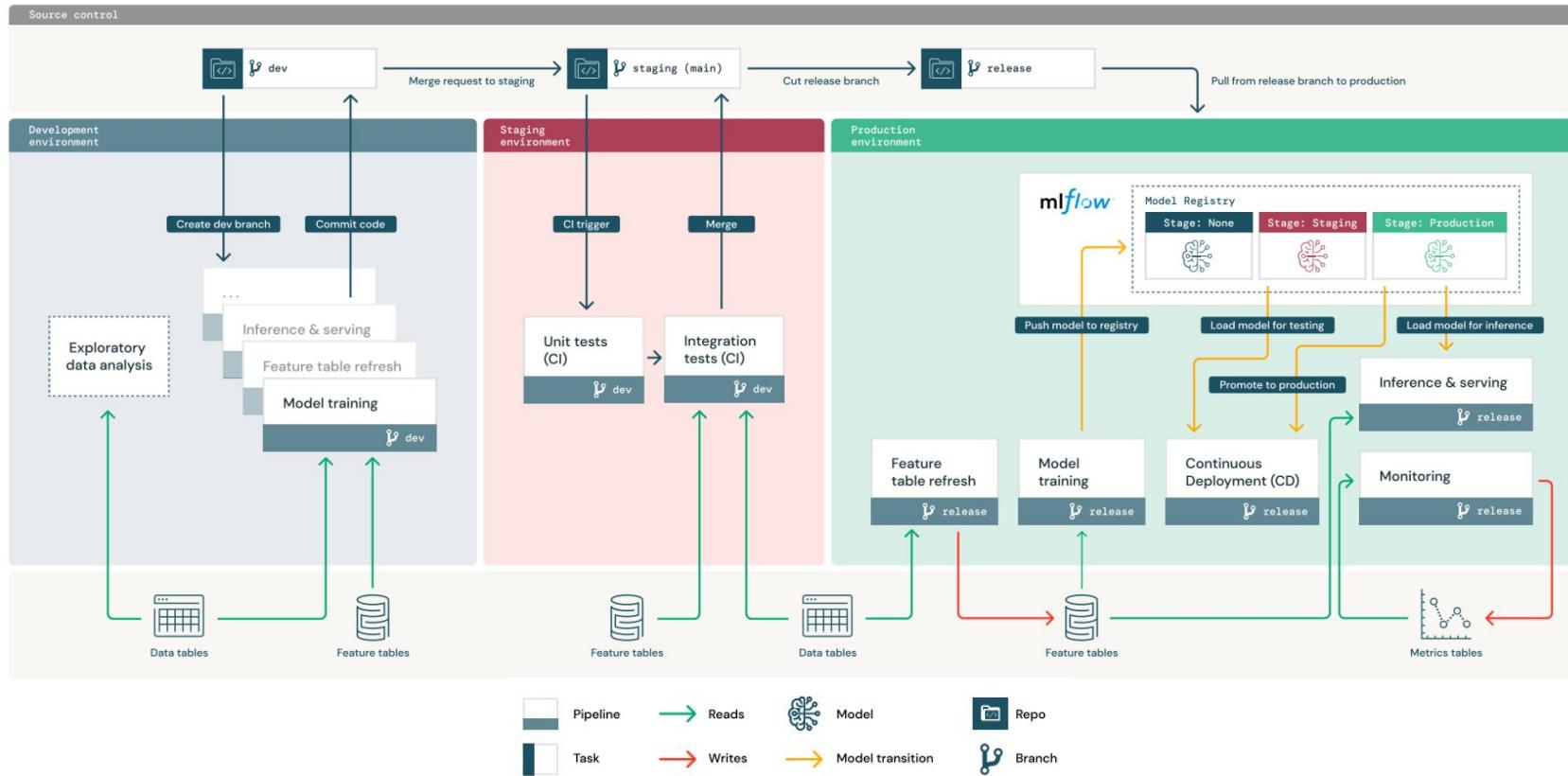
- Feature dictionary
- Cost management
- A/B testing
- Performance optimization



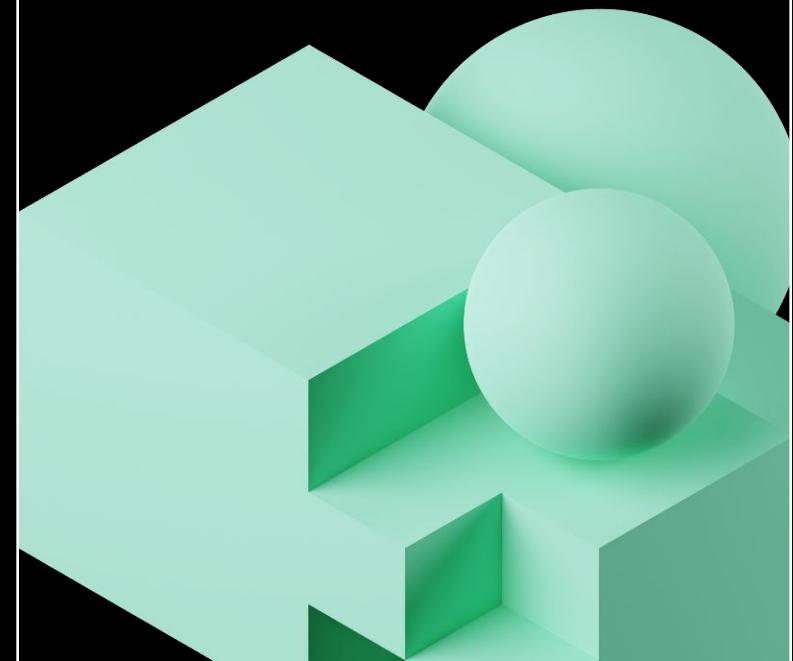
# Architecture Example



# Architecture Perspective

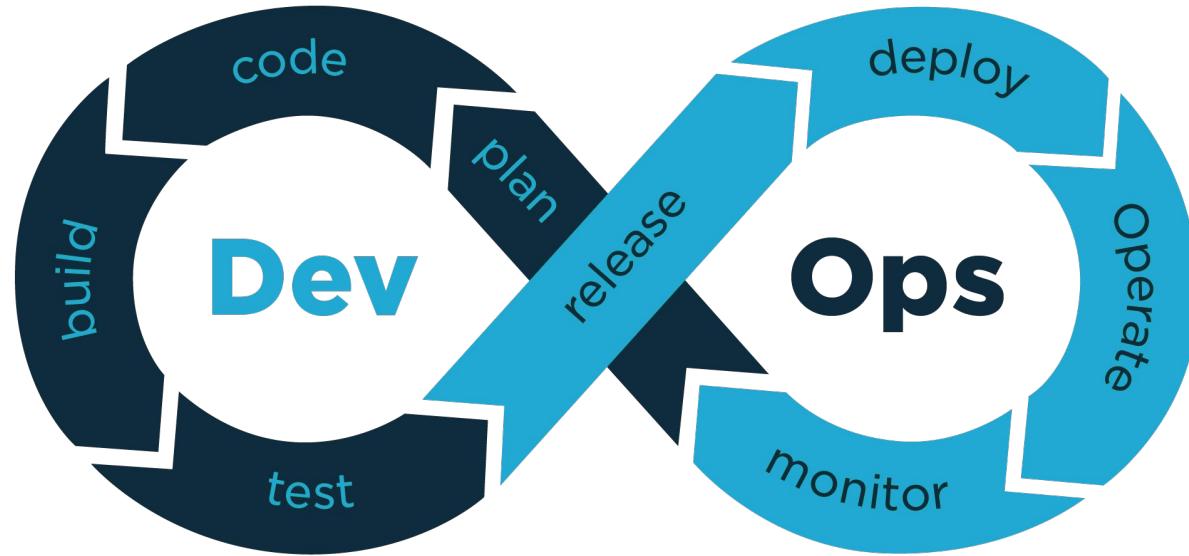


# CI/CD

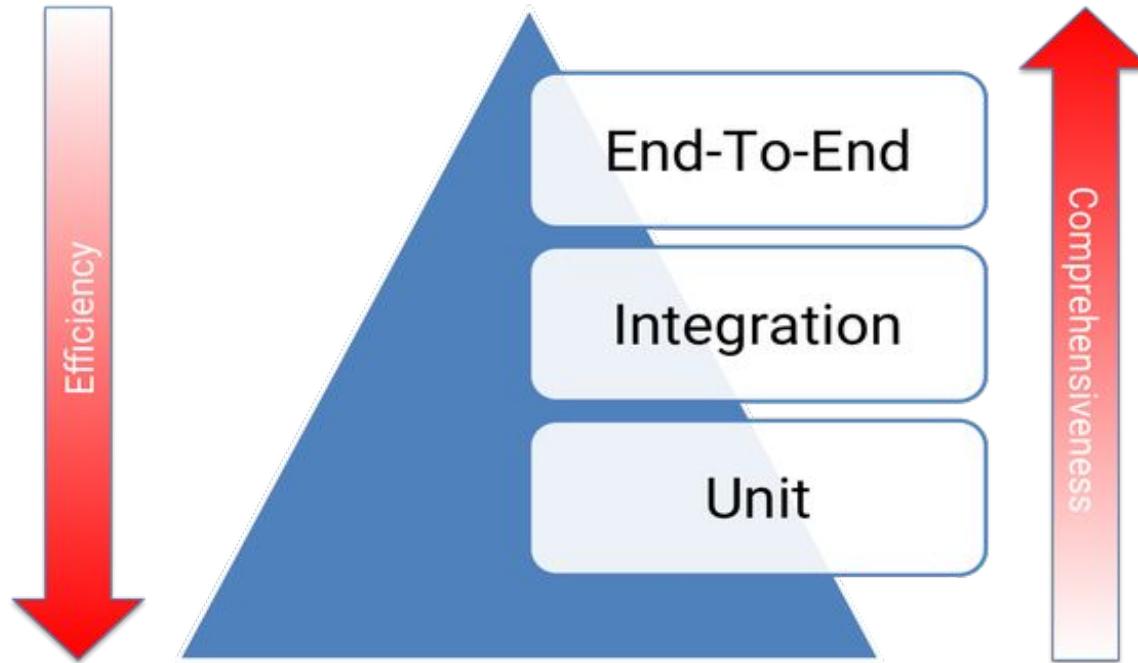


# DevOps...

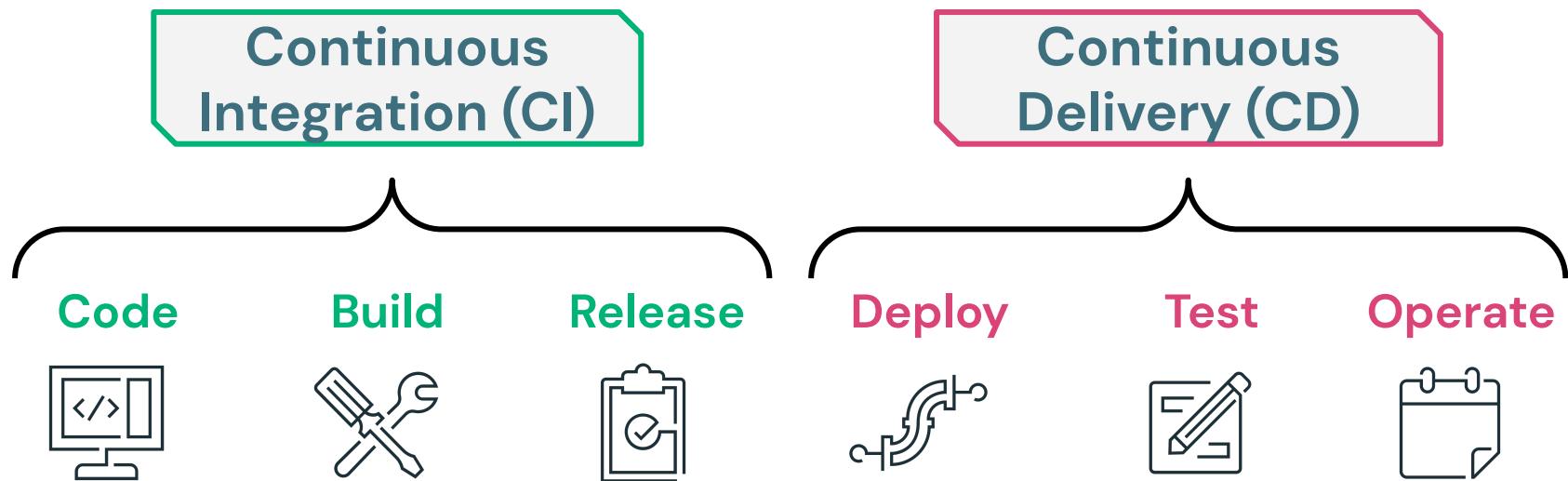
A process for continuously integrating, testing and deploying your code



# The Test Pyramid



# A Typical Databricks CI/CD Pipeline



# Types of Pipelines

- Scheduled by a workflow manager
- Types of pipelines
  - **ETL:** heavy lifting data transformation
  - **Featurization:** Generate features for model development
  - **Training:** Produce models for deployment
  - **Scoring:** Inference pipeline post deployment
  - **Monitoring:** recurring jobs post deployment with alerts set up

# Key Benefits of Continuous Integration

- Regressions are captured early by the automated tests
- Building the release is easy as all integration issues have been solved early
- Your QA team spends less time testing and can focus on significant improvements to the quality culture

# Key Benefits of Continuous Delivery

- The complexity of deploying software has been taken away
- You can release more often, thus accelerating the feedback loop with your customers
- There is much less pressure on decisions for small changes, hence encouraging iterating faster

# CI/CD Terminology

- **Build pipeline** – set of steps that build, test (usually only unit tests) & package code or notebooks
- **Release pipeline** – set of steps that performing deployment of the assets
- **Trigger** – event or action that initiates the execution of pipeline. It could be a commit to repository, or explicit execution of the pipeline
- **Asset / Artifact** – deliverable produced by build pipeline, or some other process (model training)
- **Environment** – typically:
  - Development – where development happens
  - Staging – for integration & E2E testing
  - Production – actual work

# CI/CD Technologies

	OSS Standard	Databricks	AWS	Azure	Third Party
<b>Orchestration</b>	Airflow, Jenkins	Jobs, notebook workflows	CodePipeline, CodeBuild, CodeDeploy	DevOps, Data Factor	
<b>Git Hooks</b>		MLflow Webhooks			Github Actions, Gitlab, Travis CI
<b>Artifact Management</b>	PyPi, Maven	MLflow Model Registry			Nexus
<b>Environment Management</b>	Docker, Kubernetes, Conda, pyenv		Elastic Container Repository	Container Registry	DockerHub
<b>Testing</b>	pytest				Sonar
<b>Alerting</b>		Jobs	CloudWatch	Monitor	PagerDuty, Slack integrations



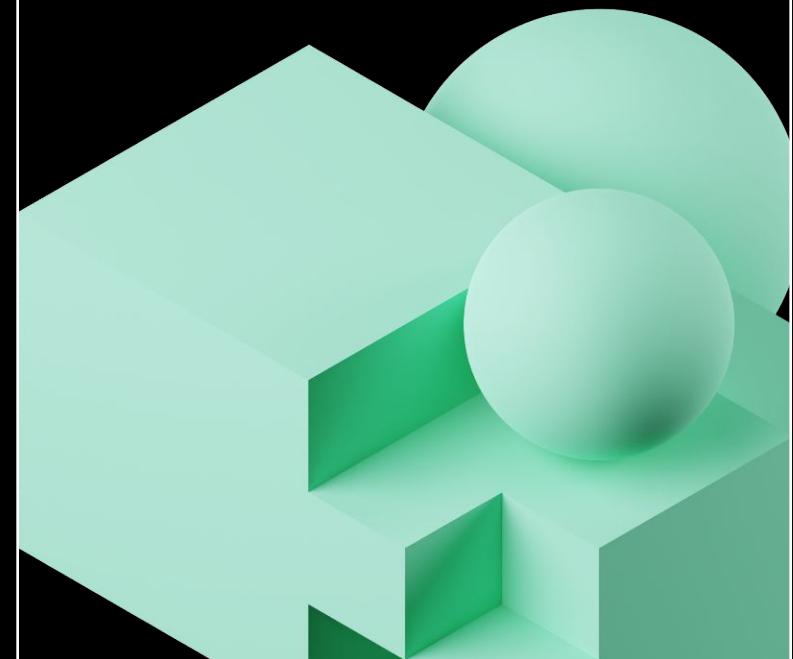
# Databricks Resources

- Implementing CI/CD on Databricks Using Databricks Notebooks and Azure DevOps [Part 1](#), [Part 2](#) (blog)
- [Continuous Integration and Delivery on Databricks using Jenkins](#) (blog)
- [Automate Continuous Integration and Continuous Delivery on Databricks using Databricks Labs CI/CD Templates](#) (blog)

# Recommended Books

- Machine Learning Engineering in Action by Ben Wilson
- Designing Data-Intensive Applications by Martin Kleppmann
- Clean Code: A Handbook of Agile Software Craftsmanship by Robert C. Martin
- Refactoring: Improving the Design of Existing Code by Martin Fowler, Kent Beck, Don Roberts
- Test-Driven Development: By Example by Kent Beck
- Code Complete, 2nd ed. by Steve McConnell
- The Pragmatic Programmer: From Journeyman to Master by Andy Hunt, Dave Thomas

# Monitoring



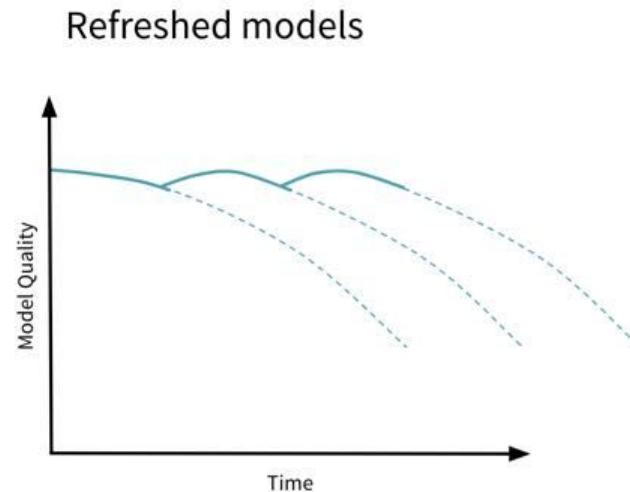
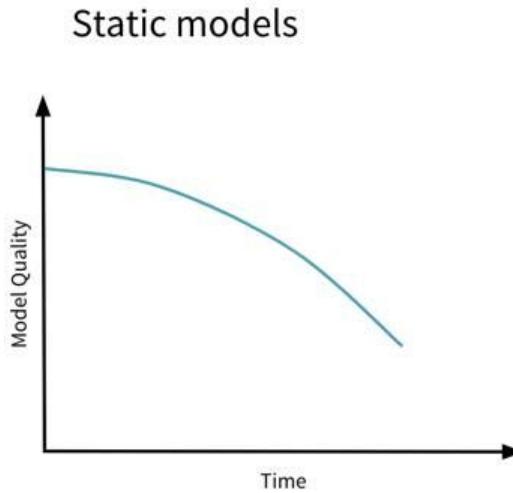
# Big Data & AI Present New Challenges...

3 changeable ingredients go into training an ML model...

Ingredient	Frequency of change	Well defined process/tools for managing change?
 Code	 ~Daily?	Yes!
 Configuration	 ~Daily/Weekly?	Emerging
 Data	 ~Every second?	No!

# Why Do ML Projects Fail in Production?

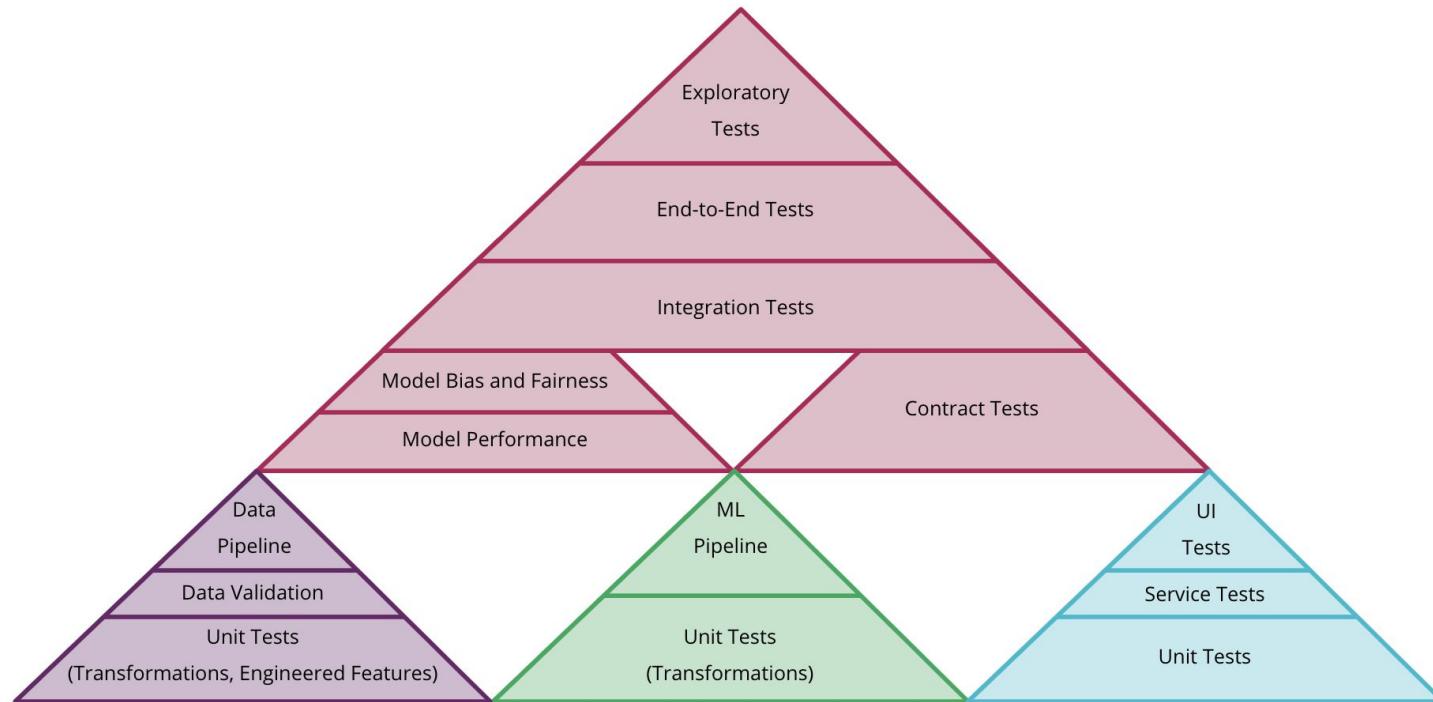
Neglect maintenance: Lack of re-training and testing



[Source](#)



# This Leads One Complicated Test Pyramid...



# MLOps: Continuous ...

CT and CM are unique to ML systems

- Continuous Integration (CI) extends the testing and validating code and components by adding testing and validating data and models
- Continuous Delivery (CD) concerns with delivery of an ML training pipeline that automatically deploys another the ML model prediction service.
- **Continuous Training (CT)** automatically retrains ML models for redeployment
- **Continuous Monitoring (CM)** concerns with monitoring production data and model performance metrics, which are bound to business metrics

# Types of Drift

Four main reasons of model degradation

## Feature Drift

Input feature(s)  
distributions deviate

## Label Drift

Label distribution  
deviates

## Prediction Drift

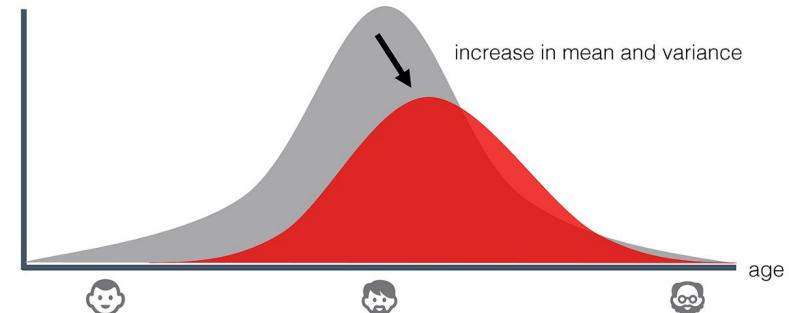
Model prediction  
distribution deviates

## Concept Drift

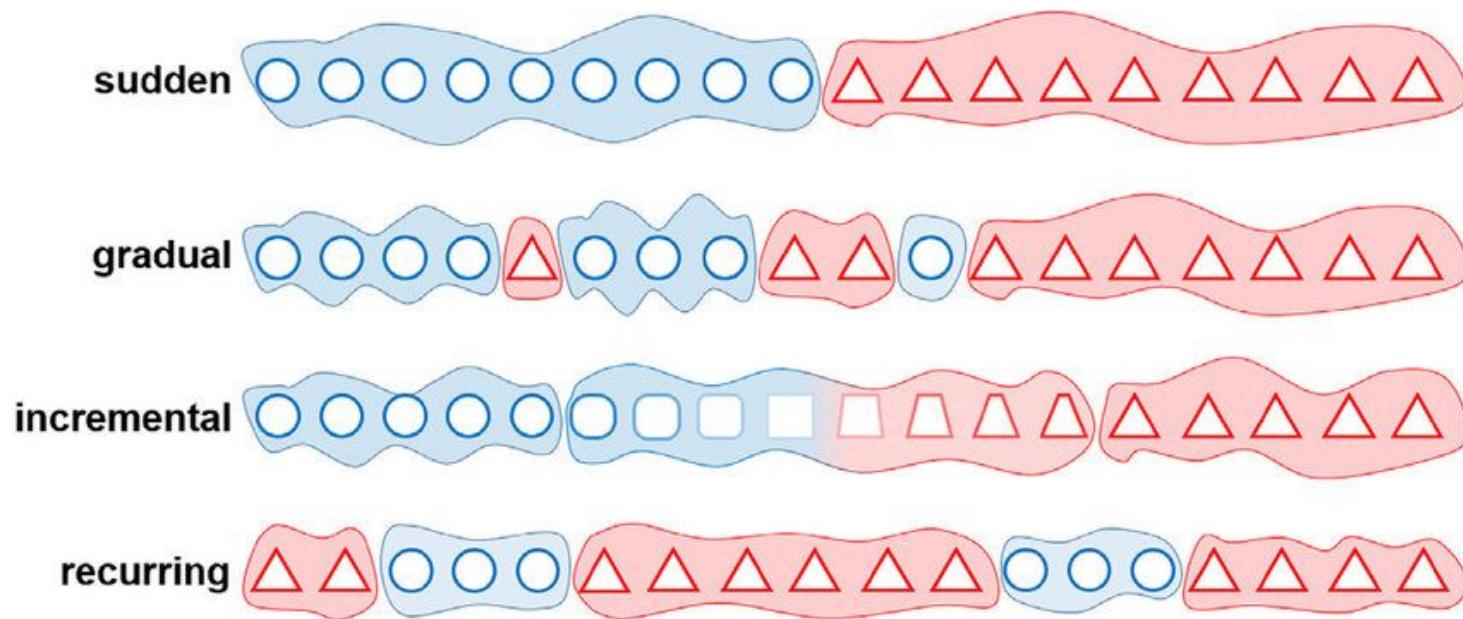
External factors  
cause the label to  
evolve

# Feature, Label, and Prediction Drift

Categories	Expected	Observed	Total
A	25	35	60
B	25	20	45
C	25	25	50
D	25	20	45
Total	100	100	200



# Concept Drift



# Drift Types and Actions to Take

Drift Type Identified	Action
Feature Drift	<ul style="list-style-type: none"><li>• Investigate feature generation process</li><li>• Retrain using new data</li></ul>
Label Drift	<ul style="list-style-type: none"><li>• Investigate label generation process</li><li>• Retrain using new data</li></ul>
Prediction Drift	<ul style="list-style-type: none"><li>• Investigate model training process</li><li>• Assess business impact of change in predictions</li></ul>
Concept Drift	<ul style="list-style-type: none"><li>• Investigate additional feature engineering</li><li>• Consider alternative approach/solution</li><li>• Retrain/tune using new data</li></ul>

# What to Monitor?

## Software Metrics

- Memory
- Compute
- Latency
- Throughput
- Server load

## Input Metrics

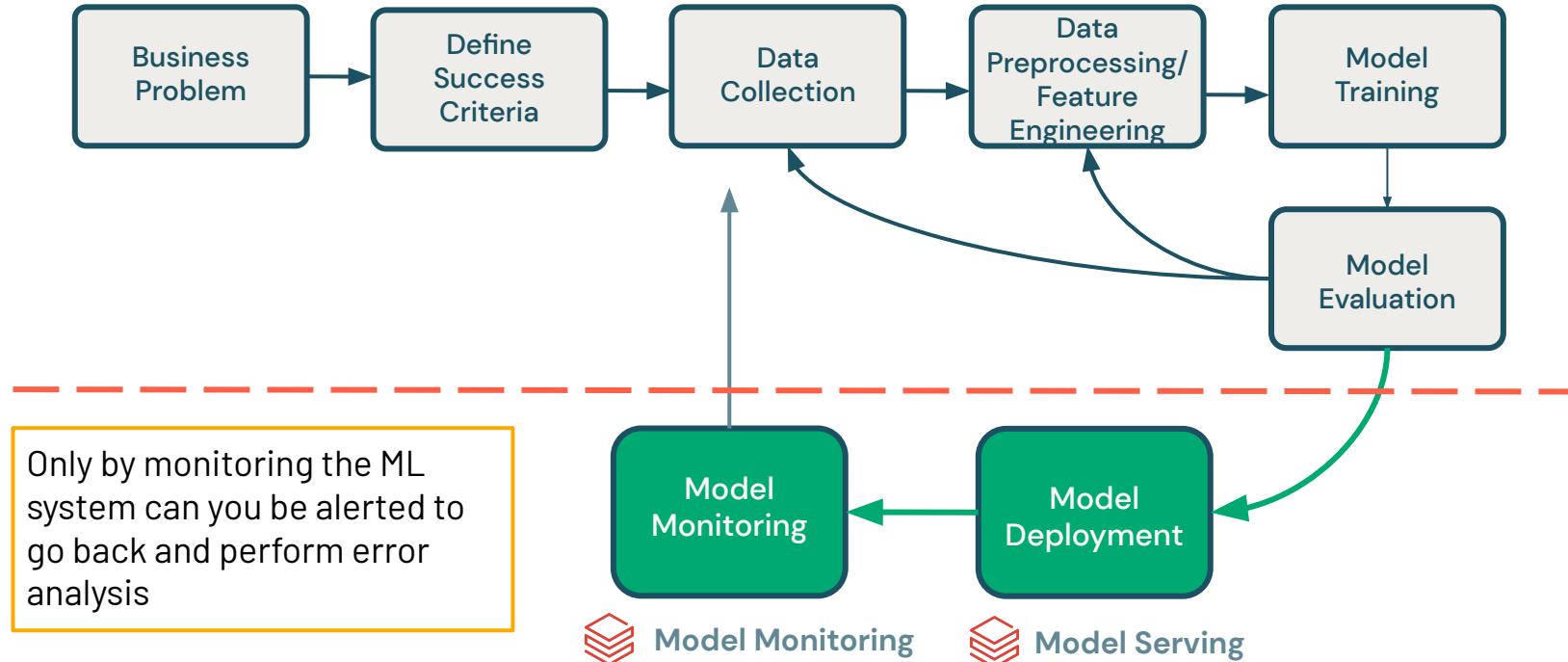
- Feature summary stats
  - Mean/Median/Min/Max
  - Num missing vals
- Statistical tests
  - KS test/Mann Whitney
  - Jensen Shannon Distance
  - Levene

## Output Metrics

- Prediction summary stats
  - Mean/Med/Min/Max
- Business metrics



# How to Action Monitoring Metrics



Model Monitoring



Model Serving

# DEMO: **MODEL MONITORING**

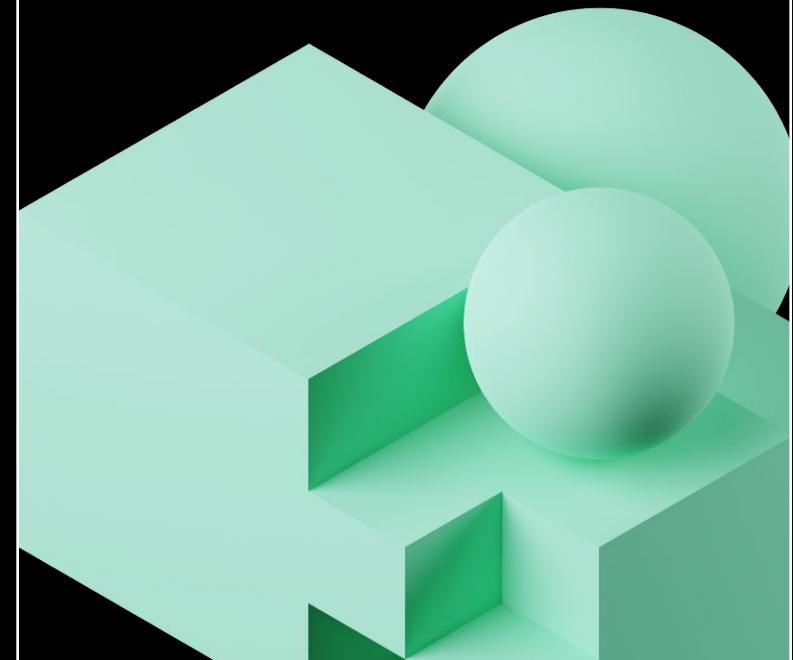
Notebook: 04-Production → 01-Monitoring



# LAB: MODEL MONITORING

Notebook: O4-Production → Labs → 01-Monitoring-Lab

# Model Rollout Strategies



# Model Rollout Strategies

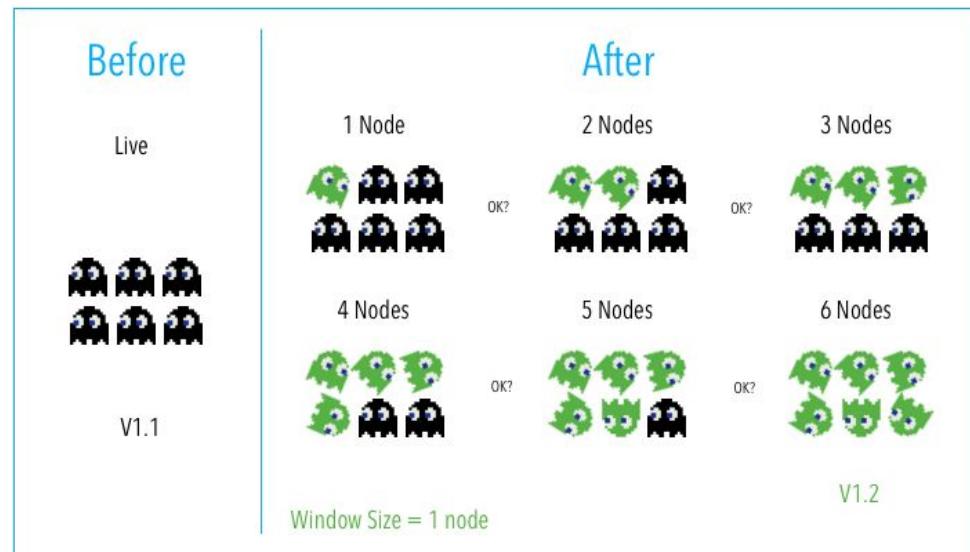
- **Shadow**
  - New deployment shows existing system but not used for decision making
- Rolling
- Blue-Green
- Canary
- A/B Testing

[Source](#)



# Model Rollout Strategies

- Shadow
- **Rolling**
  - Updates nodes in a target environment incrementally in batches with the new service version.
- Blue-Green
- Canary
- A/B Testing



[Source](#)



# Model Rollout Strategies

- Shadow
- Rolling
- **Blue-Green**
  - Utilizes two identical environments, a “blue” and a “green” environment with different versions of an application or service
- Canary
- A/B Testing

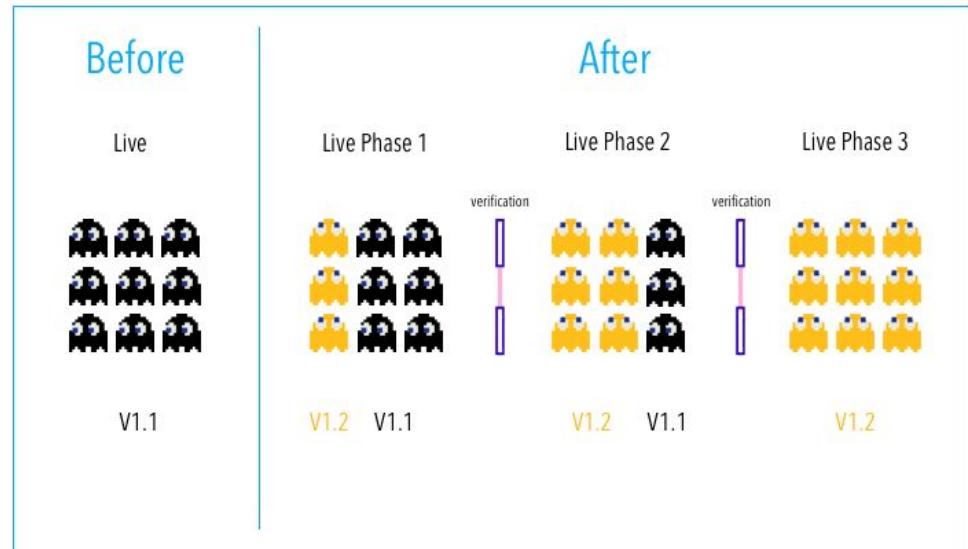


[Source](#)



# Model Rollout Strategies

- Shadow
- Rolling
- Blue-Green
- **Canary**
  - Releases an application or service incrementally to a subset of users
- A/B Testing

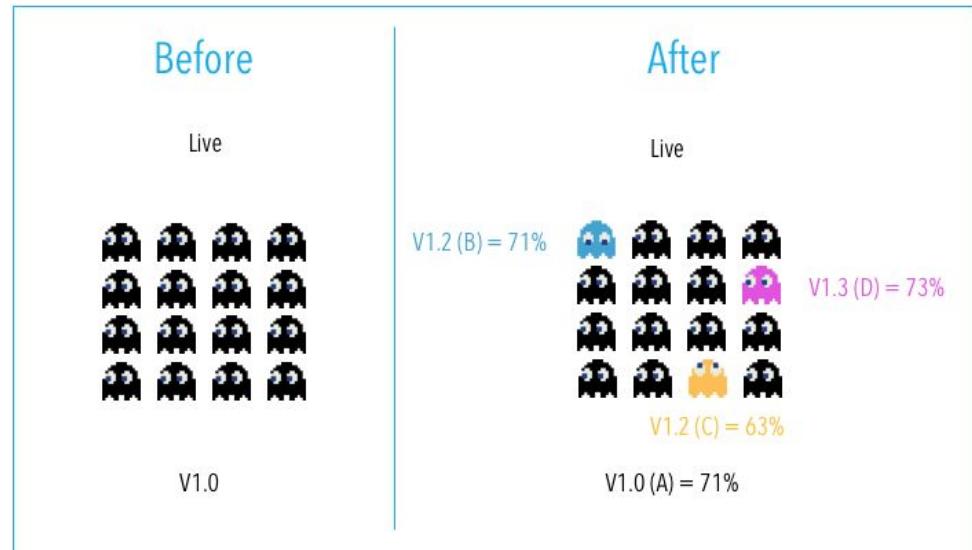


[Source](#)



# Model Rollout Strategies

- Shadow
- Rolling
- Blue-Green
- Canary
- **A/B Testing**
  - Different versions of the same service run simultaneously in the same environment for a period of time.



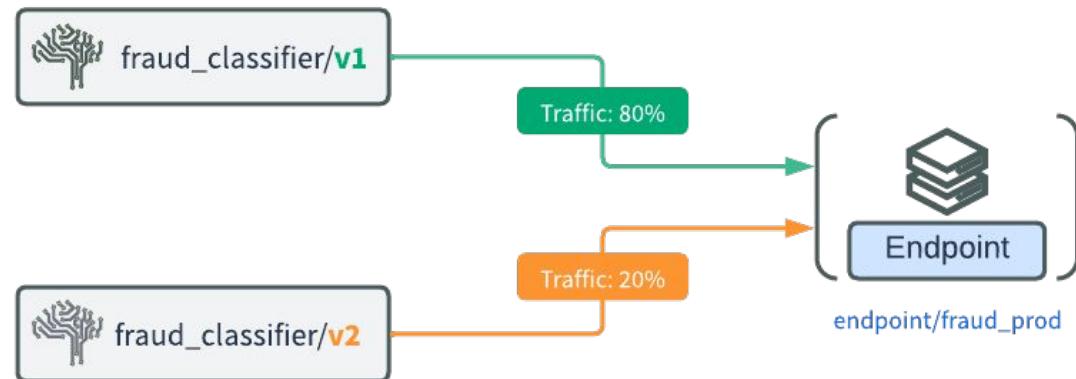
[Source](#)



# A/B Testing Framework

A/B testing is essentially an experiment where two or more variants of a deployment are shown to users at random, and statistical analysis is used to determine which variation performs better for a given conversion goal.

- Collect data
- Identify goals
- Generate hypothesis
- Create variations
- Run experiment
- Analyze results

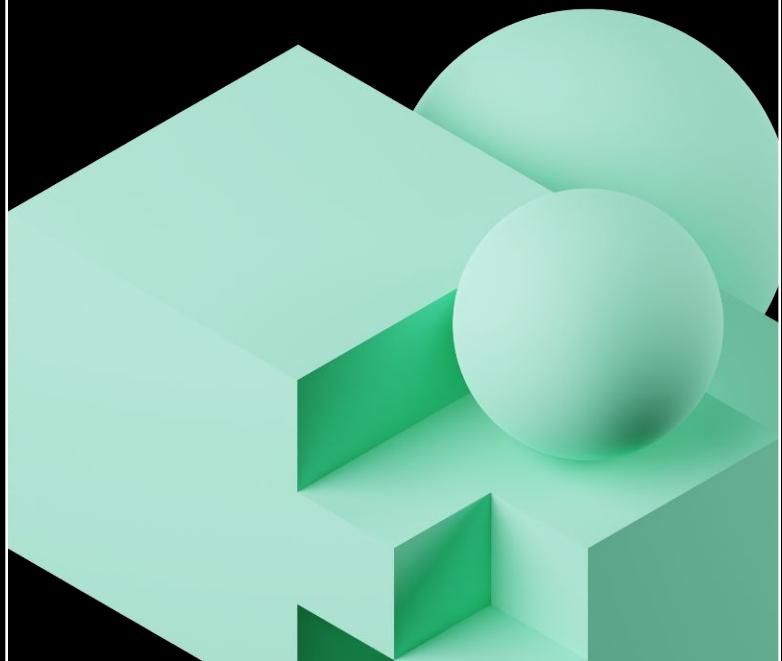


# A/B Testing Metrics

Project	ML metric	Business Metric
Fraud Detection	Area Under PR, Area Under ROC, f1	Fraud Loss \$, # Fraud Investigations
Churn Prediction	Area Under PR, Area Under ROC, f1	Recency of purchases Login events for high churn risk
Sales Forecasting	AIC, BIC, RMSE, et al.	Revenue
Sentiment Analysis	bleurt, bertscore	Number of users of tool, engagement rate
Ice Cream Coupons	MAE, MSE, RMSE	Revenue, Coupon usage

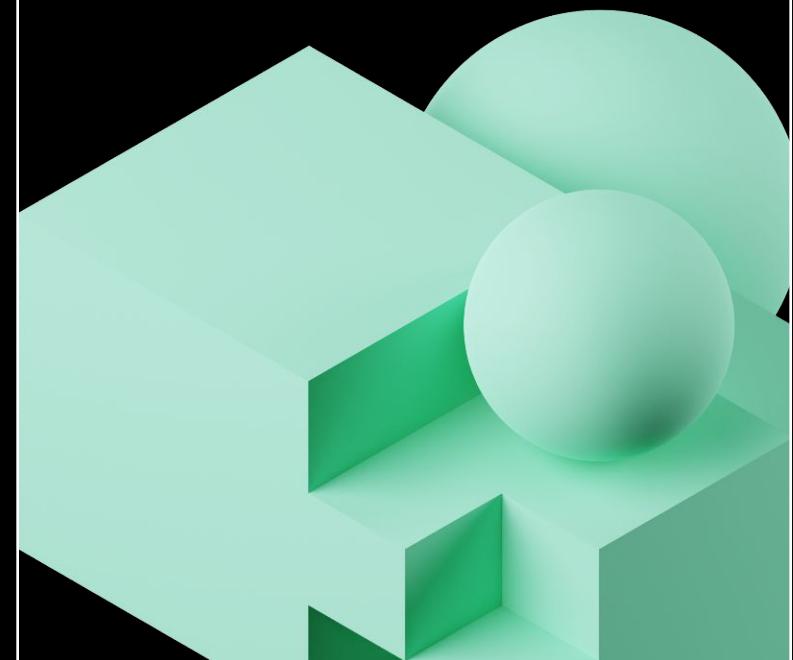


# Questions?



# Summary and Next Steps

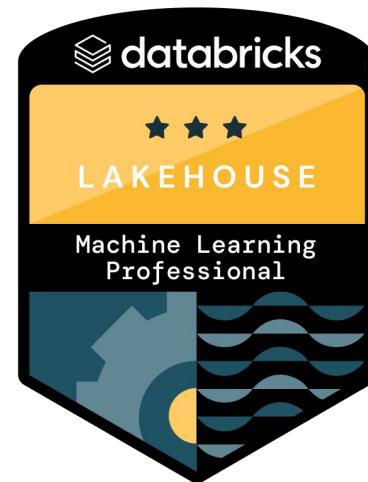
DAWD 00-1



# Databricks Certified ML Professional

Certification helps you gain industry recognition, competitive differentiation, greater productivity, and results.

- This course helps you prepare for the **Databricks Certified Machine Learning Professional exam**
- Please see the Databricks Academy for additional prep materials



For more information visit:  
[databricks.com/learn/certification](https://databricks.com/learn/certification)

# THANK YOU!



databricks