

The Modern MLOps Framework

MLOps (Machine Learning Operations) is a discipline that applies DevOps principles to machine learning systems. It ensures that ML models are not only developed but also deployed, monitored, and maintained in a scalable, automated, and reproducible manner.

2. Definition and Scope

- **MLOps:** A set of principles, practices, and tools that automate and streamline the lifecycle of ML models.
- **Scope:** Covers the entire ML lifecycle, from development to deployment and maintenance. This course focuses on the **Operations** phase post-training deployment, monitoring, and updates.

3. MLOps vs DevOps

- **DevOps:** Focuses on software development and deployment using source code.
- **MLOps:** Adds complexity by integrating data and models into the development and deployment pipeline.

Example: DevOps deploys a web application built from code. MLOps deploys a trained ML model that relies on data pipelines, model artifacts, and monitoring systems.

4. Risks of ML Without MLOps

- Manual workflows lead to inefficiency and errors.
- Lack of monitoring causes model drift and degraded performance.
- Accumulated technical debt slows innovation and increases operational cost.

Technical Debt Definition: The implied cost of additional rework caused by choosing an easy solution now instead of a better, longer-term approach.

Reference: Google's paper "Machine Learning: The High-Interest Credit Card of Technical Debt" highlights how unmanaged ML systems accumulate risk and inefficiency.

5. ML Workflow Stages

Typical ML workflows include:

1. Data Collection and Preparation
2. Data Labeling

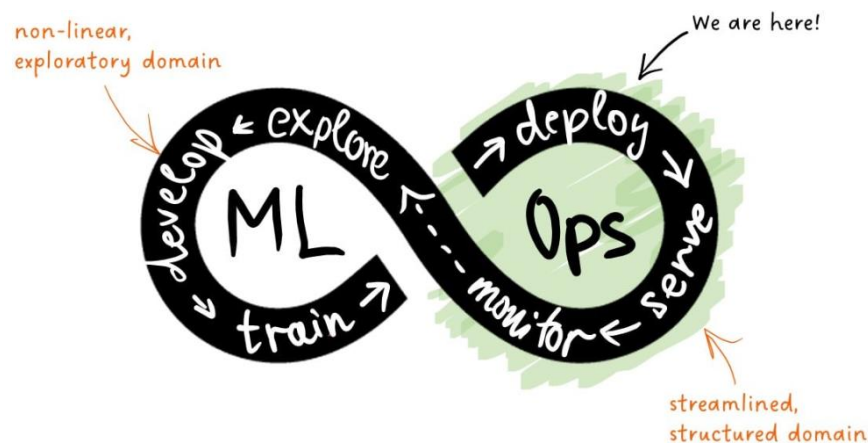
3. Model Selection
4. Model Training
5. Model Packaging
6. Model Deployment
7. Model Monitoring
8. Model Maintenance

The more these workflows are automated and integrated into the IT infrastructure, the higher the MLOps maturity of the organization.

6. Benefits of MLOps Implementation

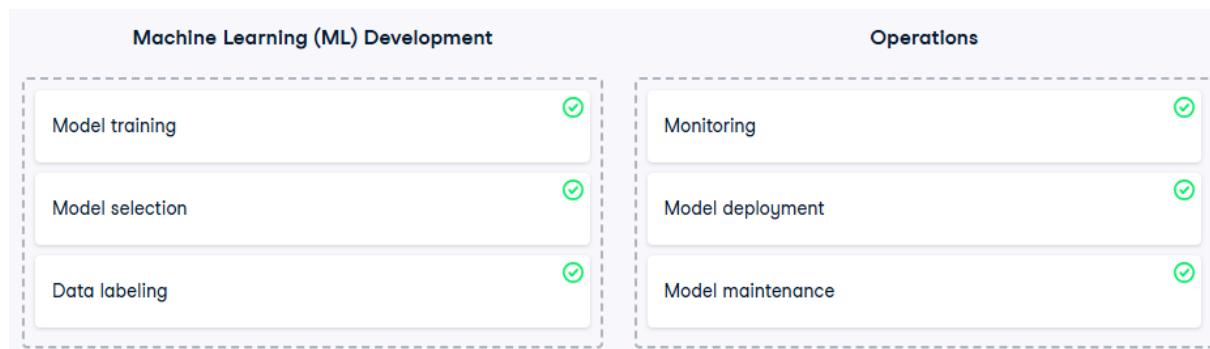
- Automation of repetitive tasks
- Reproducibility of model results
- Faster deployment cycles
- Explainability and transparency
- Scalable and maintainable ML systems
- Improved customer trust and service quality

7. Focus on the Operations Phase



- Begins after model training
- Includes deployment as a service, real-time monitoring, and maintenance
- Transitions from exploratory, nonlinear development to structured, linear production workflows

Key Insight: Once deployed, the model is exposed to users. Any error or drift is immediately visible, requiring rapid response and minimal margin for error.



Model Life Cycle in MLOps

1. Introduction to Life Cycles in Machine Learning

In machine learning, the term "life cycle" can refer to multiple layers of activity. To avoid confusion, it's important to distinguish between:

- **ML Project Life Cycle**
- **ML Application Life Cycle**
- **ML Model Life Cycle**

ML Model Life Cycle, which begins after model training and continues through deployment, monitoring, and eventual decommissioning.

2. Definitions and Relationships

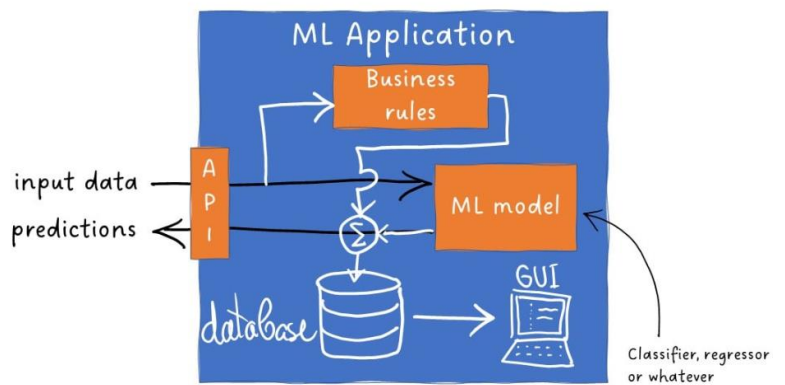
ML Project Life Cycle

- Represents the overarching effort to solve a business problem using machine learning.
- Includes problem definition, data acquisition, experimentation, and solution delivery.
- If successful, it results in the creation of ML applications and models.

ML Application vs ML Model

- **ML Model:** A trained estimator that performs a specific task (e.g., predicting daily sales).
- **ML Application:** A complete software system that includes the ML model and additional components such as:
 - Business rules (e.g., **fallback** logic for cold-start users)

- Databases (for storing features and logs)
- Graphical User Interfaces (for configuration and troubleshooting)
- APIs (for secure external communication)



Example: An ML model predicts movie ratings. An ML application uses that model, applies business rules, stores user data, and serves recommendations via an API.

3. Monolithic vs Microservice Architecture

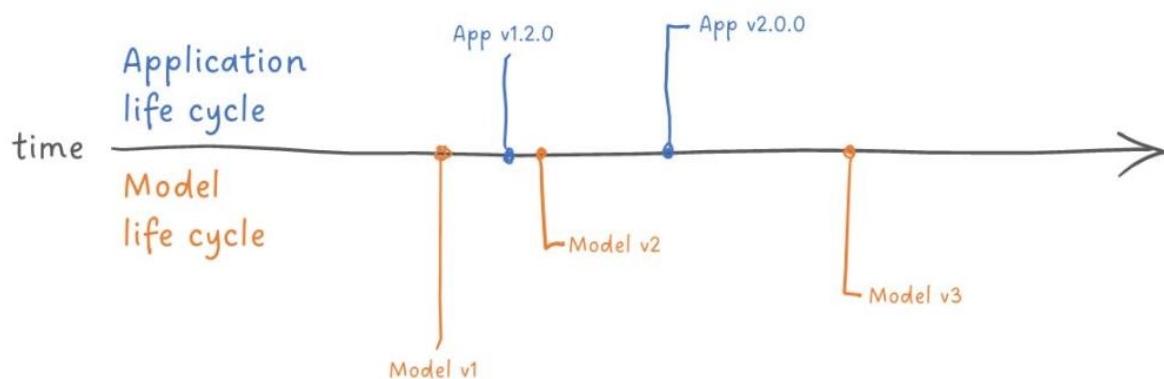
Monolithic ML Application

- The model is tightly coupled with the application.
- Difficult to update or replace components independently.

Microservice Architecture

- Model and application are decoupled.
- Each component evolves independently.
- Enables separate life cycles for the application and the model.

Analogy: Think of the ML application as a car and the ML model as its tires. The car may last decades, but the tires (models) are replaced frequently.



4. ML Model Life Cycle Stages

1. Deployment

- A trained model, along with necessary resources, forms a **deployment package**.
- Deployment marks the beginning of the model's operational life.
- The model is exposed to real-world data and users.
- Model object + deployment resources = deployment package

2. Monitoring

- Continuous observation of model behavior post-deployment.
- Ensures the model is running and performing as expected.
- Detects issues like model drift, latency, or prediction errors.

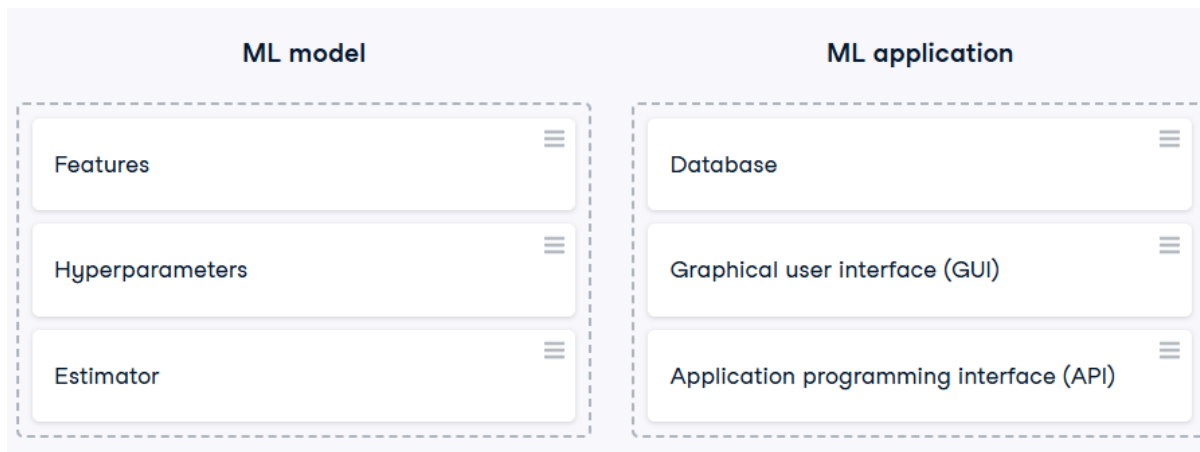
3. Decommissioning

- Triggered when the model becomes outdated or underperforms.
- Reasons may include:
 - Better model available
 - Improved features
 - Changes in the underlying business process
- The model is retired and replaced.

4. Archiving

- Essential for regulated industries.
- Archived models must be reproducible able to be reloaded and executed exactly as before.
- Enables auditability and historical analysis.

Example: A financial institution may need to reproduce a credit scoring model from five years ago to justify a loan decision during an audit.



Core Components of the MLOps Framework

MLOps extends DevOps principles to machine learning systems, enabling automated, reproducible, and scalable workflows for building, deploying, and maintaining ML models and applications. This framework includes both general software engineering concepts and ML-specific infrastructure.

2. Foundational Concepts

Workflows

- A **workflow** is a sequence of tasks that transforms inputs into outputs.
- Execution modes:
 - Manual
 - Automatic
 - Semi-automatic

Pipelines

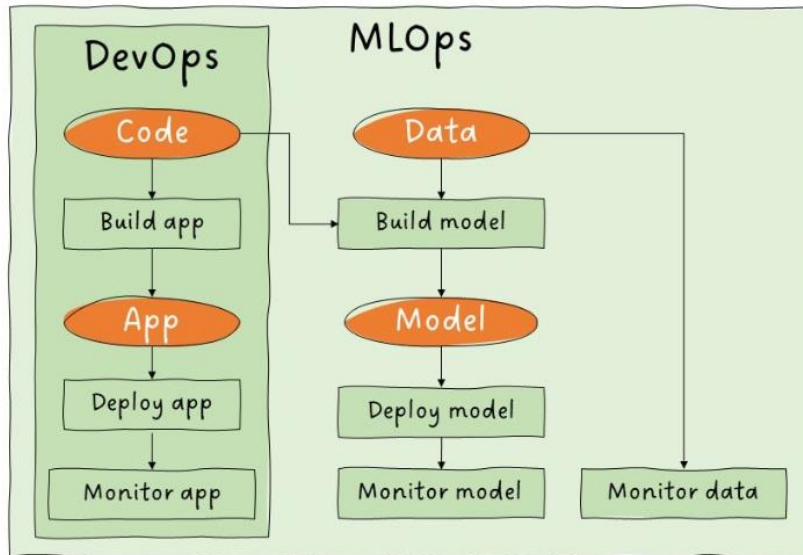
- A **pipeline** is a scripted, automated workflow.
- Common across DevOps, DataOps, and MLOps.
- Converts workflows into repeatable, programmable processes.

Artifacts

- **Artifacts** are outputs of pipelines.

- Includes compiled code, trained models, configuration files, and metadata.
- Used for deployment, versioning, and reproducibility.

3. ML-Specific Components

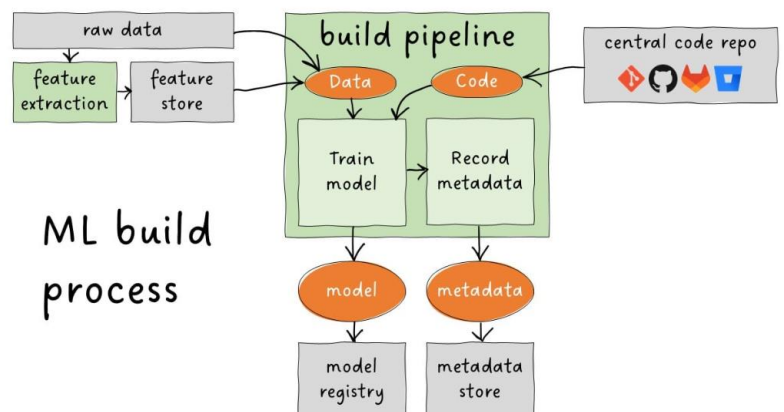


Component	Description
Feature Store	A database for storing processed variables (features) used in model training.
Model Registry	A system for storing, versioning, and managing trained ML models.
Metadata Store	Stores auxiliary data about models, including training parameters, datasets, and performance metrics.

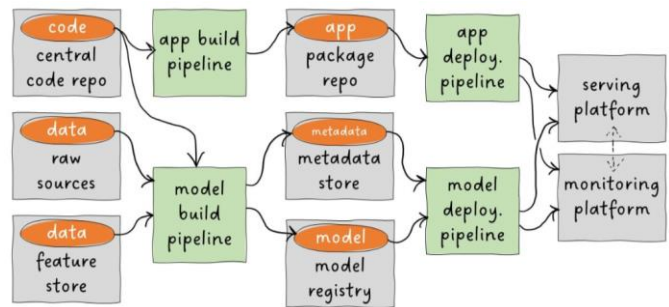
4. Build Pipelines

Model Build Pipeline (Training Pipeline)

- Transforms raw or processed data and model code into a trained model.
- Requires:
 - Model definition
 - Training script
 - Training data



- Outputs:
 - Trained model artifact
 - Model metadata
- Stores results in:
 - Model Registry
 - Metadata Store



Application Build Pipeline

- Classical DevOps process for packaging ML applications.
- Steps:
 - Pull code from central repository
 - Apply transformations (e.g., containerization)
 - Store deployable app in a repository

5. Deployment Pipeline

- Combines model and application artifacts.
- Deploys them to a target serving platform.
- Ensures the system is ready for production use.

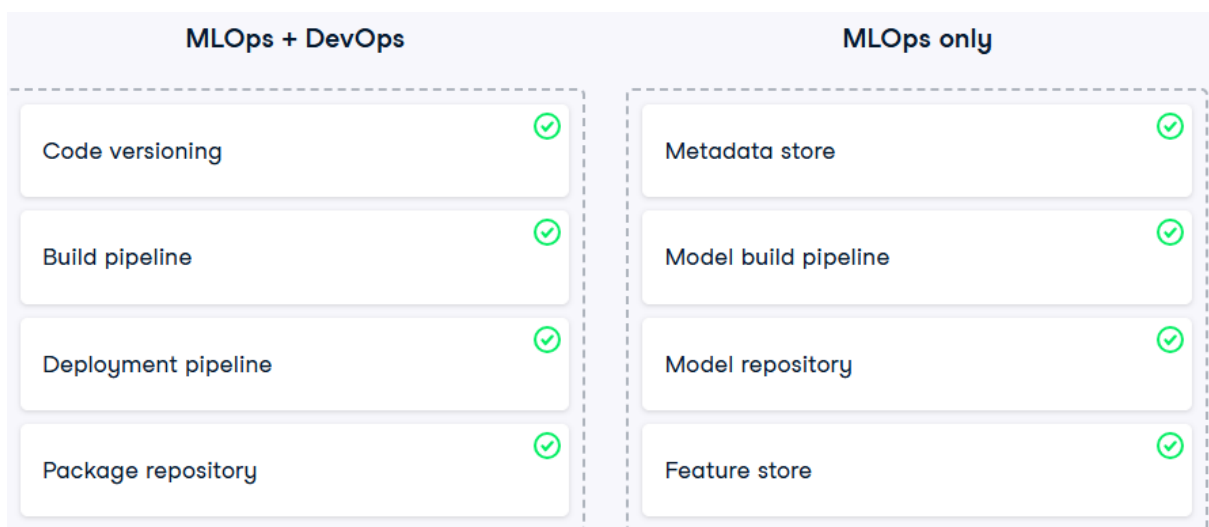
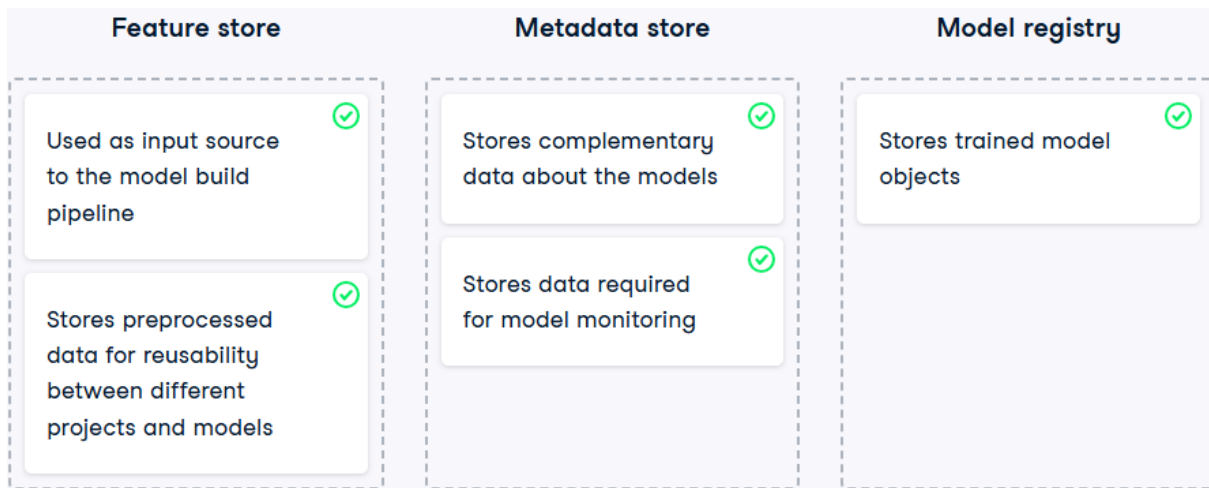
Example: The deployment pipeline may use Docker and Kubernetes to launch a REST API that serves predictions from a trained model.

6. Monitoring

- Begins post-deployment.
- Tracks:
 - Service uptime
 - Prediction accuracy
 - Latency
 - Model drift
- Ensures the model continues to perform as expected in real-world conditions.

7. High-Level MLOps Architecture Summary

Stage	Description
Workflow Design	Define tasks and dependencies
Pipeline Scripting	Automate workflows via scripts
Build Pipelines	Train models and package applications
Artifact Management	Store outputs in registries and metadata stores
Deployment Pipeline	Launch model and app on serving infrastructure
Monitoring	Continuously evaluate performance and reliability



Deployment-Driven Development in MLOps

Deployment-driven development is a mindset in MLOps that emphasizes planning for deployment from the earliest stages of model development. It ensures that models are not only accurate but also deployable, maintainable, and robust in real-world environments.

Key Principle: MLOps does not begin at deployment it begins during development. Ignoring deployment until the end risks producing models that are incompatible with production systems.

2. Analogy: The Race Car Driver

Just as a race car driver adjusts speed and steering before reaching a curve, ML engineers must anticipate deployment constraints early in the development process. This proactive approach prevents costly rework and deployment failures.

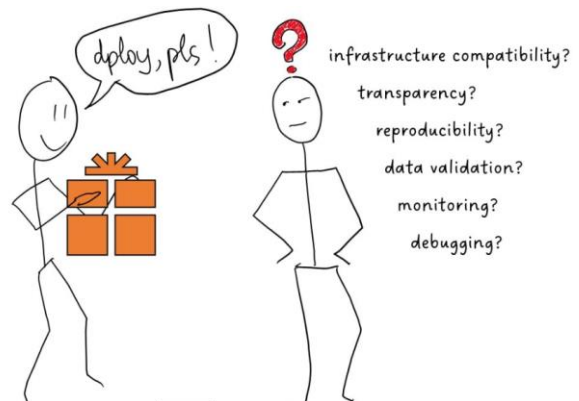
3. Deployment Scenario: Inheriting a Colleague's Model

Imagine being tasked with deploying and maintaining someone else's ML model and application. Several critical concerns immediately arise:

4. Key Deployment Concerns and Best Practices

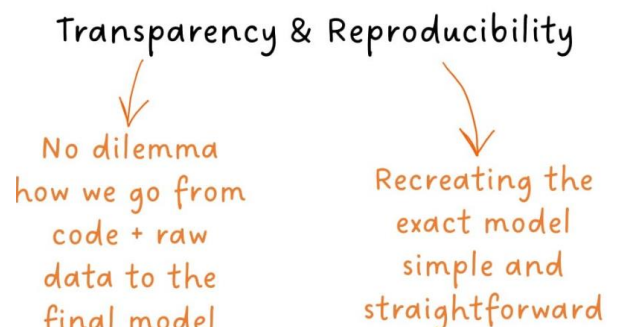
1. Infrastructure Compatibility

- **Concern:** Can the model run on the target platform?
- **Example:** A model trained on a 16-core server may be incompatible with a smartphone deployment.
- **Best Practice:** Understand deployment constraints (memory, CPU, OS) early in development.



2. Transparency and Reproducibility

- **Concern:** Is it clear who trained the model, when, and with what data and parameters?
- **Best Practice:**
 - Use versioned datasets and transparent pipelines.
 - Log all training experiments in a metadata store.



- Ensure reproducibility by documenting every step from raw data to trained model.

3. Input Data Validation

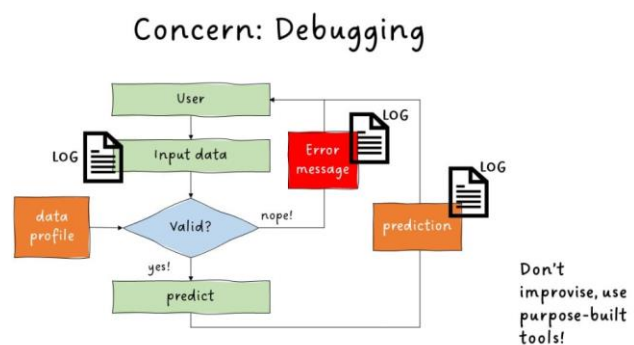
- **Concern:** How do we handle invalid user inputs?
- **Example:** A user provides a negative value where only positive values are allowed.
- **Best Practice:**
 - Define data profiles or expectations.
 - Save validation rules in model metadata during the build pipeline.

4. Performance Monitoring

- **Concern:** How do we detect if model performance deteriorates over time?
- **Best Practice:**
 - Log inputs and predictions.
 - Monitor metrics like accuracy, latency, and drift.
 - Set up alerts for performance thresholds.

5. Debugging

- **Concern:** Can we locate and fix bugs efficiently?
- **Best Practice:**
 - Implement detailed logging throughout the application.
 - Use specialized debugging tools.
 - Avoid ad hoc fixes build structured logging from the start.



6. Code Maintainability and Testing

- **Concern:** Can we safely modify the code without introducing new bugs?
- **Best Practice:**
 - Develop a comprehensive test suite:
 - Unit tests

- Integration tests
- Load tests
- Stress tests
- Deployment tests
- Maintain test coverage as the codebase evolves.

Data Profiling, Versioning, and Feature Stores in MLOps

Modern MLOps frameworks rely on robust data management practices to ensure model reliability, reproducibility, and performance in production. Three foundational components are:

- Data profiling
- Data versioning
- Feature stores

These tools and practices help validate inputs, track data lineage, and prevent training-serving mismatches.

2. Data Profiling

- Automated analysis of input data to generate high-level summaries called **data profiles** or **expectations**.
- Used for validating and monitoring data in production environments.

Purposes

- Provide feedback to users submitting invalid inputs.
- Detect data drift and trigger model retraining.
- Prevent misattribution of model errors when inputs are flawed.

Risks of Omission

- Clients may blame models for poor predictions caused by bad inputs.
- No mechanism to detect when retraining is needed due to data drift.

Best Practices

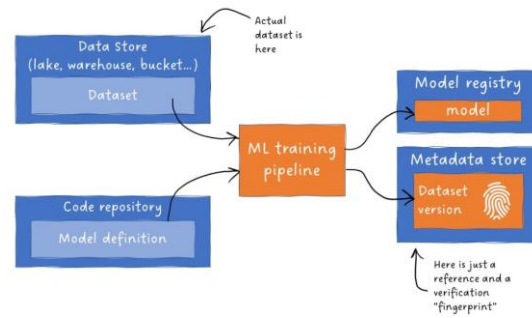
- Include a profiling step in the model training pipeline.
- Store data profiles in the **metadata store** alongside model metadata.

Tool Example

- **Great Expectations:** A popular Python-based open-source tool for data profiling and validation.

3. Data Versioning

- Recording the exact version of the dataset used to train a model.
- Ensures reproducibility without duplicating entire datasets.



Implementation

- Store datasets in a centralized location.
- Record a **pointer** and a **dataset fingerprint** in the model metadata.
- Fingerprint ensures no records have changed since training.

Bonus Practice

- Record metadata to reconstruct the exact **train-test split** used during model evaluation.

Tool Example

- **DVC (Data Version Control):** A widely used tool for tracking datasets, models, and experiments.

4. Feature Stores

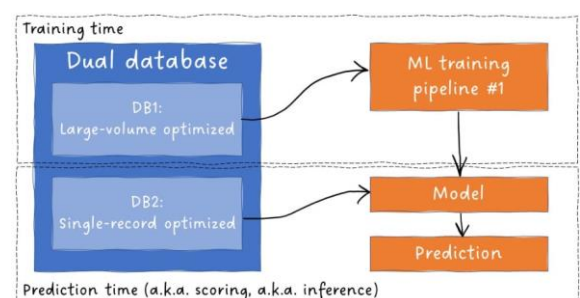
Centralized databases that store processed variables (features) for ML training and inference.

Benefits

- Reuse features across multiple models and projects.
- Reduce time spent on feature engineering.
- Prevent **training-serving skew**.

Architecture

- Often implemented as **dual databases**:
 - One optimized for bulk data retrieval during training.



- One optimized for fast, row-level access during inference.

Training-Serving Skew

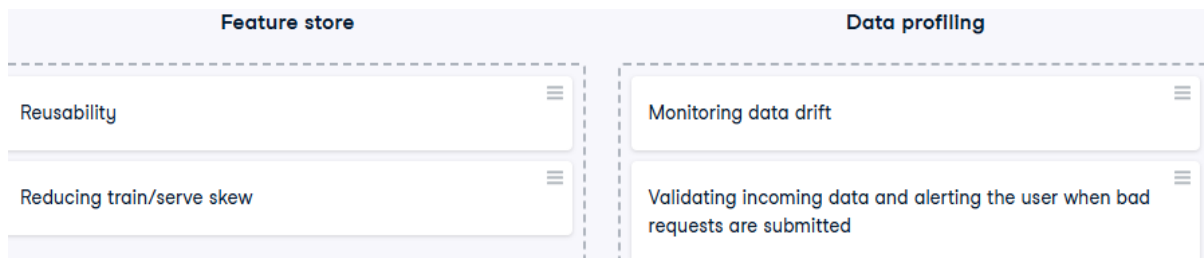
- Occurs when models perform well during training but poorly in production due to inconsistent data preprocessing like cleaning data in test but forgot to do so in production.
- Example: Training a spam filter on clean text emails, but deploying it on HTML emails.

5. Insum

Component	Purpose	Tools/Practices
Data Profiling	Validate inputs, detect drift, guide retraining	Great Expectations, metadata store
Data Versioning	Ensure reproducibility, track dataset lineage	DVC, dataset fingerprinting
Feature Stores	Centralize features, prevent training-serving skew	Dual DB architecture, reusable features

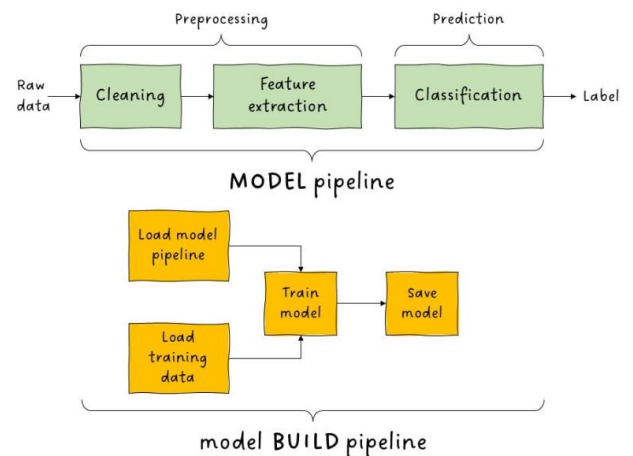
Robust data profiling, versioning, and feature management are essential for building reliable, maintainable, and auditable ML systems. These practices reduce risk, improve transparency, and ensure consistent performance across environments.

True	False
<div>A feature store is a type of a database, which stores data prepared specifically for ML models. ✓</div> <div>Feature stores help us improve efficiency and consistency, by allowing us to build a feature once, then reuse it for different models. ✓</div> <div>Advanced features stores are implemented as so-called dual databases, one for grabbing the training data and the other for making predictions. ✓</div>	<div>We can also store trained models in a feature store. ✓</div> <div>Feature stores help us monitor models deployed to production. ✓</div>



Model Build Pipelines in CI/CD for MLOps

In MLOps, model build pipelines are automated workflows that train machine learning models and prepare them for deployment. These pipelines are distinct from model pipelines (which define the data processing steps within a model) and must meet rigorous standards for deployment, reproducibility, monitoring, and CI/CD integration.



2. Distinction: Model Pipeline vs Model Build Pipeline

Term	Definition
Model Pipeline	Sequence of data processing steps (e.g., cleaning, feature extraction, prediction) executed by the model during inference.
Model Build Pipeline	Build Automated workflow that loads training data and model code, trains the model, and outputs deployment-ready artifacts.

3. Dual Build Pipelines in ML Systems

1. Application Build Pipeline

- Standard DevOps pipeline for packaging and deploying the ML application.
- Includes code compilation, dependency resolution, and containerization.

2. Model Build Pipeline

- Central to MLOps.
- Trains the model and generates a complete deployment package.

4. Requirements for an MLOps-Grade Model Build Pipeline

1. Deployment Readiness

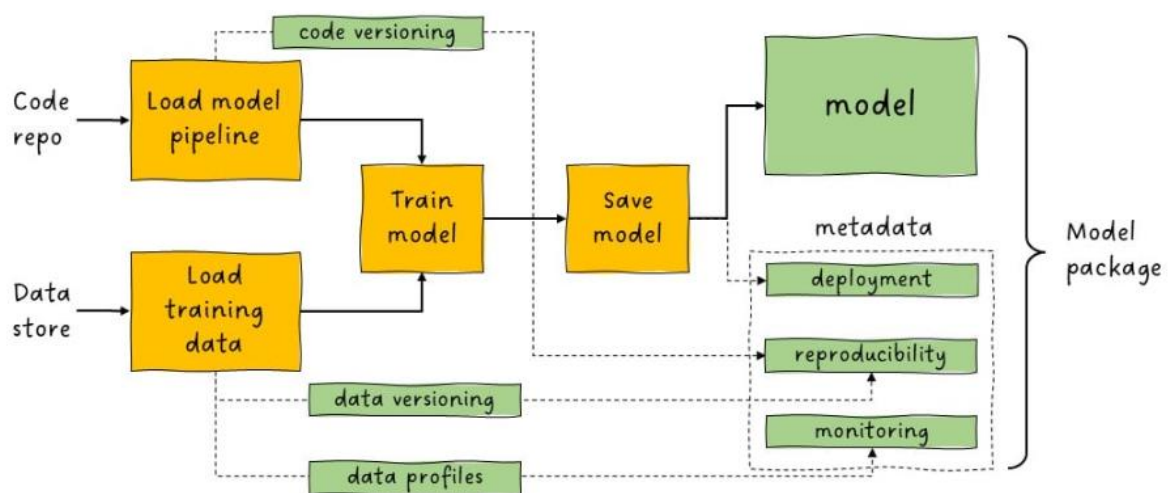
- Outputs a **complete model package**, not just the trained model.
- Includes:
 - Model object
 - Dependency specifications (e.g., environment files, Dockerfiles)
 - Configuration files
- Supports **test deployments** to verify compatibility with target infrastructure.

2. Reproducibility

- Ability to recreate the model from scratch at any time.
- Key practices:
 - **Code versioning** (e.g., Git commits)
 - **Data versioning** (e.g., DVC pointers and fingerprints)
 - Recording both in the **model metadata**

3. Monitoring Enablement

- Ensures the model behaves as expected in production.
- Requires:
 - Logging of inputs and predictions
 - Creation of **data profiles** during pipeline execution
 - Integration with monitoring tools and dashboards

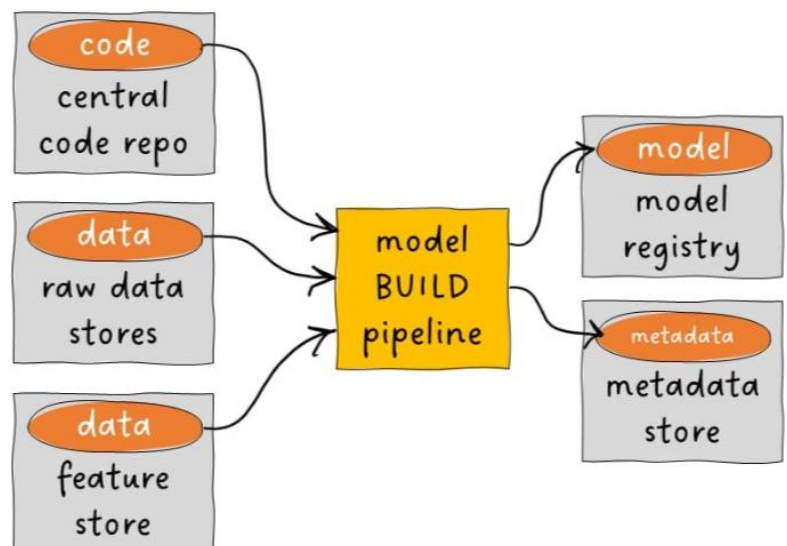


4. CI/CD Integration

- Pipeline runs within a **Continuous Integration/Continuous Deployment** framework.

- Benefits:

- Prevents use of unversioned code or data
- Enforces consistency and traceability
- Automates testing, packaging, and deployment



- Requires:

- Connection to input sources (e.g., code repo, feature store)
- Storage for generated artifacts (e.g., model registry, metadata store)

5. insum

Component	Purpose
Model Build Pipeline	Trains and packages models for deployment
Deployment Artifacts	Ensure compatibility and reproducibility

Component

Purpose

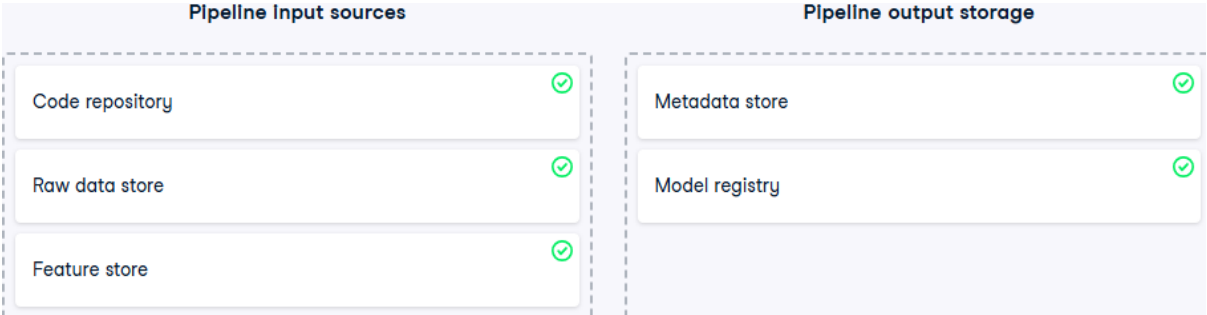
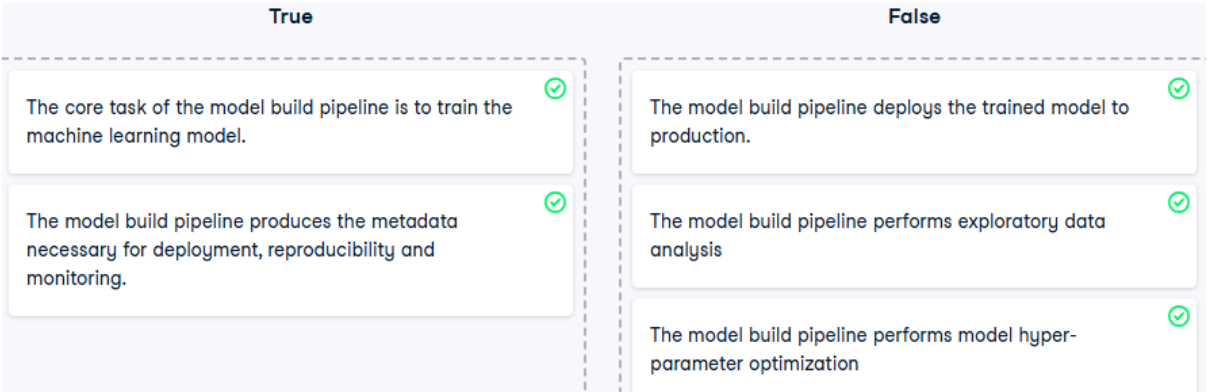
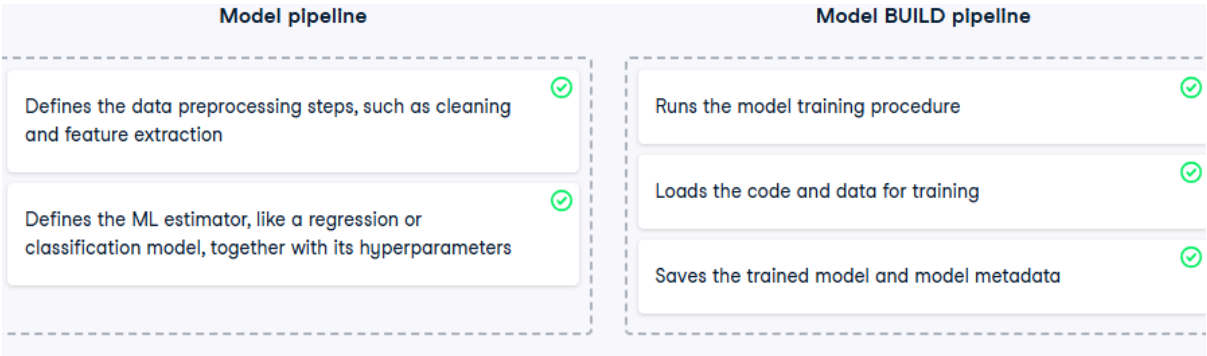
Monitoring Integration

Enables performance tracking and drift detection

CI/CD Enablement

Automates and secures the build process

A robust model build pipeline is the backbone of scalable, trustworthy MLOps. It ensures that models are not only trained effectively but also deployed reliably, monitored continuously, and reproducible under audit conditions.



Model Packaging in MLOps

Model packaging marks the transition from ML development to operations. It is the final step before deployment and must be designed to meet three core MLOps objectives:

- **Smooth Deployment**
- **Reproducibility**
- **Monitoring**

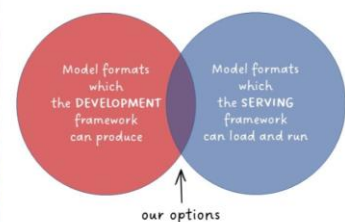
1. Model Storage Format

The trained model is the centerpiece of the package. It must be saved in a format that:

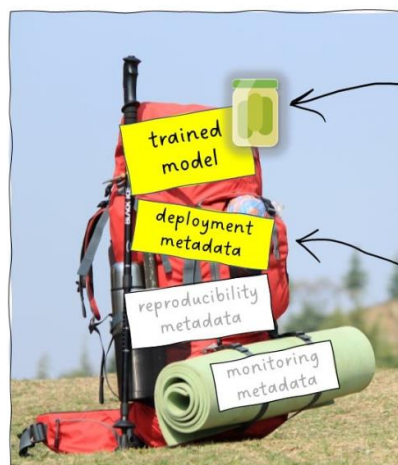
- Can be produced by the model development framework
- Can be loaded and executed by the serving framework



Model storage options



Common Format Options



IF `model_format == .pickle`
THEN:



store list of model's dependencies within model package and verify server compatibility when loading

Only choose the format after confirming the deployment compatibility

a. PMML (Predictive Model Markup Language)

- **Purpose:** Designed for cross-language and cross-platform compatibility.
- **Use Case:** Train in one language, serve in another.
- **Limitation:** Difficult to customize; less flexible for complex or custom models.
- **Best Fit:** When interoperability is critical and customization is minimal.

b. Pickle (Python Object Serialization)

- **Purpose:** Native to Python; stores almost any Python object.

- **Use Case:** Widely used in Python-based ML workflows.
- **Limitation:** Tied to Python; requires identical library versions on both training and serving sides.
- **Best Fit:** When working entirely within the Python ecosystem.

Note: If using Pickle, include a list of model dependencies in the package metadata to ensure compatibility during deployment.

2. Reproducibility

A model is reproducible if it can be rebuilt automatically at any time. This demonstrates full control over the model production process.

Required checklist for Reproducibility

- **Version Pointer to Build Pipeline Code:** Ensures the exact logic used to train the model is retrievable.
- **Versioned Dataset References:** Includes training and evaluation splits used during model development.
- **Performance Record:** Captures metrics from the test set to validate reproducibility.

3. Monitoring

Monitoring ensures the model behaves as expected in production. It can be integrated into the serving application or handled by an external service.

Monitoring Prerequisite

- **Data Profiles:** Define expectations for input and output data distributions. These must be saved within the model package to enable drift detection and alerting.

Model package		Not In model package	
Reproducibility metadata	✓	Model registry	✓
Deployment metadata	✓	Full copy of the training dataset	✓
Trained model object	✓	CI/CD pipeline	✓

Model Serving Modes in MLOps

Once a model package is complete—with the trained model and all necessary metadata—we enter the **Operations phase**. This is where the model begins its real-world lifecycle, **delivering predictions as a service**.

Model SErving

Model serving refers to the process of **providing predictions as a service**. From the end-user's perspective, the ML model behaves like any other service—similar to a food delivery app where users expect a timely response.

- Users **send a request** to the ML application
- The application **returns predictions** in response

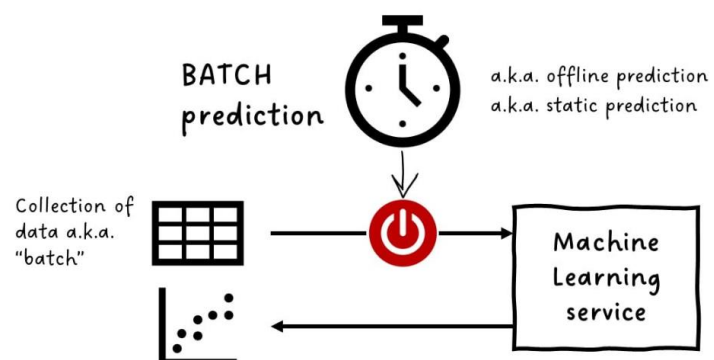
The **way** this service is implemented is called the **serving mode**, and choosing the right mode depends on the use case.

Serving Mode: When Should the Model Generate Predictions?

This is the first and most critical question. There are two primary categories:

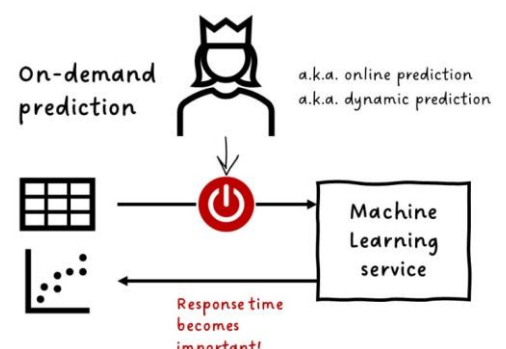
1. Batch Prediction (Offline / Static Prediction)

- **Definition:** Predictions are generated on a **scheduled basis**, not triggered by user actions.
- **Data:** Operates on a **batch** of data (e.g., all records from the past day or month).
- **Latency:** Not time-sensitive.
- **Use Case:** Monthly sales forecasts, churn prediction reports, etc.
- **Advantages:**
 - Simple to implement
 - Efficient for large-scale, non-urgent tasks



2. On-Demand Prediction (Online / Dynamic Prediction)

- **Definition:** Predictions are generated **in response to specific events or user requests**.
- **Latency:** Time becomes critical; users expect fast responses.
- **Use Case:** Chatbots, recommendation engines, fraud detection, etc.



- **Variants:**

a. Near-Real-Time Prediction (Stream Processing)

- **Latency:** Acceptable delay of a few seconds to minutes
- **Mechanism:** Data flows continuously into the model; predictions are streamed out
- **Use Case:** Social media trend analysis, sensor data monitoring

b. Real-Time Prediction

- **Latency:** Must respond in **under a second**
- **Use Case:** Credit card fraud detection, real-time bidding in ad tech
- **Constraint:** Delayed predictions are often useless

Latency Considerations

- **Latency** is the time between a **user request** and the **model's response**.
- Acceptable latency varies by use case:
 - Some users can wait minutes or hours
 - Others require instant responses

In latency-critical applications, trade-offs may be necessary:

- **Faster, simpler models** may be preferred over more accurate but slower ones
- **Edge Deployment** may be used:
 - Model is deployed **directly on the user's device**
 - Eliminates network latency
 - Common in:
 - Facial recognition for phone unlocking
 - Navigation apps
 - Real-time image filtering

API Architectural Styles

Common styles include:

- **REST (Representational State Transfer)**

- **RPC (Remote Procedure Call)**
- **SOAP (Simple Object Access Protocol)**

Each has its own conventions, but regardless of style, every API should include core functional components: input validation, output validation, authentication, and throttling.

Example: Navigation App API

Simplest Workflow

1. Client sends velocity and distance
2. Server runs the model
3. Server returns estimated time of arrival (ETA)

Input Validation

To prevent server errors and ensure robustness:

- Validate that all required fields are present
- Check that values are of expected types

Example

- Required fields: velocity, distance
- Expected types: positive real numbers

This validation is defined in the **request model**, which must be documented for API users.

Output Validation

To avoid sending faulty predictions:

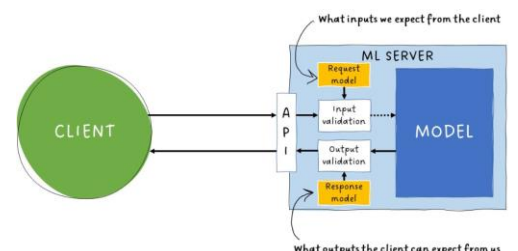
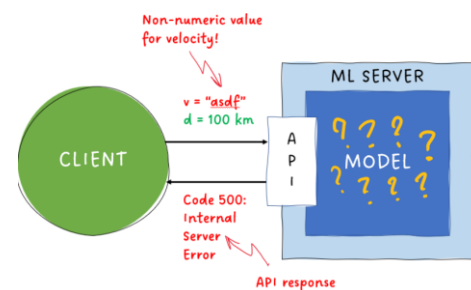
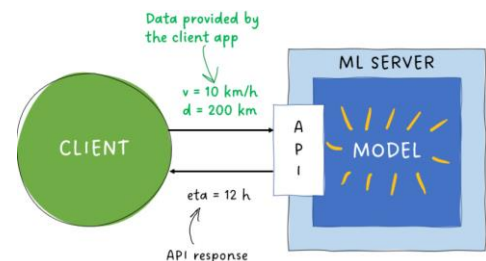
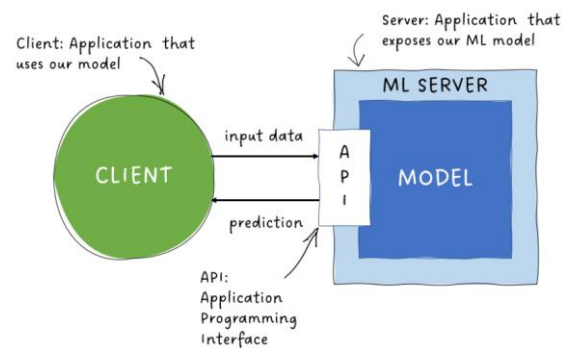
- Validate model outputs before returning them
- Detect anomalies (e.g., -7 minutes ETA but times can't be negative)

This is enforced via the **response model**, which specifies acceptable output formats and value ranges.

Authentication and Throttling

Authentication

- Restricts API access to authorized clients



- Implemented before input validation

Throttling

- Limits the number of requests per client per time unit
- Prevents abuse and ensures fair usage

Recommended Framework: FastAPI

- **FastAPI** is a Python-based, open-source framework ideal for ML API development.
- Features:
 - Built-in input/output validation
 - Automatic documentation generation
 - High performance and ease of use

Deployment Progression and Testing in ML Applications

Before deploying an ML-powered web API to production, it is essential to rigorously test the entire application not just the model's predictive performance. This ensures reliability, scalability, and safety in real-world usage.

Testing Scope

The focus is on testing the **ML application as a whole**, including:

- Database communication
- User authentication
- Logging
- API behavior
- External service integration

Environments for Testing

Environments are isolated setups used for different stages of development and testing:

Environment	Purpose
Development	Experimental coding and feature building
Test	Stable environment for running unit tests

Environment Purpose

Staging Replica of production for integration and load testing

Production Live environment serving real users

Testing Types and Progression

1. Unit Testing

- **Goal:** Verify individual functions or components behave as expected.
- **Characteristics:**
 - Fast and simple
 - Run frequently during development
- **Best Practice:** Execute in a stable **test environment**, not in the volatile development setup.

2. Integration Testing

- **Goal:** Ensure components work together and external services (e.g., databases, APIs) are accessible.
- **Environment:** Conducted in **staging**, which mirrors production setup.
- **Requirements:**
 - Same software versions as production
 - Representative subset of production data

3. Smoke Testing

- **Goal:** Confirm the application can be deployed and started without crashing.
- **Environment:** Run in staging before deeper tests.

4. Load Testing

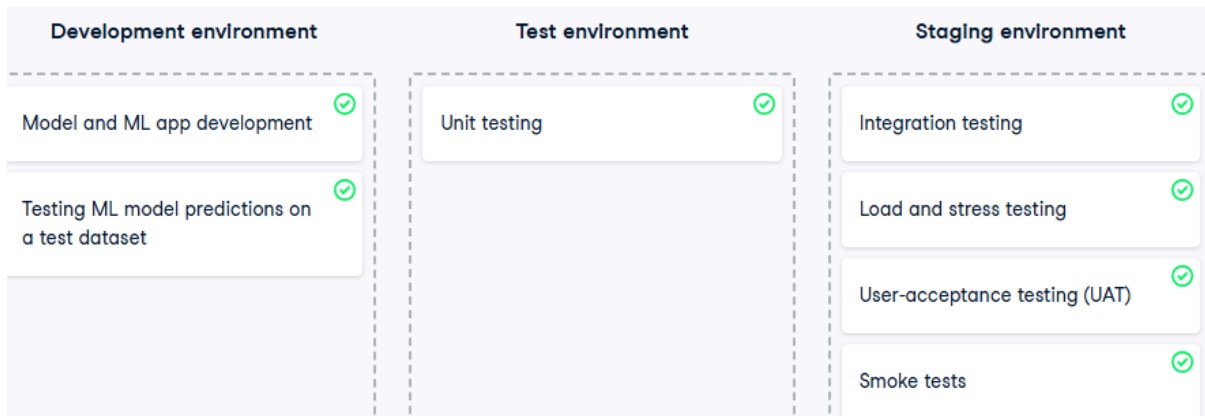
- **Goal:** Validate performance under expected user load (request frequency, input size).
- **Stress Testing:** Pushes load beyond normal limits to test system resilience.

5. User Acceptance Testing (UAT)

- **Goal:** Final validation by real users before production release.
- **Outcome:** Confirms usability and readiness for deployment.

Deployment Caution and Prioritization

- Even a single faulty component can compromise the entire system.
- Historical example: The Challenger SHUTTLE disaster 1968 was caused by a single defective sealing ring.
- **Prioritize testing critical components** over exhaustive edge-case coverage to maintain efficiency.

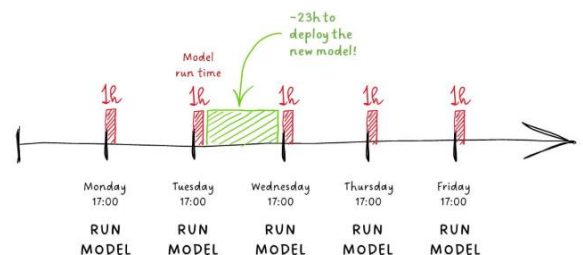


Model Deployment Strategies in MLOps

Once an ML model has passed all testing phases and is successfully deployed, the challenge shifts to **updating and replacing models** without disrupting service. Choosing the right deployment strategy depends on the prediction mode (batch vs real-time), risk tolerance, and system architecture.

1. Offline Deployment

- **Use Case:** Batch prediction services (e.g., daily or monthly runs)
- **Characteristics:**
 - Predictions are generated on a schedule
 - Model can be swapped between runs with minimal risk
- **Advantage:** Simple and low-risk
- **Limitation:** Not suitable for real-time systems



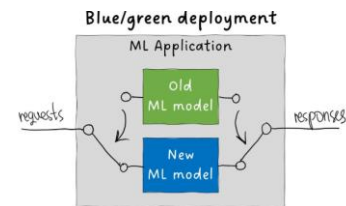
2. Blue/Green Deployment

- **Mechanism:**
 - Both old and new models are loaded in parallel

- Incoming requests are redirected from the old to the new model at once

- **Advantage:**

- Instantaneous switch with no downtime
- Easy rollback to the old model if issues arise



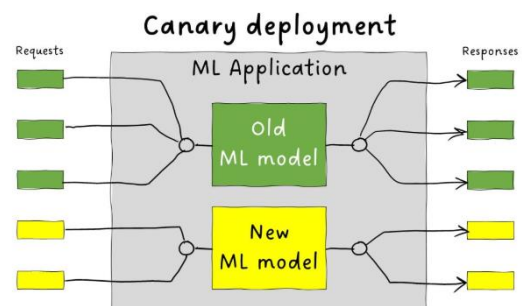
- **Risk:**

- All users are exposed to the new model immediately
- If the new model fails, it affects everyone

3. Canary Deployment

- **Mechanism:**

- Gradual rollout of the new model
- Start with a small percentage of traffic
- Increase traffic incrementally after successful validation



- **Advantage:**

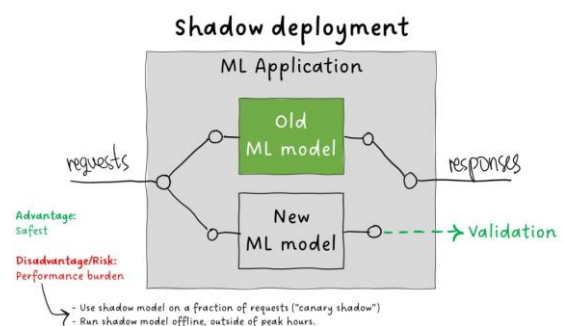
- Controlled exposure
- Easier to detect and isolate issues

- **Use Case:** Real-time systems with moderate risk tolerance

4. Shadow Deployment

- **Mechanism:**

- Requests are sent to both old and new models in parallel
- Only the old model's predictions are returned to users
- New model's outputs are logged for validation



- **Advantage:**

- Safest strategy; no user impact
- Enables performance comparison and debugging

- **Limitation:**
 - Doubles compute load
 - May affect performance if used in real-time
- **Optimization:**
 - Run shadow model on a subset of requests
 - Schedule shadow runs during off-peak hours

Summary Comparison

Strategy	Risk	Downtime	Rollback	Best For
Offline	Low	None	Easy	Batch systems
Blue/Green	Medium	None	Instant	Real-time, low-latency systems
Canary	Low	Minimal	Gradual	Real-time with cautious rollout
Shadow	Very Low	None	Not needed	Validation before full rollout

Monitoring and Maintaining ML Services

Once an ML service is deployed and actively serving predictions, especially to paying users, maintaining its quality becomes critical. Monitoring ensures that both the **technical health** and **predictive performance** of the service remain reliable over time.

1. Quality Assurance via Monitoring

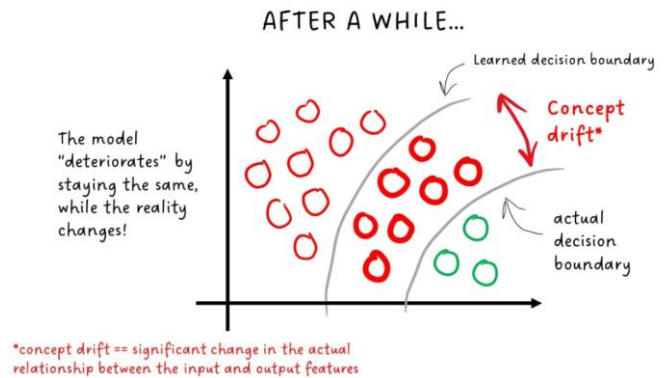
Monitoring is the first step in quality control. It involves tracking:

- **Service uptime:** Is the API running?
- **Request volume:** How many predictions are being served?
- **Success rate:** Are requests being handled correctly?
- **Latency:** Are responses delivered within acceptable timeframes?
- **Prediction quality:** Are the outputs still accurate and relevant?

2. Predictive Performance and Concept Drift

What is Concept Drift?

- A significant change in the relationship between input features and output predictions.
- **Cause:** The real-world data distribution evolves, but the model remains static.
- **Effect:** The model starts making incorrect predictions because its learned decision boundary no longer reflects reality.



Example

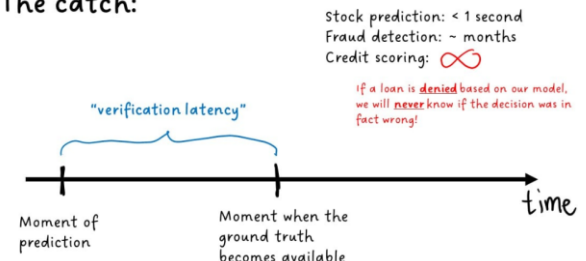
- A classifier trained to separate objects using two features learns a decision boundary.
- Over time, the properties of the classes shift, and the boundary becomes outdated.
- The model continues to apply the old boundary, leading to degraded performance.

3. Detecting Concept Drift

Ground Truth Comparison

- **Ideal method:** Compare predictions against actual outcomes.
- **Challenge:** Ground truth is often delayed or unavailable.

The catch:



Verification Latency

Even when available, ground truth can be costly and slow to obtain.

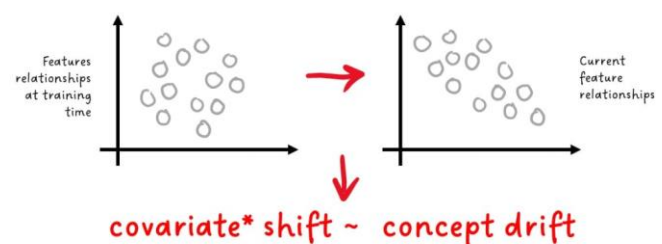
4. Alternative Monitoring Strategies

a. Input Feature Monitoring

- **Goal:** Detect changes in relationships between input features.
- **Term:** Covariate shift — a shift in feature distributions.
- **Limitation:**
 - Covariate shift \neq concept drift

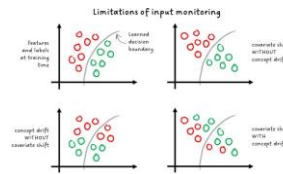
Indirect approach:

Monitor what you have: input features!



* features == "covariates" in statistical jargon

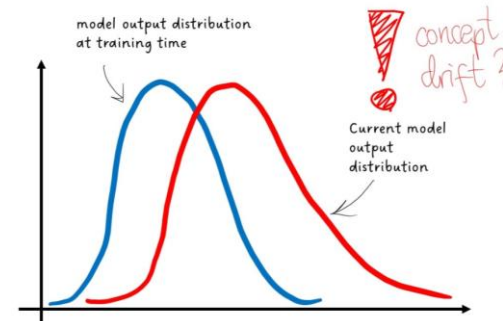
- Concept drift can occur without shift
- Both can occur simultaneously



covariate

b. Output Monitoring

- **Goal:** Track changes in the distribution of model outputs.
- **Use Case:** When ground truth is unavailable
- **Indicator:** Shifts in output patterns may signal that the model no longer reflects reality



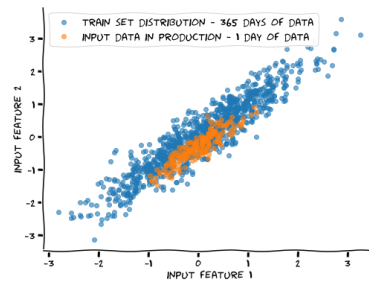
Summary

Monitoring Type	Detects	Strength	Limitation
Ground Truth Comparison	Concept drift	Direct and accurate	Latency and cost
Input Monitoring	Covariate shift	Fast and cheap	Indirect, may miss drift
Output Monitoring	Output anomalies	Useful without ground truth	Less precise

Based on the ground truth only, we might get the impression that our model has deteriorated, while, in fact, our data pipeline is broken and feeding corrupted inputs to our model. So input monitoring should *always* be included in your overall approach to monitoring models in production.

True	False
<p>Concept drift can only be confirmed using the ground truth ✓</p> <p>Monitoring input data is sufficient to detect covariate shift ✓</p> <p>Concept drift is a measure of how much the real modeled process has changed since we have trained our model ✓</p>	<p>Models deteriorate by starting to produce different outputs for the same input data ✓</p> <p>If covariate shift is detected, it is certain that concept drift has also occurred ✓</p> <p>If there is no covariate shift, there can be no concept drift ✓</p>

Imagine the following scenario: You have trained a model using a full year of labeled data and deployed it into production. One day after deployment, you decide to check your data monitoring dashboard, and you see the following picture, where the training inputs are represented with blue and the production inputs with orange dots.



Although reality is constantly changing, it rarely happens overnight. We should not jump to conclusions and first collect enough data points to be certain in our suspicions that the model requires maintenance. We can make no certain conclusion based on comparing one day of production data with a full year of training data. Moreover, the differences in the distributions are not drastic in any way.

Monitoring and Alerting in ML Services

Beyond concept drift and external data shifts, ML services are vulnerable to internal failures—bugs, data errors, and infrastructure issues. A robust monitoring and alerting system is essential to detect, diagnose, and resolve these problems quickly.

1. Internal Failures in ML Systems

ML services consist of multiple components:

- Model
- Data pipeline
- API
- Hardware and infrastructure

Each component is a potential point of failure. Monitoring must go beyond detecting that something is wrong—it must help pinpoint **what** is wrong and **where** to look.

2. Granular Logging

Detailed logging is critical for diagnosing issues:

- Track request metadata (e.g., user ID, input size, timestamps)
- Identify patterns in latency, errors, or throughput
- Example: If 5% of requests show high latency, logs should reveal which users, what data volume, and which endpoint was involved

Without granular logs, debugging becomes guesswork.

3. Data Pipeline Monitoring

Validate all stages of the data pipeline:

- Monitor each input source individually (e.g., separate tables)
- Don't rely solely on the final merged dataset
- Detect schema mismatches, missing values, and unexpected formats early

4. Data Profiles and Expectations

Use **data profiles** to define and enforce input/output expectations:

- Acceptable value ranges
- Missing value thresholds
- Feature relationships and statistical distributions

Compare incoming data against these profiles to detect anomalies.

Validation Types

Type	Description	Strength	Limitation
Hard Constraints	Lists of valid values, type checks	Fast and precise	Limited scope
Statistical Validation	Distribution shifts, correlation changes	Detects subtle drift	Can be noisy and overly sensitive

5. Alert Fatigue and Threshold Tuning

Overly sensitive alerts can lead to **alert fatigue**, where critical warnings are ignored due to excessive noise. To avoid this:

- Choose metrics carefully
- Set thresholds that balance sensitivity and relevance
- Prioritize actionable alerts

6. Alerting Mechanism

Once an incident is detected:

- Alerts must reach the right stakeholders immediately
- Use multi-channel notifications (email, Slack, dashboards)
- Include context and severity in alert payloads

7. Incident Logging and Postmortem Analysis

After resolving an incident:

- Record root cause and resolution steps
- Maintain an incident history for future reference
- Insight: In a 10-year study of one ML pipeline, over two-thirds of outages were **not ML-related** (source: USENIX OPML20)

8. Centralized Monitoring Infrastructure

For scalability and consistency:

- Use a centralized monitoring service
- Run it on dedicated infrastructure
- Monitor all ML services from one place
- Benefits:
 - Easier maintenance
 - Unified alerting
 - Higher service quality guarantees

True	False
It should preserve a log of all previous major issues and mitigation actions. ✓	It should only focus on the most important service health metrics, such as uptime and latency. ✓
It should be implemented as a centralized service for monitoring all running ML applications. ✓	It should eliminate any potential unethical bias from the model predictions. ✓
It should tell us <i>what</i> is wrong and <i>where</i> (in which application component) to look for the root cause. ✓	It should stop the ML service as soon as it detects any covariate shift. ✓

Maintaining ML Model Quality: Data-Centric vs Model-Centric Approaches

When an ML model deteriorates below acceptable performance levels, we must intervene. There are two primary strategies:

- **Improve the model itself** → Model-centric approach
- **Improve the training data** → Data-centric approach

1. Model-Centric Approach

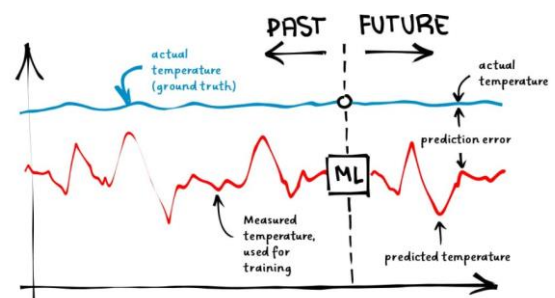
- **Definition:** Focuses on optimizing model architecture, hyperparameters, and feature engineering.
- **Common in:** ML competitions where datasets are fixed.
- **Techniques:**
 - Try different algorithms (e.g., XGBoost, neural networks)
 - Ensemble multiple models
 - Engineer new features from existing data
- **Limitation:** Performance gains are often incremental when data quality is poor.

2. Data-Centric Approach

- **Definition:** Focuses on improving the quality, richness, and relevance of training data.
- **Common in:** Real-world ML systems where data can be cleaned, enriched, or expanded.
- **Benefits:**
 - Yields better performance for the same effort compared to model-centric tuning
 - More robust and scalable improvements

Key Principles

- **Quality > Quantity:** More informative features and accurate labels matter more than dataset size.
- **Label Quality:**
 - Labels must closely reflect ground truth
 - Poor labels (e.g., noisy thermometer readings) lead to poor models



3. Labeling Tools and Efficiency

Manual labeling is error-prone and time-consuming. Modern labeling tools offer:

- Purpose-built interfaces for efficient annotation

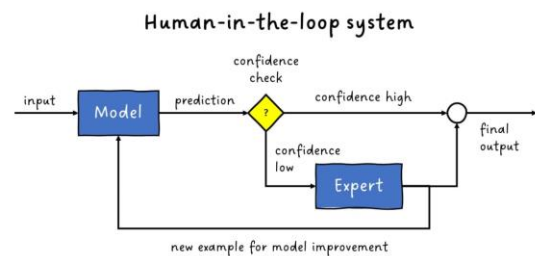
- Prioritization of high-impact samples
- Error detection and correction
- Example: Tools for image classification labeling

4. Human-in-the-Loop (HITL) Systems

- **Definition:** Humans and ML models collaborate during prediction and labeling.

- **Workflow:**

- ML model predicts with confidence score
- If confidence is high → prediction accepted
- If low → forwarded to human expert
- Each human decision becomes a new labeled example



Use Case

- Medical diagnostics: ML suggests diagnosis, doctor confirms or corrects

5. Iterative Improvement and Experiment Tracking

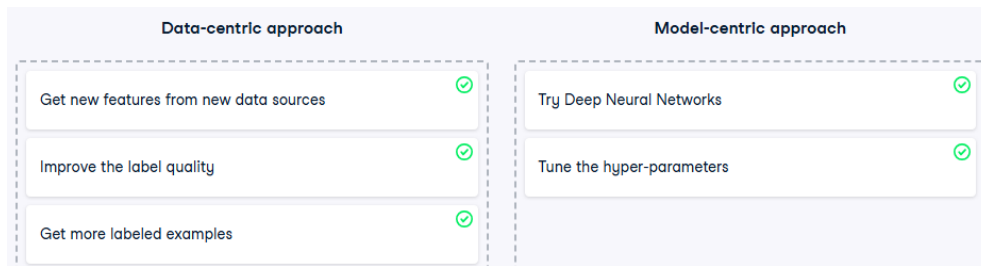
Once new labels are collected:

- Rebuild the training dataset
- Run the ML build pipeline
- Evaluate the new model

If performance improves → deploy. If not → continue iterating.

Tooling: Metadata Stores

- Use tools like **MLFlow Tracking** to:
 - Log experiments
 - Avoid redundant runs
 - Document model selection journey



Model Governance in MLOps

Model governance refers to the structured oversight of machine learning models throughout their lifecycle to ensure ethical, secure, and accountable deployment. As ML systems increasingly influence high-stakes decisions across industries, governance becomes essential to mitigate risks and uphold trust.

Why Model Governance Matters

- ML models are used to make **frequent, high-impact decisions** in sectors like finance, healthcare, and logistics.
- A flawed model (e.g., underestimating loan default risk) can lead to **catastrophic consequences**, including financial collapse and systemic disruption.
- Unlike humans, ML models can make **thousands of decisions per second**, amplifying the impact of any error.
- Therefore, **uncontrolled deployment** of models is unacceptable.

Definition

“Model governance is the overall process for how an organization controls access, implements policy, and tracks activity for models and their results.” — *DataRobot*

Effective governance minimizes:

- Financial risk
- Reputational damage
- Regulatory non-compliance

Governance Across the ML Lifecycle

1. Design Phase

- **Ethical Considerations:**

- When is it appropriate to use ML?
- Are we using private or sensitive data responsibly?
- How do we detect and mitigate bias?

2. Development Phase

- **Documentation and Traceability:**

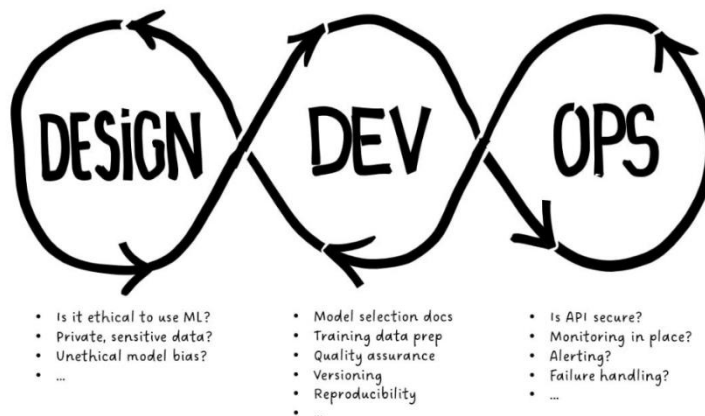
- Model selection rationale
- Data preparation and quality assurance
- Versioning of data and code
- Reproducibility of results

- **Tooling:** Metadata stores (e.g., MLflow Tracking) support traceability and auditability.

3. Pre-Production Phase

- **Security and Reliability Checks:**

- API security validation
- Monitoring and alerting systems in place
- Failure mode analysis and mitigation



- **Audit Readiness:**

- Define who performs each check
- Document outcomes for compliance and traceability

Regulatory Context

- **Self-Governance:** In low-risk domains, organizations may define their own governance standards.
- **Regulated Sectors:** In finance, healthcare, and critical infrastructure, governance is mandated by local and international laws.
- **Non-compliance:** Can lead to legal penalties, financial loss, and reputational harm.

Risk-Based Governance

Risk Level	Example Use Case	Governance Strictness
Low	Product recommendation engine	Minimal
Medium	Customer churn prediction	Moderate
High	Anti-money laundering, credit scoring	Strict

Risk is assessed based on:

- **Impact** of decisions
- **Frequency** of predictions
- **Financial and reputational stakes**

CRUX

- Governance may introduce **friction** into ML workflows, but it prevents **anarchy** and **reckless deployment**.
- The goal is not to deploy as many models as fast as possible, but to ensure they are **safe, ethical, and valuable**.
- As MLOps maturity increases, the **value of governance becomes more evident**, especially in large-scale, multi-model environments.

