

MLOps in Industrial Settings

Role of ML in Industry

- **Digital Transformation:** Modern companies are generating vast amounts of data through digitized operations.
 - **ML Applications:**
 - Building intelligent systems that analyze operational data.
 - Delivering predictions to support business goals.
 - Enhancing customer service (e.g., chatbots, recommendation engines).
 - Optimizing internal operations (e.g., predictive maintenance, inventory forecasting).
-

Example:

A logistics company uses ML to forecast delivery delays based on traffic, weather, and historical data, improving customer satisfaction and reducing operational costs.

2. Optimizing for Business Value

- **Profit Maximization:** ML use cases are selected based on their potential to generate profit.
 - **Portfolio Planning:** Companies may estimate ROI across multiple ML projects.
 - Example: Deploying 6 ML use cases could yield a projected profit of \$25M.
-

3. Cost Components in Software Development

ML systems are software solutions and inherit traditional development costs:

- **Development:** Coding, testing, integration.
 - **Project Management:** Planning, coordination, stakeholder communication.
 - **UI/UX Design:** User interface and experience optimization.
 - **Quality Assurance:** Testing, bug fixing, performance validation.
 - **Technical Debt:** Cost of future rework due to shortcuts or suboptimal decisions.
-

4. Understanding Technical Debt

- **Definition:** The future cost incurred by choosing a quick solution over a robust one.
- **Impact:** Leads to rework, instability, and reduced scalability.

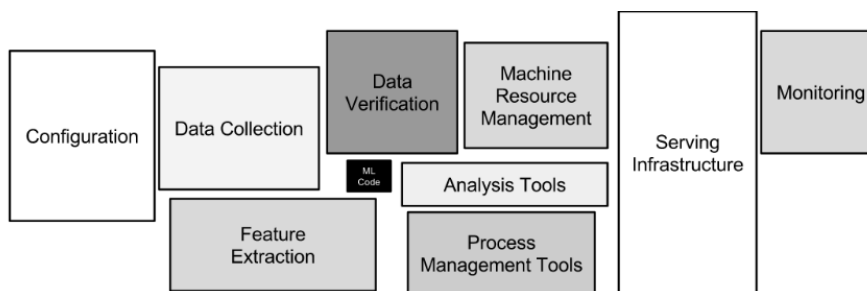
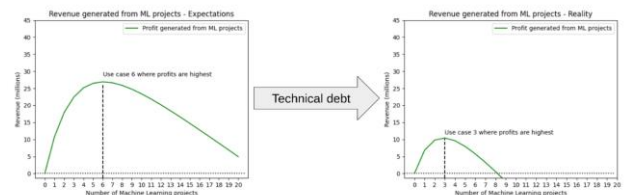


Analogy:

Choosing a shortcut in code today is like borrowing money—you'll pay interest later in the form of maintenance and refactoring.

5. Hidden Technical Debt in ML Systems

ML introduces unique forms of technical debt beyond traditional software:



• Data Dependency:

- ML behavior is learned from data, not explicitly coded.
- Poor data quality leads to unpredictable outcomes.

• Model Complexity:

- Lack of standardized practices for rigorous model development.

• Infrastructure Overhead:

- ML pipelines require complex orchestration (e.g., data ingestion, feature engineering, model serving).

• Monitoring Challenges:

- Continuous integration of new data demands robust monitoring.
- Real time prediction systems require uptime and accuracy tracking.

Quote:

“Machine Learning is the high interest credit card of technical debt.” — Google Research [【source】](#)

6. Business Impact of Hidden Costs

- **Profitability Risk:**

- Hidden technical debt can drastically reduce expected profits.
 - Example: A projected \$25M profit may shrink to \$10M due to unforeseen ML system costs.
-

7. MLOps: Strategic Debt Management

- **Definition:** MLOps (Machine Learning Operations) is the discipline of managing the ML lifecycle efficiently.

- **Goals:**

- Minimize hidden technical debt.
 - Ensure sustained business value from ML systems.
-

- **Components:**

- Data preparation
 - Model development
 - Testing and validation
 - Deployment and monitoring
-

Fully Automated MLOps

- **Automation Benefits:**

- Streamlined workflows
 - Reduced human error
 - Faster iteration and deployment
 - Reliable model updates
-

Example:

A retail company uses MLOps to automate demand forecasting. The pipeline retrains models weekly using fresh sales data, deploys updates seamlessly, and monitors prediction accuracy in real time.

Summary: Why MLOps Matters

Aspect	Without MLOps	With MLOps
Technical Debt	High, unpredictable	Minimized, manageable
Model Deployment	Manual, error prone	Automated, reliable
Business Value	Eroded by hidden costs	Sustained through efficiency
Monitoring	Reactive, fragmented	Proactive, integrated

1. Overview of the MLOps Lifecycle

The MLOps lifecycle consists of three interconnected stages:

- Designing the ML System
- ML Experimentation and Development
- ML Deployment and Operations

These stages form an iterative and incremental process. It is common to revisit earlier stages as new insights emerge or requirements evolve.

2. Stage 1: Designing the ML System

Key Activities

- **Business Understanding:**
 - Define business goals and success criteria.
 - Align ML objectives with strategic priorities.
- **Data Understanding:**
 - Explore and visualize available data.
 - Assess data quality, relevance, and alignment with business goals.
- **Solution Architecture:**
 - Design the ML system architecture.
 - Consider data security, privacy regulations, scalability, and maintainability.

Automation Opportunities

- Automate data quality checks using tools for schema validation, missing value detection, and profiling.

Human Collaboration

- Involve domain experts, business stakeholders, and developers.
- Use frameworks like CRISP DM to streamline manual processes and reduce errors.

Example

A healthcare startup designing a patient risk prediction system:

- Doctors define risk factors.
- Data scientists explore patient records.
- Developers build secure, compliant data pipelines.

3. Stage 2: ML Experimentation and Development

Key Activities

- **Proof of Concepts (PoCs):**
 - Rapid prototyping to validate feasibility.
- **Data Engineering:**
 - Clean, transform, and prepare data for modeling.
- **Model Development:**
 - Train, evaluate, and refine ML models.

Automation Opportunities

- Use templated frameworks for fast PoC iteration.
- Automate ETL pipelines with tools like Airflow or Prefect.
- Automate model development tasks:
 - Experiment tracking (e.g., MLflow)
 - Training pipelines
 - Hyperparameter tuning (e.g., Optuna, Ray Tune)

Example

An e-commerce company builds a recommendation engine:

- PoC validates feasibility using historical purchase data.
- Automated pipelines clean and prepare data.
- Models are trained and tuned using automated workflows.

4. Stage 3: ML Deployment and Operations

Key Activities

- Deploy ML systems to production.
- Implement testing, versioning, continuous delivery, and monitoring.

Automation Opportunities

- Automate deployment pipelines using CI/CD tools.
- Monitor model performance and data drift in real time.
- Enable frequent and reliable model updates.

Example

A financial institution deploys a fraud detection model:

- Automated tests validate model accuracy.
- CI/CD pipelines ensure smooth deployment.
- Monitoring tools track prediction reliability and alert on anomalies.

5. Building for Scale: Automation First

Principles

- Prioritize automation and process streamlining.
- Avoid hidden technical debt that limits scalability.
- Use reproducible frameworks like CRISP DM and Microsoft's TDSP.

Benefits

- Reduced errors
- Improved collaboration
- Reliable decision making
- Scalable ML systems

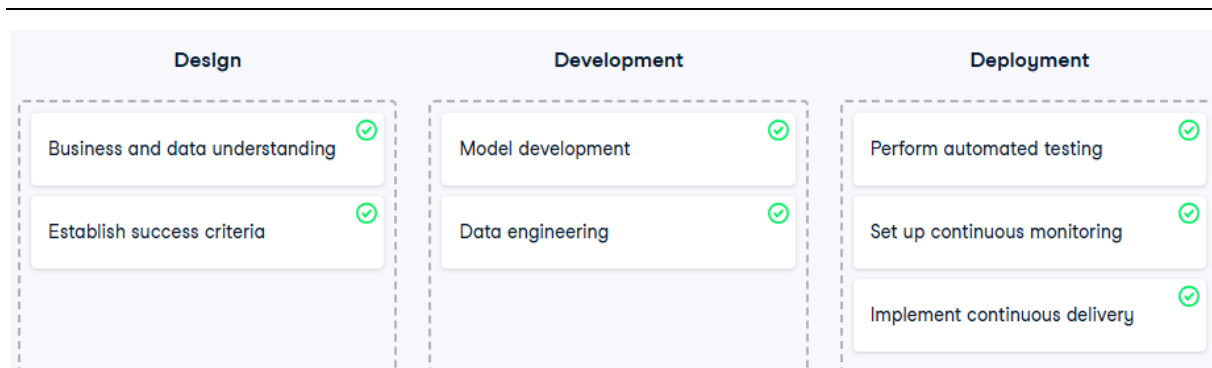
6. Best Practices

Design Phase

- Involve domain experts and business stakeholders early.
- Gather feedback from end users to refine requirements.

Development Phase

- Write clean, maintainable, and well organized code.
 - Maintain thorough documentation for reproducibility and auditability.
-



Reference architecture

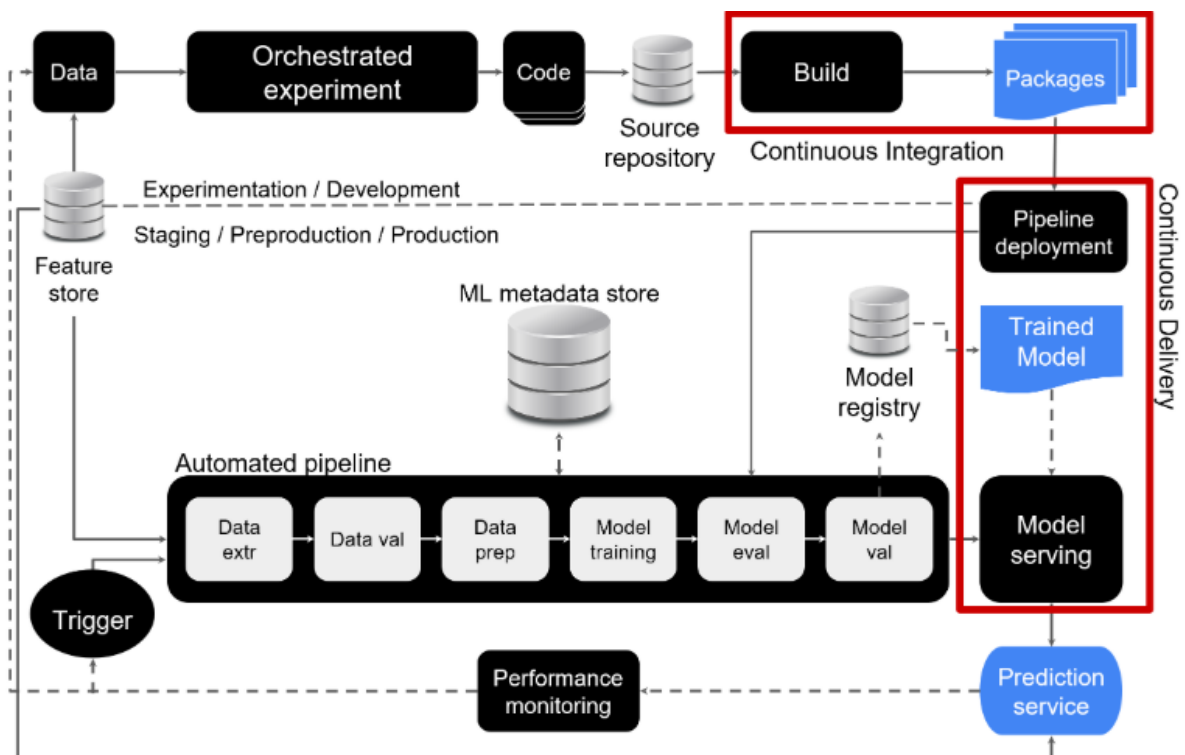
- A reference architecture is a **blueprint** for designing IT solutions, including ML systems.
- It defines **standard components, integrations, and design patterns** commonly used in industry.
- Benefits:
 - Leverages proven best practices.
 - Minimizes technical debt.
 - Promotes scalability, maintainability, and consistency.

2. Core Capabilities of Fully Automated MLOps

A fully automated MLOps system integrates the following capabilities:

- **Continuous Integration (CI):** Automated testing and packaging of source code.
- **Continuous Delivery (CD):** Automated deployment of ML artifacts to production.
- **Continuous Monitoring:** Real time tracking of model performance.
- **Continuous Training:** Automated retraining based on performance triggers.

3. Component Walkthrough of the Architecture



a. Orchestrated Experiments

- ML experiments are executed in a controlled development environment.
- Supports reproducibility and structured experimentation.

b. Source Code and CI

- Source code is pushed to a version controlled repository.
- CI pipeline builds the code and runs automated tests.
- Outputs include validated artifacts (e.g., packages, executables).

c. Artifacts and CD

- Artifacts from CI are deployed via CD pipelines.
- Ensures reliable and repeatable delivery to production environments.

d. ML Pipeline Deployment

- A new version of the ML pipeline is deployed with updated models.
- CI/CD mechanisms ensure automation and traceability.

e. Metadata Store

- Stores metadata from pipeline executions:

- Logs
- Hyperparameters
- Training configurations

f. Model Registry

- Central repository for storing trained models.
- Facilitates versioning, auditing, and rollback.

g. Prediction Services

- New models trigger automated deployment of prediction services.
- Enables **continuous model delivery**.

h. Continuous Monitoring

- Monitors prediction outputs and model performance metrics.
- Detects drift, anomalies, and degradation.

i. Automated Trigger

- Performance statistics can activate automated triggers.
- Triggers initiate pipeline execution for retraining or fallback.

j. Automated Retraining

- New models are trained automatically when triggered.
- Backup models or heuristics can be activated if needed.

4. Feature Store: A Critical Component

Purpose

- Provides **reusable and reproducible access to features** across the ML lifecycle.
- Ensures consistency between development and production environments.

Integration Points

- Feeds data to:
 - Orchestrated experiments
 - Automated ML pipelines
 - Prediction services

5. Summary of Architectural Flow

Stage	Description
Experimentation	Orchestrated ML experiments in development environment
Source Code & CI	Code pushed, built, and tested automatically
Artifacts & CD	Validated artifacts deployed to production
Pipeline Deployment	New ML pipeline version deployed with updated model
Metadata Store	Logs and training metadata recorded
Model Registry	Models versioned and stored
Prediction Services	Models deployed to serve predictions
Monitoring	Performance tracked continuously
Automated Trigger	Triggers retraining or fallback based on metrics
Feature Store	Centralized, consistent feature access across environments

☒ A reference architecture provides structures and integrations of ML elements and patterns commonly used when designing ML applications.

☒ Reference architectures provide a consistent and repeatable approach to building solutions.

☐ A reference architecture is a predefined set of rules that should not be modified.

☒ Reference architectures are designed to be customized for each specific project.

☐ Reference architectures are an experimental concept and still remain to be proven to work in the industry.

True	False
In the MLOps reference architecture, the feature store handles input data used by the ML workflows. ✓	An automated trigger component is used by the fully automated ML system to deliver predictions. ✓
The model registry in the architecture is used to manage the ML models produced by the automated workflows. ✓	A feature store is used in the architecture to store ML models produced during experimentation. ✓
The metadata store in the architecture handles the logs produced during the execution of the ML workflow. ✓	A prediction service is used in the architecture to store all features used in the execution of your ML pipelines. ✓

Why Experiment Tracking Matters in ML

- Machine learning is inherently **experimental**.
- ML engineers and data scientists frequently:
 - Apply different data transformations.
 - Try various algorithms.
 - Tune hyperparameters.
 - Evaluate models using different metrics.
- This creates a **vast space of possibilities** that must be tracked to avoid confusion, duplication, and wasted effort.

2. Challenges of Manual Tracking

- Manually logging every experiment detail is **impractical** and **error prone**.
- Without automation:
 - Experiments are hard to reproduce.
 - Results are difficult to compare.
 - Collaboration and sharing become inefficient.
 - Valuable insights may be lost or misinterpreted.

3. What Should Be Logged Automatically

Category	Examples
Code & Configuration	Source code, environment setup, config files
Data Specifications	Data source, type, format, cleaning steps, augmentation methods
Model Details	Architecture, parameters, hyperparameters
Evaluation Metrics	Accuracy, precision, recall, F1 score, AUC, loss
Execution Metadata	Timestamps, hardware used, training duration

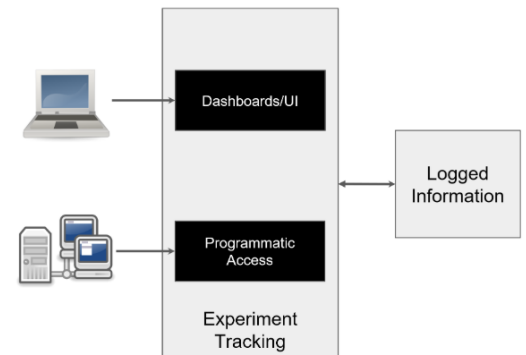
4. Goals of Logging

- **Reproducibility:** Enables others (and future you) to replicate experiments exactly.
- **Transparency:** Clarifies what was done, how, and why.
- **Performance Tracking:** Helps compare models and identify improvements.
- **Decision Support:** Informs model selection and deployment choices.

5. Automated Experiment Tracking Systems

Core Functions

- Organize and store metadata from ML experiments.
- Compare training runs and performance metrics.
- Reproduce experiments reliably.
- Interface via dashboards or APIs.

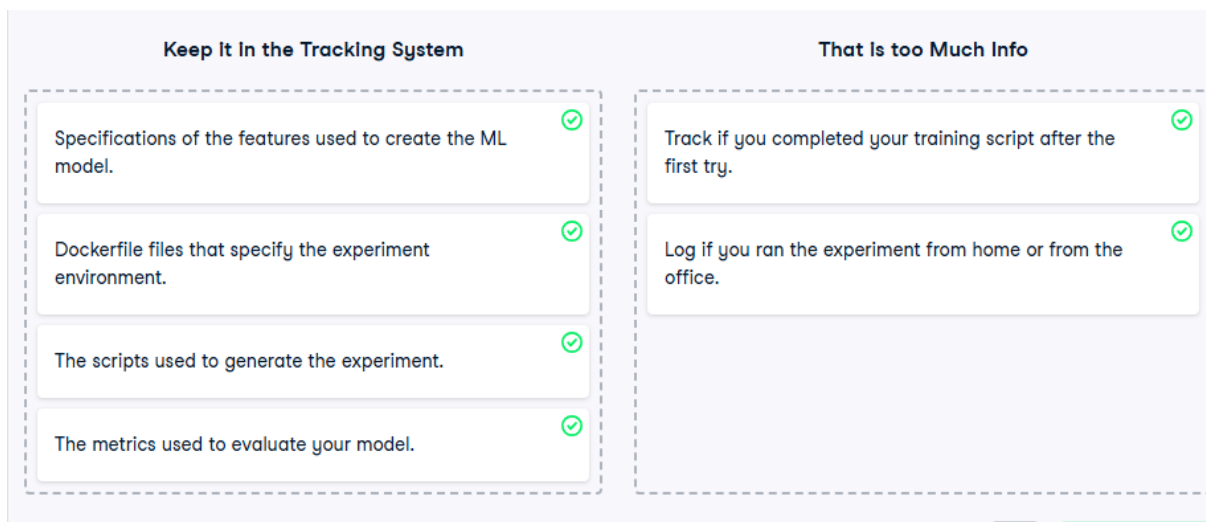


Integration with MLOps

- Supports model promotion to the **model registry**.
- Enables traceable and auditable model lifecycle management.

6. Tools Available in the Market

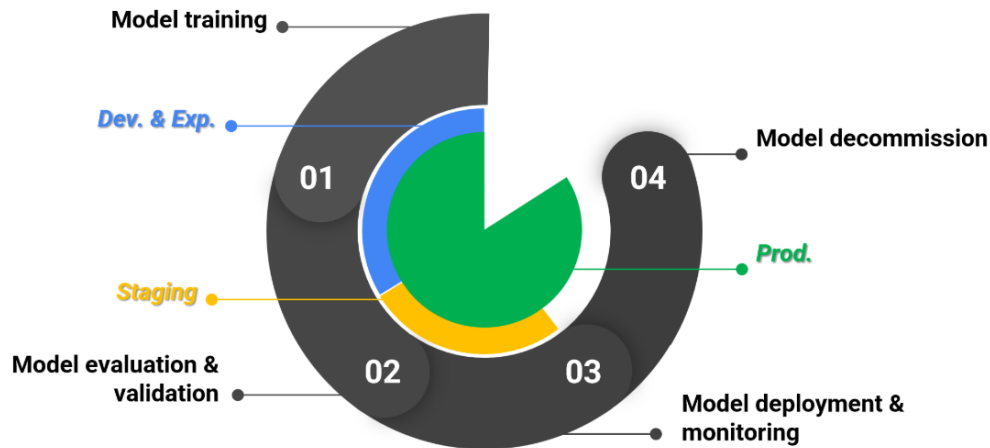
Tool Name	Type	Description
MLflow	Open source	Popular tool for experiment tracking, model registry, and deployment
Weights & Biases	Freemium	End to end platform with dashboards, collaboration, and cloud integration
Neptune.ai	Freemium	Focused on metadata management and team collaboration
Comet.ml	Freemium	Offers experiment tracking, model comparison, and visualizations
TensorBoard	Open source	Visualization tool for TensorFlow experiments



Role in the ML Model Lifecycle

- ML models follow a lifecycle from **creation to archiving**.
- Models may be produced via:
 - **Orchestrated manual pipelines** in development environments.

- **Automated pipelines** triggered by retraining mechanisms.



- Lifecycle stages typically include:
 - Development
 - Staging
 - Production

Environment Transition

- Models must be validated and tested in staging before production deployment.
- Managing these transitions requires a centralized and structured approach.

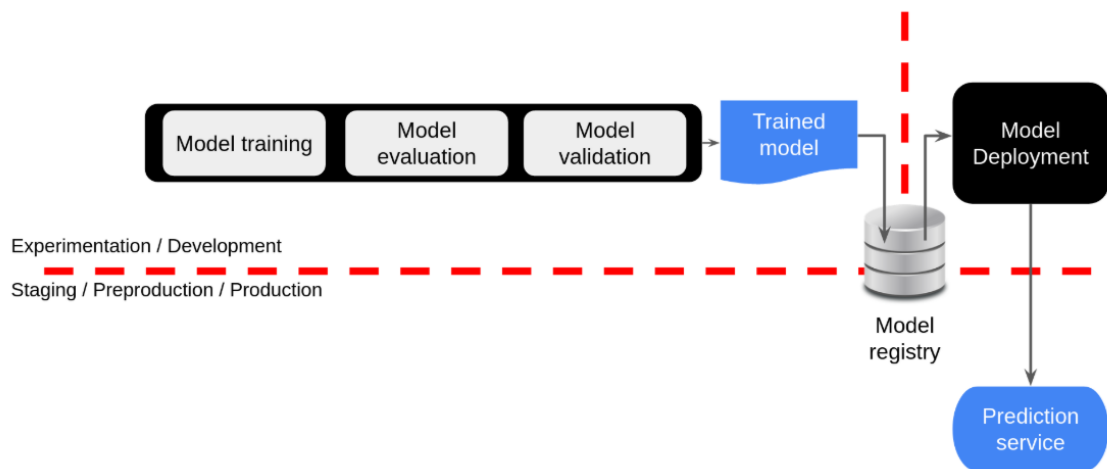
2. The Problem: "Throwing Models Over the Fence"

- In traditional workflows, a data scientist trains a model and hands it off informally to operations (e.g., via email or USB).
- This leads to:
 - Lack of reproducibility
 - Poor traceability
 - Deployment errors
 - Collaboration breakdown between Dev and Ops

3. First Step Toward Automation: Ad hoc Model Registry

- Data scientists identify the best performing model after experimentation.

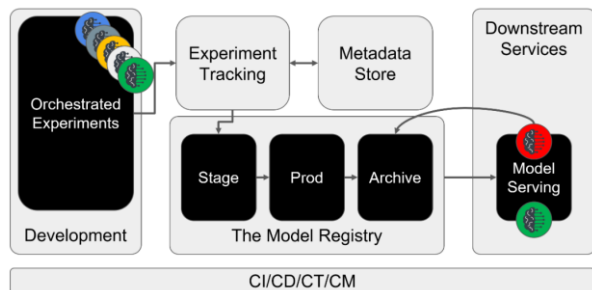
- They **register** the model in a central repository.



- Registration includes:
 - Model artifacts (e.g., weights, architecture)
 - Metadata required for deployment
- This signals readiness for staging and production testing.

4. What Is a Model Registry?

- A **model registry** is a central repository for managing production ready models.
- It is a core component of the MLOps architecture.



Functions

- Publishes and stores models and their deployment artifacts.
- Manages transitions between environments.
- Integrates with CI/CD workflows and staging environments.

5. Integration with Experiment Tracking

Component	Role
Experiment Tracking	Logs metadata from experimentation runs
Metadata Store	Stores hyperparameters, configurations, and training logs

Component	Role
Model Registry	Promotes selected models for deployment and lifecycle management <ul style="list-style-type: none"> The model registry does not store experiment metadata directly (in this design). Alternatively, it can be designed to unify experiment tracking and metadata storage.

6. Registering and Deploying Models

Steps

1. Select the best performing model from tracked experiments.
2. Register the model in the model registry.
3. Trigger automated tests in staging via CI/CD integration.
4. Upon successful validation, deploy the model to production.
5. Update downstream prediction services with the new model.

7. Model Decommissioning and Archiving

- The previous production model is **decommissioned** and **archived** in the model registry.
- This ensures:
 - Version control
 - Auditability
 - Rollback capability

8. Summary of Model Registry Functions

Function	Description
Centralized Storage	Stores models and deployment artifacts
Lifecycle Management	Tracks model status across environments
CI/CD Integration	Automates testing and deployment workflows
Collaboration Hub	Bridges Dev and Ops teams
Archiving and Versioning	Maintains historical models for rollback and audit

True	False
Typically, after experimenting and finding the model with the best performance, you can register it into the model registry. ✓	The model registry is responsible for managing the features you engineer in your system. ✓
After registering a model, the model registry can trigger automated tests to ensure that the model can be properly deployed to production. ✓	After executing plenty of experiments, the model registry allows you to compare your experiments. ✓
The model registry is a central repository that allows you to publish production-ready models. ✓	Metadata about the experiments is stored in the model registry itself. ✓

What Is Feature Engineering?

Feature engineering involves transforming raw data into meaningful inputs for machine learning models. Common transformations include:

- **Standardizing or recoding categorical variables** (e.g., one hot encoding)
- **Grouping values** (e.g., binning age ranges)
- **Creating new features** using domain knowledge (e.g., risk scores, ratios)

2. Feature Engineering in Enterprise Settings

Scenario

In a typical enterprise:

- Multiple teams develop ML models using similar internal data sources.
- Each team applies its own transformations and stores features independently.

Challenges

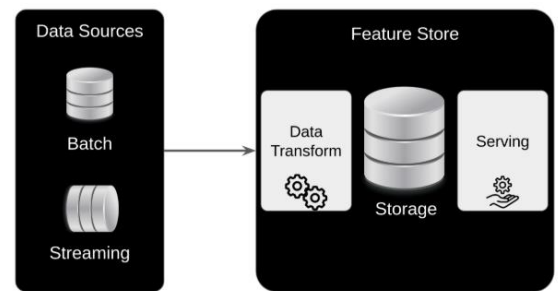
- **Duplication of effort:** Teams recreate similar features from the same data.
- **Inconsistent definitions:** Similar features may be defined differently across teams.
- **Lack of discoverability:** No centralized way to explore or reuse existing features.

3. The Feature Store: Centralized Feature Management

A **feature store** is a centralized repository that standardizes the definition, storage, and access of features across the ML lifecycle.

Core Functions

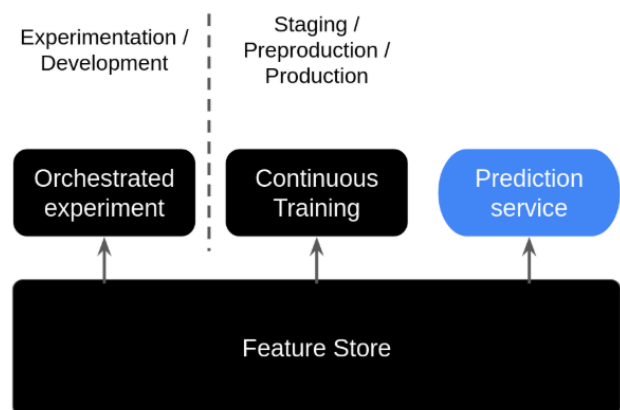
- **Ingest raw data** from batch or streaming sources.
- **Apply transformations** to generate features.
- **Store features** in a centralized, versioned repository.
- **Serve features** via APIs for:
 - High throughput batch processing
 - Low latency real time inference



4. Benefits of Using a Feature Store

a. Accelerated Experimentation

- Data scientists can extract pre engineered features for experimentation.
- Enables **reuse of feature sets** across teams and projects.
- Maintains **metadata** such as:
 - Feature definitions
 - Data versioning
 - Lineage tracking



b. Continuous Training

- Automated ML pipelines can retrieve **up to date feature values** for training tasks.
- Ensures consistency and freshness of training data.

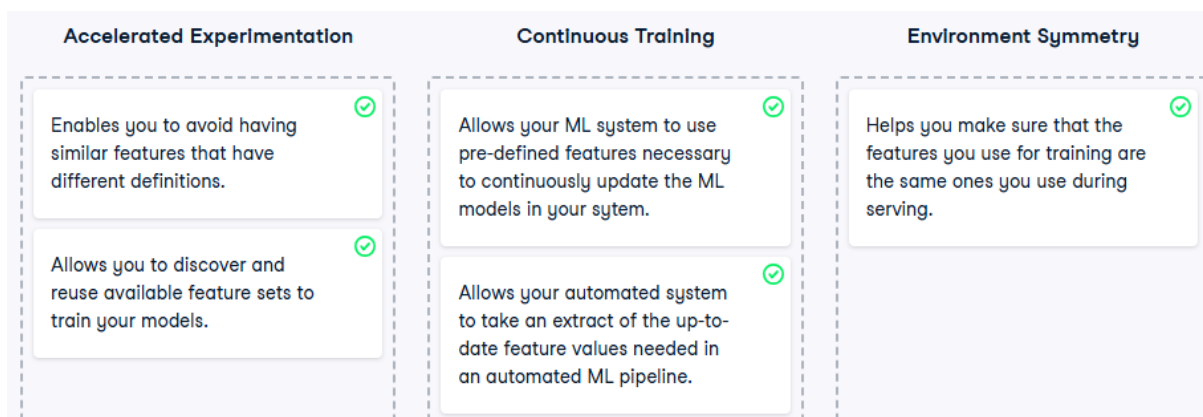
c. Online Predictions

- Prediction services can access relevant feature batches for real time inference.
- Examples include:
 - Customer segmentation features
 - Geographical attributes

- Domain specific transformations

d. Environment Symmetry

- Prevents **data skew** between training and serving environments.
- Ensures identical feature definitions and values across:
 - Experimentation
 - Training
 - Production serving



5. Summary of Feature Store Capabilities

Capability	Description
Centralized Storage	Stores features in a reusable and standardized format
Transformation Standardization	Automates and unifies feature engineering across teams
Metadata Management	Tracks feature lineage, versions, and definitions
Batch & Real time Serving	Supports both high throughput and low latency access
Environment Symmetry	Aligns training and serving environments to prevent data skew

1. What Is Metadata in MLOps?

Metadata refers to information about artifacts generated during the execution of ML pipelines. It includes:

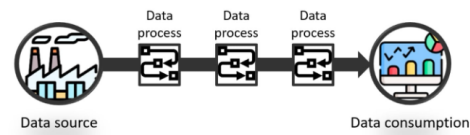
Artifact Type	Metadata Examples
Data Sources	Data version, creation/modification logs, transformation history
Model Training	Model type, version, hyperparameters, evaluation metrics
Pipeline Execution	Execution logs, timestamps, hardware utilization

This metadata is essential for tracking, auditing, and automating ML workflows.

2. Key Aspects of Metadata in ML Systems

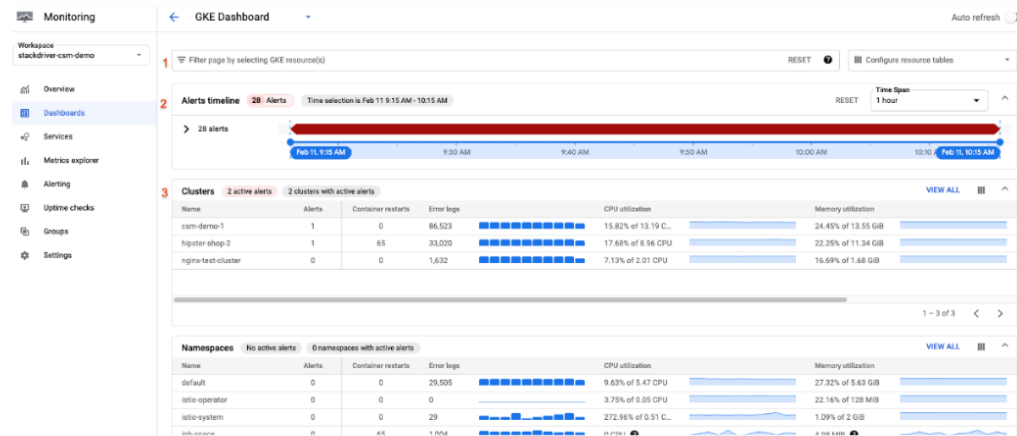
a. Data Lineage

- Tracks the lifecycle of data from creation to consumption.
- Enables visibility into transformations, aggregations, and usage.
- Example: Customer data lineage includes registration, transformation steps, and consumption points.



b. Reproducibility

Example monitoring tool



¹ <https://cloud.google.com/stackdriver/docs/solutions/gke/observing>

- Metadata ensures experiments can be replicated.
- Includes hyperparameter settings, model configurations, and training conditions.
- Builds trust and scientific rigor into ML workflows.

c. Monitoring

- Metadata enables real time tracking of pipeline execution status.
- Supports system health checks, anomaly detection, and performance diagnostics.

d. Automation Support

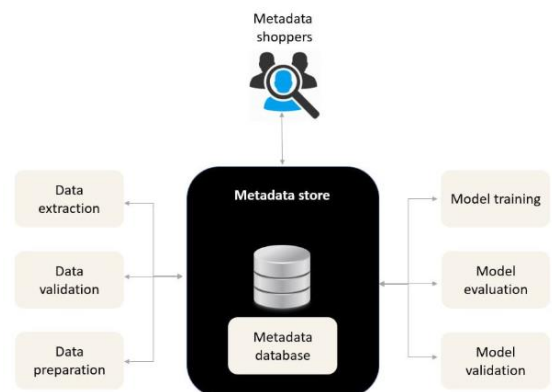
- Facilitates automated responses to system events (e.g., drift detection, error recovery).
- Enables retraining and rollback mechanisms based on logged metadata.

3. What Is a Metadata Store?

A metadata store is a centralized repository for managing metadata related to ML experiments, models, and pipelines.

Characteristics

- Stores metadata, not actual data or model artifacts.
- Provides interfaces (UI and API) for reading/writing metadata.
- Integrates with:
 - **Model Registry** (stores actual models)
 - **Feature Store** (handles actual data)



4. Role in MLOps Architecture

Component	Function
Metadata Store	Logs and manages metadata across all pipeline stages
Automated Pipeline	Reads/writes metadata during execution
Monitoring System	Uses metadata to track performance and trigger automated responses

5. Metadata Store in Fully Automated MLOps

Capabilities

- Tracks system functioning through execution logs.
- Enables:

- Continuous performance monitoring
- Automated retraining on performance decay
- Rollbacks on critical errors
- Supports uninterrupted operation until root cause analysis is completed.

Example Use Cases

- **Drift Detection:** Evaluation metrics indicate model drift -> retraining triggered.
- **Error Recovery:** Pipeline failure detected -> rollback to last stable model/data version.

You start with a fully automated MLOps system that can deliver good predictions via a prediction service and this is being constantly monitored. ✓

As your prediction service is functioning, evaluation metrics are logged to the metadata store. ✓

The evaluation metrics start to become worse over time, your system automatically detects a decay in the performance of your model. ✓

The trigger in your system starts the automated training pipeline and after its successful completion, pushes an updated model to the registry. ✓

After an updated model lands in the model registry, this is automatically deployed and the prediction service is updated to account for the drift. ✓

True	False
<p>Metadata contains information about the artifacts created during the ML workflow.</p>	<p>The metadata store allows machine learning engineers to publish prod-ready models.</p>
<p>Data versioning, the information about the data used in a training pipeline is saved to the metadata store.</p>	<p>The metadata store is a repository used to save and version trained machine learning models.</p>
<p>The metadata store interacts intensively with all the steps in the workflow, it stores logs as part of its automatic operation.</p>	<p>The metadata store is used to manage all feature-related resources (training, serving, etc.) across the ML workflow.</p>

1. Automation Across the MLOps Lifecycle

The MLOps lifecycle consists of three key stages:

- **Design**
- **Development**
- **Operations**

Each stage presents different opportunities and constraints for automation.

2. Maturity Levels of MLOps Automation

Maturity Level Description

Manual	Ad hoc experimentation, no tracking, manual handoff to operations
Semi Automated	Orchestrated experimentation, partial automation of pipeline components
Fully Automated	Continuous delivery of pipelines, automated monitoring, retraining, and deployment

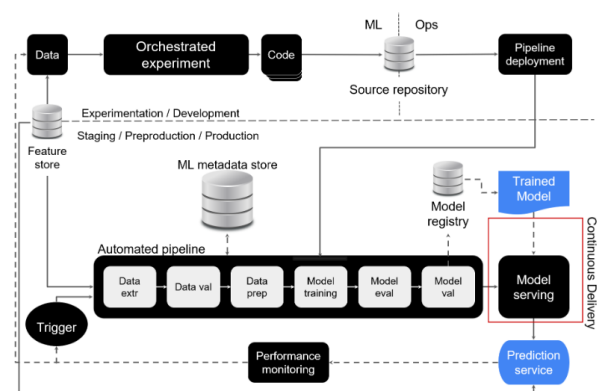
3. Manual ML Workflow

- Initial stage for most teams.
- ML experiments are run manually without tracking or logging.
- Models are handed off informally from development to operations.
- Reproducibility and traceability are limited.

4. Semi Automated ML Workflow

Key Enhancements

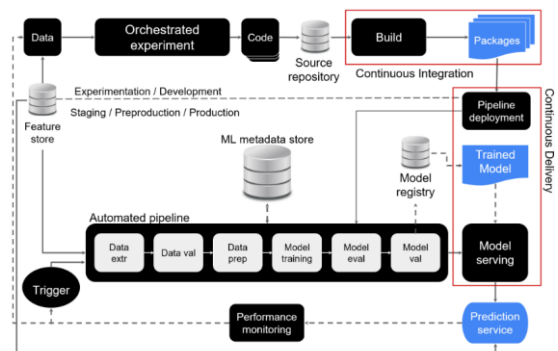
- **Orchestrated experimentation:** Structured and repeatable experiment pipelines.
- **Source code repository:** Centralized version control for ML code.
- **Feature store:** Reusable and standardized feature access across experiments.



- **Automated pipeline execution:** Triggered by code updates.
- **Metadata store:** Logs pipeline execution metadata.
- **Model registry:** Stores trained models and deployment artifacts.
- **Continuous Delivery:** Automates deployment of updated prediction services.
- **Monitoring:** Tracks prediction performance and triggers retraining.

5. Fully Automated ML Workflow

- Builds on semi automated architecture.
- Integrates automated pipeline deployment into Continuous Delivery.
- Enables continuous updates to prediction services using the latest trained models.
- Supports:
 - Faster deployment cycles
 - Reliable model updates
 - Scalable and efficient ML workflows



6. Automation in the ML Lifecycle

a. Design Phase

- Cannot be fully automated due to reliance on domain expertise and stakeholder input.
- Use reproducible frameworks like:
 - **Team Data Science Process (TDSP)**
 - Document problem definitions, success metrics, and design decisions.

b. Development Phase

- Manual effort is required, but many tasks can be automated:
 - Data preprocessing
 - Model selection
 - Hyperparameter tuning
- Best practices:

- Write clean, maintainable code.
- Use version control (e.g., Git).
- Apply orchestration tools (e.g., Airflow, Kubeflow).

c. Operations Phase

- Automation is critical and highly effective.
- Key practices:
 - Automated testing of data, models, and workflows.
 - Continuous Integration (CI)
 - Continuous Deployment (CD)
 - Continuous Training (CT)
 - Continuous Monitoring (CM)

7. Summary of Automation Opportunities

Lifecycle Stage	Automation Potential	Best Practices
Design	Low	Use reproducible frameworks, document thoroughly
Development	Medium	Automate preprocessing, tuning, use version control
Operations	High	Automate testing, deployment, monitoring, and retraining

Manual	Semi-automated	Fully automated
<p>A team of data scientists develops a state-of-the-art model, this model is delivered to an engineering team to take it to production.</p> <p>At this maturity level, the data science team manages a few models with only a couple of deployments per year.</p>	<p>Here, several workflow steps are automated except for the deployment of new pipeline implementations.</p> <p>At this maturity level, a good degree of automation is used; the testing of pipelines and their components is manual.</p>	<p>Here, ML in production refers to deploying an ML pipeline with seamless, retraining and deployment of new models end-to-end, without manual effort.</p> <p>At this level, team members advocate for automation and monitoring at all steps of ML system construction.</p>

What Is a Design Pattern?

- A design pattern is a **reusable solution** to a recurring problem in a specific domain.
- In software development, it provides a **template** for solving common challenges.
- In MLOps, design patterns help standardize solutions for reliability, scalability, and maintainability of ML systems.

2. Overview of the Pattern: Automate, Monitor, Respond

This pattern is essential for building robust MLOps systems. It includes:

- **Automation:** Deployment, operation, and maintenance of ML models and pipelines.
- **Monitoring:** Continuous tracking of model performance and system health.
- **Incident Response:** Fail safe mechanisms to recover from performance decay or data issues.

3. Key Examples of the Pattern in MLOps

1. Automated Model Retraining

Problem

- ML models degrade over time due to data drift or changing conditions.
- Declining performance can impact business outcomes.

Solution

- Use **continuous training** to maintain model performance.

Workflow

1. **Prediction Service:** Continuously serves predictions.
2. **Monitoring:** Logs performance metrics (e.g., accuracy, latency).
3. **Trigger:** Activated when performance drops below a predefined threshold.
4. **Automated Pipeline:**
 - Extracts fresh data from the feature store.

- Retrains the model.
- Registers the updated model.

5. **Deployment:** Automatically updates the prediction service with the new model.

2. Model Rollback

Problem

- Retrained models may fail validation or perform worse than expected.
- Failed updates can disrupt production services.

Solution

- Implement **model rollback** to restore the last known good model.

Workflow

1. **Validation Fail:** New model fails evaluation or testing.
2. **Rollback Trigger:** System identifies failure and activates rollback.
3. **Redeployment:** Last functional model is redeployed to restore service.
4. **Decision Criteria:**
 - Severity of performance issues.
 - Availability of fresh training data.

3. Feature Imputation

Problem

- Data quality issues (e.g., missing values) can degrade model performance.
- Manual intervention is inefficient and error prone.

Solution

- Automate **feature imputation** based on predefined thresholds.

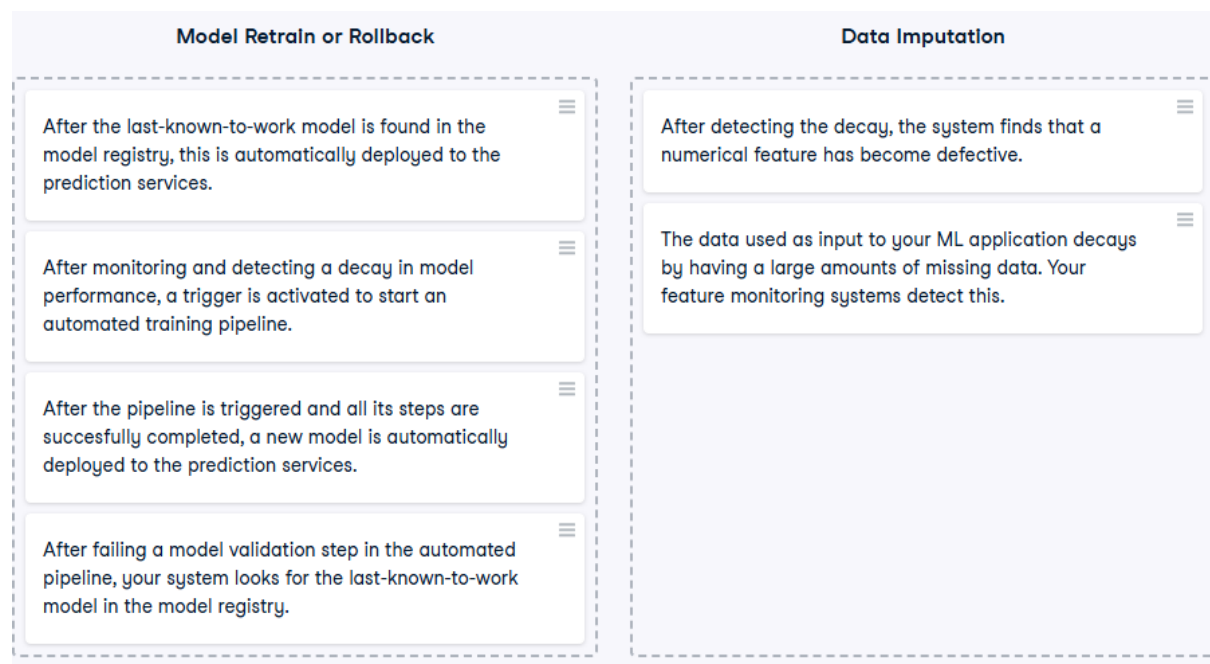
Workflow

1. **Quality Threshold:** Define acceptable limits (e.g., max 30% missing values).
2. **Monitoring:** Detect features exceeding the threshold.
3. **Trigger:** Alarm activates incident response.
4. **Imputation Methods:**
 - **Numerical Features:**

- Mean or median imputation
- KNN imputation
- **Categorical Features:**
 - Frequent category imputation
 - Add a "missing" category

4. Summary of the Pattern

Component	Function
Automation	Streamlines deployment, training, and updates
Monitoring	Tracks performance, data quality, and system health
Incident Response	Enables retraining, rollback, and data imputation
Benefits	Improves reliability, reduces downtime, and supports scalability



What Is Software Testing?

Software testing is the process of evaluating and verifying that a software product or application performs as intended. It ensures reliability, correctness, and performance.

Common Types of Software Tests

Test Type	Purpose
Unit Test	Validates individual components or functions of the application
Integration Test	Checks interactions between multiple components
End to End Test	Verifies the entire system behaves as expected from start to finish

2. Unique Nature of ML Software

Unlike traditional software, ML systems:

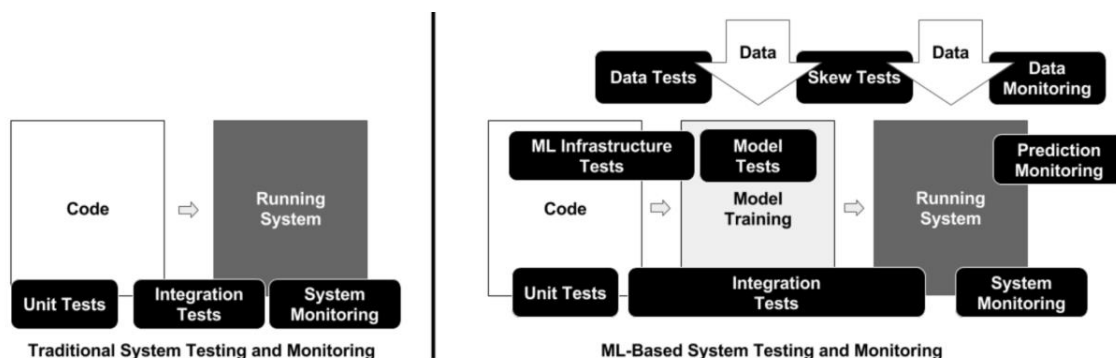
- Do not have explicitly coded behavior.
- Learn behavior from training data.
- Depend heavily on:
 - Data quality and distribution
 - Model architecture and parameters

This makes testing ML systems more complex and data dependent.

3. Testing in MLOps Systems

In addition to traditional software testing, MLOps systems require:

Test Category	Focus Area
Data Tests	Validating feature quality, distribution, and compliance
Model Tests	Evaluating model performance, generalization, and freshness
Infrastructure Tests	Ensuring pipeline reproducibility, integration, and operational stability



4. Data Testing in ML Systems

Key Practices

- **Deterministic Tests:** Validate feature values fall within expected ranges.
 - Example: Temperature values should be within realistic bounds.
- **Statistical Tests:** Check if features follow expected distributions.
- **Feature Value Assessment:** Use feature importance to justify computational cost.
- **Privacy and Compliance:** Ensure lawful use and protection of sensitive data.

5. Model Testing in ML Systems

Key Practices

- **End User Impact:** Ensure statistical improvements translate to better user experience.
- **Hyperparameter Tuning:** Test optimal configurations for performance gains.
- **Validation Techniques:** Use cross validation and holdout sets to avoid overfitting.
- **Model Freshness:**
 - Assess staleness and its impact on predictions.
 - Define update frequency based on performance decay.
- **Baseline Comparison:** Regularly test against simple models to validate complexity benefits.

6. Pipeline Testing in ML Systems

Key Practices

- **Workflow Validation:** Ensure complex pipelines execute correctly across components.
- **Reproducibility:** Training on the same data should yield the same model.
- **Integration Testing:** Include both data and model components.
- **Debugging Support:** Enable step by step inspection during training or inference.

7. Summary of Testing Scope in MLOps

Area Key Tests and Goals

Software Unit, integration, and end to end tests for application logic

Data Range checks, distribution tests, feature importance, privacy compliance

Models Accuracy, hyperparameter tuning, freshness, baseline comparison

Pipelines Reproducibility, integration, step by step debugging

Select the options that list tests related exclusively to data used in ML systems. Remember, these data tests should be specific to ML-powered applications.

✓ Answer the question 50XP

Possible Answers
Select all correct answers

☒ Feature expectations PRESS 1

☒ Value from added data justify its costs PRESS 2

☐ Business & ML metrics correlate PRESS 3

☒ Privacy controls are put in place PRESS 4

☐ All hyperparameters have been tuned PRESS 5

Data Tests	Model Tests	Pipeline Tests
<div>Tests to check that the input for your training pipelines fall within predefined, reasonable ranges. ✓</div> <div>Tests to ensure that the inputs used during model training exhibit expected statistical patterns, such as conforming to a specific distribution. ✓</div>	<div>Tests to ensure that all adjustable hyperparameters for training your machine learning systems have been optimized. ✓</div> <div>Tests that ensure the correlation between ML performance metrics and desired business outcomes. ✓</div>	<div>Tests to ensure that repeated execution of an end-to-end training process, with all variables held constant, result in predictable outcomes. ✓</div> <div>Tests to verify that the sequence of steps in an end-to-end ML workflow are properly integrated. ✓</div>

What Is a Hyperparameter?

- A **hyperparameter** is a tunable configuration value set **before** training a machine learning model.

- Unlike model parameters (e.g., weights, biases), which are **learned from data**, hyperparameters are **externally defined** and influence how the model learns.
- Hyperparameters significantly affect model performance and must be carefully selected.

Common Examples

Model Type	Hyperparameters Examples
Neural Networks	Number of layers, number of neurons, learning rate
Decision Trees	Max depth, number of branches, min samples split
Gradient Boosting	Learning rate, number of estimators, max depth
k Nearest Neighbors	Number of neighbors (k), distance metric

2. What Is Hyperparameter Tuning?

- **Hyperparameter tuning** is the process of selecting the best combination of hyperparameters to optimize model performance.
- It is a critical step in improving model accuracy and generalization.
- Tuning is typically done **before training** and is often the first step in ML experimentation.

3. Hyperparameter Tuning Methods

Method	Description
Grid Search	Exhaustively evaluates all combinations in a predefined parameter grid
Random Search	Randomly samples combinations from the parameter space
Bayesian Search	Uses probabilistic models to guide the search based on past performance

4. Automating Hyperparameter Tuning

Why Automate?

- The search space for hyperparameters can be **large and complex**.
- Manual tuning is **time consuming** and **inefficient**.
- Automation improves **scalability**, **reproducibility**, and **performance optimization**.

Steps to Automate

1. **Specify hyperparameters** to tune (e.g., learning rate, number of layers).
2. **Define search space:**
 - Discrete values (e.g., [0.01, 0.1, 1])
 - Continuous ranges (e.g., learning rate \in [0.0001, 0.1])
3. **Select performance metric** to optimize (e.g., accuracy, F1 score, AUC).
4. **Set stopping criteria:**
 - Maximum number of trials
 - Time budget
 - Early stopping rules

5. Output of Automated Tuning

- The automated tuning process returns the **best performing set of hyperparameters** from the defined search space.
- These hyperparameters are then used to train the final model.

6. Environment Symmetry in MLOps

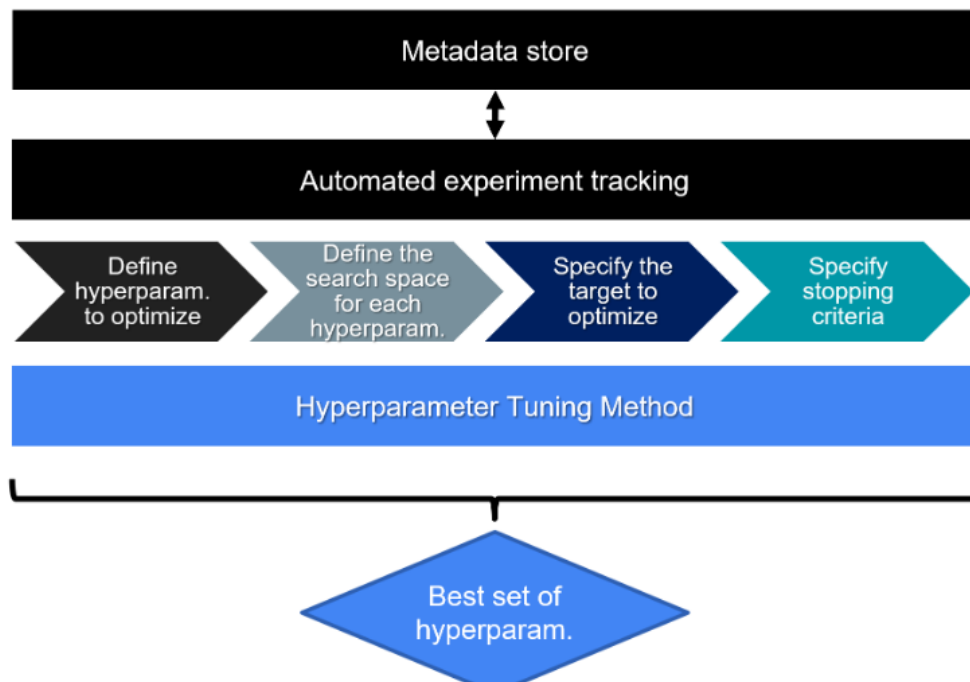
- In MLOps, it is essential to maintain **symmetry across development, staging, and production environments**.
- After tuning, the selected hyperparameters must be **consistently applied** in all environments to ensure reproducible and stable model behavior.

7. Experiment Tracking and Metadata Logging

- Hyperparameter tuning should be **tracked and logged automatically**.
- The **experiment tracking module** records:
 - All hyperparameter configurations tested
 - Corresponding performance metrics
 - Metadata such as timestamps, hardware used, and trial outcomes
- This metadata is stored in the **metadata store** for reproducibility and auditability.

8. Visualization of Hyperparameter Experiments

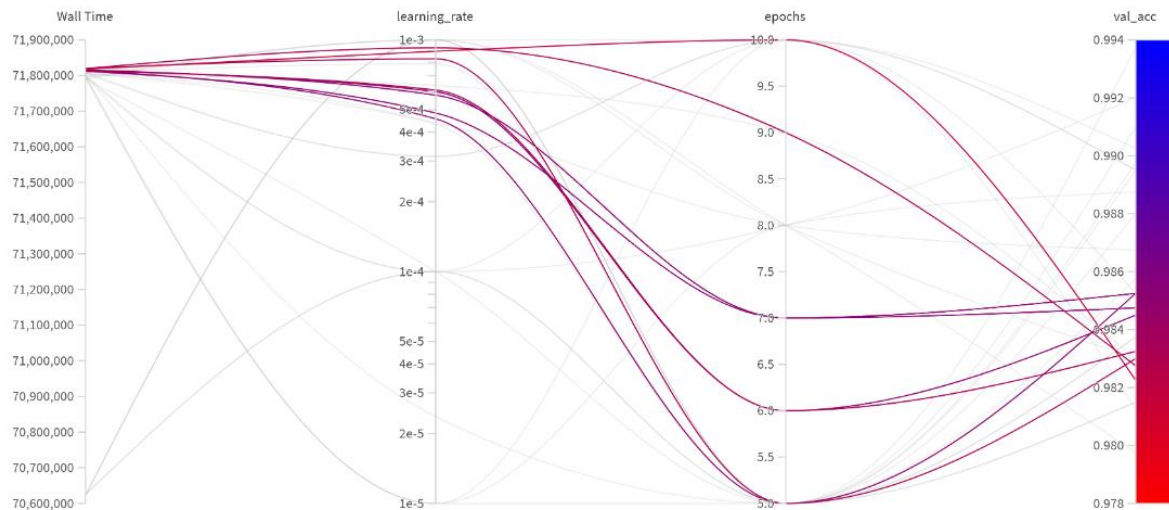
- Many experiment tracking tools (e.g., MLflow, Weights & Biases) offer



visualizations to analyze:

- How different hyperparameter values affect model performance
- Trends and interactions between hyperparameters

-



- These visualizations support **informed decision making** and **model interpretability**.

Specify which hyperparameters you want the tuning algorithm to search over and optimize.

Specify the possible values or range of values for each hyperparameter; this is known as defining a search space.

To determine the best performance by a set of hyperparameters, specify a metric that you want to optimize for.

Set a maximum number of trials, a maximum time limit, or other stopping method; these can be used as stopping criteria.

After the stopping criteria are met, automatically chose the set of hyperparameters with the best performance.

True	False
<div>Proper record-keeping and streamlining the process of choosing hyperparameters is crucial in a well-organized MLOps workflow.</div>	<div>Logging hyperparameters during hyperparameter tuning is not necessary as the results can be easily remembered and recorded manually.</div>
<div>When performing hyperparameter tuning in an MLOps system, it is important that the environments (dev/stag/prod) are identical.</div>	<div>Automated hyperparameter tuning guarantees finding the best set of hyperparameters, even if these fall outside the search space.</div>
	<div>Manually tuning hyperparameters leads to more accurate results and is preferred by experienced MLOps practitioners.</div>

Role of ML Pipelines in MLOps

- ML pipelines are foundational to **automating and operationalizing** machine learning workflows.
- In fully automated MLOps architectures, pipelines are present in both:
 - Development & Experimentation** environments
 - Production** environments

- They enable consistent, scalable, and repeatable ML processes across the lifecycle.

2. Modularity and Reusability

Characteristics

- Pipelines are designed in a **modular fashion**, where each step (e.g., data preprocessing, model training, evaluation) is an independent component.
- This modularity allows:
 - Independent development and maintenance of components
 - Easy updates and improvements without affecting the entire pipeline
 - Reusability across different datasets and ML tasks

CI/CD Integration

- Continuous Integration and Continuous Delivery (CI/CD) systems promote changes from development to production automatically and reliably.

3. Orchestration and Automation

Development & Experimentation

- Pipelines serve as a **sandbox** for testing new models and algorithms.
- Allow rapid iteration and experimentation by data scientists.
- Can be triggered manually or automatically based on code or data changes.

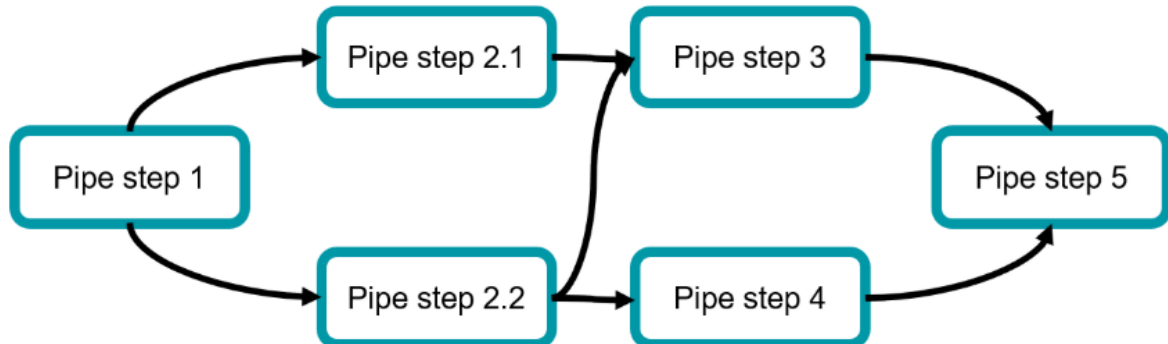
Production

- Pipelines automate the **deployment** of trained models into production environments.
- Triggered by events such as:
 - New data availability
 - Model performance degradation
 - Scheduled retraining cycles
- Benefits:
 - Reduces manual errors
 - Ensures models are always current and performant
 - Streamlines the deployment process

4. Direct Acyclic Graphs (DAGs) in MLOps

Definition

- A **DAG** is a graphical representation of the ML pipeline:



- **Nodes** represent individual tasks or steps.
- **Edges** represent dependencies between tasks.

Benefits

- Provides a **clear visualization** of the pipeline structure.
- Helps manage **complex workflows** in a controlled and reproducible manner.
- Enables debugging, optimization, and documentation of pipeline logic.

5. What Is Orchestration in MLOps?

- **Pipeline orchestration** refers to the management and automation of task execution within an ML pipeline.
- Responsibilities include:
 - Task scheduling and coordination
 - Monitoring execution status
 - Managing data dependencies
 - Ensuring correct data flow between stages

6. ML Pipelines in Development & Experimentation

Functions

- Orchestrate the full training and evaluation workflow.
- Ensure tasks are executed in the correct order.
- Record and track results for reproducibility.
- Enable **parallel execution** to accelerate experimentation.

Outcome

- Faster model development cycles
- Efficient experimentation
- Easier transition to production

7. ML Pipelines in Production

Functions

- Automate the deployment of new models.
- Manage and execute deployment steps consistently.
- Reduce human error and ensure model freshness.
- Provide centralized monitoring and control over deployment workflows.

Benefits

- Reliable and repeatable model deployment
- Easier troubleshooting and rollback
- Scalable production ML operations

8. Summary of ML Pipeline Capabilities in MLOps

Capability	Development & Experimentation	Production
Modularity	Independent component development	Reuse of validated components
Orchestration	Task sequencing and parallel execution	Automated deployment and monitoring
Reproducibility	Tracked experiments and results	Consistent model behavior across updates
Automation	Triggered by code/data changes	Triggered by events or schedules
Monitoring & Control	Experiment tracking	Centralized deployment oversight

Provides a unified environment for creating new models, allows easy access to required data and tools necessary for experiments.	Provides a framework for automating the rolling back of models in the event of issues.
Provides tools for distributing models across multiple nodes, enabling them to scale model training and evaluation.	Enables the automation of the ML re-training pipeline, ensuring that models are continuously updated.
Enables the automation of testing and evaluation of models, enabling engineers to evaluate models and compare their performance automatically.	Provides tools for managing resources such as compute, storage, and network bandwidth, ensuring the operability of the automated ML system.

True	False
DAGs are a valuable tool for defining and managing the flow of tasks in ML pipelines and improving the reproducibility of results.	DAGs are not necessary for effective orchestration in MLOps and can be replaced with other scalable manual methods.
Orchestration is essential in MLOps for ensuring the automation, scalability, and reliability of ML pipelines	The use of DAGs in MLOps only serves a visual purpose and provides no additional benefits.
Orchestration is important in both the development and experimentation phase and the automated production environment of MLOps.	
DAGs play a crucial role in visualizing and managing complex ML workflows in a controlled and reproducible manner.	

1. Prediction Services in MLOps

Prediction services are the interface through which ML models deliver predictions to end users. These services must be backed by a fully automated MLOps architecture to ensure:

- Scalability
- Reliability
- Real time performance

Modes of Prediction Serving

Mode	Description
Batch Serving	Processes large volumes of data at scheduled intervals
Streaming	Continuously processes incoming data records for real time predictions
Real Time	Processes individual records and returns predictions instantly
On the Edge	Runs models locally on edge devices (e.g., mobile, IoT) to reduce latency

2. Deployment Strategy Overview

Deployment strategies define how new models are introduced into production. The choice depends on:

- Serving mode
- Resource constraints
- Risk tolerance
- Performance requirements

Common Deployment Strategies

Strategy	Description
A/B Testing	Routes traffic between two models (A and B); performance is monitored and traffic is shifted to the better performing model
Shadow Deployment	New model runs in parallel; only the live model serves predictions; outputs are compared for evaluation
Blue/Green	Deploys new model in a separate environment (green); traffic is gradually shifted from the current (blue) environment
Canary Deployment	Not covered in this video, but typically involves rolling out the model to a small subset of users before full deployment

3. A/B Testing Deployment

Workflow

- Two models (A and B) run in parallel.
- A load balancer distributes incoming requests.
- Performance metrics are monitored continuously.
- Traffic is gradually shifted to the better performing model.

Benefits

- Dynamic adjustment based on performance
- Continuous validation of alternative models
- High prediction quality

4. Shadow Deployment

Workflow

- New model runs alongside the production model.
- Both models receive the same input, but only the production model serves predictions.
- Outputs are compared and monitored.

Benefits

- Safe evaluation without affecting users
- Enables performance benchmarking
- Smooth transition to production if validated

5. Blue/Green Deployment

Workflow

- Current model runs in the "blue" environment.
- Updated model is deployed in a replica "green" environment.
- Traffic is gradually shifted from blue to green.
- If no issues arise, full traffic is routed to green.

Benefits

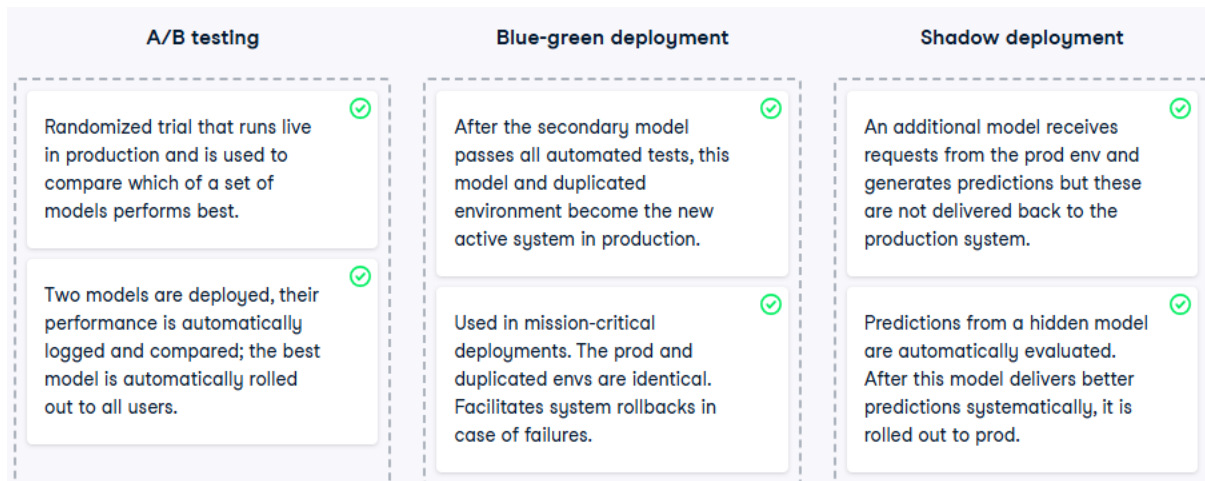
- Minimal downtime
- Easy rollback
- Controlled transition

6. Key Evaluation Dimensions

When selecting a deployment strategy, consider:

Dimension	Description
Downtime	Impact on service availability during deployment

Dimension	Description
Condition Based Triggers	Whether deployment is triggered by performance metrics or schedules
Rollback Time	Speed and ease of reverting to a previous model
Deployment Cost	Resource and operational overhead



CI/CD/CT/CM in Fully Automated MLOps

1. DevOps Foundations

- **DevOps** is a software engineering practice that promotes collaboration between development and operations teams.
- It emphasizes **automation**, **continuous delivery**, and **process optimization** to improve software quality and deployment speed.

Core DevOps Techniques

Technique	Description
CI (Continuous Integration)	Frequent code commits integrated into a shared repository with automated testing
CD (Continuous Deployment)	Automated release of validated code to production environments

2. MLOps and DevOps: Shared Principles

- MLOps extends DevOps principles to machine learning systems.

- Both focus on:
 - Automation across the lifecycle
 - Seamless integration and deployment
 - Reliability and scalability

3. ML Specific Considerations

a. Changing Data and World

- ML systems learn behavior from data, which evolves over time.
- This dynamic nature requires **Continuous Monitoring (CM)** to track:
 - Data quality
 - Model performance
 - System health

b. Model Performance Decay

- ML models degrade as data distributions shift.
- **Continuous Training (CT)** ensures models remain accurate by:
 - Retraining on new data
 - Triggering updates based on performance thresholds or schedules

Replicate your blue environment; this means creating the green environment. ✓

Deploy your new model to the newly created environment. ✓

After deploying the new model, start switching traffic to this new model. ✓

After the system delivers predictions from the green environment for long enough, switch all traffic to the green environment. ✓

After all traffic is switched to the green environment, phase out the blue environment and make the green environment the new blue environment. ✓

4. CI/CD/CT/CM in MLOps

Component	Role in MLOps
CI (Continuous Integration)	Integrates code changes with automated ML specific tests (e.g., data validation, model evaluation)
CD (Continuous Deployment)	Automates deployment of models and pipelines to production
CT (Continuous Training)	Retrains models using fresh data to maintain performance
CM (Continuous Monitoring)	Tracks system metrics, detects drift, and triggers corrective actions

5. Implementation Highlights

- **CI/CD** enables seamless integration of updates and fixes into production pipelines.
- **CT** ensures models adapt to evolving data and remain relevant.
- **CM** detects issues like data drift or model degradation and initiates automated responses.
- These components support:
 - Scalable ML operations
 - Reduced manual intervention
 - High system reliability

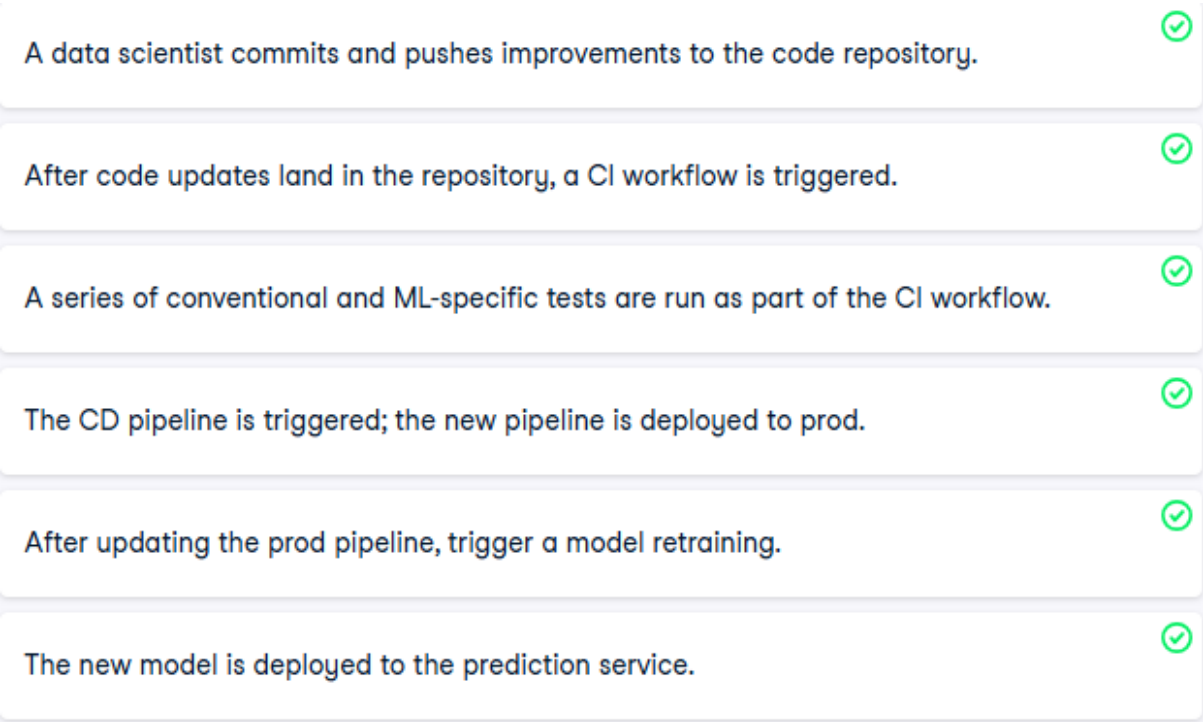
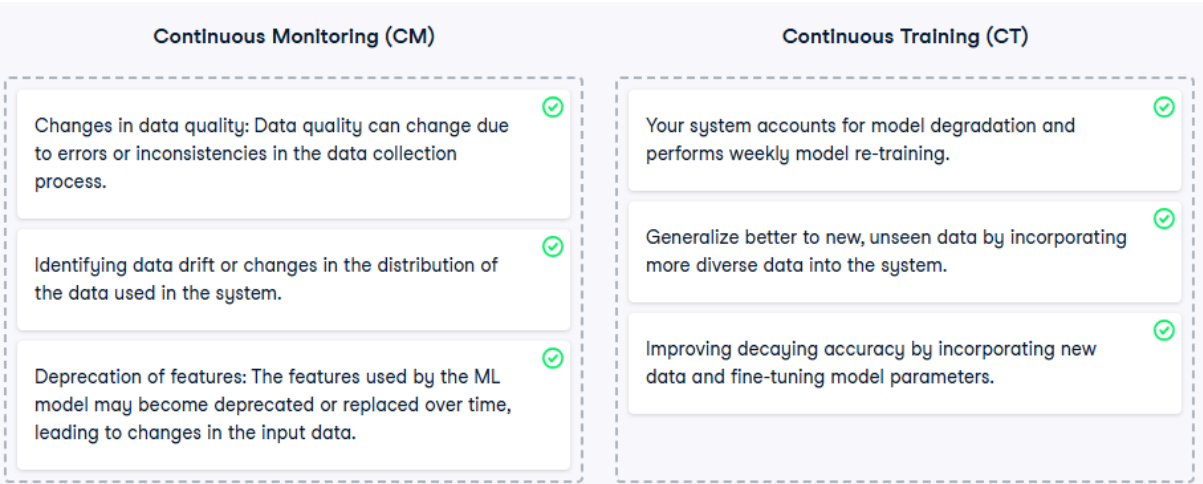
6. Supporting Patterns









- **Automation First:** Design systems to minimize manual steps across the ML lifecycle.
- **Automate Monitor Incident Response Pattern:**
 - Automate workflows
 - Monitor system health
 - Respond to incidents with retraining or rollback

7. Summary: Building Robust MLOps Systems

Together, CI/CD/CT/CM form the backbone of a fully automated MLOps system:

Capability	Benefit
Integration & Deployment	Faster updates, reduced errors
Training & Monitoring	Sustained model accuracy and reliability
Automation & Response	Scalable, resilient ML infrastructure



Design	Development	Operations
<div>How will you detect and eliminate unethical model bias?</div> <div>Is it ethical to use Machine Learning in this use case?</div> <div>Is it legally justified to use personal data in this use case?</div>	<div>Is the training data versioned and produced in a transparent manner?</div> <div>Are the experiments during model development tracked and stored?</div> <div>Are the trained models reproducible?</div>	<div>Do we perform continuous model monitoring?</div> <div>Is the API secure?</div>