# MLOps Pipeline with Kubeflow and MLflow:
## A Comprehensive Guide to Production-Ready ML Infrastructure

Vladimir Ovcharov

MLOps Engineer & ML Systems Architect

vladimir.ovcharov@highfunk.uk

July 10, 2025

### Abstract

This article presents a comprehensive guide to building production-ready machine learning pipelines using Kubeflow for orchestration and MLflow for experiment tracking. We explore the architectural patterns, implementation strategies, and best practices for creating scalable, maintainable, and robust MLOps infrastructure. The guide includes practical code examples, deployment configurations, and monitoring strategies that can be directly applied in enterprise environments.

**Keywords:** MLOps, Kubeflow, MLflow, Kubernetes, Machine Learning Pipeline, Experiment Tracking, Model Deployment

# 1 Introduction

The evolution of machine learning from research prototypes to production-ready systems represents one of the most significant challenges facing modern data science organizations. While the academic community has made tremendous strides in developing sophisticated algorithms and achieving state-of-the-art results on benchmark datasets, the transition to real-world applications introduces a complex web of operational, technical, and organizational challenges that traditional software engineering practices struggle to address effectively.

## 1.1 The Production Challenge

The journey from a Jupyter notebook experiment to a production machine learning system involves far more than simply deploying a trained model. Organizations typically encounter what has been termed the "ML production gap" – the substantial difference between the controlled environment of model development and the chaotic reality of production systems where models must operate reliably, scale dynamically, and maintain performance over time.

Consider the typical lifecycle of a machine learning project: data scientists begin with exploratory data analysis, experiment with various algorithms and feature engineering approaches, tune hyperparameters, and eventually arrive at a model that demonstrates promising performance on held-out test data. However, this process often occurs in isolation, using static datasets, controlled environments, and manual workflows that bear little resemblance to the dynamic, distributed, and automated systems required for production deployment.

The challenges multiply when we consider the operational requirements of production ML systems:

- **Data Pipeline Reliability:** Production systems must handle streaming data, data quality issues, schema evolution, and upstream system failures gracefully

- **Model Performance Monitoring:** Unlike traditional software where bugs are typically deterministic, ML models can degrade silently due to data drift, concept drift, or changing business conditions

- **Scalability Requirements:** Models must serve predictions at scale, often handling thousands of requests per second with strict latency requirements

- **Regulatory Compliance:** Many industries require explainability, auditability, and bias detection capabilities that are difficult to retrofit into existing systems

- **Continuous Learning:** Models must be retrained and updated regularly to maintain performance, requiring sophisticated automation and validation pipelines

## 1.2    The Fragmentation Problem

Traditional approaches to ML deployment have led to significant fragmentation across the ML lifecycle. Different teams often use disparate tools for data preparation, model training, validation, deployment, and monitoring. This fragmentation creates several critical issues:

**Tool Proliferation:** Organizations frequently find themselves managing dozens of different tools and platforms, each optimized for a specific phase of the ML lifecycle. Data engineers might use Apache Airflow for data pipelines, data scientists prefer Jupyter notebooks for experimentation, ML engineers deploy models using custom Docker containers, and operations teams monitor systems using traditional APM tools that lack ML-specific metrics.

**Knowledge Silos:** The fragmented toolchain often leads to knowledge silos where each team becomes expert in their specific tools but lacks understanding of the broader system. This creates bottlenecks, communication barriers, and makes it difficult to optimize the end-to-end workflow.

**Reproducibility Crisis:** Without standardized environments and workflows, reproducing experimental results becomes extremely challenging. Models that perform well in development may fail in production due to subtle differences in data preprocessing, library versions, or infrastructure configurations.

**Deployment Bottlenecks:** The handoff between data science teams and production systems often becomes a significant bottleneck. Models developed in Python on local machines must be translated into production-ready services, often requiring significant engineering effort and introducing opportunities for errors.

## 1.3    The MLOps Revolution

Machine Learning Operations (MLOps) has emerged as a discipline specifically designed to address these challenges by applying DevOps principles and practices to machine learning workflows. MLOps recognizes that ML systems are fundamentally different from traditional software applications and require specialized approaches to testing, deployment, monitoring, and maintenance.

The core principles of MLOps include:

1. **Automation:** Minimizing manual interventions through automated pipelines for data processing, model training, validation, and deployment

2. **Reproducibility:** Ensuring that every experiment, training run, and deployment can be exactly reproduced through proper versioning of code, data, models, and infrastructure

3. **Monitoring:** Implementing comprehensive monitoring that goes beyond traditional system metrics to include model-specific indicators like prediction drift, data quality, and business impact

4. **Collaboration:** Breaking down silos between data science, engineering, and operations teams through shared tools, processes, and vocabulary

5. **Governance:** Establishing controls and processes for model approval, deployment, and compliance with regulatory requirements

## 1.4   Technology Landscape and Solution Selection

The MLOps ecosystem has rapidly evolved to include numerous platforms, frameworks, and tools, each addressing different aspects of the ML production challenge. The selection of appropriate technologies requires careful consideration of organizational needs, existing infrastructure, team capabilities, and long-term strategic goals.

**Kubeflow** has emerged as a leading platform for ML workflows on Kubernetes, providing a comprehensive suite of tools that span the entire ML lifecycle. Built on the foundation of Kubernetes, Kubeflow inherits the scalability, reliability, and ecosystem benefits of the container orchestration platform while adding ML-specific capabilities:

- **Kubeflow Pipelines:** A platform for building and deploying portable, scalable ML workflows

- **Katib:** Automated hyperparameter tuning and neural architecture search

- **KServe:** Model serving platform with advanced features like canary deployments and multi-framework support

- **Notebooks:** Managed Jupyter notebook environments with resource allocation and sharing capabilities

- **Training Operators:** Distributed training support for TensorFlow, PyTorch, and other frameworks

**MLflow** complements Kubeflow by providing robust experiment tracking and model management capabilities. Originally developed by Databricks, MLflow has become a de facto standard for ML lifecycle management, offering:

- **Experiment Tracking:** Comprehensive logging of parameters, metrics, and artifacts for every experiment

- **Model Registry:** Centralized model store with versioning, staging, and annotation capabilities

- **Model Packaging:** Standardized format for packaging models with their dependencies

- **Model Serving:** Simple deployment options for various serving platforms

## 1.5   Integration Strategy and Benefits

The combination of Kubeflow and MLflow creates a powerful MLOps platform that addresses the full spectrum of production ML challenges. This integration strategy leverages the strengths of both platforms while mitigating their individual limitations:

**Kubeflow provides the infrastructure and orchestration layer**, handling the complex task of managing distributed workloads, resource allocation, and pipeline execution across Kubernetes clusters. Its native integration with Kubernetes means that ML workloads can benefit from the same scalability, reliability, and operational practices used for other cloud-native applications.

**MLflow serves as the metadata and lifecycle management layer**, providing the tracking, versioning, and governance capabilities essential for maintaining reproducibility and compliance in production environments. Its framework-agnostic approach ensures that teams can continue using their preferred ML libraries while benefiting from standardized lifecycle management.

The synergy between these platforms enables several key capabilities:

- **End-to-End Traceability:** Every model deployed in production can be traced back to its training data, code version, hyperparameters, and experimental results

- **Automated Retraining:** Pipelines can automatically detect model performance degradation and trigger retraining workflows with minimal human intervention

- **A/B Testing and Gradual Rollouts:** New model versions can be safely deployed using canary deployments and traffic splitting capabilities

- **Resource Optimization:** Kubernetes-native resource management ensures efficient utilization of computational resources across training and serving workloads

- **Multi-Environment Consistency:** The same pipeline definitions can be deployed across development, staging, and production environments with environment-specific configurations

## 1.6   Article Scope and Objectives

This comprehensive guide provides practical, implementable solutions for organizations seeking to establish robust MLOps practices using Kubeflow and MLflow. Rather than focusing on theoretical concepts, we emphasize hands-on implementation with real-world code examples, configuration files, and architectural patterns that have been proven in production environments.

The article is structured to take readers through a complete implementation journey, from initial environment setup through advanced monitoring and optimization techniques. Each section builds upon previous concepts while providing sufficient detail for independent implementation. Code examples are production-ready and include error handling, logging, and best practices gleaned from real-world deployments.

Our target audience includes ML engineers, DevOps practitioners, data scientists, and technical leaders who are responsible for moving ML systems from experimentation to production. We assume familiarity with basic ML concepts, Kubernetes fundamentals, and Python programming, but provide sufficient context for readers to understand and adapt the solutions to their specific environments.

By the end of this guide, readers will have a complete understanding of how to:

- Design and implement scalable ML pipelines using Kubeflow

- Establish comprehensive experiment tracking and model management with MLflow

- Deploy models safely and efficiently using modern serving platforms

- Monitor ML systems for performance, drift, and operational issues

- Implement automated retraining and continuous deployment workflows

- Apply security, governance, and compliance best practices

# 2 Implementation Guide

This section provides detailed, step-by-step instructions for implementing a production-ready MLOps platform using Kubeflow and MLflow. The implementation follows a progressive approach, starting with foundational infrastructure and gradually building up to advanced features. Each step includes comprehensive code examples, configuration files, and troubleshooting guidance based on real-world deployment experiences.

## 2.1 Prerequisites and Environment Preparation

Before beginning the implementation, ensure that your environment meets the necessary requirements and that all prerequisite tools are properly configured.

### 2.1.1 Infrastructure Requirements

**Kubernetes Cluster Specifications:**

- **Minimum cluster size:** 3 nodes with 4 CPU cores and 16GB RAM each

- **Recommended cluster size:** 5+ nodes with 8 CPU cores and 32GB RAM each

- **Kubernetes version:** 1.24 or later (tested up to 1.28)

- **Storage:** Dynamic volume provisioning with SSD-backed storage classes

- **Network:** CNI-compatible networking (Calico, Flannel, or cloud provider CNI)

- **Load Balancer:** Cloud provider load balancer or MetalLB for on-premises

  **Additional Infrastructure Components:**

- Object storage (AWS S3, Google Cloud Storage, Azure Blob, or MinIO)

- Container registry (Docker Hub, ECR, GCR, or Harbor)

- DNS management for custom domains and SSL certificates

- Monitoring infrastructure (Prometheus operator recommended)

### 2.1.2 Required Tools and Dependencies

Install and configure the following tools on your management workstation:

```bash
#!/bin/bash

# Install kubectl
curl -LO "https://dl.k8s.io/release/$(curl -L -s https://dl.k8s.io/release/
    stable.txt)/bin/linux/amd64/kubectl"
sudo install -o root -g root -m 0755 kubectl /usr/local/bin/kubectl

# Install Helm
curl https://raw.githubusercontent.com/helm/helm/main/scripts/get-helm-3 | bash

```

```
10  # Install kustomize
11  curl -s "https://raw.githubusercontent.com/kubernetes-sigs/kustomize/master/
        hack/install_kustomize.sh" | bash
12  sudo mv kustomize /usr/local/bin/
13
14  # Install yq for YAML processing
15  sudo wget -qO /usr/local/bin/yq https://github.com/mikefarah/yq/releases/latest
        /download/yq_linux_amd64
16  sudo chmod +x /usr/local/bin/yq
17
18  # Install kubens and kubectx for easier cluster management
19  sudo git clone https://github.com/ahmetb/kubectx /opt/kubectx
20  sudo ln -s /opt/kubectx/kubectx /usr/local/bin/kubectx
21  sudo ln -s /opt/kubectx/kubens /usr/local/bin/kubens
22
23  # Verify installations
24  echo "Verifying tool installations..."
25  kubectl version --client
26  helm version
27  kustomize version
28  yq --version
```

Listing 1: Tool Installation Script

## 2.2  Kubeflow Installation and Configuration

Kubeflow installation requires careful attention to component selection and configuration to ensure optimal performance and security.

### 2.2.1  Namespace and Security Setup

Begin by creating the necessary namespaces and security configurations:

```
1  # kubeflow-namespaces.yaml
2  apiVersion: v1
3  kind: Namespace
4  metadata:
5    name: kubeflow
6    labels:
7      control-plane: kubeflow
8      istio-injection: enabled
9  ---
10  apiVersion: v1
11  kind: Namespace
12  metadata:
13    name: kubeflow-user-example-com
14    labels:
15      control-plane: kubeflow
16      istio-injection: enabled
17      user: example@example.com
18  ---
19  apiVersion: v1
20  kind: ServiceAccount
21  metadata:
22    name: kubeflow-admin
23    namespace: kubeflow
24  ---
25  apiVersion: rbac.authorization.k8s.io/v1
26  kind: ClusterRoleBinding
27  metadata:
28    name: kubeflow-admin
29  roleRef:
```

```
30    apiGroup: rbac.authorization.k8s.io
31    kind: ClusterRole
32    name: cluster-admin
33 subjects:
34 - kind: ServiceAccount
35    name: kubeflow-admin
36    namespace: kubeflow
```

Listing 2: Namespace and RBAC Configuration

### 2.2.2  Kubeflow Manifests Installation

Install Kubeflow using the official manifests with customizations for production environments:

```bash
1  #!/bin/bash
2
3  set -e
4
5  # Configuration variables
6  export KF_VERSION="v1.7.0"
7  export KF_NAME="mlops-platform"
8  export BASE_DIR="${HOME}/kubeflow"
9  export KF_DIR="${BASE_DIR}/${KF_NAME}"
10 export CONFIG_URI="https://raw.githubusercontent.com/kubeflow/manifests/${
       KF_VERSION}/kfdef/kfctl_k8s_istio.v1.7.0.yaml"
11
12 # Create directory structure
13 mkdir -p ${KF_DIR}
14 cd ${KF_DIR}
15
16 # Download Kubeflow manifests
17 echo "Downloading Kubeflow manifests..."
18 wget -O kubeflow-manifests.tar.gz \
19   "https://github.com/kubeflow/manifests/archive/${KF_VERSION}.tar.gz"
20 tar -xzf kubeflow-manifests.tar.gz
21 cd manifests-${KF_VERSION#v}
22
23 # Apply custom configurations
24 echo "Applying custom configurations..."
25 cat > custom-config.yaml << EOF
26 apiVersion: kfdef.apps.kubeflow.org/v1
27 kind: KfDef
28 metadata:
29   name: ${KF_NAME}
30   namespace: kubeflow
31 spec:
32   applications:
33   - kustomizeConfig:
34       repoRef:
35         name: manifests
36         path: stacks/kubernetes/application/istio-1-16
37     name: istio-1-16
38   - kustomizeConfig:
39       repoRef:
40         name: manifests
41         path: stacks/kubernetes/application/cluster-local-gateway-1-16
42     name: cluster-local-gateway-1-16
43   - kustomizeConfig:
44       repoRef:
45         name: manifests
46         path: apps/pipeline/upstream/env/cert-manager/platform-agnostic-multi-
     user
47     name: kubeflow-pipelines
```

```
48     - kustomizeConfig:
49         repoRef:
50           name: manifests
51           path: apps/jupyter/jupyter-web-app/upstream/overlays/istio
52       name: jupyter-web-app
53     - kustomizeConfig:
54         repoRef:
55           name: manifests
56           path: apps/katib/upstream/installs/katib-with-kubeflow
57       name: katib
58     - kustomizeConfig:
59         repoRef:
60           name: manifests
61           path: apps/training-operator/upstream/overlays/kubeflow
62       name: training-operator
63     repos:
64     - name: manifests
65       uri: https://github.com/kubeflow/manifests/archive/${KF_VERSION}.tar.gz
66 EOF
67
68 # Install cert-manager first (required for webhooks)
69 echo "Installing cert-manager..."
70 kubectl apply -f https://github.com/cert-manager/cert-manager/releases/download
      /v1.12.0/cert-manager.yaml
71 kubectl wait --for=condition=ready pod -l app=cert-manager -n cert-manager --
      timeout=300s
72
73 # Install Kubeflow components
74 echo "Installing Kubeflow components..."
75 while ! kustomize build example | kubectl apply -f -; do
76   echo "Retrying to apply resources"
77   sleep 10
78 done
79
80 # Wait for deployments to be ready
81 echo "Waiting for Kubeflow components to be ready..."
82 kubectl wait --for=condition=ready pod -l app=istiod -n istio-system --timeout
      =300s
83 kubectl wait --for=condition=ready pod -l app=istio-proxy -n kubeflow --timeout
      =300s
84
85 echo "Kubeflow installation completed successfully!"
```

Listing 3: Kubeflow Installation Process

### 2.2.3   Kubeflow Post-Installation Configuration

Configure Kubeflow for production use with proper security and resource management:

```
1 # kubeflow-production-config.yaml
2 apiVersion: v1
3 kind: ConfigMap
4 metadata:
5   name: kubeflow-config
6   namespace: kubeflow
7 data:
8   # Pipeline configuration
9   pipeline-config.yaml: |
10     executorImage: gcr.io/ml-pipeline/api-server:2.0.0
11     cacheEnabled: true
12     cacheDatabase:
13       host: mysql.kubeflow.svc.cluster.local
14       port: "3306"
```

```
15        database: mlpipeline
16     objectStore:
17        endpoint: minio -service.kubeflow.svc.cluster.local:9000
18        bucket: mlpipeline
19        accessKey: minio
20        secretKey: minio123
21
22   # Resource limits
23   resource -limits.yaml: |
24     defaultRequests:
25        cpu: "100m"
26        memory: "128Mi"
27     defaultLimits:
28        cpu: "1000m"
29        memory: "1Gi"
30     maxRequests:
31        cpu: "8000m"
32        memory: "16Gi"
33 ---
34 apiVersion: v1
35 kind: Secret
36 metadata:
37   name: kubeflow -secrets
38   namespace: kubeflow
39 type: Opaque
40 data:
41   # Base64 encoded values
42   mysql -password: bWxwaXBlbGluZQ==   # mlpipeline
43   minio -access -key: bWluaW8=         # minio
44   minio -secret -key: bWluaW8xMjM=     # minio123
```

Listing 4: Production Kubeflow Configuration

## 2.3   MLflow Setup and Integration

MLflow serves as the experiment tracking and model registry backbone of our MLOps platform.
This section covers both standalone MLflow deployment and integration with Kubeflow.

### 2.3.1   MLflow Server Deployment

Deploy MLflow server with persistent storage and database backend:

```
1 # mlflow -deployment.yaml
2 apiVersion: apps/v1
3 kind: Deployment
4 metadata:
5   name: mlflow -server
6   namespace: kubeflow
7   labels:
8     app: mlflow -server
9 spec:
10   replicas: 2
11   selector:
12     matchLabels:
13        app: mlflow -server
14   template:
15     metadata:
16        labels:
17          app: mlflow -server
18     spec:
19        containers:
20        - name: mlflow -server
```

```
21          image: python:3.9-slim
22          ports:
23          - containerPort: 5000
24          env:
25          - name: MLFLOW_S3_ENDPOINT_URL
26            value: "http://minio-service.kubeflow.svc.cluster.local:9000"
27          - name: AWS_ACCESS_KEY_ID
28            valueFrom:
29              secretKeyRef:
30                name: mlflow-secrets
31                key: aws-access-key-id
32          - name: AWS_SECRET_ACCESS_KEY
33            valueFrom:
34              secretKeyRef:
35                name: mlflow-secrets
36                key: aws-secret-access-key
37          - name: MLFLOW_DATABASE_URI
38            valueFrom:
39              secretKeyRef:
40                name: mlflow-secrets
41                key: database-uri
42          command:
43          - /bin/bash
44          - -c
45          - |
46            pip install mlflow[extras]==2.7.1 psycopg2-binary boto3
47            mlflow server \
48              --host 0.0.0.0 \
49              --port 5000 \
50              --backend-store-uri ${MLFLOW_DATABASE_URI} \
51              --default-artifact-root s3://mlflow-artifacts/ \
52              --serve-artifacts
53          resources:
54            requests:
55              cpu: 100m
56              memory: 256Mi
57            limits:
58              cpu: 500m
59              memory: 1Gi
60          livenessProbe:
61            httpGet:
62              path: /health
63              port: 5000
64            initialDelaySeconds: 30
65            periodSeconds: 10
66          readinessProbe:
67            httpGet:
68              path: /health
69              port: 5000
70            initialDelaySeconds: 5
71            periodSeconds: 5
72 ---
73 apiVersion: v1
74 kind: Service
75 metadata:
76   name: mlflow-server
77   namespace: kubeflow
78 spec:
79   selector:
80     app: mlflow-server
81   ports:
82   - port: 5000
83     targetPort: 5000
```

```
84        protocol: TCP
85    type: ClusterIP
86  ---
87  apiVersion: v1
88  kind: Secret
89  metadata:
90    name: mlflow-secrets
91    namespace: kubeflow
92  type: Opaque
93  data:
94    aws-access-key-id: bWluaW8=  # minio
95    aws-secret-access-key: bWluaW8xMjM=  # minio123
96    database-uri:
        cG9zdGdyZXNxbDovL21sZmxvdzptbGZsb3dAcG9zdGdyZXNxbC5rdWJlZmxvdy5zdmMuY2x1c3Rlci5sb2NhbDo1Nl
        ==
```

<div align="center">Listing 5: MLflow Server Deployment</div>

### 2.3.2   Database and Storage Configuration

Set up PostgreSQL database and MinIO object storage for MLflow:

```
1   # postgresql-deployment.yaml
2   apiVersion: apps/v1
3   kind: Deployment
4   metadata:
5     name: postgresql
6     namespace: kubeflow
7   spec:
8     replicas: 1
9     selector:
10      matchLabels:
11        app: postgresql
12    template:
13      metadata:
14        labels:
15          app: postgresql
16      spec:
17        containers:
18        - name: postgresql
19          image: postgres:13
20          env:
21          - name: POSTGRES_DB
22            value: mlflow
23          - name: POSTGRES_USER
24            value: mlflow
25          - name: POSTGRES_PASSWORD
26            value: mlflow
27          - name: PGDATA
28            value: /var/lib/postgresql/data/pgdata
29          ports:
30          - containerPort: 5432
31          volumeMounts:
32          - name: postgresql-storage
33            mountPath: /var/lib/postgresql/data
34          resources:
35            requests:
36              cpu: 100m
37              memory: 256Mi
38            limits:
39              cpu: 500m
40              memory: 1Gi
41        volumes:
```

```yaml
42        - name: postgresql-storage
43          persistentVolumeClaim:
44            claimName: postgresql-pvc
45 ---
46 apiVersion: v1
47 kind: Service
48 metadata:
49   name: postgresql
50   namespace: kubeflow
51 spec:
52   selector:
53     app: postgresql
54   ports:
55   - port: 5432
56     targetPort: 5432
57 ---
58 apiVersion: v1
59 kind: PersistentVolumeClaim
60 metadata:
61   name: postgresql-pvc
62   namespace: kubeflow
63 spec:
64   accessModes:
65   - ReadWriteOnce
66   resources:
67     requests:
68       storage: 10Gi
69   storageClassName: fast-ssd
70 ---
71 # minio-deployment.yaml
72 apiVersion: apps/v1
73 kind: Deployment
74 metadata:
75   name: minio
76   namespace: kubeflow
77 spec:
78   replicas: 1
79   selector:
80     matchLabels:
81       app: minio
82   template:
83     metadata:
84       labels:
85         app: minio
86     spec:
87       containers:
88       - name: minio
89         image: minio/minio:RELEASE.2023-09-04T19-57-37Z
90        args:
91        - server
92        - /data
93        - --console-address=:9001
94        env:
95        - name: MINIO_ROOT_USER
96          value: minio
97        - name: MINIO_ROOT_PASSWORD
98          value: minio123
99        ports:
100       - containerPort: 9000
101       - containerPort: 9001
102       volumeMounts:
103       - name: minio-storage
104         mountPath: /data
```

```
105        resources:
106          requests:
107            cpu: 100m
108            memory: 256Mi
109          limits:
110            cpu: 500m
111            memory: 1Gi
112      volumes:
113      - name: minio-storage
114        persistentVolumeClaim:
115          claimName: minio-pvc
116 ---
117 apiVersion: v1
118 kind: Service
119 metadata:
120   name: minio-service
121   namespace: kubeflow
122 spec:
123   selector:
124     app: minio
125   ports:
126   - name: api
127     port: 9000
128     targetPort: 9000
129   - name: console
130     port: 9001
131     targetPort: 9001
132 ---
133 apiVersion: v1
134 kind: PersistentVolumeClaim
135 metadata:
136   name: minio-pvc
137   namespace: kubeflow
138 spec:
139   accessModes:
140   - ReadWriteOnce
141   resources:
142     requests:
143       storage: 50Gi
144   storageClassName: fast-ssd
```

Listing 6: MLflow Storage Infrastructure

## 2.4   Network Configuration and Security

Proper network configuration ensures secure communication between components while maintaining performance and accessibility.

### 2.4.1   Istio Service Mesh Configuration

Configure Istio for secure service-to-service communication:

```
1 # istio-security-policies.yaml
2 apiVersion: security.istio.io/v1beta1
3 kind: PeerAuthentication
4 metadata:
5   name: default
6   namespace: kubeflow
7 spec:
8   mtls:
9     mode: STRICT
10 ---
```

```
11  apiVersion: security.istio.io/v1beta1
12  kind: AuthorizationPolicy
13  metadata:
14    name: mlflow-access
15    namespace: kubeflow
16  spec:
17    selector:
18      matchLabels:
19        app: mlflow-server
20    rules:
21    - from:
22      - source:
23          principals: ["cluster.local/ns/kubeflow/sa/default"]
24      - source:
25          namespaces: ["kubeflow-user-example-com"]
26    - to:
27      - operation:
28          methods: ["GET", "POST", "PUT", "DELETE"]
29  ---
30  apiVersion: networking.istio.io/v1beta1
31  kind: VirtualService
32  metadata:
33    name: mlflow-vs
34    namespace: kubeflow
35  spec:
36    hosts:
37    - mlflow.example.com
38    gateways:
39    - kubeflow-gateway
40    http:
41    - match:
42      - uri:
43          prefix: /
44      route:
45      - destination:
46          host: mlflow-server.kubeflow.svc.cluster.local
47          port:
48            number: 5000
49      timeout: 300s
50  ---
51  apiVersion: networking.istio.io/v1beta1
52  kind: Gateway
53  metadata:
54    name: kubeflow-gateway
55    namespace: kubeflow
56  spec:
57    selector:
58      istio: ingressgateway
59    servers:
60    - port:
61        number: 80
62        name: http
63        protocol: HTTP
64      hosts:
65      - "*.example.com"
66      tls:
67        httpsRedirect: true
68    - port:
69        number: 443
70        name: https
71        protocol: HTTPS
72      hosts:
73      - "*.example.com"
```

```
74      tls:
75        mode: SIMPLE
76        credentialName: kubeflow-tls-secret
```

<div align="center">Listing 7: Istio Security Configuration</div>

## 2.5   Monitoring and Observability Setup

Comprehensive monitoring is essential for maintaining platform health and performance.

### 2.5.1   Prometheus and Grafana Configuration

Deploy monitoring stack with custom dashboards for MLOps metrics:

```
1  # monitoring-namespace.yaml
2  apiVersion: v1
3  kind: Namespace
4  metadata:
5    name: monitoring
6    labels:
7      name: monitoring
8  ---
9  # prometheus-config.yaml
10 apiVersion: v1
11 kind: ConfigMap
12 metadata:
13   name: prometheus-config
14   namespace: monitoring
15 data:
16   prometheus.yml: |
17     global:
18       scrape_interval: 15s
19       evaluation_interval: 15s
20
21     rule_files:
22       - "mlops_rules.yml"
23
24     scrape_configs:
25     - job_name: 'kubeflow-pipelines'
26       kubernetes_sd_configs:
27       - role: pod
28         namespaces:
29           names:
30           - kubeflow
31       relabel_configs:
32       - source_labels: [__meta_kubernetes_pod_label_app]
33         action: keep
34         regex: ml-pipeline.*
35
36     - job_name: 'mlflow-server'
37       kubernetes_sd_configs:
38       - role: pod
39         namespaces:
40           names:
41           - kubeflow
42       relabel_configs:
43       - source_labels: [__meta_kubernetes_pod_label_app]
44         action: keep
45         regex: mlflow-server
46
47     - job_name: 'model-servers'
48       kubernetes_sd_configs:
49       - role: pod
```

```
50        namespaces:
51          names:
52            - kubeflow-user-example-com
53      relabel_configs:
54      - source_labels: [
     __meta_kubernetes_pod_annotation_serving_kserve_io_inferenceservice]
55        action: keep
56        regex: .+
57
58  mlops_rules.yml: |
59    groups:
60    - name: mlops.rules
61      rules:
62      - alert: ModelServerDown
63        expr: up{job="model-servers"} == 0
64        for: 1m
65        labels:
66          severity: critical
67        annotations:
68          summary: "Model server {{ $labels.instance }} is down"
69          description: "Model server has been down for more than 1 minute"
70
71      - alert: HighModelLatency
72        expr: histogram_quantile(0.95, rate(
     model_request_duration_seconds_bucket[5m])) > 1
73        for: 2m
74        labels:
75          severity: warning
76        annotations:
77          summary: "High model inference latency detected"
78          description: "95th percentile latency is {{ $value }}s"
79
80      - alert: MLflowServerDown
81        expr: up{job="mlflow-server"} == 0
82        for: 2m
83        labels:
84          severity: critical
85        annotations:
86          summary: "MLflow server is down"
87          description: "MLflow server has been unreachable for more than 2
     minutes"
```

Listing 8: Monitoring Stack Deployment

## 2.6  Initial Platform Validation

After completing the installation, validate that all components are functioning correctly:

```bash
#!/bin/bash

set -e

echo "=== MLOps Platform Validation ==="

# Check namespace status
echo "Checking namespace status..."
kubectl get namespaces kubeflow kubeflow-user-example-com monitoring

# Check Kubeflow components
echo "Checking Kubeflow components..."
kubectl get pods -n kubeflow | grep -E "(Running|Completed)" || exit 1

# Check MLflow server
```

```
16  echo "Checking MLflow server..."
17  kubectl get pods -n kubeflow -l app=mlflow-server
18  kubectl wait --for=condition=ready pod -l app=mlflow-server -n kubeflow --
        timeout=300s
19
20  # Test MLflow API
21  echo "Testing MLflow API..."
22  MLF_HOST=$(kubectl get svc mlflow-server -n kubeflow -o jsonpath='{.spec.
        clusterIP}')
23  kubectl run test-pod --rm -i --tty --image=curlimages/curl -- \
24    curl -f http://${MLF_HOST}:5000/health || exit 1
25
26  # Check storage components
27  echo "Checking storage components..."
28  kubectl get pods -n kubeflow -l app=postgresql
29  kubectl get pods -n kubeflow -l app=minio
30
31  # Check Istio configuration
32  echo "Checking Istio configuration..."
33  kubectl get virtualservices,gateways -n kubeflow
34
35  # Test pipeline functionality
36  echo "Testing pipeline functionality..."
37  kubectl get workflows -n kubeflow-user-example-com || echo "No workflows found
        (expected for new installation)"
38
39  echo "=== Platform validation completed successfully! ==="
40  echo ""
41  echo "Access URLs (configure DNS or port-forward):"
42  echo "- Kubeflow Central Dashboard: https://kubeflow.example.com"
43  echo "- MLflow UI: https://mlflow.example.com"
44  echo "- MinIO Console: https://minio.example.com"
45  echo ""
46  echo "Next steps:"
47  echo "1. Configure DNS entries for your domain"
48  echo "2. Set up SSL certificates"
49  echo "3. Create user profiles and RBAC policies"
50  echo "4. Run your first ML pipeline"
```

Listing 9: Platform Validation Script

This implementation guide provides a solid foundation for deploying a production-ready MLOps platform. The next sections will cover advanced pipeline development, model serving configurations, and operational best practices.

# 3    Model Deployment and Serving

Model deployment and serving represent critical phases in the MLOps lifecycle where trained models transition from experimental artifacts to production systems that deliver real business value. This section provides comprehensive guidance for implementing robust, scalable, and secure model serving infrastructure using KServe, along with advanced deployment patterns that ensure safe model releases and optimal performance.

## 3.1    KServe Architecture and Integration

KServe (formerly KFServing) serves as the cornerstone of our model serving strategy, providing a Kubernetes-native platform for deploying and managing machine learning models at scale. Built on the foundation of Knative Serving, KServe inherits powerful capabilities for autoscaling, traffic management, and serverless deployment patterns while adding ML-specific functionality.

17

### 3.1.1   KServe Core Components

Understanding KServe's architecture is essential for effective model deployment and troubleshooting:

**InferenceService:** The primary custom resource that defines how models should be served, including predictor configurations, transformers, and explainers.

**Model Server:** Framework-specific serving runtimes that handle model loading, inference request processing, and response formatting. KServe supports multiple built-in servers including:

- **Scikit-learn Server:** Optimized for traditional ML models with pickle format support

- **TensorFlow Serving:** High-performance serving for TensorFlow models with batching and GPU support

- **PyTorch Server (TorchServe):** Native PyTorch model serving with custom handler support

- **XGBoost Server:** Specialized serving for gradient boosting models

- **Custom Predictors:** User-defined serving logic for complex inference pipelines

**Data Plane:** Handles actual inference requests and responses, implementing the KServe v1 and v2 inference protocols for standardized communication.

**Control Plane:** Manages the lifecycle of inference services, including deployment, scaling, and traffic routing decisions.

### 3.1.2   KServe Installation and Configuration

Deploy KServe with production-ready configurations:

```bash
#!/bin/bash

set -e

echo "Installing KServe with dependencies..."

# Install Knative Serving (required for KServe)
echo "Installing Knative Serving..."
kubectl apply -f https://github.com/knative/serving/releases/download/knative-v1.11.0/serving-crds.yaml
kubectl apply -f https://github.com/knative/serving/releases/download/knative-v1.11.0/serving-core.yaml

# Install Knative Istio controller
kubectl apply -f https://github.com/knative/net-istio/releases/download/knative-v1.11.0/net-istio.yaml

# Configure Knative Serving
kubectl patch configmap/config-network \
  --namespace knative-serving \
  --type merge \
  --patch '{"data":{"ingress-class":"istio.ingress.networking.knative.dev"}}'

# Set domain configuration
kubectl patch configmap/config-domain \
  --namespace knative-serving \
  --type merge \
  --patch '{"data":{"example.com":""}}'

# Install KServe CRDs and controllers
```

```
28  echo "Installing KServe..."
29  kubectl apply -f https://github.com/kserve/kserve/releases/download/v0.11.0/
        kserve.yaml
30
31  # Install KServe built-in ClusterServingRuntimes
32  kubectl apply -f https://github.com/kserve/kserve/releases/download/v0.11.0/
        kserve-runtimes.yaml
33
34  # Wait for KServe controller to be ready
35  echo "Waiting for KServe controller..."
36  kubectl wait --for=condition=ready pod -l control-plane=kserve-controller-
        manager -n kserve --timeout=300s
37
38  echo "KServe installation completed successfully!"
```

Listing 10: KServe Installation Script

### 3.1.3  Production KServe Configuration

Configure KServe for production environments with proper resource management and security:

```
1  # kserve-config.yaml
2  apiVersion: v1
3  kind: ConfigMap
4  metadata:
5    name: inferenceservice-config
6    namespace: kserve
7  data:
8    predictors: |
9      {
10       "tensorflow": {
11         "image": "tensorflow/serving:2.13.0",
12         "defaultImageVersion": "2.13.0",
13         "defaultGpuImageVersion": "2.13.0-gpu",
14         "supportedFrameworks": ["tensorflow"],
15         "multiModelServer": false
16       },
17       "pytorch": {
18         "image": "pytorch/torchserve:0.8.2-cpu",
19         "defaultImageVersion": "0.8.2-cpu",
20         "defaultGpuImageVersion": "0.8.2-gpu",
21         "supportedFrameworks": ["pytorch"],
22         "multiModelServer": false
23       },
24       "sklearn": {
25         "image": "kserve/sklearnserver:v0.11.0",
26         "defaultImageVersion": "v0.11.0",
27         "supportedFrameworks": ["sklearn"],
28         "multiModelServer": true
29       },
30       "xgboost": {
31         "image": "kserve/xgbserver:v0.11.0",
32         "defaultImageVersion": "v0.11.0",
33         "supportedFrameworks": ["xgboost"],
34         "multiModelServer": true
35       }
36     }
37
38    transformer: |
39      {
40       "feast": {
41         "image": "kserve/feast-transformer:v0.11.0",
42         "defaultImageVersion": "v0.11.0"
```

```
 43          }
 44        }
 45
 46    explainer: |
 47        {
 48          "alibi": {
 49            "image": "kserve/alibi-explainer:v0.11.0",
 50            "defaultImageVersion": "v0.11.0"
 51          }
 52        }
 53
 54    storageInitializer: |
 55        {
 56          "image": "kserve/storage-initializer:v0.11.0",
 57          "memoryRequest": "100Mi",
 58          "memoryLimit": "1Gi",
 59          "cpuRequest": "100m",
 60          "cpuLimit": "1000m"
 61        }
 62
 63    credentials: |
 64        {
 65          "gcs": {
 66            "gcsCredentialFileName": "gcloud-application-credentials.json"
 67          },
 68          "s3": {
 69            "s3AccessKeyIDName": "AWS_ACCESS_KEY_ID",
 70            "s3SecretAccessKeyName": "AWS_SECRET_ACCESS_KEY",
 71            "s3Endpoint": "",
 72            "s3UseHttps": true,
 73            "s3Region": "us-west-1",
 74            "s3VerifySSL": true,
 75            "s3UseVirtualBucket": false,
 76            "s3UseAnonymousCredential": false,
 77            "s3CABundle": ""
 78          }
 79        }
 80
 81    ingress: |
 82        {
 83          "ingressGateway": "kubeflow/kubeflow-gateway",
 84          "ingressService": "istio-ingressgateway.istio-system.svc.cluster.local",
 85          "localGateway": "knative-serving/knative-local-gateway",
 86          "localGatewayService": "knative-local-gateway.istio-system.svc.cluster.
     local",
 87          "ingressDomain": "example.com",
 88          "ingressClassName": "istio",
 89          "domainTemplate": "{{.Name}}-{{.Namespace}}.{{.IngressDomain}}",
 90          "urlScheme": "https",
 91          "disableIstioVirtualHost": false
 92        }
 93
 94    deploy: |
 95        {
 96          "defaultDeploymentMode": "Serverless",
 97          "progressDeadlineSeconds": 600,
 98          "defaultCpuRequest": "100m",
 99          "defaultMemoryRequest": "128Mi",
100          "defaultCpuLimit": "1000m",
101          "defaultMemoryLimit": "2Gi"
102        }
103 ---
104 apiVersion: v1
```

```
105  kind: ConfigMap
106  metadata:
107    name: kserve-logger-config
108    namespace: kserve
109  data:
110    logger.properties: |
111      # Root logger option
112      log4j.rootLogger=INFO, stdout
113
114      # Direct log messages to stdout
115      log4j.appender.stdout=org.apache.log4j.ConsoleAppender
116      log4j.appender.stdout.Target=System.out
117      log4j.appender.stdout.layout=org.apache.log4j.PatternLayout
118      log4j.appender.stdout.layout.ConversionPattern=%d{yyyy-MM-dd HH:mm:ss} %-5p
         %c{1}:%L - %m%n
119
120      # Suppress unnecessary logs
121      log4j.logger.org.apache.hadoop=WARN
122      log4j.logger.org.apache.spark=WARN
123      log4j.logger.org.eclipse.jetty=WARN
124      log4j.logger.org.apache.kafka=WARN
```

Listing 11: KServe Production Configuration

## 3.2   Model Packaging and Registry Integration

Effective model deployment requires standardized packaging and integration with model registries to ensure consistency, traceability, and reproducibility.

### 3.2.1   MLflow Model Integration

Create a comprehensive MLflow integration that automatically packages models for KServe deployment:

```python
1   import mlflow
2   import mlflow.sklearn
3   import mlflow.pytorch
4   import mlflow.tensorflow
5   import yaml
6   import json
7   import os
8   from typing import Dict, Any, Optional
9   from kubernetes import client, config
10  from datetime import datetime
11  import logging
12
13  class MLflowKServeIntegration:
14      """
15      Handles integration between MLflow model registry and KServe deployments
16      """
17
18      def __init__(self,
19                   mlflow_uri: str,
20                   namespace: str = "kubeflow-user-example-com",
21                   kserve_domain: str = "example.com"):
22          self.mlflow_uri = mlflow_uri
23          self.namespace = namespace
24          self.kserve_domain = kserve_domain
25
26          # Initialize MLflow client
27          mlflow.set_tracking_uri(mlflow_uri)
28          self.mlflow_client = mlflow.MlflowClient()
```

21

```python
29
30          # Initialize Kubernetes client
31          try:
32              config.load_incluster_config()
33          except:
34              config.load_kube_config()
35
36          self.k8s_client = client.ApiClient()
37          self.custom_client = client.CustomObjectsApi()
38
39          # Configure logging
40          logging.basicConfig(level=logging.INFO)
41          self.logger = logging.getLogger(__name__)
42
43      def get_model_info(self, model_name: str, version: Optional[str] = None) ->
        Dict[str, Any]:
44          """
45          Retrieve comprehensive model information from MLflow registry
46          """
47          try:
48              if version is None:
49                  # Get latest version
50                  latest_versions = self.mlflow_client.get_latest_versions(
51                      model_name, stages=["Production", "Staging"]
52                  )
53                  if not latest_versions:
54                      raise ValueError(f"No versions found for model {model_name}
    ")
55                  model_version = latest_versions[0]
56              else:
57                  model_version = self.mlflow_client.get_model_version(model_name
    , version)
58
59              # Get run information
60              run = self.mlflow_client.get_run(model_version.run_id)
61
62              # Extract model artifacts
63              model_uri = f"models:/{model_name}/{model_version.version}"
64
65              return {
66                  "name": model_name,
67                  "version": model_version.version,
68                  "stage": model_version.current_stage,
69                  "model_uri": model_uri,
70                  "run_id": model_version.run_id,
71                  "framework": self._detect_framework(run),
72                  "artifacts": run.data.tags.get("mlflow.log-model.history", "[]"
    ),
73                  "metrics": run.data.metrics,
74                  "params": run.data.params,
75                  "tags": run.data.tags,
76                  "creation_timestamp": model_version.creation_timestamp,
77                  "last_updated_timestamp": model_version.last_updated_timestamp
78              }
79
80          except Exception as e:
81              self.logger.error(f"Error retrieving model info: {str(e)}")
82              raise
83
84      def _detect_framework(self, run) -> str:
85          """
86          Detect ML framework from run artifacts and tags
87          """
```

```python
88          artifacts = json.loads(run.data.tags.get("mlflow.log-model.history", "
    []"))

89
90          for artifact in artifacts:
91              if "sklearn" in artifact.get("flavors", {}):
92                  return "sklearn"
93              elif "pytorch" in artifact.get("flavors", {}):
94                  return "pytorch"
95              elif "tensorflow" in artifact.get("flavors", {}):
96                  return "tensorflow"
97              elif "xgboost" in artifact.get("flavors", {}):
98                  return "xgboost"

99
100         # Fallback to run tags
101         if "mlflow.source.type" in run.data.tags:
102             source_type = run.data.tags["mlflow.source.type"].lower()
103             if "sklearn" in source_type:
104                 return "sklearn"
105             elif "pytorch" in source_type:
106                 return "pytorch"
107             elif "tensorflow" in source_type:
108                 return "tensorflow"

109
110         return "sklearn"  # Default fallback

111
112     def create_inference_service(self,
113                                  model_name: str,
114                                  model_version: Optional[str] = None,
115                                  service_name: Optional[str] = None,
116                                  resources: Optional[Dict[str, Any]] = None,
117                                  autoscaling: Optional[Dict[str, Any]] = None,
118                                  canary_percent: int = 0) -> Dict[str, Any]:
119         """
120         Create KServe InferenceService from MLflow model
121         """

122
123         # Get model information
124         model_info = self.get_model_info(model_name, model_version)

125
126         if service_name is None:
127             service_name = f"{model_name.lower().replace('_', '-')}-v{
    model_info['version']}"

128
129         # Default resource configuration
130         default_resources = {
131             "requests": {"cpu": "100m", "memory": "256Mi"},
132             "limits": {"cpu": "1000m", "memory": "2Gi"}
133         }
134         if resources:
135             default_resources.update(resources)

136
137         # Default autoscaling configuration
138         default_autoscaling = {
139             "minReplicas": 1,
140             "maxReplicas": 10,
141             "targetUtilizationPercentage": 70,
142             "scaleToZeroGracePeriod": "30s",
143             "scaleDownDelay": "0s",
144             "stableWindow": "60s"
145         }
146         if autoscaling:
147             default_autoscaling.update(autoscaling)

148
```

23

```python
149         # Create InferenceService specification
150         inference_service = {
151             "apiVersion": "serving.kserve.io/v1beta1",
152             "kind": "InferenceService",
153             "metadata": {
154                 "name": service_name,
155                 "namespace": self.namespace,
156                 "labels": {
157                     "model-name": model_name,
158                     "model-version": str(model_info['version']),
159                     "model-stage": model_info['stage'],
160                     "framework": model_info['framework'],
161                     "managed-by": "mlflow-kserve-integration"
162                 },
163                 "annotations": {
164                     "mlflow.model.uri": model_info['model_uri'],
165                     "mlflow.run.id": model_info['run_id'],
166                     "deployment.timestamp": datetime.utcnow().isoformat(),
167                     "serving.kserve.io/deploymentMode": "Serverless"
168                 }
169             },
170             "spec": {
171                 "predictor": {
172                     model_info['framework']: {
173                         "storageUri": model_info['model_uri'],
174                         "resources": default_resources,
175                         "env": [
176                             {
177                                 "name": "STORAGE_URI",
178                                 "value": model_info['model_uri']
179                             },
180                             {
181                                 "name": "MODEL_NAME",
182                                 "value": model_name
183                             }
184                         ]
185                     }
186                 }
187             }
188         }
189
190         # Add canary configuration if specified
191         if canary_percent > 0:
192             inference_service["spec"]["predictor"]["canaryTrafficPercent"] =
    canary_percent
193
194         # Add autoscaling annotations
195         inference_service["metadata"]["annotations"].update({
196             "autoscaling.knative.dev/minScale": str(default_autoscaling["
    minReplicas"]),
197             "autoscaling.knative.dev/maxScale": str(default_autoscaling["
    maxReplicas"]),
198             "autoscaling.knative.dev/target": str(default_autoscaling["
    targetUtilizationPercentage"]),
199             "autoscaling.knative.dev/scaleToZeroGracePeriod":
    default_autoscaling["scaleToZeroGracePeriod"],
200             "autoscaling.knative.dev/scaleDownDelay": default_autoscaling["
    scaleDownDelay"],
201             "autoscaling.knative.dev/window": default_autoscaling["stableWindow
    "]
202         })
203
204         # Deploy to Kubernetes
```

```python
205        try:
206            response = self.custom_client.create_namespaced_custom_object(
207                group="serving.kserve.io",
208                version="v1beta1",
209                namespace=self.namespace,
210                plural="inferenceservices",
211                body=inference_service
212            )
213
214            self.logger.info(f"Successfully created InferenceService: {
    service_name}")
215
216            # Update MLflow model stage if deploying to production
217            if model_info['stage'] != "Production":
218                self.mlflow_client.transition_model_version_stage(
219                    name=model_name,
220                    version=model_info['version'],
221                    stage="Production",
222                    archive_existing_versions=False
223                )
224
225            return {
226                "service_name": service_name,
227                "namespace": self.namespace,
228                "model_info": model_info,
229                "inference_service": response,
230                "endpoint_url": f"https://{service_name}-{self.namespace}.{self
    .kserve_domain}"
231            }
232
233        except Exception as e:
234            self.logger.error(f"Failed to create InferenceService: {str(e)}")
235            raise
236
237    def update_traffic_split(self,
238                             service_name: str,
239                             traffic_config: Dict[str, int]) -> Dict[str, Any]:
240        """
241        Update traffic splitting between model versions
242        """
243        try:
244            # Get current InferenceService
245            current_service = self.custom_client.get_namespaced_custom_object(
246                group="serving.kserve.io",
247                version="v1beta1",
248                namespace=self.namespace,
249                plural="inferenceservices",
250                name=service_name
251            )
252
253            # Update traffic configuration
254            if "canary" in traffic_config:
255                patch_body = {
256                    "spec": {
257                        "predictor": {
258                            "canaryTrafficPercent": traffic_config["canary"]
259                        }
260                    }
261                }
262
263                response = self.custom_client.patch_namespaced_custom_object(
264                    group="serving.kserve.io",
265                    version="v1beta1",
```

```python
266                    namespace=self.namespace,
267                    plural="inferenceservices",
268                    name=service_name,
269                    body=patch_body
270                )
271
272                self.logger.info(f"Updated traffic split for {service_name}: {
     traffic_config}")
273                return response
274
275        except Exception as e:
276            self.logger.error(f"Failed to update traffic split: {str(e)}")
277            raise
278
279    def rollback_deployment(self, service_name: str, target_version: str) ->
     Dict[str, Any]:
280        """
281        Rollback deployment to a previous model version
282        """
283        try:
284            # Get service metadata to extract model info
285            current_service = self.custom_client.get_namespaced_custom_object(
286                group="serving.kserve.io",
287                version="v1beta1",
288                namespace=self.namespace,
289                plural="inferenceservices",
290                name=service_name
291            )
292
293            model_name = current_service["metadata"]["labels"]["model-name"]
294
295            # Get target version info
296            target_model_info = self.get_model_info(model_name, target_version)
297
298            # Update the service to use target version
299            patch_body = {
300                "spec": {
301                    "predictor": {
302                        target_model_info['framework']: {
303                            "storageUri": target_model_info['model_uri']
304                        }
305                    }
306                },
307                "metadata": {
308                    "labels": {
309                        "model-version": str(target_model_info['version'])
310                    },
311                    "annotations": {
312                        "mlflow.model.uri": target_model_info['model_uri'],
313                        "mlflow.run.id": target_model_info['run_id'],
314                        "rollback.timestamp": datetime.utcnow().isoformat(),
315                        "rollback.target.version": target_version
316                    }
317                }
318            }
319
320            response = self.custom_client.patch_namespaced_custom_object(
321                group="serving.kserve.io",
322                version="v1beta1",
323                namespace=self.namespace,
324                plural="inferenceservices",
325                name=service_name,
326                body=patch_body
```

```
327            )
328
329            self.logger.info(f"Successfully rolled back {service_name} to
       version {target_version}")
330            return response
331
332        except Exception as e:
333            self.logger.error(f"Failed to rollback deployment: {str(e)}")
334            raise
335
336 # Usage example
337 if __name__ == "__main__":
338     # Initialize integration
339     integration = MLflowKServeIntegration(
340         mlflow_uri="http://mlflow-server.kubeflow.svc.cluster.local:5000",
341         namespace="kubeflow-user-example-com"
342     )
343
344     # Deploy a model
345     deployment_result = integration.create_inference_service(
346         model_name="fraud-detection-model",
347         resources={
348             "requests": {"cpu": "200m", "memory": "512Mi"},
349             "limits": {"cpu": "2000m", "memory": "4Gi"}
350         },
351         autoscaling={
352             "minReplicas": 2,
353             "maxReplicas": 20,
354             "targetUtilizationPercentage": 80
355         }
356     )
357
358     print(f"Model deployed successfully: {deployment_result['endpoint_url']}")
```

Listing 12: MLflow KServe Integration

## 3.3   Advanced Deployment Patterns

Production model deployment requires sophisticated strategies to minimize risk while ensuring continuous service availability. This section covers advanced deployment patterns that enable safe model releases.

### 3.3.1   Canary Deployments

Implement gradual rollout of new model versions with automated monitoring and rollback capabilities:

```
1 import asyncio
2 import logging
3 from typing import Dict, List, Optional
4 from dataclasses import dataclass
5 from datetime import datetime, timedelta
6 import numpy as np
7 from prometheus_client.parser import text_string_to_metric_families
8 import aiohttp
9 import yaml
10
11 @dataclass
12 class CanaryConfig:
13     """Configuration for canary deployment"""
14     initial_traffic_percent: int = 5
15     increment_percent: int = 10
```

```python
16      max_traffic_percent: int = 50
17      evaluation_duration_minutes: int = 10
18      success_rate_threshold: float = 0.99
19      latency_threshold_ms: float = 1000
20      error_rate_threshold: float = 0.01
21      auto_promote: bool = True
22      auto_rollback: bool = True
23
24  @dataclass
25  class DeploymentMetrics:
26      """Metrics for deployment evaluation"""
27      success_rate: float
28      avg_latency_ms: float
29      error_rate: float
30      request_count: int
31      timestamp: datetime
32
33  class CanaryDeploymentController:
34      """
35      Automated canary deployment controller with monitoring and decision making
36      """
37
38      def __init__(self,
39                   kserve_integration: MLflowKServeIntegration,
40                   prometheus_url: str,
41                   config: CanaryConfig):
42          self.kserve_integration = kserve_integration
43          self.prometheus_url = prometheus_url
44          self.config = config
45          self.logger = logging.getLogger(__name__)
46
47      async def deploy_canary(self,
48                              model_name: str,
49                              new_version: str,
50                              service_name: str) -> Dict[str, Any]:
51          """
52          Execute automated canary deployment process
53          """
54
55          deployment_log = {
56              "model_name": model_name,
57              "new_version": new_version,
58              "service_name": service_name,
59              "start_time": datetime.utcnow(),
60              "stages": [],
61              "final_status": "in_progress"
62          }
63
64          try:
65              # Stage 1: Deploy new version with minimal traffic
66              self.logger.info(f"Starting canary deployment for {model_name} v{
67  new_version}")
67
68              # Create new InferenceService for canary version
69              canary_service_name = f"{service_name}-canary"
70              canary_deployment = self.kserve_integration.
71  create_inference_service(
71                  model_name=model_name,
72                  model_version=new_version,
73                  service_name=canary_service_name,
74                  canary_percent=self.config.initial_traffic_percent
75              )
76
```

```python
77              deployment_log["stages"].append({
78                  "stage": "initial_deployment",
79                  "traffic_percent": self.config.initial_traffic_percent,
80                  "timestamp": datetime.utcnow(),
81                  "status": "success"
82              })
83
84              # Wait for deployment to stabilize
85              await asyncio.sleep(60)
86
87              # Stage 2: Gradual traffic increase with monitoring
88              current_traffic = self.config.initial_traffic_percent
89
90              while current_traffic < self.config.max_traffic_percent:
91                  # Evaluate current performance
92                  baseline_metrics = await self.get_deployment_metrics(
     service_name)
93                  canary_metrics = await self.get_deployment_metrics(
     canary_service_name)
94
95                  # Make deployment decision
96                  decision = self._evaluate_canary_performance(baseline_metrics,
     canary_metrics)
97
98                  stage_log = {
99                      "stage": "traffic_increase",
100                     "traffic_percent": current_traffic,
101                     "timestamp": datetime.utcnow(),
102                     "baseline_metrics": baseline_metrics.__dict__ if
     baseline_metrics else None,
103                     "canary_metrics": canary_metrics.__dict__ if canary_metrics
      else None,
104                     "decision": decision
105                 }
106
107                 if decision["action"] == "continue":
108                     # Increase traffic
109                     current_traffic = min(
110                         current_traffic + self.config.increment_percent,
111                         self.config.max_traffic_percent
112                     )
113
114                     await self._update_traffic_split(canary_service_name,
     current_traffic)
115                     stage_log["new_traffic_percent"] = current_traffic
116                     stage_log["status"] = "success"
117
118                     self.logger.info(f"Increased canary traffic to {
     current_traffic}%")
119
120                 elif decision["action"] == "rollback":
121                     # Automatic rollback
122                     if self.config.auto_rollback:
123                         await self._rollback_canary(service_name,
     canary_service_name)
124                         stage_log["status"] = "rollback"
125                         deployment_log["final_status"] = "failed"
126                         break
127                     else:
128                         # Manual intervention required
129                         stage_log["status"] = "requires_manual_intervention"
130                         deployment_log["final_status"] = "requires_intervention
     "
```

29

```python
131                         break
132
133                 elif decision["action"] == "hold":
134                     # Hold current traffic level for extended evaluation
135                     self.logger.info(f"Holding canary traffic at {
      current_traffic}% for extended evaluation")
136                     await asyncio.sleep(self.config.evaluation_duration_minutes
       * 60 * 2)  # Extended wait
137                     stage_log["status"] = "hold"
138
139                 deployment_log["stages"].append(stage_log)
140
141                 # Wait for evaluation period
142                 await asyncio.sleep(self.config.evaluation_duration_minutes *
      60)
143
144             # Stage 3: Final evaluation and promotion decision
145             if deployment_log["final_status"] == "in_progress":
146                 final_baseline_metrics = await self.get_deployment_metrics(
      service_name)
147                 final_canary_metrics = await self.get_deployment_metrics(
      canary_service_name)
148
149                 final_decision = self._evaluate_canary_performance(
150                     final_baseline_metrics,
151                     final_canary_metrics,
152                     final_evaluation=True
153                 )
154
155                 if final_decision["action"] == "promote" and self.config.
      auto_promote:
156                     # Promote canary to production
157                     await self._promote_canary(service_name,
      canary_service_name, new_version)
158                     deployment_log["final_status"] = "promoted"
159                     self.logger.info(f"Successfully promoted {model_name} v{
      new_version} to production")
160                 else:
161                     deployment_log["final_status"] = "requires_manual_promotion
      "
162                     self.logger.info(f"Canary deployment ready for manual
      promotion")
163
164                 deployment_log["stages"].append({
165                     "stage": "final_evaluation",
166                     "timestamp": datetime.utcnow(),
167                     "final_decision": final_decision,
168                     "status": deployment_log["final_status"]
169                 })
170
171             deployment_log["end_time"] = datetime.utcnow()
172             deployment_log["duration_minutes"] = (
173                 deployment_log["end_time"] - deployment_log["start_time"]
174             ).total_seconds() / 60
175
176             return deployment_log
177
178         except Exception as e:
179             self.logger.error(f"Canary deployment failed: {str(e)}")
180             deployment_log["final_status"] = "error"
181             deployment_log["error"] = str(e)
182
183             # Attempt cleanup
```

```python
184                try:
185                    await self._cleanup_failed_canary(canary_service_name)
186                except:
187                    pass
188
189                raise
190
191     async def get_deployment_metrics(self, service_name: str) -> Optional[
        DeploymentMetrics]:
192         """
193         Retrieve performance metrics for a deployment from Prometheus
194         """
195         try:
196             queries = {
197                 "success_rate": f'rate(kserve_request_total{{service_name="{
        service_name}",code!~"5.."}}[5m]) / rate(kserve_request_total{{service_name
        ="{service_name}"}}[5m])',
198                 "avg_latency": f'histogram_quantile(0.50, rate(
        kserve_request_duration_seconds_bucket{{service_name="{service_name}"}}[5m])
        ) * 1000',
199                 "error_rate": f'rate(kserve_request_total{{service_name="{
        service_name}",code=~"5.."}}[5m]) / rate(kserve_request_total{{service_name
        ="{service_name}"}}[5m])',
200                 "request_count": f'rate(kserve_request_total{{service_name="{
        service_name}"}}[5m]) * 300'  # 5 minute rate * 300 seconds
201             }
202
203             metrics = {}
204             async with aiohttp.ClientSession() as session:
205                 for metric_name, query in queries.items():
206                     url = f"{self.prometheus_url}/api/v1/query"
207                     params = {"query": query}
208
209                     async with session.get(url, params=params) as response:
210                         if response.status == 200:
211                             data = await response.json()
212                             result = data.get("data", {}).get("result", [])
213
214                             if result:
215                                 value = float(result[0]["value"][1])
216                                 metrics[metric_name] = value
217                             else:
218                                 metrics[metric_name] = 0.0
219                         else:
220                             self.logger.warning(f"Failed to query {metric_name
        }: {response.status}")
221                             metrics[metric_name] = 0.0
222
223             return DeploymentMetrics(
224                 success_rate=metrics.get("success_rate", 0.0),
225                 avg_latency_ms=metrics.get("avg_latency", 0.0),
226                 error_rate=metrics.get("error_rate", 0.0),
227                 request_count=int(metrics.get("request_count", 0)),
228                 timestamp=datetime.utcnow()
229             )
230
231         except Exception as e:
232             self.logger.error(f"Failed to retrieve metrics for {service_name}:
        {str(e)}")
233             return None
234
235     def _evaluate_canary_performance(self,
```

31

```python
236                                        baseline_metrics: Optional[DeploymentMetrics
    ],
237                                        canary_metrics: Optional[DeploymentMetrics],
238                                        final_evaluation: bool = False) -> Dict[str,
     any]:
239         """
240         Evaluate canary performance against baseline and thresholds
241         """
242         if not canary_metrics:
243             return {
244                 "action": "rollback",
245                 "reason": "No canary metrics available",
246                 "confidence": 1.0
247             }
248
249         # Check absolute thresholds
250         threshold_checks = {
251             "success_rate": canary_metrics.success_rate >= self.config.
    success_rate_threshold,
252             "latency": canary_metrics.avg_latency_ms <= self.config.
    latency_threshold_ms,
253             "error_rate": canary_metrics.error_rate <= self.config.
    error_rate_threshold,
254             "min_requests": canary_metrics.request_count >= 10  # Minimum
    sample size
255         }
256
257         failed_checks = [check for check, passed in threshold_checks.items() if
     not passed]
258
259         if failed_checks:
260             return {
261                 "action": "rollback",
262                 "reason": f"Failed threshold checks: {failed_checks}",
263                 "canary_metrics": canary_metrics.__dict__,
264                 "failed_checks": failed_checks,
265                 "confidence": 1.0
266             }
267
268         # Compare with baseline if available
269         if baseline_metrics and baseline_metrics.request_count >= 10:
270             relative_checks = {
271                 "success_rate_degradation": (canary_metrics.success_rate /
    baseline_metrics.success_rate) >= 0.99,
272                 "latency_increase": (canary_metrics.avg_latency_ms /
    baseline_metrics.avg_latency_ms) <= 1.2,
273                 "error_rate_increase": canary_metrics.error_rate <= (
    baseline_metrics.error_rate * 2.0 + 0.001)
274             }
275
276             failed_relative_checks = [check for check, passed in
    relative_checks.items() if not passed]
277
278             if failed_relative_checks:
279                 # Calculate confidence based on sample size and magnitude of
    degradation
280                 confidence = min(1.0, canary_metrics.request_count / 100)
281
282                 return {
283                     "action": "rollback" if confidence > 0.7 else "hold",
284                     "reason": f"Performance degradation detected: {
    failed_relative_checks}",
285                     "canary_metrics": canary_metrics.__dict__,
```

```python
286                    "baseline_metrics": baseline_metrics.__dict__,
287                    "failed_checks": failed_relative_checks,
288                    "confidence": confidence
289                }

291        # Determine action based on evaluation type
292        if final_evaluation:
293            return {
294                "action": "promote",
295                "reason": "All performance checks passed",
296                "canary_metrics": canary_metrics.__dict__,
297                "confidence": 1.0
298            }
299        else:
300            return {
301                "action": "continue",
302                "reason": "Performance within acceptable range",
303                "canary_metrics": canary_metrics.__dict__,
304                "confidence": 0.8
305            }

307    async def _update_traffic_split(self, canary_service_name: str,
    traffic_percent: int):
308        """Update traffic split for canary deployment"""
309        await self.kserve_integration.update_traffic_split(
310            canary_service_name,
311            {"canary": traffic_percent}
312        )

314    async def _rollback_canary(self, service_name: str, canary_service_name:
    str):
315        """Rollback canary deployment"""
316        try:
317            # Set canary traffic to 0
318            await self._update_traffic_split(canary_service_name, 0)

320            # Delete canary service after grace period
321            await asyncio.sleep(30)

323            self.kserve_integration.custom_client.
    delete_namespaced_custom_object(
324                group="serving.kserve.io",
325                version="v1beta1",
326                namespace=self.kserve_integration.namespace,
327                plural="inferenceservices",
328                name=canary_service_name
329            )

331            self.logger.info(f"Rolled back canary deployment: {
    canary_service_name}")

333        except Exception as e:
334            self.logger.error(f"Error during rollback: {str(e)}")

336    async def _promote_canary(self, service_name: str, canary_service_name: str
    , new_version: str):
337        """Promote canary to production"""
338        try:
339            # Get canary service configuration
340            canary_service = self.kserve_integration.custom_client.
    get_namespaced_custom_object(
341                group="serving.kserve.io",
342                version="v1beta1",
```

```
343                    namespace=self.kserve_integration.namespace,
344                    plural="inferenceservices",
345                    name=canary_service_name
346                )
347
348                # Update production service with canary configuration
349                patch_body = {
350                    "spec": canary_service["spec"],
351                    "metadata": {
352                        "labels": canary_service["metadata"]["labels"],
353                        "annotations": {
354                            **canary_service["metadata"]["annotations"],
355                            "promotion.timestamp": datetime.utcnow().isoformat(),
356                            "promoted.from": canary_service_name
357                        }
358                    }
359                }
360
361                # Remove canary-specific configurations
362                if "canaryTrafficPercent" in patch_body["spec"]["predictor"]:
363                    del patch_body["spec"]["predictor"]["canaryTrafficPercent"]
364
365                # Update production service
366                self.kserve_integration.custom_client.
    patch_namespaced_custom_object(
367                    group="serving.kserve.io",
368                    version="v1beta1",
369                    namespace=self.kserve_integration.namespace,
370                    plural="inferenceservices",
371                    name=service_name,
372                    body=patch_body
373                )
374
375                # Clean up canary service
376                await asyncio.sleep(60)  # Wait for traffic to shift
377                await self._cleanup_failed_canary(canary_service_name)
378
379                self.logger.info(f"Successfully promoted {canary_service_name} to
    production")
380
381        except Exception as e:
382                self.logger.error(f"Error during promotion: {str(e)}")
383                raise
384
385     async def _cleanup_failed_canary(self, canary_service_name: str):
386        """Clean up failed canary deployment"""
387        try:
388                self.kserve_integration.custom_client.
    delete_namespaced_custom_object(
389                    group="serving.kserve.io",
390                    version="v1beta1",
391                    namespace=self.kserve_integration.namespace,
392                    plural="inferenceservices",
393                    name=canary_service_name
394                )
395                self.logger.info(f"Cleaned up canary service: {canary_service_name}
    ")
```
\subsubsection{Blue-Green Deployments}

Implement zero-downtime deployments with instant rollback capabilities:

\begin{lstlisting}[language=python, caption=Blue-Green Deployment
    Implementation]

```python
401  class BlueGreenDeploymentManager:
402      """
403      Manages blue-green deployments for zero-downtime model updates
404      """
405
406      def __init__(self, kserve_integration: MLflowKServeIntegration):
407          self.kserve_integration = kserve_integration
408          self.logger = logging.getLogger(__name__)
409
410      async def deploy_blue_green(self,
411                                  model_name: str,
412                                  new_version: str,
413                                  service_name: str,
414                                  validation_tests: List[callable] = None) -> Dict[
     str, Any]:
415          """
416          Execute blue-green deployment with automated validation
417          """
418
419          deployment_result = {
420              "model_name": model_name,
421              "new_version": new_version,
422              "service_name": service_name,
423              "start_time": datetime.utcnow(),
424              "status": "in_progress"
425          }
426
427          try:
428              # Step 1: Deploy green environment
429              green_service_name = f"{service_name}-green"
430
431              self.logger.info(f"Deploying green environment: {green_service_name
     }")
432
433              green_deployment = self.kserve_integration.create_inference_service
     (
434                  model_name=model_name,
435                  model_version=new_version,
436                  service_name=green_service_name
437              )
438
439              # Step 2: Wait for green environment to be ready
440              await self._wait_for_service_ready(green_service_name)
441
442              # Step 3: Run validation tests against green environment
443              if validation_tests:
444                  validation_results = await self._run_validation_tests(
445                      green_service_name,
446                      validation_tests
447                  )
448
449                  if not validation_results["all_passed"]:
450                      # Cleanup and abort
451                      await self._cleanup_service(green_service_name)
452                      deployment_result["status"] = "validation_failed"
453                      deployment_result["validation_results"] =
     validation_results
454                      return deployment_result
455
456              # Step 4: Switch traffic from blue to green
457              self.logger.info("Switching traffic from blue to green")
458
459              # Get current blue service configuration
```

```python
460             blue_service = self.kserve_integration.custom_client.
        get_namespaced_custom_object(
461                 group="serving.kserve.io",
462                 version="v1beta1",
463                 namespace=self.kserve_integration.namespace,
464                 plural="inferenceservices",
465                 name=service_name
466             )
467
468             # Backup blue configuration
469             blue_backup_name = f"{service_name}-blue-backup-{int(datetime.
        utcnow().timestamp())}"
470             deployment_result["blue_backup_name"] = blue_backup_name
471
472             # Rename current service to backup
473             await self._rename_service(service_name, blue_backup_name)
474
475             # Rename green service to production
476             await self._rename_service(green_service_name, service_name)
477
478             # Step 5: Monitor post-deployment metrics
479             post_deployment_metrics = await self._monitor_post_deployment(
        service_name)
480             deployment_result["post_deployment_metrics"] =
        post_deployment_metrics
481
482             # Step 6: Cleanup old blue environment after successful deployment
483             await asyncio.sleep(300)  # Wait 5 minutes before cleanup
484             await self._cleanup_service(blue_backup_name)
485
486             deployment_result["status"] = "completed"
487             deployment_result["end_time"] = datetime.utcnow()
488
489             self.logger.info(f"Blue-green deployment completed successfully for
         {model_name}")
490
491             return deployment_result
492
493         except Exception as e:
494             self.logger.error(f"Blue-green deployment failed: {str(e)}")
495             deployment_result["status"] = "failed"
496             deployment_result["error"] = str(e)
497
498             # Attempt rollback if switch was attempted
499             if "blue_backup_name" in deployment_result:
500                 try:
501                     await self.rollback_blue_green(service_name,
        deployment_result["blue_backup_name"])
502                     deployment_result["rollback_attempted"] = True
503                 except:
504                     deployment_result["rollback_failed"] = True
505
506             raise
507
508     async def rollback_blue_green(self, service_name: str, blue_backup_name:
        str):
509         """
510         Rollback blue-green deployment to previous version
511         """
512         try:
513             self.logger.info(f"Rolling back blue-green deployment for {
        service_name}")
514
```

```python
515            # Delete current green service
516            await self._cleanup_service(service_name)
517
518            # Restore blue service
519            await self._rename_service(blue_backup_name, service_name)
520
521            self.logger.info("Blue-green rollback completed successfully")
522
523        except Exception as e:
524            self.logger.error(f"Blue-green rollback failed: {str(e)}")
525            raise
526
527    async def _wait_for_service_ready(self, service_name: str, timeout_seconds:
     int = 300):
528        """Wait for service to be ready"""
529        start_time = datetime.utcnow()
530
531        while (datetime.utcnow() - start_time).total_seconds() <
     timeout_seconds:
532            try:
533                service = self.kserve_integration.custom_client.
     get_namespaced_custom_object(
534                    group="serving.kserve.io",
535                    version="v1beta1",
536                    namespace=self.kserve_integration.namespace,
537                    plural="inferenceservices",
538                    name=service_name
539                )
540
541                # Check if service is ready
542                conditions = service.get("status", {}).get("conditions", [])
543                ready_condition = next((c for c in conditions if c["type"] == "
     Ready"), None)
544
545                if ready_condition and ready_condition["status"] == "True":
546                    self.logger.info(f"Service {service_name} is ready")
547                    return
548
549            except Exception as e:
550                self.logger.debug(f"Waiting for service readiness: {str(e)}")
551
552            await asyncio.sleep(10)
553
554        raise TimeoutError(f"Service {service_name} did not become ready within
     {timeout_seconds} seconds")
555
556    async def _run_validation_tests(self, service_name: str, validation_tests:
     List[callable]) -> Dict[str, Any]:
557        """Run validation tests against deployed service"""
558        results = {
559            "all_passed": True,
560            "test_results": [],
561            "execution_time": datetime.utcnow()
562        }
563
564        for test_func in validation_tests:
565            try:
566                test_result = await test_func(service_name)
567                results["test_results"].append({
568                    "test_name": test_func.__name__,
569                    "status": "passed" if test_result else "failed",
570                    "result": test_result
571                })
```

```
572
573                    if not test_result:
574                        results["all_passed"] = False
575
576                except Exception as e:
577                    results["test_results"].append({
578                        "test_name": test_func.__name__,
579                        "status": "error",
580                        "error": str(e)
581                    })
582                    results["all_passed"] = False
583
584            return results
585
```

\subsection{A/B Testing Framework}

Implement statistical A/B testing for model variants:

\begin{lstlisting}[language=python, caption=A/B Testing Implementation]

```python
import scipy.stats as stats
from typing import Tuple
import pandas as pd

class ABTestingFramework:
    """
    Statistical A/B testing framework for model variants
    """

    def __init__(self,
                 kserve_integration: MLflowKServeIntegration,
                 prometheus_url: str,
                 significance_level: float = 0.05,
                 minimum_sample_size: int = 1000,
                 test_duration_hours: int = 72):

        self.kserve_integration = kserve_integration
        self.prometheus_url = prometheus_url
        self.significance_level = significance_level
        self.minimum_sample_size = minimum_sample_size
        self.test_duration_hours = test_duration_hours
        self.logger = logging.getLogger(__name__)

    async def setup_ab_test(self,
                            control_model: Dict[str, str],
                            variant_model: Dict[str, str],
                            traffic_split: int = 50,
                            success_metric: str = "conversion_rate") -> Dict[str,
    Any]:
        """
        Setup A/B test between two model variants
        """

        test_config = {
            "test_id": f"ab_test_{int(datetime.utcnow().timestamp())}",
            "control_model": control_model,
            "variant_model": variant_model,
            "traffic_split": traffic_split,
            "success_metric": success_metric,
            "start_time": datetime.utcnow(),
            "status": "active"
        }

        try:
```

```
634                 # Deploy control version (if not already deployed)
635                 control_service_name = f"{control_model['name']}-control-{
       test_config['test_id']}"
636
637                 if not await self._service_exists(control_service_name):
638                     await self.kserve_integration.create_inference_service(
639                         model_name=control_model['name'],
640                         model_version=control_model['version'],
641                         service_name=control_service_name
642                     )
643
644                 # Deploy variant version
645                 variant_service_name = f"{variant_model['name']}-variant-{
       test_config['test_id']}"
646                 await self.kserve_integration.create_inference_service(
647                     model_name=variant_model['name'],
648                     model_version=variant_model['version'],
649                     service_name=variant_service_name
650                 )
651
652                 # Configure traffic splitting
653                 await self._configure_ab_traffic(
654                     control_service_name,
655                     variant_service_name,
656                     traffic_split
657                 )
658
659                 test_config["control_service"] = control_service_name
660                 test_config["variant_service"] = variant_service_name
661
662                 self.logger.info(f"A/B test setup completed: {test_config['test_id
       ']}")
663
664                 return test_config
665
666         except Exception as e:
667             self.logger.error(f"Failed to setup A/B test: {str(e)}")
668             test_config["status"] = "failed"
669             test_config["error"] = str(e)
670             raise
671
672     async def analyze_ab_test(self, test_config: Dict[str, Any]) -> Dict[str,
       Any]:
673         """
674         Analyze A/B test results and make statistical conclusions
675         """
676
677         try:
678             # Collect metrics for both variants
679             control_metrics = await self._collect_ab_metrics(
680                 test_config["control_service"],
681                 test_config["start_time"]
682             )
683
684             variant_metrics = await self._collect_ab_metrics(
685                 test_config["variant_service"],
686                 test_config["start_time"]
687             )
688
689             # Perform statistical analysis
690             analysis_result = self._perform_statistical_analysis(
691                 control_metrics,
692                 variant_metrics,
```

39

```
693                    test_config["success_metric"]
694                )
695
696                # Generate recommendations
697                recommendation = self._generate_recommendation(analysis_result,
     test_config)
698
699                result = {
700                    "test_id": test_config["test_id"],
701                    "analysis_timestamp": datetime.utcnow(),
702                    "control_metrics": control_metrics,
703                    "variant_metrics": variant_metrics,
704                    "statistical_analysis": analysis_result,
705                    "recommendation": recommendation,
706                    "test_duration_hours": (datetime.utcnow() - test_config["
     start_time"]).total_seconds() / 3600
707                }
708
709                return result
710
711        except Exception as e:
712            self.logger.error(f"Failed to analyze A/B test: {str(e)}")
713            raise
714
715    def _perform_statistical_analysis(self,
716                                      control_metrics: Dict,
717                                      variant_metrics: Dict,
718                                      success_metric: str) -> Dict[str, Any]:
719        """
720        Perform statistical significance testing
721        """
722
723        # Extract success counts and total counts
724        control_successes = control_metrics.get(f"{success_metric}_count", 0)
725        control_total = control_metrics.get("total_requests", 0)
726        variant_successes = variant_metrics.get(f"{success_metric}_count", 0)
727        variant_total = variant_metrics.get("total_requests", 0)
728
729        if control_total == 0 or variant_total == 0:
730            return {
731                "test_type": "insufficient_data",
732                "statistical_significance": False,
733                "p_value": None,
734                "confidence_interval": None,
735                "effect_size": None
736            }
737
738        # Calculate conversion rates
739        control_rate = control_successes / control_total
740        variant_rate = variant_successes / variant_total
741
742        # Perform two-proportion z-test
743        count = np.array([control_successes, variant_successes])
744        nobs = np.array([control_total, variant_total])
745
746        # Calculate z-statistic and p-value
747        z_stat, p_value = stats.proportions_ztest(count, nobs)
748
749        # Calculate confidence interval for difference
750        pooled_rate = (control_successes + variant_successes) / (control_total
     + variant_total)
751        se_diff = np.sqrt(pooled_rate * (1 - pooled_rate) * (1/control_total +
     1/variant_total))
```

```python
752        rate_diff = variant_rate - control_rate
753
754        margin_of_error = stats.norm.ppf(1 - self.significance_level/2) *
    se_diff
755        ci_lower = rate_diff - margin_of_error
756        ci_upper = rate_diff + margin_of_error
757
758        # Calculate effect size (Cohen's h)
759        effect_size = 2 * (np.arcsin(np.sqrt(variant_rate)) - np.arcsin(np.sqrt
    (control_rate)))
760
761        return {
762            "test_type": "two_proportion_z_test",
763            "control_rate": control_rate,
764            "variant_rate": variant_rate,
765            "rate_difference": rate_diff,
766            "relative_improvement": (rate_diff / control_rate) * 100 if
    control_rate > 0 else 0,
767            "z_statistic": z_stat,
768            "p_value": p_value,
769            "statistical_significance": p_value < self.significance_level,
770            "confidence_interval": (ci_lower, ci_upper),
771            "effect_size": effect_size,
772            "sample_sizes": {"control": control_total, "variant": variant_total
    }
773        }
774
775    def _generate_recommendation(self,
776                                 analysis_result: Dict[str, Any],
777                                 test_config: Dict[str, Any]) -> Dict[str, Any]:
778        """
779        Generate actionable recommendations based on test results
780        """
781
782        if analysis_result["test_type"] == "insufficient_data":
783            return {
784                "action": "continue_test",
785                "reason": "Insufficient data for statistical analysis",
786                "required_sample_size": self.minimum_sample_size
787            }
788
789        # Check minimum sample size requirement
790        min_sample_met = all(
791            size >= self.minimum_sample_size
792            for size in analysis_result["sample_sizes"].values()
793        )
794
795        if not min_sample_met:
796            return {
797                "action": "continue_test",
798                "reason": "Minimum sample size not reached",
799                "current_samples": analysis_result["sample_sizes"],
800                "required_sample_size": self.minimum_sample_size
801            }
802
803        # Check test duration
804        test_duration = (datetime.utcnow() - test_config["start_time"]).
    total_seconds() / 3600
805        if test_duration < self.test_duration_hours:
806            return {
807                "action": "continue_test",
808                "reason": f"Test duration ({test_duration:.1f}h) below minimum
    ({self.test_duration_hours}h)",
```

```
809                    "current_duration_hours": test_duration ,
810                    "required_duration_hours": self.test_duration_hours
811                }
812
813        # Make recommendation based on statistical results
814        if analysis_result["statistical_significance"]:
815            if analysis_result["rate_difference"] > 0:
816                return {
817                    "action": "deploy_variant",
818                    "reason": "Variant shows statistically significant
    improvement",
819                    "improvement": f"{analysis_result['relative_improvement
    ']:.2f}%",
820                    "confidence": f"{(1 - self.significance_level) * 100:.0f}%"
821                }
822            else:
823                return {
824                    "action": "keep_control",
825                    "reason": "Control performs significantly better than
    variant",
826                    "degradation": f"{analysis_result['relative_improvement
    ']:.2f}%",
827                    "confidence": f"{(1 - self.significance_level) * 100:.0f}%"
828                }
829        else:
830            return {
831                "action": "no_significant_difference",
832                "reason": "No statistically significant difference detected",
833                "p_value": analysis_result["p_value"],
834                "recommendation": "Consider other factors like cost, complexity
    , or business requirements"
835            }
836
837 # Example usage for A/B testing
838 async def example_ab_test():
839     """Example A/B test implementation"""
840
841     # Initialize framework
842     ab_framework = ABTestingFramework(
843         kserve_integration=integration ,
844         prometheus_url="http://prometheus.monitoring.svc.cluster.local:9090"
845     )
846
847     # Setup A/B test
848     test_config = await ab_framework.setup_ab_test(
849         control_model={"name": "recommendation -model", "version": "1.2.0"},
850         variant_model={"name": "recommendation -model", "version": "1.3.0"},
851         traffic_split=50,
852         success_metric="click_through_rate"
853     )
854
855     # Monitor test for specified duration
856     while True:
857         analysis = await ab_framework.analyze_ab_test(test_config)
858
859         if analysis["recommendation"]["action"] != "continue_test":
860             print(f"Test completed with recommendation: {analysis['
    recommendation']['action']}")
861             break
862
863         print(f"Test continuing... Current improvement: {analysis['
    statistical_analysis'].get('relative_improvement', 0):.2f}%")
```

```
864          await asyncio.sleep(3600)  # Check every hour
```

Listing 13: Automated Canary Deployment Controller

This comprehensive model deployment and serving section provides production-ready implementations for advanced deployment patterns including canary deployments, blue-green deployments, and statistical A/B testing frameworks, all integrated with KServe and MLflow for complete MLOps workflow management.

# 4    Monitoring and Observability

## 4.1    Model Performance Monitoring

Implementing comprehensive monitoring is crucial for maintaining model performance in production:

```
1  import prometheus_client
2  from prometheus_client import Counter, Histogram, Gauge
3  import numpy as np
4  from scipy import stats
5
6  class ModelMonitor:
7      def __init__(self):
8          # Prometheus metrics
9          self.prediction_counter = Counter(
10             'model_predictions_total',
11             'Total number of predictions made',
12             ['model_name', 'version']
13         )
14
15         self.prediction_latency = Histogram(
16             'model_prediction_duration_seconds',
17             'Model prediction latency',
18             ['model_name', 'version']
19         )
20
21         self.model_accuracy = Gauge(
22             'model_accuracy_score',
23             'Current model accuracy',
24             ['model_name', 'version']
25         )
26
27         self.drift_score = Gauge(
28             'model_drift_score',
29             'Data drift detection score',
30             ['model_name', 'feature']
31         )
32
33     def log_prediction(self, model_name, version, latency):
34         """Log prediction metrics"""
35         self.prediction_counter.labels(
36             model_name=model_name,
37             version=version
38         ).inc()
39
40         self.prediction_latency.labels(
41             model_name=model_name,
42             version=version
43         ).observe(latency)
44
45     def detect_data_drift(self, reference_data, current_data, feature_name):
46         """Detect data drift using Kolmogorov-Smirnov test"""
```

43

```python
47        try:
48            # Perform KS test
49            ks_statistic, p_value = stats.ks_2samp(reference_data, current_data
    )
50
51            # Update drift metric
52            self.drift_score.labels(
53                model_name="production_model",
54                feature=feature_name
55            ).set(ks_statistic)
56
57            # Alert if significant drift detected
58            if p_value < 0.05:
59                self.send_drift_alert(feature_name, ks_statistic, p_value)
60
61            return ks_statistic, p_value
62
63        except Exception as e:
64            print(f"Error detecting drift for {feature_name}: {str(e)}")
65            return None, None
66
67    def send_drift_alert(self, feature_name, ks_statistic, p_value):
68        """Send alert when data drift is detected"""
69        alert_message = f"""
70        Data Drift Alert!
71        Feature: {feature_name}
72        KS Statistic: {ks_statistic:.4f}
73        P-value: {p_value:.4f}
74        Recommendation: Review model performance and consider retraining
75        """
76        # Integration with alerting system (Slack, PagerDuty, etc.)
77        print(alert_message)
```

Listing 14: Model Monitoring Implementation

# 5    Best Practices and Recommendations

## 5.1    Pipeline Design Principles

1. **Modularity:** Design pipeline components as independent, reusable modules that can be easily tested and maintained.

2. **Reproducibility:** Ensure all experiments and deployments are fully reproducible through proper versioning of code, data, and dependencies.

3. **Scalability:** Design pipelines to handle varying workloads and data volumes without manual intervention.

4. **Monitoring:** Implement comprehensive monitoring at every stage of the pipeline to quickly identify and resolve issues.

5. **Security:** Apply security best practices including proper authentication, authorization, and data encryption.

## 5.2    Performance Optimization

- Use GPU acceleration for training intensive models

- Implement efficient data loading and preprocessing

- Optimize model serving with batching and caching strategies

- Monitor resource utilization and scale components as needed

- Implement model quantization and pruning for inference optimization

# 6 Troubleshooting Common Issues

## 6.1 Pipeline Failures

Common issues and their solutions:

| Issue | Symptoms | Solution |
|---|---|---|
| Resource limitations | Pipeline steps timing out or failing | Increase resource requests/limits |
| Data access issues | Permission denied errors | Check RBAC and storage permissions |
| Model convergence | Poor model performance | Adjust hyperparameters and data quality |
| Deployment failures | Service unavailable errors | Verify KServe configuration |

Table 1: Common Pipeline Issues and Solutions

# 7 Conclusion

This comprehensive guide demonstrates how to build production-ready MLOps pipelines using Kubeflow and MLflow. The combination of these technologies provides a robust foundation for managing the complete machine learning lifecycle, from experimentation to production deployment.

Key benefits of this approach include:

- **Standardization:** Consistent workflows across teams and projects

- **Scalability:** Kubernetes-native scaling capabilities

- **Reproducibility:** Complete experiment and deployment tracking

- **Monitoring:** Comprehensive observability and alerting

- **Collaboration:** Shared infrastructure and knowledge base

As MLOps practices continue to evolve, this foundation provides the flexibility to adapt and integrate new tools and methodologies while maintaining operational excellence.

# 8 References

1. Kubeflow Documentation. `https://www.kubeflow.org/docs/`

2. MLflow Documentation. `https://mlflow.org/docs/latest/index.html`

3. KServe Documentation. `https://kserve.github.io/website/`

4. Kubernetes Documentation. `https://kubernetes.io/docs/`

5. Sculley, D., et al. "Hidden Technical Debt in Machine Learning Systems." NIPS 2015.

# 9   About the Author

Vladimir Ovcharov is an MLOps Engineer and ML Systems Architect with 8+ years of experience in building production-ready machine learning infrastructure. He specializes in designing scalable ML pipelines, implementing robust monitoring systems, and optimizing model deployment workflows.

Contact: vladimir.ovcharov@example.com