

***Pipeline K8s + Data Workflow**

K8s Architecture:

Traffic → Ingress (routing, SSL) → Service (ClusterIP/LoadBalancer) → Deployment → ReplicaSet → Pods

• Deployment: Rolling update, Rollback, Auto-scaling (HPA) | • ConfigMap + Secret: Configuration management

Data on K8s:

• StatefulSets: PostgreSQL, MongoDB, Cassandra, Kafka | • PV/PVC: Persistent storage (dynamic provisioning)

• Operators: postgres-operator, mongodb-operator, kafka-operator | • Init Containers: Schema migration

• Backup: Velero for K8s backup | • Service Mesh: Istio for data traffic management

***Pipeline Machine Learning + MLOps Data Workflow**

ML Pipeline Stages:

Data Collection → EDA → Preprocessing → Feature Engineering → Model Building → Evaluation → Deployment → Monitoring

MLOps Data Components:

1. Feature Store

Feast, Tecton, Hopsworks

✓ Consistent features ✓ Low latency ✓ Time-travel

2. Data Validation

TFDV, Great Expectations, Pandera

✓ Schema validation ✓ Drift detection

3. Model Registry

MLflow, DVC, Weights & Biases

✓ Versioning ✓ Lineage ✓ Reproducibility

4. Model Monitoring

Evidently AI, WhyLabs, Arize

✓ Data/Model/Concept drift ✓ Performance tracking

Orchestration & Deployment:

Airflow (Batch) | Kubeflow (K8s-native) | Metaflow (Production) | Shadow → Canary → A/B → Full Rollout

CI/CD for ML: Data versioning (DVC) → Training pipeline → Validation → Registry → Deploy → Monitor → Retrain

Best Practices:

- 1. Shift-Left Security & Data Quality**
 - ✓ Early validation ✓ Pre-commit hooks ✓ Fail fast
- 2. Infrastructure as Code (IaC)**
 - ✓ Terraform/Pulumi ✓ GitOps (ArgoCD) ✓ Version everything
- 3. Observability-Driven Development**
 - ✓ Metrics/logs/traces from day 1 ✓ SLO/SLI definition
- 4. Data Contract & Schema Evolution**
 - ✓ Define contracts ✓ Backward compatibility ✓ Version schemas
- 5. Progressive Delivery**
 - ✓ Feature flags ✓ Canary ✓ Blue/Green ✓ Auto rollback