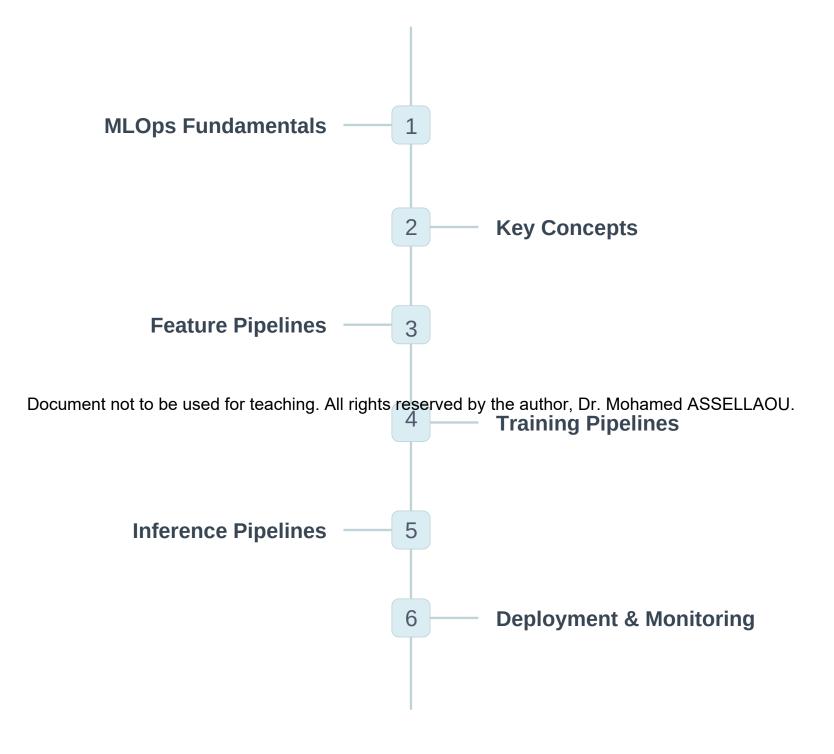


Introduction to MLOps

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Mohamed ASSELLAOU

AGENDA



MLOps overview

The bridge

MLOps connects machine learning development with models work in real-world environments.

Traditional failures

Most ML workflows fail in production due to inconsistent operational systems. It ensures environments, manual Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU. processes, and poor monitoring.

DevOps evolution

MLOps builds on DevOps principles while addressing unique challenges of machine learning systems.

What is MLOps?

Machine learning

Algorithms, models, and

experimentation

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Goal

Reliable and efficient ML lifecycle management

DevOps

Automation, CI/CD, infrastructure

Data engineering

Data pipelines, validation, storage

Why MLOps matters

90%

Failed projects

Percentage of ML projects that never reach production

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

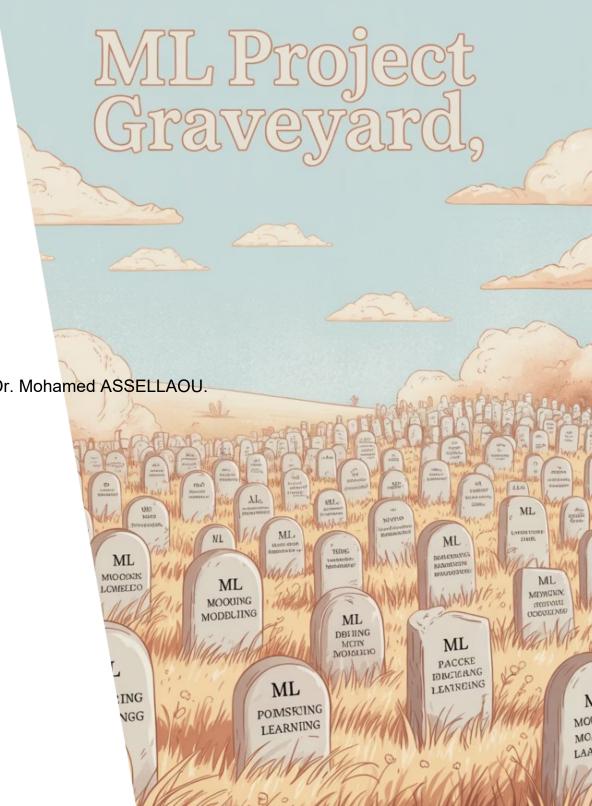
Technical debt

ML systems accumulate technical debt faster than traditional software

40%

Performance drop

Average model accuracy decrease due to data drift in 3 months



The ML lifecycle challenge



System complexity

ML systems have more moving parts than traditional software. They combine code, data, and models.

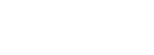


Changing data creates new failure points. Models

depend on specific data distributions.









Document a per the per the period of the per

Models require continuous Cross-functional teams need shared workflows. evaluation. Performance degrades over time without Data scientists and retraining. engineers must collaborate effectively.

MLOps maturity levels

Level 3: Full Automation

Automated end-to-end MLOps lifecycle

Level 2: CI/CD Automation

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU. Automated testing and deployment

Level 1: ML Pipeline Automation

Automated training workflows

Level 0: Manual Process

Manual experimentation and deployment

MLOps core principles

Version control

Track all components: code, data, models, and configurations. Nothing exists without versioning.

Continuous integration

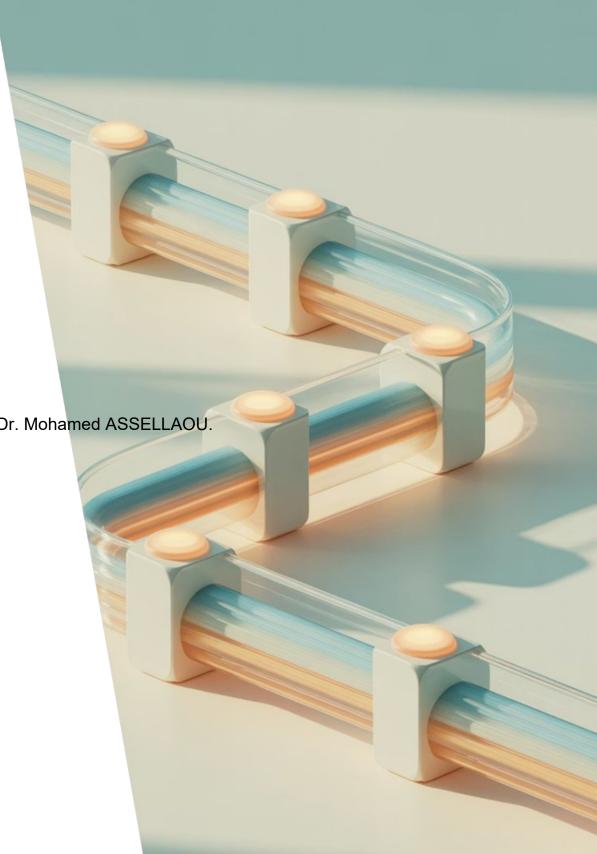
Automatically test ML components when changes occur. Validate code,
Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.
data, and models.

Continuous delivery

Automate model deployment to target environments. Ensure consistency and reliability.

Continuous training

Automatically retrain models on new data. Monitor for drift and performance degradation.





EVOLUTION Of Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Ad-hoc notebooks

Manual experimentation and deployment.

No standardization or reproducibility.

High technical debt.

Monolithic pipelines

Single end-to-end pipelines. Difficult to maintain and test. Tight coupling between components.

Modern FTI approach

Separate Feature/Training/Inference pipelines. Clear interfaces. Independent development and operation.

Monolith vs FTI pipelines

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Monolithic approach

- Hard to debug
- Difficult to test
- Tight coupling
- Team bottlenecks

FTI approach

- Easy to debug
- Testable components
- Loose coupling
- Team autonomy

FTI Pipeline Architecture Feature Training Training Inference

FTI Pipeline Architecture

Feature pipeline

Transforms raw data into ML-ready features and labels. Outputs to feature store.

cument not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Creates models using features and labels. Outputs registered models.

Inference pipeline

Generates predictions using features and models. Outputs to end users.

Feature pipeline

Data ingestion

Collect data from sources

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Feature storage

Save to feature store

Data cleaning

Handle missing values and outliers

Feature transformation

Create ML-ready features

Feature validation

Ensure quality and consistency

Feature engineering

Domain-driven creation

- Domain expertise
- Focus on business relevance

Validation & testing

- Check distributions
- Verify transformations

Document not to be used for teaching. All rights reserved by the author, the honor manner as the bound of the author, the honor manner as the bound of the author, the honor manner as the bound of the author, the honor manner as the bound of the author, the honor manner as the bound of the author of the bound of the

Feature selection

- Remove redundant features
- Measure feature importance

Consistency

- Same logic in training and inference
- Version all transformations
- Document dependencies

Feature stores

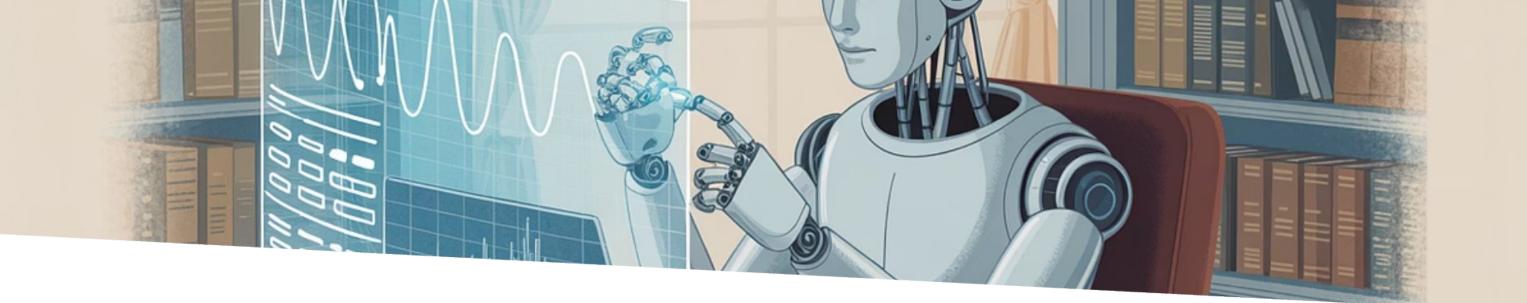
What is a feature store?

A feature store is a repository for managing and serving will rights reserved by the author or Mohamed ASSELLAOU. features. It connects feature engineering with model training and serving.

Think of it as a database specifically designed for ML features.

Key benefits

- Consistency between training and serving
- Feature versioning and lineage tracking
- Online and offline feature serving
- Reduced redundant computations



Feature pipeline debugging of the author, Dr. Mohamed ASSELLAOU.

Identify data quality Issues

Use profiling tools to find outliers, missing values, and distribution shifts in raw data.

Verify transformations

Check outputs after each transformation step. Visualize before and after distributions.



Proactive maintenance example (feature pipeline)

Define transformations

Design feature pipeline

Document not to be used for teaching. All rights reserved by the author of the hand of the series aggregations and

Create a feature pipeline for predicting equipment failures. Think about what data would help predict if a machine will break down.

- Sensor data (temperature, vibration)
- Maintenance logs and history
- **Environmental conditions**

normalization for sensor data.

Plan testing strategy

Outline how you would validate feature quality. Define expected distributions and quality thresholds for metrics.

Feature store implementation

Choose a feature store technology. Plan how features derived from sensor data and logs will be registered and accessed.

Training pipeline

Model registration

Register validated models with metadata. Prepare for deployment.

Model validation

Desument of the perdant teacher that the property of the prope

Model training

Train models with hyperparameter optimization. Log all experiments and metrics.

Feature retrieval

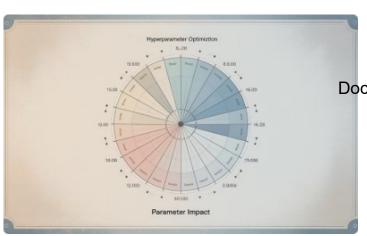
Get features from feature store. Create training datasets with proper splitting.

Model Experimentation



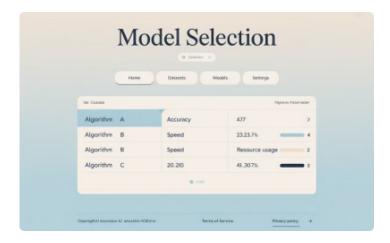
Experiment tracking

Log all model runs with parameters, metrics, and artifacts. Tools like MLflow make results searchable.



Hyperparameter tuning

Systematically search for optimal parameters. Use Bayesian optimization or grid search approaches. Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.



Model selection

Compare models based on performance metrics. Consider tradeoffs between accuracy and inference speed.

Model validation

Performance metrics

Measure accuracy, precision, recall, and business-specific KPIs. Set clear thresholds for acceptance.

Fairness assessment

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Test for bias across protected attributes. Ensure model performs equally well across all groups.

Evaluate performance across different segments.
Identify underperforming slices that need attention.

Data slicing

Validation tests

Run automated test suite on models. Validate model behavior on edge cases.



Training pipeline best practices

Deterministic training

Set random seeds for reproducibility. Document all randomness sources. Use same environment for reruns.

Extensive logging

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Log all parameters, metrics, and artifacts. Track resource usage and training time.

Controlled data splitting

Use stratified sampling for imbalanced data. Ensure test data resembles production data.

Configuration-driven

Drive experiments via configuration files. Avoid hardcoded parameters in code.

Inference pipeline overview

Feature retrieval

Get latest features from feature store

Model loading

Load appropriate model version

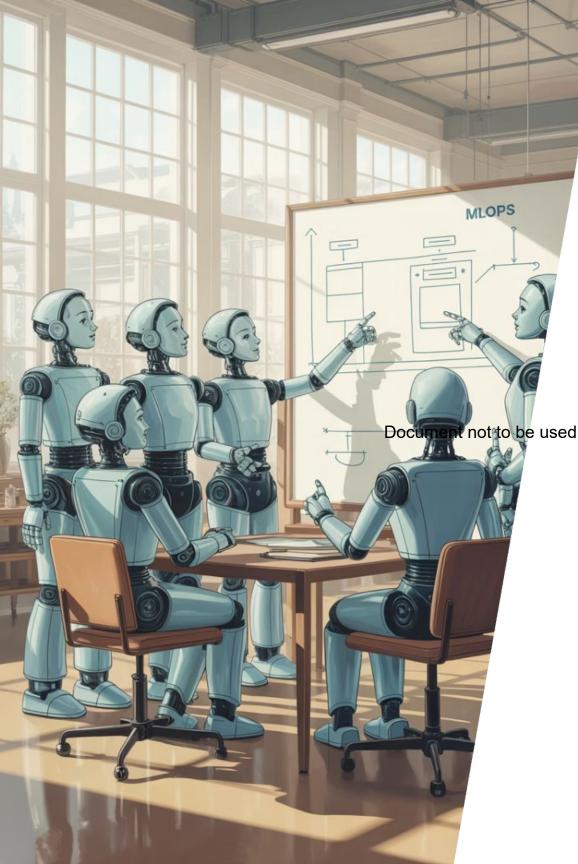
Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Post-processing

Format and enhance predictions

Prediction generation

Apply model to features



Practical exercice

Map your ML project

Take a current or planned ML project (e.g., phosphate production prediction). Map it to the FTI architecture. Identify components for each pipeline.

<mark>not to</mark> be used for teaching Ain ghis test வெடும் the author, Dr. Mohamed ASSELLAOU.

Determine clear inputs and outputs for each pipeline. Document data formats and schemas.

Assign responsibilities

Outline which team members would own each pipeline. Consider required skills and knowledge.



Batch vs. Real-time Inference

Characteristic	Batch Inference	Real-time Inference
Timing	Periodic (hourly/daily)	Immediate (milliseconds)
Volume	Large data sets	Single instances
Latency	Can be high	Must be low
Infrastructure	Optimized for throughput	Optimized for latency
Use Cases	Reports, daily recommendations	Fraud detection, search results

Model serving technologies

API frameworks

- Flask: Simple Python webDocument not to be used for tears in the faulth or the lauth or the laut framework
- FastAPI: Modern, high-performance framework
- Django: Full-featured web framework

Model servers

- TorchServe: For PyTorch models
- ONNX Runtime: For frameworkagnostic serving

Orchestration

- Kubernetes: Container orchestration
- Vertex Ai Pipeline, Kubeflow

Inference Performance Optimization

Model optimization

- Quantization: Reduce precision
- weights
- Distillation: Create smaller student models
- Graph optimization: Fuse operations

Serving optimization

- Batching: Process multiple
- Caching: Store common predictions
- Precomputing: Generate results ahead of time
- Asynchronous processing: Nonblocking APIs

Infrastructure optimization

- Horizontal scaling: Add more
- Vertical scaling: More powerful machines
- Hardware acceleration: GPUs, **TPUs**
- Deployment location: Edge vs. cloud

Deployment & Monitoring

Deployment

Safely releasing models to production

Monitoring

Tracking model health and performance

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Governance

Ensuring compliance and documentation

Retraining

Updating models with new data

Deployment strategies

Strategy	Key Benefit	Challenge
Blue/Green	Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU. Instant rollback Resource intensive	
Canary	Limited exposure	Complex traffic management

Blue/Green Deployment

Process

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOUSES

- Two parallel environments
- Blue: current version
- Green: new version

- Zero downtime
- Quick switch
- Simple rollback

- mod / toolle/ too.
 - Database migrations

Double resources

Canary deployment

Traffic shifting

Increase exposure on success

Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU.

Metrics monitoring

Watch for issues in real traffic

Gradual rollout

Start with small traffic percentage

Continuous integration/Continuous delivery

Code changes

Developers commit code to version control. Changes trigger automated workflows.

Automated testing

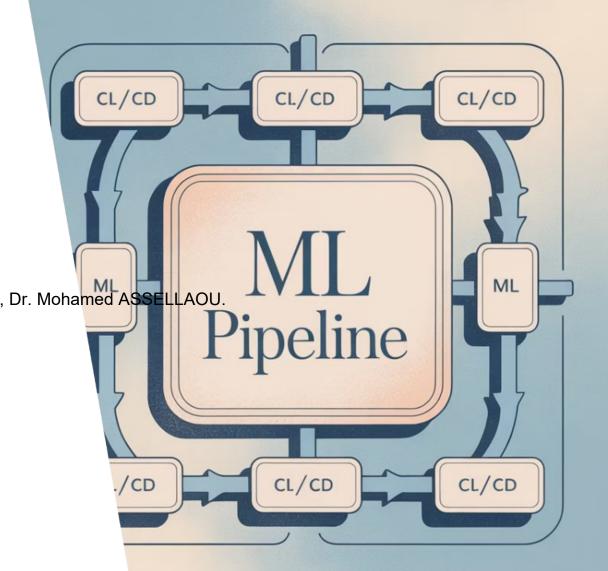
Document not to be used for teaching. All rights reserved by the author, Dr. Mohamed ASSELLAOU. Run unit, integration, and model quality tests. Verify all components work together.

Package & release

Build deployment artifacts. Create containers or packages for each pipeline.

Deployment

Safely deploy to production using blue/green or canary strategies. Monitor rollout health.



Model monitoring

Data drift detection

Compare production data ចាំទី៩ម៉ោងម្នាក់លេខ សេរីម៉ោជន៍ខ្លាំក្រាំម៉ូខែaching. All rights reserved by the author, Dr. ហ្វែងក្រាម៉ា Assetulasou. and target. data. Identify shifts in input patterns.

Performance monitoring

Track accuracy, precision, recall, and business metrics. Alert when performance degrades.

Concept drift detection

Detect changes in relationships Often requires delayed ground truth.

Alerting system

Set up automated alerts for drift thresholds. Create action plans for different alert types.

MLOps observability

Three pillars of observability

Comprehensive visibility into ML systems requires three complementary approaches. Document not to be used for teaching. All rights reservite abijution author for the abijution a

- Logs: Detailed system events
- Metrics: Quantitative performance data
- Traces: Request flows through system

ML-specific observability

Machine learning systems need additional monitoring beyond

- Feature distributions
- **Prediction distributions**
- Model performance metrics
- Data quality metrics
- **Explanation metrics**



Conclusion

Start Small

Begin with simple pipelines. Iterate and improve continuously.

Focus on Automation

Eliminate manual steps gradually. Prioritize reproducibility in all processes Document not to be used for teaching. All rights reserved by the author through the author of the beautiful and the control of the control

Invest in Tooling

Choose appropriate tools for your scale. Build or buy suitable infrastructure.

Build Capabilities

Develop cross-functional MLOps skills. Remember that MLOps is about people, process, and technology.