# *From Theory To Practice:*
# Implementing Machine Learning Solutions Safely and Effectively in the Clinical Laboratory

## Nick Spies, MD

*Co-Medical Director - Applied AI & Clinical Chemistry, ARUP Laboratories*
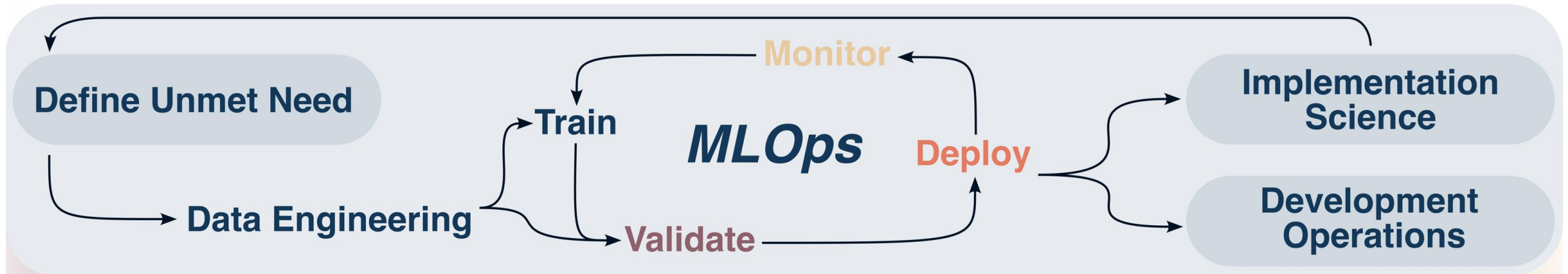
*Assistant Professor - University of Utah*
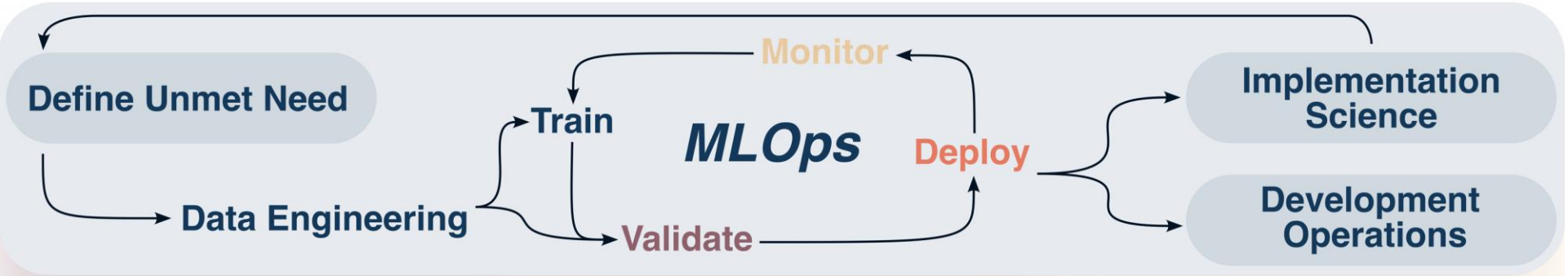
# Disclosures

- I have no relevant conflicts of interest to disclose.

# Learning Objectives

- Define key roles and responsibilities in the machine learning life-cycle.

- Explore techniques for validating, deploying, and monitoring models.

- Reinforce these concepts within a relevant, lab-based example.

# The Machine Learning Life Cycle

**Define Unmet Need** → **Data Engineering**

**Train** · **MLOps** · **Monitor** · **Deploy** · **Validate**

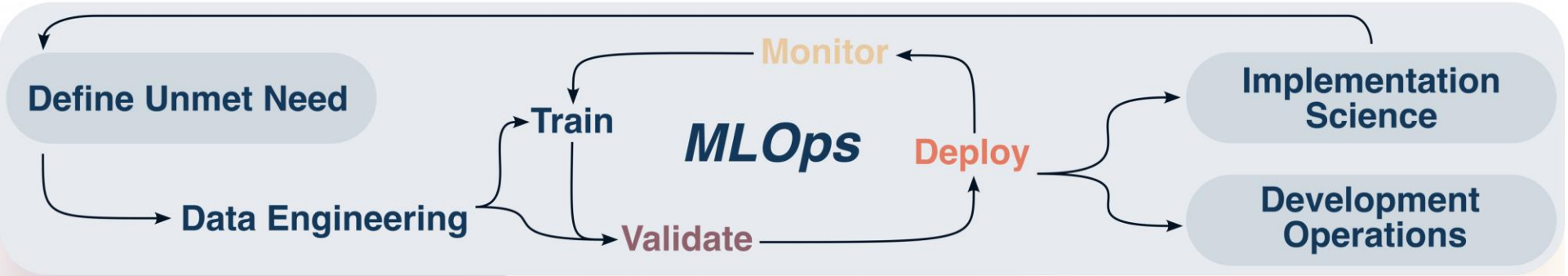**Implementation Science** · **Development Operations**

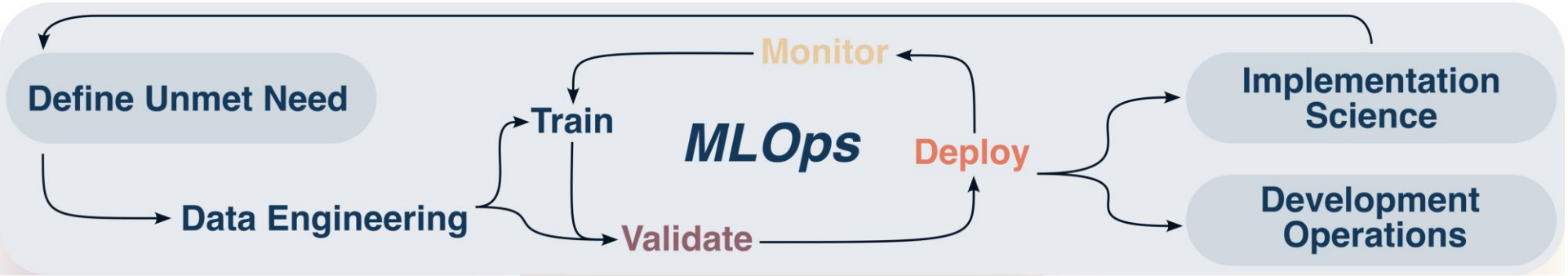## Validation

Metric Selection

Target Label Appraisal

Prediction Calibration

Generalizability & Applicability Assessment

Measuring Inequity & Algorithmic Fairness

Explainability & Interpretability

## Deployment

Production Environments & the IT Stack

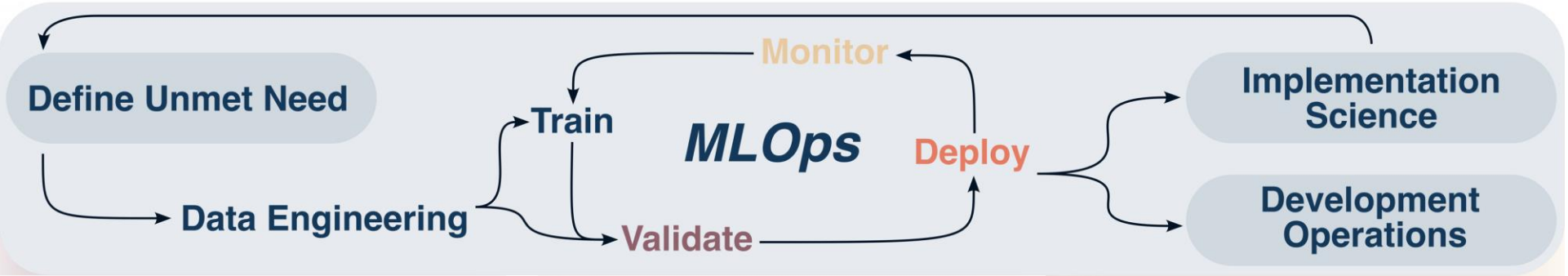Latency, Uptime, & Failure Modes Analysis

CI/CD & Logging

*Development Operations*

*Implementation Science*

Integration Domains

Human-in-the-Loop vs. Automated Inference

Governance & RACI Analysis
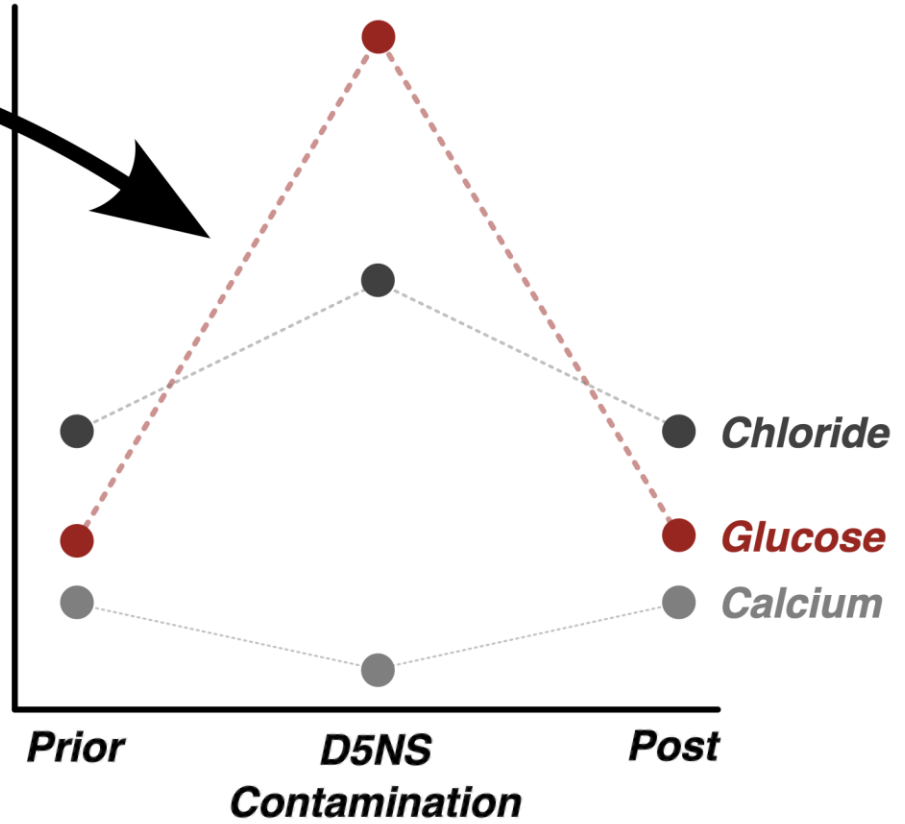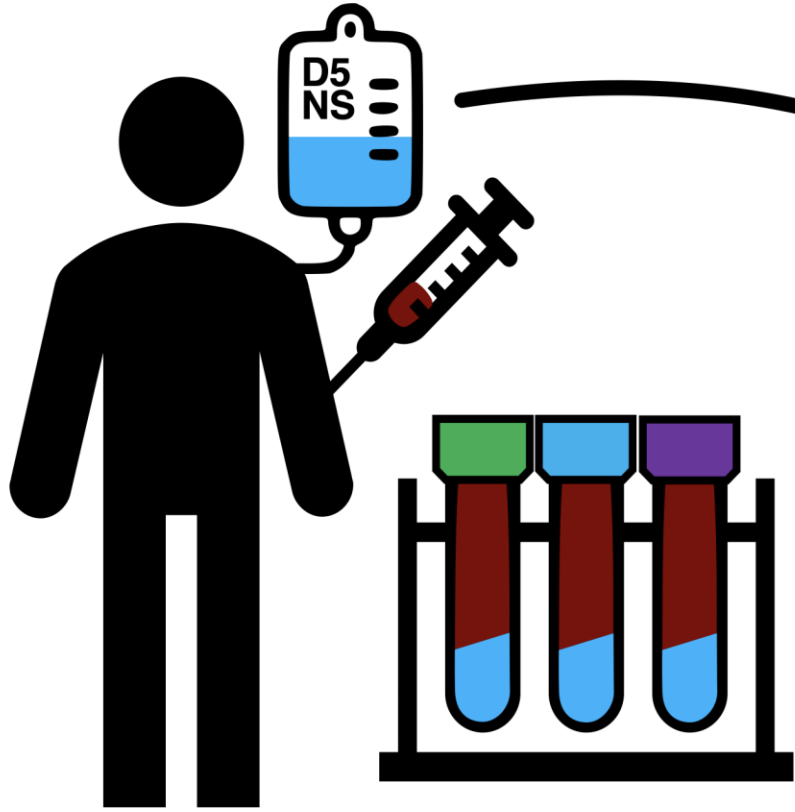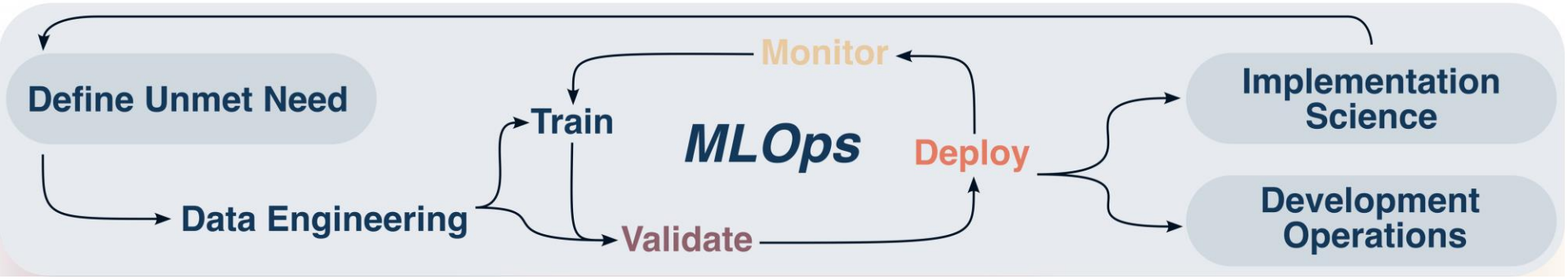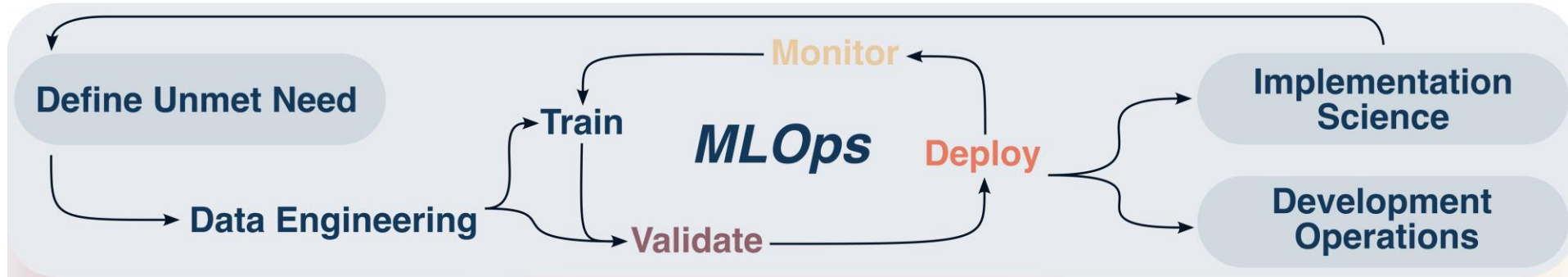
## Monitoring

Input & Prediction Drift

Prediction Impact Analysis

Online Performance Assessment

Model Updating Strategies

Algorithmic Stewardship Principles

Algorithm Inventories & Managing Conflicting Models

**Validation**

Metric Selection

Target Label Appraisal

Prediction Calibration

Generalizability & Applicability Assessment

Measuring Inequity & Algorithmic Fairness

Explainability & Interpretability

**Define Unmet Need**

**Train**

**MLOps**

**Monitor**

**Deploy**

**Implementation Science**

**Data Engineering**

**Validate**

**Development Operations**

## Deployment

Production Environments & the IT Stack

Latency, Uptime, & Failure Modes Analysis

CI/CD & Logging

*Development Operations*

*Implementation Science*

Integration Domains

Human-in-the-Loop vs. Automated Inference

Governance & RACI Analysis

**MLOps**

Define Unmet Need → Data Engineering → Train → Validate → Monitor → Deploy → Implementation Science / Development Operations

**Monitoring**

Input & Prediction Drift

Prediction Impact Analysis

Online Performance Assessment

Model Updating Strategies

Algorithmic Stewardship Principles

Algorithm Inventories & Managing Conflicting Models

# Validating, Implementing, and Monitoring Machine Learning Solutions in the Clinical Laboratory Safely and Effectively 🔓

Nicholas C Spies ✉, Christopher W Farnsworth, Sarah Wheeler, Christopher R McCudden

# Defining A Machine Learning Pipeline

# Building The Model

```r
1  # Load required libraries
2  library(tidymodels)
3  library(arrow)
4
5  # Load the data
6  train <-
7    read_feather("https://figshare.com/ndownloader/files/45407401") |>
8    select(contam_comment, bun:sodium) |>
9    mutate(contam_comment = factor(contam_comment))
10
11 # Define the feature recipe
12 recipe <- recipe(contam_comment ~ ., data = train)
13
14 # Define the model
15 model <- boost_tree(mode = "classification") |> set_engine("xgboost")
16
17 # Create the workflow
18 workflow <- workflow() |> add_recipe(recipe) |> add_model(model)
19
20 # Fit the model
21 fit <- workflow |> fit(data = train)
```

# Testing The Model

| | |
|---|---|
| *Sodium* | 147 ! |
| *Potassium* | 3.5 |
| *Chloride* | 119 ! |
| *CO2* | 17 ! |
| *Creatinine* | 0.9 |
| *BUN* | 22 |
| *Calcium* | 6.6 ! |
| *Glucose* | 86 |

**Model**



**Output**

**Output:** *0.97*
**Label:** *Positive*
**Applicable:** ✓
**Explanation:** ✓

# How can we *validate* a developed model?
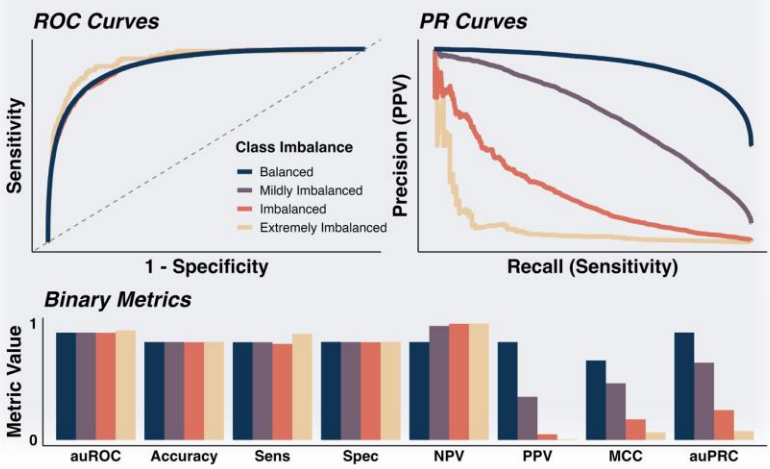
# Successful *Validation* of Machine Learning Pipelines

# Successful *Validation* of Machine Learning Pipelines

**Inputs**

**Model**

**Output**

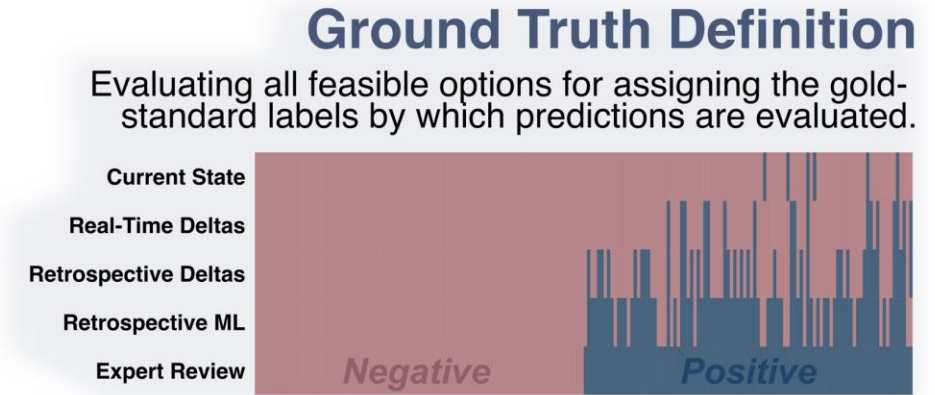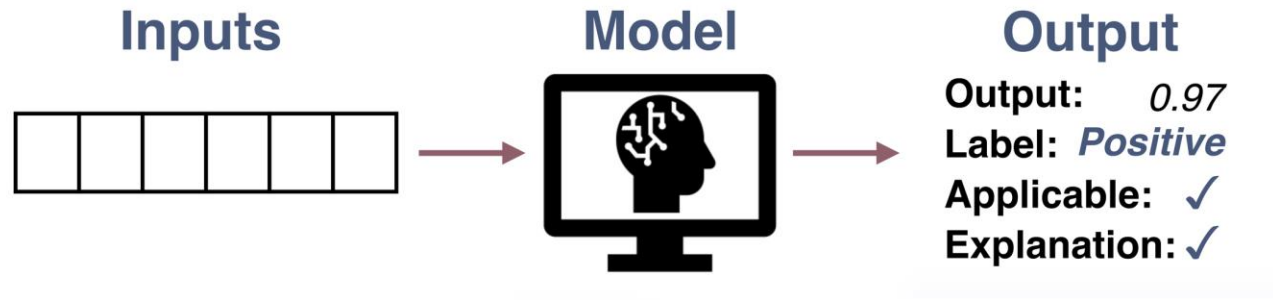Output:        *0.97*
Label: *Positive*
Applicable: ✓
Explanation: ✓

**Ground Truth Definition**

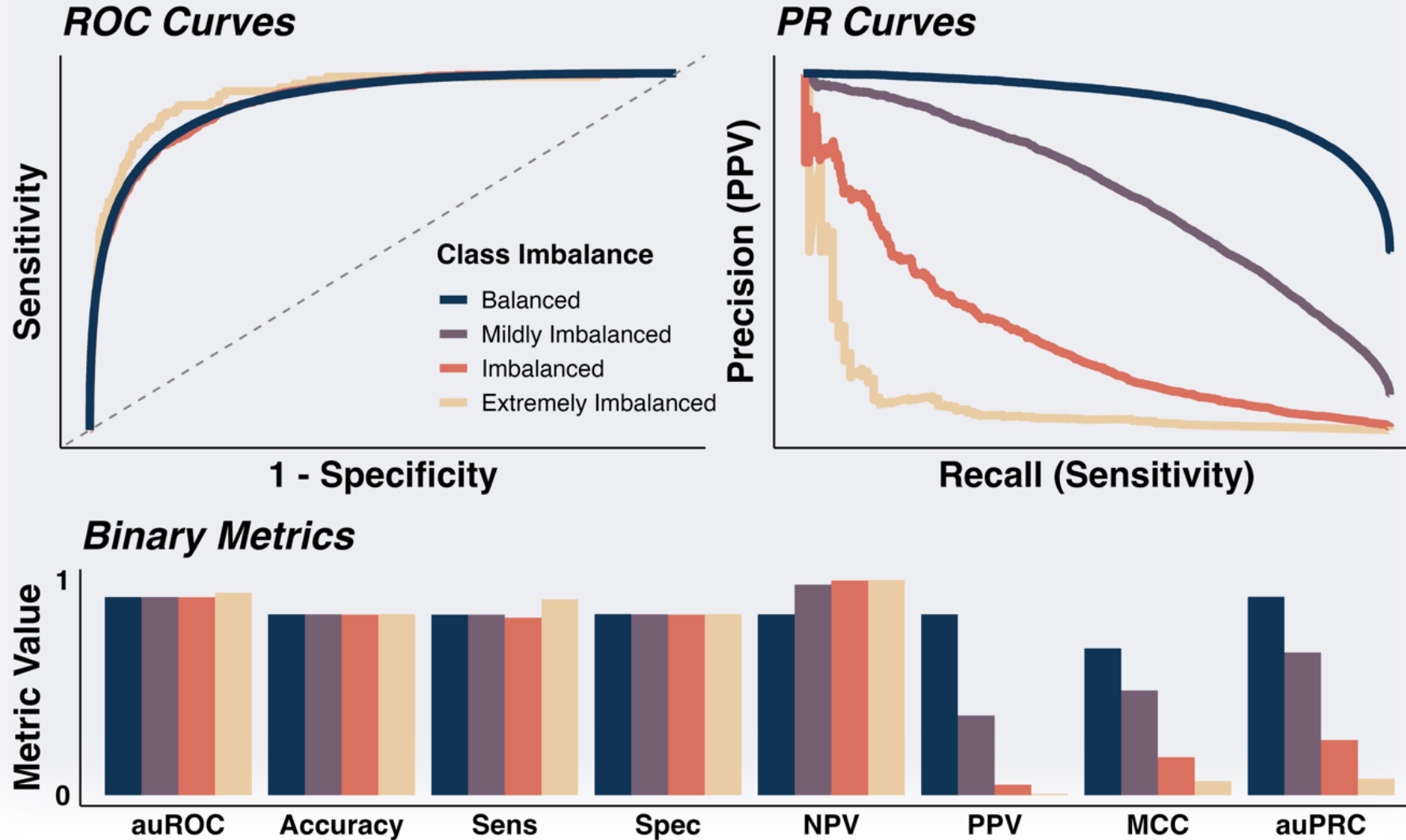Evaluating all feasible options for assigning the gold-standard labels by which predictions are evaluated.

Current State

Real-Time Deltas

Retrospective Deltas

Retrospective ML

Expert Review

*Negative*                              *Positive*

# Successful *Validation* of Machine Learning Pipelines

**Inputs**

**Model**

**Output**

**Output:** *0.97*
**Label:** *Positive*
**Applicable:** ✓
**Explanation:** ✓

## Ground Truth Definition

Evaluating all feasible options for assigning the gold-standard labels by which predictions are evaluated.

Current State
Real-Time Deltas
Retrospective Deltas
Retrospective ML
Expert Review

*Negative*          *Positive*

## Metric Selection

Appropriately measuring pipeline performance.

*ROC Curves*

Sensitivity

1 - Specificity

Class Imbalance
Balanced
Mildly Imbalanced
Imbalanced
Extremely Imbalanced

*PR Curves*

Precision (PPV)

Recall (Sensitivity)

*Binary Metrics*

Metric Value

1

0

auROC    Accuracy    Sens    Spec    NPV    PPV    MCC    auPRC

# Successful *Validation* of Machine Learning Pipelines



**Inputs**

**Model**

**Output**
Output: 0.97
Label: *Positive*
Applicable: ✓
Explanation: ✓

**Ground Truth Definition**
Evaluating all feasible options for assigning the gold-standard labels by which predictions are evaluated.

Current State
Real-Time Deltas
Retrospective Deltas
Retrospective ML
Expert Review

*Negative*  *Positive*

**Metric Selection**
Appropriately measuring pipeline performance.

*ROC Curves*

Sensitivity

Class Imbalance
Balanced
Mildly Imbalanced
Imbalanced
Extremely Imbalanced

1 - Specificity

*PR Curves*

Precision (PPV)

Recall (Sensitivity)

*Binary Metrics*

Metric Value

auROC  Accuracy  Sens  Spec  NPV  PPV  MCC  auPRC

**Threshold Optimization**
Defining the optimal decision boundaries to convert continuous outputs into class labels

*Negative*  Equivocal  *Positive*
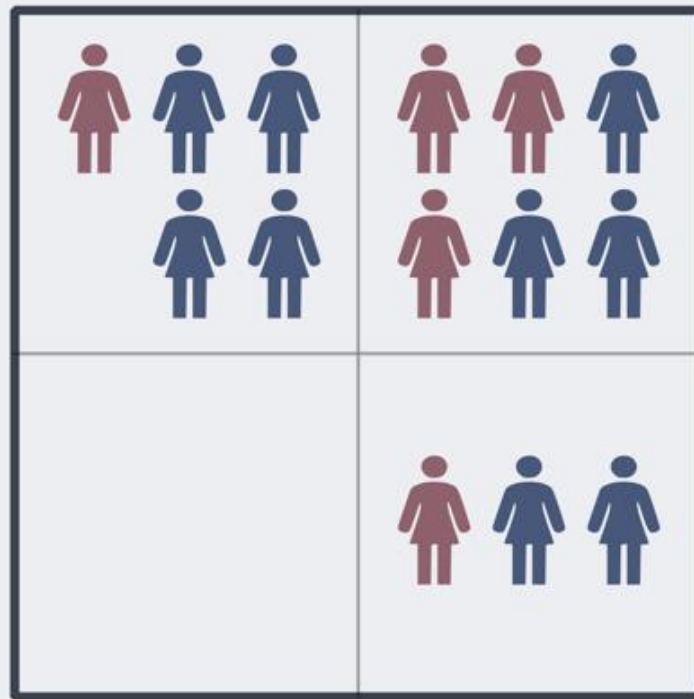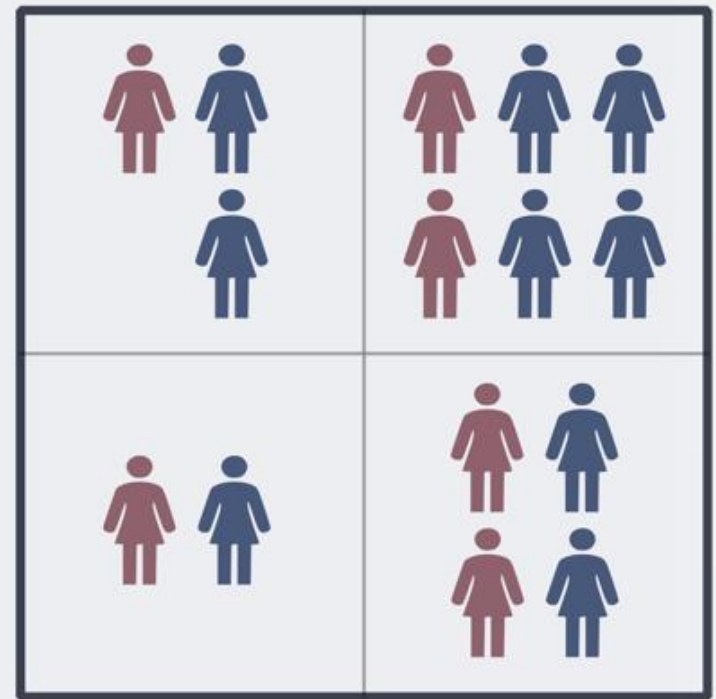
**Predicted Probability**

# Algorithmic Fairness

Ensuring the model does not introduce or exacerbate inequity across populations

**Demographic Parity**

**Equalized Odds**
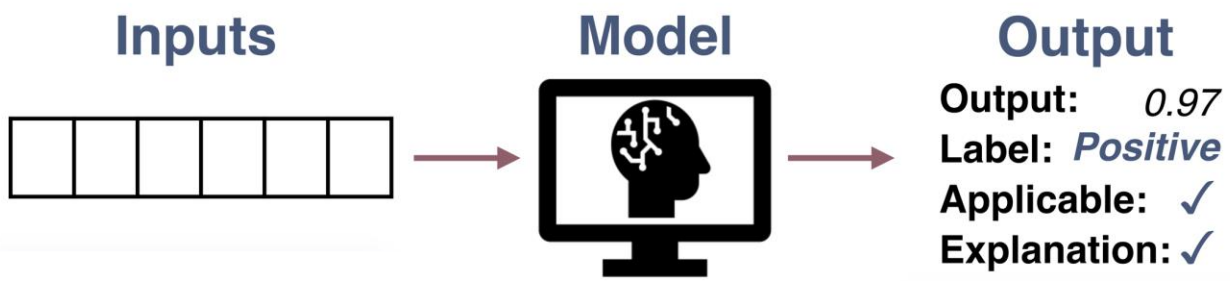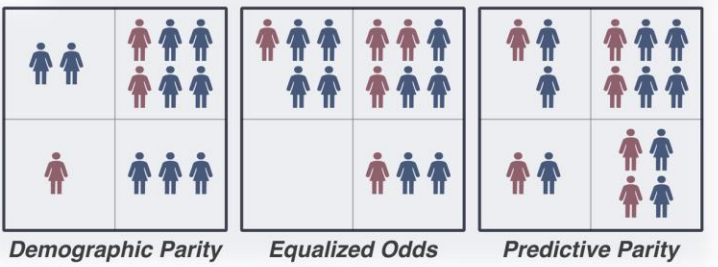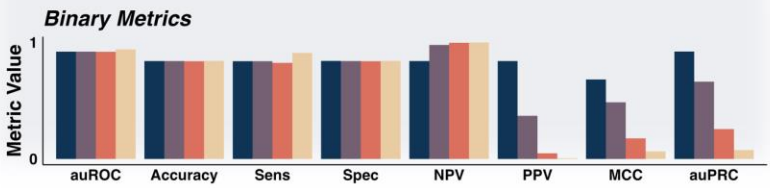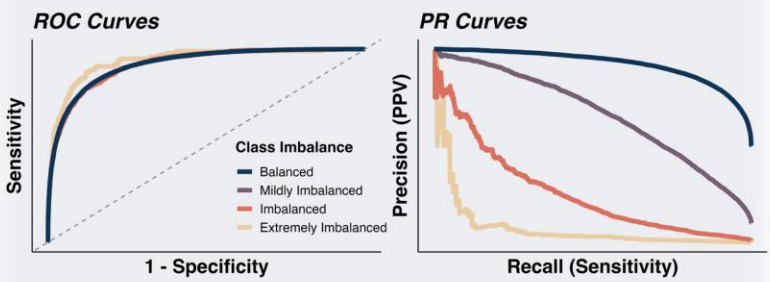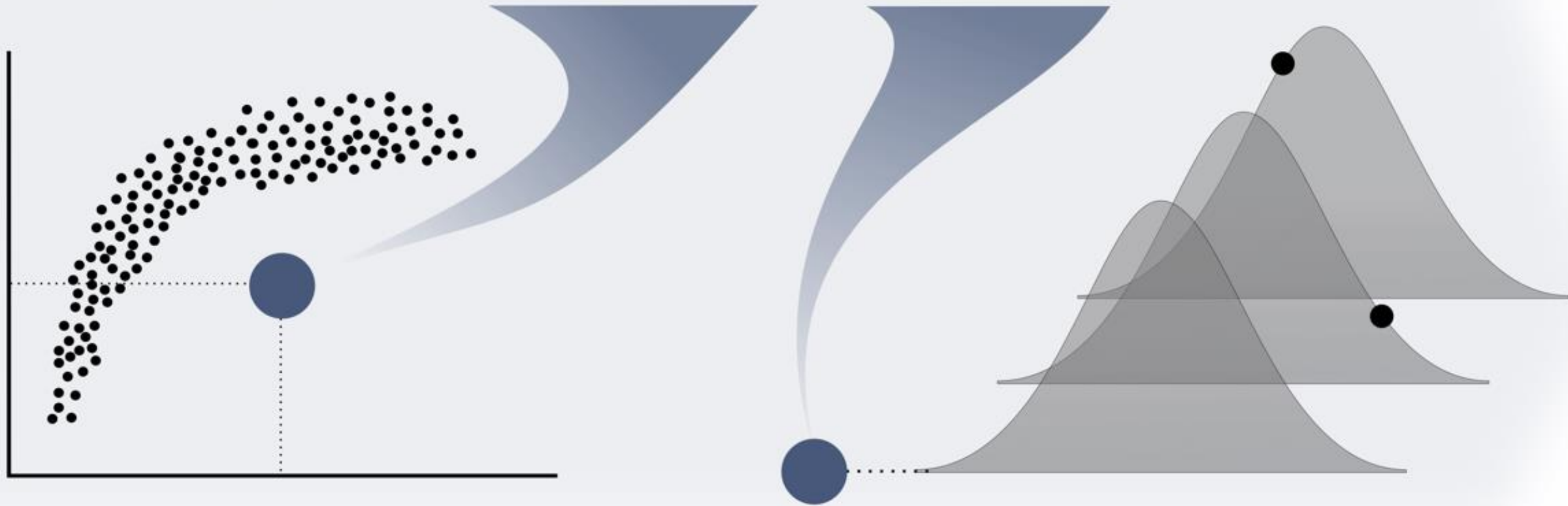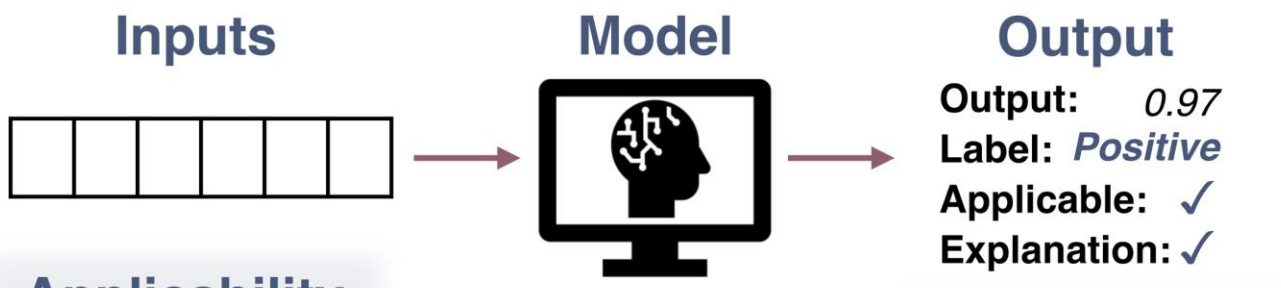
**Predictive Parity**

# Applicability

Identifying inputs that diverge from training data *across* or *within* features.

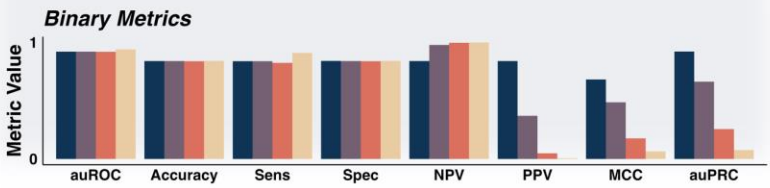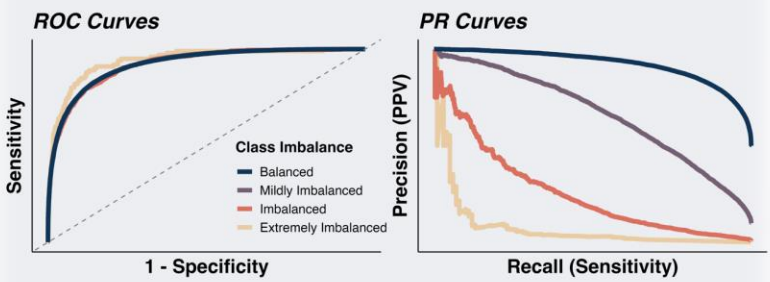# Successful *Validation* of Machine Learning Pipelines

# Explainability

Estimating the impact of each feature on the final prediction

**Feature Value**

Low          High

**Impact on Prediction**

# How do we *implement* a validated model?

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities



### *Subject Matter Experts*

- Align implementation to fit unmet clinical need.
- Evaluate failure modes and off-target effects.

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities

### *Subject Matter Experts*

- Align implementation to fit unmet clinical need.
- Evaluate failure modes and off-target effects.

### *Data Scientists & Data Engineers*

- Build and evaluate models for deployment.
- Optimize storage and retrieval of input data.

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities

### *Subject Matter Experts*

- Align implementation to fit unmet clinical need.
- Evaluate failure modes and off-target effects.

### *Data Scientists & Data Engineers*

- Build and evaluate models for deployment.
- Optimize storage and retrieval of input data.

### *Software & ML Engineers*

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.

### *MLOps*

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities

### *Subject Matter Experts*

- Align implementation to fit unmet clinical need.
- Evaluate failure modes and off-target effects.

### *Data Scientists & Data Engineers*

- Build and evaluate models for deployment.
- Optimize storage and retrieval of input data.

### *Software & ML Engineers*

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.

### *MLOps*

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.

### *Information Technology & Systems*

- Maintain interfaces for ML inputs and outputs.
- Develop infrastructure and allocate resources.

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities

### *Subject Matter Experts*

- Align implementation to fit unmet clinical need.
- Evaluate failure modes and off-target effects.

### *Data Scientists & Data Engineers*

- Build and evaluate models for deployment.
- Optimize storage and retrieval of input data.

## Terms and Technologies



Validated Model → Container

# Successful *Implementation* of Machine Learning Pipelines

## Key Roles and Responsibilities

## Terms and Technologies

Validated Model

Container

### *Deployment*

***Deploying*** a model to a ***server*** makes it **accessible** to other systems or users within a **network**.
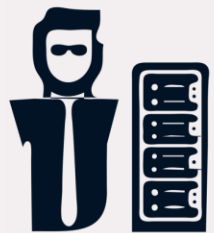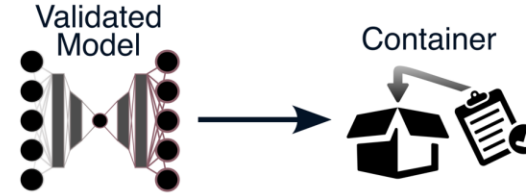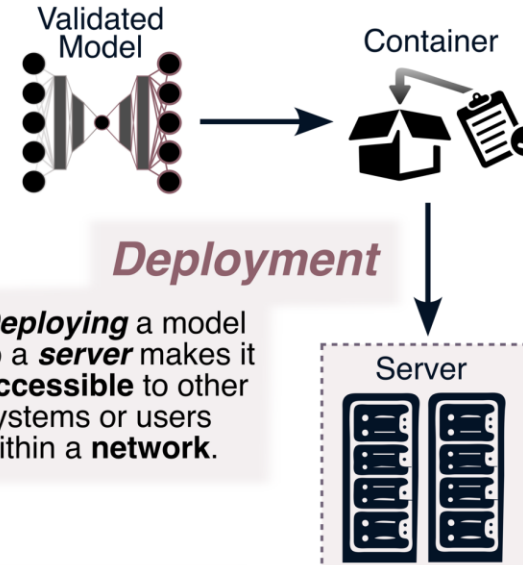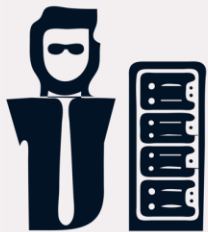
Server

### Software & ML Engineers

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.

### MLOps

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.

### Information Technology & Systems

- Maintain interfaces for ML inputs and outputs.
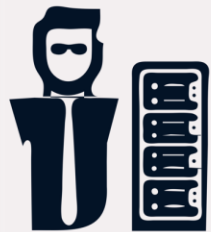- Develop infrastructure and allocate resources.

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities

## Terms and Technologies
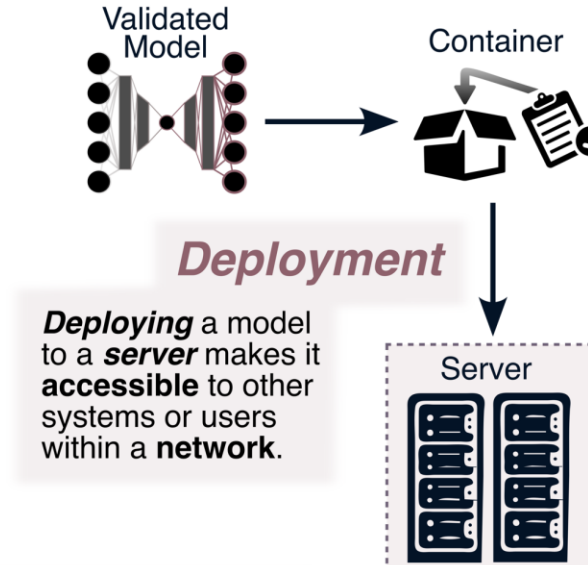
**Validated Model**

**Container**

### Development Environment

Before deployment, the **full pipeline** should be **robustly tested** in an offline "sandbox" that is **completely isolated** from clinical workflows.

### Deployment

**Deploying** a model to a **server** makes it **accessible** to other systems or users within a **network**.

**Server**

### Software & ML Engineers

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.

### MLOps

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.
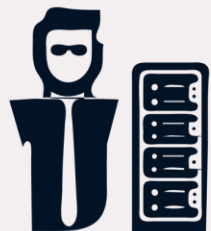
### Information Technology & Systems

- Maintain interfaces for ML inputs and outputs.
- Develop infrastructure and allocate resources.

# Successful *Implementation* of Machine Learning Pipelines

## Key Roles and Responsibilities

## Terms and Technologies

Validated Model

Container

### Development Environment

Before deployment, the **full pipeline** should be **robustly tested** in an offline "**sandbox**" that is **completely isolated** from clinical workflows.

### Deployment

*Deploying* a model to a *server* makes it **accessible** to other systems or users within a **network**.

### Software & ML Engineers

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.

### MLOps

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.

Server

API

### Logging Metrics

**Latency:** *Turn-around time for predictions.*

**Uptime:** *% of time a prediction can be made.*

**Scalability:** *Change in latency/uptime with increased volume.*

### Information Technology & Systems

- Maintain interfaces for ML inputs and outputs.
- Develop infrastructure and allocate resources.

# *Successful Implementation of Machine Learning Pipelines*

## Key Roles and Responsibilities

### Data Scientists & Data Engineers

- Build and evaluate models for deployment.
- Optimize storage and retrieval of input data.

### Software & ML Engineers

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.
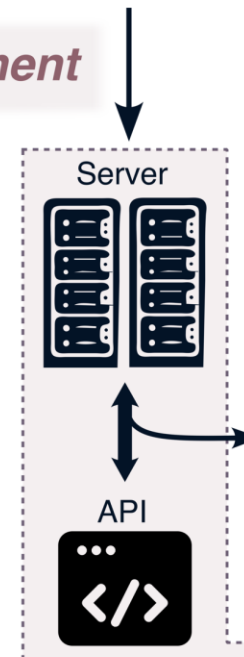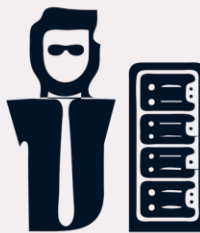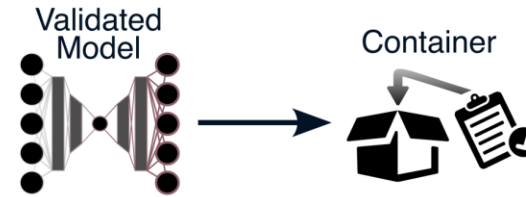
### Information Technology & Systems

- Maintain interfaces for ML inputs and outputs.
- Develop infrastructure and allocate resources.

## Terms and Technologies

Validated Model

Container

### *Deployment*

*Deploying* a model to a *server* makes it **accessible** to other systems or users within a **network**.

### *MLOps*

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.

### *Development Environment*

Before deployment, the **full pipeline** should be **robustly tested** in an offline "**sandbox**" that is **completely isolated** from clinical workflows.
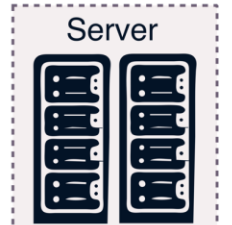
### *Logging Metrics*

**Latency:** *Turn-around time for predictions.*
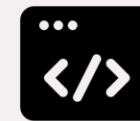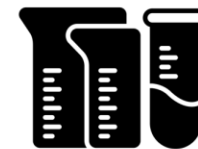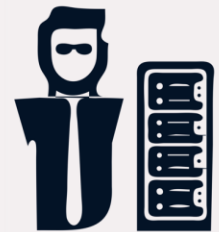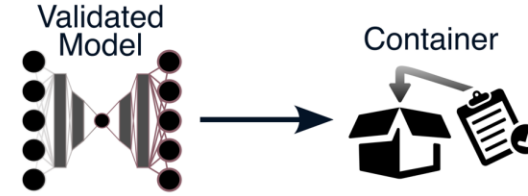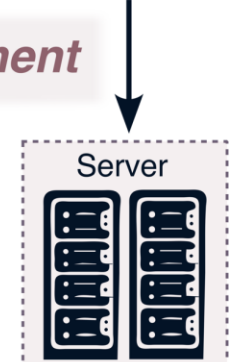
**Uptime:** *% of time a prediction can be made.*

**Scalability:** *Change in latency/uptime with increased volume.*

Server

API

### *Production Environment*

Instruments

Middleware

Client

# Successful *Implementation* of Machine Learning Pipelines

## Key Roles and Responsibilities

### *Subject Matter Experts*

- Align implementation to fit unmet clinical need.
- Evaluate failure modes and off-target effects.

### *Data Scientists & Data Engineers*

- Build and evaluate models for deployment.
- Optimize storage and retrieval of input data.

### *Software & ML Engineers*

- Build robust and secure prediction pipelines.
- Implement best practices in DevOps/MLOps.

### Information Technology & Systems

- Maintain interfaces for ML inputs and outputs.
- Develop infrastructure and allocate resources.

## Terms and Technologies

Validated Model

Container

### *Development Environment*

Before deployment, the **full pipeline** should be **robustly tested** in an offline "**sandbox**" that is **completely isolated** from clinical workflows.

### *Deployment*

*Deploying* a model to a *server* makes it **accessible** to other systems or users within a **network**.
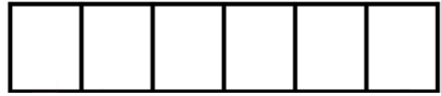
Server

### *MLOps*

The framework for **building**, **deploying**, and **monitoring** end-to-end ML solutions in live, **production environments** safely and effectively.

API

### *Logging Metrics*

**Latency:** *Turn-around time for predictions.*

**Uptime:** *% of time a prediction can be made.*

**Scalability:** *Change in latency/uptime with increased volume.*

### *Production Environment*

Instruments

Middleware

Client

Users

# How do we *monitor* an implemented model?

# Successful *Monitoring* of Machine Learning Pipelines
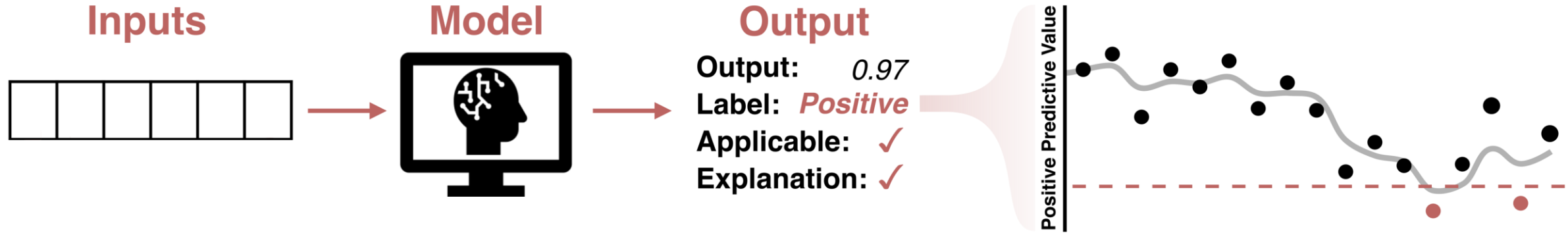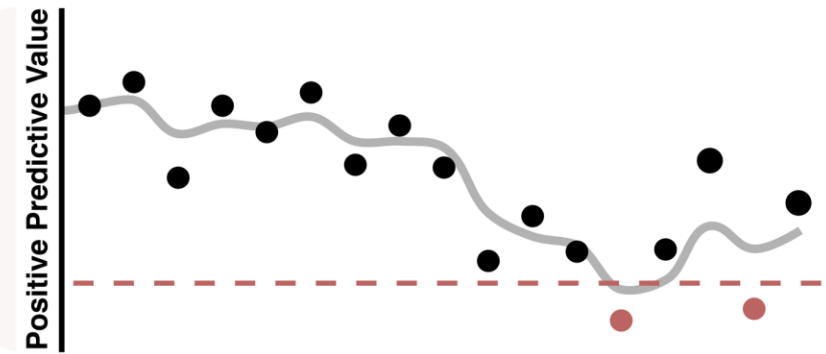
**Inputs**

**Model**

**Output**

Output: *0.97*
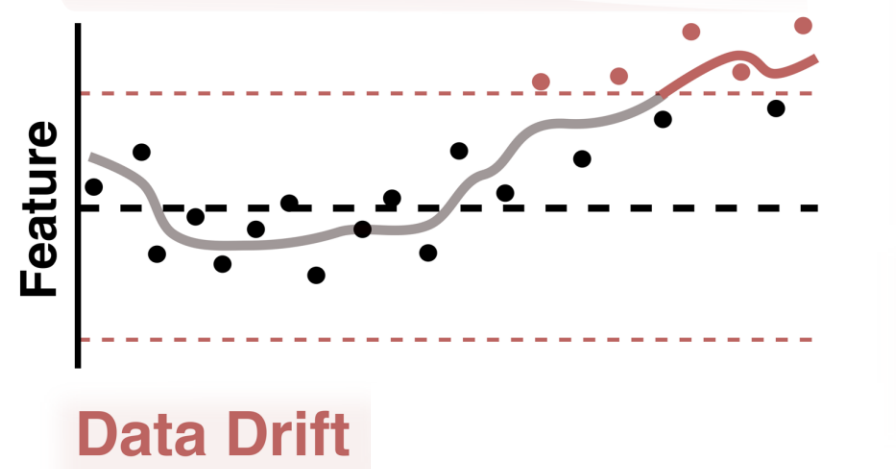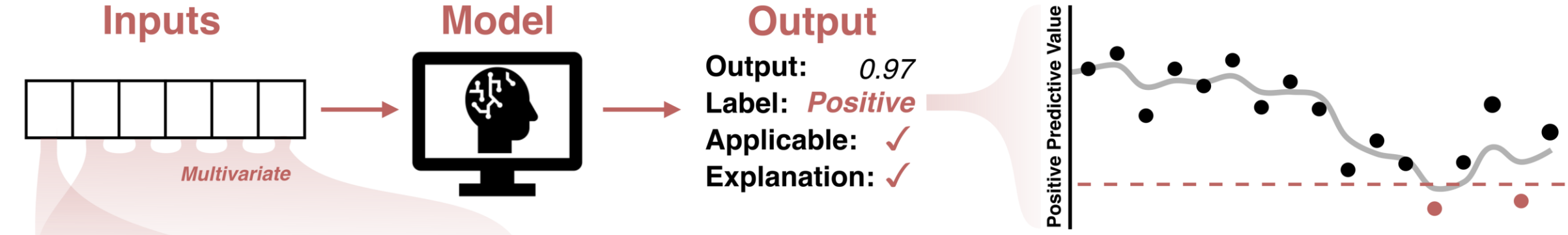Label: *Positive*
Applicable: ✓
Explanation: ✓

Positive Predictive Value

# *Successful **Monitoring** of Machine Learning Pipelines*

**Inputs**

**Model**

**Output**

**Output:** *0.97*
**Label:** *Positive*
**Applicable:** ✓
**Explanation:** ✓

Positive Predictive Value

## Performance Drift

***Data drift*** or ***concept drift*** causes ML pipelines to lose performance. Closed-loop systems are crucial for identifying when this occurs.

# Successful *Monitoring* of Machine Learning Pipelines



**Inputs**

**Model**

**Output**

Output:      *0.97*
Label:      *Positive*
Applicable:  ✓
Explanation: ✓

Positive Predictive Value

## Performance Drift

***Data drift*** or ***concept drift*** causes ML pipelines to lose performance. Closed-loop systems are crucial for identifying when this occurs.
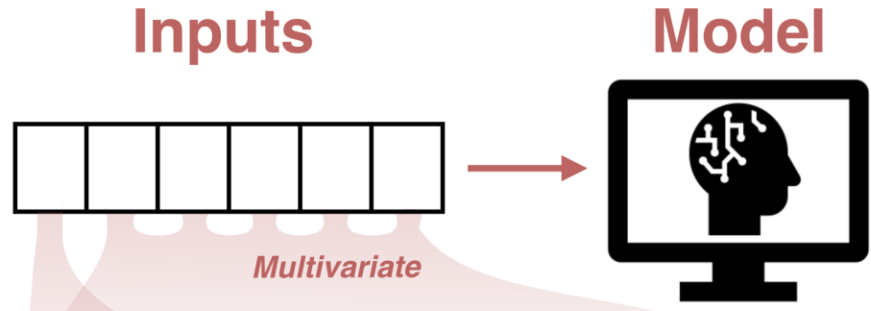
Feature

**Data Drift**

# Successful *Monitoring* of Machine Learning Pipelines



**Inputs**

**Model**

**Output**

Output: 0.97
Label: *Positive*
Applicable: ✓
Explanation: ✓

*Multivariate*

*Univariate*

Positive Predictive Value

## Performance Drift

***Data drift*** or ***concept drift*** causes ML pipelines to lose performance. Closed-loop systems are crucial for identifying when this occurs.

Feature

## Data Drift

# *Successful Monitoring of Machine Learning Pipelines*

## Inputs

**Multivariate**

*Univariate*

## Model

## Output

**Output:** 0.97
**Label:** *Positive*
**Applicable:** ✓
**Explanation:** ✓

Positive Predictive Value

Feature

## Performance Drift

***Data drift*** or ***concept drift*** causes ML pipelines to lose performance. Closed-loop systems are crucial for identifying when this occurs.

## Data Drift

| | *Detection* | *Correction* |
|---|---|---|
| *Uni-* | - Threshold Flags<br>- Moving Averages | - Input Preprocessing<br>- Analyzer Recalibration |
| *Multi-* | - Principal Components<br>- Mahalanobis Distance | - Input Transformation<br>- Model Retraining |

# *Successful Monitoring of Machine Learning Pipelines*

## Inputs

**Multivariate**

*Univariate*

## Model

## Output

**Output:** 0.97
**Label:** *Positive*
**Applicable:** ✓
**Explanation:** ✓

Positive Predictive Value

Feature

## Performance Drift

***Data drift*** or ***concept drift*** causes ML pipelines to lose performance. Closed-loop systems are crucial for identifying when this occurs.

## Concept Drift

Occurs when **real-world labels diverge** from **training** labels. Difficult to detect and correct without input from **subject-matter experts**.

## Data Drift

|        | Detection | Correction |
|--------|-----------|------------|
| *Uni-* | - Threshold Flags<br>- Moving Averages | - Input Preprocessing<br>- Analyzer Recalibration |
| *Multi-* | - Principal Components<br>- Mahalanobis Distance | - Input Transformation<br>- Model Retraining |

# *Successful Monitoring of Machine Learning Pipelines*

**Inputs**

**Model**

**Output**
**Output:** *0.97*
**Label:** *Positive*
**Applicable:** ✓
**Explanation:** ✓

*Multivariate*

*Univariate*

Positive Predictive Value

**Performance Drift**

***Data drift*** or ***concept drift*** causes ML pipelines to lose performance. Closed-loop systems are crucial for identifying when this occurs.

**Concept Drift**

Occurs when **real-world labels diverge** from **training** labels. Difficult to detect and correct without input from **subject-matter experts**.

Feature

**Data Drift**

*Uni-*

*Multi-*

| *Detection* | *Correction* |
|---|---|
| - Threshold Flags<br>- Moving Averages | - Input Preprocessing<br>- Analyzer Recalibration |
| - Principal Components<br>- Mahalanobis Distance | - Input Transformation<br>- Model Retraining |

**Updating Models**

Replacement models can be **continuously retrained and evaluated** to replace deteriorating models before they impact live workflows.

*Champion*

*Challengers*

# A Note On Regulatory Guidance

**Define Unmet Need**

Monitor

Train

**MLOps**

Deploy

**Implementation Science**

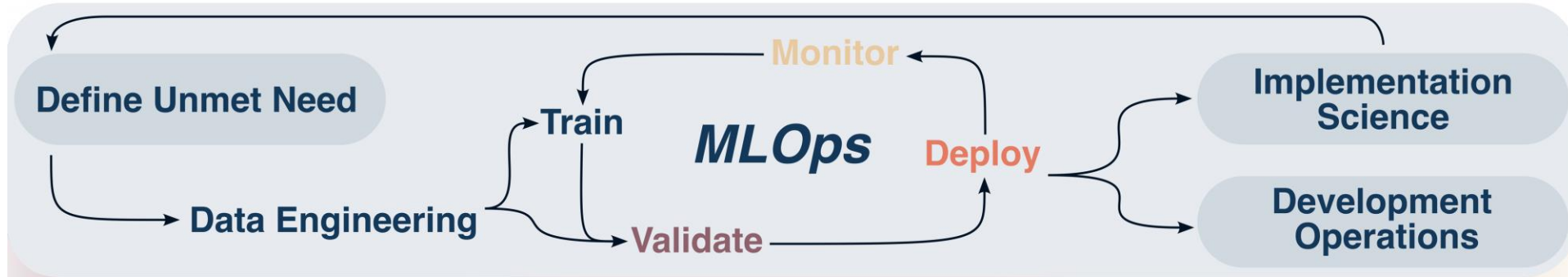Data Engineering

Validate

**Development Operations**

## Validation

Metric Selection

Target Label Appraisal

Prediction Calibration

Generalizability & Applicability Assessment

Measuring Inequity & Algorithmic Fairness

Explainability & Interpretability

## Deployment

Production Environments & the IT Stack

Latency, Uptime, & Failure Modes Analysis

CI/CD & Logging

*Development Operations*

*Implementation Science*

Integration Domains

Human-in-the-Loop vs. Automated Inference

Governance & RACI Analysis

## Monitoring

Input & Prediction Drift

Prediction Impact Analysis

Online Performance Assessment

Model Updating Strategies

Algorithmic Stewardship Principles

Algorithm Inventories & Managing Conflicting Models

# Questions?

Email: nick.spies@aruplab.com

U HEALTH
UNIVERSITY OF UTAH