# Data-Driven Classification of Canadian Crude Oils

Ryan Bulger
DTSA-5506
University of Colorado Boulder

# Problem Definition

## Assessing Current Crude Oil Classifications

### Grouping Limitations

Existing classifications of crude oils may not accurately reflect the **true similarities** among different types, leading to potential inefficiencies in management and utilization of these resources.

# Importance of Reclassifying Crude Oils

### Marketing

Improved classifications can enhance **targeted marketing strategies** effectively.

### Pricing

Accurate groupings lead to **fair pricing structures** in the market.

### Efficiency

Optimized classifications promote **efficient pipeline utilization** and resource management.

# Related Work
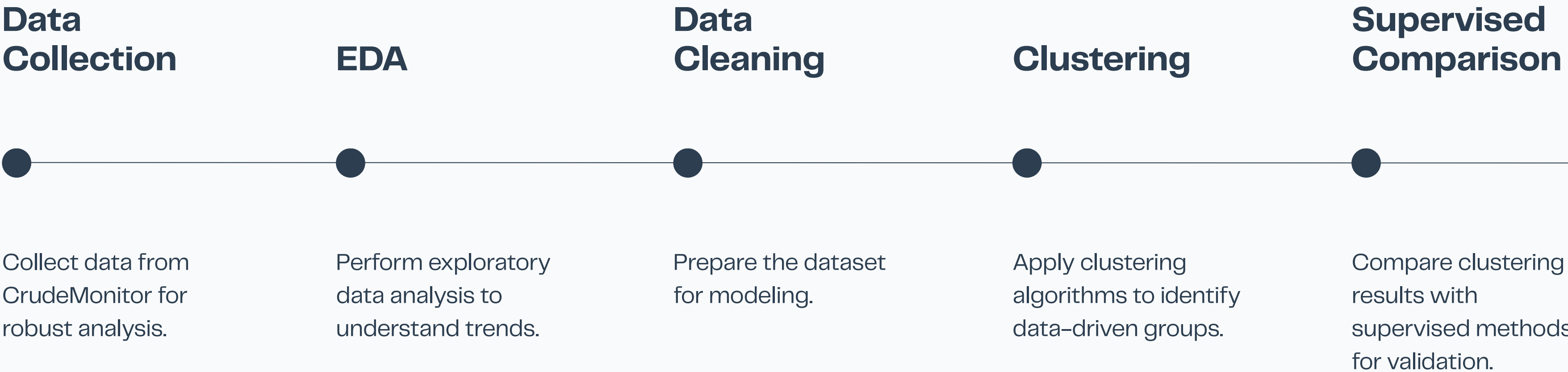
## Overview of traditional and modern methods

### Traditional Approaches

Traditional methods like API gravity and sulfur content often **fail to capture the complexities** of crude oil. These simplistic classifications may overlook significant variations in crude oil characteristics.

### Modern ML Studies

Recent machine learning approaches, such as **FTIR spectroscopy combined with clustering techniques**, have shown promise. However, they often rely on limited datasets, leading to potential inaccuracies in classifications.

# Project Steps

**Data Collection**

**EDA**

**Data Cleaning**

**Clustering**

**Supervised Comparison**

Collect data from CrudeMonitor for robust analysis.

Perform exploratory data analysis to understand trends.

Prepare the dataset for modeling.

Apply clustering algorithms to identify data–driven groups.

Compare clustering results with supervised methods for validation.

**COMPLETED**

# Data Collection

### CrudeMonitor API

Data was captured using CrudeMonitor's API endpoints.

### Dataset Shape

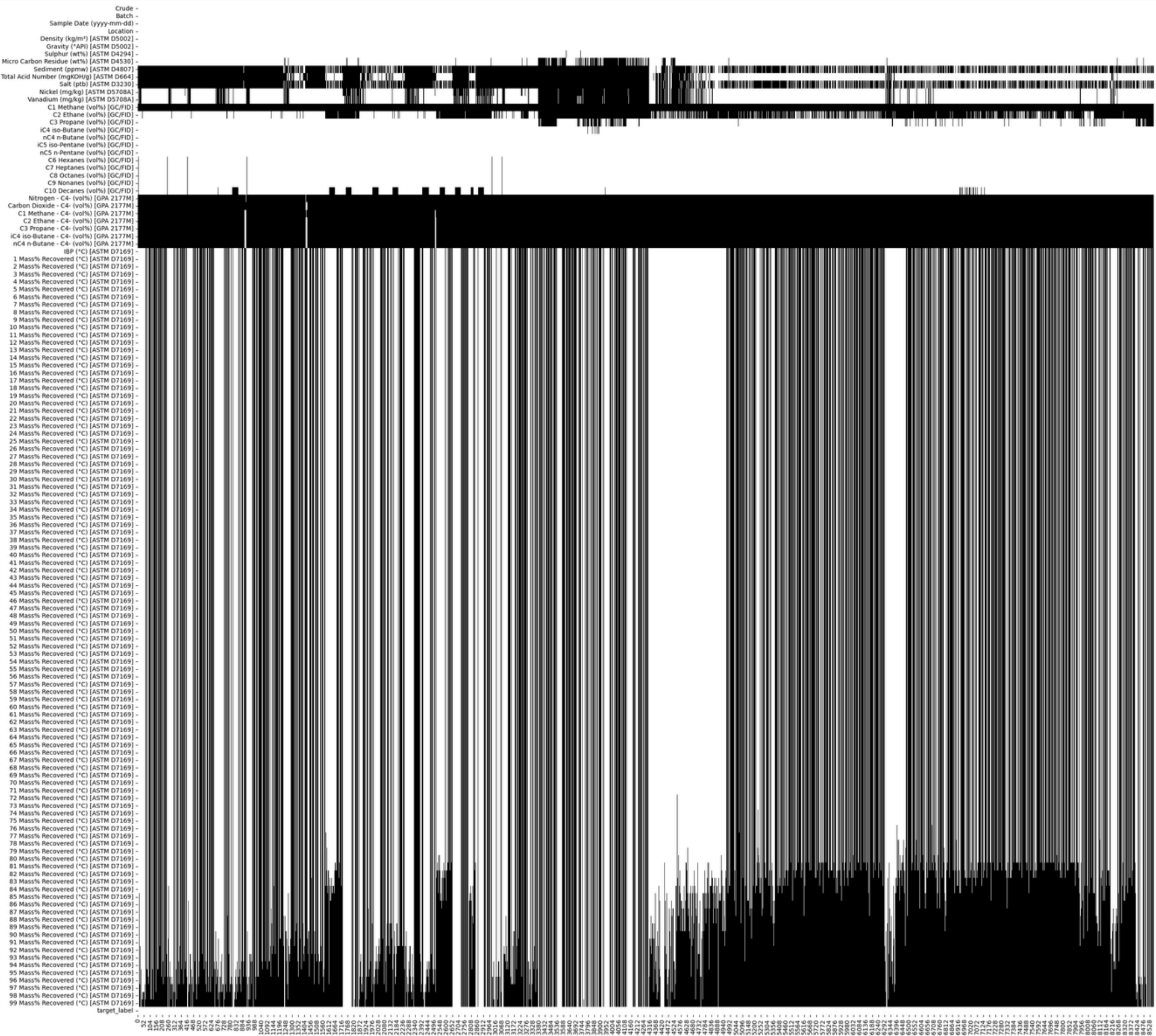The dataset contains over **8,500 rows and 133 columns**

### Data Types

The datatypes are 1 datetime, 4 categorical, and 128 numeric columns.
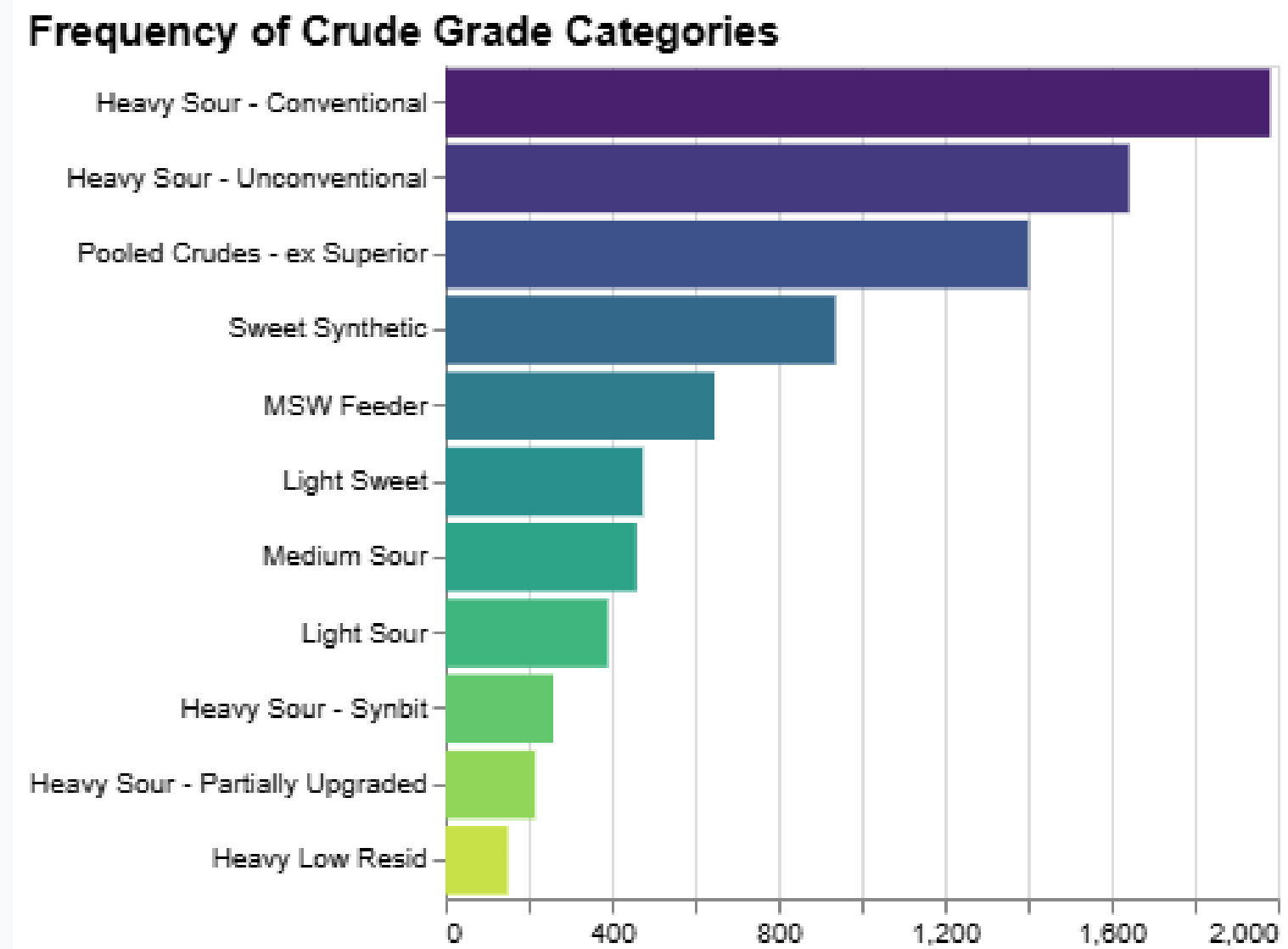
# Missing Values

## Heatmap

The heatmap color reveals **large swaths of missing data** throughout the dataset
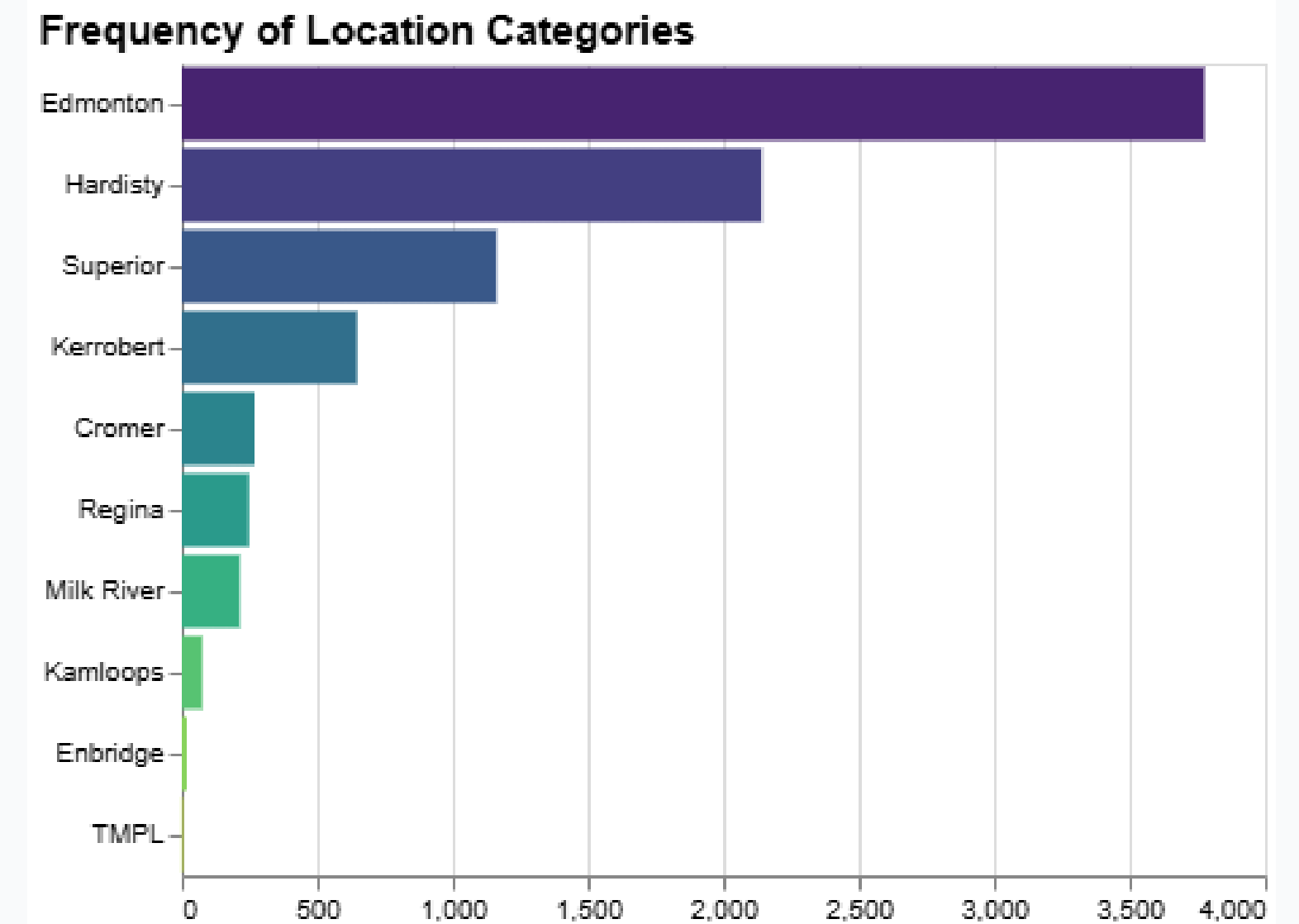
# Categorical Frequencies

## Crude Grades

Majority of the dataset's rows are related to the production of **heavy oils and their by–products**



Frequency of Crude Grade Categories

## Sample Locations

Oil production typically **flows from northern Alberta**, through Edmonton, and onto Hardisty
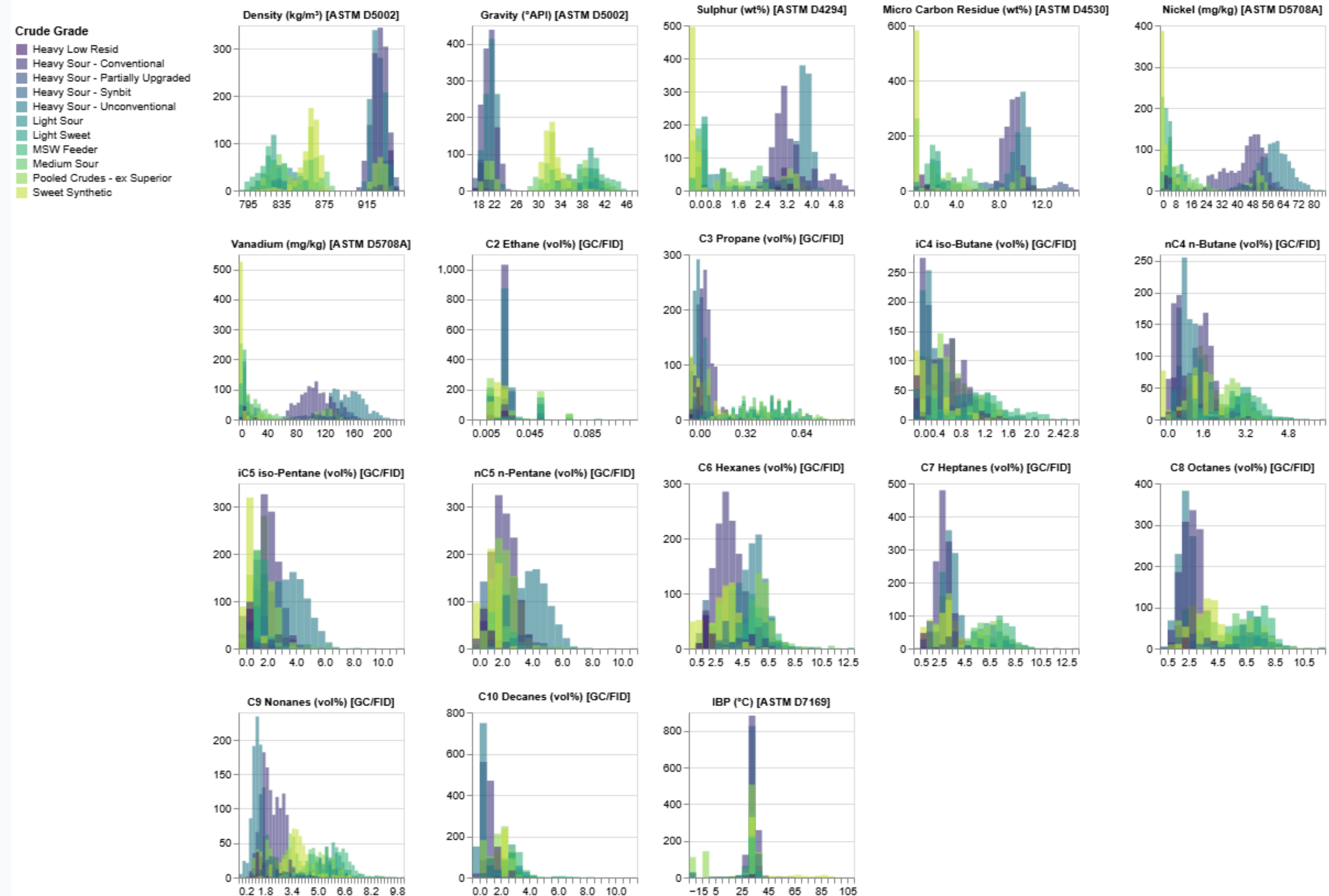


Frequency of Location Categories

# Numeric Distributions

## Crude Quality Columns

Noticeable patterns in the distributions suggest that **many groups of crudes share similar properties**
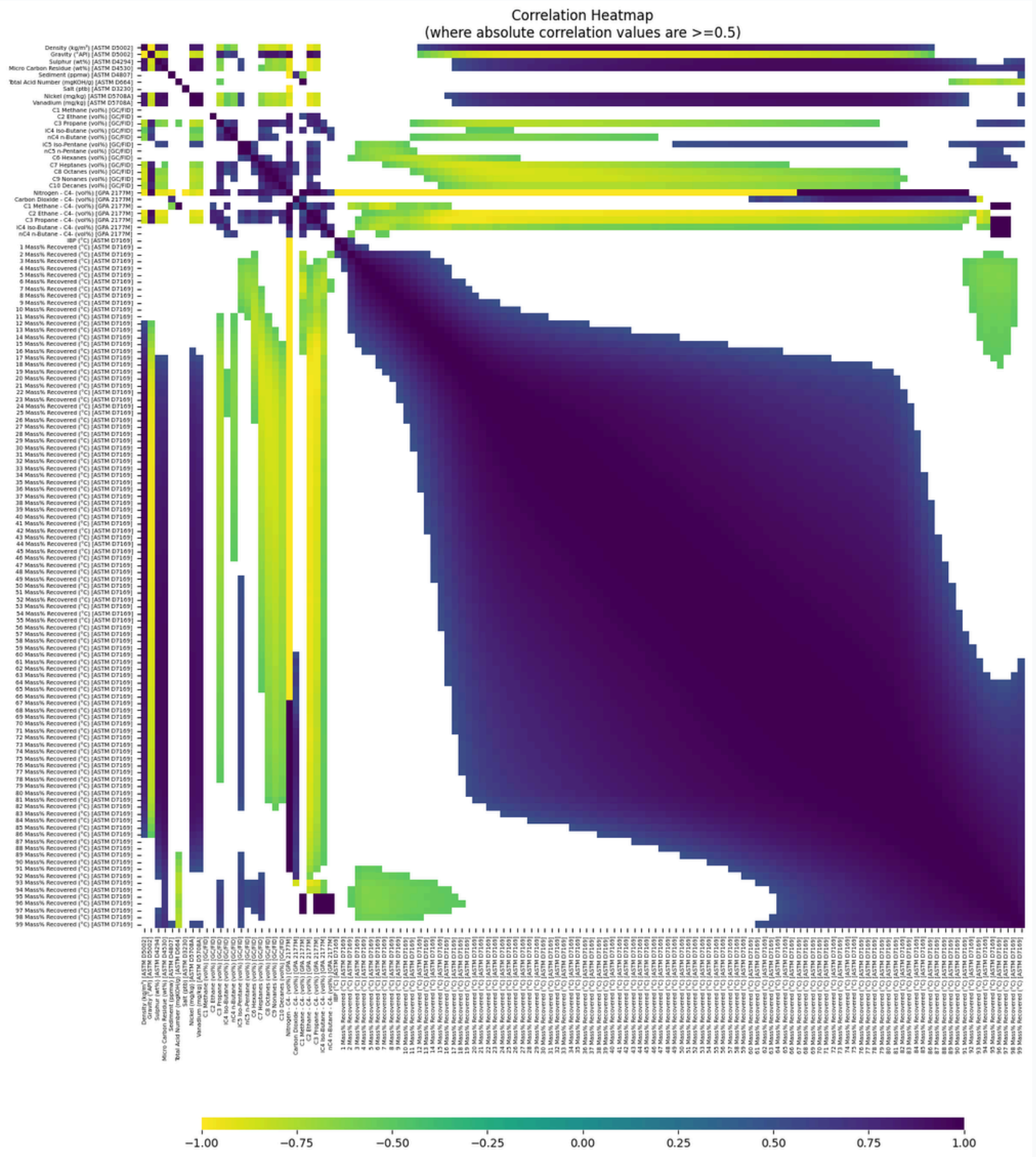


Crude Quality Columns Data Distributions

# Numeric Correlations

## Correlation Heatmap

**Strong correlations** are found throughout the crude quality and distillation temperature columns



Correlation Heatmap
(where absolute correlation values are >=0.5)

# Clustering Algorithms

## K–Means

K–Means is a **baseline clustering** method that partitions data into K distinct clusters by minimizing intra–cluster variance, allowing for efficient and straightforward classification of crude oil samples.

## Agglomerative

Agglomerative clustering is a **hierarchical approach** that builds clusters by merging similar data points iteratively, resulting in a tree–like structure that captures the relationships among crude oils.

## GMM

Gaussian Mixture Models (GMM) utilize a **probabilistic framework** to model data as a mixture of multiple Gaussian distributions, providing flexibility and capturing complex data distributions in crude oil classification.

# Model Evaluation

### Internal Metrics

Internal metrics of **Silhouette Score, DBI, and CHI** are used to assess clustering quality using statistical measures

### External Metrics

External metrics of **ARI, NMI, and FMI** are used to compare cluster assignments to known labels
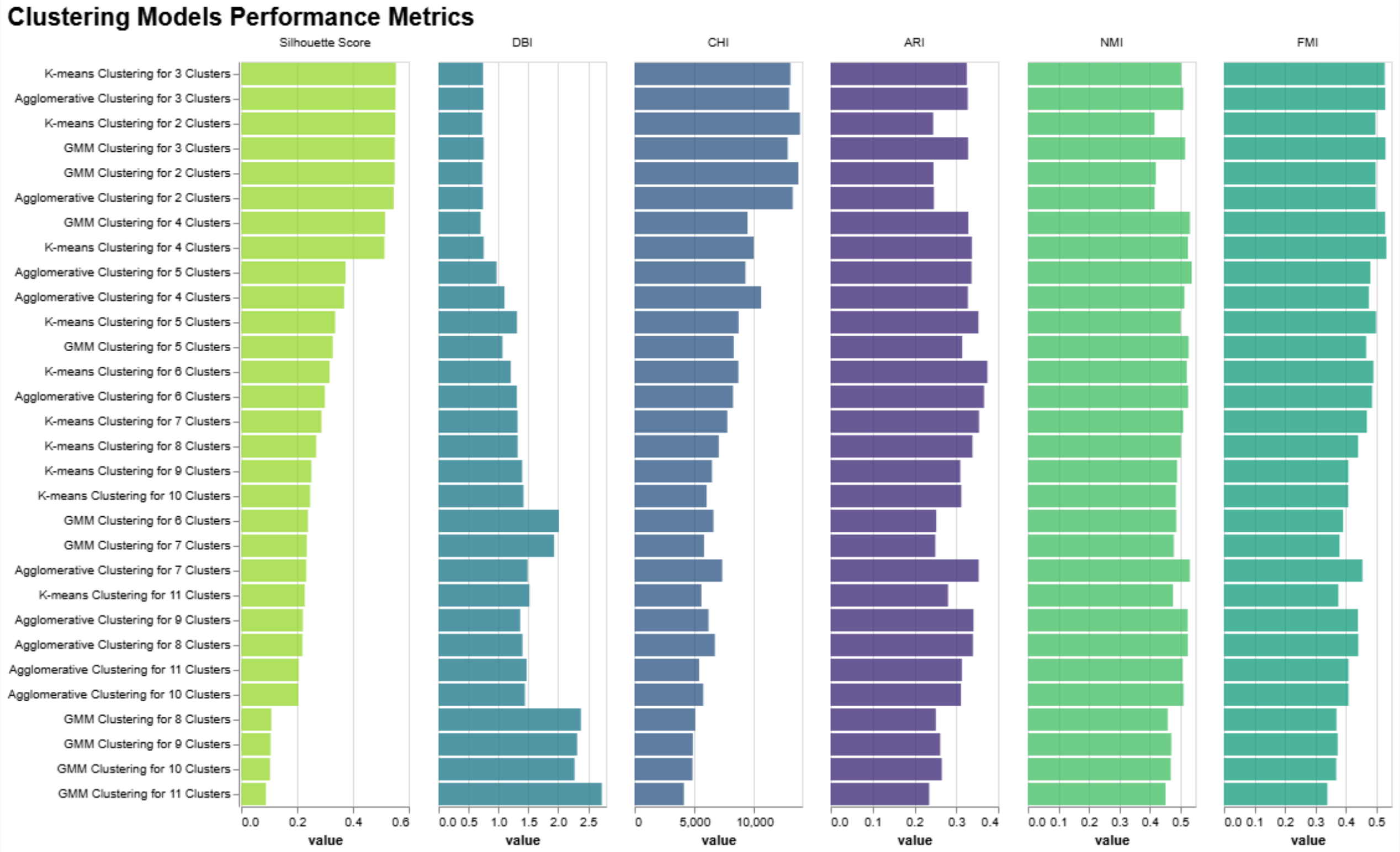
### Visualization Plots

Visualization plots are used to illustrate clustering results and distribution patterns

# Evaluation Metrics

## Best Performer

Models with **two to three clusters** showed the strongest performance, with Agglomerative Clustering and GMM performing best at three clusters
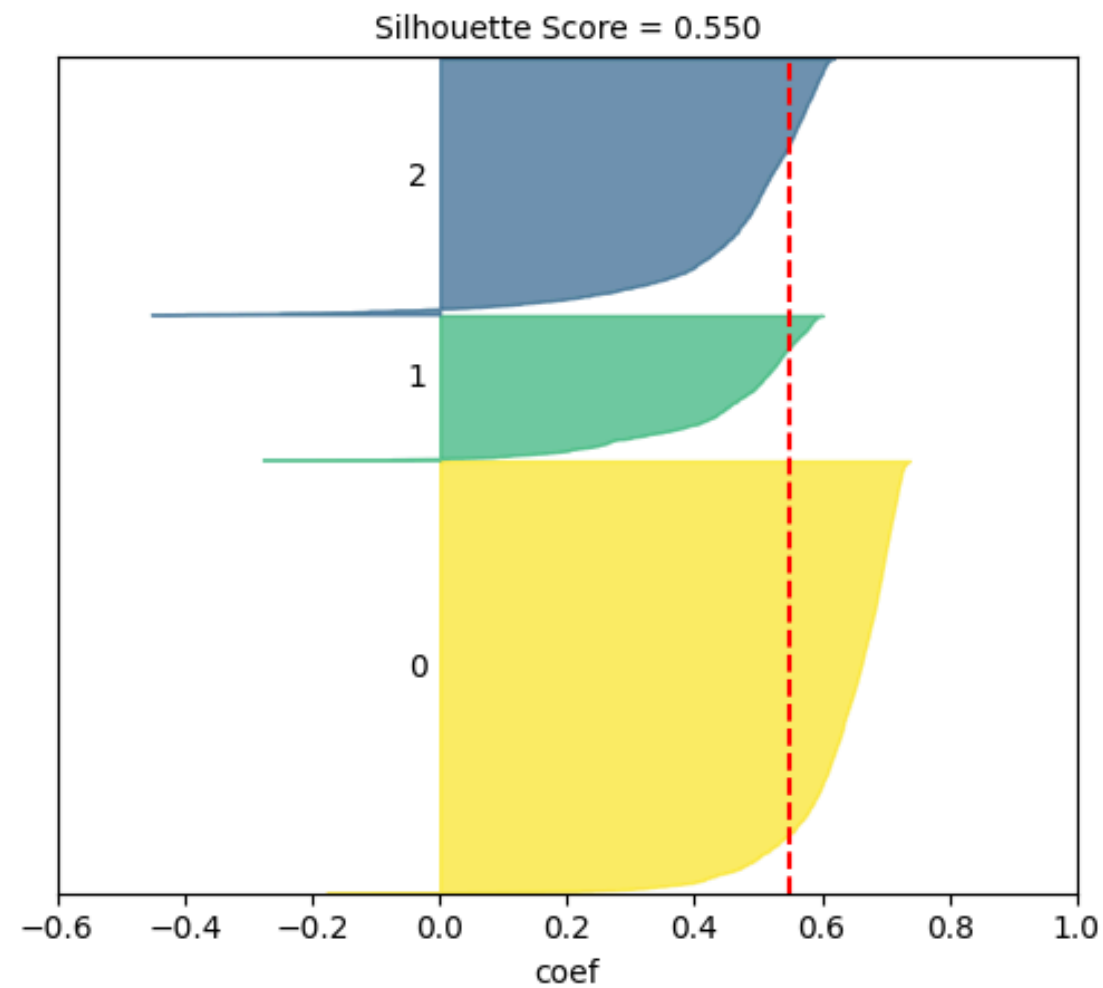


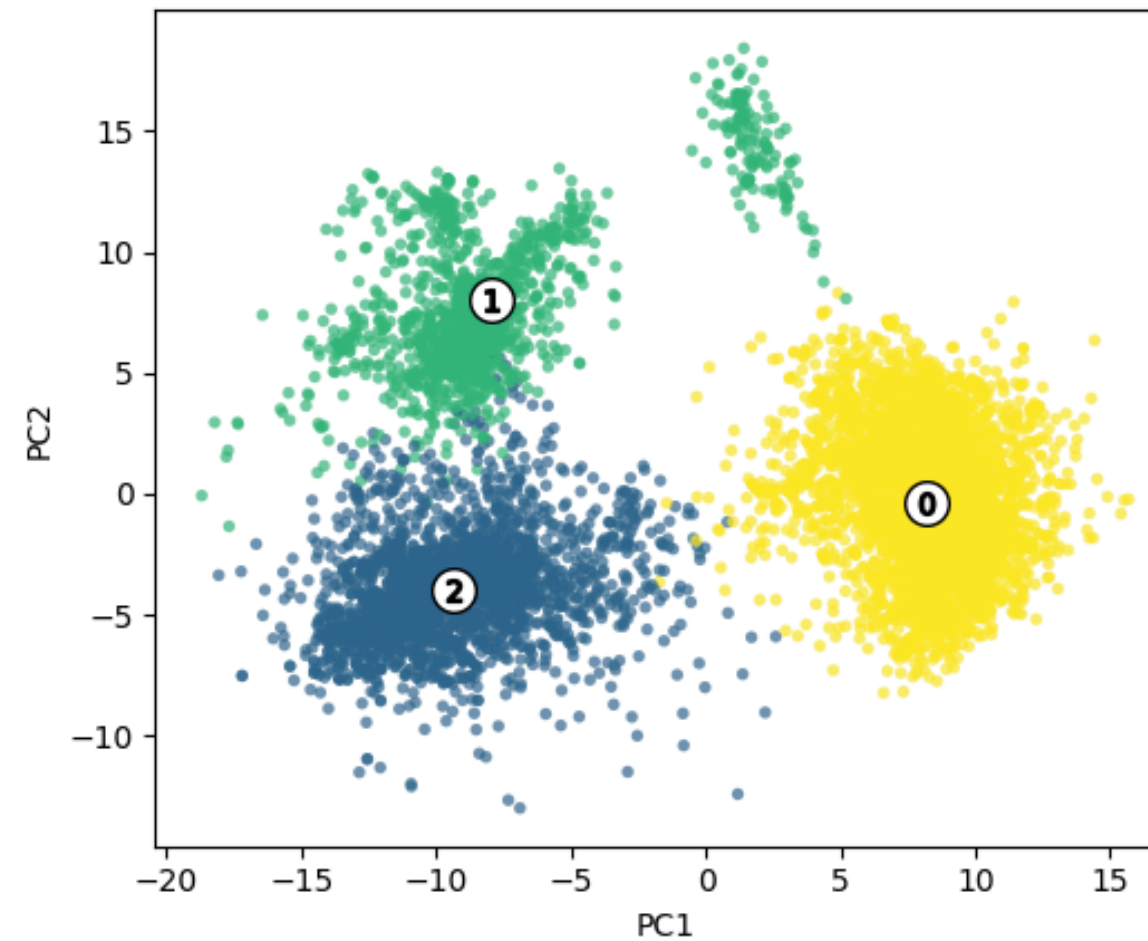Clustering Models Performance Metrics

# Evaluation Plots

## Dimensionally Reduced Clustering

The silhouette and scatter plots clearly show that **three clusters** provide the most distinct and well–separated groupings
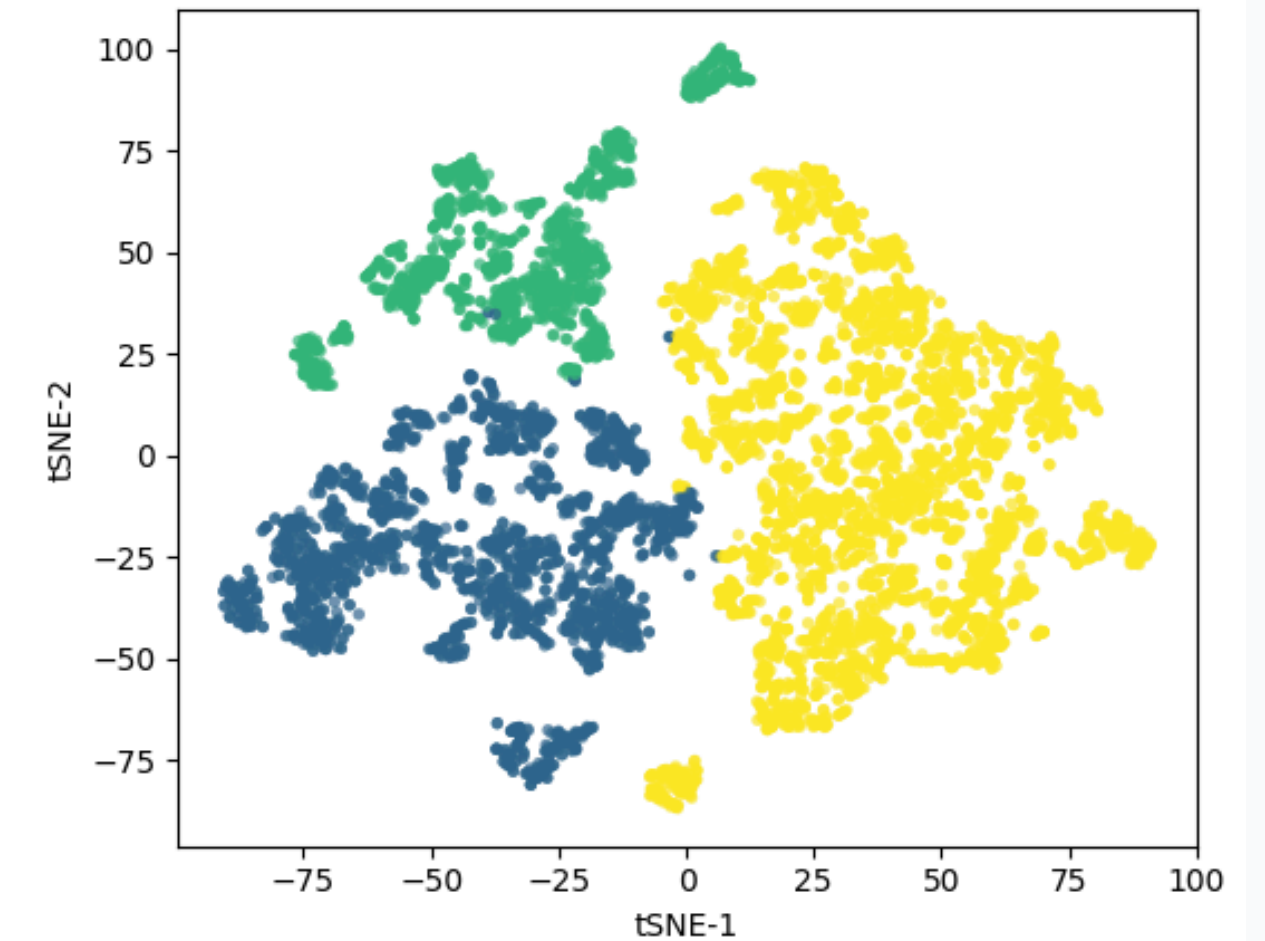


Agglomerative Clustering for 3 Clusters

# Unsupervised Confusion Matrix

Confusion Matrices of Actual vs Predicted Target Labels (for k=11 clusters)

# Supervised Benchmarking

## Comparing key machine learning methods

### Logistic Regression

Logistic Regression serves as a fundamental method, modeling the **probability of class membership** using a linear approach, effective for binary classification problems in crude oil analysis.

### Support Vector Classifier

SVC utilizes hyperplane separation to classify data points, **excelling in high–dimensional spaces**, providing robust classification through kernel trick adaptations tailored for complex crude oil datasets.

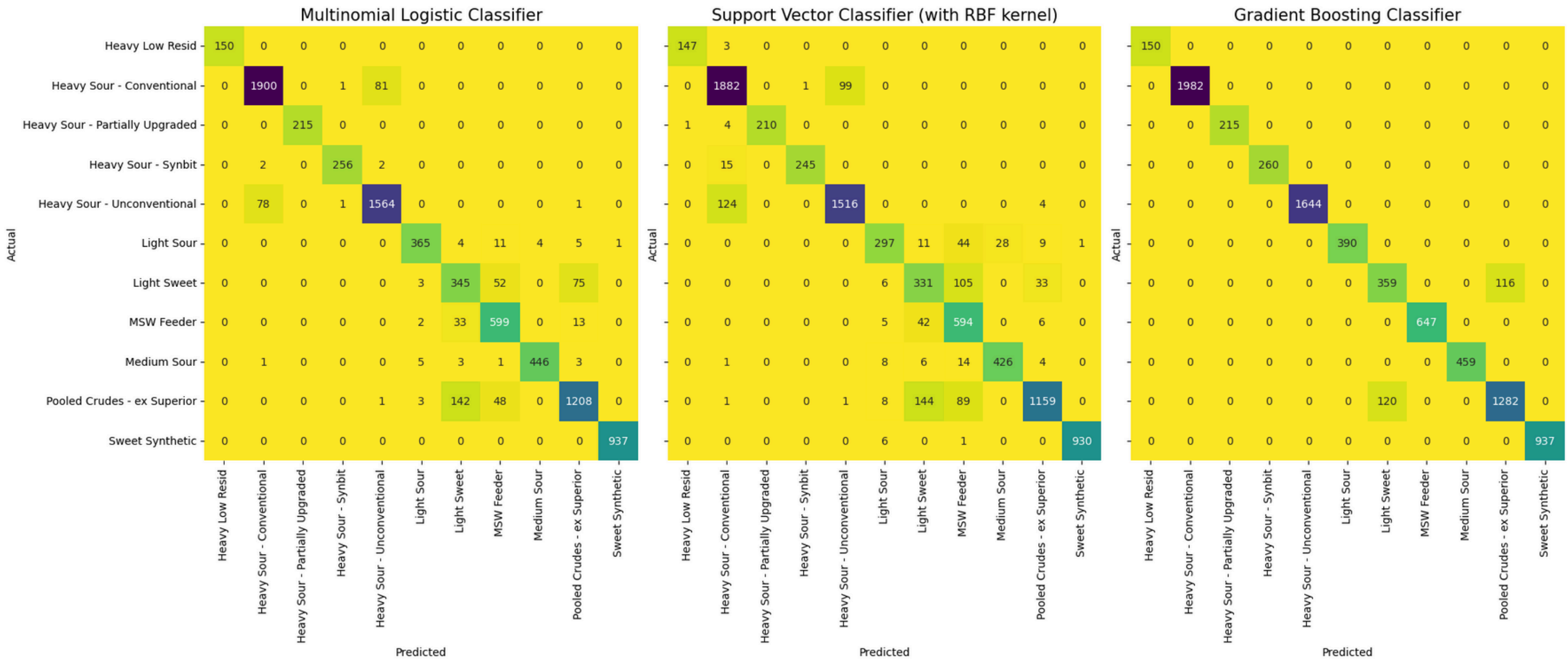### Gradient Boosting Classifier

Gradient Boosting Classifier builds an ensemble of decision trees where each model corrects the errors of the previous one, producing an **accurate and flexible classifier** that captures nonlinear relationships in complex crude oil data.

# Supervised Confusion Matrix

**Better Performance then Unsupervised**

Persistent misclassification further suggest **common groupings**



Confusion Matrices for Supervised Classification Models

# Project Timeline

**WEEK 1**

**Data Extraction**

**Data Cleaning**

**EDA**

**WEEK 2**

**Modeling**

**Evaluation**

**Final Presentation**

COMPLETED

# Conclusion

### Key Findings

Unsupervised models consistently identified **two to three natural clusters**, revealing a simpler structure than the current eleven–grade classification

### Model Agreement

K–Means, Agglomerative, and GMM **all produced similar groupings**, confirmed by supervised models showing the same overlapping crude categories

### Industry Insight

Results suggest Canadian crudes exist on a **continuous quality spectrum**, better represented by broader heavy, medium, and light/synthetic families

### Future Direction

Expanding the analysis with **economic, carbon, and yield data** could strengthen classification accuracy and industry applicability