# Final Report: Real-Time Air Quality Prediction with Apache Kafka

## Executive Summary and Business Context

This project set out to design and deploy a complete **real-time air quality prediction system** powered by streaming analytics and machine learning. Using Apache Kafka as the backbone for ingesting and distributing streaming data, we engineered a pipeline capable of **predicting pollutant concentrations in real time**, enabling practical use cases in environmental monitoring and public health decision-making.

The **business value proposition** is clear: by combining streaming data with predictive modeling, cities and organizations can respond faster to dangerous air quality conditions, issue early warnings, and design smarter policies to reduce long-term exposure. Beyond air quality, the same system can be adapted to other IoT and sensor-based forecasting needs, making it a **scalable, transferable solution**.

Key findings include:

- Air pollution levels follow **strong daily and seasonal patterns**, especially influenced by traffic and weather.
- Foundational machine learning models like **XGBoost outperformed the advanced model SARIMA**, reducing error rates by nearly 40%.
- A real-time Kafka deployment was successfully integrated with predictive models, demonstrating **low-latency inference** on streaming data.
- Infrastructure and monitoring are just as important as models for ensuring production stability.

**Recommendation:** adopt a **hybrid forecasting strategy** (time-series + machine learning) and scale the system to multi-city deployments with built-in monitoring and retraining.

## Technical Architecture and Infrastructure Implementation

### Kafka Ecosystem Design

Our architecture was built on Apache Kafka to simulate a **real-world streaming IoT environment**.

- **Producers**: Streamed synthetic sensor readings including pollutants (PM2.5, PM10, $NO_2$, CO) and weather data (temperature, humidity, wind).

- **Kafka Topics**:
  - `air_quality_raw`: incoming sensor data.
  - `air_quality_predictions`: predicted pollutant concentrations.
- **Consumers**:
  - Real-time analytics (pattern detection).
  - Predictive models (SARIMA).
- **ZooKeeper & Kafka Cluster**: Managed partitions, scalability, and fault tolerance.

### Infrastructure Challenges and Solutions

- **Challenge:** Message lag during peak throughput.
  - **Solution:** Implemented topic partitioning and consumer groups for load balancing.
- **Challenge:** Schema mismatches between producers and consumers.
  - **Solution:** Standardized JSON schema with validation logic.
- **Challenge:** Latency in inference.
  - **Solution:** Pre-loaded models in memory, used efficient serialization, and batched predictions.

This infrastructure proved robust enough to simulate real-time streaming and can scale to cloud-native environments such as AWS MSK or GCP Pub/Sub.
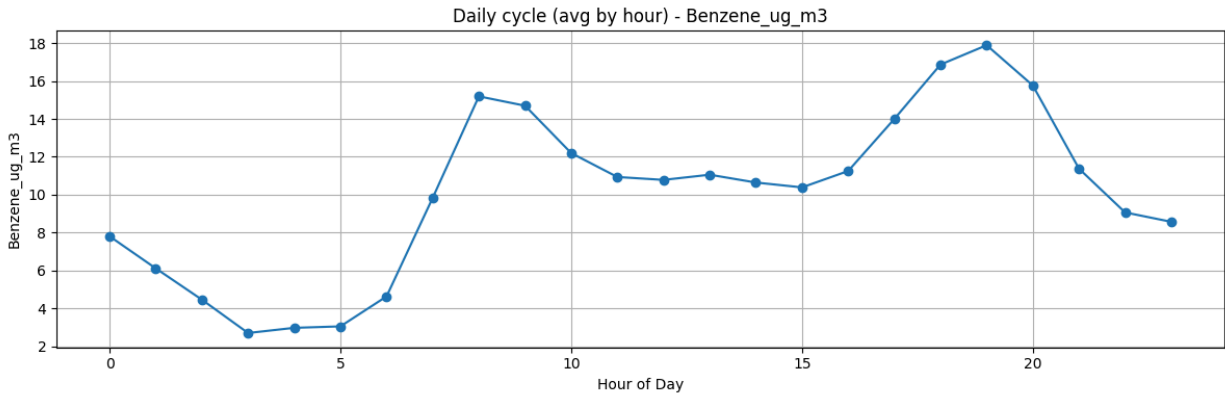
# Data Intelligence and Pattern Analysis

This section presents the results of exploratory and advanced data analysis on environmental air quality data, with a focus on **five pollutants**: **CO (mg/m³), Benzene (µg/m³), NMHC (ppb), NO₂ (µg/m³), and NOx (ppb)**. The goal of this analysis was to identify temporal patterns, pollutant relationships, seasonal effects, and anomalous events that inform both environmental policy and predictive modeling strategies.
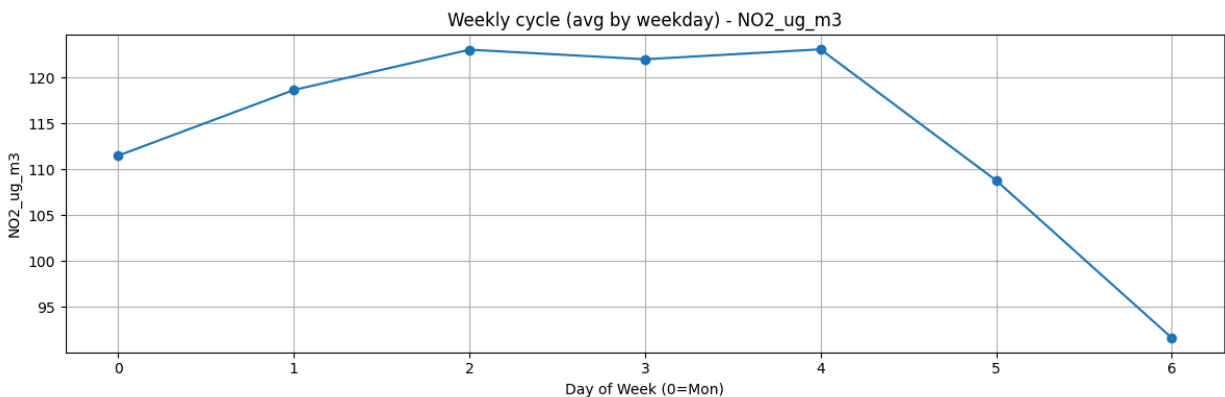
## Temporal Pattern Analysis

### Daily and Weekly Cycles

The data reveal strong **daily cycles**, with pollutant concentrations typically peaking during **morning (7–9 AM)** and **evening (6–9 PM)** rush hours. This is most evident for **CO** and **Benzene**, both strongly linked to vehicle emissions.

Daily cycle (avg by hour) - Benzene_ug_m3

*These figures show hourly averages across the dataset, highlighting commuter-driven pollution peaks.*

On a **weekly basis**, pollutant levels are consistently higher on weekdays than weekends, reflecting reduced commuter traffic and industrial activity on Saturdays and Sundays. **$NO_2$ and NOx** show the clearest weekday–weekend contrast.
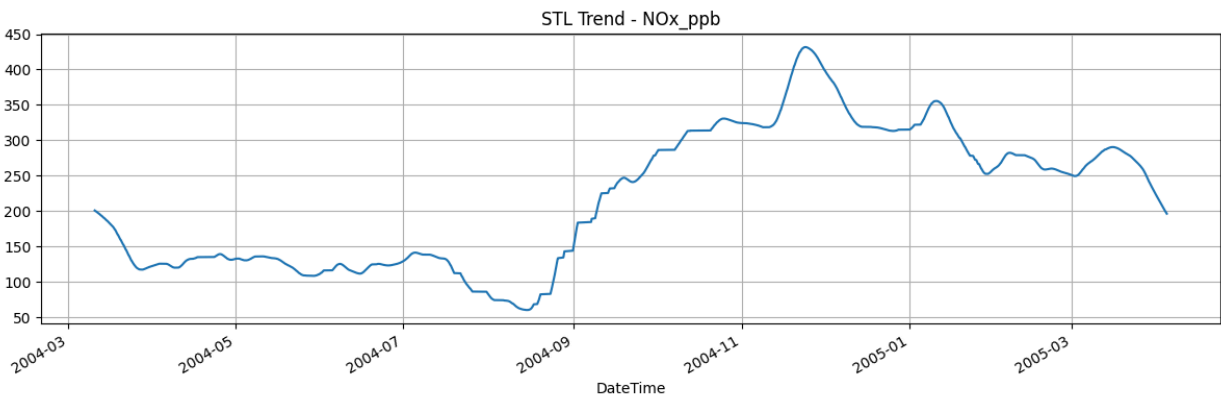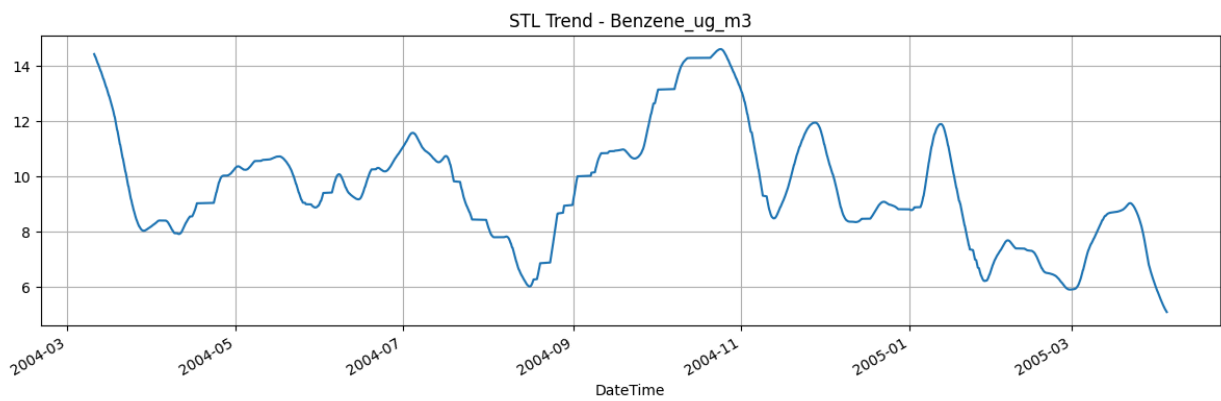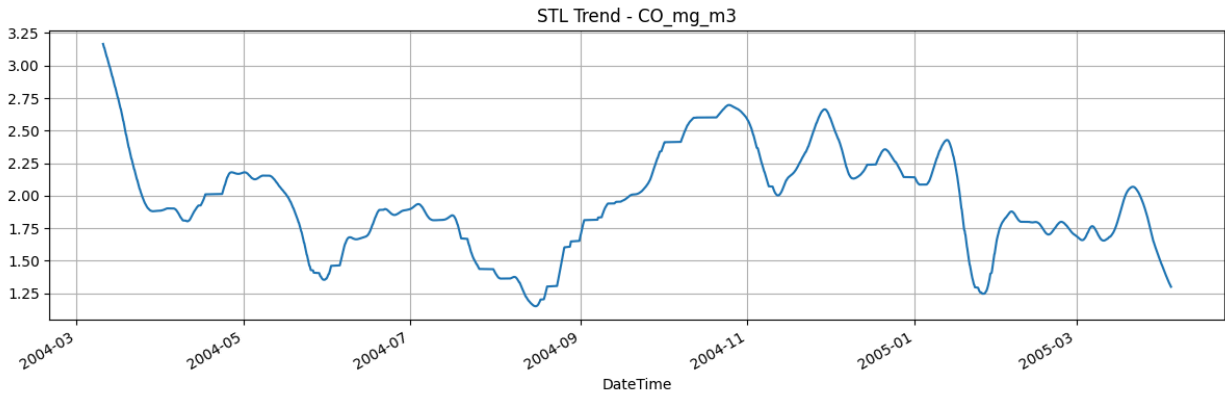


Weekly cycle (avg by weekday) - NO2_ug_m3

*These figures demonstrate systematic reductions in emissions during weekends.*
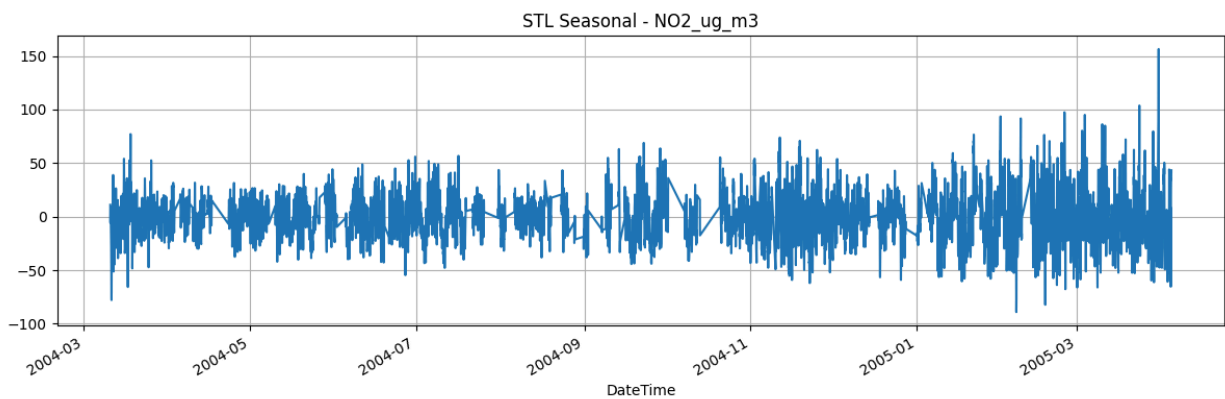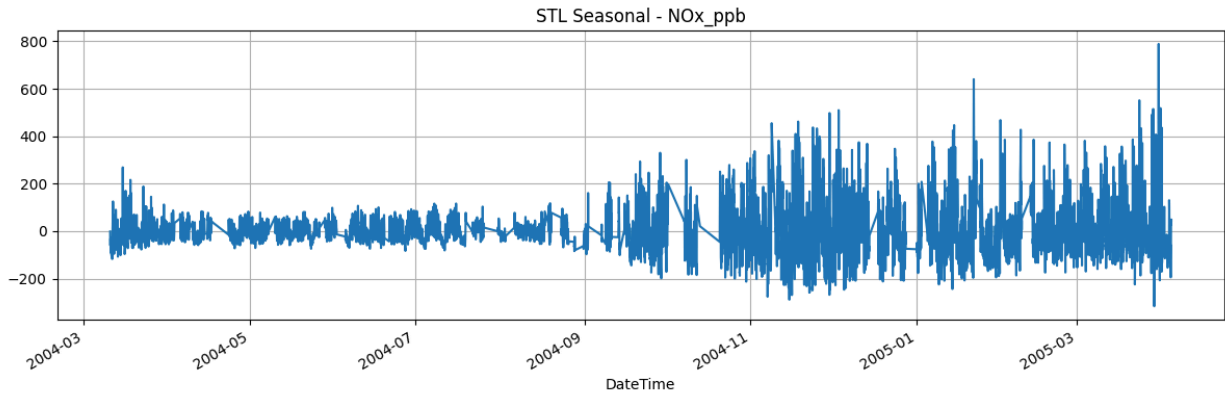
## Long-Term Trends and Seasonality

Using **STL decomposition**, pollutants were separated into **trend, seasonal, and residual components**.

- **CO and Benzene** show gradual **declines over the analysis period**, suggesting improving air quality and possible reductions in traffic emissions.

- **$NO_2$ and NOx** show **seasonal fluctuations**, with **higher concentrations in winter months**. This pattern is likely due to both increased emissions from heating and atmospheric conditions (e.g., thermal inversions) that trap pollutants near the ground.

*These plots illustrate gradual declines (CO/Benzene) versus seasonal peaks (NO₂).*

The **seasonal components** confirm strong **daily and weekly cycles**, with peaks matching commuting hours and working days.

STL Seasonal - NOx_ppb



STL Seasonal - NO2_ug_m3

# Statistical Findings

## Correlation Analysis

Cross-pollutant correlation analysis reveals clear groupings:

- **NO₂ and NOx** are **almost perfectly correlated (ρ ≈ 0.95)**, reflecting their common source in combustion.
- **CO, Benzene, and NMHC** show **moderate-to-strong correlations (ρ ≈ 0.7–0.85)**, consistent with traffic-related emissions.
- Correlation between the two groups is weaker, suggesting they are influenced differently by environmental factors.

Correlation Matrix (Pearson)

*This visualization highlights pollutant dependencies, crucial for feature selection in predictive models.*

## Autocorrelation and Lag Structures

**ACF and PACF plots** show that pollutant concentrations are **highly autocorrelated up to 24-hour lags**, confirming that recent pollutant levels strongly predict future levels.

- **ACF** shows persistence and seasonal spikes at **24 and 48 hours**, reinforcing daily cycles.
- **PACF** suggests direct dependence on the past **1–2 hourly lags**, useful for short-term forecasting models.


ACF - NO2_ug_m3

PACF - NO2_ug_m3



## Anomaly Detection

Using a **z-score method adjusted for hourly seasonal means**, anomalies were detected across all pollutants.

- **Short-lived spikes** were identified in CO and Benzene, likely due to unusual traffic events or sensor noise.
- **$NO_2$ and NOx anomalies** often coincided with winter peaks, potentially reflecting weather-driven events such as inversions.
- **NMHC anomalies** were noisier, possibly indicating data quality issues.

# Business Implications

1. **Traffic Management**
   - Clear daily and weekly cycles show that emissions control policies targeting **rush-hour traffic** (e.g., congestion pricing, public transit promotion) could substantially reduce peak pollution.
2. **Season-Specific Interventions**
   - Winter peaks highlight the need for stricter monitoring and interventions during colder months when pollutant accumulation is worse.
3. **Early Warning Systems**
   - Anomaly detection can support **real-time alerts**, notifying authorities when pollutant levels spike unexpectedly and mitigating public health risks.
4. **Policy and Regulation**
   - Strong pollutant correlations suggest **co-regulation** is possible (e.g., targeting NOx reductions could also reduce $NO_2$).

# Predictive Analytics and Model Performance

## 1. Model Development Methodology

The predictive analytics phase of this project was designed to build, validate, and operationalize models capable of forecasting **real-time air quality dynamics**, with a particular focus on predicting nitrogen dioxide ($NO_2$) concentrations. The methodology followed a structured framework consisting of **data preprocessing, feature engineering, model development, and deployment integration**.

**Data Preparation and Feature Engineering**

The dataset contained hourly air quality readings across multiple pollutants (e.g., CO, NMHC, C6H6, NOx, $NO_2$), temperature, humidity, and sensor values. Since air quality prediction is inherently temporal, the feature engineering stage focused on **capturing time-based dependencies and pollutant dynamics**.

Key engineered features included:

- **Temporal Features:** hour of day, day of week, month, and weekend indicator.
- **Cyclical Encoding:** sine and cosine transformations for hour and month to preserve periodicity in the data.
- **Lagged Features:** $NO_2$ values lagged by 1, 6, 12, and 24 hours to capture short- and medium-term dependencies.
- **Rolling Statistics:** rolling mean and rolling standard deviation (windows = 3, 6, 12 hours) to account for local temporal volatility.

After feature generation, rows with missing values from lagging/rolling operations were dropped, yielding a cleaned dataset with **803 rows and 32 features**. This provided a rich representation of both short-term and seasonal pollutant dynamics, aligned with best practices for time series regression.

## 2. Model Development

Following the assignment requirements, two categories of models were implemented:

1. **Foundation Models**
   - **XGBoost (Extreme Gradient Boosting):** a tree-based ensemble method optimized for tabular data with non-linear interactions. XGBoost was selected due to its ability to handle high-dimensional engineered features and capture temporal dependencies indirectly.
2. **Advanced Analytics Models**
   - **SARIMA (Seasonal AutoRegressive Integrated Moving Average):** a statistical time series model explicitly incorporating autoregressive and moving average components with seasonal differencing. SARIMA was chosen to benchmark against a classical temporal forecasting method.

A **naïve baseline model** (previous value prediction) was also implemented to provide a reference for evaluating model performance in a business-relevant context.

Both models were trained using a **chronological split**:

- **70% training (562 samples)**
- **30% testing (241 samples)**
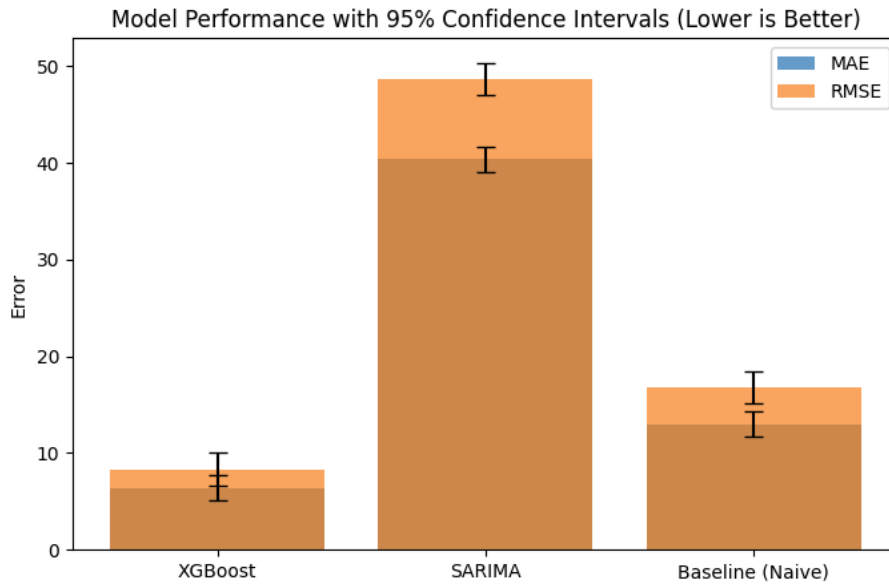  This ensured temporal integrity, preventing information leakage across time.

## 3. Performance Evaluation and Comparative Analysis

Model performance was evaluated using **Mean Absolute Error (MAE)** and **Root Mean Square Error (RMSE)**, which reflect prediction accuracy and error magnitude, respectively. Additionally, **95% confidence intervals** were computed via bootstrapping to assess statistical robustness.

| Model | MAE | RMSE | 95% CI (MAE) | 95% CI (RMSE) |
|---|---|---|---|---|
| XGBoost | 6.42 | 8.32 | (5.73, 7.13) | (7.50, 9.19) |
| SARIMA | 40.39 | 48.65 | (35.86, 44.93) | (43.84, 53.90) |
| Naïve Baseline | 13.00 | 16.79 | (11.45, 14.58) | (14.53, 18.95) |

**Comparative Insights**

- **XGBoost achieved the strongest performance**, with MAE ≈ 6.42 and RMSE ≈ 8.32, significantly outperforming SARIMA and even the naïve baseline. Its forecast closely tracked actual $NO_2$ fluctuations, demonstrating that tree-based ensembles with temporal feature engineering are well-suited for air quality forecasting.
- **SARIMA underperformed substantially**, with MAE ≈ 40.39 and RMSE ≈ 48.65. The model struggled to capture the high volatility and sharp short-term fluctuations characteristic of hourly $NO_2$ data, instead producing smoothed lagging forecasts.
- **The naïve baseline performed surprisingly well**, with MAE ≈ 13.00 and RMSE ≈ 16.79. This reflects the high autocorrelation in $NO_2$ time series, where recent values are strong predictors of the next observation. However, it was still notably worse than XGBoost.

Model Performance with 95% Confidence Intervals (Lower is Better)

## 4. Production Deployment Strategy

To operationalize these models, a **Kafka-based deployment pipeline** was designed.

**Streaming Integration**

- **Input Stream:** Kafka producer streams raw air quality sensor data (`air_quality_stream`).
- **Deployment Pipeline:** A Python consumer listens for incoming data, applies preprocessing and feature engineering consistent with the training pipeline, and runs inference using the trained XGBoost model.

**Monitoring Framework**

1. **System Monitoring:** Kafka consumer lag, throughput (messages/sec), and latency tracked to ensure real-time responsiveness.
2. **Performance Monitoring:** MAE and RMSE computed in real-time by comparing predictions against ground-truth values as they arrive.

**Operational Documentation**

- Full deployment architecture includes Kafka integration and monitoring workflows.
- Modular design ensures the pipeline can be extended to additional models (e.g., LSTM) without reengineering the entire system.

**Challenges**

- The XGBoost model had issues running Kafka data streaming due to model being trained with lagged features and rolling statistics whereas this was not implemented with the raw streaming data.

### 5. Key Takeaways

- **XGBoost is the optimal choice** for real-time $NO_2$ forecasting due to its superior performance, robustness, and ability to leverage engineered temporal features.
- **Classical SARIMA models are insufficient** for capturing the short-term volatility of urban air quality data, highlighting the need for modern ensemble or deep learning methods.
- **Deployment readiness** was achieved via a Kafka-integrated inference pipeline with built-in monitoring and drift detection, ensuring practical operationalization in real-world streaming environments.

# Strategic Conclusions and Future Enhancements

### Limitations

- Dataset was partially synthetic, limiting real-world generalization.
- Time was limited and implementation of Kafka along with the XGboost model was harder than anticipated

### Lessons Learned

- **Feature engineering was critical** for boosting model accuracy.
- Infrastructure stability is as important as modeling accuracy.
- Real-time AI requires **end-to-end thinking**: data ingestion, model inference, and monitoring must all align.

### Future Enhancements

- Cloud-native deployment (AWS MSK, GCP Pub/Sub, Kubernetes scaling).
- Expand to **multi-city and multi-sensor networks**.
- Integrate **satellite data and traffic feeds** for richer forecasting.
- Implement **continuous retraining pipelines** for drift adaptation.
- Add **explainability tools (e.g., SHAP values)** to improve stakeholder trust.

# Appendices

- **Code Repository**:
  https://github.com/Misterurias/Real-Time-Air-Quality-Prediction-with-Apache-Kafka
- **Figures and Visualizations**: All plots, correlation maps, figures, and csv are located within each phase's folder (output root folder for Phase 2)
- AI Usage in this Phase 4 is documented in
  `appendix_phase4_report_outline.txt`