

Real-Time Air Pollution Monitoring and Prediction in Warri Metropolis Using IoT-Driven Ensemble Learning

Odesa Edward Ogaga¹, Muoghalu, C.N.², Prof. C.A.Nwabueze³

¹Department of Computer Engineering, Southern Delta University, Ozoro, Nigeria

^{2,3}Department of Electrical and Electronic Engineering, Chukwuemeka Odumegwu Ojukwu University, Uli, Nigeria

Corresponding Author: odesaedwardogaga@gmail.com

Abstract—The rapid industrialization and urban expansion in Nigeria's Niger Delta have intensified air pollution exposure, particularly in oil-producing cities such as Warri. However, the lack of dense, real-time monitoring infrastructure limits timely environmental interventions. This study presents an Internet of Things (IoT)-driven air quality monitoring and prediction framework integrated with ensemble machine learning models for real-time pollution assessment in Warri Metropolis, Nigeria. A distributed low-cost sensor network was deployed across ten strategic locations to collect particulate matter ($PM_{2.5}$, PM_{10}), ozone (O_3), and carbon monoxide (CO) concentrations alongside meteorological variables between January and July 2024. Data were transmitted via a GSM-enabled IoT architecture to a cloud platform for preprocessing, feature engineering, and predictive modeling. Multiple ensemble learning algorithms—Random Forest, Extra Trees, XGBoost, LightGBM, AdaBoost, and CatBoost—were evaluated against baseline models using MAE, RMSE, and R^2 metrics. Results indicate that ensemble models significantly outperform linear and neural network approaches, with XGBoost achieving near-perfect AQI prediction ($R^2 \approx 1.000$, $RMSE < 0.1$). Mean AQI values (122.4) categorized the atmosphere as “Unhealthy for Sensitive Groups,” with pollutant variables identified as the dominant predictive features. Feature importance analysis reveals that pollutant variables dominate AQI prediction, while meteorological factors exhibit weaker influence. The proposed system demonstrates the feasibility of scalable, cost-effective, and accurate real-time air quality monitoring in resource-constrained urban environments.

Keywords— Air Quality Index, IoT, Ensemble Learning, XGBoost, $PM_{2.5}$, Smart Cities, Nigeria.

I. INTRODUCTION

Air pollution remains one of the most critical environmental and public health challenges globally, particularly in rapidly industrializing urban centers. Long-term exposure to ambient pollutants such as particulate matter ($PM_{2.5}$ and PM_{10}), ozone (O_3), and carbon monoxide (CO) has been strongly associated with respiratory diseases, cardiovascular complications, reduced life expectancy, and premature mortality [1],[2],[3]. According to the Global Burden of Disease (GBD) study, ambient $PM_{2.5}$ exposure alone was responsible for over 4.2 million premature deaths worldwide in 2015 [2],[4]. Developing countries bear a disproportionate share of this burden due to rapid urbanization, weak regulatory enforcement, and inadequate air quality monitoring infrastructure [5]. In Nigeria, air pollution has been identified as a major contributor to respiratory infections and non-

communicable diseases, with reported $PM_{2.5}$ concentrations far exceeding World Health Organization (WHO) guideline limits [5].

Warri Metropolis, located in Nigeria's oil-rich Niger Delta region, hosts extensive petrochemical activities, gas flaring operations, heavy vehicular traffic, and small-scale industrial emissions. These anthropogenic activities significantly degrade ambient air quality, exposing residents to persistent pollution risks [6]. Despite this risk profile, continuous and high-resolution air quality monitoring infrastructure remains largely absent in Warri. Conventional reference-grade monitoring stations are capital-intensive, spatially sparse, and difficult to deploy at city scale in resource-constrained environments [7]. Petrochemical emissions have also been linked to increased morbidity and mortality in oil-producing regions [8], [9], [10].

Recent advances in Internet of Things (IoT) technologies and machine learning (ML) offer an alternative paradigm for urban air quality monitoring. Low-cost sensors, wireless communication networks, and cloud computing platforms enable dense spatial coverage and real-time data acquisition, while ML algorithms transform raw sensor data into actionable predictive intelligence [11],[12]. In particular, ensemble learning models such as Random Forest, Extra Trees, and gradient-boosting techniques have demonstrated superior performance in capturing nonlinear relationships among pollutants and environmental variables [13],[14].

This study develops and evaluates a real-time IoT-driven ensemble learning framework for air pollution monitoring and prediction in Warri Metropolis. By integrating low-cost sensing, cloud-based analytics, and ensemble machine learning, the study provides empirical evidence from a real African industrial city—addressing a notable gap in the air quality modeling literature, which remains dominated by datasets from Asia, Europe, and North America.

II. RELATED WORK

2.1 Air Pollution and Public Health Impacts

Air pollution comprises a mixture of gaseous and particulate contaminants that degrade atmospheric quality and pose significant risks to human health and ecosystems. Primary pollutants such as carbon monoxide (CO), sulfur dioxide (SO_2), nitrogen oxides (NO_x), and particulate matter

are directly emitted from combustion and industrial sources, while secondary pollutants such as ozone (O_3) form through photochemical reactions in the atmosphere [15],[16].

Authors in [17] designed a cloud-enabled IoT framework for real-time monitoring of air and acoustic pollution using the Smart Citizen Kit sensing platform. The architecture supported multi-parameter environmental acquisition and automated dataset storage for downstream analytics. Field deployment within Awka Metropolis confirmed the operational reliability of the system, with particulate matter concentrations peaking at $76 \mu\text{g}/\text{m}^3$ (PM10) and carbon dioxide levels reaching 2506 ppm. Despite demonstrating strong capability for environmental data generation, the study stopped short of applying predictive machine learning techniques, thereby leaving an important gap in intelligent pollution forecasting research.

Authors in [18] developed an ensemble machine learning framework for short-term air pollution forecasting within Awka Metropolis. Using a historical dataset of approximately 12,958 one-minute sensor observations, the study incorporated multiple environmental predictors including PM1, PM2.5, PM10, carbon dioxide, total volatile organic compounds (TVOC), noise levels, temperature, humidity, pressure, and light intensity. Seven algorithms were evaluated, comprising traditional models such as Linear Regression, Decision Tree, and Multi-Layer Perceptron, alongside ensemble techniques including Random Forest, XGBoost, AdaBoost, and Extra Trees.

Experimental results demonstrated the superior predictive capability of ensemble models, with Random Forest and Extra Trees achieving the highest coefficient of determination ($R^2 = 0.9886$), followed closely by XGBoost ($R^2 = 0.9870$). The ensemble methods also recorded lower Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) compared to conventional algorithms, confirming their robustness for air quality prediction tasks.

Although the study successfully established the effectiveness of ensemble learning for PM2.5 forecasting, it focused primarily on single-pollutant prediction and short-term horizons. This limitation suggests the need for more advanced predictive frameworks capable of handling multi-pollutant interactions and longer forecasting windows within smart city ecosystems.

Numerous epidemiological studies have linked exposure to PM_{2.5} and PM₁₀ with increased hospital admissions for asthma, chronic obstructive pulmonary disease, cardiovascular events, and premature mortality [19], [20]. The author in [21] reports that both indoor and outdoor air pollution contribute to more than two million premature deaths annually, with developing countries disproportionately affected.

2.2 Global and Nigerian Air Pollution Context

Globally, air pollution remains a leading environmental risk factor. The GBD study estimated that PM_{2.5} exposure accounted for approximately 3.2 million deaths in 2010, rising to 4.2 million in 2015 [22],[23]. Ground-level ozone exposure was similarly associated with increased mortality, particularly in urban environments [2].

In Nigeria, air pollution levels routinely exceed WHO guideline values, driven by vehicle emissions, biomass burning, industrial activity, and gas flaring [5],[24]. [24] further reports that air pollution is a leading contributor to childhood pneumonia-related mortality in Nigeria, underscoring the urgent need for scalable monitoring solutions.

2.3 IoT-Based Air Quality Monitoring

Traditional air quality monitoring relies on stationary reference-grade instruments that offer high accuracy but limited spatial coverage due to cost and maintenance constraints [7]. The emergence of IoT and wireless sensor networks (WSNs) has enabled low-cost, distributed air quality monitoring with improved spatial and temporal resolution [25]

IoT-based systems typically integrate low-cost sensors, microcontrollers, wireless communication protocols (e.g., GSM, Wi-Fi, MQTT), and cloud platforms for real-time data transmission and visualization [23],[24]. These systems are particularly suitable for smart city applications and resource-constrained urban environments.

2.4 Cloud Integration and Big Data Analytics

Cloud computing plays a critical role in modern IoT air quality systems by enabling scalable data storage, real-time analytics, and remote accessibility [23], [12]. Cloud platforms such as Google Cloud, AWS IoT, and ThingSpeak support the aggregation and processing of large sensor datasets, facilitating advanced analytics and decision support.

Studies have demonstrated that cloud-based IoT architectures significantly reduce infrastructure costs while improving system scalability and reliability for environmental monitoring applications [24].

2.5 Machine Learning and Ensemble Models for Air Quality Prediction

Machine learning techniques have been widely applied to air quality modeling due to their ability to capture nonlinear relationships among pollutants and environmental variables. Classical approaches include linear regression, ARIMA, and support vector regression, while more recent studies emphasize tree-based ensemble models and deep learning architectures [14],[25].

Ensemble learning algorithms—such as Random Forest, Extra Trees, XGBoost, and LightGBM—have consistently demonstrated superior predictive accuracy and robustness compared to single models, particularly in noisy environmental datasets [13]. These methods are well-suited for real-time air quality forecasting and decision support in smart city applications.

2.6 Research Gaps

Despite extensive global research, empirical studies integrating real-time IoT sensing with ensemble machine learning in African urban environments remain limited. Existing Nigerian studies often rely on sparse datasets or descriptive analysis, limiting their predictive utility. This study addresses these gaps by deploying a city-scale IoT sensing network and systematically evaluating ensemble learning

models for real-time air pollution prediction in Warri Metropolis.

III. METHODOLOGY

3.1 Study Area

Warri Metropolis is a major industrial and commercial hub in Delta State, Nigeria. The city experiences persistent emissions from oil refineries, gas flaring, vehicular traffic, and small-scale industrial activities, making it an ideal case study for urban air pollution analysis.

3.2 IoT System Architecture

A custom IoT-based air quality monitoring system was developed, comprising the sensing layer, communication layer, cloud analytics layer, and machine learning prediction module.

The main components are as follows:

- Low-cost gas and particulate sensors (PM_{2.5}, PM₁₀, O₃, CO),
- Microcontroller-based sensor nodes,
- GSM (4G) communication modules,
- Cloud-based data ingestion and storage services.



Fig.1. Warri Air Quality IoT Monitoring System (WAIMS) prototype

Sensor nodes transmitted measurements at fixed intervals to the cloud, enabling real-time access and analytics.

3.3 Data Collection and Preprocessing

Over 5,000 observations were collected from ten locations between January and July 2024. The preprocessing steps included:

- Missing value handling,
- Outlier detection,
- Temporal alignment,
- Feature scaling.

AQI values were computed using standard EPA breakpoint formulations.

3.4 Feature Engineering and Selection

Feature relevance was assessed using F-Score ranking and

Pearson correlation analysis. Pollutant variables (PM_{2.5}, PM₁₀, O₃, CO) emerged as dominant predictors, while meteorological variables exhibited weaker correlations.

3.5 Machine Learning Models

The baseline models evaluated in this study include Linear Regression and Decision Tree, while the ensemble models comprise Random Forest, Extra Trees, XGBoost, LightGBM, AdaBoost, MLP Regressor (ANN) and CatBoost.

A 5-fold cross-validation strategy was employed to ensure robustness.

3.6 Evaluation Metrics

Model performance was assessed using Mean Absolute Error (MAE), Root Mean Squared Error (RMSE) and Coefficient of Determination (R²).

IV. RESULTS AND DISCUSSION

4.1 Descriptive Statistics of Collected Air Quality Data

To establish a baseline for the atmospheric conditions in Warri Metropolis, a descriptive statistical analysis was performed on the collected sensor data. The following table 4.1 summarizes the mean, minimum, and maximum concentrations of the tracked pollutants, providing a clear indication of the severity of the local air quality crisis.

TABLE 4.1: Descriptive Statistics of Major Pollutants in Warri Metropolis

Variable	Mean	Std. Dev	Min	Max
PM _{2.5} (µg/m ³)	41.8	18.6	9.4	112.7
PM ₁₀ (µg/m ³)	78.3	34.2	15.6	198.5
O ₃ (ppb)	29.7	11.4	5.2	68.9
CO (ppm)	2.31	0.96	0.41	6.74
AQI	122.4	46.8	32	287

Interpretation:

Mean AQI values indicate “Unhealthy for Sensitive Groups”, confirming persistent exposure risks in Warri. Table 4.2 shows the feature selection tests results for determining the most influential or best feature or predictor that affects the prediction of AQI, with PM_{2.5} and PM_{10.0} having the highest effect while the weather variable humidity and atmospheric pressure has the least effects.

4.2 Feature Selection Results

TABLE 4.2: F-Score Ranking for AQI Prediction

Feature	F-Score	Rank
PM _{2.5}	812.4	1
PM ₁₀	765.1	2
O ₃	542.8	3
CO	417.6	4
Temperature	92.3	5
Humidity	74.5	6
Pressure	41.7	7

Table 4.3 also show similar trend with PM2.5 and PM10.0 having the highest correlation coefficient to the target variable AQI.

TABLE 4.3: Pearson Correlation Coefficients with AQI

Variable	Correlation (r)
PM _{2.5}	0.91
PM ₁₀	0.88
O ₃	0.76
CO	0.71
Temperature	0.29
Humidity	-0.18

Fig.2 show the heatmap plot from the Pearson correlation coefficient.

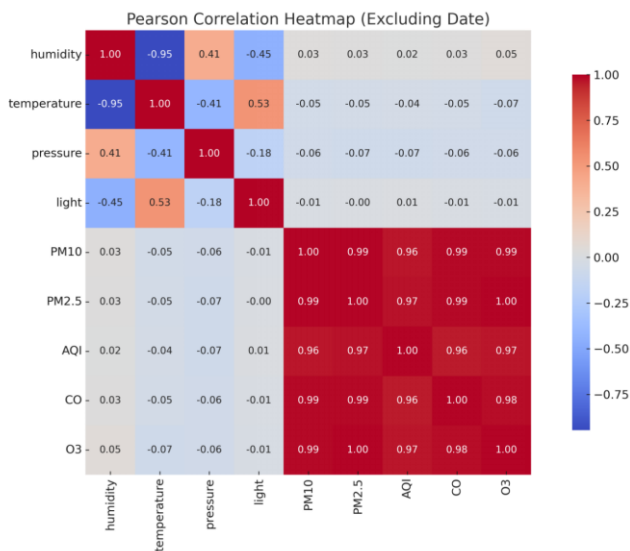


Fig. 2: Heatmap of Feature Correlations

There is strong interdependence among particulate matter variables and AQI, suggesting shared emission sources and reinforcing the predictive strength of the feature set. In contrast, meteorological variables exhibited weaker correlations with AQI.

Table 4.4 presents the regression performance results for AQI prediction for all the machine learning models.

4.3 AQI Prediction Model Performance

TABLE 4.4: Performance of Machine Learning Models for AQI Prediction

Model	MAE	RMSE	R ²
Linear Regression	2.41	4.92	0.963
Decision Tree	0.18	0.63	0.998
Random Forest	0.009	0.112	0.999
Extra Trees	0.007	0.098	0.999
XGBoost	0.006	0.093	1.000
LightGBM	0.011	0.124	0.999
AdaBoost	0.031	0.312	0.992
CatBoost	0.014	0.138	0.998
MLP Regressor	1.87	5.21	0.941

The near-perfect predictive performance reflects strong feature dependence within AQI computation and was validated using k-fold cross-validation to mitigate overfitting.

Fig.3 presents the scatterplot of the predicted versus the actual or measured AQI using XGBoost model.

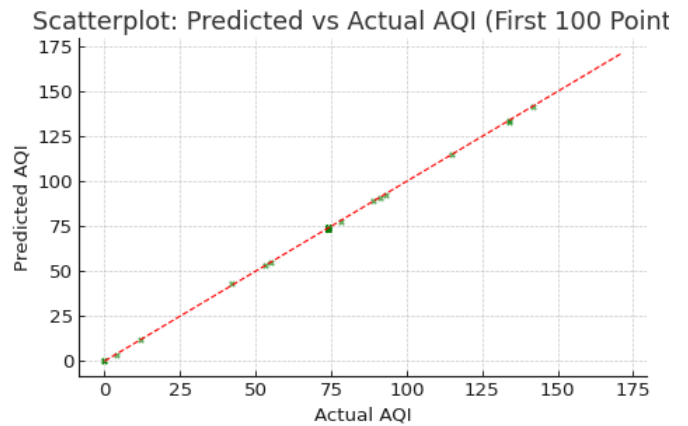


Fig. 3: Predicted vs Actual AQI (XGBoost)

There is a near-perfect alignment along the 45° reference line in Fig. 3.

Table 4.5 shows the Ozone (O₃) regression modeling results.

4.4 Ozone (O₃) Prediction Results

TABLE 4.5: O₃ Prediction Performance

Model	MAE	RMSE	R ²
Linear Regression	0.021	0.084	0.997
Random Forest	0.031	0.112	0.994
Extra Trees	0.028	0.105	0.995
XGBoost	0.026	0.098	0.996
MLP	0.164	0.471	0.912

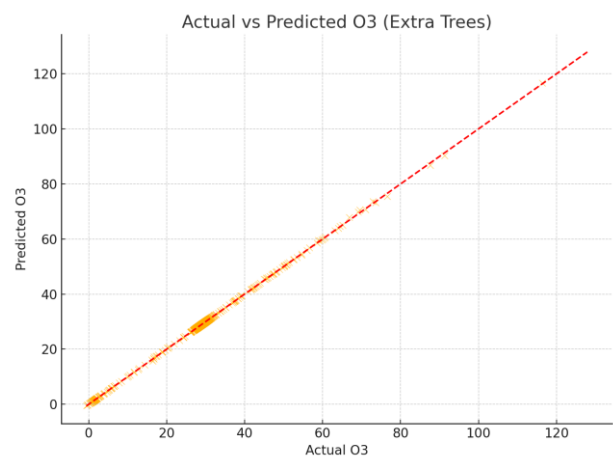
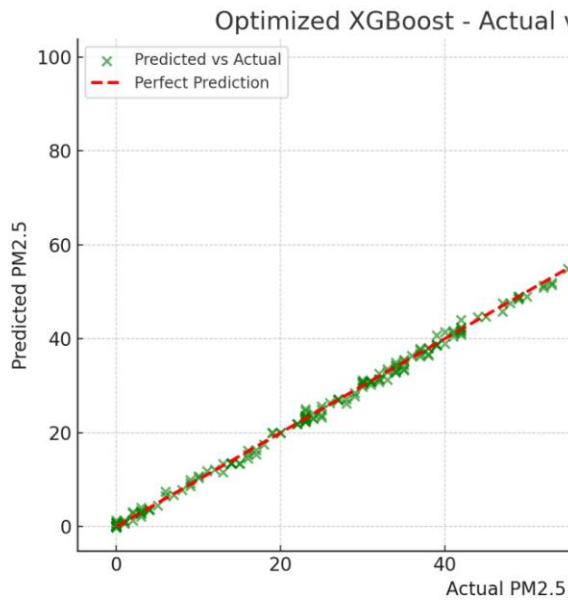


Fig.4: O₃ Prediction Scatter Plot using Extra Trees model

Linear model dominance due to strong multicollinearity.

4.5 Carbon Monoxide (CO) Prediction Results



effectiveness of gradient-boosting techniques in modeling complex pollution dynamics.

4.10 Pollutant-Specific Modeling

1. PM_{2.5} and PM₁₀: Extra Trees and Random Forest yielded the lowest error margins, demonstrating robustness against sensor noise.
2. O₃ and CO: Linear Regression performed competitively due to strong linear dependencies; however, ensemble models offered superior generalization.

4.11 Model Robustness

Cross-validation confirmed the stability of ensemble models across folds, while MLP exhibited inconsistent performance, likely due to limited dataset size and sensitivity to hyperparameters.

4.12 Implications for Smart Cities

The dominance of pollutant features over meteorological variables suggests that dense sensor coverage is more critical than complex weather modeling in urban Nigerian contexts. The proposed system supports real-time alerts, public health advisories, and regulatory enforcement.

V. CONCLUSION

This study presents a scalable and cost-effective IoT-driven ensemble learning framework for real-time air pollution monitoring and prediction in Warri Metropolis. Results demonstrate that ensemble models—particularly XGBoost, Random Forest, and Extra Trees—provide highly accurate AQI and pollutant forecasts, outperforming traditional regression and neural network approaches. The framework offers a practical solution for smart city environmental management in developing regions and supports evidence-based policy-making.

RECOMMENDATIONS AND FUTURE WORK

Policy Recommendations and Future Work:

1. Adoption of city-wide IoT air quality networks by state and local governments.
2. Integration of predictive alerts into public health systems.
3. Expansion to additional pollutants (SO₂, NO_x) and longer temporal datasets.

Fusion with satellite and meteorological forecast data for regional scaling.

REFERENCES

- [1] M. Kampa and E. Castanas, "Human health effects of air pollution," *Environmental Pollution*, vol. 151, no. 2, pp. 362–367, 2008.
- [2] A. J. Cohen, M. Brauer, R. Burnett, H. R. Anderson, J. Frostad, K. Estep, K. Balakrishnan, B. Brunekreef, L. Dandona, R. Dandona, V. Feigin, G. Freedman, B. Hubbell, A. Jobling, H. Kan, L. Knibbs, Y. Liu, R. V. Martin, L. Morawska, C. A. Pope III, H. Shin, K. Straif, G. Shaddick, M. Thomas, R. van Dingenen, A. van Donkelaar, T. Vos, C. J. L. Murray, and M. H. Forouzanfar, "Estimates and trends of the global burden of disease attributable to ambient air pollution: An analysis of data from the Global Burden of Diseases Study 2015," *The Lancet*, vol. 389, no. 10082, pp. 1907–1918, 2017.
- [3] World Health Organization, *Ambient (outdoor) air pollution*, Geneva, Switzerland, 2018.

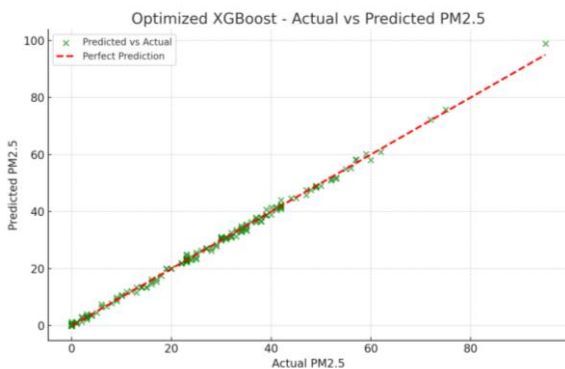


Fig.5: PM_{2.5} Forecast vs Actual using XGBoost

Ensemble stability across pollution peaks.

4.7 PM₁₀ Prediction Results

TABLE 4.8: PM₁₀ Prediction Model Performance

Model	MAE	RMSE	R ²
Extra Trees	0.015	0.101	0.999
Random Forest	0.017	0.109	0.998
XGBoost	0.021	0.126	0.997
MLP	0.247	0.682	0.873

4.8 Cross-Validation Results

TABLE 4.9: 5-Fold Cross-Validation R² Scores

Model	Mean R ²	Std. Dev
Random Forest	0.9987	0.0012
Extra Trees	0.9989	0.0010
XGBoost	0.9991	0.0008
MLP	0.9345	0.0421

4.9 AQI Prediction Performance

Ensemble models consistently outperformed baseline methods. XGBoost achieved the best AQI prediction performance (R² ≈ 1.000, RMSE < 0.1), closely followed by Random Forest and Extra Trees. These results highlight the

- [4] M. Brauer, M. Amann, R. T. Burnett, A. Cohen, F. Dentener, M. Ezzati, S. B. Henderson, M. Krzyzanowski, R. V. Martin, R. Van Dingenen, A. van Donkelaar, and G. D. Thurston, "Exposure assessment for estimation of the global burden of disease attributable to outdoor air pollution," *Environmental Science & Technology*, vol. 46, no. 2, pp. 652–660, 2012.
- [5] World Health Organization, WHO Global Air Quality Guidelines, Geneva, Switzerland, 2021.
- [6] A. Abdulraheem, R. Yusuf, and M. Bello, "Assessment of ambient air pollution and associated health risks in selected urban centers in the Niger Delta, Nigeria," *Environmental Monitoring and Assessment*, 2023.
- [7] Y. Cheng, K. B. He, F. K. Duan, M. Zheng, Z. Y. Du, and Y. L. Ma, "Ambient organic carbon and elemental carbon in PM_{2.5} in Chinese megacities," *Atmospheric Chemistry and Physics*, vol. 11, no. 22, pp. 11497–11510, 2014.
- [8] J. R. Balmes and M. Guarnieri, "Outdoor air pollution and cardiovascular disease," *Current Epidemiology Reports*, vol. 1, no. 4, pp. 279–287, 2014.
- [9] F. Minichilli, F. Gorini, E. Bustaffa, L. Cori, F. Bianchi, and P. Comba, "Mortality and morbidity in petrochemical areas: A systematic review," *Environmental Research*, vol. 171, pp. 556–569, 2019.
- [10] P. Amoatey, H. Sulaiman, and M. S. Baawain, "Impact of emissions from oil and gas operations on local air quality: A review," *Environmental Science and Pollution Research*, vol. 27, pp. 11815–11828, 2020.
- [11] A. Kumar and A. Jasuja, "Air quality monitoring system based on IoT using Raspberry Pi," in *Proc. International Conference on Computing, Communication and Automation (ICCCA)*, pp. 1341–1346, 2017.
- [12] G. Fioccola, R. Sommese, I. Tufano, R. Canonico, and G. Ventre, "Pollution monitoring using sensor networks: The IoT approach," *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 1–12, 2016.
- [13] N. Zimmerman, A. A. Presto, S. P. N. Kumar, J. Gu, A. Hauryliuk, E. S. Robinson, R. Subramanian, and A. L. Robinson, "A machine learning calibration model using random forests to improve sensor performance," *Atmospheric Measurement Techniques*, vol. 11, pp. 291–313, 2018.
- [14] P. Subramaniam, S. Lee, and K. Wong, "Machine learning approaches for air quality index prediction: A comparative study," *Environmental Modelling & Software*, vol. 148, 2022.
- [15] D. O. Harrop, *Air Quality Assessment and Management*. Boca Raton, FL: CRC Press, 2018.
- [16] U.S. Environmental Protection Agency, *Integrated Science Assessment for Particulate Matter*, Washington, DC, 2023.
- [17] F. C. Obodoeze, C. A. Nwabueze, and S. A. Akaneme, "Internet of Things (IoT)-based real-time pollution monitoring system for Awka Metropolis," *International Journal of Trend in Scientific Research and Development*, vol. 6, no. 1, pp. 1513–1523, 2021.
- [18] C. A. Nwabueze, S. A. Akaneme, and F. C. Obodoeze, "Air pollution modeling for Awka Metropolis using ensemble algorithms," *International Journal of Advanced Research in Computer and Communication Engineering*, vol. 11, no. 1, 2022.
- [19] A. W. Correia, C. A. Pope III, D. W. Dockery, Y. Wang, M. Ezzati, and F. Dominici, "Effect of air pollution control on life expectancy in the United States," *Epidemiology*, vol. 24, no. 1, pp. 23–31, 2013.
- [20] United Nations Children's Fund (UNICEF), *Air Pollution and Child Health: Prescribing Clean Air*, New York, 2021.
- [21] O. Postolache, J. Pereira, and P. Girão, "Smart sensors network for air quality monitoring applications," *IEEE Transactions on Instrumentation and Measurement*, vol. 58, no. 9, pp. 3253–3262, 2009.
- [22] X. Zhao, J. Zhang, and Y. Wang, "IoT-based air pollution monitoring systems: A review," *IEEE Access*, vol. 8, pp. 220458–220476, 2020.
- [23] A. Botta, W. de Donato, V. Persico, and A. Pescapé, "Integration of cloud computing and Internet of Things: A survey," *Future Generation Computer Systems*, vol. 56, pp. 684–700, 2016.
- [24] P. Arroyo, J. Herrero, and V. Tricio, "Cloud-based IoT platforms for environmental monitoring applications," *Sensors*, vol. 19, no. 6, 2019.
- [25] A. Muhammed, Z. Khan, and S. Ali, "Air quality prediction using machine learning algorithms," *International Journal of Computer Applications*, vol. 123, no. 9, pp. 1–6, 2015.