



# **AirSense - AI-Powered Air Quality Monitoring & Forecasting**

Literature, Technology and Data Review Submission

## **GROUP MEMBERS**

1. Tekleeyesus Munye
2. Mukitar Seid
3. Dawit Getachew
4. Helen Zelalem
5. Gelasa J.

## CONTENTS

---

<b>Literature Review.....</b>	<b>4</b>
Introduction.....	4
Thematic Organization of Literature.....	4
1. Classical Machine Learning Approaches.....	4
Survey of ML Algorithms for Air Quality Forecasting.....	4
AirNet: Web-Based ML Model.....	5
2. Deep Learning Frameworks.....	5
Hybrid DL with Spatial Autocorrelation.....	5
3. Surveys and Reviews of DL Techniques.....	5
Deep Learning for Air Quality Forecasts: A Review.....	5
Summary and Synthesis.....	6
Comparative Insights.....	6
Conclusion.....	7
Key Takeaways.....	7
Contribution of AirSense	
Building upon these insights, AirSense will:.....	7
References to Research Papers.....	8
<b>Data Research Report: AirSense - AI-Powered Air Quality Monitoring &amp; Forecasting.....</b>	<b>8</b>
Introduction.....	8
Importance of the Research Questions.....	8
Organization of Data Research.....	9
Data Description.....	9
Dataset Overview.....	9
Why This Dataset?.....	10
Data Preprocessing Needs.....	10
Data Analysis and Insights.....	10
Descriptive Statistics.....	10
Visualizations.....	10
Key Findings.....	11
Conclusion.....	11
Summary of Insights.....	11
Relevance to AirSense.....	11
Citations.....	11
<b>Comprehensive Technology Review: AI/ML Tools for Air Quality Monitoring &amp; Forecasting.....</b>	<b>12</b>
Introduction.....	12
Importance of the Technology Review.....	12
Relevance to AirSense's Goals.....	13

Technology Overview.....	13
Machine Learning Frameworks.....	13
(A) Scikit-learn.....	13
(B) TensorFlow/Keras.....	14
(C) PyTorch.....	14
Data Processing Tools.....	14
Visualization & Deployment.....	15
• Matplotlib/Plotly:.....	15
• Flask/Django:.....	15
Relevance to AirSense.....	15
Addressing Project Challenges.....	15
Why LSTMs + Scikit-learn?.....	15
Comparison and Evaluation.....	16
ML/DL Framework Comparison.....	16
Data Tools Comparison.....	16
Use Cases & Case Studies.....	17
Real-World Implementations.....	17
1. IBM Green Horizons.....	17
2. OpenAQ Platform.....	17
Lessons for AirSense.....	17
Gaps and Research Opportunities.....	17
Identified Limitations.....	17
Future Enhancements.....	17
Conclusion.....	18
Citations.....	18

# Literature Review

## Introduction

Air pollution poses significant risks to human health and urban sustainability, contributing to respiratory illnesses, cardiovascular diseases, and premature mortality. Monitoring and forecasting key pollutants—such as fine particulate matter (PM<sub>2.5</sub>) and nitrogen dioxide (NO<sub>2</sub>)—are essential for issuing timely health advisories and informing policy interventions. While traditional deterministic models (e.g., chemical transport models) provide valuable insights, they are computationally intensive and often lack the granularity needed for real-time urban applications. Consequently, data-driven approaches leveraging machine learning (ML) and deep learning (DL) have emerged as powerful alternatives, offering efficient, scalable, and accurate predictions of air quality indices. A systematic review of existing literature is therefore necessary to synthesize current methodologies, identify gaps, and position the proposed **AirSense** system within this evolving landscape.

## Thematic Organization of Literature

### 1. Classical Machine Learning Approaches

#### Survey of ML Algorithms for Air Quality Forecasting

Méndez *et al.* (2023) provide a comprehensive survey of classical ML techniques, including Random Forests, Support Vector Regression, and Gradient Boosting Machines—for forecasting air quality indices across diverse urban environments. They report that ensemble methods (e.g., Random Forest) often outperform single learners due to their robustness against overfitting and capacity to capture nonlinear pollutants–meteorology relationships [link.springer](#).

#### **AirNet: Web-Based ML Model**

Rahman *et al.* (2024) introduce **AirNet**, a predictive ML framework with a user-friendly web interface that forecasts PM<sub>2.5</sub> and NO<sub>2</sub> levels using historical pollutant concentrations and meteorological features. AirNet employs feature selection via recursive elimination and optimizes ensemble regressors, achieving root mean square error (RMSE) improvements of 12–18% over baseline linear models [AirNet](#). The study highlights the importance of interpretability for stakeholder adoption, aligning with AirSense's goal of providing actionable insights to city officials.

## 2. Deep Learning Frameworks

### DeepAir: Convolutional LSTM for Spatiotemporal Patterns

Alléon *et al.* (2020) propose **PlumeNet**, a convolutional LSTM architecture that jointly forecasts  $\text{NO}_2$ ,  $\text{O}_3$ ,  $\text{PM}_{2.5}$ , and  $\text{PM}_{10}$  over a  $0.5^\circ$  grid by integrating ground-monitor data, weather forecasts, and outputs from physical–chemical models. PlumeNet achieves up to 25% lower mean absolute error (MAE) for four-day forecasts compared to persistence and standard LSTM baselines [arxiv](#). Its design underscores the value of capturing both spatial autocorrelation (via convolutions) and temporal dependencies (via LSTMs), a strategy AirSense may adapt for high-resolution urban deployments.

### Hybrid DL with Spatial Autocorrelation

Zhao *et al.* (2023) develop a hybrid DL model combining graph convolutional networks (GCNs) with LSTM layers to account for spatial autocorrelation among monitoring stations during the COVID-19 period. Their framework outperforms pure LSTM and GCN models by 8–15% in RMSE for  $\text{PM}_{2.5}$  prediction, demonstrating that embedding spatial topology significantly enhances forecasting accuracy [nature](#). This insight informs AirSense’s potential integration of graph-based modules to model intra-city pollutant dispersion.

## 3. Surveys and Reviews of DL Techniques

### Deep Learning for Air Quality Forecasts: A Review

Qi *et al.* (2019) review DL architectures applied to air quality forecasting, including autoencoders for data gap filling, CNNs for spatial feature extraction, and attention mechanisms for dynamic temporal weighting. They identify challenges, such as data sparsity, model interpretability, and transferability across regions—and recommend hybridizing DL with physical models to improve generalization [researchgate](#). This survey underscores the necessity of rigorous preprocessing (e.g., imputation, scaling) and domain knowledge incorporation, both of which are integral to AirSense’s preprocessing pipeline.

## Summary and Synthesis

Study	Methodology	Key Findings	Contribution
Méndez et al. (2023)	Survey of Random Forest, SVR, GBM	Ensemble methods yield superior accuracy; feature importance analysis guides pollutant drivers	Benchmark of classical ML approaches
Rahman et al. (2024) – AirNet	Recursive feature elimination + ensemble regressors	RMSE reduction of 12–18% over linear models; web interface for stakeholder engagement	Demonstrates interpretability and usability
Alléon et al. (2020) – PlumeNet	ConvLSTM integrating ground data, weather forecasts, AQPCM outputs	25% lower MAE for 4-day forecasts vs. persistence/LSTM baselines	Joint spatiotemporal modeling at continental scale
Zhao et al. (2023)	GCN + LSTM hybrid model	RMSE improvements of 8–15% over pure DL models; spatial autocorrelation boosts accuracy	Highlights spatial graph integration
Qi et al. (2019)	Review of DL: autoencoders, CNNs, attention mechanisms	Identifies data sparsity, interpretability, and transferability as core challenges; advocates hybrid DL–physical modeling	Roadmap for future DL applications

## Comparative Insights

- **Interpretability vs. Accuracy:** Classical ML (e.g., Random Forest) offers interpretability through feature importance, whereas DL models (e.g., PlumeNet, hybrid GCN-LSTM) deliver higher predictive accuracy but at the cost of transparency.
- **Spatial Modeling:** Incorporating spatial dependencies—either via convolutional layers (PlumeNet) or graph structures (Zhao *et al.*)—consistently enhances performance, suggesting AirSense should embed spatial modules for intra-urban forecasts.
- **Hybrid Architectures:** Merging physical model outputs with DL (as in PlumeNet) or blending GCNs with LSTMs (Zhao *et al.*) yields robust forecasts, pointing to the benefit of multimodal data fusion.

## Conclusion

### Key Takeaways

- Ensemble ML methods remain competitive for short-term, low-dimensional forecasting tasks, offering ease of interpretation and rapid deployment.
- Deep learning frameworks, particularly those integrating spatial autocorrelation (via CNNs or GCNs) and temporal sequence modeling (via LSTMs), achieve superior accuracy for multi-pollutant, multi-day forecasts.
- Hybrid models that fuse deterministic physical outputs with data-driven learning provide a promising pathway to balance accuracy, scalability, and computational efficiency.

### Contribution of AirSense

Building upon these insights, **AirSense** will:

1. **Data Preprocessing:** Employ rigorous missing-value imputation, feature scaling, and time-series windowing on the Beijing Multi-Site Air-Quality Data Set.
2. **Modeling Approach:** Start with interpretable ML baselines (Random Forest, Gradient Boosting) and progressively integrate spatiotemporal DL modules (e.g., graph-enhanced LSTM) to capture urban dispersion patterns.
3. **Actionable Alerts:** Translate continuous forecasts into health advisories aligned with WHO guidelines, empowering residents and officials to mitigate exposure risks.

By synthesizing classical and advanced methodologies, AirSense aims to deliver an end-to-end, AI-powered air quality monitoring and forecasting platform that advances SDG 3 (Good Health and Well-Being) and SDG 11 (Sustainable Cities and Communities).

## References to Research Papers

1. Méndez, M., Merayo, M. G., & Núñez, M. (2023). *Machine learning algorithms to forecast air quality: a survey*. Artificial Intelligence Review, 56, 10031–10066.  
<https://doi.org/10.1007/s10462-023-10424-4> [link.springer](#)
2. Rahman, M. M., et al. (2024). *AirNet: predictive machine learning model for air quality forecasting using web interface*. Environmental Systems Research, 13, Article 44.

<https://doi.org/10.1186/s40068-024-00378-z>

[environmentalsystemsresearch.springeropen](https://environmentalsystemsresearch.springeropen)

3. Alléon, A., Jauvion, G., Quennehen, B., & Lissmyr, D. (2020). *PlumeNet: Large-Scale Air Quality Forecasting Using A Convolutional LSTM Network*. arXiv:2006.09204. <https://arxiv.org/abs/2006.09204> [arxiv](https://arxiv.org/abs/2006.09204)
4. Zhao, Z., Wu, J., Cai, F., Zhang, S., & Wang, Y. G. (2023). *A hybrid deep learning framework for air quality prediction with spatial autocorrelation during the COVID-19 pandemic*. Scientific Reports, 13, 1015. <https://doi.org/10.1038/s41598-023-28287-8> [nature](https://doi.org/10.1038/s41598-023-28287-8)
5. Qi, L., et al. (2019). *Deep Learning for Air Quality Forecasts: a Review*. ResearchGate. [https://www.researchgate.net/publication/344086668\\_Deep\\_Learning\\_for\\_Air\\_Quality\\_Forecasts\\_a\\_Review](https://www.researchgate.net/publication/344086668_Deep_Learning_for_Air_Quality_Forecasts_a_Review) [researchgate](https://www.researchgate.net/publication/344086668_Deep_Learning_for_Air_Quality_Forecasts_a_Review)

# Data Research Report: AirSense - AI-Powered Air Quality Monitoring & Forecasting

## Introduction

### Importance of the Research Questions

Air pollution is a leading environmental health risk, contributing to 7 million premature deaths annually (WHO, 2022). The AirSense project seeks to answer critical research questions:

- How can historical air quality data predict future pollution levels (PM<sub>2.5</sub>, NO<sub>2</sub>)?
- What machine learning models perform best for time-series air quality forecasting?
- How can real-time data improve public health decision-making?

A thorough exploration of data is essential because:

- ✓ Data quality directly impacts model accuracy (e.g., handling missing sensor readings).
- ✓ Identifying trends (e.g., seasonal pollution spikes) informs model selection.
- ✓ Public health interventions require reliable, interpretable forecasts.



# Organization of Data Research

This report is structured thematically:

1. Data Description: Source, format, and relevance.
2. Data Analysis: Key insights, visualizations, and statistics.
3. Conclusion: Summary of findings and project alignment.

## Data Description

### Dataset Overview

Attribute	Details
Dataset Name	Beijing Multi-Site Air-Quality Data
Source	Repository
Format	CSV (12 files, one per monitoring station)
Size	~70 MB (2014–2017 hourly data)
Variables	PM2.5, PM10, NO <sub>2</sub> , SO <sub>2</sub> , temperature, pressure, humidity

### Why This Dataset?

- Relevance: Covers critical pollutants (PM2.5, NO<sub>2</sub>) aligned with AirSense’s goals.
- Geographical Focus: Beijing’s air quality challenges mirror urban areas globally.
- Temporal Granularity: Hourly data enables high-resolution forecasting.

Data Preprocessing Needs

- Missing Values: ~15% of PM2.5 readings (addressed via interpolation).
- Feature Engineering: Time-lagged features for trend capture.
- Normalization: Min-max scaling for neural networks.

Data Analysis and Insights

Descriptive Statistics

Pollutant	Mean	Max	Std Dev	Correlation with PM2.5
PM2.5	98 µg/m³	898 µg/m³	88	1.00
NO₂	52 µg/m³	340 µg/m³	32	0.72
Temperature	12°C	40°C	10	-0.31

Key Observations:

- High PM2.5 variability (std dev = 88) indicates frequent pollution spikes.
- NO₂ strongly correlates with PM2.5, suggesting co-emission sources (e.g., traffic).
- Temperature inversely correlates with PM2.5 (cold weather traps pollutants).

Visualizations

(A) PM2.5 Trends Over Time

(B) Pollutant Correlation Matrix

Key Findings

1. Seasonality: PM2.5 levels surge 3× higher in winter than summer.
2. Diurnal Patterns: NO₂ peaks during rush hours (8 AM, 6 PM).
3. Missing Data: Gaps concentrated in 2014 (requires imputation).

## Conclusion

### Summary of Insights

- ❖ The Beijing dataset validates the feasibility of forecasting PM2.5 using ML.
- ❖ NO<sub>2</sub> and temperature are critical auxiliary features.
- ❖ Time-series gaps must be addressed to avoid model bias.

### Relevance to AirSense

- ❖ Model Selection: LSTMs will capture seasonal/diurnal trends.
- ❖ Health Alerts: Real-time dashboards can highlight high-risk periods.
- ❖ SDG Alignment: Data-driven insights support SDG 3 (Health) and SDG 11 (Cities).

## Citations

1. WHO. (2022). *Air Pollution and Health*. [Link]
2. Zheng, Y., et al. (2015). *Beijing Multi-Site Air-Quality Data*. UCI.
3. Liang, X., et al. (2021). *DeepAir: Forecasting PM2.5 with LSTMs*. IEEE.

# Comprehensive Technology Review: AI/ML Tools for Air Quality Monitoring & Forecasting

*(AirSense - AI-Powered Air Quality Monitoring & Forecasting System)*

## Introduction

### Importance of the Technology Review

Air pollution is responsible for an estimated **7 million premature deaths annually** (WHO, 2022). The **AirSense** project aims to mitigate this crisis by developing an AI-driven system that:

- **Forecasts** PM2.5, NO<sub>2</sub>, and other pollutants using historical and real-time data.
- **Issues health alerts** to urban residents and policymakers.
- **Supports SDG 3 (Good Health & Well-being)** and **SDG 11 (Sustainable Cities)** by enabling data-driven interventions.

This review evaluates **machine learning (ML) and deep learning (DL) technologies** to determine the optimal tools for:

- ✓ **Time-series forecasting** (LSTMs, Gradient Boosting)
- ✓ **Real-time data processing** (Apache Kafka, Dask)
- ✓ **Visualization & public alerts** (Plotly, Flask)

### Relevance to AirSense's Goals

Project Need	Technology Solution
Accurate PM2.5 forecasting	LSTMs (for temporal patterns)

Handling missing sensor data	Pandas + Scikit-learn imputation
Scalable real-time processing	Apache Kafka (if deploying IoT sensors)
Public-friendly health alerts	Plotly Dash/Flask web dashboard

## Technology Overview

### Machine Learning Frameworks

#### (A) Scikit-learn

- **Purpose:** Prototyping traditional ML models (Random Forest, Gradient Boosting).
- **Key Features:**
  - Fast training for structured data.
  - Interpretability (feature importance analysis).
  - Integrates with Pandas for preprocessing.
- **Common Use Cases:**
  - Baseline AQI prediction (e.g., *Beijing Air Quality Dataset*).
  - IBM's early air quality models.

#### (B) TensorFlow/Keras

- **Purpose:** Deep learning for complex time-series forecasting.
- **Key Features:**

- o Built-in LSTM/GRU layers for sequential data.
  - o GPU acceleration for large datasets.
- **Common Use Cases:**
  - o **DeepAir** (LSTM-based PM2.5 forecasting).
  - o NASA's pollution trend modeling.

### (C) PyTorch

- **Purpose:** Research-focused DL with dynamic computation graphs.
- **Key Features:**
  - o Flexible architecture experimentation.
  - o Preferred for cutting-edge NN research.
- **Common Use Cases:**
  - o Hybrid models (e.g., CNN-LSTM for spatial-temporal data).

### Data Processing Tools

Tool	Role in AirSense
<b>Pandas</b>	Clean, impute, and scale the Beijing dataset.
<b>Dask</b>	Parallelize preprocessing for scalability.
<b>Apache Kafka</b>	Stream real-time sensor data (future phase).

## Visualization & Deployment

- **Matplotlib/Plotly:**
  - Generate interactive maps of pollution hotspots.
  - Example: *Plotly Dash* for real-time AQI dashboards.
- **Flask/Django:**
  - Deploy forecasts as a web app for public alerts.

## Relevance to AirSense

### Addressing Project Challenges

Challenge	Technology Solution
Missing data in Beijing dataset	Pandas (interpolation) + Scikit-learn impute
Capturing long-term trends	LSTMs (Keras)
Explaining model decisions	SHAP (for Scikit-learn interpretability)

### Why LSTMs + Scikit-learn?

- **LSTMs:** Outperform ARIMA and Random Forest in **temporal dependency** tasks (*DeepAir, 2021*).
- **Scikit-learn:** a **baseline** for benchmarking DL model performance.

## Comparison and Evaluation

### ML/DL Framework Comparison

Criterion	Scikit-learn	LSTM (Keras)	PyTorch
Ease of Use	★★★★★	★★★★★	★★★★
Interpretability	High (feature importance)	Medium (attention layers)	Low (research-focused)
Scalability	Good (CPU-friendly)	Excellent (GPU support)	Excellent (GPU)
Accuracy	Moderate (for trends)	High (for sequences)	High (customizable)

### Data Tools Comparison

Tool	Best For	Limitations
Pandas	Data cleaning & wrangling	Struggles with >10GB datasets
Dask	Distributed preprocessing	Requires cluster setup
Apache Kafka	Real-time IoT data streaming	Overkill for static datasets



## Use Cases & Case Studies

### Real-World Implementations

#### 1. IBM Green Horizons

- **Tech Stack:** LSTMs + Kafka.
- **Outcome:** 30% improvement in Beijing's PM2.5 forecasts.

#### 2. OpenAQ Platform

- **Tech Stack:** Scikit-learn + Plotly.
- **Outcome:** Global real-time AQI visualization.

### Lessons for AirSense

- **LSTMs** are proven for pollution forecasting but require **large datasets**.
- **Hybrid models** (e.g., CNN-LSTM) may improve sudden spike predictions.

## Gaps and Research Opportunities

### Identified Limitations

- **Data Quality:** Beijing dataset has missing entries (requires advanced imputation).
- **Model Explainability:** DL models are "black boxes" (SHAP/LIME can help).

### Future Enhancements

- **Generative Models:** GANs to synthesize missing sensor data.
- **Edge AI:** Deploy lightweight models on IoT devices.

## Conclusion

For **AirSense**, we recommend:

1. **Scikit-learn** for baseline models (interpretability).
2. **LSTMs (Keras)** for final high-accuracy forecasts.
3. **Plotly Dash** for public-facing alerts.

This stack balances **accuracy, scalability, and usability** while aligning with SDG 3 and 11.

## Citations

1. WHO. (2022). *Air Pollution and Health*. [Link]
2. Zheng, Y., et al. (2015). *Beijing Multi-Site Air-Quality Data*. UCI.
3. Liang, X., et al. (2021). *DeepAir: LSTM Forecasting*. IEEE.
4. IBM. (2016). *Green Horizons Initiative Case Study*.