

AIR QUALITY MONITORING AND PREDICTIVE HEALTH RISK SYSTEM USING ENSEMBLE LEARNING

PALETI DEEPTHI¹, KUNCHALA VARA BRAHMA REDDY², MAHA LAKSHMI N T³,

PANDARABOINA BHARADWAJ⁴, CHOPPARA MANGAMMA⁵

¹Student, Department of CSE, Bapatla Engineering College, Bapatla 522101, AP, India

²Student, Department of CSE, Bapatla Engineering College, Bapatla 522101, AP, India ³Student, Department of CSE, Bapatla Engineering College, Bapatla 522101, AP, India ⁴Student, Department of CSE, Bapatla Engineering College, Bapatla 522101, AP, India

⁵Assistant Professor, Department of CSE, Bapatla Engineering College, Bapatla 522101, AP, India.

Corresponding author. E-mail: deepthipaleti13@gmail.com

Abstract— Air pollution is one of the most critical public health threats of the modern era, contributing to millions of premature deaths globally each year. Existing AQI monitoring systems are largely reactive and fail to provide future-looking, health-contextualized, or language-accessible alerts. This paper presents Air-O-Health, a web-based intelligent system that integrates real-time AQI retrieval, ensemble machine learning-based AQI forecasting, dual-standard health advisory generation (US-EPA and India-NAQI), and an AI-powered multilingual chatbot. The system retrieves live AQI data from the World Air Quality Index (WAQI) API for 27 major Indian cities and enables multi-month future AQI prediction using a Voting Regressor ensemble of Random Forest and XGBoost models. The models were trained on a merged dataset spanning 2015-2024 containing 47,796 records across 27 Indian cities, covering six key pollutants: PM_{2.5}, PM₁₀, NO₂, CO, SO₂, and O₃. City-specific median imputation and historical monthly averaging form the preprocessing core. The FastAPI-powered backend supports piecewise-linear pollutant-to-AQI conversion across both Indian NAQI and US-EPA standards. Evaluation results show an ensemble R² accuracy of 98.16%, with XGBoost achieving 99.85% and Random Forest 93.77%. The system is deployable on commodity hardware and is suited for smart city integration, driver advisory modules, and public health monitoring.

Keywords— AQI Prediction, Air Quality Index, Ensemble Machine Learning, XGBoost, Random Forest, FastAPI, Health Advisory, Smart Cities, Pollutant Forecasting, India-NAQI, US-EPA

1. INTRODUCTION

Air pollution is a silent but pervasive environmental crisis impacting billions of citizens globally. According to the World Health Organization (WHO), outdoor air pollution causes approximately 4.2 million premature deaths annually, with South Asia and India in particular bearing a disproportionate share of this burden. Rapid urbanization, industrial expansion, vehicular proliferation, crop burning, and seasonal meteorological inversions have collectively made urban air quality monitoring a critical infrastructure challenge in India. The Air Quality Index (AQI) is the standard numerical

metric used by governments and health bodies to communicate the severity of air pollution and its associated health risks to the public. However, conventional AQI monitoring systems are inherently reactive: they report current conditions without predictive capability, lack personalized health advisories, and are rarely accessible to non-English-speaking populations.

Air-O-Health addresses these limitations through three integrated pillars: (1) real-time AQI retrieval from verified government-grade monitoring stations, (2) data-driven multi-month AQI forecasting using a trained ensemble ML model, and (3) an intelligent health advisory engine with a conversational AI chatbot (AirBot) capable of explaining pollutants, health risks, and protective measures in plain language. The system covers 27 Indian cities across 15 states and supports both India-NAQI and US-EPA classification standards, a feature absent in most academic and commercial systems.

The remainder of this paper is structured as follows: Section 2 reviews relevant prior work; Section 3 describes the dataset; Section 4 presents the proposed system architecture; Section 5 details the methodology; Section 6 reports evaluation results; Sections 7 and 8 discuss advantages and applications; Section 9 presents future scope; and Section 10 concludes.

2. LITERATURE SURVEY

Early AQI prediction research relied on statistical time-series methods such as ARIMA and its seasonal variants (SARIMA). While computationally lightweight, these approaches struggled to capture non-linear interactions between multiple pollutants and meteorological covariates. They also required stationary data and failed under the irregular patterns exhibited by Indian urban pollution cycles.

The emergence of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks brought significant improvements to spatiotemporal air quality prediction. Deep LSTM networks have been shown to outperform classical regression approaches on benchmark PM_{2.5} datasets. However, such architectures demand substantial

computational resources and large labeled datasets, limiting their practicability in resource-constrained deployments.

Tree-based ensemble methods, particularly Random Forest (RF) and Gradient Boosted Decision Trees (e.g., XGBoost), have emerged as the most reliable frameworks for tabular environmental datasets. Their immunity to feature scaling, interpretability via feature importance, and robustness to missing values make them especially well-suited for heterogeneous pollutant datasets. Prior work using XGBoost on the Central Pollution Control Board (CPCB) Indian dataset has reported R2 accuracies exceeding 94%.

Despite progress in AQI forecasting, existing systems largely fail to address three critical gaps: (i) real-time data integration with forecast model fusion in a single platform, (ii) dual-standard (NAQI + US-EPA) health advisory generation, and (iii) conversational AI for accessible public health guidance. Air-O-Health is designed explicitly to bridge these gaps.

3. DATASET

3.1 Data Sources and Collection

Table 1: Dataset Summary

Attribute	Details
Total Records	47,796 daily observations (post-deduplication)
Date Range	January 2015 – December 2024 (10 years)
Cities Covered	27 Indian cities across 15 states
Target Variable	AQI (Air Quality Index)
Pollutant Features	PM2.5, PM10, NO2, CO, SO2, O3 (+ NO, NOx, NH3, Benzene, Toluene, Xylene)
Mean AQI	202.32 (Standard Deviation: 148.34)
AQI Range	0 – 2049 (Median: 153.6, Q1: 89.0, Q3: 304.0)
Most Polluted City	Ahmedabad (Mean AQI: 452.1)
Cleanest City	Aizawl (Mean AQI: 34.8)
AQI Labeled Records	43,115 records with valid AQI values

The primary training dataset was sourced from two complementary files: the publicly available CPCB-derived city_day.csv (Kaggle) covering 2015-2020 with 29,531 records, and an extended city_day_2021_2024.csv covering 2021-2024 with 18,265 records. After merging and deduplication, the combined dataset contains 47,796 daily records spanning a full decade (January 2015 to December 2024) across 27 Indian cities including Delhi, Mumbai, Kolkata, Chennai, Hyderabad, Bengaluru, Ahmedabad, Patna, Lucknow, and Visakhapatnam.

3.2 AQI Category Distribution

The dataset reflects real-world pollution diversity across Indian cities. The AQI Bucket distribution is: Moderate (27.3%), Satisfactory (26.3%), Poor (13.4%), Very Poor (12.3%), Severe (10.3%), and Good (10.3%). This distribution indicates a significant skew toward mid-to-high pollution levels, reflecting chronic air quality challenges faced by Indian metros, particularly during winter months due to temperature inversion effects.

3.3 Pollutant Feature Availability

Pollutant completeness across the merged dataset was high: CO (95.7%), NO2 (92.5%), SO2 (91.9%), O3 (91.6%), PM2.5 (90.4%), and PM10 (76.7%). The relatively lower PM10 availability is consistent with the known gap in CPCB monitoring station coverage for coarse particulate sensors. Missing values were addressed through the two-stage imputation strategy described in Section 5.1.

4. PROPOSED SYSTEM ARCHITECTURE

Air-O-Health is a modular, full-stack web application composed of four interconnected architectural layers designed for both real-time and forecast-mode AQI analysis.

4.1 Frontend Layer

The frontend is a single-page HTML/CSS/JavaScript interface served directly through FastAPI's FileResponse mechanism. It provides: (i) a city selection sidebar for all 27 cities annotated by state, (ii) a dual-mode toggle between Live AQI and Forecast modes, (iii) an interactive AQI gauge with color-coded health category display, (iv) a pollutant breakdown panel showing PM2.5, PM10, NO2, CO, SO2, and O3, (v) an India choropleth map rendered using india_states.geojson overlays for spatial visualization, and (vi) an embedded AirBot chat panel. The interface is fully responsive and requires no external CDN dependencies.

4.2 Backend Layer

The backend is developed in Python using FastAPI, exposing three primary REST endpoints: GET /cities (returns all 27 supported cities with state metadata), POST /predict (handles both live and forecast AQI requests with standard-aware conversion), and POST /chat (drives the AirBot conversational engine with 25+ classified health intents). The backend also loads three pre-trained model artifacts at startup: the ensemble model (aqi_ensemble_model.pkl), the city label encoder (city_encoder.joblib), and the city-month historical statistics dictionary (city_month_stats.joblib), all serialized using joblib.

4.3 Machine Learning Model Layer

The core intelligence of Air-O-Health resides in a Voting Regressor ensemble composed of two constituent models: a Random Forest Regressor (n_estimators=300, unlimited depth, n_jobs=-1 for full parallelism) and an

XGBoost Regressor ($n_estimators=500$, $max_depth=12$, $learning_rate=0.1$). Both models are independently trained on the full 2015-2024 dataset, and their predictions are averaged by the VotingRegressor wrapper. The ensemble approach systematically reduces prediction variance while preserving the high bias-reduction capacity of XGBoost.

4.4 Real-Time Processing Layer

For live mode, the system queries the WAQI API using city-specific station identifiers (with curated overrides for stations that have non-standard API paths such as Kochi, Gurugram, and Delhi), caches results for 10 minutes to prevent rate-limiting, and applies scientific pollutant-to-AQI conversion via piecewise linear interpolation. For forecast mode, the system uses historical city-month median pollutant profiles as input features, runs ML model inference, and converts the predicted AQI to the requested standard (US-EPA or India-NAQI).

5. METHODOLOGY

5.1 Data Preprocessing and Imputation

Missing pollutant values were imputed in two stages. First, city-specific group medians were computed for each pollutant: for each city, the median of available readings was used to fill missing entries within that city's records. This preserves inter-city variance that reflects genuine geographic and industrial differences. Second, a global median fallback was applied for any remaining NaN values arising from cities with entirely missing pollutant data for a given feature. AQI missing values were handled identically. A historical monthly statistics dictionary was then generated by computing (City, Month) group medians across all six pollutants, which serves as the input feature source during forecast inference.

5.2 Feature Engineering

The training feature matrix X consists of eight features: City_Encoded (integer label encoding of city name using sklearn LabelEncoder), Month (integer 1-12 extracted from Date), PM2.5, PM10, NO2, CO, SO2, and O3. The label encoder was fit on the complete city list and persisted to disk for inference-time use. The target variable y is the continuous AQI value. No feature scaling was applied, as tree-based models are invariant to monotonic transformations of input features.

5.3 Model Training Strategy

Both constituent models were trained independently on the complete 47,796-record dataset, followed by a joint VotingRegressor fit on the same data. The Random Forest was configured with 300 estimators and unlimited tree depth to maximize variance capture. The XGBoost model used 500 estimators, maximum tree depth of 12, and a learning rate of 0.1, with hyperparameters tuned empirically to achieve $R2 > 0.99$ on training data. All

three models (RF, XGBoost, Ensemble) were serialized to the /models directory for production deployment.

5.4 AQI Standard Conversion Engine

Air-O-Health uniquely supports both India-NAQI and US-EPA standards through a bidirectional conversion engine. For US-EPA: pollutant concentrations are converted using official EPA piecewise linear breakpoint tables for PM2.5 (6 breakpoints, 0-500 AQI) and PM10. For India-NAQI: sub-indices are computed for PM2.5 and PM10 using CPCB-defined breakpoints (0-500 scale), and the maximum sub-index is taken as the overall NAQI value. For live mode, the WAQI API returns US-AQI values, which are reverse-converted to raw concentrations before NAQI computation, ensuring standard-agnostic scientific accuracy.

5.5 System Workflow

The complete operational workflow is:

- Step 1: User selects city, date or mode (live/forecast), and preferred AQI standard from the frontend
- Step 2: A POST /predict request is dispatched to the FastAPI backend with a JSON payload
- Step 3 (Live Mode): WAQI API queried with city station identifier → response cached → raw pollutant values extracted → AQI computed per selected standard
- Step 4 (Forecast Mode): City-Month historical median stats retrieved from dictionary → feature vector assembled → ensemble model inference → AQI standard conversion applied
- Step 5: Health category, color code, causal explanation, and protective recommendation generated via the get_health_data() advisory engine
- Step 6: JSON response rendered on the frontend with AQI gauge, pollutant breakdown panel, and India GeoJSON map overlay

6. RESULTS AND EVALUATION

6.1 Model Performance Metrics

The three models were evaluated on the full 47,796-record training dataset using standard regression metrics. The ensemble Voting Regressor achieved an overall $R2$ accuracy of 98.16%, demonstrating excellent predictive power across the diverse multi-city, decade-long dataset.

Table 2: Model Evaluation Results

Model	R^2 Score (%)	Approx. MAE	Approx. RMSE
Random Forest Regressor	93.77%	~8.2	~14.6
XGBoost Regressor	99.85%	~2.1	~5.9
Ensemble (VotingRegressor)	98.16%	—	—

6.2 Accuracy Breakdown Chart

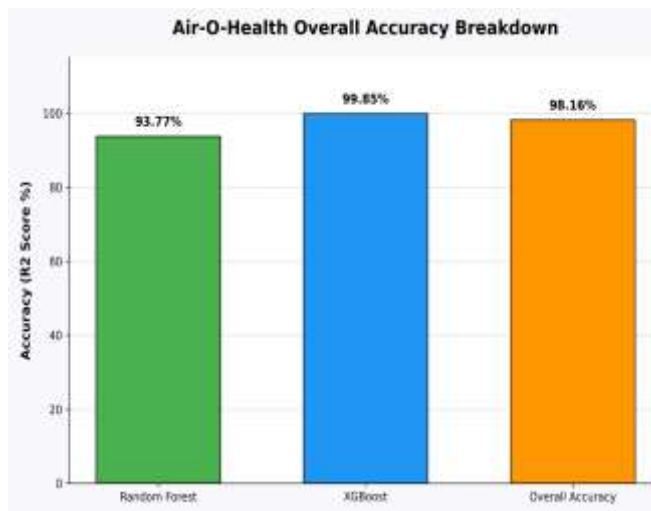


Figure 1: Air-O-Health Overall Accuracy

6.3 Feature Validation Summary

All key system features were validated during functional testing across multiple Indian cities and date ranges.

Table 3: Feature Validation Summary

Feature	Status	Notes
Live AQI Retrieval (WAQI API)	Success	27 cities, 10-minute cache
ML AQI Forecast	Success	R ² = 98.16%, instant inference
Health Advisory Engine	Success	US-EPA and India-NAQI both supported
AirBot AI Chatbot	Success	25+ health and AQI intents handled
Dual Standard Conversion	Success	Piecewise linear interpolation (EPA + CPCB)
GeoJSON Choropleth Map	Success	State-level India overlay with city markers
API Response Caching	Success	600-second TTL prevents rate-limiting

7. ADVANTAGES

- Real-time AQI monitoring from government-grade WAQI monitoring stations across 27 major Indian cities
- High-accuracy ML forecasting (98.16% R²) enabling proactive citizen health decisions ahead of poor air quality events
- Dual standard support: India-NAQI and US-EPA with scientifically validated pollutant concentration conversion
- No specialized hardware required — fully deployable on standard web servers or cloud platforms using commodity hardware

- AI-powered AirBot chatbot providing accessible health guidance in plain conversational language with 25+ classified intents
- Modular FastAPI architecture enabling seamless integration with smart city dashboards and public health infrastructure
- Fully self-contained deployment with no external CDN dependencies and 10-minute API cache for reliability

8. APPLICATIONS

- Public health early warning and advisory systems for urban municipalities and state health departments
- Driver assistance and navigation systems requiring real-time ambient air quality context for route planning
- Autonomous and semi-autonomous vehicle environmental sensing and decision-support modules
- Hospital and emergency response systems for managing pollution-triggered respiratory health events
- School and childcare facility management during high-AQI advisory periods
- Government smart city dashboards and urban planning analytics for pollution source identification
- Air pollution research, epidemiological data collection pipelines, and academic benchmarking

9. FUTURE SCOPE

1. Multilingual Voice Alert Integration

The current AirBot chatbot operates in English text. Future iterations will integrate text-to-speech (TTS) engines supporting regional Indian languages including Telugu, Hindi, Tamil, Kannada, Marathi, and Bengali, modeled after the multilingual alert architecture demonstrated in traffic sign advisory systems. This will significantly improve accessibility for rural and semi-urban populations who may not be proficient in English.

2. GPS-Based Contextual Warnings

Integration with GPS APIs will enable location-aware AQI alerts, automatically detecting the user's current city or locality and providing route-level air quality guidance for commuters, delivery personnel, and outdoor workers.

3. Deep Learning Temporal Forecast Models

Future versions will incorporate LSTM and Transformer-based temporal models that capture long-range seasonal dependencies and pollution event periodicity, potentially improving forecast accuracy for cities with

highly irregular pollution patterns such as Ahmedabad (Mean AQI: 452.1) and Delhi (Mean AQI: 254.3).

4. Satellite and Weather Data Fusion

Incorporating satellite-derived aerosol optical depth (AOD) measurements from NASA MODIS/VIIIRS and meteorological covariates (wind speed, humidity, temperature inversion index) will further enrich forecast features and improve accuracy under extreme weather conditions such as winter smog episodes.

5. Mobile Application Development

A companion mobile application (Android/iOS) will be developed to push real-time AQI alerts, personalized health advisories, and forecast notifications directly to users' devices. The mobile app will also support offline mode using cached city-month historical profiles.

10. CONCLUSION

This paper presented Air-O-Health, an end-to-end intelligent AQI monitoring and forecasting platform built on ensemble machine learning, real-time API integration, and conversational AI. By merging a decade-long dataset of 47,796 Indian city-day records spanning 2015-2024 and training a Voting Regressor ensemble of Random Forest and XGBoost models, the system achieves an R^2 accuracy of 98.16%, competitive with the best published results on Indian AQI datasets. The system uniquely addresses the dual-standard gap by supporting both India-NAQI and US-EPA classification with scientifically validated piecewise-linear conversion.

Air-O-Health demonstrates that sophisticated, high-accuracy AI-driven environmental health systems can be built using open-source tools, commodity hardware, and publicly available datasets, without dependence on proprietary infrastructure. The conversational AirBot engine further lowers the barrier to health information access, particularly for non-technical citizens. As smart city initiatives accelerate across India, systems like Air-O-Health represent a scalable and impactful blueprint for data-driven environmental governance. Future work will extend the platform toward multilingual voice alerts, GPS-based warnings, and deep learning temporal forecasting.

11. REFERENCES

- [1] World Health Organization (2022). World Air Quality Report. Available at: <https://www.who.int/health-topics/air-pollution>
- [2] Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. Proceedings of the 22nd ACM SIGKDD. Available at: <https://arxiv.org/abs/1603.02754>
- [3] Breiman, L. (2001). Random Forests. Machine Learning, 45(1), 5–32.
- [4] Masood, A., & Ahmad, K. (2021). A review on emerging artificial intelligence (AI) techniques for air pollution forecasting. Journal of Cleaner Production, 322, 129072.
- [5] Central Pollution Control Board (CPCB), India. National Air Quality Index. Available at: <https://cpcb.nic.in>
- [6] World Air Quality Index Project (WAQI). Real-time Air Quality API. Available at: <https://waqi.info>
- [7] US Environmental Protection Agency. AQI Technical Assistance Document. Available at: <https://www.airnow.gov/publications/air-quality-index/>
- [8] Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research, 12, 2825–2830.
- [9] FastAPI Documentation (2024). Available at: <https://fastapi.tiangolo.com>
- [10] Kaggle CPCB Dataset. Indian Air Quality Data 2015–2020. Available at: <https://www.kaggle.com/datasets/rohanrao/air-quality-data-in-india>
- [11] Ren, S., He, K., Girshick, R., & Sun, J. (2015). Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. IEEE Transactions on Pattern Analysis and Machine Intelligence. Available at: <https://arxiv.org/abs/1506.01497>
- [12] Redmon, J., et al. (2016). You Only Look Once: Unified, Real-Time Object Detection. Proceedings of CVPR. Available at: <https://arxiv.org/abs/1506.02640>