



1 Research of Low-cost Air Quality Monitoring Models with Different 2 Machine Learning Algorithms

3 Gang Wang^{1, 2, 3}, Chunlai Yu^{1, 3}, Kai Guo², Haisong Guo^{1, 3}, Yibo Wang²

4 ¹Huanghe Science and Technology College, Zhengzhou 450063, China

5 ²Hanwei Electronics Group Corporation, Zhengzhou 450001, China

6 ³Zhengzhou Key Laboratory of Intelligent Measurement Techniques and Applications, Zhengzhou, 450063, China

7 Correspondence to: Gang Wang (wywangang163@163.com)

8 **Abstract.** To improve the prediction for the future air quality trends, the demand for low-cost sensor-based air quality grid
9 monitoring is growing gradually. In this study, a low-cost multi-parameter air quality monitoring system (LCS) based on different
10 machine learning algorithm is proposed. The LCS can measure particulate matter (PM_{2.5} and PM₁₀) and gas pollutants (SO₂, NO₂,
11 CO and O₃) simultaneously. The multi-dimensional multi-response prediction model is developed based on the original signals of
12 the sensors, ambient temperature (*T*) and relative humidity (*RH*), and the measurements of the reference instrumentations. The
13 performance of the different algorithms (RF, MLR, KNN, BP, GA-BP) with the parameters such as determination coefficient R^2
14 and Root Mean Square Error (RMSE) are compared and discussed. Using these methods, the R^2 of the algorithms (RF, MLR, KNN,
15 BP, GA-BP) for the PM is in the range 0.68 - 0.99; the mean RMSE values of PM_{2.5} and PM₁₀ are within 3.96 - 16.16 $\mu\text{g m}^{-3}$ and
16 7.37 - 28.90 $\mu\text{g m}^{-3}$, respectively. The R^2 of the algorithms (RF, MLR, KNN, BP, GA-BP) for the gas pollutants (O₃, CO and NO₂)
17 is within 0.70 - 0.99; the mean RMSE values for these pollutants are 4.06 - 16.07 $\mu\text{g m}^{-3}$, 0.04 - 0.15 mg m^{-3} , 3.25 - 13.90 $\mu\text{g m}^{-3}$,
18 respectively. The R^2 of the algorithms (RF, KNN, BP, GA-BP, except for MLR) for SO₂ is within 0.27 - 0.97, and the mean RMSE
19 value is in the range 1.05 - 3.22 $\mu\text{g m}^{-3}$. These measurements are consistent with the national environmental protection standard
20 requirement of China, and the LCS based on the machine learning algorithms can be used to predict the concentrations of PM and
21 gas pollution.

22 1 Introduction

23 The development along with increased population and urbanization brings disadvantages, such as decreasing air quality and impact
24 on public and individual health (Ioannis et al., 2020; Khreis et al., 2022; Singh et al., 2021). Among the atmospheric pollutants,
25 the primary pollutant is fine particulate matter, which affects the respiratory system and cardiac activity of humans. The secondary
26 pollutants are SO₂, CO, NO_x, and O₃, which also induce disease or chronic poisoning. To improve the understanding of air pollution
27 exposure and to predict future air quality trends (Zimmerman et al., 2018), air quality assessment and forecasting are the essentials.
28 The conventional air quality monitoring instrumentations are high cost, which has limited the spatial coverage of the monitoring
29 stations (Zimmerman et al., 2018). The development and applications of the low-cost commercially available sensor-based air
30 quality monitoring system (LCS) would considerably reduce both installation and maintenance costs (Spinelle et al., 2017). The
31 larger spatial density of the air quality grid monitoring network becomes possible, which would play an important role in
32 monitoring pollution trend, locating of pollution source, supporting environmental management (Zhao et al., 2019) and support
33 better epidemiological models (Khreis et al., 2022; Zimmerman et al., 2018). These demands promote the LCS growing
34 gradually (Cui et al., 2021; Wang et al., 2016).



1 The LCS typically utilizes the electrochemical, metal oxide or light scattering sensors for gas-phase or particulate pollutants
2 measurement, such as sulfur dioxide (SO₂), nitrogen oxide (NO₂), carbon monoxide (CO), ozone (O₃), total volatile organic
3 compounds (VOCs), and particulate matters (PM). These electrochemical and metal oxide sensors have intrinsic problems such as
4 uneven quality, signal drift, temperature and humidity impacts, and gaseous cross-sensitivities (Spinelle et al., 2015, 2017;
5 Zimmerman et al., 2018) (Brilli et al., 2020; Guo et al., 2020; Jiao et al., 2016; Magi et al., 2020). For example, limited by the poor
6 selection performance, the NO₂ electrochemical sensor also undergo redox reactions in the presence of O₃ gaseous pollutants. The
7 diffusion coefficient of the electrochemical sensor can be affected by temperature and relative humidity(Hitchman et al., 1997;
8 Masson et al., 2015). The reagent of the electrochemical sensor is consumed over time, which affects the stability of the sensor.
9 These features of the sensors have historically been poorly addressed by laboratory calibrations, limiting the utility for air quality
10 monitoring (Zimmerman et al., 2018).

11 The de-convolving of cross-sensitivity effect and stability on sensor performance is complex, and this has been increasing interest
12 in multifarious algorithms for low-cost sensor calibration. The linear or multivariate linear calibration models (Alexopoulos, 2010;
13 Khreis et al., 2022; Zoest et al., 2019), have been developed and studied to overcome these sensors weaknesses. The R^2 performance
14 when measuring NO₂ and PM_{2.5} is less than 0.58(Khreis et al., 2022) with mixed results after the calibration and generally
15 deteriorating performance. The accurate and precise calibration of low cost gas sensors still represents a challenge in ambient air
16 quality monitoring. High-dimensional multi-response calibration models(Alexopoulos, 2010) are built for CO, NO, NO₂, and O₃,
17 with the determination coefficient R^2 within 0.39–0.88 between the models result and the reference monitors results. The artificial
18 neural network (ANN) calibration model has the intelligence to process nonlinear data with self-learning and self-memory(Spinelle
19 et al., 2015), which does not need an accurate mathematical model and has been widely utilized in data analysis to meet
20 computational intelligence requirements(Amuthadevi et al., 2021; Janabi et al., 2021). The ANN has been used in calibration
21 models for measuring ozone or nitrogen oxide(Esposito et al., 2016; Spinelle et al., 2015). For example, the ANN calibration model
22 was used to calibrate O₃ and the uncertainty could meet the European data quality objectives; however, meeting these objectives
23 for NO₂ remains a challenge. Dynamic neural network calibrations of NO₂ sensors was demonstrated with the mean absolute error
24 less than 2 ppb; however, the same performance for O₃ was not observed. Furthermore, the model for calibration performance for
25 measuring NO, NO₂, and O₃ was tested on 4 weeks of data. Random-forest-based machine learning algorithm was also used to
26 improve the calibration strategies of low-cost sensors(Zimmerman et al., 2018). The average mean absolute error on the testing
27 data set from the random forest models was 38 ppb for CO (14% relative error), 10 ppm for CO₂ (2 % relative error), 3.5 ppb for
28 NO₂ (29 % relative error), and 3.4 ppb for O₃ (15 % relative error), and Pearson r versus the reference monitors exceeded 0.8 for
29 most units. Multiple linear regression(Ionascu et al., 2021) based temperature and humidity correction and ANN based calibration
30 shown the potential for significant further improvement for leave one out cross validation(Ali et al., 2021). An integrated genetic
31 programming dynamic neural network model was used to accurately estimate the carbon monoxide and nitrogen dioxide pollutant
32 concentrations from the multi-sensor measurement data(Davut et al., 2022). However, these calibrations have only been tested on
33 a short measurement period and small number of sensor matrix, each containing one sensor per pollutant (Cross et al., 2017;
34 Esposito et al., 2016; Spinelle et al., 2015), not have been utilized to evaluate and predict the concentration values of multi
35 pollutants simultaneously, such as PM_{2.5}, PM₁₀, SO₂, NO₂, CO and O₃.

36 In this work, the LCS is developed to measure PM_{2.5}, PM₁₀, SO₂, NO₂, CO and O₃ simultaneously. Taking the original electronic
37 signals of the sensors as input and measurements obtained by the reference instrumentations as output, five calibration strategies
38 are applied and contrasted. The calibration model used are multivariate linear regression (MLR)(Alexopoulos, 2010), genetic-
39 algorithm-back-propagation neural (GA-BP) network (Ning et al., 2019; Wang et al., 2019), BP neural network(Xu et al., 2021),
40 K Nearest Neighbor (KNN)(Zhao et al., 2021) and random-forest (RF)(Breiman, 2001; Liu et al., 2012) based machine learning



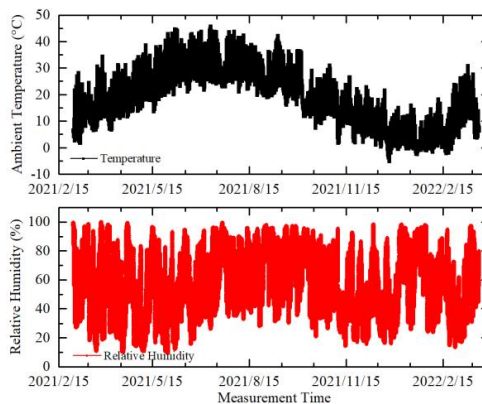
1 algorithm. The measurement is implemented in the real-world conditions almost a 12-month period (1 March 2021 and 28 February
2 2022) spanning multiple seasons and a wide range of meteorological conditions to ensure calibration model robustness. The
3 performance of the different algorithms with the parameters such as R^2 and Root Mean Square Error (RMSE) are compared and
4 discussed.

5 The rest of this paper is organized as follows. The measurement setup is described in section 2. The principles of the calibration
6 strategies are presented in section 3. The results and discussion are shown in section 4. The conclusion is drawn in section 5.

7 **2 Measurement setup**

8 This section describes the measurement site and data collection, schematic block of the LCS, and the reference instrumentation.
9 The low-cost here is defined as below 150 dollars per pollutant, commercial availability and low maintenance. The sensors typically
10 utilize electro-chemical signal and scattering light intensity for gas-phase pollutants (SO_2 , NO_2 , CO and O_3) and particle pollutants
11 ($\text{PM}_{2.5}$, PM_{10}) measurement.

12 **2.1 Measurement site and data collection**

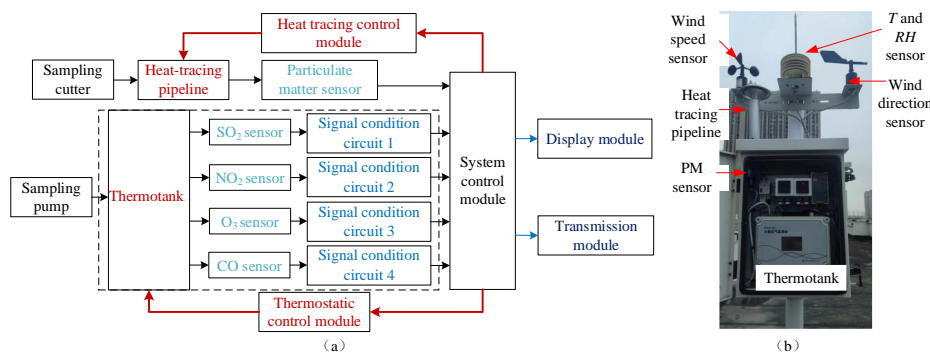


13
14 **Figure 1. Temperature and relative humidity ranges during the measurement period (1 March 2021 and 28 February 2022)**

15
16 Measurements for $\text{PM}_{2.5}$, PM_{10} , CO , SO_2 , NO_2 and O_3 were made continuously between 1 March 2021 and 28 February 2022,
17 which were used as the start and end dates for the analyses. The measurement was made on the 30 Yaochang Street, Zhongyuan
18 District, Zhengzhou City, Henan Province, where there was an independent reference monitoring system for $\text{PM}_{2.5}$, PM_{10} , CO , SO_2 ,
19 NO_2 and O_3 measurement. The LCS was mounted at a consistent height with the reference monitoring system. The data collection
20 interval of the LCS and reference instruments was 5 minutes. During the measurement period, the ranges of the ambient
21 temperature and relative humidity separately were -5°C to 50°C and 10% to 98%, shown in Figure 1. The ambient temperature
22 increased, decreased and fluctuated separately within 1 March 2021 and 30 June 2021, 1 July 2021 and 31 October 2021, 1
23 November 2021 and 28 February 2022, dividing the whole measurement period into three segments.



1 2.2 Schematic block of LCS



2

3 **Figure 2. Schematic block and site photo of the LCS. The left panel (a) is the schematic block of the LCS. The system control module can**
4 **ensure the temperature stability of the heat tracing pipeline and thermo-tank through the heat tracing control module and thermostatic**
5 **control module, respectively. The sampling cutter is used to filter particles larger than 10 μm . The sampling pump is utilized to deliver**
6 **ambient air to the surface of the sensors. The right panel (b) is the site photo of the LCS.**

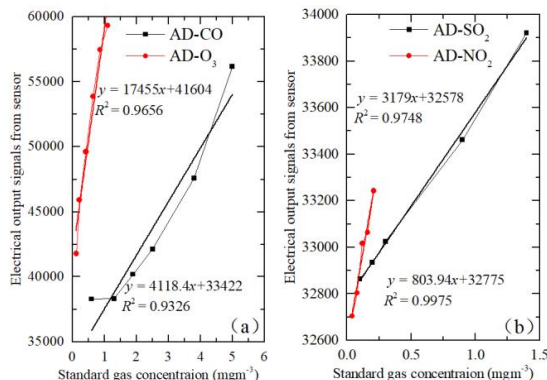
7

8 In this study, the LCS is developed by Hanwei Electronics Group Corporation, and its schematic block diagram is shown in Figure
9 2. The LCS uses the commercially available particulate matter sensor and electrochemical sensors from Cubic Ltd and Alphasense
10 Ltd, respectively. The particulate matter sensor device is the laser diode (LD) based particle sensor, using a spectrophotometer to
11 measure the particle scattering light intensity. The PM sensor device (PM3006) can measure size dependent $\text{PM}_{2.5}$ and PM_{10}
12 concentration of the particles in the size range of 0.3 to 10 μm . The gas pollution (SO_2 \|\ NO_2 \|\ O_3 \|\ CO) sensor used are with 4
13 electrodes (i.e. reference, worker, counter and auxiliary electrodes), where the auxiliary electrode is not exposed to the target
14 analyte to account for changes in the sensor baseline signal under different meteorological conditions (Mead et al., 2013).

15 The electrochemical sensor outputs are measured using electronic circuitry designed by Hanwei and optimized for signal stability.
16 The circuitry is developed with custom electronics to drive the device, multiple stages of filtering circuitry for specific noise
17 signatures, and an analog-to-digital converter for measurement of the conditioned signal.

18 Due to the redox reaction on the anode and the cathode of the electrochemical sensor, the movement of charge between the
19 electrodes produces a current proportional to the analyte reaction rate, which can be used to determine the analyte
20 concentration (Mead et al., 2013) and the sensor whether working effectively. The linearity of the gas sensors of SO_2 , NO_2 , CO and
21 O_3 is examined in laboratory to evaluate the performance of the sensors before used in real-world conditions.

22 In laboratory, the sensors were tested under steadily increased different single standard gas concentration, which was from 0 - 5
23 mgm^{-3} for CO sensor, 0 - 0.2 mgm^{-3} for NO_2 , 0 - 1.1 mgm^{-3} for O_3 and 0 - 1.4 mgm^{-3} for SO_2 . The electrical output signals of the
24 gas pollution sensors, proportional to the concentrations of the single standard substances, shown in Figure 3, verified the sensor
25 working properly and effectively and could be applied to the LCS.



1 **Figure 3. Electrical output signals versus single standard gas concentration in laboratory condition. The left panel (a) and right panel (b)**
2 **represent the proportional relations between CO and O₃ sensors, SO₂ and NO₂ sensors, respectively. The duration of each measurement**
3 **is about 30 minutes.**
4

5
6 However, even with an auxiliary electrode, electrochemical sensors may insufficiently account for the impacts of temperature and
7 relative humidity. With the purpose to eliminate the influence of the external environment on the sensor as much as possible, the
8 particles flow through a sampling cutter and heat-tracing pipeline to the particulate matter sensor, and the gaseous pollutants are
9 pumped to the electrochemical sensors, which are secured in a thermo-tank. The temperature values of the heat-tracing pipeline
10 and thermo-tank can be maintained at $60\text{ }^{\circ}\text{C} \pm 2\text{ }^{\circ}\text{C}$ to reduce the influence of relative humidity and $25\text{ }^{\circ}\text{C} \pm 2\text{ }^{\circ}\text{C}$ (Wei et al., 2018)
11 to keep the sensor operating at a stable temperature, respectively.

12 The measurement results of particulate matter sensor and gas pollution sensors, transmitted to the system control module through
13 the data buses, are directly displayed on the local display module and wirelessly transmitted to the corresponding online server
14 through the transmission module. As the uni-variate linear models does not in-corporate any cross-sensitivities to other pollutants
15 or any nonlinearities in the response, we attempt to using the sensor electronic results as the input and the reference measurements
16 as the output, to build multi-dimensional multi-response prediction models to de-convolve the effects of cross-sensitivity and
17 stability on sensor performance utilizing MLR, RF, KNN, BP and GA-BP calibration models.

18 2.3 Reference instrumentation

19 In order to reduce the adsorption effect on particle matter and gaseous pollutants, the reference measurements are made on ambient
20 air continuously drawn through Teflon fluorinated ethylene propylene (FEP) (Wei et al., 2018) tubing with a six-port stainless steel
21 manifold for flow distribution to the gas analyzers and particulate monitors (Mead et al., 2013). It should be pointed out that the
22 LCS was mounted at a consistent height with the reference monitoring system during the measurement period.

23 The reference ambient particulate monitor 5014i, which uses beta attenuation of the ambient particulate deposited onto a filter tape,
24 is applied to measure the mass concentration of suspended and refined particulates. The reference NO-NO₂-NO_x monitor 42i, using
25 the linear proportional of the chemi-luminescence reaction of NO and O₃ after NO₂ transformed into NO, is utilized to measure the
26 NO₂ concentration. The SO₂ reference analyzer is 43i using the ultraviolet light (which is emitted as the excited SO₂ molecules
27 decay to lower energy states) intensity proportional to the SO₂ concentration. The CO reference monitor is 48i utilizing the principle
28 that CO absorbs infrared radiation at a wavelength of 4.6 μm and the infrared absorption can be transformed to be proportional to
29 the CO concentration. The 49i O₃ analyzer operates on the principle that O₃ molecules absorb UV light at a wavelength of 254 nm,
30 and the absorption intensity of the UV light is directly related to the ozone concentration. All these reference monitors are produced



1 by Thermo Fisher Scientific Inc. The time interval for all reference measurements is 5 minutes. The reference gas and particulate
 2 analyzers are checked and calibrated weekly using calibration gas mixtures to correct the baseline drift.

3 Principles

4 This section describes the principles of the calibration methods, such as MLR, BP, GA-BP, KNN and RF, and the metrics for
 5 performance evaluation. The calibration models are constructed with the sensors (i.e., PM_{2.5}, PM₁₀, CO, SO₂, NO₂ and O₃ sensors.)
 6 electronic results as the input and the reference measurements as the output.

3.1 Calibration methods

3.1.1 Multiple linear regression model

9 After the data collected by the LCS, the raw data should be preprocessed. The PM3006 particulate matter sensor can measure
 10 particle counters in the size range of 0.3 to 10 μm. By using the particle counters $x_{0.5}$, $x_{1.0}$, $x_{2.5}$, $x_{5.0}$ and $x_{10.0}$ of the sensors, listed in
 11 Table 1, the measured particle number concentration is converted to PM mass concentrations in the PM_{2.5} and PM₁₀ size fractions.

12
 13 **Table 1. Size range of the particulate matter sensor. The sensor can measure particles with the size range of 0.3–0.5 μm, 0.5–1.0 μm,**
 14 **1.0–2.5 μm, 2.5–5.0 μm and 5.0–10 μm, simultaneously. The corresponding particle counters are expressed as $x_{0.5}$, $x_{1.0}$, $x_{2.5}$, $x_{5.0}$ and $x_{10.0}$,**
 15 **respectively.**

Range (μm)	0.3–0.5	0.5–1.0	1.0–2.5	2.5–5.0	5.0–10.0
Particle counter	$x_{0.5}$	$x_{1.0}$	$x_{2.5}$	$x_{5.0}$	$x_{10.0}$

16
 17 Taking the particle counters, listed in Table 1, as input and the concentrations $Y_{pm2.5}$ and Y_{pm10} of PM_{2.5} and PM₁₀ measured by
 18 5014i as output, the multivariate linear regression (MLR) models (Alexopoulos, 2010; Zoest et al., 2019) is built. Due to the
 19 previously established influence of ambient temperature (T) and relative humidity (RH) on sensor response (Jiao et al., 2016;
 20 Masson et al., 2015), the multi-dimensional multi-response preprocessing and prediction models can be written as Eq. (1).

$$21 \begin{cases} Y_{pm2.5} = w_{1_pm2.5} \cdot x_{0.5} + w_{2_pm2.5} \cdot x_{1.0} + w_{3_pm2.5} \cdot x_{2.5} + w_{4_pm2.5} \cdot T + w_{5_pm2.5} \cdot RH + b_{pm2.5} \\ Y_{pm10} = w_{1_pm10} \cdot x_{0.5} + w_{2_pm10} \cdot x_{1.0} + w_{3_pm10} \cdot x_{2.5} + w_{4_pm10} \cdot x_{5.0} + w_{5_pm10} \cdot x_{10.0} + w_{6_pm10} \cdot T + w_{7_pm10} \cdot RH + b_{pm10} \end{cases} \quad (1)$$

22 The equation (1) can be simplified as,

$$23 \begin{cases} Y_{pm2.5} = W_{pm2.5} \cdot X_{pm2.5} + b_{pm2.5} \\ Y_{pm10} = W_{pm10} \cdot X_{pm10} + b_{pm10} \end{cases} \quad (2)$$

24 Where $W_{pm2.5} = [w_{1_pm2.5}, w_{2_pm2.5}, w_{3_pm2.5}, w_{4_pm2.5}, w_{5_pm2.5}]$ and $W_{pm10} = [w_{1_pm10}, w_{2_pm10}, w_{3_pm10}, w_{4_pm10}, w_{5_pm10}, w_{6_pm10}, w_{7_pm10}]$
 25 are the corresponding weight coefficients; the $X'_{pm2.5} = [x_{0.5}, x_{1.0}, x_{2.5}, T, RH]$ and $X'_{pm10} = [x_{0.5}, x_{1.0}, x_{2.5}, x_{5.0}, x_{10.0}, T, RH]$ represent
 26 the input particle counters and values obtained from the PM sensor, the temperature sensor and humidity sensor; the $b_{pm2.5}$ and
 27 b_{pm10} are the intercept values of the model.

28 Due to the poor selection performance and cross interference of the electro-chemical sensors response, the output values from each
 29 sensor using net sensor response to the target analyte, such as O₃, NO₂, SO₂, concentration measured by the inference monitor are
 30 used to build the MLR model. The CO gaseous pollution is also one of the criteria pollutants, which is must to be measured in
 31 China. Thus, the multi-dimensional multi-response preprocessing and prediction model for the 4 gas pollutions, T and RH can be
 32 written as Eq. (3).



$$1 \quad \begin{cases} Y_{\text{SO}_2} = w_{11} \cdot x_{\text{SO}_2} + w_{12} \cdot x_{\text{NO}_2} + w_{13} \cdot x_{\text{CO}} + w_{14} \cdot x_{\text{O}_3} + w_{15} \cdot T + w_{16} \cdot RH + b_{\text{SO}_2} \\ Y_{\text{NO}_2} = w_{21} \cdot x_{\text{SO}_2} + w_{22} \cdot x_{\text{NO}_2} + w_{23} \cdot x_{\text{CO}} + w_{24} \cdot x_{\text{O}_3} + w_{25} \cdot T + w_{26} \cdot RH + b_{\text{NO}_2} \\ Y_{\text{CO}} = w_{31} \cdot x_{\text{SO}_2} + w_{32} \cdot x_{\text{NO}_2} + w_{33} \cdot x_{\text{CO}} + w_{34} \cdot x_{\text{O}_3} + w_{35} \cdot T + w_{36} \cdot RH + b_{\text{CO}} \\ Y_{\text{O}_3} = w_{41} \cdot x_{\text{SO}_2} + w_{42} \cdot x_{\text{NO}_2} + w_{43} \cdot x_{\text{CO}} + w_{44} \cdot x_{\text{O}_3} + w_{45} \cdot T + w_{46} \cdot RH + b_{\text{O}_3} \end{cases} \quad (3)$$

2 The equation (3) can be simplified as,

$$3 \quad Y_{[\text{SO}_2, \text{NO}_2, \text{CO}, \text{O}_3]} = W_{\text{gas}} \cdot X_{\text{gas}} + B_{\text{gas}}, \quad (4)$$

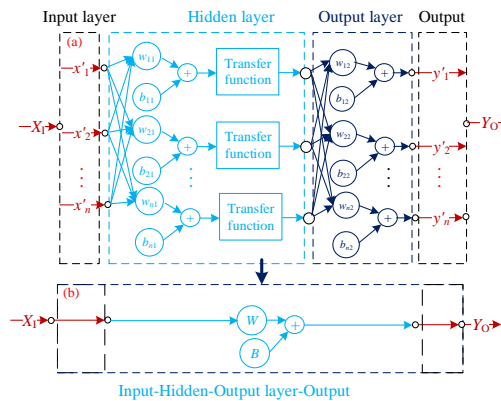
4 Where $W_{\text{gas}} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} & w_{15} & w_{16} \\ w_{21} & w_{22} & w_{23} & w_{24} & w_{25} & w_{26} \\ w_{31} & w_{32} & w_{33} & w_{34} & w_{35} & w_{36} \\ w_{41} & w_{42} & w_{43} & w_{44} & w_{45} & w_{46} \end{bmatrix}$ is the corresponding weight coefficient; the $X_{\text{gas}} = [x_{\text{SO}_2}, x_{\text{NO}_2}, x_{\text{CO}}, x_{\text{O}_3}, T,$

5 $RH]$ is the convertor output values of the sensors through the electronic circuitries; the $B_{\text{gas}}' = [b_{\text{SO}_2}, b_{\text{NO}_2}, b_{\text{CO}}, b_{\text{O}_3}]$ is the
 6 intercept value of the model.

7 Hereto, the multi MLR models for the gas sensor and PM sensor are separately developed. The training data is used to calculate
 8 the model regression coefficient and intercept values, and the withheld testing data is utilized to evaluate the performance of the
 9 model performance.

10 3.1.2 BP neural network model

11 The BP neural network algorithm is one of the most widely used ANN models. It is a multi-layer feed-forward network trained
 12 through an error back propagation algorithm by constantly adjusting the weight and intercept of the network. The feed-forward
 13 topological structure of the BP neural network model, shown in Figure 4, includes the input layer, hidden layer and output layer.
 14 With the purpose to avoid the numerical problems caused by the extreme values of polarization, eliminate the misleading effects
 15 for feature extraction and obtain the accurate estimation of pollutant concentrations (Janabi et al., 2021), the collected input sensor
 16 date X_I and output date Y_O should be respectively normalized with min-max normalization to limit values in each dimension
 17 between 0 and 1 (Hande et al., 2021).



18 **Figure 4. Topological structure of BP neural network model.** The up panel (a) is the feed-forward topological structure. The X_I and Y_O
 19 are the input data and output data, respectively. The X'_i and Y'_i separately indicate the normalized items of X and Y . The w_{j1} and b_{j1}
 20 and b_{j2} separately represent the weight value and intercept value of the hidden layer and output layer. The down panel (b) is equivalent
 21 to panel (a) to simplify the formulas.
 22

23
 24 After the normalization process, the BP network can be established. To optimize the best parameters of the network, the number
 25 of hidden layer, the transfer functions of the layers, and the end conditions should be determined. If the parameters are inappropriate,
 26 the BP-model will be over trained or insufficient. In this study, a shallow structure with a single hidden layer is chosen, as extensive



1 testing did not show any noticeable improvement in calibration performance with deeper structure consisting of multiple hidden
2 layers(Ali et al., 2021). This also reduced the complexity and the training time.

3 **3.1.3 Genetic algorithm-BP model**

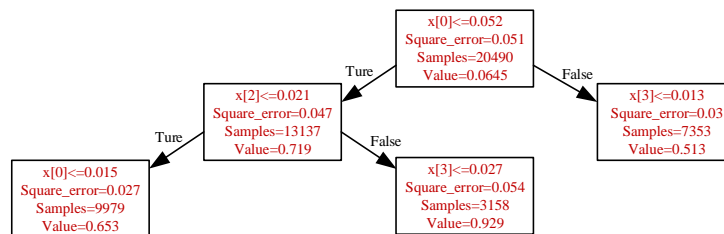
4 In the traditional BP neural network, the initial weights and thresholds are randomly generated. The results often fall into a local
5 minimum rather than a global minimum, and would lead to the distortion of the prediction result. In addition, the convergence
6 speed of the BP neural network is usually slow. To solve these problems, the genetic algorithm (GA) (Liang et al., 2018) with BP
7 algorithm is also used to avoid the inherent defects of BP algorithm. The GA method is essentially a direct search method that does
8 not rely on specific problems and gradient information. It follows the survival and elimination rule of biological evolution,
9 generates the following hypotheses by mutating and reconstructing the best existed hypothesis and makes it possible to solve the
10 problem(Ning et al., 2019). Generally, the GA is used to find an optimal initial weight and a threshold value for the model, so that
11 the model could converge in the direction of minimum value(Wang et al., 2019). The GA-BP hybrid algorithm is used to reduce
12 the time for the BP neural network to adjust the weight and threshold itself and achieve the goal of improving work efficiency.

13 **3.1.4 K nearest neighbor model**

14 The k nearest neighbor (KNN) is also one of the simplest method for classification as well as regression problem(Kumar, 2015;
15 Zhao and Lai, 2021). The KNN is a supervised method that uses estimation based on values of neighbors, which can automatically
16 adapt to the supervised learning problems with arbitrary Bayes decision boundaries(Zhao and Lai, 2021). From the supervisor
17 dataset, the KNN solution utilizes the values of given dependent variable y_i to approximate the dependent variable y^* , which is
18 closest with respect to distance between their corresponding model parameters. For regression problem, the mean of the observed
19 labels of k nearest neighbors of independent variable X is assigned to be the predicted label. In this study, the k is set to 10 with
20 the performance having no obvious difference from other numbers.

21 **3.1.5 Random forest model**

22 The random forest (RF) model is used for solving regression or classification problems(Breiman, 2001; Liu et al., 2012). It works
23 by constructing an ensemble of decision trees using a training data set; the mean value from that ensemble of decision trees is then
24 used to predict the value for new input data(Zimmerman et al., 2018). With the purpose to establish a RF model, the maximum
25 number of decision trees of the forest should be specified. Each tree is constructed using a bootstrapped random sample from the
26 training data set. By considering a random subset of the possible explanatory variables with the strongest predictor of the response,
27 the origin node of the decision tree can be split into sub-nodes. The node splitting process is repeated until a terminal node is
28 reached. The terminal node can be specified using the maximum number of sub-nodes or the minimum number of data points in
29 the node. To illustrate the method, consider building a random forest model for one LCS using a single decision tree and a subset
30 of 20490 data points to build a calibration model, shown in Figure 5. The RF model can predict data with variable parameters
31 within the training range. Therefore, a larger and more variable training data set should create a better final model. To avoid missing
32 any spikes during the training window, a 5-fold cross-validation approach(Zimmerman et al., 2018) is also used to maximize
33 utilization of the training data set. This approach helps to minimize bias in training data selection when predicting new data and
34 ensures that every point in the training window is used to build the model.



1
 2 **Figure 5. Simplified illustration of the RF with a single decision tree and a subset. The $x[0]$, $x[2]$, $x[3]$ represent the CO, SO₂ and O₃**
 3 **pollutants. At the first split, points with normalized CO sensor signal ≤ 0.052 are sent to a terminal node; the remaining points go to the**
 4 **other splitting node. The Samples is the number of data points in each terminal node. The Value is the average in each terminal node.**

5 3.2 Metrics for performance evaluation

6 To quantitatively compare the MLR, BP, GA-BP, KNN and RF model output to the reference monitor concentrations, the
 7 determination coefficient R^2 and root mean square error (RMSE) (Janabi et al., 2021) are utilized. Where, the R^2 is obtained as
 8 Eq. (5).

$$9 R^2 = 1 - \frac{\sum_{i=1}^n (y_i - f_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (5)$$

10 Where y_i and \bar{y} are the real-time data and mean data obtained by reference instrumentations, respectively. The f_i is the model
 11 output data according to the model algorithm. The R^2 reflects the fit degree between the model output data and the reference
 12 monitor measurement. The measurement results should meet the requirements of environmental standards of China (Jiao et al.,
 13 2016). The RMSE measures how much error there is between the predicted values and the reference measurements, which is
 14 calculated by the following Eq. (6).

$$15 RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_{i,p} - y_i)^2}, \quad (6)$$

16 Where $y_{i,p}$ and y_i represent the i th model output data from the algorithm-based LCS system and the reference data from the
 17 reference instrumentations, the n is the number of the measurement data in the dataset.

18 4 Result and discussion

19 Following the model building, the goodness of regression and root mean square error between the model output concentrations
 20 and the reference monitor concentrations are evaluated for all calibration model approaches. The plots for the PM_{2.5}, PM₁₀, O₃, CO,
 21 NO₂ and SO₂ illustrating the time series and goodness of fit of the models are provided in the Figure 9 - Figure 14. The R^2 and
 22 RMSE values are listed in Table 2 – Table 5.

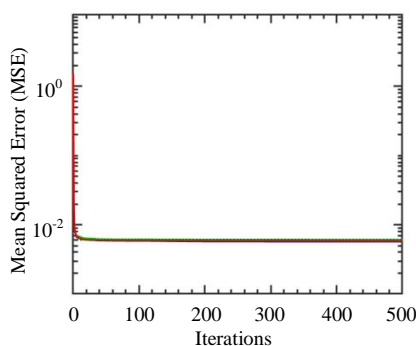
23 4.1 Parameters of the model

24 For the BP and GA-BP models, the parameters are the functions for the hidden layer and output layer, the type of the hidden layer,
 25 the number of iteration times, and the number of the nerve units (Wang et al., 2013). The functions for the hidden layer and output
 26 layer in this study respectively are the default tansig and the purelin functions. With the more complex type of the hidden layer
 27 and less obvious improvement, the hidden layer is single type to achieve the goal of work efficiency.

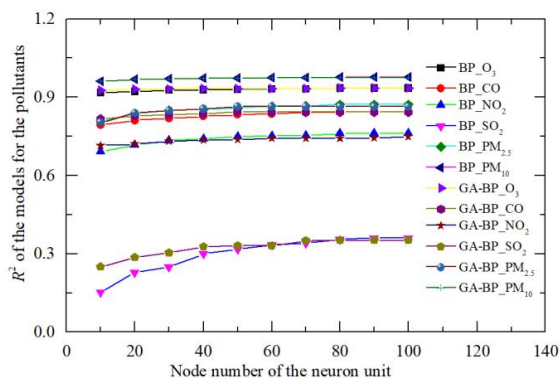
28 To determine the best number of iteration times and nerve units, the measurement from the LCS and reference monitor between 1
 29 March 2021 and 30 June 2021 is used. The number of iteration time is optimized using the mean squared error (MSE) between the
 30 model value from the model and the reference monitor output value. The tendency of the MSE is shown in Figure 6. The training
 31 is performed for 500 iterations. It is observed that the MSE decreases with the number of iteration time increasing, the rate of



1 decrease and the variation of the MSE is negligible beyond 100 iterations. More iterations incur higher computational cost for the
2 training and small performance improvement. There is also the risk of overtraining resulting in poor generalization capability.



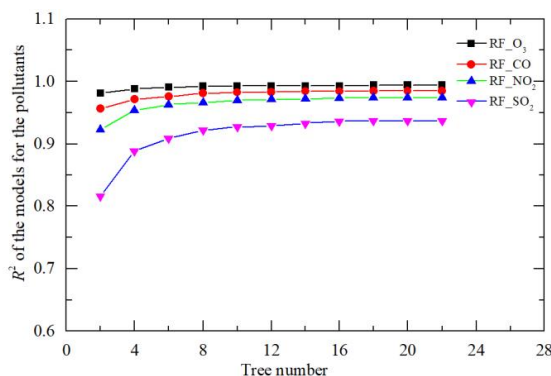
3
4 **Figure 6. The MSE with the number of iterations.**



5
6 **Figure 7. The R^2 with different node number of the neuron for the pollutants.**

7 The node number of the nerve units is determined by the contrast results of determination coefficient R^2 for different gas and PM
8 pollutants within 1 March 2021 and 30 June 2021. The results are shown in Figure 7. The R^2 is improved as the number of nerve
9 units increasing. The rate of increase and the variation of R^2 is negligible beyond 70 units. More units incur higher computational
10 cost and time for the training and small performance improvement.

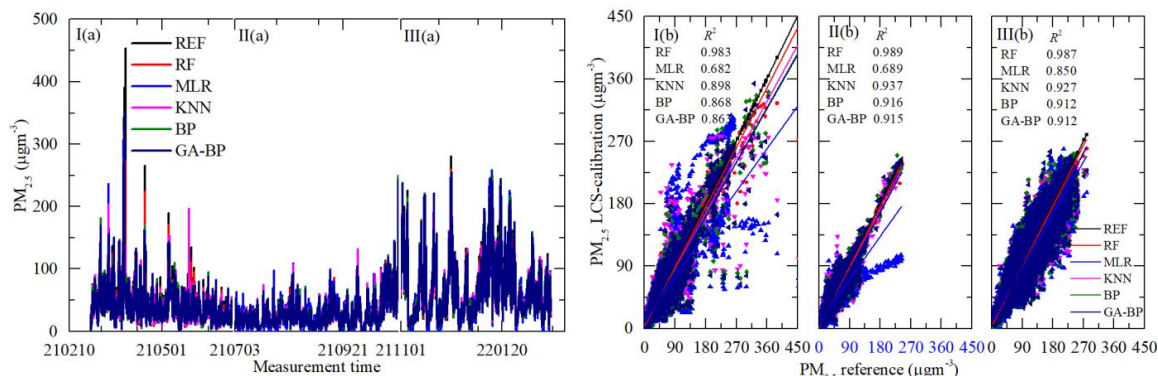
11 For the RF model, the number of tree is determined by the contrast results of determination coefficient R^2 for different gas pollutants
12 within 1 March 2021 and 30 June 2021. The results are shown in Figure 8. The R^2 is improved as the number of tree increasing.
13 The rate of increase and the variation of R^2 is negligible beyond 20. The terminal node is specified using a maximum number of
14 sub-node points per node. The R^2 is also improved as the number of sub-node increasing under the same tree number. The rate of
15 increase and the variation of R^2 is negligible beyond 100. More number of the tree or the sub-node incur higher computational cost
16 and time for the training and small performance improvement.



1
 2 **Figure 8. The R^2 with different tree number of the RF model for the pollutants.**

3 **4.2 Measurement results of PM**

4 With the results from 1 March 2021 to 28 February 2022 and according the trend of the ambient temperature, shown in Figure 1,
 5 the whole data is divided into three segments. The three segments (I, II, and III) separately are within 1 March 2021 and 30 June
 6 2021, 1 July 2021 and 31 October 2021, 1 November 2021 and 28 February 2022. The time series data and regressions of five
 7 modes for PM from reference monitor and LCS calibration output are shown in Figure 9 and Figure 10. With the purpose of
 8 avoiding over-fit in the five models, the randomly divide parameters of train ratio and test ratio are 80% and 20%, respectively.



9
 10 **Figure 9. Time series and regressions comparing the reference monitor $PM_{2.5}$ data (black) to five calibration model $PM_{2.5}$ results. Where**
 11 **red, blue, magenta, olive and navy represent RF, MLR, KNN, BP, GA-BP, respectively. The left panel (a) shows the whole time series**
 12 **data of the measurement period. The right panel (b) shows the regressions of the five calibration models.**

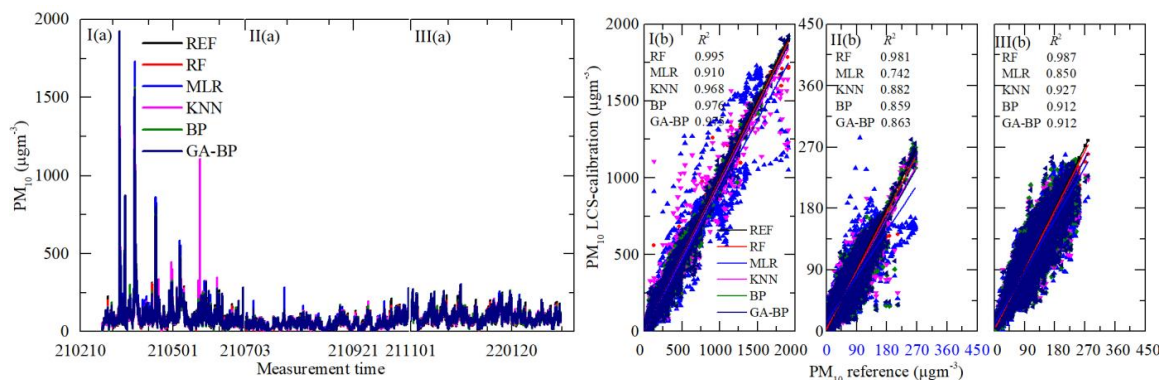
13 As shown in Figure 9 (a) and Figure 10 (a), the general tendency of the data fluctuation between the reference monitor and the RF,
 14 MLR, KNN, BP, GA-BP based algorithm of the LCS are consistent with each other. The best performance is RF model, the next
 15 are KNN, BP and GA-BP, the worst is MLR. The regressions consistence between the reference data and the five model data are
 16 also shown in Figure 9 (b) - Figure 10 (b), and listed in Table 2.

17 The R^2 of RF for the PM is better than 0.98. The R^2 of MLR for the PM is less than 0.91, and even less than 0.7. The R^2 of the
 18 other three model are within 0.86 and 0.98. The performance of different calibration models for the PM against reference monitor
 19 is also evaluated using RMSE, and the results are listed in Table 3.

20 Using the data listed in Table 3, the RMSE values from the first (I) and third (III) stages are large than the one from the second (II)
 21 stage, the main reason maybe the large fluctuation range of the PM for the climatic factors in winter and spring resulting in the
 22 poor model fit. The mean RMSE values of $PM_{2.5}$ between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms
 23



1 data are calculated as $3.96 \mu\text{g m}^{-3}$, $16.16 \mu\text{g m}^{-3}$, $9.48 \mu\text{g m}^{-3}$, $10.67 \mu\text{g m}^{-3}$ and $10.74 \mu\text{g m}^{-3}$, respectively. The mean RMSE values
 2 of PM_{10} between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms data are calculated as $7.37 \mu\text{g m}^{-3}$,
 3 $28.90 \mu\text{g m}^{-3}$, $18.50 \mu\text{g m}^{-3}$, $18.04 \mu\text{g m}^{-3}$ and $18.17 \mu\text{g m}^{-3}$, respectively.



4 **Figure 10.** Time series and regressions comparing the reference monitor PM_{10} data (black) to five calibration model PM_{10} results. Where
 5 red, blue, magenta, olive and navy represent RF, MLR, KNN, BP, GA-BP, respectively. The left panel (a) shows the whole time series
 6 data of the measurement period. The right panel (b) shows the regressions of the five calibration models.
 7

8 **Table 2.** Performance of different calibration models for the $\text{PM}_{2.5}$ and PM_{10} against reference monitor. The determination coefficient R^2
 9 (higher is better, maximum of 1) of different calibration models (RF, MLR, KNN, BP, GA-BP) versus reference monitor.
 10

R^2 Model	$\text{PM}_{2.5}$			PM_{10}		
	I	II	III	I	II	III
RF	0.983	0.989	0.987	0.995	0.981	0.987
MLR	0.682	0.689	0.850	0.910	0.742	0.850
KNN	0.898	0.937	0.927	0.968	0.882	0.927
BP	0.868	0.916	0.912	0.976	0.859	0.912
GA-BP	0.863	0.915	0.912	0.975	0.863	0.912

11 **Table 3.** Performance of different calibration models for the $\text{PM}_{2.5}$ and PM_{10} against reference monitor. The RMSE errors (lower is better)
 12 of different calibration models (RF, MLR, KNN, BP, GA-BP) versus reference monitor.
 13

RMSE Model	$\text{PM}_{2.5} (\mu\text{g m}^{-3})$				$\text{PM}_{10} (\mu\text{g m}^{-3})$			
	I	II	III	MEAN	I	II	III	MEAN
RF	4.03	2.36	5.49	3.96	10.37	4.55	7.19	7.37
MLR	17.18	12.63	18.68	16.16	45.05	16.43	25.22	28.90
KNN	9.73	5.67	13.05	9.48	27.08	11.14	17.29	18.50
BP	11.09	6.56	14.35	10.67	23.10	12.15	18.88	18.04
GA-BP	11.27	6.61	14.35	10.74	23.65	11.99	18.87	18.17

15 4.3 Measurement results of gas pollution

16 With the results from 1 March 2021 to 28 February 2022 and according the trend of the ambient temperature, shown in Figure 1,
 17 the whole data is also divided into three same segments as section 4.2. The time series data and regressions of five modes for gas
 18 pollution from reference monitor and LCS calibration output are shown in Figure 11 - Figure 14. With the same purpose of avoiding
 19 over-fit in the five models, the randomly divide parameters of train ratio and test ratio are also 80% and 20%, respectively.

20 As shown in Figure 11 (a) - Figure 14 (a), the general tendency of the data fluctuation between the reference monitor and the RF,
 21 MLR, KNN, BP, GA-BP based algorithm of the LCS are consistent with each other. The best performance is RF model, the next
 22 are KNN, BP and GA-BP, the worst is MLR. The regressions consistence between the reference data and the five model data are
 23 also shown in Figure 11 (b) - Figure 14 (b), and listed in Table 4. The RMSE values between the reference data and the five model
 24 data are listed in Table 5.



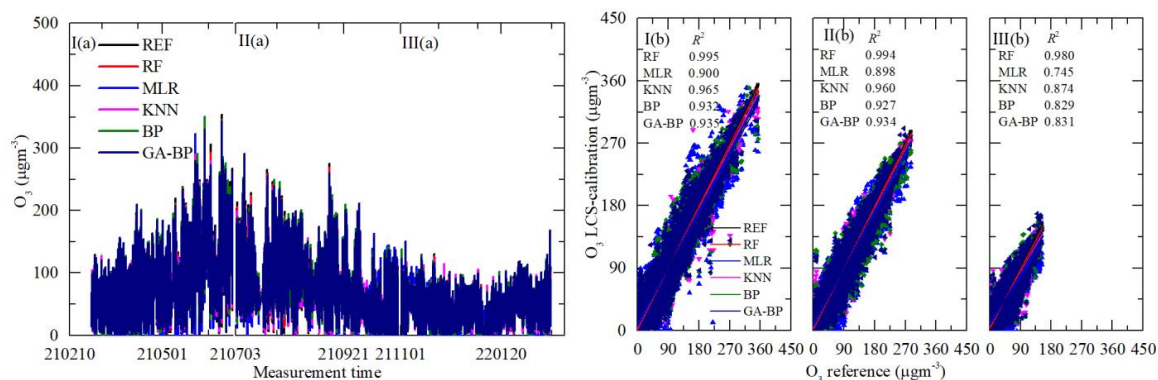
1

2 **Table 4. Performance of different calibration models for the gaseous pollutant (SO₂, CO, NO₂, O₃) against reference monitor. The**
 3 **determination coefficient R^2 (higher is better, maximum of 1) of different calibration models (RF, MLR, KNN, BP, GA-BP) versus**
 4 **reference monitor.**

Model	O ₃			CO			NO ₂			SO ₂		
	I	I	I	II	III	I	II	III	I	II	III	
RF	0.995	0.994	0.980	0.989	0.978	0.981	0.981	0.967	0.962	0.969	0.939	
MLR	0.900	0.898	0.745	0.729	0.807	0.710	0.456	0.530	0.570	0.065	0.333	
KNN	0.965	0.960	0.874	0.921	0.934	0.861	0.866	0.878	0.786	0.686	0.797	
BP	0.932	0.927	0.829	0.837	0.858	0.815	0.756	0.775	0.716	0.332	0.609	
GA-BP	0.935	0.934	0.831	0.841	0.871	0.816	0.742	0.782	0.708	0.341	0.622	

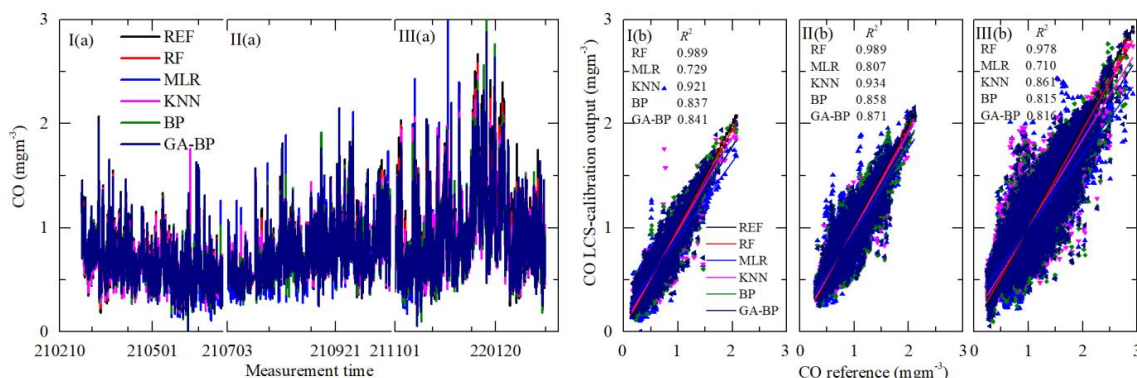
5

6 For the O₃ model, the R^2 of RF is better than 0.98. The R^2 of MLR is less than 0.90, and even less than 0.8. The R^2 of the other
 7 three models are within 0.82 and 0.97. The performance of different calibration models for the O₃ against reference monitor is also
 8 evaluated using RMSE, and the results are listed in Table 5. Using the data listed in Table 5, the RMSE values from the first (I) and
 9 third (III) stages have little difference with the one from the second (II) stage, indicating the O₃ electrochemical sensor suitable for
 10 the ambient ozone measurement. The mean RMSE values of O₃ between the reference data and the RF, MLR, KNN, BP, GA-BP-
 11 based algorithms data are calculated as 4.06 $\mu\text{g}\text{m}^{-3}$, 16.07 $\mu\text{g}\text{m}^{-3}$, 10.23 $\mu\text{g}\text{m}^{-3}$, 13.35 $\mu\text{g}\text{m}^{-3}$ and 13.0 $\mu\text{g}\text{m}^{-3}$, respectively.



12

13 **Figure 11. Time series and regressions comparing the reference monitor O₃ data (black) to five calibration model O₃ results. Where red,**
 14 **blue, magenta, olive and navy represent RF, MLR, KNN, BP, GA-BP, respectively. The left panel (a) shows the whole time series data of the**
 15 **measurement period. The right panel (b) shows the regressions of the five calibration models.**



16

17 **Figure 12. Time series and regressions comparing the reference monitor CO data (black) to five calibration model CO results. Where red,**
 18 **blue, magenta, olive and navy represent RF, MLR, KNN, BP, GA-BP, respectively. The left panel (a) shows the whole time series**
 19 **data of the measurement period. The right panel (b) shows the regressions of the five calibration models.**

20

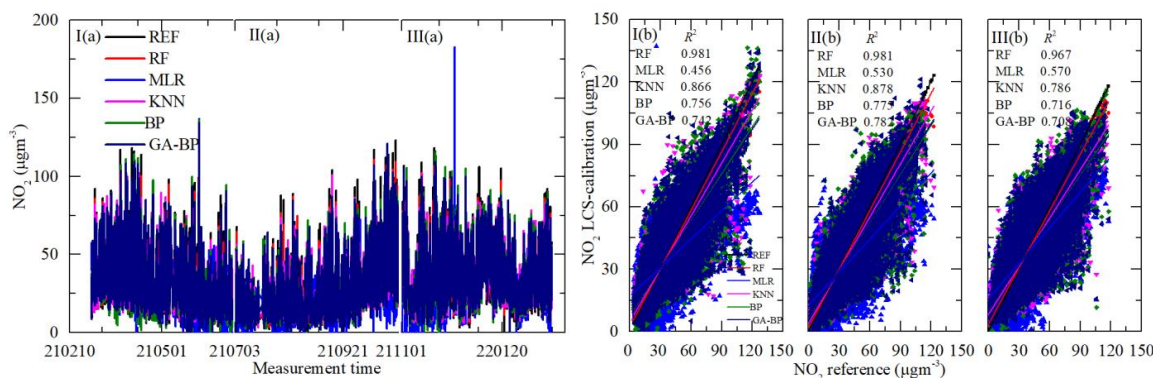


1 For the CO model, the R^2 of RF is better than 0.97. The R^2 of MLR is less than 0.81, and even less than 0.7. The R^2 of the other
 2 three model are within 0.81 and 0.94. The performance of different calibration models for the CO against reference monitor is also
 3 evaluated using RMSE, and the results are listed in Table 5. The RMSE values from the first (I) and third (III) stages have little
 4 difference with the one from the second (II) stage, indicating the CO electrochemical sensor also suitable for the ambient CO
 5 measurement. The mean RMSE values of CO between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms
 6 data are calculated as 0.04 mgm^{-3} , 0.15 mgm^{-3} , 0.10 mgm^{-3} , 0.12 mgm^{-3} and 0.12 mgm^{-3} , respectively.

7
 8 **Table 5. Performance of different calibration models for the gaseous pollutant (SO_2 , CO, NO_2 and O_3) against reference monitor. The**
 9 **RMSE errors (lower is better) of different calibration models (RF, MLR, KNN, BP, GA-BP) versus reference monitor.**

RMSE Model	O_3 (μgm^{-3})				CO (mgm^{-3})				NO_2 (μgm^{-3})				SO_2 (μgm^{-3})			
	I	II	III	MEAN	I	II	III	MEAN	I	II	III	MEAN	I	II	III	MEAN
RF	4.05	4.06	4.08	4.06	0.02	0.03	0.06	0.04	2.88	2.88	3.99	3.25	0.83	0.64	1.68	1.05
MLR	17.79	16.42	14.00	16.07	0.12	0.12	0.23	0.15	14.54	13.54	13.61	13.90	3.53	2.69	5.37	3.86
KNN	10.57	10.28	9.84	10.23	0.06	0.07	0.16	0.10	7.25	6.93	9.61	7.93	2.06	1.49	4.05	2.53
BP	14.67	13.91	11.46	13.35	0.09	0.10	0.18	0.12	9.75	9.37	11.07	10.06	2.98	2.06	4.63	3.22
GA-BP	14.40	13.19	11.41	13.00	0.09	0.10	0.18	0.12	10.02	9.21	11.21	10.15	2.97	2.03	4.60	3.20

10
 11 For the NO_2 model, the R^2 of RF is better than 0.96. The R^2 of MLR is less than 0.60, and even less than 0.5. The R^2 of the other
 12 three model are within 0.70 and 0.90. The performance of different calibration models for the NO_2 against reference monitor is
 13 also evaluated using RMSE, and the results are listed in Table 5. The RMSE values from the first (I) and third (III) stages have little
 14 difference with the one from the second (II) stage, indicating the NO_2 electrochemical sensor still suitable for the ambient NO_2
 15 measurement. The mean RMSE values of NO_2 between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms
 16 data are calculated as 3.25 μgm^{-3} , 13.90 μgm^{-3} , 7.93 μgm^{-3} , 10.06 μgm^{-3} and 10.15 μgm^{-3} , respectively.



17
 18 **Figure 13. Time series and regressions comparing the reference monitor NO_2 data (black) to five calibration model NO_2 results. Where**
 19 **red, blue, magenta, olive and navy represent RF, MLR, KNN, BP, GA-BP, respectively. The left panel (a) shows the whole time series**
 20 **data of the measurement period. The right panel (b) shows the regressions of the five calibration models.**

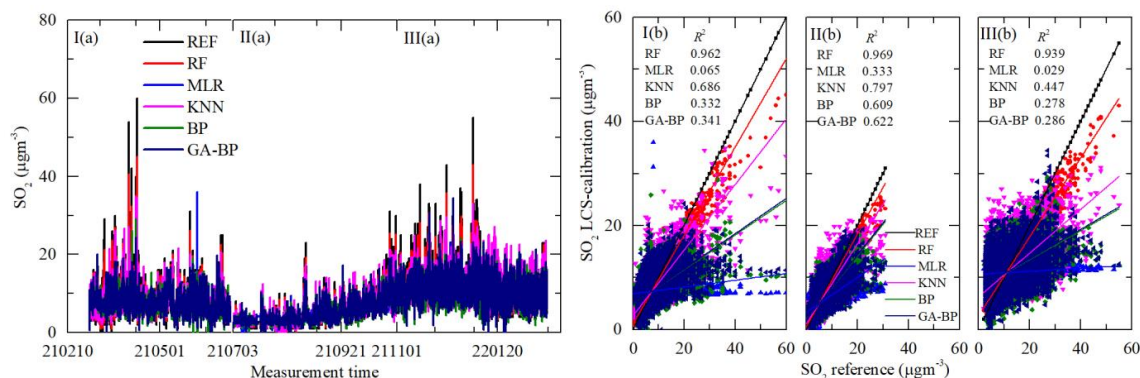


Figure 14. Time series and regressions comparing the reference monitor SO₂ data (black) to five calibration model SO₂ results. Where red, blue, magenta, olive and navy represent RF, MLR, KNN, BP, GA-BP, respectively. The left panel (a) shows the whole time series data of the measurement period. The right panel (b) shows the regressions of the five calibration models.

For the SO₂ model, the R² of RF is better than 0.93. The R² of MLR is less than 0.40, and even less than 0.1. The R² of the other three model are within 0.27 and 0.80. The performance of different calibration models for the SO₂ against reference monitor is also evaluated using RMSE, and the results are listed in Table 5. The RMSE values from the first (I) and third (III) stages have little difference with the one from the second (II) stage, indicating the SO₂ electrochemical sensor with the RF calibration can be used for the ambient SO₂ measurement. The mean RMSE values of SO₂ between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms data are calculated as 1.05 μgm⁻³, 3.86 μgm⁻³, 2.53 μgm⁻³, 3.22 μgm⁻³ and 3.20 μgm⁻³, respectively.

As shown in Figure 11 - Figure 14 and listed in Table 4 – Table 5, the models (except the RF model) perform poorly for SO₂, especially during the spring and winter. There maybe three reasons for this phenomenon. The first one is the cross interference effect from NO₂ and O₃. The second one is the reaction products (SO₃) dissolved in the electrolyte, which can affect the ion concentration. The last one is the stability of electrode material, which is easily affected by external interference.

5 Conclusion

A low-cost air quality monitoring system based on RF, MLR, KNN, BP, GA-BP algorithms are proposed. The system can measure PM_{2.5}, PM₁₀, SO₂, NO₂, CO and O₃, simultaneously. The PM prediction model is proposed by taking the particle counters $x_{0.5}$, $x_{1.0}$, $x_{2.5}$, $x_{5.0}$, $x_{10.0}$ of the sensors, ambient temperature (T) and relative humidity (RH) as input and the concentrations $Y_{pm2.5}$ and Y_{pm10} of PM_{2.5} and PM₁₀ measured by the reference instrumentation as output. The gas pollutant predictions model is also proposed by taking results of the electro-chemical sensors, T and RH as input and the measurements by the reference monitors as output. The experimental results show that the R² of RF for the PM is better than 0.98; the R² of MLR for the PM is less than 0.91; the R² of the other three model are within 0.86 and 0.98. The mean RMSE values of PM_{2.5} and PM₁₀ between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms data are calculated as 3.96 μgm⁻³, 16.16 μgm⁻³, 9.48 μgm⁻³, 10.67 μgm⁻³, 10.74 μgm⁻³, and 7.37 μgm⁻³, 28.90 μgm⁻³, 18.50 μgm⁻³, 18.04 μgm⁻³, 18.17 μgm⁻³, respectively. For the gas pollutants (SO₂, NO₂, CO and O₃), the R² of RF for is better than 0.93; the R² of KNN, BP and GA-BP for the gas pollutants (SO₂, NO₂, CO and O₃) is within 0.27 to 0.97; the R² of MLR for the NO₂, CO and O₃ is within 0.46 to 0.90, but for SO₂ less than 0.40, and even less than 0.1. The mean RMSE values of O₃, CO, NO₂ and SO₂ between the reference data and the RF, MLR, KNN, BP, GA-BP-based algorithms data are calculated as 4.06 μgm⁻³, 16.07 μgm⁻³, 10.23 μgm⁻³, 13.35 μgm⁻³, 13.0 μgm⁻³; 0.04 mgm⁻³, 0.15 mgm⁻³, 0.10 mgm⁻³, 0.12 mgm⁻³, 0.12 mgm⁻³; 3.25 μgm⁻³, 13.90 μgm⁻³, 7.93 μgm⁻³, 10.06 μgm⁻³, 10.15 μgm⁻³; and 1.05 μgm⁻³, 3.86 μgm⁻³, 2.53 μgm⁻³, 3.22 μgm⁻³ and 3.20 μgm⁻³, respectively. These measurements are consistent with the national environmental protection standard requirement



1 of China. Therefore, the low-cost multi-parameter air quality monitoring system, based on RF, MLR, KNN, BP, GA-BP algorithms,
2 can be used to predict the concentrations of PM and gas pollution. In the next research, we will focus on improving the
3 generalization of the algorithms in more applications, and the performance of the SO₂ sensor.

4 **Competing interests**

5 The contact author has declared that none of the authors has any competing interests.

6 **Acknowledgment**

7 Funding for this study was provided by the National Key Research and Development Program of China (Assistance Agreement
8 No. 2021YFB3200403) and Zhengzhou Education Department (No. 23B413006). The authors also wish to thank Minghui, Li,
9 Hongbiao Liu, and Jinlong Wang for helpful conversations.

10 **References**

- 11 Alexopoulos, E. C.: Introduction to Multivariate Regression Analysis, Hippokratia, 14(1), 23-28, 2010.
- 12 Ali, S., Glass, T., Parr, B., Potgieter, J., and Alam, F.: Low Cost Sensor With IoT LoRaWAN Connectivity and Machine Learning-
13 Based Calibration for Air Pollution Monitoring, IEEE Transactions on Instrumentation and Measurement, 70, 5500511,
14 doi:10.1109/TIM.2020.3034109, 2021.
- 15 Amuthadevi, C., Vijayan, D. S., and Ramachandran, V.: Development of air quality monitoring (AQM) models using different
16 machine learning approaches, Journal of Ambient Intelligence and Humanized Computing, 33(2022), doi:10.1007/s12652-020-
17 02724-2, 2021.
- 18 Breiman, L.: Random Forests, Machine Learning, 45(1), 5-32, doi:10.1023/A:1010933404324, 2001.
- 19 Brilli, L., Berton, A., Carotenuto, F., Gioli, B., Gualtieri, G., Martelli, F., and Profeti, S.: Innovative low-cost air quality stations
20 as a supporting means for road traffic regulations in urban areas, 15th International Conference on Atmospheric Sciences and
21 Applications to Air Quality (ASSAQ), 489(2020), 012023, doi:10.1088/1755-1315/489/1/012023, 2020.
- 22 Cross, E. S., Williams, L. R., Lewis, D. K., Magoon, G. R., Onasch, T. B., Kaminsky, M. L., and Worsnop, D. R.: Use of
23 electrochemical sensors for measurement of air pollution: correcting interference response and validating measurements, atmos.
24 Meas. Tech., 10(9), 3575-3588, doi:10.5194/amt-10-3575-2017, 2017.
- 25 Cui, H., Zhang, L., Li, W., Yuan, Z., Wu, M., Wang, C., and Ma, J.: A new calibration system for low-cost Sensor Network in air
26 pollution monitoring, Atmospheric Pollution Research, 12(5), 101049, doi:10.1016/j.apr.2021.03.012, 2021.
- 27 Davut, A., and Baykant, A. B.: An effective integrated genetic programming and neural network model for electronic nose
28 calibration of air pollution monitoring application, Neural Computing and Applications, 34(15), 12633-12652,
29 doi:10.1007/s00521-022-07129-0, 2022.
- 30 Esposito, E., De, V. S., Salvato, M., Bright, V., Jones, R. L., and Popoola, O.: Dynamic neural network architectures for on field
31 stochastic calibration of indicative low cost air quality sensing systems, Sensors Actuators B-Chemical, 231, 701-713,
32 doi:10.1016/j.snb.2016.03.038, 2016.
- 33 Guo, Y., Gao, J., Wang, W., Qin, X., Ren, J., and Ni, D.: Research on outdoor environmental performance evaluation of low-cost
34 atmospheric particulate sensors, China Environmental Science, 40(12), 5133-5141, doi:10.19674/j.cnki.issn1000-
35 6923.2020.0565, 2020.
- 36 Hande, B., and Selda, G.: Estimation of Concentration Values of Different Gases Based on Long Short-Term Memory by Using
37 Electronic Nose Biomedical Signal Processing and Control, 69, doi:10.1016/j.bspc.2021.102908, 2021.
- 38 Hitchman, M. L., Cade, N. J., Gibbs, T. K., and Hedley, N. J. M.: Study of the factors affecting Mass Transport in Electrochemical
39 Gas Sensors, Analyst, 122(11), 1411-1417, doi:10.1039/a703644b, 1997.
- 40 Ioannis, M., Elisavet, S., Agathangelos, S., and Eugenia, B.: Environmental and Health Impacts of Air Pollution: A Review,
41 Frontiers in Public Health, 8(14), 1-13, doi:10.3389/fpubh.2020.00014, 2020.
- 42 Ionascu, M. E., Castell, N., Boncalo, O., Schneider, P., Darie, M., and Marcu, M.: Calibration of CO, NO₂, and O₃ Using Airify:
43 A Low-Cost Sensor Cluster for Air Quality Monitoring, Sensors, 21(23), 7997, doi:10.3390/s21237977, 2021.
- 44 Janabi, S. A., Alkaim, A., Al-Janabi, E., Aljeboree, A., and Mustafa, M.: Intelligent forecaster of concentrations (PM_{2.5}, PM₁₀,
45 NO₂, CO, O₃, SO₂) caused air pollution (IFCsAP), Neural Computing and Applications, 33(21), 14199-14229,
46 doi:10.1007/s00521-021-06067-7, 2021.



- 1 Jiao, W., Hagler, G., Williams, R., Sharpe, R., Brown, R., Garver, D., and Judge, R.: Community Air Sensor Network
2 (CAIRSENSE) project: evaluation of low-cost sensor performance in a suburban environment in the southeastern United States,
3 *atmos. Meas. Tech.*, 9(11), 5281-5292, doi:10.5194/amt-9-5281-2016, 2016.
- 4 Khreis, H., Johnson, J., Jack, K., Dadashova, B., and Park, E. S.: Evaluating the Performance of Low-Cost Air Quality Monitors
5 in Dallas, Texas, *International Journal of Environmental Research and Public Health*, 19(3), 1647, doi:10.3390/ijerph19031647,
6 2022.
- 7 Kumar, T.: Solution of Linear and Non Linear Regression Problem by K Nearest Neighbour Approach: By Using Three Sigma
8 Rule, 2015 IEEE International Conference on Computational Intelligence & Communication Technology, 197-201,
9 doi:10.1109/CICT.2015.110, 2015.
- 10 Liang, Y., Ren, C., Wang, H., Huang, Y., and Zheng, Z.: Research on soil moisture inversion method based on GA-BP neural
11 network model, *International Journal of Remote Sensing*, 40(5-6), 2087-2103, doi:10.1080/01431161.2018.1484961, 2018.
- 12 Liu, Y., Wang, Y., and Zhang, J.: New machine learning algorithm: Random forest, *Information Computing and Applications*.
13 *Proceedings of the Third International Conference, ICICA 2012*, 246-252, doi:10.1007/978-3-642-34062-8_32, 2012.
- 14 Magi, B. I., Cupini, C., Francis, J., Green, M., and Hauser, C.: Evaluation of PM_{2.5} measured in an urban setting using a low-cost
15 optical particle counter and a Federal Equivalent Method Beta Attenuation Monitor, *Aerosol Science & Technology*, 54(2),
16 147-159, doi:10.1080/02786826.2019.1619915, 2020.
- 17 Masson, N., Piedrahita, R., and Hannigan, M.: Quantification method for electrolytic sensors in long-term monitoring of ambient
18 air quality, *Sensors*, 15(10), 27283-27302, doi:10.3390/s151027283, 2015.
- 19 Mead, M. I., Popoola, O. A. M., Stewart, G. B., and Landshoff, P.: The use of electro-chemical sensors for monitoring urban air
20 quality in Low-cost, high-density networks, *Atmos. Environ.*, 70, 186-203, doi:10.1016/J.ATMOSENV.2012.11.060, 2013.
- 21 Ning, M., Guan, J., Liu, P., Zhang, Z., and O'Hare, G. M. P.: GA-BP Air Quality Evaluation Method Based on Fuzzy Theory,
22 *CMC-Computers Materials & Continua*, 58(1), 215-227, doi:10.32604/cmc.2019.03763, 2019.
- 23 Singh, A., Ng'ang'a, D., Gatari, M. J., Kidane, A. W., Alemu, Z. A., Derrick, N., and Webster, M. J.: Air quality assessment in
24 three East African cities using calibrated low-cost sensors with a focus on road-based hotspots, *Environmental Research*
25 *Communications*, 3(7), 075007, doi:10.1088/2515-7620/ac0e0a, 2021.
- 26 Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Field calibration of a cluster of low-cost available
27 sensors for air quality monitoring. Part A: Ozone and nitrogen dioxide, *Sensors Actuators, B Chem.*, 215, 249-257,
28 doi:10.1016/j.snb.2015.03.031, 2015.
- 29 Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Field calibration of a cluster of low-cost
30 commercially available sensors for air quality monitoring. Part B: NO, CO and CO₂, *Sensors and Actuators B-Chemical*, 238,
31 706-715, doi:10.1016/j.snb.2016.07.036, 2017.
- 32 Wang, C., Pan, B., Wu, X., Song, Y., Zhang, L., Ma, J., and Sun, K.: Research on Quality Control of Atmospheric Grid Monitoring
33 Based on Large Data Analysis, *Environmental Monitoring in China* 32(6), 1-6, doi:10.19316/j.issn.1002-6002.2016.06.01, 2016.
- 34 Wang, S., Li, L., Ma, W., and Chen, X.: Trajectory analysis for on-demand services: A survey focusing on spatial-temporal demand
35 and supply patterns, *Transportation Research Part C: Emerging Technologies*, 108, 74-99, doi:10.1016/j.trc.2019.09.007, 2019.
- 36 Wang, X., Shi, F., Yu, L., and Li, Y. (2013). *MATLAB neural network analysis*: Beihang University Press.
- 37 Wei, P., Ning, Z., Ye, S., Sun, L., and Yang, F.: Impact Analysis of Temperature and Humidity Conditions on Electrochemical
38 Sensor Response in Ambient Air Quality Monitoring, *Sensors*, 18(2), doi:10.3390/s18020059, 2018.
- 39 Xu, X., Guo, H., and Fan, J.: Water quality evaluation of Xiaoshan water quality station in eastern Zhejiang Water Diversion
40 Project Based on BP network, *Journal of Physics: Conference Series*, 1732, 012035, doi:10.1088/1742-6596/1732/1/012035,
41 2021.
- 42 Zhao, C., Wang, Y., Shi, X., Zhang, D., Wang, C., Jiang, J., and Zhang, Q.: Estimating the Contribution of Local Primary
43 Emissionsto Particulate Pollution Using High - Density Station Observations, *J. Geophys. Res.: Atmospheres*, 124, 1648-1661,
44 doi:10.1029/2018JD028888, 2019.
- 45 Zhao, P., and Lai, L.: Minimax Rate Optimal Adaptive Nearest Neighbor Classification and Regression, *IEEE Transactions on*
46 *Information Theory*, 67(5), 3155-3182, doi:10.48550/arXiv.1910.10513, 2021.
- 47 Zimmerman, N., Presto, A. A., Kumar, S. P. N., Gu, J., Haurlyliuk, A., Robinson, E. S., and Robinson, A. L.: A machine learning
48 calibration model using random forests to improve sensor performance for lower-cost air quality monitoring, *atmos. Meas.*
49 *Tech.*, 11(1), 291-313, doi:10.5194/amt-11-291-2018, 2018.
- 50 Zoest, V. v., Osei, F. B., Stein, A., and Hoek, G.: Calibration of low-cost NO₂ sensors in an urban air quality network, *Atmospheric*
51 *Environment*, 210, 66-75, doi:10.1016/j.atmosenv.2019.04.048, 2019.
- 52
53