



# Machine Learning Calibration of Low-Cost Air Quality Gas Sensors

Giannis Ioannidis<sup>1</sup>, Vincent Langat<sup>1</sup>, Roubina Papaconstantinou<sup>1</sup>, Spyros Bezantakos<sup>1</sup>,  
Prashant Kumar<sup>2,3</sup>, and George Biskos<sup>1,4</sup>

<sup>1</sup>Climate and Atmosphere Research Centre (CARE-C), The Cyprus Institute, Nicosia 2121, Cyprus

5 <sup>2</sup>Global Centre for Clean Air Research (GCARE), School of Sustainability, Civil and Environmental Engineering, Faculty of Engineering and Physical Sciences, University of Surrey, Guildford GU2 7XH, United Kingdom

<sup>3</sup>Institute for Sustainability, University of Surrey, Guildford GU2 7XH, United Kingdom

<sup>4</sup>Faculty of Civil Engineering and Geosciences, Delft University of Technology, Delft, The Netherlands

*Correspondence to:* George Biskos (g.biskos@cyi.ac.cy, g.biskos@tudelft.nl)

10 **Abstract.** Low cost sensors (LCSs) for measuring the concentrations of gaseous pollutants hold great promises for air quality monitoring (AQM) as they can improve the spatio-temporal resolution of observational networks. However, the performance of LCSs is affected by a number of factors including temperature and relative humidity of ambient air, as well as cross-sensitivities with gaseous species other than the target gas, thereby deteriorating the quality of their measurements. To address these issues, data from LCSs can be calibrated against reference instruments using machine learning (ML) algorithms. Here, we have  
15 evaluated the performance of a number of ML algorithms for calibrating measurements from CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub> LCSs against respective reference measurements. The best model is then used to determine (1) the influence of temporal resolution of the measurements to the calibration performance, (2) the minimum fraction of data needed for model training while maintaining the quality of calibrated measurements within acceptable levels, and (3) the ideal calibration frequency with collocated reference measurements. We found that the quality of LCS measurements improve significantly for all sensors  
20 after ML calibration, with Random Forest (RF) being the best performing algorithm, corroborating previous works. By varying the temporal resolution of the training data from 1 h to 2 min, the performance of the RF model in terms of the normalized root mean squared error and the relative expanded uncertainty calculated at maximum observed concentration improves by 11-21%. The results also suggest that the minimum fraction of data required for training the ML models depends on the frequency of carrying out collocated measurements with reference instruments and using the resulting datasets for  
25 training the calibration model. If the calibrations are carried out on a monthly basis, ca. 50% of the period is needed for collecting data to train the RF algorithm and qualify the LCSs for indicative measurements as defined by the EU directive (2008/50/EC). If the training is carried out every 3 or 6 months by sampling the training data continuously, then ca. 60% of the measuring period is required for collecting training data. In those cases, if the sampling of the training data is made over specific periods every month, but the entire training dataset is used to calibrate the measurements over 3 or 6 months, the  
30 amount of data required for qualifying the LCSs for indicative measurements can significantly reduce to 22%. However, this would require that the measurements from the LCSs be calibrated retrospectively, which for specific applications is not such of a problem.



## 1 Introduction

35 The adverse effects of air pollution on human health and the environment warrant for continuous monitoring of air quality  
(Heal et al., 2012). Traditionally, air quality monitoring for regulatory purposes is carried out at a number of observational  
stations equipped with reference-grade air quality monitors. The number of stations that can be operated within a given  
geographical area, however, is limited by the high installation and operation cost of the instruments (Kumar and Sahu, 2021;  
Arroyo et al., 2021). This has spurred the development of low-cost air quality gas sensors (Kumar et al., 2015) that are  
40 significantly less expensive compared to their reference-grade instrument counterparts, motivating research towards improving  
their very sensing nanomaterials to meet requirements in air quality monitoring (Baranwal et al., 2022; Isaac et al., 2022). In  
fact, low-cost sensors (LCSs) are increasingly being used to complement the reference instruments in measuring the  
concentration of air pollutants, as they can significantly increase spatio-temporal resolution of air quality monitoring (AQM)  
networks (Lewis et al., 2018; Chen et al., 2018; Zuidema et al., 2021; Nowack et al., 2021; Zimmerman, 2022) and capture  
45 spatial and temporal gradients (Papaconstantinou et al., 2026). Their low cost and ease of operation also allow personal use  
for monitoring air quality in the indoor and outdoor environments, providing great means for assessing the impacts of air  
pollutants on human health (Schäfer et al., 2021; Patra et al., 2021).

Although LCSs hold great promises for expanding existing AQM networks, they have a number of technical limitations  
including high limits of detection, low precision and accuracy, and signal drift, which at the moment prohibit their widespread  
50 use (Papaconstantinou et al., 2023; Rai et al., 2015). Some of these technical limitations are related to environmental conditions  
they are operated under (i.e., temperature and Relative Humidity; RH), and cross-sensitivities to other gaseous pollutants. As  
a result, when LCSs are deployed for field measurements at ambient conditions where these factors vary substantially, the  
deviation between the measurements they provide and those reported by reference instruments can become large, failing to  
meet the requirements defined by environmental regulatory agencies (Samad et al., 2020; Schäfer et al., 2021). According to  
55 the European Union (EU) directive on ambient and cleaner air (2008/50/EC), for air quality measurements to qualify as  
reference measurements, they should exhibit relative expanded uncertainties (REUs) lower than 15%; i.e., something that is  
achieved by reference-grade instruments typically used in AQ monitoring. If the REUs are higher than 15% and lower than  
25% for CO, NO<sub>2</sub> and SO<sub>2</sub>, and 30% for O<sub>3</sub>, the measurements can be qualified as indicative (EU-directive, 2008; Equivalence,  
2010), which can be useful for a number of applications; e.g., for identifying pollution sources and hotspots in cities, and  
60 determining the spatio-temporal variability in the concentration of air pollutants (Schäfer et al., 2021).

The poor agreement between LCS and reference-grade instrument measurements is not a surprise considering that the former  
typically come calibrated by the manufacturers at laboratory conditions, which do not capture the complexity encountered in

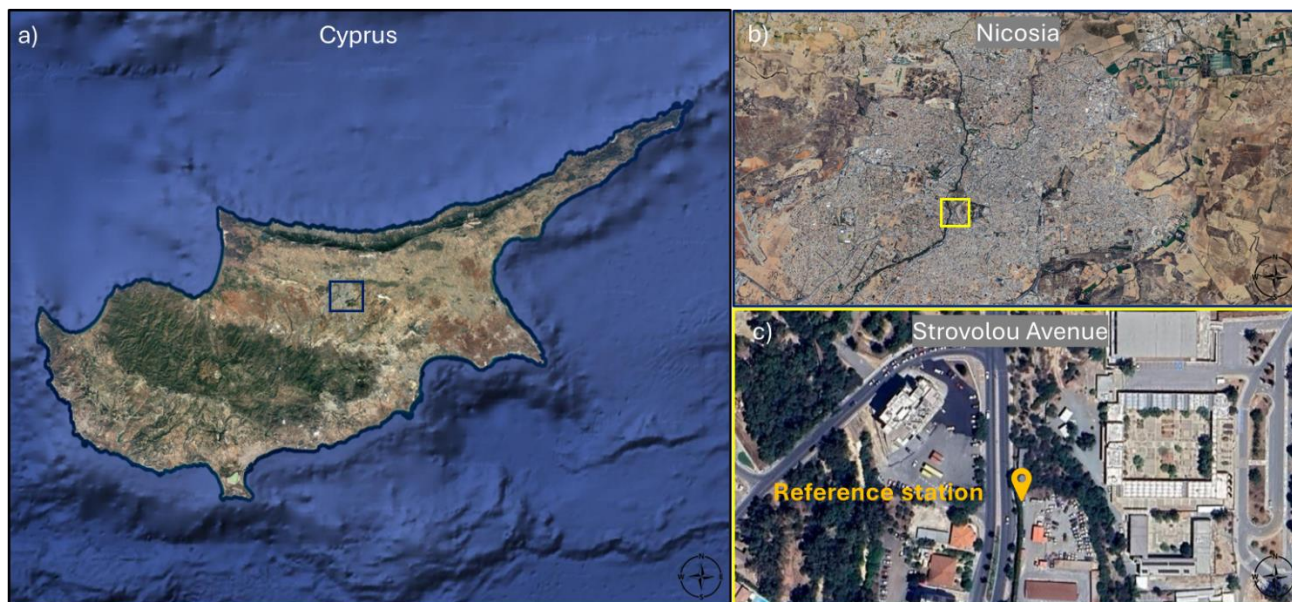


the field (Zuidema et al., 2021; Koziel et al., 2024; Garbagna et al., 2025). To address this gap, previous efforts have focused on the post calibration of LCS measurements using machine learning (ML) models that take into account the variability of temperature, relative humidity and interfering pollutants, as well as other factors that may influence measurements carried out in the field (Zimmerman et al., 2018; Ferrer-Cid et al., 2020; Mahajan and Kumar, 2020; Okafor and Delaney, 2020; Song et al., 2020; Nowack et al., 2021; Patra et al., 2021; Kumar and Sahu, 2021; Vajs et al., 2021; Podder et al., 2024; Sousan et al., 2025). To do so, one needs to collocate the LCSs with respective reference instruments for a certain amount of time in order to gather data for training the calibration models.

In this work, we use concentration measurements of atmospheric pollutants recorded by LCSs and reference instruments over a period of 6 months in a traffic station in the city of Nicosia, Cyprus, in order to train and evaluate the performance of five ML algorithms; namely Linear Regression (LR; Kumar and Sahu 2021); Support Vector Regression (SVR; Kumar and Sahu 2021); Random Forest Regressor (RF; Zimmerman et al. 2018); Artificial Neural Network (ANN; Spinelle et al. 2015); and Extreme Gradient Boosting (XGBoost; Chen and Guestrin 2016). The best model is then used to determine the effects of a number of practical parameters including the temporal resolution of the training data, how the training data is sampled and the frequency of training/calibration on the performance of the algorithm, as well as the minimum fraction of data needed for training without compromising significantly their quality (i.e., the accuracy of the resulted post calibrated measurements).

## 2 Experimental

In this work, we employed four electrochemical sensors manufactured by Alphasense to measure the concentrations of CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub>. More information about the technical specifications of each sensor is provided by Papaconstantinou et al. (2023) and in Table S1 of the Supplement. Each sensor provides two raw analogue voltage signals: one from the working electrode and the other from an auxiliary electrode that serves as the zero background signal against which the signal of the working electrode is compared (Masic et al., 2018; ANN803-05, 2019; Arroyo et al., 2021). These signals are then converted to corresponding concentrations, expressed in ppb, using calibration equations provided by the manufacturer (ANN803-05, 2019). For our analysis, we will refer to the concentrations calculated based on such equations as laboratory (LAB) calibrated concentrations.



90 **Figure 1:** Satellite image showing the island of Cyprus (Figure 1a). Figure 2b illustrates the greater area of the capital of Cyprus, Nicosia and Figure 3c indicates the location of the traffic station where the measurements were carried out. The station is located next to one of the busiest avenues, i.e., Strovolou Avenue, in Nicosia, Cyprus, with coordinates 35°09'07.2" N and 33°20'52.0" E. Base map from Google Earth, imagery © Google, Maxar Technologies, accessed February 2026.

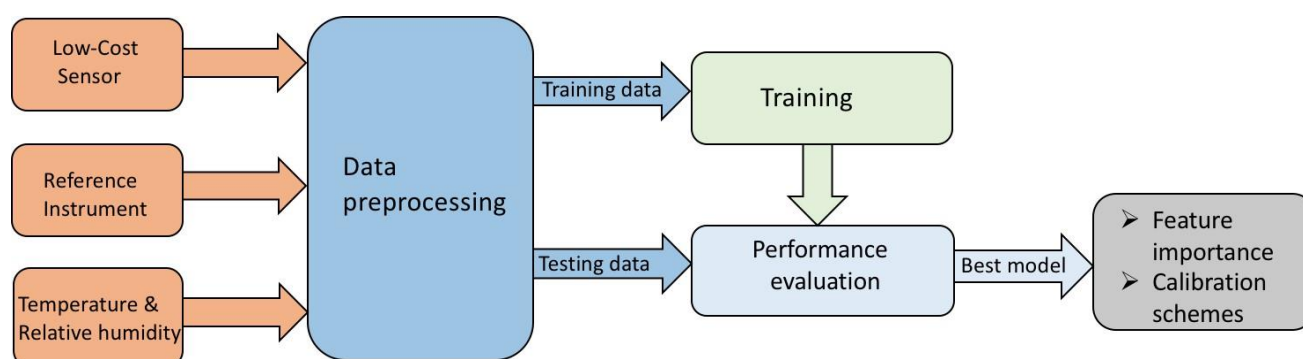
The measurements were carried out between 2 October 2019 and 31 March 2020 at one of the national regulatory air quality monitoring stations in Nicosia. As shown in Figure 1c, the station is located approximately 10 m away from one of the busiest  
95 avenues and is equipped with reference-grade instruments for measuring the concentration of CO (Ecotech Serinus 30), NO<sub>2</sub> (Ecotech Serinus 40), O<sub>3</sub> (Thermo Scientific 49i) and SO<sub>2</sub> (Ecotech Serinus 50). These instruments are calibrated at least once every 3 months, and after maintenance. More detailed information on the technical specifications of these instruments is provided in the Supplement Table S2.

## 2.1 ML calibration procedure

100 Figure 2 shows the steps followed for ML calibration of the LCS measurements. The data from the reference instruments, as well as the temperature and RH recorded by the sensors located close to the reference station, had a time resolution of 2 min, whereas the measurements from the LCSs were recorded every 2 s. To align the dataset, the LCS measurements were averaged every 2 min and concatenated with the respective reference data, temperature and RH. Subsequently, data cleaning was done by dropping all rows with the missing values and those with negative net sensor signals (i.e., obtained by subtracting the signals  
105 reported by the auxiliary electrode from those reported by working electrode). From the cleaned dataset, the net sensor signals (NSS), temperature, relative humidity (RH), month, day of week and hour were selected as input variables (also referred to as



features) for each algorithm, whereas the respective reference concentrations were used as response (also referred to as dependent) variables. The temporal variables (month, day of week, and hour) were included as input to unravel the importance of monthly, weekly and diurnal variabilities. For the ML calibration of the NO<sub>2</sub> and O<sub>3</sub> LCSs we also included the O<sub>3</sub> and NO<sub>2</sub> reference concentrations, respectively, as input variables in order to account for cross-sensitivities. According to manufacturer and literature reports, the O<sub>3</sub> LCS is affected by NO<sub>2</sub> (ANN803-05, 2019; Pang et al., 2017), and vice versa (Maag et al., 2016).



**Figure 2:** Flow diagram showing the stages of the ML calibration process employed in this work. First, we merge the measurements from the LCSs, the reference instruments and the meteorological sensors. The resulting dataset are then pre-processed and segmented into training and testing sets. The training sets are used to train the ML models while the testing sets are used to evaluate their performance. The best-performing model is then used to assess (1) the significance of input variables to model performance, and (2) how different calibration practices affect the fraction of training data required.

LCS calibration was performed using five ML models as mentioned in the introduction: LR, SVR, RF, ANN and XGBoost. All calculations were carried out in the anaconda environment using python version 3.8.6. To train and evaluate the performance of these models, the 6-month datasets were split into training and testing sets. The training sets were derived from the first 80% of the data of each month and used to train and tune the parameters of the models. The testing sets, which were used to evaluate the performance of the models, comprised of the remaining 20% of the data (see Fig. S1 in the Supplement).

The hyper-parameters for the SVR and ANN ML models were tuned through a 5-fold cross-validation and grid search, whereas those of the RF and XGBoost algorithms were auto-tuned through the AutoML library (FLAML) developed by Microsoft (Wang et al., 2021). The optimal values for the main hyper-parameters obtained for each model are provided in Table S3 in the Supplement. The other parameters for each model were kept at their default values. The calibration performance of the tuned models were then evaluated based on a number of statistical indicators and the best-performing model was used to determine the importance of each feature in the model and assess the effect of (1) the temporal resolution of the training data, (2) how the training data are sampled, and (3) the calibration frequency on the fraction of data required for training the ML



algorithms. We should note that the importance of each input variable in the model was determined using the permutation feature importance function within the Scikit-learn ML library.

## 2.2 ML calibration procedure

135 Performance evaluation of the models was done in two stages. In the first stage we evaluated the accuracy of their calibration based on Pearson correlation coefficient ( $r$ ), coefficient of determination ( $R^2$ ) and normalized root mean squared error (NRMSE), given respectively as:

$$r = \frac{\sum_{t=1}^n (\text{Cal}_t - \overline{\text{Cal}})(\text{Ref}_t - \overline{\text{Ref}})}{\sqrt{\sum_{t=1}^n (\text{Cal}_t - \overline{\text{Cal}})^2 \sum_{t=1}^n (\text{Ref}_t - \overline{\text{Ref}})^2}}, \quad (1)$$

140

$$R^2 = 1 - \frac{\sum_{t=1}^n (\text{Cal}_t - \text{Ref}_t)^2}{\sum_{t=1}^n (\text{Ref}_t - \overline{\text{Ref}})^2}, \quad (2)$$

$$\text{NRMSE} = \frac{\sqrt{\frac{\sum_{t=1}^n (\text{Cal}_t - \text{Ref}_t)^2}{n}}}{\overline{\text{Ref}}}. \quad (3)$$

145 Here  $\text{Cal}_t$ ,  $\overline{\text{Cal}}$ ,  $\text{Ref}_t$ ,  $\overline{\text{Ref}}$  and  $n$  denote respectively the calibrated concentration at time  $t$ , the mean of calibrated concentrations, the reference concentration at time  $t$ , the mean of the reference concentrations, and the total number of data points. The Pearson correlation coefficient,  $r$ , provides information on the strength of associations between the calibrated concentrations and their corresponding reference concentrations.  $R^2$  is a measure of the extent to which the input variables used in the models explain the variation in the dependent variable, whereas NRMSE expresses the error between the calibrated and reference  
150 concentrations normalized by the mean of the reference concentrations.

In the second stage, we created target diagrams to infer the level of bias and variance in the model, which are key parameters for diagnosing whether a model is over-fitting or under-fitting the data. An under-fitted model has high bias and low variance, whereas the opposite is true for an over-fitted model (Kumar and Sahu, 2021; Yu et al., 2006). For LCS calibration, the root mean square error (RMSE) of the calibration model captures both the bias and the variance in the data provided by LCSs, and  
155 is typically used to create target diagrams (Jolliff et al., 2009) and to visualize the performance of different models. The RMSE is determined as:

$$\text{RMSE}^2 = \text{MBE}^2 + \text{CRMSE}^2 \quad (4)$$



160 Here, MBE and CRMSE are the mean bias error and the centred root mean square error (Zimmerman et al., 2018; Thunis et al., 2012), expressed as:

$$\text{MBE} = \frac{1}{n} \sum_{t=1}^n (\text{Cal}_t - \text{Ref}_t) \quad (5)$$

$$\text{CRMSE} = \sqrt{\sigma_{\text{Cal}}^2 + \sigma_{\text{Ref}}^2 - 2r\sigma_{\text{Cal}}\sigma_{\text{Ref}}} \quad (6)$$

165

where  $\sigma_{\text{Cal}}$  and  $\sigma_{\text{Ref}}$  are the standard deviation of calibrated and reference concentrations, respectively.

### 2.3 Data quality assessments

The goal of calibrating measurements from LCSs using ML models is to improve their accuracy to meet the data quality levels required for specific applications. Here we use the EU directive 2008/50/EC as a guideline to assess the quality of the data produced by the LCSs and the different calibration schemes we tried. According to the directive, the LCSs can primarily be used for indicative measurements (i.e., measurements with less stringent uncertainty requirements in comparison to regulatory measurements), in which case the uncertainty limits are set to 25% for CO, NO<sub>2</sub> and SO<sub>2</sub>, and 30% for O<sub>3</sub>. These uncertainty requirements have to be met at specific limit values for CO, NO<sub>2</sub> and SO<sub>2</sub>, and at the target value for O<sub>3</sub>. Limit values are specific concentrations below which the harmful effects of gaseous pollutants on human health and the environment are reduced. Similarly, target values are fixed levels below which the long-term effects of air pollutants on human health and the environment are minimized (Equivalence, 2010). Limit or target values are determined for specific periods over which the measurements are averaged depending on the pollutant. For example, the limit values for NO<sub>2</sub> and SO<sub>2</sub> are determined on an hourly basis and have values of 200 and 350 ppb, respectively, whereas the limit value for CO and the target value for O<sub>3</sub> are determined based on an 8-hour averaging, having respective values of 10 ppm and 125 ppb (EU-directive, 2008). We should note here that the concentrations for CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub> observed at the station used in this study were lower than the above-mentioned limit or target values. Considering that, the uncertainties reported in this work are calculated at the maximum concentrations reported by the reference instruments.

The EU directive uses the relative expanded uncertainty (REU) as a data quality indicator (DQI), calculated based on the guidelines provided by Walker and Schneider (2020). According to those, the calibrated LCS concentrations are assumed to have a linear relationship with the reference concentrations given as:

$$\text{Cal} = \theta_0 + \theta_1 \cdot \text{Ref} \quad (7)$$

where Cal and Ref are calibrated and reference concentrations, respectively, whereas  $\theta_0$  and  $\theta_1$  are regression parameters obtained by a two-step adjusted orthogonal regression. The relative expanded uncertainty at the maximum concentration is then



190 calculated at 95% confidence interval as:

$$REU_{\max} = \frac{2 \sqrt{\frac{RSS}{n-2} + (\lambda - (\theta_1 - 1)^2) \cdot U^2_{\text{Ref}_{\max}} + \theta_0 + (\theta_1 - 1) \cdot \text{Ref}_{\max}}}{\text{Ref}_{\max}} \cdot 100\% , \quad (8)$$

where,

$$195 \quad RSS = \sum_{t=1}^n ((\text{Cal}_t - \theta_0 - \theta_1 \cdot \text{Ref}_t)^2) - (\theta_1^2 + \lambda) \cdot U^2_{\text{Ref}_{\max}}, \text{ and} \quad (9)$$

$$\lambda = \frac{(U_{\text{Cal}_{\max}})^2}{(U_{\text{Ref}_{\max}})^2} \quad (10)$$

Here,  $RSS$ ,  $REU_{\max}$ ,  $\text{Ref}_{\max}$ ,  $\text{Cal}_{\max}$ ,  $U_{\text{Cal}_{\max}}$  and  $U_{\text{Ref}_{\max}}$  are the residual sum of squares, the relative expanded uncertainty at  
200 the maximum concentration observed, the maximum reference concentration, the maximum calibrated concentration, the  
random uncertainty for maximum calibrated concentration, and the random uncertainty for maximum reference concentration,  
respectively. For  $\text{CO}$ ,  $\text{NO}_2$  and  $\text{SO}_2$ , the  $U_{\text{Ref}_{\max}}$  values were obtained from the technical specification of each reference in-  
strument as provided by the manufacturers (see Table S2 in the Supplement). For the  $\text{O}_3$  reference measurements, for which  
no information about the accuracy of the reference instrument is provided by the manufacturer, we assumed  $U_{\text{Ref}_{\max}}$  to be 0.1%.  
205 For  $U_{\text{Cal}_{\max}}$ , which is difficult to determine considering that its value depends on uncertainty of the reference concentrations,  
we assumed  $U_{\text{Cal}_{\max}} = U_{\text{Ref}_{\max}}$  and thus  $\lambda = 1$  in our calculations following recommendation provided by Walker and Schneider  
(2020).

### 3 Results and discussion

Table 1 provides summary statistics of the concentration measurements reported by the reference instruments and the  
210 LAB calibrated LCSs. The differences in the mean and standard deviation between the measurements obtained by the  
LCSs and the reference instruments for the four pollutants range between ca. 10 and 1300% and ca. 30 and 2600%,  
respectively. To evaluate the statistical significance of these differences, we first checked normality of all LCS and  
reference measurements by running Shapiro-Wilk tests, which returned p-values less than 0.05, indicating that all the  
LCS and reference measurements were normally distributed. Next, we performed t-tests for all the pollutants to assess  
215 the statistical significance of the difference between the mean of LCS and reference concentrations (Tiku and Akkaya,  
2004). To evaluate the statistical significance of the differences in variances of the LCS and reference concentration  
measurements, we performed the Fligner and Killeen tests (F-K tests; Si et al. 2020) for all the cases.



220 **Table 1:** Summary statistics of concentration measurements of CO, NO<sub>2</sub>, SO<sub>2</sub> and O<sub>3</sub> recorded by the LCS, using the laboratory (LAB) calibrations, and the reference (REF) instruments.

	CO (REF)	CO (LAB)	NO <sub>2</sub> (REF)	NO <sub>2</sub> (LAB)	O <sub>3</sub> (REF)	O <sub>3</sub> (LAB)	SO <sub>2</sub> (REF)	SO <sub>2</sub> (LAB)
Number of data points	111,025	111,025	102,427	102,427	94,520	94,520	44,520	44,520
Average (ppb)	470	430	18	33	18	62	1.6	27
Standard deviation (ppb)	288	262	13	26	15	39	0.98	23
Minimum (ppb)	2.87	1.79	0.26	0.030	0.002	0.070	0.001	0.170
Maximum (ppb)	3,573	3,110	99	1,365	64	1,789	96	252

The t-tests showed p-values  $< 10^{-5}$  for all four pollutants, which is lower than the set significance level of 0.05, suggesting that the difference in mean values of the laboratory calibrated and reference concentrations are statistically significant. Similarly, the F-K tests returned p-values  $< 10^{-10}$ , indicating that the differences in variances of laboratory calibrated LCSs and reference concentration measurements for all the pollutants are statistically significant considering that a significance level of 0.05 was used in the calculations. We should note here that the difference between LCS and reference measurements was more pronounced for SO<sub>2</sub>. This can be explained by the fact that the concentration of SO<sub>2</sub> at the observational site was on average less than 5ppb, which was lower than the limit of detection (LoD) of the SO<sub>2</sub> LCS (see Table S1 in the Supplement).

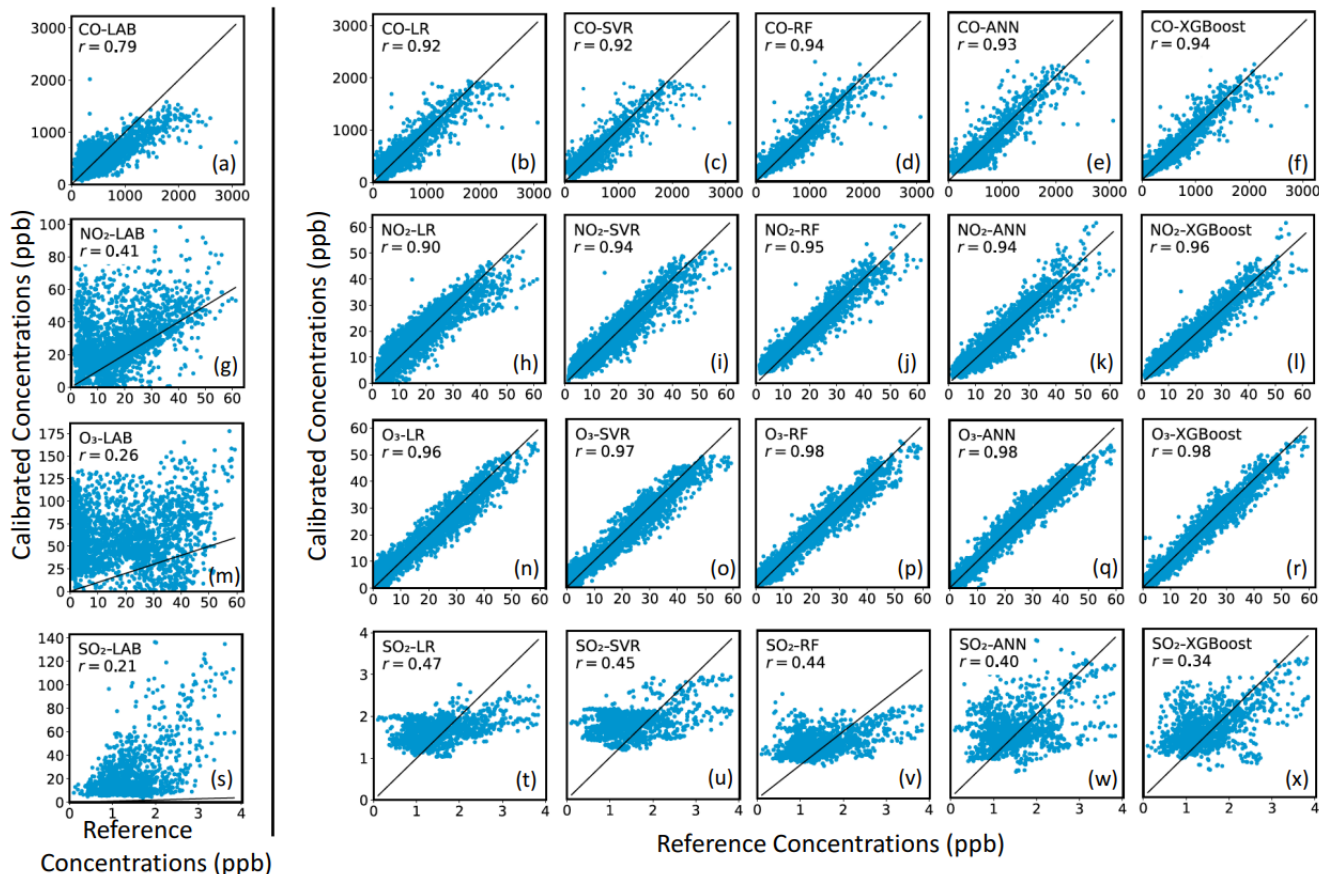
230 The statistically significant difference between the measurements recorded with the LCSs and the respective reference instruments is not surprising considering that the laboratory calibrations carried out by the manufacturer only account for the effect of the temperature at room conditions, which does not fully capture the meteorological variabilities in ambient conditions. Moreover, the calibration equations provided by the manufacturer do not take into account the influence of other factors, such as RH which can significantly affect the performance of the LCSs (Papaconstantinou et al., 2023; Samad et al., 235 2020; Gonzalez et al., 2019; Pang et al., 2018).

### 3.1 Evaluation of ML algorithms

Figure 3 shows the correlation between the calibrated and reference concentrations for the CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub> for all calibration approaches we tested. We should note that for all the tests we used 80 % of the entire dataset for training and 20% for validation, whereas all the measurements had a 2-min time resolution. Evidently, we observe a good agreement between the measurements from the CO LCS and the corresponding reference measurements, even without using any of the ML calibrations (see Fig. 3a). The rest of the LCSs, however, exhibit moderate to low correlation with their respective reference measurements (see Fig. 3g, 3m and 3s). The measurements reported by the SO<sub>2</sub> LCS highly overestimate the SO<sub>2</sub> concentrations as compared to those reported by the reference instrument. Overall, we observe an improvement in the agreement between measurements reported by the LCS and the respective reference instruments after ML calibration. This



245 improvement is more pronounced for CO, NO<sub>2</sub> and O<sub>3</sub> measurements, as reflected by the high *r* values that are larger than 0.9, following ML calibration. In contrast, for SO<sub>2</sub>, despite that there was a general improvement in agreement between the SO<sub>2</sub> LCS and reference measurements when using ML calibration, they exhibited *r* values that were hardly above 0.5 (see Fig. 3s-x). As mentioned above, this can be attributed to the low concentrations of SO<sub>2</sub> observed at the measuring site, which on average were far lower (i.e., 1.6 ± 1 ppb; see Table 1) than the LoD of the SO<sub>2</sub> LCS (i.e., 5 ppb; see Table S1 in the Supplement and Alphasense-SO<sub>2</sub>-B4 2019).

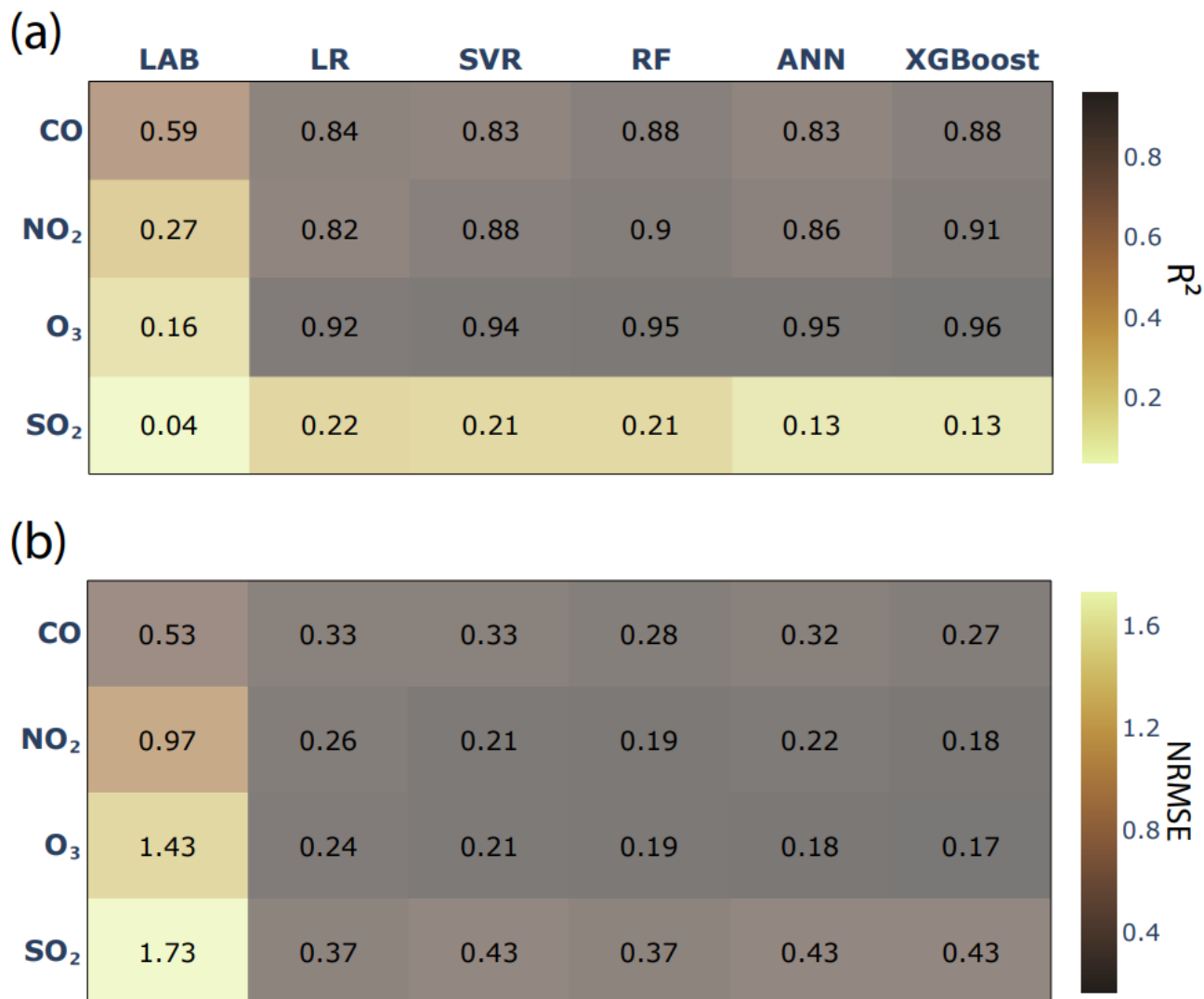


**Figure 3:** Correlation between calibrated LCS and reference concentration measurements for CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub>. The black dotted correspond to 1:1 relation.

Figure 4 provides heat maps that show the performance of LAB and ML calibrations based on R<sup>2</sup> and NRMSE. With the exception of the SO<sub>2</sub> LCS, all the ML calibrations performed reasonably well in terms of R<sup>2</sup> as indicated in Figure 4a. Overall, the RF model exhibited the highest R<sup>2</sup> values for all pollutants, followed by the ANN and the XGBoost models. The good performance of the ML calibrations was also reflected by the relatively low NRMSE values (see Fig. 4b), which exhibit a significant improvement compared to the values of their corresponding LAB calibration. Another key point to note here is



260 that, even though the CO LCS exhibited a good performance without ML calibration, it is evident that performance of ML calibrations for NO<sub>2</sub> and O<sub>3</sub> measurements are generally better than that of the CO. This is explained by the fact that, for the NO<sub>2</sub> and O<sub>3</sub> ML calibrations we accounted for cross-sensitivities by including O<sub>3</sub> and NO<sub>2</sub> reference concentrations, respectively, as input variables.

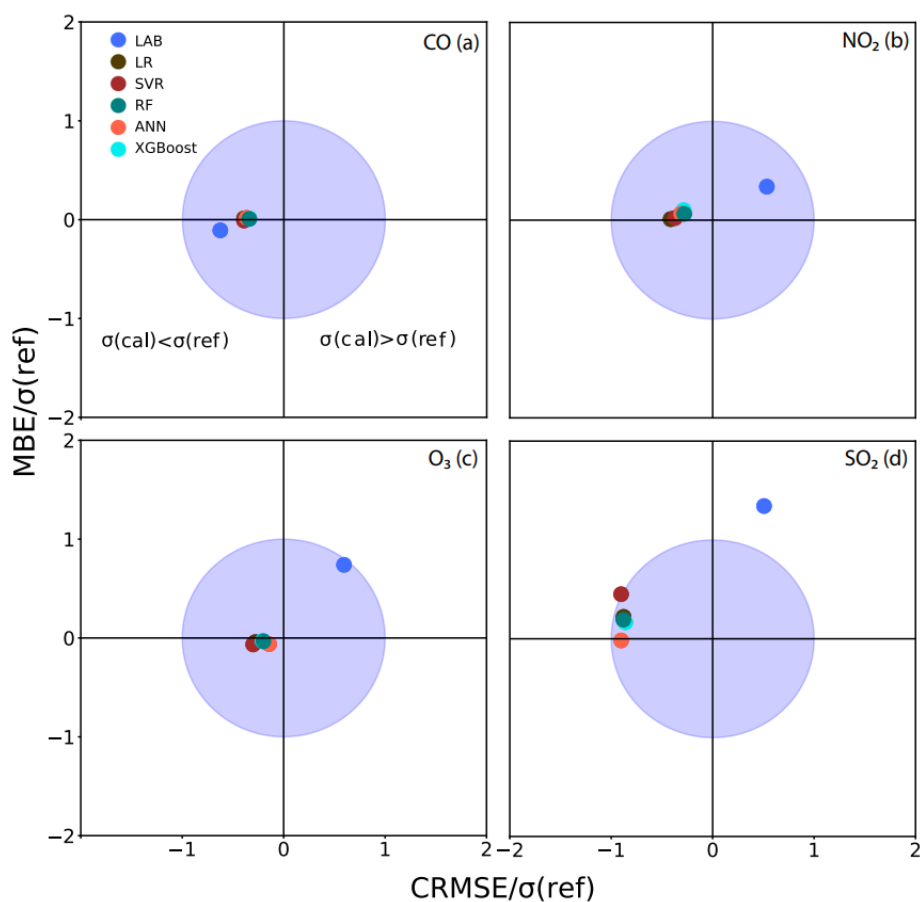


265 **Figure 4:** Heat maps showing the performance of the LCSs in terms of (a) coefficient of determination ( $R^2$ ), and (b) normalized root mean squared error (NRMSE) after calibration using the laboratory tests carried out by the manufacturer (LAB) and the five ML algorithms employed in this work.



270 Figure 5 shows target diagrams that help to further assess the performance of the ML models in terms of the bias and variance, and to compare the standard deviations of their predictions with those of their respective reference concentrations. In each diagram the normalised MBE is plotted against the normalized CRMSE. The MBE and CRMSE were normalised by the standard deviation of their corresponding reference concentrations to facilitate comparison of the calibration performance of the models across different pollutants. Consequently, the distance of each point from the centre of the diagrams corresponds to the RMSE normalized by the standard deviation of the respective reference concentrations. We will denote this with nRMSE to distinguish it from the one described in Eq. (3). As shown in the diagrams, the ML calibrations generally exhibited lower nRMSE values compared to corresponding laboratory calibrations. The RF, ANN and XGBoost models generally performed better in terms of nRMSE when compared to the LR and SVR models.

275



**Figure 5:** Target diagrams showing the overall performance of the laboratory and ML learning calibrations of the data from the LCS tested in this work. The light blue circles show the region where the NRMSE is less than unity and the variance of residuals of the calibrated concentrations are lower than those of their corresponding reference concentrations.

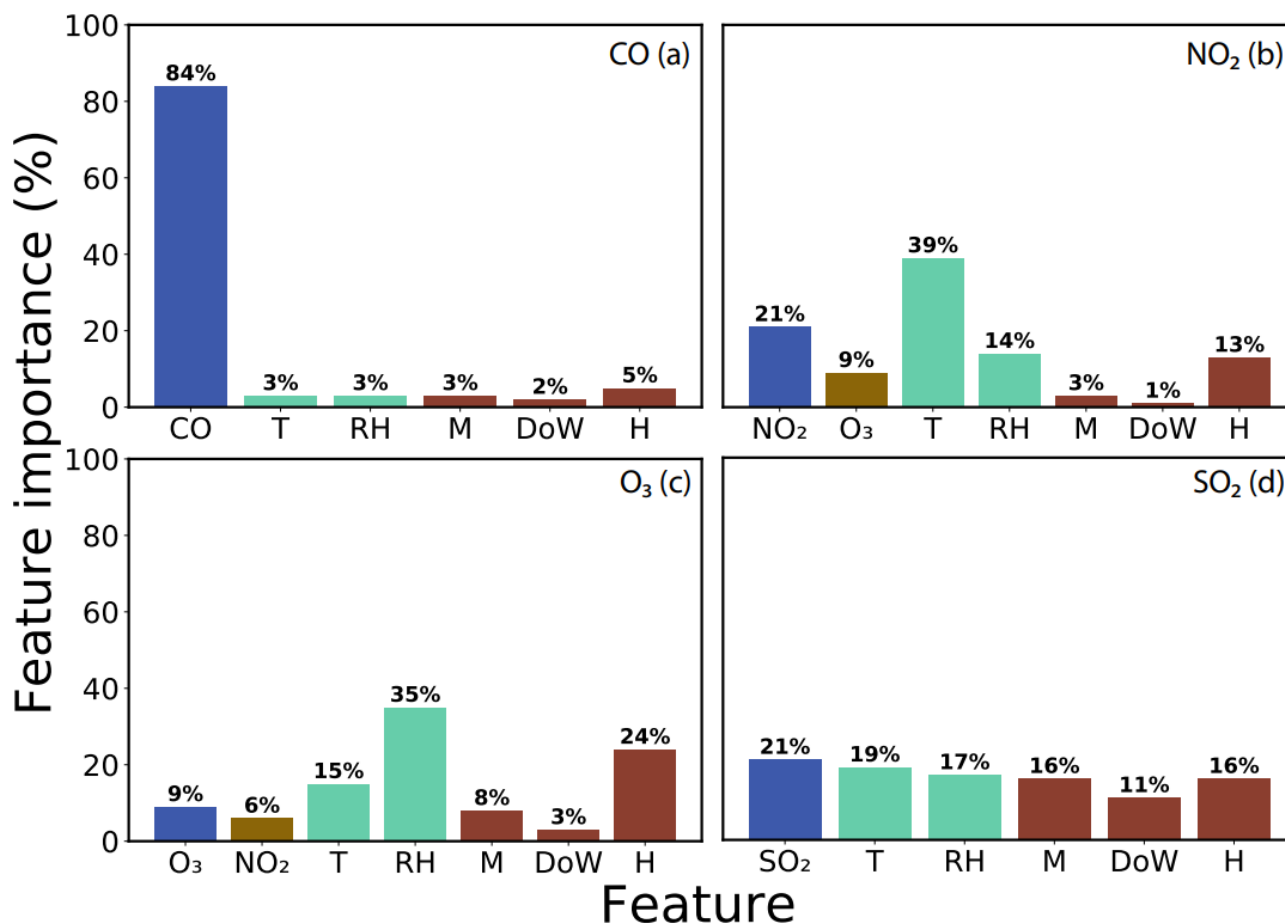


280 The horizontal lines passing through the centers of the circles in Figure 5 correspond to cases with zero bias. Any point above  
or below these lines indicate that the corresponding algorithm overestimates or underestimates the values reported by the  
reference instruments, respectively. The CO LAB calibration exhibit a small bias and thus corroborating the good correlation  
between those, and the respective reference measurements as shown in Figure 3. This is not the case for the other sensors. As  
indicated in Figure 5b-d, the LAB calibrations for the measurements with the NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub> LCSs overestimate the  
285 respective reference measurements. With the exception of the SO<sub>2</sub> measurements, the calibrations by the ML models have a  
bias close to zero, indicating a significant improvement in the agreement between the LCS and reference concentration  
measurements.

When the standard deviation of the calibrated LCS data is lower than that of their respective reference measurements, the point  
associated with the model used for calibration would be illustrated on the left quadrant of the respective target diagram and  
290 vice versa. Therefore, the standard deviations of LAB calibrated measurements from all the LCSs, except the one for CO, are  
higher than those of their respective reference concentrations, whereas ML calibrated concentrations generally have standard  
deviations which are lower than those of the corresponding reference concentrations. All the points corresponding to ML  
models generally lie inside the unit circles, indicating that the variance of residuals of their calibrated concentrations are lower  
than those of the respective reference concentrations and thus that they do not over-fit on the training data.

### 295 **3.2 Feature importance**

To quantify the influence of a number of features that affect the performance of the CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub> LCSs, we carried  
out sensitivity analysis using the RF model, which exhibited the best performance as discussed in the previous sections. To do  
that, the model is trained with the entire 6-month dataset by using reference concentrations, temperature, relative humidity,  
month, day of week and hour as input/independent variables, and the NSS as dependent variables. As stated earlier, the NO<sub>2</sub>  
300 LCS is highly cross-sensitive to O<sub>3</sub> and vice versa (ANN803-05, 2019). To assess the effect of cross-sensitivity on model  
performance, O<sub>3</sub> and NO<sub>2</sub> reference concentrations were included as additional input variables in the training of the RF model  
for the NO<sub>2</sub> and the O<sub>3</sub> LCSs, respectively. The influence of a given input variable to model performance was then calculated  
by permuting its values and determining the change in the performance. For each variable, 20 random permutations were  
carried out and the average change in R<sup>2</sup> was calculated and expressed as a percentage of the total contribution from all the  
305 variables.



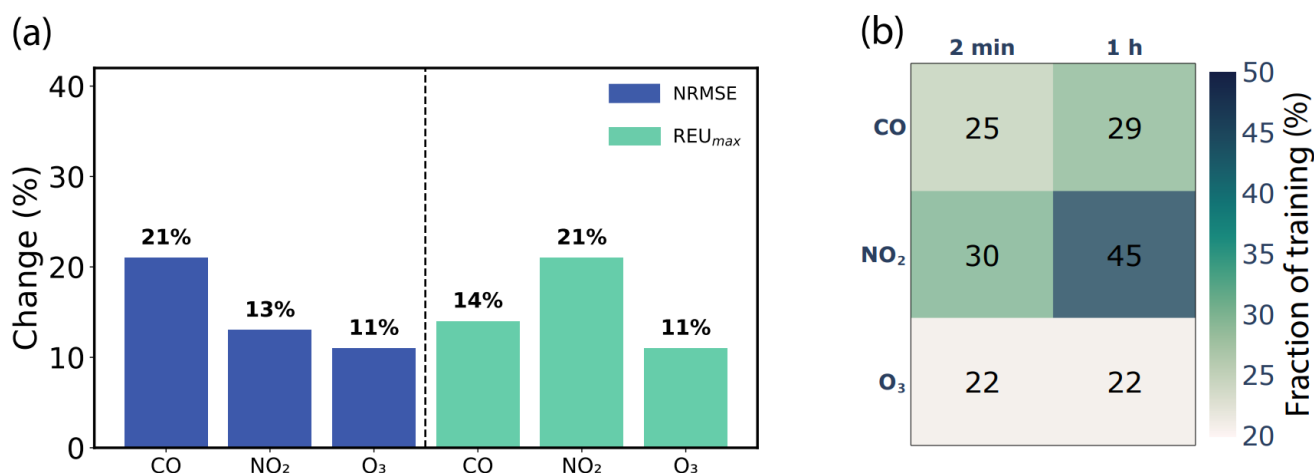
**Figure 6:** Estimated importance of each input feature in determining the variability of the signals reported by CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub> LCSs, determined for each sensor using the RF model. For CO and SO<sub>2</sub> the model is trained using respective reference concentrations of the target gases, temperature, relative humidity, month, day of week and hour as input variables, and the NSS as dependent variable. For NO<sub>2</sub> and O<sub>3</sub>, the training of the RF algorithm also includes respectively, the O<sub>3</sub> and NO<sub>2</sub> reference concentrations as input variables in order to assess cross-sensitivity.

For the CO (see Fig. 6a), the target gas (i.e., the CO reference concentration) has the highest degree of importance in the performance of the model relative to the importance of other variables. This implies that the CO sensor is not affected much by the other variables used as input to the model, and thus responds well to the fluctuations of the CO concentration in ambient air. This is not the case for NO<sub>2</sub> and O<sub>3</sub>. For these two cases, the influence of target gases on the performance of the model seems to be low, with temperature and RH having strong influence. By accounting for cross-sensitivity, the performance of the models further improved by 6-9% (see Fig. 6c and 6d). For the SO<sub>2</sub>, the performance of the model is almost equally affected by the target gas, temperature and RH as well as some of the temporal parameters (see Fig. 6d). This is expected considering that the SO<sub>2</sub> LoD is higher compared to what is observed.



### 320 3.3 Assessing calibration practices

For the evaluation of the models described in subsection 3.1, we have used 80% of the dataset for training the models. While it is a good practice in ML to use higher fraction of data (e.g., > 70 %) for training, it is not practical for calibrating the data from the LCSs as that will require collocating LCSs with corresponding reference instruments for long periods of time, adding significantly to the cost of the measurements. Using RF model, here we investigate how the amount of training data can be minimized, while not significantly sacrificing their quality, by varying practical parameters such as (1) the temporal resolution of the training data, (2) how the training is sampled, and (3) the frequency of calibration.



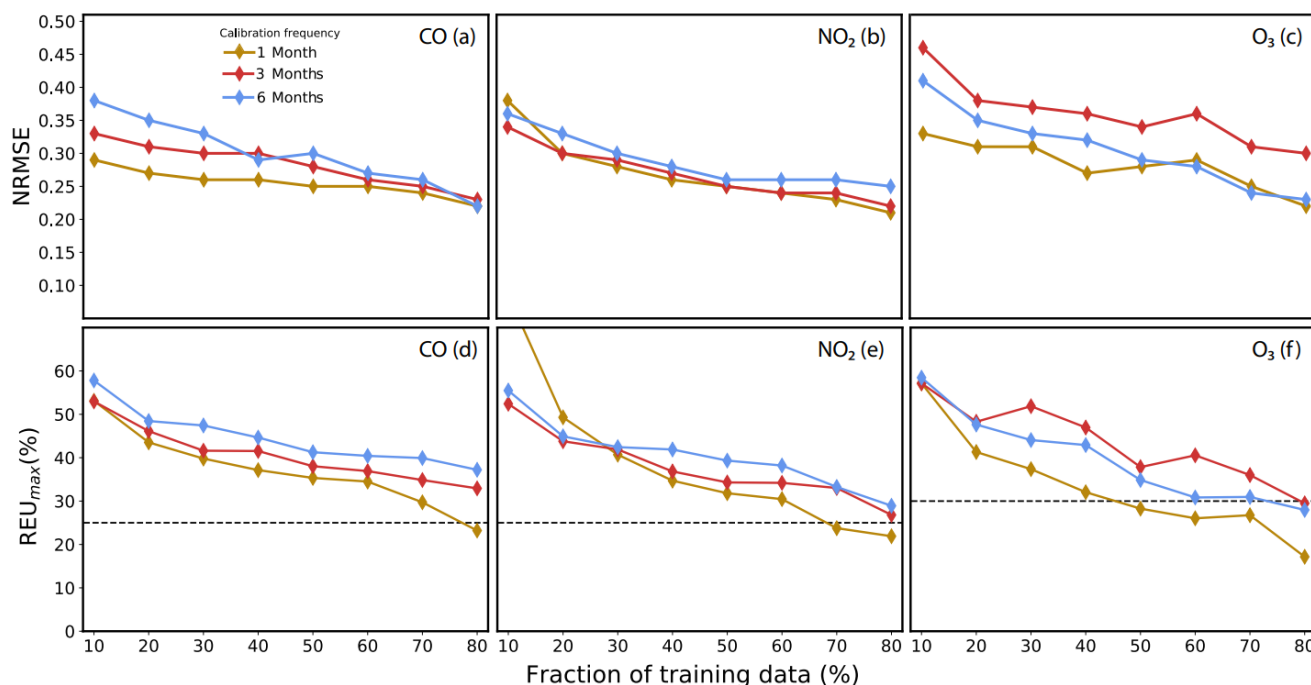
**Figure 7:** (a) Percentage change in the normalized root mean squared error (NRMSE) and the average relative expanded uncertainty calculated at maximum concentration observed (REU<sub>max</sub>) when the resolution of the data used to train the RF model is increase from 1 h to 2 min, and (b) the minimum fraction for 2 min and 1 h measurements required to train the RF model that will guarantee its calibrations to meet the data quality objective for indicative measurements defined by the EU directive.

Figure 7 shows the percentage change in (1) the normalized root mean squared error ( $\Delta$ NRMSE), and (2) the relative expanded uncertainty calculated at maximum concentration observed ( $\Delta$ REU<sub>max</sub>), when the resolution of the data used to train the RF models increases from 1 h to 2 min.  $\Delta$ NRMSE and  $\Delta$ REU<sub>max</sub> were calculated respectively as  $\Delta$ NRMSE =  $\frac{|NRMSE(2min) - NRMSE(1h)|}{NRMSE(1h)} \times 100\%$  and  $\Delta$ REU<sub>max</sub> =  $\frac{|REU_{max}(2min) - REU_{max}(1h)|}{REU_{max}(1h)} \times 100\%$ . The results show that by increasing the temporal resolution of the training data from 1 h to 2 min (which is the base case here), the performance of the RF model in terms of NRMSE and REU<sub>max</sub> improves by 11-21%. As a result, the minimum fraction of data required for training the RF models that will qualify their calibrations for indicative measurements as defined by EU directive reduces by an average of 6.3% (see Fig. 7b)

To investigate the effect of how the training data is sampled and calibration frequency on the amount of data required for training, we performed calibrations using different fractions of the dataset every 1, 3 and 6 months. For 1 month calibration,

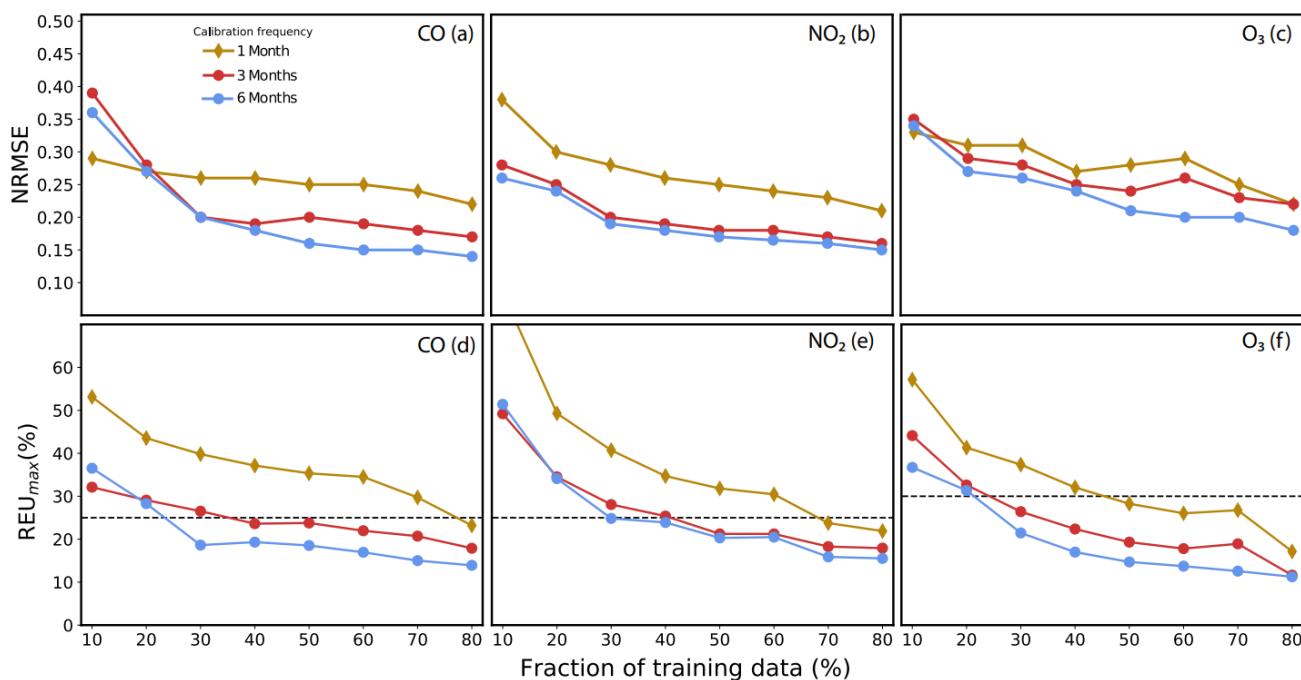


the training data is obtained through continuously sampling, i.e., training data sampled continuously at the start of each month. For 3- and 6-month calibrations, in addition to continuous data sampling, we also considered interceptive data sampling, where the training datasets were sampled from all the months. Note that for the latter case, the post-calibration of the LCSs by applying the trained ML algorithm and the evaluation of the algorithm are conducted retrospectively at the end each period (i.e., every 3 or 6 months for the 3- and 6- month periods, respectively).



**Figure 8:** Average normalized root mean squared error (NRMSE; a-c), and the average relative expanded uncertainty calculated at maximum concentration observed (REU<sub>max</sub>: d-f) for 1, 3 and 6 calibrations obtained by training RF following continuous data sampling (i.e., continuously at the beginning of each period). The dotted lines in the bottom show the EU DQOs for indicative measurements, which are set to 25% for CO and NO<sub>2</sub>, and 30% for O<sub>3</sub>.

Figure 8 shows the average NRMSE and REU<sub>max</sub> for 1-, 3- and 6-months calibrations for CO, NO<sub>2</sub> and O<sub>3</sub> obtained through continuous data sampling and varying the training data from 10 to 80%. The NRMSE and the REU<sub>max</sub> decreases with increase in the fraction of training data. The results also show that under the continuous data sampling scheme, the monthly calibration performing better in terms of NRMSE and REU<sub>max</sub> values compared to calibration after 3- or 6-month. This is explained by the fact that the RF algorithm lacks the power to extrapolate on training data, and thus will perform better in monthly calibration where the seasonal variabilities of the training and testing datasets are lower compared to calibration after 3 or 6 months. It is also evident that under the continuous data sampling, a larger fraction of data (at least 70% for CO and NO<sub>2</sub>, and 50% for O<sub>3</sub>) is needed to train the models for the quality of calibrated data to meet the EU directive DQOs for indicative measurements (see Fig. 8d-f).



**Figure 9:** Average normalized root mean squared error (NRMSE: a-c), and the average relative expanded uncertainty calculated at maximum concentration observed (REU<sub>max</sub>: d-f) for 1, 3 and 6 month calibration periods obtained by training the RF following interceptive data sampling for 3 and 6 month calibrations and comparing with monthly calibration carried out through continuous data sampling. The dotted lines in the bottom show the EU DQOs for indicative measurements, which are set to 25 % for CO and NO<sub>2</sub>, and 30 % for O<sub>3</sub>.

Figure 9 shows the average NRMSE and REU<sub>max</sub> for CO, NO<sub>2</sub> and O<sub>3</sub> obtained by performing 3- and 6-months through interceptive sampling. Note that the results of monthly calibrations are added here for comparison. As shown in the figure, performing calibration after 3 or 6 months using interceptive data sampling scheme requires less training data for achieving a data quality meeting the requirements of EU directive DQOs. For instance, performing 6-month calibration under this scheme will reduce the fraction of training data to as low as 22% without deteriorating the quality of the data to levels below those that will qualify them for indicative measurements. This reduces substantially the cost of these measurements, as significantly less amount of time (from 70-80% down to 22%) is needed for collocated measurements with expensive reference grade instruments.

#### 4 Conclusion

We have evaluated the performance of various ML algorithms for calibrating LCSs data from CO, NO<sub>2</sub>, O<sub>3</sub> and SO<sub>2</sub>, and identified that the Random Forest model provides the best results. This model is then used (1) to investigate how different input variables affect the calibration performance, and (2) to investigate the extent to which practical parameters such as the



temporal resolution of the measurements, how the training data is sampled, and the frequency of calibration reduces the fraction of data needed for training, while keeping the quality of calibrated data within the accepted levels. Our results show that the  
380 CO LCS responds well to the target gas and is not affected much by the other variables such as temperature and RH. This is not the case for the NO<sub>2</sub> and O<sub>3</sub> LCSs. For these two cases, the influence of target gases on the performance of the model seems to be low, while temperature and RH having strong influence on their performance. By increasing the temporal resolution of the training data from 1 h to 2 min, the minimum fraction of data required for training the RF models that will qualify their calibrations for indicative measurements as defined by EU directive reduces by an average of 6.3%. The results  
385 also suggest that the minimum fraction of data required for training the ML models depends on the frequency of carrying out collocated measurements with reference instruments and using the resulting datasets for training the calibration model. If the calibrations are carried out monthly, ca. 50 % of the period is needed for collecting data to train the RF algorithm and qualify the LCSs for indicative measurements as defined by the EU directive (2008/50/EC). If the training is carried out every 3- or 6-months by following continuous data sampling, then ca. 60% of the measuring period is required for collecting training data.  
390 In those cases, if the sampling of the training data is collected over specific periods every month, but the entire training dataset is used to calibrate the measurements over 3 or 6 months, then the amount of data required for qualifying the LCSs for indicative measurements can reduce to 22 %. This reduces substantially the costs, as the amount of time needed for collocated measurements with reference grade instruments significantly reduce from 70-80% to 22%. This not only improves the quality of measurements obtained by LCS networks but assures quality in the most cost-effective manner.

#### 395 **Code and data availability**

Datasets used for this publication are available to the community here: [10.5281/zenodo.18629746](https://zenodo.org/doi/10.5281/zenodo.18629746)

#### **Author contributions**

Conceptualization: V.L., S.B., P.K., G.B.; Methodology: V.L., P.K., S.B., G.B; Software: V.L., R.P., G.I., S.B.;  
Investigation: V.L., G.I., S.B., G.B.; Data curation: V.L., R.P., G.I.; Writing – original draft preparation: V.L., G.I., P.K., G.B;  
400 Writing – review & editing: V.L., G.I., S.B., P.K., G.B; Supervision: P.K., G.B.; All authors have agreed to the version of the manuscript.

#### **Competing interests**

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



## 405 Acknowledgements

PK acknowledges the support received from the UKRI (EPSRC, NERC, AHRC) funded RECLAIM Network Plus ([EP/W034034/1](https://doi.org/10.1039/C9EM00343A)) and the NERC-funded Green Cities project ([NE/X002799/1](https://doi.org/10.1039/C9EM0027991)) projects.

## References

- Alphasense-SO2-B4: SO2-B4 Sulfur Dioxide Sensor Technical Specification 4-Electrode,  
410 <https://www.alphasense.com/wp-content/uploads/2019/09/SO2-B4.pdf>, 2019.
- ANN803-05: Correcting for background currents in four electrode toxic gas sensors (2019). Alphasense Application Note AAN 803-05, 2019.
- Arroyo, P., Gómez-Suárez, J., Suárez, J. I., and Lozano, J.: Low-Cost Air Quality Measurement System Based on Electrochemical and PM Sensors with Cloud Connection, *Sensors*, 21, 6228, 2021.
- 415 Baranwal, J., Barse, B., Gatto, G., Broncova, G., and Kumar, A.: Electrochemical sensors and their applications: A review, *Chemosensors*, 10, 363, 2022.
- Chen, C.-C., Kuo, C.-T., Chen, S.-Y., Lin, C.-H., Chue, J.-J., Hsieh, Y.-J., Cheng, C.-W., Wu, C.-M., and Huang, C.-M.: Calibration of low-cost particle sensors by using machine-learning method, in: 2018 IEEE Asia Pacific Conference on Circuits and Systems (APCCAS), pp. 111–114, IEEE, 2018.
- 420 Chen, T. and Guestrin, C.: Xgboost: A scalable tree boosting system, in: Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining, pp. 785–794, 2016.
- Equivalence: Guide to the demonstration of equivalence of ambient air monitoring methods, 2010.
- EU-directive: Directive 2008/50/EC of the European Parliament and of the Council of 21 May 2008 on ambient air quality and cleaner air for Europe, *Official Journal of the European Union*, 2008.
- 425 Ferrer-Cid, P., Barcelo-Ordinas, J. M., Garcia-Vidal, J., Ripoll, A., and Viana, M.: Multisensor data fusion calibration in IoT air pollution platforms, *IEEE Internet of Things Journal*, 7, 3124–3132, 2020.
- Garbagna, L., Babu Saheer, L., & Maktab Dar Oghaz, M. (2025). AI-driven approaches for air pollution modelling: A comprehensive systematic review. In *Environmental Pollution* (Vol. 373). Elsevier Ltd. <https://doi.org/10.1016/j.envpol.2025.125937>
- 430 Gonzalez, A., Boies, A., Swason, J., and Kittelson, D.: Field calibration of low-cost air pollution sensors, *Atmospheric Measurement Techniques Discussions*, pp. 1–17, 2019.
- Heal, M. R., Kumar, P., & Harrison, R. M. (2012). Particles, air quality, policy and health. *Chemical Society Reviews*, 41(19), 6606–6630. <https://doi.org/10.1039/c2cs35076a>



- 435 Isaac, N., Pikaar, I., and Biskos, G.: Metal oxide semiconducting nanomaterials for air quality gas sensors: operating principles, performance, and synthesis techniques, *Microchimica Acta*, 189, 196, 2022.
- Jolliff, J. K., Kindle, J. C., Shulman, I., Penta, B., Friedrichs, M. A., Helber, R., and Arnone, R. A.: Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment, *Journal of Marine Systems*, 76, 64–82, 2009.
- Koziel, S., Pietrenko-Dabrowska, A., Wojcikowski, M., & Pankiewicz, B. (2024). High-performance machine-learning-based calibration of low-cost nitrogen dioxide sensor using environmental parameter differentials and global data scaling. *Scientific Reports*, 14(1). <https://doi.org/10.1038/s41598-024-77214-y>
- 440 Kumar, P., Morawska, L., Martani, C., Biskos, G., Neophytou, M., Di Sabatino, S., Bell, M., Norford, L., and Britter, R.: The rise of low-cost sensing for managing air pollution in cities, *Environment international*, 75, 199–205, 2015.
- Kumar, V. and Sahu, M.: Evaluation of nine machine learning regression algorithms for calibration of low-cost PM<sub>2.5</sub> sensor, *Journal of Aerosol Science*, 157, 105–110, 2021.
- 445 Lewis, A., Peltier, W. R., and von Schneidemesser, E.: Low-cost sensors for the measurement of atmospheric composition: overview of topic and future applications, 2018.
- Maag, B., Saukh, O., Hasenfratz, D., and Thiele, L.: Pre-Deployment Testing, Augmentation and Calibration of Cross-Sensitive Sensors., in: *EWSN*, pp. 169–180, 2016.
- Mahajan, S., & Kumar, P. (2020). Evaluation of low-cost sensors for quantitative personal exposure monitoring. *Sustainable Cities and Society*, 57. <https://doi.org/10.1016/j.scs.2020.102076>
- 450 Masic, A., Bibic, D., Pikula, B., and Razic, F.: NEW APPROACH OF MEASURING TOXIC GASES CONCENTRATIONS: PRINCIPLE OF OPERATION, *Annals of DAAAM & Proceedings*, 29, 2018.
- Nowack, P., Konstantinovskiy, L., Gardiner, H., and Cant, J.: Machine learning calibration of low-cost NO<sub>2</sub> and PM<sub>10</sub> sensors: non-linear algorithms and their impact on site transferability, *Atmospheric Measurement Techniques*, 14, 5637–5655, 2021.
- 455 Okafor, N. U. and Delaney, D. T.: Application of Machine Learning Techniques for the Calibration of Low-cost IoT Sensors in Environmental Monitoring Networks, in: *2020 IEEE 6th World Forum on Internet of Things (WF-IoT)*, pp. 1–3, IEEE, 2020.
- Pang, X., Shaw, M. D., Gillot, S., and Lewis, A. C.: The impacts of water vapour and co-pollutants on the performance of electrochemical gas sensors used for air quality monitoring, *Sensors and Actuators B: Chemical*, 266, 674–684, 2018.
- 460 Pang, X., Shaw, M. D., Lewis, A. C., Carpenter, L. J., and Batchellier, T.: Electrochemical ozone sensors: A miniaturised alternative for ozone measurements in laboratory experiments and air-quality monitoring, *Sensors and Actuators B: Chemical*, 240, 829–837, 2017.



- 465 Papaconstantinou, R., Bezantakos, S., Pikridas, M., Parolin, M., Stylianou, M., Savvides, C., Sciare, J., & Biskos, G. (2026). Comparing spatial and temporal variabilities between the Vaisala AQT530 monitor and reference measurements. *Atmospheric Measurement Techniques*, 19(1), 63–78. <https://doi.org/10.5194/amt-19-63-2026>
- Papaconstantinou, R., Demosthenous, M., Bezantakos, S., Hadjigeorgiou, N., Costi, M., Stylianou, M., Symeou, E., Savvides, C., and Biskos, G.: Field evaluation of low-cost electrochemical air quality gas sensors under extreme temperature and relative humidity conditions, *Atmospheric Measurement Techniques*, 16, 3313–3329, 2023.
- 470 Patra, S. S., Ramsisaria, R., Du, R., Wu, T., and Boor, B. E.: A machine learning field calibration method for improving the performance of low-cost particle sensors, *Building and Environment*, 190, 107 457, 2021.
- Podder, I., Fischl, T., & Bub, U. (2024). Smart calibration and monitoring: leveraging artificial intelligence to improve MEMS-based inertial sensor calibration. *Complex and Intelligent Systems*, 10(6), 7451–7474. <https://doi.org/10.1007/s40747-024-01531-y>
- 475 Rai, A. C., Kumar, P., Pilla, F., Skouloudis, A. N., Di Sabatino, S., Ratti, C., Yasar, A., & Rickerby, D. (2017). End-user perspective of low-cost sensors for outdoor air pollution monitoring. In *Science of the Total Environment (Vols. 607–608, pp. 691–705)*. Elsevier B.V. <https://doi.org/10.1016/j.scitotenv.2017.06.266>
- Samad, A., Obando Nuñez, D. R., Solis Castillo, G. C., Laquai, B., and Vogt, U.: Effect of relative humidity and air temperature on the results obtained from low-cost gas sensors for ambient air quality measurements, *Sensors*, 20, 5175, 480 2020.
- Schäfer, K., Lande, K., Grimm, H., Jenniskens, G., Gijsbers, R., Ziegler, V., Hank, M., and Budde, M.: High-Resolution Assessment of Air Quality in Urban Areas? A Business Model Perspective, *Atmosphere*, 12, 595, 2021.
- Si, M., Ying, X., Du, S., and Du, K.: Evaluation and Calibration of a Low-cost Particle Sensor in Ambient Conditions Using Machine Learning Technologies 2, 2020.
- 485 Song, J., Han, K., and Stettler, M. E.: Deep-MAPS: Machine-Learning-Based Mobile Air Pollution Sensing, *IEEE Internet of Things Journal*, 8, 7649–7660, 2020.
- Sousan, S., Wu, R., Popoviciu, C., Fresquez, S., & Park, Y. M. (2025). Advancing low-cost air quality monitor calibration with machine learning methods. *Environmental Pollution*, 374. <https://doi.org/10.1016/j.envpol.2025.126191>
- Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Field calibration of a cluster of low-cost 490 available sensors for air quality monitoring. Part A: Ozone and nitrogen dioxide, *Sensors and Actuators B: Chemical*, 215, 249–257, 2015.
- Thunis, P., Pederzoli, A., and Pernigotti, D.: Performance criteria to evaluate air quality modeling applications, *Atmospheric Environment*, 59, 476–482, 2012.
- Tiku, M. L. and Akkaya, A. D.: Robust estimation and hypothesis testing, New Age International, 2004.



- 495 Vajs, I., Drajić, D., Gligorić, N., Radovanović, I., and Popović, I.: Developing Relative Humidity and Temperature Corrections for Low-Cost Sensors Using Machine Learning, *Sensors*, 21, 3338, 2021.
- Walker, S.-E. and Schneider, P.: A study of the relative expanded uncertainty formula for comparing low-cost sensor and reference measurements, NILU rapport, 2020.
- Wang, C., Wu, Q., Weimer, M., and Zhu, E.: FLAML: A fast and lightweight automl library, *Proceedings of Machine Learning and Systems*, 3, 434–447, 2021.
- 500 Yu, L., Lai, K. K., Wang, S., and Huang, W.: A bias-variance-complexity trade-off framework for complex system modeling, in: *International Conference on Computational Science and Its Applications*, pp. 518–527, Springer, 2006.
- Zimmerman, N., Presto, A. A., Kumar, S. P., Gu, J., Hauryliuk, A., Robinson, E. S., Robinson, A. L., and Subramanian, R.: A machine learning calibration model using random forests to improve sensor performance for lower-cost air quality monitoring, *Atmospheric Measurement Techniques*, 11, 291–313, 2018.
- 505 Zimmerman, N.: Tutorial: Guidelines for implementing low-cost sensor networks for aerosol monitoring, *Journal of Aerosol Science*, 159, 105 872, 2022.
- Zuidema, C., Schumacher, C. S., Austin, E., Carvlin, G., Larson, T. V., Spalt, E. W., Zusman, M., Gasset, A. J., Seto, E., Kaufman, J. D., et al.: Deployment, Calibration, and Cross-Validation of Low-Cost Electrochemical Sensors for Carbon Monoxide, Nitrogen Oxides, and Ozone for an Epidemiological Study, *Sensors*, 21, 4214, 2021.
- 510