# Performance evaluation of air quality sensors for environmental epidemiology

Miriam Chacón-Mateos [a,1], Héctor García-Salamero [a], Bernd Laquai [a], Ulrich Vogt [a]

[a] University of Stuttgart, Institute of Combustion and Power Plant Technology, Department of Flue Gas Cleaning and Air
5   Quality Control, Stuttgart, 70569, Germany

[1] permanent address: German Aerospace Center, Institute of Combustion Technology, Stuttgart, 70569, Germany

**Correspondence**: Miriam Chacón-Mateos (miriam.chacon-mateos@ifk.uni-stuttgart.de)

**Abstract.** Over the past few decades, the study and the use of air quality sensors have significantly increased, leading to a wealth of experience and a deeper understanding of their strengths and limitations. This study aimed to transcend the limitations
10   by developing and evaluating a methodology for $PM_{2.5}$ and $NO_2$ sensors to enhance sensor accuracy to a level suitable for epidemiological studies, where ensuring data quality is paramount. The performance evaluation of indoor and outdoor sensors was carried out during the co-location phase with reference instruments (RIs), by calculating common error metrics, target diagrams and the relative expanded uncertainties (REUs) stated in the EU Air Quality Directive 2008/50/EC and the recently published EU Directive 2024/2881, before the deployment of the air quality sensor systems (AQSSs) in the houses of patients
15   suffering from chronic obstructive pulmonary disease (COPD) or asthma in Stuttgart (Germany). Regression and machine learning models for sensor calibration were tested during the co-location. Moreover, an original methodology was designed and evaluated to validate the sensor data during the epidemiological study. The study found that indoor sensor calibration using artificially generated $NO_2$ and aerosols does not ensure model transferability, emphasizing the need for training data that matches the intended deployment environment in terms of real patterns of concentration, particle composition and
20   environmental conditions. Integrating metadata such as activity logs, window status, and data from official monitoring stations, proved essential for data validation and interpretation during the sensor deployment in the houses of the participants. Despite limitations at low pollutant levels, calibrated AQSSs are a promising tool to increase the ubiquity of epidemiological studies for low- and middle-income countries or regions where higher air pollutant concentrations are expected.

**Keywords** Low-cost sensors; Air quality; $PM_{2.5}$; $NO_2$; Epidemiological studies; Measurement uncertainty

25   ## 1 Introduction

The World Health Organization (WHO) updated its global air quality guidelines in September 2021. The new air quality recommendations proposed by the WHO resulted from the findings based on recent epidemiological studies. The increase in evidence on the adverse health effects of air pollution has been possible thanks to the advances in technology for air pollution

monitoring and personal exposure (WHO, 2021). A major air pollutant is particulate matter (PM), especially the fraction $PM_{2.5}$,

30  which can cause respiratory and cardiovascular diseases, reproductive and central nervous system dysfunctions, and cancer (Manisalidis et al., 2020). In a meta-analysis, Braithwaite et al. (2019) also found statistically significant associations between long-term $PM_{2.5}$ exposure and mental illnesses such as depression and anxiety. Another air pollutant of special interest is $NO_2$, which has been associated with higher morbidity for vulnerable groups such as asthma and chronic obstructive pulmonary disease (COPD) patients (Hoffmann et al., 2022). Moreover, a recent review paper has shown that both short- and long-term

35  exposure to $PM_{2.5}$ or $NO_2$ adjusted for $NO_2$ and $PM_{2.5}$, respectively, revealed a synergistic effect appearing as higher mortality from respiratory diseases (Mainka and Żak, 2022).

Exposure measurements are carried out using direct or indirect approaches. The direct approaches measure the exposure levels by using personal passive samplers (Piechocki-Minguy et al., 2006; Shirdel et al., 2019; Samon et al., 2022) or mobile monitors (Rea et al., 2001; Koehler et al., 2019) that must be worn by the person during the campaign. In recent years more studies have

40  deployed air quality sensors allowing multi-pollutant exposure assessment (Piedrahita et al., 2014; Chatzidiakou et al., 2020; Novak et al., 2021). This methodology is considered the most accurate estimate of a person's 'true' exposure. However, this type of personal exposure assessment is only adequate for short-term exposure (Steinle et al., 2013). The main challenges of these studies are the complexity of the data integration including the time-activity-location profiles (Chatzidiakou et al., 2022), and the measurement uncertainty due to the position of the sampling inlet, which may be largely affected by the

45  perihuman/personal cloud effect (Licina et al., 2017; Pantelic et al., 2020). In theory, the sampling inlet should be placed close to the breathing zone, but this is in reality not always feasible, especially for multi-pollutant devices. Additional factors, such as vibrations, static electricity (Shirdel et al., 2019) and movement (e.g. isokinetic sampling of PM cannot not be guaranteed), have also an influence on the accuracy of the instrumentation. Moreover, other external factors like the accuracy of the GPS signal, the accelerometer, etc. may be crucial to characterize the true exposure.

50  The indirect approaches measure air quality at fixed monitoring sites or are based on modelling (Goldman et al., 2012; Beloconi and Vounatsou, 2020; Huang et al., 2021) which can also integrate satellite data (Hang et al., 2022). Among the indirect approaches, some studies rely on outdoor measurements at fixed-site monitoring stations (Harré et al., 1997; Meng et al., 2013). This has been the cause of exposure misclassification in the past (Shaw et al., 2018), as outdoor monitoring stations fail to capture the real concentrations in the different microenvironments an individual is exposed to (Krause, 2021). Moreover,

55  strong correlations among the ambient pollutants can lead to biased health effect estimates due to confounding (Sarnat et al., 2001). Other indirect approaches are based on static measurements in the most visited microenvironments of the participants (Scott Downen et al., 2022). The main advantage of this methodology is the lower effort required of the participant which allows longer measurement periods, making it the ideal candidate for long-term exposure assessment (Steinle et al., 2013).

The use of air quality sensors for environmental epidemiology has many advantages, for instance, the decrease in the bias of

60  exposure estimations when compared with fixed outdoor monitoring stations (Chatzidiakou et al., 2019). Another benefit of using sensors is the possibility of increasing the number of participants with the same fixed budget, which helps to ensure adequate statistical power of the study. Moreover, sensors allow time resolutions in the order of seconds, making possible the

exposure assessment on movement and the correlation of pollution patterns with personal behaviour when this information also exists (Jerrett et al., 2017; Novak et al., 2022).

65   On the other side, some characteristics of the sensors have kept them away from applications where high accuracy is required. One of them is the influence of meteorological conditions such as temperature (T) and relative humidity (RH) and cross-sensitivities in the sensor signal (Samad et al., 2020; Venkatraman Jagatha et al., 2021; Zamora et al., 2022). That makes the calibration of the sensors more complex than traditional monitoring devices, as the correction algorithms should account also for those influences, and that limits the transferability of the calibration models when moving the sensor to a different location

70   (Zauli-Sajani et al., 2021; Diez et al., 2024). Another parameter that affects the sensor accuracy for long-term measurements is the signal drift caused by the sensor degradation (Tancev, 2021; deSouza et al., 2023). Last but not least, the unit-to-unit variability poses a challenge when it comes to calibrating many units at the same time, as is the case for epidemiological studies (Gäbel et al., 2022).

Some recent studies have shown that the above-mentioned concerns can be overcome and that getting highly personalized air

75   pollution exposure outweighs the measurement uncertainty of the air quality sensors. The AIRLESS study (Effects of AIR pollution on cardiopuLmonary disEaSe in urban and peri-urban reSidents in Beijing) demonstrated that sensing technologies can revolutionise health studies and address scientific, health and policy questions in a way that has not been possible before (Krause, 2021). The results of the AIRLESS project have been well-documented (Chatzidiakou et al., 2020; Krause, 2021) and are a prove of the potential use of sensing technologies for epidemiological studies in very different environments, i.e.

80   high- and middle-income countries like London (Evangelopoulos et al., 2021) and Beijing (Han et al., 2020; Han et al., 2021) but also low-income countries like Kenya (Krause, 2021).

This work aims to evaluate the performance of $NO_2$ and $PM_{2.5}$ sensors for their use in health research. We present an approach to calibrate the sensors based on co-location with reference instruments (RIs) and assess the reliability of the different correction models before and during deployment. The sensors were deployed in the houses of seven COPD and asthma patients.

85   The measurements were conducted in two microenvironments per participant representing the outdoor and the indoor levels of exposure for one month. The multiple linear regression (MLR) and three machine learning (ML) models (random forest, support vector regressor as well as artificial neural network) have been evaluated for indoor and outdoor calibration of $NO_2$ sensors. A univariate linear regression (ULR) to correct the data of $PM_{2.5}$ sensors for indoor measurements as well as the ULR and a low-cost dryer attached to the inlet of outdoor $PM_{2.5}$ sensors have been investigated. The performance evaluation has

90   been carried out using common error metrics, REUs according to the European data quality objectives (DQOs), and target diagrams. Finally, we discuss the capabilities as well as the limitations of the proposed methodology.

## 2 Methodology

### 2.1 Participant recruitment and study protocol

The participants consisted of seven patients suffering from COPD or asthma. All the participants lived in Stuttgart (Germany)
95 (see Fig. S1 in the Supplement) and agreed to perform the measurements in their homes for 30 days. One participant agreed to have two outdoor air quality sensor systems (AQSSs) instead of one to study the difference between street-side and garden-side concentrations. For this participant, the measuring campaign was reduced to 19 days.

The study protocol was evaluated and approved by the Ethics Committee of the Medical Association of the State of Baden-Württemberg (reference number F-2019-105) and by the data protection officer of the University of Stuttgart. Before the
100 beginning of the measurements, participants were informed about the study and requested to provide written consent. The participants are referred to by a patient identification number from P1 to P7. An environmental questionnaire in the German language was designed to characterize the living area, the house, and the habits, and was completed prior to the measurements with the help of a worker from the University of Stuttgart. Participants also completed a spirometry test, a health survey assessing their symptoms, and a logbook with hourly information on indoor activities every day. At the end of the
105 measurements, we asked the participants for written feedback. Participants who started the study before March 2020 received the study indications at their homes. However, those who started the study after the COVID-19 outbreak performed the interview by phone, and the contact between the participants and the university staff was kept to a minimum. A detailed description of the data collected and the further analysis to determine the feasibility of using the developed sensor systems and methodology for exposure assessment and indoor source apportionment can be read in Chacón-Mateos et al. (2024).

110 ### 2.2 Indoor and outdoor air quality monitors

Two different AQSSs for indoor and outdoor measurements were designed for this study (see Fig. 1), each one containing one electrochemical sensor for $NO_2$ (Alphasense, UK, model B43F), and one optical particle counter for $PM_{2.5}$ (Alphasense, UK, model OPC-R1). The sensor selection was based on our own tests of different sensors in the laboratory. Another important factor that was considered was the price, being 150 Euro the maximum possible price per sensor. Additionally, a T and RH
115 sensor was included (IST AG, Switzerland, model HTY221). The microcontroller Arduino UNO was used to control and save the data every two seconds on an SD card. During the deployment, participants did not have access to the data in order to avoid behavioural changes.

As an outdoor AQSS must be weather resistant, we selected an enclosure made of glass fibre-reinforced polyester with the following dimensions: 200×300×150 mm. For the indoor AQSS, a polypropylene box with the dimensions 240×195×112 mm
120 was chosen. The cost of the materials amounted to a total of 540 and 460 Euro for the outdoor and indoor AQSSs, respectively. To counteract the effect of the high RH in the PM sensor readings, a low-cost dryer was designed for the outdoor PM sensor. The main advantage of using a low-cost dryer is that it allows the use of the same correction models independently of the location of the PM sensor. Other techniques based on the κ-Köhler theory or machine learning have shown incorrect results

when moving the sensor to another location, as the particle composition may differ from the one in the co-located site (Di

125 Antonio et al., 2018; deSouza et al., 2022). The dryer consists of a 50 cm brass tube with a resistive wire wound around its surface. The wire is heated when the RH is higher than 70 % using 12 V and 10 W. The T is controlled by using the internal T sensor of the OPC-R1. A detailed description and evaluation of the low-cost dryer can be read in Chacón-Mateos et al. (2022).
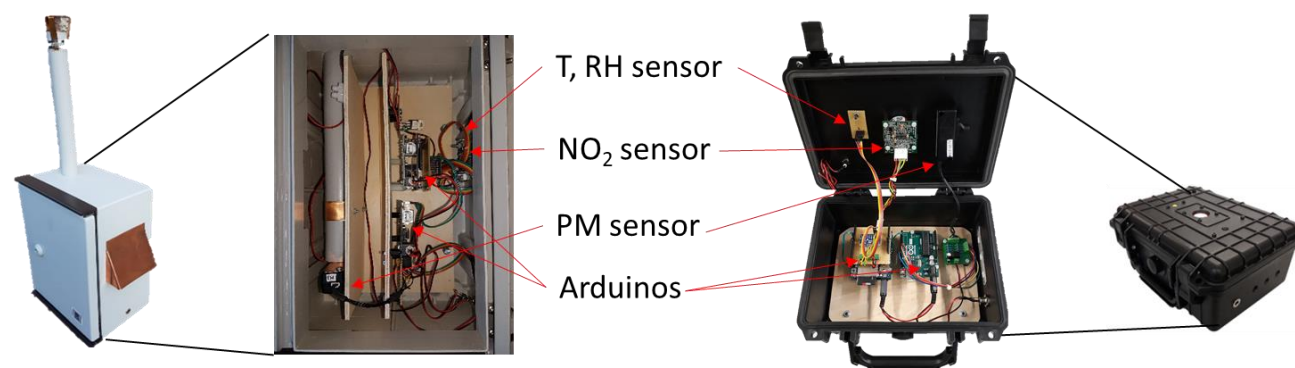


130
**Figure 1.** Designed AQSSs for outdoor (left) and indoor (right) measurements (Chacón-Mateos et al., 2024).

## 2.2 Study design and quality assurance

The measurements in the houses of the patients took place in Stuttgart region (Germany) between 20 December 2019 and 28 May 2020. Figure S1 shows the approximate locations of the participants' homes, the governmental outdoor air quality

135 monitoring stations in Stuttgart and the monitoring stations of the University of Stuttgart. The calibration of the indoor and outdoor AQSSs took place discontinuously starting on 7 November 2019 and finishing on 5 June 2020. The sensor calibration was done weeks before the individual deployment in houses or immediately after it. A general overview of the measuring campaign showing the periods where the calibration and deployment of the AQSSs took place can be seen in Fig. S2. In the following sub-sections, a detailed description of the methodology used to verify and assess the quality of the data before and

140 after the deployment in the homes of the patients is described.

### 2.3.1 Sensor calibration before deployment

The calibration procedures for both indoor and outdoor AQSSs were conducted in distinct locations to replicate real-world environmental conditions as closely as possible. Likewise, the methodology was tailored to address the specific conditions encountered in indoor as well as outdoor environments. The main objective was to cover the maximum range of possible

145 concentrations, T and RH that could be found later in the indoor and outdoor locations. A summary of the different procedures can be seen in Table 1.

Before deployment, the $NO_2$ sensors for indoor measurements were co-located in the laboratory for a minimum of seven days and a maximum of 34 days, depending on the availability of the AQSSs. A chemiluminescence device (MLU, Austria, model

200A) was used as RI for NO$_2$. Due to the low NO$_2$ concentrations measured in the laboratory, it was necessary to generate

150    higher concentrations using a Gas Phase Titration system (GPT) (Ecotech, Australia, model Serinus Cal 3000). For this

purpose, the indoor AQSSs were placed inside a sealed box made of inert glass with gas supply connections. The dimensions

of the box were as follows: $310 \times 525 \times 375$ mm. The sensors were exposed to the following pyramid of NO$_2$ concentrations:

0 - 50 - 100 - 50 - 0 ppb. Each stage was maintained for 3 hours and the pyramid was repeated at least twice in different days.

The changes in the T and RH were forced using an infrared lamp close to the calibration box and an air humidifier, respectively.

155    Moreover, natural changes in the room conditions were simulated by opening and closing the windows several times a day.

**Table 1.** Calibration methodology of the NO$_2$ and PM$_{2.5}$ sensors.

| Pollutant | Indoor AQSSs | Outdoor AQSSs |
|---|---|---|
| NO$_2$ | Co-location in the laboratory:<br>- Low concentrations: indoor air.<br>- High concentrations (up to 180 µg m$^{-3}$): artificial generated NO$_2$.<br>- Changes in T using an infrared lamp.<br>- Changes in the RH by manually opening and closing the windows and using an air humidifier. | Co-location at Hauptstätter Street monitoring station. |
| PM$_{2.5}$ | Co-location:<br>- Low concentrations: indoor air.<br>- High concentrations (up to 150 µg m$^{-3}$): calibration powders in particle chamber. | Co-location at Hauptstätter Street monitoring station. |

The calibration of PM$_{2.5}$ sensors was performed in a particle chamber. High particle concentrations (up to 150 µg m$^{-3}$) with a

160    peak concentration at an aerodynamic diameter of less than 3 µm were dispersed using an aerosol generator and liquid paraffin.

The complete indoor sensor system was placed inside the chamber to correct the possible influence of the material on the

particle concentration measured by the PM sensor. The RI was a light-scattering device (Grimm, Germany, model 1.108). The

experiments in the particle chamber took about an hour and were repeated at least twice. More information about this

calibration setup can be read in Laquai and Saur (2017). After the calibration in the particle chamber, the sensors were co-

165    located in the laboratory for some days to expose the sensors to real PM indoors.

The calibration of the outdoor AQSSs required less effort as it consisted of seven to 34 days of co-location in the hotspot

monitoring station at Hauptstätter Street (48°45´55.8936´´ N, 9°10´12.9396´´ E), Stuttgart. The average co-location time was

15 days. As RI for NO$_2$, the model 405 nm from the company 2B Technologies (USA) was used. An EDM 180 from the

company Grimm GmbH (Germany) was used as an RI for the PM measurements. The RI for NO$_2$ was calibrated once a month

170    and the measurements of the Grimm EDM 180 were corrected against gravimetric measurements at the beginning of the

campaign.

During the measurement campaign and after having analysed the first results, we decided to experiment with a new calibration

strategy: for patient P7 an outdoor AQSS (B03) calibrated in the Hauptstätter Street monitoring station was used for indoor air

175 quality measurements. The reason for that was the high deviation of the indoor $NO_2$ concentrations modelled by the support vector regressor (SVR) and random forest regressor (RFR) models when compared to the results of the passive samples located in the same place (See Section 3.2.1).

### 2.3.2 Data validation during deployment

Due to the data reliability problems that air sensors have, it is vital to be able to identify if the AQSSs are working properly during the deployment in the houses of the patients. In an ideal case, having an RI co-located would be the best option.

180 However, this is usually not possible for epidemiological studies with a lot of participants. For that reason, we present here a methodology that can be used in epidemiological studies having a high number of participants. This novel approach has been summarized in Table 2.

**Table 2.** Validation of the $NO_2$ and $PM_{2.5}$ sensors during deployment.

| Pollutant | Indoor AQSSs | Outdoor AQSSs |
|---|---|---|
| $NO_2$ | Comparison with passive samples (quantitative). | Comparison with passive samples (quantitative). Comparison with outdoor air quality monitoring stations less than 6 km apart (qualitative). |
| $PM_{2.5}$ | Identification of possible sources of peak concentrations using the activity log. | Comparison with outdoor air quality monitoring stations less than 6 km apart (qualitative). |

185

To have a reference $NO_2$ concentration value in the houses, $NO_2$ passive samples (diffusion tubes) from the company Passam (Switzerland) were attached to the indoor and outdoor AQSSs to perform discontinuous measurements. In this technique, $NO_2$ is absorbed in a metal mesh that has been treated with triethanolamine. After 14 days of exposure time, the passive samples were collected and analysed in the laboratory as described in VDI 2453 Part 1 (1990). Thanks to this, it was possible to compare

190 the sensor data to the concentration measured by the passive samples to identify which models were under or over-estimating the $NO_2$ concentrations during deployment. For patients P2 and P4 only one period was collected of 14 and 19 days, respectively. Additionally, the data of the nearest governmental monitoring station was also collected to qualitatively compare the trend of the official monitoring stations to the $NO_2$ sensors used outdoors.

The $PM_{2.5}$ concentrations measured by the sensors indoors were validated using the activity log to identify the possible sources

195 of $PM_{2.5}$. In the case of the outdoor $PM_{2.5}$ sensors, a comparison with a governmental monitoring station nearby was carried out.

### 2.3 Sensor correction models

In this section, the correction models used for $PM_{2.5}$ and $NO_2$ sensors are described. For PM sensors the univariate linear regression (ULR) shown in Eq. 1 was used,

$$PM_{2.5,corrected} = \beta_0 + \beta_1 PM_{2.5,raw}$$

(1)

7

200    where $\beta_0$ is the calibration constant and $\beta_1$ the calibration factor of the linear fitting between the PM$_{2.5}$ concentrations of the sensor and the RI. The use of a low-cost dryer prevents the outdoor PM sensor readings from the influence of hygroscopic growth of PM when the RH is higher than 70 %. The indoor PM sensor was also calibrated using ULR and it did not include the low-cost dryer as RH higher than 70 % indoors was not expected. During the deployment, we measured indoor RH between 18 to 58 %.

205    For NO$_2$ sensors, different parametric and non-parametric models were investigated to take into account the influence of the RH and the T in the sensor signal: multiple linear regression (MLR), random forest regressor (RFR), support vector regressor (SVR) and artificial neural network (ANN). These models have been already investigated to correct the data of air quality sensors with promising results (Esposito et al., 2016; Topalović et al., 2019; Zimmerman et al., 2018; Bigi et al., 2018) but literature about how these models perform when the sensor is transferred to a new location is scarce.

210    The explanatory variables (also called features in ML models) used for all the models were data from the working (WE) and the auxiliary (AE) electrodes, and the T and RH of the HYT221 sensor. The MLR model shown in Eq. 2 is applied to correct the NO$_2$ sensor data. In Eq. 2, $\alpha_0$ is the intercept and $\alpha_n$ are the coefficients that applied to each explanatory variable.

$$NO_{2,corrected} = \alpha_0 + \alpha_1 WE + \alpha_2 AE + \alpha_3 T + \alpha_4 RH \qquad (2)$$

The RFR is an ML algorithm based on ensembles of decision trees (Breiman, 2001). The main characteristics are that it randomizes both the selection of the data points used to build the trees and the explanatory variables at each node to determine

215    the split. Thus, leading to each decision tree being built on a slightly different dataset with a different subset of features (Müller and Guido, 2017). During prediction, the RFR calculates the average of the predicted values from all the decision trees, resulting in a more accurate prediction than a single decision tree. The RFR is known for its ability to handle noisy and complex data while reducing overfitting and improving model performance.

The SVR models come originally from support vector machine algorithms, which are usually used for classification purposes

220    (Boser et al., 1992). In SVR, instead of trying to minimize the residuals between the predicted values and the actual values using the conventional sum of the squared residuals of a linear fitting, the goal is to find a margin that includes as many data points as possible within a certain distance, also called epsilon (ε), from the predicted values. To achieve that, a hyperplane in a high-dimensional feature space, i.e. a function, must be found, so that the threshold distance of the ε-tube between the hyperplane and the support vectors is maximized while the errors of the predicted values are minimized. The support vectors

225    are the data points that lie either on the edge of the ε-tube or violate the margin constraints (Awad and Khanna, 2015). This model is very robust in handling outliers.

The ANN is an ML algorithm inspired by the connections among the cells of the nervous system (McCulloch and Pitts, 1943). In this model, the training data containing the explanatory variables are inserted as input nodes in the network. This input is used in the first step, called forward propagation, to estimate the value of the parameters (biases and weights). These parameters

230    connect the neurons in the hidden layer/s using the selected nonlinear function (so-called activation function) so that a first prediction of the output node, which is in this case the NO$_2$ concentration, can be estimated. As the output from the forward propagation may not be correct, in the second step, the so-called backpropagation, the biases, and weights are optimized to

minimize the residual sum of squares between the observed values ($NO_2$ concentration of the RI) and the predicted values using gradient descent. In order to avoid wrong predictions caused by local minimums, a parameter called learning rate ($\alpha$)

235    should be as small as possible. Note that the smaller the learning rate, the longer the computational time so an optimum must be found (Bishop, 2006; Awad and Khanna, 2015).

The hyperparameter tuning for the ML models was carried out in Python (version 0.22.1) using the *RandomizedGridSearchCV* optimizer provided by the Scikit-learn library. Additionally, Keras and TensorFlow libraries for ANN models were used. In order to avoid overfitting, a 5-fold cross-validation was used. Some of the preliminary hyperparameter values were based on

240    the literature (Wei et al., 2020; Spinelle et al., 2015; Pedregosa et al., 2011) whereas others were manually tested by means of observing how the learning curves react (Géron, 2019). The grid of parameters for each model is shown in Table S1-S3 in the Supplement. Among the whole calibration dataset, 75 % of the data was used for training and the other 25 % for validation. Both datasets were randomly selected. The hyperparameters were tuned for each sensor individually. All the ML models were built using the Scikit-learn library in Python. A total of 217 simulations were run, 96 % of which were completed in less than

245    15 minutes on a single 2.50 GHz Intel i7-6500U CPU.

**2.4 Data pre-processing**

Firstly, in order to identify and remove outliers, data cleaning was carried out using an unsupervised learning algorithm, the Density-Based Spatial Clustering of Applications and Noise (DBSCAN) (Ester et al., 1996), prior to the training of the models. The warming-up period of $NO_2$ sensors was observed to range from some hours to some days. Therefore, unsuitable data from

250    the warming-up phase were manually removed after visual inspection of the data.

For $PM_{2.5}$ sensors, the data for calibration of the indoor sensors were averaged every 1 min whereas the data of outdoor sensors were averaged every 30 min. In the case of the calibration of the $NO_2$ sensors, we evaluated the impact of the averaging time in the model performance by using 1-, 5-, 10- and 15-minute averages. In theory, the higher the number of data points, the better the model performs. However, the $NO_2$ sensor data contained a lot of noise so an optimum between the number of

255    training data points and the reduction of the noise in the sensor signal must be found. During the deployment in the houses of the patients, hourly and daily averages were used for the analysis.

For ANN and SVR models, the data of the explanatory variables were normalised from 0 to 1 using Eq. 3,

$$X_N = \frac{X_i - X_{min}}{X_{max} - X_{min}} \tag{3}$$

where $X_N$ is the normalized value, $X_i$ is the feature value ($i$) to be normalized, and $X_{min}$ and $X_{max}$ are the minimum and maximum values of the feature, respectively. After the prediction, the results were transformed back to the real values.

260    **2.5 Model evaluation**

Various goodness-of-fit indexes were used to assess the performance of the models including root-mean-square error (RMSE), centred root-mean-square error (CRMSE), mean bias error (MBE), mean absolute error (MAE), the coefficient of

determination ($R^2$), Person correlation coefficient (r), model efficiency (MEF) and fractional bias (FB). The respective formulas and ideal values are summarized in Table S4 of the Supplement.

265 The results of the $PM_{2.5}$ and $NO_2$ sensors were also evaluated using target diagrams. A target diagram is built using the CRMSE and the MBE of the testing set as the x-axis and y-axis, respectively, both normalised by the standard deviations of the RI ($\sigma_{ref}$). As the values of CRMSE are always positive, the model predictions are plotted in the left quadrants if their standard deviation is lower than the standard deviation of the RI (Zimmerman et al., 2018). The outermost circle of the diagram corresponds to the performance criteria, set as 1, whereas the inner circle represents the performance goal which has been

270 defined for this study as 0.5, that is, 50 % more stringent. In general, the performance of the model is higher the closer the attained performance score is to the target diagram's origin (Thunis et al., 2012).

Following the recommendation of the CEN/TS 17660-1:2021 and the CEN/TS 17660-2:2025, the relative expanded uncertainty (REU) has been calculated to determine whether the sensor data fulfil the DQOs as defined in the Directive 2008/50/EC. On November 24, 2024, the EU Directive 2024/2881 was published, establishing stricter limit values to be

275 achieved by January 2030. The new directive also specifies in Annex V new DQOs for indicative measurements (I.M.) and objective estimation (O.E.).

In the CEN/TS 17660-1 (2021) and CEN/TS 17660-2 (2025), three different classes for sensors are defined. Sensors fulfilling the DQO required for indicative measurements belong to class 1 whereas sensors in class 2 fulfil the DQO for objective estimations. A third class, which is less strict and is not formally associated with the Directive, has also been defined. Class 3

280 is not object of study of this work, as it is not formally linked to a binding DQO.

In Table 3 the DQOs for both Directives, the 2008/50/EC and 2024/2881 are presented. More information about how to calculate the REU can be read in the Supplement. As shown in Table 3, the DQO of the objective estimation for hourly $NO_2$ values has changed from 75 % in Directive 2008/50/EC to 80 % in Directive 2024/2881 whereas the DQO for daily $PM_{2.5}$ mean concentrations has changed from 100 % in Directive 2008/50/EC to 85 % in Directive 2024/2881. For indicative

285 measurements, only the DQO of daily mean concentrations of $PM_{2.5}$ has been re-defined from 50 to 35 %.

Another aspect that should be noted is the average time. For epidemiological studies, especially those using portable monitors, 24 h average or even 1 h average may still not yield enough resolution, as higher time resolutions are needed to detect real-time pollution peaks, especially in mobile measurements. Moreover, long co-location periods are not always possible during the exposure assessment campaigns and consequently, the use of a 1-hour average can decrease considerably the available data

290 to train the models and reduce the range of T and RH, as well as the pollution concentration range used to train the model. Therefore, in this work, we present the REUs of the $NO_2$ models for different averaging times, that is, 1, 5, 10, and 15 min and thus, an evaluation of the REUs at the limit value is not applicable. Similarly, co-location measurements of indoor $PM_{2.5}$ sensors in a particle chamber with high particle concentrations lasted an average of 2 to 3 hours. Therefore, the uncertainties were calculated for a 1-min average. For outdoor $PM_{2.5}$ sensors where more data points are available, a 30-min average was

295 used so that neither REU for $PM_{2.5}$ measurements for indoor or outdoor are applicable in the region of the limit values.

**Table 3.** DQOs specified as the largest REU for short-term concentrations (Directive 2024/2881, 2024; Directive 2008/50/EC, 2008).

| Air pollutant | DQO I.M. | | DQO O.E. | |
|---|---|---|---|---|
| | 2008/50/EC | 2024/2881 | 2008/50/EC | 2024/2881 |
| $NO_2$ (1-hour) | 25 % | 25 % | 75 % | 80 %[a] |
| $PM_{2.5}$ (24-hours[b]) | 50 % | 35 % | 100 % | 85 %[c] |

[a] Calculated as the maximum ratio (3.2) over the uncertainty of indicative measurements (see Annex V of EU Directive 2024/2881).

[b] The EU Directives do not include uncertainty for $PM_{2.5}$ hourly values.

[c] According to Annex V of EU Directive 2024/2881: "The uncertainty of objective estimation shall not exceed the uncertainty for indicative measurements by more than the applicable maximum ratio and shall not exceed 85 %".

## 3 Results

### 3.1 Sensor data validation before deployment

#### 3.1.1 Relative expanded uncertainty

The REU of the testing data for the indoor and outdoor $PM_{2.5}$ sensors before the deployment in the houses of the patients can be seen in Fig. 2. The DQOs of the EU Directive 2008/50/EC and the new EU Directive 2024/2881 for both objective estimation and indicative measurements of $PM_{2.5}$ are also indicated. As shown in Fig. 2 (a), the unit-to-unit variability of indoor $PM_{2.5}$ sensors is significant. Specifically, the $PM_{2.5}$ sensor in B04-P3 meets the DQO for indicative measurements up to 2 µg m⁻³ and 3 µg m⁻³ under Directives 2008/50/EC and 2024/2881, respectively. In contrast, the $PM_{2.5}$ sensor in B01-P4 meets the DQO for objective estimation only for the Directive 2008/50/EC and concentrations higher than approximately 36 µg m⁻³. Three out of six indoor sensors fulfil the DQO for objective estimation set in the Directives 2008/50/EC and 2024/2881 at 12 µg m⁻³ and 14 µg m⁻³, respectively, and meet the DQO for indicative measurements for $PM_{2.5}$ concentration higher than 24 µg m⁻³ and 35 µg m⁻³ for the same directives respectively.

As can be observed in Fig. 2 (b), the unit-to-unit variability of outdoor calibrated sensors is less pronounced, with some sensors reaching the DQO for indicative measurements for concentrations higher than 5 to 6 µg m⁻³ (B06-P4, B06-P7_end) for both Directives. Four out of nine calibrated sensors fail to fulfil the DQO for indicative measurements of the new Directive 2024/2881 in contrast to only two that do not achieve the DQO for indicative measurements contemplated in the Directive 2008/50/EC. For the latter Directive, most sensors reach the mentioned DQO at concentrations higher than 16 µg m⁻³.

Similar to the indoor AQSSs, the results for outdoor sensors present data from different testing datasets for the same AQSS. For instance, the AQSS B05 was used by two patients (P2 and P4) and therefore calibrated twice before each deployment. The AQSS B03 was used in the houses of three patients but calibrated four times, including an additional co-location period after the deployment in the house of patient P7. In contrast to indoor calibrated sensors, outdoor sensors exhibit generally consistent REU across different deployments, as observed by the overlapping points. This consistency suggests that the calibration method may influence the REU, possibly because the aerosol (liquid paraffin) used in the particle chamber for calibration does

325 not have the same composition as the urban dust. The OPC-R1 sensor has been designed for ambient aerosol monitoring, assuming a density of 1.65 g/ml and a refractive index of 1.5+i0 for the calculation of the PM fraction. Additional details regarding the calibration conditions, the $PM_{2.5}$ concentration range and the calibration coefficients can be read in Table S5 in the Supplement.
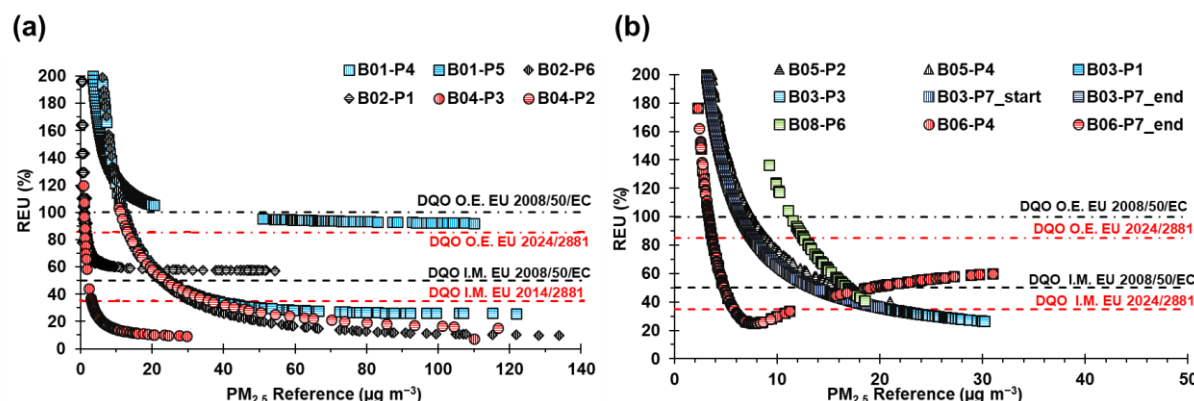


330 **Figure 2.** REU for (a) indoor and (b) outdoor $PM_{2.5}$ sensors. The dashed lines indicate the DQOs for indicative measurements while the dash-dot lines represent the DQOs for objective estimation (black for EU Directive 2008/50/EC and red for EU Directive 2024/2881).

Examples to illustrate the REUs of indoor and outdoor $NO_2$ sensors are shown in Fig. 3, which contains the results of the tested models (MLR, SVR, RFR, and ANN) as well as the influence of the averaging time of the training dataset, for 1, 5, 10, and
335 15 min on the REU. The DQOs of both Directives 2008/50/EC and 2024/2881 for objective estimation and indicative measurements of $NO_2$ are also indicated. Note that both directives have the same DQO for indicative measurements (25 %). The y-axis has been limited to 110 % so that the difference among the models can be distinguished. In Figures S3 and S4 the diagrams for all the other indoor and outdoor AQSSs are shown, respectively.

In general, the longer the averaging time used for training the data, the lower the REU. However, the longer the averaging
340 time, the smaller the dynamic range of the input variables, which can lead to higher uncertainties due to data extrapolation. Thus, an optimum averaging time shall be used. In our study, we found a 10-minute average time a good compromise between training the models with enough data points and achieving low uncertainties.

For the sensors calibrated in indoor conditions, SVR and RFR usually perform better than ANN and MLR. The MLR trained using data averaged 1-min performs in most of the cases the worst. This could be due to the signal noise, not only from the
345 sensor but also from the data of the RI used for the training. Results show that the DQO for indicative measurements (25 %) is achieved with a 10- or 15-minute average and mixing ratios larger than about 5 - 22 ppb for indoor and 10 - 25 ppb for outdoor AQSSs. The low REU that is achieved during the calibration of AQSSs in indoor conditions may be due to the controlled conditions, as the $NO_2$ gas was given stepwise and kept constant for 3.5 hours, as well as the controlled changes of the T and the RH. Other authors have also observed better results when the sensors are calibrated in control conditions as

350    compared to outdoor calibrations but they fail later during the field deployment (Castell et al., 2017). This creates the challenge of calibrating indoor AQSSs for a wide range of $NO_2$ concentrations and meteorological parameters.
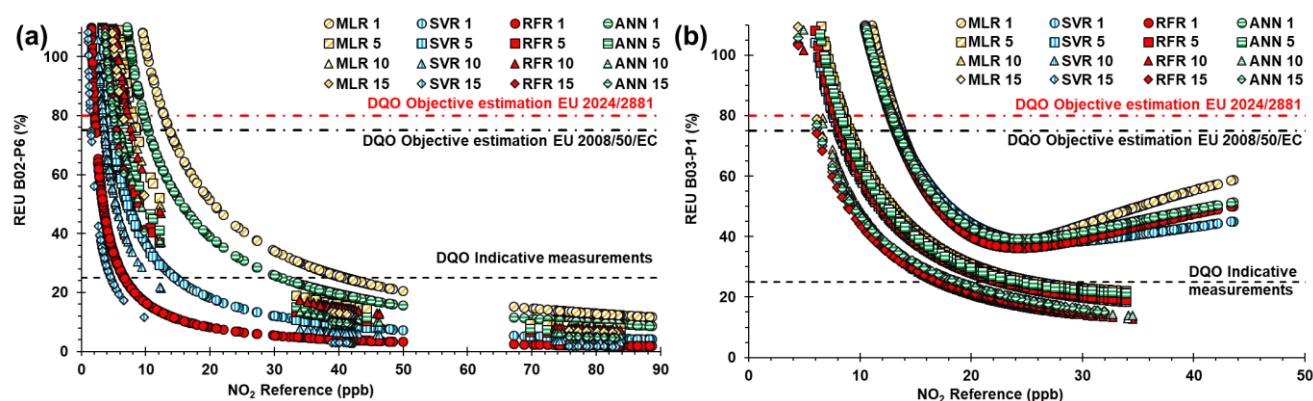


**Figure 3.** Example of REU for (a) indoor and (b) outdoor $NO_2$ sensors for the tested models (MLR, SVR, RFR and ANN) at different
355    averaging times (1 min, 5 min, 10 min and 15 min) against reference concentrations. The dashed line indicates the DQO for indicative measurements while the dash-dot lines represent the DQOs for objective estimation (black for EU Directive 2008/50/EC and red for EU Directive 2024/2881).

### 3.1.2 Target diagrams

360    The target diagrams for the testing data of the indoor and outdoor $PM_{2.5}$ sensors are shown in Fig. 4. Two main results can be inferred from these diagrams: (i) Different outcomes are obtained with the same sensor for each calibration period, as indicated by the symbols with the same form and colour and (ii) the results of indoor $PM_{2.5}$ sensors remain within the unit circumference, being most of them even within the inner circle, which is 50 % more stringent. In contrast, four out of seven outdoor $PM_{2.5}$ sensors do not perform well enough to reach the inner circle, and most of them remain outside the unit circumference. The
365    differences between the indoor and the outdoor sensors' performances can be attributed to the same factors discussed in Section 3.1.1. Other researchers have obtained similar results, with $PM_{2.5}$ sensors falling within and without the target circle without specific patterns (Borrego et al., 2016). The question of whether the prototype of the dryer unit helped to improve the performance of the $PM_{2.5}$ sensors of the outdoor AQSSs may arise after analysing this outcome. In Chacón-Mateos et al. (2022) the weaknesses and strengths of the thermal dryer used for this study were discussed in detail. In that study, it was concluded
370    that the dryer was causing an excess of heating and therefore an underestimation of $PM_{2.5}$ concentrations compared to the RI. In this regard, we have developed a new prototype to keep the air temperature inside the dryer at less than 40 °C.
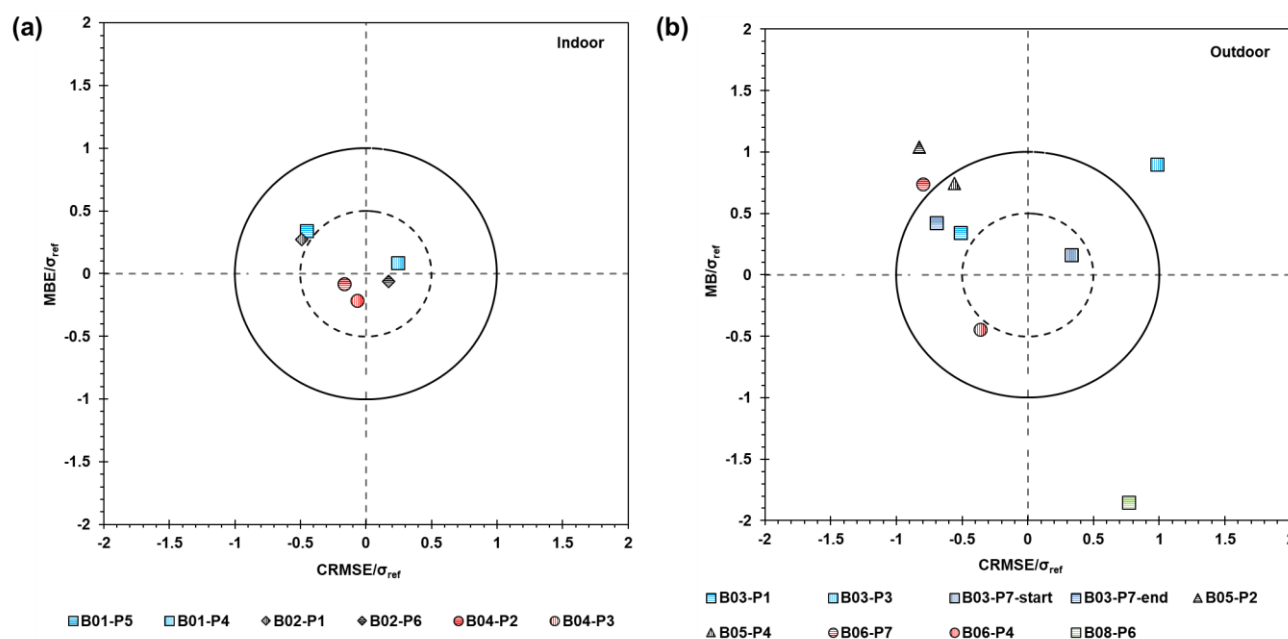
13

**Figure 4.** Target diagrams for (a) indoor and (b) outdoor PM$_{2.5}$ sensors.

375

Figure 5 illustrates two examples of target diagrams for the tested models for indoor and outdoor NO$_2$ sensors. The remaining results for indoor and outdoor NO$_2$ sensors are available in Figures S5 and S6 respectively. All the indoor sensors fall within the performance goal ($\pm$ 0.5) independently of the average time and the model used, indicating high accuracy (low mean bias or systematic error) and high precision (low CRMSE or random error) for all the models.

380 The models for correcting NO$_2$ sensor readings outdoors show more discrepancies among the models and averaging times. Models trained using 1-min averaging time show the worst performance, followed by the 5-min average. For most of the models, the results of the target diagrams for 1- and 5-min averages do not reach the performance target ($\pm$ 0.5). Higher averaging periods like 10-, 15- or 1-hour usually reach the inner circle. In terms of models, SVR and RFR tend to outperform MLR and ANN achieving higher accuracy and precision. In all the cases, the results are situated on the left side, indicating

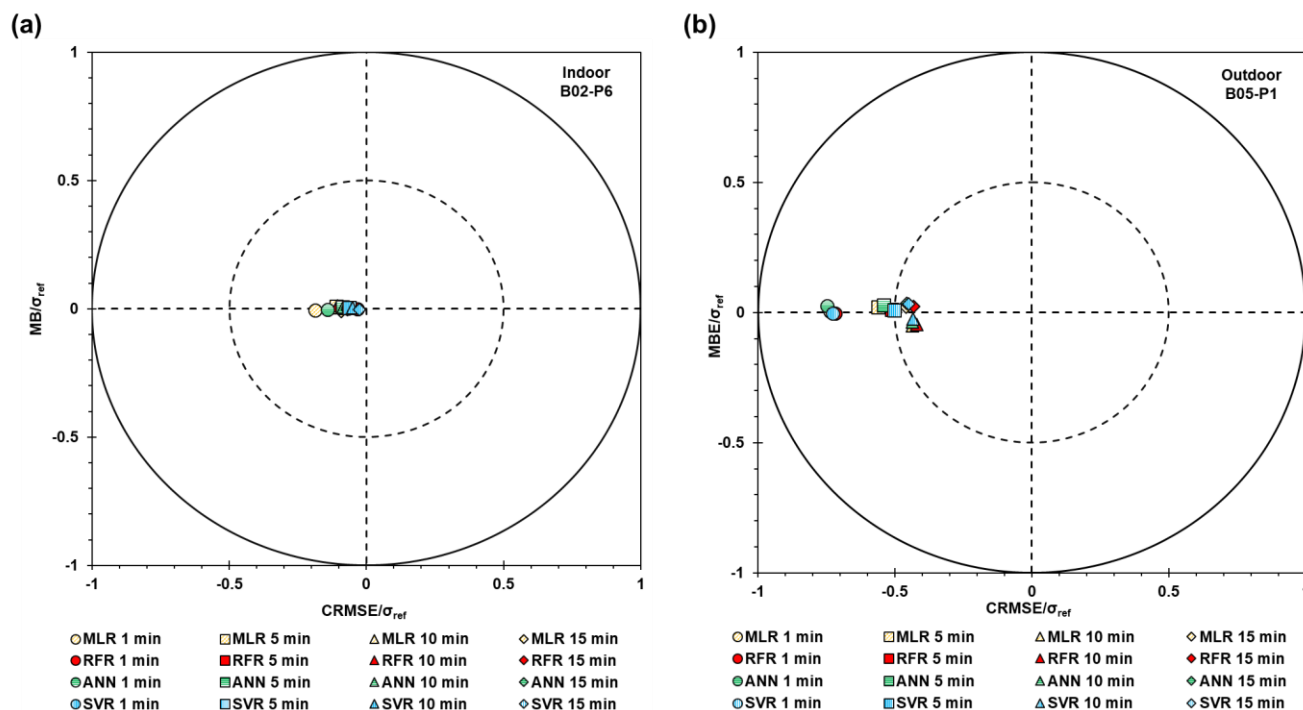385 that the standard deviation of the sensors was lower than the standard deviation of the RI.

**Figure 5.** Example of target plots for indoor (left) and outdoor (right) NO₂ sensors for the tested models and different averaging times.

### 3.1.3 Performance metrics

Figure 6 presents the statistical results for various metrics (orthogonal slope and intercept, model efficiency, MAE, and Pearson

390    correlation coefficient) of the models tested for indoor and outdoor NO₂ sensors at different averaging times. Consistent with
previous findings, the indoor models outperform the outdoor models, likely due to the more controlled laboratory conditions.
Notably, the model efficiency for all indoor models is nearly 1, indicating an almost perfect match to the RI data. When
comparing different time aggregations, it is evident that higher aggregation intervals result in the orthogonal slope approaching
1 and the orthogonal intercept approaching 0 for all the tested models. This is attributed to the reduction in sensor noise and

395    increased data stability with higher time aggregation. However, when comparing the MEF for 10- to 15-minute time
aggregations, no improvement is observed; instead, there is a decrease in performance across all models. This decline is likely
due to the excessive reduction in the number of training data points, with approximately 35 % fewer data points (see Fig. S7).
This trend is also observed in the MAE, which decreases from an average of 10 ppb across all models with 1-min averaging
time to 5 ppb using 10- and 15-min averaging times for outdoor NO₂ sensors. The improvement in the indoor NO₂ sensors is

400    less notable. The Pearson correlation coefficient shows an improvement between 1-min and 5-min averaging time but remains
stable thereafter for both indoor and outdoor sensors. In general, MLR shows the worst performance across the tested models.
SVR and RFR exhibit the best performance, closely followed by ANN.
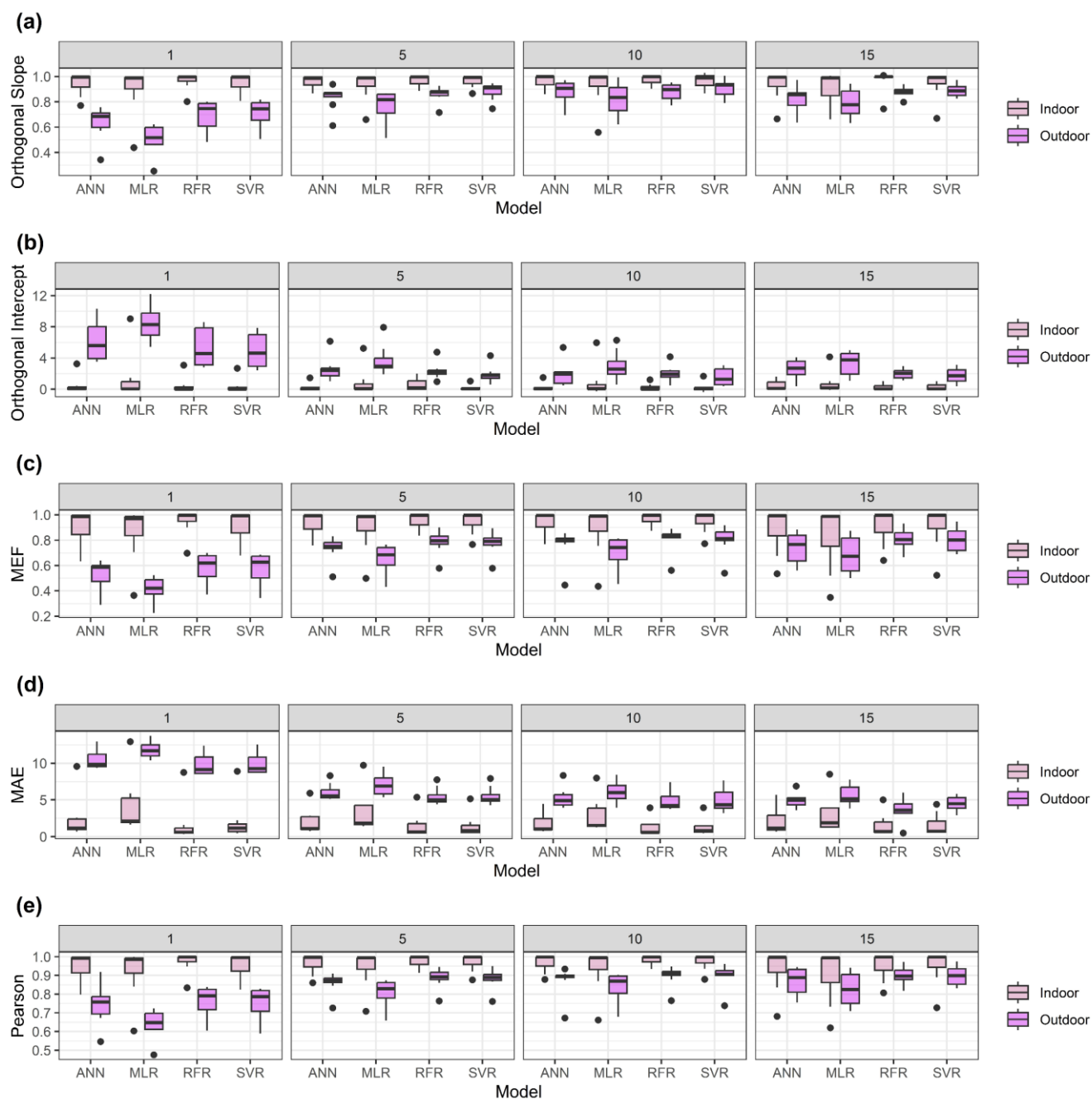
15

**Figure 6.** Boxplots of various performance evaluation metrics: (a) orthogonal slope, (b) orthogonal intercept (in ppb), (c) model efficiency
(MEF), (d) MAE (in ppb) and (e) Pearson correlation coefficient, for different tested models (ANN, MLR, RFR and SVR) for the different
time aggregations (1-, 5-, 10- and 15-min) applied to the testing data for indoor and outdoor $NO_2$ sensors.

Figure 7 presents the performance evaluation metrics for the indoor and outdoor PM$_{2.5}$ sensors. The calibration factor ($\beta_1$) and calibration constant ($\beta_0$) for the indoor sensors are closer to 1 and 0, respectively, compared to the outdoor sensors. Notably,
410 almost all the sensors exhibit a calibration constant greater than zero ($\beta_0 > 0$). This constant deviation, or displacement error, may be attributed to the different limits of detection of the OPC-R1 (0.35 µm) compared to 0.30 µm of the RI. As mentioned in Section 2.3.1, the indoor sensors were calibrated using an aerosol generator and liquid paraffin. However, these particles do not accurately represent the heterogeneity of the particles present in the indoor air. This discrepancy likely explains why the indoor sensors perform better across most metrics except for the MAE, as higher concentrations (median 124 µg m$^{-3}$) were
415 generated during the calibration. In contrast, the highest median PM$_{2.5}$ concentration measured during the outdoor calibration is 35 µg m$^{-3}$. Overall, the calibrated indoor and outdoor sensors exhibit a median FB of less than 0.3, which is within the acceptable limits, and Pearson correlation coefficients of more than 0.75.
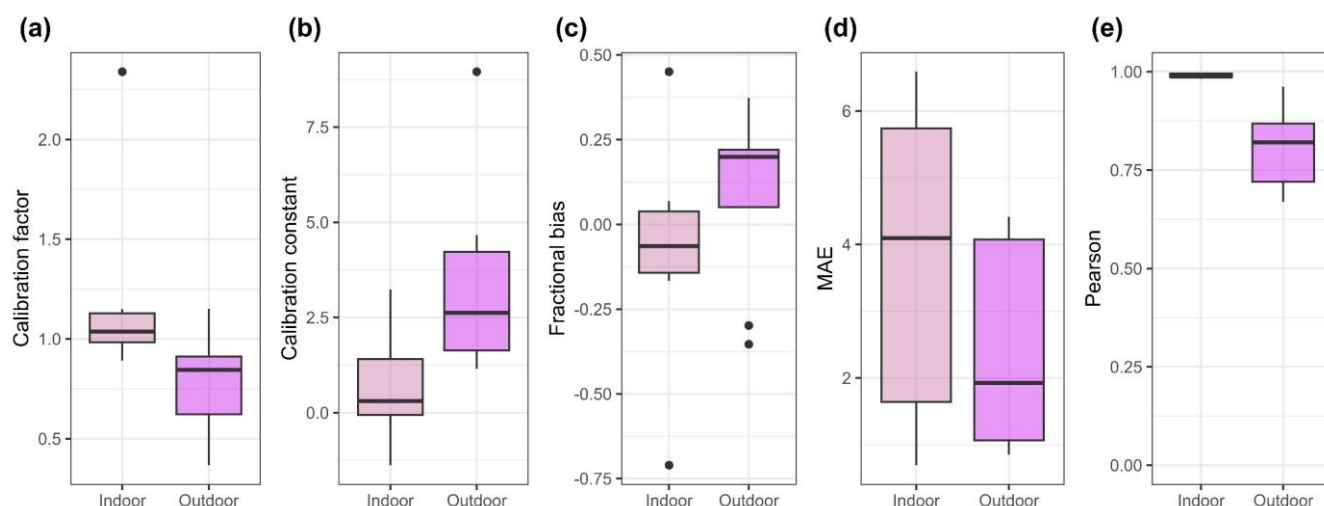


**Figure 7.** Boxplots of various performance evaluation metrics: (a) calibration factor, (b) calibration constant (in µg m$^{-3}$), (c) fractional bias,
420 (d) MAE (in µg m$^{-3}$) and (e) Pearson correlation coefficient, for indoor and outdoor PM$_{2.5}$ sensors.

## 3.2 Sensor data validation during deployment

### 3.2.1 Passive samples for NO$_2$

Figure 8 presents the results of the discontinuous NO$_2$ measurements using passive samples for the indoor and outdoor microenvironments during the deployment in the houses of the patients, compared to the results of the tested sensor calibration
425 models. Each sampling period spans 14 or 15 days, except for patient P4, whose period extended to 19 days. Considering the passive samples as the "true value", it is evident from Fig. 8 that the SVR model predicts indoor NO$_2$ poorly, with concentrations higher than 18 µg m$^{-3}$ in all the cases. This occurs despite achieving similar levels of uncertainty and better performance metrics as other models for the same averaging time during the testing period (see Section 3.1). RFR tends to overestimate the results, particularly for the indoor concentration measured in the house of patient P6 (average of both periods

430   35 µg m$^{-3}$ of NO$_2$ compared to 8 µg m$^{-3}$ measured with passive samples). These discrepancies suggest that SVR and RFR overfitted the training data. The negative average values of the MLR model deployed in the house of patient P6 indicate a signal drift. Both SVR and RFR also tend to overestimate outdoor NO$_2$ concentrations, although this tendency is less pronounced compared to indoor predictions. The MLR model sometimes overestimates and sometimes underestimates the concentrations. ANN appears to be the most robust model for both indoor and outdoor sensors, even though it occasionally

435   overestimates the actual NO$_2$ concentrations (up to 5 µg m$^{-3}$ more than the passive samples).

Figure 8 (a) also shows the results of the AQSS calibrated outdoors but used indoors in the house of patient P7. When analysing closely the outcomes, we can observe that the ML models are overestimating the results compared to the passive sample value but for SVR and RFR less than the indoor results of the other patients. The ANN is the model that better agrees with the results, showing 2 and 3 µg m$^{-3}$ more than the results of the passive samples for the first and the second period, respectively. The MLR

440   underestimates the NO$_2$ concentration the first period and overestimates in the second period. It should be noted that the warming-up period of the NO$_2$ sensor was in this case three days, longer than usual.

Overall, this comparison underscores the importance of not relying solely on pre-deployment performance evaluations. Reference values during deployment are crucial for verifying sensor performance. In this context, passive samples have proven to be a simple and effective tool to verify calibrated sensor data.
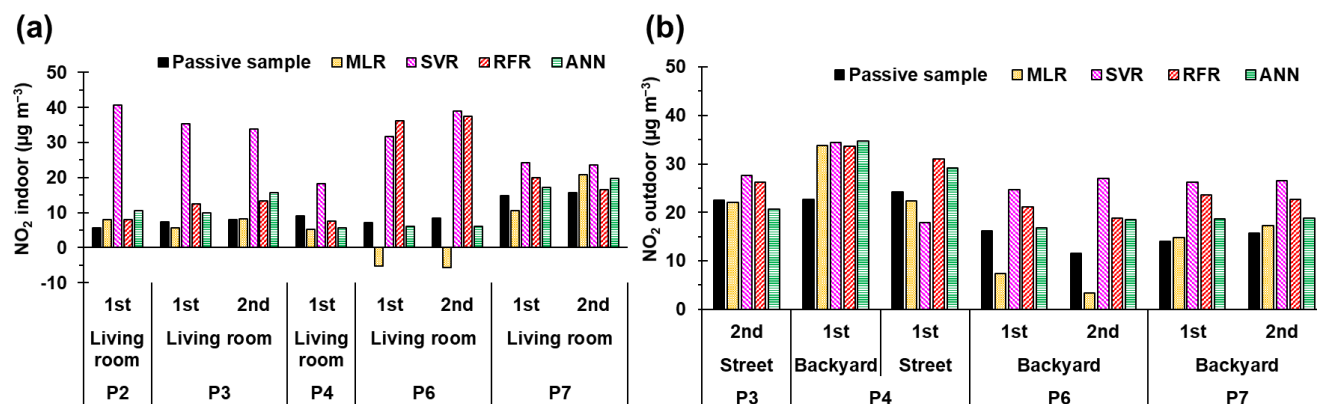


445

**Figure 8.** Comparison of the NO$_2$ models with the concentration measured by the passive samples (two-week period) for (a) indoor and (b) outdoor sensors. Models were trained with data averaged every 10 min.

### 3.2.2 Comparison of outdoor sensors with air quality monitoring stations

450   As part of the data validation process, the measurements from the outdoor AQSSs were compared with NO$_2$ and PM$_{2.5}$ data from the governmental air quality monitoring stations in the city and our measurement station at Hauptstätter Street. Figure 9 presents the results of the deployment of the AQSS placed outside the window of patient P1 and the nearest monitoring station.

18

Additional results are provided in the Supplement (Figures S8 - S13). The correction models for $NO_2$ sensors were trained with 10-min time aggregations.

455 The monitoring station shown in Fig. 9 is located at Arnulf-Klett-Platz, 1.1 km from the AQSS location, near a busy road. In contrast, the AQSS was installed at the window of a second-floor apartment adjacent to a secondary road. Due to the different locations, comparisons should be approached with caution, although similar temporal patterns in the pollution concentration are expected due to the shared urban and rural background concentrations.

Different trends in the $NO_2$ concentrations of the tested models are shown in Fig. 9 (a). Notably, the RFR model underperforms,

460 exhibiting excessively constant $NO_2$ levels over extended periods. This suggests our RFR model is not suitable for sensor data correction. Conversely, the SVR model fails to detect $NO_2$ concentrations below $20\,\mu g\,m^{-3}$, likely due to its limited extrapolation capability. The ANN model generally demonstrates satisfactory performance. Both the ANN and MLR models display trends that closely match the expected concentration trends. However, for other patients, MLR prediction reaches negative peaks up to $-100\,\mu g\,m^{-3}$ (see Fig. S11). The negative peaks occurred when the T was above 25 °C. This is a clear

465 example of the effect of the T and RH on the electrochemical sensors. The calibration period covered a T and a RH range of 2 - 25 °C and 40.8 - 77.4 %, respectively. However, during the measurement campaign in the house of patient P6, the $NO_2$ sensor was exposed to T up to 31 °C and RH as low as 8 %, which were far beyond the ranges covered during the calibration period. The MLR model must be used cautiously for T above 25 °C, as the influence of the T and the RH on the sensor signal is not linear (Samad et al., 2020).

470 Figure 9 (b) shows that the $PM_{2.5}$ sensor equipped with a low-cost dryer and calibrated using ULR closely follows the trend of the nearby reference station. A detailed examination reveals that the $PM_{2.5}$ readings were more accurate at the beginning of the deployment period compared to the end when the calibrated sensor reported higher concentrations than those from the reference station. Although initially unexpected, this discrepancy could be attributed to the highly localized nature of particulate matter concentrations. The placement of the AQSS in a building corner, which disrupts airflow, and its proximity to a tram line and

475 the entrance of a hospital parking, might result in higher concentrations. If there is one field where sensors have proven valuable, that is in identifying new pollution hotspots (deSouza et al., 2022).
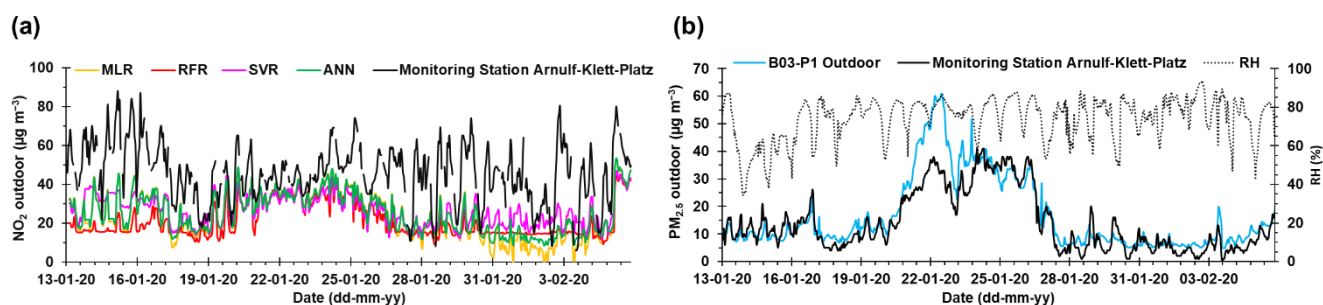


**Figure 9.** Time series of hourly outdoor (a) $NO_2$ and (b) $PM_{2.5}$ concentrations during deployment in the house of patient P1.

19

480  A highlight from Fig. 9 and Fig. S13 is that all the PM$_{2.5}$ sensors show similar trends compared to the monitoring station at Hauptstätter Street ($0.40 < R² < 0.93$), even the hours when the RH is higher than 70 %. The typical overestimation of the PM concentration by the sensors at high RH is avoided thanks to the thermal dryer.

### 3.2.3 Metadata for qualitative sensor data validation

In this section, we present an example of how the use of metadata, specifically the activity and window status log, contributes

485  to the validation of sensor data. Figure 10 shows the indoor NO$_2$ and PM$_{2.5}$ concentrations during the second week of deployment in the house of patient P2. Additionally, the different activities on an hourly basis and the status of the windows in the living room where the AQSS was located are shown. The NO$_2$ sensor readings have been corrected using the ANN model using 10-min aggregation time.

As illustrated in Fig. 10, pollution peaks can be correlated with specific activities at home. The information collected in the

490  logbook is invaluable for interpreting sensor data. It allows for the detection of anomalies and helps in understanding the source of pollution peaks. For instance, in Fig. 10 (c), there is a noticeable decrease in PM$_{2.5}$ concentration during sleeping hours and an increase during activities like exercising (on January 24th) and cooking (on January 24th and 27th). For NO$_2$, the activity log is especially useful when considering window status, as NO$_2$ typically originates from outdoor sources. This is evident in Fig. 10 (b), where some peaks occur when the window is open or tilted. A deeper analysis of the information acquired in the

495  log books and the relationship with the indoor air quality in the houses of the patients can be read in Chacón-Mateos et al. (2024).
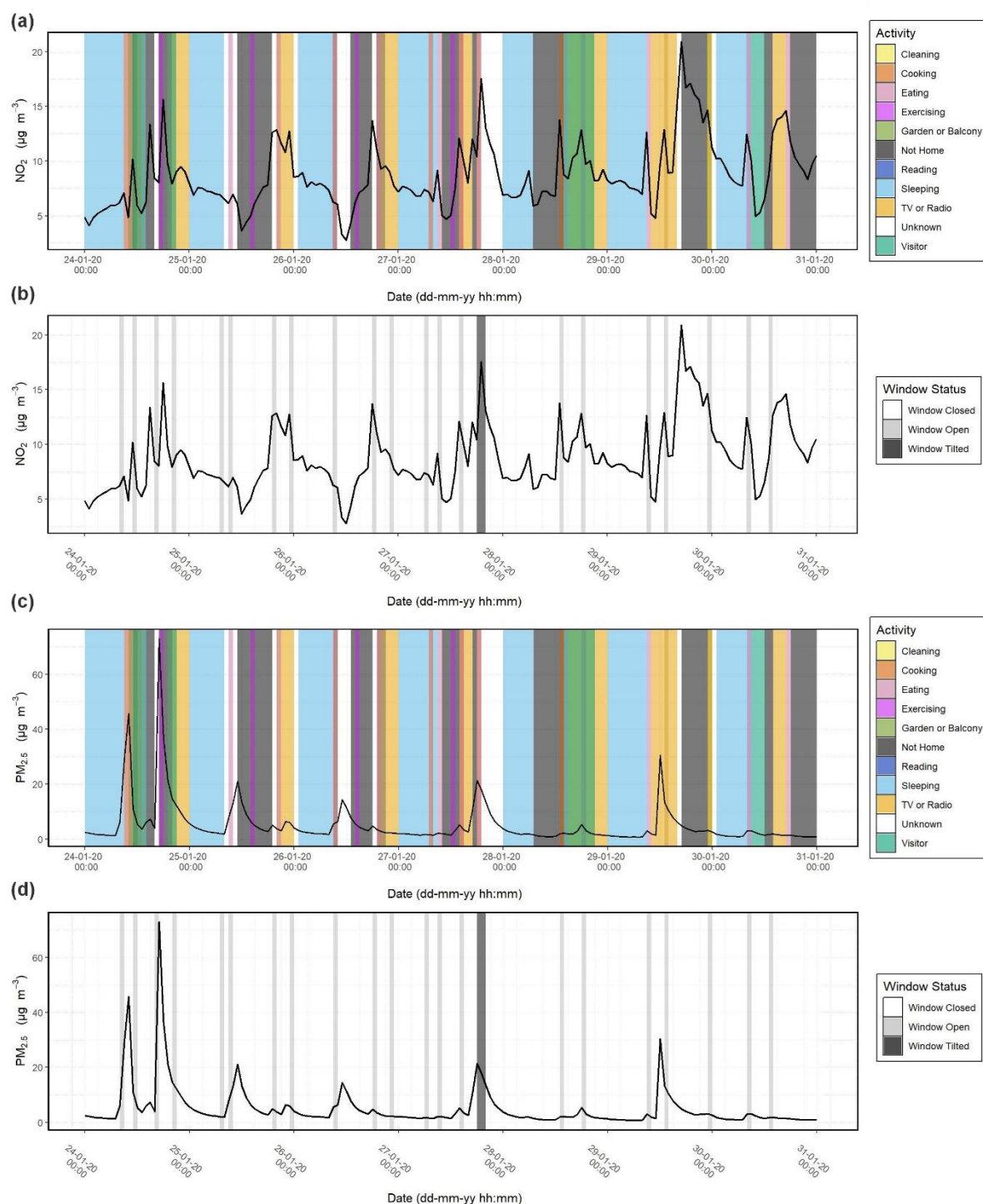
**Figure 10.** Hourly times series of (a) indoor NO$_2$ concentration and activities, (b) indoor NO$_2$ concentration and window status, (c) indoor PM$_{2.5}$ concentration and activities and (d) indoor PM$_{2.5}$ concentration and window status during one week in the house of patient P2.

## 4 Discussion

### 4.1 Evaluation of the NO₂ sensors

#### 4.1.1 Indoor NO₂ sensors

505 The results of this study indicate that using indoor co-location and artificially generated $NO_2$ to correct the signal of electrochemical sensors for $NO_2$ may not be effective for the tested models. Calibrating indoor sensors presents particular challenges due to two main factors: the low concentrations and the need to train models across different spans of concentrations, T, and RH. Although the testing data showed nearly perfect results after training the sensor data with artificially generated $NO_2$ and controlled changes of T and RH, applying the model to the sensor data deployed in the patients' homes yielded significantly different outcomes. Some models like SVR and RFR struggled to accurately predict the $NO_2$

510 concentrations in the new indoor environment.

We conducted an experiment to test whether an AQSS for indoor use could be calibrated in outdoor co-location to better learn real $NO_2$ concentration and meteorology patterns. However, in this case the machine learning models consistently overestimated concentrations compared to the passive samples. This discrepancy arises because in Stuttgart outdoor $NO_2$ concentrations are generally higher than indoor concentration, and the models are not adept at extrapolating to lower

515 concentrations.

A promising solution may be a hybrid calibration, for instance, the Enhanced Ambient Sensing Environment (EASE), which combines the advantages of laboratory calibration with the increased accuracy of field calibration (Russell et al., 2022). To date, this approach has only been tested with multilinear regression models. Further research is needed to determine whether it is suitable for training machine learning algorithms. Another possible solution is the calibration of the sensors in occupied

520 homes (Suriano and Penza, 2022) or exposing the sensors to cooking events (Tryner et al., 2021). However, these studies did not deal with the re-colocation of the monitors after the calibration in a new environment. Therefore, further research is needed to expand our knowledge of calibration transfer in indoor environments for electrochemical sensors.

#### 4.1.2 Outdoor NO₂ sensors

The calibration of $NO_2$ sensors through a co-location with RI outdoors is at this moment a common procedure (Karagulian et

525 al., 2019). Many studies have tested different regression and ML models (Spinelle et al., 2015; Cordero et al., 2018; Zimmerman et al., 2018; Malings et al., 2019). Our results on the performance evaluation for the outdoor $NO_2$ sensors are similar to the outcomes from Bigi et al. (2018) obtained using 10-minute averages for MLR, SVR and RFR and to those from Apostolopoulos et al. (2023) for the ANN model (note that their results are based on hourly values). One limitation of our study was the lack of ozone data. It has been demonstrated that the sensor B43F has cross-sensitivity to ozone despite having

530 an ozone filter and that the influence of ozone increases as the filter saturates (Li et al., 2021). Knowing this, we experimented with adding ozone data from the air quality monitoring station located at Marienplatz (Stuttgart) to train the models of data acquired in the houses of patients P1 and P7, representing the cold and the warm months, respectively. The results of the error

metrics are shown in Fig. S14. Even though the results of the $R^2$ and RMSE seem to improve in most of the cases adding ozone data, the results of the MAE show the opposite trend in the cold months (P1). Moreover, the difference in the RMSE results

535 during the warm month (P7) is minimal for most of the models except for MLR. Therefore, we did not further investigate the addition of ozone for the rest of the data. Furthermore, ozone concentrations are higher in summer months and our measurement campaign ran from December to May, i.e. mainly in winter months when ozone concentrations are lower. In addition, studies have shown that the performance degradation of the ozone filter starts 200 days after sensor unpacking (Li et al., 2021), which is approximately the number of days that our campaign lasted. Nevertheless, for future studies, we recommend adding an

540 ozone sensor so the cross-sensitivity can be corrected for all seasons.

### 4.1.3 Evaluating averaging times

The choice of the temporal resolution significantly affects the data's noise level and the quality of training data for $NO_2$ sensor models. High-resolution data (e.g., 1-min intervals) often contain more noise, which can degrade training quality. Conversely, using lower resolutions (e.g., hourly averages) may excessively reduce the number of training data points available and the

545 concentration range covered in the calibration. Our study demonstrated that using a 10-min averaging period over a two-week calibration phase resulted in lower MAE for $NO_2$ sensors. Although some researchers have employed hourly averages (Cai et al., 2009; Wei et al., 2020; Goulier et al., 2020), others have also identified a 10-minute average as the optimal (Paas et al., 2017; Bigi et al., 2018). In contrast, Sahu et al. (2021) in their analysis of the effect of temporal data averaging, found that data averaged every 5-min provided better results. All in all, the selection of an appropriate averaging time depends largely on the

550 quantity of available training data and must be carefully selected.

### 4.1.4 Evaluating calibration models for $NO_2$ sensors

The results of our study show that ANN is the most robust model for transferring the calibration parameters of the sensors to be used in another place, either indoors or outdoors, using the proposed calibration methodology. Even though RFR and SVR show better results for the metrics RMSE, MAE and Pearson correlation coefficient during the validation phase prior to

555 deployment, the comparison with the passive samples during the measurement campaign in the houses of the patients showed that SVR and RFR overestimated in most cases the $NO_2$ concentrations. The MLR showed the worst performance among the tested models.

### 4.2 Evaluation of the PM$_{2.5}$ sensors

### 4.2.1 Indoor PM$_{2.5}$ sensors

560 Validating PM$_{2.5}$ sensor measurements indoors presents significant challenges. While activity logs provide invaluable information regarding events that might cause elevated PM$_{2.5}$ concentrations, it remains unclear without having a RI in the houses, whether sensors can accurately quantify these peaks. This uncertainty may be particularly problematic in short-term

exposure studies, where precise measurement of peaks is critical. However, in long-term studies, short-duration peaks contribute less to the overall concentration average, thus presenting a lesser concern.

565 The results of this study suggest that using test aerosols like liquid paraffin in a particle chamber may not be an optimal technique for PM sensor data correction. This is likely due to discrepancies between the density assumed in the sensor's internal algorithm and the actual density of the generated particles.

Our study also explored the use of an outdoor calibrated AQSS intended for indoor deployment. However, due to the lack of an RI indoors during deployment, we cannot conclusively determine if this method outperformed indoor calibrations carried

570 out in the laboratory. Previous research, such as the study by Koehler et al. (2023), suggests that calibrations using ambient outdoor air data can enhance the data quality of indoor sensors compared to using manufacturer-provided calibrations. Nonetheless, the composition and concentration ranges of PM indoors can significantly differ from that of outdoor air, which may affect the correct performance of the sensor calibration. Further research is necessary to evaluate various calibration methods for indoor sensors and to understand how different PM compositions influence sensor performance.

575 **4.2.2 Outdoor PM$_{2.5}$ sensors**

One of the biggest concerns about PM sensor measurements outdoors is the effect of hygroscopic growth or fog. The use of either physical air preconditioning or data correction approaches is a must in regions where high relative humidity and hygroscopic aerosols are expected, as it is the case of Stuttgart. For this project, a low-cost dryer unit was designed to avoid the overestimation of PM$_{2.5}$ concentrations.

580 The results of the comparison of sensor data with data from local monitoring stations in Stuttgart in the vicinity of the houses of the patients showed that the PM$_{2.5}$ sensors showed a similar trend even when the RH was higher than 70 %. Given the fact that a simple linear regression applied to the outdoor PM$_{2.5}$ sensors with a dryer shows plausible results when compared to the nearest measurement stations, this method can be used to simplify the models for PM data correction. However, it is important to control the drying temperature as temperatures higher than 40 °C could evaporate semi-volatile organic compounds and

585 trigger the underestimation of the PM mass concentration (Chacón-Mateos et al., 2022).

**4.3 Do the sensors fulfil the Data Quality Objectives?**

Previous studies have indicated that while commercially available sensor systems often meet the criteria for indicative measurements of PM$_{2.5}$, NO$_2$ sensors frequently struggle to fulfil the DQO (Castell et al., 2017). This challenge prompted the design and evaluation of our AQSSs. However, the rapid advancement in sensor technology outpaces scientific literature,

590 making it difficult to keep up with the latest developments.

Regarding NO$_2$ sensor units, many researchers have applied calibration models that account for parameters such as RH, T, and ozone data. These models have demonstrated that the DQO for indicative measurements can be achieved for NO$_2$ concentrations above 20 ppb (Spinelle et al., 2015; Bigi et al., 2018; D'Elia et al., 2024). Our findings align with these results, showing that outdoor NO$_2$ sensors meet the DQOs of both EU Directive 2008/50/EC and 2024/2881 for indicative

595  measurements between 10 and 25 ppb, depending on the specific sensor unit and the averaging time used. Sensors calibrated in indoor conditions performed even better, achieving the DQOs at even lower concentrations. However, we have also argued that the use of a GPT system to generate specific $NO_2$ concentrations may not be appropriate for training ML models intended for deployment in indoor environments.

It is evident that even after calibration, the "hardware" of electrochemical sensors has not reached enough maturity yet for

600  applications requiring low measurement uncertainty, especially for low concentrations, making the measurement very dependent on the "software" used to correct the data (regression models, ML, etc.). Recent advancements in sensor units include onboard temperature monitoring near the electrical cell, which appears highly promising to improve the accuracy of the signal correction.

Our research also highlights the impact of the averaging time on the REUs of calibrated sensors. Generally, longer averaging

605  times improve the likelihood of meeting the DQO at lower concentrations, though this often reduces the concentration range covered during calibration. Moreover, ML models may not predict accurately outside the concentration range for which they were trained. Therefore, we recommend adjusting the averaging time based on the available data to potentially enhance model performance.

For $PM_{2.5}$ measurements, both DQOs are in the new EU Directive 2024/2881 stricter, from 50 to 35 % for indicative

610  measurements and from 100 to 85 % for objective estimation. Considering that, the DQO for indicative measurements after an indoor sensor calibration is typically achieved at concentrations above 23 and 35 µg m$^{-3}$ for the Directives 2008/50/EC and 2024/2881, respectively. After field calibration of the outdoor units, the DQO for indicative measurements is achieved at concentrations higher than 16 µg m$^{-3}$ under EU Directive 2008/50/EC. However, four out of nine sensors fail to meet the DQO criteria of EU Directive 2024/2881. Moreover, a significant unit-to-unit variability exists. This variability has been noted in

615  previous studies, such as those on the SDS011 sensor (Liu et al., 2019).

In summary, while the tested sensor units generally fulfil the DQOs for higher concentrations, their performance at lower concentrations may not be good enough for certain applications such as health studies. This limitation suggests that these sensors may be more suitable for environments with expected high pollution levels, where more epidemiological studies are needed (Amegah, 2018). Nevertheless, it is important to acknowledge that even reference-grade monitors are not free from

620  uncertainties (Diez et al., 2024). Regular quality control is essential for all air quality monitoring devices, whether they are gold standard, reference-grade, or sensor-based.

**4.4 Considerations on the concept of "low-cost" sensors**

In this study, we designed two AQSSs costing approximately 400 euros for indoor and 500 euros for outdoor applications, excluding labour costs. Despite the relatively low acquisition cost compared to a reference-grade air monitor, the

625  implementation and maintenance of the sensor systems are not necessarily low-cost. Moreover, the term "low-cost" varies significantly by region, and we have intentionally avoided its use in this manuscript. Even though we acknowledge that the term "low-cost" or the abbreviation "LCS" has helped to differentiate them from traditional air monitors and form a

recognizable community, we recommend that future publications also refrain from using "low-cost" or "LCS" and instead use "air quality sensors" or "AQS".

## 5 Conclusion

In this study, we investigated the sensor models OPC-R1 and the B43F sensors to measure $PM_{2.5}$ and $NO_2$, respectively, for their use in health studies covering indoor as well as outdoor microenvironments. The performance evaluation was carried out using common error metrics, target diagrams and the DQOs described in the EU Directive 2008/50/EC as well as in the recently published EU Directive 2024/2881. The co-location phase was conducted before the measurement campaign, where the data from reference-grade instruments were used to calibrate the $PM_{2.5}$ sensors and test regression and ML models to correct the $NO_2$ sensor data. Moreover, a methodology to validate the sensor data during the deployment of the sensor systems in the houses of the patients was designed and evaluated.

Even though ML is a promising tool in the field of AQS, special care during the co-location is crucial to assure the quality of the training data. The performance evaluation has shown that indoor sensor calibration using test aerosols for PM and artificially generated $NO_2$ is not appropriate to train the models for sensor re-location as model transferability cannot be guaranteed. The concentration range covered and the similarity of the environment (T and RH ranges) in the training are key factors in achieving reliable data after the transfer of the sensor to another location. Moreover, the integration of metadata, such as activity logs, window status, data from official monitoring stations and diffusive samples, was proved essential for validating and interpreting sensor data.

Although this study focused on specific sensor models for $NO_2$ and $PM_{2.5}$, the findings of this study can be generalized to enhance the understanding of the opportunities and limitations of AQSSs for environmental epidemiology. Additional research is necessary, particularly on sensor re-location following co-location with reference-grade instruments, especially in indoor settings, and on drift. While sensors currently face limitations at low pollutant concentrations, they are a promising tool to enhance statistical power and expand epidemiological studies in low- and middle-income countries, as well as in moderately to highly polluted regions.

## Data availability

The data of this study are available from the authors upon request.

## Author contributions

Conceptualization, M.C.-M., U.V. and B.L.; data curation, M.C.-M. and H.G.S.; formal analysis, M.C.-M.; funding acquisition, U.V. and B.L.; investigation, M.C.-M., B.L. and H.G.S.; methodology, M.C.-M., B.L. and U.V.; project

administration, U.V.; resources, B.L. and U.V.; software, H.G.S.; supervision, B.L. and U.V.; validation, M.C.-M.; visualization, M.C.-M. and H.G.S.; writing—original draft, M.C.-M.; writing—review and editing, M.C.-M., B.L. and U.V. All authors have read and agreed to the published version of the manuscript.

**Competing interests**

660 The authors declare that they have no conflict of interest.

**References**

Amegah, A. K.: Proliferation of low-cost sensors. What prospects for air pollution epidemiologic research in Sub-Saharan
670    Africa?, Environ. Pollut., 241, 1132–1137, https://doi.org/10.1016/j.envpol.2018.06.044, 2018.

Apostolopoulos, I. D., Fouskas, G., and Pandis, S. N.: Field Calibration of a Low-Cost Air Quality Monitoring Device in an Urban Background Site Using Machine Learning Models, Atmosphere (Basel), 14, 368, https://doi.org/10.3390/atmos14020368, 2023.

Awad, M. and Khanna, R.: Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System
675    Designers, Apress; Imprint, Berkeley, CA, 300 pp., 2015.

Beloconi, A. and Vounatsou, P.: Bayesian geostatistical modelling of high-resolution NO2 exposure in Europe combining data from monitors, satellites and chemical transport models, Environ. Int., 138, 105578, https://doi.org/10.1016/j.envint.2020.105578, 2020.

Bigi, A., Mueller, M., Grange, S. K., Ghermandi, G., and Hueglin, C.: Performance of NO, $NO_2$ low cost sensors and three
680    calibration approaches within a real world application, Atmos. Meas. Tech., 11, 3717–3735, https://doi.org/10.5194/amt-11-3717-2018, 2018.

Bishop, C. M.: Pattern recognition and machine learning, Information science and statistics, Springer, New York, NY, 738 pp., 2006.

Borrego, C., Costa, A. M., Ginja, J., Amorim, M., Coutinho, M., Karatzas, K., Sioumis, T., Katsifarakis, N., Konstantinidis,
685     K., Vito, S. de, Esposito, E., Smith, P., André, N., Gérard, P., Francis, L. A., Castell, N., Schneider, P., Viana, M.,
        Minguillón, M. C., Reimringer, W., Otjes, R. P., Sicard, O. von, Pohle, R., Elen, B., Suriano, D., Pfister, V., Prato, M.,
        Dipinto, S., and Penza, M.: Assessment of air quality microsensors versus reference methods: The EuNetAir joint
        exercise, Atmos. Environ., 147, 246–263, https://doi.org/10.1016/j.atmosenv.2016.09.050, 2016.

Boser, B. E., Guyon, I. M., and Vapnik, V. N.: A training algorithm for optimal margin classifiers, in: Proceedings of the
690     fifth annual workshop on Computational Learning Theory, Pittsburgh Pennsylvania USA, July 27 - 29, 144–152, 1992.

Braithwaite, I., Zhang, S., Kirkbride, J. B., Osborn, D. P. J., and Hayes, J. F.: Air Pollution (Particulate Matter) Exposure and
        Associations with Depression, Anxiety, Bipolar, Psychosis and Suicide Risk: A Systematic Review and Meta-Analysis,
        Environ. Health Perspect., 127, 126002, https://doi.org/10.1289/EHP4595, 2019.

Breiman, L.: Random Forests, Machine Learning, 45, 5–32, https://doi.org/10.1023/A:1010933404324, 2001.

695 Cai, M., Yin, Y., and Xie, M.: Prediction of hourly air pollutant concentrations near urban arterials using artificial neural
        network approach, Transportation Research Part D: Transport and Environment, 14, 32–41,
        https://doi.org/10.1016/j.trd.2008.10.004, 2009.

Castell, N., Dauge, F. R., Schneider, P., Vogt, M., Lerner, U., Fishbain, B., Broday, D., and Bartonova, A.: Can commercial
        low-cost sensor platforms contribute to air quality monitoring and exposure estimates?, Environ. Int., 99, 293–302,
700     https://doi.org/10.1016/j.envint.2016.12.007, 2017.

CEN/TS 17660-1: Air quality - Performance evaluation of air quality sensor systems - Part 1: Gaseous pollutants in ambient
        air, 2021.

CEN/TS 17660-2: Air quality - Performance evaluation of air quality sensor systems - Part 2: Particulate matter in ambient
        air, 2025.

705 Chacón-Mateos, M., Remy, E., Liebers, U., Heimann, F., Witt, C., and Vogt, U.: Feasibility Study on the Use of NO2 and
        PM2.5 Sensors for Exposure Assessment and Indoor Source Apportionment at Fixed Locations, Sensors (Basel), 24,
        5767, https://doi.org/10.3390/s24175767, 2024.

Chacón-Mateos, M., Laquai, B., Vogt, U., and Stubenrauch, C.: Evaluation of a low-cost dryer for a low-cost optical particle
        counter, Atmos. Meas. Tech., 15, 7395–7410, https://doi.org/10.5194/amt-15-7395-2022, 2022.

710 Chatzidiakou, L., Krause, A., Kellaway, M., Han, Y., Li, Y., Martin, E., Kelly, F. J., Zhu, T., Barratt, B., and Jones, R. L.:
        Automated classification of time-activity-location patterns for improved estimation of personal exposure to air pollution,
        Environ. Health, 21, 125, https://doi.org/10.1186/s12940-022-00939-8, 2022.

Chatzidiakou, L., Krause, A., Han, Y., Chen, W., Yan, L., Popoola, O. A. M., Kellaway, M., Wu, Y., Liu, J., Hu, M., Barratt,
        B., Kelly, F. J., Zhu, T., and Jones, R. L.: Using low-cost sensor technologies and advanced computational methods to

715      improve dose estimations in health panel studies: results of the AIRLESS project, J. Expo. Sci. Environ. Epidemiol., 30, 981–989, https://doi.org/10.1038/s41370-020-0259-6, 2020.

Chatzidiakou, L., Krause, A., Popoola, O. A. M., Di Antonio, A., Kellaway, M., Han, Y., Squires, F. A., Wang, T., Zhang, H., Wang, Q., Fan, Y., Chen, S., Hu, M., Quint, J. K., Barratt, B., Kelly, F. J., Zhu, T., and Jones, R. L.: Characterising low-cost sensors in highly portable platforms to quantify personal exposure in diverse environments, Atmos. Meas.

720      Tech., 12, 4643–4657, https://doi.org/10.5194/amt-12-4643-2019, 2019.

Cordero, J. M., Borge, R., and Narros, A.: Using statistical methods to carry out in field calibrations of low cost air quality sensors, Sens. Actuators B Chem., 267, 245–254, https://doi.org/10.1016/j.snb.2018.04.021, 2018.

D'Elia, G., Ferro, M., Sommella, P., Ferlito, S., Vito, S. de, and Di Francia, G.: Concept Drift Mitigation in Low-Cost Air Quality Monitoring Networks, Sensors (Basel), 24, 2786, https://doi.org/10.3390/s24092786, 2024.

725      deSouza, P., Barkjohn, K., Clements, A., Lee, J., Kahn, R., Crawford, B., and Kinney, P.: An analysis of degradation in low-cost particulate matter sensors, Environ. Sci. Atmos., 3, 521–536, https://doi.org/10.1039/d2ea00142j, 2023.

deSouza, P., Kahn, R., Stockman, T., Obermann, W., Crawford, B., an Wang, Crooks, J., Li, J., and Kinney, P.: Calibrating networks of low-cost air quality sensors, Atmos. Meas. Tech., 15, 6309–6328, https://doi.org/10.5194/amt-15-6309-2022, 2022.

730      Di Antonio, A., Popoola, O. A. M., Ouyang, B., Saffell, J., and Jones, R. L.: Developing a Relative Humidity Correction for Low-Cost Sensors Measuring Ambient Particulate Matter, Sensors, 18, 2790, https://doi.org/10.3390/s18092790, 2018.

Diez, S., Lacy, S., Coe, H., Urquiza, J., Priestman, M., Flynn, M., Marsden, N., Martin, N. A., Gillott, S., Bannan, T., and Edwards, P. M.: Long-term evaluation of commercial air quality sensors: an overview from the QUANT (Quantification of Utility of Atmospheric Network Technologies) study, Atmos. Meas. Tech., 17, 3809–3827,

735      https://doi.org/10.5194/amt-17-3809-2024, 2024.

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, in: Official Journal of the European Union, 2008.

Directive 2024/2881 of the European Parliament and of the Council of 23 October 2024 on ambient air quality and cleaner air for Europe (recast), in: Official Journal of the European Union, 2024.

740      Esposito, E., Vito, S. de, Salvato, M., Bright, V., Jones, R. L., and Popoola, O.: Dynamic neural network architectures for on field stochastic calibration of indicative low cost air quality sensing systems, Sens. Actuators B Chem., 231, 701–713, https://doi.org/10.1016/j.snb.2016.03.038, 2016.

Ester, M., Kriegel, H. P., Sander, J., and Xiaowei, X.: A density-based algorithm for discovering clusters in large spatial databases with noise, in: Proceedings / Second International Conference on Knowledge Discovery & Data Mining,

745      AAAI Press, Menlo Park, Calif., 226–231, 1996.

Evangelopoulos, D., Chatzidiakou, L., Walton, H., Katsouyanni, K., Kelly, F. J., Quint, J. K., Jones, R. L., and Barratt, B.: Personal exposure to air pollution and respiratory health of COPD patients in London, Eur. Respir. J., 58, https://doi.org/10.1183/13993003.03432-2020, 2021.

Gäbel, P., Koller, C., and Hertig, E.: Development of Air Quality Boxes Based on Low-Cost Sensor Technology for
750     Ambient Air Quality Monitoring, Sensors (Basel, Switzerland), 22, 3830, https://doi.org/10.3390/s22103830, 2022.

Géron, A.: Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build
        intelligent systems, Second Edition, O'Reilly, Sebastopol, CA, 2019.

Goldman, G. T., Mulholland, J. A., Russell, A. G., Gass, K., Strickland, M. J., and Tolbert, P. E.: Characterization of
        Ambient Air Pollution Measurement Error in a Time-Series Health Study using a Geostatistical Simulation Approach,
755     Atmos. Environ., 57, 101–108, https://doi.org/10.1016/j.atmosenv.2012.04.045, 2012.

Goulier, L., Paas, B., Ehrnsperger, L., and Klemm, O.: Modelling of Urban Air Pollutant Concentrations with Artificial
        Neural Networks Using Novel Input Variables, Int. J. Environ. Res. Public Health, 17,
        https://doi.org/10.3390/ijerph17062025, 2020.

Han, Y., Chatzidiakou, L., Yan, L., Chen, W., Zhang, H., Krause, A., Xue, T., Chan, Q., Liu, J., Wu, Y., Barratt, B., Jones,
760     R., Zhu, T., and Kelly, F. J.: Difference in ambient-personal exposure to PM2.5 and its inflammatory effect in local
        residents in urban and peri-urban Beijing, China: results of the AIRLESS project, Faraday Discuss., 226, 569–583,
        https://doi.org/10.1039/d0fd00097c, 2021.

Han, Y., Chen, W., Chatzidiakou, L., Krause, A., Yan, L., Zhang, H., Chan, Q., Barratt, B., Jones, R., Liu, J., Wu, Y., Zhao,
        M., Zhang, J., Kelly, F. J., and Zhu, T.: Effects of AIR pollution on cardiopuLmonary disEaSe in urban and peri-urban
765     reSidents in Beijing: protocol for the AIRLESS study, Atmos. Chem. Phys., 20, 15775–15792,
        https://doi.org/10.5194/acp-20-15775-2020, 2020.

Hang, Y., Meng, X., Li, T., Wang, T., Cao, J., Fu, Q., Dey, S., Li, S., Huang, K., Liang, F., Kan, H., Shi, X., and Liu, Y.:
        Assessment of long-term particulate nitrate air pollution and its health risk in China, iScience, 25, 104899,
        https://doi.org/10.1016/j.isci.2022.104899, 2022.

770     Harré, E. S., Price, P. D., Ayrey, R. B., Toop, L. J., Martin, I. R., and Town, G. I.: Respiratory effects of air pollution in
        chronic obstructive pulmonary disease: a three month prospective study, Thorax, 52, 1040–1044,
        https://doi.org/10.1136/thx.52.12.1040, 1997.

Hoffmann, C., Maglakelidze, M., Schneidemesser, E. von, Witt, C., Hoffmann, P., and Butler, T.: Asthma and COPD
        exacerbation in relation to outdoor air pollution in the metropolitan area of Berlin, Germany, Respir. Res., 23, 64,
775     https://doi.org/10.1186/s12931-022-01983-1, 2022.

Huang, C., Hu, J., Xue, T., Xu, H., and Wang, M.: High-Resolution Spatiotemporal Modeling for Ambient PM2.5 Exposure
        Assessment in China from 2013 to 2019, Environ. Sci. Technol., 55, 2152–2162,
        https://doi.org/10.1021/acs.est.0c05815, 2021.

Jerrett, M., Donaire-Gonzalez, D., Popoola, O., Jones, R., Cohen, R. C., Almanza, E., Nazelle, A. de, Mead, I., Carrasco-
780     Turigas, G., Cole-Hunter, T., Triguero-Mas, M., Seto, E., and Nieuwenhuijsen, M.: Validating novel air pollution
        sensors to improve exposure estimates for epidemiological analyses and citizen science, Environ. Res., 158, 286–294,
        https://doi.org/10.1016/j.envres.2017.04.023, 2017.

Karagulian, F., Barbiere, M., Kotsev, A., Spinelle, L., Gerboles, M., Lagler, F., Redon, N., Crunaire, S., and Borowiak, A.:
    Review of the Performance of Low-Cost Sensors for Air Quality Monitoring, Atmosphere (Basel), 10, 506,
785    https://doi.org/10.3390/atmos10090506, 2019.

Koehler, K., Wilks, M., Green, T., Rule, A. M., Zamora, M. L., Buehler, C., Datta, A., Gentner, D. R., Putcha, N., Hansel, N.
    N., Kirk, G. D., Raju, S., and McCormack, M.: Evaluation of Calibration Approaches for Indoor Deployments of
    PurpleAir Monitors, Atmos. Environ. (1994), 310, https://doi.org/10.1016/j.atmosenv.2023.119944, 2023.

Koehler, K., Good, N., Wilson, A., Mölter, A., Moore, B. F., Carpenter, T., Peel, J. L., and Volckens, J.: The Fort Collins
790    commuter study: Variability in personal exposure to air pollutants by microenvironment, Indoor air, 29, 231–241,
    https://doi.org/10.1111/ina.12533, 2019.

Krause, A.: Using novel portable air quality monitors to improve personal exposure and dose estimations for health studies,
    Apollo - University of Cambridge Repository, 2021.

Laquai, B. and Saur, A.: Development of a Calibration Methodology for the SDS011 Low-Cost PM-Sensor with respect to
795    Professional Reference Instrumentation,
    https://www.researchgate.net/publication/322628807_Development_of_a_Calibration_Methodology_for_the_SDS011_
    Low-Cost_PM-Sensor_with_respect_to_Professional_Reference_Instrumentation, last access: 2 December 2024, 2017.

Li, J., Hauryliuk, A., Malings, C., Eilenberg, S. R., Subramanian, R., and Presto, A. A.: Characterizing the Aging of
    Alphasense NO2 Sensors in Long-Term Field Deployments, ACS Sens., 6, 2952–2959,
800    https://doi.org/10.1021/acssensors.1c00729, 2021.

Licina, D., Tian, Y., and Nazaroff, W. W.: Emission rates and the personal cloud effect associated with particle release from
    the perihuman environment, Indoor air, 27, 791–802, https://doi.org/10.1111/ina.12365, 2017.

Liu, H.-Y., Schneider, P., Haugen, R., and Vogt, M.: Performance Assessment of a Low-Cost PM2.5 Sensor for a near Four-
    Month Period in Oslo, Norway, Atmosphere (Basel), 10, 41, https://doi.org/10.3390/atmos10020041, 2019.

805 Mainka, A. and Żak, M.: Synergistic or Antagonistic Health Effects of Long- and Short-Term Exposure to Ambient NO2
    and PM2.5: A Review, Int. J. Environ. Res. Public Health, 19, https://doi.org/10.3390/ijerph192114079, 2022.

Malings, C., Tanzer, R., Hauryliuk, A., Kumar, S. P. N., Zimmerman, N., Kara, L. B., and Presto, A. A.: Development of a
    general calibration model and long-term performance evaluation of low-cost sensors for air pollutant gas monitoring,
    Atmos. Meas. Tech., 12, 903–920, https://doi.org/10.5194/amt-12-903-2019, 2019.

810 Manisalidis, I., Stavropoulou, E., Stavropoulos, A., and Bezirtzoglou, E.: Environmental and Health Impacts of Air
    Pollution: A Review, Front. Public Health, 8, 14, https://doi.org/10.3389/fpubh.2020.00014, 2020.

McCulloch, W. S. and Pitts, W.: A logical calculus of the ideas immanent in nervous activity, Bulletin of Mathematical
    Biophysics, 5, 115–133, https://doi.org/10.1007/BF02478259, 1943.

Meng, X., Wang, C., Cao, D., Wong, C.-M., and Kan, H.: Short-term effect of ambient air pollution on COPD mortality in
815    four Chinese cities, Atmos. Environ., 77, 149–154, https://doi.org/10.1016/j.atmosenv.2013.05.001, 2013.

Müller, A. C. and Guido, S.: Introduction to machine learning with Python: A guide for data scientists, O'Reilly, Beijing, 2017.

Novak, R., Robinson, J. A., Kanduč, T., Sarigiannis, D., and Kocman, D.: Assessment of Individual-Level Exposure to Airborne Particulate Matter during Periods of Atmospheric Thermal Inversion, Sensors (Basel, Switzerland), 22, 7116,

820     https://doi.org/10.3390/s22197116, 2022.

Novak, R., Petridis, I., Kocman, D., Robinson, J. A., Kanduč, T., Chapizanis, D., Karakitsios, S., Flückiger, B., Vienneau, D., Mikeš, O., Degrendele, C., Sáňka, O., García Dos Santos-Alves, S., Maggos, T., Pardali, D., Stamatelopoulou, A., Saraga, D., Persico, M. G., Visave, J., Gotti, A., and Sarigiannis, D.: Harmonization and Visualization of Data from a Transnational Multi-Sensor Personal Exposure Campaign, Int. J. Environ. Res. Public Health, 18, 11614,

825     https://doi.org/10.3390/ijerph182111614, 2021.

Paas, B., Stienen, J., Vorländer, M., and Schneider, C.: Modelling of Urban Near-Road Atmospheric PM Concentrations Using an Artificial Neural Network Approach with Acoustic Data Input, Environments, 4, 26, https://doi.org/10.3390/environments4020026, 2017.

Pantelic, J., Liu, S., Pistore, L., Licina, D., Vannucci, M., Sadrizadeh, S., Ghahramani, A., Gilligan, B., Sternberg, E.,

830     Kampschroer, K., and Schiavon, S.: Personal CO2 cloud: laboratory measurements of metabolic CO2 inhalation zone concentration and dispersion in a typical office desk setting, J. Expo. Sci. Environ. Epidemiol., 30, 328–337, https://doi.org/10.1038/s41370-019-0179-5, 2020.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Müller, A., Nothman, J., Louppe, G., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot,

835     M., and Duchesnay, É.: Scikit-learn: Machine Learning in Python, J. Mach. Learn. Res., 12, 2825–2830, 2011.

Piechocki-Minguy, A., Plaisance, H., Schadkowski, C., Sagnier, I., Saison, J. Y., Galloo, J. C., and Guillermo, R.: A case study of personal exposure to nitrogen dioxide using a new high sensitive diffusive sampler, Sci. Total Environ., 366, 55–64, https://doi.org/10.1016/j.scitotenv.2005.08.009, 2006.

Piedrahita, R., Xiang, Y., Masson, N., Ortega, J., Collier, A., Jiang, Y., Li, K., Dick, R. P., Lv, Q., Hannigan, M., and Shang,

840     L.: The next generation of low-cost personal air quality sensors for quantitative exposure monitoring, Atmos. Meas. Tech., 7, 3325–3336, https://doi.org/10.5194/amt-7-3325-2014, 2014.

Rea, A. W., Zufall, M. J., Williams, R. W., Sheldon, L., and Howard-Reed, C.: The influence of human activity patterns on personal PM exposure: a comparative analysis of filter-based and continuous particle measurements, J. Air Waste Manag. Assoc., 51, 1271–1279, https://doi.org/10.1080/10473289.2001.10464351, 2001.

845  Russell, H. S., Frederickson, L. B., Kwiatkowski, S., Emygdio, A. P. M., Kumar, P., Schmidt, J. A., Hertel, O., and Johnson, M. S.: Enhanced Ambient Sensing Environment-A New Method for Calibrating Low-Cost Gas Sensors, Sensors (Basel), 22, 7238, https://doi.org/10.3390/s22197238, 2022.

Sahu, R., Nagal, A., Dixit, K. K., Unnibhavi, H., Mantravadi, S., Nair, S., Simmhan, Y., Mishra, B., Zele, R., Sutaria, R., Motghare, V. M., Kar, P., and Tripathi, S. N.: Robust statistical calibration and characterization of portable low-cost air

850    quality monitoring sensors to quantify real-time $O_3$ and $NO_2$ concentrations in diverse environments, Atmos. Meas.
        Tech., 14, 37–52, https://doi.org/10.5194/amt-14-37-2021, 2021.

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
        (Basel, Switzerland), 20, 5175, https://doi.org/10.3390/s20185175, 2020.

855    Samon, S. M., Hammel, S. C., Stapleton, H. M., and Anderson, K. A.: Silicone wristbands as personal passive sampling
        devices: Current knowledge, recommendations for use, and future directions, Environ. Int., 169, 107339,
        https://doi.org/10.1016/j.envint.2022.107339, 2022.

Sarnat, J. A., Schwartz, J., Catalano, P. J., and Suh, H. H.: Gaseous pollutants in particulate matter epidemiology:
        confounders or surrogates?, Environ. Health Perspect., 109, 1053–1061, https://doi.org/10.1289/ehp.011091053, 2001.

860    Scott Downen, R., Dong, Q., Chorvinsky, E., Li, B., Tran, N., Jackson, J. H., Pillai, D. K., Zaghloul, M., and Li, Z.: Personal
        NO2 sensor demonstrates feasibility of in-home exposure measurements for pediatric asthma research and management,
        J. Expo. Sci. Environ. Epidemiol., 32, 312–319, https://doi.org/10.1038/s41370-022-00413-0, 2022.

Shaw, P. A., Deffner, V., Keogh, R. H., Tooze, J. A., Dodd, K. W., Küchenhoff, H., Kipnis, V., and Freedman, L. S.:
        Epidemiologic analyses with error-prone exposures: review of current practice and recommendations, Ann. Epidemiol.,
865    28, 821–828, https://doi.org/10.1016/j.annepidem.2018.09.001, 2018.

Shirdel, M., Bergdahl, I. A., Andersson, B. M., Wingfors, H., Sommar, J. N., and Liljelind, I. E.: Passive personal air
        sampling of dust in a working environment-A pilot study, J. Occup. Environ. Hyg., 16, 675–684,
        https://doi.org/10.1080/15459624.2019.1648814, 2019.

Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Field calibration of a cluster of low-cost
870    available sensors for air quality monitoring. Part A: Ozone and nitrogen dioxide, Sens. Actuators B Chem., 215, 249–
        257, https://doi.org/10.1016/j.snb.2015.03.031, 2015.

Steinle, S., Reis, S., and Sabel, C. E.: Quantifying human exposure to air pollution--moving from static monitoring to spatio-
        temporally resolved personal exposure assessment, Sci. Total Environ., 443, 184–193,
        https://doi.org/10.1016/j.scitotenv.2012.10.098, 2013.

875    Suriano, D. and Penza, M.: Assessment of the Performance of a Low-Cost Air Quality Monitor in an Indoor Environment
        through Different Calibration Models, Atmosphere (Basel), 13, 567, https://doi.org/10.3390/atmos13040567, 2022.

Tancev, G.: Relevance of Drift Components and Unit-to-Unit Variability in the Predictive Maintenance of Low-Cost
        Electrochemical Sensor Systems in Air Quality Monitoring, Sensors (Basel, Switzerland), 21, 32908,
        https://doi.org/10.3390/s21093298, 2021.

880    Thunis, P., Georgieva, E., and Pederzoli, A.: A tool to evaluate air quality model performances in regulatory applications,
        Environ. Model. Softw., 38, 220–230, https://doi.org/10.1016/j.envsoft.2012.06.005, 2012.

Topalović, D. B., Davidović, M. D., Jovanović, M., Bartonova, A., Ristovski, Z., and Jovašević-Stojanović, M.: In search of
        an optimal in-field calibration method of low-cost gas sensors for ambient air pollutants: Comparison of linear,

multilinear and artificial neural network approaches, Atmos. Environ., 213, 640–658,

885    https://doi.org/10.1016/j.atmosenv.2019.06.028, 2019.

Tryner, J., Phillips, M., Quinn, C., Neymark, G., Wilson, A., Jathar, S. H., Carter, E., and Volckens, J.: Design and Testing
of a Low-Cost Sensor and Sampling Platform for Indoor Air Quality, Build. Environ., 206,
https://doi.org/10.1016/j.buildenv.2021.108398, 2021.

VDI 2453 Part 1: Gaseous air pollution measurement; determination of nitrogen dioxide concentration; photometric manual

890    standard method (Saltzmann), 1990.

Venkatraman Jagatha, J., Klausnitzer, A., Chacón-Mateos, M., Laquai, B., Nieuwkoop, E., van der Mark, P., Vogt, U., and
Schneider, C.: Calibration Method for Particulate Matter Low-Cost Sensors Used in Ambient Air Quality Monitoring
and Research, Sensors, 21, 3960, https://doi.org/10.3390/s21123960, 2021.

Wei, P., Sun, L., Anand, A., Zhang, Q., Huixin, Z., Deng, Z., Wang, Y., and Ning, Z.: Development and evaluation of a

895    robust temperature sensitive algorithm for long term NO2 gas sensor network data correction, Atmos. Environ., 230,
117509, https://doi.org/10.1016/j.atmosenv.2020.117509, 2020.

WHO: WHO global air quality guidelines. Particulate matter (PM2.5 and PM10), ozone, nitrogen dioxide, sulfur dioxide and
carbon monoxide, Geneva, 2021.

Zamora, M. L., Buehler, C., Lei, H., Datta, A., Xiong, F., Gentner, D. R., and Koehler, K.: Evaluating the Performance of

900    Using Low-Cost Sensors to Calibrate for Cross-Sensitivities in a Multipollutant Network, ACS ES&T engineering, 2,
780–793, https://doi.org/10.1021/acsestengg.1c00367, 2022.

Zauli-Sajani, S., Marchesi, S., Pironi, C., Barbieri, C., Poluzzi, V., and Colacci, A.: Assessment of air quality sensor system
performance after relocation, Atmos. Pollut. Res., 12, 282–291, https://doi.org/10.1016/j.apr.2020.11.010, 2021.

Zimmerman, N., Presto, A. A., Kumar, S. P. N., Gu, J., Hauryliuk, A., Robinson, E. S., and Robinson, A. L.: A machine

905    learning calibration model using random forests to improve sensor performance for lower-cost air quality monitoring,
Atmos. Meas. Tech., 11, 291–313, https://doi.org/10.5194/amt-11-291-2018, 2018.