

1 **Leveraging Machine Learning Algorithms to Advance Low-Cost Air**
2 **Sensor Calibration in Stationary and Mobile Settings**

3

4

5 An Wang¹, Yuki Machida¹, Priyanka deSouza², Simone Mora^{1,3*}, Tiffany Duhl⁴, Neelakshi
6 Hudda⁴, John L. Durant⁴, Fabio Durate¹, and Carlo Ratti¹

7 1. Senseable City Lab, Department of Urban Studies and Planning, Massachusetts Institute of
8 Technology

9 2. Department of Urban and Regional Planning, University of Colorado Denver

10 3. Department of Computer Science, Norwegian University of Science and Technology

11 4. Department of Civil and Environmental Engineering, Tufts University

12 * corresponding author

13

1 **ABSTRACT**

2 Low-cost air sensing is changing the paradigm of ambient air quality management research
3 and practices. However, consensus on a structured low-cost sensor calibration and performance
4 evaluation framework is lacking. Our study aims to devise a standardized low-cost sensor
5 calibration protocol and evaluate the performance of various calibration algorithms. Extensive
6 collocation data were collected in stationary and mobile settings in two American cities, New York
7 and Boston. We trained the calibration models using stationary data aggregated at various intervals
8 to examine the performance of several commonly used calibration algorithms described in the
9 literature. Linear models provide consistently satisfactory calibration results, indicating linear
10 responses from the low-cost sensors in our stationary test environment. Its simplicity is
11 recommended for citizen science and education usages. Models that can account for non-linear
12 relationships, especially random forest, perform well and transfer between sensors better than
13 generalized linear regression models for PM_{2.5} calibration, which should be adopted for regulatory
14 and scientific purposes. Data collected in a mobile validation campaign in Boston were passed
15 through the best-performing calibration models to assess their transferability. The results indicate
16 that models trained with data from a different urban environment and season in the stationary
17 setting did not transfer well to a mobile setting. It is recommended that low-cost sensors should be
18 calibrated more often than suggested in Environmental Protection Agency’s air sensor
19 performance evaluation guidelines and used in an environment that is as similar as possible to the
20 calibration environment.

21 Keywords: low-cost sensor calibration, PM_{2.5}, NO₂, machine learning, mobile monitoring,
22 environmental justice

23

1 1 INTRODUCTION

2 Lacking the resources to deploy high-quality monitors, low- and middle-income countries
3 are held back from making effective measurements to inform air pollution management. The rapid
4 development of low-cost sensors (< \$2500(US EPA, 2014)) in recent years provides a unique
5 opportunity to shift the current air quality monitoring paradigm. While low-cost sensors have been
6 widely adopted in high-income regions to supplement traditional air quality monitoring and
7 mapping of local air quality (Castell et al., 2017; Crawford et al., 2021; Gressent et al., 2020;
8 Miskell et al., 2018), their applications in low-income countries are limited (Brauer et al., 2019;
9 deSouza et al., 2020; SM et al., 2019). The benefits of incorporating low-cost sensors in the
10 existing monitoring network are multifold, including improving public awareness of air pollution,
11 pushing for environmentally just decision-making and reducing *data colonialism* in air quality
12 monitoring, where government agencies and high-income countries claim ownership of collected
13 air quality data (Duarte & deSouza, 2020).

14 Despite the advantages in cost efficiency, flexibility, and ease of use, low-cost sensors
15 suffer from constant data quality and stability issues; therefore, sensor collocation and calibration
16 are of utmost importance prior to field deployment. Collocation is the process of deploying sensors
17 side-by-side with reference monitors, and calibration involves adjusting raw sensor readings using
18 collocation data and mathematical methods. We summarized seven issues that affect low-cost
19 sensor performance from low-cost sensor calibration literature within the last five years: inter-
20 sensor variability, intra-sensor variability, drift, aging, response time, cross-sensitivity, and
21 sensitivity to environmental factors. Table S1 presents detailed findings from each study.

22 Inter-sensor variability refers to the variability in measurements using multiple identical
23 sensors under the same testing environment. The calibration for inter-sensor variability is crucial
24 for low-cost sensors as it is the foundation of large-scale sensor deployment and data transferability.
25 Intra-sensor variability describes the variability in consecutive measurements made by a given
26 sensor under a similar testing environment. Drift is the gradual change in sensor response over
27 time. Aging refers to the continuous deterioration of sensor performance over time. Unlike
28 reference sensors, low-cost sensors are prone to drifting and have a much shorter lifetime. Thus,
29 they require routine sensor calibration and replacement. Response time reflects the lag before
30 sensors reach stable readings in the test environment, parity between response time and temporal

1 scale of change in signal is critical for time-resolved sampling. Cross-sensitivity, an issue exclusive
2 to gas sensors, denotes a sensor's false response to gases other than the target gas. Finally,
3 sensitivity to environmental factors, including temperature, humidity, wind, barometric pressure,
4 and particle composition, is ubiquitous in both low-cost and reference sensors. These
5 environmental factors are commonly identified as the main explanatory features in low-cost sensor
6 calibration models. Acknowledging that there might be no single solution to all seven issues in
7 given applications, it is essential to characterize the performance of various calibration algorithms
8 and identify the proper context for their usage.

9 In light of low-cost sensor performance issues and the surging interest in using sensors for
10 air quality monitoring in areas where regulatory monitors are absent, the US Environmental
11 Protection Agency (EPA) has been engaging local communities, subject matter experts, and air
12 sensor manufacturers to develop guidelines for low-cost sensor usage. EPA published two
13 performance target reports for O₃ and particulate matter (PM) sensors to support low-cost sensor
14 non-regulatory supplemental and informational monitoring (NSIM) applications (US EPA, 2021a,
15 2021b). These reports provide low-cost sensor collocation, calibration, performance evaluation,
16 and deployment protocols. The reports suggest using linear regression to calibrate low-cost sensor
17 data against the reference station for ease of use; however, linear calibration may not meet
18 performance targets in many cases (Kelly et al., 2017; Malings et al., 2019).

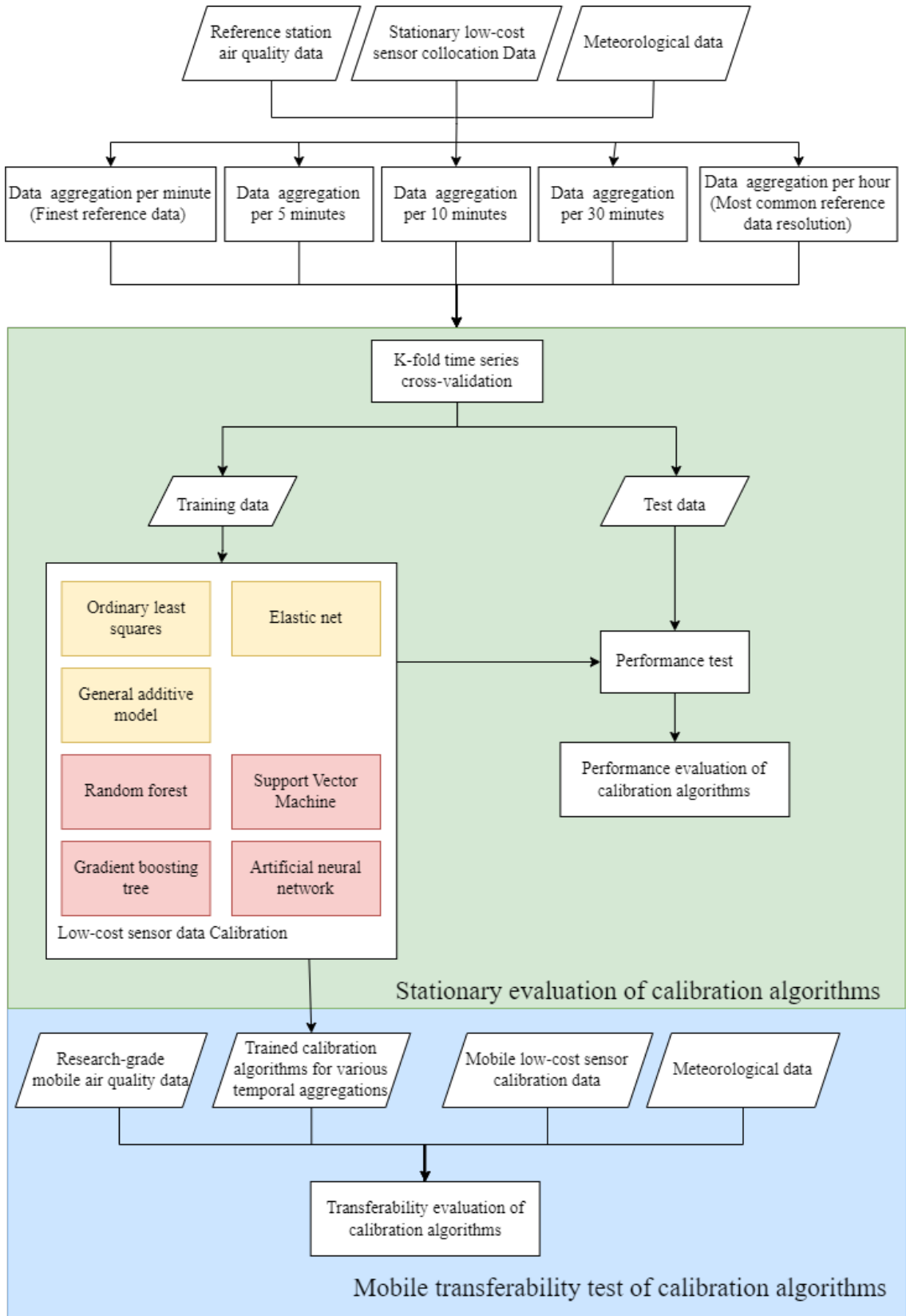
19 The three main categories of low-cost sensor calibration algorithms currently in use include
20 physical mechanism-based models, parametric models such as linear and non-linear regression
21 models, and machine learning models (Liang, 2021). To date, physical mechanism-based models
22 have targeted the problem that relative humidity can change PM hygroscopic size and refractive
23 index, influencing conversion from PM number count to mass concentrations. Crilley et al. (2018)
24 and Zheng et al. (2018) applied κ -Köhler theory and scattering efficiency, respectively, to adjust
25 relative humidity correction factors for better low-cost PM sensor calibration. Parametric models,
26 including linear and non-linear regression models, are more widely used in the literature, as shown
27 in Table S1. They can consider more features affecting low-cost sensor readings, while physical
28 mechanism-based methods mainly calibrate for relative humidity. As state-of-the-art calibration
29 algorithms, machine learning approaches can provide higher accuracy in data-rich environments

1 (Liang, 2021; Malings et al., 2019; Wang et al., 2021). Nonetheless, some machine learning
2 algorithms lack interpretability and transferability, which hinders their broader implementation.

3 Given that there is no all-in-one solution to low-cost sensor calibration, it is necessary to
4 explore various calibration approaches in different applications and compare their performance.
5 To this end, we tested multiple commonly used low-cost sensor calibration algorithms in stationary
6 and mobile settings. The calibration models trained with stationary collocation data were further
7 examined with mobile collected data for transferability tests. We were motivated to simulate a
8 common scenario deploying sensors calibrated in one location in a different location that does not
9 have high-quality air instruments to enable local collocation and calibration. This is especially
10 useful for sensor applications in low-income regions and countries to address the existing air
11 quality data gap and related environmental justice issues. This study is among the first to
12 incorporate the standardized EPA air sensor evaluation frameworks and the effects of data
13 aggregation in stationary and mobile settings. Beyond simply comparing model performance by
14 investigating algorithm behaviors and the impact factors of sensor performance, this study explores
15 the appropriate application scenarios for multiple calibration algorithms to account for non-
16 linearity in low-cost sensor performance. Finally, our study provides valuable information on
17 algorithm transferability by exploiting two large air monitoring datasets collected in different
18 locations under stationary and mobile conditions.

19 **2 METHODOLOGY**

20 Our study included two main components: developing calibration models and evaluating
21 the transferability of the models. As shown in Figure 1, reference station data, raw low-cost sensing
22 data, and meteorological data during the collocation were first aggregated at various temporal
23 aggregations to test the effects of data aggregation on the algorithms. We employed a k-fold time
24 series cross-validation to avoid overfitting and performance overestimation, where past
25 observations were used to predict future observations but not the other way around. Two major
26 categories of calibration algorithms were tested: the generalized linear regressions (yellow) and
27 the machine learning algorithms (red). The best-performing model of each calibration algorithm
28 was passed to the transferability test using mobile data collected with research-grade mobile
29 sensors.



1

2

Figure 1 Methodology flow chart

1 2.1 Low-Cost Sensors

2 The mobile low-cost sensing platform we evaluate is developed as part of the City Scanner
3 initiative at the Senseable City Lab, Massachusetts Institute of Technology. It is designed to
4 address big data needs for urban development and planning while helping advance environmental
5 justice in global cities. For simplicity, the sensing platform is referred to as the ‘City Scanner’ or
6 ‘CS’ hereafter. Each CS unit is equipped with an Alphasense OPC-N3 and an Alphasense NO₂-
7 A43F that measure particle counts and nitrogen dioxide, respectively. The Alphasense OPC-N3 is
8 an optical particle counter with a nominal particle monitoring size range from 0.35 to 40
9 micrometers, which has been widely adopted in multiple previous studies (Bezantakos et al., 2018;
10 Crilley et al., 2018, 2020; Sousan et al., 2016). It provides particle counts in 24 size bins and
11 estimates PM₁, PM_{2.5}, and PM₁₀ mass concentrations based on particle shape and density
12 presumptions. Alphasense NO₂-A43F is an electro-chemical gas sensor that can operate in a
13 variety of ambient environments but is known for its cross-sensitivity issue with NO and CO₂,
14 while its cross-sensitivity with ozone is minimized with an ozone scrubber. The NO₂ readings are
15 given as electric signals (millivoltages) and need to be converted to parts per billion by volume.
16 The low-cost sensors are integrated into the CS, which is a standalone sensing platform with data
17 storage and remote monitoring capability. The platform is powered by a battery that can be charged
18 through a power supply or a solar panel. The data from the sensors are stored on an onboard SD
19 card, which can be accessed remotely through an LTE connection for routine checks on the data
20 from sensors and the device’s status.

21 2.2 Stationary and Mobile Data Collection

22 As per EPA’s air sensor calibration guidelines (US EPA, 2021a, 2021b), the stationary
23 collocation lasted four weeks, from August 13 to September 10, 2021. With permission from the
24 New York State Department of Environmental Conservation, five CS units were placed next to a
25 reference station collecting PM and NO₂ data every five seconds. The reference station (Site ID:
26 IS 52/MS 302) is located on the roof of New York City Department of Education Public School
27 52/Middle School 302 in South Bronx. The station measures criteria pollutants, particulate matter
28 speciation, hydrocarbons, and volatile organic chemicals continuously or intermittently. We
29 obtained continuous minute-by-minute measurements of PM_{2.5} and NO₂ at the station from

1 NYSDEC, measured by Teledyne T640 at 5.0 liter per minute (Federal Equivalent Method 236)
2 and Thermo Environmental Instruments 42C (Federal Reference Method 074), respectively.

3 In addition to the stationary collocation, a mobile validation comparison campaign was
4 carried out using the Tufts Air Pollution Laboratory (TAPL), which measured real-time particulate
5 matter and NO_x concentrations (Hudda et al., 2020; Padró-Martínez et al., 2012). TAPL is housed
6 in an electric vehicle to avoid self-pollution and equipped with research-grade instruments,
7 including a DustTrak™ DRX Aerosol Monitor (Model 8533, TSI, Shoreview, MN) for size-
8 resolved mass fraction concentrations and a NO-NO₂-NO_x Analyzer (Model 42i, Thermo Fisher,
9 Waltham, MA) for NO/NO₂/NO_x concentrations. TAPL instruments were synchronized with CS
10 units using GPS time, reporting PM_{2.5} and NO₂ readings every second and every ten seconds,
11 respectively. TAPL Dusttrak was calibrated against a regulatory site in Boston managed by the
12 Massachusetts Department of Environmental Protection (Nubian Square, site ID 250250042).
13 Detailed information on this calibration is documented in the SI. Two CS units were mounted on
14 the TAPL for mobile validation in Boston. The mobile validation campaign occurred in
15 neighborhoods north of Boston Logan International Airport and lasted from February 11 to April
16 8, 2022, with data collected on 16 different days. It is worth noticing that between the New York
17 collocation and Boston mobile validation with TAPL, from September 13 to December 16, we
18 continuously used these CS units for opportunistic mobile air quality for four months (unattended
19 in an open environment) on municipal service vehicles. This study does not report these data as
20 they are irrelevant to sensor calibration.

21 2.3 Data Processing and Calibration

22 2.3.1 *Data quality assurance and aggregation*

23 In total, during the stationary collocation, 1.6 million unique 5-second resolution data
24 points were collected by the five CS units in New York. During the mobile validation, 130,000
25 unique 5-second-interval data points were recorded by two CS units in Boston. The CS data were
26 then filtered by two criteria before they were aggregated. Fifteen percent of raw CS data from
27 stationary collocation were excluded in the process, where 12% percent were removed due to high
28 relative humidity (> 90% or raining), while the other 3% due to extreme readings (< 1 µg/m³
29 or >1000 µg/m³ for PM_{2.5}, < 200 mv or > 900 mv for NO₂ electro-signal).

1 The filtered raw CS data were aggregated into five temporal resolutions, including minute-
2 level, 5-minute-level, 10-minute-level, 30-minute-level, and hour-level. Minute-by-minute
3 reference air quality data were aggregated in the same intervals. Hourly meteorological data were
4 gathered from the nearest LaGuardia Airport, about 4 kilometers away. Raw CS data, reference
5 data, and meteorology were synchronized and matched, yielding five datasets of various temporal
6 aggregations to explore data aggregation's effect on sensor calibration. Similarly, mobile
7 validation CS data were cleaned using the same criteria as in the stationary collocation. The filtered
8 CS and TAPL data were matched and aggregated to the five temporal resolutions. Hourly
9 meteorological data were obtained from the Boston Logan International Airport (within 5 km of
10 the driving area) and matched with air quality data.

11 *2.3.2 Development of calibration models*

12 Our study evaluated seven distinct algorithms in two categories, the generalized linear
13 regressions and machine learning algorithms, which frequently appear in the existing literature.
14 The algorithms were selected in light of a study that assessed the air quality prediction performance
15 by Kerckhoffs et al. (2019). Ordinary least squares (OLS) is the most widely adopted and
16 straightforward algorithm in sensor calibration, which is also recommended by the EPA air sensor
17 performance evaluation guidelines. It is an effective tool for sensor calibration providing intuitive
18 results and straightforward interpretation. Elastic net is essentially an extension of linear regression
19 with penalty terms (Lasso and/or Ridge regularization) for correlated explanatory features. While
20 we acknowledge that there is great variability and capability in linear regression algorithms, we
21 only adopted the simplest form of these algorithms as recommended by the EPA guidelines for
22 citizen science and educational purposes without considering feature transformations. In addition,
23 feature transformation and smoothing are taken into account in the generalized additive model
24 (GAM), which estimates the target feature with linear combinations of smooth functions of the
25 features. GAM models in this study did not consider interactions between the features for better
26 interpretability.

27 Regarding machine learning algorithms, supporter vector regression (SVR) is a commonly
28 used regression method that uses a kernel function to transform the data into a higher dimension
29 and finds the hyperplane that fits the data with a desirable error margin. Random forest and Light
30 Gradient Boosting Machine (LightGBM) belong to the bigger ensemble modeling family but

1 employ bagging and boosting techniques, respectively. Specifically, LightGBM is considered the
2 state-of-the-art gradient boosting framework originally developed by Microsoft, which has faster
3 training speed, higher efficiency, and lower memory usage while yielding better accuracy (Ke et
4 al., 2017). Artificial neural network (ANN) is popular in regression for prediction purposes, yet its
5 black box nature hinders its wider adoption in explanatory tasks.

6 For each temporally aggregated dataset, we adopted a k-fold time series cross-validation
7 approach to evaluate calibration models' performance (Pedregosa et al., 2011). The approach is
8 designed to avoid predicting past observations from future ones, limiting data leakage and
9 performance overestimation. Each dataset is first ranked temporally and divided into five
10 sequential sections with the same number of observations. In this case, the first k folds are used to
11 train the models, and the k+1 fold split is used to validate and test the model. The train-test split is
12 illustrated in Figure 2. In the training and test sets, the CS readings for PM_{2.5} and NO₂, and
13 meteorological factors, including temperature in Celsius, relative humidity in percentage, dew
14 point in Celsius, and air pressure in kilopascal, were treated as explanatory features, and the
15 reference pollutant readings were the target feature. Other meteorological factors, such as wind
16 speed, wind direction, and feel temperature, were also tested in the models but excluded due to
17 statistical insignificance. Detailed feature descriptions and model-tuning information are
18 documented in the SI. All features were tested for their normality and transformed if not following
19 a normal distribution. It is worth noting that unlike the PM sensor, which provides mass
20 concentration, the NO₂ sensor yields electric signals instead of parts per billion by volume as the
21 raw readings. Conventionally, one converts electro-signals to gas concentrations using sensitivity
22 factors provided by the manufacturer from factory calibration, which is conducted for every batch
23 of sensors. Then, the converted gas concentrations are calibrated with reference concentrations
24 collected in the field. Our approach directly mapped the electric signals to reference concentrations
25 rather than using factory-calibrated sensitivity factors, which can be different between sensor
26 batches and taking average within each batch.

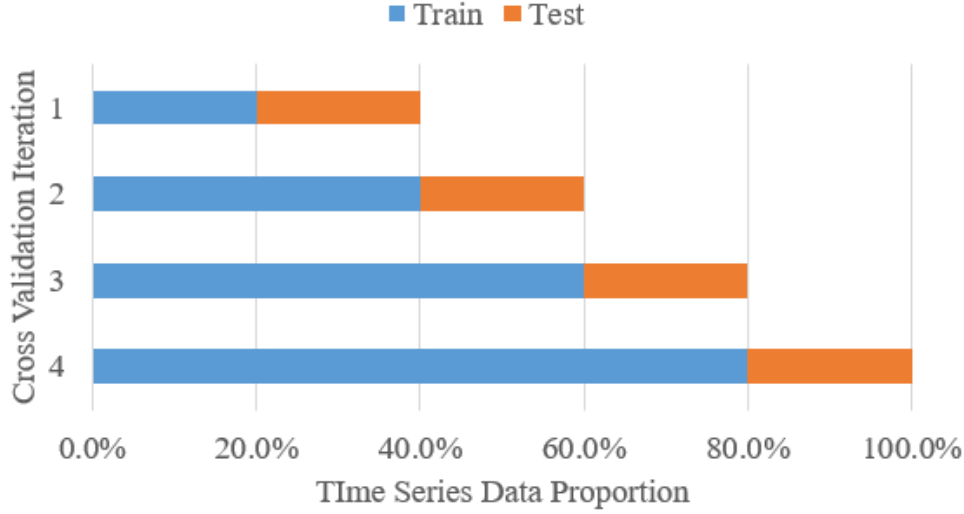


Figure 2. Illustration of k-fold time series cross-validation splitCalibration Algorithm Performance Evaluation and Transferability tests

2.3.3 Performance measures

In each fold of the time series cross-validation, the trained model was assessed against the test dataset using two performance measures, r^2 and the root mean squared error (RMSE) as shown in Equations 1 and 2. r^2 simply represents the square of the Pearson correlation coefficient between predictions and observations. It measures the correlation between calibrated CS readings and reference readings, while RMSE quantifies the absolute difference between them. It is worth noting that even though we train models with log-transformed CS and reference air pollutant readings, both performance measures are calculated using back-transformed predictions and observations.

$$r^2 = \left(\frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}} \right)^2 \quad \text{Eq.1}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad \text{Eq.2}$$

Where:

y_i and \hat{y}_i are observed and predicted target features in the test dataset;

\bar{y} and $\bar{\hat{y}}$ are mean values of observed and predicted target features in the test dataset;

n is the size of the test dataset.

1 The best-performing models and temporally aggregated dataset were selected for each
2 algorithm (7 algorithms \times 5 temporal resolutions = 35 models). Mobile CS data collected by TAPL
3 were passed to the models, whose transferability was evaluated using the same performance
4 measures.

5 2.3.4 *Model interpretation*

6 In addition to performance evaluation, we explored the impact of meteorological factors
7 on sensor and model performance. We focus the model interpretation on the ensemble algorithms,
8 as they are widely adopted in the existing calibration literature to account for non-linear
9 relationships with an established global interpretation method. We employed the Shapley Additive
10 exPlanations (SHAP) method (Lundberg et al., 2018; Lundberg & Lee, 2017) to interpret the
11 trained random forest and LightGBM model. The method is based on Shapley values, a commonly
12 used metric to measure contributions from players in cooperative game theory. It treats each
13 explanatory feature as a player and calculates its importance by comparing the model performance
14 with and without the feature. Another useful visualization, the SHAP summary plot, sorts the
15 features by their global impact and calculates each feature’s impact on the model output of each
16 test sample, illustrating complex associations between the target and explanatory features.

17 **3 RESULTS**

18 3.1 Descriptive Analysis of Stationary Collocation and Mobile Validation

19 We first examined the raw data collected by CS units compared to the reference. As shown
20 in Table 1, the average PM_{2.5} concentrations during the sampling periods were 7.9 and 4.7 $\mu\text{g}/\text{m}^3$,
21 respectively, in New York City and Boston. All five CS units underestimated PM_{2.5} concentrations
22 in stationary collocation by 62.1% on average in New York City. During the mobile validation in
23 Boston, CS Unit 3 overestimated PM_{2.5} concentrations, while CS Unit 5 still underestimated PM_{2.5}.
24 The discrepancy is mainly caused by the different particle characteristics in New York and Boston
25 and inter-sensor sensing variability for these particles. The inter-sensor variability of CS units is
26 exhibited, which indicates the importance of developing unit-specific calibration models. The
27 specifically high inter-sensor variability in the Boston campaign is mainly attributable to the
28 continuous deployment in an unattended and open environment in NYC. The average reference
29 NO₂ concentrations were 10.2 and 14.3 $\mu\text{g}/\text{m}^3$, respectively, in New York City and Boston, which
30 is not included in Table 1 as CS provides only raw millivoltage readings.

1 Table 1. Raw PM_{2.5} levels in µg/m³ collected by City Scanner and reference monitors during
 2 stationary collocation in New York City and mobile comparison in Boston

	New York City (µg/m ³)			Boston (µg/m ³)		
	25 th percentile	50 th percentile	75 th percentile	25 th percentile	50 th percentile	75 th percentile
CS Unit 1	1.6	2.6	4.1	-	-	-
CS Unit 2	2.8	4.1	6.2	-	-	-
CS Unit 3	1.7	2.8	4.4	3.9	5.4	8.1
CS Unit 4	2.3	3.6	5.6	-	-	-
CS Unit 5	2.2	3.6	5.5	0.2	0.5	1.1
Reference	5.6	7.3	10.3	1.2	2.2	3.0

3 The weather conditions were different during the NYC stationary collocation and Boston
 4 mobile validation. Table 2 demonstrates the meteorological factors that were employed in the
 5 calibration algorithms. While relative humidity and air pressure were relatively consistent
 6 regionally across the two monitoring campaigns, air temperature and dew point were about 20 °C
 7 lower during the Boston mobile campaign than NYC.

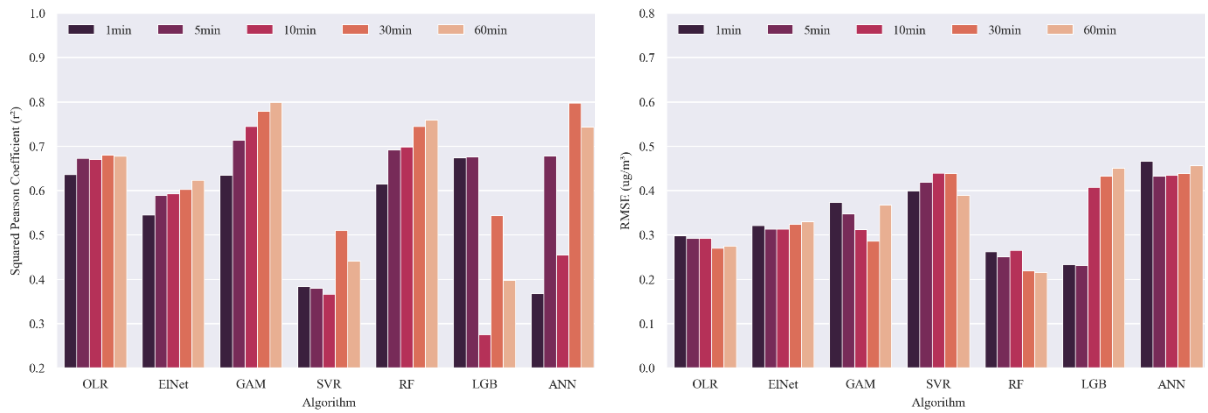
8 Table 2. Meteorological conditions during the stationary collocation and mobile validation

Duration and meteorology	New York average (Standard deviation)	Boston average (Standard deviation)
Duration	August 13 - September 10, 2021	February 11 - April 8, 2021
Temperature (Celsius)	25.6 (±3.5)	4.1 (±6.4)
Relative humidity (%)	64.5 (±14.6)	58.7 (±21.8)
Dew point (Celsius)	17.9 (±3.5)	-4.1 (±8.6)
Air pressure (kPa)	101.5 (±0.5)	101.6 (±1.0)

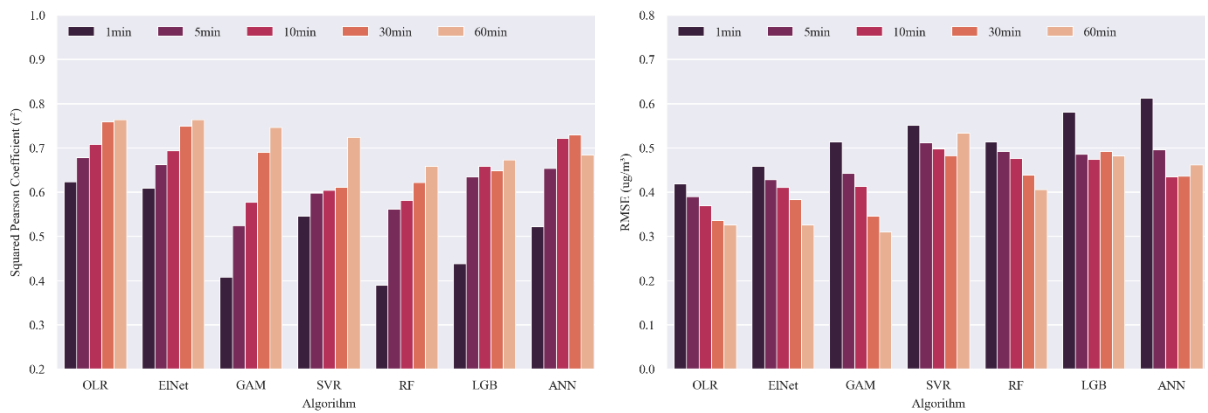
9 3.2 Algorithm Performances in a Stationary Setting

10 We summarized the calibration algorithm performance (r^2 and RMSE) of CS Unit 1 for
 11 PM_{2.5} and NO₂ in Figure 3. Performance summary plots for other CS units can be found in the
 12 Supporting Information. After calibration, all five CS units present good agreement with the
 13 reference data. In general, the algorithms provide good calibration results for both PM_{2.5} and NO₂,
 14 while PM calibration models have slightly better performance than NO₂ ones across all CS units.
 15 Best PM models' r^2 values reach over 0.8, compared to those of NO₂ models around 0.7. Still, it
 16 indicates the feasibility of calibrating CS NO₂ readings directly from millivoltage signals to

1 reference concentrations. Among all algorithms, ordinary least squares and elastic net have
 2 demonstrated consistently satisfactory performance in PM_{2.5} calibration, r^2 ranging from 0.6 to 0.8.
 3 This phenomenon is also observed in NO₂ calibration, indicating good linear responses between
 4 the low-cost sensors and reference sensors. RMSE values are very low after calibration, mostly
 5 below 0.5 $\mu\text{g}/\text{m}^3$ and 0.8 ppb for PM_{2.5} and NO₂, respectively, considering that the median
 6 concentrations for PM_{2.5} and NO₂ are 7.3 $\mu\text{g}/\text{m}^3$ and 9.0 ppb in the stationary collocation. Machine
 7 learning models can also deliver moderate calibration results but sometimes fail due to overfitting
 8 the training data. Across all units and algorithms, we do not observe an obvious effect of temporal
 9 aggregation interval on model performance.



(a) City Scanner Unit 1 PM_{2.5} performance

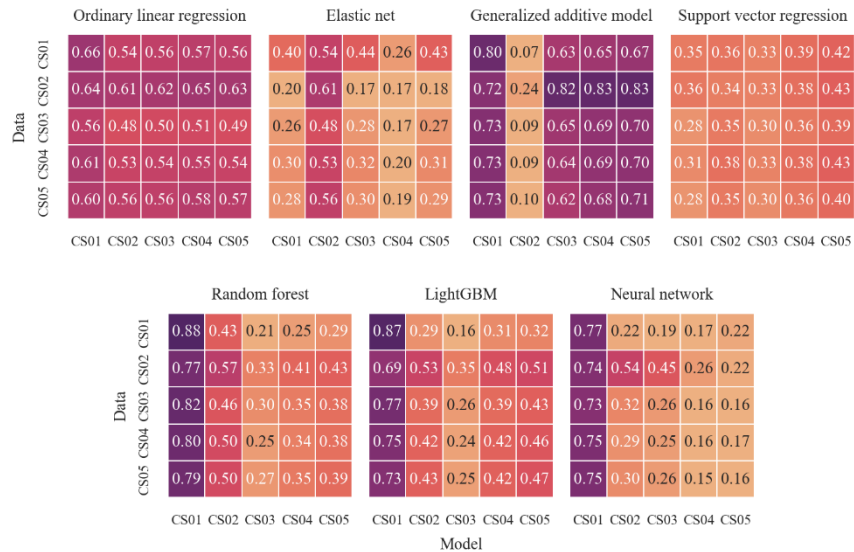


(b) City Scanner Unit 1 NO₂ performance

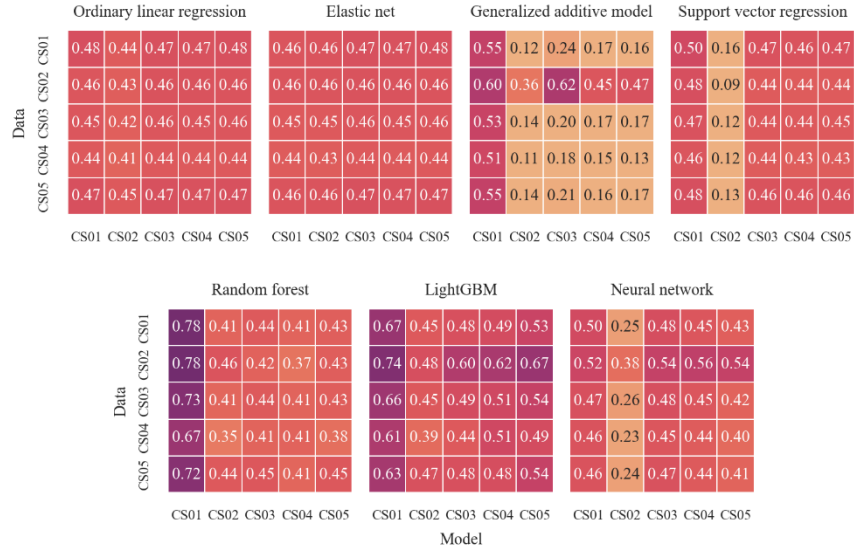
10 Figure 3. K-fold cross validation model performances (a) PM_{2.5} and (b) NO₂ in a stationary
 11 setting for City Scanner Unit 1.

12 Figure 4 tests the inter-sensor algorithm transferability for the best-performing PM_{2.5} and
 13 NO₂ models. Test data from a CS unit were passed through models developed using data from

1 another CS unit, which process was repeated for all five CS units and all calibration algorithms.
 2 1-min aggregated data and models were used as they provide the worst-case scenario in
 3 transferability and are most prone to overfitting data from a specific sensor. It is hypothesized that
 4 linear regressions should have the best transferability across sensors as they are the least likely to
 5 overfit. Nevertheless, for PM_{2.5} models, the generalized additive model and random forest are
 6 observed to have the best inter-sensor transferability, while the transferability of linear regressions
 7 is moderate. For NO₂, ordinary least squares and elastic net models perform, on average, as well
 8 as random forest and LightGBM models. But the latter two's transferability still outperforms
 9 generalized linear models developed for CS Unit 1. It is worth noting that models developed for
 10 CS Unit 2 have the worst transferability, and those developed for CS Unit 1 have the best
 11 transferability for both pollutants and all other algorithms.



(a) PM_{2.5} inter-sensor test



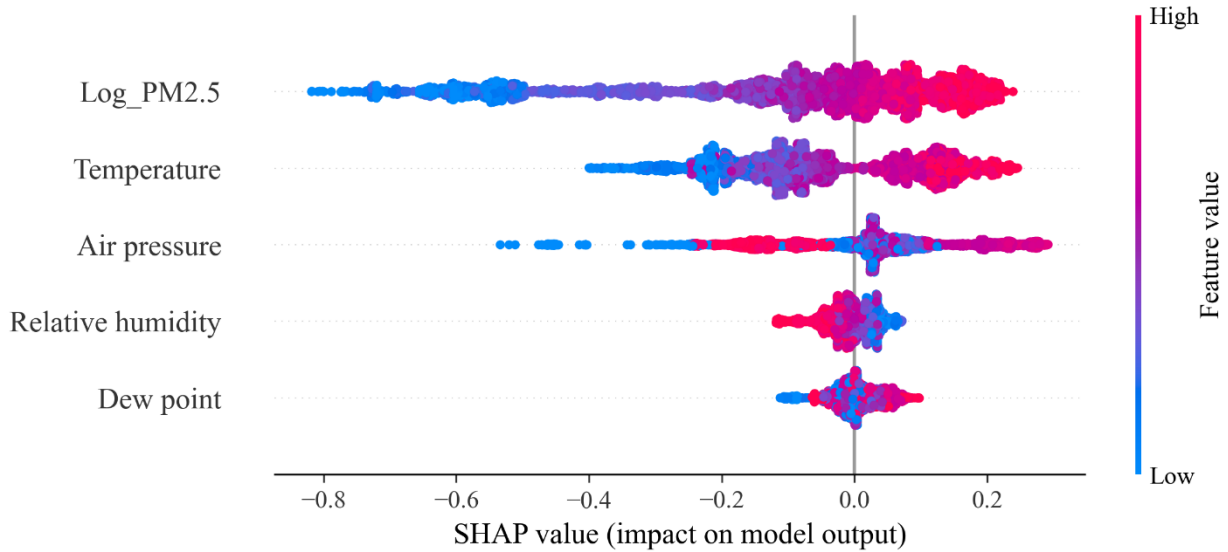
(b) NO_2 inter-sensor test

1 Figure 4. Inter-sensor variability tests of best-performing 1-min (a) $\text{PM}_{2.5}$ and (b) NO_2 models.
 2 Values in the tables are r^2 values calculated from passing sensor-specific data (y-axis) to sensor-
 3 specific models (x-axis); darker color represents higher values.

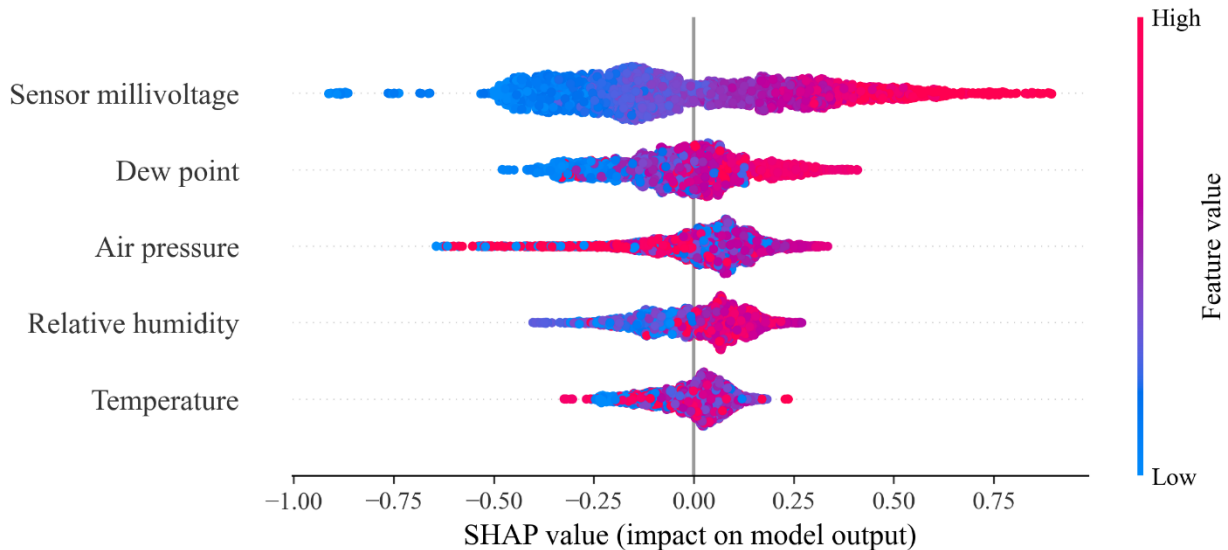
4 Finally, we examine the model behaviors in the stationary setting from the best-performing
 5 random forest models for $\text{PM}_{2.5}$ and NO_2 . The results are visualized using SHAP summary plots,
 6 which consist of marginal contributions of all data points of each explanatory feature, represented
 7 by their SHAP values. Positive values denote positive effects, and higher absolute values denote
 8 larger effects on the output and vice versa. Data points are color-coded from blue to red, indicating
 9 low to high values. Data points of each feature are stacked vertically if different values have the
 10 same impact as in a violin plot. Features are ranked by the absolute average SHAP values of all
 11 observations from high to low, indicating high to low feature importance. In Figure 5, log-
 12 transformed $\text{PM}_{2.5}$ readings from CS units have the highest feature importance. Moreover, its
 13 impact on the predictions is monotonic and positive as the SHAP value (x-axis) increases as log
 14 $\text{PM}_{2.5}$ values increase with a smooth color change from blue to red (low to high feature values). A
 15 similar phenomenon is found in the NO_2 calibration model, even though we mapped the electro-
 16 signal directly to reference readings. It is desirable that low-cost sensor readings can reflect air
 17 quality changes. Other than CS unit readings, the temperature is the second most important that
 18 positively correlates with reference $\text{PM}_{2.5}$ readings, while in the NO_2 model, the dew point feature
 19 plays the same role. The effect of relative humidity is not significant in the $\text{PM}_{2.5}$ and NO_2 models,
 20 where it is positively correlated with NO_2 readings but negatively correlated with $\text{PM}_{2.5}$. This

1 observation is opposite to previous literature (Crilley et al., 2018; di Antonio et al., 2018), which
2 can be attributed to a difference in particle hygroscopicity from city to city, especially in the
3 collocation region. We observe that the relationship between air pressure and model output is non-
4 linear. It is worth noting that air pressure can be merely a good predictor rather than a good
5 explainer for sensor readings. It is in line with our main purpose to develop high-accuracy
6 calibration models. Dew point in the PM_{2.5} model and temperature in the NO₂ model demonstrate
7 similar behaviors, proving the necessity of using non-linear models.

8



(a) PM_{2.5} best performing 1-min model



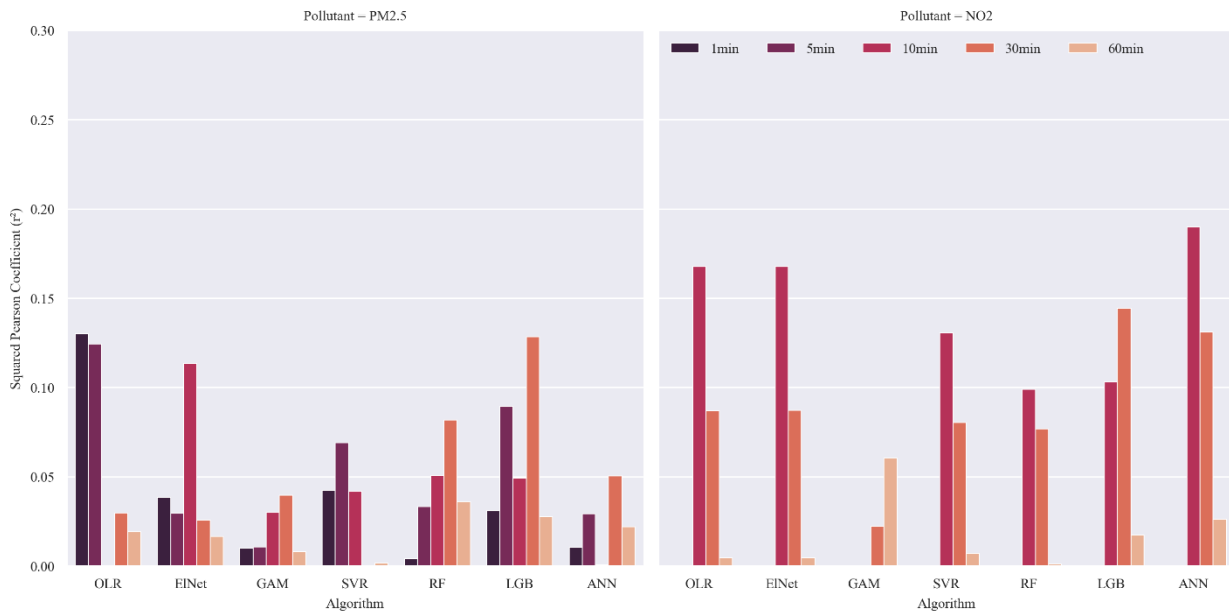
(b) NO₂ best performing 1-min model

1 Figure 5. SHAP summary plots of best performing 1-min random forest models on a City
 2 Scanner unit for (a) PM_{2.5} and (b) NO₂. Dots represent explanatory feature values and color-
 3 coded from low to high values in blue to red; data points with the same SHAP value are stacked
 4 vertically; all features are ranked from high to low feature importance on the y-axis.

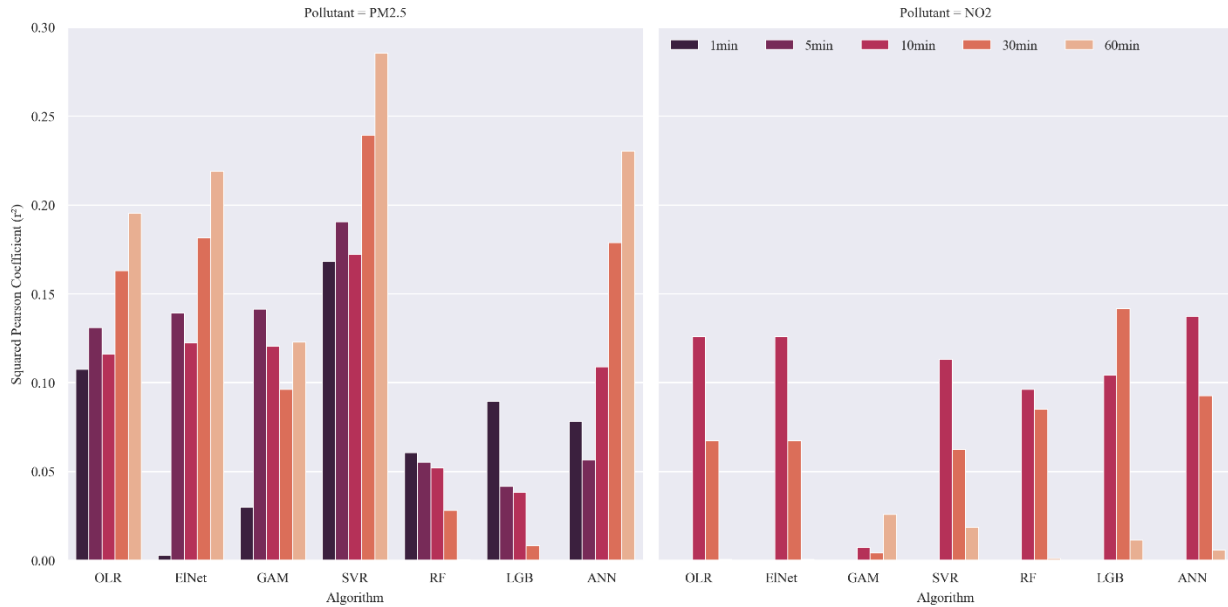
5 3.3 Algorithm Transferability in a Mobile Setting

6 We further transfer the best-performing models trained with stationary collocation data in
 7 NYC to predict mobile validation observations in Boston. We observe in Figure 6 significantly
 8 deteriorated performance across all algorithms for both CS units used in Boston (Units 3 and 5).
 9 The r^2 values are low in the range of 0.1 to 0.3. Multiple possible reasons lead to unsatisfactory

1 model transferability. On the one hand, mobile validation has brought extra challenges to sensor
 2 performance. In mobile validation, the sensors can be exposed to a variety of emission sources at
 3 different distances in a short amount of time. Such drastic changes in particle properties and
 4 pollutant concentrations are difficult to account for in stationary calibration. Also, the sensors
 5 being mobile validated were used for another mobile measurement campaign, which has led to
 6 accelerated drifting and aging due to vibration than promised in the manual. On the other hand,
 7 some problems exist ubiquitously for low-cost sensor calibration models in both stationary and
 8 mobile settings. First, the weather in NYC during stationary collocation and in Boston during the
 9 mobile validation are drastically different, especially the temperature and dew points, which play
 10 important roles in the PM_{2.5} and NO₂ models, respectively. Some weather feature values are unseen
 11 in the training data, thus, hard to extrapolate when transferred to mobile data prediction. Second,
 12 while we use reference stationary data to train our model, the target data to predict during the
 13 Boston mobile validation are from research-grade instruments rather than reference-grade ones.
 14 Third, it is common that the performance of low-cost sensors degrades, and their behaviors change
 15 after long-time outdoor usage, given that the stationary collocation and mobile validation were
 16 conducted more than six months apart.



(a) City Scanner Unit 3



(b) City Scanner Unit 5

1 Figure 6. Model performances in a mobile setting for two City Scanner units and five data
 2 aggregations. Box colors illustrate 1 min to 60 min temporal aggregations from dark to light.

3

4 4 DISCUSSION

5 This study systematically evaluated the effects of temporal aggregation, calibration
 6 algorithms, and meteorological conditions on low-cost sensor calibration in both stationary and
 7 mobile settings.

8 We observed that linear regression models, including ordinary least squares and elastic net,
 9 perform consistently well in both PM_{2.5} and NO₂ calibrations. Meanwhile, algorithms that consider
 10 non-linear relationships, including random forest and GAM, provide good prediction accuracy and
 11 better transferability between CS units than those that do not when calibrating PM_{2.5}. This
 12 phenomenon is not found in NO₂ calibration models, where linear regression transfers as well as
 13 machine learning models. It indicates that Alphasense NO₂-A43F sensor provides consistent linear
 14 responses in most circumstances with little inter-sensor variability. It is recommended to employ
 15 algorithms that can account for non-linear interactions between meteorology and low-cost sensor
 16 behaviors, especially for PM calibration. For most educational and citizen science projects,
 17 generalized linear regressions are sufficient for low-cost sensor calibration as they are easy to
 18 apply and interpret. Robust results can be reached if we properly follow the procedures

1 recommended in the EPA guidelines for field tests. We also observed that random forest and GAM
2 models are highly transferable between PM_{2.5} sensors, second by linear regression models. It
3 supports adopting low-cost air sensors at a large scale to obtain big air quality data in urban areas.

4 By interpreting the best PM_{2.5} and NO₂ model behavior using SHAP plots, we confirmed
5 the non-linear relationships between meteorological features and model outputs, especially for PM
6 calibration. Moreover, we developed all models using only the CS readings and four
7 meteorological features without further feature engineering. It indicates that capturing non-
8 linearity and feature interactions might be more effective in improving calibration models than
9 incorporating and compositing more features.

10 We observed that calibration model performance deteriorates when transferred from
11 stationary to mobile use. This is mainly caused by several reasons, some of which are ubiquitous
12 to all low-cost sensors, while others are unique to mobile measurement. The ubiquitous ones
13 include weather differences between stationary collocation and mobile validation periods, using
14 research-grade instead of reference-grade monitors in the mobile validation, and performance
15 degradation of the low-cost sensors, including drifting and aging. Drifting and aging are common
16 among all types of low-cost sensors (Kim et al., 2018; Malings et al., 2019) and can happen after
17 2 to 6 months of usage (Van Zoest et al., 2019). Moreover, mobile deployment between the
18 stationary collocation and mobile validation reported in this study has led to accelerated aging and
19 drifting. It is suggested to calibrate low-cost sensors before and validate after each deployment in
20 an environment as similar to the actual monitoring as possible if it lasts more than six months and
21 is unsupervised, as in our case. Regular collocations and calibrations are recommended for low-
22 cost sensors performed no more than six months apart. Our work also demonstrates the risk of
23 deploying low-cost sensors in an environment different from the calibration environment. The lack
24 of local reference instruments to perform calibration can constrain low-cost sensor applications in
25 low-income regions and countries; alternatively, the viability of calibration (spatial and temporal
26 differences in calibration and deployment locations) should be acknowledged in data analysis.

27 We acknowledge several limitations in our study. First, the PM_{2.5} and NO₂ levels were low
28 during the stationary collocation and mobile validation periods ($< 20 \mu\text{g}/\text{m}^3$ and $< 50 \text{ ppb}$,
29 respectively). Ideally, an as wide as possible air quality range should be included in the calibration
30 process to account for the variability of air pollution in mobile on-road deployments. Kelly et al.

1 (2017) reported that some low-cost sensors begin to exhibit a non-linear response only when in
2 high PM concentrations ($>40 \mu\text{g}/\text{m}^3$). Second, not all CS units are tested in the mobile comparison
3 due to limited space and resources. Lastly, while it is more rigorous to calibrate low-cost sensors
4 against Federal Reference Methods instruments, our study used data from Federal Equivalent
5 Methods instruments as they are available minute-by-minute, which can better serve the purpose
6 of our study.

7 **5 CONCLUDING REMARKS**

8 Our work presents two major takeaways regarding effective and robust low-cost calibration
9 with a special focus on mobile air quality monitoring. We first demonstrated that PM calibration
10 should consider the complex relationship between sensor responses, pollutant concentrations, and
11 meteorological factors. Therefore, algorithms that consider non-linear relationships should be
12 adopted in this case.

13 Models trained in the stationary setting can hardly be transferred to the mobile setting in a
14 different urban environment and climate; additionally, sensors degrade after long-term outdoor
15 mobile deployment. It is necessary to calibrate low-cost sensors in an environment that is similar
16 to real-world deployment. It is recommended to calibrate or assess comparability both before and
17 after mobile deployments of durations exceeding six months. Our recommendation exceeds
18 existing EPA guidelines that are primarily focused on stationary monitoring.

19 The findings in our study are important to citizen scientists, air quality researchers, and
20 practitioners interested in advancing low-cost sensor applications, such as raising awareness and
21 community engagement in scientific analysis. Further research is needed on other commonly used
22 low-cost particulate matter and gas sensors and low-cost sensor performance in temperature,
23 humidity, and concentration range not included in this study. Finally, there is a need to assess low-
24 cost sensor performance in mobile settings given the rapidly increasing use of and demands for
25 mobile environmental sensing.

26

1 **REFERENCE**

- 2 Bezantakos, S., Schmidt-Ott, F., & Biskos, G. (2018). Performance evaluation of the cost-effective
3 and lightweight Alphasense optical particle counter for use onboard unmanned aerial vehicles.
4 *Aerosol Science and Technology*, 52(4), 385–392.
5 [https://doi.org/10.1080/02786826.2017.1412394/SUPPL_FILE/UAST_A_1412394_SM0236.ZI](https://doi.org/10.1080/02786826.2017.1412394/SUPPL_FILE/UAST_A_1412394_SM0236.ZIP)
6 P
- 7 Brauer, M., Guttikunda, S. K., K A, N., Dey, S., Tripathi, S. N., Weagle, C., & Martin, R. v. (2019).
8 Examination of monitoring approaches for ambient air pollution: A case study for India.
9 *Atmospheric Environment*, 216, 116940. <https://doi.org/10.1016/J.ATMOSENV.2019.116940>
- 10 Castell, N., Dauge, F. R., Schneider, P., Vogt, M., Lerner, U., Fishbain, B., Broday, D., &
11 Bartonova, A. (2017). Can commercial low-cost sensor platforms contribute to air quality
12 monitoring and exposure estimates? *Environment International*, 99, 293–302.
- 13 Crawford, B., Hagan, D. H., Grossman, I., Cole, E., Holland, L., Heald, C. L., & Kroll, J. H. (2021).
14 Mapping pollution exposure and chemistry during an extreme air quality event (the 2018 Kīlauea
15 eruption) using a low-cost sensor network. *Proceedings of the National Academy of Sciences*,
16 118(27). <https://doi.org/10.1073/PNAS.2025540118>
- 17 Crilley, L. R., Shaw, M., Pound, R., Kramer, L. J., Price, R., Young, S., Lewis, A. C., & Pope, F.
18 D. (2018). Evaluation of a low-cost optical particle counter (Alphasense OPC-N2) for ambient air
19 monitoring. *Atmospheric Measurement Techniques*, 11(2), 709–720.
20 <https://doi.org/10.5194/AMT-11-709-2018>
- 21 Crilley, L. R., Singh, A., Kramer, L. J., Shaw, M. D., Alam, M. S., Apte, J. S., Bloss, W. J.,
22 Hildebrandt Ruiz, L., Fu, P., Fu, W., Gani, S., Gatari, M., Ilyinskaya, E., Lewis, A. C., Ng’ang’a,
23 D., Sun, Y., Whitty, R. C. W., Yue, S., Young, S., & Pope, F. D. (2020). Effect of aerosol
24 composition on the performance of low-cost optical particle counter correction factors.
25 *Atmospheric Measurement Techniques*, 13(3), 1181–1193. [https://doi.org/10.5194/AMT-13-](https://doi.org/10.5194/AMT-13-1181-2020)
26 1181-2020
- 27 deSouza, P., Kahn, R. A., Limbacher, J. A., Marais, E. A., Duarte, F., & Ratti, C. (2020).
28 Combining low-cost, surface-based aerosol monitors with size-resolved satellite data for air

1 quality applications. *Atmospheric Measurement Techniques*, 13(10), 5319–5334.
2 <https://doi.org/10.5194/AMT-13-5319-2020>

3 Duarte, F., & deSouza, P. (2020). Data Science and Cities: A Critical Approach. *Harvard Data*
4 *Science Review*, 2(3). <https://doi.org/10.1162/99608F92.B3FC5CC8>

5 Gressent, A., Malherbe, L., Colette, A., Rollin, H., & Scimia, R. (2020). Data fusion for air quality
6 mapping using low-cost sensor observations: Feasibility and added-value. *Environment*
7 *International*, 143, 105965. <https://doi.org/10.1016/J.ENVINT.2020.105965>

8 Hua, J., Zhang, Y., Foy, B. de, Mei, X., Shang, J., Zhang, Y., Sulaymon, I. D., & Zhou, D. (2021).
9 Improved PM2.5 concentration estimates from low-cost sensors using calibration models
10 categorized by relative humidity. <https://doi.org/10.1080/02786826.2021.1873911>, 55(5), 600–
11 613. <https://doi.org/10.1080/02786826.2021.1873911>

12 Hudda, N., Simon, M. C., Patton, A. P., & Durant, J. L. (2020). Reductions in traffic-related black
13 carbon and ultrafine particle number concentrations in an urban neighborhood during the COVID-
14 19 pandemic. *Science of The Total Environment*, 742, 140931.
15 <https://doi.org/10.1016/J.SCITOTENV.2020.140931>

16 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., & Liu, T.-Y. (2017). LightGBM:
17 A Highly Efficient Gradient Boosting Decision Tree. *31st Conference on Neural Information*
18 *Processing Systems (NIPS 2017)*.
19 <https://proceedings.neurips.cc/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf>

20 Kelly, K. E., Whitaker, J., Petty, A., Widmer, C., Dybwad, A., Sleeth, D., Martin, R., & Butterfield,
21 A. (2017). Ambient and laboratory evaluation of a low-cost particulate matter sensor.
22 *Environmental Pollution*, 221, 491–500. <https://doi.org/10.1016/J.ENVPOL.2016.12.039>

23 Kerckhoffs, J., Hoek, G., Portengen, L., Brunekreef, B., & Vermeulen, R. C. H. (2019).
24 Performance of Prediction Algorithms for Modeling Outdoor Air Pollution Spatial Surfaces.
25 *Environmental Science & Technology*, 53(3), 1413–1421. <https://doi.org/10.1021/acs.est.8b06038>

26 Liang, L. (2021). Calibrating low-cost sensors for ambient air monitoring: Techniques, trends, and
27 challenges. *Environmental Research*, 197, 111163.
28 <https://doi.org/10.1016/J.ENVRES.2021.111163>

1 Lundberg, S. M., Erion, G. G., & Lee, S.-I. (2018). *Consistent Individualized Feature Attribution*
2 *for Tree Ensembles*.

3 Lundberg, S. M., & Lee, S.-I. (2017). *A Unified Approach to Interpreting Model Predictions*.
4 *Section 2*, 1–10.

5 Malings, C., Tanzer, R., Hauryliuk, A., Kumar, S. P. N., Zimmerman, N., Kara, L. B., Presto, A.
6 A., & Subramanian, R. (2019). Development of a general calibration model and long-term
7 performance evaluation of low-cost sensors for air pollutant gas monitoring. *Atmospheric*
8 *Measurement Techniques*, *12*(2), 903–920. <https://doi.org/10.5194/AMT-12-903-2019>

9 Miskell, G., Salmond, J. A., & Williams, D. E. (2018). Use of a handheld low-cost sensor to
10 explore the effect of urban design features on local-scale spatial and temporal air quality variability.
11 *Science of The Total Environment*, *619–620*, 480–490.
12 <https://doi.org/10.1016/J.SCITOTENV.2017.11.024>

13 Padró-Martínez, L. T., Patton, A. P., Trull, J. B., Zamore, W., Brugge, D., & Durant, J. L. (2012).
14 Mobile monitoring of particle number concentration and other traffic-related air pollutants in a
15 near-highway neighborhood over the course of a year. *Atmospheric Environment*, *61*, 253–264.
16 <https://doi.org/10.1016/J.ATMOSENV.2012.06.088>

17 Pedregosa, F., Michel, V., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Vanderplas, J.,
18 Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, É. (2011). Scikit-learn: Machine Learning
19 in Python. *Journal of Machine Learning Research*, *12*(85), 2825–2830.
20 <http://jmlr.org/papers/v12/pedregosa11a.html>

21 SM, S. N., Reddy Yasa, P., MV, N., Khadirnaikar, S., & Pooja Rani. (2019). Mobile monitoring
22 of air pollution using low cost sensors to visualize spatio-temporal variation of pollutants at urban
23 hotspots. *Sustainable Cities and Society*, *44*, 520–535.

24 Sousan, S., Koehler, K., Hallett, L., & Peters, T. M. (2016). Evaluation of the Alphasense optical
25 particle counter (OPC-N2) and the Grimm portable aerosol spectrometer (PAS-1.108). *Aerosol*
26 *Science and Technology*, *50*(12), 1352–1365.
27 https://doi.org/10.1080/02786826.2016.1232859/SUPPL_FILE/UAST_A_1232859_SM0343.ZI
28 P

1 US EPA. (2014). *Air Sensor Guidebook*.
2 https://cfpub.epa.gov/si/si_public_record_Report.cfm?Lab=NERL&dirEntryId=277996

3 US EPA. (2021a). *Performance Testing Protocols, Metrics, and Target Values for Fine*
4 *Particulate Matter Air Sensors: Use in Ambient, Outdoor, Fixed Sites, Non-Regulatory*
5 *Supplemental and Informational Monitoring Applications*. <https://doi.org/EPA/600/R-20/280>

6 US EPA. (2021b). *Performance Testing Protocols, Metrics, and Target Values for Ozone Air*
7 *Sensors: Use in Ambient, Outdoor, Fixed Site, Non-Regulatory and Informational Monitoring*
8 *Applications*. <https://doi.org/EPA/600/R-20/279>

9 Wang, S., Ma, Y., Wang, Z., Wang, L., Chi, X., Ding, A., Yao, M., Li, Y., Li, Q., Wu, M., Zhang,
10 L., Xiao, Y., & Zhang, Y. (2021). Mobile monitoring of urban air quality at high spatial resolution
11 by low-cost sensors: Impacts of COVID-19 pandemic lockdown. *Atmospheric Chemistry and*
12 *Physics*, 21(9), 7199–7215. <https://doi.org/10.5194/acp-2020-1169>

13 Zheng, T., Bergin, M. H., Johnson, K. K., Tripathi, S. N., Shirodkar, S., Landis, M. S., Sutaria, R.,
14 & Carlson, D. E. (2018). Field evaluation of low-cost particulate matter sensors in high-and low-
15 concentration environments. *Atmospheric Measurement Techniques*, 11(8), 4823–4846.
16 <https://doi.org/10.5194/AMT-11-4823-2018>

17