# 1 Calibrating Networks of Low-
# 2 Cost Air Quality Sensors

3 Priyanka deSouza[1*], Ralph Kahn[2], Tehya Stockman[3,4], William Obermann[3], Ben Crawford[5], An
4 Wang[6], James Crooks[7,8], Jing Li[9], Patrick Kinney[10]

5
6 1: Department of Urban and Regional Planning, University of Colorado Denver, 80202
7 2: NASA Goddard Space Flight Center, Greenbelt MD
8 3: Denver Department of Public Health and Environment, USA
9 4: Department of Civil, Environmental, and Architectural Engineering, University of Colorado
10 Boulder, Boulder, Colorado 80309, United States
11 5: Department of Geography and Environmental Sciences, University of Colorado Denver, 80202
12 6: Senseable City Lab, Massachusetts Institute of Technology, Cambridge 02139
13 7: Division of Biostatistics and Bioinformatics, National Jewish Health, 2930
14 8: Department of Epidemiology, University of Colorado at Denver - Anschutz Medical Campus,
15 129263
16 9: Department of Geography and the Environment, University of Denver, Denver, CO, USA
17 10: Boston University School of Public Health, Boston, MA, USA
18
19 *: priyanka.desouza@ucdenver.edu

## 20 Abstract

21 Ambient fine particulate matter (PM$_{2.5}$) pollution is a major health risk. Networks of low-cost
22 sensors (LCS) are increasingly being used to understand local air pollution variation. However,
23 measurements from LCS have uncertainties which can act as a potential barrier for effective
24 decision-making. LCS data thus need to be calibrated to obtain better quality PM$_{2.5}$ estimates. In
25 order to develop correction factors, LCS are typically co-located with gold-standard reference
26 monitors. A calibration equation is then developed that relates the raw output of the LCS as closely
27 as possible to measurements from the reference monitor. This calibration algorithm is then
28 typically *transferred* to measurements from monitors in the network. Calibration algorithms tend to
29 be evaluated based on their performance at co-location sites. It is often implicitly assumed that the
30 conditions at the relatively sparse co-location sites are representative of the LCS network, overall.
31 Little work has been done to explicitly evaluate the sensitivity of the LCS network hotspot
32 detection, and spatial and temporal PM$_{2.5}$ trends to the correction method applied. This paper
33 provides a first look at how transferable different calibration methods are using a dense network of
34 Love My Air LCS monitors in Denver. It offers a series of transferability metrics that can be
35 applied to other networks and offers suggestions for which calibration method would be most
36 useful for different end goals. Finally, it develops a set of best practice suggestions on calibrating
37 LCS networks.
38

Atmospheric
Measurement
Techniques
Open Access

Discussions

EGU

39  **Key words**: low-cost sensors, PM$_{2.5}$, calibration, Love My Air

# 1 Introduction

41  Poor air quality is currently the single largest environmental risk factor to human health in the
42  world, with ambient air pollution responsible for 6.7 million premature deaths every year (State of
43  Global Air, 2020). Accurate air quality data is crucial for tracking long-term trends in air quality
44  levels, and for the development of effective pollution management plans. Levels of fine particulate
45  matter (PM$_{2.5}$), a criterion pollutant that poses more of danger to human health than other
46  widespread pollutants (Kim et al., 2015), can vary over distances as small as ~ 10's of meters in
47  complex urban environments (Brantley et al., 2019; deSouza et al., 2020a). Therefore, dense
48  monitoring networks are often needed to capture relevant spatial variations. Due to their costliness,
49  EPA air quality reference monitoring networks, the gold standard for measuring air pollutants, are
50  sparsely positioned across the US (Apte et al., 2017; Anderson and Peng, 2012).
51
52  Low-cost sensors (LCS) (<USD \$2500 as defined by the US EPA Air Sensor Toolbox) (Williams
53  et al., 2014) have the potential to capture concentrations of PM in previously unmonitored
54  locations and democratize air pollution information (Castell et al., 2017; Kumar et al., 2015;
55  Morawska et al., 2018; Snyder et al., 2013; deSouza and Kinney, 2021; deSouza, 2022). However,
56  LCS measurements have several sources of uncertainty (Bi et al., 2020; Giordano et al., 2021;
57  Liang, 2021).
58
59  Most low-cost PM sensors rely on optical measurement techniques. Optical instruments face
60  several inherent challenges that introduce potential differences in mass estimations compared to
61  reference methods (Barkjohn et al., 2021; Crilley et al., 2018; Giordano et al., 2021; Malings et al.,
62  2020):
63
64  1. Optical methods do not directly measure mass concentrations; rather, they estimate mass based
65  on calibrations that convert light scattering data to particle number and mass. LCS come with
66  factory-supplied calibrations, but in practice must be re-calibrated in the field to ensure accuracy,
67  due to variations in ambient particle characteristics.
68
69  2. High relative humidity (RH) can produce hygroscopic particle growth, leading to mass
70  overestimation if the particles are not dessicated by the instrument.
71
72  3. The inability to detect particles with diameters below a specific size, which is determined by
73  the wavelength of laser light within each device, and is generally in the vicinity of 0.3 μm, whereas
74  the peak in pollution particle size distributions is typically smaller than 0.3 μm.
75
76  4. The physical and chemical parameters of the aerosol (particle size distribution, shape, indices
77  of refraction, hygroscopicity, volatility etc.) which might vary significantly across different
78  microenvironments with diverse sources impact light scattering, which in turn affects the aerosol
79  mass concentrations reported by these instruments.
80

81    The need for field calibration to correct LCS measurements is particularly important. This is
82    typically done by co-locating a small number of LCS with a reference monitor at a representative
83    monitoring location or locations. The co-location could be carried out for a brief period before
84    and/or after the actual study or may continue at a small number of sites for the duration of the
85    study. In either case, the co-location provides data from which a calibration equation is then
86    developed that relates the raw output of the LCS as closely as possible to the desired quantity as
87    measured by the reference monitor. Thereafter, the calibration equation is transferred to other LCS
88    in the network, based upon the presumption that ongoing sampling conditions are within the same
89    range as those during the calibration period.

91    Calibration models typically correct for 1) systematic error in LCS by adjusting for bias using
92    reference monitor measurements, and 2) the dependence of LCS measurements on environmental
93    conditions affecting the ambient particle properties such as relative humidity (RH), temperature
94    (T), and/or dew-point (D). Correcting for RH, T and D is carried out through either a) a physics-
95    based approach that accounts for aerosol hygroscopic growth given particle composition using $\kappa$-
96    kohler's theory, or b) empirical models, such as regression and machine learning techniques. In
97    this paper, we will focus on the latter, as it is the most widely used (Barkjohn et al., 2021).
98    Previous work has also shown that the two approaches yield comparable improvements in the case
99    of $PM_{2.5}$ LCS (Malings et al., 2020).

101    Prior studies have used multivariate regressions, piecewise linear regressions, or higher-order
102    polynomial models to account for RH, T and D in these calibration equations (Holstius et al., 2014;
103    Magi et al., 2020; Zusman et al., 2020). More recently, machine learning techniques such as
104    random forests, neural networks, and gradient boosted decision trees have been used (Considine et
105    al., 2021; Liang, 2021; Zimmerman et al., 2018). Researchers have also started including
106    additional covariates in their models besides what is directly measured by the LCS, such as time of
107    day, seasonality and site-type, which have been shown to yield significantly improved results
108    (Considine et al., 2021).

110    Past research has shown that there are several important decisions, in addition to the choice of
111    statistical model, that need to be made during calibration and can impact the results (Bean, 2021;
112    Giordano et al., 2021; Hagler et al., 2018). These include a) the kind of reference air quality
113    monitor used, b) the time-interval (e.g., hour/day) over which to average measurements used when
114    developing the calibration algorithm, c) how cross-validation (e.g., leave one site out/10-fold cross
115    validation) is carried out, and d) how long the co-location experiment takes place.

117    Calibration algorithms are evaluated based on how well the corrected measurements agree with
118    those from the reference monitor. A commonly used metric is the coefficient of determination, $R^2$,
119    which quantifies the strength of the association. However, it might be a mis-leading indicator of
120    sensor performance when measurements are observed close to the level of detection of the
121    instrument. Therefore, Root Mean Square Error (RMSE) is also often used in practice. Neither of
122    these metrics captures how well the calibration method developed at the co-located sites *transfers*
123    to the rest of the network.

Atmospheric
Measurement
Techniques
Discussions

124

125 If the conditions at the calibration site (meteorological conditions, pollution source mix) are the
126 same as at the rest of the network, the calibration function developed at the co-location site can be
127 assumed to be transferable to the rest of the network. In order to ensure that the sampling
128 conditions of the co-location site are representative of sampling conditions of the network, most
129 researchers tend to deploy monitors in the same general sampling area as the network (Zusman et
130 al., 2020). However, it is difficult to definitively test if the co-location site is representative of the
131 locations of all monitors in the network; ambient PM concentrations can vary on scales as small as
132 a few meters. Furthermore, LCS are often deployed specifically in areas where the air pollution
133 conditions are poorly understood, meaning that representativeness cannot be assessed ahead-of-
134 time.

135

136 Where multiple co-location sites exist, one way to address this challenge is to leave out one or
137 another co-location site to test if the calibration algorithm is transferable to the left-out site. This
138 method was used in recent work evaluating the feasibility of developing a US-wide correction to
139 the PurpleAir low-cost sensor network (Barkjohn et al., 2021; Nilson et al., 2022). Although this
140 approach helps, co-location sites are sparse relative to other sites in the network. Even in the
141 PurpleAir network (which is one of the densest low-cost networks in the world) there were only 39
142 co-location sites in 16 US states, a small fraction of the several thousand PurpleAir sites overall
143 (Barkjohn et al., 2021). It is thus important to test how sensitive the spatial and temporal trends of
144 pollution derived from the network are to the calibration algorithm used.

145

146 Examining the reliability of calibration methods is timely because, as mentioned earlier, more
147 researchers are opting to use machine learning calibration models. Although in most cases, such
148 models have yielded better results than traditional linear regressions, it is important to examine if
149 these models are overfitted to conditions at the co-location sites, and how transferable they are to
150 the rest of the network. Indeed, because of concerns of overfitting, some researchers have
151 explicitly eschewed employing machine learning calibration models altogether (Nilson et al.,
152 2022). It is important to test if these concerns are warranted.

153

154 This paper uses a dense low-cost air quality monitoring network deployed in Denver, termed
155 "Love My Air" network, to quantify the uncertainty in the spatial and temporal trends of the
156 network to the calibration algorithm used, as well as to ask the question: How much do we have to
157 worry about the transferability of different calibration functions across a $PM_{2.5}$ network in a
158 relatively small area in a single city? The methodology proposed in this paper to evaluate the
159 transferability of calibration adjustments can be applied to other low-cost sensor networks, with the
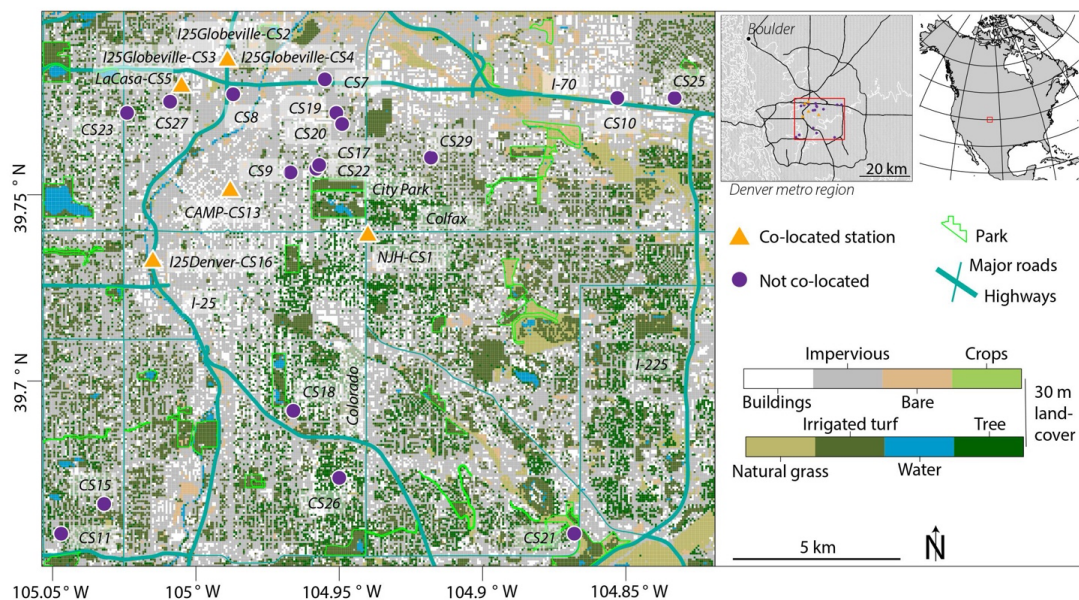160 understanding that the actual results will vary with study region.

## 2 Data and Methods

### 2.1 Data Sources

163 Between January 1 and September 30, 2021, Denver's Love My Air sensor network collected data
164 from 24 low-cost sensors deployed across the city outside of public schools and at reference

165     monitor locations (**Figure 1, Table 1**). The Love My Air sensors are Canary-S models equipped
166     with a Plantower 5003, made by Lunar Outpost Inc. The Canary-S sensors detect $PM_{2.5}$, T, and
167     RH, and upload minute-resolution measurements to an online platform via cellular data network.



168

169     *Figure 1: Locations of all 24 Love My Air Sensors. Sensors displayed with an orange triangle*
170     *indicate that they were co-located with a reference monitor. The labels of the co-located sensors include the*
171     *name of the corresponding reference monitor. The base map of land cover was obtained from*
172     *https://drcog.org/services-and-resources/data-maps-and-modeling/regional-land-use-land-cover-*
173     *project, last accessed April 2021.*

174

175     After removing missing values in the $PM_{2.5}$, T and RH data, RH < 0 (unrealistic values), T ≤ -30⁰C
176     (unrealistically low), and $PM_{2.5}$ values above 1,500 μg/m³ (outside the operational range of the
177     Plantower sensors used) from the Canary-S sensors (Considine et al., 2021), we were left with
178     8,809,340 measurements. We calculated hourly averages and obtained a total 147,101
179     measurements. From inspection, one of the monitors, CS13, worked intermittently in January and
180     February, before resuming continuous measurement in March (**Figure S1** *in Supplementary*
181     *Information*). When CS13 worked intermittently, large spikes in the measurements were observed,
182     likely due to power surges. We thus only retained measurements taken after March 1, 2021, for this
183     monitor. The total number of hourly measurements was thus reduced to 146,583.

184

185     Love My Air sensors were co-located with FEM reference monitors at La Casa (Sensor ID: CS5),
186     CAMP (Sensor ID: CS13), I25 Globeville (Sensor ID: CS2, CS3, CS4), I25 Denver (Sensor ID:
187     CS16), and NJH (Sensor ID: CS1) for the entire period of the experiment. Three Love My Air
188     sensors were co-located with the I25 Globeville Monitor, whereas there were single Love My Air
189     sensors at the other co-location sites. We obtained high-quality hourly $PM_{2.5}$ measurements from
190     the five reference monitors for the duration of the experiment. We joined hourly averages from
191     each of the co-located Love My Air monitors with the corresponding FEM monitor. We had a total
192     of 35,593 co-located measurements for which we had data for both the Love My Air sensor and the

193    corresponding reference monitor. **Figure S2** displays time-series plots of PM$_{2.5}$ from all co-located
194    Love My Air sensors. **Figure S3** displays time-series plots of PM$_{2.5}$ from the corresponding
195    reference monitors.
196
197    *Table 1*: Site location of each Love My Air sensor, as well as summary statistics of hourly
198    *measurements from each sensor*

| | | | | | PM$_{2.5}$ (µg/m³) | | | Temperature (⁰C) | RH (%) | Dewpoint (⁰C) |
|---|---|---|---|---|---|---|---|---|---|---|
| Sensor ID | Co-location Information | Latitude | Longitude | Hours operational | Mean | Median | Min-Max | Mean | Mean | Mean |
| CS1 | Co-located at NJH | 39.739 | -104.940 | 5,478 | 13 | 8 | 0 - 121 | 14.9 | 57.4 | 4.4 |
| CS2 | Co-located at I25 Globeville | 39.786 | -104.989 | 5,818 | 14 | 9 | 0 - 142 | 16.4 | 63.6 | 7.6 |
| CS3 | Co-located at I25 Globeville | 39.786 | -104.989 | 2,490 | 18 | 13 | 0 - 159 | 9.3 | 62.5 | 0.1 |
| CS4 | Co-located at I25 Globeville | 39.786 | -104.989 | 5,765 | 12 | 8 | 0 - 137 | 15.8 | 67.6 | 8.0 |
| CS5 | Co-located at La Casa | 39.779 | -105.005 | 5,761 | 12 | 8 | 0 - 129 | 13.4 | 69.6 | 6.0 |
| CS7 | - | 39.781 | -104.955 | 6,540 | 13 | 8 | 0 - 136 | 16.5 | 55.6 | 5.0 |
| CS8 | - | 39.777 | -104.987 | 6,282 | 13 | 8 | 0 - 133 | 17.3 | 38.3 | 0.0 |
| CS9 | - | 39.756 | -104.967 | 6,552 | 12 | 8 | 0 - 115 | 15.3 | 62.8 | 6.1 |
| CS10 | - | 39.776 | -104.853 | 6,552 | 12 | 7 | 0 - 142 | 17.9 | 32.6 | -2.4 |
| CS11 | - | 39.659 | -105.047 | 6,548 | 12 | 7 | 0 - 127 | 15.0 | 58.2 | 4.5 |
| CS13 | Co-located at CAMP | 39.751 | -104.988 | 4,449 | 13 | 8 | 0 - 115 | 21.9 | 54.7 | 10.2 |
| CS15 | - | 39.667 | -105.032 | 6,552 | 10 | 6 | 0 - 106 | 17.0 | 34.6 | -1.5 |
| CS16 | Co-located at I25 Denver | 39.732 | -105.015 | 5,832 | 12 | 9 | 0 - 100 | 17.4 | 33.6 | -2.2 |
| CS17 | - | 39.757 | -104.958 | 6,527 | 12 | 7 | 0 - 149 | 17.1 | 35.1 | -1.3 |
| CS18 | - | 39.692 | -104.966 | 6,552 | 12 | 7 | 0 - 115 | 16.9 | 36.3 | -1.0 |
| CS19 | - | 39.772 | -104.951 | 1,749 | 11 | 5 | 0 - 66 | 3.4 | 40.0 | -11.1 |
| CS20 | - | 39.769 | -104.949 | 6,551 | 10 | 6 | 0 - 105 | 17.9 | 34.2 | -1.2 |
| CS21 | - | 39.659 | -104.868 | 6,551 | 12 | 6 | 0 - 129 | 15.2 | 39.2 | -1.2 |
| CS22 | - | 39.758 | -104.957 | 6,551 | 12 | 7 | 0 - 118 | 17.5 | 35.4 | -0.9 |
| CS23 | - | 39.772 | -105.024 | 6,552 | 14 | 9 | 0 - 139 | 16.5 | 34.6 | -2.0 |
| CS25 | - | 39.776 | -104.833 | 6,551 | 12 | 7 | 0 - 135 | 16.2 | 35.8 | -1.8 |
| CS26 | - | 39.674 | -104.950 | 6,552 | 12 | 7 | 0 - 115 | 15.9 | 36.9 | -1.2 |

Atmospheric
Measurement
Techniques
Open Access

Discussions

EGU

| CS27 | - | 39.775 | -105.009 | 6,552 | 12 | 7 | 0 - 115 | 16.4 | 35.6 | -1.4 |
| CS29 | - | 39.760 | -104.918 | 6,552 | 11 | 7 | 0 - 114 | 15.7 | 37.5 | -1.2 |

199

200  The three Love My Air sensors co-located at the I25 Globeville sites (CS2, CS3, CS4) agreed well
201  with each other (correlation = 0.98) (**Figures S4** and **Figure S5**). To ensure that our co-located
202  dataset was well balanced across sites, we only retained measurements from CS2 at the I25
203  Globeville site. We were left with a total of 27,338 co-located measurements that we used to
204  develop a calibration algorithm. **Figure S6** displays the time-series plots of PM$_{2.5}$ from all other
205  Love My Air sensors in the network.

206

207  Reference monitors at La Casa, CAMP, I25 Globeville and I25 Denver, also reported minute-level
208  PM$_{2.5}$ concentrations between April 23 11:16 and September 30, 22:49. We joined minute-level
209  Love my Air concentrations with minute-level reference data at these sites. We had a total of
210  1,062,141 co-located minute-level measurements during this time period. As with the hourly-
211  averaged data, we only retained data from one of the Love My Air sensors at the I25 Globeville
212  site and were thus left with 815,608 measurements.  **Table S1** has information on the minute-level
213  co-located measurements. **Figure S7** displays the time-series plot of minute-level data from the
214  LCS at the four co-location sites. As can be seen, the data at the minute-level displays more
215  variation and peaks in PM$_{2.5}$ concentrations than the hourly-averaged measurements, likely due to
216  the impact of passing sources. It is also important to mention that minute-level reference data may
217  have some additional uncertainties given the time resolution. Unless explicitly referenced, we will
218  be reporting results from using hourly-averaged measurements.

219

220  We found that RH and T reported by the Love My Air sensors were well correlated with that
221  reported by the reference monitoring stations. We used the Love My Air T and RH measurements
222  in our calibration models as they most closely represent the conditions experienced by the sensors.

223

224  We derived dew-point (D) from T and RH reported by the Love My Air sensors using the
225  *weathermetrics* package in the programming language R (Anderson and Peng, 2012), as D has
226  been shown to be a good proxy of particle hygroscopic growth in previous research (Clements et
227  al., 2017; Malings et al., 2020). Some previous work has also used a nonlinear correction for RH in
228  the form of $RH^2/(1-RH)$, that we also calculated for this study.

229

230  We extracted hour, weekend, and month variables from the Canary-S sensors and converted hour
231  and month into cyclic values to capture periodicities in the data by taking the cosine and sine of
232  hour*$2\pi/24$ and month*$2\pi/12$, which we designate as cos_time, sin_time, cos_month and
233  sin_month, respectively. Sinusoidal corrections for seasonality have been shown to improve
234  accuracy of PM$_{2.5}$ measurements in machine learning models(Considine et al., 2021).

## 2.2 Statistical Modeling

236  The goal of the calibration algorithm is to predict, as accurately as possible, the 'true' PM$_{2.5}$
237  concentrations given the concentrations reported by the Love My Air sensors. At the co-located

238    sites, the FEM $PM_{2.5}$ measurements, which we take to be the "true" $PM_{2.5}$ concentrations, are the
239    dependent variable in the models. We tested 21 increasingly complex models that included T, RH,
240    D as well as metrics that captured the time-varying patterns of $PM_{2.5}$ to correct the Love My Air
241    $PM_{2.5}$ measurements (**Table 2**).

243    Sixteen models were multivariate models that were used in a recent paper (Barkjohn et al., 2021) to
244    calibrate another network of low-cost sensors: the PurpleAir, that rely on the same $PM_{2.5}$ sensor
245    (Plantower) as the Canary-S monitors in this study. As T, RH and D are not independent (**Figure
246    S8**), the 16 linear regression models include adding the meteorological conditions considered as
247    interaction terms, instead of additive terms. The remaining 5 relied on machine learning
248    techniques.

250    Machine learning models can capture more complex nonlinear effects (for instance, unknown
251    relationships between additional spatial and temporal variables). We opted to use the following
252    machine learning techniques that have been widely used in calibrating LCS:

254    1. *Random forest (RF)*: RF is a decision-tree-based machine learning algorithm that has been
255    shown to perform well in air quality predictions. Briefly, to generate a random forest model, the
256    user specifies the maximum number of trees that make up the forest. Each tree is constructed using
257    a bootstrapped random sample from the training data set. The origin node of the decision tree is
258    split into sub-nodes by considering a random subset of the possible explanatory variables. Trees
259    are split based on which of the explanatory variables in each subset is the strongest predictor of the
260    outcome. This process of node splitting is repeated until a terminal node is reached (Zimmerman et
261    al., 2018). For our random forest models, the terminal node was specified using a minimum node
262    size of five data points per node.

264    2. *Neural Network (NN)*: NN consists of interconnected neurons organized in layers. Each neuron
265    or unit passes received information through an activation function and produces output values that
266    are then processed by neurons in the next layer. The NN training process is based on updating the
267    weights of neurons via supervised learning (Spinelle et al., 2014). A simple single hidden layer
268    neural network with a linear transfer function was chosen in this study.

270    3. *Gradient Boosting (GB)*: GB is a decision-tree-based approach that uses 'boosting' methods to
271    improve model performance. 'Boosting' sequentially combines many 'weak' models (learners)
272    into a final, improved model. The final model is built in an additive forward stagewise manner
273    where at each step a new learner is added that minimizes the negative gradient using a least squares
274    approach. The residuals of the current model are then used as the input for the next tree allowing
275    the model to 'learn' from the errors of the previous models (Johnson et al., 2018).

277    4. *SuperLearner (SL)*: SL is an ensemble-based machine learning algorithm, which allows for the
278    simultaneous evaluation (by cross-validation) of a library of plausible machine learning algorithms
279    to determine which models are most appropriate for the data, based on minimizing a least squares

280 loss function, and then averages over these chosen models to produce a composite model (Van der
281 Laan et al., 2007).
282
283 All machine learning models were run using the *caret* package in R (Kuhn, 2015).

### 2.2.1 Types of Corrections

285 For each of the 21 models considered, we developed four main corrections:
286
287 (C1) Developed using training data for the entire period of co-location.
288 (C2) Developed using all data for the same week of the measurement.
289 (C3) Developed using co-located data collected for a brief period (2 weeks) at the beginning of the
290 study (Jan 1 - Jan 14, 2021).
291 (C4) Developed using co-located data collected for two 2-week periods in different seasons (Jan 1
292 - Jan 14, 2021, and May 1 - May 14, 2021).
293
294 Although models developed using co-located data over the entire time period (C1) tend to be more
295 accurate over the entire spatiotemporal data set, it is inefficient to re-run large models frequently
296 (incorporating new data). On-the-fly corrections (such as C2) can help characterize short-term
297 variation in air pollution and sensor characteristics. The duration of calibration is a key question
298 that remains unanswered (Liang, 2021). We opted to test corrections C3 and C4 as many low-cost
299 sensor networks rely on developing calibration models based on relatively short co-location
300 periods (deSouza et al., 2020b; West et al., 2020; Singh et al., 2021).

### 2.2.2 Cross-Validation techniques to avoid overfitting in the machine learning models

302 We used a Leave-One-Site (I25 Globeville, I25 Denver, La Casa, CAMP)-Out (LOSO) approach
303 for cross validation (CV) to prevent overfitting in our machine learning models (Models 17 - 21 in
304 **Table 2**). Briefly, we split the data into four groups, with each group excluding data from a single
305 reference monitoring site. In each cross-validation iteration, we selected each group in turn to fit
306 the model and made predictions at the left-out site. This CV approach was used to tune the hyper
307 parameters in the machine learning models adopted in this study using correction approaches: C1,
308 C2, C3 and C4.
309
310 For the correction conducted on the complete archived dataset (C1), we also conducted a leave-
311 out-by-date (LOBD) CV for the machine learning models considered (**Table 3**). For the LOBD
312 model validation method, the project time period was split into 3-week periods. Each period
313 contained between ~ 700 and 900 hourly data points, with typically more sensors running
314 continuously during later chunks as more sensors were deployed and came online over time.
315 Thirteen periods were available in total, and, for each test-train set, 12 periods were used to train
316 the correction model, whereas the remaining interval was selected to test the correction model. By
317 eliminating, using data from the same calendar week, where measurements are likely to be
318 correlated, we eliminate the possibility of obtaining overly optimistic model performance summary
319 statistics.
320

321    Models were generated for all combinations of training and test data. To summarize: each of the 21
322    calibration models considered was tested under four potential correction schemes (C1, C2, C3 and
323    C4). For C1, the machine-learning algorithms were trained using two CV approaches: LOSO and
324    LOBD, separately. For C2, C3 and C4 only LOSO was conducted, as model application is already
325    being performed on a different time period from the training. Note that for simple linear
326    regressions, overfitting is not an issue, and no CV is required.
327
328    Zusman et al., (2020) have reported that for more than 3 co-location sites, a LOSO CV is preferred,
329    as it replicates our ultimate objective of applying the calibration developed to other sites in the
330    network. However, in this case, due to the high correlation across co-located sites (**Figure S5,**
331    **Figure S6**), a LOBD CV is likely to produce more robust results.
332
333    Overall, we test 89 models (26 (C1) + 21 x 3 (C2, C3, C4) = 89) listed in **Tables 2** and **3**.

334    **2.2.3 Evaluating the correction models at the co-location sites**
335    **Figure S9** displays the PM$_{2.5}$ concentrations from the reference monitors and the corresponding
336    levels from the co-located Love My Air sensors by RH. Uncorrected Love My Air measurements
337    tend to be biased upwards by an average of ~12%.
338
339    We evaluate the performance of the corrections across the range of PM$_{2.5}$ concentrations for the
340    entire time period of co-location in our sample using the following metrics: R (Pearson correlation
341    coefficient), and RMSE (**Tables 2** and **3**). We also evaluated calibrations using corrections C3 and
342    C4 only for the time-period over which the calibration algorithm was developed, which was Jan 1 -
343    Jan 14, 2021, for C3 and Jan 1 - Jan 14, 2021, and May 1 - May 14, 2021 for C4 (**Table S2**).
344
345    Mean PM$_{2.5}$ concentrations from the reference monitors between Jan 1 - Jan 14, 2021, was 9 µg/m$^3$
346    (Median: 7 µg/m$^3$, Min:0 µg/m$^3$, Max: 79 µg/m$^3$). Nineteen measurements were > 30 µg/m$^3$.  Mean
347    PM$_{2.5}$ concentrations from the reference monitors between May 1 - May 14 was 6 µg/m$^3$ (Median:
348    5 µg/m$^3$, Min: 1 µg/m$^3$, Max: 22 µg/m$^3$). Zero measurements were > 30 µg/m$^3$.
349
350    We evaluated model performance for true/reference PM$_{2.5}$ concentrations > 30 µg/m$^3$ and ≤ 30
351    µg/m$^3$, as these concentrations account for the greatest differences in health and air pollution
352    avoidance behavior impacts (Nilson et al., 2022). Further, lower concentrations (PM$_{2.5}$ ≤ 30 µg/m$^3$)
353    represent most measurements observed in our network; better performance at these levels will
354    ensure better day-to-day functionality of the correction. In order to compare errors observed in the
355    two different concentration ranges, in addition to reporting R and RMSE of the calibration
356    approaches, we also report the normalized RMSE (normalized by the mean of the true
357    concentrations) (**Table S3**).
358
359    One of the key advantages of LCS is that they report high frequency measurements of pollution.
360    As reference monitoring stations provide hourly, or daily average pollution values, most often the
361    calibration algorithm is developed using hourly averaged data and then applied to the high
362    frequency LCS measurements. We applied the calibration algorithms described in **Tables 2** and **3**

363    developed using hourly-averaged co-located measurements on minute-level measurements from

364    the co-located LCS described in **Table S1**. We evaluated the performance of the corrected high-

365    frequency measurements against the 'true' measurements from the corresponding reference

366    monitor using the metrics R and RMSE (**Tables 4** and **5**).

367

368    ***Table 2****: Performance of the calibration models as captured using root mean square error*

369    *(RMSE), and Pearson correlation (R). LOSO CV was used to prevent overfitting in the machine*

370    *learning models. All corrected values were evaluated over the entire time-period (Jan 1 -*

371    *September 30, 2021)*

| ID | Name | Equation | C1 *Correction developed on data during the entire period of network operation* | | C2 *On-the-fly correction developed using data for the same week of measurement* | | C3 *Correction developed using measurements made in the first two weeks of January* | | C4 *Correction developed using measurements from the first two weeks of January and the first two weeks in May* | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | R | RMSE ($\mu g/m^3$) | R | RMSE ($\mu g/m^3$) | R | RMSE ($\mu g/m^3$) | R | RMSE ($\mu g/m^3$) |
| | **Raw Love My Air measurements** | | | | | | | | | |
| 0 | Raw | | 0.927 | 6.469 | - | - | - | - | - | - |
| | **Multivariate Regression (LOSO CV)** | | | | | | | | | |
| 1 | Linear | $PM_{2.5, corrected} = PM_{2.5}$ x s1 + b | 0.927 | 3.421 | 0.944 | 3.008 | 0.927 | 3.486 | 0.927 | 3.424 |
| 2 | +RH | $PM_{2.5, corrected} = PM_{2.5}$ x $s_1$ + RH x $s_2$ + b | 0.929 | 3.379 | 0.948 | 2.904 | 0.928 | 3.618 | 0.929 | 3.462 |
| 3 | +T | $PM_{2.5, corrected} = PM_{2.5}$ x $s_1$ + T x $s_2$ + b | 0.928 | 3.409 | 0.949 | 2.896 | 0.925 | 3.948 | 0.928 | 3.460 |
| 4 | +D | $PM_{2.5, corrected} = PM_{2.5}$ x $s_1$ + D x $s_2$ + b | 0.928 | 3.417 | 0.947 | 2.934 | 0.917 | 3.713 | 0.925 | 3.470 |
| 5 | +RH x T | $PM_{2.5, corrected} = PM_{2.5}$ x $s_1$ + RH x $s_2$ + T x $s_3$ + RH x T x $s_4$ + b | 0.934 | 3.260 | 0.953 | 2.782 | 0.931 | 3.452 | 0.933 | 3.344 |
| 6 | +RH x D | $PM_{2.5, corrected} = PM_{2.5}$ x $s_1$ + RH x $s_2$ + D x $s_3$ + RH x D x $s_4$ + b | 0.930 | 3.361 | 0.953 | 2.785 | 0.911 | 3.973 | 0.929 | 3.461 |
| 7 | +D x T | $PM_{2.5, corrected} = PM_{2.5}$ x $s_1$ + D x $s_2$ + T x $s_3$ + D | 0.928 | 3.409 | 0.952 | 2.798 | 0.888 | 5.698 | 0.921 | 3.720 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | $\times\,T\times s_4 + b$ | | | | | | | | |
| 8 | +RH x T x D | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + RH\times s_2 + T\times s_3 + D\times s_4 + RH\times T\times s_5 + RH\times D\times s_6 + T\times D\times s_7 + RH\times T\times D\times s_8 + b$ | 0.935 | 3.246 | 0.955 | 2.724 | 0.779 | 7.077 | 0.926 | 3.625 |
| 9 | PM x RH | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + RH\times s_2 + RH\times PM_{2.5}\times s_3 + b$ | 0.930 | 3.362 | 0.950 | 2.854 | 0.925 | 3.949 | 0.925 | 3.767 |
| 10 | PM x D | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + D\times s_2 + D\times PM_{2.5}\times s_3 + b$ | 0.932 | 3.324 | 0.950 | 2.871 | 0.883 | 4.460 | 0.913 | 3.777 |
| 11 | PM x T | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + T\times s_2 + T\times PM_{2.5}\times s_3 + b$ | 0.930 | 3.365 | 0.952 | 2.809 | 0.906 | 6.509 | 0.928 | 3.466 |
| 12 | PM x nonlinear RH | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + \frac{RH^2}{(1-RH)}\times s_2 + \frac{RH^2}{(1-RH)}\times PM_{2.5}\times s_3 + b$ | 0.934 | 3.277 | 0.948 | 2.900 | 0.931 | 3.510 | 0.932 | 3.403 |
| 13 | PM x RH x T | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + RH\times s_2 + T\times s_3 + PM_{2.5}\times RH\times s_4 + PM_{2.5}\times T\times s_5 + RH\times T\times s_6 + PM_{2.5}\times RH\times T\times s_7 + b$ | 0.938 | 3.165 | 0.956 | 2.672 | 0.891 | 6.220 | 0.928 | 3.497 |
| 14 | PM x RH x D | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + RH\times s_2 + D\times s_3 + PM_{2.5}\times RH\times s_4 + PM_{2.5}\times D\times s_5 + RH\times D\times s_6 + PM_{2.5}\times RH\times D\times s_7 + b$ | 0.933 | 3.288 | 0.957 | 2.663 | 0.879 | 7.289 | 0.917 | 4.033 |
| 15 | PM x T x D | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + T\times s_2 + D\times s_3 + PM_{2.5}\times T\times s_4 + PM_{2.5}\times D\times s_5 + T\times D\times s_6 + PM_{2.5}\times T\times D\times s_7 + b$ | 0.932 | 3.315 | 0.957 | 2.665 | 0.734 | 6.302 | 0.905 | 4.574 |
| 16 | PM x RH x T x D | $PM_{2.5,\,corrected} = PM_{2.5}\times s_1 + RH\times s_2 + T\times s_3 + D\times s_4 + PM_{2.5}\times RH\times s_5 + PM_{2.5}\times T\times s_6 + T\times RH\times s_7 + PM_{2.5}\times D\times s_8 + D\times RH\times s_9 + D\times T\times s_{10} + PM_{2.5}\times RH\times T\times s_{11} + PM_{2.5}\times RH\times$ | 0.940 | 3.115 | 0.960 | 2.557 | 0.324 | 32.951 | 0.765 | 6.746 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | $D \times s_{12} + PM_{2.5} \times D \times T \times s_{13} + D \times RH \times T \times s_{14} + PM_{2.5} \times RH \times T \times D \times s_{15} + b$ | | | | | | | | |
| | **Machine Learning (LOSO CV)** | | | | | | | | | |
| 17 | Random Forest | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.983 | 1.713 | 0.988 | 1.450 | 0.913 | 3.926 | 0.911 | 3.824 |
| 18 | Neural Network (One hidden layer) | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.933 | 3.286 | 0.948 | 2.916 | 0.932 | 3.550 | 0.913 | 4.725 |
| 19 | Gradient Boosting | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.950 | 2.870 | 0.964 | 2.452 | 0.910 | 3.854 | 0.909 | 3.834 |
| 20 | SuperLearner | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.950 | 2.855 | 0.970 | 2.236 | 0.910 | 3.917 | 0.923 | 3.582 |
| 21 | Random Forest | For C1: $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, cos\_time, cos\_month, sin\_month)$<br><br>For C2, C3, C4 $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, cos\_time)$ | 0.987 | 1.475 | 0.990 | 1.289 | 0.870 | 5.032 | 0.884 | 4.617 |

372

***Table 3***: *Performance of the calibration models using the C1 correction as captured using root mean square error (RMSE), normalized RMSE, and Pearson correlation (R) LOBD CV was used to prevent overfitting in the machine learning models*

| ID | Machine Learning (LOBD CV) | | R | RMSE ($\mu g/m^3$) |
|---|---|---|---|---|
| 17 | Random Forest | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.983 | 1.710 |
| 18 | Neural Network (One hidden layer) | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.933 | 3.285 |
| 19 | Gradient Boosting | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.953 | 2.759 |
| 20 | SuperLearner | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.956 | 2.692 |
| 21 | Random Forest | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, cos\_time, cos\_month, sin\_month)$ | 0.987 | 1.480 |

Atmospheric
Measurement
Techniques

Open Access

EGU

Discussions

### 2.3 Evaluating transferability

### 2.3.1 Evaluating the representativeness of meteorological conditions at the co-location sites of the entire network

We first evaluated if meteorological conditions (T and RH) at the co-location sites corresponding to measurements used to construct calibration models were representative of conditions of operation for the rest of the network by comparing distributions of these parameters across sites.

### 2.3.2 Evaluating transferability at the co-location sites

To evaluate how transferable the calibration technique developed at the co-located sites was to the rest of the network, we ran the models proposed in **Tables 2** and **3**, after leaving out each one of the 5 co-located sites in turn. We report the distribution of RMSE from each model across the different test datasets using boxplots (**Figure 2**).

We compare statistically the errors in predictions on each test dataset with errors in predictions from using all sites in our main analysis. Such an approach is useful to understand how well the proposed correction can transfer to other areas in the Denver region. To compare statistical difference between errors, t-tests were used to compare normally distributed datasets (as determined by Shapiro–Wilk), and Wilcoxon tests were used for nonparametric datasets with a significance value of 0.05.

We have only 5 co-location sites in the network. Although evaluating the transferability among these sites is useful, as we know the true $PM_{2.5}$ concentrations at these sites, we also evaluated the transferability of these models in the larger network by predicting $PM_{2.5}$ concentrations using the models proposed in **Tables 2** and **3** at each of the 24 sites in the Love My Air network. For each site, we display time series plots of corrected $PM_{2.5}$ measurements in order to visually compare the ensemble of corrected values at each site.

### 2.3.3 Evaluating the sensitivity of hotspot detection across the network of sensors to the calibration method

One of the key use-cases of low-cost sensors is hotspot detection. We report the labels of sites that are the most polluted using corrected measurements from the 89 different models using hourly data. We repeat this process for daily, weekly and monthly-averaged corrected measurements. We ignore missing measurements from the network when calculating time averaged values for the different time periods considered. We report the mean number of sensors that are ranked 'most polluted' across the different correction functions for the different averaging periods.

### 2.3.4 Evaluating sensitivity of the spatial and temporal trends of the low-cost sensor network to the method of calibration

We compared the differences in corrected $PM_{2.5}$ using similar methods to that in (Jin et al., 2019; deSouza et al., 2022) by calculating:

Atmospheric
Measurement
Techniques
Discussions

414      (1) The spatial root mean square difference (RMSD) between any two corrected exposures at

415      the same site: $SRMSD_{h,d} = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(Conc_{hi} - Conc_{di})^2}$, where $Conc_{hi}$ and $Conc_{di}$ are

416      Jan 1- September 30, 2021 averaged $PM_{2.5}$ concentrations estimated from correction h and

417      d for site i. N is the total number of sites.

418      (2) The temporal RMSD between pairs of exposures: $TRMSD_{h,d} =$

419      $\sqrt{\frac{1}{M}\sum_{t=1}^{M}(Conc_{ht} - Conc_{dt})^2}$, where $Conc_{ht}$ and $Conc_{dt}$ are hourly corrected $PM_{2.5}$

420      concentrations averaged over all operational Love My Air sites estimated from correction h

421      and d for time t. M is the total number of hours of operation of the network.

422      (3) The spatial pearson correlation coefficient: $R_S =$

423      $\frac{\sum_{i=1}^{N}(Conc_{hi}-\underline{Conc_h})(Conc_{di}-\underline{Conc_d})}{\sqrt{\sum_{i=1}^{N}(Conc_{hi}-\underline{Conc_h})^2 \sum_{i=1}^{N}(Conc_{di}-\underline{Conc_d})^2}}$, where $\underline{Conc_h}$ and $\underline{Conc_d}$ are the average (across

424      all sites and times) corrected $PM_{2.5}$ concentrations estimated from corrections h and d

425      respectively.

426      (4) The temporal pearson correlation coefficient: $R_T =$

427      $\frac{\sum_{t=1}^{M}(Conc_{ht}-\underline{Conc_h})(Conc_{dt}-\underline{Conc_d})}{\sqrt{\sum_{t=1}^{M}(Conc_{ht}-\underline{Conc_h})^2 \sum_{i=1}^{N}(Conc_{dt}-\underline{Conc_d})^2}}$

428

429 We characterized the uncertainty in the 'corrected' $PM_{2.5}$ estimates at each site across the different

430 models using two metrics: a normalized range (NR) and uncertainty. NR for a given site represents

431 the spread of $PM_{2.5}$ across the different correction approaches.

432      (5) $NR = \frac{1}{M}\sum_{t=1}^{M} \frac{max_{k \in K}\, C_{kt} - min_{k \in K}\, C_{kt}}{\underline{C_t}}$

433 $C_{kt}$ is the $PM_{2.5}$ concentration at hour t from the kth model from the ensemble of K (which in this

434 case is 89) correction approaches. $\underline{C_t}$ represents the ensemble mean across the K different products

435 at hour t. M is the total number of hours in our sample for which we have $PM_{2.5}$ data for the site

436 under consideration.

437

438 For our sample (K = 89), we assume the variations in $PM_{2.5}$ across multiple models follows the t-

439 statistical distribution with the mean being the ensemble average. The confidence interval (CI) for

440 the ensemble mean at a given time t is:

441

442      (6) $CI_t = \underline{C_t} + t^* \frac{SD_t}{\sqrt{K}}$

443 Where $\underline{C_t}$ represents the ensemble mean at time t; t* is the upper (1-CI)/2 critical value for the t-

444 distribution with K-1 degrees of freedom. For K=89, t* for the 95% double tailed confidence

445 interval is 1.99. $SD_t$ is the sample standard deviation at time t.

446      (7) $SD_t = \sqrt{\frac{\sum_{k=1}^{K}(C_{k,t}-\underline{C_t})^2}{K-1}}$

447

448      We define an overall estimate of uncertainty as follows:

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

449    (8) $uncertainty = \frac{1}{M}\sum_{t=1}^{M} t^* \frac{SD_t}{C_t\sqrt{K}}$ , which can also be expressed as

450    (8) $uncertainty = \frac{1}{M}\sum_{t=1}^{M} \frac{CI_t - C_t}{C_t}$

## 3 Results

### 3.1 Evaluating the correction models at the co-location sites

When we evaluated each of the 21 correction models proposed on the entire co-location dataset
(**Tables 2** and **3**), we found that the C2 correction performed better overall than the C1, C3 and C4
corrections.

We also found that for corrections C3 and C4, more complex models yielded a better performance
(for example the RMSE for Model 16: 2.813 μg/m$^3$, RMSE for Model 2: 0.915 μg/m$^3$ generated
using the C3 correction) when evaluated during the period of co-location, alone (**Table S2**).
However, when models generated using the C3 and C4 correction were transferred to the entire
time period of co-location, we find that more complex multivariate regression models (Models 13-
16) and the machine learning model (Model 21) that include cos_time, performed significantly
worse than the simpler models. In some cases, these models performed worse even than the
uncorrected measurements. For example, applying Model 16 generated using C3 on the entire
dataset resulted in an RMSE of 32.951 μg/m$^3$ compared to 6.469 μg/m$^3$ for the uncorrected
measurements. Including data for another season in the training sample (C4), resulted in
significantly increased performance of the calibration over the entire dataset compared to C3,
although it did not result in an improvement in performance for all models compared to the
uncorrected measurements. For example, Model 16 generated using C4 yielded an RMSE of 6.746
μg/m$^3$. Among the multivariate regression models, we found that models of the same form that
corrected for RH instead of T or D did best. The best performance was observed for models that
included the nonlinear correction for RH (Model 12) or included an RH X T term (Model 5)
(**Tables 2** and **3**).

For corrections C1 and C2, we found that an increase in complexity of model form resulted in a
decreased RMSE. Overall, Model 21 yielded the best performance (RMSE = 1.281 μg/m$^3$ when
using the C2 correction, and 1.475 μg/m$^3$ when using the C1 correction with a LOSO CV and
1.480 μg/m$^3$ when using a LOBD correction). In comparison, the simplest model that corrected for
bias yielded an RMSE of 3.421 μg/m$^3$ for the C1 correction, and 3.008 μg/m$^3$ when using the C2
correction.

For correction C1, using a LOBD CV with the machine learning models resulted in better
performance than using a LOSO CV, except for Model 21 which is an RF model with additional
time-of-day and month covariates, for which performance using the LOSO was slightly better
(RMSE: 1.475 μg/m$^3$ versus 1.480 μg/m$^3$).

When we evaluated how well the models performed at high PM$_{2.5}$ concentrations (> 30 μg/m$^3$)
versus lower concentrations (≤ 30 μg/m$^3$), we found that multivariate regression models generated

489    using the C1 correction did not perform well in capturing spikes in $PM_{2.5}$ concentrations
490    (normalized RMSE > 25%). Multivariate regression models generated using the C2 correction
491    performed better (normalized RMSE ~ 20 -25 %). Machine learning algorithms generated using
492    both C1 and C2 corrections captured $PM_{2.5}$ spikes well (C1: normalized RMSE ~ 10 - 25%, C2:
493    normalized RMSE ~ 10 - 20%). Specifically, the C2 RF model (Model 21) yielded the lowest
494    RMSE values (4.180 μg/m$^3$, normalized RMSE: 9.8%), of all models considered. Machine learning
495    models generated using the C1 corrected that were tuned using LOBD CV instead of LOSO
496    performed better in both $PM_{2.5}$ concentration regimes. Models generated using C3 and C4 had the
497    worst performance in both concentration regimes and yielded poorer agreement with reference
498    measurements than even the uncorrected measurements. As in the case with the entire dataset,
499    more complex multivariate regression models and machine learning models generated using C3
500    and C4 performed worse than more simple models in both $PM_{2.5}$ concentration intervals (**Tables**
501    **S3** and **S4**).
502
503    We then evaluated how well the models generated using C1, C2, C3 and C4 corrections performed
504    when applied to minute-level LCS data at co-located sites. We found that the machine learning
505    models generated using C1 and C2 improved the performance of the LCS (Model 21 (CV=LOSO)
506    generated using C1 yielded an RMSE of 15.482 μg/m$^3$ compared to 16.409 μg/m$^3$ obtained from
507    the uncorrected measurements.) The more complex multivariate regression models yielded a
508    significantly worse performance across all corrections. (Model 16 generated using C1 yielded an
509    RMSE of 41.795 μg/m$^3$.) As in the case with the hourly-averaged measurements, using correction
510    C1, LOBD CV instead of LOSO for the machine learning models resulted in better model
511    performance except for Model 21. Few models generated using C3 and C4 resulted in improved
512    performance when applied to the minute-level measurements (**Tables 4** and **5**).
513
514    ***Table 4****: Performance of the calibration models developed using the co-located hourly*
515    *measurements to the minute-level data as captured using root mean square error (RMSE), and*
516    *Pearson correlation (R). LOSO CV was used to prevent overfitting in the machine learning models.*
517    *All corrected values were evaluated over the entire time period (April 23 - September 30, 2021).*

| ID | Name | Equation | C1 *Correction developed on data during the entire period of network operation* | | C2 *On-the-fly correction developed using data for the same week of measurement* | | C3 *Correction developed using measurements made in the first two weeks of January* | | C4 *Correction developed using measurements from the first two weeks of January and the first two weeks in May* | |
|----|------|----------|------|------|------|------|------|------|------|------|
| | | | **R** | **RMSE (μg/m$^3$)** | **R** | **RMSE (μg/m$^3$)** | **R** | **RMSE (μg/m$^3$)** | **R** | **RMSE (μg/m$^3$)** |
| | **Raw Love My Air measurements** | | | | | | | | | |

| 0 | Raw | | 0.497 | 16.409 | - | - | - | - | - | - |
|---|---|---|---|---|---|---|---|---|---|---|
| | **Multivariate Regression (LOSO CV)** | | | | | | | | | |
| 1 | Linear | $PM_{2.5, corrected} = PM_{2.5} \times s1 + b$ | 0.497 | 15.667 | 0.498 | 15.646 | 0.497 | 15.657 | 0.497 | 15.663 |
| 2 | +RH | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + b$ | 0.495 | 15.678 | 0.500 | 15.618 | 0.492 | 15.721 | 0.494 | 15.686 |
| 3 | +T | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + T \times s_2 + b$ | 0.496 | 15.670 | 0.500 | 15.621 | 0.493 | 15.822 | 0.495 | 15.671 |
| 4 | +D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + D \times s_2 + b$ | 0.497 | 15.663 | 0.498 | 15.640 | 0.491 | 15.805 | 0.495 | 15.693 |
| 5 | +RH x T | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + T \times s_3 + RH \times T \times s_4 + b$ | 0.499 | 15.634 | 0.500 | 15.621 | 0.495 | 15.669 | 0.498 | 15.640 |
| 6 | +RH x D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + D \times s_3 + RH \times D \times s_4 + b$ | 0.496 | 15.671 | 0.500 | 15.622 | 0.477 | 15.892 | 0.494 | 15.684 |
| 7 | +D x T | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + D \times s_2 + T \times s_3 + D \times T \times s_4 + b$ | 0.470 | 15.928 | 0.014 | 323.684 | 0.018 | 257.153 | 0.032 | 135.647 |
| 8 | +RH x T x D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + T \times s_3 + D \times s_4 + RH \times T \times s_5 + RH \times D \times s_6 + T \times D \times s_7 + RH \times T \times D \times s_8 + b$ | 0.138 | 33.817 | 0.041 | 111.569 | 0.029 | 160.447 | 0.027 | 160.963 |
| 9 | PM x RH | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + RH \times PM_{2.5} \times s_3 + b$ | 0.494 | 15.688 | 0.501 | 15.615 | 0.485 | 15.896 | 0.486 | 15.844 |
| 10 | PM x D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + D \times s_2 + D \times PM_{2.5} \times s_3 + b$ | 0.498 | 15.644 | 0.499 | 15.630 | 0.477 | 16.145 | 0.491 | 15.820 |
| 11 | PM x T | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + T \times s_2 + T \times PM_{2.5} \times s_3 + b$ | 0.495 | 15.675 | 0.501 | 15.610 | 0.483 | 17.172 | 0.495 | 15.675 |
| 12 | PM x nonlinear RH | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + \frac{RH^2}{(1-RH)} \times s_2 + \frac{RH^2}{(1-RH)} \times PM_{2.5} \times s_3 + b$ | 0.496 | 15.659 | 0.497 | 15.650 | 0.494 | 15.705 | 0.495 | 15.681 |
| 13 | PM x RH x T | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + T \times s_3 + PM_{2.5} \times RH \times s_4 + PM_{2.5} \times T \times s_5 + RH \times T \times s_6 + PM_{2.5} \times RH \times T \times s_7 + b$ | 0.501 | 15.611 | 0.502 | 15.601 | 0.462 | 17.111 | 0.489 | 15.732 |
| 14 | PM x RH x D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + D \times s_3 + PM_{2.5} \times RH \times s_4 + PM_{2.5} \times D \times s_5 + RH \times D \times s_6 + PM_{2.5} \times RH \times D \times s_7 + b$ | 0.496 | 15.657 | 0.502 | 15.602 | 0.460 | 17.710 | 0.479 | 15.948 |

18

| ID | Model | Equation | R | RMSE | R | RMSE | R | RMSE | R | RMSE |
|---|---|---|---|---|---|---|---|---|---|---|
| 15 | PM x T x D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + T \times s_2 + D \times s_3 + PM_{2.5} \times T \times s_4 + PM_{2.5} \times D \times s_5 + T \times D \times s_6 + PM_{2.5} \times T \times D \times s_7 + b$ | 0.134 | 35.196 | 0.020 | 217.684 | 0.012 | 178.589 | 0.044 | 114.530 |
| 16 | PM x RH x T x D | $PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + T \times s_3 + D \times s_4 + PM_{2.5} \times RH \times s_5 + PM_{2.5} \times T \times s_6 + T \times RH \times s_7 + PM_{2.5} \times D \times s_8 + D \times RH \times s_9 + D \times T \times s_{10} + PM_{2.5} \times RH \times T \times s_{11} + PM_{2.5} \times RH \times D \times s_{12} + PM_{2.5} \times D \times T \times s_{13} + D \times RH \times T \times s_{14} + PM_{2.5} \times RH \times T \times D \times s_{15} + b$ | 0.112 | 41.795 | 0.029 | 159.921 | 0.010 | 482.333 | 0.019 | 203.714 |
| | **Machine Learning (LOSO CV)** | | | | | | | | | |
| 17 | Random Forest | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.505 | 15.565 | 0.510 | 15.527 | 0.489 | 15.863 | 0.488 | 15.821 |
| 18 | Neural Network (One hidden layer) | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.496 | 15.669 | 0.501 | 15.611 | 0.495 | 15.699 | 0.477 | 16.202 |
| 19 | Gradient Boosting | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.500 | 15.625 | 0.502 | 15.604 | 0.485 | 15.779 | 0.486 | 15.765 |
| 20 | SuperLearner | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.500 | 15.622 | 0.503 | 15.591 | 0.483 | 15.805 | 0.490 | 15.719 |
| 21 | Random Forest | For C1: $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, cos\_time, cos\_month, sin\_month)$<br><br>For C2, C3, C4: $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, cos\_time)$ | 0.514 | 15.482 | 0.512 | 15.502 | 0.481 | 16.349 | 0.481 | 16.185 |

518

519 **Table 5**: *Performance of the calibration models developed using the co-located hourly*
520 *measurements to the minute-level data as captured using root mean square error (RMSE), and*
521 *Pearson correlation (R). LOBD CV was used to prevent overfitting in the machine learning*
522 *models. All corrected values were evaluated over the entire time period (April 23 - September 30,*
523 *2021)*

| ID | Machine Learning (LOBD CV) | | R | RMSE ($\mu g/m^3$) |
|---|---|---|---|---|

| 17 | Random Forest | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.506 | 15.561 |
|---|---|---|---|---|
| 18 | Neural Network (One hidden layer) | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.496 | 15.666 |
| 19 | Gradient Boosting | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.501 | 15.610 |
| 20 | SuperLearner | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH)$ | 0.503 | 15.594 (1.326) |
| 21 | Random Forest | $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, cos\_time, cos\_month, sin\_month)$ | 0.510 | 15.516 |

## 3.1 Evaluating the representativeness of meteorological conditions at the co-location sites of the entire network

Temperature at the co-located sites across the entire period of the experiment during the development of C1 were similar to those at the rest of Love My Air network (**Figure S10**). The sensor CS19 is the only one that recorded lower temperatures than those at any of the other sites. Relative humidity at the co-located sites appears to be larger than at the other sites in the network (**Figure S11**).

We also compared meteorological conditions during the development of corrections C3 (Jan 1 - Jan 14, 2021) and C4 (Jan 1 - Jan 14, 2021 and May 1 - May 14, 2021), to those measured during the duration of network operation (C3: **Figures S12** and **S13**; C4: **Figures S14** and **S15**). Temperatures at the co-located sites during the development of C3 were on average lower than those reported during the operation of the network. Temperatures at the co-located sites during the development of C4 were more representative of the network than C3, although they too are smaller than the average temperatures experienced by the network. RH values during C3 and C4 tend to be on the higher side and are not representative of conditions experienced by some Love My Air sensors.

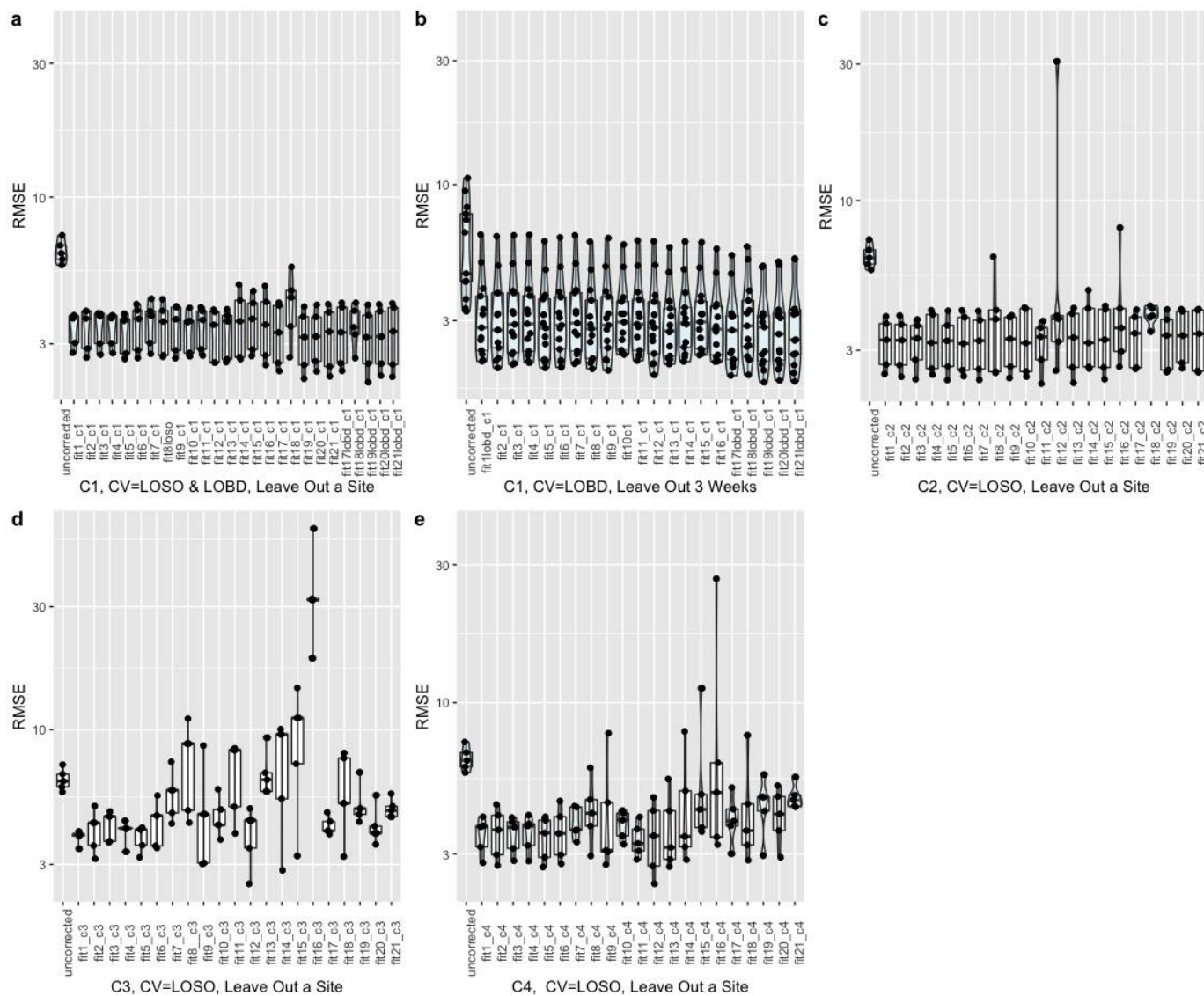We then evaluated the transferability of the corrections developed.

## 3.2 Evaluating transferability at the co-location sites

**Figure 2** shows the performance (RMSE) of corrected Love My Air PM$_{2.5}$ data by generating corrections based on the 21 models previously proposed using the (a) C1 correction, CV= LOSO and CV = LODB for Models 17 - 21, when leaving out a test site (**Figure 2a**). Also shown is the result using the C1 correction when leaving out a three week period of data at a time and generating corrections based on the data from the remaining time periods across each site and using CV = LOBD for Models 17 - 21 (**Figure 2b**). Finally, **Figures 2c, 2d** and **2e illustrate** using the C2, C3 and C4 corrections, respectively, (CV= LOSO for Models 17 - 21) when leaving out a test site.

Large reductions in RMSE are observed when applying simple linear corrections (*Models 1 - 4*) to the uncorrected data across all corrections. Increasing the complexity of the model does not result

555 in marked changes in correction performance on different test sets for C1 and C2. Although the
556 performance of the corrected datasets did improve on average for some of the complex models
557 considered (*Model 17, 20, 21* for example, vis-a-vis simple linear regressions when using the C1
558 correction) (**Figures 2a, 2b**), this was not the case for *all* test datasets considered, as evinced by the
559 overlapping distributions of RMSE performances (e.g., Model 11 using the C2 correction resulted
560 in a worse fit for one of the test datasets). For C3 and C4, the performance of corrections was
561 worse across all datasets for the more complex multivariate model formulations (**Figures 2d**, **2e**),
562 indicating that using uncorrected data is better than using these corrections and calibration models.
563
564 Wilcoxon tests and t-tests (based on whether Shapiro-Wilk tests revealed that the distribution of
565 RMSEs was normal) revealed significant improvements in the distribution of RMSEs for all
566 corrected test sets vis-a-vis the uncorrected data. There was no significant difference in the
567 distribution of RMSE values from applying C1 and C2 corrections to the test sets, across the
568 different models. For corrections C3 and C4, we found significant differences in the distribution of
569 RMSEs obtained from running different models on the data, implying that the choice of model has
570 a significant impact on transferability of the calibration models to other monitors.

Atmospheric
Measurement
Techniques
Discussions
Open Access
EGU

Figure 2: *Performance (RMSE) of corrected Love My Air PM$_{2.5}$ data by generating corrections based on the 21 models previously proposed using (a) Correction C1 when leaving out a co-location site in turn and then running the generated correction on the test site (Note that for machine learning models (Models 17- 21), we performed CV using a LOSO CV as well as a LOBD CV approach), (b) Correction C1 when leaving out 3 week periods of data at a time and generating corrections based on the data from the remaining time periods across each site, and evaluating the performance of the developed corrections on the held out 3 weeks of data (Note that for machine learning models (Models 17- 21), we performed CV using a LOBD CV approach), (c) Correction C2 when leaving out a co-location site in turn and then running the generated correction on the test site, (c) Correction C3 when leaving out a co-location site in turn and then running the generated correction on the test site, (c) Correction C4 when leaving out a co-location site in turn and then running the generated correction on the test site. Each point represents the*

22

584 *RMSE for each test dataset permutation. The distribution of RMSEs is displayed using boxplots*
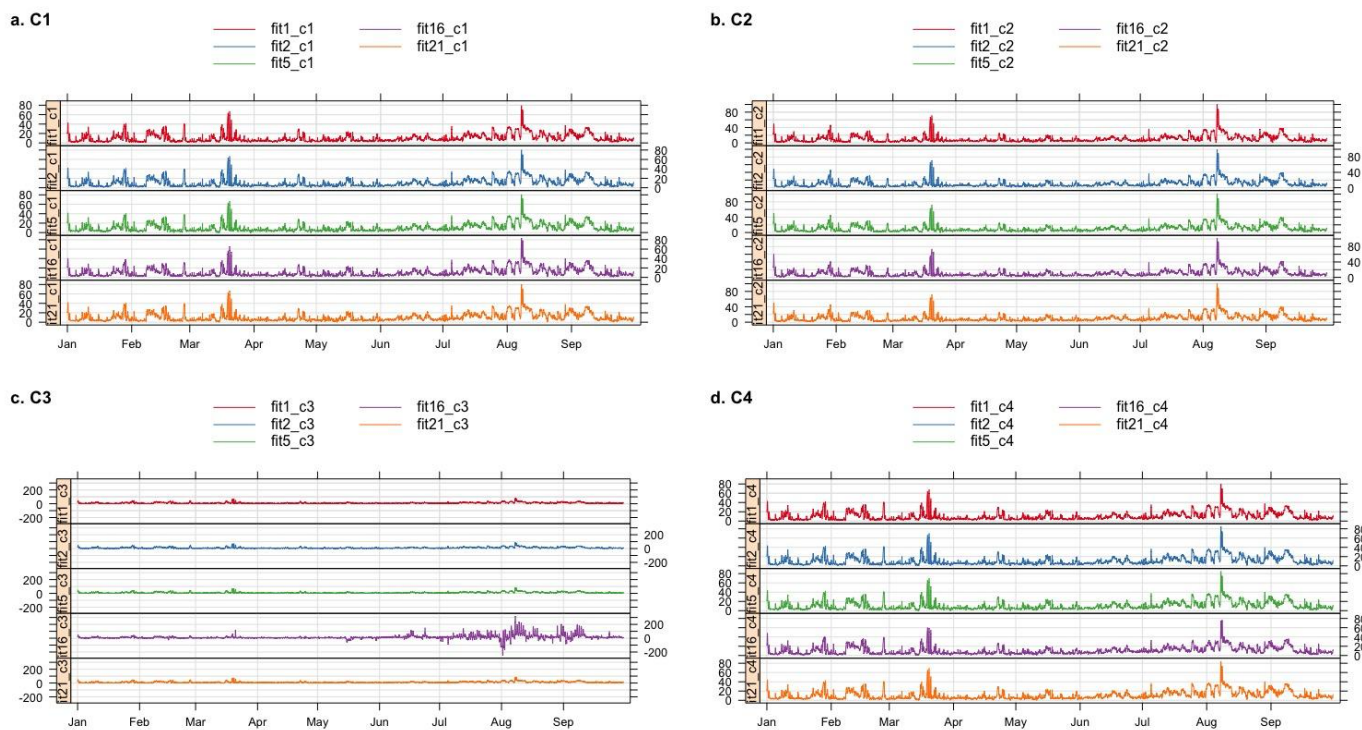585 *and violinplots*
586
587 The time-series of corrected PM$_{2.5}$ values for Models 1, 2, 5, 16, and 21 (RF using additional
588 variables) (using CV = LOSO for the machine learning Models 17 and 21) for corrections
589 generated using C1, C2, C3 and C4 are displayed in **Figure 3** for Love My Air sensor CS 1. These
590 subsets of models were chosen as they cover the range of model forms considered in this analysis.
591
592 From **Figure 3**, we note that although the different corrected values from C1 and C2 track each
593 other well, there are small systematic differences between the different corrections. Peaks in
594 corrected values using on-the-fly data tend to be higher than those using archived data. Peaks in
595 corrected values using machine learning methods on the archived data are higher than those
596 generated from multivariate regression models. There are marked differences in the corrected
597 values from C3 and C4. Specifically Model 16 yields peaks in the data that corrections using the
598 other models do not generate. This pattern was consistent when applying this suite of corrections to
599 other Love My Air sensors.
600



602 **Figure 3**: *Time-series of the different PM$_{2.5}$ corrected values for Models 1, 2, 5, 16 and 21 across*
603 *corrections (a) C1, (b) C2, (c )C3 and (d) C4 for the Love My Air monitor CS1*
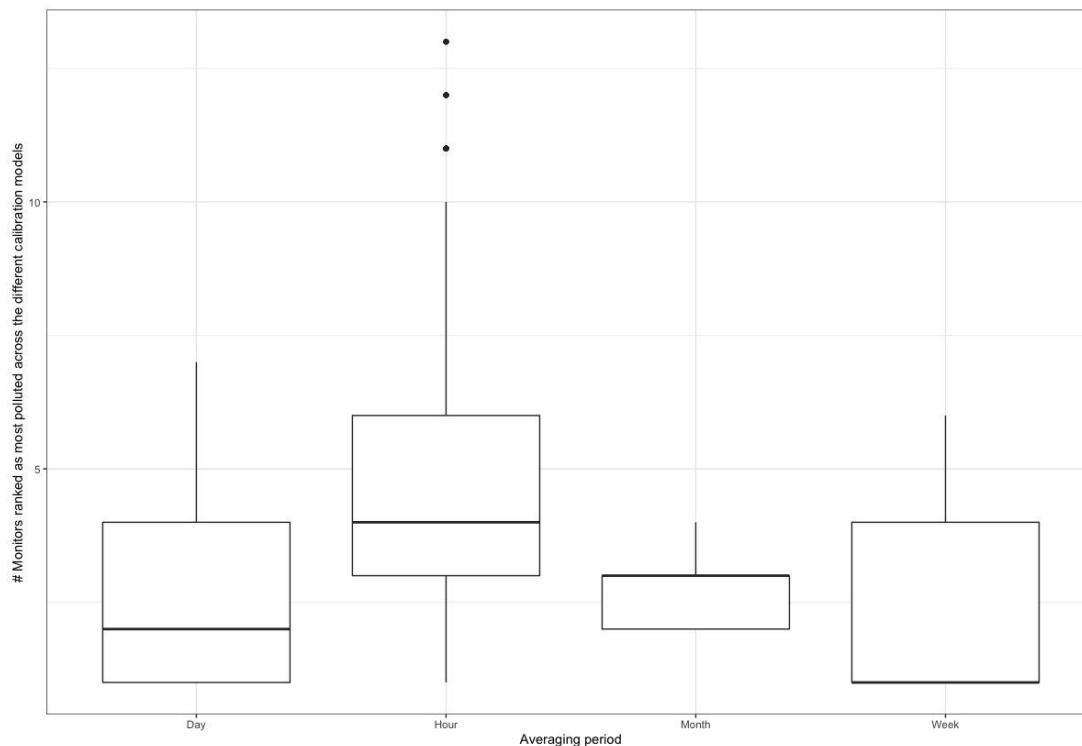
Atmospheric
Measurement
Techniques
Open Access
Discussions
EGU

### 3.3 Evaluating the sensitivity of hotspot detection across the network of sensors to the calibration method

604
605

606  Mean (95% CI) PM$_{2.5}$ concentrations across the different models (as well as CV technique) and
607  corrections (26 (C1) + 21 x 3 (C2, C3, C4) = 89 listed in **Tables 2** and **3**) at each Love My Air site
608  for the duration of the experiment (Jan 1 - September 30, 2021) are displayed in **Figure S16**. Due
609  to overlap between the different corrected measurements across sites, identification of the most
610  polluted site is dependent on the correction algorithm used. We examined the sensitivity of the
611  'most polluted site' at different time-intervals

612

613  Every hour, we ranked the different monitors for each of the 89 different corrections. We found
614  that there were on average 4.4 (median = 5) monitors that were ranked most polluted. When this
615  calculation was repeated using daily-averaged corrected data, there were on average 2.5 (median =
616  2) monitors that were ranked the most polluted. The corresponding value for weekly-corrected data
617  was 2.4 (median = 1), and for monthly data was 3 (median = 3) (**Figure 4**).
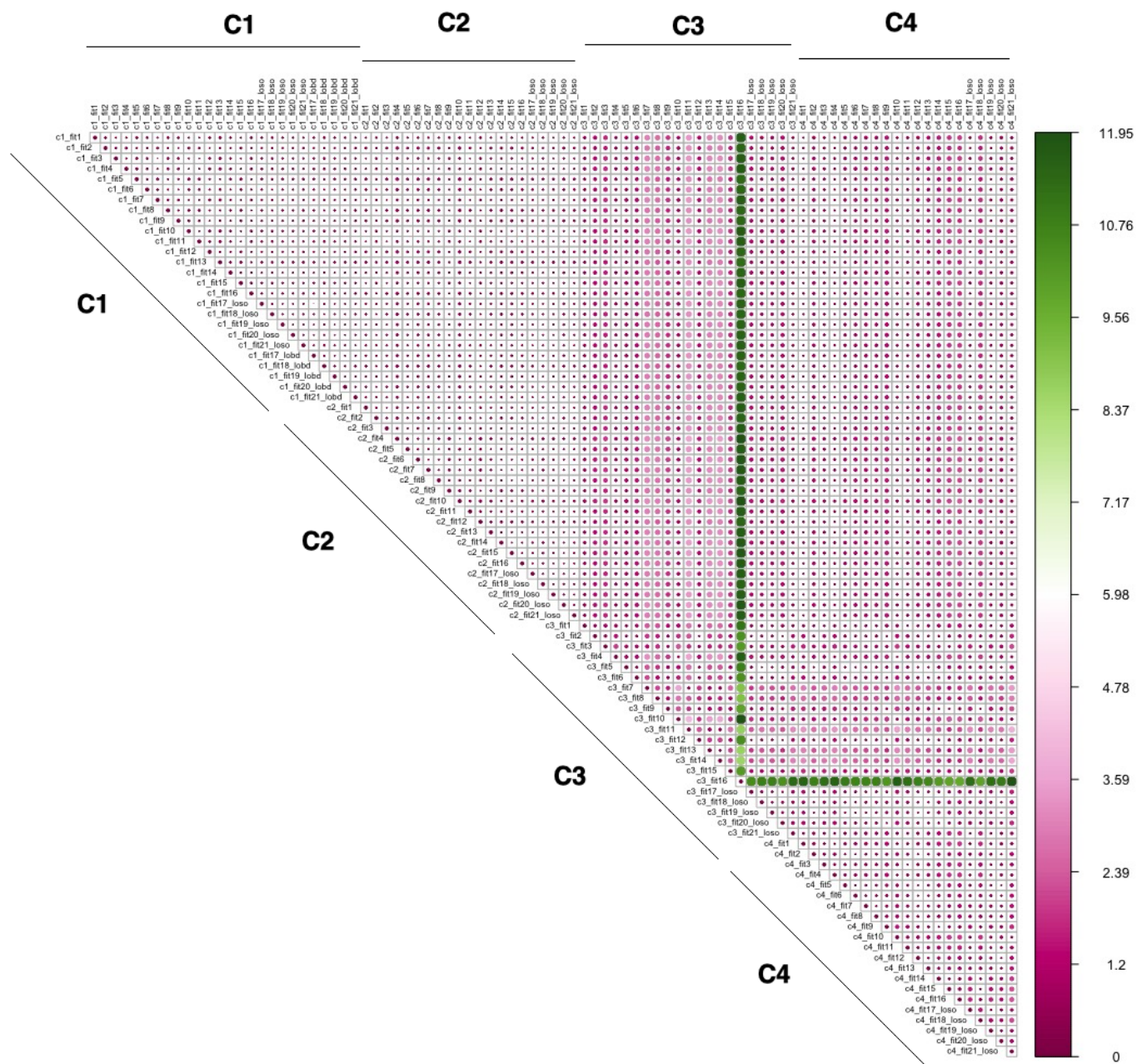


618
619  *Figure 4: Variation in the number of sites that were ranked as 'most polluted' across the 89*
620  *different corrections for different time-averaging periods displayed using boxplots*

### 3.4 Evaluating sensitivity of the spatial and temporal trends of the low-cost sensor network to the method of calibration

621
622

623  The spatial and temporal RMSD values between corrected values generated from applying each of
624  the 89 models using the four different correction approaches across all monitoring sites in the Love
625  My Air network are displayed **Figures 5** and **6**, respectively. It appears that there is larger temporal

626    variation (max 32.79 μg/m$^3$), in comparison to spatial variations displayed across corrections (max:
627    11.95 μg/m$^3$). Model 16 generated using the C3 correction has the greatest spatial and temporal
628    RMSD in comparison with all other models. Models generated using the C3 and C4 corrections
629    displayed the greatest spatial and temporal RMSD vis-a-vis C1 and C2. **Figures S17- S20** display
630    spatial RMSD values between all models corresponding to corrections C1-C4, respectively.
631    **Figures S21- S24** display temporal RMSD values between all models corresponding to corrections
632    C1-C4, respectively. Across all corrections the temporal RMSD between models is greater than the
633    spatial RMSD.
634
635    Spatial and temporal correlation coefficients between corrected measurements generated from
636    applying all 89 models using the four different correction approaches across the entire network are
637    displayed in **Figures S25** and **S26**, respectively. The spatial correlations are lower than temporal
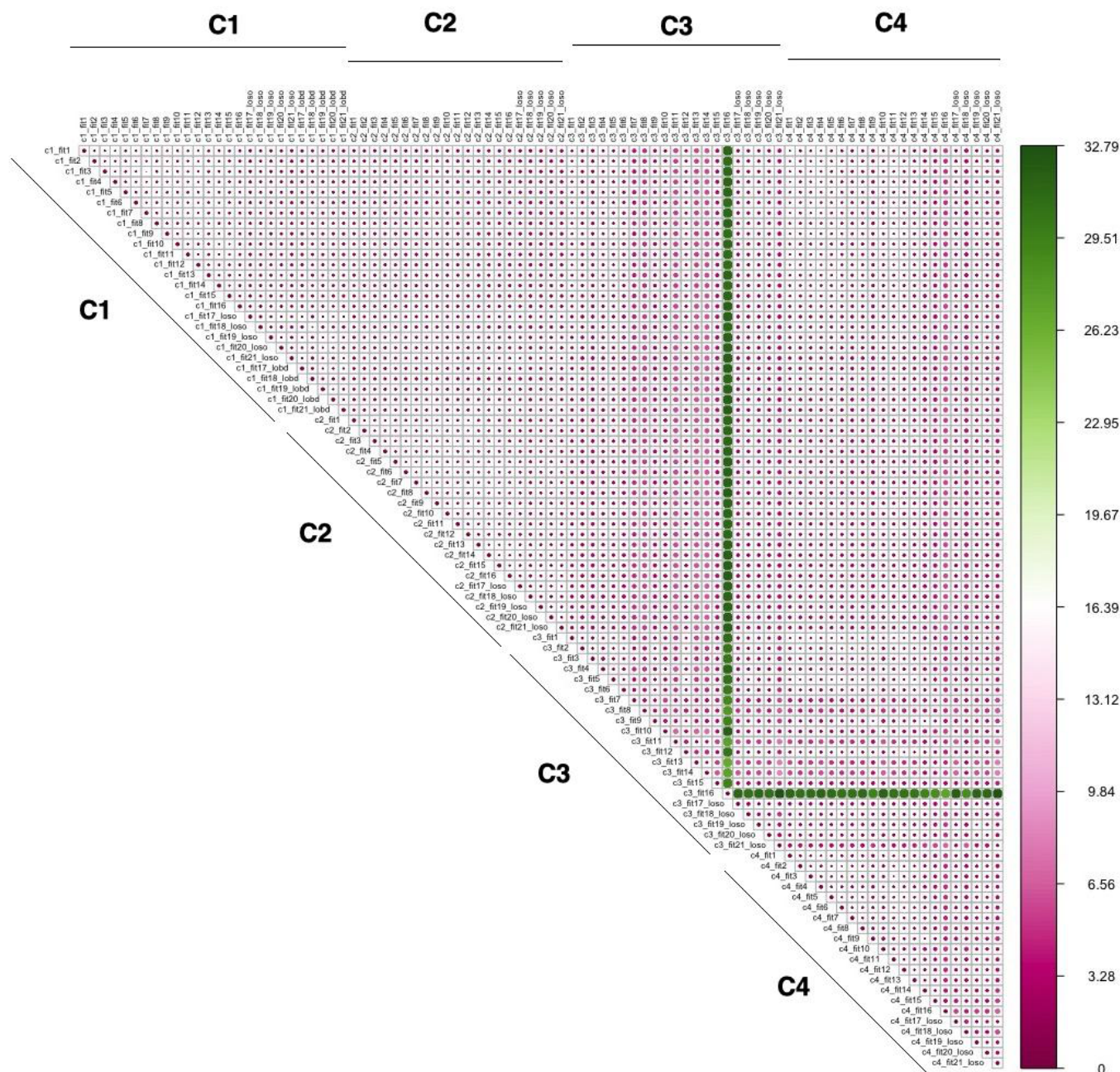638    correlations between corrected measurements.

**Figure 5**: *Spatial RMSD (μg/m$^3$) calculated using the method detailed in section 2.3.4 from applying each of the 89 models using the four different correction approaches to all monitoring sites in the Love My Air network*
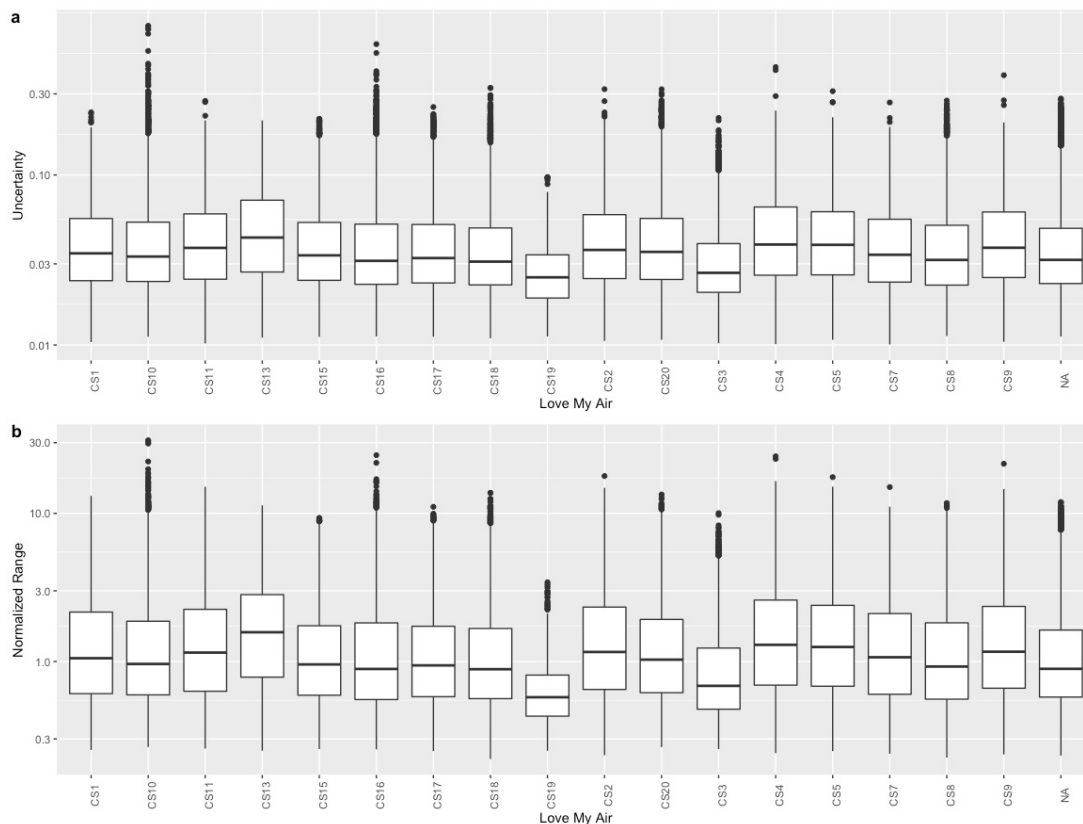
*Figure 6: Temporal RMSD (µg/m³) calculated using the method detailed in section 2.3.4 from applying each of the 89 models using the four different correction approaches to all monitoring sites in the Love My Air network*

The distribution of uncertainty and the NR in hourly corrected measurements over the 89 models by monitor are displayed in **Figure 7**. Overall, there are small differences in uncertainties and NR of the exposure assessment across sites. The average NR and uncertainty across all sites are 1.554

27

Atmospheric
Measurement
Techniques
Discussions

653    (median: 0.9768) and 0.044 (median: 0.033) , respectively.  We note that although the uncertainties
654    in the data are small, the average normalized range tends to be quite large.

655



656
657    **Figure 7**: *Distribution of (a) uncertainty and (b) normalized range (NR) in hourly-corrected*
658    *measurements across all 89 correction models at each site using the methodology described in*
659    *Section 2.3.4*

## 4 Discussion and Conclusions

661    In our analysis of how transferable the correction algorithms developed at the Love My Air co-
662    location sites are to the rest of the network, we found that for C1 and C2 corrections, more
663    complex models yielded a better fit at the co-location sites. When examining the C3 and C4
664    corrections, we found that although these corrections appeared to significantly improve LCS
665    measurements for the time period of model development (**Table S2**), when applied to the entire
666    time period of operation they did not perform well. Many of the models, especially the more
667    complex multivariate regression models, performed significantly worse than even the uncorrected
668    measurements. This indicates that calibration models generated during short-term time periods,
669    even if the time periods correspond to different seasons, may not necessarily transfer well to other
670    times, likely due to changes in the aerosol composition, and differences in meteorological
671    conditions, among other potential factors. This suggests the need for calibration models to be
672    developed over longer time periods that better capture different LCS operating conditions. For C3

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

673   and C4, we found models that relied on nonlinear formulations of RH, that serve as proxies for
674   hygroscopic growth, yielded the best performance, as compared to more complex models. This
675   suggests that physics-based calibrations are potentially an alternative approach when relying on
676   short co-location periods and need to be explored further.

678   When evaluating how transferable the calibration models using the different correction approaches
679   were to the rest of the network, we found that for C1 and C2, more complex models that appeared
680   to perform well at the co-location sites did not necessarily transfer best to the rest of the network.
681   Specifically, when we tested these models on a co-located site that was left out when generating
682   the correction, we found that some of the more complex models run using the C2 correction
683   yielded a significantly worse performance at some test sites (**Figure 2**). If the corrected data were
684   going to be used to make site-specific decisions then such corrections would lead to important
685   errors. When evaluating C3 and C4 correction approaches we observed a large distribution of
686   RMSE values across sites. For several of the more complex models developed using C3 and C4
687   corrections, the RMSE values were larger than observed for the uncorrected data, suggesting that
688   certain calibration models could result in even more error-prone data than using uncorrected
689   measurements.

691   For C1 and C2, we found that there were no significant differences in the distribution of the
692   performance metric: RMSE of corrected measurements from simpler models in comparison to
693   those derived from more complex corrections at test sites (**Figure 2**). For C3 and C4, we found
694   significant differences in the distribution of RMSE across test sites, which indicates that these
695   models are likely site-specific and not easily transferable to other sites in the network. This
696   suggests that less complex models might be preferred when short-term co-locations are carried out
697   for sensor calibration.

699   Our findings reinforce the idea that evaluating calibration models at all co-location sites using
700   overall metrics like RMSE should not be seen as the only/best way to determine how to calibrate a
701   network of LCS. Instead, approaches like LOSO, LOBD, or a combination of these, as
702   demonstrated should be used to evaluate calibration transferability.

704   We also found that the calibration models yielded different performance results at different $PM_{2.5}$
705   concentration ranges. Machine learning models developed using C1, and models developed using
706   C2 were better than multivariate regression models generated using C1 at capturing peaks in
707   pollution (> 30 μg/m$^3$). All models using C3 and C4 yielded poor performance results across both
708   concentration ranges ($PM_{2.5} > 30$ μg/m$^3$ and $PM_{2.5} \le 30$ μg/m$^3$).

710   When evaluating how well the calibration models translated to minute-level data (**Tables 4** and **5**),
711   we observed that machine learning models generated using C1 and C2, improved the LCS
712   measurements. More complex multivariate regression models performed poorly. All C3 and C4
713   models also performed poorly. This suggests that caution needs to be exercised when transferring
714   models developed at a particular time scale to another (**Tables S3** and **S4**).

715

716     Our findings thus far indicate that different calibration approaches are required for different end
717     purposes. There may not be a single one-size-fits-all calibration approach.
718

719     We found that the 'most polluted' site in the Love My Air network was dependent on the
720     calibration algorithm used on the network. We found that for the Love My Air network, the
721     detection of the most polluted site was sensitive to the duration of time-averaging of the corrected
722     measurements (**Figure 4**). Hotspot detection was most robust using weekly-averaged
723     measurements. Such an analysis thus reveals the most robust temporal scale for decision-making
724     related to evaluating hotspots.
725

726     We found that the temporal RMSD (**Figure 6**) was greater than the spatial RMSD (**Figure 5**) for
727     the ensemble of 47 corrected exposure assessments developed for the Love My Air network. One
728     of the reasons this may be the case is that $PM_{2.5}$ concentrations across the different Love My Air
729     sites in Denver are highly correlated (**Figure S5**), indicating that the contribution of local sources
730     to $PM_{2.5}$ concentrations in Denver is small. Due to the low variability in $PM_{2.5}$ concentrations
731     across sites, it makes sense that the variations in the corrected $PM_{2.5}$ concentrations will be seen in
732     time rather than space. The largest pairwise temporal RMSD were all seen between corrections
733     derived from complex models using the C3 correction.
734

735     However, we note that the temporal correlation coefficients (**Figure S26**) for all-pairwise
736     correction models were higher than the corresponding spatial coefficients (**Figure S25**). This
737     implies that although the corrections generated from all models considered tended to track each
738     other (except for a few models using C3) some corrected values were biased low, whereas some
739     were biased high. It's important to understand under what conditions these biases occur. One of the
740     ways this can be determined is by evaluating the performance of the calibrated data under different
741     conditions, such as in different pollution regimes as demonstrated in this paper (**Tables S3** and **S4**).
742

743     Finally, we observed that the uncertainty in $PM_{2.5}$ concentrations across the ensemble of
744     corrections was consistently small for the Love My Air Denver network. The normalized range in
745     the corrected measurements, on the other hand, was large, indicating that the corrections yield a
746     large range of corrected measurements; however, most of the corrected measurements fall within a
747     relatively small interval. Thus, deciding which calibration algorithm to pick has important
748     consequences for decision-makers using data from this network.
749

750     In summary: this paper makes the case that it is not enough to evaluate calibration algorithms
751     based on metrics of performance at co-located sites, alone. We need to:
752

753     1) Evaluate models under different conditions (e.g., pollution concentrations) to evaluate the
754     circumstances under which different calibration algorithms do well to determine which model to
755     use for which use-case.
756

Atmospheric
Measurement
Techniques
Open Access
EGU
Discussions

757   2) Determine how well calibration adjustments can be transferred to other locations. Specifically,
758   although we found that in Denver some corrections performed well at co-location sites, they could
759   result in large errors at specific sites that would create difficulties for site-specific decision making.
760
761   3) Examine how well calibration adjustments can be transferred to other time periods. In this study
762   we found that models developed using the C3 correction were not transferable to other time
763   periods because the conditions during the co-location were not representative of broader operating
764   conditions in the network.
765
766   4) Evaluate how well calibration algorithms developed for a specific time-scale transfer to
767   measurements at other time intervals.
768
769   5) Use a variety of approaches to quantify transferability, both focusing on co-location sites (using
770   a LOSO and/or LOBD cross-validation scheme) and looking at the wider low-cost sensor network
771   (e.g., with spatio-temporal correlations and RMSD). The metrics proposed in this paper to evaluate
772   model transferability can be used in other networks.
773
774   6) Investigate how adopting a certain timescale for averaging measurements could mitigate the
775   uncertainty induced by the calibration process. Namely, we found that in the Love My Air
776   network, hotspot identification was more robust to using daily-averaged data than hourly-averaged
777   data.
778
779   In this work, the Love My Air network under consideration is located over a fairly small area in a
780   single city. In this network, for the time period considered, $PM_{2.5}$ seems to be mainly a regional
781   pollutant and the contribution of local sources is small. More work needs to be done to evaluate
782   model transferability in networks in other settings. Concerns about model transferability are likely
783   to be even more key when thinking about larger networks that span different cities and should be
784   considered in future research.

## Author Contributions

786   PD conceptualized the study, developed the methodology, carried out the analysis and wrote the first draft.
787   TS and WO provided PD with access to the data. PD and BC obtained funding for this study. BC produced
788   Figure 1. All authors helped in refining the methodology and editing the draft.

## Acknowledgements

## Data Availability

796   The data used in this study can be obtained from the author on request
797

798 **Competing Interests**

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

800 **References**

801 Anderson, G. and Peng, R.: weathermetrics: Functions to convert between weather metrics (R package),
802 2012.
803
804 State of Global Air: https://www.stateofglobalair.org/, last access: 18 June 2020.
805
806 Apte, J. S., Messier, K. P., Gani, S., Brauer, M., Kirchstetter, T. W., Lunden, M. M., Marshall, J. D., Portier,
807 C. J., Vermeulen, R. C. H., and Hamburg, S. P.: High-Resolution Air Pollution Mapping with Google Street
808 View Cars: Exploiting Big Data, Environ. Sci. Technol., 51, 6999–7008,
809 https://doi.org/10.1021/acs.est.7b00891, 2017.
810
811 Barkjohn, K. K., Gantt, B., and Clements, A. L.: Development and application of a United States-wide
812 correction for PM$_{2.5}$ data collected with the PurpleAir sensor, Atmospheric Meas. Tech., 14, 4617–4637,
813 https://doi.org/10.5194/amt-14-4617-2021, 2021.
814
815 Bean, J. K.: Evaluation methods for low-cost particulate matter sensors, Atmospheric Meas. Tech., 14,
816 7369–7379, https://doi.org/10.5194/amt-14-7369-2021, 2021.
817
818 Bi, J., Wildani, A., Chang, H. H., and Liu, Y.: Incorporating Low-Cost Sensor Measurements into High-
819 Resolution PM2.5 Modeling at a Large Spatial Scale, Environ. Sci. Technol., 54, 2152–2162,
820 https://doi.org/10.1021/acs.est.9b06046, 2020.
821
822 Brantley, H. L., Hagler, G. S. W., Herndon, S. C., Massoli, P., Bergin, M. H., and Russell, A. G.:
823 Characterization of Spatial Air Pollution Patterns Near a Large Railyard Area in Atlanta, Georgia, Int. J.
824 Environ. Res. Public. Health, 16, 535, https://doi.org/10.3390/ijerph16040535, 2019.
825
826 Castell, N., Dauge, F. R., Schneider, P., Vogt, M., Lerner, U., Fishbain, B., Broday, D., and Bartonova, A.:
827 Can commercial low-cost sensor platforms contribute to air quality monitoring and exposure estimates?,
828 Environ. Int., 99, 293–302, https://doi.org/10.1016/j.envint.2016.12.007, 2017.
829
830 Clements, A. L., Griswold, W. G., Rs, A., Johnston, J. E., Herting, M. M., Thorson, J., Collier-Oxandale,
831 A., and Hannigan, M.: Low-Cost Air Quality Monitoring Tools: From Research to Practice (A Workshop
832 Summary), Sensors, 17, 2478, https://doi.org/10.3390/s17112478, 2017.
833
834 Considine, E. M., Reid, C. E., Ogletree, M. R., and Dye, T.: Improving accuracy of air pollution exposure
835 measurements: Statistical correction of a municipal low-cost airborne particulate matter sensor network,
836 Environ. Pollut., 268, 115833, https://doi.org/10.1016/j.envpol.2020.115833, 2021.
837
838 Crilley, L. R., Shaw, M., Pound, R., Kramer, L. J., Price, R., Young, S., Lewis, A. C., and Pope, F. D.:
839 Evaluation of a low-cost optical particle counter (Alphasense OPC-N2) for ambient air monitoring,
840 Atmospheric Meas. Tech., 11, 709–720, https://doi.org/10.5194/amt-11-709-2018, 2018.
841
842 deSouza, P. and Kinney, P. L.: On the distribution of low-cost PM 2.5 sensors in the US: demographic and
843 air quality associations, J. Expo. Sci. Environ. Epidemiol., 31, 514–524, https://doi.org/10.1038/s41370-
844 021-00328-2, 2021.
845
846 deSouza, P., Anjomshoaa, A., Duarte, F., Kahn, R., Kumar, P., and Ratti, C.: Air quality monitoring using
847 mobile low-cost sensors mounted on trash-trucks: Methods development and lessons learned, Sustain. Cities
848 Soc., 60, 102239, https://doi.org/10.1016/j.scs.2020.102239, 2020a.

849
850   deSouza, P., Lu, R., Kinney, P., and Zheng, S.: Exposures to multiple air pollutants while commuting:
851   Evidence from Zhengzhou, China, Atmos. Environ., 118168,
852   https://doi.org/10.1016/j.atmosenv.2020.118168, 2020b.
853
854   deSouza, P. N.: Key Concerns and Drivers of Low-Cost Air Quality Sensor Use, Sustainability, 14, 584,
855   https://doi.org/10.3390/su14010584, 2022.
856
857   deSouza, P. N., Dey, S., Mwenda, K. M., Kim, R., Subramanian, S. V., and Kinney, P. L.: Robust
858   relationship between ambient air pollution and infant mortality in India, Sci. Total Environ., 815, 152755,
859   https://doi.org/10.1016/j.scitotenv.2021.152755, 2022.
860
861   Giordano, M. R., Malings, C., Pandis, S. N., Presto, A. A., McNeill, V. F., Westervelt, D. M., Beekmann,
862   M., and Subramanian, R.: From low-cost sensors to high-quality data: A summary of challenges and best
863   practices for effectively calibrating low-cost particulate matter mass sensors, J. Aerosol Sci., 158, 105833,
864   https://doi.org/10.1016/j.jaerosci.2021.105833, 2021.
865
866   Hagler, G. S. W., Williams, R., Papapostolou, V., and Polidori, A.: Air Quality Sensors and Data
867   Adjustment Algorithms: When Is It No Longer a Measurement?, Environ. Sci. Technol., 52, 5530–5531,
868   https://doi.org/10.1021/acs.est.8b01826, 2018.
869
870   Holstius, D. M., Pillarisetti, A., Smith, K. R., and Seto, E.: Field calibrations of a low-cost aerosol sensor at
871   a regulatory monitoring site in California, Atmospheric Meas. Tech., 7, 1121–1131,
872   https://doi.org/10.5194/amt-7-1121-2014, 2014.
873
874   Jin, X., Fiore, A. M., Civerolo, K., Bi, J., Liu, Y., Donkelaar, A. van, Martin, R. V., Al-Hamdan, M., Zhang,
875   Y., Insaf, T. Z., Kioumourtzoglou, M.-A., He, M. Z., and Kinney, P. L.: Comparison of multiple PM 2.5
876   exposure products for estimating health benefits of emission controls over New York State, USA, Environ.
877   Res. Lett., 14, 084023, https://doi.org/10.1088/1748-9326/ab2dcb, 2019.
878
879   Johnson, N. E., Bonczak, B., and Kontokosta, C. E.: Using a gradient boosting model to improve the
880   performance of low-cost aerosol monitors in a dense, heterogeneous urban environment, Atmos. Environ.,
881   184, 9–16, https://doi.org/10.1016/j.atmosenv.2018.04.019, 2018.
882
883   Kim, K.-H., Kabir, E., and Kabir, S.: A review on the human health impact of airborne particulate matter,
884   Environ. Int., 74, 136–143, https://doi.org/10.1016/j.envint.2014.10.005, 2015.
885
886   Kuhn, M.: caret: Classification and Regression Training, Astrophys. Source Code Libr., ascl:1505.003,
887   2015.
888
889   Kumar, P., Morawska, L., Martani, C., Biskos, G., Neophytou, M., Di Sabatino, S., Bell, M., Norford, L.,
890   and Britter, R.: The rise of low-cost sensing for managing air pollution in cities, Environ. Int., 75, 199–205,
891   https://doi.org/10.1016/j.envint.2014.11.019, 2015.
892
893   Liang, L.: Calibrating low-cost sensors for ambient air monitoring: Techniques, trends, and challenges,
894   Environ. Res., 197, 111163, https://doi.org/10.1016/j.envres.2021.111163, 2021.
895
896   Magi, B. I., Cupini, C., Francis, J., Green, M., and Hauser, C.: Evaluation of PM2.5 measured in an urban
897   setting using a low-cost optical particle counter and a Federal Equivalent Method Beta Attenuation Monitor,
898   Aerosol Sci. Technol., 54, 147–159, https://doi.org/10.1080/02786826.2019.1619915, 2020.
899
900   Malings, C., Tanzer, R., Hauryliuk, A., Saha, P. K., Robinson, A. L., Presto, A. A., and Subramanian, R.:
901   Fine particle mass monitoring with low-cost sensors: Corrections and long-term performance evaluation,
902   Aerosol Sci. Technol., 54, 160–174, https://doi.org/10.1080/02786826.2019.1623863, 2020.

903
904  Morawska, L., Thai, P. K., Liu, X., Asumadu-Sakyi, A., Ayoko, G., Bartonova, A., Bedini, A., Chai, F.,
905  Christensen, B., Dunbabin, M., Gao, J., Hagler, G. S. W., Jayaratne, R., Kumar, P., Lau, A. K. H., Louie, P.
906  K. K., Mazaheri, M., Ning, Z., Motta, N., Mullins, B., Rahman, M. M., Ristovski, Z., Shafiei, M.,
907  Tjondronegoro, D., Westerdahl, D., and Williams, R.: Applications of low-cost sensing technologies for air
908  quality monitoring and exposure assessment: How far have they gone?, Environ. Int., 116, 286–299,
909  https://doi.org/10.1016/j.envint.2018.04.018, 2018.
910
911  Nilson, B., Jackson, P. L., Schiller, C. L., and Parsons, M. T.: Development and Evaluation of Correction
912  Models for a Low-Cost Fine Particulate Matter Monitor, Atmospheric Meas. Tech. Discuss., 1–16,
913  https://doi.org/10.5194/amt-2021-425, 2022.
914
915  Singh, A., Ng'ang'a, D., Gatari, M. J., Kidane, A. W., Alemu, Z. A., Derrick, N., Webster, M. J.,
916  Bartington, S. E., Thomas, G. N., Avis, W., and Pope, F. D.: Air quality assessment in three East African
917  cities using calibrated low-cost sensors with a focus on road-based hotspots, Environ. Res. Commun., 3,
918  075007, https://doi.org/10.1088/2515-7620/ac0e0a, 2021.
919
920  Snyder, E. G., Watkins, T. H., Solomon, P. A., Thoma, E. D., Williams, R. W., Hagler, G. S. W., Shelow,
921  D., Hindin, D. A., Kilaru, V. J., and Preuss, P. W.: The Changing Paradigm of Air Pollution Monitoring,
922  Environ. Sci. Technol., 47, 11369–11377, https://doi.org/10.1021/es4022602, 2013.
923
924  Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Calibration of a cluster of
925  low-cost sensors for the measurement of air pollution in ambient air, in: 2014 IEEE SENSORS, 2014 IEEE
926  SENSORS, 21–24, https://doi.org/10.1109/ICSENS.2014.6984922, 2014.
927
928  Van der Laan, M. J., Polley, E. C., and Hubbard, A. E.: Super learner, Stat. Appl. Genet. Mol. Biol., 6,
929  2007.
930
931  West, S. E., Buker, P., Ashmore, M., Njoroge, G., Welden, N., Muhoza, C., Osano, P., Makau, J., Njoroge,
932  P., and Apondo, W.: Particulate matter pollution in an informal settlement in Nairobi: Using citizen science
933  to make the invisible visible, Appl. Geogr., 114, 102133, https://doi.org/10.1016/j.apgeog.2019.102133,
934  2020.
935
936  Williams, R., Kilaru, V., Snyder, E., Kaufman, A., Dye, T., Rutter, A., Russel, A., and Hafner, H.: Air
937  Sensor Guidebook, US Environmental Protection Agency, Washington, DC, EPA/600/R-14/159 (NTIS
938  PB2015-100610), 2014.
939
940  Zimmerman, N., Presto, A. A., Kumar, S. P. N., Gu, J., Hauryliuk, A., Robinson, E. S., Robinson, A. L.,
941  and R. Subramanian: A machine learning calibration model using random forests to improve sensor
942  performance for lower-cost air quality monitoring, Atmospheric Meas. Tech., 11, 291–313,
943  https://doi.org/10.5194/amt-11-291-2018, 2018.
944
945  Zusman, M., Schumacher, C. S., Gassett, A. J., Spalt, E. W., Austin, E., Larson, T. V., Carvlin, G., Seto, E.,
946  Kaufman, J. D., and Sheppard, L.: Calibration of low-cost particulate matter sensors: Model development
947  for a multi-city epidemiological study, Environ. Int., 134, 105329,
948  https://doi.org/10.1016/j.envint.2019.105329, 2020.