

# Calibrating Networks of Low-Cost Air Quality Sensors

Priyanka deSouza<sup>1\*</sup>, Ralph Kahn<sup>2</sup>, Tehya Stockman<sup>3,4</sup>, William Obermann<sup>3</sup>, Ben Crawford<sup>5</sup>, An Wang<sup>6</sup>, James Crooks<sup>7</sup>, Jing Li<sup>8</sup>, Patrick Kinney<sup>9</sup>

1: Department of Urban and Regional Planning, University of Colorado Denver, 80202

2: NASA Goddard Space Flight Center, Greenbelt MD

3: Denver Department of Public Health and Environment, USA

4: Department of Civil, Environmental, and Architectural Engineering, University of Colorado Boulder, Boulder, Colorado 80309, United States

5: Department of Geography and Environmental Sciences, University of Colorado Denver, 80202

6: Senseable City Lab, Massachusetts Institute of Technology, Cambridge 02139

7: Division of Biostatistics and Bioinformatics, National Jewish Health, 2930

8: Department of Geography and the Environment, University of Denver, Denver, CO, USA

9: Department of Epidemiology, University of Colorado at Denver - Anschutz Medical Campus, 129263

10: Boston University School of Public Health, Boston, MA, USA

\*: [priyanka.desouza@ucdenver.edu](mailto:priyanka.desouza@ucdenver.edu)

## Abstract

Ambient fine particulate matter (PM<sub>2.5</sub>) pollution is a major health risk. Networks of low-cost sensors (LCS) are increasingly being used to understand local-scale air pollution variation. However, measurements from LCS have uncertainties that can act as a potential barrier to effective decision-making. LCS data thus need adequate calibration to obtain good quality PM<sub>2.5</sub> estimates. In order to develop calibration factors, one or more LCS are typically co-located with reference monitors for short- or long periods of time. A calibration model is then developed that characterizes the relationships between the raw output of the LCS and measurements from the reference monitors. This calibration model is then typically *transferred* from the co-located sensors to other sensors in the network. Calibration models tend to be evaluated based on their performance only at co-location sites. It is often implicitly assumed that the conditions at the relatively sparse co-location sites are representative of the LCS network overall, and that the calibration model developed is not overfitted to the co-location sites. Little work has explicitly evaluated how transferable calibration models developed at co-location sites are to the rest of an LCS

37 network, even after appropriate cross-validation. Further, few studies have evaluated the  
38 sensitivity of key LCS use-cases such as hotspot detection to the calibration model  
39 applied. Finally, there has been a dearth of research on how the duration of co-location  
40 (short-term/long-term) can impact these results. This paper attempts to fill these gaps  
41 using data from a dense network of LCS monitors in Denver deployed through the city's  
42 Love My Air program. It offers a series of transferability metrics for calibration models that  
43 can be used in other LCS networks and some suggestions as to which calibration model  
44 would be most useful for achieving different end goals.

45

46 **Key words:** low-cost sensors, PM<sub>2.5</sub>, calibration, LoveMyAir

## 47 **1 Introduction**

48 Poor air quality is currently the single largest environmental risk factor to human health in  
49 the world, with ambient air pollution responsible for approximately 6.7 million premature  
50 deaths every year (State of Global Air, 2020). Having accurate air quality measurements  
51 is crucial for tracking long-term trends in air pollution levels, identifying hotspots, and for  
52 developing effective pollution management plans. The dry-mass concentration of fine  
53 particulate matter (PM<sub>2.5</sub>), a criterion pollutant that poses more of danger to human health  
54 than other widespread pollutants (Kim et al., 2015), can vary over distances as small as ~  
55 10's of meters in complex urban environments (Brantley et al., 2019; deSouza et al.,  
56 2020a). Therefore, dense monitoring networks are often needed to capture relevant  
57 spatial variations. Due to their costliness, Environmental Protection Agency (EPA) air  
58 quality reference monitoring networks, the gold standard for measuring air pollutants, are  
59 sparsely positioned across the US (Apte et al., 2017; Anderson and Peng, 2012).

60

61 Low-cost sensors (LCS) (<USD \$2500 as defined by the US EPA Air Sensor Toolbox)  
62 (Williams et al., 2014) have the potential to capture concentrations of PM in previously  
63 unmonitored locations and to democratize air pollution information (Castell et al., 2017;  
64 Crawford et al., 2021; Kumar et al., 2015; Morawska et al., 2018; Snyder et al., 2013;  
65 deSouza and Kinney, 2021; deSouza, 2022). However, LCS measurements have several  
66 sources of greater uncertainty than reference monitors (Bi et al., 2020; Giordano et al.,  
67 2021; Liang, 2021).

68

69 Most low-cost PM sensors rely on optical measurement techniques. Optical instruments  
70 face inherent challenges that introduce potential differences in mass estimates compared  
71 to reference methods (Barkjohn et al., 2021; Crilley et al., 2018; Giordano et al., 2021;  
72 Malings et al., 2020):

73

74 1. Optical methods do not directly measure mass concentrations; rather, they estimate  
75 mass based on calibrations that convert light scattering data to particle number and mass.  
76 LCS come with factory-supplied calibrations, but in practice must be re-calibrated in the  
77 field to ensure accuracy, due to variations in ambient particle characteristics and  
78 instrument drift.

- 79
- 80 2. High relative humidity (RH) can produce hygroscopic particle growth, leading to dry  
81 mass overestimation unless particle hydration can accurately be taken into account or the  
82 particles are desiccated by the instrument.
- 83
- 84 3. LCS are not able to detect particles with diameters below a specific size, which is  
85 determined by the wavelength of laser light within each device, and is generally in the  
86 vicinity of 0.3  $\mu\text{m}$ , whereas the peak in pollution particle number size distribution is  
87 typically smaller than 0.3  $\mu\text{m}$ .
- 88
- 89 4. The physical and chemical parameters describing the aerosol (particle size  
90 distribution, shape, indices of refraction, hygroscopicity, volatility etc.), that might vary  
91 significantly across different microenvironments with diverse sources, impact light  
92 scattering; this in turn affects the aerosol mass concentrations reported by these  
93 instruments.

94

95 The need for field calibration to correct LCS measurements is particularly important. This  
96 is typically done by co-locating a small number of LCS with one or a few reference  
97 monitors at a representative monitoring location or locations. The co-location could be  
98 carried out for a brief period before and/or after the actual study or may continue at a  
99 small number of sites for the duration of the study. In either case, the co-location provides  
100 data from which a calibration model is developed that relates the raw output of the LCS as  
101 closely as possible to the desired quantity as measured by the reference monitor.  
102 Thereafter, the calibration model is transferred to other LCS in the network, based upon  
103 the presumption that ongoing sampling conditions are within the same range as those at  
104 the collocation site(s) during the calibration period.

105

106 Calibration models typically correct for 1) systematic error in LCS by adjusting for bias  
107 using reference monitor measurements, and 2) the dependence of LCS measurements  
108 on environmental conditions affecting the ambient particle properties such as relative  
109 humidity (RH), temperature (T), and/or dew-point (D). Correcting for RH, T and D is  
110 carried out through either a) a physics-based approach that accounts for aerosol  
111 hygroscopic growth given particle composition using  $\kappa$ -Köhler's theory, or b) empirical  
112 models, such as regression and machine learning techniques. In this paper, we focus on  
113 the latter, as it is currently the most widely used (Barkjohn et al., 2021). Previous work  
114 has also shown that the two approaches yield comparable improvements in the case of  
115  $\text{PM}_{2.5}$  LCS (Malings et al., 2020).

116

117 Prior studies have used multivariate regressions, piecewise linear regressions, or higher-  
118 order polynomial models to account for RH, T and D in these calibration models (Holstius  
119 et al., 2014; Magi et al., 2020; Zusman et al., 2020). More recently, machine learning  
120 techniques such as random forests, neural networks, and gradient boosted decision trees  
121 have been used (Considine et al., 2021; Liang, 2021; Zimmerman et al., 2018).

122 Researchers have also started including additional covariates in their models besides  
123 what is directly measured by the LCS, such as time of day, seasonality, wind direction,  
124 and site-type, which have been shown to yield significantly improved results (Considine et  
125 al., 2021).

126  
127 Past research has shown that there are several important decisions, in addition to the  
128 choice of calibration model, that need to be made during calibration and that can impact  
129 the results (Bean, 2021; Giordano et al., 2021; Hagler et al., 2018). These include a) the  
130 kind of reference air quality monitor used, b) the time-interval (e.g., hour/day) over which  
131 to average measurements used when developing the calibration model, c) how cross-  
132 validation (e.g., leave one site out/10-fold cross-validation) is carried out, and d) how long  
133 the co-location experiment takes place.

134  
135 Calibration models are typically evaluated based on how well the corrected data agree  
136 with measurements from reference monitors at the corresponding co-location site. A  
137 commonly used metric is the Pearson correlation coefficient,  $R$ , which quantifies the  
138 strength of the association. However, it is a misleading indicator of sensor performance  
139 when measurements are observed close to the limit of detection of the instrument.  
140 Therefore, Root Mean Square Error (RMSE) is often included in practice. Unfortunately,  
141 neither of these metrics captures how well the calibration method developed at the co-  
142 located sites *transfers* to the rest of the network in both time and space.

143  
144 If the conditions at the co-location sites (meteorological conditions, pollution source mix)  
145 for the period of co-location are the same as for the rest of the network during the total  
146 operational period, the calibration model developed at the co-location sites can be  
147 assumed to be transferable to the rest of the network. In order to ensure that the sampling  
148 conditions at the co-location site are representative of sampling conditions across the  
149 network, most researchers tend to deploy monitors in the same general sampling area as  
150 the network (Zusman et al., 2020). However, it is difficult to definitively test if the co-  
151 location site during the period of co-location is representative of conditions at all monitors  
152 in the network; ambient PM concentrations can vary on scales as small as a few meters.  
153 Furthermore, LCS are often deployed specifically in areas where the air pollution  
154 conditions are poorly understood, meaning that representativeness cannot be assessed in  
155 advance.

156  
157 In order to evaluate whether calibration models are transferable in time, we test if models  
158 generated using typical short-term co-locations at specific co-location sites perform well  
159 during other time periods at all co-location sites. Where multiple co-location sites exist,  
160 one way to evaluate how transferable calibration models are in space is to leave out one  
161 or more co-location sites and test if the calibration model is transferable to the left-out  
162 sites. This method was used in recent work evaluating the feasibility of developing a US-  
163 wide calibration model for the PurpleAir low-cost sensor network (Barkjohn et al., 2021;  
164 Nilson et al., 2022).

165  
166 Although these approaches are useful, co-location sites are sparse relative to other sites  
167 in the network. Even in the PurpleAir network (which is one of the densest low-cost  
168 networks in the world) there were only 39 co-location sites in 16 US states, a small  
169 fraction of the several thousand PurpleAir sites overall (Barkjohn et al., 2021). It is thus  
170 important to develop metrics to test how *sensitive* the spatial and temporal trends of  
171 pollution derived from the entire network are to the calibration model applied. Finally, a  
172 key use-case of LCS networks is to identify hotspots. It is important to also evaluate how  
173 sensitive the hotspot identified in an LCS network is to the calibration model applied.  
174

175 Examining the reliability of calibration models is timely because more researchers are  
176 opting to use machine learning models. Although in most cases, such models have  
177 yielded better results than traditional linear regressions, it is important to examine if these  
178 models are overfitted to conditions at the co-location sites, even after appropriate cross-  
179 validation, and how transferable they are to the rest of the network. Indeed, because of  
180 concerns of overfitting, some researchers have explicitly eschewed employing machine  
181 learning calibration models altogether (Nilson et al., 2022). It is important to test under  
182 what circumstances such concerns might be warranted.  
183

184 This paper uses a dense low-cost PM<sub>2.5</sub> monitoring network deployed in Denver, the  
185 “Love My Air” network deployed primarily outside the city’s public schools, to evaluate the  
186 transferability of different calibration models in space and time across the network. To do  
187 so, new metrics are proposed to quantify the Love My Air network spatial and temporal  
188 trend uncertainty due to the calibration model applied. Finally, for key LCS network use-  
189 cases such as hotspot detection, tracking high pollution events and evaluating pollution  
190 trends at a high temporal resolution, the sensitivity of the results to the choice of  
191 calibration model is evaluated. The methodologies and metrics proposed in this paper can  
192 be applied to other low-cost sensor networks, with the understanding that the actual  
193 results will vary with study region.

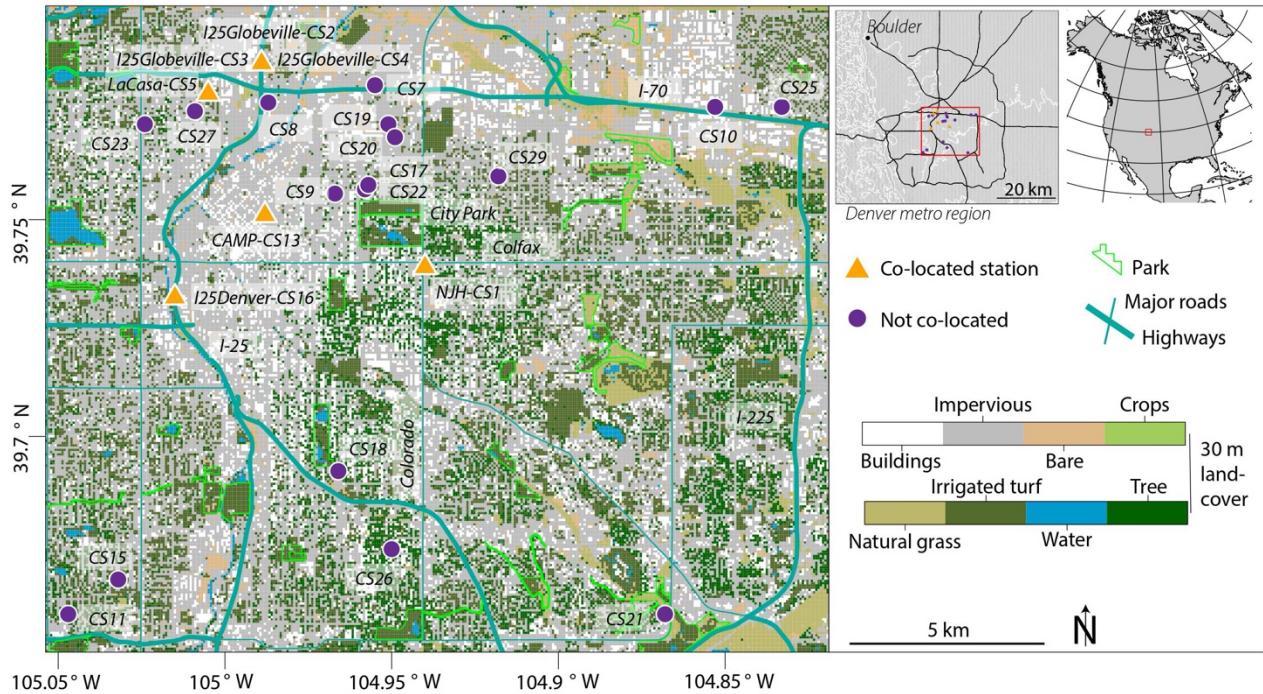
## 194 **2 Data and Methods**

### 195 **2.1 Data Sources**

196 Between Jan 1 and Sep 30, 2021, Denver’s Love My Air sensor network collected minute-  
197 level data from 24 low-cost sensors deployed across the city outside of public schools and  
198 at 5 federal equivalent method (FEM) reference monitor locations (**Figure 1**). The Love  
199 My Air sensors are Canary-S models equipped with a Plantower 5003, made by Lunar  
200 Outpost Inc. The Canary-S sensors detect PM<sub>2.5</sub>, T, and RH, and upload minute-  
201 resolution measurements to an online platform via cellular data network.  
202

203 We found that RH and T reported by the Love My Air sensors were well correlated with  
204 that reported by the reference monitoring stations. We used the Love My Air LCS T and

205 RH measurements in our calibration models as they most closely represent the conditions  
 206 experienced by the sensors.



207  
 208 **Figure 1:** Locations of all 24 Love My Air sensors. Sensors displayed with an orange  
 209 triangle indicate that they were co-located with a reference monitor. The labels of the co-  
 210 located sensors include the name of the reference monitor with which they were co-  
 211 located after a hyphen.

### 212 2.1.1 Data cleaning protocol for measurements from the Love My Air network

213 A summary of the data cleaning and data preparation steps carried out on the Love My  
 214 Air data from the entire network are listed below:

- 215
- 216 1) Removed data for time-steps where key variables:  $PM_{2.5}$ , T and RH measurements  
 217 were missing
- 218 2) Removed unrealistic RH and T values ( $RH < 0$  and  $T \leq -30^{\circ}C$ )
- 219 3) Removed  $PM_{2.5}$  values above  $1,500 \mu g/m^3$  (outside the operational range of the  
 220 Plantower sensors used) from the Canary-S sensors (Considine et al., 2021)
- 221 4) We were left with 8,809,340 minute-level measurements and then calculated  
 222 hourly-average  $PM_{2.5}$ , T, and RH measurements for each sensor. We had a total of  
 223 147,101 hourly-averaged measurements
- 224 5) From inspection, one of the monitors, CS13, worked intermittently in Jan and Feb,  
 225 before resuming continuous measurement in March (**Figure S1 in Supplementary**  
 226 **Information**). When CS13 worked intermittently, large spikes in the measurements  
 227 were observed, likely due to power surges. We thus retained measurements taken  
 228 after March 1, 2021 for this monitor. The total number of hourly measurements was  
 229 thus reduced to 146,583.

230  
231 Love My Air sensors (indicated by Sensor ID) were co-located with FEM reference  
232 monitors from which we obtained high quality hourly  $PM_{2.5}$  measurements at (**Table 1**):  
233 1) La Casa (Sensor ID: CS5)  
234 2) CAMP (Sensor ID: CS13)  
235 3) I25 Globeville (Sensor ID: CS2, CS3, CS4)  
236 4) I25 Denver (Sensor ID: CS16)  
237 5) NJH (Sensor ID: CS1) for the entire period of the experiment

### 238 **2.1.2 Data preparation steps for preparing a training dataset used to develop** 239 **the various calibration models**

240 A summary of the data preparation steps for preparing a training dataset used to develop  
241 the various calibration models are described below:  
242

- 243 1) We joined hourly averages from each of the seven co-located Love My Air  
244 monitors with the corresponding FEM monitor. We had a total of 35,593 co-located  
245 hourly measurements for which we had data for both the Love My Air sensor and  
246 the corresponding reference monitor.  
247 **Figure S2** displays time-series plots of  $PM_{2.5}$  from all co-located Love My Air  
248 sensors. **Figure S3** displays time-series plots of  $PM_{2.5}$  from the corresponding  
249 reference monitors.
- 250 2) The three Love My Air sensors co-located at the I25 Globeville sites (CS2, CS3,  
251 CS4) agreed well with each other (correlation = 0.98) (**Figures S4** and **Figure S5**).  
252 To ensure that our co-located dataset was well balanced across sites, we only  
253 retained measurements from CS2 at the I25 Globeville site. We were left with a  
254 total of 27,338 co-located hourly measurements that we used to develop a  
255 calibration model. **Figure S6** displays the time-series plots of  $PM_{2.5}$  from all other  
256 Love My Air sensors in the network.  
257

258 Reference monitors at La Casa, CAMP, I25 Globeville and I25 Denver, also reported  
259 minute-level  $PM_{2.5}$  concentrations between April 23 11:16 and Sep 30, 22:49. We also  
260 joined minute-level Love My Air  $PM_{2.5}$  concentrations with minute-level reference data at  
261 these sites. We had a total of 1,062,141 co-located minute-level measurements during  
262 this time period. As with the hourly-averaged data, we only retained data from one of the  
263 Love My Air sensors at the I25 Globeville site and were thus left with 815,608 minute-level  
264 measurements from one LCS at each of the four co-location sites.  
265

266 **Table S1** has information on the minute-level co-located measurements. The data at the  
267 minute-level displays more variation and peaks in  $PM_{2.5}$  concentrations than the hourly-  
268 averaged measurements (**Figure S7**), likely due to the impact of passing sources. It is  
269 also important to mention that minute-level reference data may have some additional  
270 uncertainties introduced due to instrument error given the finer time resolution. We will

271 use the minute-level data in supplementary analyses, only. Thus, unless explicitly  
 272 referenced, we will be reporting results from hourly-averaged measurements.

273 **2.1.3 Deriving additional covariates**

274 We derived dew-point (D) from T and RH reported by the Love My Air sensors using the  
 275 *weathermetrics* package in the programming language R (Anderson and Peng, 2012), as  
 276 D has been shown to be a good proxy of particle hygroscopic growth in previous research  
 277 (Barkjohn et al., 2021; Clements et al., 2017; Malings et al., 2020). Some previous work  
 278 has also used a nonlinear correction for RH in the form of  $RH^2/(1-RH)$ , that we also  
 279 calculated for this study (Barkjohn et al., 2021).

280  
 281 We extracted hour, weekend, and month variables from the Canary-S sensors and  
 282 converted hour and month into cyclic values to capture periodicities in the data by taking  
 283 the cosine and sine of  $hour * 2\pi/24$  and  $month * 2\pi/12$ , which we designate as *cos\_time*,  
 284 *sin\_time*, *cos\_month* and *sin\_month*, respectively. Sinusoidal corrections for seasonality  
 285 have been shown to improve accuracy of PM<sub>2.5</sub> measurements in machine learning  
 286 models (Considine et al., 2021).

287  
 288 **Table 1: Site location of each Love My Air sensor, as well as summary statistics of hourly**  
 289 **measurements from each sensor**

Sensor ID	Co-location Information	Latitude	Longitude	Hours operational	PM <sub>2.5</sub> (µg/m <sup>3</sup> )			Temperature (°C)	RH (%)	Dewpoint (°C)
					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
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

## 290 2.2 Defining the Calibration Models Used

291 The goal of the calibration model is to predict, as accurately as possible, the ‘true’ PM<sub>2.5</sub>  
292 concentrations given the concentrations reported by the Love My Air sensors. At the co-  
293 located sites, the FEM PM<sub>2.5</sub> measurements, which we take to be the “true” PM<sub>2.5</sub>  
294 concentrations, are the dependent variable in the models.

295  
296 We evaluated 21 increasingly complex models that included T, RH, D as well as metrics  
297 that captured the time-varying patterns of PM<sub>2.5</sub> to correct the Love My Air PM<sub>2.5</sub>  
298 measurements (**Tables 2** and **3**).

299  
300 Sixteen models were multivariate regression models that were used in a recent paper  
301 (Barkjohn et al., 2021) to calibrate another network of low-cost sensors: the PurpleAir,  
302 that rely on the same PM<sub>2.5</sub> sensor (Plantower) as the Canary-S sensors in the current  
303 study. As T, RH, and D are not independent (**Figure S8**), the 16 linear regression models  
304 include adding the meteorological conditions considered as interaction terms, instead of  
305 additive terms. The remaining five calibration models relied on machine learning  
306 techniques.

307  
308 Machine learning models can capture more complex nonlinear effects (for instance,  
309 unknown relationships between additional spatial and temporal variables). We opted to  
310 use the following machine learning techniques: Random Forest (RF), Neural Network  
311 (NN), Gradient Boosting (GB), SuperLearner (SL) that have been widely used in  
312 calibrating LCS. A description of each technique is described in detail in **section S1** in  
313 *Supplementary Information*. All machine learning models were run using the *caret*  
314 package in R (Kuhn, 2015).

315

316 We used both Leave-One-Site-Out (LOSO) (**Table 2**) and Leave-Out-By-Date, where we  
317 left out a 3-weeks period of data at a time at all sites (LOBD) (**Table 3**) cross-validation  
318 (CV) methods to avoid overfitting in the machine learning models. For more details on the  
319 cross-validation methods used to avoid overfitting in the machine learning models refer to  
320 **section S2** in *Supplementary Information*.

### 321 **2.2.1 Corrections generated using different co-location time periods (long-** 322 **term, on-the-fly, short-term)**

323 As described earlier, co-location studies in the LCS literature have been conducted over  
324 different time periods. Some studies co-locate one or more LCS for brief periods of time  
325 before or after an experiment, whereas others co-locate a few LCS for the entire duration  
326 of the experiment. These studies apply calibration models generated using the co-located  
327 data to measurements made by the entire network over the entire duration of the  
328 experiment. We attempt to replicate these study designs in our experiment to evaluate the  
329 transferability of calibration models across time by generating four different corrections:  
330

331 (C1) *Entire data set correction*: The 21 calibration models were developed using data at  
332 all co-location sites for the entire period of co-location.

333 (C2) *On the fly correction*: The 21 calibration models to correct a measurement during a  
334 given week were developed using data across all co-located sites for the same week of  
335 the measurement.

336 (C3) *2-week winter correction*: The 21 calibration models were developed using co-  
337 located data collected for a brief period (2 weeks) at the beginning of the study (Jan 1 -  
338 Jan 14, 2021). They were then applied to measurements from the network during the rest  
339 of the period of operation.

340 (C4) *2-week winter + 2-week spring*: The 21 calibration models were developed using co-  
341 located data collected for two 2-week periods in different seasons (Jan 1 - Jan 14, 2021  
342 and May 1 - May 14, 2021). They were then applied to measurements from the network  
343 during the rest of the period of operation.

344  
345 Although models developed using co-located data over the entire time period (C1) tend to  
346 be more accurate over the entire spatiotemporal data set, it is inefficient to re-run large  
347 models frequently (incorporating new data). On-the-fly corrections (such as C2) can help  
348 characterize short-term variation in air pollution and sensor characteristics. The duration  
349 of calibration is a key question that remains unanswered (Liang, 2021). We opted to test  
350 corrections C3 and C4 as many low-cost sensor networks rely on developing calibration  
351 models based on relatively short co-location periods (deSouza et al., 2020b; West et al.,  
352 2020; Singh et al., 2021). Each of the 21 calibration models considered was tested under  
353 four potential correction schemes (C1, C2, C3 and C4).

354  
355 For C1, the five machine-learning models were trained using two CV approaches: LOSO  
356 and LOBD, separately. For C2, C3 and C4 only LOSO was conducted, as model

357 application is already being performed on a different time period from the training (for  
358 more details refer to **section S2**).

359  
360 Overall, we test 89 calibration models (21 (C1, CV=LOSO) + 5 (C1, CV=LOBD) + 21 × 3  
361 (C2, C3, C4) = 89) listed in **Tables 2** and **3**.

## 362 **2.3 Evaluating the calibration models developed under the four** 363 **different correction schemes**

364 Uncorrected Love My Air measurements tend to be biased upwards from the  
365 corresponding reference PM<sub>2.5</sub> levels by an average of ~12% (**Figure S9**). We first  
366 evaluate:

- 367 1) Were meteorological conditions at the co-location sites representative of network  
368 operating conditions?
- 369 2) How well do different calibration models perform when using the traditional method  
370 of model evaluation at co-location sites, during the period of co-location?

371  
372 We then evaluate transferability of the calibration models in time and space by evaluating:

- 373 1) How well do calibration models developed during short-term co-locations  
374 (corrections: C3 and C4) perform when transferred to long-term network  
375 measurements?
- 376 2) How well do calibration models developed at a small number of co-locations sites  
377 transfer in space to other sites, even after appropriate cross-validation to prevent  
378 overfitting?
- 379 3) Different metrics to quantify the uncertainty in spatial and temporal trends in PM<sub>2.5</sub>  
380 reported by the LCS network to the calibration model applied.

381  
382 Finally, we evaluate the impact of the choice of calibration model on key LCS network  
383 use-cases, such as hotspot detection, or detection of the most-polluted site. In  
384 supplementary analyses, we also evaluate how much the calibration model impacts the  
385 following additional use-cases:

- 386 1) LCS are increasingly used to evaluate pollution trends on increasingly short  
387 timescales. We evaluated how well calibration models developed using hourly  
388 aggregated data to minute-level LCS measurements
- 389 2) LCS have been deployed to track smoke from fires. We evaluate how well different  
390 calibration models perform at high PM<sub>2.5</sub> concentrations.

### 391 **2.3.1 Evaluating the representativeness of meteorological conditions at the** 392 **co-location sites of the entire network**

393 LCS measurements are impacted by T and RH. We thus, first evaluated if meteorological  
394 conditions (T and RH) at the co-location sites during time-periods used to construct the  
395 calibration models were representative of conditions of operation for the rest of the  
396 network by comparing distributions of these parameters across sites (**Figure 2**).

397 **2.3.2 Traditional Evaluation of the different Calibration Models**

398 We evaluated the performance of the calibration models for the time period of co-location  
399 in our sample using: R (Pearson correlation coefficient), and RMSE (**Tables 2** and **3**).

400 **2.3.3 Evaluating transferability of short-term calibrations developed to the**  
401 **entire period of operation of the network**

402 We evaluated calibration models using corrections C3 and C4 only for the time-period  
403 over which the calibration models were developed, which was Jan 1 - Jan 14, 2021, for  
404 C3 and Jan 1 - Jan 14, 2021, and May 1 - May 14, 2021, for C4 (**Table S2**) and compared  
405 the performance with applying these models to the entire time period of the network  
406 (**Table 2**).

407 **2.3.4 Evaluating whether the calibration models are overfitted to the co-**  
408 **location sites even after appropriate cross-validation**

409 To evaluate how transferable the calibration technique developed at the co-located sites  
410 was to the rest of the network, even after conducting LOSO CV, we left out each of the  
411 five co-located sites in turn and using data from the remaining sites ran the models  
412 proposed in **Tables 2** and **3**. We then applied the models generated to the left-out site.  
413 We report the distribution of RMSE from each calibration model considered at the left-out  
414 sites using box-plots (**Figure 3**). For correction C1, we also left out a three-week period of  
415 data at a time and generated the calibration models based on the data from the remaining  
416 time periods at each site. For the machine learning models (Models 17 – 21), we used CV  
417 = LOBD. We plotted the distribution of RMSE from each model considered for the left-out  
418 three week period (**Figure 3**).

419  
420 We statistically compared the errors in predictions on each test dataset with errors in  
421 predictions from using all sites in our main analysis. Such an approach is useful to  
422 understand how well the proposed correction can transfer to other areas in the Denver  
423 region. To compare statistical differences between errors, we used t-tests if the  
424 distribution of errors were normally distributed (as determined by a Shapiro–Wilk test),  
425 and Wilcoxon signed rank tests, if not, using a significance value of 0.05 (**Section 3.1.4**).

426  
427 We have only five co-location sites in the network. Although evaluating the transferability  
428 among these sites is useful, as we know the true  $PM_{2.5}$  concentrations at these sites, we  
429 also evaluated the transferability of these models in the larger network by predicting  $PM_{2.5}$   
430 concentrations using the models proposed in **Tables 2** and **3** at each of the 24 sites in the  
431 Love My Air network. For each site, we display time series plots of corrected  $PM_{2.5}$   
432 measurements in order to visually compare the ensemble of corrected values at each site  
433 (**Figure 3**).

434 **2.3.5 Evaluating sensitivity of the spatial and temporal trends of the low-cost**  
 435 **sensor network to the method of calibration**

436 We evaluate the spatial and temporal trends in the PM<sub>2.5</sub> concentrations corrected using  
 437 the 89 different calibration models using similar methods to that described in (Jin et al.,  
 438 2019; deSouza et al., 2022) by calculating:

439  
 440 (1) The spatial root mean square difference (RMSD) (**Figure 5**) between any two  
 441 corrected exposures at the same site:  $SRMSD_{h,d} = \sqrt{\frac{1}{N} \sum_{i=1}^N (Conc_{hi} - Conc_{di})^2}$ ,

442 where  $Conc_{hi}$  and  $Conc_{di}$  are Jan 1- Sep 30, 2021 averaged PM<sub>2.5</sub> concentrations  
 443 estimated from correction  $h$  and  $d$  for site  $i$ .  $N$  is the total number of sites.

444 (2) The temporal RMSD (**Figure 6**) between pairs of exposures:  $TRMSD_{h,d} =$   
 445  $\sqrt{\frac{1}{M} \sum_{t=1}^M (Conc_{ht} - Conc_{dt})^2}$ , where  $Conc_{ht}$  and  $Conc_{dt}$  are hourly corrected PM<sub>2.5</sub>  
 446 concentrations averaged over all operational Love My Air sites estimated from  
 447 correction  $h$  and  $d$  for time  $t$ .  $M$  is the total number of hours of operation of the  
 448 network.

449 (3) The spatial Pearson correlation coefficient (**Figure 7**):  $R_S =$   
 450  $\frac{\sum_{i=1}^N (Conc_{hi} - \overline{Conc_h})(Conc_{di} - \overline{Conc_d})}{\sqrt{\sum_{i=1}^N (Conc_{hi} - \overline{Conc_h})^2 \sum_{i=1}^N (Conc_{di} - \overline{Conc_d})^2}}$ , where  $\overline{Conc_h}$  and  $\overline{Conc_d}$  are the average

451 (across all sites and times) corrected PM<sub>2.5</sub> concentrations estimated from  
 452 corrections  $h$  and  $d$  respectively.

453 (4) The temporal Pearson correlation coefficient (**Figure 8**):  $R_T =$   
 454  $\frac{\sum_{t=1}^M (Conc_{ht} - \overline{Conc_h})(Conc_{dt} - \overline{Conc_d})}{\sqrt{\sum_{t=1}^M (Conc_{ht} - \overline{Conc_h})^2 \sum_{t=1}^M (Conc_{dt} - \overline{Conc_d})^2}}$

455  
 456 We characterized the uncertainty in the ‘corrected’ PM<sub>2.5</sub> estimates at each site across the  
 457 different models using two metrics: a normalized range (NR) (**Figure 9a**) and uncertainty,  
 458 calculated from the 95% confidence interval (CI) assuming a t-statistical distribution  
 459 (**Figure 9b**). NR for a given site represents the spread of PM<sub>2.5</sub> across the different  
 460 correction approaches.

461 (5)  $NR = \frac{1}{M} \sum_{t=1}^M \frac{\max_{k \in K} C_{kt} - \min_{k \in K} C_{kt}}{\bar{C}_t}$

462  $C_{kt}$  is the PM<sub>2.5</sub> concentration at hour  $t$  from the  $k$ th model from the ensemble of  $K$  (which  
 463 in this case is 89) correction approaches.  $\bar{C}_t$  represents the ensemble mean across the  $K$   
 464 different products at hour  $t$ .  $M$  is the total number of hours in our sample for which we  
 465 have PM<sub>2.5</sub> data for the site under consideration.

466  
 467 For our sample ( $K = 89$ ), we assume the variations in PM<sub>2.5</sub> across multiple models  
 468 follows the Student-t distribution with the mean being the ensemble average. The  
 469 confidence interval ( $CI$ ) for the ensemble mean at a given time  $t$  is:

470

471 
$$(6) CI_t = \bar{C}_t + t^* \frac{SD_t}{\sqrt{K}}$$

472 Where  $\bar{C}_t$  represents the ensemble mean at time  $t$ ;  $t^*$  is the upper  $\frac{(1-CI)}{2}$  critical value for  
 473 the t-distribution with  $K-1$  degrees of freedom. For  $K=89$ ,  $t^*$  for the 95% double tailed  
 474 confidence interval is 1.99.  $SD_t$  is the sample standard deviation at time  $t$ .

475 
$$(7) SD_t = \sqrt{\frac{\sum_{k=1}^K (C_{k,t} - \bar{C}_t)^2}{K-1}}$$

476  
 477 We define an overall estimate of uncertainty as follows:

478 
$$(8) \text{uncertainty} = \frac{1}{M} \sum_{t=1}^M t^* \frac{SD_t}{\bar{C}_t \sqrt{K}}, \text{ which can also be expressed as}$$

479 
$$(8) \text{uncertainty} = \frac{1}{M} \sum_{t=1}^M \frac{CI_t - \bar{C}_t}{\bar{C}_t}$$

### 480 **2.3.6 Evaluating the sensitivity of hotspot detection across the network of** 481 **sensors to the calibration method**

482 One of the key use-cases of low-cost sensors is hotspot detection. We report the labels of  
 483 sites that are the most polluted using calibrated measurements from the 89 different  
 484 models using hourly data (**Section 3.1.6**) We repeat this process for daily, weekly and  
 485 monthly-averaged calibrated measurements. We ignore missing measurements from the  
 486 network when calculating time averaged values for the different time periods considered.  
 487 We report the mean number of sensors that are ranked ‘most polluted’ across the  
 488 different correction functions for the different averaging periods (**Figure 10**). We do this  
 489 to identify if the choice of the calibration model impacts the hotspot identified by the  
 490 network (i.e. depending on the calibration model different sites show up as the most  
 491 polluted).

### 492 **2.3.7 Supplementary Analysis: Evaluating transferability of calibration** 493 **models developed in different pollution regimes**

494 We evaluated model performance for true/reference  $PM_{2.5}$  concentrations  $> 30 \mu\text{g}/\text{m}^3$  and  
 495  $\leq 30 \mu\text{g}/\text{m}^3$ , as Nilson et al. (2022) has shown that calibration models can have different  
 496 performances in different pollution regimes. We chose to use  $30 \mu\text{g}/\text{m}^3$  as the threshold,  
 497 as these concentrations account for the greatest differences in health and air pollution  
 498 avoidance behavior impacts (Nilson et al., 2022). Lower concentrations ( $PM_{2.5} \leq 30$   
 499  $\mu\text{g}/\text{m}^3$ ) represent most measurements observed in our network; better performance at  
 500 these levels will ensure better day-to-day functionality of the correction. High  $PM_{2.5}$  ( $> 30$   
 501  $\mu\text{g}/\text{m}^3$ ) concentrations in Denver typically occur during fires. Better performance of the  
 502 calibration models in this regime will ensure that the LCS network can accurately capture  
 503 pollution concentrations under smoky conditions. In order to compare errors observed in  
 504 the two different concentration ranges, in addition to reporting R and RMSE of the  
 505 calibration approaches, we also report the normalized RMSE (normalized by the mean of  
 506 the true concentrations) (**Tables S3 and S4**).

507 **2.3.8 Supplementary Analysis: Evaluating transferability of calibration**  
508 **models developed across different time aggregation intervals**

509 One of the key advantages of LCS is that they report high frequency (time scales shorter  
510 than an hour) measurements of pollution. As reference monitoring stations provide hourly  
511 or daily average pollution values, most often the calibration model is developed using  
512 hourly averaged data and then applied to the unaggregated, high-frequency LCS  
513 measurements. We applied the calibration models described in **Tables 2** and **3** developed  
514 using hourly-averaged co-located measurements on minute-level measurements from the  
515 co-located LCS described in **Table S1**. We evaluated the performance of the corrected  
516 high-frequency measurements against the ‘true’ measurements from the corresponding  
517 reference monitor using the metrics R and RMSE (**Tables S5** and **S6**).

518 **3 Results**

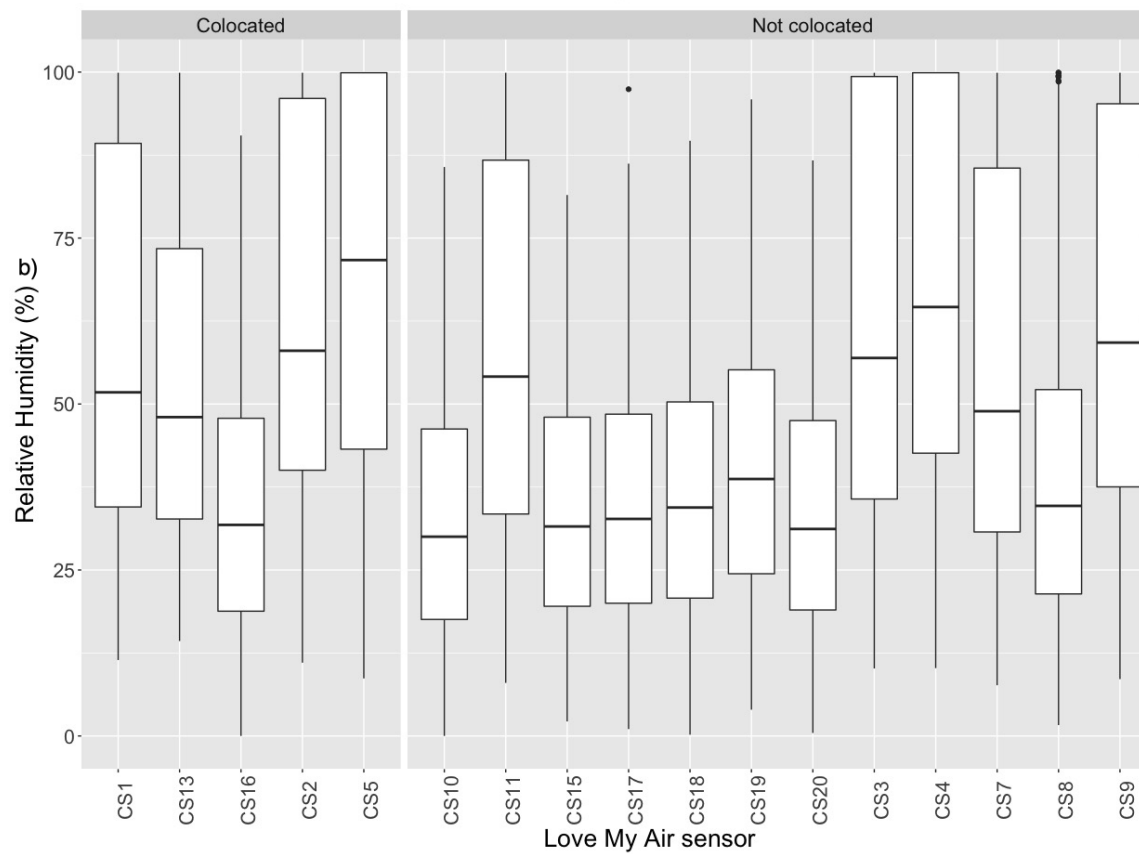
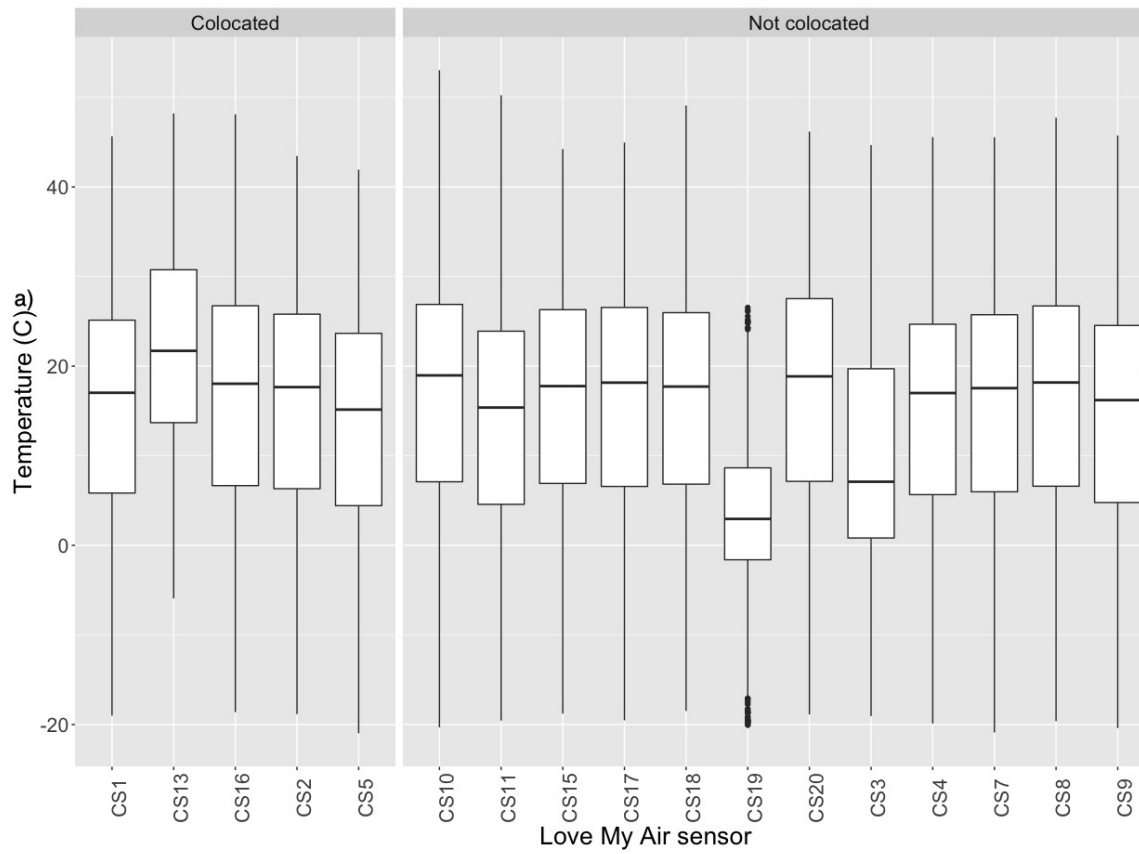
519 **3.1 Evaluating the correction models at the co-location sites**

520 **3.1.1 Evaluating the representativeness of meteorological conditions at**  
521 **the co-location sites of the entire network**

522 Temperature at the co-located sites across the entire period of the experiment (from Jan 1  
523 – Sep 30, 2021) were similar to those at the rest of Love My Air network (**Figure 2a**). The  
524 sensor CS19 is the only one that recorded lower temperatures than those at any of the  
525 other sites. Relative humidity at the co-located sites (three of the four co-located sites  
526 have a median RH close to 50 % or higher) is higher than at the other sites in the network  
527 (7 of the 12 other sites have a median RH < 50%) (**Figure 2b**).

528  
529 We also compared meteorological conditions during the development of corrections C3  
530 (Jan 1 - Jan 14, 2021) and C4 (Jan 1 - Jan 14, 2021, and May 1 - May 14, 2021), to those  
531 measured during the duration of network operation (C3: **Figures S10** and **S11**; C4:  
532 **Figures S12** and **S13**). Unsurprisingly, temperatures at the co-located sites during the  
533 development of C4 were more representative of the network than C3, although they were  
534 on average lower (median temperatures ~ 10 - 17<sup>0</sup>C) than the average temperatures  
535 experienced by the network (median temperatures ~ 5 - 23<sup>0</sup>C). RH values at co-located  
536 sites during C3 and C4 tend to be higher than conditions experienced by some Love My  
537 Air sensors.

538





540 **Figure 2:** (a) Distribution of temperature recorded by each Love My Air sensor, (b)  
 541 Distribution of RH recorded by each Love My Air sensor. The distribution of temperature  
 542 and RH recorded by co-located LCS is shown on the left. The distribution of temperature  
 543 and RH recorded by all LCS not used to construct the calibration models are displayed on  
 544 the right  
 545

### 546 3.1.2 Traditional Evaluation of the different Calibration Models

547 When we evaluated the performance of applying each of the 89 calibration models on all  
 548 co-located data, we found that based on R and RMSE values, the on-the-fly C2 correction  
 549 performed better overall than the C1, C3 and C4 corrections for most calibration model  
 550 forms (**Tables 2 and 3**).

551  
 552 Within corrections C1 and C2, we found that an increase in complexity of model form  
 553 resulted in a decreased RMSE. Overall, Model 21 yielded the best performance (RMSE =  
 554 1.281  $\mu\text{g}/\text{m}^3$  when using the C2 correction, 1.475  $\mu\text{g}/\text{m}^3$  when using the C1 correction with  
 555 a LOSO CV and 1.480  $\mu\text{g}/\text{m}^3$  when using a LOBD correction). In comparison, the simplest  
 556 model yielded an RMSE of 3.421  $\mu\text{g}/\text{m}^3$  for the C1 correction, and 3.008  $\mu\text{g}/\text{m}^3$  when  
 557 using the C2 correction.

558  
 559 For correction C1, using a LOBD CV (**Table 3**) with the machine learning models resulted  
 560 in better performance than using a LOSO CV (**Table 2**), except for Model 21 which is an  
 561 RF model with additional time-of-day and month covariates, for which performance using  
 562 the LOSO CV was marginally better (RMSE: 1.475  $\mu\text{g}/\text{m}^3$  versus 1.480  $\mu\text{g}/\text{m}^3$ ).

563  
 564 **Table 2:** Performance of the calibration models as captured using root mean square error  
 565 (RMSE), and Pearson correlation (R). LOSO CV was used to prevent overfitting in the  
 566 machine learning models. All corrected values were evaluated over the entire time-period  
 567 (Jan 1 - Sep 30, 2021)

ID	Name	Model	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 Jan	C4 Correction developed using measurements from the first two weeks of Jan and the first two weeks in May

			R	RMSE ( $\mu\text{g}/\text{m}^3$ )	R	RMSE ( $\mu\text{g}/\text{m}^3$ )	R	RMSE ( $\mu\text{g}/\text{m}^3$ )	R	RMS E ( $\mu\text{g}/\text{m}^3$ )
<b>Raw Love My Air measurements</b>										
0	Raw		0.927	6.469	-	-	-	-	-	-
<b>Multivariate Regression (LOSO CV)</b>										
1	Linear	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + b$	0.927	3.421	0.944	3.008	0.927	3.486	0.927	3.424
2	+RH	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{RH} \times s_2 + b$	0.929	3.379	0.948	2.904	0.928	3.618	0.929	3.462
3	+T	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{T} \times s_2 + b$	0.928	3.409	0.949	2.896	0.925	3.948	0.928	3.460
4	+D	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{D} \times s_2 + b$	0.928	3.417	0.947	2.934	0.917	3.713	0.925	3.470
5	+RH x T	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{RH} \times s_2 + \text{T} \times s_3 + \text{RH} \times \text{T} \times s_4 + b$	0.934	3.260	0.953	2.782	0.931	3.452	0.933	3.344
6	+RH x D	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{RH} \times s_2 + \text{D} \times s_3 + \text{RH} \times \text{D} \times s_4 + b$	0.930	3.361	0.953	2.785	0.911	3.973	0.929	3.461
7	+D x T	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{D} \times s_2 + \text{T} \times s_3 + \text{D} \times \text{T} \times s_4 + b$	0.928	3.409	0.952	2.798	0.888	5.698	0.921	3.720
8	+RH x T x D	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{RH} \times s_2 + \text{T} \times s_3 + \text{D} \times s_4 + \text{RH} \times \text{T} \times s_5 + \text{RH} \times \text{D} \times s_6 + \text{T} \times \text{D} \times s_7 + \text{RH} \times \text{T} \times \text{D} \times s_8 + b$	0.935	3.246	0.955	2.724	0.779	7.077	0.926	3.625
9	PM x RH	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{RH} \times s_2 + \text{RH} \times \text{PM}_{2.5} \times s_3 + b$	0.930	3.362	0.950	2.854	0.925	3.949	0.925	3.767
10	PM x D	$\text{PM}_{2.5, \text{ corrected}} = \text{PM}_{2.5} \times s_1 + \text{D} \times s_2 + \text{D} \times \text{PM}_{2.5} \times s_3$	0.932	3.324	0.950	2.871	0.883	4.460	0.913	3.777

		+ b								
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	$PM_{2.5, corrected} = PM_{2.5} \times s_1 + RH \times s_2 + T \times s_3 + D \times s_4 + PM_{2.5} \times RH \times s_5 + PM_{2.5} \times T \times s_6 + T \times RH \times s_7 + PM_{2.5} \times D \times s_8 + D \times RH \times s_9 + D \times T \times s_{10} + PM_{2.5} \times RH \times T \times s_{11} + PM_{2.5} \times RH \times D \times s_{12} + PM_{2.5} \times D \times T \times s_{13} + D \times RH \times T \times s_{14} + PM_{2.5} \times RH \times T \times D \times s_{15} + b$	0.940	3.115	0.960	2.557	0.324	32.951	0.765	6.746
<b>Machine Learning (LOSO CV)</b>										
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: $PM_{2.5, corrected} = f(PM_{2.5}, T, RH, D, \cos\_time, \cos\_month, \sin\_month)$  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

568

569 **Table 3:** Performance of the calibration models using the C1 correction as captured using  
570 root mean square error (RMSE), and Pearson correlation (R) LOBD CV was used to  
571 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

572 **3.1.3 Evaluating transferability of short-term calibrations developed to the**  
573 **entire period of operation of the network**

574 We also found that for corrections of short-term calibrations, C3 and C4, more complex  
575 models yielded a better performance (for example the RMSE for Model 16: 2.813  $\mu g/m^3$ ,  
576 RMSE for Model 2: 3.110  $\mu g/m^3$  generated using the C3 correction) when evaluated

577 during the period of co-location, alone (**Table S3**). However, when models generated  
578 using the C3 and C4 corrections were transferred to the entire time period of co-location,  
579 we find that more complex multivariate regression models (Models 13-16) and the  
580 machine learning model (Model 21) that include *cos\_time*, performed significantly worse  
581 than the simpler models (**Table 2**). In some cases, these models performed worse than  
582 the uncorrected measurements. For example, applying Model 16 generated using C3 on  
583 the entire dataset resulted in an RMSE of 32.951  $\mu\text{g}/\text{m}^3$  compared to 6.469  $\mu\text{g}/\text{m}^3$  for the  
584 uncorrected measurements.

585  
586 Including data from another season, spring in addition to winter, in the training sample  
587 (C4), resulted in significantly increased performance of the calibration over the entire  
588 dataset compared to C3 (winter), although it did not result in an improvement in  
589 performance for all models compared to the uncorrected measurements. For example,  
590 Model 16 generated using C4 yielded an RMSE of 6.746  $\mu\text{g}/\text{m}^3$ . Among the multivariate  
591 regression models, we found that models of the same form that corrected for RH instead  
592 of T or D did best. The best performance was observed for models that included the  
593 nonlinear correction for RH (Model 12) or included an  $RH \times T$  term (Model 5) (**Table 2**).

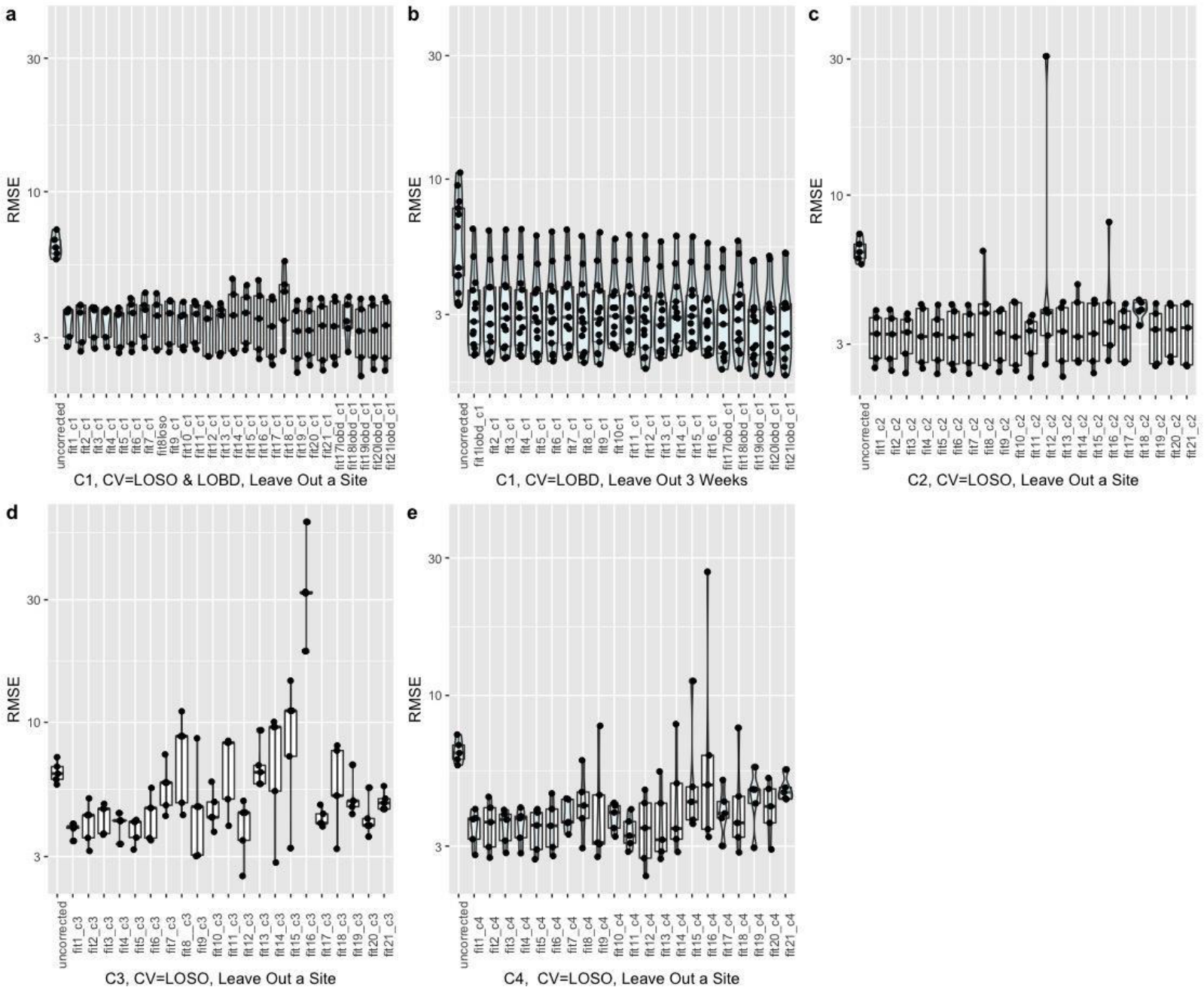
### 594 **3.1.4 Evaluating if the calibration models are overfitted to the co-location** 595 **sites even after appropriate cross-validation**

596 **Figure 3** shows the performance (RMSE) of corrected Love My Air  $\text{PM}_{2.5}$  data by  
597 generating corrections based on the 21 models previously proposed using the C1  
598 correction, CV= LOSO and CV = LOBD for Models 17 - 21, when leaving out a test site  
599 (**Figure 3a**). Also shown is the result using the C1 correction when leaving out a three  
600 week period of data at a time and generating calibration models based on the data from  
601 the remaining time periods across each site, using CV = LOBD for Models 17 – 21, and  
602 applying the models to the remaining three-week period (**Figure 3b**). Finally, **Figures 3c,**  
603 **3d** and **3e** illustrate using the C2, C3 and C4 corrections, respectively, (CV= LOSO for  
604 Models 17 - 21) when leaving out a test site.

605  
606 Large reductions in RMSE are observed when applying simple linear corrections (Models  
607 1 - 4) to the uncorrected data across C1, C2, C3 and C4. Increasing the complexity of the  
608 model does not result in marked changes in correction performance on different test sets  
609 for C1 and C2. Although the performance of the corrected datasets did improve on  
610 average for some of the complex models considered (Model 17, 20, 21 for example, vis-a-  
611 vis simple linear regressions when using the C1 correction) (**Figures 3a, 3b**), this was not  
612 the case for *all* test datasets considered, as evinced by the overlapping distributions of  
613 RMSE performances (e.g., Model 11 using the C2 correction resulted in a worse fit for  
614 one of the test datasets). For C3 and C4, the performance of corrections was worse  
615 across all datasets for the more complex multivariate model formulations (**Figures 3d,**  
616 **3e**), indicating that using uncorrected data is better than using these corrections and  
617 calibration models.

618

619 Wilcoxon tests and t-tests (based on whether Shapiro-Wilk tests revealed that the  
 620 distribution of RMSEs was normal) revealed significant improvements in the distribution of  
 621 RMSEs for all corrected test sets vis-a-vis the uncorrected data. There was no significant  
 622 difference in the distribution of RMSE values from applying C1 and C2 corrections to the  
 623 test sets, across the different models. For corrections C3 and C4, we found significant  
 624 differences in the distribution of RMSEs obtained from running different models on the  
 625 data, implying that the choice of model has a significant impact on transferability of the  
 626 calibration models to other monitors.



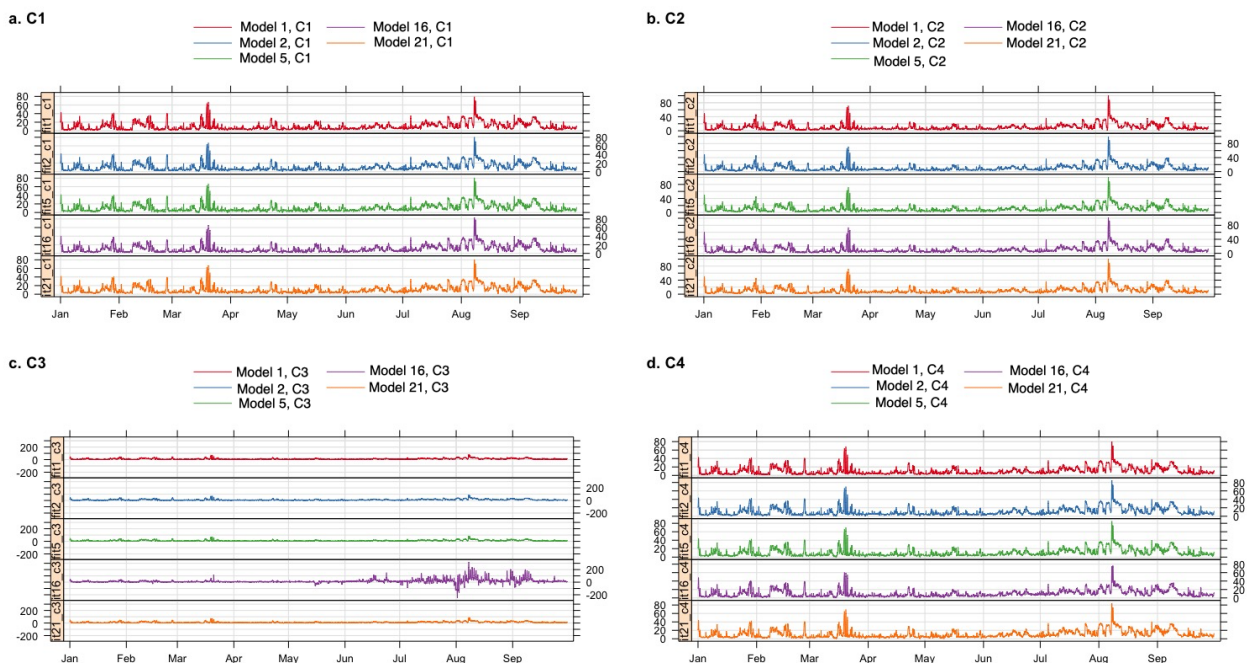
628 **Figure 3:** Performance (RMSE) of corrected Love My Air PM<sub>2.5</sub> data by generating  
 629 corrections based on the 21 models (designated as fit) previously proposed using (a)  
 630 Correction C1 when leaving out a co-location site in turn and then running the generated  
 631 correction on the test site (Note that for machine learning models (Models 17- 21), we

632 performed CV using a LOSO CV as well as a LOBD CV approach), **(b)** Correction C1  
 633 when leaving out 3 week periods of data at a time and generating corrections based on  
 634 the data from the remaining time periods across each site, and evaluating the  
 635 performance of the developed corrections on the held out 3 weeks of data (Note that for  
 636 machine learning models (Models 17- 21), we performed CV using a LOBD CV  
 637 approach), **(c)** Correction C2 when leaving out a co-location site in turn and then running  
 638 the generated correction on the test site, **(d)** Correction C3 when leaving out a co-location  
 639 site in turn and then running the generated correction on the test site, **(e)** Correction C4  
 640 when leaving out a co-location site in turn and then running the generated correction on  
 641 the test site. Each point represents the RMSE for each test dataset permutation. The  
 642 distribution of RMSEs is displayed using box-plots and violin-plots.

643  
 644 The time-series of corrected PM<sub>2.5</sub> values for Models 1, 2, 5, 16, and 21 (RF using  
 645 additional variables) (using CV = LOSO for the machine learning Models 17 and 21) for  
 646 corrections generated using C1, C2, C3 and C4 are displayed in **Figure 4** for Love My Air  
 647 sensor CS1. These subsets of models were chosen as they cover the range of model  
 648 forms considered in this analysis.

649  
 650 From **Figure 4**, we note that although the different corrected values from C1 and C2 track  
 651 each other well, there are small systematic differences between the different corrections.  
 652 Peaks in corrected values using C2 tend to be higher than those using C1. Peaks in  
 653 corrected values using machine learning methods using C1 are higher than those  
 654 generated from multivariate regression models. **Figure 4** also shows marked differences  
 655 in the corrected values from C3 and C4. Specifically Model 16 yields peaks in the data  
 656 that corrections using the other models do not generate. This pattern was consistent  
 657 when applying this suite of corrections to other Love My Air sensors.

658



659

660 **Figure 4:** Time-series of the different  $PM_{2.5}$  corrected values for Models 1, 2, 5, 16 and 21  
661 across corrections (a) C1, (b) C2, (c) C3 and (d) C4 for the Love My Air monitor CS1.  
662 Note that the scales are the same for C1, C2 and C4, but not for C3.

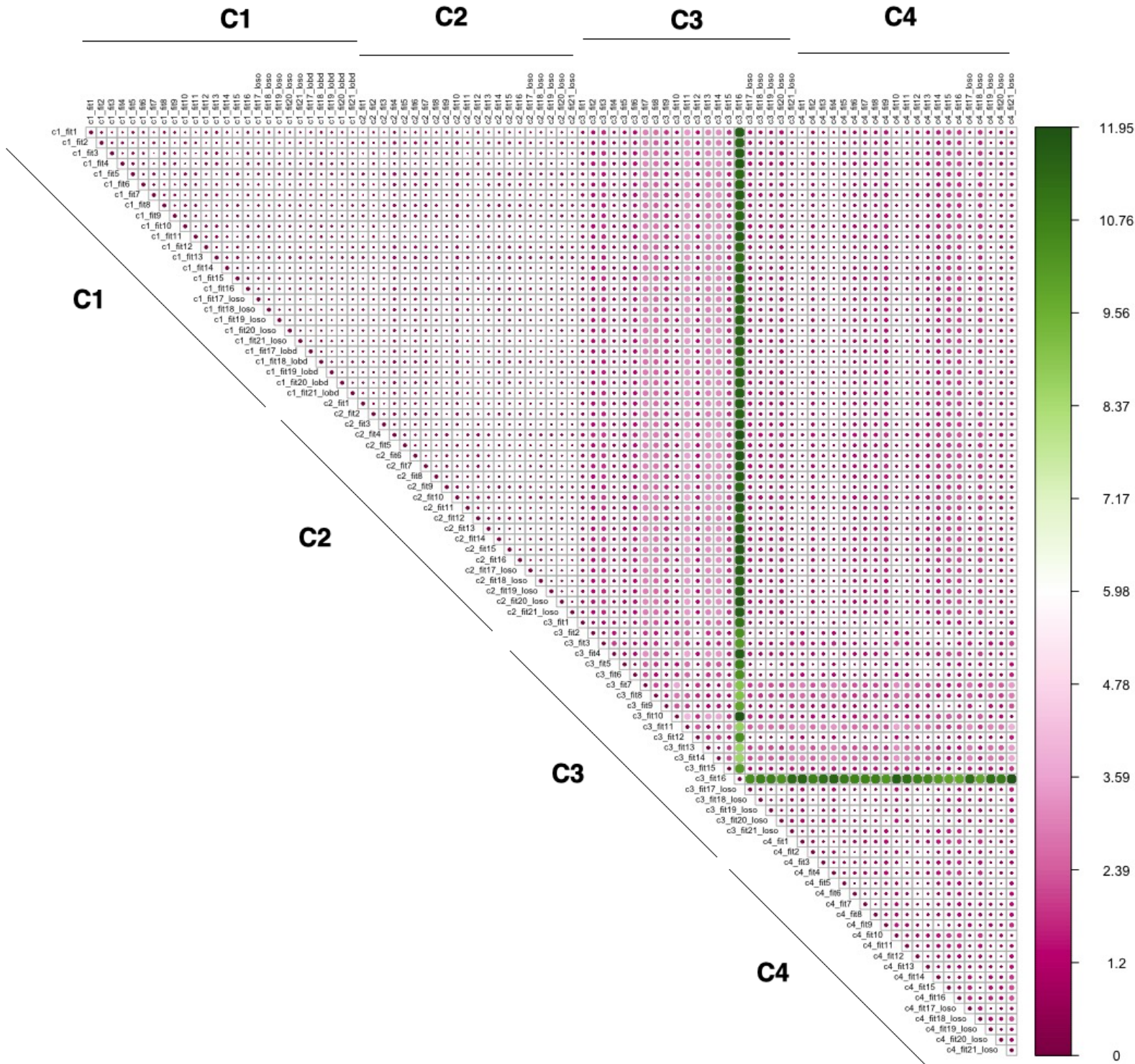
### 663 **3.1.5 Evaluating sensitivity of the spatial and temporal trends of the low-cost** 664 **sensor network to the method of calibration**

665 The spatial and temporal RMSD values between corrected values generated from  
666 applying each of the 89 models using the four different correction approaches across all  
667 monitoring sites in the Love My Air network are displayed **Figures 5 and 6**, respectively.  
668 There is larger temporal variation (max  $32.79 \mu\text{g}/\text{m}^3$ ), in comparison to spatial variations  
669 displayed across corrections (max:  $11.95 \mu\text{g}/\text{m}^3$ ). Model 16 generated using the C3  
670 correction has the greatest spatial and temporal RMSD in comparison with all other  
671 models. Models generated using the C3 and C4 corrections displayed the greatest spatial  
672 and temporal RMSD vis-a-vis C1 and C2.

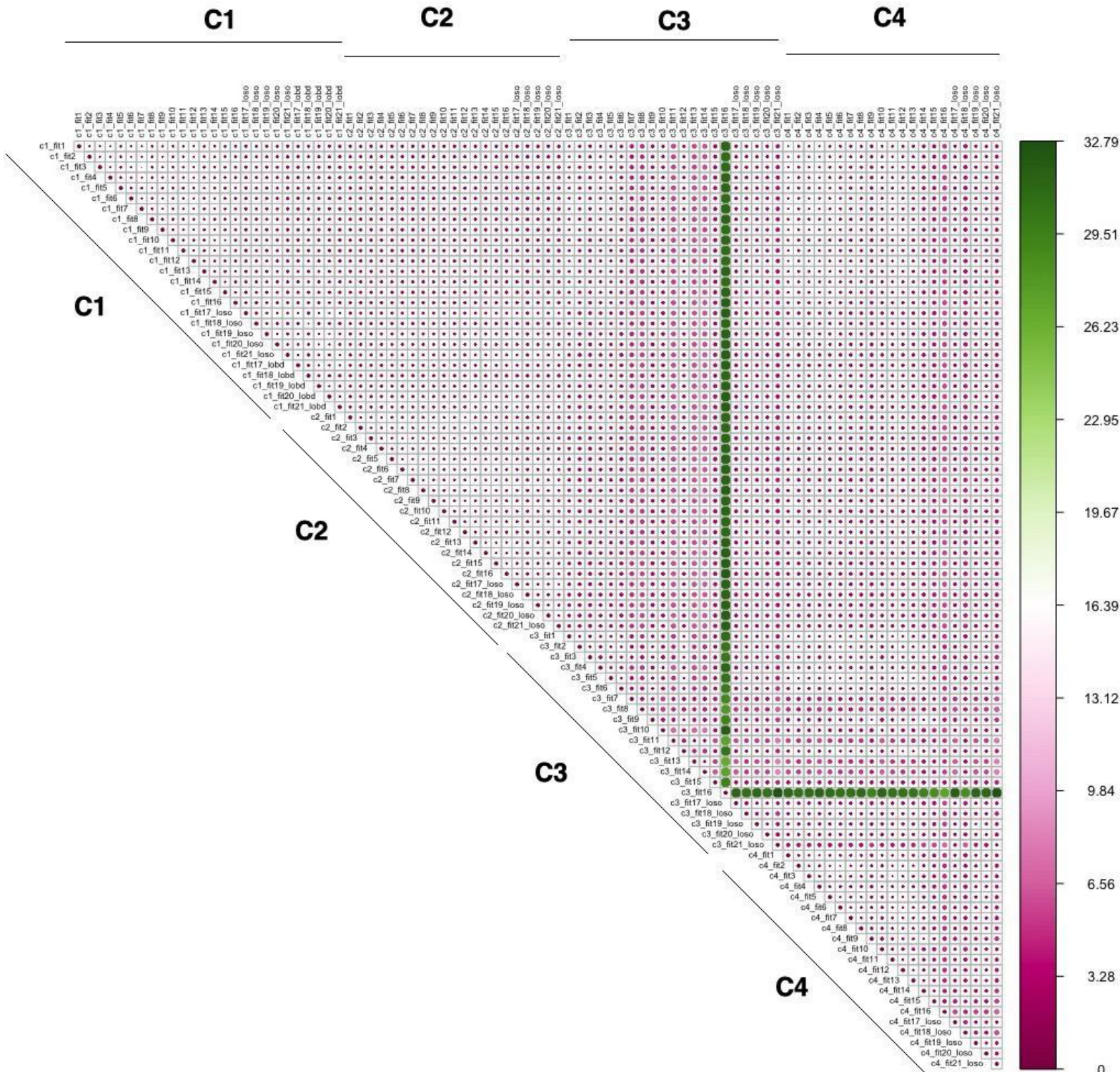
673  
674 **Figures S14- S17** display spatial RMSD values between all models corresponding to  
675 corrections C1-C4, respectively, to allow for a zoomed in view of the impact of the  
676 different model forms for the 4 corrections. Similarly, **Figures S18- S21** display temporal  
677 RMSD values between all models corresponding to corrections C1-C4, respectively.  
678 Across all models the temporal RMSD between models is greater than the spatial RMSD.

679  
680 Spatial and temporal correlation coefficients between corrected measurements generated  
681 from applying all 89 models using the four different correction approaches across the  
682 entire network are displayed in **Figures 7 and 8**, respectively. The spatial correlations are  
683 lower than temporal correlations between corrected measurements.

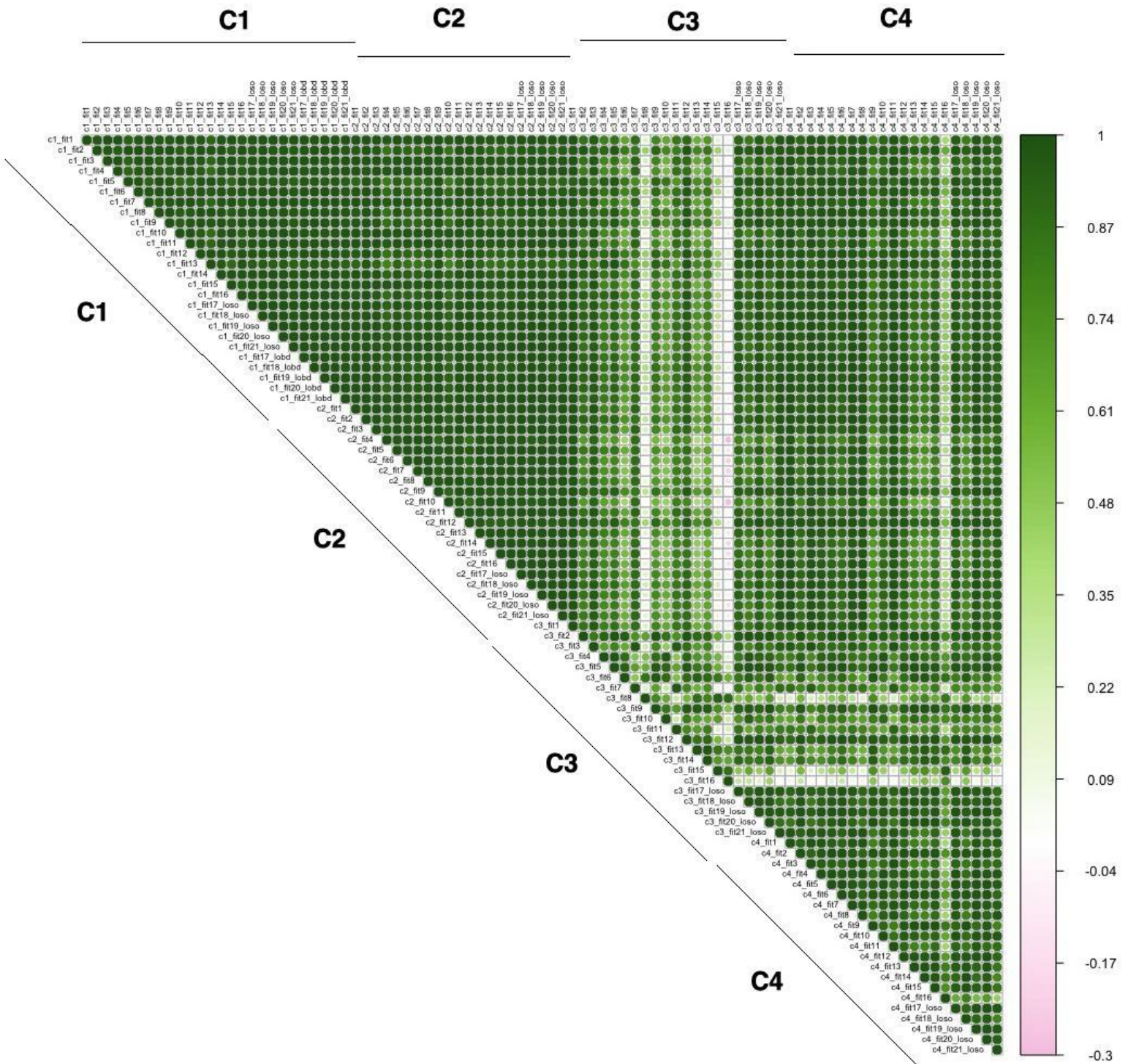




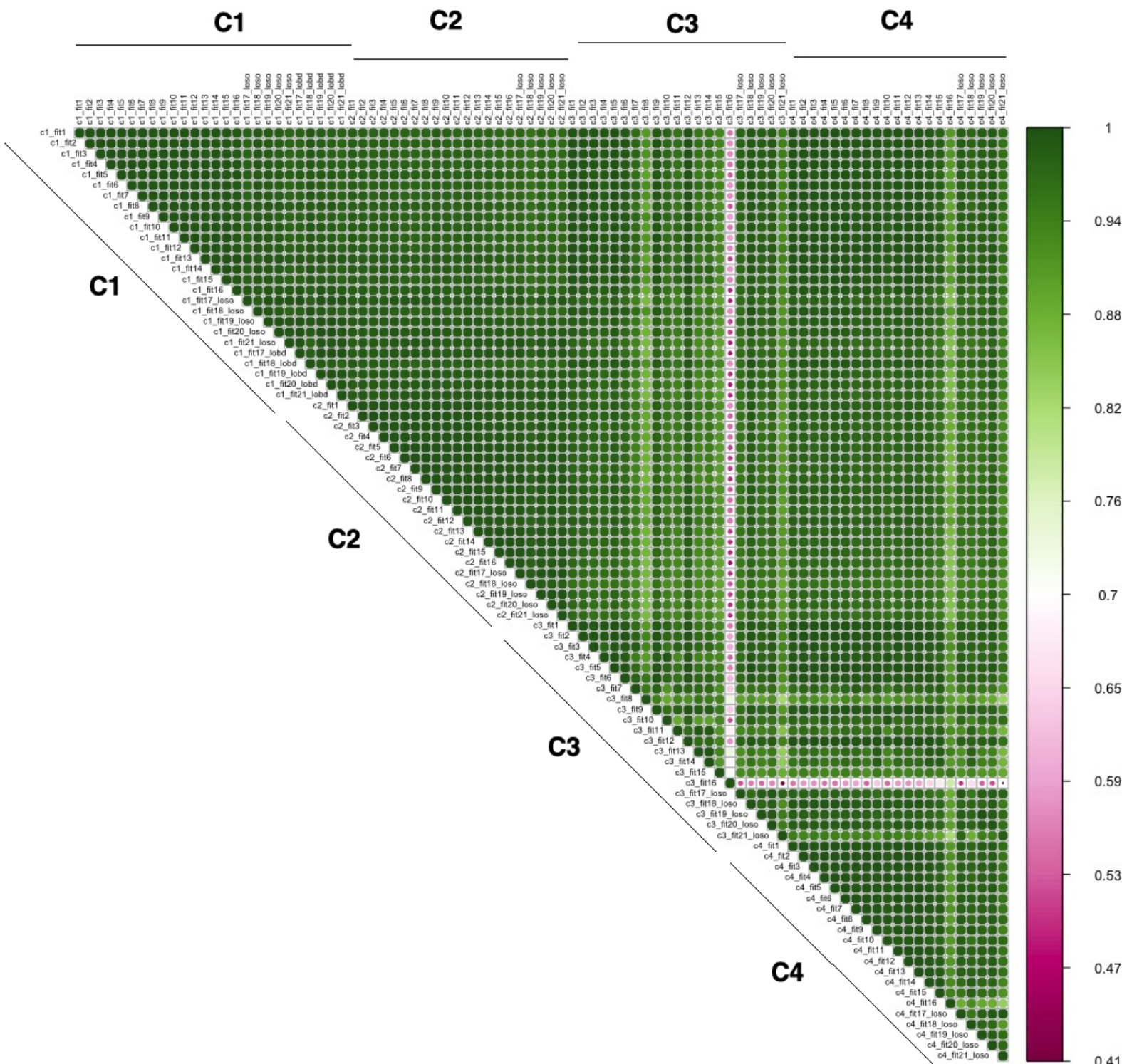
685  
 686 **Figure 5:** Spatial RMSD ( $\mu\text{g}/\text{m}^3$ ) calculated using the method detailed in section 2.3.5  
 687 from applying each of the 89 calibration models using the four different correction  
 688 approaches to all monitoring sites in the Love My Air network.  
 689



691 **Figure 6:** Temporal RMSD ( $\mu\text{g}/\text{m}^3$ ) calculated using the method detailed in section 2.3.5  
 692 from applying each of the 89 calibration models using the four different correction  
 693 approaches to all monitoring sites in the Love My Air network.  
 694



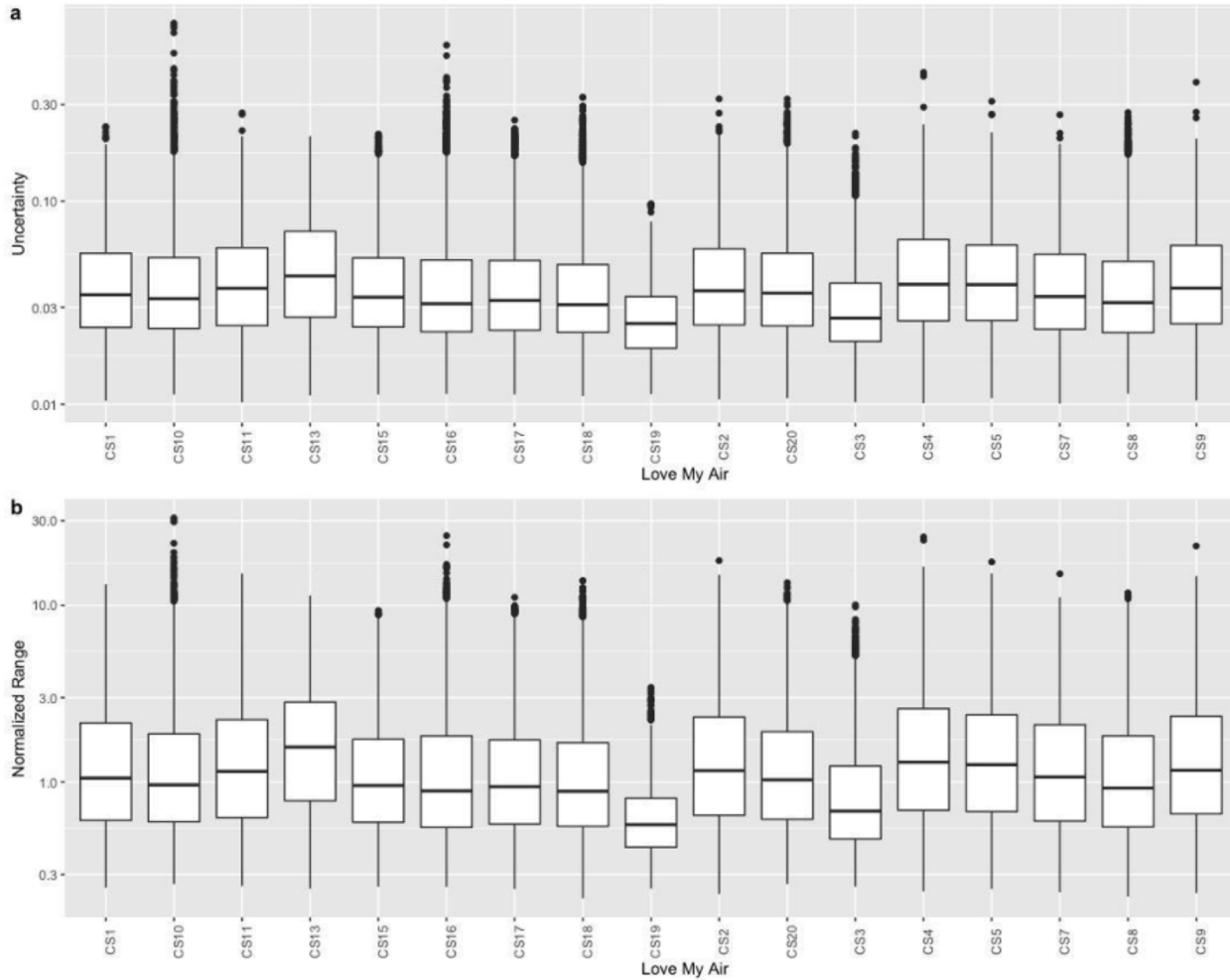
696  
 697 **Figure 7:** Spatial Correlations from applying each of the 89 calibration models using  
 698 corrections C1-C4 to all monitoring sites in the Love My Air network calculated using the  
 699 method described in section 2.3.5.  
 700



702 **Figure 8:** Temporal Correlations from applying each of the 89 calibration models using  
 703 corrections C1-C4 approaches to all monitoring sites in the Love My Air network  
 704 calculated using the method described in section 2.3.5.

705  
 706 The distribution of uncertainty and the NR in hourly-calibrated measurements over the 89  
 707 models by monitor are displayed in **Figure 9**. Overall, there are small differences in  
 708 uncertainties and NR of the calibrated measurements across sites. The average NR and

709 uncertainty across all sites are 1.554 (median: 0.9768) and 0.044 (median: 0.033),  
710 respectively. We note that although the uncertainties in the data are small, the average  
711 normalized range tends to be quite large.  
712



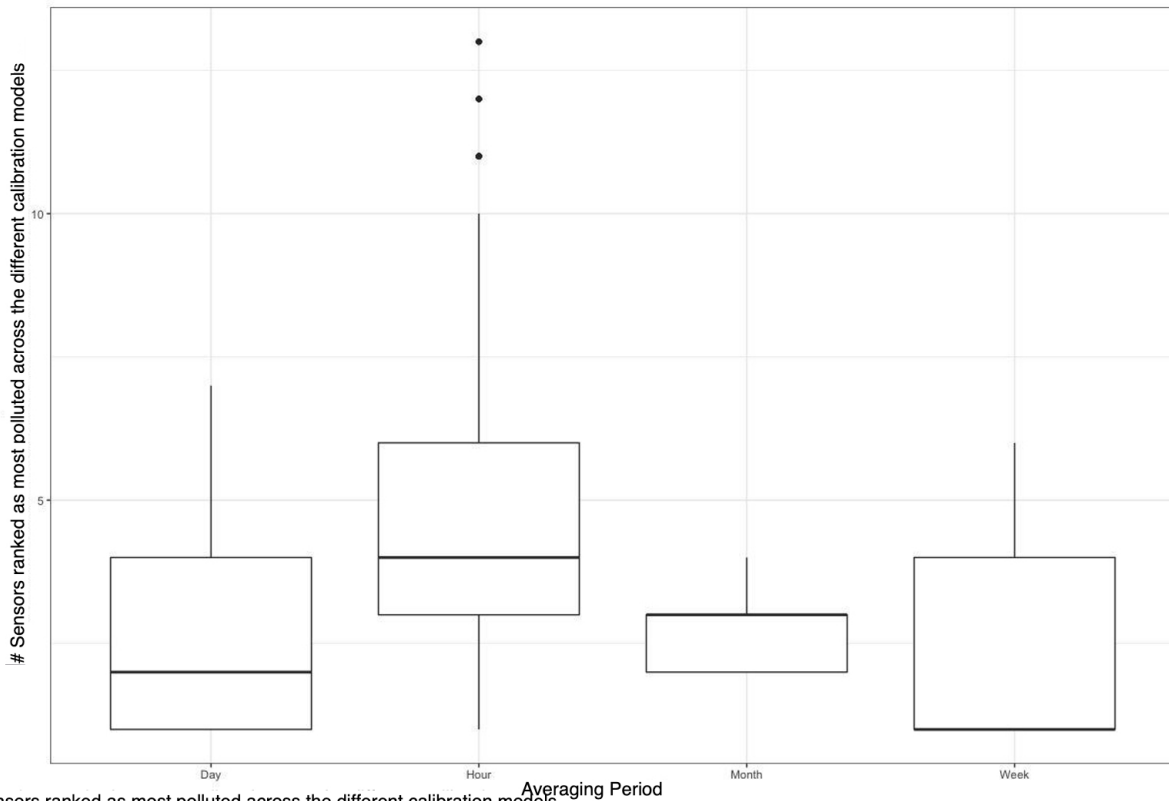
713  
714 **Figure 9:** Distribution of (a) uncertainty and (b) normalized range (NR) in hourly-calibrated  
715 measurements across all 89 calibration models at each site using the methodology  
716 described in Section 2.3.5.

### 717 3.1.6 Evaluating the sensitivity of hotspot detection across the network of 718 sensors to the calibration method

719 Mean (95% CI) PM<sub>2.5</sub> concentrations across the 89 different calibration models listed in  
720 **Tables 1 and 2**) at each Love My Air site for the duration of the experiment (Jan 1 - Sep  
721 30, 2021) are displayed in **Figure S22**. Due to overlap between the different calibrated  
722 measurements across sites, the ranking of sites based on pollutant concentrations is  
723 dependent on the calibration model used.

724  
725 Every hour, we ranked the different monitors for each of the 89 different calibration  
726 models, in order to evaluate how sensitive pollution hotspots were to the calibration model  
727 used. We found that there were on average 4.4 (median = 5) sensors that were ranked

728 most polluted. When this calculation was repeated using daily-averaged calibrated data,  
729 there were on average 2.5 (median = 2) sensors that were ranked the most polluted. The  
730 corresponding value for weekly-calibrated data was 2.4 (median = 1), and for monthly  
731 data was 3 (median = 3) (Figure 9).



732 # Sensors ranked as most polluted across the different calibration models  
733 **Figure 9:** Variation in the number of sites that were ranked as 'most polluted' across the  
734 89 different calibration models for different time-averaging periods displayed using box-  
735 plots

### 736 3.1.7 Supplementary Analysis: Evaluating transferability of calibration 737 models developed in different pollution regimes

738 When we evaluated how well the models performed at high PM<sub>2.5</sub> concentrations (> 30  
739 µg/m<sup>3</sup>) versus lower concentrations (≤ 30 µg/m<sup>3</sup>), we found that multivariate regression  
740 models generated using the C1 correction did not perform well in capturing peaks in PM<sub>2.5</sub>  
741 concentrations (normalized RMSE > 25%) (Tables S3 and S4).

742  
743 Multivariate regression models generated using the C2 correction performed better than  
744 those generated using C1 (normalized RMSE ~ 20 -25 %). Machine learning models  
745 generated using both C1 and C2 corrections captured PM<sub>2.5</sub> peaks well (C1: normalized  
746 RMSE ~ 10 - 25%, C2: normalized RMSE ~ 10 - 20%). Specifically, the C2 RF model  
747 (Model 21) yielded the lowest RMSE values (4.180 µg/m<sup>3</sup>, normalized RMSE: 9.8%), of all  
748 models considered. The performance of models generated using C1 and C2 corrections  
749 in the low-concentration regime was the same as that over the entire dataset. This is  
750 because most measurements made were < 30 µg/m<sup>3</sup>.

751  
752 Models generated using C3 and C4 had the worst performance in both concentration  
753 regimes and yielded poorer agreement with reference measurements than even the  
754 uncorrected measurements. As in the case with the entire dataset, more complex  
755 multivariate regression models and machine learning models generated using C3 and C4  
756 performed worse than more simple models in both PM<sub>2.5</sub> concentration intervals (**Tables**  
757 **S3** and **S4**).

### 758 **3.1.8 Supplementary Analysis: Evaluating transferability of calibration** 759 **models developed across different time aggregation intervals**

760 We then evaluated how well the models generated using C1, C2, C3 and C4 corrections  
761 performed when applied to minute-level LCS data at co-located sites (**Tables S5** and **S6**).  
762 We found that the machine learning models generated using C1 and C2 improved the  
763 performance of the LCS. Model 21 (CV=LOSO) generated using C1 yielded an RMSE of  
764 15.482 µg/m<sup>3</sup> compared to 16.409 µg/m<sup>3</sup> obtained from the uncorrected measurements.

765  
766 The more complex multivariate regression models yielded a significantly worse  
767 performance across all corrections. (Model 16 generated using C1 yielded an RMSE of  
768 41.795 µg/m<sup>3</sup>). As in the case with the hourly-averaged measurements, using correction  
769 C1, LOBD CV instead of LOSO for the machine learning models resulted in better model  
770 performance except for Model 21. Few models generated using C3 and C4 resulted in  
771 improved performance when applied to the minute-level measurements (**Tables S5** and  
772 **S6**).

## 773 **4 Discussion and Conclusions**

774 In our analysis of how transferable the correction models developed at the Love My Air  
775 co-location sites are to the rest of the network, we found that for C1 and C2, more  
776 complex model forms yielded better predictions (higher R, lower RMSE) at the co-located  
777 sites. This is likely because the machine learning models were likely best able to capture  
778 complex, non-linear relationships between the LCS measurements, meteorological  
779 parameters and reference data. Model 21, which included additional covariates intended  
780 to capture periodicities in the data, such as seasonality yielded the best performance,  
781 suggesting that in this study the relationship between LCS measurements and reference  
782 data varies over time. One possible reason for this could be the impact of changing  
783 aerosol composition in time which has been shown to impact the LCS calibration function  
784 (Malings et al., 2020).

785  
786 When examining the short-term, C3 and C4 corrections, we found that although these  
787 corrections appeared to significantly improve LCS measurements during the time period  
788 of model development (**Table S2**), when transferred to the entire time period of operation  
789 they did not perform well (**Table 2**). Many of the models, especially the more complex  
790 multivariate regression models, performed significantly worse than even the uncorrected  
791 measurements. This indicates that calibration models generated during short time

792 periods, even if the time periods correspond to different seasons, may not necessarily  
793 transfer well to other times, likely due to changes in the aerosol composition, and  
794 differences in meteorological conditions, among other potential factors. This indicates the  
795 need for statistical calibration models to be developed over longer time periods that better  
796 capture different LCS operating conditions. For C3 and C4, we did however find models  
797 that relied on nonlinear formulations of RH, that serve as proxies for hygroscopic growth,  
798 yielded the best performance, as compared to more complex models (**Table 2**). This  
799 suggests that physics-based calibrations are potentially an alternative approach,  
800 especially when relying on short co-location periods and need to be explored further.

801  
802 When evaluating how transferable different calibration models were to the rest of the  
803 network, we found that for C1 and C2, more complex models that appeared to perform  
804 well at the co-location sites did not necessarily transfer best to the rest of the network.  
805 Specifically, when we tested these models on a co-located site that was left out when  
806 generating the calibration models, we found that some of the more complex models using  
807 the C2 correction yielded a significantly worse performance at some test sites (**Figure 3**).  
808 If the corrected data were going to be used to make site-specific decisions then such  
809 corrections would lead to important errors. For C3 and C4, we observed a large  
810 distribution of RMSE values across sites. For several of the more complex models  
811 developed using C3 and C4 corrections, the RMSE values at some left-out sites were  
812 larger than observed for the uncorrected data, suggesting that certain calibration models  
813 could result in even more error-prone data than using uncorrected measurements. As the  
814 meteorological parameters for the duration of the C3 and C4 co-locations are not  
815 representative of overall operating conditions of the network, it is likely that the more  
816 complex models were overfit to conditions during the co-location, leading to them not  
817 performing well over the network operations.

818  
819 For C1 and C2, we found that there were no significant differences in the distribution of  
820 the performance metric RMSE of corrected measurements from simpler models in  
821 comparison to those derived from more complex corrections at test sites (**Figure 3**). For  
822 C3 and C4, we found significant differences in the distribution of RMSE across test sites,  
823 which indicates that these models are likely site-specific and not easily transferable to  
824 other sites in the network. This suggests that less complex models might be preferred  
825 when short-term co-locations are carried out for sensor calibration, especially when  
826 conditions during the short-term co-location are not representative of that of the network.

827  
828 We found that the temporal RMSD (**Figure 6**) was greater than the spatial RMSD (**Figure**  
829 **5**) for the ensemble of corrected measurements developed by applying the 89 different  
830 calibration models to the Love My Air network. One of the reasons this may be the case is  
831 that  $PM_{2.5}$  concentrations across the different Love My Air sites in Denver are highly  
832 correlated (**Figure S5**), indicating that the contribution of local sources to  $PM_{2.5}$   
833 concentrations in the Denver neighborhoods in which Love My Air was deployed is small.  
834 Due to the low variability in  $PM_{2.5}$  concentrations across sites, it makes sense that the



835 variations in the corrected PM<sub>2.5</sub> concentrations will be seen in time rather than space.  
836 The largest pairwise temporal RMSD were all seen between corrections derived from  
837 complex models using the C3 correction.

838  
839 However, we note that the temporal correlation coefficients (**Figure 8**) for all-pairwise  
840 correction models was higher than the corresponding spatial coefficient (**Figure 7**). This  
841 implies that although the corrections generated from all models considered tended to  
842 track each other (except for a few models using C3) some corrected values were biased  
843 low, whereas some were biased high. It's important for future work to be done to  
844 characterize under what conditions these biases occur.

845  
846 Finally, we observed that the uncertainty in PM<sub>2.5</sub> concentrations across the ensemble of  
847 89 calibration models (**Figure 9**) was consistently small for the Love My Air Denver  
848 network. The normalized range in the corrected measurements, on the other hand, was  
849 large; however, the uncertainty (95% CI) in the corrected measurements fall within a  
850 relatively small interval. Thus, deciding which calibration model to pick has important  
851 consequences for decision-makers when using data from this network.

852  
853 Our findings reinforce the idea that evaluating calibration models at all co-location sites  
854 using overall metrics like RMSE should not be seen as the only/best way to determine  
855 how to calibrate a network of LCS. Instead, approaches like the ones we have  
856 demonstrate, and metrics like the ones we have proposed should be used to evaluate  
857 calibration transferability.

858  
859 We found that the detection of the 'most polluted' site in the Love My Air network (an  
860 important use-case of LCS networks) was dependent on the calibration model used on  
861 the network. We also found that for the Love My Air network, the detection of the most  
862 polluted site was sensitive to the duration of time-averaging of the corrected  
863 measurements (**Figure 10**). Hotspot detection was most robust using weekly-averaged  
864 measurements. A possible reason for this is that temporal variation in PM<sub>2.5</sub> in Denver  
865 varied primarily on a weekly-scale, and therefore analysis conducted using weekly-values  
866 resulted in the most robust results. Such an analysis thus provides guidance on the most  
867 useful temporal scale for decision-making related to evaluating hotspots in the Denver  
868 network.

869  
870 In supplementary analyses, when we evaluated the sensitivity of other LCS use-cases to  
871 the calibration model applied such as tracking high pollution concentrations during fire or  
872 smoke-events, we found that different models yielded different performance results in  
873 different pollution regimens. Machine learning models developed using C1, and models  
874 developed using C2 were better than multivariate regression models generated using C1  
875 at capturing peaks in pollution (> 30 µg/m<sup>3</sup>). All models using C3 and C4 yielded poor  
876 performance results in tracking high pollution events (**Tables S3** and **S4**). This is likely  
877 because PM<sub>2.5</sub> concentrations during the C3 and C4 co-location tended to be low. The

878 calibration model developed thus did not transfer well to other concentrations. When  
879 evaluating how well the calibration models developed using hourly-aggregated  
880 measurements translated to high-resolution minute-level data (**Tables S5** and **S6**), we  
881 observed that machine learning models generated using C1 and C2, improved the LCS  
882 measurements. More complex multivariate regression models performed poorly. All C3  
883 and C4 models also performed poorly. This suggests that caution needs to be exercised  
884 when transferring models developed at a particular time scale to another. Note that in this  
885 paper, because pollution concentrations did not show much spatial variation, we focus on  
886 evaluating transferability across time-scales, only.

887

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

890

891 1) *Determine how well calibration adjustments can be transferred to other locations.*

892 Specifically, although we found that in Denver some calibration models performed well at  
893 co-location sites, the models could result in large errors at specific sites that would create  
894 difficulties for site-specific decision making.

895

896 2) *Examine how well calibration adjustments can be transferred to other time periods.* In  
897 this study we found that models developed using the short-term C3 and C4 corrections  
898 were not transferable to other time periods because the conditions during the co-location  
899 were not representative of broader operating conditions in the network.

900

901 3) *Use a variety of approaches to quantify transferability of calibration models in the*  
902 *overall network* (e.g., with spatio-temporal correlations and RMSD). The metrics proposed  
903 in this paper to evaluate model transferability can be used in other networks.

904

905 4) *Investigate how adopting a certain time-scale for averaging measurements could*  
906 *mitigate the uncertainty induced by the calibration process for specific use-cases.*

907 Namely, we found that in the Love My Air network, hotspot identification was more robust  
908 to using daily-averaged data than hourly-averaged data. Our analyses also revealed  
909 which models performed best when needing to transfer the calibration model developed  
910 using hourly-averaged data to higher-resolution data, and which models best captured  
911 peaks in pollution during fire- or smoke- events.

912

913 In this work, the Love My Air network under consideration is located over a fairly small  
914 area in a single city. In this network, for the time period considered, PM<sub>2.5</sub> seems to be  
915 mainly a regional pollutant and the contribution of local sources is small. More work needs  
916 to be done to evaluate model transferability in networks in other settings. Concerns about  
917 model transferability are likely to be even more pressing when thinking about larger  
918 networks that span different cities and should be considered in future research. In this  
919 study, we present a first attempt to demonstrate the importance of considering the

920 transferability of calibration models. In future work, we also aim to explore the physical  
921 factors that drive concerns about transferability to generalize our findings more broadly.

## 922 **Author Contributions**

923 PD conceptualized the study, developed the methodology, carried out the analysis and wrote the  
924 first draft. PD and BC obtained funding for this study. BC produced Figure 1. All authors helped in  
925 refining the methodology and editing the draft.

## 926 **Acknowledgements**

927 PD and BC gratefully acknowledge a CU Denver Presidential Initiative grant that  
928 supported their work. The authors are grateful to the Love My Air team for setting up and  
929 maintaining the Love My Air network. The authors are also grateful to Carl Malings for  
930 useful comments

## 931 **Competing Interests**

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

## 933 **References**

934 Anderson, G. and Peng, R.: weathermetrics: Functions to convert between weather metrics (R  
935 package), 2012.

936

937 [State of Global Air: https://www.stateofglobalair.org/](https://www.stateofglobalair.org/), last access: 18 June 2022.

938

939 Apte, J. S., Messier, K. P., Gani, S., Brauer, M., Kirchstetter, T. W., Lunden, M. M., Marshall, J.  
940 D., Portier, C. J., Vermeulen, R. C. H., and Hamburg, S. P.: High-Resolution Air Pollution Mapping  
941 with Google Street View Cars: Exploiting Big Data, *Environ. Sci. Technol.*, 51, 6999–7008,  
942 <https://doi.org/10.1021/acs.est.7b00891>, 2017.

943

944 Barkjohn, K. K., Gantt, B., and Clements, A. L.: Development and application of a United States-  
945 wide correction for PM<sub>2.5</sub> data collected with the PurpleAir sensor, *Atmospheric Meas. Tech.*, 14,  
946 4617–4637, <https://doi.org/10.5194/amt-14-4617-2021>, 2021.

947

948 Bean, J. K.: Evaluation methods for low-cost particulate matter sensors, *Atmospheric Meas.*  
949 *Tech.*, 14, 7369–7379, <https://doi.org/10.5194/amt-14-7369-2021>, 2021.

950

951 Bi, J., Wildani, A., Chang, H. H., and Liu, Y.: Incorporating Low-Cost Sensor Measurements into  
952 High-Resolution PM<sub>2.5</sub> Modeling at a Large Spatial Scale, *Environ. Sci. Technol.*, 54, 2152–2162,  
953 <https://doi.org/10.1021/acs.est.9b06046>, 2020.

954

955 Brantley, H. L., Hagler, G. S. W., Herndon, S. C., Massoli, P., Bergin, M. H., and Russell, A. G.:  
956 Characterization of Spatial Air Pollution Patterns Near a Large Railyard Area in Atlanta, Georgia,  
957 *Int. J. Environ. Res. Public Health*, 16, 535, <https://doi.org/10.3390/ijerph16040535>, 2019.

958

959 Castell, N., Dauge, F. R., Schneider, P., Vogt, M., Lerner, U., Fishbain, B., Broday, D., and

960 Bartonova, A.: Can commercial low-cost sensor platforms contribute to air quality monitoring and  
961 exposure estimates?, *Environ. Int.*, 99, 293–302, <https://doi.org/10.1016/j.envint.2016.12.007>,  
962 2017.

963

964 Clements, A. L., Griswold, W. G., Rs, A., Johnston, J. E., Herting, M. M., Thorson, J., Collier-  
965 Oxandale, A., and Hannigan, M.: Low-Cost Air Quality Monitoring Tools: From Research to  
966 Practice (A Workshop Summary), *Sensors*, 17, 2478, <https://doi.org/10.3390/s17112478>, 2017.

967

968 Considine, E. M., Reid, C. E., Ogletree, M. R., and Dye, T.: Improving accuracy of air pollution  
969 exposure measurements: Statistical correction of a municipal low-cost airborne particulate matter  
970 sensor network, *Environ. Pollut.*, 268, 115833, <https://doi.org/10.1016/j.envpol.2020.115833>,  
971 2021.

972

973 Crawford, B., Hagan, D.H., Grossman, I., Cole, E., Holland, L., Heald, C.L. and Kroll, J.H., 2021.  
974 Mapping pollution exposure and chemistry during an extreme air quality event (the 2018 Kīlauea  
975 eruption) using a low-cost sensor network. *Proceedings of the National Academy of Sciences*,  
976 118(27), p.e2025540118.

977

978 Crilley, L. R., Shaw, M., Pound, R., Kramer, L. J., Price, R., Young, S., Lewis, A. C., and Pope, F.  
979 D.: Evaluation of a low-cost optical particle counter (Alphasense OPC-N2) for ambient air  
980 monitoring, *Atmospheric Meas. Tech.*, 11, 709–720, <https://doi.org/10.5194/amt-11-709-2018>,  
981 2018.

982

983 deSouza, P. and Kinney, P. L.: On the distribution of low-cost PM 2.5 sensors in the US:  
984 demographic and air quality associations, *J. Expo. Sci. Environ. Epidemiol.*, 31, 514–524,  
985 <https://doi.org/10.1038/s41370-021-00328-2>, 2021.

986

987 deSouza, P., Anjomshoaa, A., Duarte, F., Kahn, R., Kumar, P., and Ratti, C.: Air quality  
988 monitoring using mobile low-cost sensors mounted on trash-trucks: Methods development and  
989 lessons learned, *Sustain. Cities Soc.*, 60, 102239, <https://doi.org/10.1016/j.scs.2020.102239>,  
990 2020a.

991

992 deSouza, P., Lu, R., Kinney, P., and Zheng, S.: Exposures to multiple air pollutants while  
993 commuting: Evidence from Zhengzhou, China, *Atmos. Environ.*, 118168,  
994 <https://doi.org/10.1016/j.atmosenv.2020.118168>, 2020b.

995

996 deSouza, P. N.: Key Concerns and Drivers of Low-Cost Air Quality Sensor Use, *Sustainability*, 14,  
997 584, <https://doi.org/10.3390/su14010584>, 2022.

998

999 deSouza, P. N., Dey, S., Mwenda, K. M., Kim, R., Subramanian, S. V., and Kinney, P. L.: Robust  
1000 relationship between ambient air pollution and infant mortality in India, *Sci. Total Environ.*, 815,  
1001 152755, <https://doi.org/10.1016/j.scitotenv.2021.152755>, 2022.

1002

1003 Giordano, M. R., Malings, C., Pandis, S. N., Presto, A. A., McNeill, V. F., Westervelt, D. M.,  
1004 Beekmann, M., and Subramanian, R.: From low-cost sensors to high-quality data: A summary of  
1005 challenges and best practices for effectively calibrating low-cost particulate matter mass sensors,  
1006 *J. Aerosol Sci.*, 158, 105833, <https://doi.org/10.1016/j.jaerosci.2021.105833>, 2021.

1007  
1008 Hagler, G. S. W., Williams, R., Papapostolou, V., and Polidori, A.: Air Quality Sensors and Data  
1009 Adjustment Algorithms: When Is It No Longer a Measurement?, *Environ. Sci. Technol.*, 52, 5530–  
1010 5531, <https://doi.org/10.1021/acs.est.8b01826>, 2018.  
1011  
1012 Holstius, D. M., Pillarisetti, A., Smith, K. R., and Seto, E.: Field calibrations of a low-cost aerosol  
1013 sensor at a regulatory monitoring site in California, *Atmospheric Meas. Tech.*, 7, 1121–1131,  
1014 <https://doi.org/10.5194/amt-7-1121-2014>, 2014.  
1015  
1016 Jin, X., Fiore, A. M., Civerolo, K., Bi, J., Liu, Y., Donkelaar, A. van, Martin, R. V., Al-Hamdan, M.,  
1017 Zhang, Y., Insaf, T. Z., Kioumourtzoglou, M.-A., He, M. Z., and Kinney, P. L.: Comparison of  
1018 multiple PM 2.5 exposure products for estimating health benefits of emission controls over New  
1019 York State, USA, *Environ. Res. Lett.*, 14, 084023, <https://doi.org/10.1088/1748-9326/ab2dcb>,  
1020 2019.  
1021  
1022 Johnson, N. E., Bonczak, B., and Kontokosta, C. E.: Using a gradient boosting model to improve  
1023 the performance of low-cost aerosol monitors in a dense, heterogeneous urban environment,  
1024 *Atmos. Environ.*, 184, 9–16, <https://doi.org/10.1016/j.atmosenv.2018.04.019>, 2018.  
1025  
1026 Kim, K.-H., Kabir, E., and Kabir, S.: A review on the human health impact of airborne particulate  
1027 matter, *Environ. Int.*, 74, 136–143, <https://doi.org/10.1016/j.envint.2014.10.005>, 2015.  
1028  
1029 Kuhn, M.: caret: Classification and Regression Training, *Astrophys. Source Code Libr.*,  
1030 ascl:1505.003, 2015.  
1031  
1032 Kumar, P., Morawska, L., Martani, C., Biskos, G., Neophytou, M., Di Sabatino, S., Bell, M.,  
1033 Norford, L., and Britter, R.: The rise of low-cost sensing for managing air pollution in cities,  
1034 *Environ. Int.*, 75, 199–205, <https://doi.org/10.1016/j.envint.2014.11.019>, 2015.  
1035  
1036 Liang, L.: Calibrating low-cost sensors for ambient air monitoring: Techniques, trends, and  
1037 challenges, *Environ. Res.*, 197, 111163, <https://doi.org/10.1016/j.envres.2021.111163>, 2021.  
1038  
1039 Magi, B. I., Cupini, C., Francis, J., Green, M., and Hauser, C.: Evaluation of PM2.5 measured in  
1040 an urban setting using a low-cost optical particle counter and a Federal Equivalent Method Beta  
1041 Attenuation Monitor, *Aerosol Sci. Technol.*, 54, 147–159,  
1042 <https://doi.org/10.1080/02786826.2019.1619915>, 2020.  
1043  
1044 Malings, C., Tanzer, R., Hauryliuk, A., Saha, P. K., Robinson, A. L., Presto, A. A., and  
1045 Subramanian, R.: Fine particle mass monitoring with low-cost sensors: Corrections and long-term  
1046 performance evaluation, *Aerosol Sci. Technol.*, 54, 160–174,  
1047 <https://doi.org/10.1080/02786826.2019.1623863>, 2020.  
1048  
1049 Morawska, L., Thai, P. K., Liu, X., Asumadu-Sakyi, A., Ayoko, G., Bartonova, A., Bedini, A., Chai,  
1050 F., Christensen, B., Dunbabin, M., Gao, J., Hagler, G. S. W., Jayaratne, R., Kumar, P., Lau, A. K.  
1051 H., Louie, P. K. K., Mazaheri, M., Ning, Z., Motta, N., Mullins, B., Rahman, M. M., Ristovski, Z.,  
1052 Shafiei, M., Tjondronegoro, D., Westerdahl, D., and Williams, R.: Applications of low-cost sensing  
1053 technologies for air quality monitoring and exposure assessment: How far have they gone?

1054 Environ. Int., 116, 286–299, <https://doi.org/10.1016/j.envint.2018.04.018>, 2018.  
1055  
1056 Nilson, B., Jackson, P. L., Schiller, C. L., and Parsons, M. T.: Development and Evaluation of  
1057 Correction Models for a Low-Cost Fine Particulate Matter Monitor, *Atmospheric Meas. Tech.*  
1058 *Discuss.*, 1–16, <https://doi.org/10.5194/amt-2021-425>, 2022.  
1059  
1060 Singh, A., Ng'ang'a, D., Gatari, M. J., Kidane, A. W., Alemu, Z. A., Derrick, N., Webster, M. J.,  
1061 Bartington, S. E., Thomas, G. N., Avis, W., and Pope, F. D.: Air quality assessment in three East  
1062 African cities using calibrated low-cost sensors with a focus on road-based hotspots, *Environ.*  
1063 *Res. Commun.*, 3, 075007, <https://doi.org/10.1088/2515-7620/ac0e0a>, 2021.  
1064  
1065 Snyder, E. G., Watkins, T. H., Solomon, P. A., Thoma, E. D., Williams, R. W., Hagler, G. S. W.,  
1066 Shelow, D., Hindin, D. A., Kilaru, V. J., and Preuss, P. W.: The Changing Paradigm of Air Pollution  
1067 Monitoring, *Environ. Sci. Technol.*, 47, 11369–11377, <https://doi.org/10.1021/es4022602>, 2013.  
1068  
1069 Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Calibration of a  
1070 cluster of low-cost sensors for the measurement of air pollution in ambient air, in: 2014 IEEE  
1071 SENSORS, 2014 IEEE SENSORS, 21–24, <https://doi.org/10.1109/ICSENS.2014.6984922>, 2014.  
1072  
1073 Van der Laan, M. J., Polley, E. C., and Hubbard, A. E.: Super learner, *Stat. Appl. Genet. Mol.*  
1074 *Biol.*, 6, 2007.  
1075  
1076 West, S. E., Buker, P., Ashmore, M., Njoroge, G., Welden, N., Muhoza, C., Osano, P., Makau, J.,  
1077 Njoroge, P., and Apondo, W.: Particulate matter pollution in an informal settlement in Nairobi:  
1078 Using citizen science to make the invisible visible, *Appl. Geogr.*, 114, 102133,  
1079 <https://doi.org/10.1016/j.apgeog.2019.102133>, 2020.  
1080  
1081 Williams, R., Kilaru, V., Snyder, E., Kaufman, A., Dye, T., Rutter, A., Russel, A., and Hafner, H.:  
1082 Air Sensor Guidebook, US Environmental Protection Agency, Washington, DC, EPA/600/R-  
1083 14/159 (NTIS PB2015-100610), 2014.  
1084  
1085 Zimmerman, N., Presto, A. A., Kumar, S. P. N., Gu, J., Hauryliuk, A., Robinson, E. S., Robinson,  
1086 A. L., and R. Subramanian: A machine learning calibration model using random forests to improve  
1087 sensor performance for lower-cost air quality monitoring, *Atmospheric Meas. Tech.*, 11, 291–313,  
1088 <https://doi.org/10.5194/amt-11-291-2018>, 2018.  
1089  
1090 Zusman, M., Schumacher, C. S., Gasset, A. J., Spalt, E. W., Austin, E., Larson, T. V., Carvlin, G.,  
1091 Seto, E., Kaufman, J. D., and Sheppard, L.: Calibration of low-cost particulate matter sensors:  
1092 Model development for a multi-city epidemiological study, *Environ. Int.*, 134, 105329,  
1093 <https://doi.org/10.1016/j.envint.2019.105329>, 2020.