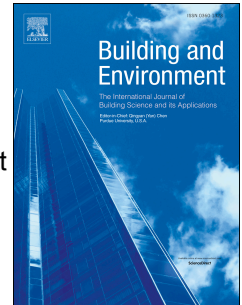


Journal Pre-proof

A machine learning field calibration method for improving the performance of low-cost particle sensors

Satya S. Patra, Rishabh Ramsisaria, Ruihang Du, Tianren Wu, Brandon E. Boor



PII: S0360-1323(20)30824-6

DOI: <https://doi.org/10.1016/j.buildenv.2020.107457>

Reference: BAE 107457

To appear in: *Building and Environment*

Received Date: 25 August 2020

Revised Date: 27 October 2020

Accepted Date: 9 November 2020

Please cite this article as: Patra SS, Ramsisaria R, Du R, Wu T, Boor BE, A machine learning field calibration method for improving the performance of low-cost particle sensors, *Building and Environment* (2020), doi: <https://doi.org/10.1016/j.buildenv.2020.107457>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier Ltd.

A Machine Learning Field Calibration Method for Improving the Performance of Low-Cost Particle Sensors

Satya S. Patra^{a,b}, Rishabh Ramsisaria^{b,c}, Ruihang Du^{b,c}, Tianren Wu^{a,b} and Brandon E. Boor^{a,b*}

^aLyles School of Civil Engineering, Purdue University, West Lafayette, Indiana, 47907, United States.

^bRay W. Herrick Laboratories, Center for High Performance Buildings, Purdue University, West Lafayette, Indiana, 47907, United States.

^cDepartment of Computer Science, Purdue University, West Lafayette, Indiana, 47907, United States.

*Corresponding author: e-mail: bboor@purdue.edu, phone: +1-765-496-0576.

Abstract

Measurements of airborne particles in buildings with low-cost optical particle counters (OPCs) are often inaccurate and subject to uncertainties. This study aims to provide a methodology to improve the performance of low-cost OPCs in measuring indoor particles through machine learning. A two-month field measurement campaign was conducted in an occupied net-zero energy house. The studied OPCs (OPC-N2, Alphasense Ltd.) report size fractionated concentrations from 0.38–17.5 μm . Co-located reference instrumentation included a scanning mobility particle sizer (SMPS: 0.01–0.30 μm) and an optical particle sizer (OPS: 0.30–10 μm). The machine learning field calibration method applies Gaussian Process Regression (GPR) and includes two components: (1.) correction of the size-resolved OPC counting efficiency from 0.38–10 μm and (2.) prediction of volume size distributions (mass proxy) below the 0.38 μm detection limit of the OPC. The field calibration method is applicable to OPCs that report size fractionated concentrations in different size bins. In (1.), a GPR function was used to correct the size-resolved counting efficiency of the OPCs between 0.38–10 μm using the OPS as reference. In (2.), a second GPR function was used to predict the volume size distribution below 0.38 μm using the SMPS/OPS as reference. This was done given the significant contribution of sub-0.38 μm particles to volume concentrations in the accumulation mode. The machine learning field calibration method resulted in a significant improvement in the accuracy of size-integrated volume concentrations ($PV_{2.5}$, PV_{10}) reported by the OPCs as compared to the SMPS/OPS. Improvements were seen in the Pearson coefficient (before correction: 0.59–0.83; after correction: 0.98–0.99); coefficient of determination (before correction: 0.35–0.69; after correction: 0.97–0.98); and mean absolute percentage error (before correction: 35–69%; after correction: 19–25%).

Keywords

indoor air quality; particulate matter; sensors; aerosol size distributions; cooking emissions; optical particle counters

1. Introduction

Airborne particles are an important contributor to indoor air pollution [1,2]. Epidemiological studies have shown that prolonged exposure to elevated particle concentrations can adversely affect respiratory and cardiovascular health [3–8]. Since people spend approximately 90% of their time indoors, it is important to characterize indoor exposure to particles of indoor and outdoor origin through routine particle monitoring in buildings [6,9]. However, traditional particle monitoring equipment is expensive, has a large footprint, and often requires external pumps [10–12]. Thus, widespread deployment of such instrumentation in buildings remains a challenge. Recent advances in sensor technology has driven the emergence of a new particle monitoring paradigm based largely on low-cost optical particle counters (OPCs) for real-time measurement of accumulation (D_p : 0.10 to 2.5 μm) and coarse (D_p : 2.5 to 10 μm) mode particles [11,13,14].

Rai et al. [13] defined a low-cost air quality sensor as one that costs significantly less than the development cost of a sophisticated laboratory-grade instrument. Such sensors have low operating and manufacturing costs and are often associated with short response times for pollutant detection [15]. However, the data provided by low-cost OPCs is often prone to error and standard evaluation criteria are lacking [16]. As a result, recent research has focused on evaluating the performance of low-cost OPCs; a brief summary is provided in **Table 1** [4,7,12,17–30]. A notable limitation of low-cost OPCs is that they often have a low detection efficiency for particles that are smaller than 0.30 μm in diameter [31,32]. This is because sub-0.30 μm particles do not scatter enough light to be detected by most photodetectors [32]. For example, one study found the detection efficiency for a low-cost OPC to be 0.04% at $D_p = 0.10 \mu\text{m}$ [33]. Poor detection of sub-0.30 μm particles is also common in optical-based reference particle instruments, which are commonly used to evaluate the performance of low-cost OPCs [7,20,26]. The sole use of optical-based references offers limited potential to characterize the performance of low-cost OPCs in measuring particles below $D_p = 0.30 \mu\text{m}$.

Low detection efficiency for $D_p \leq 0.30 \mu\text{m}$ by both low-cost OPCs and optical-based references is a concern as the dominant peak for indoor and urban particle volume and mass size distributions is often observed in the accumulation mode [34,35]. As reported by previous studies [34–38], sub-0.30 μm particles can contribute meaningfully to volume and mass concentrations in the accumulation mode. The measurement of size-integrated particle volume ($PV_{2.5}$, PV_{10} , $\mu\text{m}^3 \text{cm}^{-3}$) and mass ($PM_{2.5}$, PM_{10} , $\mu\text{g m}^{-3}$) concentrations must account for the sub-0.30 μm fraction. Holstius et al. [22] evaluated a low-cost OPC in reporting $PM_{2.5}$ against both optical ($0.30 \mu\text{m} < D_p < 2.5 \mu\text{m}$) and non-optical ($D_p < 2.5 \mu\text{m}$) references. High R^2 values (0.64 to 0.95) were found for comparison to optical references, whereas lower R^2 values (0.55 to 0.60) were observed for comparison to non-optical references. Evaluation of low-cost OPCs with optical-based reference instruments with lower detection limits of $D_p = 0.30 \mu\text{m}$ limits the range of assessment to $D_p = 0.30$ to $10 \mu\text{m}$; hence, evaluating $PV_{0.3-2.5}/PM_{0.3-2.5}$ or $PV_{0.3-10}/PM_{0.3-10}$, and not the true $PV_{2.5}/PM_{2.5}$ or PV_{10}/PM_{10} , respectively. In order to develop improved calibration methods for low-cost OPCs in reporting $PV_{2.5}/PM_{2.5}$ or PV_{10}/PM_{10} , it is desirable to use a combination of reference instruments that can measure the full accumulation and coarse modes, from $D_p = 0.10$ to $10 \mu\text{m}$.

Calibration techniques for low-cost OPCs have been improved through application of machine learning algorithms. Researchers have used multiple linear regression models [17, 18] and reported improvement in the performance of low-cost OPCs. While these studies report an improvement in the OPC's performance, the calibration function often introduces a parametric assumption of linear dependency between the OPC and the reference instrument. This may not hold true in many practical applications. Thus, non-parametric machine learning algorithms may be preferred. Si et al. [39] used two non-parametric machine learning algorithms, XGBoost and feedforward neural network, to calibrate a low-cost OPC for ambient air quality monitoring. Several other studies have documented the use of non-parametric machine learning algorithms to develop calibration models for low-cost OPCs [4,40–42]. A disadvantage of these calibration models, however, is that the calibration functions developed are applied directly to size-integrated concentrations (e.g. $PM_{2.5}$) in order to obtain the corrected value. Such calibrations tend to lose essential size-resolved particle information after calibration [43].

Therefore, the literature lacks a robust calibration regime that could improve the performance of low-cost OPCs in forecasting size-integrated concentrations while retaining basic particle information. In this regard, the current study integrates knowledge of basic particle size distribution functions and machine learning algorithms to establish a new field calibration methodology. Two machine-learning non-parametric gaussian process regression (GPR) functions were developed that worked towards addressing the measurement limitations of low-cost OPCs used in indoor environments. One of these developed GPR functions corrected the counting efficiency of coarse mode particles, while the other

predicted the size distribution of accumulation mode particles below the detection limit of the sensor. To the best of the authors' knowledge, no such methodology exists in the literature. The developed machine learning field calibration method was applied to a two-month measurement campaign conducted in an occupied net-zero energy building – the Purdue Retrofit Net-zero: Energy, Water, and Waste (ReNEWW) House. Overall, the paper first presents an overview of the OPCs, reference instruments, field campaign, and the calibration methods, followed by a detailed assessment of how the machine learning calibration approach improves the performance of low-cost OPCs in measuring indoor particles from $D_p = 0.10$ to $10 \mu\text{m}$.

2. Materials and Methods

2.1. Particle Instrumentation and Indoor Air Field Measurements

2.1.1. Description of Low-Cost OPC

The low-cost OPC examined in this study (OPC-N2, Alphasense Ltd.) operates under the principle of light scattering for particle detection and counting [44]. A detailed description of the examined OPC can be found at [11] and [20]. The OPC provides particle counts in 16 discrete size fractions, or bins, from $D_p = 0.38$ to $17.5 \mu\text{m}$. These number concentration outputs are directly converted into size-integrated mass concentrations ($PM_{2.5}$, PM_{10}) using an on-board factory calibration in compliance with European Standard EN481 [44].

The 16 bins of the OPC were defined as ($D_{bin,lower}$ to $D_{bin,upper}$): $0.38\text{--}0.54 \mu\text{m}$, $0.54\text{--}0.78 \mu\text{m}$, $0.78\text{--}1.05 \mu\text{m}$, $1.05\text{--}1.34 \mu\text{m}$, $1.34\text{--}1.59 \mu\text{m}$, $1.59\text{--}2.07 \mu\text{m}$, $2.07\text{--}3 \mu\text{m}$, $3\text{--}4 \mu\text{m}$, $4\text{--}5 \mu\text{m}$, $5\text{--}6.5 \mu\text{m}$, $6.5\text{--}8 \mu\text{m}$, $8\text{--}10 \mu\text{m}$, $10\text{--}12 \mu\text{m}$, $12\text{--}14 \mu\text{m}$, $14\text{--}16 \mu\text{m}$, and $16\text{--}17.5 \mu\text{m}$. The mean diameters for each bin ($D_{bin,mean}$) were: $0.46 \mu\text{m}$, $0.66 \mu\text{m}$, $0.915 \mu\text{m}$, $1.195 \mu\text{m}$, $1.465 \mu\text{m}$, $1.83 \mu\text{m}$, $2.535 \mu\text{m}$, $3.5 \mu\text{m}$, $4.5 \mu\text{m}$, $5.75 \mu\text{m}$, $7.25 \mu\text{m}$, $9 \mu\text{m}$, $11 \mu\text{m}$, $13 \mu\text{m}$, $15 \mu\text{m}$, and $16.75 \mu\text{m}$, respectively. To download the collected data, the serial peripheral interface (SPI) of the OPC was used to connect it to a single-board computer (SBC) (Raspberry Pi 3 Model B+, Raspberry Pi Fdn.) [45]. A Python code was written using the "py-opc" library [46] to log the output of the OPC into a text file in real-time. Three OPCs were used throughout the field measurement campaign in this study. However, one of the OPCs encountered periodic issues with data storage and output and was excluded from the analysis. The remaining two OPCs are referred to as OPC 1 and OPC 2.

2.1.2. Reference Particle Instrumentation

Reference particle instrumentation included an optical particle sizer (OPS) (Model 3330, TSI Inc.) and a scanning mobility particle sizer (SMPS) (Model 3938NL88, TSI Inc.) with a Kr-85 bi-polar charger (370 MBq, Model 3077 A, TSI Inc.), a long Differential Mobility Analyzer (long-DMA, Model 3081, TSI Inc.), and a water-based Condensation Particle Counter (wCPC, Model 3788, TSI Inc.) (Fig. 1). The OPS is a particle spectrometer that measures particle number size distributions from $D_p = 0.30$ to $10 \mu\text{m}$ in optical equivalent diameter [47]. The OPS bin widths were adjusted to match the bin widths of the OPC across its detection range. As the OPS is based on optical size classification, the raw particle data requires correction for the refractive index (R_i) of the sample particle population. The default R_i for the OPS is $1.5 - 0i$. As the field measurements were conducted in a residential indoor environment, three particle categories were considered based on emission activities logged by the residents: background (no documented indoor sources, particles primarily of outdoor origin), candle emissions, and cooking emissions; the latter two were found to be the most common sub-micron indoor particle sources during the measurement period. The R_i values used to correct the raw OPS data for each category are: $1.53 - 0.008i$ for background [48], $1.55 - 0.09i$ for candle emissions [49], and $1.53 - 0.1i$ for cooking emissions [50]. The SMPS operates based on electrical mobility size classification with single particle counting [51-53] and measured particle number size distributions from $D_p = 0.01$ to $0.30 \mu\text{m}$ in electrical mobility diameter.

2.1.3. Indoor Field Measurement Site: Purdue ReNEWW House

To develop and evaluate the machine learning field calibration method, the OPCs were co-located with the reference OPS and SMPS (Fig. 1) in an occupied net-zero energy residence, the Purdue ReNEWW House, located at Purdue University in West Lafayette, Indiana [54,55]. The indoor particle measurements were conducted from November 15, 2018 to January 24, 2019. During this period, the house was occupied by three adult residents, aside from parts of Thanksgiving vacation (November 21 to 24, 2018) and the entirety of winter break (December 15, 2018 to January 05, 2019). The OPCs, OPS, and SMPS were positioned on a table adjacent to the kitchen, a location agreed upon with the residents, with sample inlets 1 m above the floor. The OPS and SMPS were not operational during the winter break period. Indoor combustion (via candle) and cooking events (via electric induction cooktop) were documented. Supplemental

measurements of the indoor air temperature, relative humidity, and operational status of the heating, ventilation, and air conditioning (HVAC) system were made. The HVAC system included an air handling unit (AHU) with a MERV 11 filter and an energy recovery ventilator (ERV), the latter of which delivered outdoor air to the AHU. The AHU/ERV remained operational during the occupied and unoccupied periods.

2.2. Data Processing

The machine learning field calibration method is focused on improving the performance of the OPCs in reporting size-integrated particle volume concentrations ($PV_{2.5}$, PV_{10}) across the full accumulation and coarse modes. Particle volume was selected over particle mass as the latter requires information on size-resolved particle effective density functions [34], which are poorly characterized for indoor particles from $D_p = 0.10$ to $10 \mu\text{m}$. Particle volume and mass size distributions are generally similar in shape [34]. Before implementation of the calibration method described in Section 2.3, the uncalibrated size-integrated particle volume concentrations determined by the OPCs (OPC 1, OPC 2) were compared to the reference OPS and SMPS data. Two approaches were used to determine the uncalibrated volume concentrations as measured by the OPCs. The first evaluation is based on the $PV_{2.5}$ and PV_{10} calculated manually using the raw number concentration output from the OPCs. The second evaluation is based on the direct output (firmware) of $PV_{2.5}$ and PV_{10} from the OPCs. While both approaches were used to evaluate the uncalibrated performance of the OPCs, only the raw number concentration output, and not the firmware output, was used in the machine learning field calibration method.

To calculate the PV manually for the first approach, the raw bin number concentration from the OPC is first converted to a bin volume concentration using Equation 1:

$$dV = \frac{\pi}{6} \times (D_{bin,mean})^3 \times dN \quad (1)$$

where dV is the bin volume concentration ($\mu\text{m}^3 \text{cm}^{-3}$), dN is the raw bin number concentration measured by the OPC (cm^{-3}), and $D_{bin,mean}$ is the mean diameter of the bin (μm), given by Equation 2:

$$D_{bin,mean} = \frac{D_{bin,lower} + D_{bin,upper}}{2} \quad (2)$$

where $D_{bin,lower}$ and $D_{bin,upper}$ are the lower and upper cutoffs for each bin (μm), respectively. $D_{bin,mean}$, $D_{bin,lower}$, and $D_{bin,upper}$ are defined in Section 2.1.1. After the discrete bin volume concentrations are obtained, they are normalized by the bin width as $d\log D_p = \log(\frac{D_{bin,upper}}{D_{bin,lower}})$. Finally, the size-integrated PV is computed via Equation 3:

$$PV = \int \frac{dV}{d\log D_p} \times d\log D_p \quad (3)$$

Two assumptions were made in this approach. First, the particles are assumed to be spheres with a dynamic shape factor (χ) of $\chi = 1$; this is a common assumption when size-resolved variations in χ are not known [56]. Second, all particles in a bin are assumed to have the same diameter, equal to $D_{bin,mean}$; this assumption helps to transform the discrete bin volume distributions into continuous volume distributions and is valid as the width of the bins are small. The uncalibrated $PV_{2.5}$ and PV_{10} were calculated through integration of the continuous volume distributions ($dV/d\log D_p$) from the lower OPC cutoff of $0.38 \mu\text{m}$ to $2.5 \mu\text{m}$ and $10 \mu\text{m}$, respectively. The uncalibrated $PV_{2.5}$ and PV_{10} determined through this approach are henceforth referred to as “calculated $PV_{2.5}$ and PV_{10} .”

For the second approach, the firmware $PV_{2.5}$ and PV_{10} are obtained using the processed firmware $PM_{2.5}$ and PM_{10} outputs, respectively, from the OPC. To convert the firmware PM to PV , the particle effective density assumed by the manufacturer of the low-cost OPC, 1.65 g cm^{-3} [44] (consistent with [34]), is used following [57]:

$$PV_{firmware} = \frac{PM_{firmware}}{1.65 \text{ g cm}^{-3}} \quad (4)$$

The uncalibrated $PV_{2.5}$ and PV_{10} calculated through this approach are henceforth referred to as “firmware $PV_{2.5}$ and PV_{10} .”

To establish the reference dataset, particle volume concentrations were calculated from the number concentrations measured by the OPS and SMPS following Equation 1. The data from the OPS and SMPS were merged to create a continuous volume distribution ($dV/d\log D_p$) from $D_p = 0.01$ to $10 \mu\text{m}$. As previously noted, the SMPS data from $D_p = 0.01$ to $0.30 \mu\text{m}$ is defined by an electrical mobility diameter and the OPS data from $D_p = 0.30$ to $10 \mu\text{m}$ is defined by an optical equivalent diameter. The merged OPS and SMPS data were integrated following Equation 3 to derive size-integrated particle volume concentrations as $PV_{2.5}$: 0.01 to $2.5 \mu\text{m}$ and PV_{10} : 0.01 to $10 \mu\text{m}$. The $D_p = 0.01$ to $0.10 \mu\text{m}$ fraction contributed negligibly to volume concentrations but was included for correctness. The time-series data from the OPCs, OPS, and SMPS were processed using a moving-average method with a 30 min window and a step size of 2 min.

2.3. Machine Learning Field Calibration Method

The development of the non-parametric, size-resolved machine learning field calibration method is motivated by the limitations of particle detection and sizing with low-cost OPCs. Two such limitations were identified. First, the low-cost OPCs are marked with errors while counting the number of particles within its detection range of $D_p = 0.38$ to $17.5 \mu\text{m}$. Second, they are restricted to a lower limit of detection of $D_p = 0.38 \mu\text{m}$ due to insufficient scattering of incident light. The OPC cannot see these small particles, which can contribute significantly to particle volume and mass in the accumulation mode. Thus, to address these limitations and account for the missing sub- $0.38 \mu\text{m}$ particles, the field calibration method includes two components (**Fig. 1**): (1.) correction of the size-resolved counting efficiency of the OPC from $D_p = 0.38$ to $10 \mu\text{m}$ (upper limit set by OPS range of detection) and (2.) prediction of particle volume size distributions ($dV/d\log D_p$) below $D_p = 0.38 \mu\text{m}$. Both (1.) and (2.) are carried out using a non-parametric machine learning algorithm, Gaussian Process Regression (GPR). In (1.), a GPR function is used to correct the raw particle counts of the OPCs from $D_p = 0.38$ to $10 \mu\text{m}$ using the OPS as reference. In (2.), a second GPR function is used to predict the particle volume size distribution below $D_p = 0.38 \mu\text{m}$ using the available data in the OPC's detection range with the OPS and SMPS as reference.

2.3.1. Correction of OPC Size-Resolved Counting Efficiencies

The size-resolved particle counting efficiency is represented as the mean ratio of the bin volume concentration of the OPC to the corresponding bin volume concentration of the OPS. Mathematically it is expressed via Equation 5 as:

$$\text{Counting efficiency} = \frac{1}{n} \sum_{i=1}^n \frac{\text{Bin volume concentration}_{\text{OPC}}}{\text{Bin volume concentration}_{\text{OPS}}} \quad (5)$$

A value greater than 1 indicates an overestimation of particle counts and less than 1 indicates an underestimation of particle counts. Low-cost OPCs are commonly associated with errors in size-resolved counting efficiency [11,20]. This results in an inaccurate estimation of size-integrated volume concentrations. Thus, to correct this, a GPR correction model is proposed. In this model, the bin concentrations of the OPC are corrected against the corresponding bin concentrations of the OPS.

GPR is a non-parametric Bayesian approach to regression problems [58]. Let there be n set of observations Y for n set of inputs X , where $Y = \{y_1, y_2, y_3 \dots y_n\}$, and $X = \{x_1, x_2, x_3 \dots x_n\}$. The regression finds a function $f(x)$ such that:

$$\begin{Bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{Bmatrix} \sim \begin{Bmatrix} (y_1) \\ (y_2) \\ \vdots \\ (y_n) \end{Bmatrix} \quad (6)$$

GPR constructs this function $f(x)$ using a Gaussian Process (GP) in such a way that the distribution connecting outputs of any two or more points in its domain will form a multivariate joint Gaussian distribution [59]. This is done using a mean ($\mu(x)$) and covariance function ($k(x, x')$). Mathematically, $f(x)$ in GPR can be written as Equation 7:

$$f(x) \sim GP(\mu(x), k(x, x')) \quad (7)$$

where $\mu(x)$ is the mean function which gives the expected value at input x . The prior mean (function before the prediction) is often set to zero in order to avoid complicated calculations [60]. Next, the covariance function $k(x, x')$ models the association between the function values at different input points x and x' . This function is usually termed the kernel of the GP [61]. The choice of a suitable kernel is based on assumptions such as smoothness and probable data patterns. A

typical assumption is that the correlation between two points decreases with the distance between them. This means that closer points should behave more similarly than distant points. One of the principal kernels to meet this assumption is the radial base function kernel [60], which takes the form of Equation 8:

$$k(x, x') = \sigma_f^2 e^{-\frac{(x-x')^2}{2l^2}} \quad (8)$$

where σ_f and l are the hyper-parameters, optimized according to the data, which determines the smoothness of the GPR model. By definition of GPR, the observations Y and the function $f(x)$ form a multivariate joint normal distribution. This distribution takes the form of Equation 9 [60]:

$$\begin{pmatrix} Y \\ f \end{pmatrix} \sim N \left(0, \begin{pmatrix} K(X_i, X_i) + \sigma^2 I & K(X_i, X_p) \\ K(X_p, X_i) & K(X_p, X_p) \end{pmatrix} \right) \quad (9)$$

Here, X_i is the set of observed input values, Y_i are the corresponding output values, and X_p is the set of inputs where a prediction is to be made. $K(X_i, X_i)$ is the covariance matrix for all observed points. This is calculated using the covariance function as given in Equation 10 [59,60]:

$$K(X_i, X_i) = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_i) \\ \vdots & \ddots & \vdots \\ k(x_i, x_1) & \cdots & k(x_i, x_i) \end{bmatrix} \quad (10)$$

In Equation 9, σ is the variance of the error and I is the identity matrix. $K(X_p, X_p)$ is the covariance matrix for inputs where a prediction is to be made. $K(X_i, X_p)$ is the covariance matrix between the observed input points and the points where a prediction is to be made. $K(X_p, X_i)$ is the covariance matrix between the points where a prediction is to be made and the observed input points. All these matrices are calculated in a similar manner as shown in Equation 10.

The best estimate of a prediction for a point x in the set X_p is the mean of the posterior joint distribution at x [59], given by Equation 11 [60,62]:

$$\text{Predicted value at } x = K(X, X_i)[K(X_i, X_i) + \sigma^2 I]^{-1} Y_i \quad (11)$$

This is the theory for prediction using GPR. For correcting the size-resolved counting efficiencies of the OPCs, their 30 min averaged raw bin volume concentrations formed the input set, and the corresponding 30 min averaged reference bin volume concentrations from the OPS constituted the output set. As stated before, both the OPCs and OPS had the same bin widths. The GP function mapping these two sets formed the corrective function.

One point to note here is that the indoor air temperature and relative humidity are not included in this corrective function. The temperature was excluded because it remained stable during the two-month measurement period at the ReNEW House (mean: 19.3°C, median: 19.8°C, stdev: 1.7°C). Similarly, the relative humidity at the ReNEW House remained below 50% (mean: 42.0%, median: 41.7%, stdev: 6.8%). Hygroscopic growth of particles is often negligible at relative humidities below 85% [11], hence, the relative humidity was excluded as it was not expected to influence the performance of the OPCs as is common for outdoor measurements at elevated relative humidities.

2.3.2. Prediction of Particle Volume Size Distributions Below the Lower Detection Limit of the OPC

The particle volume size distribution can be expressed mathematically as a multi-modal lognormal distribution function [57]:

$$\frac{dV}{d \log D_p} = \sum_{i=1}^n \frac{A_i}{(2\pi)^{1/2} \log(c_i)} e^{-\frac{(\log D_p - \log b_i)^2}{2 \log^2(c_i)}} \quad (12)$$

where A_i is the volume concentration ($\mu\text{m}^3 \text{ cm}^{-3}$), b_i is geometric mean diameter (μm), and c_i is the geometric standard deviation (-) for each mode (i); and D_p is the particle diameter (μm), here the $D_{bin,mean}$ for each bin. In general, there exists two dominant modes (or peaks) for particle volume size distributions for indoor and urban air [57]: one mode in the $D_p = 0.10$ to $2.5 \mu\text{m}$ fraction (termed the accumulation mode) and one mode in the $D_p = 2.5$ to $10 \mu\text{m}$ fraction (termed the coarse mode). As the OPCs have a lower limit of detection of $D_p = 0.38 \mu\text{m}$, they cannot resolve the full extent of the accumulation mode. A second GPR function was developed to predict the particle volume size distribution below the $D_p = 0.38 \mu\text{m}$ detection limit of the OPC using its calibrated concentration data (following Section 2.3.1) above $D_p = 0.38 \mu\text{m}$.

Reference particle volume size distributions as measured by the OPS and SMPS were used to develop the second GPR function. The OPS and SMPS data were fit to the multi-modal lognormal distribution function based on a nonlinear least-squares curve fitting function in MATLAB (The MathWorks, Inc.). The fitting provided the scalar parameters A_i , b_i , c_i for the accumulation and coarse modes. A GPR training architecture was created whereby the OPS and SMPS reference volume concentrations for $D_p > 0.38 \mu\text{m}$ formed the input set and the corresponding accumulation mode scalar fitting parameters for the fitted curve formed the output set. This architecture was then trained for mapping the volume concentration with the accumulation mode scalar parameters. The formed GP function predicts the scalar fitting parameters for the accumulation mode of the particle volume size distribution using the measured volume concentration data for $D_p > 0.38 \mu\text{m}$.

Once the training was completed, the calibrated bin volume concentrations for the OPCs (following Section 2.3.1) were fed into it as input, the output of which were the scalar fitting parameters for the accumulation mode. The fitting parameters were used to draw the sub- $0.38 \mu\text{m}$ fraction of the particle volume size distribution curve for the OPCs. This created a continuous particle volume size distribution across the accumulation and coarse modes, from $D_p = 0.10$ to $10 \mu\text{m}$. The $D_p = 0.01$ to $0.10 \mu\text{m}$ (termed the Aitken mode) was included for completeness, however, it is present as the tail of the volume size distribution. The particle volume size distribution was integrated following Equation 3 to derive calibrated and size-integrated particle volume concentrations ($PV_{2.5}$, PV_{10}) for each OPC (OPC 1, OPC 2). The calibrated $PV_{2.5}$ and PV_{10} calculated through this approach are henceforth referred to as “corrected $PV_{2.5}$ and PV_{10} .”

The computation and optimization of the GPR functions for correction of the size-resolved counting efficiency of the OPC and prediction of particle volume size distributions below $D_p = 0.38 \mu\text{m}$ was executed in the regression learner toolbox in MATLAB. For training and testing the GPR functions, the collected data from the measurement campaign at the ReNEW House was divided into two sets: a training set and a testing set. The two GPR functions were built using the data from the training set, which consisted of approximately 6 weeks of data. Once the correction functions were trained, they were evaluated on the testing set over one week.

2.4. Evaluation Metrics for OPC Performance

Two features used to evaluate the performance of low-cost sensors are accuracy and precision. Precision is calculated for a set of similar test sensors and accuracy is computed between a test sensor and a reference instrument. There are several metrics to quantify them. The coefficient of variation (CV) is an indicator of precision among a set of similar sensors [11]. It is calculated using Equation 15 as:

$$CV = \frac{\sigma}{\mu} \quad (15)$$

where σ is the standard deviation and μ is the mean of the intra-OPCs' observation for each 30 min time-averaged interval. The average CV for both the calculated and firmware $PV_{2.5}$ and PV_{10} as reported by OPC 1 and OPC 2 was determined. This quantity represents the variation in the measurements given by OPC 1 and OPC 2. According to the U.S. Environmental Protection Agency (EPA), CV values up to 10% are acceptable [20]. This threshold for CV is mentioned in CFR Part 58-Ambient Air Quality Surveillance (Subchapter C) [33].

For accuracy evaluation, the calculated, firmware, and corrected $PV_{2.5}$ and PV_{10} from the two OPCs were compared against the reference values from the OPS and SMPS. For each comparison, statistical measures, such as the Pearson coefficient (r), coefficient of determination (R^2), and the slope between the OPC values and OPS/SMPS values were determined. The U.S. EPA recommends a $r \geq 0.97$ and a slope of 1 ± 0.1 [20] between the test sensor and the reference

instrument. The U.S. EPA criteria is often used to evaluate low-cost OPCs [11,20,33]. R^2 is another statistical measure reflecting the proportion of variance between two variables. Thus, values closer to 1 are desirable.

This study uses another evaluation metric for sensor accuracy, the mean absolute percentage error (MAPE) [63]. MAPE is calculated using Equation 16 as:

$$MAPE = \left(\frac{1}{n} \sum_{i=1}^n \frac{|PV_{OPC} - PV_{OPS/SMPS}|}{PV_{OPS/SMPS}} \right) \times 100 \quad (16)$$

MAPE values reflect the overall bias of the OPCs. MAPE was calculated for the $PV_{2.5}$ and PV_{10} obtained from OPC 1 and OPC 2 and the reference instruments. Lower MAPE values indicate better sensor output.

A statistical test, a paired t -test, was also performed in this evaluation. A paired t -test is a statistical analysis used to evaluate the significance of the difference between two separate measurements of the same subject [64]. Here, this test is used to assess the difference between the $PV_{2.5}$ and PV_{10} given by the two OPCs and the reference instruments. To conduct this test, the difference between the two observations in each pair is first calculated. Thereafter, the mean and the standard deviation of the difference is estimated. A t -statistic is then determined using Equation 17:

$$t - statistic = \frac{\text{mean of the difference}}{\text{standard deviation of the difference}} \quad (17)$$

Finally, corresponding to the value of the t -statistic, using the t -distribution table, enables determination of the p -value for the paired t -test. The null hypothesis of the paired t -test assumes that the statistic follows a t -distribution. If the p -value of the paired t -test is less than 0.05, it is concluded that the mean difference between the paired observations is significantly different [65]. The measurements are therefore only acceptable if the p -value for t -testing between the OPCs and reference instruments is greater than 0.05.

3. Results and Discussion

The following sections present the performance evaluation of the low-cost OPCs prior to, and after, implementation of the machine learning field calibration method for the two-month measurement campaign at the Purdue ReNEW House. First, the uncalibrated performance of the low-cost OPCs is compared against the reference instruments (OPS and SMPS). The assessment for the same is discussed as per the criteria mentioned in CFR Part 58-Ambient Air Quality Surveillance (Subchapter C) [33]. Thereafter, the results of the two components of the field calibration method are discussed. Finally, the cumulative results of the proposed calibration method are compared with the reference measurements in the test dataset to determine if the low-cost OPCs met the U.S. EPA's criteria after the field calibration.

3.1. Performance Evaluation of the OPC Prior to Field Calibration

Fig. 2 illustrates the comparison between the uncalibrated $PV_{2.5}$ and PV_{10} (calculated and firmware) time-series as measured by OPC 1 with the $PV_{2.5}$ and PV_{10} time-series as measured by the reference instruments (OPS and SMPS) for the entirety of the measurement campaign at the ReNEW House. For the occupied sampling periods at the ReNEW House, the uncalibrated OPCs' calculated mean $PV_{2.5}$ concentrations were: $0.52 \mu\text{m}^3 \text{cm}^{-3}$ for OPC 1 and $0.54 \mu\text{m}^3 \text{cm}^{-3}$ for OPC 2. For the same duration and time-averaging window, the mean $PV_{2.5}$ as measured by the OPS/SMPS was $2.65 \mu\text{m}^3 \text{cm}^{-3}$. This suggests an underestimation in $PV_{2.5}$ as calculated using the raw number concentrations from both OPCs. Similar findings were observed for the calculated PV_{10} . The underestimations in calculated PV concentrations are attributed to limitations in particle detection and sizing by the low-cost OPCs, including errors in counting efficiency across their detection range and lack of accounting for the volume contribution of sub- $0.38 \mu\text{m}$ particles. The proprietary self-calibration in the OPC firmware used to report the firmware PM outputs [20] resulted in an improvement over the calculated PV values when compared to the reference instruments. The mean firmware $PV_{2.5}$ concentrations were: OPC 1: $1.91 \mu\text{m}^3 \text{cm}^{-3}$ and OPC 2: $2.02 \mu\text{m}^3 \text{cm}^{-3}$. Despite the improvements over the calculated PV values, the firmware outputs also underestimated the size-integrated volume concentrations when compared to the reference values (**Fig. 2**).

For the unoccupied sampling periods at the ReNEW House, mean calculated PV values were: OPC 1 $PV_{2.5}$: $0.23 \mu\text{m}^3 \text{cm}^{-3}$, OPC 2 $PV_{2.5}$: $0.22 \mu\text{m}^3 \text{cm}^{-3}$; OPC 1 PV_{10} : $0.71 \mu\text{m}^3 \text{cm}^{-3}$, OPC 2 PV_{10} : $0.68 \mu\text{m}^3 \text{cm}^{-3}$. The mean firmware PV values were: OPC 1 $PV_{2.5}$: $0.85 \mu\text{m}^3 \text{cm}^{-3}$, OPC 2 $PV_{2.5}$: $0.81 \mu\text{m}^3 \text{cm}^{-3}$; OPC 1 PV_{10} : $1.98 \mu\text{m}^3 \text{cm}^{-3}$, OPC 2 PV_{10} : $2.12 \mu\text{m}^3 \text{cm}^{-3}$.

These values were significantly less (p -value $\ll 0.05$) than the corresponding mean PV concentrations during the occupied periods. This is due in part to a greater prevalence of human-associated indoor sources of accumulation and coarse mode particles during occupied periods [66]. As the AHU/ERV remained operational during the unoccupied periods, the measured particles are likely of outdoor origin as outdoor air was continually introduced into the AHU by the ERV. Despite issues in the accuracy of the uncalibrated OPCs in reporting PV concentrations, they were capable of detecting changes in concentrations due to variations in house occupancy patterns and associated emission events.

To quantify the difference in PV concentrations reported by the two uncalibrated OPCs (OPC 1 and OPC 2) with each other, and the reference instruments (OPS and SMPS), the precision and accuracy metrics discussed in Section 2.4 were estimated. Each accuracy metric was calculated for the PV concentration data and for the period when the ReNEW House was occupied. Firstly, the CV values for both OPC outputs were estimated to evaluate the variability in measurements between OPC 1 and OPC 2. The mean CV values for the calculated and firmware $PV_{2.5}$ concentrations were 4.82% and 4.07%, respectively. For PV_{10} , the mean CV values for the calculated and firmware concentrations were 4.94% and 8.73%, respectively. For all comparisons, the CV values were within the limits recommended by the U.S. EPA ($\leq 10\%$). This indicates that both OPCs were consistent with respect to each other in terms of output PV concentrations.

Accuracy metrics, including r , R^2 , slope, $MAPE$, and p -values of the t -test, for comparing the OPCs with the reference instruments (OPS and SMPS) were calculated for the uncalibrated $PV_{2.5}$ and PV_{10} (calculated and firmware). **Table 2** summarizes the results. The slope and r values, relative to the reference concentrations, obtained for both calculated and firmware PV values given by both OPCs did not meet the U.S. EPA's performance criteria. The $MAPE$ for the calculated $PV_{2.5}$ relative to the reference $PV_{2.5}$ was observed to be 79.81% for OPC 1 and 78.97% for OPC 2. For the calculated PV_{10} , the $MAPE$ was found to be 67.67% for OPC 1 and 67.59% for OPC 2. Such high values in $MAPE$ suggest high bias in the sensors' output. In addition, lower R^2 values were observed (**Table 2**). The p -values for the t -tests were less than 0.05. Therefore, at 95% confidence, it is statistically concluded that the calculated data given by both OPCs differed significantly from the reference concentrations measured by the OPS and SMPS. The performance metrics for the firmware PV outputs by both OPCs were better than the PV outputs calculated using the raw number concentrations (**Table 2**). However, despite the proprietary self-calibration in the firmware, the error values remained as high as 40 to 50% and there still exists a statistical disagreement between the firmware values and the reference concentrations.

Despite pronounced differences in the PV concentrations reported by the uncalibrated OPCs and the OPS/SMPS, they tend to follow a similar pattern over time (**Fig. 2**). In order to better characterize this, the average diurnal variation in the PV concentrations during occupied periods at the ReNEW House were determined for the two OPCs and OPS/SMPS (**Fig. 3**). A consistent numerical difference in PV concentrations given by the OPCs and the reference instruments can be observed. However, the shape of the diurnal trend is consistent, suggesting the uncalibrated OPCs are capable of detecting changes in accumulation and coarse mode particle concentrations due to variations in resident activity patterns. Lower PV levels are observed during the evening hours (21:00 to 06:00) due to a lower prevalence of indoor particle sources. PV concentrations gradually increase during the day due in part to human-associated indoor particle sources. Diurnal variations in outdoor PV concentrations are also a contributing factor, due to outdoor air delivery to the AHU by the ERV, however, outdoor particles were not measured during the campaign. Such observations are consistent with prior measurements of diurnal trends in indoor particle concentrations [66,67]. Another observation from **Fig. 3** is that the diurnal trends reported by both OPCs agree well with respect to each other (for both calculated and firmware PV values). This corroborates the low CV values obtained.

3.2. Performance Evaluation of the OPC After Field Calibration

3.2.1. Correction of OPC Size-Resolved Counting Efficiencies

The first component of the machine learning field calibration method is correction of the size-resolved counting efficiency of the OPCs from $D_p = 0.38$ to $10 \mu m$ through use of a GPR function with the OPS as reference (Section 2.3.1). Once the GPR function was trained during the training set, the raw bin volume concentration data from the testing set was corrected using the developed function. **Fig. 4** shows the comparison in the mean size-resolved counting efficiencies for both OPCs before (raw) and after (corrected) calibration. The counting efficiency is reported for the mean particle diameter of each bin, $D_{bin,mean}$. The raw counting efficiencies are calculated for the entire measurement campaign (training and testing sets), while the corrected counting efficiencies are calculated for the testing set. Raw and corrected size-resolved counting efficiencies were similar for OPC 1 and OPC 2.

Before correction with the trained GPR function, counting efficiencies for $D_p < 0.80 \mu\text{m}$ were less than unity for OPC 1 and OPC 2, suggesting that sub- $0.80 \mu\text{m}$ PV concentrations were under-estimated by both OPCs. The counting efficiency reached a maxima at $D_p = 1.195 \mu\text{m}$. For particles in the range of $D_p = 0.80$ to $1.6 \mu\text{m}$, the OPCs over-estimate the raw particle counts. Above $D_p = 1.6 \mu\text{m}$, the OPCs begin to underestimate again, reaching a minimum value at $D_p = 1.83 \mu\text{m}$. However, the counting efficiency values begin to improve with an increase in particle size above $D_p = 2.0 \mu\text{m}$. Therefore, it is evident that prior to correction, there are periodic size fractions where raw counts were under- and over-estimated when compared to the reference instrument (OPS). This inconsistency in size-resolved particle detection is one of the key reasons for the poor performance of the OPCs described in Section 3.1. Following correction with the trained GPR function, a substantial improvement in the size-resolved counting efficiencies of both OPCs can be seen (**Fig. 4**). Values near unity were found for all OPC bins, demonstrating that the proposed correction scheme was effective in improving the particle detection efficiency of both OPCs from $D_p = 0.38$ to $10 \mu\text{m}$.

3.2.2. Prediction of Particle Volume Size Distributions Below the Lower Detection Limit of the OPC

Following correction of the size-resolved counting efficiency of both OPCs with the first GPR function, the corrected PV concentrations from $D_p = 0.38$ to $10 \mu\text{m}$ were used in the second GPR function to obtain estimates of the sub- $0.38 \mu\text{m}$ particle volume size distributions (Section 2.3.2). As previously noted, sub- $0.38 \mu\text{m}$ particles can contribute meaningfully to particle volume concentrations in the accumulation mode [34,35]. To confirm this observation, the contribution of sub- $0.38 \mu\text{m}$ particles to indoor $PV_{2.5}$ and PV_{10} as measured with the OPS and SMPS at the ReNEW House were calculated for the entire field campaign (**Fig. 5**). The average contribution of sub- $0.38 \mu\text{m}$ particles to $PV_{2.5}$ and PV_{10} were found to be approximately 40 and 27%, respectively. Therefore, the uncalibrated $PV_{2.5}$ and PV_{10} reported by the OPCs in Section 3.1 do not account for a significant fraction of the indoor particle volume. Accounting for the missing sub- $0.38 \mu\text{m}$ particle volume is a critical element in improving the performance of OPCs in reporting $PV_{2.5}$ and PV_{10} .

Application of the trained GPR function for predicting the particle volume size distribution below $D_p = 0.38 \mu\text{m}$ is illustrated in **Fig. 6**. Three example volume distributions from $D_p = 0.01$ to $10 \mu\text{m}$ as measured during the testing set are shown. For $D_p > 0.38 \mu\text{m}$, the red curve indicates the corrected OPC data using the first GPR function (for counting efficiency). For $D_p < 0.38 \mu\text{m}$, the blue curve indicates the predicted OPC data using the second GPR function (missing sub- $0.38 \mu\text{m}$ particle volume). The second GPR function enables prediction of volume size distributions below $D_p = 0.10 \mu\text{m}$, however, Aitken mode particles contributed negligibly to volume concentrations (**Fig. 6**).

Fig. 7 illustrates the reference particle volume size distribution time-series as measured by the OPS and SMPS and the calibrated particle volume size distribution time-series as reported by OPC 1 and OPC 2 after implementation of the two GPR functions for the entirety of the testing set (January 15 to 23, 2019). It is evident that the correction functions result in good agreement between the calibrated OPC distributions and the reference distributions between $D_p = 0.01$ to $10 \mu\text{m}$. Notably, the non-parametric, size-resolved machine learning field calibration method enables for an effective OPC measurement range beyond its traditional detection range. **Fig. 7** demonstrates that the calibrated OPCs can reliably capture temporal variations in the shape and magnitude of the indoor $dV/d\log D_p$ at the ReNEW House. The diurnal trend in indoor particle volume distributions as measured by the OPCs and OPS/SMPS is similar to that observed for the size-integrated volume concentrations in **Fig. 3**. Values for $dV/d\log D_p$ in the coarse mode are especially pronounced during the day and peak during evening hours, due in part to human activity-driven emissions [68]. The use of a candle on the evening of January 19 resulted in a spike in particle concentrations across all size fractions (**Fig. 7**), similar to prior studies [69,70].

3.2.3. Evaluation of the Field Calibration Method

The corrected OPC particle volume size distributions during the testing set (**Fig. 7**) were integrated following Equation 3 to obtain corrected size-integrated volume concentrations ($PV_{2.5}$ and PV_{10}). A comparison in the corrected $PV_{2.5}$ and PV_{10} time-series as measured by OPC 1 and OPC 2 with the reference instruments is shown in **Fig. 8**. It is evident that after the calibration, the agreement in $PV_{2.5}$ and PV_{10} as measured by the two OPCs with the OPS/SMPS has substantially improved. The mean corrected $PV_{2.5}$ and PV_{10} were $2.64 \mu\text{m}^3 \text{cm}^{-3}$ and $6.43 \mu\text{m}^3 \text{cm}^{-3}$ for OPC 1, and $2.31 \mu\text{m}^3 \text{cm}^{-3}$ and $6.21 \mu\text{m}^3 \text{cm}^{-3}$ for OPC 2, respectively. The mean reference $PV_{2.5}$ and PV_{10} were $2.97 \mu\text{m}^3 \text{cm}^{-3}$ and $6.69 \mu\text{m}^3 \text{cm}^{-3}$, respectively. The corrected values were close to the reference concentrations. Diurnal variations in the corrected $PV_{2.5}$ and PV_{10} as reported by OPC 1 and OPC 2 during the testing set is shown in **Fig. 9**. The diurnal trend is similar to that

observed in **Fig. 3**, however, the corrected output from the two OPCs more closely follows the reference values in both shape and magnitude.

To quantify the performance of the calibration methodology, the accuracy metrics, r , R^2 , slope, $MAPE$, and p -values of the t -test, between the corrected OPC PV and reference PV concentrations for the testing set were calculated (**Table 3**). To evaluate the level of improvement achieved after the field calibration was applied, the accuracy metrics were also evaluated for the uncalibrated calculated and firmware PV concentrations reported by both OPCs during the testing set (**Table 3**). Before correction, high $MAPE$ (57.1 to 69.2%) and low R^2 (0.41 to 0.44) values were observed for both OPCs for calculated $PV_{2.5}$ and PV_{10} . The r and the slope values did not meet the U.S. EPA's criteria. In addition, the p -values showed a statistical difference between the calculated PV concentrations by both OPCs and the reference concentrations. The firmware $PV_{2.5}$ and PV_{10} before correction agreed better with the reference instruments ($MAPE$: 35.4 to 63.3%; R^2 : 0.35 to 0.69), however, the values still failed to meet the U.S. EPA's recommended values for r and slope. In addition, it was observed that the firmware $PV_{2.5}$ values had a greater number of outliers than the firmware PV_{10} values. For the same reason, the firmware $PV_{2.5}$ were reported with low r , R^2 , and slope values (linear dependency decreased due to more outliers). Nevertheless, the overall improvements in the $MAPE$ values for the firmware $PV_{2.5}$ was greater than that of the firmware PV_{10} .

The corrected PV concentrations for both OPCs after implementation of the field calibration methodology showed a significant reduction in $MAPE$ (19.2 to 25.1%) as compared to the uncalibrated OPC data. High R^2 values were also obtained (> 0.97). The r and the slope values for the corrected PV concentrations met the U.S. EPA's criteria. The p -values for the t -tests between the corrected OPC PV concentrations and the reference concentrations are > 0.05 . This indicates that after the corrections were applied to the OPC data, there was no statistical evidence of significant disparities between the OPC and reference $PV_{2.5}$ and PV_{10} concentrations.

Correlation plots for the calculated and firmware (uncalibrated) and corrected $PV_{2.5}$ and PV_{10} concentrations as reported by OPC 1 and OPC 2 against the reference $PV_{2.5}$ and PV_{10} concentrations as measured by the OPS/SMPS for the testing set are shown in **Fig. 10**. Points nearest to the 1:1 line are considered as an acceptable measurement. The calculated and firmware (uncalibrated) $PV_{2.5}$ and PV_{10} for OPC 1 and OPC 2 deviate significantly from the 1:1 line. The firmware PV values are heteroscedastic and the calculated PV values consistently underestimate the reference PV concentrations. Conversely, the corrected PV values are homoscedastic and close to the 1:1 line, thus supporting the results given in **Table 3**. The results presented in **Fig. 6-10** and **Table 3** demonstrate that the machine learning field calibration method was successful in addressing the limitations in particle detection and sizing of the OPCs, improving their performance in measuring accumulation and coarse mode particles in residential indoor environments.

3.2.4. Training Dataset Size Sensitivity on Model Results

Being a data-driven approach, the developed field calibration methodology may depend on the size of the training dataset. This is because a larger dataset can better model the uncertainties and interactions between the predicted and response variables. However, as the training data size increases, the complexity of the model also increases. Therefore, to understand the effect of the training dataset size on the performance of the developed field calibration model, it was evaluated by varying the size of the training dataset. 8 sets of training datasets were formed with 3, 3.5, 4, 4.5, 5, 5.5, 6 and 6.5 weeks of training data, respectively. The remaining data in each set (out of seven weeks of data) formed the testing data. In each of the developed sets, the correction functions were developed using the data from the training set, and the $MAPE$ was calculated applying the corrective functions in the testing set (as discussed in equation 16). **Fig. 11** presents the results of the $MAPE$ values obtained for different sizes of training dataset ($PV_{2.5}$ and PV_{10}).

As expected, the errors in the calibration model have reduced with an increase in the size of the training datasets for both of the size integrated concentrations. However, when the training sets were smaller, the reduction was steeper. The reduction of errors was negligible after six weeks of training data, suggesting reasonable sensor drift capture after this timeframe.

3.2.5. Advantages and Limitations of the Field Calibration Method

There are several advantages of using this field calibration method for low-cost OPCs for indoor air measurement. First, the correction functions used are non-parametric, hence, they are expected to perform better than parametric regression

models. For example, Holstius et al. [22] used a linear regression model for a low-cost OPC and achieved a maximum R^2 value of 0.72. In another study, Magi et al. [18] used a multiple linear regression model and achieved a R^2 of 0.60. The use of GPR models in this study resulted in R^2 values in the range of 0.97 to 0.98. This signifies that the developed correction model can explain up to 97 to 98% of variation in the data. The use of a GPR model throughout the correction is inspired by the fact that it is a non-parametric machine learning regression model. Thus, any prior information about the relationship between concentration values from the OPCs and the reference instruments is not needed, making the correction function free from any prior assumptions. Other non-parametric machine learning models are available that can do the same; however, GPR has the advantage of being simple and learning better from the data [71].

Another advantage of using this field calibration method is that it enables the user to obtain an estimated particle volume size distribution across a wide size range ($D_p = 0.01$ to $10\ \mu\text{m}$) using a single low-cost OPC. It is important to know the full extent of the volume size distribution prior to calculating size-integrated volume concentrations. This is especially true for the accumulation mode, where sub- $0.38\ \mu\text{m}$ particle contribute significantly to indoor particle volume concentrations (Fig. 5). Interestingly, the manual for the OPC (OPC-N2, Alphasense Ltd.) states: "The OPC-N2 calculations of particle mass assume a negligible contribution from particles below approximately $0.38\ \mu\text{m}$ " [44].

There are several limitations of the field calibration method. First, the results presented in this study are expressed in terms of particle volume and not mass. Although the two are typically analogs to one another, PM is much more commonly used by the air quality research community. Integration of measured size-resolved particle effective densities are needed to expand the calibration method from PV to PM . Another limitation is that the method is limited to sampling of particles in conditioned indoor environments, similar to the ReNEW House. The method cannot be applied to outdoor particle measurement with the OPCs as air temperature and relative humidity should be included in the corrective functions. Furthermore, the proposed field calibration method is data driven. Therefore, a large training dataset (at least 5.5 to 6 weeks of training data) is required to effectively develop the correction GPR functions as discussed in the study. Lastly, due to its data-driven nature, the developed corrective functions are site-specific. Changing the indoor field site, or a significant change in particle sources, would require retraining of the functions. However, it may be noted that, as suggested by Wang et al. [15], low-cost OPCs are sensitive to particle sources. Therefore, it is advisable to periodically calibrate the sensors whenever particle source changes are expected. The current study aimed to provide a novel calibration regime that could possibly be used by the manufacturers or end-users of OPCs, similar to the Alphasense OPC-N2. The advantage of using this methodology is that it would improve the low-cost OPC's performance in forecasting size integrated concentrations and enable the OPC to forecast reliable size distributions even below its detection range.

4. Conclusion

This study discussed the development of a machine learning field calibration method to improve the accuracy of low-cost OPCs in measuring accumulation and coarse mode particles in residential buildings. The evaluated low-cost OPC was co-located with two reference instruments (OPS and SMPS) in a net-zero energy house. The machine learning field calibration method addressed measurement limitations of OPCs by developing two GPR functions. First, errors in the size-resolved counting efficiencies of the OPCs from $D_p = 0.38$ to $10\ \mu\text{m}$ were corrected using a GPR function. Then, a second GPR function was used to predict the volume size distribution below the $D_p = 0.38\ \mu\text{m}$ detection limit of the OPCs. The final corrected OPC-based size-integrated volume concentrations ($PV_{2.5}$ and PV_{10}) exhibited a significant improvement in several accuracy metrics (R^2 , $MAPE$) over the uncalibrated data. The r and the slope values for the corrected OPC data relative to the reference data were within the limits prescribed by the U.S. EPA. In addition, statistical tests proved insignificant differences between the calibrated OPCs and reference instruments. This suggests that the proposed methodology was successful in correcting the low-cost OPCs for forecasting reliable size-integrated particle volume concentrations in indoor environments. Furthermore, the proposed field calibration method enabled the low-cost OPC to estimate particle volume size distribution across a wide size range, with size distributions comparable to those measured by the OPS and SMPS.

Acknowledgements

This work was supported by the National Science Foundation (CBET-1805804) and the Purdue University Instructional Innovation Program. The authors thank the residents of the Purdue ReNEW House and Jason Schneemann of Whirlpool Corporation for their help and support.

References

- [1] P.N. Breyse, G.B. Diette, E.C. Matsui, A.M. Butz, N.N. Hansel, M.C. McCormack, Indoor air pollution and asthma in children, *Proc. Am. Thorac. Soc.* 7 (2010) 102–106. <https://doi.org/10.1513/pats.200908-083RM>.
- [2] S. and M. Hegde Kyeong T. and Moore, James and Lundrigan, Philip and Patwari, Neal and Collingwood, Scott and Balch, Alfred and Kelly, Kerry E., Indoor Household Particulate Matter Measurements Using a Network of Low-cost Sensors, *Aerosol Air Qual. Res.* 20 (2020) 381–394. <https://doi.org/10.4209/aaqr.2019.01.0046>.
- [3] J.O. Anderson, J.G. Thundiyil, A. Stolbach, Clearing the Air: A Review of the Effects of Particulate Matter Air Pollution on Human Health, *J. Med. Toxicol.* 8 (2012) 166–175. <https://doi.org/10.1007/s13181-011-0203-1>.
- [4] M. Badura, P. Batog, A. Drzeniecka-Osiadacz, P. Modzel, Evaluation of Low-Cost Sensors for Ambient PM_{2.5} Monitoring, *J. Sensors.* 2018 (2018) 5096540. <https://doi.org/10.1155/2018/5096540>.
- [5] P. de Prado Bert, E.M.H. Mercader, J. Pujol, J. Sunyer, M. Mortamais, The Effects of Air Pollution on the Brain: a Review of Studies Interfacing Environmental Epidemiology and Neuroimaging, *Curr. Environ. Heal. Reports.* 5 (2018) 351–364. <https://doi.org/10.1007/s40572-018-0209-9>.
- [6] S.S. Patra, Prediction of Indoor PM_{2.5} concentrations using Support Vector Regression, *Int. J. Adv. Res. IDEAS Innov. Technol.* 5 (2019) 187–190.
- [7] A.L. Northcross, R.J. Edwards, M.A. Johnson, Z.-M. Wang, K. Zhu, T. Allen, K.R. Smith, A low-cost particle counter as a realtime fine-particle mass monitor, *Environ. Sci. Process. Impacts.* 15 (2013) 433–439. <https://doi.org/10.1039/C2EM30568B>.
- [8] P.J. Dacunto, K.-C. Cheng, V. Acevedo-Bolton, R.-T. Jiang, N.E. Klepeis, J.L. Repace, W.R. Ott, L.M. Hildemann, Real-time particle monitor calibration factors and PM_{2.5} emission factors for multiple indoor sources, *Environ. Sci. Process. Impacts.* 15 (2013) 1511–1519. <https://doi.org/10.1039/C3EM00209H>.
- [9] K. Isiugo, R. Jandarov, J. Cox, P. Ryan, N. Newman, S.A. Grinshpun, R. Indugula, S. Vesper, T. Reponen, Indoor particulate matter and lung function in children, *Sci. Total Environ.* 663 (2019) 408–417. <https://doi.org/10.1016/j.scitotenv.2019.01.309>.
- [10] F.D. Pope, M. Gatari, D. Ng, A. Poynter, R. Blake, Airborne particulate matter monitoring in Kenya using calibrated low-cost sensors, *Atmos. Chem. Phys.* 18 (2018) 15403–15418. <https://doi.org/10.5194/acp-18-15403-2018>.
- [11] L.R. Crilley, M. Shaw, R. Pound, L.J. Kramer, R. Price, S. Young, A.C. Lewis, F.D. Pope, Evaluation of a low-cost optical particle counter (Alphasense OPC-N2) for ambient air monitoring, *Atmos. Meas. Tech.* 11 (2018) 709–720. <https://doi.org/10.5194/amt-11-709-2018>.
- [12] N. Zikova, M. Masiol, D. Chalupa, D. Rich, A. Ferro, P. Hopke, Estimating Hourly Concentrations of PM_{2.5} across a Metropolitan Area Using Low-Cost Particle Monitors, *Sensors.* 17 (2017) 1922. <https://doi.org/10.3390/s17081922>.
- [13] A.C. Rai, P. Kumar, F. Pilla, A.N. Skouloudis, S. Di Sabatino, C. Ratti, A. Yasar, D. Rickerby, End-user perspective of low-cost sensors for outdoor air pollution monitoring, *Sci. Total Environ.* 607–608 (2017) 691–705. <https://doi.org/10.1016/j.scitotenv.2017.06.266>.
- [14] Z. Wang, W.W. Delp, B.C. Singer, Performance of low-cost indoor air quality monitors for PM_{2.5} and PM₁₀ from residential sources, *Build. Environ.* 171 (2020) 106654. <https://doi.org/10.1016/j.buildenv.2020.106654>.
- [15] Y. Wang, J. Li, H. Jing, Q. Zhang, J. Jiang, P. Biswas, Laboratory Evaluation and Calibration of Three Low-Cost Particle Sensors for Particulate Matter Measurement, *Aerosol Sci. Technol.* 49 (2015) 1063–1077. <https://doi.org/10.1080/02786826.2015.1100710>.
- [16] A.K. Amegah, S. Agyei-Mensah, Urban air pollution in Sub-Saharan Africa: Time for action, *Environ. Pollut.* 220 (2017) 738–743. <https://doi.org/10.1016/j.envpol.2016.09.042>.
- [17] H.-Y. Liu, P. Schneider, R. Haugen, M. Vogt, Performance Assessment of a Low-Cost PM_{2.5} Sensor for a near Four-Month Period in Oslo, Norway, *Atmosphere (Basel).* 10 (2019) 41. <https://doi.org/10.3390/atmos10020041>.
- [18] B.I. Magi, C. Cupini, J. Francis, M. Green, C. Hauser, Evaluation of PM_{2.5} measured in an urban setting using a low-cost optical particle counter and a Federal Equivalent Method Beta Attenuation Monitor, *Aerosol Sci. Technol.* 54 (2020) 147–159. <https://doi.org/10.1080/02786826.2019.1619915>.
- [19] L. Bai, L. Huang, Z. Wang, Q. Ying, J. Zheng, X. Shi, J. Hu, Long-term field Evaluation of Low-cost Particulate Matter Sensors in Nanjing, *Aerosol Air Qual. Res.* 20 (2020) 242–253. <https://doi.org/10.4209/aaqr.2018.11.0424>.
- [20] S. Sousan, K. Koehler, L. Hallett, T.M. Peters, Evaluation of the Alphasense optical particle counter (OPC-N2) and the Grimm portable aerosol spectrometer (PAS-1.108), *Aerosol Sci. Technol.* 50 (2016) 1352–1365. <https://doi.org/10.1080/02786826.2016.1232859>.
- [21] N. Zimmerman, A.A. Presto, S.P.N. Kumar, J. Gu, A. Hauriuk, E.S. Robinson, A.L. Robinson, A machine learning calibration model using random forests to improve sensor performance for lower-cost air quality monitoring, *Atmos. Meas. Tech.* 11 (2018) 291–313. <https://doi.org/10.5194/amt-11-291-2018>.
- [22] D.M. Holstius, A. Pillarisetti, K.R. Smith, E. Seto, Field calibrations of a low-cost aerosol sensor at a regulatory monitoring site in California, *Atmos. Meas. Tech.* 7 (2014) 1121–1131. <https://doi.org/10.5194/amt-7-1121-2014>.
- [23] M. Gao, J. Cao, E. Seto, A distributed network of low-cost continuous reading sensors to measure spatiotemporal variations of PM_{2.5} in Xi'an, China, *Environ. Pollut.* 199 (2015) 56–65. <https://doi.org/10.1016/j.envpol.2015.01.013>.
- [24] M. Jovašević-Stojanović, A. Bartonova, D. Topalović, I. Lazović, B. Pokrić, Z. Ristovski, On the use of small and

- cheaper sensors and devices for indicative citizen-based monitoring of respirable particulate matter, *Environ. Pollut.* 206 (2015) 696–704. <https://doi.org/10.1016/j.envpol.2015.08.035>.
- [25] S. Steinle, S. Reis, C.E. Sabel, S. Semple, M.M. Twigg, C.F. Braban, S.R. Leeson, M.R. Heal, D. Harrison, C. Lin, H. Wu, Personal exposure monitoring of PM_{2.5} in indoor and outdoor microenvironments, *Sci. Total Environ.* 508 (2015) 383–394. <https://doi.org/10.1016/j.scitotenv.2014.12.003>.
- [26] M. Alvarado, F. Gonzalez, A. Fletcher, A. Doshi, Towards the Development of a Low Cost Airborne Sensing System to Monitor Dust Particles after Blasting at Open-Pit Mine Sites, *Sensors*. 15 (2015) 19667–19687. <https://doi.org/10.3390/s150819667>.
- [27] L.R. Crilley, A. Singh, L.J. Kramer, M.D. Shaw, M.S. Alam, J.S. Apte, W.J. Bloss, L. Hildebrandt Ruiz, P. Fu, W. Fu, S. Gani, M. Gatari, E. Ilyinskaya, A.C. Lewis, D. Ng'ang'a, Y. Sun, R.C.W. Whitty, S. Yue, S. Young, F.D. Pope, Effect of aerosol composition on the performance of low-cost optical particle counter correction factors, *Atmos. Meas. Tech.* 13 (2020) 1181–1193. <https://doi.org/10.5194/amt-13-1181-2020>.
- [28] I. Han, E. Symanski, T.H. Stock, Feasibility of using low-cost portable particle monitors for measurement of fine and coarse particulate matter in urban ambient air, *J. Air Waste Manage. Assoc.* 67 (2017) 330–340. <https://doi.org/10.1080/10962247.2016.1241195>.
- [29] K.E. Kelly, J. Whitaker, A. Petty, C. Widmer, A. Dybwad, D. Sleeth, R. Martin, A. Butterfield, Ambient and laboratory evaluation of a low-cost particulate matter sensor, *Environ. Pollut.* 221 (2017) 491–500. <https://doi.org/10.1016/j.envpol.2016.12.039>.
- [30] B. Feenstra, V. Papapostolou, S. Hasheminassab, H. Zhang, B. Der Boghossian, D. Cocker, A. Polidori, Performance evaluation of twelve low-cost PM_{2.5} sensors at an ambient air monitoring site, *Atmos. Environ.* 216 (2019) 116946. <https://doi.org/10.1016/j.atmosenv.2019.116946>.
- [31] A. Thomas, J. Gebhart, Correlations between gravimetry and light scattering photometry for atmospheric aerosols, *Atmos. Environ.* 28 (1994) 935–938. [https://doi.org/10.1016/1352-2310\(94\)90251-8](https://doi.org/10.1016/1352-2310(94)90251-8).
- [32] K.A. Koehler, T.M. Peters, New Methods for Personal Exposure Monitoring for Airborne Particles, *Curr. Environ. Heal. Reports*. 2 (2015) 399–411. <https://doi.org/10.1007/s40572-015-0070-z>.
- [33] S. Sousan, K. Koehler, G. Thomas, J.H. Park, M. Hillman, A. Halterman, T.M. Peters, Inter-comparison of low-cost sensors for measuring the mass concentration of occupational aerosols, *Aerosol Sci. Technol.* 50 (2016) 462–473. <https://doi.org/10.1080/02786826.2016.1162901>.
- [34] T. Wu, B.E. Boor, Urban Aerosol Size Distributions: A Global Perspective, *Atmos. Chem. Phys. Discuss.* (2020) 1–83. <https://doi.org/10.5194/acp-2020-92>.
- [35] J. Zhao, W. Birmili, B. Wehner, A. Daniels, K. Weinhold, L. Wang, M. Merkel, S. Kecorius, T. Tuch, U. Franck, T. Hussein, A. Wiedensohler, Particle Mass Concentrations and Number Size Distributions in 40 Homes in Germany: Indoor-to-Outdoor Relationships, Diurnal and Seasonal Variation, *Aerosol Air Qual. Res.* (2020). <https://doi.org/10.4209/aaqr.2019.09.0444>.
- [36] T. Fazli, Y. Zeng, B. Stephens, Fine and ultrafine particle removal efficiency of new residential HVAC filters, *Indoor Air*. 29 (2019) 656–669. <https://doi.org/10.1111/ina.12566>.
- [37] C.M. Long, H.H. Suh, P.J. Catalano, P. Koutrakis, Using Time- and Size-Resolved Particulate Data To Quantify Indoor Penetration and Deposition Behavior, *Environ. Sci. Technol.* 35 (2001) 2089–2099. <https://doi.org/10.1021/es001477d>.
- [38] W.J. Riley, T.E. McKone, A.C.K. Lai, W.W. Nazaroff, Indoor Particulate Matter of Outdoor Origin: Importance of Size-Dependent Removal Mechanisms, *Environ. Sci. Technol.* 36 (2002) 200–207. <https://doi.org/10.1021/es010723y>.
- [39] M. Si, Y. Xiong, S. Du, K. Du, Evaluation and calibration of a low-cost particle sensor in ambient conditions using machine-learning methods, *Atmos. Meas. Tech.* 13 (2020) 1693–1707. <https://doi.org/10.5194/amt-13-1693-2020>.
- [40] B.G. Loh, G.H. Choi, Calibration of Portable Particulate Matter–Monitoring Device using Web Query and Machine Learning, *Saf. Health Work.* 10 (2019) 452–460. <https://doi.org/10.1016/j.shaw.2019.08.002>.
- [41] C. Chen, C. Kuo, S. Chen, C. Lin, J. Chue, Y. Hsieh, C. Cheng, C. Wu, C. Huang, Calibration of Low-Cost Particle Sensors by Using Machine-Learning Method, in: 2018 IEEE Asia Pacific Conf. Circuits Syst., 2018: pp. 111–114. <https://doi.org/10.1109/APCCAS.2018.8605619>.
- [42] Y. Wang, Y. Du, J. Wang, T. Li, Calibration of a low-cost PM_{2.5} monitor using a random forest model, *Environ. Int.* 133 (2019) 105161. <https://doi.org/10.1016/j.envint.2019.105161>.
- [43] A. Di Antonio, O.A.M. Popoola, B. Ouyang, J. Saffell, R.L. Jones, Developing a relative humidity correction for low-cost sensors measuring ambient particulate matter, *Sensors (Switzerland)*. 18 (2018) 2790. <https://doi.org/10.3390/s18092790>.
- [44] Alphasense Ltd, Alphasense User Manual OPC-N2 Optical Particle Counter Issue 3, 2015. www.alphasense.com (accessed May 4, 2020).
- [45] RaspberryPi, Raspberry Pi 3 Model B+ – Raspberry Pi, (n.d.). <https://www.raspberrypi.org/products/raspberry-pi-3-model-b-plus/> (accessed May 4, 2020).
- [46] D.H. Hagan, A. Tolmie, J. Trochim, py-opc: operate the Alphasense OPC-N2 from a raspberry pi or other popular microcontrollers/microcomputers, *J. Open Source Softw.* 3 (2018) 782. <https://doi.org/10.21105/joss.00782>.
- [47] TSI, TSI Optical Particle Sizer 3330 (n.d.). <https://tsi.com/products/particle-sizers/particle-size-spectrometers/optical-particle-sizer-3330/> (accessed May 5, 2020).

- [48] O. Dubovik, B. Holben, T.F. Eck, A. Smirnov, Y.J. Kaufman, M.D. King, D. Tanré, I. Slutsker, Variability of Absorption and Optical Properties of Key Aerosol Types Observed in Worldwide Locations, *J. Atmos. Sci.* 59 (2002) 590–608. [https://doi.org/10.1175/1520-0469\(2002\)059<0590:VOAAOP>2.0.CO;2](https://doi.org/10.1175/1520-0469(2002)059<0590:VOAAOP>2.0.CO;2).
- [49] J. Pagels, A. Wierzbicka, E. Nilsson, C. Isaxon, A. Dahl, A. Gudmundsson, E. Swietlicki, M. Bohgard, Chemical composition and mass emission factors of candle smoke particles, *J. Aerosol Sci.* 40 (2009) 193–208. <https://doi.org/10.1016/j.jaerosci.2008.10.005>.
- [50] R. Beltman, R. van der Bruggen, G.C. Dol, S. Fritsche, Cooking Fat Product With Improved Spattering Behaviour, US20080305237A1, 2004.
- [51] TSI, TSI Scanning Mobility Particle Sizer Spectrometer 3938 (n.d.). <https://tsi.com/products/particle-sizers/particle-size-spectrometers/scanning-mobility-particle-sizer-spectrometer-3938/> (accessed May 5, 2020).
- [52] M.R. Stolzenburg, P.H. McMurry, Method to assess performance of scanning mobility particle sizer (SMPS) instruments and software, *Aerosol Sci. Technol.* 52 (2018) 609–613. <https://doi.org/10.1080/02786826.2018.1455962>.
- [53] TSI, Scanning Mobility Particle Sizer (SMPS) Spectrometer Model 3938 Operation And Service Manual, 2016.
- [54] S.L. Caskey, E.A. Groll, Hybrid air-hydrionic HVAC performance in a residential net-zero energy retrofit, *Energy Build.* 158 (2018) 342–355. <https://doi.org/10.1016/j.enbuild.2017.10.003>.
- [55] S.L. Caskey, E.J. Bowler, E.A. Groll, Analysis on a net-zero energy renovation of a 1920s vintage home, *Sci. Technol. Built Environ.* 22 (2016) 1060–1073. <https://doi.org/10.1080/23744731.2016.1216226>.
- [56] Y. Zou, M. Young, M. Wickey, A. May, J.D. Clark, Response of eight low-cost particle sensors and consumer devices to typical indoor emission events in a real home (ASHRAE 1756-RP), *Sci. Technol. Built Environ.* 26 (2020) 237–249. <https://doi.org/10.1080/23744731.2019.1676094>.
- [57] J.H. Seinfeld, S.N. Pandis, *Atmospheric Chemistry and Physics: From Air Pollution to Climate Change*, 3rd ed., WILEY, 2016.
- [58] S.J. Gershman, D.M. Blei, A tutorial on Bayesian nonparametric models, *J. Math. Psychol.* 56 (2012) 1–12. <https://doi.org/10.1016/j.jmp.2011.08.004>.
- [59] C.E. Rasmussen, Gaussian Processes in Machine Learning BT - Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia, February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures, in: O. Bousquet, U. von Luxburg, G. Rätsch (Eds.), Springer Berlin Heidelberg, Berlin, Heidelberg, 2004: pp. 63–71. https://doi.org/10.1007/978-3-540-28650-9_4.
- [60] E. Schulz, M. Speekenbrink, A. Krause, A tutorial on Gaussian process regression: Modelling, exploring, and exploiting functions, *BioRxiv.* (2017) 95190. <https://doi.org/10.1101/095190>.
- [61] F. Jäkel, B. Schölkopf, F.A. Wichmann, A tutorial on kernel methods for categorization, *J. Math. Psychol.* 51 (2007) 343–358. <https://doi.org/10.1016/j.jmp.2007.06.002>.
- [62] C.E. Rasmussen, H. Nickisch, Gaussian Processes for Machine Learning (GPML) Toolbox, *J. Mach. Learn. Res.* 11 (2010) 3011–3015. <http://www.kyb.tuebingen.mpg.de/bs/people/carl/code/minimize/>. (accessed May 15, 2020).
- [63] S. Kim, H. Kim, A new metric of absolute percentage error for intermittent demand forecasts, *Int. J. Forecast.* 32 (2016) 669–679. <https://doi.org/10.1016/j.ijforecast.2015.12.003>.
- [64] D. Dinu, M. Fayolas, M. Jacquet, E. Leguy, J. Slavinski, N. Houel, Accuracy of Postural Human-motion Tracking Using Miniature Inertial Sensors, in: *Procedia Eng.*, Elsevier Ltd, 2016: pp. 655–658. <https://doi.org/10.1016/j.proeng.2016.06.266>.
- [65] F. Schoonjans, A. Zalata, C.E. Depuydt, F.H. Comhaire, MedCalc: a new computer program for medical statistics, *Comput. Methods Programs Biomed.* 48 (1995) 257–262. [https://doi.org/10.1016/0169-2607\(95\)01703-8](https://doi.org/10.1016/0169-2607(95)01703-8).
- [66] T. Kapwata, B. Language, S. Piketh, C. Wright, Variation of Indoor Particulate Matter Concentrations and Association with Indoor/Outdoor Temperature: A Case Study in Rural Limpopo, South Africa, *Atmosphere (Basel)*. 9 (2018) 124. <https://doi.org/10.3390/atmos9040124>.
- [67] A.J. Wheeler, L.A. Wallace, J. Kearney, K. Van Ryswyk, H. You, R. Kulka, J.R. Brook, X. Xu, Personal, Indoor, and Outdoor Concentrations of Fine and Ultrafine Particles Using Continuous Monitors in Multiple Residences, *Aerosol Sci. Technol.* 45 (2011) 1078–1089. <https://doi.org/10.1080/02786826.2011.580798>.
- [68] M. Braniš, P. Řezáčová, M. Domasová, The effect of outdoor air and indoor human activity on mass concentrations of PM₁₀, PM_{2.5}, and PM₁ in a classroom, *Environ. Res.* 99 (2005) 143–149. <https://doi.org/10.1016/j.envres.2004.12.001>.
- [69] C. Isaxon, A. Gudmundsson, E.Z. Nordin, L. Lönnblad, A. Dahl, G. Wieslander, M. Bohgard, A. Wierzbicka, Contribution of indoor-generated particles to residential exposure, *Atmos. Environ.* 106 (2015) 458–466. <https://doi.org/10.1016/j.atmosenv.2014.07.053>.
- [70] L. Wallace, Indoor Sources of Ultrafine and Accumulation Mode Particles: Size Distributions, Size-Resolved Concentrations, and Source Strengths, *Aerosol Sci. Technol.* 40 (2006) 348–360. <https://doi.org/10.1080/02786820600612250>.
- [71] Y. Li, B. Hannaford, Gaussian Process Regression for Sensorless Grip Force Estimation of Cable-Driven Elongated Surgical Instruments, *IEEE Robot. Autom. Lett.* 2 (2017) 1312–1319. <https://doi.org/10.1109/LRA.2017.2666420>.

Tables

Table 1. Summary of selected studies evaluating low-cost OPCs.

Study	Low-Cost OPC(s) Evaluated	Reference Instrument(s) Used	Summary of Results
Northcross et al. [7]	Dylos 1700 (Dylos Corp.) (modified by the authors)	DustTrak 8520 (TSI Inc.)	Different particle types were added into an experimental chamber and the performance of the Dylos and DustTrak were compared. R^2 of 0.99 was reported for PSL spheres and $(\text{NH}_4)_2\text{SO}_4$. Woodsmoke reported R^2 of 0.97-0.98 and ambient particle sampling reported R^2 of 0.81-0.99.
Holstius et al. [22]	PPD42NS (Shinyei Tech. Co.)	BAM 1020 (Met One Instruments Inc.), DustTrak 8530 (TSI Inc.), PAS 1.108 (GRIMM GmbH), Dylos 1700 (Dylos Corp.)	At a 1 h scale, the R^2 reported ranged from 0.55-0.60, 0.87-0.92, 0.90-0.94, and 0.64-0.80 for comparison against BAM, Dylos, PAS 1.108, and DustTrak, respectively. Improvements were observed for 24 h scale. Linear corrections explained 60% variance in 1 h and 72% variance in 24 h data.
Gao et al. [23]	PPD42NS (Shinyei Tech. Co.)	APS 3321 (TSI Inc.)	The performance of the Shinyei PPD42NS against the APS was evaluated with PSL spheres and ASHRAE Test Dust. For concentrations less than $50 \mu\text{g m}^{-3}$, a linear model captured the sensor response, and for higher concentrations, a non-linear function captured the same.
Jovašević-Stojanović et al. [24]	Dylos 1700 (Dylos Corp.)	OPS 3330 (TSI Inc.), PAS 1.108 (GRIMM GmbH)	The Dylos reported R^2 values between 0.88-0.99 for indoor air sampling. The outdoor air evaluations reported lower R^2 values, ranging from 0.74-0.84.
Steinle et al. [25]	Dylos 1700 (Dylos Corp.)	TEOM-FDMS (Thermo Fisher Scientific)	The performance of the Dylos was validated against the TEOM at ambient monitoring sites. The R^2 values reported were 0.9 and 0.7 at rural and urban sites, respectively.
Alvarado et al. [26]	GP2Y10 (Sharp), DMS501A (Samyoung Elec. Co.)	DustTrak 8520 (TSI Inc.)	The R^2 for the Sharp GP2Y10 against the DustTrak ranged from 0.92-0.98. However, the Samyoung DSM501A showed relatively lower values of R^2 (approx. 0.5).
Zikova et al. [12]	DMS501A (Samyoung Elec. Co.)	PAS 1.109 (GRIMM GmbH)	The DMS501A reported R^2 values ranging from 0.07-0.29 for 1 min data. When the sampling interval in the study was increased, higher R^2 values were observed (0.15-0.46).
Crilley et al. [11]	OPC-N2 (Alphasense Ltd.)	TEOM-FDMS (Thermo Fisher Scientific), PAS 1.108 (GRIMM GmbH)	The R^2 values for $\text{PM}_{2.5}$ ranged from 0.7-0.74 and 0.71-0.74 for evaluation against TEOM and PAS 1.108, respectively. For PM_{10} , the values ranged from 0.64-0.67 and 0.66-0.72. Linear corrections were applied for $\text{RH} < 85\%$, and κ -Köhler corrections were applied for $\text{RH} > 85\%$.
Han et al. [28]	Dylos 1700 (Dylos Corp.)	Mini-LAS 11-R (GRIMM GmbH)	The R^2 for overall correlation between the Dylos and the Mini-LAS was 0.78. The concentration values reported by both were closer when the $\text{RH} < 60\%$, compared to elevated RH .
Badura et al. [4]	SDS011 (Nova Fitness), ZH03A (Winsen), PMS7003 (Plantower), OPC-N2 (Alphasense Ltd.)	TEOM 1400a (Thermo Fisher Scientific)	The R^2 values ranged from 0.79-0.90, 0.70-0.89, 0.80-0.93, and 0.43-0.69 for SDS011, ZH03A, PMS7003, and OPC-N2, respectively. Overestimation in the sensors were observed for $\text{RH} > 80\%$.
Liu et al. [17]	SDS011 (Nova Fitness)	TEOM 1405 (Thermo Fisher Scientific)	The R^2 values reported for SDS011 against the TEOM were 0.55-0.71. When the $\text{RH} > 80\%$, it negatively impacted the OPC's response. The inclusion of RH and temperature in the calibration increased the R^2 values.
Magi et al. [18]	PA-II (Purple Air)	BAM 1022 (Met One Instruments Inc.)	The reported R^2 for comparison of the PA-II against the BAM was 0.54. A multiple linear regression model was used to correct the sensor. The results showed a 27-57% improvement in the accuracy of the sensor.
Bai et al. [19]	PPD42NS (Shinyei Tech. Co.)	BAM 1020 (Met One Instruments Inc.)	The PPD42NS was evaluated against BAM for long-term outdoor use. R^2 values ranged from 0.71-0.84.

Table 2. Statistical summary of the comparison of calculated and firmware $PV_{2.5}$ and PV_{10} concentrations reported by OPC 1 and OPC 2 with the reference instruments (OPS and SMPS) during occupied periods at the Purdue ReNEW House.

Concentration Type	Slope		r		R^2		MAPE (%)		p -value	
	OPC 1	OPC 2	OPC 1	OPC 2	OPC 1	OPC 2	OPC 1	OPC 2	OPC 1	OPC 2
Calculated $PV_{2.5}$	0.124	0.123	0.79	0.79	0.63	0.64	79.81	78.97	3.76E-32	6.20E-25
Calculated PV_{10}	0.229	0.211	0.77	0.78	0.59	0.60	67.67	67.59	5.64E-24	2.50E-25
Firmware $PV_{2.5}$	0.511	0.502	0.77	0.77	0.60	0.59	41.37	40.02	2.53E-03	2.67E-03
Firmware PV_{10}	0.437	0.447	0.75	0.75	0.57	0.56	47.50	52.53	4.59E-12	3.21E-20

Table 3. Statistical summary of the comparison of calculated and firmware (uncalibrated) and corrected $PV_{2.5}$ and PV_{10} concentrations reported by OPC 1 and OPC 2 with the reference instruments (OPS and SMPS) during the testing set.

Concentration Type	Slope		r		R^2		MAPE (%)		p-value	
	OPC 1	OPC 2	OPC 1	OPC 2	OPC 1	OPC 2	OPC 1	OPC 2	OPC 1	OPC 2
Calculated $PV_{2.5}$	0.081	0.086	0.66	0.67	0.44	0.44	57.3	57.1	1.34E-05	1.21E-05
Calculated PV_{10}	0.172	0.171	0.64	0.65	0.41	0.42	66.3	69.2	2.40E-06	3.20E-06
Firmware $PV_{2.5}$	0.308	0.331	0.59	0.60	0.35	0.36	37.7	35.4	1.79E-02	9.10E-03
Firmware PV_{10}	0.654	0.683	0.83	0.83	0.69	0.69	60.5	63.3	2.82E-06	2.66E-06
Corrected $PV_{2.5}$	0.982	0.984	0.99	0.99	0.98	0.98	23.4	25.1	6.50E-01	3.20E-01
Corrected PV_{10}	0.978	0.966	0.99	0.98	0.98	0.97	19.2	20.6	8.20E-01	6.40E-01

Figures

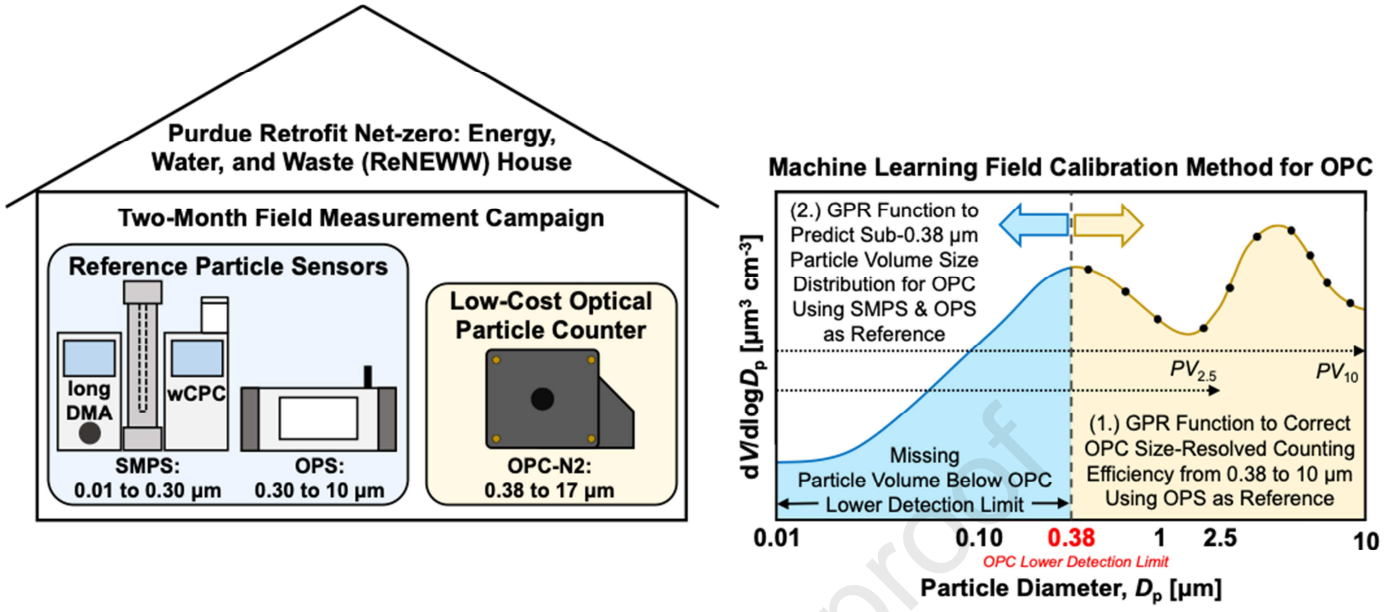


Figure 1. (left) Overview of indoor particle measurements at the Purdue ReNEWW House with SMPS, OPS, and low-cost OPCs and (right) illustration of two-component machine learning field calibration method for low-cost OPCs.

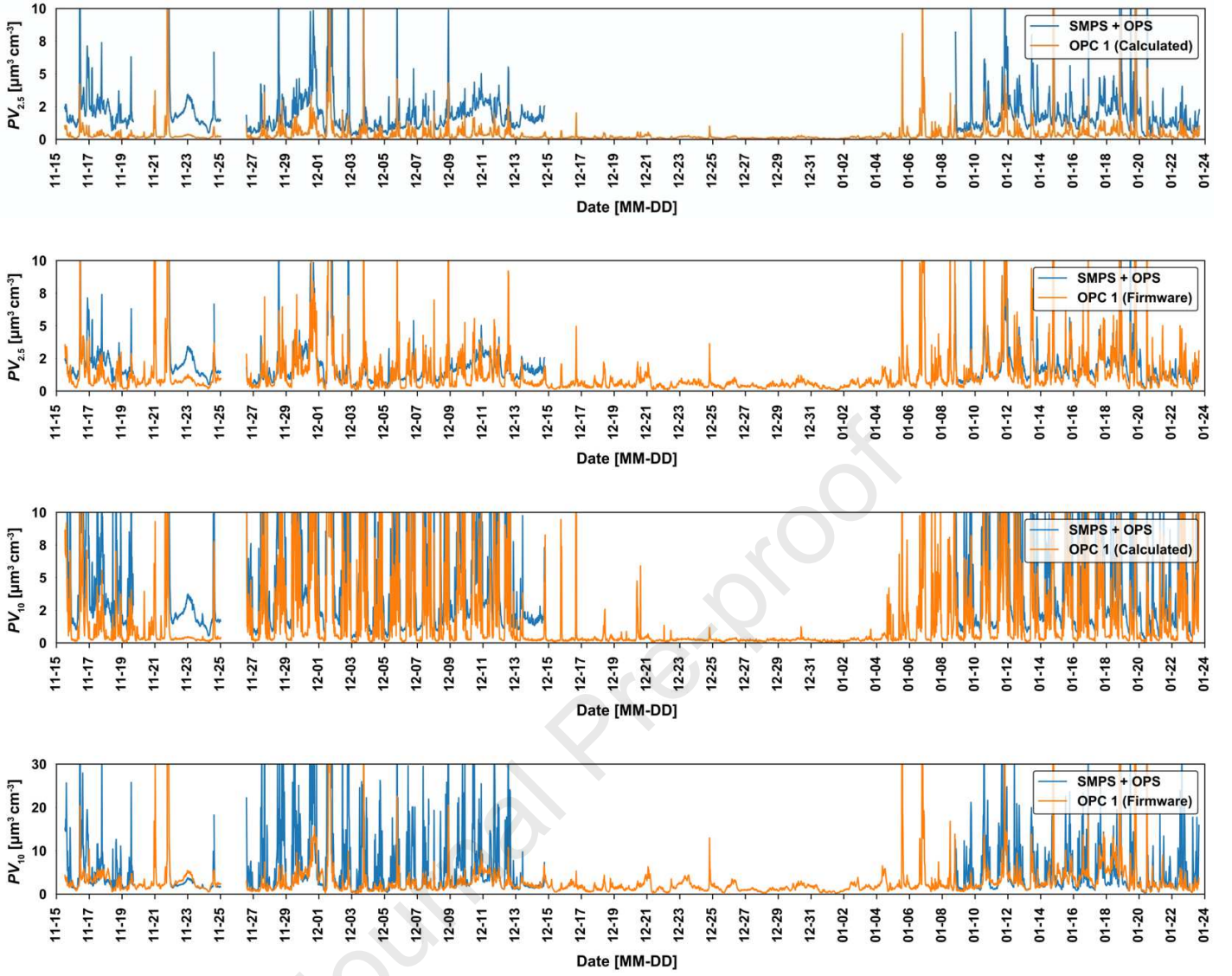


Figure 2. Comparison of $PV_{2.5}$ and PV_{10} concentration time-series reported by OPC 1 (both calculated and firmware) with $PV_{2.5}$ and PV_{10} concentration time-series reported by the reference instruments (OPS and SMPS). The Purdue ReNEW House remained unoccupied for parts of Thanksgiving vacation (November 21 to 24, 2018) and the entirety of winter break (December 15, 2018 to January 05, 2019; OPS and SMPS not used). Data for all instruments was not recorded from November 25 to 26, 2018 due to data logging and power outage issues and from December 08 to 15, 2018 for OPC 2 due to a data logging issue.

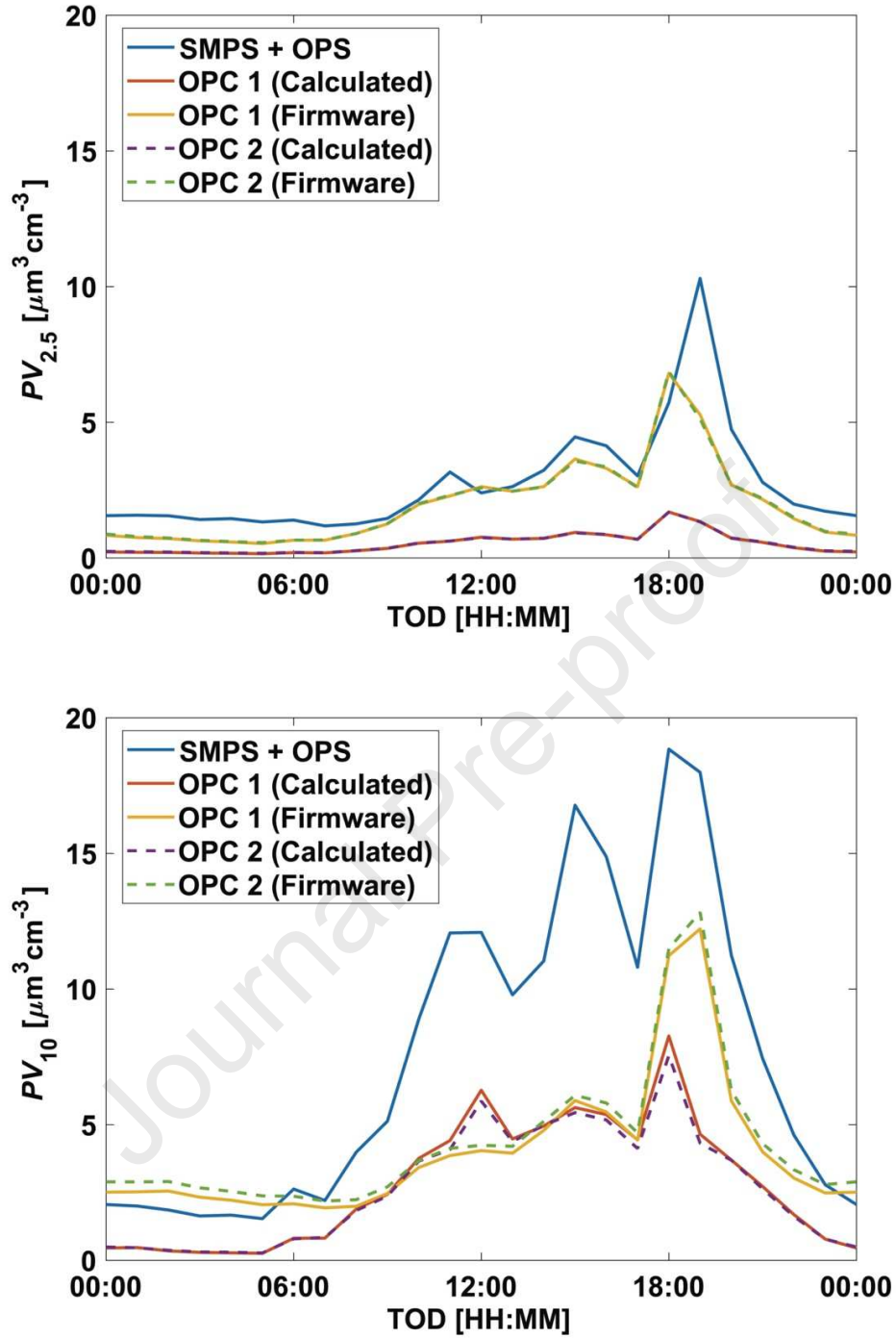


Figure 3. Comparison of diurnal profiles in the mean $PV_{2.5}$ and PV_{10} concentrations as reported by OPC 1 and OPC 2 (both calculated and firmware) with those reported by the reference instruments (OPS and SMPS) during occupied periods at the Purdue ReNEW House.

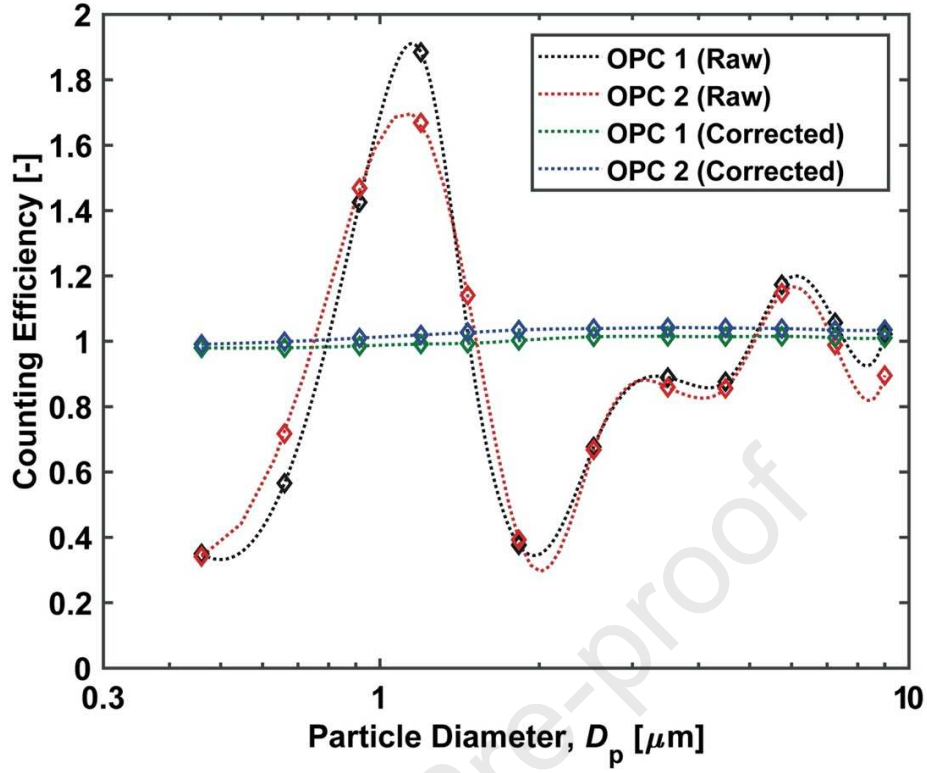


Figure 4. Raw and corrected size-resolved counting efficiencies for OPC 1 and OPC 2. Each marker represents the mean counting efficiency per bin, shown at the mean diameter of the bin ($D_{bin,mean}$). The raw size-resolved counting efficiencies are calculated for the entire measurement campaign (training and testing sets) and the corrected size-resolved counting efficiencies are calculated for the testing set.

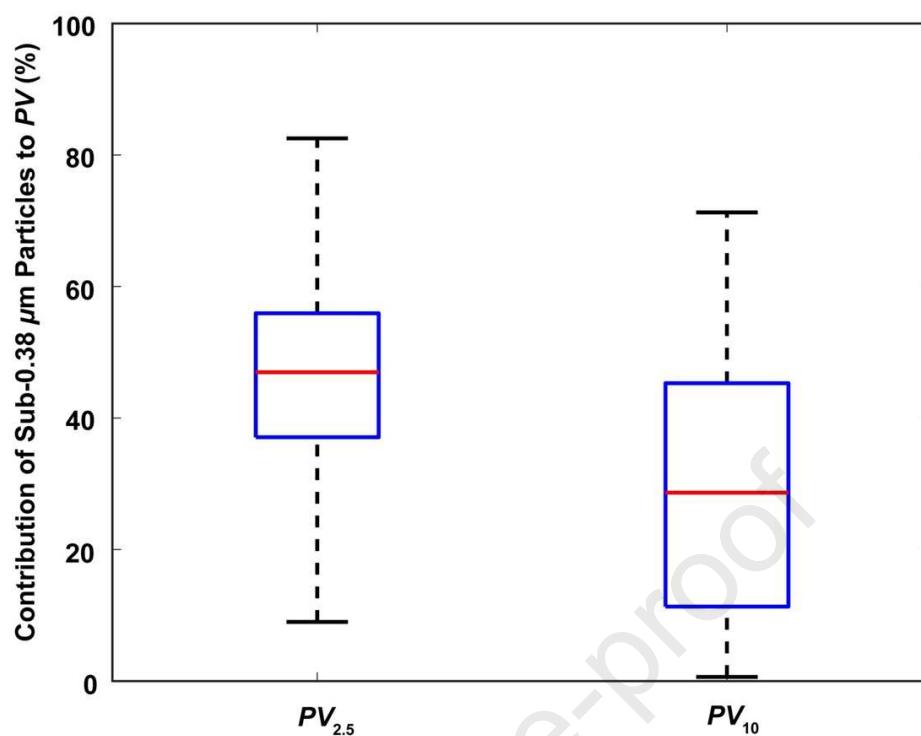


Figure 5. The contribution of sub-0.38 µm particles to $PV_{2.5}$ and PV_{10} concentrations as measured by the OPS and SMPS over the entirety of the measurement campaign at the Purdue ReNEW House.

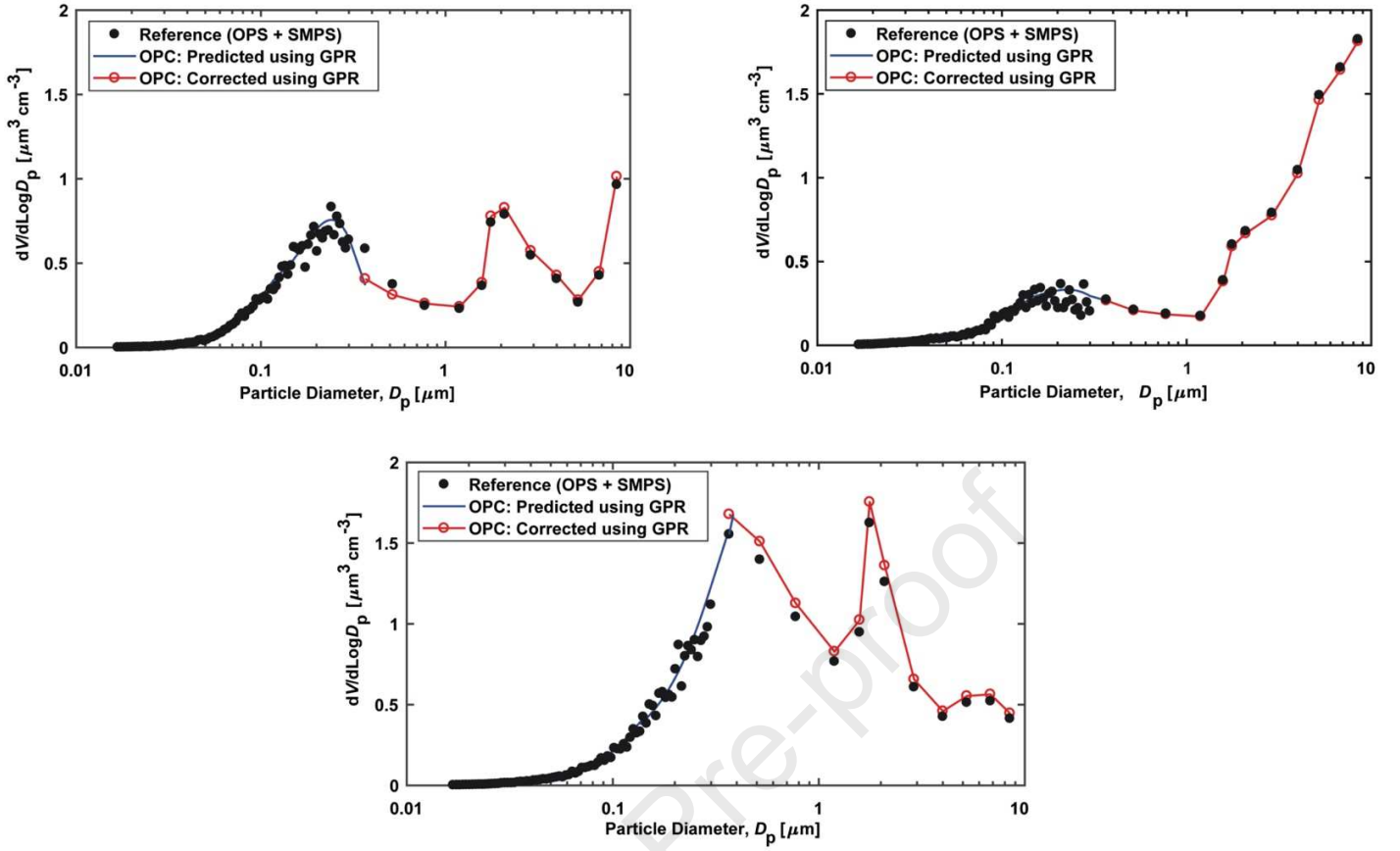


Figure 6. Comparison of reference particle volume size distributions as measured by the OPS and SMPS with the estimated particle volume size distributions obtained by the OPCs following implementation of the machine learning calibration method. The top two distributions are from OPC 1 and the bottom distribution is taken from OPC 2. The red line represents the counting efficiency correction with the first GPR function and the blue line represents the sub-0.38 μm volume prediction with the second GPR function. The randomly selected volume distributions are determined for 30 min windows on January 22, 2019 at 23:30 (top left), January 20, 2019 at 16:00 (top right), and January 16, 2019 at 10:00 (bottom).

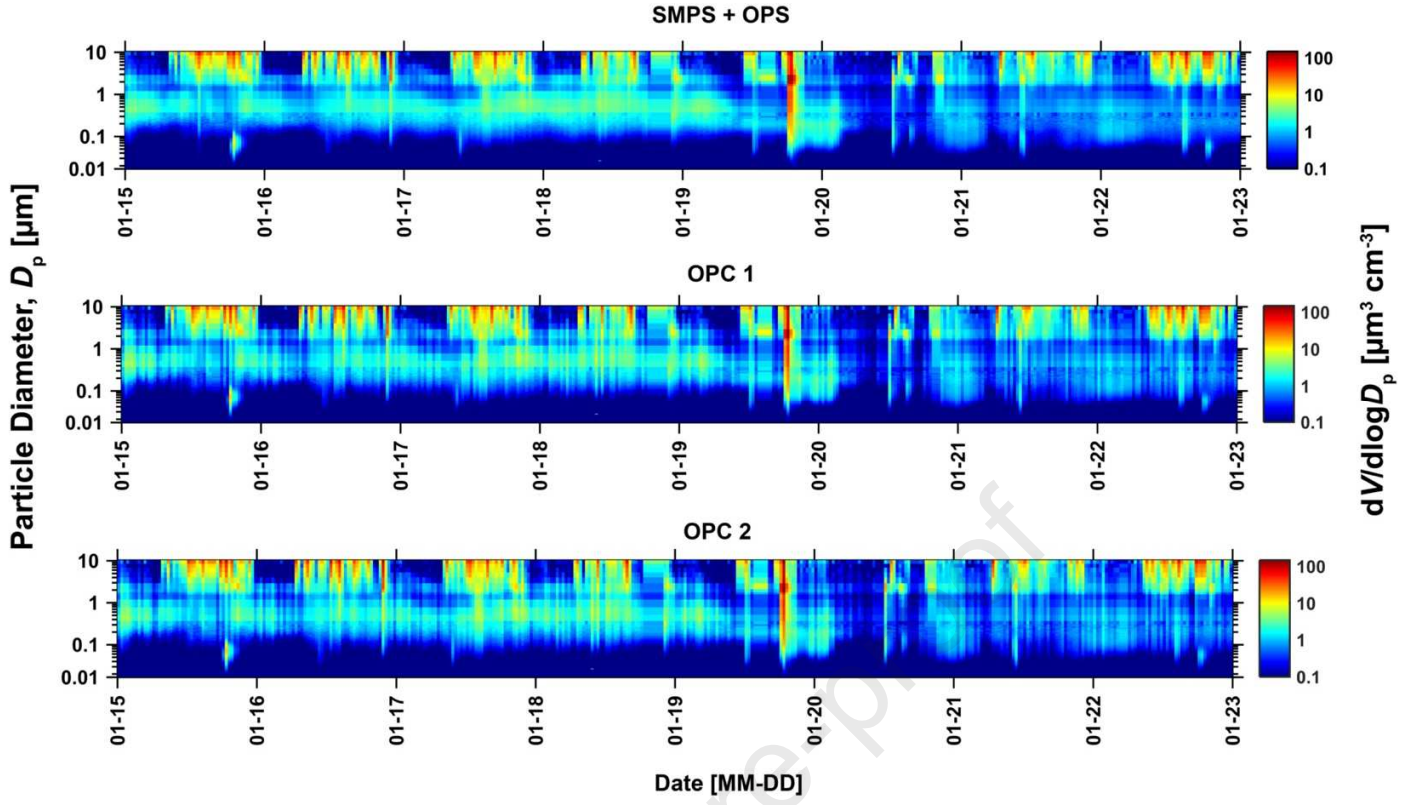


Figure 7. Comparison of reference particle volume size distribution time-series as measured by the OPS and SMPS with the estimated particle volume size distribution time-series obtained by OPC 1 and OPC 2 following implementation of the machine learning calibration method for the testing set.

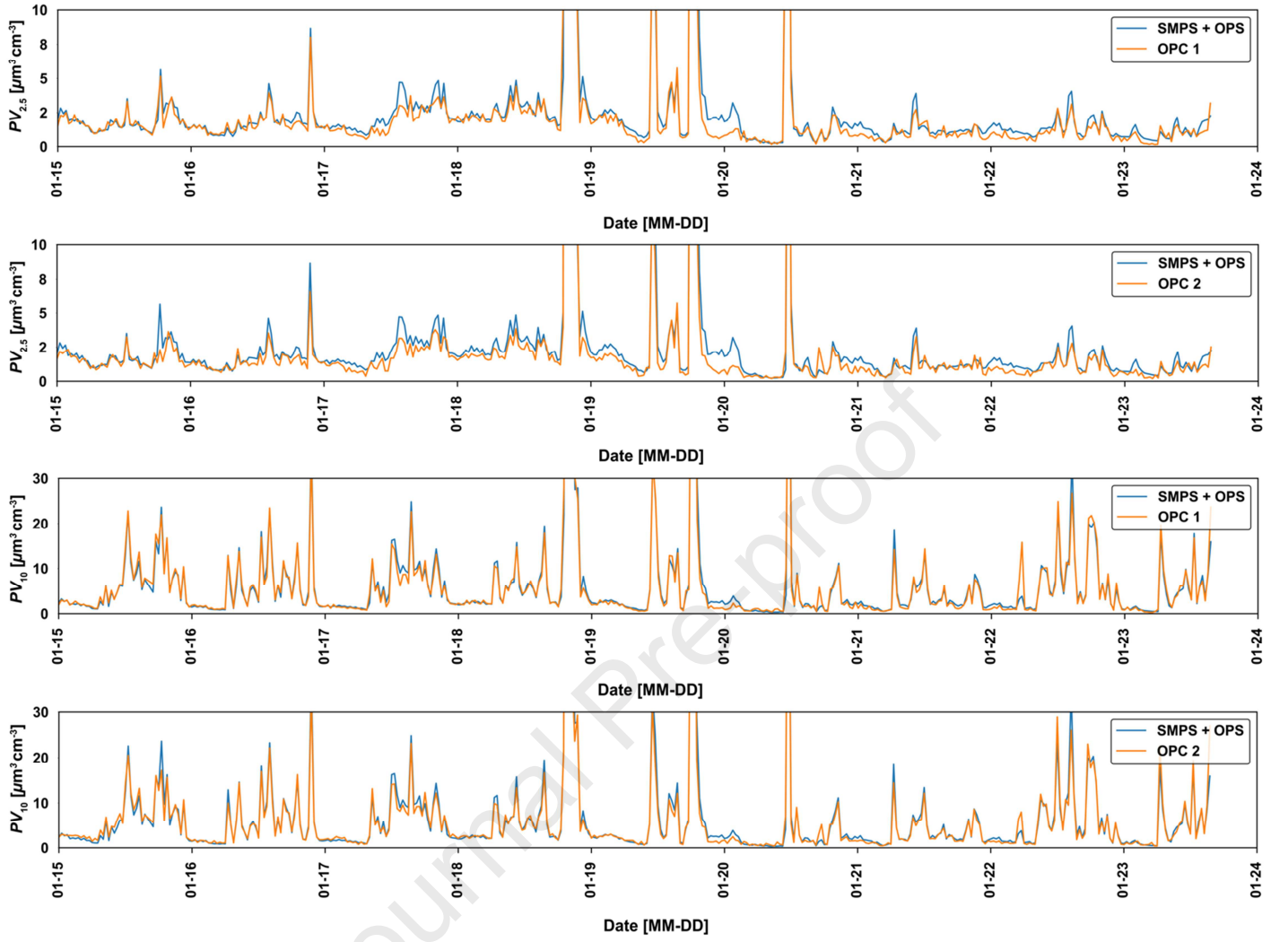


Figure 8. Comparison of the corrected $PV_{2.5}$ and PV_{10} concentration time-series reported by OPC 1 and OPC 2 with $PV_{2.5}$ and PV_{10} concentration time-series reported by the reference instruments (OPS and SMPS) during the testing set.

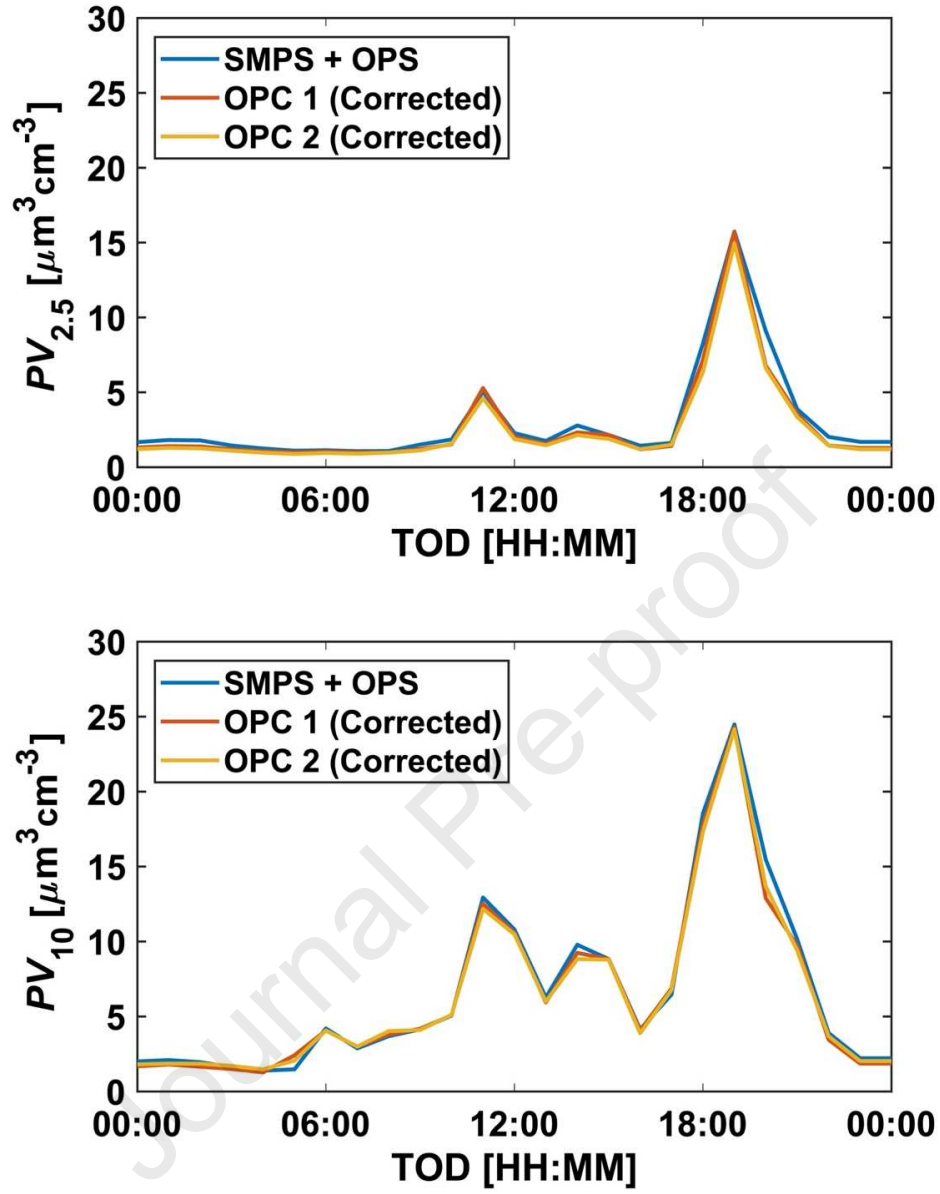


Figure 9. Comparison of diurnal profiles in the mean corrected $PV_{2.5}$ and PV_{10} concentrations as reported by OPC 1 and OPC 2 with those reported by the reference instruments (OPS and SMPS) during the testing set.

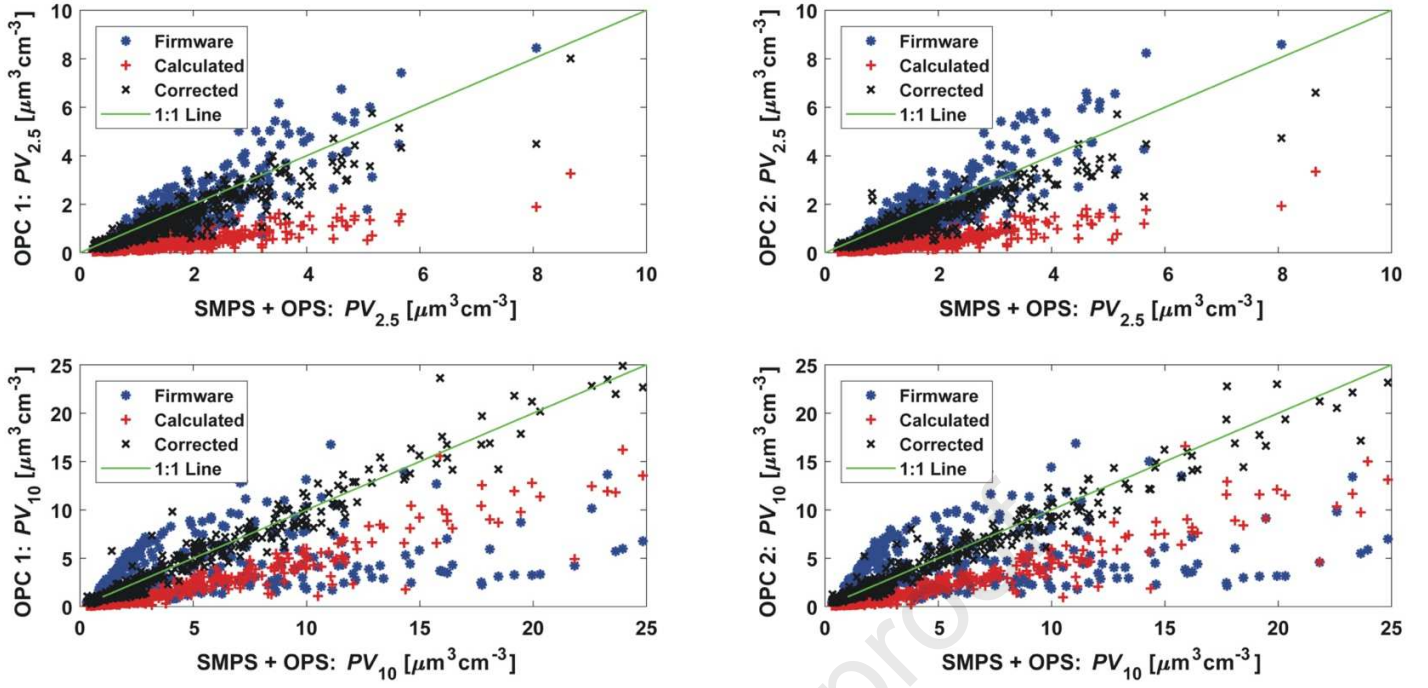


Figure 10. Correlation plots of the calculated and firmware (uncalibrated) and corrected $PV_{2.5}$ and PV_{10} concentrations reported by OPC 1 and OPC 2 with those reported by the reference instruments (OPS and SMPS) during the testing set.

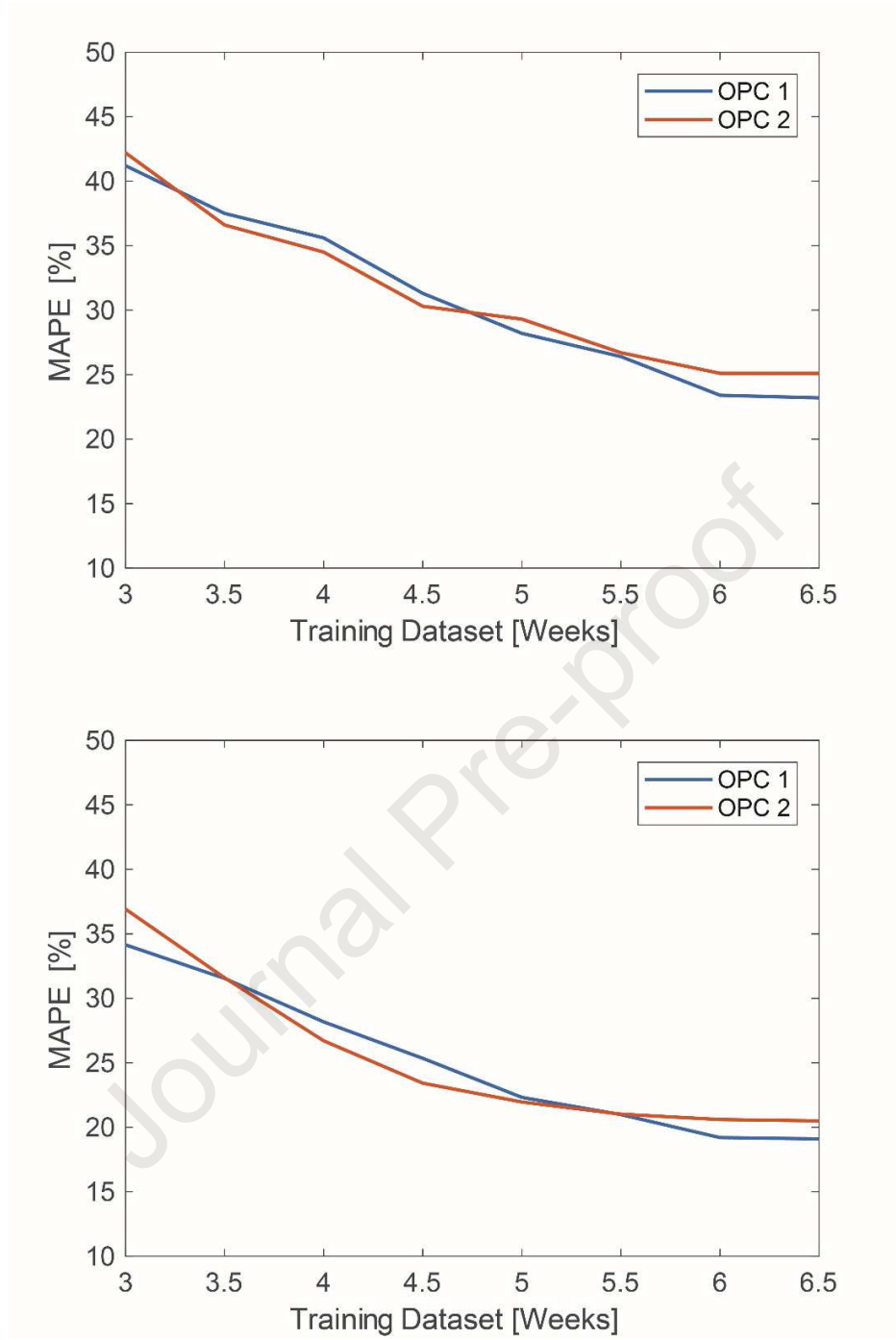


Figure 11. Variation of $MAPE$ for $PV_{2.5}$ (top) and PV_{10} (bottom) concentrations with varying size of the training dataset.

Highlights

- Developed a new machine learning field calibration method for improving the performance of low-cost OPCs.
- A two-month measurement campaign was conducted in a net-zero energy house to evaluate the calibration method.
- Calibration method improves the accuracy with which OPCs report indoor particle volume size distributions.
- Method is first to account for the volume contribution of the “missing” sub-0.38 μm particles not measured by OPCs.
- A non-parametric, size-resolved approach to OPC calibration targets limitations in particle detection/sizing by OPCs.

Declaration of Interest

The authors declare that they have no conflict of interest.