

1 Predicting PM<sub>2.5</sub> in well-mixed indoor air for a large  
2 office building using regression and artificial neural  
3 network models

4 *Brent Lagesse<sup>†</sup>, Shuoqi Wang<sup>‡</sup>, Timothy V. Larson<sup>‡</sup>, and Amy A. Kim<sup>\*‡</sup>*

5 <sup>†</sup>Division of Computing and Software Systems, University of Washington Bothell, Bothell,  
6 Washington 98011, United States

7 <sup>‡</sup>Department of Civil and Environmental Engineering, University of Washington, Seattle,  
8 Washington 98195, United States

9 ABSTRACT

10 Although the exposure to PM<sub>2.5</sub> has serious health implications, indoor PM<sub>2.5</sub> monitoring is not  
11 a widely applied practice. Regulations on indoor PM<sub>2.5</sub> level and measurement schemes are not  
12 well-established. Compared to other indoor settings, PM<sub>2.5</sub> prediction models for large office  
13 buildings are particularly lacking. In response to these challenges, statistical models were  
14 developed in this paper to predict the PM<sub>2.5</sub> concentration in well-mixed indoor air in a commercial  
15 office building. The performance of different modeling methods, including multiple linear  
16 regression (MLR), partial least squares regression (PLS), distributed lag model (DLM), least  
17 absolute shrinkage selector operator (LASSO), simple artificial neural networks (ANN), and long-

18 short term memory (LSTM), were compared. Various combinations of environmental and  
19 meteorological parameters were used as predictors. The root mean square error (RMSE) of the  
20 predicted hourly  $PM_{2.5}$  was  $1.73 \mu\text{g}/\text{m}^3$  for the LSTM model and in the range of  $2.20\sim 4.71 \mu\text{g}/\text{m}^3$   
21 for the other models when regulatory ambient  $PM_{2.5}$  data were used as predictors. The LSTM  
22 models outperformed other modeling approaches across the used performance metrics by learning  
23 the predictors' temporal patterns. Even without any ambient  $PM_{2.5}$  information, the developed  
24 models still demonstrated relatively high skill in predicting the  $PM_{2.5}$  levels in well-mixed indoor  
25 air.

26

## 27 1. INTRODUCTION

28 The exposure to ambient fine particulate matter (PM, with an aerodynamic diameter smaller than  
29  $2.5 \mu\text{m}$ ), or  $PM_{2.5}$ , and its health implication has been studied extensively<sup>1-6</sup>. The monitoring  
30 network of ambient  $PM_{2.5}$  is now well established in the United States, which provides essential  
31 evidence for ambient air quality regulations. However, indoor  $PM_{2.5}$  monitoring is not a widely  
32 applied practice, and regulation of indoor  $PM_{2.5}$  is lacking. Previous studies have shown that  
33 people in developed countries spent up to 90% of their time indoors<sup>7</sup> and the American worker  
34 spends eight hours a day on average at the workplace<sup>8</sup>. While long term measurement is available  
35 for ambient  $PM_{2.5}$  through various agencies<sup>9</sup>, indoor  $PM_{2.5}$  data are usually scarce.

36 To enable the indoor air quality (IAQ) assessment where direct measurement is not feasible,  
37 researchers are seeking to develop prediction models using other environmental variables that are  
38 readily available. Most recently, Wei et al.<sup>10</sup> conducted a review of studies using machine learning  
39 and statistical models for predicting the IAQ in various types of buildings and found that artificial  
40 neural networks (ANN) and regression were the most popular techniques. The results also showed

41 that just five out of the 37 reviewed studies were carried out in offices, and the models of these  
42 five studies were all developed using different types of feed-forward ANNs<sup>11-15</sup>. Only one of the  
43 five office studies focused on predicting indoor PM<sub>2.5</sub> using ambient PM<sub>2.5</sub> measurements<sup>11</sup>. In that  
44 study, Challoner et al.<sup>11</sup> predicted indoor PM<sub>2.5</sub> using ambient PM<sub>2.5</sub> concentrations and  
45 meteorological data in a mechanically ventilated office building with ANN and reported large  
46 errors ranging from -8.09 to 4.93 μg/m<sup>3</sup>. However, the ambient PM<sub>2.5</sub> concentration was calculated  
47 using a personal-exposure activity location model instead of measured directly at the building site,  
48 which might cause the large errors in the predictions<sup>11</sup>.

49 Other regression models, including multiple linear regression (MLR), stepwise regression,  
50 partial least squares regression (PLS), and principal component regression (PCR), have been  
51 applied in dwellings, schools, and subway stations but not in offices<sup>10</sup>. It is unclear whether these  
52 regression models could be used for predicting the PM<sub>2.5</sub> in offices and how well they perform  
53 compared to the ANN model.

54 A commercial office often consists of various types of regularly occupied spaces, e.g., open  
55 workstations, conference rooms, and common areas. There is no definitive method in the  
56 placement of monitors to assess the overall PM<sub>2.5</sub> level in the entire space. Predicting the spatial  
57 variation of PM<sub>2.5</sub> in the large area is also potentially complex. Therefore, developing a prediction  
58 model for spatial-averaged PM<sub>2.5</sub> could be the first step. The exhaust air is a well-mixed sample of  
59 the return air from different indoor locations and could serve as a representation of the spatially  
60 averaged condition.

61 The objective of this paper is to develop statistical models to predict the PM<sub>2.5</sub> concentration in  
62 well-mixed indoor air inside a commercial office building using MLR, PLS, a simple ANN, and a  
63 specific type of ANN known as a Long Short-Term Memory (LSTM) neural network. Several

64 recent studies have explored the use of the LSTM neural network in predicting ambient  $PM_{2.5}$   
65 concentration<sup>16-20</sup>, given that it is better suited for long time-series predictions than simple ANN  
66 models. Yet, its use in an indoor office setting has not been investigated as a comparison to a  
67 simple ANN model. Regression models that are capable of handling time-series data and  
68 evaluating delayed effects, i.e., the distributed lag model (DLM), the least absolute shrinkage  
69 selector operator (LASSO), as well as PLS with lagged predictors (PLS-Lag), are also considered.  
70 The dependent variable in this paper was limited to indoor  $PM_{2.5}$ , while other pollutants, such as  
71 chemicals emitted by the occupants, were not included. The independent variables considered  
72 included meteorological variables (e.g., wind speed, wind direction, air temperature, and relative  
73 humidity), publicly available ambient  $PM_{2.5}$  concentration from other locations, number of  
74 occupants at the study site, and building operational data. Various combinations of the independent  
75 variables were tested to evaluate the prediction accuracy with or without ambient  $PM_{2.5}$   
76 concentration (whether from publicly available monitoring sites or measured directly at the study  
77 site). The performance of the various models was compared using several performance indicators,  
78 including the normalized absolute error (NAE), the root mean square error (RMSE), the coefficient  
79 of determination ( $R^2$ ), and the index of agreement (IA).

## 80 2. EXPERIMENTAL SECTION

### 81 2.1. Prediction Variables

82 The variables used to predict hourly indoor  $PM_{2.5}$  ( $PM_E$ ) included hourly outdoor  $PM_{2.5}$ , relative  
83 humidity (RH), air temperature (T), and wind speed. We also included building air intake damper  
84 opening fraction on an hourly basis as a measure of outdoor air intake. Filtration of the outdoor air  
85 is discussed in detail in Section 1.1 of the Supporting Information (SI). Occupancy level on an  
86 hourly basis was also introduced as a predictor to account for the impact of indoor human activity.

87 The damper opening and occupancy on an hourly basis are rarely used in other studies due to the  
88 lack of such information.

89 A list of all the predictor variables and relevant descriptions is given in Table 1.

90

91 Table 1. Descriptions of the prediction variables.

No.	Variable	Unit	Description
1	$PM_{LYN}$	$\mu\text{g}/\text{m}^3$	Ambient $PM_{2.5}$ measured at the PSCAA Lynnwood site.
2	$PM_B$	$\mu\text{g}/\text{m}^3$	Ambient $PM_{2.5}$ measured at the PSCAA Bellevue site.
3	$PM_W$	$\mu\text{g}/\text{m}^3$	Ambient $PM_{2.5}$ measured at the PSCAA TW site.
4	$PM_{LFP}$	$\mu\text{g}/\text{m}^3$	Ambient $PM_{2.5}$ measured at the PSCAA LFP site.
5	$PM_D$	$\mu\text{g}/\text{m}^3$	Ambient $PM_{2.5}$ measured at the PSCAA Duwamish site.
6	$PM_O$	$\mu\text{g}/\text{m}^3$	Ambient $PM_{2.5}$ measured at location 4.
7	$T_S$	$^{\circ}\text{C}$	Air temperature of the supply air (location 1).
8	$T_E$	$^{\circ}\text{C}$	Air temperature of the exhaust air (location 2).
9	$T_F$	$^{\circ}\text{C}$	Air temperature of the floor air (location 3).
10	$T_O$	$^{\circ}\text{C}$	Air temperature of the ambient air logged by the BCS.
11	$RH_S$	%	Relative humidity of the supply air (location 1).
12	$RH_E$	%	Relative humidity of the exhaust air (location 2).
13	$RH_F$	%	Relative humidity of the floor air (location 3).
14	$O_F$	-	Relative occupancy of the floor (see SI).
15	$D$	-	Air intake damper opening fraction logged by the BCS.
16	$WD$	$^{\circ}$	Wind direction recorded on the ATG rooftop.
17	$WS$	m/s	Wind speed recorded on the ATG rooftop.

92

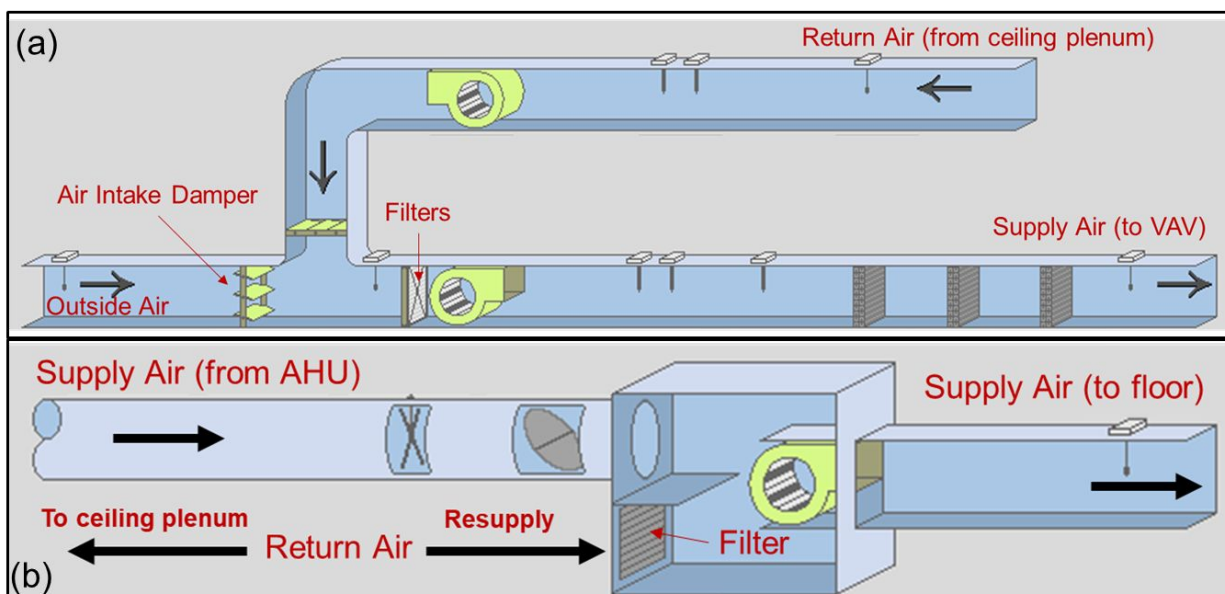
## 93 2.2. Ambient $PM_{2.5}$ Measurements

94 The hourly average ambient  $PM_{2.5}$  concentrations were collected from five monitoring sites  
95 managed by the Puget Sound Clean Air Agency (PSCAA)<sup>21</sup> as part of Washington's air monitoring  
96 network<sup>22</sup>. The selected five sites were all within 16 km of the UW Tower building. A general  
97 description of the environment near each site is given in Table S1. The wind speed and wind

98 direction records were obtained from a weather station on the rooftop of the Atmospheric Sciences-  
99 Geophysics (ATG) Building<sup>23</sup> on the UW campus, approximately 0.9 km from the UW Tower.  
100 Figure S1 shows the location of UW Tower in relation to the PSCAA sites and the weather station.

### 101 2.3. Indoor Measurements

102 The indoor measurements were recorded every five minutes on floor O-3 in the University of  
103 Washington (UW) Tower building in Seattle. The building schematic and floor plan of the selected  
104 O-3 office space is shown in Figure S2. A detailed description of the ventilation system operation  
105 can also be found in Section 1.1 of the SI. The building control system (BCS) manages the building  
106 ventilation and logged the air intake damper opening fraction and ambient air temperature every  
107 five minutes. Figure 1 shows the configuration of the air handling unit and variable air volume  
108 boxes on floor O-3. Since the supply air is a mixture of outside air and return air, the air intake  
109 damper opening is a better indicator of the amount of outdoor air brought into the space than  
110 ventilation rate. The RH and T were measured at three locations on floor O-3, as shown in Figure  
111 S2, using three units of Particles Plus 7302-AQM air quality monitors (AQM)<sup>24</sup>. The PM<sub>2.5</sub>  
112 concentration in the exhaust air was measured at location 2. The exhaust air was considered a well-  
113 mixed air sample representing the spatial average of the indoor air on the O-3 floor. A Radiance  
114 M903 nephelometer<sup>25</sup> was used at location 4 outdoors to record the concentration of ambient PM<sub>2.5</sub>  
115 adjacent to the building. The particle mass concentration calibration process is detailed in Section  
116 1.2 of the SI.



117  
 118 Figure 1. Configuration of the ventilation system on floor O-3. (a) air handling unit (AHU) located  
 119 in the mechanical room; (b) one of the variable air volume (VAV) boxes located in the ceiling  
 120 plenum space on the floor.

121 In addition to the AQM, an occupancy sensor was installed at location 3 in the open office space  
 122 on floor O-3 (see Figure S2). It estimated relative occupancy by counting the number of Media  
 123 Access Control (MAC) addresses that communicate during the five-minute sampling period. The  
 124 working theory of the sensor is explained in Section 1.3 of the SI.

#### 125 2.4. Data Processing

126 The measurements were conducted from August 2<sup>nd</sup> through November 13<sup>th</sup> in 2019. Due to  
 127 various technical issues (e.g., repair of power supply, failure to start logging), missing values and  
 128 measurement gaps existed for some of the variables. The duration of the measurements for each  
 129 variable is shown in Figure S3. Only a subset of the observations was used to construct the analysis  
 130 dataset. The collected data were transformed into 1-hour averages and merged to create a time  
 131 series matrix. The  $PM_E$  data were also checked for outliers. The first quantile ( $Q1$ ), third quantile

132 ( $Q3$ ), and the interquartile range ( $IQR$ ) were calculated, and the data points above the upper bound  
133 (defined as  $Q3 + 1.5IQR$ ) or below the lower bound (defined as  $Q1 - 1.5IQR$ ) were labeled as  
134 outliers. In total, 21 outliers were identified for  $PM_E$  out of 1,535 observations (1.4%) and  
135 removed. After removing missing observations from all the other variables, the remaining matrix  
136 with complete data contained 670 observations (hourly averaged values) for a total of 17 variables.

### 137 3. MODELS

#### 138 3.1. Multiple Linear Regression (MLR)

139 MLR models are the most commonly used for IAQ predictions as summarized in Wei et al.<sup>10</sup>  
140 The MLR approach enables the use of various independent variables (e.g., ambient  $PM_{2.5}$ ,  
141 meteorological, occupancy) to predict the outcome (i.e., the  $PM_{2.5}$  in well-mixed indoor air). The  
142 coefficient of each independent variable reflects the effect of the variable in predicting the  
143 outcome. All 17 independent variables were included at the beginning and a stepwise selection  
144 was conducted to find the final model that has the lowest Akaike Information Criterion (AIC)  
145 value. Multicollinearity is a known issue with MLR models when high correlations exist between  
146 independent variables<sup>26, 27</sup>. Therefore, the variance inflation factors (VIF) of the independent  
147 variables in the final model were examined. Variables with large VIF values were excluded to  
148 ensure low correlations among all the predictors in the final model<sup>27</sup>.

#### 149 3.2. Partial Least Squares (PLS)

150 PLS is a regression method with dimension reduction capability. For building a prediction model  
151 with many potentially correlated predictors, applying the PLS technique allows the user to  
152 transform the predictors into a reduced set of orthogonal latent variables (or components), which  
153 are a linear combination of the original predictors<sup>28</sup>. Compared to MLR, the PLS method has



154 shown its capability in building robust  $PM_{2.5}$  prediction models by coping with the  
155 multicollinearity issue present with a large number of predictors<sup>29-32</sup>.

### 156 3.3. Artificial Neural Network (ANN)

157 An ANN is a collection of algorithms used to learn patterns from data and then use those patterns  
158 for predicting or classifying new data that has not previously been seen. This paper utilized a  
159 simple ANN (additional background information is included in Section 1.4 of the SI), and a type  
160 of recurrent neural network (RNN) called LSTM. An RNN is a neural network where some of the  
161 connections propagate backward in addition to forward. RNNs are commonly used in applications  
162 that involve time-series data since the feedback connections help it learn temporal sequences.  
163 LSTMs are a refinement on RNNs that mitigate the vanishing gradient problem where the model  
164 unintentionally learns to ignore the feedback connections. Handling missing data is an active area  
165 of research in the machine learning community<sup>33</sup>. The models described in this paper ignore  
166 observations with missing features and the LSTM uses the two most recent observations, even if  
167 they are not the two preceding hours, but other approaches from emerging research could be  
168 examined in future work.

169 The term “ANN” is used hereafter to only refer to the simple ANN model. Hyperparameters for  
170 each model were chosen through a grid search approach after eliminating hyperparameters that  
171 never resulted in competitive models. The hyperparameters used for the ANN and LSTM can be  
172 found in Table S2, while Table S3 lists the options for hyperparameters that were selected.

### 173 3.4. Time Series Regression

174 Given LSTM’s capability in handling time-series data and learning temporal patterns, two  
175 additional regression techniques, i.e., DLM and LASSO, were employed as a way to control for  
176 autocorrelation. Both DLM and LASSO can model the delayed effects of the independent time-

177 series variables on the dependent time-series variable including different time lags. The PLS model  
178 can also be modified to include lagged independent variables as predictors (PLS-Lag). To facilitate  
179 a fair comparison with LSTM, additional ANN models (ANN-Lag) were evaluated where the two  
180 most recent observations were used in the learning process similar to LSTM.

### 181 3.5. Model Training and Testing

182 The dataset of 670 observations with 17 columns of the independent variables and one column  
183 of the dependent variable was separated into a training set and a testing set. The time sequence  
184 structure of the dataset was maintained. The first 546 observations (81%) were kept in the training  
185 set and the rest in the testing set. This split was done due to the fact that a large time gap existed  
186 between observations 546 and 547 because of missing data. The training set was further split into  
187 a training subset and validation subset for the implementation of cross-validation (CV). A rolling  
188 forecasting origin technique as discussed by Hyndman and Athanasopoulos<sup>34</sup> was used to create  
189 10 resamples of the training and validation subsets while maintaining time sequence. The details  
190 of the CV scheme are illustrated in Figure S4. During CV, the final model was selected based on  
191 RMSE.

### 192 3.6. Model Implementation

193 A consistent three-phase model framework was implemented across the four modeling  
194 approaches, i.e., MLR, PLS, ANN, and LSTM:

- 195 • During Phase 1, for PLS, ANN, and LSTM models, all 17 predictors were included. For  
196 the MLR model, because some predictors could be highly correlated, a bidirectional  
197 stepwise regression was run to determine the best subset of predictors for the full model.  
198 Remaining predictors in the subset with high VIF values were removed.
- 199 • During Phase 2, the Phase 1 models were re-evaluated after removing predictor  $PM_{10}$ .

- 200           • During Phase 3, the Phase 2 models were re-evaluated after removing all other PM<sub>2.5</sub>  
201           independent variables.

202       The three phases were designed in accordance with the potential difficulties of obtaining PM<sub>2.5</sub>  
203       measurements. On-site ambient PM<sub>2.5</sub> measurement requires high-grade instruments to provide  
204       accurate readings regardless of the ambient weather condition, which may not be feasible for some  
205       buildings due to economic or operational constraints. Ambient PM<sub>2.5</sub> records from government  
206       agencies are publicly available, but some labor costs may be involved to acquire and process the  
207       data. In addition, regulatory monitoring sites may not exist in the target city or even nearby cities.  
208       Therefore, by evaluating models without certain PM<sub>2.5</sub> predictors, users are given the option to  
209       choose the best approach based on each building's unique condition.

### 210       3.7. Model Evaluation

211       Several indicators, i.e., NAE, RMSE, R<sup>2</sup>, and IA, are used to compare the performance of the  
212       predictive models. The use of these indicators was also demonstrated in Elbayoumi et al<sup>26</sup>. The  
213       NAE and RMSE are smaller-the-better metrics that measure the existing error of the model, while  
214       the R<sup>2</sup> and IA are larger-the-better metrics that measure the accuracy of the model. Calculation of  
215       each indicator is given in Equations (1) through (4):

$$216 \quad NAE = \frac{\sum_{i=1}^N |P_i - O_i|}{\sum_{i=1}^N O_i} \quad (1)$$

$$217 \quad RMSE = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (P_i - O_i)^2} \quad (2)$$

$$R^2 = \left( \frac{\sum_{i=1}^N (P_i - \bar{P})(O_i - \bar{O})}{N \cdot S_p \cdot S_o} \right)^2 \quad (3)$$

218

$$IA = 1 - \left[ \frac{\sum_{i=1}^N (P_i - O_i)^2}{\sum_{i=1}^N (|P_i - \bar{O}| + |O_i - \bar{O}|)^2} \right] \quad (4)$$

219 where  $N$  is the number of observations;  $P_i$  and  $O_i$  are the predicted and observed values of the  $i^{\text{th}}$   
 220 observation;  $\bar{P}$  and  $\bar{O}$  are the averages of the predicted and observed values;  $S_p$  and  $S_o$  are the  
 221 standard deviations of the predicted and observed values.

## 222 4. RESULTS

### 223 4.1. Descriptive Statistics

224 A summary of the descriptive statistics of the dataset is provided in Table S4. The dataset  
 225 contained observations from October 10<sup>th</sup> through November 13<sup>th</sup> in 2019. As shown in Table S4,  
 226 the mean hourly averaged indoor  $\text{PM}_{2.5}$  was  $5.68 \mu\text{g}/\text{m}^3$  while the mean hourly ambient  $\text{PM}_{2.5}$   
 227 measured at the UW Tower was  $4.07 \mu\text{g}/\text{m}^3$ . The ambient  $\text{PM}_{2.5}$  at the five PSCAA monitoring  
 228 sites were also within acceptable range per the National Primary and Secondary Ambient Air  
 229 Quality Standards<sup>35</sup> most of the time (see Figure S5). Pearson's correlation coefficients of the  
 230 predictors are summarized in Table S5.

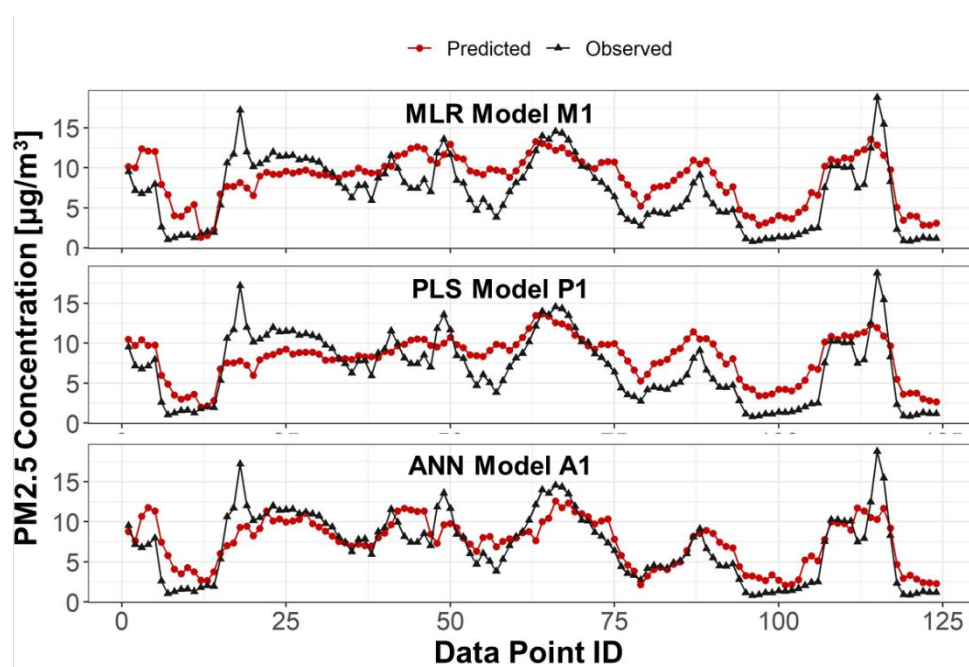
### 231 4.2. Model Performance

232 Using the three-phase implementation framework and the three modeling approaches, several  
 233 models were evaluated for their performance in predicting indoor  $\text{PM}_{2.5}$ . Considering that  
 234 regulatory ambient  $\text{PM}_{2.5}$  records from nearby monitoring sites should be relatively easy to obtain  
 235 for most commercial office buildings located in urban centers, only the Phase 2 model results are  
 236 presented here to be succinct. Results for Phase 1 and Phase 3 models are included in the SI (See

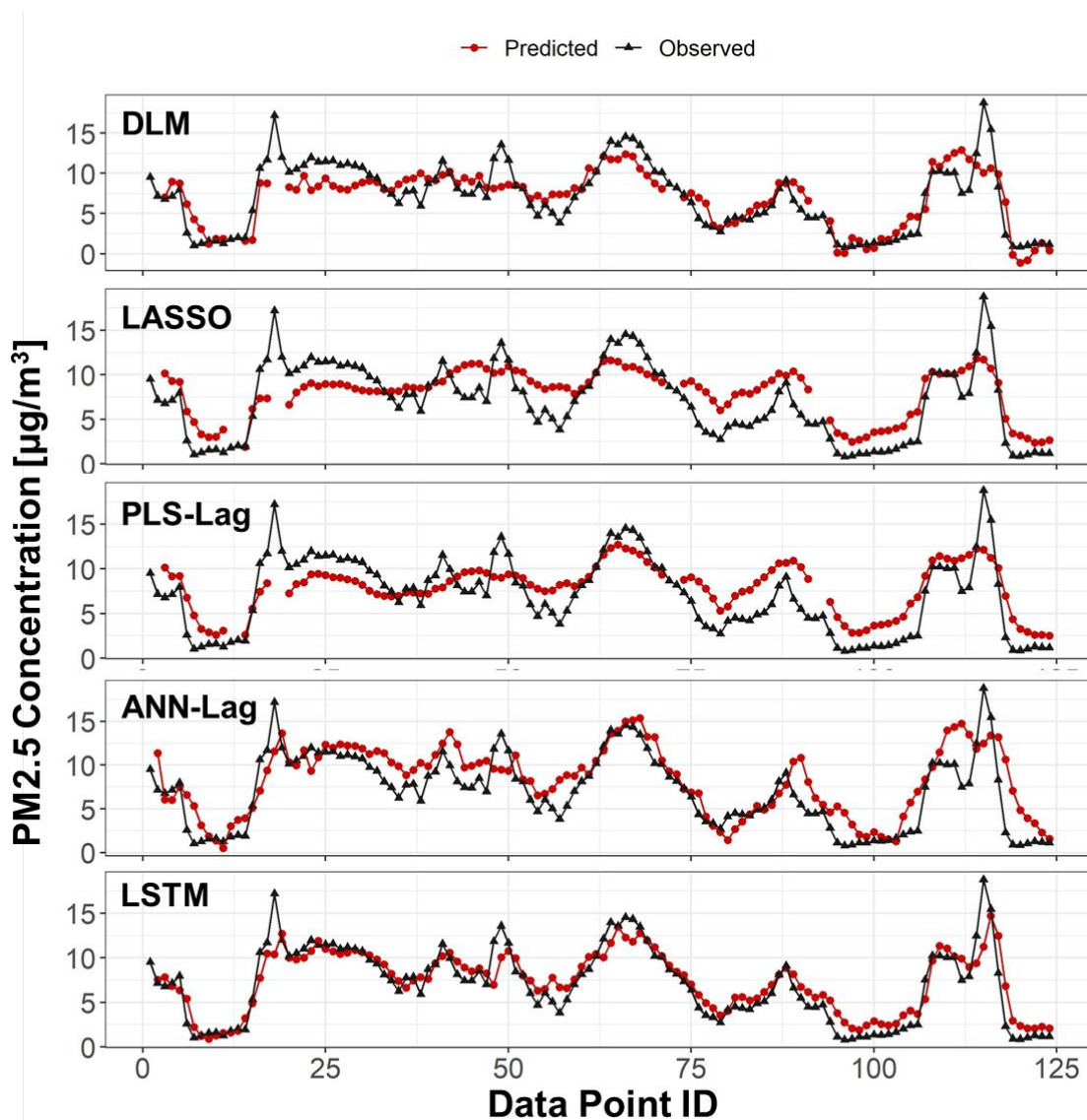
237 Figure S6, Table S6, and Table S7). The values of calculated Phase 2 model performance indicators  
 238 are summarized in Table 2. Figures 2 and 3 show the predicted values versus observations for  
 239 different models using the testing dataset without and with the temporal information considered.  
 240 Details of the results obtained for each modeling approach are discussed in the following sections.  
 241  
 242 Table 2. Performance indicators of the Phase 2 models.

Temporal Information	Model No.	Method	Predictors	NAE	RMSE	R <sup>2</sup>	IA
Not considered	M1	MLR	10	0.37	3.07	0.60	0.82
	P1	PLS	16	0.35	2.89	0.60	0.83
	A1	ANN	16	0.26	2.38	0.67	0.88
Considered	DL1	DLM	10	0.24	2.20	0.71	0.91
	LA1	LASSO	20	0.33	2.65	0.65	0.85
	PL1	PLS-Lag	30	0.54	4.71	0.02	0.50
	AL1	ANN-Lag	16	0.29	2.63	0.60	0.89
	L1	LSTM	16	0.18	1.73	0.83	0.94

243



245 Figure 2. Plots of predicted and observed values for Phase 2 MLR, PLS, and ANN models using  
 246 the testing set.



247  
 248 Figure 3. Time series plots of predicted and observed values for DLM, LASSO, PLS-Lag, ANN-  
 249 Lag and LSTM models.

#### 250 4.2.1. MLR Modeling Results

251 The MLR was conducted in R<sup>36</sup> using the “caret”<sup>37</sup> package for cross-validation and the  
 252 “MASS”<sup>38</sup> package for bidirectional stepwise regression. As discussed in Mansfield and Helms<sup>39</sup>,

253 multicollinearity is not a problem if the VIFs are not unusually larger than 1.0. The VIFs of the ten  
 254 predictors in Model M1 are in the range of 1.22-2.92, and the air intake damper opening  $D$  appears  
 255 to have a significant effect on the indoor  $PM_{2.5}$ , as shown in Table 3. Similar results can be  
 256 observed for the Phase 1 and Phase 3 models (see Section 1.5.1 and Table S8 in the SI). The  
 257 exclusion of on-site ambient  $PM_{2.5}$  predictor  $PM_O$  led to a slight increase of RMSE (12%) in Model  
 258 M1 compared to Model M2 while the ambient  $PM_{2.5}$  from other locations were kept in the model.  
 259 Removal of all the ambient  $PM_{2.5}$  predictors, led to an increase of RMSE of 18% in Model M3  
 260 compared to Model M2. A modified version of Model M2 was evaluated by swapping  $PM_W$ ,  $P$   
 261  $M_{LFP}$ , and  $PM_D$  with on-site  $PM_O$  data, and the results were similar (RMSE decreased by 9%). It  
 262 shows that the inclusion of some ambient  $PM_{2.5}$  information, not necessarily measured on-site,  
 263 could improve the prediction accuracy of the model.

264

265 Table 3. Summary of Model M1 results.

	Intercept	$PM_W$	$PM_{LFP}$	$PM_D$	$T_S$	$T_F$	$RH_E$	$O_F$	$D$	$WD$	$WS$
VIF	-	2.71	2.21	2.92	1.84	1.44	1.32	1.26	1.43	1.22	1.49
Coef	-7.28	0.08	0.14	0.17	-0.10	0.17	0.17	0.01	-6.19	0.004	-0.96

266

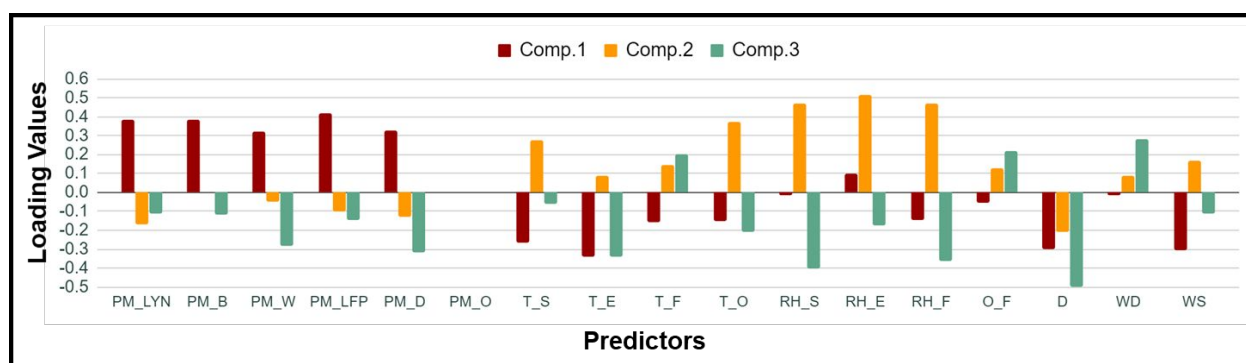
#### 267 4.2.2. PLS Modeling Results

268 The PLS regression was conducted in R<sup>36</sup> using the “pls”<sup>40</sup> package and the kernel algorithm<sup>41</sup>.

269 The optimal number of components in the model was determined using a randomization test  
 270 approach<sup>42</sup>, which checked whether the squared prediction errors of the models with fewer  
 271 components were significantly larger than in the reference model and selected the smallest model  
 272 not significantly worse than the reference model. Figure S7 shows the cross-validation plots and  
 273 the determined number of components for each model.

274 The percentage of variance explained by each component of the PLS model for both the  
 275 predictors and outcome is summarized in Table S9. Component 1 of all three PLS models appears  
 276 to make the most contribution (49.54%~62.01%) in explaining the variance of the outcome  
 277 variable. For each component of Model P1, the loading value of each predictor is shown in Figure  
 278 4. It can be seen that all of the ambient PM<sub>2.5</sub> predictors carried large positive loading values in  
 279 Component 1 when they were included in the model. A similar effect can be observed for Models  
 280 P2 and P3 (see Figure S8). Additional discussion regarding other predictors can be found in Section  
 281 1.5.2 of SI.

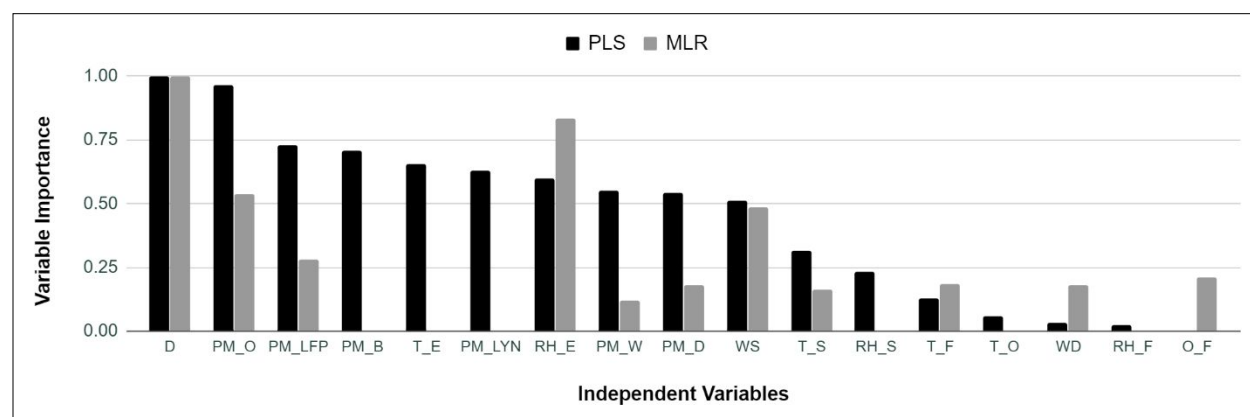
282



283

284 Figure 4. Loading values of each variable for each component of Model P1.

285



286

287 Figure 5. The variable importance of all the predictors using PLS and MLR approaches.



288

289 The variable importance (VIP) was calculated using the “caret”<sup>37</sup> package in R for both the MLR  
290 and PLS Phase 1 models which included all the predictors, as shown in Figure 5. It can be seen  
291 that the critical predictors are relatively consistent in either approach. The damper opening,  
292 ambient PM<sub>2.5</sub> level (on-site as well as from other locations), well-mixed air temperature and RH,  
293 and wind speed appear to be the variables that affect the prediction to a greater extent.

#### 294 4.2.3. ANN and LSTM Modeling Results

295 The results reported in Table 2 represent the average values for each metric across all folds of  
296 the 10-fold cross-validation run using that model (see Table S6 for Phase 1 and Phase 3 models).  
297 The model was selected based on the best average RMSE value. Table S10 shows the  
298 hyperparameters that were used in each of these models. Note that many of the models built with  
299 alternate hyperparameter configurations also performed nearly as well in RMSE and some  
300 performed better in the other metrics.

301 The ANN was tested both with and without temporal information included in the training. The  
302 model performance did not always improve when the previous two hours’ data were considered in  
303 the training (RMSE degraded by 0.89 and 0.25  $\mu\text{g}/\text{m}^3$  for Phase 1 and Phase 2 models, and  
304 improved by 0.52  $\mu\text{g}/\text{m}^3$  for the Phase 3 model). The LSTM model, which learns temporal patterns,  
305 was also trained and tested with inputs from the previous two hours in addition to the current values  
306 for every variable excluding the dependent variable,  $PM_E$ . As shown in Table 2, the LSTM model  
307 outperformed the other models across the four metrics.

#### 308 4.2.4. Time Series Regression Results

309 The DLM regression was conducted in R<sup>36</sup> using the “dLagM”<sup>43</sup> package. The built-in  
310 autoregressive distributed lag (ARDL) bounds testing function was used to compute the optimal

311 lag structure for the independent variables of Model DL1. The maximum lag period considered  
312 was two hours (same as in LSTM) and only ambient  $PM_{2.5}$ , occupancy, and damper opening  
313 variables were included in the lag structure. The results suggested a 2-hour lag for  $PM_D$ , 1-hour  
314 lag for damper opening and zero lags were used for all the other variables (as listed in Table S7).  
315 By using the lagged variables, the delayed effect of these predictors was included in the DL1  
316 model.

317 The LASSO regression was ran using the “glmnet”<sup>44</sup> package in R. The optimal tuning parameter  
318  $\lambda$  which controls the overall strength of the penalty was selected by rolling forecast origin CV as  
319 discussed in Section 3.5. Both 1-hour and 2-hour lagged  $PM_W$ ,  $PM_{LFP}$ ,  $PM_D$ ,  $O$ , and  $D$  were  
320 included in LASSO as well as the unlagged versions. Other variables were included without any  
321 lags. The PLS-Lag model was based on Model P1. Similar to Model LA1, both 1-hour and 2-hour  
322 lagged  $PM_{LYN}$ ,  $PM_B$ ,  $PM_W$ ,  $PM_{LFP}$ ,  $PM_D$ ,  $O$  and  $D$  were included in PLS-Lag as well as their  
323 unlagged versions. Other variables were kept without any lags.

324 The results show better performance of Model DL1 compared to LA1 and PL1 as listed in Table  
325 2. However, from Figure 3, the three models are outperformed by LSTM. Notice that gaps exist in  
326 the predicted time series of DL1, LA1, and PL1 models, due to missing data of in some of the  
327 predictors. As mentioned in Section 3.3, LSTM was able to use the most recent two observations  
328 regardless of missing data, no gaps exist for the LSTM predictions.

## 329 5. DISCUSSION

330 The development and comparison of the various predictive models have shown that the indoor  
331  $PM_{2.5}$  in the well-mixed air in this office space could be estimated by using readily available  
332 variables. In general, when temporal information is not considered, the performance of the models  
333 developed using the MLR, PLS, and ANN methods were comparable in terms of their NAE,

334 RMSE,  $R^2$ , and IA, as shown in Table 2 and Table S6. With temporal information included in the  
335 model, the LSTM method outperformed the DLM, LASSO, PLS-Lag, and ANN-Lag presumably  
336 because the LSTM took into consideration the lagged effect as well as the rate of change in the  
337 predictor variables. For example, the LSTM model may learn that the value of a variable from 2  
338 hours ago has an effect on  $PM_E$  at the current time. In addition, it may also learn that the rate of  
339 change of a variable over the past 2 hours has an effect on  $PM_E$  at the current time. The models  
340 with a large number of independent variables appear to provide a marginally better prediction for  
341 regression models (M2 and P2), but not for the neural network models (A2, AL2, and L2). The  
342 reduced models with fewer predictors are still capable of making accurate predictions. The results  
343 suggest that although on-site measurement of ambient  $PM_{2.5}$  could aid in predicting the indoor  
344 level, using measurements from other publicly available monitors instead has minimal impact on  
345 the model performance. By including some form of ambient  $PM_{2.5}$  measurements (not necessarily  
346 on-site), there is a significant improvement in the model results.

347 The RMSE values of the regression and ANN models developed in this paper for the office space  
348 are in the range of 2.05~4.71  $\mu\text{g}/\text{m}^3$  while the values of the LSTM models are in the range of  
349 1.73~1.93  $\mu\text{g}/\text{m}^3$ . In comparison, as summarized in Wei et al.<sup>10</sup>, the RMSE values for regression  
350 models developed for indoor  $PM_{2.5}$  in schools<sup>26</sup> and private dwellings<sup>45-47</sup> were in the range of  
351 0.45~1.7  $\mu\text{g}/\text{m}^3$ . The ANN type models have been used to predict indoor  $PM_{2.5}$  in subway stations<sup>48</sup>,  
352 <sup>49</sup>, dwellings<sup>50</sup>, and schools<sup>26</sup>. Similarly to the regression models, the reported error values were  
353 small for dwellings and schools (1~3  $\mu\text{g}/\text{m}^3$  RMSE) but large for the subway stations (RMSE over  
354 10  $\mu\text{g}/\text{m}^3$ ). This paper shows that regression and simple ANN models are quite capable of  
355 predicting indoor  $PM_{2.5}$  in offices. Using an LSTM to account for time trends in the predictor

356 variables further provides a significant improvement of the prediction performance over time  
357 series regression methods.

358 Unlike schools, dwellings, and ambient air, a mechanically ventilated office has a relatively  
359 consistent indoor environment controlled by the ventilation system. With the air filtration in place,  
360 the correlation of indoor and outdoor  $PM_{2.5}$  exists but is not as high as in a naturally ventilated  
361 building. Therefore, in addition to the usual meteorological variables, i.e., air temperature, RH,  
362 and wind speed, other building-related variables, i.e., damper opening and occupancy, also  
363 appeared to be useful. Due to the difficulties in obtaining these building-related variables, few  
364 studies have included them in the prediction models. As shown in Figure 5, the damper opening  
365 was the most important variable in both MLR and PLS models. Since the studied office space did  
366 not have any operable windows, the ventilation system was the main route through which ambient  
367 PM entered the indoor environment while infiltration and tracking remained secondary routes. The  
368 high variable importance of the damper opening and ambient PM variables in Figure 5 show that  
369 the ambient condition has a major influence on the indoor environment. In addition, the air  
370 temperature and relative humidity of the well-mixed air, as well as the outdoor wind speed, also  
371 appear to carry large weights in the prediction model. As discussed in Gundel and Destailats<sup>51</sup>,  
372 the ambient particles go through a phase change when entering the building via ventilation or  
373 infiltration due to the change of temperature and relative humidity conditions.

374 The results presented in this paper have some practical implications. First, a rooftop ambient  
375  $PM_{2.5}$  monitor is not always necessary to produce a fairly good prediction of well-mixed indoor  
376  $PM_{2.5}$  unless nearby regulatory monitors do not exist. From a building management perspective,  
377 this is encouraging as there is always cost associated with conducting on-site  $PM_{2.5}$  measurements,  
378 including material and labor costs for sensor procurement, installation, and data analysis. The other

379 predictors, such as indoor air temperature, relative humidity, and wind speed, are all commonly  
380 monitored environmental parameters in existing commercial buildings or could be obtained from  
381 nearby weather stations. The air intake damper opening could be difficult to obtain if the building  
382 control system lacks the capability of continuous monitoring of the damper position. In this case,  
383 substitute parameters, e.g., outdoor airflow rate, could be used if such monitoring is easier to  
384 implement. Building occupancy monitoring also requires a specialized sensor for data collection.  
385 Nevertheless, the importance of occupancy in the models is relatively low compared to the other  
386 predictors, as shown in Figure 5.

387 Some limitations exist in this paper. The analysis dataset only contained measurements from  
388 October and November in 2019 when the outdoor weather was relatively mild in the Seattle area.  
389 Extensive data collection is needed to evaluate the model performance in different seasons. The  
390 ambient  $PM_{2.5}$  in the Seattle area was maintained at a healthy level during the measurement. It is  
391 unknown whether the degradation of ambient air quality (e.g., during wildfire events) could affect  
392 the predictive capability of the models. The outcome  $PM_{2.5}$  variable was measured in the  
393 exhaust/well-mixed air, and it may not be the same as the  $PM_{2.5}$  measured at other locations in the  
394 building. How well these models predict the  $PM_{2.5}$  level at other indoor locations is not in the  
395 scope of this paper. It is also recognized that the models were trained and tested using data from  
396 one floor in a single building. Their applicability at other building sites with different ambient air  
397 condition and building characteristics is rather limited. Buildings with natural ventilation and  
398 operable windows could have very different set of significant predictors than discovered in this  
399 paper which in turn would affect the model performance. This issue has also been raised in  
400 Challoner et al.<sup>11</sup> and Wei et al.<sup>10</sup> Future field studies covering various climate regions and  
401 building types could validate and improve the results obtained from small-scale investigations. As

402 the low-cost PM sensors become more reliable and widely used, it would require less effort to  
403 conduct these large-scale investigations to obtain more generalized findings. IAQ simulation tools  
404 such as CONTAM<sup>52</sup> could also serve as another avenue for validating the predictive models from  
405 a physical and mechanical perspective. When the simulation is coupled with commercial reference  
406 buildings<sup>53</sup>, the results could be applicable to the most common commercial buildings.

407 In summary, this paper shows that it is feasible to develop a relatively accurate indoor PM<sub>2.5</sub>  
408 prediction model for well-mixed air in a mechanically ventilated office space using some readily  
409 available meteorological and building-related variables. A straightforward indoor PM<sub>2.5</sub> prediction  
410 model could provide the building owner, facility manager, and occupants insight into the average  
411 air quality of the space and empower the stakeholders to make informed decisions related to the  
412 management of the indoor environment.

413 In the future, researchers should continue to explore not just prediction models, but also how to  
414 optimize the cost and accessibility of prediction relative to accuracy. The quantity, quality, and  
415 placement of sensors augmented by external information and machine learning models are critical  
416 to widespread access to such systems. Furthermore, research in this field should progress from  
417 prediction to active management, where predictive models such as those presented in this paper  
418 are used to actively improve the efficiency of building operations and the quality of life of the  
419 building's inhabitants.

## 420 ASSOCIATED CONTENT

### 421 **Supporting Information (SI)**

422 The Supporting Information is available free of charge at <http://pubs.acs.org>.

423 Texts: building characteristics and sampling locations; particle mass concentration calculation;  
424 occupancy sensor background; ANN background; additional model results.

425 Figures: location of the UW Tower and other monitoring sites; overview of the study site and  
426 sampling locations; measurement duration for all of the variables; data splitting scheme for the  
427 cross-validation; time series plot of the indoor and outdoor PM<sub>2.5</sub> measurements; plots of predicted  
428 and observed values for Phase 1 and 3 models; cross-validation plots of the PLS models; loading  
429 values of each variable in Models P2 and P3; boxplots of PM<sub>2.5</sub> level at different hours.

430 Tables: descriptions of the surrounding environment of each ambient PM<sub>2.5</sub> monitoring site; list  
431 of hyperparameters used for the ANN and LSTM; hyperparameter search space; descriptive  
432 statistics of the variables; Pearson's correlation coefficients between independent variables;  
433 performance indicators of Phase 1 and 3 models; independent variables included in Phase 1 and 3  
434 models; summary of MLR models M2 and M3; percentage of variance explained by each  
435 component of the PLS models; hyperparameters of the best ANN and LSTM models; recursive  
436 model results.

## 437 AUTHOR INFORMATION

### 438 **Corresponding Author**

439 \*Phone: +1 (206) 685-0228; E-mail: [amyakim@uw.edu](mailto:amyakim@uw.edu).

### 440 **Author Contributions**

441 The contributions of each author are: conceptualization, B.L., T.L., and A.K.; methodology, B.L.,  
442 S.W., and T.L.; data curation, B.L. and S.W.; writing – original draft preparation, B.L., S.W. and  
443 A.K.; writing – review and editing, T.L.; visualization, S.W.; and supervision, A.K.

### 444 **Notes**

445 The authors declare no competing financial interest.

## 446 ACKNOWLEDGMENT

447 The authors would like to thank Troy Swanson, UW Tower Facility Manager; Kim Lokan, UW-  
448 IT Facilities Services Manager. This material is based upon work supported by the National  
449 Science Foundation under Grant No. 1852995.

## 450 REFERENCES

- 451 1. Burnett, R.; Chen, H.; Szyszkowicz, M.; Fann, N.; Hubbell, B.; Pope, C. A.; Apte, J.  
452 S.; Brauer, M.; Cohen, A.; Weichenthal, S.; Coggins, J.; Di, Q.; Brunekreef, B.; Frostad, J.;  
453 Lim, S. S.; Kan, H.; Walker, K. D.; Thurston, G. D.; Hayes, R. B.; Lim, C. C.; Turner, M. C.;  
454 Jerrett, M.; Krewski, D.; Gapstur, S. M.; Diver, W. R.; Ostro, B.; Goldberg, D.; Crouse, D. L.;  
455 Martin, R. V.; Peters, P.; Pinault, L.; Tjepkema, M.; van Donkelaar, A.; Villeneuve, P. J.;  
456 Miller, A. B.; Yin, P.; Zhou, M.; Wang, L.; Janssen, N. A. H.; Marra, M.; Atkinson, R. W.;  
457 Tsang, H.; Quoc Thach, T.; Cannon, J. B.; Allen, R. T.; Hart, J. E.; Laden, F.; Cesaroni, G.;  
458 Forastiere, F.; Weinmayr, G.; Jaensch, A.; Nagel, G.; Concini, H.; Spadaro, J. V., Global  
459 estimates of mortality associated with long-term exposure to outdoor fine particulate matter.  
460 *Proceedings of the National Academy of Sciences* **2018**, *115* (38), 9592-9597.
- 461 2. Cesaroni, G.; Badaloni, C.; Gariazzo, C.; Stafoggia, M.; Sozzi, R.; Davoli, M.;  
462 Forastiere, F., Long-Term Exposure to Urban Air Pollution and Mortality in a Cohort of More than  
463 a Million Adults in Rome. *Environmental Health Perspectives* **2013**, *121* (3), 324-331.
- 464 3. Cohen, A. J.; Brauer, M.; Burnett, R.; Anderson, H. R.; Frostad, J.; Estep, K.;  
465 Balakrishnan, K.; Brunekreef, B.; Dandona, L.; Dandona, R.; Feigin, V.; Freedman, G.;  
466 Hubbell, B.; Jobling, A.; Kan, H.; Knibbs, L.; Liu, Y.; Martin, R.; Morawska, L.; Pope, C. A.,  
467 III; Shin, H.; Straif, K.; Shaddick, G.; Thomas, M.; van Dingenen, R.; van Donkelaar, A.; Vos,  
468 T.; Murray, C. J. L.; Forouzanfar, M. H., Estimates and 25-year trends of the global burden of  
469 disease attributable to ambient air pollution: an analysis of data from the Global Burden of  
470 Diseases Study 2015. *The Lancet* **2017**, *389* (10082), 1907-1918.
- 471 4. Li, T.; Zhang, Y.; Wang, J.; Xu, D.; Yin, Z.; Chen, H.; Lv, Y.; Luo, J.; Zeng, Y.; Liu,  
472 Y.; Kinney, P. L.; Shi, X., All-cause mortality risk associated with long-term exposure to ambient  
473 PM<sub>2.5</sub> in China: a cohort study. *The Lancet Public Health* **2018**, *3* (10), e470-e477.
- 474 5. Pope III, C. A.; Lefler, J. S.; Ezzati, M.; Higbee, J. D.; Marshall, J. D.; Kim, S.-Y.;  
475 Bechle, M.; Gilliat, K. S.; Vernon, S. E.; Robinson, A. L.; Burnett, R. T., Mortality Risk and  
476 Fine Particulate Air Pollution in a Large, Representative Cohort of U.S. Adults. *Environmental*  
477 *Health Perspectives* **2019**, *127* (7), 077007.
- 478 6. Schraufnagel, D. E.; Balmes, J. R.; Cowl, C. T.; De Matteis, S.; Jung, S.-H.; Mortimer,  
479 K.; Perez-Padilla, R.; Rice, M. B.; Riojas-Rodriguez, H.; Sood, A.; Thurston, G. D.; To, T.;  
480 Vanker, A.; Wuebbles, D. J., Air Pollution and Noncommunicable Diseases: A Review by the  
481 Forum of International Respiratory Societies' Environmental Committee, Part 1: The Damaging  
482 Effects of Air Pollution. *CHEST* **2019**, *155* (2), 409-416.
- 483 7. Klepeis, N. E.; Nelson, W. C.; Ott, W. R.; Robinson, J. P.; Tsang, A. M.; Switzer, P.;  
484 Behar, J. V.; Hern, S. C.; Engelmann, W. H., The National Human Activity Pattern Survey



- 485 (NHAPS): a resource for assessing exposure to environmental pollutants. *Journal of Exposure*  
486 *Science & Environmental Epidemiology* **2001**, *11*, 231-252.
- 487 8. Bureau of Labor Statistics Graphics for Economic News Releases.  
488 <https://www.bls.gov/charts/american-time-use/activity-by-sex.htm#> (accessed November 29,  
489 2019).
- 490 9. EPA Air Quality System (AQS). <https://www.epa.gov/aqs> (accessed January 27, 2020).
- 491 10. Wei, W.; Ramalho, O.; Malingre, L.; Sivanantham, S.; Little, J. C.; Mandin, C., Machine  
492 learning and statistical models for predicting indoor air quality. *Indoor Air* **2019**, *29* (5), 704-726.
- 493 11. Challoner, A.; Pilla, F.; Gill, L., Prediction of Indoor Air Exposure from Outdoor Air  
494 Quality Using an Artificial Neural Network Model for Inner City Commercial Buildings.  
495 *International Journal of Environmental Research and Public Health* **2015**, *12* (12), 15233-15253.
- 496 12. Putra, J. C. P.; Safrilah, S.; Ihsan, M. In *The prediction of indoor air quality in office room*  
497 *using artificial neural network*, Proceedings of the 4th International Conference on Engineering,  
498 Technology, and Industrial Application, Surakarta, Indonesia, AIP Conference Proceedings:  
499 Surakarta, Indonesia, 2018; p 020040.
- 500 13. Sofuoglu, S. C., Application of artificial neural networks to predict prevalence of building-  
501 related symptoms in office buildings. *Building and Environment* **2008**, *43* (6), 1121-1126.
- 502 14. Xie, H.; Ma, F.; Bai, Q. In *Prediction of Indoor Air Quality Using Artificial Neural*  
503 *Networks*, 2009 Fifth International Conference on Natural Computation, Tianjin, China, Aug 14-  
504 16; Wang, H.; Low, K. S.; Wei, K.; Sun, J., Eds. IEEE: Tianjin, China, 2009; pp 414-418.
- 505 15. Chen, X.; Zheng, Y.; Chen, Y.; Jin, Q.; Sun, W.; Chang, E.; Ma, W.-Y. In *Indoor air*  
506 *quality monitoring system for smart buildings*, UbiComp '14: Proceedings of the 2014 ACM  
507 International Joint Conference on Pervasive and Ubiquitous Computing, Seattle, WA, USA,  
508 Association for Computing Machinery: Seattle, WA, USA, 2014; pp 471-475.
- 509 16. Bai, Y.; Zeng, B.; Li, C.; Zhang, J., An ensemble long short-term memory neural network  
510 for hourly PM<sub>2.5</sub> concentration forecasting. *Chemosphere* **2019**, *222*, 286-294.
- 511 17. Kim, H. S.; Park, I.; Song, C. H.; Lee, K.; Yun, J. W.; Kim, H. K.; Jeon, M.; Lee, J.;  
512 Han, K. M., Development of a daily PM<sub>10</sub> and PM<sub>2.5</sub> prediction system using a deep long short-  
513 term memory neural network model. *Atmospheric Chemistry and Physics* **2019**, *19* (20), 12935-  
514 12951.
- 515 18. Li, X.; Peng, L.; Yao, X.; Cui, S.; Hu, Y.; You, C.; Chi, T., Long short-term memory  
516 neural network for air pollutant concentration predictions: Method development and evaluation.  
517 *Environmental Pollution* **2017**, *231* (Part 1), 997-1004.
- 518 19. Qi, Y.; Li, Q.; Karimian, H.; Liu, D., A hybrid model for spatiotemporal forecasting of  
519 PM<sub>2.5</sub> based on graph convolutional neural network and long short-term memory. *Science of The*  
520 *Total Environment* **2019**, *664*, 1-10.
- 521 20. Zhao, J.; Deng, F.; Cai, Y.; Chen, J., Long short-term memory - Fully connected (LSTM-  
522 FC) neural network for PM<sub>2.5</sub> concentration prediction. *Chemosphere* **2019**, *220*, 486-492.
- 523 21. Puget Sound Clean Air Agency Air Quality. <https://www.pscleanair.gov/27/Air-Quality>  
524 (accessed 2020 July 29).
- 525 22. Washington State Department of Ecology Washington's Air Monitoring Network.  
526 <https://fortress.wa.gov/ecy/enwiwa/> (accessed November 30, 2019).
- 527 23. University of Washington Rooftop Observations - ATG building UW.  
528 [https://a.atmos.washington.edu/cgi-bin/list\\_uw.cgi](https://a.atmos.washington.edu/cgi-bin/list_uw.cgi) (accessed 2020 July 29).

- 529 24. Particles Plus Inc. 7301-AQM and 7302-AQM Remote Air Quality and Environmental  
530 Monitor. <https://particlesplus.com/7301-iaq-remote-particle-counter/> (accessed December 1,  
531 2019).
- 532 25. NOAA Radiance Research Nephelometer.  
533 [https://www.esrl.noaa.gov/gmd/aero/instrumentation/RR\\_neph.html](https://www.esrl.noaa.gov/gmd/aero/instrumentation/RR_neph.html) (accessed December 1,  
534 2019).
- 535 26. Elbayoumi, M.; Ramli, N. A.; Yusof, N. F. F. M., Development and comparison of  
536 regression models and feedforward backpropagation neural network models to predict seasonal  
537 indoor PM<sub>2.5-10</sub> and PM<sub>2.5</sub> concentrations in naturally ventilated schools. *Atmospheric Pollution*  
538 *Research* **2015**, *6* (6), 1013-1023.
- 539 27. Elbayoumi, M.; Ramli, N. A.; Yusof, N. F. F. M.; Yahaya, A. S. B.; Al Madhoun, W.;  
540 Ul-Saufie, A. Z., Multivariate methods for indoor PM10 and PM2.5 modelling in naturally  
541 ventilated schools buildings. *Atmospheric Environment* **2014**, *94*, 11-21.
- 542 28. Marengo, E.; Bobba, M.; Robotti, E.; Liparota, M. C., Modeling of the Polluting  
543 Emissions from a Cement Production Plant by Partial Least-Squares, Principal Component  
544 Regression, and Artificial Neural Networks. *Environmental Science and Technology* **2006**, *40* (1),  
545 272-280.
- 546 29. Huang, T.; Yu, Y.; Wei, Y.; Wang, H.; Huang, W.; Chen, X., Spatial-seasonal  
547 characteristics and critical impact factors of PM2.5 concentration in the Beijing-Tianjin-Hebei  
548 urban agglomeration. *PLoS ONE* **2018**, *13* (9), e0201364.
- 549 30. Kim, M.; SankaraRao, B.; Kang, O.; Kim, J.; Yoo, C., Monitoring and prediction of  
550 indoor air quality (IAQ) in subway or metro systems using season dependent models. *Energy and*  
551 *Buildings* **2012**, *46*, 48-55.
- 552 31. Lee, S.; Kim, M. J.; Kim, J. T.; Yoo, C. K., In search for modeling predictive control of  
553 indoor air quality and ventilation energy demand in subway station. *Energy and Buildings* **2015**,  
554 *98*, 56-65.
- 555 32. Lim, J.; Kim, Y.; Oh, T.; Kim, M.; Kang, O.; Kim, J. T.; Kim, I.-W.; Kim, J.-C.; Jeon,  
556 J.-S.; Yoo, C., Analysis and prediction of indoor air pollutants in a subway station using a new key  
557 variable selection method. *Korean Journal of Chemical Engineering* **2012**, *29* (8), 994-1003.
- 558 33. Che, Z.; Purushotham, S.; Cho, K.; Sontag, D.; Liu, Y., Recurrent Neural Networks for  
559 Multivariate Time Series with Missing Values. *Scientific Reports* **2018**, *8*, 6085.
- 560 34. Hyndman, R. J.; Athanasopoulos, G., *Forecasting: principles and practice*. OTexts:  
561 Melbourne, Australia, 2018.
- 562 35. EPA, National Primary and Secondary Ambient Air Quality Standards. U.S.  
563 Environmental Protection Agency, Ed. 2016; Vol. 40 CFR Part 50.
- 564 36. R Core Team *R: A Language and Environment for Statistical Computing*, R Foundation  
565 for Statistical Computing: Vienna, Austria, 2019.
- 566 37. Kuhn, M. *caret: Classification and Regression Training*, R Foundation for Statistical  
567 Computing: Vienna, Austria, 2019.
- 568 38. Ripley, B. *MASS: Support Functions and Datasets for Venables and Ripley's MASS*, R  
569 Foundation for Statistical Computing: Vienna, Austria, 2019.
- 570 39. Mansfield, E. R.; Helms, B. P., Detecting Multicollinearity. *The American Statistician*  
571 **1982**, *36* (3 Part 1), 158-160.
- 572 40. Mevik, B.-H.; Wehrens, R.; Liland, K. H. *pls: Partial Least Squares and Principal*  
573 *Component Regression*, R Foundation for Statistical Computing: Vienna, Austria, 2019.

- 574 41. Dayal, B. S.; MacGregor, J. F., Improved PLS algorithms. *Journal of Chemometrics* **1997**,  
575 *11* (1), 73-85.
- 576 42. van der Voet, H., Comparing the predictive accuracy of models using a simple  
577 randomization test. *Chemometrics and Intelligent Laboratory Systems* **1994**, *25* (2), 313-323.
- 578 43. Demirhan, H. *dLagM: Time Series Regression Models with Distributed Lag Models*, R  
579 Foundation for Statistical Computing: Vienna, Austria, 2020.
- 580 44. Friedman, J.; Hastie, T.; Tibshirani, R.; Narasimhan, B.; Tay, K.; Simon, N.; Qian, J.  
581 *glmnet: Lasso and Elastic-Net Regularized Generalized Linear Models*, R Foundation for  
582 Statistical Computing: Vienna, Austria, 2020.
- 583 45. Jafta, N.; Barregard, L.; Jeena, P. M.; Naidoo, R. N., Indoor air quality of low and middle  
584 income urban households in Durban, South Africa. *Environmental Research* **2017**, *156*, 47-56.
- 585 46. Yuchi, W. Modelling Fine Particulate Matter Concentrations inside the Homes of Pregnant  
586 Women in Ulaanbaatar, Mongolia. Master of Science, Simon Fraser University, Burnaby, British  
587 Columbia, Canada, 2017.
- 588 47. Yuchi, W.; Gombojav, E.; Boldbaatar, B.; Galsuren, J.; Enkhmaa, S.; Beejin, B.;  
589 Naidan, G.; Ochir, C.; Legtseg, B.; Byambaa, T.; Barn, P.; Henderson, S. B.; Janes, C. R.;  
590 Lanphear, B. P.; McCandless, L. C.; Takaro, T. K.; Venners, S. A.; Webster, G. M.; Allen, R.  
591 W., Evaluation of random forest regression and multiple linear regression for predicting indoor  
592 fine particulate matter concentrations in a highly polluted city. *Environmental Pollution* **2019**, *245*,  
593 746-753.
- 594 48. Kim, M.; Kim, Y.; Sung, S.; Yoo, C. In *Data-driven prediction model of indoor air quality*  
595 *by the preprocessed recurrent neural networks*, Proceedings of the ICROS-SICE International  
596 Joint Conference 2009, Fukuoka City, Japan, IEEE: Fukuoka City, Japan, 2009; pp 1688-1692.
- 597 49. Loy-Benitez, J.; Vilela, P.; Li, Q.; Yoo, C., Sequential prediction of quantitative health  
598 risk assessment for the fine particulate matter in an underground facility using deep recurrent  
599 neural networks. *Ecotoxicology and Environmental Safety* **2019**, *169*, 316-324.
- 600 50. Das, P.; Shrubsole, C.; Jones, B.; Hamilton, I.; Chalabi, Z.; Davies, M.; Mavrogianni,  
601 A.; Taylor, J., Using probabilistic sampling-based sensitivity analyses for indoor air quality  
602 modelling. *Building and Environment* **2014**, *78*, 171-182.
- 603 51. Gundel, L. A.; Destailats, H., Aerosol Chemistry and Physics - An Indoor Perspective. In  
604 *Aerosols handbook : measurement, dosimetry, and health effects*, 2nd ed.; Ruzer, L. S.; Harley, N.  
605 H., Eds. CRC Press: Boca Raton, FL, 2013.
- 606 52. NIST CONTAM. <https://www.nist.gov/services-resources/software/contam> (accessed  
607 August 2).
- 608 53. DOE Commercial Reference Buildings.  
609 <https://www.energy.gov/eere/buildings/commercial-reference-buildings> (accessed August 2).

610