Closing the gap on lower cost air quality monitoring: machine learning calibration models to improve low-cost sensor performance

Naomi Zimmerman¹, Albert A. Presto¹, Sriniwasa P.N. Kumar¹, Jason Gu², Aliaksei Hauryliuk¹, Ellis S. Robinson¹, Allen L. Robinson¹, R. Subramanian¹

⁵ ¹Center for Atmospheric Particle Studies, Carnegie Mellon University, Pittsburgh, 15213, USA ²Sensevere LLC, Pittsburgh, 15222, USA

Correspondence to: R. Subramanian (subu@cmu.edu)

Supplemental.





10 Figure S1. Average reference monitor concentrations during training and testing windows for each RAMP



Figure S2: Choosing collocation length for training. Optimization assessed on a subset of three RAMPs (RAMP #2, #4, and #14). Ultimately a 4-week collocation period was chosen as being the best across all four species.

Table S1: Change in model performance if one consecutive 4-week colocation at the beginning of the study is conducted vs spacing out training data throughout the study in 8 half week increments.

	Impact of using one consecutive four week training window vs. distributed half week colocatio				
Metric	CO	CO ₂	NO ₂	O 3	
RMSE	+11.9 ppb	+3.2 ppm	+0.4 ppb	+1.6 ppb	
MAE	+12.1 ppb	+1.8 ppm	+0.4 ppb	+1.6 ppb	
Pearson r	-0.01	-0.03	-0.08	-0.05	

Table S2. Metrics used for comparing sensor data. M indicates a value measured by one of the sensors participating in the experiment and O indicates the observations from the reference measurements.

Statistic	Abbrev.	Formula	Characteristics	
Mean Bias	MBE	$MBE = \overline{M} - \overline{O}$	Estimation of the magnitude of differences (bias) between sensors estimation and reference values	
Error			averaged over the whole sampling period	
Mean		$1 \sum_{n=1}^{n}$		
Absolute	MAE	MAE = $\frac{1}{n} \sum M_i - O_i $	• Indicates the average of the magnitude of the errors.	
Error			• Sensitive to outliers.	
Pearson		$\sum_{i=1}^{n} (M_i - \overline{M}) (O_i - \overline{O})$	Manner the strength and the line time of a line	
Correlation	r	$r = \frac{1}{\sqrt{2}}$	relationship between two variables.	
Coefficient		$\sqrt{\sum_{i=1}^{n} (M_i - M)^2 (O_i - O)^2}$		
Root Mean		$1\sum_{n=1}^{n}$		
Square Error	RMSE	$RMSE = \left \frac{1}{n} \sum (M_i - O_i)^2 \right $	• Magnitude of the error and retains the variable's unit	
-		$\sqrt{i=1}$	• Sensitive to extreme values and to outliers	
Centred Root				
Mean Square	CRMSE	$CRMSE = \sqrt{RMSE^2 - MBE^2}$	RMSE corrected for bias	
Error			• Measure of random error	

5 s

Plots of Goodness of Fit (Results from Training, Figures S3-S6)



Figure S3: Goodness of fit on training data across all 19 RAMPs for CO using random forests.



Figure S4: Goodness of fit on training data across all 19 RAMPs for NO2 using random forests.



Figure S5: Goodness of fit on training data across all 19 RAMPs for O3 using random forests.



Figure S6: Goodness of fit on training data across all 19 RAMPs for CO₂ using random forests.

Plots of Random Forest Performance by RAMP (Results from Testing, Figures S7-S10)



Figure S7: Model performance on testing data across 16 RAMPs for CO using random forests.



Figure S8: Model performance on testing data across 10 RAMPs for NO2 using random forests.



Figure S9: Model performance on testing data across 16 RAMPs for O3 using random forests.



Figure S10: Model performance on testing data across 15 RAMPs for CO2 using random forests.

US EPA Air Sensor Guidebook Precision and Bias Estimators:

Precision:

5 The precision estimator (CV) is the upper bound of the 90% confidence interval on the coefficient of variation.

$$CV = \sqrt{\frac{n \cdot \sum_{i=1}^{n} d_i^2 - (\sum_{i=1}^{n} d_i)^2}{n(n-1)}} \cdot \sqrt{\frac{n-1}{\chi^2_{(0.1,n-1)}}}$$

Where n is the number of data points, $\chi^2_{(0.1,n-1)}$ is the 10th percentile of a chi-squared distribution with n-1 degrees of freedom, and d_i is equal to:

$$d_i = \frac{RAMP - reference}{reference} \cdot 100\%$$

Bias:

15

The bias estimator is the upper bound of the 95% confidence interval on the mean absolute value of the percent difference between the RAMPs and the reference monitor.

$$|\text{Bias}| = \text{AB} + t_{0.95,n-1} \cdot \frac{\text{AS}}{\sqrt{n}}$$

Where $t_{0.95,n-1}$ is the 95th quartile of a t-distibution with n-1 degrees of freedom and AB is the mean of the absolute values of the d_i's.

$$AB = \frac{1}{n} \sum_{i=1}^{n} |\mathbf{d}_i|$$

And AS is the standard deviation of the absolute value of the d_i's:

$$AS = \sqrt{\frac{n \cdot \sum_{i=1}^{n} |d_i|^2 - (\sum_{i=1}^{n} |d_i|)^2}{n(n-1)}}$$

25