## Supporting Information

# Precursor apportionment of atmospheric oxygenated organic molecules using a machine learning method

Xiaohui Qiao<sup>1</sup>, Xiaoxiao Li<sup>1</sup>, Chao Yan<sup>2,3,4</sup>, Nina Sarnela<sup>3</sup>, Rujing Yin<sup>1</sup>, Yishuo Guo<sup>4</sup>, Lei Yao<sup>3,4</sup>, Wei Nie<sup>2</sup>, Dandan Huang<sup>5</sup>, Zhe Wang<sup>6</sup>, Federico Bianchi<sup>3,4</sup>, Yongchun Liu<sup>4</sup>, Neil M. Donahue<sup>7,8</sup>, Markku Kulmala<sup>3,4</sup>, Jingkun Jiang<sup>1,\*</sup>

<sup>1</sup> State Key Joint Laboratory of Environment Simulation and Pollution Control, School of Environment, Tsinghua University, 100084 Beijing

<sup>2</sup> Joint International Research Laboratory of Atmospheric and Earth System Research, School of Atmospheric Sciences, Nanjing University, Nanjing, China

<sup>3</sup> Institute for Atmospheric and Earth System Research/Physics, Faculty of Science, University of Helsinki, 00014 Helsinki, Finland

<sup>4</sup> Aerosol and Haze Laboratory, Beijing Advanced Innovation Center for Soft Matter Science and Engineering, Beijing University of Chemical Technology, 100029 Beijing, China

<sup>5</sup> State Environmental Protection Key Laboratory of Formation and Prevention of Urban Air

Pollution Complex, Shanghai Academy of Environmental Sciences, Shanghai, China

<sup>6</sup>Division of Environment and Sustainability, The Hong Kong University of Science and Technology, Hong Kong SAR, China

<sup>7</sup>Center for Atmospheric Particle Studies, Carnegie Mellon University, Pittsburgh, PA, USA

<sup>8</sup> Department of Chemistry, Carnegie Mellon University, Pittsburgh, PA, USA

#### **Evaluation of precursor apportionment model**

Table S1 is a confusion matrix for the prediction results. It is taken as an example to explain the evaluation of the precursor apportionment model. TP, FP, or TN is a description of the reality and predictions, in which T is for true, F is for false, P is for positive, and N is for negative. Aromatic was taken as an example in the table. Thus, when aromatics are predicted as aromatics, the result is TP; when aromatics are predicted as aliphatics and monoterpene, the result is TN; when aliphatics and monoterpene are predicted as aromatics, the result is FP. In one word, the first letter is for reality and the latter is for prediction.

Label\Prediction	Aromatics	Aliphatics	Monoterpenes
Aromatics	TP (₽)	TN (원)	TN (印)
Aliphatics	NP (원)		
Monoterpenes	NP (원)		

Table S1. An example of the confusion matrix for a prediction.

*Precision, recall*, and *F1-score* are generally used to evaluate the performance of machine learning models. They are calculated as below,

$$Precision = \frac{TP}{TP + FP}$$
(1)

$$\frac{TP}{Recall = \overline{TP + TN}}$$
(2)

$$F1 \ score = \frac{2 * Precision * Recall}{Precision + Recall}$$
(3)

#### **Cross-validation**

The model was cross-validated using 10-fold cross-validation, which randomly partitioning the training dataset into 10 sets, then reserving each set as validation data, and training the model using the other k-1 sets. Moreover, in this study, we evaluated the 10 models above with the same testing dataset, the 30% of original dataset that was reserved at the beginning, to present a more contrasting result (Figure 2 & S3).

Table S2 shows two examples of the determination of overlapping molecules. In example 1, ten trees give unbalanced votes, while in example 2, ten trees give balanced votes. According to the rules of voting, the molecule in example 1 is attributed to monoterpene, while the molecule in example 2 is attributed to an overlapping OOM oxidized from aromatic, aliphatic, or monoterpene.

	Example 1	Example 2
Aromatic OOMs	3	3
Aliphatic OOMs	2	4
Monoterpene OOMs	4	3
Isoprene OOMs	1	0
Mean+Standard deviation	2.5+1.29=3.79	2.5+1.73=4.23
Mean-Standard deviation	2.5-1.29=1.21	2.5-1.73=0.77
Final	Monoterpene	Aromatic Aliphatic Monoterpene

Table S2. Examples of the voting strategy.



**Figure S1.** The Venn-plot of oxygenated organic molecules in Table 1. It should be noted that these aliphatic OOMs are without expansion.



**Figure S2.** The apportionment result for oxygenated organic molecules in urban Beijing using the workflow method<sup>1</sup>.

### **Expansion of aliphatic OOMs peaklist**

In this study, we test the case that artifically construct molecules homogenous to the studied aliphatic OOMs. As shown in Table S2, we add or minus a integer number of -CH2 to those studied aliphatic OOMs (63 compounds) to obtain an enlarged peaklist of aliphatics (346 compounds). It should be mentioned that the reported aliphatic OOMs are mainly C6 and C10 compounds, we processed every molecule in the same way. Finally, there were some duplicate molecules, and only one of them was retained.

minus C <sub>n</sub> H <sub>2n</sub> (n=2,3)	minus CH <sub>2</sub> (C9)	Studied OOMs (C10)	add CH <sub>2</sub> (C11)	add $C_nH_{2n}$ (n=2,3,4)
$\begin{array}{c} C_{10\text{-n}}H_{20\text{-}2n}O_2\\ C_{10\text{-n}}H_{21\text{-}2n}O_2\\ C_{10\text{-n}}H_{20\text{-}2n}O_3\\ C_{10\text{-n}}H_{20\text{-}2n}O_3\\ \end{array}$	$\begin{array}{c} C_{9}H_{18}O_{2}\\ C_{9}H_{19}O_{2}\\ C_{9}H_{18}O_{3}\\ C_{9}H_{19}O_{3} \end{array}$	$\begin{array}{c} C_{10}H_{20}O_2\\ C_{10}H_{21}O_2\\ C_{10}H_{20}O_3\\ C_{10}H_{21}O_3 \end{array}$	$\begin{array}{c} C_{11}H_{22}O_2\\ C_{11}H_{23}O_2\\ C_{11}H_{22}O_3\\ C_{11}H_{23}O_3 \end{array}$	$\begin{array}{c} C_{10+n}H_{20+2n}O_2\\ C_{10+n}H_{21+2n}O_2\\ C_{10+n}H_{20+2n}O_3\\ C_{10+n}H_{21+2n}O_3\\ \end{array}$
	C <sub>9</sub> H <sub>19</sub> NO <sub>3</sub>	C <sub>10</sub> H <sub>21</sub> NO <sub>3</sub>	C <sub>11</sub> H <sub>23</sub> NO <sub>3</sub>	

Table S3. Examples of aliphatic O	OOMs peaklist expansion
-----------------------------------	-------------------------

#### The impact of data expansion

There is an increase in the overall accuracy of decision trees when comparing the evaluation result of the model trained with the original dataset (Figure S5) and that trained with the expanded dataset (Figure S2a). The *F1-score* of aliphatic OOMs increased from ~0.3 to ~0.5. For aromatic OOMs and monoterpene OOMs, however, it did not change evidently. Although the accuracy of identifying isoprene OOMs has decreased, the overall accuracy for these four precursors has improved generally.

Moreover, the application results showed that the proportion of monoterpene OOMs in Beijing decreased  $\sim 1\%$  and that of aliphatic OOMs increased by  $\sim 4\%$  when using the model trained with the expanded dataset in comparison to that trained with the original dataset. For Hyytiälä, monoterpene OOMs increased 3%, aromatic OOMs decreased 3% when using the model trained with the expanded dataset. The proportion of aliphatic OOMs increased and that of isoprene OOMs decreased.



**Figure S3.** The overall accuracy of the decision tree model in precursor apportionment of atmospheric oxygenated organic molecules using laboratory dataset without the expansion of aliphatic OOMs.



**Figure S4.** The proportion of the identified oxygenated organic molecules that are oxidized from monoterpenes, aliphatics, aromatics, and isoprene in Beijing and Hyytiälä predicted by model trained with laboratory dataset without the expansion of aliphatic OOMs: (a) and (c) are the proportions of OOM species; (b) and (d) are the proportions of OOM number concentration.



Figure S5. The *F1-score* of models consisting of different number of trees.

Along with the increasing number of trees, the *F1-score* of the trained model fluctuates without significant improvement. Thus, ten trees was adopted in this study. The performance of models with more trees may be limited by the imbalance training dataset, which can be improved by supplementing with more comprehensive laboratory data.