



Supplement of

Constraining the response factors of an extractive electrospray ionization mass spectrometer for near-molecular aerosol speciation

Dongyu S. Wang et al.

Correspondence to: Dongyu S. Wang (dongyu.wang@psi.ch), Imad El Haddad (imad.el-haddad@psi.ch), Jay G. Slowik (jay.slowik@psi.ch), and David M. Bell (david.bell@psi.ch)

The copyright of individual parts of the supplement might differ from the article licence.

Section S1. Condensation sink estimation

The condensation sink, CS (Lehtinen et al., 2003; Dal Maso et al., 2002) in s⁻¹ is calculated as

$$CS = 2\pi D \int_0^\infty d_p \cdot \beta_M(d_p) \cdot N(d_p) \cdot dd_p \qquad \qquad \text{Eq. (S1)}$$

where *D* is the vapor diffusivity, d_P is the particle diameter, the $N(d_P)$ is the number of particle of diameter d_P , and $\beta_M(d_P)$ is the Fuchs-Sutugin correction factor for gas-phase diffusion over particles in the transition regime. Using a discrete particle size distribution as measured by the SMPS, we calculate *CS* using an approximation of the integral, namely

$$CS = 2\pi D \sum_{i} \beta_{i} d_{pi} N_{i}$$
 Eq. (S2)

The lifetime for gaseous condensation in the presence of a *CS* is (Markku Kulmala and Wagner 2001)

$$\tau_{cond} = \frac{1}{CS}$$
 Eq. (S3)

As an approximation, *D* can be assumed to be 6 to 7 x 10^{-6} m² s⁻¹ for condensable organic vapors (Palm et al., 2016; Krechmer et al., 2017). A more nuanced estimation is described below. The Fuchs-Sutugin correction factor β , is calculated

$$\beta = \frac{K_n + 1}{0.377 \cdot K_n + \frac{4}{3} \cdot \alpha^{-1} \cdot K_n + \frac{4}{3} \cdot \alpha^{-1} \cdot K_n^2 + 1}$$
Eq. (S4)

where α is the mass accommodation coefficient. In lieu of empirical values, unity is assumed for α (Markku et al., 2001). Recent experimental results support this unity assumption (Krechmer et al., 2017; Liu et al., 2019). K_n is the Knudsen number,

$$K_n = \frac{2\lambda}{d_p}$$
 Eq. (S5)

where the particle radius $(d_p/2)$ is used as the characteristic length; λ is the effective free mean path of vapor molecules. The mean free path in dry air varies slightly in the literature, e.g. 6.53 to 6.673 x 10⁻⁸ m (Jennings 1988). The mean free path of any organic compound can be calculated if its gas-phase diffusion coefficient (at the bath gas pressure), D_{Pr} , and average molecular speed *c*, are known,

$$\lambda = \frac{3D_{pr}}{c}$$
Eq. (S6)
$$c = \sqrt{\frac{8RT}{\pi MW}}$$
Eq. (S7)

where *R* is the ideal gas constant (8.314 J mol⁻¹ K⁻¹), T is the temperature in K, and *MW* is the molar mass (kg mol⁻¹). Note that D_{pr} is a function of bath gas pressure, *P* (Torr), and the gas diffusivity, *D* (Torr cm⁻² s⁻¹)

$$D_{pr} = \frac{D}{P}$$
 Eq. (S8)

For a trace gas *A* in a bath gas *B*, the gas diffusivity could be estimated using Fuller's method (Fuller et al., 1966; Tang et al., 2015),

$$D(A,B) = \frac{1.0868 \times T^{1.75}}{\sqrt{m(A,B)} \left(\sqrt[3]{V_A} + \sqrt[3]{V_B}\right)^2}$$
 Eq. (S9)

where V_A and V_B are dimensionless diffusion volumes of A, and B; m(A,B) is the reduced mass of the A-B pair and can be calculated based on the molecular masses (g mol⁻¹) of A and B, m_A and m_B , respectively

$$m(A,B) = \frac{2}{(1/m_A + 1/m_B)}$$
 Eq. (S10)

 V_A may be estimated from the molecular formula of the trace gas

$$V = \sum n_i V_i$$
 Eq. (S11)

where n_i is the number of atoms with diffusion volume of V_i , which is 15.9 for C, 2.31 for H, 6.11 for O, and 4.54 for N (Reid et al., 1987). Subtracting 18.3 from the total diffusion volume accounts for the effect of the aromatic ring. For compounds containing multiple aromatic rings, it maybe be best to correct only for independent aromatic rings, based on limited experimental data (Tang et al., 2015). Alicyclic rings are not expected to have an effect on the diffusion volume (Tang et al., 2015). Diffusion volumes of common bath gasses are known instead of estimated: N₂ (18.5), O₂ (19.7), H₂O (13.1). For inorganic and slightly oxygenated organic compounds, the mean free path of condensable vapors may be quite uniform (within 20%), where the Knudsen number can be estimated based on pressure and particle diameter alone (Tang et al., 2015),

$$K_n = \frac{2}{d_P} \times \frac{\lambda_P}{P}$$
 Eq. (S12)

where *P* is the pressure of air in atm, and λ_P is the pressure normal mean free path equal to 100 nm atm. The deviation of K_n estimated using Eq. S12 for a 100 nm particle (i.e. $K_n = 2$) with respect to that estimated using Eq. S5 for selected compounds is shown in Table S1 with the corresponding gas diffusivity *D*, estimated using Eq. S9. All compounds are assumed to be non-aromatic unless indicated otherwise. For C₅ to C₁₀ VOCs (e.g. isoprene, monoterpenes) and their oxidation products (e.g. C₅ to C₁₀ monomers and C₂₀ dimers), the estimated diffusivities differ less than a factor of 2 from $6.5 \cdot 10^{-6}$ cm² s⁻¹. Diffusion volume correction for (single) aromatic rings results in minor differences (< 5%) of the estimated *D* values. The estimated Knudsen numbers agree within 15%, as do the estimated Fuchs-Sutugin correction factors, β , between the simplified and the more rigorous estimation methods, assuming either a mass accommodation coefficient of 1 (3.08 · 10⁻¹ for all compounds) or 0.1 (3.67 · 10⁻² for all compounds), estimated using Eq. S4 and Eq. S12.

Gas	Kn	%Diff ^a	Diffusivity (m ² s ⁻¹)	$\beta (\alpha = 1)^{b}$	$\beta (\alpha = 0.1)^{b}$
C ₃ H ₆	1.83	9.30	1.18 x 10 ⁻⁵	3.29 x 10 ⁻¹	4.00 x 10 ⁻²
$C_3H_6O_2$	2.05	-2.44	9.97 x 10 ⁻⁶	3.02 x 10 ⁻¹	3.58 x 10 ⁻²
$C_3H_6O_4$	2.20	-9.26	8.96 x 10 ⁻⁶	2.85 x 10 ⁻¹	3.34 x 10 ⁻²
$C_3H_6O_6$	2.32	-13.77	8.27 x 10 ⁻⁶	2.73 x 10 ⁻¹	3.18 x 10 ⁻²
C_5H_8	1.77	13.02	8.98 x 10 ⁻⁶	3.38 x 10 ⁻¹	4.13 x 10 ⁻²
$C_5H_8O_2$	1.94	3.00	8.13 x 10 ⁻⁶	3.15 x 10 ⁻¹	3.78 x 10 ⁻²
$C_5H_8O_4$	2.07	-3.55	7.55 x 10 ⁻⁶	2.99 x 10 ⁻¹	3.54 x 10 ⁻²
$C_5H_8O_6$	2.18	-8.19	7.12 x 10 ⁻⁶	2.88 x 10 ⁻¹	3.38 x 10 ⁻²
$C_5H_8O_8$	2.26	-11.67	6.77 x 10 ⁻⁶	2.79 x 10 ⁻¹	3.25 x 10 ⁻²
$C_7H_8O_1^c$	1.94	2.88	7.83 x 10 ⁻⁶	3.14 x 10 ⁻¹	3.77 x 10 ⁻²
$C_7H_8O_1$	1.83	9.41	7.36 x 10 ⁻⁶	3.30 x 10 ⁻¹	4.00 x 10 ⁻²
$C_7H_8O_2$	1.89	5.55	7.12 x 10 ⁻⁶	3.21 x 10 ⁻¹	3.87 x 10 ⁻²
$C_7H_8O_4$	2.01	-0.45	6.73 x 10 ⁻⁶	3.06 x 10 ⁻¹	3.65 x 10 ⁻²
$C_7H_8O_6$	2.10	-4.91	6.42 x 10 ⁻⁶	2.96 x 10 ⁻¹	3.49 x 10 ⁻²
$C_7H_8O_8$	2.18	-8.37	6.16 x 10 ⁻⁶	2.87 x 10 ⁻¹	3.37 x 10 ⁻²
$C_7H_8O_{10}$	2.25	-11.12	5.93 x 10 ⁻⁶	2.80 x 10 ⁻¹	3.27 x 10 ⁻²
$C_9H_{12}^d$	1.81	10.46	6.92 x 10 ⁻⁶	3.32 x 10 ⁻¹	4.04 x 10 ⁻²
$C_{9}H_{12}$	1.72	16.08	6.58 x 10 ⁻⁶	3.44 x 10 ⁻¹	4.24 x 10 ⁻²
$C_9H_{12}O_2$	1.84	8.65	6.25 x 10 ⁻⁶	3.28 x 10 ⁻¹	3.98 x 10 ⁻²
$C_9H_{12}O_4$	1.94	3.13	5.98 x 10 ⁻⁶	3.15 x 10 ⁻¹	3.78 x 10 ⁻²
$C_9H_{12}O_6$	2.02	-1.13	5.76 x 10 ⁻⁶	3.05 x 10 ⁻¹	3.63 x 10 ⁻²
$C_9H_{12}O_8$	2.10	-4.54	5.57 x 10 ⁻⁶	2.97 x 10 ⁻¹	3.51 x 10 ⁻²
$C_9H_{12}O_{10}$	2.16	-7.33	5.40 x 10 ⁻⁶	2.90 x 10 ⁻¹	3.41 x 10 ⁻²
$C_{10}H_{16}$	1.70	17.36	6.11 x 10 ⁻⁶	3.47 x 10 ⁻¹	4.29 x 10 ⁻²
$C_{10}H_{16}O_2$	1.81	10.41	5.85 x 10 ⁻⁶	3.32 x 10 ⁻¹	4.04 x 10 ⁻²
$C_{10}H_{16}O_4$	1.90	5.13	5.63 x 10 ⁻⁶	3.20 x 10 ⁻¹	3.85 x 10 ⁻²
$C_{10}H_{16}O_{6}$	1.98	0.96	5.44 x 10 ⁻⁶	3.10 x 10 ⁻¹	3.70 x 10 ⁻²
$C_{10}H_{16}O_8$	2.05	-2.42	5.28 x 10 ⁻⁶	3.02 x 10 ⁻¹	3.58 x 10 ⁻²
$C_{10}H_{16}O_{10}$	2.11	-5.21	5.13 x 10 ⁻⁶	2.95 x 10 ⁻¹	3.48 x 10 ⁻²
$C_{10}H_{16}O_{12}$	2.16	-7.57	5.00 x 10 ⁻⁶	2.89 x 10 ⁻¹	3.40 x 10 ⁻²
$C_{10}H_{16}O_{14}$	2.21	-9.58	4.88 x 10 ⁻⁶	2.84 x 10 ⁻¹	3.33 x 10 ⁻²
$C_{10}H_{16}O_{16}$	2.26	-11.32	4.77 x 10 ⁻⁶	2.80 x 10 ⁻¹	3.26 x 10 ⁻²
$C_{20}H_{32}O_{6}$	1.83	9.56	3.98 x 10 ⁻⁶	3.30 x 10 ⁻¹	4.01 x 10 ⁻²
$C_{20}H_{32}O_8$	1.87	6.86	3.92 x 10 ⁻⁶	3.24 x 10 ⁻¹	3.91 x 10 ⁻²
$C_{20}H_{32}O_{10}$	1.91	4.49	3.85 x 10 ⁻⁶	3.18 x 10 ⁻¹	3.83 x 10 ⁻²
$C_{20}H_{32}O_{12}$	1.95	2.38	3.80 x 10 ⁻⁶	3.13 x 10 ⁻¹	3.75 x 10 ⁻²
$C_{20}H_{32}O_{14}$	1.99	0.49	3.74 x 10 ⁻⁶	3.09 x 10 ⁻¹	3.69 x 10 ⁻²
$C_{20}H_{32}O_{16}$	2.02	-1.21	3.69 x 10 ⁻⁶	3.10 x 10 ⁻¹	3.63 x 10 ⁻²

Table S1. Knudsen number and gas diffusivity

⁽a). Percent difference of K_n =2, estimated using Eq. S12, with respect to the K_n estimated using Eq. S5. (b). Fuchs-Sutugin correction factors estimated using Eq. S4 assuming different values for mass accommodation coefficients; the K_n used here was estimated using Eq. S5 (c) *o*-Cresol (d) 1,2,4-trimethylbenzene

Section S2. Oxidation flow reactor schematic

A schematic of the experiment setup is shown in Figure S1 along with the physical dimensions of the oxidation flow reactor (OFR). VOC precursor and seed particles are injected near the entrance region of the OFR, whereas O_3 is injected coaxially in the direction of the flow through a 6 mm outer diameter stainless-steel tubing about 61 cm downstream of the entrance region. Instruments sampled from near the exit region of the OFR. The cross-sectional area of the OFR is approximately $4.3 \cdot 10^{-3}$ m². At 12 L min⁻¹, the plug flow velocity is roughly $4.65 \cdot 10^{-2}$ m s⁻¹. The residence time within the oxidation region (i.e. 39 cm) is roughly 8.38s, or an effective dilution rate of 0.12 s⁻¹.



Figure S1. Flow tube dimension.



Section S3. Vocus-PTR Calibration

Figure S2. Vocus-PTR calibration

(a) The mass transmission efficiency curve for Vocus-PTR is fitted using a lognormal function., Eq. S13. Calibration of the mass transmission efficiency for PTR is described in details elsewhere (Holzinger et al., 2019). (b) Measured sensitivity as a function of the kinetic capture rate, k_{MH} for compounds in a multicomponent calibration tank. The linear regression line with forced 0 intercept was used to estimate sensitivities for additional uncalibrated compounds.

$$MT = 0.20 + 0.96 \times \exp\left(-\frac{\ln\left(\frac{m/z}{95.98}\right)}{0.58^2}\right)$$
 Eq. (S13)

Section S4. AMS Vaporizer artifact correction

The high-resolution aerosol mass spectrometer (AMS) determines the aerosol composition in terms of NO₃, NH₄, SO₄, Chl, and Organics (OA). All experiments were conducted under low-NO_x conditions using NH₄NO₃ seed particles. Therefore, all NH₄⁺ and NO₃⁻ observed are attributed to NH₄NO₃. Due to the high inorganic concentrations used (up to 11.6 mg m⁻³), caution needs to be taken to account for vaporizer artifacts, where NO_x⁺ ions generated from nitrate particles during the electron impact ionization process could oxidize organic residues on the vaporizer surface, producing CO₂⁺ ions that are falsely attributed to organic aerosols (Pieber et al., 2016). The extent of this artifact is determined by injecting NH₄NO₃ seed particles into the OFR in the absence of any organic oxidation products. As shown in Figure S3a below, the correlation of the organic vaporizer artifact, *Orgartifact*, can be described by an exponential function of the NH₄NO₃ (i.e. combined mass concentrations of NO₃⁻ and NH₄⁺). This correlation is used to correct for *Orgartifact* for all runs, as shown in Figure S3d. Note that this correlation could change with the vaporizer history (Pieber et al., 2016). Here, the vaporizer artifact was characterized in the midst of the campaign.





(a) Artefact organics concentration observed by the AMS when sampling nebulized NH_4NO_3 in the absence of any organic oxidation products. An exponential function of NH_4NO_3 concentration is used to estimate the organic signal attributable to the vaporizer artifact. The

organic concentrations with and without applying this correction are shown in (b) for limonene ozonolysis, in (c) for the OH oxidation of *o*-cresol, and in (d) for the OH oxidation of 1,2,4-trimethylbenzene. The correlation between condensed organics and NH₄NO₃ seed concentrations can be roughly described by a double exponential function.

Section S5. Oxidation flow reactor model

The organic vapor wall loss may be estimated from the OFR dimension and the gasdiffusivities as proposed by McMurry and Grosjean (1995),

$$k_{wall} = \frac{1}{\tau_{wall}} = \frac{A}{V} \cdot \frac{2}{\pi} \sqrt{k_e D}$$
 Eq. (S14)

when the vapor wall accommodation coefficient is greater than 10^{-5} , i.e. eddy diffusion dominates. This is the case for oxidation flow reactors (OFR) of similar dimensions to the one used in this study (Brune 2019; George et al., 2007). *A* and *V* are the surface area (1.02 x 10^{-1} m²) and volume (1.72 x 10^{-3} m³) of the OFR, respectively. k_e is the coefficient of Eddy diffusion, which may be estimated as a function of the enclosure volume (Krechmer et al., 2016),

$$k_e(s^{-1}) = 0.004 + (5.6 \times 10^{-3})(V)^{0.74}$$
 Eq. (S15)

which is 4.05 x 10^{-3} s⁻¹. Due to their relatively small enclosure volume (relative to that of a typical smog chamber, k_e would be close to $4 \cdot 10^{-3}$ s⁻¹ for most OFR designs. For estimated gas diffusivity, *D* ranging from 3.69 $\cdot 10^{-6}$ (C₂₀H₃₂O₁₆) to $1.18 \cdot 10^{-5}$ (C₃H₆) m² s⁻¹, the corresponding k_{wall} ranges from 4.60 $\cdot 10^{-3}$ s⁻¹ to 8.22 $\cdot 10^{-3}$ s⁻¹, resulting in a wall loss timescale, τ_{wall} between 122 and 218 s. Two different vapor wall loss experiments conducted using a PTR-TOF and an acetate atmospheric pressure interface chemical ionization TOF-MS indicate a 50% vapor wall loss rate at 10 L min⁻¹ flow rate, which suggest a τ_{wall} similar to that of the dilution lifetime, i.e. 27 seconds, meaning that the actual k_{wall} is close to $3.7 \cdot 10^{-2}$ s⁻¹, roughly 4 to 8 times higher than Eq. S14 and Eq. S15 would suggest. For simplicity, a k_w value of 0.04 s⁻¹ is used as the base case scenario. The effects of higher k_w (i.e. 0.4 s⁻¹) and lower k_w (i.e. 0.04 s⁻¹) values on the gas- and particle-phase concentrations are simulated and shown in Figure 3a-c for generic oxidation products of differing saturation vapor concentrations ranging from 10^{-2} to $10^6 \,\mu g \,m^{-3}$. The OFR wall is also assumed to be a perfect sink for organic vapors, i.e. no back-partitioning of organic vapor from the wall to the gas-phase is considered.

The remaining gas-phase concentration, G_{remain} and the condensed particle-phase concentration, P_{cond} during seed injection are expressed in relative terms with respect to the steady gas-phase concentration prior to the seed injection, G_{ss} (e.g. G_{remain}/G_{ss} and P_{cond}/G_{ss}). So that they are not dependent on the absolute value of G_{ss} , and vice versa on the actual production rate, provided that the production rate is not affected by the seed injection.

The modeled gas-particle partitioning is shown below in Figure S4. A sensitivity analysis was performed by varying the organic aerosol concentration (OA), the condensation sink (CS), or the wall loss rate (k_w) from the base condition (20 µg m⁻³ OA, 1 s⁻¹ CS, and 0.04 s⁻¹ k_w) in Figure S4a-c. The observed OA and CS values were used to simulate the partitioning behaviors as shown in Figure S4d-i. For each VOC system, the observed OA concentration and CS roughly followed a linear correlation. Figure S4d shows the P_{cond} normalized to the maximum value as a function of CS, and suggests that it may be possible to infer the saturation vapor concentration, C^* of semi-volatile compounds based on the uptake trend without the knowledge of near-molecular particle-phase sensitivity or gas-phase concentration (as long as G_{SS} remains constant in this case). However, compounds of different C^* may exhibit similar

trends, i.e. high inter-correlations, which cannot be numerically resolved due to noise. Visually, this is obvious for compounds with $\log(C^*) > 2$ or < -1 as shown in Figure S4d.

To determine the range of $log(C^*)$ that could be in theory numerically resolved from the P_{cond} behaviors alone, we modeled the normalized P_{cond} for compounds with $\log(C^*)$ ranging from -2 to 6 using OA and CS values observed for each system. The lower C^* threshold is set at the point beyond which all compounds with lower C^* would exhibit normalized P_{cond} trends with intercorrelation (R^2 value from linear regression between the normalized P_{cond} values corresponding to any pair of C* values, i.e. any two "vertical slices" from Figure S4e and S4f) above 0.99. The decision to set the cutoff at $R^2 = 0.99$ is arbitrary. The upper C^{*} threshold is similarly defined in Figure S4g-i. The experimentally constrainable $log(C^*)$ ranges based on the uptake behavior alone are narrow: 1.25 to 2.02 for the cresol system, 1.18 to 2.09 for the TMB system, and 0.57 to 1.85 for the limonene system. The span of the constrainable C^* range is wider for the limonene system due to the higher maximum CS range explored experimentally $(>2 \text{ s}^{-1} \text{ as compared to } <1 \text{ s}^{-1} \text{ for the anthropogenic systems})$. The upper constrainable C^* range for limonene system (i.e. $\log(C^*) = 1.85$) is lower compared to that for either cresol (i.e. $\log(C^*)$ =2.02) or TMB system (i.e. $log(C^*)$ =2.09) due to the lower maximum OA uptake as a function of CS for the limonene system as compared to the anthropogenic systems. All else being equal, the constrainable range of $log(C^*)$ increases with the experimental CS range, which is limited by the maximum particle concentrations the instruments could accommodate before clogging or signal depletion becomes too severe.







(a-c) Expected distribution of organic oxidation products of differing volatilities between the gas- and particle-phase during the seed injection period for a hypothetical base case scenario of 20 µg m⁻³ organic aerosol concentration (*OA*), 1 s⁻¹ condensation sink (*CS*), and 0.04 s⁻¹ vapor wall loss rate (k_w). Alternative scenarios assume higher or lower *OA*, *CS*, and k_w . (d-i) Modeled ratio of P_{cond} to G_{ss} for compounds of varying $log(C^*)$ under observed OA and CS conditions. (a) Ratio of condensed organic material during seed injection, P_{cond} to the steady-state gas-phase concentration prior to seed injection, G_{SS} . The ratio can exceed 1 under high *CS* conditions. (b) Ratio of P_{cond} to the sum of P_{cond} with the gas-phase concentration during the seed injection period, G_{remain} . Partitioning between P_{cond} and G_{remain} is invariant with respect

to k_w . (c) Ratio of P_{cond} to the sum of P_{cond} and G_{SS} . (d) Normalized P_{cond} relative to the maximum expected value, $P_{cond,max}$ as a function of CS for compounds of different volatility. Note again that observed CS and OA values from the anthropogenic experiments are used to simulate the uptake behavior shown in (d), whereas hypothetical CS, OA, and k_w conditions are used to simulate the behaviors shown in (a-c). (e) Ratio of P_{cond} to G_{SS} for compounds of varying $\log(C^*)$ at different CS for the cresol and TMB systems, which exhibited similar intercorrelations between observed OA and CS. (f) Ratio of P_{cond} to G_{SS} for compounds of varying $\log(C^*)$ at different CS for the limonene system. (g) Inter-correlation of the expected normalized P_{cond} , similar to those shown in (d), for compounds of varying $\log(C^*)$ under the uptake conditions in the cresol system. (h) Inter-correlation of the expected normalized P_{cond} , for compounds of varying $\log(C^*)$ under the uptake conditions in the TMB system. (i) Intercorrelation of the expected normalized P_{cond} , for compounds of varying $\log(C^*)$ under the uptake conditions in the limonene system. (j) Similar to g, but with the maximum CS range extrapolated to 2 s^{-1} from ~0.8 s⁻¹ to examine its effect on constrainable log(C^{*}) range. Regions with R^2 values exceeding 0.99 are shown in white in (g-j), where the log(C^*) empirically determined from the normalized P_{cond} is considered as highly uncertain due to experimental noise and high intercorrelations of the normalized P_{cond} behavior with compounds of different $\log(C^*)$. Behaviors of compounds with $\log(C^*)$ below -1 or above 4 are not shown, as they are indistinguishable per our definition based on the intercorrelation value R^2 .





Figure S5. Average particle-phase composition

Ion intensity of $[M+Na]^+$ adducts observed during (a) OH-oxidation of cresol, (b) OH-oxidation of TMB, and (c) ozonolysis of limonene. For each VOC and oxidant system, the average composition over all seed injection / organic aerosol uptake events is shown. Ion intensities are grouped by their carbon number (#*C*) and further distinguished by the oxygen number as shown in the legend.



Intensities of selected $[M+Na]^+$ adducts observed by the EESI-TOF for the particle-phase are shown for (a) C₇ OH + cresol oxidation products, (c) C₉ OH + TMB oxidation products, and (e) C₁₀ limonene + O₃ oxidation products. Intensities of selected $[M+H]^+$ ions observed by the Vocus-PTR in the gas-phase are shown for (b) C₇ OH + cresol oxidation products, (d) C₉ OH + TMB oxidation products, and (f) C₁₀ limonene + O₃ oxidation products. Average particle-phase signals over all uptake events are shown in (a), (c), and (e). Average steadystate gas-phase concentrations prior to each uptake event are shown in (b), (d), and (f). Note that the color scales are only consistent within each of the (a-b), (c-d), and (e-f) pairs. Ion

intensities are grouped by the number of hydrogens (#H) and further distinguished by the number of oxygen as indicated in the legends. The *o*-cresol is not included in (a) and (b) because it is the VOC precursor and not an oxidation product.

Section S7. Parameterization and Model Validation

The EESI-TOF response factor in ions s⁻¹ ppb⁻¹, RF_x^* , can be estimated by performing a linear regression of I_x (in ions s⁻¹) as a function of $P_{cond,x}$ (in ppb⁻¹) as described in the main text, and taking the slope. Because ordinary least square regression (OLS) minimizes only the vertical (i.e. I_x on the y-axis) distance of the dependent variable, it cannot account for uncertainties in the explanatory variable (i.e. $P_{cond,x}$ on the x-axis) during error propagation. Propagation of uncertainties in the explanatory and dependent variables can be achieved by performing an orthogonal distance regression (ODR). The slope values obtained using either method agree within a factor of 2, as shown in Figure S7.



Figure S7. Comparison of the response factor (RF_x^*) values determined using ordinary least square (OLS) and orthogonal distance regression (ODR). Uncertainties in the explanatory and response variables are taken into consideration by ODR during fitting. Vertical and horizontal error bars shown represent the standard deviation of the fitted slope of EESI-TOF vs. Vocus-PTR measurements, i.e. RF_x^*

Based on the elemental formulae measured by the EESI-TOF and the Vocus-PTR, several additional features could be derived from the number of carbon (n_C), hydrogen (n_H), and oxygen (n_O), including the exact molecular mass (*MW*), the mass defect (Δm), the hydrogen-to-carbon ratio (*H*:*C*), the oxygen-to-carbon ratio (*O*:*C*), the double bond equivalent (*DBE*), and the double bond equivalent per carbon (*DBEpC*)

$$DBE = 1 + \frac{1}{2}(2C - H + N + P)$$
 Eq. (S16)

The aromaticity index (AI) can be calculated as

$$AI = \frac{DBE_{AI}}{C_{AI}} = \frac{1 + C - O - S - 0.5H}{C - O - S - N - P}$$
Eq. (S17)

Which has been reported to underestimate the aromaticity compared to the aromaticity equivalent (X_c) proposed by Yassine et al. (Yassine et al., 2014)

$$X_C = \frac{C - (H - C)}{DBE} + 1$$
 Eq. (S18)

where, if $DBE \le 0$, X_c is set to 0. Note that for CHO compounds, Eq. S18 simplifies to

$$X_C = 3 - \frac{2}{DBE}$$
 Eq. (S19)

In addition, the carbon-oxygen non-ideality (NI_{CO}) from Eq. (7) itself is an interaction term between the product of the number of carbon and oxygen atoms (P_{CO}) and the inverse of the sum of carbon and oxygen atoms (I_{CO}),

$$NI_{CO} = \frac{n_C n_O}{n_C + n_O} = P_{CO} \times I_{CO}$$
 Eq. (S20)

In addition to the aforementioned features, the log of effective saturation vapor concentration, $log(C^*)$ is included as a feature.

Preliminary ordinary least square (OLS) regressions of the near-molecular EESI-TOF response factor, RF_x^* (which was obtained with ODR) as a function of n_C , n_O , MW, NI_{CO} , P_{CO} , or I_{CO} , are shown in Figure S8a-f for each of the three VOC systems studied. The RF_x^* values estimated for cresol and TMB oxidation products appear to increase as the molecules increase in size (i.e. positive correlation with MW and n_C) and/or become more functionalized (i.e. positive correlation with n_O). The correlations also appear to be steeper for the TMB system than for the cresol system. In contrast, the RF_x^* values estimated for limonene oxidation products do not appear to be well correlated with n_C , n_O , MW, P_{CO} , I_{CO} , or NI_{CO} . The discrepancies observed between the aromatic systems and the biogenic system are likely due to differences in the structure of the oxidation products as discussed in the main text.





OLS regression analysis of the log of RF_x^* with respect to (a) the number of carbon, n_C , (b) the number of oxygen, n_O , (c) the molecular weight, MW, (d) the carbon-oxygen non-ideality, NI_{CO} , (e) the product of n_C and n_O , P_{CO} , and (f) the inverse of the sum of n_C and n_O , I_{CO} . The red, blue, and green dashed lines correspond to the linear fitting lines for the log(RF_x^*) values of TMB, cresol, and LMN oxidation products, respectively. The coefficient of determination, R^2 of ordinary linear regression for the log(RF_x^*) as a function of the feature is shown in brackets after the corresponding VOC label.

The full regression analysis was performed on two types of datasets: The log of measured EESI sensitivity in ions s⁻¹ ppb⁻¹, $log(RF_x^*)$ from (1) the TMB system alone, or (2) all three VOC systems. Two approaches were taken for the combined dataset: (2a) the precursor

VOC identity was not included as a feature or (2b) the VOC identity was one-hot encoded and included as a feature, i.e. limonene, TMB, and cresol products would have attributes of [VOC_{LMN}:1, VOC_{TMB}: 0, VOC_{Cresol}: 0], [VOC_{LMN}:0, VOC_{TMB}: 1, VOC_{Cresol}: 0], and [VOC_{LMN}:0, VOC_{TMB}: 0, VOC_{Cresol}: 1], respectively.

First, an exhaustive search over the feature space was performed to determine the optimal set of features for each regressor using their respective default hyperparameter values. Leave-one-out (LOO) cross-validation was used to evaluate the model performance in terms of the coefficient of determination, R^2

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$
 Eq. (S21)

where y_i and \hat{y}_i are the true and the predicted value for the *i*-th sample among a total of n samples, and \overline{y} is the mean of the *n* samples. If the model always predicts \overline{y} , the R^2 will be 0, e.g. a naive model where all values are predicted to equal that of the sample mean regardless of input. The R^2 can be negative if it performs worse than this naive model, i.e. assuming the mean value regardless of model input produces on average better results. For a dataset of size *n*, LOO involves setting aside each data point (y_i) in turn as the test sample while the remaining (*n*-1) data points are used to train the model and make a prediction, \hat{y}_i . y_i and \hat{y}_i are then used to estimate the R^2 using Eq. (S21). LOO can be considered as performing a K-fold crossvalidation where the number of K is equal to the number of data points. Compared to the Kfold cross-validation method, LOO is more computationally intensive to perform, but is nonetheless appropriate given the small size of the dataset used here $(n_{sample} = 28 \text{ for case } 1 \text{ and}$ 70 for 2a and 2b). During cross-validation, a portion of the dataset is used to train the model (i.e. "train" set), while the remaining dataset is withheld to validate against the model predictions (i.e. "test" set). For each train-test set, the training feature values $(n = n_{sample} - 1)$ were standardized, which involves subtracting by their mean and dividing by their standard deviation. The same transformation was then applied to the feature values from the test set (n= 1), which was not included in deriving the transformation required for the standardization to prevent information leak between the training and test sets.

The results of the feature optimization are shown in Figure S8 in terms of the best R^2 vs. the number of features used. The optimal feature sets are shown in Table S2. In addition to OLS, linear ridge regression ("Ridge") and Bayesian ridge regression (BayRR) are included. Both Ridge and BayRR implement L₂ regularization, making them more resilient against overfitting and feature co-linearity. Support vector regression (SVR) with linear kernel is also included as a linear regression model for comparison. Exploratory analysis using SVR with radial basis functions (rbf) yielded better R^2 , but the relative feature importance was not easily interpretable when rbf was used, hence the choice of linear kernel. Lastly, nonparametric regressions such as random forest regressor (RFR) and gradient boosting regressor (GBR) were included, as the RF_x^* is likely not a linear function of features already included. While it is possible that RF_x^* could be well-described by a linear combination of engineered features, it is not feasible to explore all nonlinear (e.g. n_c^2) or interaction ($n_c n_H$) feature terms, hence the necessity of nonparametric regressors. For the purpose of feature selection and later hyperparameter tuning, the random state (which controls the permutation of features at each split within the decision tree) for RFR and GBR are fixed (i.e. given a seed value of 0) so that the models generate producible outputs.



Figure S9. Feature selection

The best R^2 from LOO cross-validation test for each regressor using different permutations of features as a function of the number of features included for (a) Case 1, where only the TMB dataset is used, (b) Case 2a, where data from all the VOC systems were used without providing the digitized VOC identity as one of the input features, or (c) Case 2b, where data from all the VOC systems were used with the one-hot encoded VOC identity provided as one of the input features, hence the one extra feature over cases 1 and 2a.

Case #	OLS	Ridge	BKK	SVK	KFK	GBK
1	0.46	0.40	0.39	0.46	0.59	0.71
	no,	NIco	NIco	NIco,	no,	no,
	Ico,			mW	n _H ,	n _H ,
	$\log(C^*)$			DBEpC, X _C	H:C	H:C
2a	0.33	0.32	0.32	0.35	0.19	0.21
	Xc,	Xc,	Xc,	Xc,	no	nc,
	H:C,	H:C,	H:C,	H:C,	H:C,	X _C ,
	$\log(C^*),$	$\log(C^*),$	$\log(C^*),$	$\log(C^*),$	NIco	NIco
	NI _{CO}	NI _{CO}	NI _{CO}	NI _{CO}		
2b	0.51	0.47	0.47	0.48	0.40	0.49
	H:C,	NI _{CO} ,	NI _{CO} ,	NI _{CO} ,	H:C,	H:C,
	no,	Xc,	DBE,	DBE,	no,	O:C,
	Δm,	$\log(C^*),$	X _C ,	Xc,	Δm,	P _{CO} ,
	Xc	DBE,	$\log(C^*),$	$\log(C^*),$	VOC	VOC
	VOC	VOC	VOC	VOC		

Table S2. Best R² scores and their corresponding feature combination obtained using leave-one-out cross-validation with default model hyperparameters

Note that in some cases (e.g. SVR in case 1 and 2a), the optimal feature set selected does not correspond to the set with the highest R^2 , but rather one with slightly lower R^2 score but also (sometimes substantially) lower total number of features used. The feature abbreviations used are as followed: Carbon-oxygen non-ideality (*NI*_{CO}), product of the number of oxygen and carbon numbers (*P*_{CO}), the inverse of the sum of the number of oxygen and carbon numbers (*I*_{CO}), logarithm of saturation vapor concentration (log(C^*)), aromaticity (*X*_C), double bond equivalent per carbon (*DBEpC*), number of oxygen atoms (*n*_O), number of hydrogen atoms (*n*_H), molecular weight (*MW*), mass defect (Δm), hydrogen-to-carbon ratio (*H*:*C*), oxygen to carbon ratio (*O*:*C*), one-hot encoded precursor VOC label (VOC).

For Case 1, NI_{CO} and n_O are identified as essential features in predicting the (log of) EESI-TOF response factor. For Case 2a, all model performances degrade due to the lack of knowledge of the VOC identity. The feature selection results that NI_{CO} (or relatedly n_O) and X_C (or relatedly H:C) are essential features to include. For Case 2b, inclusion of the VOC label as a feature results in substantial increase in R^2 for all regressors, as shown in Figure S9 and Table S3 below. As the number of features increase beyond 5, regressor performances do not show any substantial improvement and may even deteriorate. As a trade-off between R² and model feature complexity, the optimal number of features for linear models (OLS, Ridge, BRR, and SVR) is set to 5 and the optimal number of features for nonparametric models (RFR and GBR) is set to 4.

			VI I			
Feature #	OLS	Ridge	BRR	SVR	RFR	GBR
1	0.12	0.12	0.12	0.14	0.09	0.06
	I _{CO}	Ico	Ico	mW	H:C	H:C
2	0.20	0.23	0.22	0.22	0.21	0.31
	NI _{CO} ,	NI _{co} ,	NI _{CO} ,	NI _{CO} ,	mW,	mW,
	$\log(C^*)$	VOC	VOC	$\log(C^*)$	VOC	VOC
3	0.29	0.33	0.33	0.35	0.38	0.39
	NI _{CO} ,	NI _{CO} ,	NI _{CO} ,	NI _{CO} ,	H:C,	H:C,
	H:C,	Χс,	Χс,	X _C ,	no,	mW,
	Ico	VOC	VOC	VOC	VOC	VOC
4	0.40	0.40	0.38	0.39	0.40	0.49
	I _{CO} ,	NI _{CO} ,	NI _{CO} ,	NI _{CO} ,	H:C,	H:C,
	Χс,	X _C ,	X _C ,	Χс,	n _O ,	O:C,
	DBE,	DBE,	H:C,	$\log(C^*),$	Δm,	P _{CO} ,
	VOC	VOC	VOC	VOC	VOC	VOC
5	0.51	0.47	0.47	0.48	0.41	0.50
	H:C,	NI _{CO} ,	NI _{CO} ,	NI _{CO} ,	H:C,	H:C,
	n _O ,	X _C ,	X _C ,	Χс,	n _O ,	O:C,
	Δm ,	$\log(C^*),$	$\log(C^*),$	$\log(C^*),$	Δm,	P _{CO} ,
	X_{C}	DBE,	DBE,	DBE,	X_{C}	I _{CO} ,
	VOC	VOC	VOC	VOC	VOC	VOC

Table S3. Best R² scores for different feature combinations obtained using leave-one-out cross-validation with default model hyperparameters for Case 2b

Having identified the optimal feature sets, we then performed grid search to find the optimal model hyperparameters using R^2 from *LOO* as the metric. The hyperparameter spaces explored for each regressor are listed in Table S4a-c below, along with the *LOO* R^2 obtained using the default vs. the optimal model hyperparameters.

Regressor	Hyperparameter Space	Optimal	R ² (Optimal)	R ² (Default)
RFR	n_estimator: [10, 20, 30, 40, 50, 100]	20	0.71	0.61
	min_samples_split: [2, 3, 4, 5]	3		
	max_features: [<u>"auto"</u> , "sqrt", "log2"]	"sqrt"		
	bootstrap: [<u>True</u> , False]	False		
SVR	C: [0.1, 0.2, 0.5, <u>1</u> , 2, 10, 100]	5	0.49	0.46
	epsilon: [<u>0.1</u> , 0.2, 0.5, 1, 10, 100]	<u>0.1</u>		
GBR	n_estimator: [5, 10, 20, 30, 40, 50, 100, 200]	200	0.83	0.71
	loss: [<u>"ls"</u> , "lad", "huber"]	"lad"		
	learning_rate: [0.05, <u>0.1</u> , 0.2, 0.5, 0.7]	0.7		
	subsample: [0.3, 0.5, 0.7, <u>1]</u>	0.7		
	max_features: [<u>"auto"</u> , "sqrt", "log2"]	<u>"auto"</u>		
	min_samples_split: [2, 3, 4, 5]	4		
BRR	n_iter: [100, 200, <u>300</u> , 500, 1000]	100	0.39	0.39
	alpha_1: [10 ⁻⁴ , 10 ⁻⁵ , <u>10⁻⁶</u> , 10 ⁻⁷ , 10 ⁻⁸]	10-4		
	alpha_2: [10 ⁻⁴ , 10 ⁻⁵ , <u>10⁻⁶</u> , 10 ⁻⁷ , 10 ⁻⁸]	10-8		
	lambda_1: [10 ⁻⁴ , 10 ⁻⁵ , <u>10⁻⁶</u> , 10 ⁻⁷ , 10 ⁻⁸]	10-8		
	lambda_2: [10 ⁻⁴ , 10 ⁻⁵ , <u>10⁻⁶</u> , 10 ⁻⁷ , 10 ⁻⁸]	10-4		
Ridge	alpha: [0.1, 0.2, 0.5, <u>1</u> , 2, 10, 100]	<u>1</u>	0.40	0.40

Table S4a. Regressor	hyperparameter	grid search	results for Case 1
----------------------	----------------	-------------	--------------------

Note: Optimal hyperparameter values that are identical to the default values are <u>underlined</u>. The seed values used to generate the random state for decision tree-type ensemble models, i.e. RFR and GBR, are fixed during hyperparameter grid search to ensure that the models give reproducible outputs for a given set of inputs and model hyperparameters.

Table	S4b.	Regressor	hyper	parameter	grid	search	results	for	Case	2a
					8					

Regressor	Hyperparameter Space	Optimal	R ² (Optimal)	R ² (Default)
RFR	n_estimator: [10, 20, 30, 40, 50,100]	10	0.22	0.19
	min_samples_split: [2, 3, 4, 5]	4		
	<pre>max_features: ["auto", "sqrt", "log2"]</pre>	<u>"auto"</u>		
SVR	C: [0.1, 0.2, 0.5, 1, 2, 10, 100]	2	0.37	0.35
	epsilon: [0.1, 0.2, 0.5, 1, 10, 100]	<u>0.1</u>		
GBR	n_estimator: [5, 10, 20, 30, 40, 50, 100, 200]	10	0.34	0.21
	loss: ["ls", "lad", "huber"]	"lad"		
	learning_rate: [0.05, 0.1, 0.2, 0.5]	0.5		
	subsample: [0.3, 0.5, 0.7, 1]	<u>1</u>		
	<pre>max_features: ["auto", "sqrt", "log2"]</pre>	"log2"		
	min_samples_split: [2, 3, 4, 5]	<u>2</u>		
BRR	n_iter: [100, 200, 500, 1000]	100	0.32	0.32
	alpha_1: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-4		
	alpha_2: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-8		
	lambda_1: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-8		
	lambda_2: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-4		
Ridge	alpha: [0.1, 0.2, 0.5, 1, 2, 10, 100]	0.2	0.33	0.32

Regressor	Hyperparameter Space	Optimal	R ² (Optimal)	R ² (Default)
RFR	n_estimator: [10, 20, 30, 40, 50,100]	50	0.42	0.40
	min_samples_split: [2, 3, 4, 5]	3		
	<pre>max_features: ["auto", "sqrt", "log2"]</pre>	"sqrt"		
	bootstrap: [True, False]	False		
SVR	C: [0.1, 0.2, 0.5, 1, 2, 10, 100]	100	0.52	0.48
	epsilon: [0.1, 0.2, 0.5, 1, 10, 100]	0.1		
GBR	n_estimator: [5, 10, 20, 30, 40, 50, 100, 200,	350	0.52	0.49
	250, 300, 350, 400]			
	loss: ["ls", "lad", "huber"]	<u>ls</u>		
	learning_rate: [0.05, 0.1, 0.2, 0.5, 0.7]	0.05		
	subsample: [0.3, 0.5, 0.7, 1]	<u>1</u>		
	max_features: ["auto", "sqrt", "log2"]	<u>"auto"</u>		
	min_samples_split: [2, 3, 4, 5]	4		
BRR	n_iter: [100, 200, 500, 1000]	100	0.47	0.47
	alpha_1: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-4		
	alpha_2: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-8		
	lambda_1: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-8		
	lambda_2: [10 ⁻⁴ , 10 ⁻⁵ , 10 ⁻⁶ , 10 ⁻⁷ , 10 ⁻⁸]	10-4		
Ridge	alpha: [0.1, 0.2, 0.5, 1, 2, 10, 100]	0.1	0.50	0.47

Table S4c. Regressor hyperparameter grid search results for Case 2b

The log(RF_x^*) predicted from the LOO cross-validation test (see discussion around Eq. S21) by the linear ridge regressor (LRR) and the gradient boosting regressor (GBR) using their respective optimal features sets and hyperparameters for Cases 1, 2a, and 2b are shown in Figure S10 and compared to the measured log(RF_x^*). Performance and fitting coefficients of all regressors used are summarized in Table S5. For a single VOC system (Case 1), the predicted and measured RF_x^* values mostly agree within a factor of 5 using LRR or a factor of 2 using GBR. When dealing with compounds from multiple VOC systems, where the VOC precursor identities are unknown (i.e. Case 2a), the predictions fare slightly better than simply assuming a uniform response factor equal to that of the sample mean, as shown in Figure S10b. If the VOC precursor identity is used as one of the features, GBR and LRR can produce reasonable predictions that agree with the measured values within a factor of 2-5, as shown in Figure S10c. Much of the scatter was related to the limonene dataset, which did not appear to have a clear predictor for log(RF_x^*), as we have also shown during our preliminary regression analysis in Figure S8.





Comparison of the log of the measured response factor, $log(RF_x)$ with those predicted using the leave-one-out cross-validation method by the linear ridge regressor (LRR) and the gradient boosting regressor (GBR) using their optimal feature sets and hyperparameters for (a) Case 1, (b) Case 2a, and (c) Case 2b. The VOC identity was made available to the regression models to use as a potential feature for Case 2b in (c), but not for Case 2a in (b). The 1-to-1 line is shown in solid black. The darker shaded region represents a factor of 2 deviation from the 1-to-1 line.

			1			
Case #	OLS	Ridge	BRR	SVR	RFR	GBR
1	0.46	0.40	0.39	0.49	0.71	0.83
	no: 2.53	NI _{co} : 0.42	NI _{CO} : 0.43	NI _{co} : 1.31	no: 0.52	n ₀ : 0.30
	Ico: -0.95			mW: -1.12	n _H : 0.23	n _H : 0.25
	$\log(C^*): 2.96$			DBEpC: -0.79	H:C: 0.25	H:C: 0.44
	-			Xc: 0.72		
2a	0.33	0.33	0.32	0.37	0.22	0.34
	X _C : 0.63	X _C : 0.62	X _C : 0.58	X _C : 0.64	n ₀ : 0.11	X _C : 0.38
	H:C: 0.57	H:C: 0.56	H:C: 0.51	H:C: 0.58	H:C: 0.54	nc: 0.17
	NIco: 0.61	NIco: 0.55	NIco: 0.39	NIco: 0.43	NIco: 0.35	NIco: 0.45
	$\log(C^*): 0.84$	$\log(C^*): 0.77$	$log(C^*): 0.61$	$\log(C^*): 0.64$		
2b	0.49	0.50	0.47	0.52	0.42	0.52
	NIco: 1.92	NIco: 1.13	NIco: 0.88	NIco: 1.23	H:C: 0.32	H:C: 0.33
	n _H : 0.71	DBE: -0.50	DBE: -0.43	DBE: -0.59	no: 0.22	O:C: 0.17
	Xc: 0.57	Xc: 0.60	Xc: 0.56	X _C : 0.68	Δm: 0.28	Pco: 0.28
	mW: -1.93	$\log(C^*): 0.81$	$\log(C^*): 0.57$	$\log(C^*): 0.91$	VOCLMN: 0.08	VOCLMN: 0.15
	VOCLMN: 0.18	VOCLMN: 0.16	VOCLMN: 0.16	VOCLMN: 0.16	VOC _{TMB} : 0.04	VOC _{TMB} : 0.02
	VOC _{TMB} : -0.02	VOC _{TMB} : -0.01	VOC _{TMB} : -0.01	VOC _{TMB} : -0.03	VOC _{Cresol} : 0.06	VOC _{Cresol} : 0.05
	VOC _{Cresol} : -0.18	VOC _{Cresol} : -0.17	VOC _{Cresol} : -0.18	VOC _{Cresol} : -0.14		

Table S5. R^2 for each regressor using their optimal features and model hyperparameters, and the weights/importance of fitted features.

The R^2 determined from the leave-one-out (LOO) cross-validation test is shown. For ordinary least square (OLS) regression, linear ridge regression (LRR), Bayesian ridge regression (BRR), and support vector regression (SVR), the weight for each feature is shown. For random forest (RFF) and gradient boosting regression (GBR), the importance is shown, which is a measure of the usefulness of a feature in constructing the decision tree. VOC_{LMN}, VOC_{TMB}, and VOC_{Cresol} are the one-hot encoded representation of the VOC identity.

Note that if we were to use the entire dataset to train and validate the model, the resulting R^2 would be overly optimistic, as shown in Figure S11 especially for those obtained using the nonparametric regressors.





Comparison of the predicted $\log(RF_x^*)$ using the entire dataset with VOC label included as one of the features using (a) linear regression models and (b) nonparametric regression models. The optimal feature sets and hyperparameters used for each model are identical to those used for Figure S10 and Table S5, except that now each model was trained with the entire dataset to predict the entire dataset, instead of following the LOO procedure. The 1-to-1 line is shown in solid black. The darker shaded region represents a factor of 2 deviation from the 1-to-1 line. The lighter shaded region represents a factor of 5 deviation from the 1-to-1 line.

For typical ambient measurements or chamber experiments with complex precursor mixtures, the VOC precursor identity is often not known without additional constraints (e.g. ion mobility or gas chromatography measurements supported with chemical reaction box models). The prediction capability of the regression model for an unknown VOC is examined in Figures S12a and S12b, using the TMB dataset as the "known" VOC system to predict the $log(RF^*_x)$ for the "unknown" cresol and limonene (LMN) systems. As shown in Figure S12a, while the regression models trained with TMB dataset tend to overestimate the $log(RF^*_x)$ for the cresol system, the predictions and observations are qualitatively consistent in terms of the relative $log(RF^*_x)$, likely due to the structural similarity of cresol and TMB, which would be reflected to varying degrees in their respective oxidation products. In contrast, regression models trained with the TMB dataset are unfit to predict the $log(RF^*_x)$ for the limonene oxidation products, as shown in Figure S12b.

The effect of the VOC precursor on the predicted $\log(RF_x^*)$ values, using the model trained in Case 2b (all data with digitized VOC label), for all CHO molecular formulae used for EESI-TOF spectral fitting is shown in Figures S12c and S12d. In general, the predicted $\log(RF_x^*)$ trend in the same direction for all VOCs. The predicted effect of VOC precursor is distinct when a linear regressor is used, as shown in Figure S12d, where $\log(RF_x^*)$ is treated as a linear combination of features, one of which is the digitized VOC precursor identity. When a decision-tree type regressor is used, the VOC precursor identity effect is not as simple, as shown in Figure S12c. Lastly, the combination of dataset from multiple VOC systems also affects the predicted $\log(RF_x^*)$, as shown in Figure S12e and S12f for the TMB system. Linear models trained with the combined dataset (i.e. Case 2b) appear to (severely) underestimate the $\log(RF_x^*)$ as compared to the models trained with a single VOC dataset (i.e. Case 1). Furthermore, regressors that performed reasonably well (e.g. LRR for Case 1) for the training dataset with a limited number of features (e.g. NI_{CO}) may be ill-equipped when predicting for a more diverse set of compounds, whose variabilities are only reflected in other features (e.g. optimal features for LRR in Case 2b, see Table S5).





(a) Comparison of the observed $\log(RF_x^*)$ for cresol oxidation products with that predicted using gradient boosting regression (GBR) and the linear ridge regression (LRR) models trained with the TMB dataset (b) Same as (a) but for the limonene (LMN) system. (c) Comparison of the $\log(RF_x^*)$ for all molecular formulae used for EESI-TOF MS fitting predicted using the GBR model from Case 2b for different VOC systems, i.e. all feature values used during

prediction were identical expect for that of the digitized VOC precursor identity. (d) Same as (c), but with the LRR model from Case 2b. (e) Comparison of the $\log(RF_x^*)$ for all molecular formulae used for EESI-TOF MS fitting predicted using the GBR model from Case 1 and Case 2b for TMB system only. (f) Same as (e), but with the LRR model from Case 1 and Case 2b. The optimal feature sets and hyperparameters used for each model are listed in Table S5. The 1-to-1 line is shown in solid black. The darker shaded region represents a factor of 2 deviation from the 1-to-1 line. The lighter shaded region represents a factor of 5 deviation from the 1-to-1 line. The case number indicated on the axis legend and in annotations indicate the how the model was trained as described throughout Tables S2-5.



Figure S13. Comparison of estimated and observed OA concentration

Comparison of the observed organic aerosol (OA) as measured by the AMS with the OA concentration estimated using EESI-TOF measurements converted from ions s⁻¹ to μ g m⁻³ using the RF_x^* (ions s⁻¹ ppb⁻¹) predicted using the gradient boosting regression (GBR) model. Conversion of ppb to molecules cm⁻³ is performed under standard conditions, i.e. $2.46 \cdot 10^{10}$ molecules cm⁻³ per ppb. The 1-to-1, 2-to-1, and 3-to-1 lines are shown in solid black. Two versions of the regression models are used to predicted the RF_x^* for TMB, one trained with single VOC dataset (Case 1) and one trained with combined VOC datasets where the VOC precursor identity is used as a training feature (Case 2b).

References

Brune, W. H.: The Chamber Wall Index for Gas-Wall Interactions in Atmospheric Environmental Enclosures, Environ. Sci. Technol., 53(7), 3645–3652, doi:10.1021/acs.est.8b06260, 2019.

Dal Maso, M., Kulmala, M., Lehtinen, K. E. J., *Mkelä*, J. M., Aalto, P. and O'Dowd, C. D.: Condensation and coagulation sinks and formation of nucleation mode particles in coastal and boreal forest boundary layers, J. Geophys. Res. Atmos., 107(19), doi:10.1029/2001JD001053, 2002.

Fuller, E. N., Schettler, P. D. and Giddings, J. C.: A new method for prediction of binary gasphase diffusion coefficients, Ind. Eng. Chem., 58(5), 18–27, doi:10.1021/ie50677a007, 1966.

George, I. J., Vlasenko, A., Slowik, J. G., Broekhuizen, K. and Abbatt, J. P. D.: Heterogeneous oxidation of saturated organic aerosols by hydroxyl radicals: Uptake kinetics, condensed-phase products, and particle size change, Atmos. Chem. Phys., 7(16), 4187–4201, doi:10.5194/acp-7-4187-2007, 2007.

Holzinger, R., Joe Acton, W. F., Bloss, W. W., Breitenlechner, M., Crilley, L. L., Dusanter, S., Gonin, M., Gros, V., Keutsch, F. F., Kiendler-Scharr, A., Kramer, L. L., Krechmer, J. J., Languille, B., Locoge, N., Lopez-Hilfiker, F., Materi, D., Moreno, S., Nemitz, E., Quéléver, L. L., Sarda Esteve, R., Sauvage, S., Schallhart, S., Sommariva, R., Tillmann, R., Wedel, S., Worton, D. D., Xu, K. and Zaytsev, A.: Validity and limitations of simple reaction kinetics to calculate concentrations of organic compounds from ion counts in PTR-MS, Atmos. Meas. Tech., 12(11), 6193–6208, doi:10.5194/amt-12-6193-2019, 2019.

Jennings, S. G.: The mean free path in air, J. Aerosol Sci., 19(2), 159–166, doi:10.1016/0021-8502(88)90219-4, 1988.

Krechmer, J. E., Pagonis, D., Ziemann, P. J. and Jimenez, J. L.: Quantification of Gas-Wall Partitioning in Teflon Environmental Chambers Using Rapid Bursts of Low-Volatility Oxidized Species Generated in Situ, Environ. Sci. Technol., 50(11), 5757–5765, doi:10.1021/acs.est.6b00606, 2016.

Krechmer, J. E., Day, D. A., Ziemann, P. J. and Jimenez, J. L.: Direct Measurements of Gas/Particle Partitioning and Mass Accommodation Coefficients in Environmental Chambers, Environ. Sci. Technol., 51(20), 11867–11875, doi:10.1021/acs.est.7b02144, 2017.

Kulmala, M. and Wagner, P. E.: Mass accommodation and uptake coefficients - A quantitative comparison, J. Aerosol Sci., 32(7), 833–841, doi:10.1016/S0021-8502(00)00116-6, 2001.

Lehtinen, K. E. J., Korhonen, H., Dal Maso, M. and Kulmala, M.: On the concept of condensation sink diameter, Boreal Environ. Res., 8(4), 405–411 [online] Available from: http://www.borenv.net/BER/pdfs/ber8/ber8-405.pdf (Accessed 22 May 2014), 2003.

Liu, X., Day, D. A., Krechmer, J. E., Brown, W., Peng, Z., Ziemann, P. J. and Jimenez, J. L.: Direct measurements of semi-volatile organic compound dynamics show near-unity mass accommodation coefficients for diverse aerosols, Commun. Chem., 2(1), 98, doi:10.1038/s42004-019-0200-x, 2019.

Palm, B. B., Campuzano-Jost, P., Ortega, A. M., Day, D. A., Kaser, L., Jud, W., Karl, T., Hansel, A., Hunter, J. F., Cross, E. S., Kroll, J. H., Peng, Z., Brune, W. H. and Jimenez, J. L.:

In situ secondary organic aerosol formation from ambient pine forest air using an oxidation flow reactor, Atmos. Chem. Phys., 16(5), 2943–2970, doi:10.5194/acp-16-2943-2016, 2016.

Pieber, S. M., El Haddad, I., Slowik, J. G., Canagaratna, M. R., Jayne, J. T., Platt, S. M., Bozzetti, C., Daellenbach, K. R., Fröhlich, R., Vlachou, A., Klein, F., Dommen, J., Miljevic, B., Jiménez, J. L., Worsnop, D. R., Baltensperger, U. and Prévôt, A. S. H.: Inorganic Salt Interference on CO2+ in Aerodyne AMS and ACSM Organic Aerosol Composition Studies, Environ. Sci. Technol., 50(19), 10494–10503, doi:10.1021/acs.est.6b01035, 2016.

Tang, M. J., Shiraiwa, M., Pöschl, U., Cox, R. A. and Kalberer, M.: Compilation and evaluation of gas phase diffusion coefficients of reactive trace gases in the atmosphere: Volume 2. Diffusivities of organic compounds, pressure-normalised mean free paths, and average Knudsen numbers for gas uptake calculations, Atmos. Chem. Phys., 15(10), 5585–5598, doi:10.5194/acp-15-5585-2015, 2015.

Yassine, M. M., Harir, M., Dabek-Zlotorzynska, E. and Schmitt-Kopplin, P.: Structural characterization of organic aerosol using Fourier transform ion cyclotron resonance mass spectrometry: aromaticity equivalent approach, Rapid Commun. Mass Spectrom., 28(22), 2445–2454, doi:10.1002/rcm.7038, 2014.