

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

SUPPLEMENTARY INFORMATION

Novel Insights on New Particle Formation Derived from a Pan-European Observing System

M. Dall'Osto, D.C.S. Beddows, A. Asmi, L. Poulain, L. Hao, E. Freney, J. D. Allan, M. Canagaratna, M. Crippa, F. Bianchi, G. de Leeuw, A. Eriksson, H. C. Hansson, J. S. Henzing, C. Granier, P. Laj, T. Onasch, A. Prevot, J. P. Putaud, K. Sellegri, E. Swietlicki, M. Vidal, A. Virtanen, K. Zemankova, R. Simo, D. Worsnop, C. O'Dowd, M. Kulmala and Roy. M. Harrison

1 **1. SMPS inspection : Aerosol contribution from urban local sources**

2 The average annual variation and diurnal profiles of the four k-means clusters are shown
3 in Figure S1d. Cluster 3 (nucleation) shows an annual trend peaking in spring (April) and
4 fall (September), and a clear diurnal profile peaking during daylight time (noon). However,
5 cluster 4 was also found to peak during day time.

6

7 After inspection of SMPS data, it was found that cluster 4 was also partially composed of
8 aerosol size distributions attributed to late stages of new particle formation growth events.
9 Indeed, the diurnal profile of cluster 4 (Fig S1d) shows a shift in time (peaking at 3pm)
10 relative to cluster 3 (first stage of new particle formation events). Figure S1c also shows a
11 growing mode of the new particle formation mode.

12

13 Part of the definition of a remote region is that it is not influenced by traffic or local sources.
14 We wanted to investigate if cluster 3 was affected by ultrafine anthropogenic emissions not
15 related to new particle formation. In order to do so, we selected three monitoring sites
16 representative of monitoring stations situated in polluted continental background areas
17 including: Ispra (ISP, Po Valley, Italy, Southern Europe); K-Puszt (KPO, Hungary,
18 Central Europe) and Kosetice (OBK, Czech Republic, Eastern Europe).

19

20 The chart presented in Figure S1e shows the four K-means clusters for two different sets
21 of monitoring stations: "polluted" represents the average of three stations (ISP, KPO, OBK)
22 relative to the average of all sites (ALL). Figure S1e clearly shows that the diurnal profiles
23 differ for polluted, except for Cluster 3 (nucleation). In particular, cluster 2 shows an

1 enhancement in polluted sites in the 6-9am time interval, suggesting it is the major one
2 affected by anthropogenic emissions. In summary, we conclude that cluster 3 is not
3 affected by local contamination, and can be attributed to new particle formation events.

4

5 **2. PMF solution of the AMS-SMPS combined dataset**

6 Number Size Distributions (NSDs) data were obtained from Ref. 28. Although the
7 instruments within the 24-site network of SMPS/DMPS devices used several different size
8 ranges, all the data collected were harmonised into one large matrix by interpolating the
9 data onto a common size bin scale; 121 size bins spanning 1 to 1000 nm with 40 channels
10 per decade were used. AMS data were used from Ref 54. Aerosol mass spectrometer
11 (AMS) measurements were carried out during 26 field campaigns at 17 different sites.
12 Only five monitoring stations were overlapping with AMS and NSD data. Particle time of
13 flight (PToF) AMS data were obtained for nitrate, ammonium, sulphate and organics. Five
14 equally spaced bins were obtained, 20-38nm, 38-72nm, 72-137nm, 137-262nm, 262-
15 500nm. More information can be found in Ref. 51.

16

17 PMF analysis was applied to the AMS-NSD dataset, following the same approach recently
18 described in Ref. 66. Compared to cluster analysis, which groups similar data together,
19 Positive Matrix Factorisation is used to identify the common 'building blocks' within the
20 data. PMF solves the general receptor modelling problem using constrained, weighted,
21 least-squares applied to the input data x which represent a matrix of concentrations, albeit
22 particle or PM, measured at specific intervals during the study⁶⁷. The general model
23 assumes there are p factors F which are interpreted as fixed emission source profiles
24 (NSD and AMS spectra) and impact the receptor site by various amounts - represented by

1 the scores G (time series of the source profiles) - during the measurement. PMF
2 determines the profiles of these factors and calculates their contribution G such that the
3 sum of linear combinations $G \times F$ of closely matches the measured concentration.
4 Mathematically, the observation x_{ij} , at the receptor is represented in the matrix equation X
5 $= G \times F + E$ whose elements are,

$$x_{ij} = \sum_{k=1}^p g_{ik} \cdot f_{kj} + e_{ij} \quad (1)$$

6 The measurements (AMS or NSD concentrations) are indexed by the integer j for the j^{th}
7 time step (hour or day). The term g_{ik} is the contribution of the k^{th} factor to the receptor site
8 on the j^{th} hour/day, f_{kj} is the fraction of the k^{th} factor (AMS or NSD concentrations) that
9 contributes to measurement j . Matrix E , comprises of elements e_{ij} which are the residual
10 e_{ij} values between the measurement and model for the j^{th} measurement on the j^{th} hour

11

12 In PMF, only x_{ij} are known and the goal is to estimate the contributions (g_{ik}) and the
13 fractions (f_{ij}). It is assumed that the contributions and number fractions are all non-
14 negative, hence the “constrained” part of the least-squares. Furthermore, PMF uses
15 uncertainties measured for each of the x_{ij} size-bin. Measurements with high uncertainty are
16 not allowed to influence the estimation of the contribution and fractions as much as those
17 with small uncertainty, thus giving the “weighted” part of the least squares.

18

19 Given the above, it is task of PMF to minimise the sum of the squares Q calculated using
20 equation 2.

$$Q = \sum_{i=1}^n \sum_{j=1}^m \left(\frac{e_{ij}}{s_{ij}} \right)^2 \quad (2)$$

1 where s_{ij} is the uncertainty in the j^{th} measurement for hour/day i and PMF can be operated
 2 in a robust mode, meaning that “outliers” are also not allowed to influence the fitting of the
 3 contributions and profiles.

4

5 The elements of the matrix S , are derived from the uncertainties entered by the user and
 6 these can be entered directly as a matrix using the `X_std-dev` file. The method chosen for
 7 the `X_std-dev` file values is based on the method used in Ref 68 and Ref. 69. In this, S is
 8 calculated using equation 3 (used in PMF2 when selecting EM=-14),

$$s_{ij} = t_{ij} + v_{ij} \max(|x_{ij}|, |y_{ij}|) \quad (3)$$

9 Where x_{ij} are the actual data values and y_{ij} are the equivalent data values fitted by PMF.
 10 Matrices t_{ij} (sometimes thought of as the the estimated detection limit – 10-20% of X for the
 11 NSD data and 21-233% of X for the AMS data) and v_{ij} (comparable to the relative
 12 uncertainty of the measurement - ~6% for the NSD and ~21% for AMS data) are given by

$$t_{ij} = T(x_{ij} + \overline{x_j}) \quad (4)$$

13

$$v_{ij} = V \quad (5)$$

14

15 We chose the empirical values of T and V by trial and error until their calculated Q value
 16 was the closest to the theoretical value it could be, i.e. $Q/Q_{\text{theory}} \sim 1$. For our 3 factor
 17 solution, $T=V=0.1$, implying that the uncertainty is 10% of the maximum of the fitted and

1 actual values of X and the detection limit of each value is taken as 10% of the sum of that
2 value and the mean of all of the values in the same column. For this case, this formulation
3 gives an uncertainty values between 16 and 27 % for the NSD data 42 and 260 % for the
4 AMS data. Q_{theory} is taken as the difference in the sum of the elements in the input matrix
5 X and output matrices G and F.

6

7 In addition to optimising Q/Q_{theory} for the whole data set according to V and T, the Q/Q_{theory}
8 value for both the AMS and NSD data were both adjusted to ~ 1 by optimising two
9 additional multipliers D_{ams} and D_{nsd} . D_{ams} and D_{nsd} were used to scale the uncertainties of
10 the AMS and NSD matrices respectively with the aim to 'balance' the PMF model. The
11 model is balanced when $Q_{\text{NSD}}/Q_{\text{t_NSD}} \sim 1$ and $Q_{\text{AMS}}/Q_{\text{t_AMS}} \sim 1$,

12

13 So for example, for a our preferred 3 factor solution the closest we could get to a balanced
14 solution was whilst still maintaining $Q/Q_{\text{theory}} \sim 1$ was when using $T=0.1$; $V=0.1$; $D_{\text{ams}}=0.2.1$;
15 $D_{\text{ams}}=0.6..$ These values are presented in Table S2 for each solution.

16

17 As with Cluster Analysis, optimum settings need to be derived (e.g. Number of Factors) for
18 the model. For a given uncertainty matrix, metrics derived from the residual matrix can be
19 used to give an approximation of the ideal number of factors^{68,69}. When deciding on the
20 number of factors a useful constraint is to determine when 'factor-splitting' occurs. This is
21 when a factor which fits a source within the data is forced to divide and can be detected by
22 strong linear relationship between the G scores for the divided factors. Inspection of the
23 scaled residuals and the ratio of Q/Q_{theory} (~ 1) are also accepted indicators of a good fit
24 where Q_{theory} approximately equal to the number of entries in your data array. However,

1 more often than not, it is more intuitive to make the final decisions based on how well the
2 model fits the science being explained by the model.

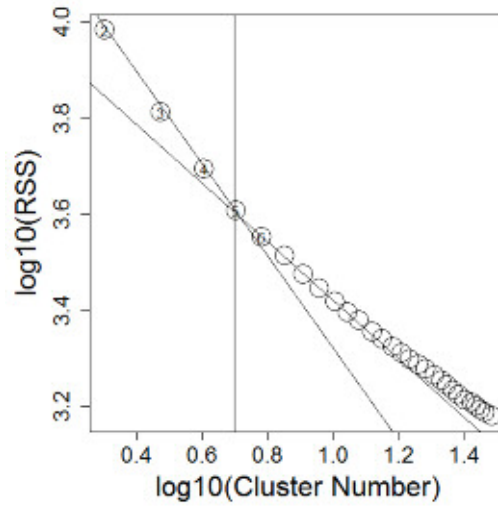
3

4 Table S2 shows the settings used for the initial investigations. For each factor number, the
5 t and v values (see equations 3-5) were adjusted until the ratios of Q/Q_{theory} , Q_{NSD}/Q_{t_NSD}
6 and Q_{AMS}/Q_{t_AMS} , were ~ 1 and after 6 factors evidence of factor splitting became stronger.
7 Hence 3 to 6 factor solutions were considered.

8

9 An estimation of how well the model is fitting to the data for each factor setting can be
10 judged by looking at the scaled residual matrix R and in particular calculating $IM =$
11 $(\text{colMeans}(R))$ the *maximum mean column value of R* and $IS = \max(\text{apply}(R,2,\text{stdev}))$ the
12 *maximum column standard deviation value of R* . These give a idea of the spread of the
13 residuals after the model has been fitted for each factor number and a smaller number
14 indicates a preferred fit. For the fitted models presented in Table S2, the values of IM
15 decrease by about 23 % when the factor number is increased from 3 to 9, whereas IS
16 increases by 16 %. The ratio of IS to IM , remains around 3.5 from 3 to 5 factors before
17 almost doubling when using 8 factors. On increasing the factor number, the expected
18 trend is for IS and IM to decrease sharply with from 2-5 factors before shallowing off over
19 the higher factors. In comparison the maximum rotation value $\max(R)$ increases and using
20 the two a minimum and maximum factor number can be selected. In this case, because
21 we optimise each model for Q/Q_{theory} , only a maximum factor number can be selected, in
22 this case 5. Hence from the study of IM , IS and $\max(\text{Rot})$, factor numbers between 3 and
23 5 can be chosen as being suitable.

1

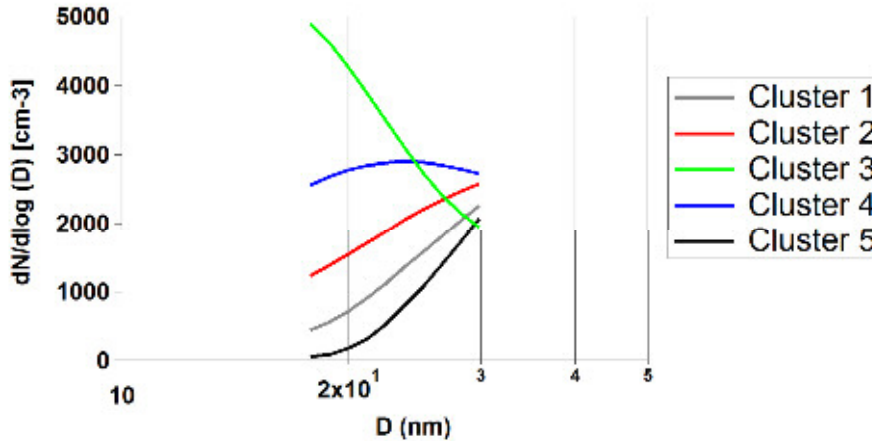


2

3 **Figure S1a.** Plot of Residual Sum of Squares (RRS) against cluster number to identify the
4 main cluster groups within the data. Clusters number 2 to 6 are marked and cluster 5 has
5 been identified as the optimum number of clusters (where the 'knee' in the curve is) using
6 the two tangents marked through the lower 5 and middle 5 to 8 points.

7

8

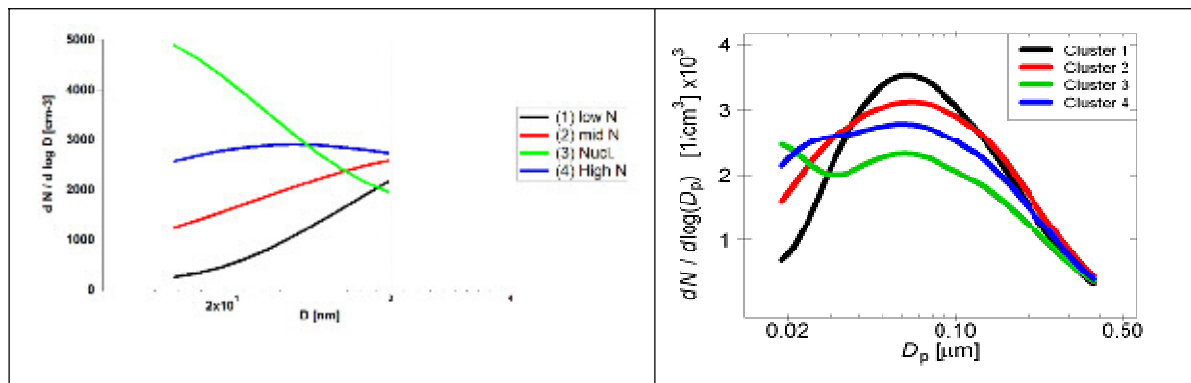


9

10 **Figure S1b.** The aerosol size distributions of the optimum number of clusters. Cluster 1
11 and 5 were found similar to each other and not relevant for high particle concentrations in
12 the studied size range, so they were merged.

13

14



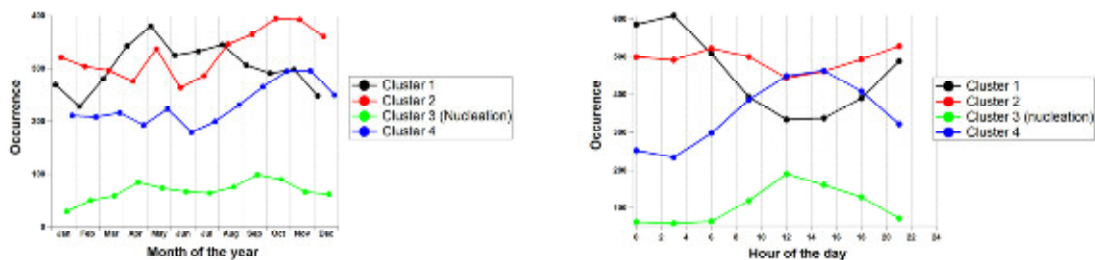
1

2 **Figure S1c.** Aerosol size distributions of the four clusters in the range 17-30 nm resulting
 3 from K-means clustering. On the right side the extracted full aerosol size distribution (17-
 4 500nm) can be seen. The resultant four factors (1-4) contributed 32±20%, 34±9%, 7±4%
 5 and 27±15% respectively of the total size distribution spectra.

6

7

8

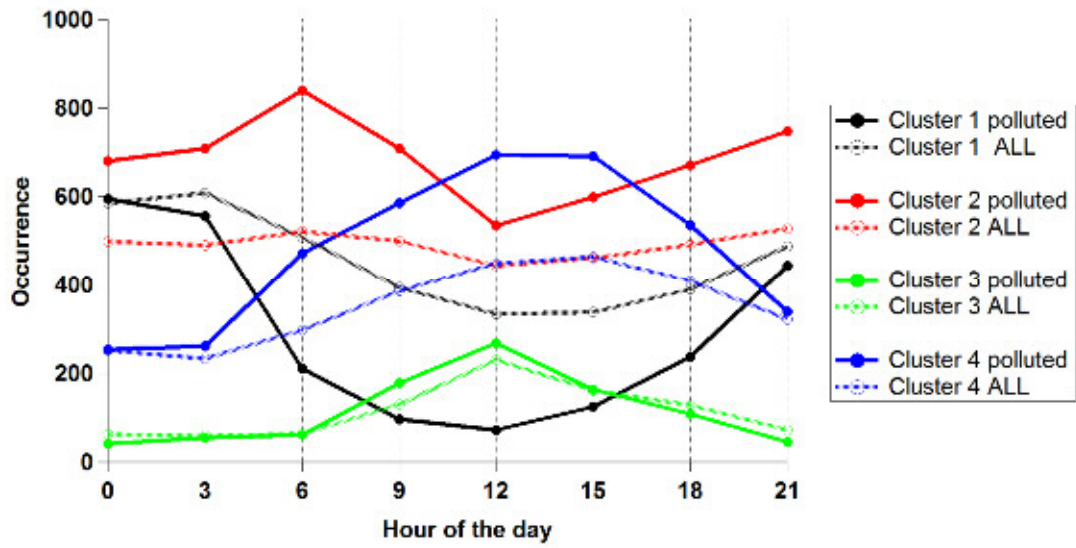


9

10 **S1d.** Annual variation and diurnal profiles of K-means clusters.

11

12



1

2 **Figure S1e.** Monitoring sites influenced by traffic (remote continental stations of ISP,
 3 KPO and MPZ) relative to average of all (ALL).

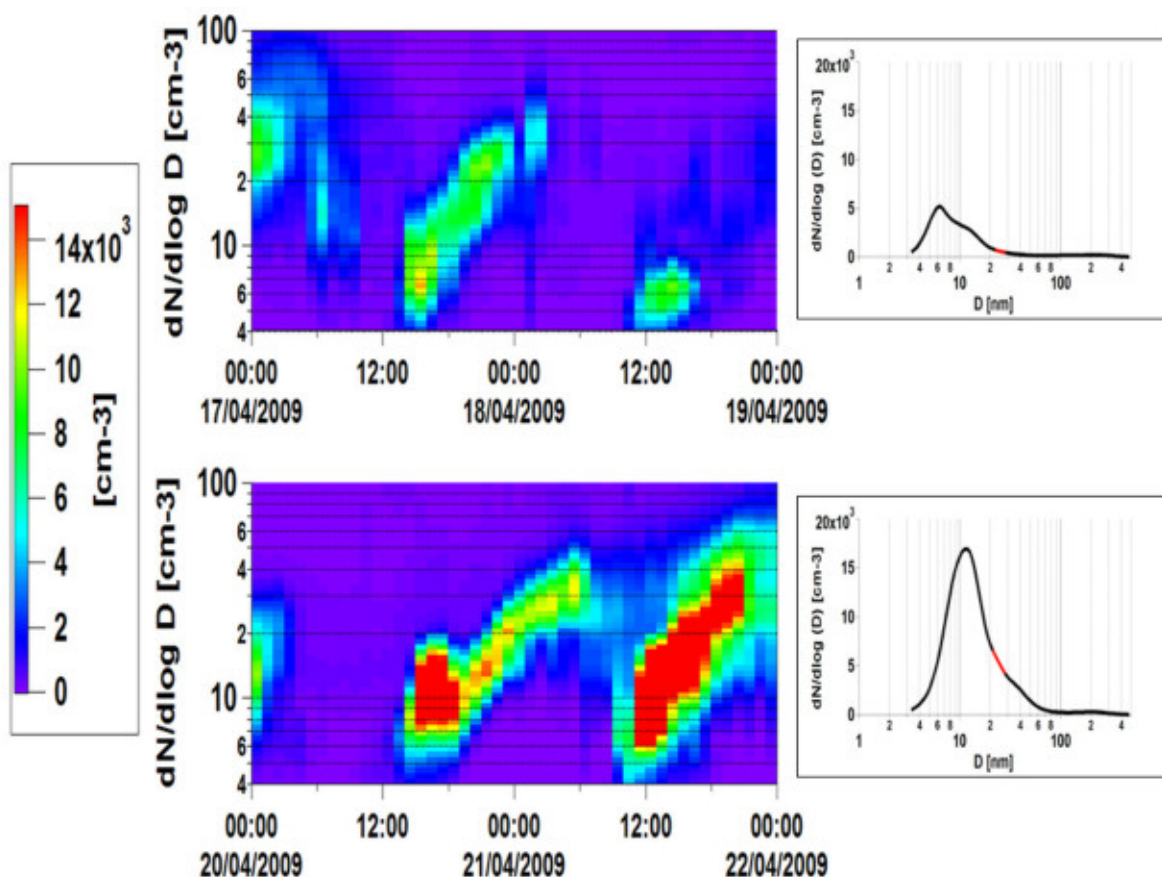
4

5

6

7

8



1

2 **Figure S1f.** Results from this study were compared with those reported by Manninen et
 3 al. (2010)³¹. It is worth mentioning that the present study covers a broader number of
 4 stations (24 vs 10) and a wider measurement period (24 months vs 14). About 55% of the
 5 nucleation events classified in the previous study³¹ were also detected with our
 6 methodology. As an example, four nucleation days classified in Manninen et al. (2010)³¹
 7 from the month of April 2008 (SMR station) are reported. Weak nucleation events in
 8 particle number concentration and growth (17-18/04/2009, Figure S1g top) are not
 9 classified in our study, whereas strong ones (20-21/04/2009, Figure S1g bottom) are. The
 10 right hand panels show aerosol size distributions during the beginning of the event (12:00-
 11 16:00), and in red the aerosol size bins (17-30nm) used in the K-means clustering used in
 12 this study.

13

14

15

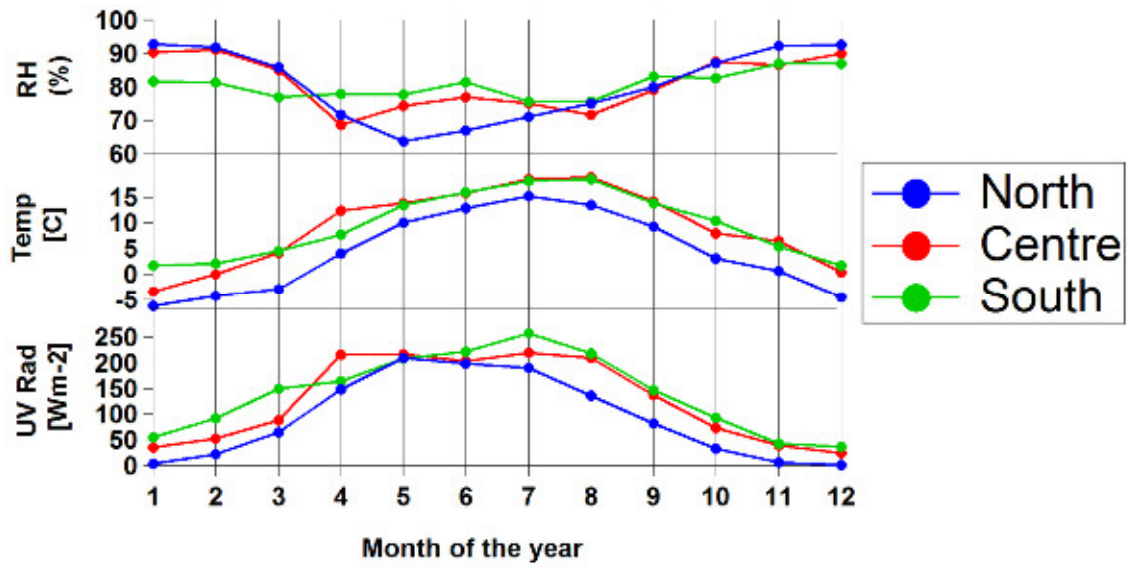


Figure S2a. Average meteorological data across Europe.

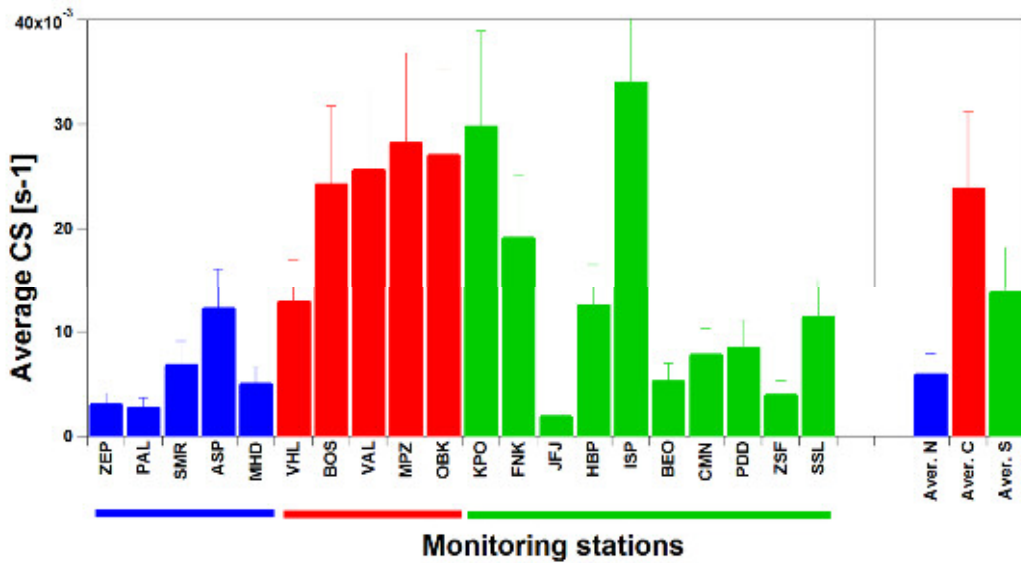
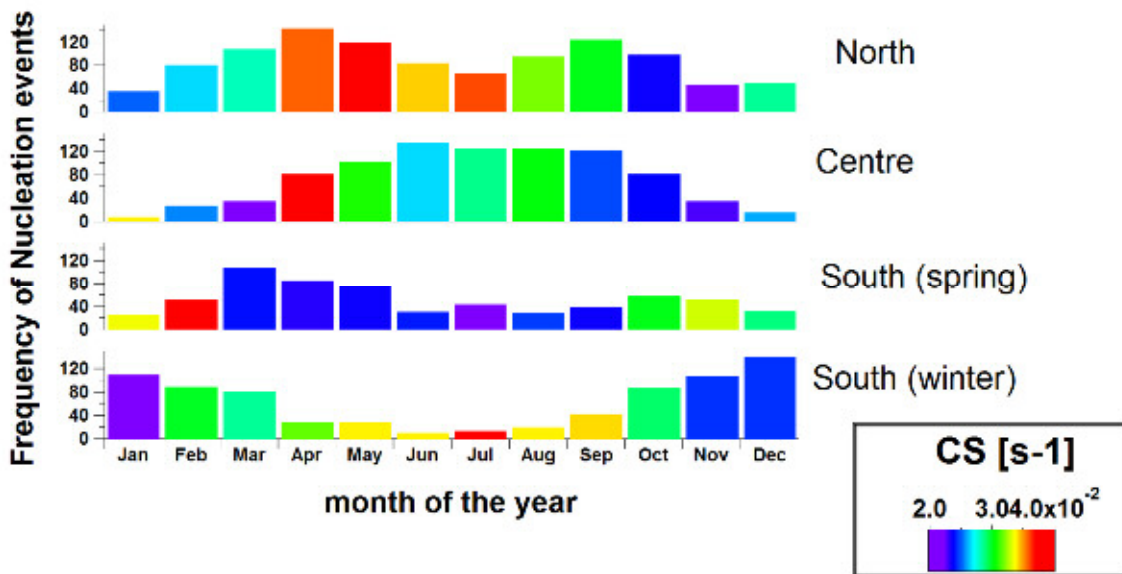
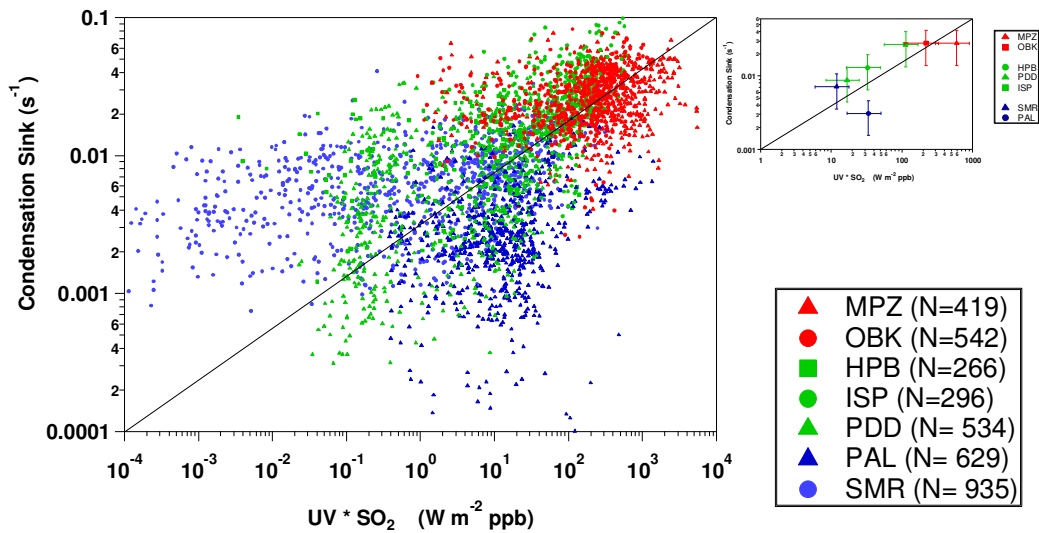


Figure S2b. Average Condensation Sink among the monitoring stations (and averages on the right side, Blue: North; Red: Centre; Green: South).



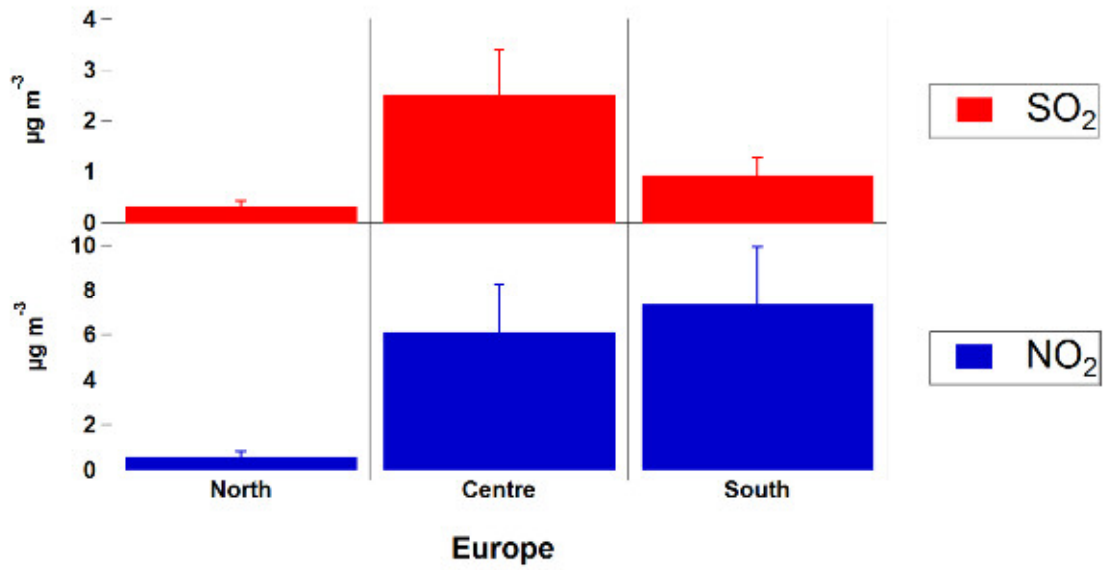
1

2 **Figure S2c.** Annual variation of the CS across the four different European regions.



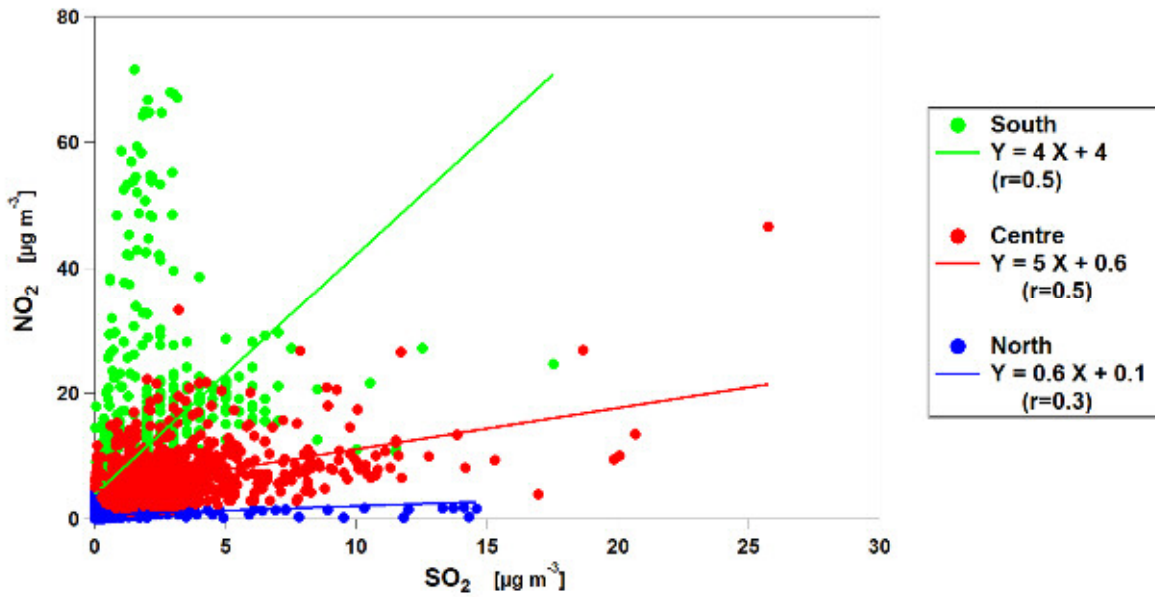
3

4 **Figure 2d:** Relationship of nucleation events at seven European sites with UV * SO₂
 5 (surrogate for H₂SO₄ formation) and condensation sink.



1
2
3
4

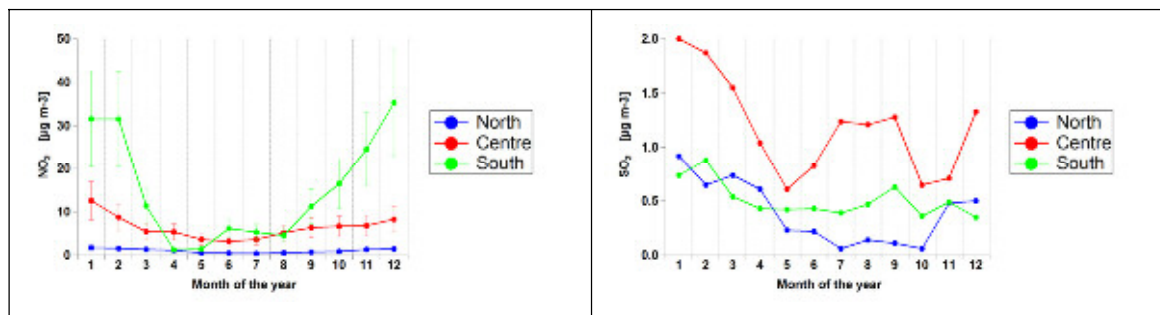
Figure S3a. SO₂ and NO₂ average concentrations for North, Centre and South Europe.



5
6
7
8
9

Figure S3b. NO₂-SO₂ correlation plot for North, Centre and South Europe.

1

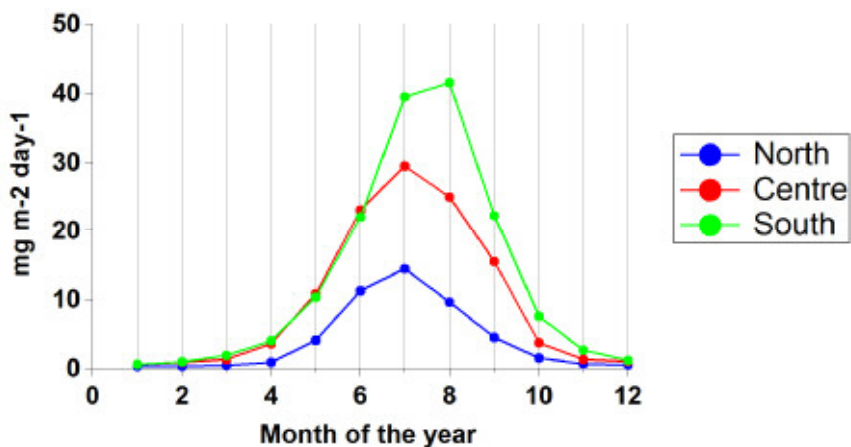


2

3 **Figure S3c.** Annual variation of SO₂ and NO₂ concentrations in North, Centre and South
4 Europe.

5

6



7

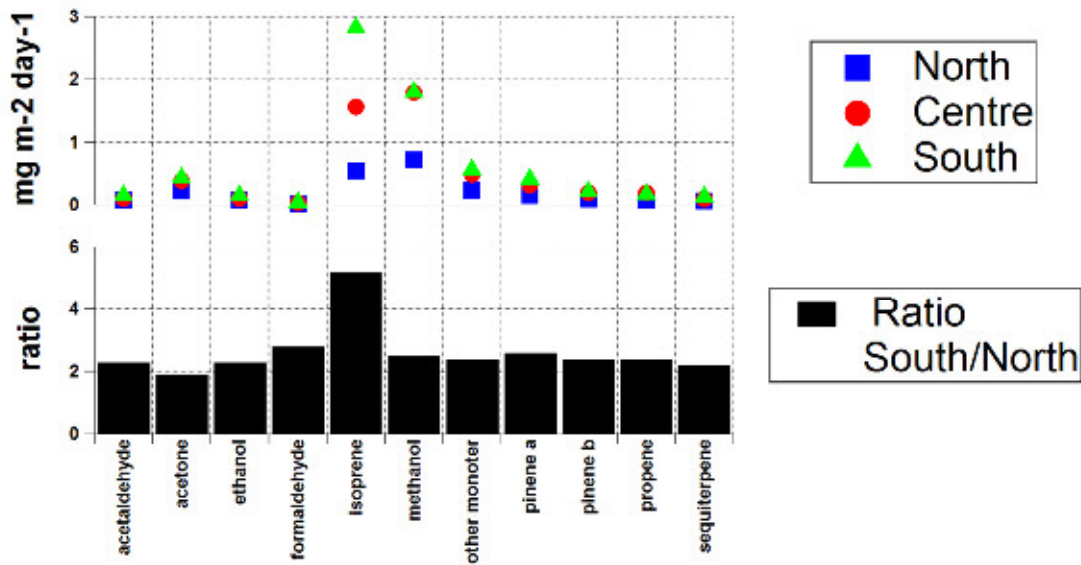
8

9 **Figure S3d.** Annual trend of sum of BVOC fluxes across North, Centre, and South
10 Europe. The annual variation of biogenic VOC (acetaldehyde, acetone, ethanol,
11 formaldehyde, isoprene, methanol, other monoterpenes, α -pinene, β -pinene, propene,
12 sesquiterpene) around Europe does not change dramatically across the different stations,
13 with the months of June-July-August having the highest concentrations.

14

15

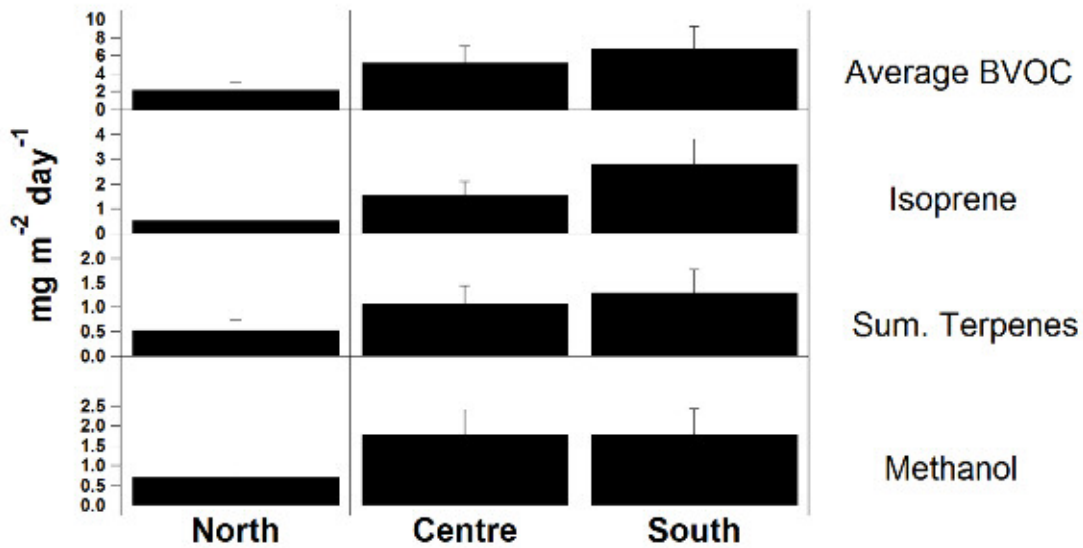
16



1

2 **Figure S3e.** Average fluxes for BVOC across North, Centre, South Europe, and gradient
 3 between South/North Europe. When we take the 24 stations and we plot the average
 4 monthly concentrations for the three different regions, we see that a gradient South-
 5 Centre-North is often observed, especially for isoprene.

6

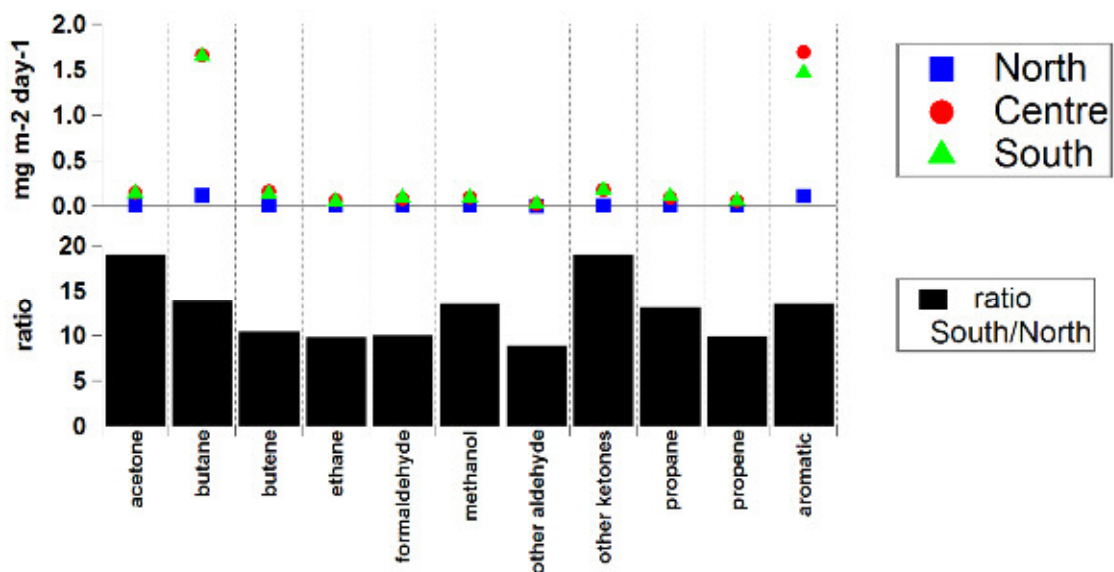


7

8 **Figure S3f.** Average fluxes for selected BVOC across North, Centre, South Europe.

9

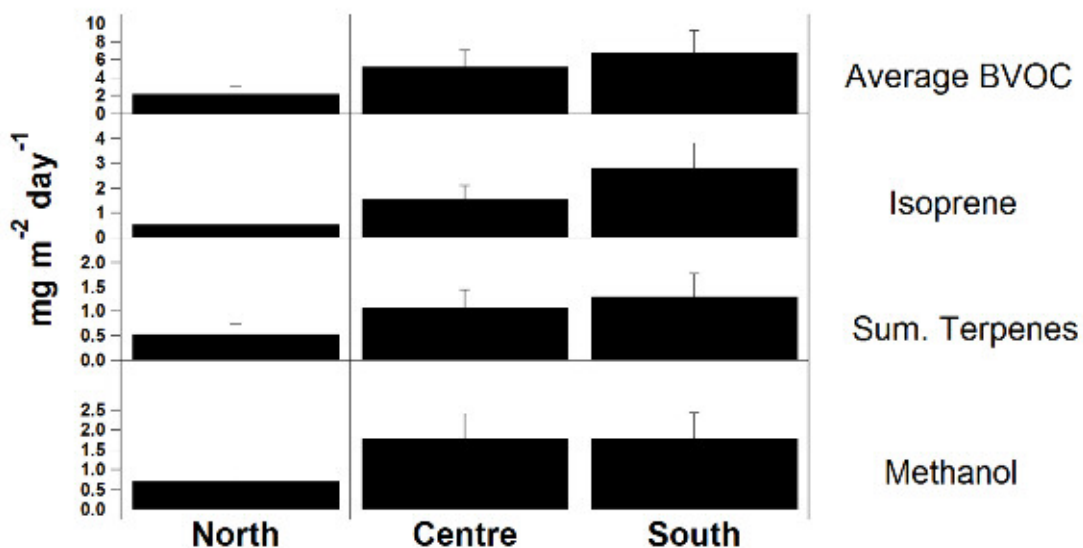
10



1

2 **Figure S3g.** Average fluxes for VOC across North, Centre, South Europe, and gradient
 3 between South/North Europe. The Anthropogenic VOC present very different trends
 4 relative to the Biogenic VOC reported before, and describing such variation is beyond the
 5 scope of this work. Moreover, the ratio South/North for anthropogenic VOC is higher
 6 (about 10-20) than biogenic VOC (about 2), as shown in Figure S4g. Moreover, a clear
 7 difference between centre and south Europe can also be seen. Formaldehyde and other
 8 aldehydes have a North South gradient, whereas aromatic compounds and butane are
 9 higher in the centre part of Europe.

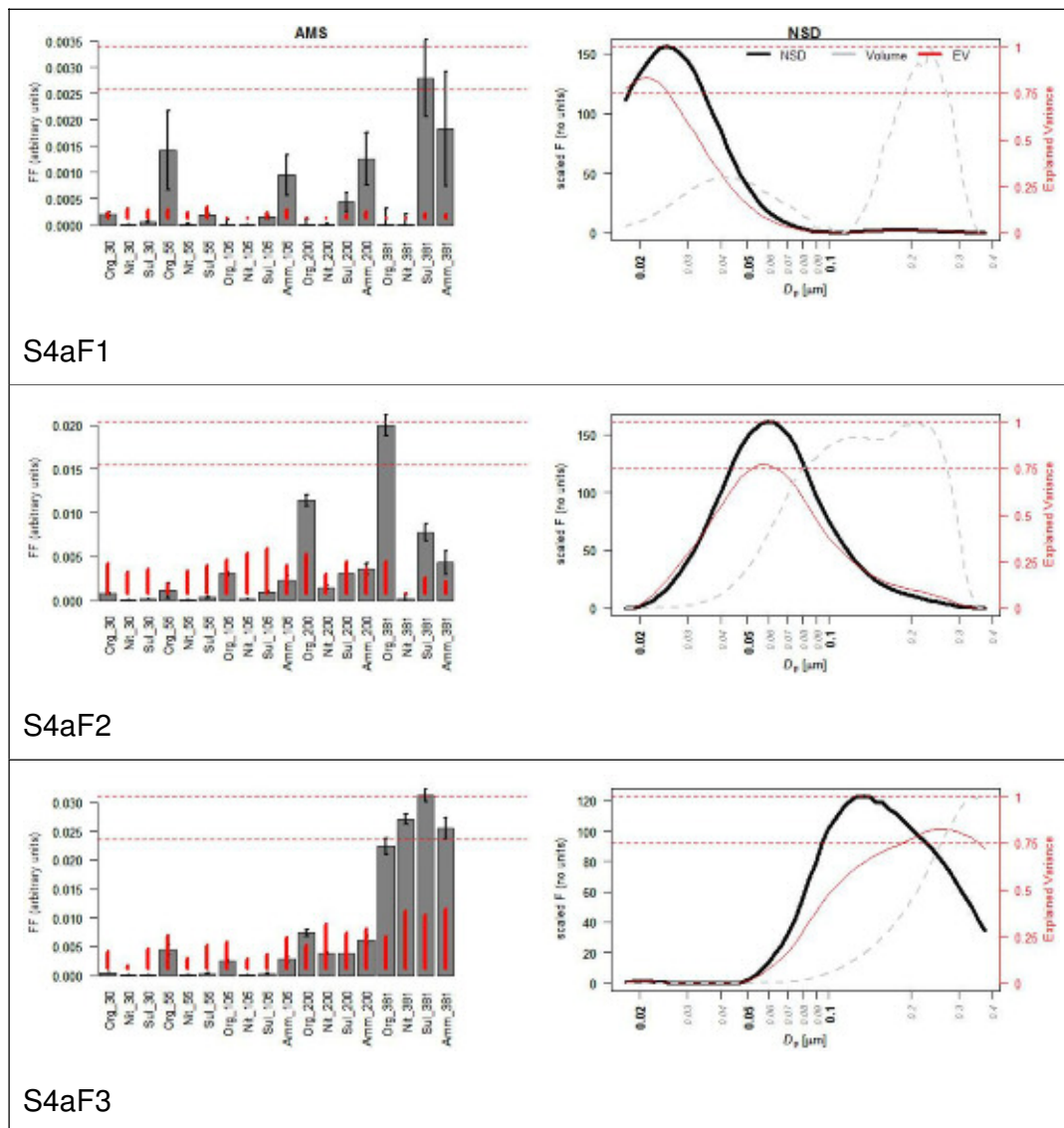
10



11

12 **Figure S3h.** Average fluxes for selected VOC across North, Centre, and South Europe.

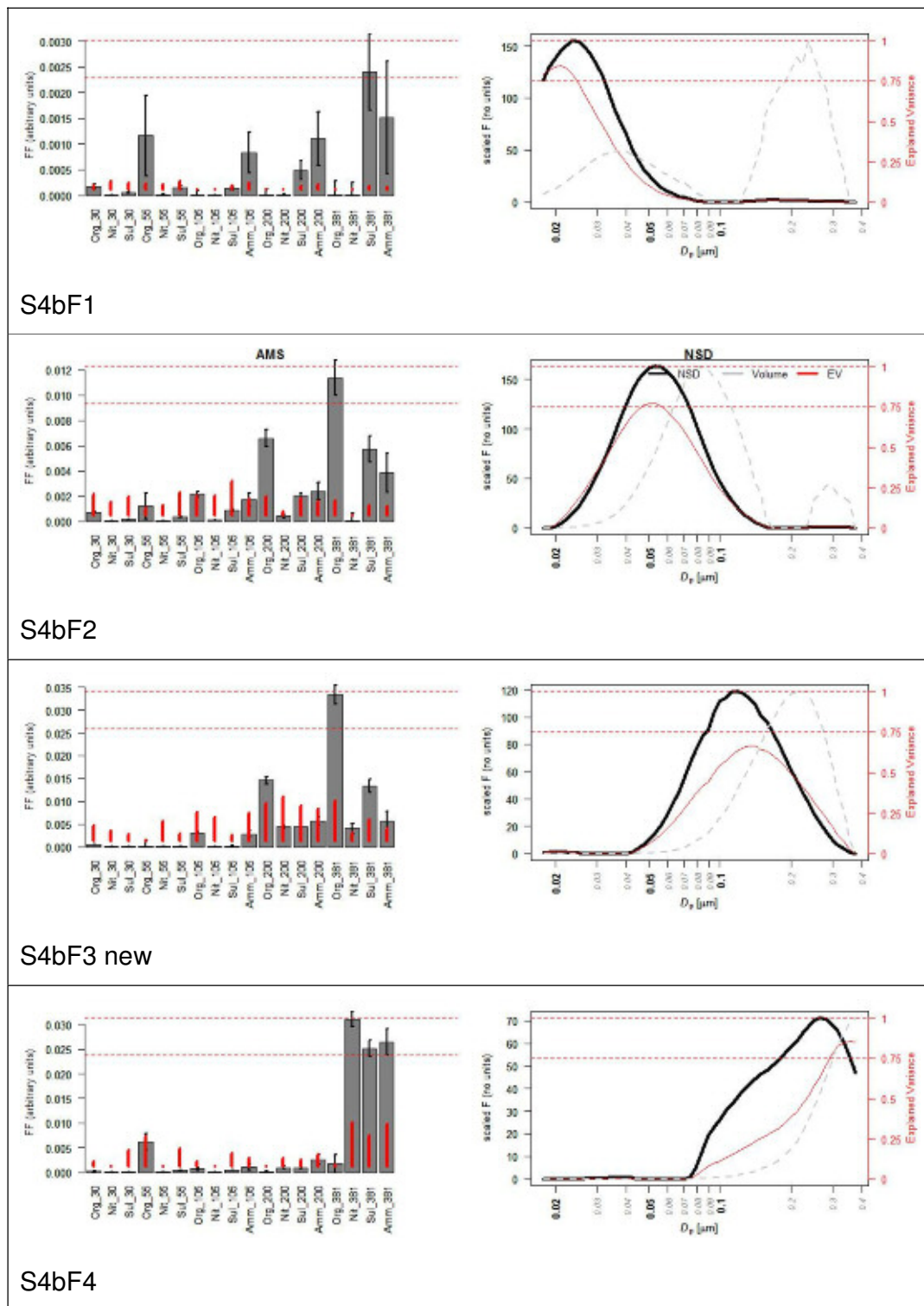
1



2 **Figure S4a.** Three factor PMF solution of combined ToF HR-AMS and SMPS data. The
 3 left-hand panels show AMS analytes according to their size bin (30-381 nm). The right-
 4 hand panels show the SMPS number and volume size distributions and the explained
 5 variation.

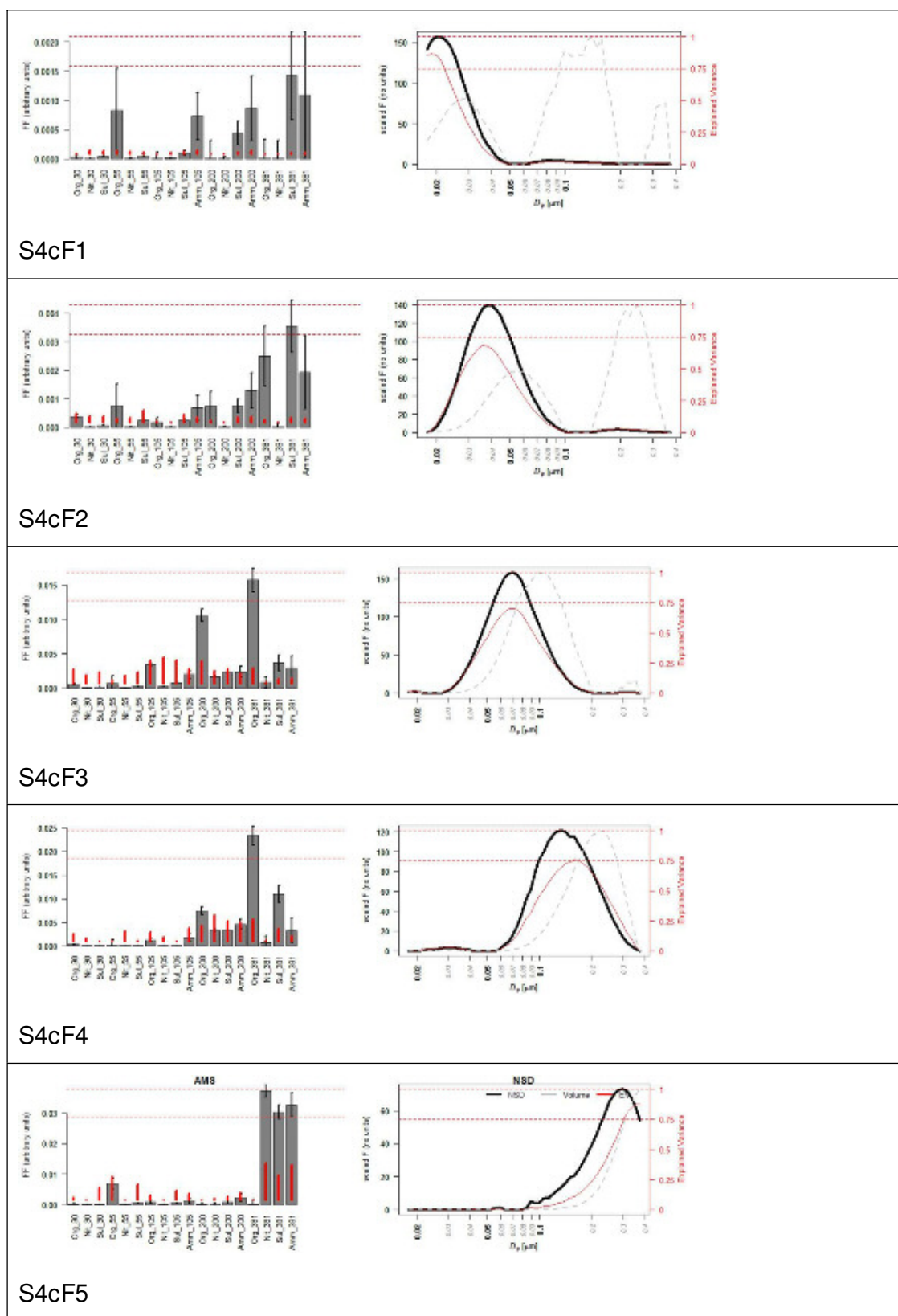
6

1



2 **Figure S4b.** Four factor PMF solution of combined ToF HR-AMS and SMPS data.

3



1 **Figure S4c.** Five factor PMF solution of combined ToF HR-AMS and SMPS data.

1 **Table S1.** Locations and names of stations used in the data analysis. The site altitudes are given
 2 in reference to standard sea level. The areas are grouped by European sub-divisions using
 3 definitions from Central Intelligence Agency (2009). Country codes are given in the ISO 3166
 4 standard. Right columns are for SMPS, AMS and SO₂ availability at the stations considered.

Station name	Station code	Country	Coordinates, altitude (lat., lon., height.)	SMPS	AMS	Gas SO ₂
<i>Nordic and Baltic</i>						
Aspvreten	ASP	SE	58°48'N, 17° 23'E, 30 m	V		
Birkenes	BIR	NO	58°23'N, 8°15'E, 190 m	V		
Pallas	PAL	FI	67° 48'N, 24°7'E, 560 m	V		V
Preila	PLA	LT	55°55'N, 21°0'E, 5 m	V		
SMEAR II	SMR	FI	61°51'N, 24°17'E, 181 m	V	V	V
Vavihil	VHL	SE	56°1'N, 13°9'E, 172 m	V	V	
<i>Central Europe</i>						
Bösel	BOS	DE	53°7'N, 57°57'E, 16 m	V		
K-Pusztá	KPO	HU	46°58'N, 19°19'E, 125 m	V	V	
Melpitz	MPZ	DE	51°32'N, 12°12'E, 87 m	V	V	V
Kosetice	OBK	CZ	49°35'N, 15°15'E, 534 m	V		V
Hohenpeissenberg	HPB	DE	47°48'N, 11°11'E, 988 m	V		V
Waldhof	WAL	DE	52°31'N, 10°46'E, 70 m	V		
<i>Western Europe</i>						
Cabauw	CBW	NL	51°18'N, 4°55'E, 60 m	V		
Harwell	HWL	UK	51°34'N, 1°19'W, 60 m	V		
Mace Head	MHD	IE	53°19'N, 9°53'W, 5 m	V		
<i>Mediterranean</i>						
Finokalia	FKL	GR	35°20'N, 25°40'E, 250 m	V		
JRC-Ispra	ISP	IT	45°49'N, 8°38'E, 209 m	V		V
<i>Arctic</i>						
Zeppelin	ZEP	NO	78°55'N, 11°54'E, 474 m	V		
<i>High Altitude sites (over 1000 msl)</i>						
<i>Western Europe</i>						
Puy de Dôme	PDD	FR	45°46'N, 2°57'E, 1465 m	V	V	V
<i>Central Europe</i>						
Schauinsland	SCH	DE	47°55'N, 7°55'E, 1210 m	V		
Zugspitze	ZSF	DE	47°25'N, 10°59'E, 2650 m	V		
Jungfraujoch	JFJ	CH	46°32'N, 7°59'E, 3580 m	V		
<i>Balkans</i>						
BEO Moussala	BEO	BG	42°10'N, 23°35'E, 2971 m	V		
<i>Mediterranean</i>						
Monte Cimone	CMN	IT	44°11'N, 10°41'E, 2165 m	V		

1
2
3
4
5

Table S2. Initial settings used to investigate the data. For each factor number F , t and v were adjusted until $Q/Q_{\text{theory}} \approx 1$; $Q_{\text{NDS}}/Q_{\text{NDS_theory}} \approx 1$; $Q_{\text{AMS}}/Q_{\text{AMS_theory}} \approx 1$.

	F	3	4	5	6	7	8	9
t	0.100	0.100	0.075	0.100	0.075	0.100	0.075	0.075
v	0.100	0.075	0.100	0.050	0.075	0.015	0.050	0.050
D_{nsd}	0.6	0.475	0.38	0.28	0.235	0.19	0.171	0.171
D_{ams}	2.1	2.4	2.75	2.95	3.3	3.8	4.3	4.3
<i>Factor Splitting (Yes/No)</i>	N	N	N	N	Y	Y	Y	Y
IM	0.55	0.52	0.51	0.47	0.47	0.42	0.42	0.42
IS	1.90	1.92	1.72	1.98	1.84	2.57	2.22	2.22
$Max(ROT)$	0.0038	0.0027	0.0015	0.005	0.004	0.0034	0.0394	0.0394

6
7