

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

**RECEPTOR MODELLING OF BOTH PARTICLE COMPOSITION
AND SIZE DISTRIBUTION FROM A BACKGROUND SITE IN
LONDON, UK – THE TWO STEP APPROACH**

David C.S. Beddows and Roy M. Harrison*†

**National Centre for Atmospheric Science
School of Geography, Earth and Environmental Sciences
University of Birmingham
Edgbaston, Birmingham B15 2TT
United Kingdom**

*To whom correspondence should be addressed.

Tele: +44 121 414 3494; Fax: +44 121 414 3708; Email: r.m.harrison@bham.ac.uk

† Also at: Department of Environmental Sciences / Center of Excellence in Environmental Studies, King Abdulaziz University, PO Box 80203, Jeddah, 21589, Saudi Arabia

28 **ABSTRACT**

29 Some air pollution datasets contain multiple variables with a range of measurement units,
30 and combined analysis by Positive Matrix Factorization (PMF) can be problematic, but can
31 offer benefits from the greater information content. In this work, a novel method is devised
32 and the source apportionment of a mixed unit data set (PM₁₀ mass and Number Size
33 Distribution NSD) is achieved using a novel two-step approach to PMF. In the first step the
34 PM₁₀ data is PMF analysed using a source apportionment approach in order to provide a
35 solution which best describes the environment and conditions considered. The time series
36 G values (and errors) of the PM₁₀ solution are then taken forward into the second step where
37 they are combined with the NSD data and analysed in a second PMF analysis. This results
38 in NSD data associated with the apportioned PM₁₀ factors. We exemplify this approach
39 using data reported in the study of Beddows et al. (2015), producing one solution which
40 unifies the two separate solutions for PM₁₀ and NSD data datasets together. We also show
41 how regression of the NSD size bins and the G time series can be used to elaborate the
42 solution by identifying NSD factors (such as nucleation) not influencing the PM₁₀ mass.

43 **Keywords:** PM₁₀; London; PMF; source apportionment; receptor modelling

44

45 **1. INTRODUCTION**

46 It is unquestionable that worldwide, the scientific vista of air quality is expanding; whether it
47 is the increasing number of observatories or the refinement of information mined from the
48 increasing sophistication of measurements often incorporated in campaign work. The
49 number of metrics being measured has increased from simple measurements of PM mass
50 and gas concentrations, and we can now probe the composition of the PM mass and the
51 size distributions with mass spectrometers, mobility analysers and optical devices.

52
53 Studies using PMF as a tool for source apportionment of particle mass using
54 multicomponent chemical analysis data are published frequently using datasets from around
55 the world. However, they do not always provide consistent outcomes (Pant and Harrison,
56 2012), and one means by which source resolution and identification can be improved is by
57 inclusion of auxiliary data, such as gaseous pollutants (Thimmaiah et al., 2009), particle
58 number count (Masiol et al., 2017) or particle size distribution (Beddows et al., 2015; Ogulei
59 et al., 2006; Leoni et al., 2018).

60
61 Harrison et al. (2011), analysed NSD data (merged SMPS and APS data) with PMF using
62 auxiliary data (meteorology, gas concentration, traffic counts and speed). The study used
63 particle size distribution data collected at the Marylebone Road supersite in London in the
64 autumn of 2007 and put forward a 10 factor solution comprised of roadside and background
65 particle source factors. Sowlat et al., 2016 carried out a similar analysis on number size
66 distribution (13nm - 10µm) data combined with several auxiliary variables collected in Los
67 Angeles. These included BC, EC/OC, PM mass, gaseous pollutants, meteorological, and
68 traffic flow data. A six-factor solution was chosen comprising of: nucleation, 2 x traffic, an
69 urban background aerosol, a secondary aerosol and a soil factor. The two traffic sources
70 contributed up to above 60% of the total number concentrations combined. Nucleation was

71 also observed as a major factor (17%). Urban background aerosol, secondary aerosol, and
72 soil, with relative contributions of approximately 12, 2.1, and 1.1%, respectively, overall
73 accounted for approximately 15% of PM number concentrations, although these factors
74 dominated the PM volume and mass concentrations, due mainly to their larger mode
75 diameters. Chan et al. (2011) considered extracting more source information from an
76 aerosol composition dataset by including data on other air pollutants and wind data in the
77 analysis of a small but comprehensive dataset from a 24-hourly sampling programme
78 carried out during June 2001 in an industrial area in Brisbane. They chose multiple types of
79 composition data (aerosols, VOCs and major gaseous pollutants) and wind data in source
80 apportionment of air pollutants and found it to result in better defined source factors and
81 better fit diagnostics, compared to when non-combined data were used. Likewise, Wang et
82 al. (2017) report an improvement in source profiles when coupling the PMF model with ^{14}C
83 data to constrain the PMF run as *a priori* information.

84

85 However, while combining, for example, particle chemical composition and size distribution
86 data in a single PMF analysis may assist source resolution, difficulties arise if the two
87 datasets have different and/or ambiguous rotations (discussed in Section 2). This tends to
88 result in factors with either mass contributions and small number contributions or number
89 contributions and small mass contributions and rarely a meaningful contribution from both
90 data types. Experimental design can of course circumnavigate this problem, for instance,
91 using chemical data which is already size segregated, measured using a cascade impactor
92 (Contini et al., 2014). Such an approach is attractive by view of the fact that there is no
93 question as to whether both datasets sufficiently overlap across the size bins. However,
94 cascade impactors do not offer the high time resolution of particle counting instruments, with
95 individual measurements lasting hours or days. Even so, for the case where two or more

96 instruments are available in a campaign to measure two or more different metrics, e.g. PM
97 mass and particle number (PN), then a combined data analysis is useful. Emami and Hopke
98 (2017) have shown that the effect of adding variables as auxiliary data (with potentially
99 different units) to a NSD data set is to decrease the rotational ambiguity of a solution from a
100 1-step PMF analysis.

101

102 In this study, we present a method for analysing simultaneously collected PM₁₀ composition
103 and NSD data. In the work of Beddows et al. (2015), both particle composition and number
104 size distribution (NSD) data from a background site in London (2011 and 2012) was
105 analysed using Positive Matrix Factorization. As part of the methodology development, it
106 was concluded that it was preferable not to combine these two data types in a single analysis
107 but to conduct separate PMF analyses for PM₁₀ mass and particle number. This yielded a
108 6 factor solution for the PM₁₀ data (Diffuse Urban; Marine; Secondary; Non-Exhaust
109 Traffic/Crustal (NET/Crustal)); Fuel Oil; and Traffic. Factors described as Diffuse Urban;
110 Secondary; and Traffic were identified in the 4 factor solution for the NSD data, together with
111 a Nucleation factor not seen in the PM₁₀ mass data analysis (see Figure 1). When combining
112 the PM₁₀ and NSD data in a single PMF analysis, Diffuse Urban; Nucleation; Secondary;
113 Aged Marine and Traffic Factors were identified but the factors were not as clearly separated
114 from each other as the factors derived from the separate datasets. For example, Fuel Oil
115 was now mixed in with Marine and called Aged Marine. This is summarized in Figure 1.
116 However, it would still be useful to obtain a number size distribution for each of the 6 PM₁₀
117 factors and/or a chemical composition for the 4 NSD factors. As a continuation of this work,
118 we present an alternative method for analysing the combined dataset in a so called, two-
119 step methodology. In the first step, we analyse the mass data (PM₁₀; units: $\mu\text{g}/\text{m}^3$) according
120 to the methodology of Beddows et al. (2015). This results in a time series factor G which is

121 carried forward into a second PMF analysis of a combined dataset consisting of the G time
122 series and an auxiliary data set (i.e. NSD; units: $1/\text{cm}^3$). The first step identifies sources and
123 apportions the G factors to their contribution to mass and in the second step, an FKEY matrix
124 is chosen such that G 'drives' the model and the NSD data 'follow'. This means that we
125 have PM_{10} factors each of which is augmented by its number size distribution. Furthermore,
126 we also consider linear regression as a second step in a PMF-LR analysis to show that
127 although the initial analysis is biased toward mass by analysing PM_{10} factors only, unseen
128 factors influencing the NSD data (e.g. nucleation) can be identified in the data.

129

130 **2. EXPERIMENTAL**

131 With a population of 8.5 million in 2014 (ONS, 2017), the UK city of London is the focus of
132 study in this work where the London *North Kensington* (NK) Site ($LAT = 51^\circ : 31' : 15.780''$
133 N and $LONG = 0^\circ : 12' : 48.571'' W$) was considered. NK is part of both the London Air
134 Quality Network and the national Automatic Urban and Rural Network and is owned and
135 part-funded by the Royal Borough of Kensington and Chelsea. The facility is located within
136 a self contained cabin within the grounds of Sion Manning School. The nearest road, St.
137 Charles Square, is a quiet residential street approximately 5 metres from the monitoring site
138 and the surrounding area is mainly residential. The nearest heavily trafficked roads are the
139 B450 (~100 m East) and the very busy A40 (~400 m South). For a detailed overview of the
140 air pollution climate at North Kensington, the reader is referred to Bigi and Harrison (2010).

141

142 **2.1 Data**

143 As alluded to, this work is a continuation of the study carried out by Beddows et al (2015),
144 which analysed NSD and PM_{10} chemical composition data collected at the London NK
145 receptor site. Number Size Distribution (NSD) data were collected continuously every 15

146 min using a Scanning Mobility Particle Sizer (SMPS) consisting of a CPC (TSI model 3775)
147 combined with an electrostatic classifier (TSI model 3080) and air dried according to the
148 EUSAAR protocol (Wiedensohler et al., 2012). The particle sizes covered were 51 size bins
149 ranging from 16 nm to 604 nm and the 15 min distributions were aggregated up to hourly
150 averages (where there were at least 3 x 15 min samples per hour) and all missing values
151 were replaced using a value calculated using the method of Polissar et al. (1998). Further
152 details of the SMPS settings are given in Table S1 and the reader is also referred to
153 Beccaceci et al. (2013a,b) for an extensive account of how the NSD data was collected and
154 quality assured.

155

156 Accompanying the NSD data from the study of Beddows et al. (2015) was the PMF output
157 from the analysis of PM₁₀ chemical composition data. The latter data consisted of 24h air
158 samples taken daily over a 2-year period (2011 and 2012) using a Thermo Partisol 2025
159 sampler fitted with a PM₁₀ size selective inlet. These filters were analysed for total metals
160 PM_{metals} (Al, Ba, Ca, Cd, Cr, Cu, Fe, K, Mg, Mo, Na, Ni, Pb, Sn, Sb, Sr, V, and Zn), using a
161 Perkin Elmer/Sciex ELAN 6100DRC following HF acid digestion of GN-4 Metrical membrane
162 filters. Water-soluble ions PM_{ions} (Ca²⁺, Mg²⁺, K, NH₄⁺, Cl⁻, NO₃⁻ and SO₄²⁻) were measured
163 using a near-real-time URG-9000B (hereafter URG) ambient ion monitor (URG Corp). The
164 data capture over the 2 years ranged from 48 to 100% as different sampling instruments
165 varied in reliability. Data gaps were filled by measurements made on daily PM₁₀ filter
166 samples collected continuously at this site using a Partisol 2025; laboratory-based ion
167 chromatography measurements were made for anions on Tissuquartz™2500 QAT-UP
168 filters). No cation measurements were available from these filters, and this resulted in a
169 lower data capture for the cations. Again, all missing data were replaced using a value
170 calculated using the method of Polissar et al. (1998). A woodsmoke metric, CWOD, was

171 also included. This was derived as PM Woodsmoke from the methodology of Sandradewi
172 et al. (2008) utilising Aethalometer and EC/OC data, as described in Fuller et al. (2014).
173 Samples were also collected using a Partisol 2025 with a PM₁₀ size selective inlet and
174 concentrations of elemental carbon (EC) and organic carbon (OC) were measured by
175 collection on quartz filters (Tissuquartz™ 2500 QAT-UP) and analysis using a Sunset
176 Laboratory thermal–optical analyser according to the QUARTZ protocol (which gives results
177 very similar to EUSAAR 2: Cavalli et al., 2010) (NPL, 2013). We refer to CWOD, EC and
178 OC as PM_{carbon}. In addition, particle mass was determined on samples collected on Teflon-
179 coated glass fibre filters (TX40HI20WW) with a Partisol sampler and PM₁₀ size-selective
180 inlet.

181

182 This aforementioned PM₁₀ data was represented in this work as the PMF solution for PM₁₀-
183 only data, derived in Beddows et al. (2015) and consisting of 6 sources, namely: Diffuse
184 Urban; Marine; Secondary; Non-Exhaust Traffic/Crustal; Fuel Oil; and Traffic. The Diffuse
185 Urban factor had a chemical profile indicative of contributions mainly from both woodsmoke
186 (CWOD) and road traffic (Ba, Cu, Fe, Zn). The Marine factor explained much of the variation
187 in the data for Na, Cl⁻ and Mg²⁺, and the Secondary factor was identified from a strong
188 association with NH₄⁺, NO₃⁻, SO₄²⁻ and organic carbon. For the Traffic emissions, the PM
189 did not simply reflect tailpipe emissions, as it also included contributions from non-exhaust
190 sources, i.e. resuspension of road dust and primary PM emissions from brake, clutch and
191 tyre wear. The Non-Exhaust Traffic/Crustal factor explained a high proportion of the variation
192 in the Al, Ca²⁺ and Ti measurements consistent with particles derived from crustal material,
193 derived either from wind-blown or vehicle-induced resuspension. There was also a
194 significant explanation of the variation in elements such as Zn, Pb, Mn, Fe, Cu and Ba, which
195 had a strong association with non-exhaust traffic emissions. As there was a strong

196 contribution of crustal material to particles resuspended from traffic this likely reflected the
197 presence of particulate matter from resuspension and traffic-polluted soils. The last factor
198 was attributed to Fuel Oil, characterised by a strong association with V and Ni together with
199 significant SO_4^{2-} . This output comprised the first-step solution in the 2-step analysis of PM_{10}
200 and NSD data and in this study we concentrate on the analysis of the NSD data in the
201 second PMF step with the aim of assigning a NSD to each of the 6 PM_{10} factors.

202

203 **2.2 Methods**

204 **2.2.1 PMF**

205 Positive Matrix Factorization (PMF) is a well-established multivariate data analysis method
206 used in the field of aerosol science. PMF can be described as a least-squares formulation
207 of factor analysis developed by Paatero (Paatero and Tapper, 1994). It assumes that the
208 ambient aerosol concentration X (represented by $n \times m$ matrix of n observations and m PM_{10}
209 constituents or NSD size bins), measured at one or more sites, can be explained by the
210 product of a source profile matrix F and source contribution matrix G whose elements are
211 given by equation 1:

$$x_{ij} = \sum_{k=1}^p g_{ik} \cdot f_{kj} + e_{ij} \quad i=1\dots n; j=1\dots m \quad (1)$$

212 where the j^{th} PM constituent (element, size bin, or auxiliary measurement) on the i^{th}
213 observation (i.e. hour) is represented by x_{ij} . The term g_{ik} is the contribution of the k^{th} factor
214 to the receptor on the i^{th} hour, f_{kj} is the fraction of the j^{th} PM constituent in the k^{th} factor, and
215 e_{ij} is the residual for the j^{th} measurement on the i^{th} hour. The residuals (i.e. difference
216 between measured and reconstructed concentrations) are accounted for in matrix E and the

217 two matrices G and F are obtained by an iterative algorithm which minimises the object
218 function Q (see equation 2).

219

220 Using the data and uncertainty matrices for the model, equation 1 is optimised in the PMF
221 algorithm by minimising the Q value (equation 2),

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

222

223 where s_{ij} is the uncertainty in the j^{th} measurement for hour i . All analyses were carried out
224 in Robust mode which reduces the impact of outliers (Paatero, 2002).

225

226 PMF is a weighted technique and the value of Q , and hence the model fit, is determined by
227 the input variables with the lowest values of uncertainty, s_{ij} , thus giving their variables a
228 higher weighting in the analysis. Input variables with low weight have little effect upon the
229 value of Q , even when their residuals are large. This can be used to the advantage of the
230 operator, e.g. when apportioning total PM mass in a conventional one-step PMF, the total
231 PM concentrations are normally input with artificially high uncertainty, so that they are
232 essentially passive in the PMF analysis and do not influence its outcome. By doing so, the
233 chemical composition data determine the apportionment of PM mass to the source-related
234 factors identified by the PMF. A similar approach can be followed in the PMF analysis of a
235 combined dataset where higher weightings can be applied to the main dataset of interest
236 such that it “drives” the analysis and the auxillary data set “follows”, i.e. the uncertainties are
237 chosen such that the balance of total weights from the two data sets is *tipped* towards the
238 measurement of interest and highest reliability in regards of rotational unambiguity.

239 To assess the PMF model, the Q value is outputted by PMF and compared to a theoretical
240 value Q_{theory} which is approximately the difference between the product of the dimensions of
241 X and the product of the number of factors and the sum of dimensions of X (i.e. $n \times m - p(n$
242 $+ m)$) $p \times m$. For a given number of factors, the whole uncertainty matrix is scaled by a
243 factor b_{scale} until the ratio between Q and Q_{theory} is approximately one (rQ value = $Q/Q_{\text{theory}} =$
244 1 ± 0.02).

245

246 With regards to the final output from PMF, a scaling has to be applied in order to achieve
247 quantitative results. This is done by scaling either G or F to unity such that the units from X
248 are carried over to either F or G respectively to complete the apportionment. However,
249 different routes have to be considered depending on whether X has homogeneous or
250 heterogeneous units.

251

252 **2.2.2** *1-Step method using data in the same units - homogeneous units*

253 Given a PMF input data matrix X , a solution $GF + E$ can be computed where G represents
254 the time series of the source profiles F , with a residual matrix E . Often X comprises columns
255 of PM_{10} component concentrations (e.g. ICPMS values measured from acid-digested filters
256 collected with a Partisol 2025) and it is common practice to also include a Total variable
257 (e.g. column of PM_{10} , measured using a TEOM) in the data matrix. The resulting PM_{10}
258 profile element value can then be used to scale G and F such that G carries the units of X
259 with F unitless. Note that neither G or F is scaled to unity in this approach. Instead, scaling
260 is done after the analysis using a constant a_k , determined by the time series of a Total
261 variable (e.g. PM_{10}), down weighted by applying a high uncertainty, within the input data.

262

$$x_{ij} = \sum_{k=1}^p (a_k g_{ik}) \left(\frac{f_{kj}}{a_k} \right) \quad (3)$$

263

264 The resulting value for the PM₁₀ contribution for each factor within the F matrix is then used
 265 as a scaling constant a_k in equation 3. Such scaling results in unitless factors F which
 266 describe the characteristics of the sources and time series G with units of $\mu\text{g}/\text{m}^3$.
 267 Apportionment can then be carried out by averaging the G values for each source factor, or
 268 a fully quantified time series of each factor can be presented, e.g. in Bivariate plots. Of
 269 course, the G and F can be normalized such that G is unitless and F carries units; an
 270 approach necessary when X contains heterogeneous units. This approach however,
 271 requires each column of G to be scaled to unity, by using the PMF setting Mean IGI = 1.

272

273 **2.2.3** *1-Step method using data with different units - heterogeneous units*

274 If the analysis of X was to be enhanced by the inclusion of data from a second instrument
 275 with different units, then a different approach to the *1-Step method with homogeneous units*
 276 would be required to analyse the joint data matrix $[X,Z] = G[X,Z] F[X,Z] + E[X,Z]$. If the
 277 previous method was applied where F was normalized, then it would not be clear what units
 278 to assign to G, whether the units from X or Z. To get around this problem, G is scaled to
 279 unity. This results in a unitless time series G and a quantified F matrix. For each source
 280 profile the sum of the species associated with either data type gives the average total
 281 apportionment, e.g. of PM₁₀ or number concentration PN. Of course, this requires the
 282 complete mass or number closure of the elements making up either PM₁₀ or PN respectively,
 283 although inclusion of measurements of total PM₁₀ or PN can be used instead, if available.

284

285 In the ideal case, if the individually computed factors for both data sets result in $G(X)$ and
286 $G(Z)$ being identical, then a straightforward joint model $[X,Z]$ is successful and $G[X,Z] = G(X)$
287 $= G(Z)$. However, if $G(X)$ and $G(Z)$ are significantly different then the joint model will fail,
288 identified by a too large Q value. A solution to this problem is to set the total weights of the
289 better dataset X significantly higher than the total weights of the auxiliary data set Z such
290 that X will “drive the model” and $G[X,Z]$ will be approximately equal to $G(X)$ and a reasonable
291 Q value is obtained for the Z . However, care is required to ensure that X or Z do not contain
292 rotational ambiguity because such rotation for X may not be suitable for Z . For such cases,
293 equal total weights for both X and Z are applied in the hope that the best rotation for both X
294 and Z can be found.

295

296 **2.2.4** *2-Step method using data with different units - heterogeneous units*

297 The method proposed in this work separates the analysis of the two data sets X and Z into
298 two different PMF analyses. Dataset X is first analysed and an unambiguous rotation is
299 selected which gives computed factors $G(X)$. These are then carried over into a second
300 PMF step in which $G(X)$ are combined with Z to form a joint matrix for analysis. By using
301 FKEY (described below) factors, $G(X,Z)$ are forced to be equal to $G(X)$ from step 1. So for
302 example, if in the first step we analyse PM_{10} data and carry forward the output $G(PM_{10})$ into
303 a second step combined with the NSD data, i.e. $[G(PM_{10}),NSD]$ this results in profiles
304 $F[G(PM_{10}),NSD]$. In other words, we force out of the NSD data source profiles which have
305 the same G factors as the PM_{10} data and extend the list of components of the sources
306 identified in the first step and thus improve characterisation of the source. Note that this is
307 equivalent to non-negative weighted regression of matrix Z by columns of matrix G for which
308 other tools exist. Furthermore, by using a two step method, we can continue to use the
309 scaling method described in Section 2.2.2 to apportion the sources using a quantified time

310 series $G(X)$ rather than normalising the $G(X,Z)$ matrix sums to 1 and relying on the
311 summation of the elements in the rows of $F(X,Z)$ to give the apportionment of X and Z. **2.2.5.**

312

313 ***Application of PMF***

314 Positive Matrix Factorization was carried out in this work using the DOS based executable
315 file PMF2 v4.2 compiled by Pentti Paatero and released on Feb 11, 2010 (downloaded from
316 www.helsinki.fi/~paatero/PMF/). This is used by the author in preference to a GUI version of
317 PMF (e.g. US EPA PMF 5.0, Norris et al., 2014) because of the ease with which it can be
318 incorporated into a Cran R procedure script using shell commands, thus facilitating
319 automation of the analysis and any optimisation. R-script can be written to manipulate and
320 organise input data for PMF2, run PMF2, collect the output and produce the necessary
321 output for consideration as text, table or plot. The main strength for this approach is to
322 improve the repeatability and transference of a method between practitioners within our
323 group.

324

325 The two step method is shown schematically in Figure 2. Matrix X yields factors 1G and 1F
326 in the first step. The timeseries 1G matrix is carried through to the second step where it is
327 combined with an auxiliary data set Z, to give the a step 2 input matrix [1G Z]. This in turn
328 is analysed to produce factors 2G and 2F . In the current example, the dataset of Beddows
329 et al. (2015) is used as a starting matrix X and comprises the PM_{10} chemical composition
330 dataset. This yields timeseries 1G and source profile 1F and the reader is referred to
331 Beddows et al. (2015) for a description of the analysis and output. Figure 1 shows the output
332 from the first step which was found to be the optimum solution after considering 3 to 8 factor
333 solutions. The normalised timeseries matrix 1G from this analysis was combined with the

334 NSD data - concurrently measured with the PM₁₀ data - to form the input matrix [¹GZ], for
335 step 2. The uncertainties of the ¹G1 matrix, ¹ΔG are transferred from the output of the first
336 step and entered as input uncertainties for the second step. The hourly NSD data was
337 aggregated into daily values to match the daily ¹G factors outputted from the PMF analysis
338 of the daily PM₁₀ data sampled. This reduced the data matrix down to 590 rows by 57
339 columns (¹G1...¹G6, NSD₁^{16nm}...NSD₅₁^{640nm}) for which we have a Q_{theory} value of 29748
340 for a 6 factor solution. For the NSD data, the uncertainties are taken as the NSD values
341 multiplied by the value of an arbitrary parameter *b*_{scale} (see Figure 2). Initially, *b*_{scale} was set
342 to 4 to ensure that the model was weighted such that it was driven by the PM₁₀ data.
343 However, this operation becomes somewhat redundant by the use of the FKEY matrix
344 discussed in the next section. However, in order to find the optimal NSD uncertainties the
345 value of the parameter *b*_{scale} (typically, 0.2) was optimised in Cran R so that the ratio of
346 Q/Q_{theory} = 1 ± 0.02, indicating an relative percentage uncertainty in the region of 20%. In
347 retrospect – by taking into account the decrease in reliability of the size bin counts towards
348 the edges of the size bin range - an improvement would be to gradually increase the
349 uncertainties from 5% in the middle range of sizes to a pre defined larger value, e.g. 50%,
350 over the lower and upper size bins. The uncertainties were entered directly into the model
351 using PMF matrix T with U and V redundant.

352

353 **2.2.6 Pulling down with GKEY and FKEY**

354 GKEY and FKEY are matrices with the same dimensions as G and F respectively, for
355 incorporating *a priori* information into a PMF analysis. They are used in the second step of
356 the PMF analysis to “pull” elements of the source profiles to zero. GKEY and FKEY indicate
357 the location of suspected zeros in source profiles ²F or contributions ²G (Figure S1). Since
358 we are concerned with the profiles, this information is given in the form of integer values in

359 an FKEY. The greater the certainty that an element of a source profile is zero, the larger the
360 integer value that is specified. In this case, in the second step for the input dataset [¹G
361 NSD], it is certain that only one unique contribution will be strong for each row of the profile
362 ²F, outputted from the second PMF analysis, e.g. only ¹G1 and not ¹G2.. ¹G6 will contribute
363 the to (¹G1, ²F₁) position in output factor ²F₁. (Figure S1). All 'non-zero' elements within
364 the output of ²F take a FKEY value of zero whereas all elements of ²F which are pulled to
365 zero take a non-zero value of *fkey*₁. This leads to a FKEY matrix which can be understood
366 in two parts. The first part is a square matrix of dimension equal to the number of columns
367 of ¹G with all its entries equal to *fkey*₁ except for the leading diagonal; this part ensures that
368 ¹G is the same as ²G. The second part of the matrix consist of all the elements as zero and
369 represents the NSD input data. An *fkey*₁ value of 7 to 9 is considered a medium to strong
370 pull, and in this work, we used a value of 24 which in comparison is very aggressive ensuring
371 only one rotational solution is available ensuring ¹G ≈ ²G.

372

373 To extend the analysis from 6 factors to 7 factors an extra row was added to FKEY. This
374 was done in order to investigate any factors missed in the NSD data which the first analysis
375 using PM₁₀ would not be sensitive to. For example, a nucleation mode would be detected
376 in NSD data but not PM₁₀ data. In order to give the model freedom to factorise out a
377 nucleation factor, the 7th row of of FKEY values consisted {*fkey*₁, *fkey*₂... *fkey*₆, *nsd*₁, *nsd*₂...
378 *nsd*₅}. This ensured that all the ²G contributions were allocated to the first 6 factors only
379 leaving the 7th factor to account for the remaining unfactorised NSD data. There is no reason
380 why more than 7 factors could not be used to investigate possible unresolved NSD factors.
381 However, we constrained the scope of our investigation to reidentifying those in Figure 1.

382

383

384 **2.3 Regression**

385 As an alternative to using PMF in the second step, a regression was carried out. Each
386 column of data for each of the 51 size bins j within the NSD was regressed against the six
387 1G time series using Equation 4

$$NSD_j = \alpha_{0,j} + \alpha_{1,j} \ ^1G_1 + \alpha_{2,j} \ ^1G_2 + \dots + \alpha_{6,j} \ ^1G_6 \quad (4)$$

388

389 where α_0 is the population intercept and α_{1-6} are the populations slope coefficients. This
390 results in a 7 by 51 matrix of values. Each column represents a size bin of the NSD data
391 and each row represents the slope coefficients associated with 6 of the factors (giving an
392 indication of how each size bin scales with each of the 6 factors) and an intercept. When
393 $\alpha_{1-6,j}$ is plotted against the size bin, 6 plots showing the dependence of each size bin j on
394 each of the 6 PM_{10} factors are produced. It is also assumed that these (referred to here as
395 NSD regression source profiles) will be comparable to the actual NSD PMF source profile.
396 Similarly, the $\alpha_{0,j}$ values are expected to give a background value due possibly to noise;
397 however, it is more likely to yield a source (such nucleation) to which the PM_{10} mass analysis
398 is insensitive.

399

400 **2.4 Peak Fitting**

401 If it is assumed that the factors derived from the daily NSD data are the same as those
402 present in the hourly data, i.e. the factors are conserved when averaging the data from
403 hourly to daily data before PMF analysis, then daily NSD profiles can be fitted to the hourly
404 NSD spectra to recover a diurnal cycle for the factors. However, it is worth noting that the
405 process of aggregating hourly data to daily NSD data may cause loss of information implying
406 that minor factors (e.g. due to event episodes) might well be averaged out of the data.

407 Given the j^{th} size bin in the i^{th} number size distribution $NSD_{i,j}$ (of dimensions $M \times N$), the
 408 factors can be fitted using equation (5).

$$D_i = \sum_{i=1}^M d_i \quad (5)$$

409 which is the i^{th} sum D_i of the difference (d_i give by equation 6) across the size bins of the i^{th}
 410 NSD_i and the linear sum of the p NSD source profiles ($p = 7$ in this case) scaled with respect
 411 to the scalar values c_{ik} , representing the timeseries of each fitted NSD source profile.

$$d_i = \sum_{j=1}^N \left\{ NSD_{ij} - \sum_{k=0}^p c_{ik} \times f_{kj} \right\}, \quad c_{ik} \geq 0 \quad (6)$$

$$1 \times 10^{10}, \quad c_{ik} < 0$$

412

413 The Cran R package Non-Linear Minimization (nlm) (R Core Team, 2018) was used to
 414 minimise the value of D_i with respect to the scalar values c_{ik} with a non-negative constraint
 415 on c_{ik} placed in the function. If a negative value is returned by any of the c_k values then D
 416 returns an excessively large value. Furthermore, in order to extract an apportionment to
 417 number concentration ($1/\text{cm}^3$) the fitted values were scaled using a scalar β_k . Seven values
 418 were derived for β_k by regressing the total particle number (total hourly SMPS) against each
 419 of the fitted values c_k (equation 7).

$$PN = \beta_0 + \beta_1 c_1 + \beta_2 c_2 + \dots + \beta_7 c_7 \quad (7)$$

420 The resulting scaled-fitted values were then used to calculate the PN concentration for each
 421 of the regression source profiles (equation 8) allowing subsequent plotting of the 7 diurnal
 422 cycles.

$$PN_k = \beta_k c_k \quad (8)$$

423

424

425 2.5 Bivariate Plot

426 Identification of the sources responsible for the factors outputted from PMF can be assisted
427 by meteorological data. Time series of the k_{th} factor (or g_k values) can be plotted against
428 wind direction and wind speed using either the polarPlot or polarAnnulus functions provided
429 in the Openair package. Polar Plots are simply used for plotting the factor contribution on a
430 polar coordinate plot with North, East, South and West axes. Mean concentrations are
431 calculated for wind speed-direction 'bins' (e.g. 0-1, 1-2 m/s,... and 0-10, 10-20 degrees etc.)
432 and smoothed using a generalized additive model. Each bin concentration is plotted as a
433 group of pixels (coloured according to a concentration-colour scale) and positioned a
434 distance away from the origin according to the magnitude of wind speed and along an angle
435 from the North axis according to the wind direction. Such plots are useful when identifying
436 the nature of the source. A diffuse source will tend to have its highest concentration showing
437 as a *hotspot* at the origin of the polar plot, whereas a point source will cause a *hotspot* both
438 away from the origin and in the direction pointing towards the source. On the other hand
439 wind blown sources tend to be recognised by their relation to wind speed and hence do not
440 necessarily produce *hotspots*. Instead, they produce a minimum to maximum gradual
441 gradient of colour from the origin, spreading radially out towards the edge of the plot in the
442 direction of the source, e.g. for a marine source. Likewise, Annulus Plots plot the mean
443 factor concentration on a colour scale by wind direction and as a function of hour-of-the-day
444 as an annulus, represented by the distance of the coloured pixels from the origin. The
445 function is good for visualising how concentrations of pollutants vary by wind direction and
446 hour of the day. For example, for the North Kensington site – positioned West of the city
447 centre – we might well expect most of the anthropogenic sources (traffic, diffuse urban, etc)
448 to show an Easterly direction with the appropriate diurnal cycle (e.g. rush hour traffic
449 patterns). Similarly, we might expect cleaner air (Marine, Nucleation, etc) to occur from a

450 Westerly direction and at times of the day when the solar strength is highest.

451

452 **3. RESULTS AND DISCUSSION**

453 The aim of this work has been to show how a given PMF result can be complemented with
454 concurrently measured auxillary data. We exemplify this using PM₁₀ and NSD data collected
455 from the North Kensington receptor site in London and start with the premise that we are
456 completely satisfied with the PM₁₀ analysis and are using a rotation which gives quantified
457 factors (quantified G and scaled F) which best represent the urban atmosphere sampled,
458 i.e. the output from Beddows et al. (2015). For each PM₁₀ factor we wish to assign a NSD
459 distribution. Rather than repeat the PMF analysis using a combined PM₁₀+NSD dataset
460 which can be complicated if the rotations of the individual PMF analyses of PM₁₀ and NSD
461 data are mismatched or ambiguous, we can carry out a a second PMF analysis or a
462 regression.

463

464 Furthermore, by the nature of any factor analysis, we also have to make the assumption that
465 each source chemical profile and size distribution not only remain unchanged between
466 source and receptor but that it remains constant throughout the measurement campaign.
467 This of course limits our capacity to fully understand the aerosol within the atmosphere we
468 are considering. Chemical reactions during the transit of the air masses will of course modify
469 the chemical composition. It might be assumed that a fully aged aerosol remains unchanged
470 and is identified as a background component, but for example we would expect progressive
471 chlorine depletion within a fresh marine aerosol passing over a city. Likewise, we also have
472 to appreciate that different particle sizes will have different atmospheric transit efficiencies
473 with large particles settling out of the air mass before smaller ones. Similarly, particles
474 nucleate and grow from 1 nm up to 20-30 nm over a short time period of time. It is these

475 finer details which are missed when making an overall assessment of the chemical and
476 physical composition of an air mass measured over a long period (e.g. 2 years) dataset
477 using PMF.

478

479 **3.1 2-Step PMF-PMF Analysis**

480 Figure 3 presents the profiles 1F_k and 2F_k from the first and second PMF analysis
481 respectively. The plots of 1F_k were carried over from Beddows et al. (2015) to complete the
482 assignment of the source profiles.

483

484 The time series 1G_k and uncertainties ${}^1\Delta G_k$ from the first PMF analysis of PM_{10} data were
485 carried over into the second step where they are combined with the NSD data for PMF
486 analysis (Figure 2). The uncertainties of the NSD data are taken as an optimised multiple
487 of the NSD values themselves ($\sim 5\%$ uncertainty, yielding a Q value of 30,333 in the robust
488 mode; see Table S2 for PMF settings). Also in order to encourage 2G_k to be proportional
489 to 1G_k for $k = 1-6$ (see Table S4), the FKEY matrix is applied to pull elements in the source
490 matrix to zero as described in section 2.3.3. This ensured that the PMF analysis of the NSD
491 data was driven by the 1G time series and resulted in a 6 factor output in which there were
492 unique contributions from the k^{th} factor 1G_k from the first analysis to the k^{th} factor 2F_k in the
493 second analysis. This is mainly due to the aggressive pulling of the factor element in 2F
494 applied using FKEY.

495

496 When inspecting Figure 3 it is notable that the source profiles are surprisingly similar to
497 those calculated for the just-NSD and PM_{10} +NSD data in Beddows et al. (2015). The Diffuse
498 Urban factor has a modal-diameter just below $0.1\ \mu\text{m}$ which is comparable to the same

499 factor in the just-NSD analysis. Marine is comparable to the Aged Marine factor derived
500 from the PM₁₀+NSD analysis. The Secondary factor is again the factor with the largest modal
501 diameter (between 0.4 and 0.5 μm) and traffic has as expected a modal diameter between
502 30 and 40 nm. The Fuel Oil factor appears to be a combination of a nucleation factor and a
503 mode comparable to diesel exhaust seen in the Traffic factor.

504

505 **3.2 2-Step PMF-LR Analysis**

506 Figure S2 shows the results of the linear regression of the NSD data plotted against the
507 PM₁₀ ¹G_k scores and again what is remarkable is the similarity between these regression
508 source profiles and both the factors derived in Beddows et al. (2015) and those from the 2-
509 step PMF-PMF analysis.

510

511 This PMF-LR analysis was carried out using daily averaged data and to obtain hourly
512 information - and thus obtain the diurnal patterns (Figure S2) - the resulting regression
513 source profiles were re-fitted to the original NSD data. On inspection of these source profiles
514 and diurnal plots, the negative values make interpretation a struggle reinforcing one of the
515 4 conditions (Hopke, 1991) in the analysis if it is to make sense. We can however fit non-
516 negative gradients using non-negative regression. However, the surprising consequence of
517 applying this constraint is that the same profiles are derived but they are clipped so that all
518 negative values are replaced by zero values – hence, information is lost by doing this. One
519 interpretation of the negative values is that these are particle sinks but this contradicts the
520 PMF-PMF findings and hence it is concluded that the PMF-LR analysis only serves as an
521 indication of how the PM₁₀ factors are augmented by the NSD data. If all profiles are shifted
522 to above the zero line then comparisons to the PMF-PMF data can be made. However,

523 what is interesting to note in this result is the intercept NSD which is comparable in profile
524 and diurnal pattern to the nucleation mode identified in Beddows et al. (2015). This is a
525 seventh regression source profile, in addition to the 6 PM₁₀ factors and suggests that
526 although the PMF analysis of the PM₁₀ data alone misses a Nucleation factor, this can be
527 recovered in a second analysis as a remainder or bias in the data. Furthermore, this result
528 indicates that the composition of the Nucleation NSD factor has no link to the chemical PM₁₀
529 composition and cannot be used to infer a composition. This is unsurprising given the very
530 small mass contributed by the nucleation mode particles.

531

532 Returning to the PMF-PMF analysis and extending the analysis from 6 factors to 7 factors,
533 an extra row in the FKEY matrix was added to pull all of the ¹G₇ contributions to ²F₇ to zero
534 in the solution (Figure S1). The same FKEY matrix of *fkey*₁ and 0 values was used but this
535 time it was augmented with a 7th row of *fkey*₂ and zero values. In this case, the *fkey*₂ values
536 were set to a value of 20.

537

538 The same 6 factor solution is obtained with the additional 7th factor (Figure 4 and Figure S3)
539 and as expected, this seventh factor was a Nucleation factor. It was suspected that in the
540 6 factor solution, the Nucleation factor was combined with the Fuel-Oil factor. This does not
541 suggest any link between the Nucleation and Fuel-Oil factor other than there was an
542 insufficient number of factors within the model for the two to factorise out of the data giving
543 the Fuel-Oil NSD profile a more reasonable modal peak between 50 and 60 nm rather than
544 20, 30 and 60 nm.

545

546 Beddows et al. (2015), applied a 1-step analysis to three different datasets: PM₁₀-only; NSD-

547 only and PM₁₀+NSD. The analyses of the PM₁₀-only and NSD-only – both with
548 homogeneous units - produced quantitative timeseries G. This was unlike the analysis of
549 the PM₁₀+NSD with heterogeneous units which could not apportion its 5 factors using G but
550 was able to factorise out a Nucleation factor from the data, seen also in the 4 sources in the
551 PMF solution for the NSD-only data. A PM₁₀-only seven factor solution did not reveal this
552 factor, presumably because the mass associated with nucleation mode particles is too small
553 to affect composition significantly. Furthermore, Fuel Oil was not factorised out of the
554 PM₁₀+NSD data and was more likely divided across all 5 factors.

555

556 Another interesting observation is that although only 4 factors were derived from the PMF
557 analysis of NSD-alone (Diffuse Urban; Secondary; Traffic and Nucleation), when extra
558 information is included from the PMF analysis of the PM₁₀ data, more information can be
559 extracted from the PMF analysis of the NSD data in the form of the Marine; Fuel Oil and
560 NET & Crustal factors. The Nucleation factor is only revealed when performing a regression
561 between the NSD size bins and the G scores of the PM₁₀ PMF analysis which leads to
562 increasing the factor number from 6 to 7 which yields the Nucleation profile. It is also
563 reassuring that the bivariate plots for the 7 factors (discussed in the next section) correspond
564 to the bivariate plots given in Beddows et al. (2015). Also note that there is no reason why
565 any further investigation might not explore using more than 7 factors. In fact the Nucleation
566 factor appears at first sight to be multimodal. However, we restricted our analysis to 7
567 factors, considering it complete in terms of identifying the sources obtained by Beddows et
568 al. (2015).

569

570

571 3.3 Diurnal and Bivariate Plots

572 The original PMF was carried out on daily PM₁₀ data and in order to make diurnal and
573 bivariate plots, a higher time resolution is desirable. It is assumed that the factors derived
574 in the hourly NSD data are the same as those derived from the daily averaged data, i.e. the
575 factors are conserved when averaging the data from hourly to daily data before PMF
576 analysis. Then the hourly NSD data can be fit with the PMF profiles derived from the daily
577 data (see Section 2.4). Figure 5 shows the resulting diurnal profiles. The diurnal trends of
578 the parameter c_k (equation 7), required to fit the 7 daily NSD factors to the hourly NSD data
579 are shown. These have been scaled to PN (measured in 1/cm³) using the integral of the
580 NSD (equation 8). The Nucleation factor diurnal trend behaves as expected rising to a
581 maximum during the day and then falling back down to a minimum at night. This
582 corresponds to the intensity of the sun during the day and the increased likelihood of
583 nucleation on clean days when there is sufficient precursor material to form particles with a
584 low particle condensation sink. The Marine factor is also high during the day presumably
585 due to higher wind speeds. Diffuse Urban, NET & Crustal, and Traffic all follow a trend which
586 is synchronised to the daily cycle of anthropogenic activity and traffic as influenced by
587 greater atmospheric stability at night. The Secondary factor shows a small diurnal range.
588 Fuel Oil is highest during the evening and night and may correspond to home heating rather
589 than shipping emissions. The particle size distributions associated with the Marine and NET
590 & Crustal sources are of limited value as these sources are dominated by coarse particles,
591 beyond the range of the SMPS data, although there is a sharp increase in the volume of the
592 particles above 0.5 µm in the Marine factor. As pointed out in Beddows et al. (2015), the
593 Marine factor is identified by its chemical profile of sodium and chloride and is accompanied
594 by an aged nucleation mode at around 30nm. This can be either viewed simply as clean
595 marine air being 'polluted' by traffic emission and/or as the consequence of nucleation

596 occurring over at city in clean maritime air masses (Brines et al. 2015). The key point here
597 is that the factors derived in this work are comparable to those factorised in Beddows et al.
598 (2015) using the combined dataset and the advantage of the 2-step approach is that now
599 we have quantified hourly timeseries G.

600

601 The hourly contributions are aggregated into daily values and plotted as bivariate plots in
602 Figure 5 to assist comparison with the daily plots in Beddows et al. (2015). In that work, the
603 same PMF analysis of the NSD data yielded 4 factors which are named identically to those
604 in the bivariate plots. The similarity of both of the polar and annular plots for each of the 4
605 factors supports our previous factor identification. The Secondary and Diffuse Urban are
606 background sources with strongest contributions in the evening and morning. Traffic is
607 strongest for all wind speeds from the East which makes sense since North Kensington is
608 to the West of the city centre of London where traffic is expected to be most dense.
609 Nucleation is also seen to be strongest for those wind direction from the West which are
610 expected to be cleaner, and have a lower condensation sink. NET & Crustal and Fuel Oil
611 are similar to Diffuse Urban suggesting a similar predominant source location in the centre
612 of London. Marine is observed to be strongest for elevated wind speeds for all wind
613 directions which is consistent with the expected strong contribution for all high wind speeds
614 from the South West, as observed in the daily polar plots in Beddows et al. (2015).

615

616 **3.4 Composition associated with the Nucleation Factor**

617 The Nucleation factor was extracted from the two-step PMF-PMF analysis which included
618 pulling the 1G_1 - 1G_6 to zero of factor 2F_7 . It might be reasonable to suggest that if the two-
619 step PMF-PMF analysis is repeated and the order of analysis of PM_{10} and NSD datasets

620 reversed that it would be possible to derive the chemical conditions within the atmosphere
621 which were conducive to nucleation. For this, the time series of the 4 NSD factors (1G_1 - 1G_4)
622 reported in Beddows et al. (2015) were combined with the PM_{10} data. We again assume
623 that the first PMF step has been carried out and that we are satisfied with how the final
624 solution represents the urban environment of the receptor site and that there are no
625 rotational ambiguities. We then carry out the second step PMF analysis on the 34×591
626 input matrix ($[{}^1G_1 \dots {}^1G_4]$, $PM_{10}[PM, PM_{\text{carbon}}, PM_{\text{ions}}, PM_{\text{metals}}]$). The hourly output
627 uncertainties from the first PMF analysis of the NSD data ${}^1\Delta G_1 \dots {}^1\Delta G_4$ were carried forward
628 into the second PMF analysis by adding them *in quadrature* to give daily uncertainties. As
629 with the analysis of the auxiliary data in the PM_{10} -NSD data, the measurement uncertainties
630 for the PM_{10} data (this time the auxiliary data) was naively taken as 4 times the PM_{10} matrix.
631 Extra care could have been taken in assigning the PM_{10} uncertainties but since we force the
632 output using FKEY a simpler approach was taken. In fact, the FKEY consisted of a 4×4
633 diagonal matrix of zero values with an $fkey_1$ of 20 for all the off-diagonal positions joined to
634 a 4×30 matrix of zeros. Furthermore, the uncertainty values of the PM_{10} were scaled until
635 $Q/Q_{\text{theory}} = 0.99$ using parameter $b_{\text{scale}} = 0.35$ (see Table S3 for more details).

636

637 Ideally, the chemical data would be limited to the composition of the particles in the same
638 size range as the SMPS data. However, when since we are using the PM_{10} composition
639 data we can at best describe the composition of the aerosol which accompanied each factor
640 (Figure S4). For the NSD Secondary factor with its strongest contribution (indicated by the
641 Explained Variation) ~ 400 nm, we have a strong contribution to PM_{10} and $PM_{2.5}$ together
642 with nitrate, sulphate and ammonium. Diffuse Urban, with its strongest contribution at 100
643 nm is accompanied by contributions from elemental carbon and wood smoke indicative of
644 traffic and recreational wood burning. There are also contributions from barium, chromium,

645 iron, molybdenum, antimony and vanadium, all indicative of non-exhaust traffic emissions
646 and the burning of fuel oil. Similarly, the Traffic factor has a modal diameter at roughly 30
647 nm which is indicative of exhaust emissions and this is accompanied by contributions to
648 aluminum, barium, calcium, copper, iron, manganese, titanium and various other metals
649 attributed to vehicles, albeit from tyre or brake wear or resuspension.

650

651 The Nucleation factor with its peak ~20 nm, was associated with marine air as indicated by
652 the strong contributions to Na, Cl and Mg (Figure S4). There are also traces of V, Cr, Ni
653 and a high contribution to PM₁₀ mass which are all associated with marine air. This is
654 explained by an association with the south-westerly wind sector which brings strong winds
655 and marine aerosol rather than reflecting the composition of the nucleation particles
656 themselves. Marine air is considered to provide the conditions required of an air mass
657 conducive to nucleation, i.e. cleaner air with particles with a low condensation sink. As these
658 air masses pass over the land and eventually into London, anthropogenic precursor gases
659 are added to this air which then nucleate particles seen at the receptor site as a nucleation
660 mode. This also goes some way to explain the earlier observation of aged nucleation
661 particles observed in the marine factor in Figure S3. There are also strong contributions to
662 vanadium which is most likely from an unresolved Fuel Oil source being mixed into the
663 Marine and Diffuse Urban factors.

664

665 **4. CONCLUSIONS**

666 A two-step PMF analysis method is presented whereby existing PMF profiles can be extend
667 to incorporate auxillary data concurrently measured and having different units. This is
668 exemplified using PM₁₀ and NSD data.

669

670 When analysing PM₁₀ data, the inclusion of auxillary data such as meteorological, gas and
671 particle number data has proved to give a clearer separation of factors. However, for a
672 successful output, there must be no rotational ambiguity in either the PM₁₀ data or in the
673 auxillary data. In the ideal case, the individually computed factors G(X), G(Z) and G(X,Z)
674 need to be similar if the joint model is to be successful and not produce large residuals and
675 hence a too large Q value. In the best case, the total weight of the PM₁₀ data can be set
676 higher than the auxillary data so that the PM₁₀ data drives the analysis. In this work, we
677 present an alternative method called the 2-step PMF method. In the first step the PM₁₀ data
678 is PMF analysed using the standard approach without the inclusion of additional data. An
679 appropriate solution is derived using the methods described in the literature in order to give
680 an initial separation of source factors. The time series G (and errors) of the PM₁₀ solution
681 are then taken forward into the second step where they are combined with the NSD data.
682 The PMF analysis is then repeated using the combined and mixed unit G time series and
683 NSD dataset. In order to ensure that unique factors are obtained for the G scores, FKEY is
684 used to pull off-diagonal values to zero thus driving the NSD data. This ensures that the
685 NSD factors are specific to the PM₁₀ solution and the PM₁₀ analysis is not affected by any
686 rotational ambiguity of the NSD data. For our demonstration using the Beddows et al. (2015)
687 analysis, this results in 6 PM₁₀ factors whose time series are not only apportioned in mass
688 but the source profiles are identified for the NSD data. Comparisons of both the factor
689 profiles, diurnal trends and bivariate plots to those of Beddows et al. (2015), show that this
690 technique produces one solution linking the two separate solutions for PM₁₀ and NSD data
691 datasets together. This generates confidence that the NSD and PM₁₀ factors ascribed to
692 one source are in fact attributable to that same source.

693

694 Hence, the process starts with a dataset which produces a solution which is sensitive to
695 mass but the factors more sensitive to number can be accessed using a second step.
696 Furthermore, by exploring a higher number of factors, NSD factors which are insensitive to
697 PM₁₀ mass can be identified as in the case of the Nucleation factor. This information can
698 also be extracted using a linear regression PMF-LR where the size bins of the NSD data are
699 regressed against the PM₁₀ PMF time series. For this dataset, the Nucleation factor profile
700 is identified as an intercept within the fitted model leading to an increase in the number of
701 PMF factors from 6 to 7.

702

703 **5. ACKNOWLEDGEMENTS**

704 The National Centre for Atmospheric Science is funded by the U.K. Natural Environment
705 Research Council. Figures were produced using CRAN R and Openair (R Core Team, 2016;
706 Carslaw and Ropkins, 2012).

707

708

709 **REFERENCES**

- 710 Beccaceci, S., Mustoe, C., Butterfield, D., Tompkins, J., Sarantaridis, D., Quincey, D.,
711 Brown, R., Green, D., Grieve, A., Jones, A.: Airborne Particulate Concentrations and
712 Numbers in the United Kingdom (phase 3), Annual Report 2011, NPL Report as 74, 2013a,
713 [https://uk-](https://uk-air.defra.gov.uk/assets/documents/reports/cat05/1306241448_Particles_Network_Annual_Report_2011_(AS74).pdf)
714 [air.defra.gov.uk/assets/documents/reports/cat05/1306241448_Particles_Network_Annual](https://uk-air.defra.gov.uk/assets/documents/reports/cat05/1306241448_Particles_Network_Annual_Report_2011_(AS74).pdf)
715 [Report 2011 \(AS74\).pdf](https://uk-air.defra.gov.uk/assets/documents/reports/cat05/1306241448_Particles_Network_Annual_Report_2011_(AS74).pdf).
716
- 717 Beccaceci, S., Mustoe, C., Butterfield, D., Tompkins, J., Sarantaridis, D., Quincey, D.,
718 Brown, R., Green, D., Fuller, G., Tremper, A., Priestman, M., Font, A. F., Jones, A.: Airborne
719 Particulate Concentrations and Numbers in the United Kingdom (phase 3), Annual Report
720 2012, NPL Report as 74, 2013b, [https://uk-](https://uk-air.defra.gov.uk/assets/documents/reports/cat05/1312100920_Particles_Network_Annual_report_2012_AS_83.pdf)
721 [air.defra.gov.uk/assets/documents/reports/cat05/1312100920_Particles_Network_Annual](https://uk-air.defra.gov.uk/assets/documents/reports/cat05/1312100920_Particles_Network_Annual_report_2012_AS_83.pdf)
722 [report 2012 AS 83.pdf](https://uk-air.defra.gov.uk/assets/documents/reports/cat05/1312100920_Particles_Network_Annual_report_2012_AS_83.pdf).
723
- 724 Beddows, D. C. S., Harrison, R. M., Green, D. C., and Fuller, G. W.: Receptor modelling of
725 both particle composition and size distribution from a background site in London, UK, *Atmos.*
726 *Chem. Phys.*, 15, 10107-10125, 2015.
727
- 728 Bigi, A., and Harrison, R. M.: Analysis of the air pollution climate at a central urban
729 background site, *Atmos. Environ.*, 44, 2004-2012, 2010.
730
- 731 Brines, M., Dall'Osto, M, Beddows, D. C. S., Harrison, R. M., Gómez-Moreno, F., Núñez, L.,
732 Artíñano, B., Costabile, F., Gobbi, G. P., Salimi, F., Morawska, L., Sioutas, C., and Querol,
733 X., Traffic and nucleation events as main sources of ultrafine particles in high-insolation
734 developed world cities, *Atmos. Chem. Phys.*, 15, 5929-5945, 2015.
735
- 736 Carslaw, D. C., and Ropkins, K.: openair - an R package for air quality data analysis,
737 *Environ. Model Softw.* 27-28, 52-61, 2012.
738
- 739 Cavalli, F., Viana, M., Yttri, K. E., Genberg, J., and Putaud, J.-P.: Toward a standardised
740 thermal-optical protocol for measuring atmospheric organic and elemental carbon: the
741 EUSAAR protocol, *Atmos. Meas. Tech.*, 3, 79-89, 2010.
742
- 743 Chan, Y.-C., Hawas, O., Hawker, D., Vowles, P., Cohen, D. D., Stelcer, E., Simpson, R.,
744 Golding, G., and Christensen E.: Using multiple type composition data and wind data in
745 PMF analysis to apportion and locate sources of air pollutants, *Atmos. Environ.*, 45, 439-
746 449, 2011.
747
- 748 Contini D., Cesari D., Genga, A., Siciliano, M., Ielpo, P., Guascito, M. R., and Conte, M.:
749 Source apportionment of size-segregated atmospheric particles based on the major water-
750 soluble components in Lecce (Italy), *Sci. Tot. Environ.*, 472, 248-261, 2014.
751
- 752 Emami, F., and Hopke, P. K.: Effect of adding variables on rotational ambiguity in positive
753 matrix factorization solutions, *Chemometr. Intell. Lab.*, 162, 198-202, 2017.
754
- 755 Fuller, G. W., Tremper, A. H., Baker, T. D., Yttri, K. E., and Butterfield, D.: Contribution of
756 wood burning to PM 10 in London, *Atmos. Environ.*, 87, 87-94, 2014.
757

758 Harrison, R. M., Beddows, D. C. S., and Dall'Osto, M.: PMF analysis of wide-range particle
759 size spectra collected on a major highway, *Environ.Sci.Technol.*, 45, 5522-5528, 2011.
760
761 Hopke, P. K.: A guide to Positive Matrix Factorization, *J. Neuroscience*, 2, 1-16, 1991.
762
763 Leoni, C., Pokorna, P., Hovorka, J., Masiol, M., Topinka J., Zhao, Y., Krupal, K., Cliff, S.,
764 Mikuska, P., and Hopke, P. K.: Source apportionment of aerosol particles at a European air
765 pollution hot spot using particle number size distributions and chemical composition,
766 *Environ. Pollut.*, 234, 45-154, 2018.
767
768 Masiol, M., Hopke, P. K., Felton, H. D., Frank, B. P., Rattigan, O. V., Wurth, M. J., and
769 LaDuke, G. H.: Source apportionment of PM_{2.5} chemically speciated mass and particle
770 number concentrations in New York City, *Atmos. Environ.*, 148, 215-229, 2017.
771
772 Norris, G., Duvall, R., Brown, S., and Bai, S.: EPA Positive Matrix Factorization (PMF) 5.0
773 Fundamentals and User Guide, U.S. Environmental Protection Agency, Washington, DC,
774 EPA/600/R-14/108 (NTIS PB2015-105147), 2014.
775
776 Ogulei, D., Hopke, P. K., Zhou, L., Pancras, J. P., Nair, N., and Ondov, J.M.: Source
777 apportionment of Baltimore aerosol from combined size distribution and chemical
778 composition data, *Atmos. Environ.*, 40, S396-S410, 2006.
779
780 Paatero, P.: User's Guide to Positive Matrix Factorization Programs PMF2 and PMF3, Part
781 2, 2002.
782
783 Pant, P., and Harrison, R. M.: Critical review of receptor modelling for particulate matter: A
784 case study of India, *Atmos. Environ.*, 49, 1-12, 2012.
785
786 Polissar, A. V., Hopke, P. K., and Paatero, P.: Atmospheric aerosol over Alaska – 2.
787 Elemental composition and sources, *J. Geophys. Res.-Atmos.*, 103, 9045-19057,1998,
788 doi:10.1029/98JD01212.
789
790 R Core Team. R: A language and environment for statistical computing. R Foundation for
791 Statistical Computing, Vienna, Austria, 2018. Available at: <https://www.r-project.org/>.
792
793 R Core Team. R: A language and environment for statistical computing. R Foundation for
794 Statistical Computing, Vienna, Austria, 2016. Available at: <https://www.R-project.org/>.
795
796
797
798 Sandradewi, J., Prevot, A. S. H., Weingartner, E., Schmidhauser, R., Gysel, M., and
799 Baltensperger, U.: A study of wood burning and traffic aerosols in an Alpine valley using a
800 multi-wavelength Aethalometer, *Atmos. Environ.*, 42, 101-112, 2008.
801
802 Sowlat M., H., Hasheminassab, S., and Sioutas, D.: Source apportionment of ambient
803 particle number concentrations in central Los Angeles using positive matrix factorization
804 (PMF), *Atmos. Chem. Phys.*, 16, 4849-4866, 2016.
805
806 Thimmaiah, D., Hovorka, J., and Hopke, P. K.: Source apportionment of winter submicron
807 Prague aerosols from combined particle number size distribution and gaseous composition

808 data, *Aerosol Air Qual. Res.*, 9, 209-236, 2009.

809

810 Wang, X., Zong, Z., Tian, C., Chen, Y., Luo, C., Li, J., Zhang, G., and Luo, Y.: Combining
811 Positive Matrix Factorization and radiocarbon measurements for source apportionment of
812 PM_{2.5} from a national background site in north China, *Sci. Rep.*, 7, 10648, 2017, doi:
813 10.1038/s41598-017-10762-8.

814

815 Wiedensohler, A., Birmili, W., Nowak, A., Sonntag, A., Weinhold, K., Merkel, M., Wehner,
816 B., Tuch, T., Pfeifer, S., Fiebig, M., Fjaraa, A. M., Asmi, E., Sellegri, K., Depuy, R., Venzac,
817 H., Villani, P., Laj, P., Aalto, P., Ogren, J. A., Swietlicki, E., Williams, P., Roldin, P., Quincey,
818 P., Hüglin, C., Fierz-Schmidhauser, R., Gysel, M., Weingartner, E., Riccobono, F., Santos,
819 S., Gruning, C., Faloon, K., Beddows, D., Harrison, R. M., Monahan, C., Jennings, S. G.,
820 O'Dowd, C. D., Marinoni, A., Horn, H.-G., Keck, L., Jiang, J., Scheckman, J., McMurry, P.
821 H., Deng, Z., Zhao, C. S., Moerman, M., Henzing, B., de Leeuw, G., Loschau, G., and
822 Bastian S.: Mobility particle size spectrometers: Harmonization of technical standards and
823 data structure to facilitate high quality long-term observations of atmospheric particle
824 number size distributions, *Atmos. Meas. Tech.*, 5, 657-685, 2012.

825

826

827 **FIGURE LEGENDS:**

828

829 **Figure 1.** Venn Diagram showing the summary of the findings of Beddows et al. (2015);
830 applying PMF to PM₁₀-only, NSD-only and PM₁₀+NSD datasets. Table shows the
831 apportionment of PM₁₀ and NSD taken from Beddows et al. (2015).

832

833 **Figure 2.** Flow diagram showing the flow of data through the 2-step PMF-PMF analysis.
834 The PMF analyses of single data set X are considered in step 1 and output indicated by
835 factors/uncertainties ¹G, ¹ΔG, ¹F and ¹ΔF. The second PMF analysis is carried out on the
836 joint data set [¹GZ] and yields factors/uncertainties ²G, ²ΔG, ²F and ²ΔF. In our analysis,
837 X and ¹G are the PM₁₀ and resulting time series from the analysis of Beddows et al. (2015)
838 and Z is the auxillary NSD data concurrently measured using a SMPS.

839

840 **Figure 3.** Source profiles ¹F and ²F from both the first and second PMF step using 6
841 factors. [Grey bars and black line indicate the values of F; red lines and dots indicate the
842 explained variations; and grey dotted line indicates the dV/dlogDp.]

843

844 **Figure 4.** Nucleation and Fuel Oil factors derived when extending the second PMF analysis
845 from the 6 factors (shown in Figure 3) to 7 factors. Source profiles ²F₁ to ²F₆ are given in
846 Figure S3. Each plot is divided into 2 showing the output ¹F_k and ²F_k. [Grey bars and black
847 line indicate the values of F; red lines and dots indicate the explained variations; and grey
848 dotted line indicates the dV/dlogDp.]

849

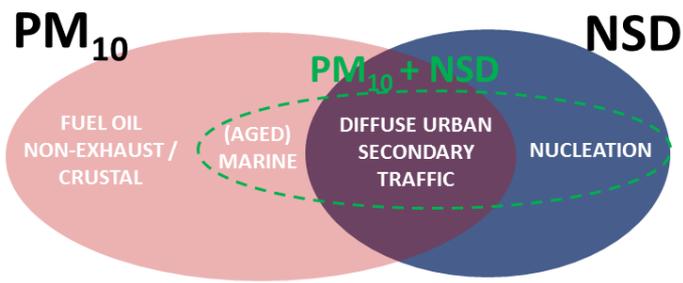
850 **Figure 5.** Diurnal cycles derived *PN_k* calculated by the fitting of the daily PMF factor profiles
851 to the hourly NSD data fitted (see equation 8 and Section 2.4). [Left-left column – diurnal
852 trends of *PN_k*; left-middle column – bivariate plot of *PN_k*; middle-right – annular plot *PN_k*;
853 right-right – bivariate plot of *PN_k*, plotted using the Openair program. Polar plots show a
854 point coloured according to the key, the number concentration at that point on the plot whose
855 distance from the origin represents wind speed and angle wind direction. Likewise for the
856 angular plots the number concentration represent wind direction at an hour of the day
857 between 0 and 23 hrs.]. Note that the diurnal plots do not start at zero.

858

859

860

861
862



	PM ₁₀ [μg/m ³]	NSD [1/cm ³]
Diffuse Urban	4.1	2370
Marine	2.6	-
Secondary	4.4	243
NET / Crustal	4.3	-
Fuel Oil	1.0	-
Traffic	0.8	2460
Nucleation	-	430
Total	17.2	5512

Figure 1. Venn Diagram showing the summary of the findings of Beddows et al. (2015); applying PMF to PM₁₀-only, NSD-only and PM₁₀+NSD datasets. Table shows the apportionment of PM₁₀ and NSD taken from Beddows et al. (2015).

863

864

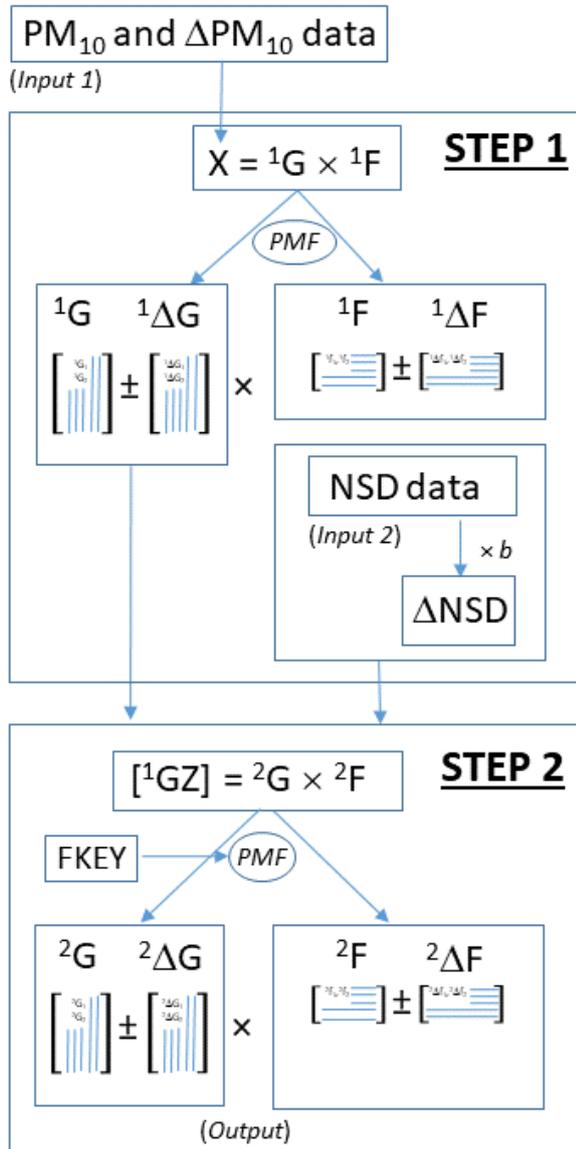


Figure 2. Flow diagram showing the flow of data through the 2-step PMF-PMF analysis. The PMF analyses of single data set X are considered in step 1 and output indicated by factors/uncertainties ¹G, ¹ΔG, ¹F and ¹ΔF. The second PMF analysis is carried out on the joint data set [¹GZ] and yields factors/uncertainties ²G, ²ΔG, ²F and ²ΔF. In our analysis, X and ¹G are the PM₁₀ and resulting time series from the analysis of Beddows et al. (2015) and Z is the auxiliary NSD data concurrently measured using a SMPS.

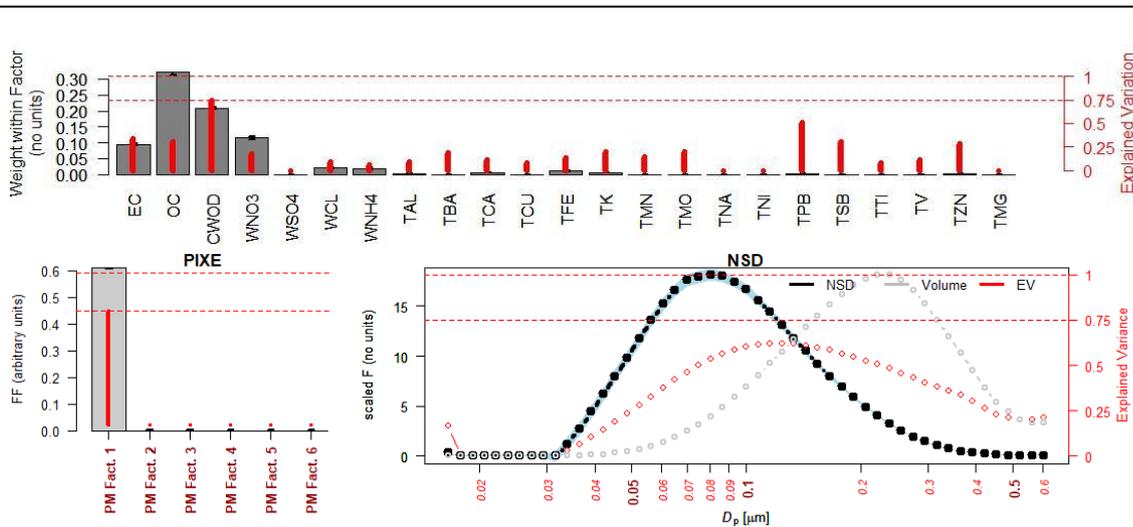
866

867

868

Diffuse Urban

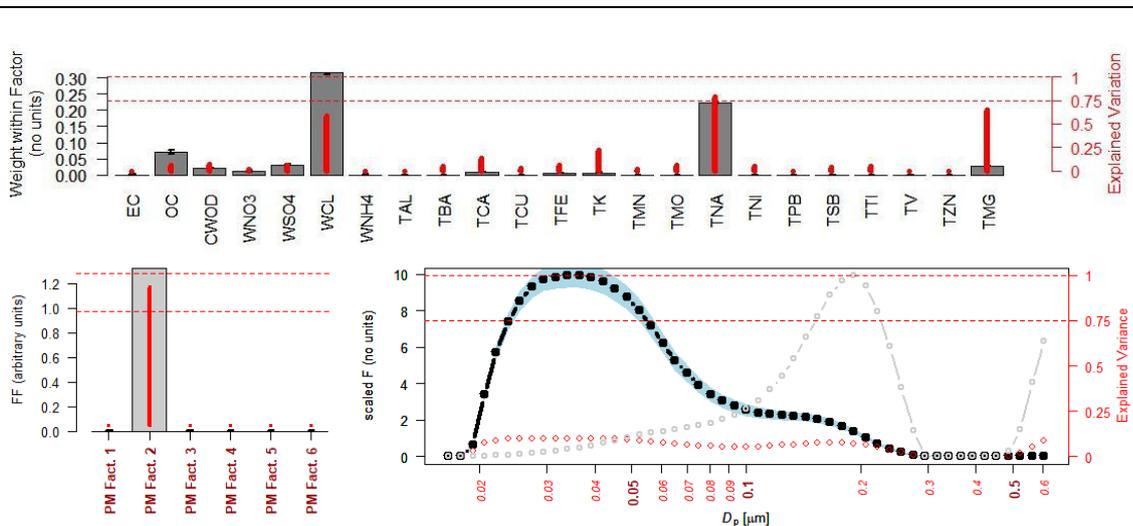
$1F_1$



$2F_1$

Marine

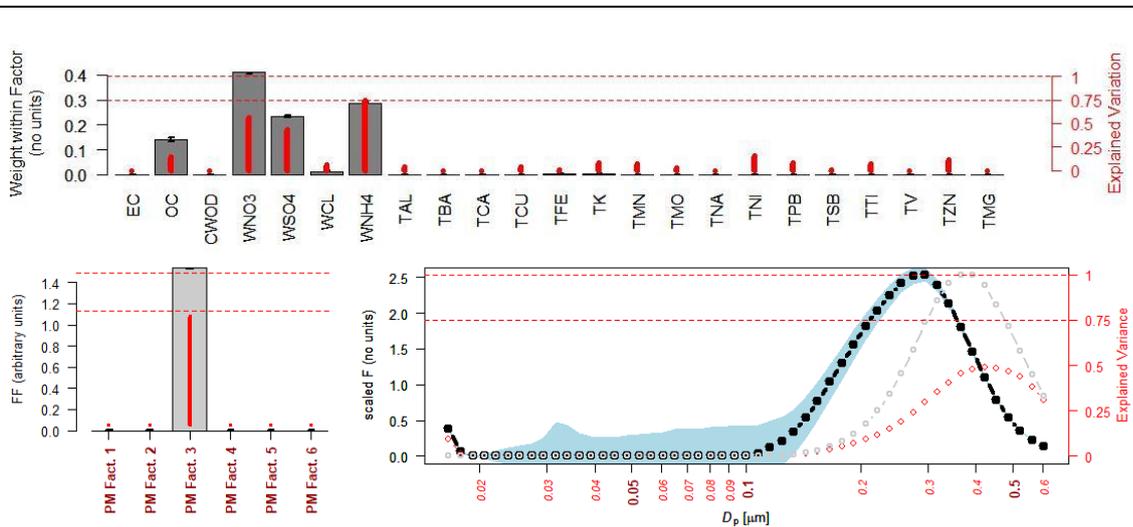
$1F_2$



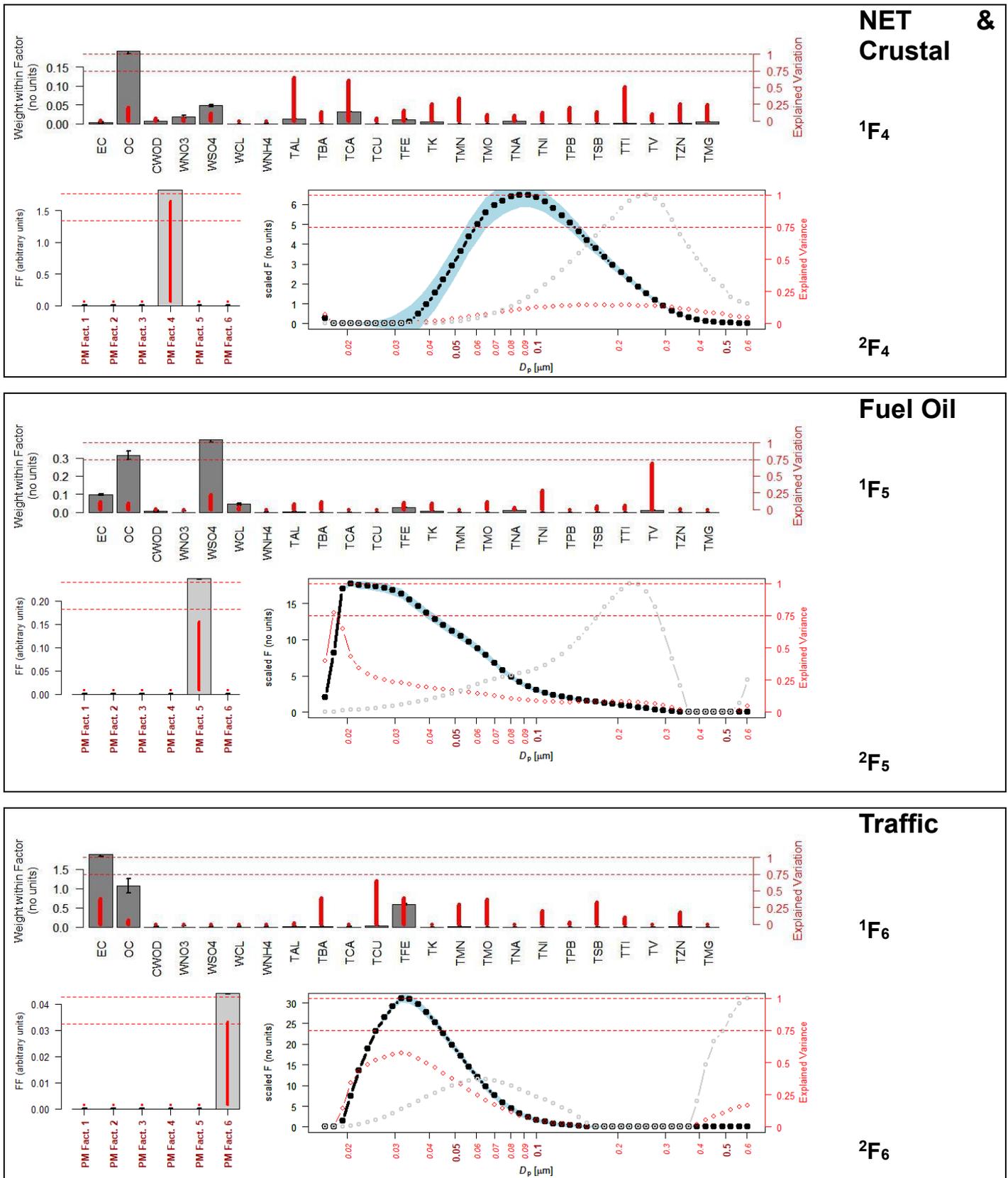
$2F_2$

Secondary

$1F_3$



$2F_3$



870 **Figure 3.** Source profiles 1F and 2F from both the first and second PMF step using 6 factors.
 871 [Grey bars and black line indicate the values of F; red lines and dots indicate the explained
 872 variations; and grey dotted line indicates the $dV/d\log D_p$.]

873

874

875

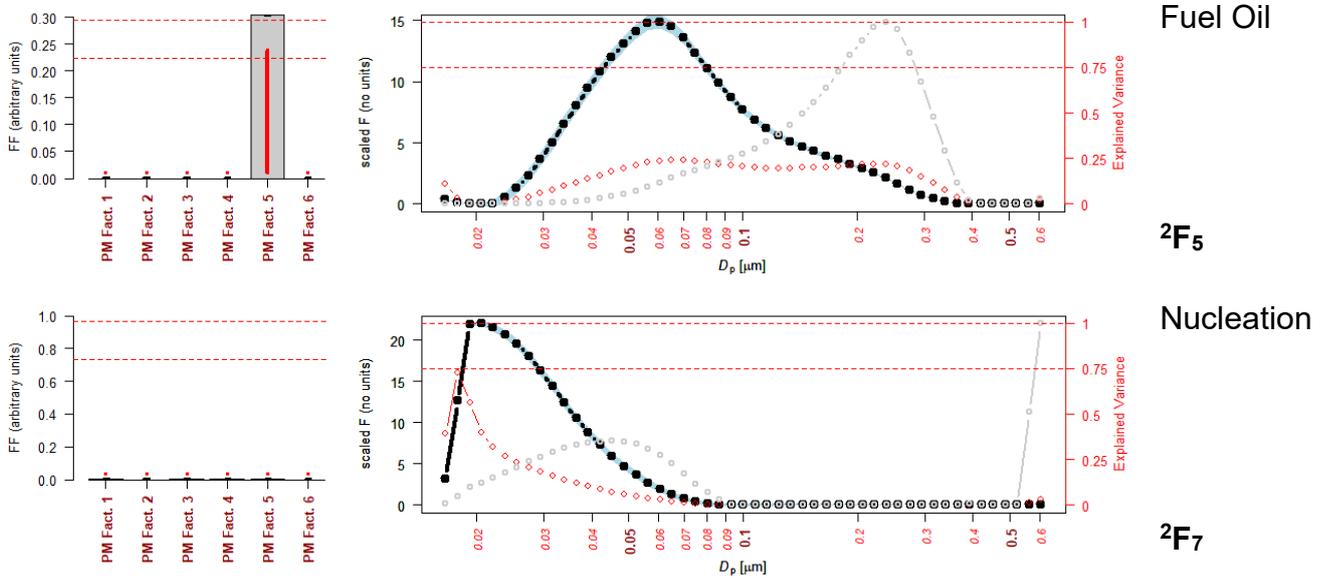


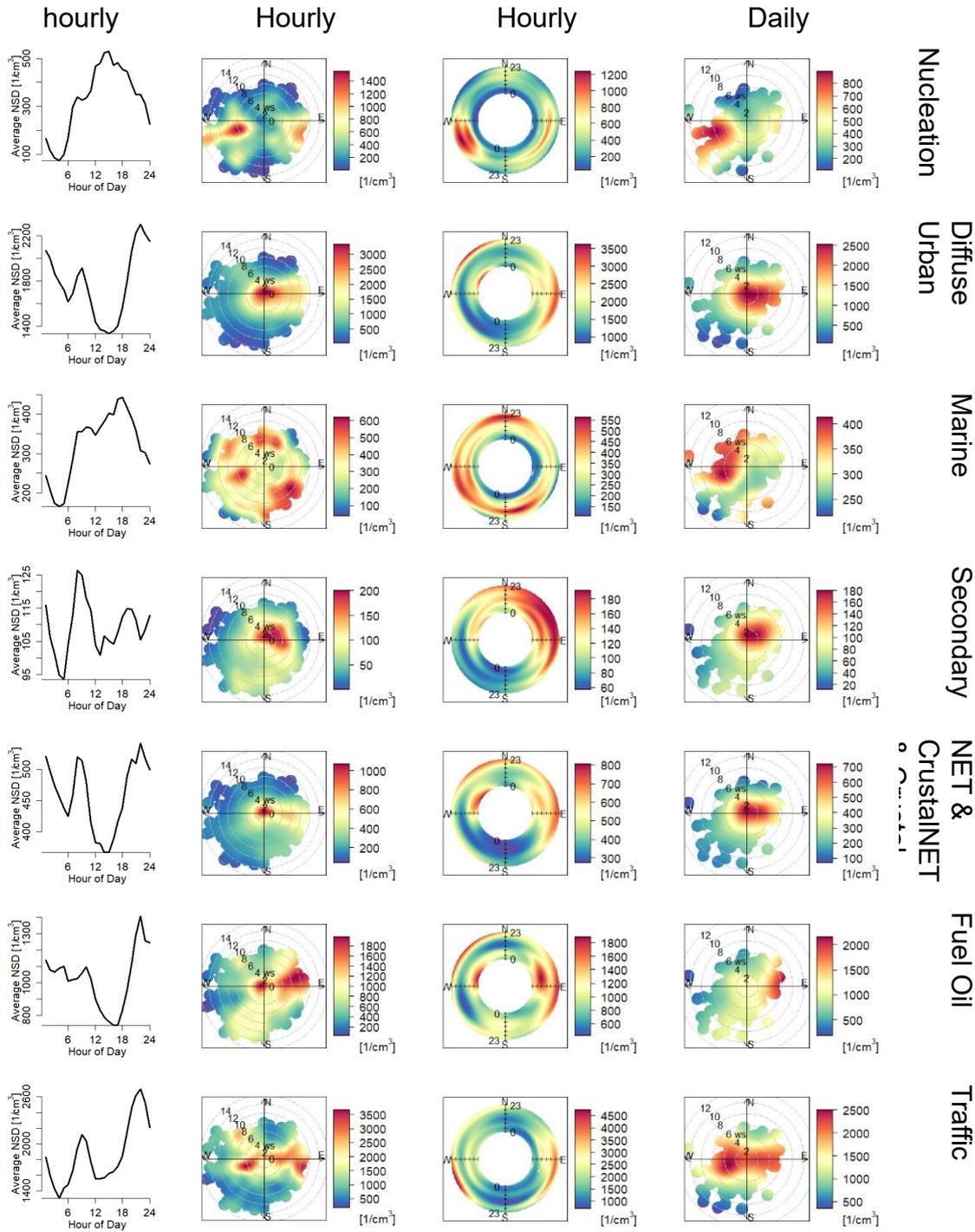
Figure 4. Nucleation and Fuel Oil factors derived when extending the second PMF analysis from the 6 factors (shown in Figure 3) to 7 factors. Source profiles 2F_1 to 2F_6 are given in Figure S3. Each plot is divided into 2 showing the output 1F_k and 2F_k . [Grey bars and black line indicate the values of F; red lines and dots indicate the explained variations; and grey dotted line indicates the $dV/d\log D_p$.]

876

877

878

879



881

882

883

884

885

886

887

888

889

890

Figure 5. Diurnal cycles derived PN_k calculated by the fitting of the daily PMF factor profiles to the hourly NSD data fitted (see equation 8 and Section 2.4). [Left-left column – diurnal trends of PN_k ; left-middle column – bivariate plot of PN_k ; middle-right – annular plot PN_k ; right-right – bivariate plot of PN_k , plotted using the Openair program. Polar plots show a point coloured according to the key, the number concentration at that point on the plot whose distance from the origin represents wind speed and angle wind direction. Likewise for the angular plots the number concentration represent wind direction at an hour of the day between 0 and 23 hrs.]. Note that the diurnal plots do not start at zero.