

1 Chemical sparsity in Bayesian receptor models for aerosol source 2 apportionment

3 Marta Via¹, Jure Demšar², Yufang Hao³, Manousos Manousakas^{3,4}, Anton Rusanen⁵, Jianhui Jiang⁶,
4 Stuart K. Grange⁷, Jean-Luc Jaffrezo⁸, Vy Ngoc Thuy Dinh⁸, Gaëlle Uzu⁸, Griša Močnik¹, and Kaspar
5 R. Daellenbach³

6
7 ¹Center for Atmospheric Research, University of Nova Gorica, Ajdovščina 5270, Slovenia

8 ²Faculty of Computer and Information Science, Tržaška Cesta 25, 1000 Ljubljana, Slovenia

9 ³Laboratory of Atmospheric Chemistry, Paul Scherrer Institute, 5232 Villigen PSI, Switzerland

10 ⁴Environmental Radioactivity Aerosol Tech. for Atmospheric Climate Impacts, INRaSTES, National Centre of Scientific
11 Research “Demokritos”, Ag. Paraskevi, 15310, Greece

12 ⁵Atmospheric Composition Research, Finnish Meteorological Institute, 00101 Helsinki, Finland

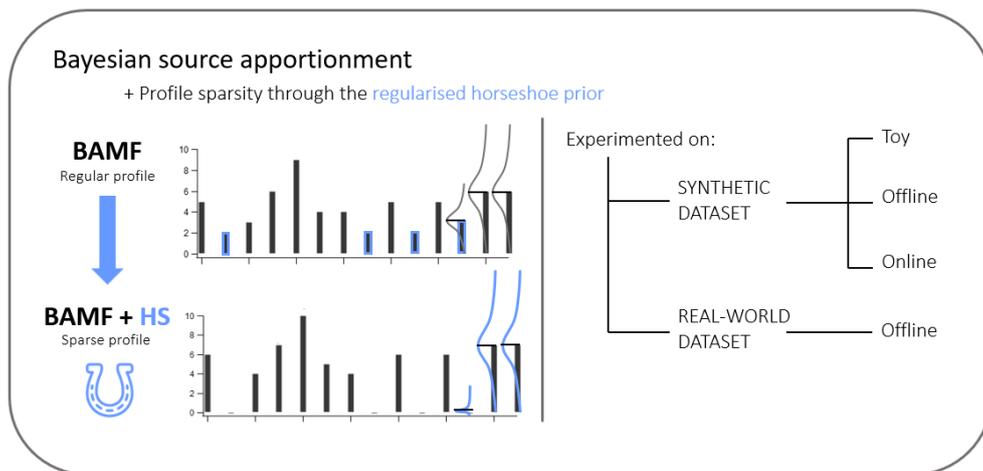
13 ⁶School of Ecological and Environmental Sciences, East China Normal University, 200241, Shanghai, China

14 ⁷Climate and Environmental Physics, Physics Institute, University of Bern, Sidlerstrasse 5, 3012 Bern, Switzerland

15 ⁸University of Grenoble Alpes, CNRS, INRAE, IRD, Grenoble INP, IGE, Grenoble 38000, France

16
17 *Correspondence to:* Marta Via (marta.viagonzalez@ung.si) and Kaspar R. Daellenbach (kaspar.daellenbach@psi.ch)

18 **Abstract.** Aerosol source apportionment is a key tool for understanding the origins of atmospheric particulate matter and for
19 guiding effective air quality management strategies. However, source apportionment techniques still struggle to properly
20 separate highly correlated sources without relying on restrictive *a priori* information, possibly skewing the solution and adding
21 subjective operator input, with varying degrees of benefit. This study introduces sparsity into the Bayesian Autocorrelated
22 Matrix Factorisation (BAMF) model with the aim of removing non-essential species contribution in the unconstrained profiles,
23 which is expected to improve the separation of factors [compared to BAMF](#). The regularised horseshoe prior (HS) has been
24 added to BAMF (BAMF+HS) to promote composition matrix F sparsity, shrinking low-signal contributions to the solutions.
25 BAMF+HS was evaluated using three synthetic datasets designed to reflect increasing levels of data complexity (Toy,
26 [representing a highly simplified dataset](#); Offline, [representing a filter dataset](#); and Online, [representing an Aerosol Chemical
27 Speciation Monitor \(ACSM\)-like dataset](#)), and a real-world multi-site filter dataset. The results demonstrate that BAMF+HS
28 effectively enforces sparsity in offline datasets and that this improves accuracy in reconstructing source profiles and time series
29 compared to BAMF and Positive Matrix Factorisation (PMF). However, its application to higher-complexity ACSM datasets
30 revealed sensitivity to sampling instability hindering sparsification. With that, even though sparsity was not achieved, the
31 quality of the BAMF+HS solution metrics were not deprecated compared to BAMF. Overall, this work underscores the value
32 of incorporating profile sparsity as a solution property in Bayesian source apportionment, and positions BAMF+HS as a
33 promising model for source apportionment.



35

36 **1. Introduction**

37 Particulate Matter (PM) adversely affects human health through both short and sustained exposures (Pope and Dockerty, 1999,
38 Yang et al. 2019). The observed relationship between decreasing PM concentrations and increased life expectancy (Keuken et
39 al. 2011; Zheng et al. 2022) highlights the importance of developing mitigation plans grounded in detailed knowledge of PM
40 sources composition and concentrations. Moreover, because some proxies for aerosol toxicity, among them oxidative potential,
41 are the toxicity of PM can be highly dependent on its sources (Daellenbach et al. 2020), implementing source-specific
42 mitigation measures contributes to more quantitative and efficient abatement and a more effective protection of the population.

43
44 Source apportionment is the process of identifying and quantifying PMF sources by using information about its chemical
45 composition, and is commonly conducted through receptor models (RMs) which differentiate PM sources according to the
46 distinctness of their chemical composition and time series characteristics. The most widely used RM is the Positive Matrix
47 Factorisation model (PMF, Paatero and Tapper, 1994), which deconvolutes/decomposes the input chemical composition into
48 the product of composition and time series matrices (**F** and **G**, respectively), and minimises the residuals of the fit through the
49 weighted least squares loss. The factorisation equation, hence, is written as

$$50 \quad X = G \cdot F + E, \quad (1)$$

51 where X is the input matrix, a $n \times m$ matrix of n timepoints and m species, which is decomposed into G and F matrices of
52 dimensions $n \times p$ and $p \times m$, respectively, where p is the number of factors, and E is the residuals matrix of dimensions $n \times m$.

Formatted: Font color: Blue

Formatted: Right

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue, Highlight

Formatted: Font: Italic, Font color: Blue, Highlight

Formatted: Font color: Blue, Highlight

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

53 where X is the $\mathbb{R}^{n \times m}$ input matrix of n timepoints and m species, which is decomposed into $G \in \mathbb{R}^{n \times p}$, and $F \in \mathbb{R}^{p \times m}$, where
54 p is the number of factors, and $E \in \mathbb{R}^{n \times m}$ is the residuals matrix.

55
56 Unconstrained PMF, although even if it can lead to robust results, is usually insufficient when the sources are highly correlated
57 or have very similar source profiles. In such cases, guiding the model by introducing a priori knowledge (common practice
58 known as constraining the model) has been proven beneficial for the source deconvolution/disentanglement (Lingwall and
59 Christensen et al. 2007, Belis et al. 2014, Dinh et al. 2025). However it can still introduce substantial bias in the solution
60 (Via et al. 2022). A very strongly-constrained RM/The extreme case of RM guidance is the Chemical Mass Balance model
61 (CMB), which factorises the initial matrices with a totally fixed G or F . Globally, the RMs cover the whole range of prior
62 knowledge required (Viana et al. 2009a, Belis et al. 2013).

63
64 Bayesian models represent a probabilistic alternative to the PMF framework. The first application of Bayesian models in
65 atmospheric source apportionment was in introduced in Park et al. (2001, 2002) for Volatile Organic Compounds (VOC)
66 source apportionment. In this approach, the mass closure condition was taken to the Bayesian framework and an autocorrelation
67 prior AR(1) (the first order autoregression) was applied, improving the solution given assuming independent G components.
68 The autocorrelation prior importance was later reinforced in Rusanen et al. (2024) with a differently formulated autocorrelation
69 prior. The latter shows the added value of the Bayesian Autocorrelated Matrix Factorisation model (BAMF) in comparison to
70 PMF in different kinds of spectrometry-based PM synthetic datasets. The Bayesian Multivariate receptor modelling software
71 BNFA and bayesMRM (Park and Oh, 2021) were developed to provide user-friendly tools for Bayesian source apportionment.

72
73 However, studies using the Bayesian Matrix Factorisation (BMF) framework are still scarce. Some examples are Oh and Park
74 (2022), which employed a Bayesian RM to conduct approach multi-site source apportionment, and Zhang et al. (2023), which
75 performed NH_4^+ source apportionment through the Bayesian SIMMR package (Govan et al. 2023). Bayesian models have also
76 been used as a complement to standard RMs, as in Balachandran et al. (2013) where a Bayesian model processing ensemble
77 solutions of a chemical transport model and solutions of three RMs are produced to then use it in CMB for production of final
78 results. The Bayesian model focused on attributing the proper weight to each of the ensemble components and improved the
79 correlation of sources with their markers compared to the traditional approach. Bayesian inference has also been used in Park
80 et al. (2002) and Dai et al. (2024) to generate spatially resolved source apportionment solutions adjusting the weights of each
81 location solution in a multi-site data scheme.

82
83 Thus, Bayesian Matrix Factorisation has become an effective and powerful tool for aerosol source apportionment. However,
84 to the authors knowledge, little attention has been given to improving the accuracy of chemical composition profiles, i.e. F
85 components. This highlights the fundamental challenge in receptor modelling of obtaining chemically distinct and interpretable
86 source profiles from complex and mixed overlapping emissions sources. Moreover, it has been shown in Rusanen et al. (2024)

Formatted: Font color: Blue

87 that in BAMF, slight ~~differences of F differences~~ can severely compromise the quality of **G** (Figure S2 in the mentioned
88 article), hence, steps towards **F** refining should result in overall source apportionment method improvement. ~~In this context,~~
89 ~~sparsity, defined as the property of a dataset, model or solution in which only a limited number of elements are substantial~~
90 ~~contributions while most are zero or close to zero, could be favourable for this problem.~~The enforcement of sparsity, defined
91 ~~as the property of a dataset, model, or solution, where most elements are insignificant or inactive often represented as zeros or~~
92 ~~values near zero, could be favourable for this problem.~~ The accomplishment of sparse source fingerprints could represent
93 “cleaner” emission sources, ~~with less mixing among resolved factor profiles, since substituting non-significant contributions~~
94 ~~in a factor by zeros might allow allocating more importance to the actually relevant contributions of species in factors, less~~
95 ~~entangled among themselves.~~ This work aims to implement sparsity on chemical fingerprints in BAMF aiming for a more
96 accurate source apportionment. We introduce sparsity with the regularised horseshoe prior (Piironen and Vehtari, 2017), which
97 unlike other sparsity priors, enables regularisation of the sparsity strength, and compare it with other sparsity priors, such as
98 Lasso (Tibshirani et al. 2015) and Spike-and-slab (Andersen et al. 2014). This model is then tested on three synthetic datasets
99 with different complexity degrees and one real-world dataset to depict the impact of sparsity and potential benefits of its
100 implementation.

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

101 2. Methodology

102 2.1 Bayesian Matrix Factorisation

103
104 Bayesian Matrix factorisation models, like other RMs, are based on the chemical mass balance equation (Eq. 1(a)). Bayesian
105 modeling approaches this problem probabilistically and bases the determination of the matrices, **F** and **G**, the main parameters
106 to determine, upon the assumptions imposed on the model, i.e. priors. Bayesian factorisation forces the decomposition through
107 modelling the **X** matrix components as a Gaussian with center on the “noise-free data matrix” **Z** (matrix of dimensions p, m) Z
108 ($\in \mathbb{R}^{n \times m}$) and a standard deviation given by the positively-defined uncertainty matrix σ (positive matrix of dimensions $p, m \in$
109 $\mathbb{R}_{\geq 0}^{n \times m}$) (Eq. 2)-(2). The matrix **Z** is, in turn, the product of the time series and profiles submatrices, **G** and **F**, respectively,
110 and (a) is rewritten as:

$$111 X \sim N(Z, \sigma) = N(G \cdot F, \sigma) \quad (2)$$

112 where N represents the normal distribution. Whilst **G** is not given any prior meaning and is sampled then by default from a
113 uniform distribution, **F** is modelled as a Dirichlet distribution to ensure positivity, with the sum of its components being equal
114 to 1 (2):

$$115 F_k \sim \text{Dirichlet}(1_m) \quad (3)$$

116 With these **F** requirements, profiles represent the normalised contribution to the spectra of one source. Usual notation for
117 indices used hereinafter here are i, j, k for elements in the range $(1, \dots, n)$, $(1, \dots, m)$, and $(1, \dots, p)$, for the timestamps, species,
118 and factors, respectively. ~~It is worth noting that PMF applies the normalisation of profiles after a **F**, **G**, solution is found, not~~

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue, Highlight

Formatted: Font color: Blue, Highlight

Formatted: Font: Italic, Font color: Blue, Highlight

Formatted: Font: Bold, Font color: Blue

Formatted: Font: (Default) Cambria Math

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue, Highlight

Formatted: Font color: Blue, Highlight

Formatted: Font: Italic, Font color: Blue, Highlight

Formatted: Font color: Blue

Formatted: Right

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font: (Default) Cambria Math, Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

as a model prior as done in BAMF. The PMF generates mass-loaded \mathbf{F}_* , \mathbf{G}_* solution matrices, which are reweighted to provide a normalised \mathbf{F}_* and a mass-loaded \mathbf{G}_* . In the Bayesian models used in this study, the normalisation of \mathbf{F}_* is inherent to the model by design. The different formulations eventually provide normalised \mathbf{F}_* and mass-weighted \mathbf{G}_* with unlikely affectations due to the normalisation procedure. The model configuration given by (2) and (3) will be referred to as Bayesian Factorisation model (BMF) and represents the analog of PMF in the Bayesian framework. All models used in this manuscript are outlined in Table 1.

On top of this structure, Rusanen et al (2024) proposed an autocorrelation prior for \mathbf{G} which should account for the inherent autocorrelation of air pollutant sources in time. The formulation of the autocorrelation prior for \mathbf{G} is given by (4) and includes two more modelling parameters, α (positive vector or dimension p) ($\in \mathbb{R}_{>0}^p$) and β (positive vector or dimension p) ($\in \mathbb{R}_{>0}^p$), which regulate the similarity of one \mathbf{G} component with the previous one as follows:

$$G_{i+1,k} \sim C^+(G_{ik}, \alpha_k \cdot \Delta t_i + \beta_k) \quad (4)$$

where $i \in \{2, \dots, n-1\}$, C represents the Cauchy distribution and $+$ represents positive real numbers $\mathbb{R}_{>0}$. This prior centers the $i+1$ 'th component distribution in the i th component with a distribution width that linearly depends on the temporal gap between these two timestamps. Hence, the more temporally-separated two consecutive points are, the less correlated they are expected to be. The Cauchy distribution was chosen due to its heavier tails which enable more probable jumps between consecutive i 's than a Gaussian distribution (Gelman et al., 2013). This flexibility could be convenient for real-world datasets which are affected by measurement gaps. The coefficients α and β are source dependent to allow for source-dependent correlation degrees. The model which introduces this prior to BMF is called Bayesian Autocorrelation Matrix Factorisation model (BAMF, Rusanen et al. 2024).

2.1.1 The horseshoe prior

The introduction of sparsity in BAMF involves entails the addition of several hyperpriors in the \mathbf{F} prior to implement the shrinkage mechanism. In this study, we used the regularised horseshoe prior (Piironen and Vehtari, 2017), which is a global-local complex of hyperpriors, i.e. the shrinkage power is both regulated globally source-wise and \mathbf{F} -component-wise. The idea behind this prior is that species with very small contributions to a factor are shrunk toward zero through an automatic shrinkage mechanism, whereas species with substantial support from the data are largely unaffected. The regularised horseshoe (HS) prior implemented in \mathbf{F} in the BAMF scheme as

$$F_{kj} = \mu_{kj} \cdot \tilde{\lambda}_{kj} \cdot \tau_j, \quad (5)$$

where μ (matrix of dimensions $\mu \cdot m$) ($\in \mathbb{R}_{>0}^{n \cdot m}$) represents the \mathbf{F} matrix without the horseshoe prior. μ , in turn is defined as a standard Cauchy distribution prior as

$$\mu_{kj} \sim C^+(0, 1), \quad (6)$$

where τ (vector of dimension p) ($\in \mathbb{R}_{>0}^m$) represents the global shrinkage parameters and

- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Bold, Font color: Blue
- Formatted: Font color: Blue
- Formatted: Font: Not Bold, Not Italic, Font color: Blue
- Formatted: Font: Not Bold, Font color: Blue
- Formatted: Font: Not Bold, Not Italic, Font color: Blue
- Formatted: Font: Not Bold, Not Italic, Font color: Blue
- Formatted: Font: Not Bold, Font color: Blue
- Formatted: Font: Not Bold, Not Italic, Font color: Blue
- Formatted: Right
- Formatted: Font color: Blue

- Formatted: Font color: Blue
- Formatted: Right
- Formatted: Font color: Blue
- Formatted: Font: Italic, Font color: Blue
- Formatted: Right
- Formatted: Font: Not Bold
- Formatted: Font: Not Bold, Font color: Blue
- Formatted: Font: Not Bold, Italic, Font color: Blue
- Formatted: Font: Not Bold

$$\tau \sim C^+(0, \tau_0 \cdot \sigma_{HS}), \quad (7)$$

where the parameter τ_0 can be regulated by the user to regulate the overall shrinkage power and σ_{HS} is sampled from an uniform distribution. The hyperparameter $\tilde{\lambda}_{kj}$ applies the local shrinkage to λ ($\in \mathbb{R}^{n \cdot m}$) as

$$\tilde{\lambda}_{kj} = \sqrt{\frac{c^2 \cdot \lambda^2}{c^2 + \lambda^2 \cdot \tau^2}}, \quad (8)$$

where

$$c^2 \sim \Gamma^{-1}(0.5 \cdot \text{slab_df}, 0.5 \cdot \text{slab_df}) \quad (9)$$

$$\lambda \sim \text{slab_scale} \cdot C^+(0, 1) \quad (10)$$

both combined providing the characteristic shrinking horseshoe shape. Here, λ is a model parameter of dimensions n, m which after regularisation becomes is denoted as $\tilde{\lambda}$. Further description of the horseshoe implementation on BAMF can be found in Section S2.1, and the prior derivation and details, in Piironen and Vehtari (2018). The distribution parameters τ_0 , σ_{HS} , and slab_df , slab_scale were tested and results did not show significant sensitivity to their variations, so we keep the defaults as provided in Piironen and Vehtari (2018) as can be found in the available shared codes. The models with the horseshoe (HS) priors are hereinafter marked with "+HS". Figure S1 shows a schematic diagram of the matrix decomposition through BAMF.

In order to assess the amount of sparsity of a dataset or a solution, we used the Gini coefficient (Gini et al., 1936), which assesses the inequality over a distribution as follows:

$$\text{Gini} = \frac{\sum_{i=1}^n (2i - n - 1) \cdot x_i}{2 \cdot n \cdot \sum_{i=1}^n x_i} \quad (11)$$

where x values are sorted in ascending order and n is the number of elements in x . Since it quantifies the inequality/unevenness, it can be a proxy for sparsity: if some values are high and the others are zero, $\text{Gini} \approx 1$, if all values are equal, $\text{Gini} = 0$. Also, the solution-to-truth Gini values ratio will be discussed throughout the analysis, referred to as "Gini ratio". To evaluate if the sparsity is enforced precisely where it should, an additional metric has been applied called "zero truth sum". This metric sums up the modelled contributions of the null species in the truth profiles.

2.1.2 Alternative factorisation methodologies

BAMF-AR1. There is an alternative formulation for the autocorrelation prior as introduced in Bayesian models by Park et al. (2001). The AR(1) autocorrelation prior is the first degree polynomial expansion of the autoregressive models and it proposes a linear progression of $G_{i+1,k}$ from $G_{i,k}$. We introduce AR(1) in the Bayesian framework as

$$G_{i+1,k} \sim N(\alpha_k \cdot G_{ik} + \beta_k, \gamma_k) \quad (12)$$

Formatted: Right

Formatted: Right

Formatted: Right

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font: Italic, Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Right

181 In this formulation, the $i+1$ -th point stems from a Gaussian distribution centered linear combination on the i -th point with
182 source-dependent slopes (α) and intercept (β), and width γ . Although, unlike (4), it disregards the decrease of correlation
183 between gapped consecutive points, this prior allows for source-specific time series trends, which would be beneficial for
184 certain source description. The model which introduces this prior to BMF will be called BAMF-AR1.

185

186 **BAMF-GS.** Another formulation is introduced, switching the Dirichlet distribution to the matrix \mathbf{G} instead of \mathbf{F} (6). This
187 swap should allow \mathbf{F} to retain the \mathbf{X} matrix mass and could potentially help [deconvolutingdisentangling](#) profiles due to the
188 upweighting of the chemical profiles.

$$189 G_k \sim \text{Dirichlet}(\mathbf{1}_n) \quad (13)$$

190 Thus, the \mathbf{G} presents two priors, the dirichlet distribution and the autocorrelation prior, whilst \mathbf{F} is sampled from the default
191 uniform distribution. This model will be called hereinafter BAMF-GS as short from BAMF-G simplex, since a simplex is the
192 set of positive vectors that sum to one hence it is the natural geometric structure for the Dirichlet distribution to sit on. This
193 model structure, nevertheless, does not allow for a horseshoe prior application, since due to the factors mass now incorporated
194 in \mathbf{F} , the coefficients will be very distinct from zero and the horseshoe prior will not perceive them as potential signals to
195 sparsify.

196

197 **CMB.** Lastly, a Bayesian formulation of the CMB model was employed in order to test the horseshoe prior capacities with the
198 most proper factorisation possible. This model is mainly analogous to CMB in the Bayesian framework, but the \mathbf{G} matrix was
199 fixed with the truth time series. Hence the model only had to determine the \mathbf{F} components distributions to match the
200 factorisation condition (2) given the truth \mathbf{G} .

201 2.1.3 Solver and Hamiltonian-Monte Carlo Markov Chain

202 All Bayesian models were compiled and run in STAN (Carpenter et al. 2017), a probabilistic programming language developed
203 for Bayesian modelling. STAN solves Bayesian inference through the Hamiltonian Monte Carlo (HMC) algorithm based on
204 Markov Chain Monte Carlo methods (MCMC). HMC uses an approximate Hamiltonian dynamics simulation with the
205 Metropolis acceptance/rejection criterion and a no-U-turn sampler (NUTS, Hoffman and Gelman, 2014). For the sake of
206 brevity, we present only the essential concepts here, directing readers to Carpenter et al. (2017), Gelman et al. (2014), STAN
207 Manual (2025) and references therein for comprehensive information.

208 The parameters of the model, primarily \mathbf{F} and \mathbf{G} but also all the other defined hyperparameters ($\tau, \lambda, \alpha, \beta$), are sampled from
209 their posterior distributions, constructed from the priors and the data introduced. In the Hamiltonian analogy, the evolution of
210 these parameters across samples is computed as the trajectory of a fictitious particle. This particle moves through the parameter
211 space driven by random momentum in all directions. [This approach avoids the random-walk behavior of simpler sampling
212 methods and enables faster convergence.](#) The trajectory is hence simulated using a discretized approximation, and candidate

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

positions are accepted or rejected according to the Metropolis criterion (Metropolis et al. 1953). Accepted positions correspond to plausible parameter values given both the model assumptions and the data. This process provides a distribution over samples of possible solutions from which confidence intervals of each of the model (hyper)parameters can be extracted. A set of samples is called a chain, each of them initialised with a different seed in order to explore the solution space more broadly. In order to initialise the model parameters more effectively, the maximum a posteriori (MAP) point parameters solution estimated by STAN is used through the LBFGS algorithm (Liu and Nocedal, 1989). Even if this approach makes the parameter sampling process much more efficient, solutions might have multiple local maxima and MAP will initialise the models based only on one of those. This highlights the importance of using different seeds to explore the solution space more widely. Since the early iterations of each Markov chain are typically influenced by the starting values and may not represent samples from the true posterior distribution, we discarded the first half of the samples from each chain. Different settings were used according to the type of experiment, shown in Table S1. The number of chains is consistent with standard practice in Bayesian modeling, and the number of samples was increased beyond commonly adopted values (e.g., 1000) in order to improve solution stability. As seen in Table S1, the more complex the datasets are, the more time BAMF+HS takes to run. Since the BAMF+HS running times are high at this development stage, BAMF+HS might currently be more adequate for exhaustive source apportionment refinement than real-time monitoring.

Formatted: Font color: Blue

In order to evaluate the convergence of a solution to the target posterior distribution, the potential scale reduction factor (\hat{R} , Gelman and Rubin, 1992) is used. This coefficient compares the variance within chains and between chains of the \mathbf{Z} matrix, hence if chains converge, $R=1$, and values of $\hat{R} \gg 1$ imply chain divergence and values of $\hat{R} \ll 1$ imply sampling divergence in chains. The convergence of all runs has been assessed using standard Bayesian diagnostics, including visual inspection of trace plots, and the effective sample size and \hat{R} statistics, and all experiments shown in the manuscript fall within satisfactory stability ranges for these criteria.

Formatted: Font color: Blue

2.2 PMF

In PMF (Paatero and Tapper, 1994), equation (1) is solved through the ME-2 solver (ME2, Paatero, 1999) based on the weighted means squares minimization of the quantity:

$$Q = \sum_{j=1}^m \sum_{i=1}^n \frac{e_{ij}^2}{\sigma_{ij}^2} \quad (12)$$

PMF was implemented on all datasets unconstrainedly through the Source Finder software (SoFi version 9.5, Canonaco et al. 2013) with 100 sorted runs which are posteriorly sorted as in BAMF. The number of runs may seem compromising the PMF quality in comparison to the 6000 samples per chain used in Bayesian models. However, this comparison is misleading, since the factorisation space is indeed better explored by PMF, with 100 different sampling seeds, while only 4 seeds (chains) were used in BAMF-like models as usual procedure in Bayesian modelling for the sake of computational resources.

244 **2.3 Pre- and post- processing for all models**

245 ~~Before~~Previously to model running, \mathbf{X} and σ are normalised to use consistent scales of all priors and posteriors for Bayesian
246 models. The normalisation is based on ensuring a mean of $\mathbf{X} = 1$.

247 $X^* = X / f_{norm}$, $\sigma^* = \sigma / f_{norm}$ where $f_{norm} = \sum_{i,j} X_{ij} / (n \cdot m)$
248 (13)

249 After the factorisation this normalisation is reverted converting the normalised matrices, hereinafter referred as \underline{G} , \underline{E} , to the
250 properly-scaled \mathbf{G} , \mathbf{F} matrices.

251
252 The factor ordering in the matrices is random in the model results, hence, ~~the solution factors~~the sorting process is, hence, the
253 ~~solutions factors~~ must be sorted. Here, as in Rusanen et al. (2024), we used the Hungarian algorithm (Kuhn, 1955) to sort the
254 \underline{Z}_k components ($\underline{Z}_k = \underline{G}_{ik} \cdot \underline{F}_{kj}$, i.e. each factor's normalised \mathbf{Z} submatrix). The metric to sort the components is the Manhattan
255 distance (i.e. the sum of the absolute differences of two Cartesian coordinates). All factors in each chain of samples are then
256 reordered upon the factor order of a small group of samples of that chain (the last 5, arbitrarily chosen) and, subsequently, one
257 all-samples-averaged F_k and G_k are retrieved for each of the chains. Then, the order of factors of each of the chains is sorted
258 again in the same way in relation to the truth F , to have all sources equally sorted in all chains. Median and quantiles are
259 computed over samples and chains to produce the final solutions and uncertainties. This sorting process is also used for the
260 PMF solution despite not being its usual sorting approach for the sake of homogeneity in comparison to the Bayesian models.

261
262 The last step of the experimental process ~~was~~ to assess the model performance on the given dataset. The evaluation of the
263 performance should be based on: i. reconstruction performance, or the difference between \mathbf{X} and \mathbf{Z} ; ii. similarity to truth, or
264 environmental sensibility based on the apportionment of source tracers in case the truth is not available; iii. computational
265 performance. The reconstruction performance ~~was~~will be assessed by checking the cell-wise correlation between \mathbf{X} and \mathbf{Z}
266 and checking the median and maximum of the absolute value of relative deviations of \mathbf{Z} and \mathbf{X} with respect to the measurement
267 uncertainty matrix σ ($(\mathbf{X}-\mathbf{Z})/\sigma$). The similarity to truth, when available, is tackled by comparing the median ratio between
268 modelled \mathbf{G} and truth $(\underline{G}/\underline{G}_0)$, the Pearson correlation for the \mathbf{G} matrix (\underline{G}, r) and the Spearman correlation for \mathbf{F} amongst
269 models (\underline{F}, ρ) . The Spearman correlation coefficient for the factor profiles was chosen due to the expected non-linearity of the
270 comparison and likely presence of outliers. These comparisons, and especially when the ground truth is not available, need to
271 be accompanied by visual inspection of the solution quality, looking for resemblance with known environmental sources. The
272 models accounting for sparsity will be also compared upon the aforementioned Gini metric and, when truth is available, the
273 Gini ratio with truth and the “zero-truth sum”. ~~Computational performance assessment will be based on the metrics of~~
274 ~~convergence metrics of the Hamiltonian-Montecarlo Markov chain methods embedded in STAN software (e.g. \hat{R}).~~

- Formatted: Font: Bold
- Formatted: Font: Bold
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font: Bold, Font color: Custom Color(79;129;189)
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font: Bold, Font color: Custom Color(79;129;189)
- Formatted: Font color: Custom Color(79;129;189), Subscript
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font: Bold, Font color: Custom Color(79;129;189)
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font: Bold, Font color: Custom Color(79;129;189)
- Formatted: Font color: Custom Color(79;129;189)
- Formatted: Font color: Blue

276 2.4 Datasets

277 The datasets created for model experimentation can be divided into synthetic and real-world datasets. Synthetic datasets are
278 artificially created with the purpose of knowing the \mathbf{F} , \mathbf{G} , to test model accuracy retrieving these matrices with respect to the
279 *truth* and these have been widely used for source apportionment validation in the last decades (Park et al. 2002, Brinkman et
280 al. 2006; Belis et al. 2015; Via et al. 2022; Rusanen et al. 2024). In order to challenge the models gradually, we created three
281 synthetic datasets with increasing degrees of complexities (toy, offline, online ACSM synthetic datasets). Additionally, a real-
282 world chemically sparse dataset was also used to test the results. AlthoughDespite the truth's factorisation is unknown and
283 results cannot be directly verified, the model's factorisation can be assessed environmentally or based on indicators on the
284 goodness of the fit. The different datasets have different levels of sparsity, as can be seen in Table 2, that the models with the
285 horseshoe prior should aim to replicate. The time resolution of modelled OA sources, used both in the chemically-sparse toy
286 dataset and the chemically less sparse datasets, is 1 hour. The time resolution of offline datasets, used in the chemically sparse
287 synthetic offline dataset and the chemically-sparse real-world offline dataset is 1 day.

288 2.3.1 2.4.1 Chemically-sparse toy dataset

289 A simplistic synthetic toy dataset was designed as a deliberately simplified test case to perform basic control and performance
290 tests. rather than to reproduce any realistic atmospheric scenario. It was devised by creating three very simple and sparse
291 profiles and using three time series (HOA, SOAbio, BBOA) from modelled source time series of the city of Zurich (Rusanen
292 et al. 2024) in order to test how sparsity priors act on very uneven species contribution. Although it is based on ACSM-like
293 time series and therefore reflects some of the temporal properties of such measurements, the three included sources do not
294 represent combinations that would be expected in a real-world environment since this toy dataset is intended solely for
295 methodological testing purposes. In addition, the source profiles were intentionally designed to be highly simplified in order
296 to facilitate an immediate visual assessment of the model fitting. For these reasons, the extracted components were not assigned
297 environmental labels, but were instead referred to generically as Factor 1, Factor 2, and Factor 3.

298 Then, \mathbf{F} and \mathbf{G} were multiplied to generate \mathbf{Z} , and some gaussian error with standard deviation σ was added to each component
299 to generate a realistic \mathbf{X} matrix. The uncertainties matrix σ was designed as a sixth of the X values plus Gaussian noise. With
300 this arrangement, the models can be applied conventionally to the \mathbf{X} , σ matrices and the modelled \mathbf{F} and \mathbf{G} , can be compared
301 to the original truth \mathbf{F} , \mathbf{G} , which will be referred hereinafter as \mathbf{F}_0 and \mathbf{G}_0 , displayed in Figure S2.

302 2.3.2 2.4.2 Chemically-sparse synthetic offline dataset

303 We created a synthetic offline filter dataset, mimicking the filter-based measurements input matrices, in order to test the
304 accuracy of the modelsmodels-accuracy in these kinds of datasets. This dataset mimics the concentrations on the coarse fraction
305 ($\text{PM}_{10} - \text{PM}_{2.5}$) as collected by a high-volume sampler on the Zurich-Kaserne site (Grange et al. 2021) including the following
306 chemical species: OC, Al, Na, Mg, Cl, K, Ca, S, Fe, Cu, Zn, Mn, Sb, Ba, mannitol, arabitol. In the original real-world dataset,

Formatted: Font color: Blue

Formatted: Font: (Default) Arial, 11 pt, Not Bold, Font color: Black

Formatted: No bullets or numbering

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: No bullets or numbering

Formatted: Font: (Default) Arial, 11 pt, Not Bold, Font color: Black

307 data obtained with two series of samples (PM_{10} and $PM_{2.5}$) were subtracted in order to focus on the coarse source
308 apportionment, since the main emission sources of these elements and organic species stem from mechanical processes leading
309 to major coarse models. It was created by crafting first the **F** and **G**, then multiply them and creating **X** and σ . The **F** matrix
310 was slightly modified from that proposed in Manousakas et al. (2025), making the chemical profiles slightly sparse by zeroing
311 the non-relevant species in each of the factors (dust, traffic, salt, coarse biological). The **G** matrix was composed of the time
312 series of:

- 313 - Dust: modelled PM_{10} dust (Vasilakos et al. in prep.) converted to coarse with the $Al_{PM_{10}}$ vs. PM_{10} ratio from Grange
314 et al. (2021).
- 315 - Traffic: modelled PM_{10} copper (Upadhyay et al. 2025) converted to coarse with the $Cu_{PM_{10}}$ vs. PM_{10} ratio from
316 Grange et al. (2021).
- 317 - Salt: coarse Na+Cl (Grange et al. 2021) converted to PM concentrations and multiplied by an arbitrary number (3 in
318 this case match the concentrations of the sea salt factor in the original dataset).
- 319 - Coarse biological: coarse Arabitol+Mannitol (Grange et al. 2021) converted to PM concentrations and multiplied by
320 3, similarly as for the salt factor.

321 This dataset will be called “offline synthetic dataset”. Another more simplistic dataset was prepared similarly but using Al and
322 Cu for dust and traffic factors, respectively, in the same way as in the salt or coarse biological factors, i.e. omitting the use of
323 modelled data. This dataset will be hereinafter named “Purely-measurement-based offline synthetic dataset” and its modelling
324 results will be described in section 3.2.

325 Once the **F** and **G** matrices were created, **X** was calculated by their multiplication and the addition of Gaussian noise with
326 amplitude σ . The uncertainties matrix σ was generated as in Grange et al. (2021) multiplied by 2 to balance the signal-to-noise
327 ratio to the datasets in Manousakas et al. (2025). The matrices **F**, **G** of this dataset are displayed in Figure S3.

328

329 ~~2.3.3~~ 2.4.3 Chemically sparse real-world offline dataset

330 A real-world dataset was employed to test the current models applicability in campaign measurements. This dataset was
331 originally used for source apportionment in Manousakas et al. (2025) and Grange et al. (2021) and consists of $PM_{10-2.5}$ samples
332 at five Swiss National Air Pollution Monitoring Network (NABEL): Basel, Bern, Magadino, Payerne, and Zurich. The
333 measurements were taken in the June 2018 - July 2019 period every fourth day and using Digitel high-volume samplers. During
334 the sampling campaign PM_{10} and $PM_{2.5}$ were collected and the respective concentrations were subtracted to generate the coarse
335 ($PM_{10-2.5}$) concentrations. These samples include: i. OC concentrations, measured through the thermal optical transmission
336 (TOT) EN16909 method with the EUSAAR2 temperature protocol; ii. elemental concentrations (Al, Fe, Cu, Zn, Mn, Sb, Ba,
337 Sr, Bi, Pb) measured by inductively coupled plasma atomic emission spectrometry (ICP-AES) and inductively coupled plasma
338 mass spectrometry (ICP-MS); iii. water soluble inorganic ion concentrations (Ca^+ , Cl^+ , Mg^+ , K^+ , Na^+), determined by ion
339 chromatography (IC); iv. Organic species (mannitol, arabitol) determined by a high-performance liquid chromatographic

Formatted: No bullets or numbering

Formatted: Font: (Default) Arial, 11 pt, Not Bold, Font color: Black

340 method followed by pulsed amperometric detection (HPLC-PAD). The uncertainties of these species were calculated as in
341 Grange et al. (2021).

342 ~~2.3.4~~ 2.4.4 Chemically less sparse synthetic online ACSM datasets

343 With the aim of recreating more complex real-world datasets to test the models, we generated 6 datasets for four European
344 cities: Krakow, Milan, Paris, and Zurich. The objective ~~was~~ to recreate OA matrices as given by a mass spectrometer
345 instrument like Q-ACSM, for which there are plenty of real-world source apportionment studies in the literature. The **G** matrix
346 was created from OA sources time series generated through the regional air quality model CAMx (Comprehensive Air Quality
347 Model with Extensions) as previously published by Jiang et al. (2019). The five sources of these datasets were hydrocarbon-
348 like OA (HOA), related to traffic emissions, biomass burning OA (BBOA), biogenic SOA (SOA_{bio}), biomass burning SOA
349 (SOA_{bb}), and traffic SOA (SOA_{tr}). To ensure seasonal representativity while keeping computational costs low, datasets
350 included the first two weeks of every second month of 2011 (January, March, ...). The relative concentrations of these datasets
351 are shown in Figure S5. This figure shows the highest seasonal OA variation for the city of Milan and the lowest for Zurich.
352 In terms of sources, the most seasonally stable sources, overall, are HOA and SOA_{tr} in contrast to the remarkable variability
353 of BBOA and SOA_{bio}. The profiles used to create the species matrix **F** were those in Table S2 for primary sources (HOA,
354 BBOA). For secondary sources, the profiles from the European megacity dataset presented in Rusanen et al. (2024) were used
355 for the Zurich city, which were slightly perturbed for the other cities due to the limited availability of these sources' profiles
356 in the literature.

357
358 The **X** matrix was obtained by multiplying the **F** and **G** submatrices and adding Gaussian noise. The procedure to calculate the
359 error matrix for such datasets is described in Via et al. (2022) and the dataset used to calculate the error matrix is that from the
360 Zurich site, which ranges from February 2011 until December 2011.

361
362 Lastly, a sensitivity analysis was carried out by ~~slightly~~ modifying the original **F**, **G** matrices slightly upon which the **X**, **σ**
363 matrices were subsequently created. The first Zurich dataset (period 01/09/2011 - 14/09/2011) was used for this purpose and
364 we chose to perturbate one factor only (HOA). The **F**, **G** submatrices were perturbed independently upon the expression:

$$365 \quad G_{HOA}' = G_{HOA} \cdot N(1, \sigma') \quad F_{HOA}' = F_{HOA} \cdot N(1, \sigma') \quad (14)$$

366 where we used $\sigma' = [0, 0.1, 0.2, 0.3, 0.4, 0.5]$ to create different degrees of perturbation. The profiles in **F** were normalised
367 after that process. It must be noted that the perturbation is more relevant on **F** than in **G** since a given σ' in the aforementioned
368 range is more comparable and impactful on the profile contributions, bounded to 1, than on the unbounded time series
369 timepoints. Consequently, ~~with~~ in this framework, we obtained 6 **G**-perturbed and 6 **F**-perturbed input matrices. Both BAMF
370 and BAMF+HS models were run with all these input matrices and their subsequent HOA results were compared to the original
371 truth in order to ~~comprehend~~ ~~grasp~~ the sensitivity of the models upon time series and profile perturbations.

Formatted: No bullets or numbering

Formatted: Font: (Default) Arial, 11 pt, Not Bold, Font color: Black

372 3. Results

373 3.1 Chemically sparse synthetic toy dataset

374 Here, we introduce the evaluated models relying on unrealistically simplified toy data with the purpose of showcasing the
375 performance of the horseshoe prior introduction to BAMF (Figure S2) and the alternative factorisation methodologies, which
376 are discussed in SI Section C.1.

377
378 In the first evaluation step, we assess the performance of the horseshoe prior under the assumption that the source matrix \mathbf{G} is
379 known, in order to isolate its effect on the estimation of \mathbf{F} . Figure 1 shows the distribution of each \mathbf{F} component for CMB with
380 and without the horseshoe prior (CMB, CMB+HS, respectively, Table 1). The distributions shown account for all the
381 variability across samples of each \mathbf{F} component for both models, and the truth is shown as a marker in the x-axis since it is a
382 point value to be compared to the centers of the distributions. The presentation of the CMB and CMB+HS distributions aims
383 to demonstrate the sparsity-inducing role of the horseshoe prior, which enforces shrinkage of the \mathbf{F} component toward zero;
384 this effect is more readily discernible when a strongly guided \mathbf{G} matrix is used to isolate the evidence of sparsity. Figure 1
385 showcases the horseshoe prior power to generate sparsity in \mathbf{F} components, shrinking more strongly the lowest signals to zero
386 than CMB and as a consequence, enlarging the most prominent signals. Table 3 shows how the Gini metric is consistently
387 higher for CMB+HS with respect to CMB, supported by a higher Gini ratio and lower zero truth metric reflecting the
388 sparsification of profiles and higher similarity to truth. The RMSE compared to the truth for the profiles improved with the
389 horseshoe prior applied for all three factors (for CMB and CMB+HS, respectively: 1.2e-04, 3.8e-05 for F1; 1.86e-04, 5e-05
390 for F2; 3.3e-05, 1e-05 for F3). Hence, the sparsity introduced in \mathbf{F} through the regularised horseshoe prior successfully
391 improved the profile description of the solution.

392
393 In the next evaluation step, we test the various models assuming no prior knowledge. Figure 2 shows the results of PMF,
394 BAMF, and BAMF+HS models on the toy dataset and Table 3 shows their factorisation performance and comparison to truth
395 metrics/performances. In terms of factorisation, median relative errors are better for BAMF+HS and BAMF than for PMF, but
396 their maximum errors are higher and the Pearson coefficients slightly lower, all this entailing comparable factorisation
397 performances. All models generally adapt well to the truth features, but they present non-negligible differences. PMF results
398 better resemble the truth in terms of \mathbf{G} R^2 , but it is the model whose G/G_0 differs from 1 the most, accumulating the greatest
399 error (2.64), followed by BAMF (2.10), while BAMF+HS exhibits the smallest deviation (0.81), indicating the highest overall
400 accuracy, factors, both in concentrations and feature adaptation. However, H In terms of profiles, the BAMF+HS model is the
401 closest to the truth both in terms of ρ and R^2 , especially for the second and third factors for which the sparsity introduction
402 results are advantageous with respect to BAMF results. Consistently, the Gini ratios of the inferred solutions relative to the
403 truth are markedly closer to unity for BAMF+HS (range 0.40–0.93) than for PMF (0.45–0.64). This is confirmed by the Gini
404 ratio with truth which is the largest for BAMF+HS and with the lower truth zero metric, indicating a higher \mathbf{F} similarity to

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font: (Default) Times New Roman, Font color: Blue

Formatted: Font: (Default) Times New Roman, Font color: Blue, Superscript

Formatted: Font: (Default) Times New Roman, Font color: Blue

Formatted: Font: (Default) Times New Roman, Font color: Blue, Subscript

Formatted: Font: (Default) Times New Roman, Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue, Superscript

Formatted: Font color: Blue

405 truth. The sparsity effects can also be seen in Figure S5, in which the horseshoe shrinkage is evident for the low m/zs allowing
406 in turn the larger m/zs to retain more mass, hence resembling better the truth profiles. Taken together, these results indicate
407 that BAMF+HS not only promotes sparsity, but does so in a chemically consistent manner, leading to a more accurate mass
408 apportionment across factors, despite a slightly reduced time-series correlation for the third factor. However, the BAMF+HS
409 could not shrink down the lowest signals in Factor 1, likely because their contribution estimated by the mass balance and the
410 autocorrelation restrictions of this model made it unclear for the horseshoe to shrink them down completely. With this result,
411 this toy dataset depicts the capacities and limitations of the horseshoe implementation on BAMF: it is capable to sparsify
412 effectively only the signals which are close enough to zero as given by the restrictions of the BAMF model.

Formatted: Font color: Blue

414 While other sparsity priors exist (e.g. Lasso and Spike-and-slab priors (Figure S6, Table 3, Table S3)), our tests show that the
415 BAMF+HS model is most effective in shrinking unnecessary contributors to F. Hence this prior will be used onwards. This is
416 evidently portrayed by the Gini ratio, for which neither Lasso nor Spike-Slab achieve the signal shrinkage that the BAMF+HS
417 does. Also, neither BAMF+Lasso nor BAMF+Spike-and-slab managed to sparsify the first factor. Additionally, different
418 autocorrelation formulations were implemented ~~tried~~ with and without the horseshoe prior, showing worse performance than
419 BAMF or BAMF+HS, respectively, as discussed in section SI C.1. This supports using the BAMF autocorrelation prior instead
420 of the alternative AR(1) prior, G simplex formulation or lack of autocorrelation prior models, although these models are also
421 tried on the other datasets to further highlight this ~~prove this further~~.

422 3.2 Chemically sparse synthetic offline dataset

423 This synthetic offline dataset was used ~~is intended~~ to assess the performance of different models on a proxy representation of
424 atmospheric aerosol data, while maintaining the verifiability property inherent to synthetic datasets as described in Section
425 2.4.2. We performed source apportionment of the **X** matrix through the aforementioned Bayesian models and PMF, obtaining
426 4 factors fingerprints and time series. The dataset used in this source apportionment is expected to be much more sparse than
427 ACSM-like datasets, hence it could better expose the capabilities and added value of the sparsity prior.

429 To avoid initialisation failure, BAMF was run by initialising **F** as a normal distribution to ensure a more sturdy sampling,
430 Model initialisation fails when no set of initial parameter values satisfying the model result in valid Bayesian solutions, and
431 are usually solved by imposing more informative priors constraints on the model parameters, to make the sampling more sturdy
432 to avoid the initialization failure otherwise. A t-test was run comparing the **F**, **G** factors from this slightly modified model and
433 BAMF to ensure their similarity. Its results passed the t-test for all factors except for one factor, although it presented a
434 $R^2=0.9990$ correlation and only a 20% of quantitative difference with that BAMF factor. Hence, one can assume that the model
435 provides an acceptable level of agreement with BAMF, capturing the essential structure of the factors with only very minor
436 deviations.

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

437

438 Figure 3 presents the (a) time series (b) auto-correlation (c) profiles of the source apportionment solution for PMF, BAMF,
439 BAMF+HS, (d) additional comparison to truth metrics, and Figure 4 shows the histograms of the models \mathbf{F} components
440 estimation. The time series and autocorrelation show only slight differences between the models, the PMF being the most
441 different to the truth in all factors except the salt one, as supported in Figure 3 (d). Amongst factors, the coarse biological
442 source is the most poorly reconstructed. If accounting for the sum of all factors \mathbf{G} R^2 s and \mathbf{G}/\mathbf{G}_0 , in the last row of Figure 3(d),
443 the most accurate model is the BAMF+HS, followed by PMF and then BAMF. In terms of profiles, the best overall model
444 performance depends on the metric, \mathbf{F} Spearman correlation coefficient being highest for BAMF and R^2 and cosine similarity
445 correlation coefficients for BAMF+HS. This fact, accompanied by Gini being the highest for BAMF+HS and the closest to 1
446 Gini ratio, indicates that the extreme values of the profile (i.e. maximum and zeros species contributions) are closer to truth
447 for BAMF+HS, whose extreme contributions would be less relevant in the Spearman correlation coefficient. Considering the
448 Truth \mathbf{F} zeros sum metric, the horseshoe shrinkage is visibly sparsifying most of the low signals whilst BAMF and PMF
449 present non-zero contributions for species whose contribution in this factor is null. Hence, the BAMF+HS model would
450 effectively promote the profiles sparsity which it was intended for.

451

452 However, the favourable results of BAMF+HS in comparison to the other models could be a dataset-dependent finding, related
453 to the properties of the created synthetic dataset. The purely-measurement-based offline synthetic dataset, whose performance
454 statistics are shown in Table S4, shows that PMF overperforms BAMF+HS, presenting slightly higher \mathbf{F} and \mathbf{G} R^2 and better
455 \mathbf{G}/\mathbf{G}_0 . This could indicate that the optimal model selection might be dataset dependent. However, the source time series of this
456 very simplistic dataset are fully correlated with some species time series, since they are used to generate factor time series,
457 which makes it a very redundant dataset. In this scenario, the source apportionment comparison might still be valid, but it is
458 not the perfect showcase for RMs testing due to the excessive source correlation with species. ~~W~~However, we found it valuable
459 to present different model performances on different datasets, which in atmospheric measurements, can suffer from artefacts
460 complicating the behaviour of some models.

461

462 In the same way, the alternative autocorrelation priors models were also tried and will be thoroughly discussed in Section SI
463 C.II. However, overall, the BAMF+HS model is the one providing the best source apportionment results for this offline dataset,
464 taking advantage of the sparsity to upgrade both profiles and time series accuracy.

465

466 3.3 Real-world offline dataset

467 To test the models on real-world data and identify their limitations for more complex datasets, we tested the models in the real-
468 world offline PM₁₀-PM_{2.5} ~~offline~~ dataset described in Section 2.4.3. Since the truth is not accessible, the model performance
469 can only be assessed upon environmental, factorisation-related, and computational criteria. For this dataset, BAMF and

470 BAMF-AR1 models presented initialisation issues ~~preventingineapacitating~~ them ~~fromto~~ properly launch~~ing~~ the models. ~~To~~
471 ~~avoid this issue and make the model more robust, we implemented a prior in F, so that its components are drawn from~~~~To avoid~~
472 ~~this issue, we made them more robust by initialising F components as~~ Gaussian distributions centered at zero and with a
473 standard deviation of 1 so that we restrict values to be bounded to 1. This modification was not needed for the other models,
474 which did not present initialisation issues.

475
476 Source apportionment results for PMF, BAMF, and BAMF+HS are shown in Figure 5 and Table 4. Figure 6 shows the **F**
477 distributions for these models, as a detail of Figure S9 (a). Figure S6, S9 (a) display very similar results for PMF, BAMF,
478 BAMF+HS both in terms of **F**, **G**, and reconstruction metrics, and only some differences can be perceived for PMF, while
479 BAMF and BAMF+HS histograms are almost overlapping in Figure 6. However, the BAMF+HS profiles present a remarkable
480 difference in terms of sparsity as seen in the **F** Gini metric, which is mostly the highest for BAMF+HS or equal, except for the
481 biological factor for which PMF is slightly higher. For some species, the relative **F** components apportionment is more strongly
482 suppressed by BAMF+HS than by BAMF or PMF, hence, their contribution on other profiles can be larger. This is clearly
483 visible, for instance, for OC, Mg⁺, K⁺, S⁺, or mannitol, which are zeroed in the Salt factor and consequently are larger on the
484 factors where these species are relevant. This is more evidently depicted in Figure S9 (a) and Figure 6, where the distribution
485 of **F** components is shown. For the aforementioned species, the horseshoe effect can be seen in the distribution, whilst BAMF
486 and PMF are further from zero. This result thus highlights the potential benefits of sparsity introduction in matrix factorisation.

487
488 The application of other autocorrelation priors was not advantageous with respect to the regular BAMF autocorrelation and
489 even worsened the shrinkage power of the horseshoe prior as discussed in SI C.III.

490 3.4 Chemically less sparse synthetic online ACSM datasets

491 The next step ~~was to test these models on more realistic synthetic datasets~~~~to test these models is their trial on more realistic~~
492 ~~synthetic datasets~~. For that purpose, 6 datasets for 4 European cities (a total of 24 datasets) were designed with 5 factors in
493 each of them (section 2.4). We applied the 8 models under discussion (PMF, BMF, BMF+HS, BAMF, BAMF+HS, BAMF-
494 AR1, BAMF-GS) to the 24 synthetic datasets and computed the summary statistics (the median of the ratios of **G** over the
495 truth **G**, **G/G**₀, the Pearson correlation of **G** with truth, **G** ρ , and the Spearman correlation of **F** with truth, **F** ρ). All metrics
496 over cities, datasets and sources are presented in Table S5, ~~and an example for one site (Zurich) and one dataset (dataset 0,~~
497 ~~from 01/01/2019 to 14/01/2019) is shown in Figure S11 as an example of the results obtained by the three models in 1 out of~~
498 ~~the 24 datasets.-~~

499
500 Figure 7 shows the model summary statistics over the 6 generated datasets for the four cities and Figure S11+2 shows the factor-
501 dependent statistics. In this case, the Pearson correlation coefficient (not squared) was used to compare the results of these
502 ACSM-like datasets more easily to those presented in Rusanen et al. (2024), ~~which used this metric~~~~in which this metric was~~

Formatted: Font color: Blue

Formatted: Font: Bold, Font color: Blue

Formatted: Font: (Default) NimbusRomNo9L-Medi, Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

503 used. Broadly speaking, Figure 7 shows a good agreement between models and the truth, with ~~most solutions~~ ~~mostly all~~
504 ~~solutions~~ with correlations with truth for \mathbf{F} and \mathbf{G} above 0.7, similarly to Rusanen et al. (2024). However, there are clear
505 differences amongst models and cities. PMF is performing worse in comparison to the Bayesian models, including BMF, the
506 Bayesian analog to PMF in all datasets except for Milan. As shown in Table S5, PMF presents the highest $|\mathbf{Z}-\mathbf{X}|/\sigma$, the highest
507 overestimations of \mathbf{G} , and correlations of \mathbf{G} and \mathbf{F} are the lowest in comparison to other models except for the Milan dataset.
508 In terms of \mathbf{G}/\mathbf{G}_0 , the model providing the best results are BAMF, BAMF-AR1, BAMF-GS, followed by their horseshoe
509 versions. The BAMF+HS, presents slightly lower $F\rho$, $F R^2$, and the sparsity Gini metric ratio is not close to one, entailing the
510 horseshoe prior did not successfully implement sparsity and the \mathbf{F} accuracy did not improve. In terms of correlations with \mathbf{F}
511 and \mathbf{G} , the models including the horseshoe prior present higher dispersion within a city with respect to the models without
512 sparsity terms. Considering all the parameters, the models with the best overall performance are BAMF, BAMF-GS, and
513 BAMF+HS.

514
515 Figure S132 shows the autocorrelation for lags 0-168 h (half of the monthly measurement period) for all the sources and sites,
516 displaying the cyclicity of the selected sources. In all cases, the short-term lags present very high autocorrelation, entailing
517 that the similarity on adjacent timestamps is very high and decays over longer periods. Typically, and as presented on the
518 figure, the autocorrelation of primary sources, with more marked daily cycles, decays faster than secondary sources, which
519 evolve more steadily due to their slower reaction to emissions. Whilst HOA and BBOA present a very steady intraday
520 structure, with one or two maxima per day, the biogenic SOA presents one peak per day and the other two secondary sources
521 may or may not present marked daily cycles. This different intra- and inter-daily structure amongst sources certainly challenges
522 the models to resolve the source-dependent characteristic.

523
524 Figure 8 shows the autocorrelation from truth and the model outputs correlate (Pearson coefficient of determination) for each
525 model and source in the 4 cities. Figure 8 shows the modelled vs. truth Pearson correlation coefficients for the autocorrelation
526 of each source in the 4 cities. Each dot represents one of the 6 datasets for each site, and colors represent the different sources.
527 The results show that all models present very high Pearson coefficient ranges for \mathbf{G} autocorrelations in comparison to truth
528 except for PMF, which struggles with this dataset aspect due to the lack of accounting of self-correlation. In general terms, the
529 best captured correlation by all models is that of SOA_{Bio} , with the most regular cyclical patterns. The SOA_{BB} and SOA_{T}
530 autocorrelations seem to challenge the models further due to more irregular patterns, and for some datasets, their
531 autocorrelation is poorly modeled. POA sources are generally accurately modeled, with HOA patterns slightly better captured
532 than those from BBOA. Regarding models, the ones with better performance are BAMF-GS, BAMF, and BAMF+HS, with
533 only slight differences between the last two. This observation suggests that the horseshoe prior addition does not significantly
534 reduce the autocorrelation power of the BAMF.

Formatted: Font color: Blue

536 Regarding sparsity, Figure S134 depicts the lack of sparsity both for input and modelled data. This figure shows the truth's 5
537 lowest m/z components as well as BAMF, BAMF+HS outcomes. The reference (truth) profiles do not present zeros but very
538 small signals, as do many ACSM-like profiles in the AMS spectral database (Ulbrich et al. 2009). Both BAMF, BAMF+HS
539 reflect this lack of sparsity, however, it could be expected that BAMF+HS would decrease the contributions of the lowest
540 components. However, the sparsity introduction was not achieved as seen before in the lack of improvement of the Gini ratios.
541 This lack of sparsity despite the enforcement through the horseshoe prior can be explained by the complexity of the data, which
542 due to chain divergence, hinders the models performance. Figure S154 shows the model \hat{R} , a typical Bayesian metric to
543 evaluate the precision of Hamiltonian chains, computing the ratio between inter- and intra-chain variabilities. In any case
544 results are very close to the ideal value, 1, so the validity of all models' solutions is assured. However, this plot reflects the
545 ~~solution~~ deprecation of the solution with models when the horseshoe prior is applied. The horseshoe prior adds more
546 complexity to the \mathbf{F} with three more parameters compared to non-sparsity models which could be the cause of the increased
547 model instability across chains.

548
549 Finally, a sensitivity analysis was run for the first Zurich dataset perturbing independently the original \mathbf{F} , \mathbf{G} matrices to
550 different degrees, monitoring the correlation of the modelled \mathbf{F} , \mathbf{G} matrices to the original truth (Figure 9). Subfigures (a) and
551 (b) show how both in the case of the original \mathbf{F} and \mathbf{G} perturbations, the \mathbf{F} accuracy drops immediately and analogously for
552 both models, with a more sudden decay for \mathbf{G} perturbations. Contrarily, the affectations in \mathbf{G} (subplots (c) and (d)) are different
553 for both models, with a steady decay for BAMF with \mathbf{G} perturbations and a non-clear trend for \mathbf{F} perturbations, whilst
554 BAMF+HS correlation rests insensitive to \mathbf{F} , \mathbf{G} perturbations with an increasing/decreasing erratic behaviour. This result
555 shows the reduced precision in \mathbf{G} of BAMF+HS in comparison to BAMF due to the chain divergence issue, which, in any
556 case, does not severely compromise its accuracy. This finding also explains justifies the bigger variations for BAMF+HS with
557 respect to BAMF in all the metrics shown in Figure 8. Additionally, it showcases the general strong sensitivity of \mathbf{F}
558 determination opposite to the general robustness of \mathbf{G} upon general \mathbf{X} matrix perturbation.

559 **3.4. Discussion**

560 This study aims to explore further BAMF capabilities and the benefits introduced through additional priors and/or
561 modifications of the current model structure as given by Rusanen et al. (2024). The introduction of sparsity in source
562 apportionment models was of particular interest to provide more distinct and concise source profiles which can, in turn,
563 improve the time series accuracy. However, in real-world applications, it may also remove small but relevant signals along
564 with noise. Therefore, comparison with BAMF results is recommended, leaving it to the user to decide whether the method's
565 use is appropriate for their case.

Formatted: Font: Bold

Formatted: Outline numbered + Level: 1 + Numbering Style: 1, 2, 3, ... + Start at: 3 + Alignment: Right + Aligned at: 0 cm + Indent at: 1.27 cm

567 Firstly, the use of the simplistic toy dataset highlighted the added value of the sparsity introduction through the horseshoe prior
568 in the totally constrained experiment. In this controlled setting, the ground truth structure is well defined, allowing the effect
569 of sparsity to be clearly isolated and the method performance validated. However, for an unconstrained experiment, sparsity
570 was proven remarkably advantageous, but subject to the underlying matrix factorisation results. That is, the horseshoe prior in
571 BAMF+HS effectively suppresses weak signals of F contributions as determined by BAMF, yet it fails to guide the model
572 toward a more accurate or sparser solution when the initial BAMF estimate is suboptimal. Other sparsity priors, like Lasso and
573 Spike-and-slab, were tried out but did not improve the regularised horseshoe performance.

Formatted: Font color: Blue

574
575 The regularised horseshoe prior introduction in BAMF improved apportionment of offline synthetic and real-world datasets
576 with respect to BAMF, promoting sparser profiles. The synthetic dataset comparison to truth was maximal for BAMF+HS,
577 with sparser profiles and consequently better G accuracy. Its application also proved advantageous for the real-world dataset,
578 despite not being able to be compared to the truth. In this case, improvements are assessed through increased profile
579 distinctness and internal consistency rather than absolute accuracy. The results show a sparsity effect which provides more
580 distinct profiles in comparison to PMF and BAMF. This result encourages the usage of the horseshoe prior for sparsity
581 introduction in datasets whose solutions are expected to be strongly sparse, in datasets of expected strong sparsity, such as
582 elemental datasets.

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue

583
584 Subsequently, in the more realistic European datasets, the sparsity introduction could not be effectively enforced. Although
585 solution quality was not substantially compromised, the profiles remained non-sparse after applying the prior. This is likely
586 due to model instability arising from the higher complexity of these datasets, which is further aggravated by the addition of
587 the horseshoe prior, as it requires sampling a larger number of parameters. Moreover, the inherent nature of ACSM datasets—
588 characterized by highly correlated species—might also contribute to this limitation, since the model struggles to disentangle
589 overlapping sources when variables are strongly interdependent. The higher chain divergence found for the horseshoed models
590 causes a drop in solution precision due to different landings on the solution space depending on the chain. This issue could be
591 reduced by selecting chains a-posteriori upon user-defined criterion as is practiced in PMF. The higher chain divergence found
592 for the horseshoed models, of lower precision due to different landings on the solution space depending on the chain, which
593 could be reduced by selecting chains a-posteriori upon user-defined criterion as is practiced in PMF. This is further confirmed
594 by the insensitivity to G or F perturbations that are visible for BAMF+HS but not for BAMF ~~the BAMF+HS suffers in~~
595 ~~comparison to BAMF~~. Nonetheless, given that ACSM-like factor profiles exhibit low sparsity in the literature, ~~hence~~ the use
596 of sparsity priors in these datasets is less justified. Also, because usually ACSM profiles obtained in chamber or ambient
597 experiments are not usually sparse, as seen in Ulbrich et al. (2009), the BAMF+HS is not as pertinent in these kinds of datasets
598 as for filter-based datasets.

Formatted: Font: Italic

Formatted: Font color: Blue

600

601 The sparsity conceptual framework could also be brought into PMF through the pulling equations, which can shrink down
602 manually the expectedly low signals in a factor. However, this methodology requires that the user indicates the species that
603 are intended to be zeroed, which introduces user-subjectivity to the problem. The BAMF+HS method, contrarily, acts globally,
604 shrinking those species with lowest signals in favour of the matrix factorisation, hence no user intervention is needed. This
605 makes the approach more objective but also less targeted, returning the factorization optimisation agency to the model.
606 However, if the purpose were to enforce a shrinkage of a certain species as in the PMF case, this feature could also be
607 implemented through the horseshoe method with minimal code modification.

Formatted: Font color: Blue

608
609 The results of the other models tested ~~The other tried-out models~~ (BAMF-AR1, BAMF-GS) did not show a significant
610 improvement with respect to BAMF. The BAMF-AR1 contains another autocorrelation to parametrisation (STAN Team,
611 2025) which should allow for trend consideration, although this matter was not tackled in the current work and remains to be
612 validated in future studies. The BAMF-GS seemed to capture slightly better the G variability in comparison to BAMF in the
613 online datasets, but led to worse correlation to truth in the offline synthetic dataset. Nonetheless, it does not support enforcing
614 sparsity in F, thereby reducing its effectiveness for profile adjustments.

Formatted: Font color: Blue

615 4.3. Conclusions

616 This study presents a sparsity introduction technique for the Bayesian Autocorrelated Matrix Factorisation model (BAMF)
617 which intends to ~~condense~~essentialise source apportionment profiles removing noisy signals. The regularised horseshoe prior
618 is introduced in BAMF (BAMF+HS) in order to narrow down the lowest signals in factor profiles while keeping the most
619 significant ones regularised. The BAMF+HS model is built in STAN, an open-source framework for statistical modelling with
620 Hamiltonian-Montecarlo Markov Chain MCMC-sampling. In order to test the capabilities of the developed model, we
621 generated three kinds of synthetic datasets to compare the model factorisation outputs to the truth factors, namely Toy, offline,
622 and online synthetic datasets, each representing a progressively increasing level of complexity. Likewise, to confirm its
623 usability to real-world data, BAMF+HS was also applied to a multi-site filter dataset. Given the opportunity to explore source
624 apportionment with different types of datasets, we also tested other receptor models such as Positive Matrix Factorisation
625 (PMF) and other BAMF-like Bayesian models. In the Bayesian framework, we ~~tested~~essayed a different formulation of the
626 autocorrelation term (BAMF-AR1) and a permutation on the factorisation matrix logic (BAMF-GS).

Formatted: Font color: Blue

627 The main result highlights can be summarised as:

- 628 - BAMF+HS has been shown to be advantageous to introduce sparsity in factor profiles for offline datasets and to not
629 deprecate the solution for the more complex datasets mimicking Aerosol Chemical Speciation Monitor (ACSM)
630 dataACSM-like datasets. Other sparsifying priors tried out were not as effective in low-signal shrinkage.

Formatted: Font color: Blue

- 631 - The BAMF+HS performance towards truth profile reconstruction was higher than for BAMF and PMF in the toy and
632 offline synthetic datasets. Improving \mathbf{F} typically led to a more accurate determination of \mathbf{G} , highlighting the strong
633 interdependence between the two factorisation matrices.
- 634 - The real-world dataset also shows a better description of sources through BAMF+HS in terms of matrix factorisation
635 metrics and profile sparsification achievement.
- 636 - As shown in the toy dataset, the introduction of sparsity did not solve factorisation issues inherent to the underlying
637 factorisation model.
- 638 - The BAMF+HS model does not create sparsity in ACSM-like datasets, which are originally, indeed, non-sparse.
639 BAMF+HS is more unstable than BAMF for these more complex datasets as a result the higher chain divergence
640 during Hamiltonian-Montecarlo Markov Chain sampling as suggested by the \hat{R} metric. However, the effects of
641 the horseshoe prior do not affect the overall performance of BAMF or its autocorrelation accuracy.
- 642 - The alternative formulations for BAMF, BAMF-AR1 and BAMF-GS, did not show a significant improvement with
643 respect to BAMF.

Formatted: Font color: Blue

644 With all that, profile sparsity has been shown to substantially enhance the accuracy of source apportionment analyses,
645 improving the separation of the chemical composition of sources. The BAMF+HS model succeeds in incorporating this
646 property in profile fingerprints, especially in filter-based datasets. Using BAMF+HS in such datasets, the solutions reflect the
647 sparsity of filter-based chemical profiles, hence, this newly introduced method is encouraged when source fingerprints are
648 expected to be substantially sparse. However, for ACSM-like datasets, the sparsity is not fully achieved due to converge issues,
649 although the quality of the solution is not substantially deprecated with respect to BAMF. With the aim of improving further
650 source apportionment techniques, future research should be directed to enhance the robustness and generalisability of the
651 BAMF+HS model across diverse data types. Moreover, continued exploration of the underlying properties of solution spaces
652 — such as profiles sparsity, time series autocorrelation — may provide valuable insights into disentangling complex source
653 contributions through receptor modelling. In this regard, the Bayesian source apportionment framework offers a particularly
654 suitable foundation, allowing for the integration of prior knowledge and uncertainty quantification in the inference process.

Formatted: Font color: Blue

655 **Code and data availability**

656 The models and datasets can be found at <https://github.com/martavia0/BAMF-horseshoe.git>

657 **Author contribution**

658 MV: Conceptualisation, data curation, formal analysis, funding acquisition, investigation, methodology, project
659 administration, resources, software, validation, visualisation, writing (original draft preparation). YH: Formal analysis,
660 investigation, software; JD: investigation, resources, software, validation. MM: Data curation. AR: Data curation, formal

661 analysis, methodology, investigation, resources, software. JJ: Data curation. SKG: Data curation. J-LJ: Data curation; VNTD:
662 Data curation. GU: Data curation. GM: conceptualisation, funding acquisition, investigation, supervision, validation. KRD:
663 Conceptualisation, data curation, formal analysis, funding acquisition, investigation, methodology, supervision, validation. All
664 co-authors participated in the revision and edition of the manuscript.

665

666 **Competing interests**

667 The authors declare that they have no conflict of interest.

668

669 **Disclaimer**

670

671 Co - funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not
672 necessarily reflect those of the European Union or European Research Exacutive Agency. Neither the European Union nor the
673 granting authority can be held responsible for them.

674 **Acknowledgements**

675 This publication is co- funded by/ has received funding from/ the European Union's Horizon Europe research and innovation
676 program under the Marie Skłodowska-Curie COFUND Postdoctoral Programme grant agreement No.101081355- SMASH
677 and by/from /the Republic of Slovenia and the European Union from the European Regional Development Fund. Work was
678 partly supported by ARIS program P1-0385. All the chemical measurements on the PM10 and PM2.5 Swiss samples used in
679 this study were performed on the Air O Sol analytical plateau at IGE (France). K.R.D. acknowledges support by the SNSF
680 Ambizione grant PZPGP2_201992. The authors thank Jimeng Wu and Imad El Haddad for their input related to
681 autocorrelation.

682

683 **References**

684 Andersen, M. R., Winther, O., & Hansen, L. K. (2014). Bayesian inference for structured spike and slab priors. *Advances in*
685 *Neural Information Processing Systems*, 27.

686 Belis, C. A., Karagulian, F., Larsen, B. R., & Hopke, P. K. (2013). Critical review and meta-analysis of ambient particulate
687 matter source apportionment using receptor models in Europe. *Atmospheric Environment*, 69, 94-108.

688 Belis, C. A., Larsen, B. R., Amato, F., El Haddad, I., Favez, O., Harrison, R. M., ... & Viana, M. (2014). European guide on
689 air pollution source apportionment with receptor models. JRC reference reports EUR26080 EN.

Formatted: Font color: Auto

Formatted: Font color: Auto

690 Belis, C., Pernigotti, D., Karagulian, F., Pirovano, G., Larsen, B., Gerboles, M., and Hopke, P.: A new methodology to assess
691 the performance and uncertainty of source apportionment models in intercomparison exercises, *Atmospheric Environment*,
692 119, 35–44, 2015.

693 Brinkman, G., Vance, G., Hannigan, M. P., and Milford, J. B.: Use of synthetic data to evaluate positive matrix factorization
694 as a source apportionment tool for PM_{2.5} exposure data, *Environmental science & technology*, 40, 1892–1901, 2006.

695 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., ... & Riddell, A. (2017). Stan: A
696 probabilistic programming language. *Journal of statistical software*, 76, 1-32.

697 Crippa, M., DeCarlo, P. F., Slowik, J. G., Mohr, C., Heringa, M. F., Chirico, R., ... & Baltensperger, U. (2013). Wintertime
698 aerosol chemical composition and source apportionment of the organic fraction in the metropolitan area of Paris. *Atmospheric
699 Chemistry and Physics*, 13(2), 961-981.

700 Dai, T., Dai, Q., Yin, J., Chen, J., Liu, B., Bi, X., ... & Feng, Y. (2024). Spatial source apportionment of airborne coarse
701 particulate matter using PMF-Bayesian receptor model. *Science of The Total Environment*, 917, 170235.

702 Daellenbach, K. R., Uzu, G., Jiang, J., Cassagnes, L. E., Leni, Z., Vlachou, A., ... & Prévôt, A. S. (2020). Sources of particulate-
703 matter air pollution and its oxidative potential in Europe. *Nature*, 587(7834), 414-419.

704 Daellenbach, K. R., Manousakas, M., Jiang, J., Cui, T., Chen, Y., El Haddad, I., ... & Prévôt, A. S. H. (2023). Organic aerosol
705 sources in the Milan metropolitan area—Receptor modelling based on field observations and air quality modelling. *Atmospheric
706 Environment*, 307, 119799.

707 Dinh, V. N. T., Uzu, G., Dominutti, P., Sauvage, S., Elazzouzi, R., Darfeuil, S., Voiron, C., Samaké, A., Zhang, S., Socquet,
708 S., Favez, O., and Jaffrezo, J.-L.: Toolbox for accurate estimation and validation of PMF solutions in PM source apportionment,
709 *EGUsphere* [preprint], <https://doi.org/10.5194/egusphere-2025-1968>, 2025.

710 Elser, M., Huang, R. J., Wolf, R., Slowik, J. G., Wang, Q., Canonaco, F., ... & Prévôt, A. S. (2016). New insights into PM 2.5
711 chemical composition and sources in two major cities in China during extreme haze events using aerosol mass spectrometry.
712 *Atmospheric Chemistry and Physics*, 16(5), 3207-3225.

713 Gelman, A., & Rubin, D. B. (1992). Inference from Iterative Simulation Using Multiple Sequences. *Statistical Science*, 7(4),
714 457–472. <https://doi.org/10.1214/ss/1177011136>

715 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B.: *Bayesian data analysis*, 3rd edn., CRC
716 Press, ISBN 9781439898208, 2014.

717 Gini, Corrado (1936). "On the Measure of Concentration with Special Reference to Income and Statistics", Colorado College
718 Publication, General Series No. 208, 73–79

719 Govan, E., Jackson, A. L., Inger, R., Bearhop, S., & Parnell, A. C. (2023). *simmr*: A package for fitting stable isotope mixing
720 models in R. *arXiv preprint arXiv:2306.07817*.

721 Grange, S. K., Fischer, A., Zellweger, C., Alastuey, A., Querol, X., Jaffrezo, J. L., ... & Hueglin, C. (2021). Switzerland's
722 PM₁₀ and PM_{2.5} environmental increments show the importance of non-exhaust emissions. *Atmospheric environment: X*,
723 12, 100145.

724 Hoffman, M. D. and Gelman, A.: The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo, *J.*
725 *Mach. Learn. Res.*, 15, 1593–1623, 2014.

726 Jiang, J., Aksoyoglu, S., El-Haddad, I., Ciarelli, G., Denier van der Gon, H. A., Canonaco, F., ... & Prévôt, A. S. (2019).
727 Sources of organic aerosols in Europe: A modelling study using CAMx with modified volatility basis set scheme. *Atmospheric*
728 *Chemistry and Physics Discussions*, 2019, 1-35.

729 Keuken, M. P., Moerman, M., Voogt, M., Blom, M., Weijers, E. P., Röckmann, T., & Dusek, U. (2013). Source contributions
730 to PM_{2.5} and PM₁₀ at an urban background and a street location. *Atmospheric Environment*, 71, 26-35.

731 Kuhn, H. W. (1955). The Hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2), 83-97.

732 Lingwall, J. W., & Christensen, W. F. (2007). Pollution source apportionment using a priori information and positive matrix
733 factorization. *Chemometrics and Intelligent Laboratory Systems*, 87(2), 281-294.

734 Liu, D. C. and Nocedal, J.: On the Limited Memory BFGS Method for Large Scale Optimization, *Math. Program.*, 45, 503–
735 528, <https://doi.org/10.1007/BF01589116>, 1989.

736 Manousakas, M. I., Rausch, J., Jaramillo Vogel, D., Schneider-Beltran, K., Alastuey, A., Jaffrezo, J. L., ... & Dällenbach, K.
737 R. Comparison of PM Source Profiles Identified by Different Techniques and the Potential of Utilizing Single-Particle
738 Analysis Data in Source Apportionment. Available at SSRN 5323830.

739 Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., & Teller, E. (1953). Equation of state calculations by
740 fast computing machines. *The Journal of Chemical Physics*, 21(6), 1087–1092. <https://doi.org/10.1063/1.1699114>

741 Oh, M. S., & Park, C. K. (2022). Regional source apportionment of PM_{2.5} in Seoul using Bayesian multivariate receptor
742 model. *Journal of Applied Statistics*, 49(3), 738-751.

743 Paatero, P., & Tapper, U. (1994). Positive matrix factorization: A non-negative factor model with optimal utilization of error
744 estimates of data values. *Environmetrics*, 5(2), 111-126.

745 Paatero, P. (1999). The multilinear engine—a table-driven, least squares program for solving multilinear problems, including
746 the n-way parallel factor analysis model. *Journal of Computational and Graphical Statistics*, 8(4), 854-888.

747 Park, E. S., Guttorp, P., & Henry, R. C. (2001). Multivariate receptor modeling for temporally correlated data by using
748 MCMC. *Journal of the American Statistical Association*, 96(456), 1171-1183.

749 Park, E. S., Spiegelman, C. H., & Henry, R. C. (2002). Bilinear estimation of pollution source profiles and amounts by using
750 multivariate receptor models. *Environmetrics*, 13(7), 775-798.

751 Park, E. S., Lee, E. K., & Oh, M. S. (2021). Bayesian multivariate receptor modeling software: BNFA and
752 bayesMRM. *Chemometrics and Intelligent Laboratory Systems*, 211, 104280.

753 Piironen, J. and Vehtari, A.: Sparsity information and regularization in the horseshoe and other shrinkage priors, *Electronic*
754 *Journal of Statistics*, 11, 5018–5051, <https://doi.org/10.1214/17-EJS1337SI>, 2017.

755 Rasmussen, M. A., & Bro, R. (2012). A tutorial on the Lasso approach to sparse modeling. *Chemometrics and Intelligent*
756 *Laboratory Systems*, 119, 21-31.

757 Pope III, C. A., & Dockery, D. W. (1999). Epidemiology of particle effects. In *Air pollution and health* (pp. 673-705).
758 Academic Press.

759 Rusanen, A., Bjorklund, A., Manousakas, M. I., Jiang, J., Kulmala, M. T., Puolamaki, K., and Daellenbach, K. R.: A novel
760 probabilistic source apportionment approach: Bayesian auto-correlated matrix factorization, *Atmospheric Measurement*
761 *Techniques*, 17, 1251–1277, <https://doi.org/10.5194/amt-17-1251-2024>, 2024.

762 Sage, A. M. (2007). Evolving mass spectra of the oxidized component of organic aerosol mass spectrometer analysis of aged
763 diesel emissions. *Atmos. Chem. Phys. Discuss.*, 7, 10065-10096.

764 STAN Development Team. (2025). *Stan user's guide* (Version 2.36). <https://mc-stan.org/docs/stan-users-guide/>. Accessed
765 July 2025.

766 Tobler, A. K., Skiba, A., Canonaco, F., Močnik, G., Rai, P., Chen, G., ... & Prevot, A. S. (2021). Characterization of non-
767 refractory (NR) PM 1 and source apportionment of organic aerosol in Kraków, Poland. *Atmospheric chemistry and physics*,
768 21(19), 14893-14906.

769 Tibshirani, R., & Wasserman, L. (2015). Sparsity and the lasso. *Statistical machine learning*, 1-15.

770 Ulbrich, I. M., Handschy, A., Lechner, M., and Jimenez, J. L.: AMS Spectral Database, [http://cires.colorado.edu/jimenez-](http://cires.colorado.edu/jimenez-group/AMSsd/)
771 [group/AMSsd/](http://cires.colorado.edu/jimenez-group/AMSsd/), accessed: 2025-05-19, n.d.

772 Upadhyay, A., Jiang, J., Cheng, Y., Vasilakos, P., Chen, Y., Banos, D. T., ... & El-Haddad, I. (2025). High-resolution modelling
773 of particulate matter chemical composition over Europe: brake wear pollution. *Environment International*, 109615.

774 Via, M., Minguillón, M. C., Reche, C., Querol, X., & Alastuey, A. (2021). Increase in secondary organic aerosol in an urban
775 environment. *Atmospheric chemistry and physics*, 21(10), 8323-8339.

776 Via, M., Chen, G., Canonaco, F., Daellenbach, K. R., Chazeau, B., Chebaicheb, H., Jiang, J., Keernik, H., Lin, C., Marchand,
777 N., et al.: Rolling vs. seasonal PMF: real-world multi-site and synthetic dataset comparison, *Atmospheric measurement*
778 *techniques*, 15, 5479–5495, 2022.

779 Viana, M., Kuhlbusch, T. A., Querol, X., Alastuey, A., Harrison, R. M., Hopke, P. K., ... & Hitzenberger, R. (2008). Source
780 apportionment of particulate matter in Europe: a review of methods and results. *Journal of aerosol science*, 39(10), 827-849.

781 Yang, Y., Ruan, Z., Wang, X., Yang, Y., Mason, T. G., Lin, H., & Tian, L. (2019). Short-term and long-term exposures to fine
782 particulate matter constituents and health: A systematic review and meta-analysis. *Environmental pollution*, 247, 874-882.

783 Zhang, Y., Albinet, A., Petit, J. E., Jacob, V., Chevrier, F., Gille, G., ... & Favez, O. (2020). Substantial brown carbon emissions
784 from wintertime residential wood burning over France. *Science of the Total Environment*, 743, 140752.

785 Zhang, Y., Ma, X., Tang, A., Fang, Y., Misselbrook, T., & Liu, X. (2023). Source Apportionment of Atmospheric Ammonia
786 at 16 Sites in China Using a Bayesian Isotope Mixing Model Based on $\delta^{15}\text{N-NH}_x$ Signatures. *Environmental Science &*
787 *Technology*, 57(16), 6599-6608.

788 Zheng, Y., Xue, T., Zhao, H., & Lei, Y. (2022). Increasing life expectancy in China by achieving its 2025 air quality
789 target. *Environmental science and ecotechnology*, 12, 100203.

790

791

792

793 **Figures**794 **Table 1. Priors used in the models tested in this manuscript.**

795

Model	G priors	F priors
BMF	None	None
BMF+HS	None	Regularised horseshoe
BAMF	Rusanen et al. (2024)	None
BAMF+HS	Rusanen et al. (2024)	Regularised horseshoe
BAMF-AR1	AR(1)	None
BAMF-AR1+HS	AR(1)	Regularised horseshoe
BAMF-GS	Rusanen et al. (2024)	None
PMF	None	None
CMB	Fixed a-priori.	None
CMB+HS	Fixed a-priori.	Regularised horseshoe

796

797 **Table 2. Sparsity of the F factors of the synthetic datasets.**

798

Dataset	Factor	F Gini	% zeros
Chemically-sparse synthetic toy dataset	Factor 1	0.67	75.0
	Factor 2	0.67	75.0
	Factor 3	0.5	25.0
Chemically-sparse synthetic offline dataset	Dust	0.74	37.5
	Traffic	0.86	12.5

	Salt		0.78	37.5
	Coarse biological		0.88	25.0
Less chemically-sparse synthetic online ACSM datasets	HOA	Krakow	0.74	5.0
		Milan	0.67	0.0
		Paris	0.68	0.0
		Zurich	0.68	0.0
	BBOA	Krakow	0.47	1.2
		Milan	0.72	0.0
		Paris	0.74	0.0
		Zurich	0.58	0.0
	SOA _{bio}	Krakow	0.52	0.0
		Milan	0.50	0.0
		Paris	0.50	0.0
		Zurich	0.67	0.0
	SOA _{BB}	Krakow	0.55	0.0
		Milan	0.53	0.0
		Paris	0.53	0.0
		Zurich	0.73	0.0

	SOA _{TR}	Krakow	0.45	0.0
		Milan	0.42	0.0
		Paris	0.46	0.0
		Zurich	0.60	0.0

799

800 **Table 3. Toy experiment statistics of (a) Factorisation performance. (b) Comparison to truth. Green sequential**
801 **colorscales represent variables whose larger value leans to a better performance and the blue-to-red divergent**
802 **colorscales (centered at 1, in white) represent G/G₀ divergence with respect to 1. Red bars in (a) depict deviations from**
803 **the ideal 0 value.**

804

(a)

Model	Factorisation		
	R ²	Median(Z-X /sigma)	Max(Z-X /sigma)
CMB	0.9985	0.2542	0.5755
CMB+HS	0.9996	0.2648	0.4774
PMF	0.9979	0.2658	0.5847
BMF	0.9650	0.2125	1.2917
BMF+HS	0.9689	0.2107	1.2863
BAMF	0.9818	0.1222	0.7247
BAMF+HS	0.9820	0.1208	0.7180
BAMF-AR1	0.9806	0.0947	1.0257
BAMF-AR1+HS	0.9790	0.1171	1.0297
BMF-GS	0.9657	0.2024	1.3031
BAMF-GS	0.9818	0.1576	0.9082

(b)

Model	Source	G			F			
		G/G ₀	R ²	ρ	R ²	Gini ratio	Zeros sum	
CMB	Source1	1.00	1.00	0.82	1.00	0.61	0.91	0.06
	Source2	1.00	1.00	0.82	1.00	0.60	0.90	0.08
	Source3	1.00	1.00	0.96	1.00	0.44	0.92	0.03
CMB+HS	Source1	1.00	1.00	0.82	1.00	0.63	0.95	0.03
	Source2	1.00	1.00	0.82	1.00	0.63	0.95	0.04
	Source3	1.00	1.00	0.96	1.00	0.47	0.96	0.02
PMF	Source1	1.00	0.96	0.82	0.95	0.35	0.52	0.36
	Source2	3.29	0.97	0.82	0.88	0.30	0.45	0.42
	Source3	0.65	0.74	0.79	0.81	0.31	0.64	0.2
BMF	Source1	1.08	0.87	0.82	0.76	0.25	0.37	0.48
	Source2	4.10	0.84	0.51	0.43	0.16	0.23	0.58
	Source3	0.60	0.46	0.79	0.80	0.21	0.44	0.25
BMF+HS	Source1	1.11	0.88	0.82	0.77	0.25	0.38	0.48
	Source2	3.92	0.85	0.82	0.50	0.15	0.23	0.57
	Source3	0.60	0.44	0.79	0.84	0.22	0.46	0.23
BAMF	Source1	1.66	0.98	0.82	0.78	0.27	0.41	0.46
	Source2	1.99	0.95	0.82	0.97	0.33	0.50	0.37
	Source3	0.54	0.50	0.96	1.00	0.38	0.79	0.08
BAMF+HS	Source1	1.96	0.98	0.82	0.76	0.27	0.40	0.47
	Source2	1.28	0.95	0.82	0.99	0.58	0.86	0.11
	Source3	0.51	0.48	0.96	1.00	0.44	0.93	0.02
BAMF-AR1	Source1	1.70	0.98	0.82	0.79	0.27	0.41	0.46
	Source2	3.05	0.92	0.82	0.89	0.37	0.55	0.36
	Source3	0.42	0.50	0.96	0.94	0.43	0.90	0.08
BAMF-AR1+HS	Source1	1.92	0.99	0.82	0.77	0.26	0.39	0.46
	Source2	2.58	0.89	0.82	0.90	0.48	0.72	0.36
	Source3	0.38	0.49	0.96	0.97	0.49	1.03	0.08
BMF-GS	Source1	0.88	0.92	0.82	0.63	0.21	0.31	0.52
	Source2	1.12	0.86	0.82	0.71	0.20	0.30	0.53
	Source3	1.02	0.29	0.79	0.91	0.22	0.46	0.22
BAMF-GS	Source1	0.89	0.94	0.82	0.66	0.21	0.32	0.53
	Source2	1.12	0.95	0.82	0.70	0.19	0.28	0.53
	Source3	1.02	0.36	0.79	0.90	0.22	0.46	0.22
BAMF-Lasso	Source1	1.76	0.98	0.82	0.76	0.27	0.40	0.47
	Source2	2.02	0.95	0.82	0.95	0.34	0.50	0.37
	Source3	0.49	0.47	0.96	1.00	0.41	0.86	0.05
BAMF-Spike-Slab	Source1	1.57	0.97	0.82	0.77	0.26	0.40	0.47
	Source2	1.78	0.95	0.82	0.99	0.34	0.52	0.34
	Source3	0.62	0.50	0.96	0.98	0.33	0.69	0.12

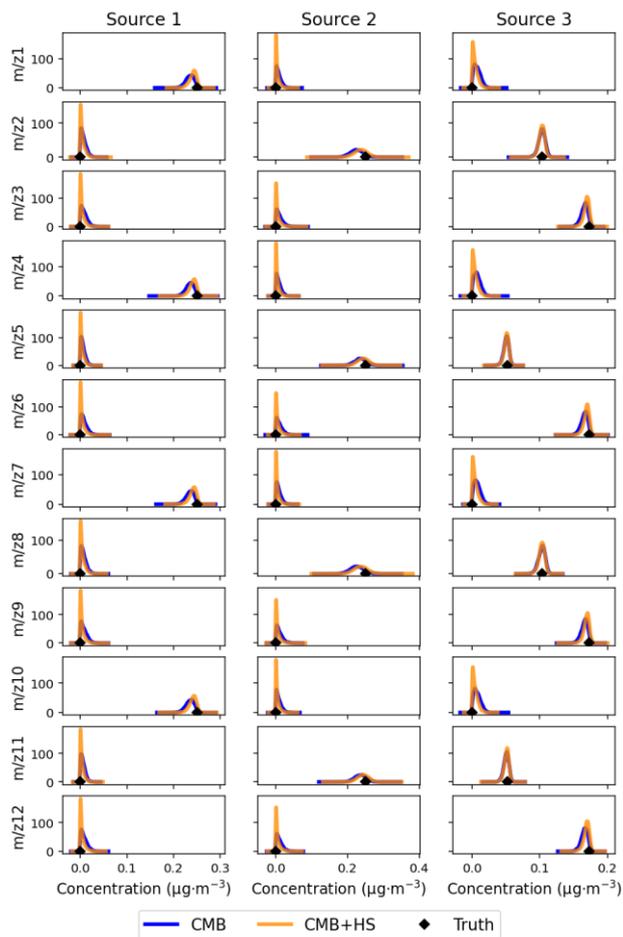


Figure 1. F-matrix components distributions for CMB and CMB+HS (solid lines) compared to truth (markers). F-components distributions for CMB and CMB+HS compared to truth.

Formatted: Font color: Auto

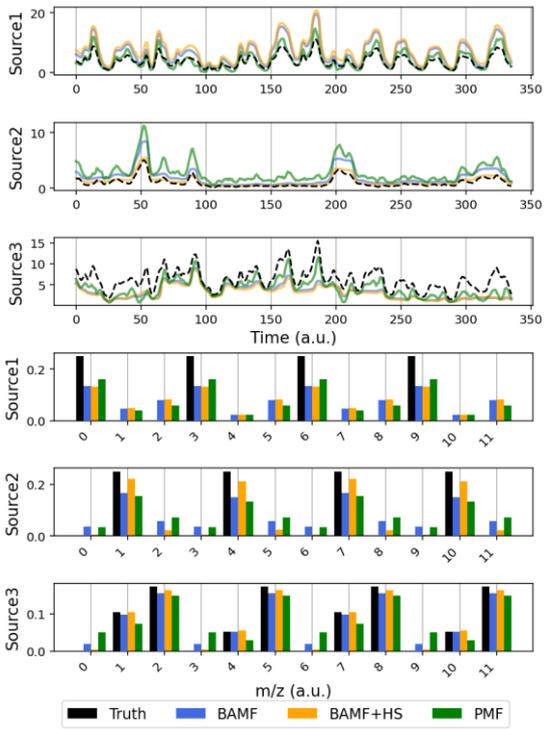
Formatted: Font color: Blue

Formatted: Font color: Auto

Formatted: Font color: Blue

Formatted: Font color: Auto

Formatted: Font color: Blue



811
812

Figure 2. Source apportionment results for the toy dataset obtained using PMF, BAMP, and BAMP+HS, compared against the true solution (black bars). (a) Factor time series. (b) Factor profiles. Toy dataset source apportionment results for PMF, BAMP, and

Formatted: Font: 10 pt, Not Bold, Font color: Auto

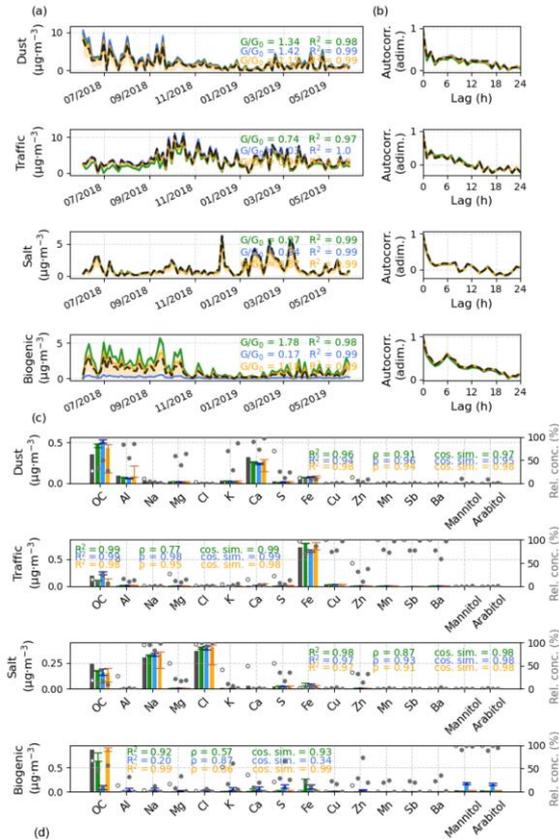
Formatted: Font color: Auto

Formatted: Font color: Blue

Formatted: Font color: Auto

Formatted: Font color: Blue

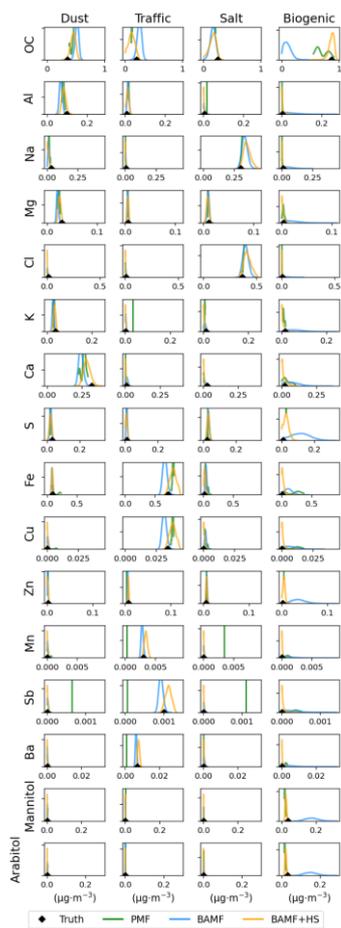
Formatted: Font color: Auto



Factor	Model	F constr. R ²	F constr. p	F Gini	F Gini ratio (*)	Truth zeros sum
Dust	PMF	0.95	0.84	0.79	0.90	0.0035
	BAMF	0.77	0.90	0.80	1.08	0.0025
	BAMF+HS	0.94	0.90	0.79	1.07	0.0011
Traffic	PMF	0.99	0.90	0.89	1.19	0.0000
	BAMF	0.93	0.94	0.85	0.99	0.0005
	BAMF+HS	0.995	0.93	0.88	1.02	0.0004
Salt	PMF	0.95	0.84	0.79	0.90	0.0012
	BAMF	0.95	0.88	0.80	1.01	0.0118
	BAMF+HS	0.98	0.87	0.82	1.04	0.0031
Biogenic	PMF	0.87	0.43	0.85	1.08	0.0035
	BAMF	0.53	0.51	0.50	0.56	0.0118
	BAMF+HS	0.96	0.84	0.90	1.01	0.0007
Σ ₂ (*)	PMF	3.80	3.04	3.34	0.43	0.0198
	BAMF	3.18	3.23	2.95	0.53	0.0172
	BAMF+HS	3.88	3.84	3.39	0.15	0.0031

814 **Figure 3. Synthetic offline dataset source apportionment results for PMF, BAMF, and BAMF+HS models. (a) Time Series. (b)**
 815 **Autocorrelation. (c) Profiles. (d) Table with additional metrics for comparison to truth. Bold numbers reflect the highest value**
 816 **amongst models. F contr. represents here the percentage of each factor into a given species. The sum row reflects the**
 817 **overall performance of the model for all sources for each statistic metric except for the ones marked with (*), in which**
 818 **the difference to 1 in absolute value is summed up.**

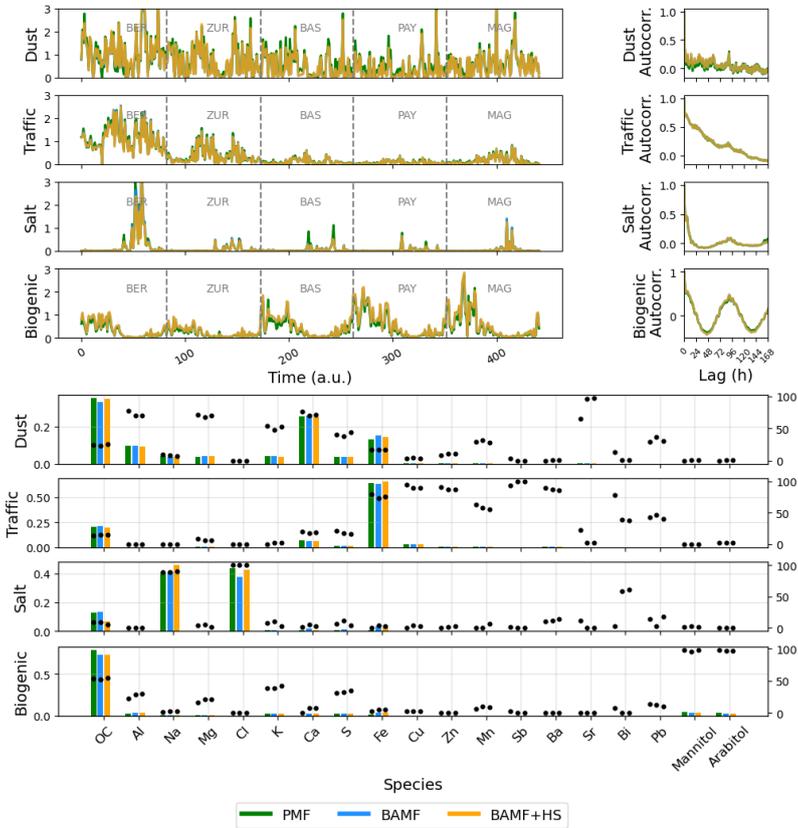
819



821
822
823
824
825

Figure 4. Profile components distribution for PMF, BAMF, BAMF+HS (solid colored lines) in comparison to the truth (markers) on the real-world filters dataset. Rows represent the species of the source apportionment and columns represent sources. Profile components distribution for PMF, BAMF, BAMF+HS on the real-world filters dataset

Formatted: Font color: Blue
Formatted: Font color: Blue
Formatted: Font color: Blue



826
827
828

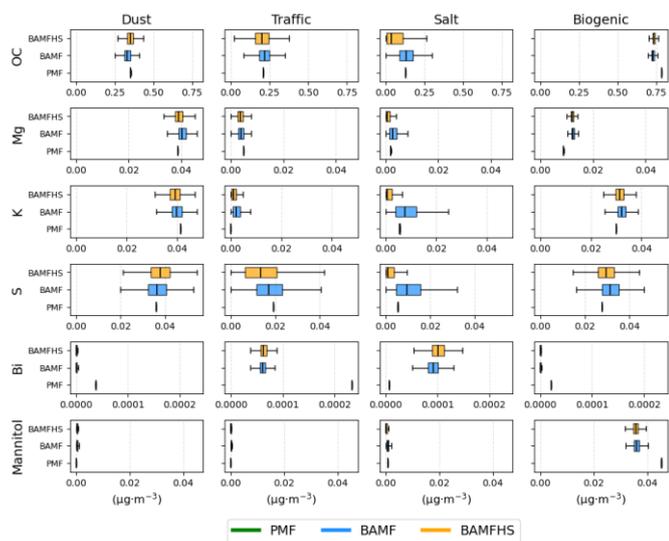
Figure 5. Comparison of PMF, BAMF, BAMF+HS for the real-world filters dataset. From left to right and top to bottom: time series, autocorrelation, and profile plots. The dots in the profiles (right axis) show the contribution of each species to the source.

829
830

Table 4. Offline real-world dataset reconstruction and sparsity statistics. Bold numbers reflect the highest value amongst models.

Model	$R^2(Z, X)$	Median $ X-Z /\sigma$	Median $ X-Z /\sigma$	Factor	F Gini
PMF	0.68	0.77	10.52	Dust	0.77
				Traffic	0.87
				Salt	0.86
				Biogenic	0.87
BAMF	0.67	0.75	11.13	Dust	0.76
				Traffic	0.87
				Salt	0.84
				Biogenic	0.83
BAMF+HS	0.67	0.75	11.12	Dust	0.77
				Traffic	0.87
				Salt	0.87
				Biogenic	0.83

831
832
833
834
835



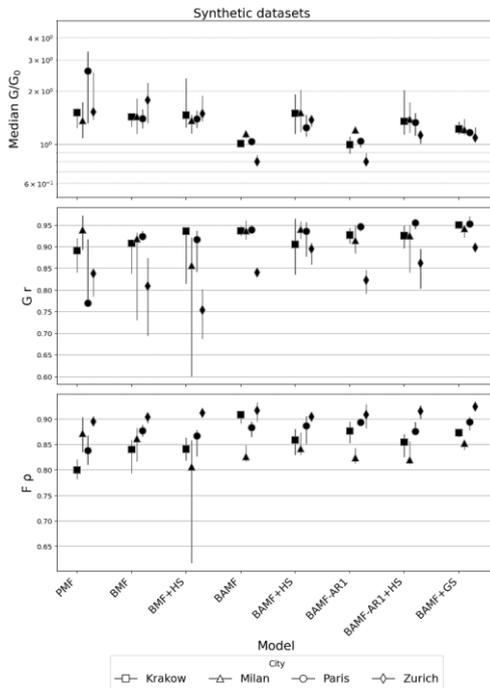
836

837

838

Figure 6. Boxplot Distributions of individual profile components derived from PMF, BAMF, and BAMF+HS analyses for the real-world filter dataset. A complete comparison of all profiles is presented in Figure S9.

839



840
841
842
843

Figure 7. European cities synthetic datasets summary statistics; from top to bottom, median ratio time series with truth (G/G_0), Pearson correlation coefficient of G with truth ($G r$), Spearman correlation coefficient of F with truth ($F \rho$). The axis of the G/G_0 plot is in logarithmic scale.

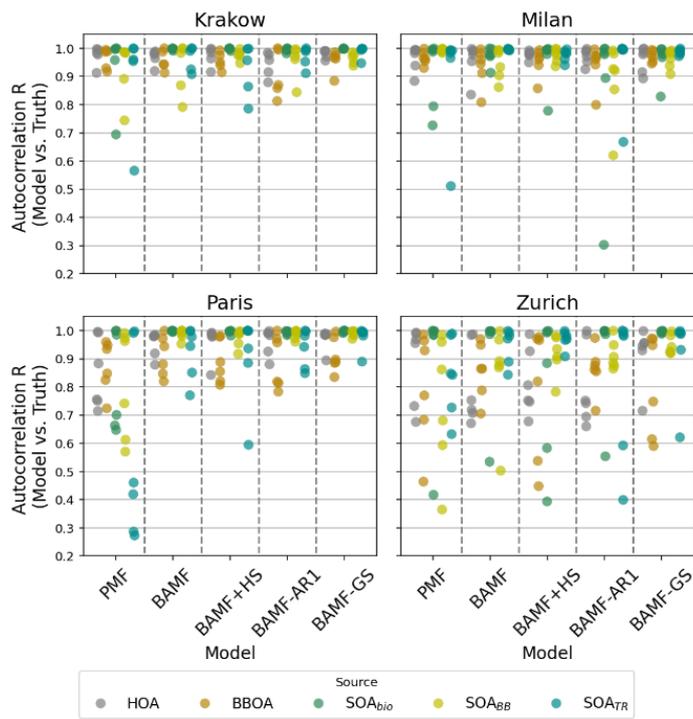
Formatted: Font color: Blue

Formatted: Font color: Blue, Subscript

Formatted: Font color: Blue

Formatted: Font color: Blue

Formatted: Font color: Blue



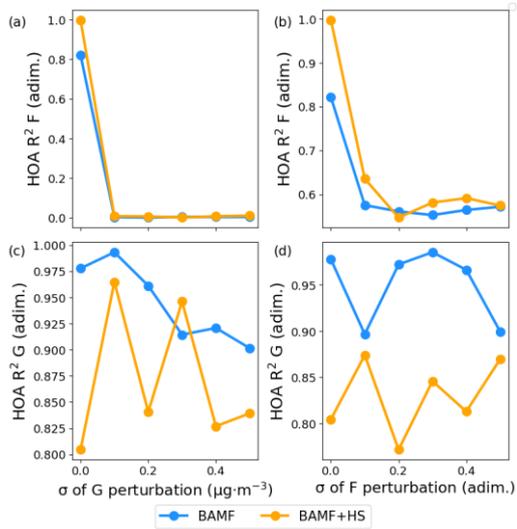
844

845

Figure 8. Pearson correlation of the autocorrelations of model solutions with the truth for all factors and all cities.

846

847



848

849

850

Figure 9. Squared Pearson coefficient of F, G matrix with original truth F, G matrices of the BAMF, BAMF+HS models with the degrees of perturbation in F and G.

1 **Supplementary information for: Sparsity introduction in Bayesian**
2 **models for aerosol source apportionment through the regularised**
3 **horseshoe prior**

4 Marta Via¹, Jure Demšar², Yufang Hao³, Manousos Manousakas^{3,4}, Anton Rusanen⁵, Jianhui Jiang⁶,
5 Stuart K. Grange⁷, Jean-Luc Jaffrezo⁸, Vy Ngoc Thuy Dinh⁸, Gaëlle Uzu⁸, Griša Močnik¹, and Kaspar
6 R. Daellenbach³

7
8 ¹Center for Atmospheric Research, University of Nova Gorica, Ajdovščina 5270, Slovenia

9 ²Faculty of Computer and Information Science, Tržaška Cesta 25, 1000 Ljubljana, Slovenia

10 ³Laboratory of Atmospheric Chemistry, Paul Scherrer Institute, 5232 Villigen PSI, Switzerland

11 ⁴Environmental Radioactivity Aerosol Tech. for Atmospheric Climate Impacts, INRaSTES, National Centre of Scientific
12 Research “Demokritos”, Ag. Paraskevi, 15310, Greece

13 ⁵Atmospheric Composition Research, Finnish Meteorological Institute, 00101 Helsinki, Finland

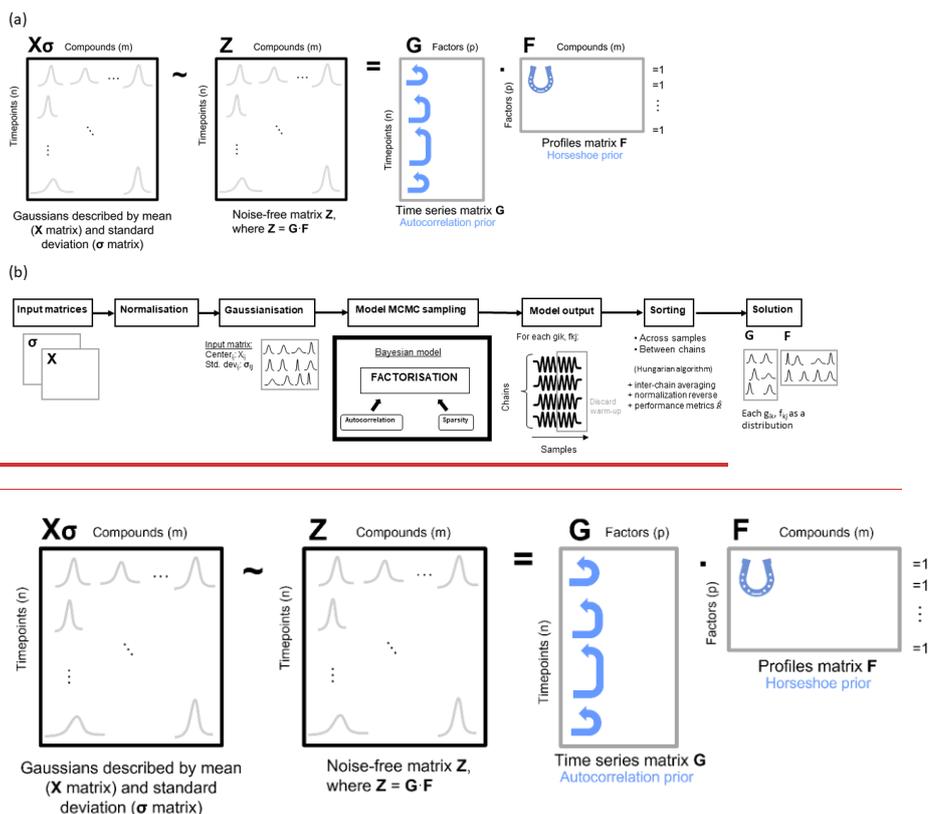
14 ⁶School of Ecological and Environmental Sciences, East China Normal University, 200241, Shanghai, China

15 ⁷Climate and Environmental Physics, Physics Institute, University of Bern, Sidlerstrasse 5, 3012 Bern, Switzerland

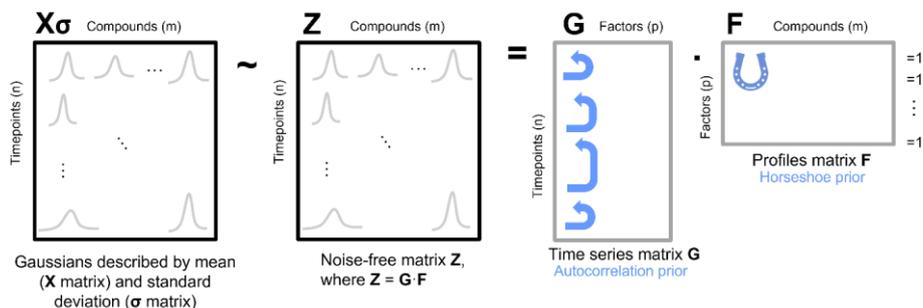
16 ⁸University of Grenoble Alpes, CNRS, INRAE, IRD, Grenoble INP, IGE, Grenoble 38000, France

17
18 *Correspondence to:* Marta Via (marta.viagonzalez@ung.si)

19 A) Supplementary figures



20



21

22 **Figure S1. (a) Bayesian matrix factorisation for source apportionment sketch with autocorrelation and horseshoe priors. (b)**
 23 **Workflow diagram showing the aerosol source apportionment stages with the BAMF+HS model. This figure portrays the**
 24 **previous and posterior processes to the MCMC sampling, and both starting and end points. Concepts as chains and**
 25 **samples, mentioned in the text, are schematised in the workflow as well.**

26 **Bayesian matrix factorisation sketch with autocorrelation and horseshoe priors.**

27 **Table S1. Hamiltonian MonteCarlo sampling parameters for all the conducted experiments.**

28

Formatted: Font color: Auto

Formatted: Font color: Blue

Formatted: Font color: Auto

Formatted: Font color: Blue

Formatted: Font color: Auto

Formatted: Font: 10 pt, Font color: Blue

Formatted: Font color: Auto

Formatted: Font color: Auto

Experiment	Total number of samples	Number of warm-up samples	Number of chains	Time *
Toy dataset	4000	2000	4	~1h
European synthetic datasets	12000	6000	4	~7h
Filters synthetic dataset	6000	3000	4	~3h
Filters real-world dataset	6000	3000	4	~1.5h

* The time here shows the BAMF+HS running time plus the post-processing time (for sorting, averaging, etc.), the latter having a nearly neglectable contribution. Model runtimes were measured as wall-clock time on a high-performance computing cluster. Jobs were executed on compute nodes equipped with dual-socket AMD EPYC 7502 processors (64 physical cores, 2.5 GHz) and approximately 256 GB of RAM, running Linux. All BAMF+HS runs used four parallel chains using approximately 4 physical CPU cores (8 logical threads).

Experiment	Total number of samples	Number of warm-up samples	Number of chains
Toy dataset	4000	2000	4
European synthetic datasets	12000	6000	4
Filters synthetic dataset	6000	3000	4
Filters real-world dataset	6000	3000	4

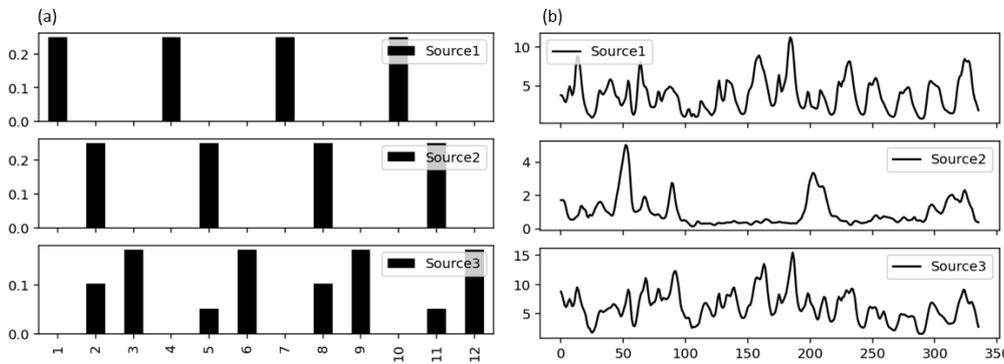
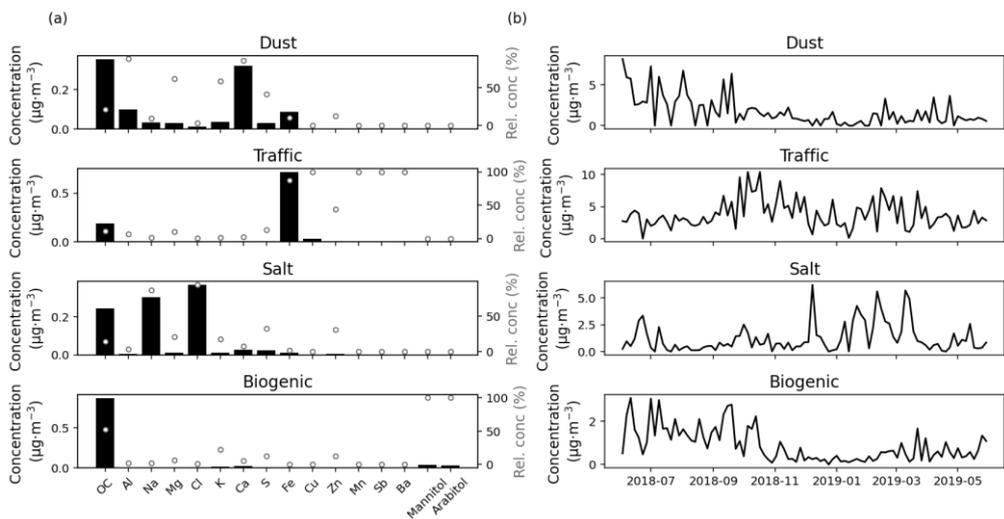
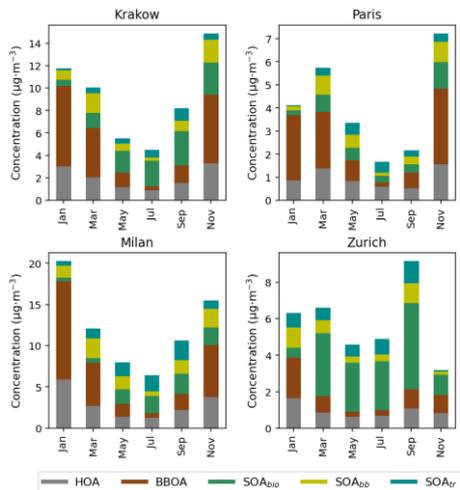


Figure S2. Toy dataset (a) Profiles. (b) Time series.

Formatted



36
37 **Figure S3. Synthetic offline dataset (a) Profiles. (b) Time series.**



38
39 **Figure S4. Generated synthetic datasets source mean concentrations for 4 European cities.**

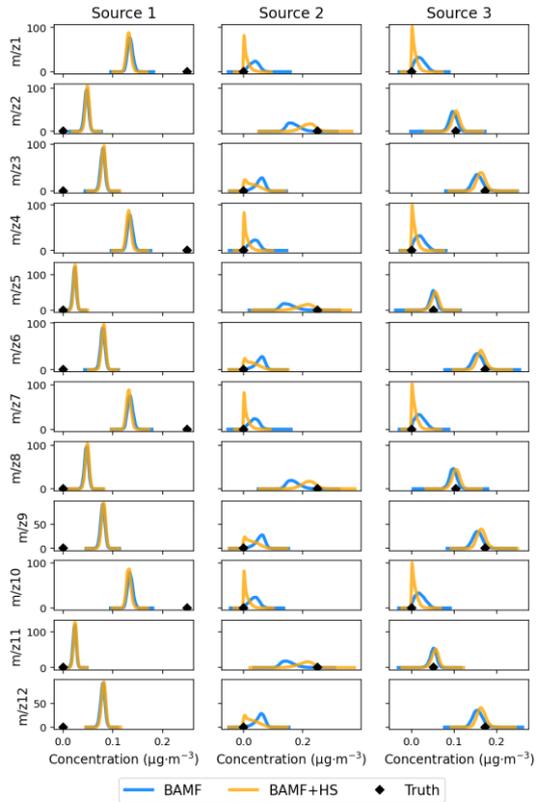
40

41 **Table S2. Profiles employed for the 4-cities synthetic datasets.**

City	Source	Citation
Krakow	HOA	Mohr et al. (2012)
	BBOA	Tobler et al. (2021)
Milan	HOA	Via et al. (2021)
	BBOA	Daellenbach et al. (2021)
Paris	HOA	Crippa et al. (2011)
	BBOA	Zhang et al. (2022)
Zurich	HOA	Elser et al. (2016)
	BBOA	Ulbrich et al. (2002), Ulbrich et al. (2022)
	SOA _{bio}	Daellenbach et al. (2017)
	SOA _{bb}	Ulbrich et al. (2002)
	SOA _{tr}	Sage et al. (2007)

42

43

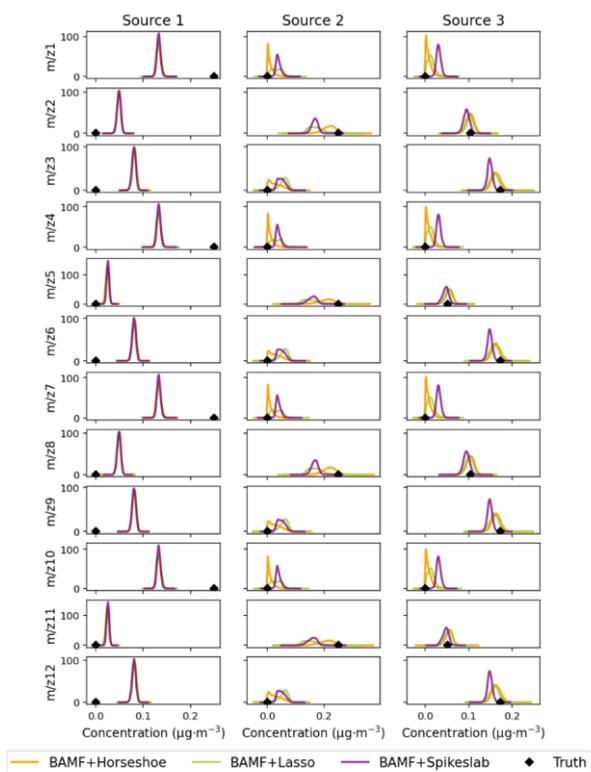


44
45 **Figure S5. Toy dataset distribution of F components for BAMF, BAMF+HS.**

46 **Table S3. Statistics of the toy dataset reconstruction and comparison to truth for the different shrinkage alternatives.**

Model	Sources	X			G			F		
		R2	Median(Z-X /sigma)	Max(Z-X /sigma)	G/G ₀	R ²	ρ	R ²	Spars. ratio	Gini ratio
BAMF+HS	Source1				1.9633	0.9850	0.8193	0.7551	0.0	0.40
	Source2				1.2786	0.9499	0.8193	0.9940	0.5	0.86
	Source3	0.9689	0.2107	1.2863	0.5129	0.4781	0.9608	0.9994	1.0	0.93
BAMF-Lasso	Source1				1.7579	0.9771	0.8193	0.7610	0.0	0.40
	Source2				2.0186	0.9491	0.8193	0.9513	0.0	0.50
	Source3	0.9816	0.1239	0.6964	0.4910	0.4609	0.9608	0.9996	0.0	0.86
BAMF-Spike-Slab	Source1				1.5725	0.9707	0.8193	0.7746	0.0	0.40
	Source2				1.7830	0.9490	0.8193	0.9936	0.0	0.52
	Source3	0.9814	0.1262	0.7094	0.6215	0.4954	0.9608	0.9829	0.0	0.69

47

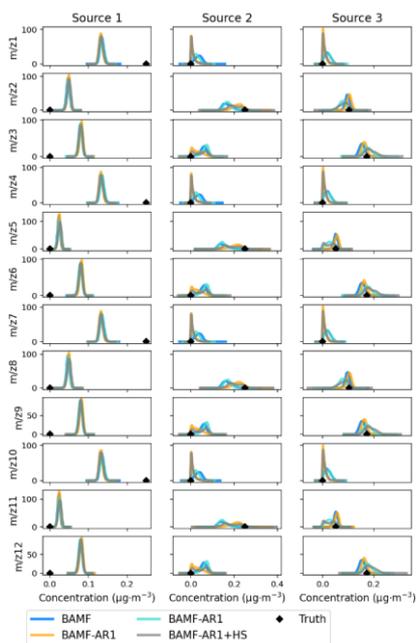
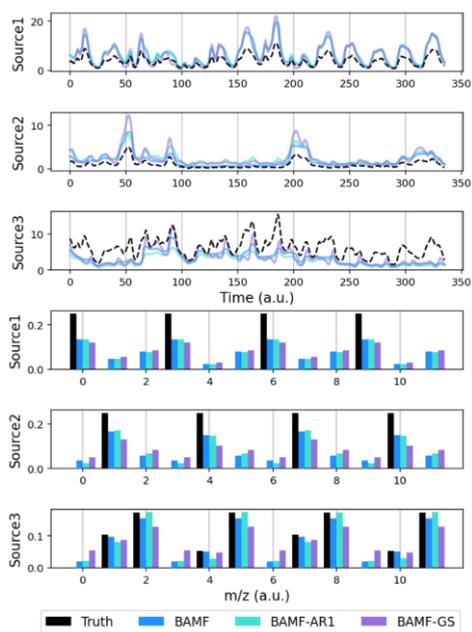


48

49

Figure S6. Comparison of toy dataset distribution of F components for BAMF+HS, BAMF+Lasso, BAMF+Spike-and-slab.

50

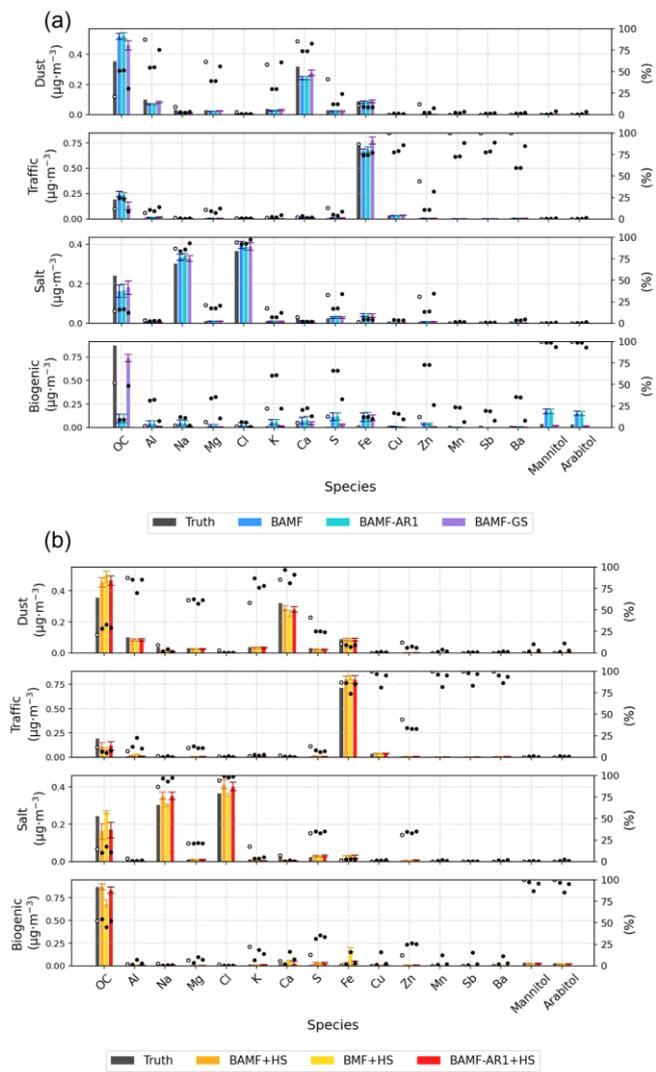


51

52

53

Figure S7. Left: Toy dataset time series (top) and profiles (bottom) results with BAMF, BAMF-AR1 and BAMF-GS. Right: Toy dataset distribution of F components for BAMF, BAMF+HS, BAMF-AR1, BAMF-AR1+HS.



54

55 **Figure S8. Synthetic offline dataset comparison between (a) BAMF, BAMF-AR1, BAMF-GS (b) BAMF+HS, BMF+HS, BAMF-**
 56 **AR1+HS.**

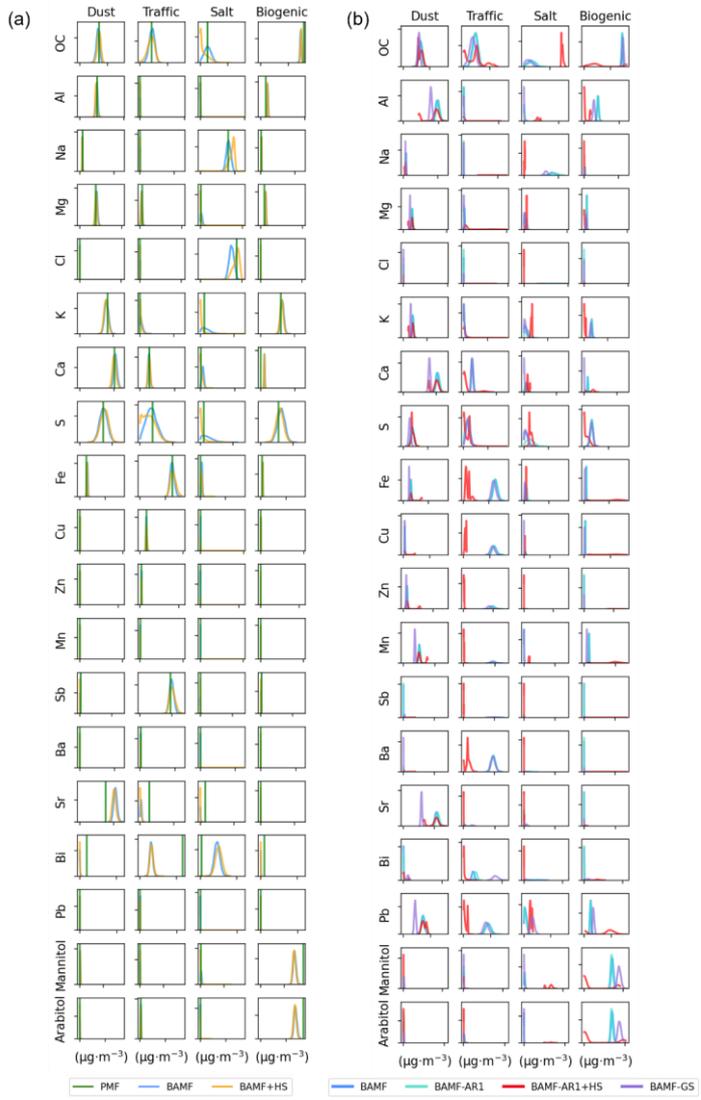
57

58 **Table S4. Purely-measurement-based offline synthetic dataset factorisation statistics. The last three rows correspond to the sum of**
 59 **the metrics for all all factors, except for the G/G₀ metric column, which shows the sum of the factors G/G₀ deviation to 1.**
 60

Model	Factor	F R ²	F ρ	F contr. R ²	F contr. ρ	F Gini	G R ²	G/G ₀ (*)
Dust	PMF	0.98	0.90	0.93	0.90	0.77	0.95	1.00
	BAMF	0.99	0.97	0.78	0.93	0.74	0.92	0.98
	BAMF+HS	0.96	0.96	0.96	0.86	0.77	0.93	1.02
Traffic	PMF	0.99	0.67	0.87	0.88	0.89	0.94	0.83
	BAMF	0.97	0.77	0.46	0.75	0.81	0.93	1.02
	BAMF+HS	0.96	0.85	0.85	0.94	0.88	0.94	0.75
Salt	PMF	0.90	0.60	0.95	0.63	0.83	0.77	0.92
	BAMF	0.90	0.92	0.71	0.63	0.63	0.30	0.11
	BAMF+HS	0.78	0.82	0.92	0.55	0.86	0.52	0.86
Biogenic	PMF	0.96	0.27	0.77	0.16	0.88	0.96	0.68
	BAMF	0.03	0.85	0.56	0.21	0.51	0.96	0.14
	BAMF+HS	0.999	0.63	0.85	0.41	0.89	0.96	0.65
\sum_k (*)	PMF	3.85	2.44	3.52	2.57	3.37	3.64	0.61
	BAMF	2.89	3.51	2.51	2.52	2.69	3.11	1.79
	BAMF+HS	3.70	3.26	3.58	2.76	3.40	3.35	0.78

61

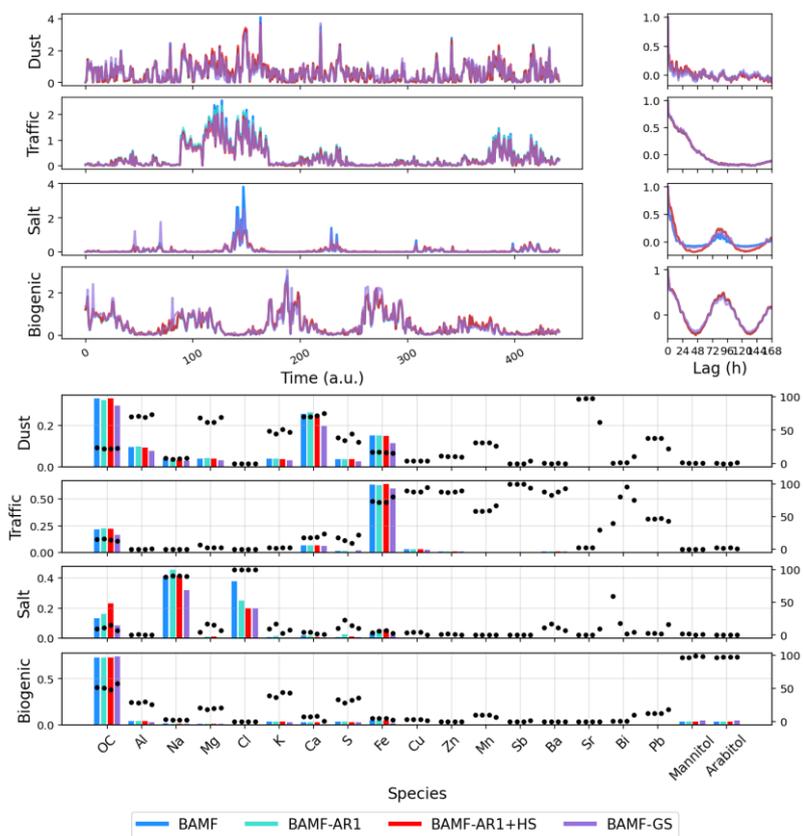
62
63



64

65
66
67
68
69

Figure S9. Filters real-world filters dataset profiles components distributions for (a) PMF, BAMF, BAMF+HS; (b) BAMF, BAMF-AR1, BAMF-AR1+HS, BAMF+GS.



70
71
72
73
74
75

Figure S10. Comparison of BAMF, BAMF-AR1, BAMF-AR1+HS, BAMF-GS on the real-world filters dataset. Plots from left to right and from the top to the bottom are: Time series, autocorrelation, profiles.

Table S5. Online European synthetic datasets reconstruction and factorisation statistics. Statistics contain data from all sites, sources, and datasets.

Statistic	Model	Mean	Median	Min	Max
X - Z /σ	PMF	4.24	4.34	1.45	7.17
	BAMF	0.64	0.66	0.04	2.04
	BAMF+HS	2.92	2.26	0.66	8.45
G/G ₀	PMF	2.82	1.10	0.39	18.45
	BAMF	1.07	0.92	0.02	7.21
	BAMF+HS	2.00	1.11	0.19	46.86
G R ²	PMF	0.78	0.91	0.01	1.00
	BAMF	0.84	0.95	0.04	1.00
	BAMF+HS	0.83	0.92	0.05	1.00
F ρ	PMF	0.85	0.87	0.53	1.00
	BAMF	0.88	0.91	0.27	0.99
	BAMF+HS	0.87	0.90	0.53	1.00
F R ²	PMF	0.84	0.92	0.04	1.00
	BAMF	0.90	0.96	0.02	1.00
	BAMF+HS	0.88	0.95	0.16	1.00
F sparsity	PMF	2.84	2.69	1.69	5.14
	BAMF	2.61	2.45	1.53	4.53
	BAMF+HS	2.63	2.45	1.65	4.53
F Gini	PMF	0.59	0.59	0.45	0.72
	BAMF	0.57	0.56	0.05	0.75
	BAMF+HS	0.57	0.58	0.41	0.83

F Gini ratio	PMF	0.99	0.96	0.75	1.46
	BAMF	0.97	0.99	0.09	1.15
	BAMF+HS	0.98	0.97	0.72	1.40

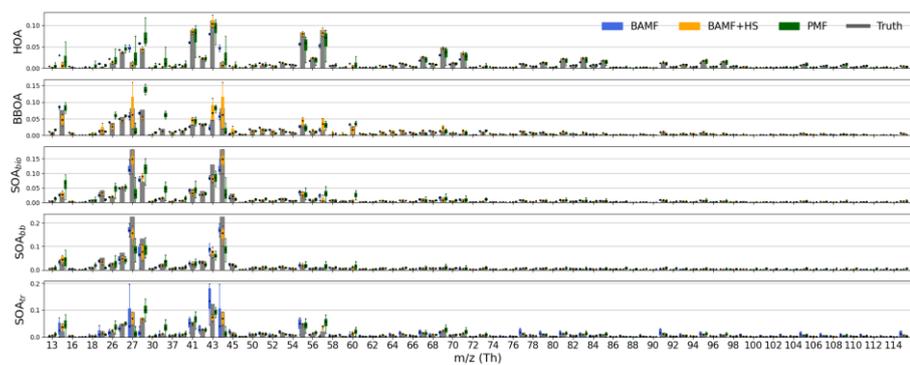
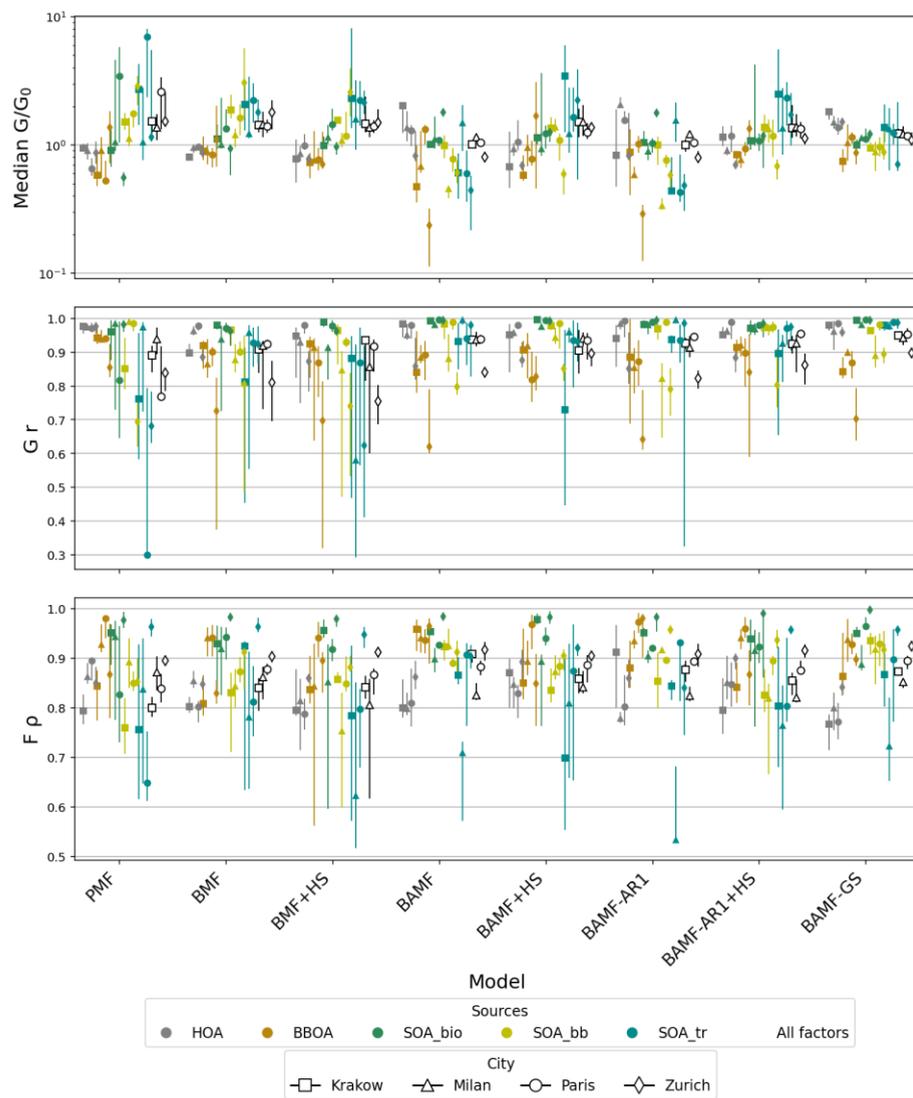


Figure S11. Example results obtained with BAMF, BAMF+HS, and PMF for one of the 24 datasets included in the ACSM-like experiment. The results correspond to Zurich, covering the period from 1 January 2019 to 14 January 2019 (dataset 0).

Formatted: Space After: 0 pt, Line spacing: 1.5 lines
 Border: Top: (No border), Bottom: (No border), Left: (No border), Right: (No border), Between : (No border)

Formatted: Font: 10 pt



81

82 **Figure S124. Summary metrics for all synthetic datasets.**

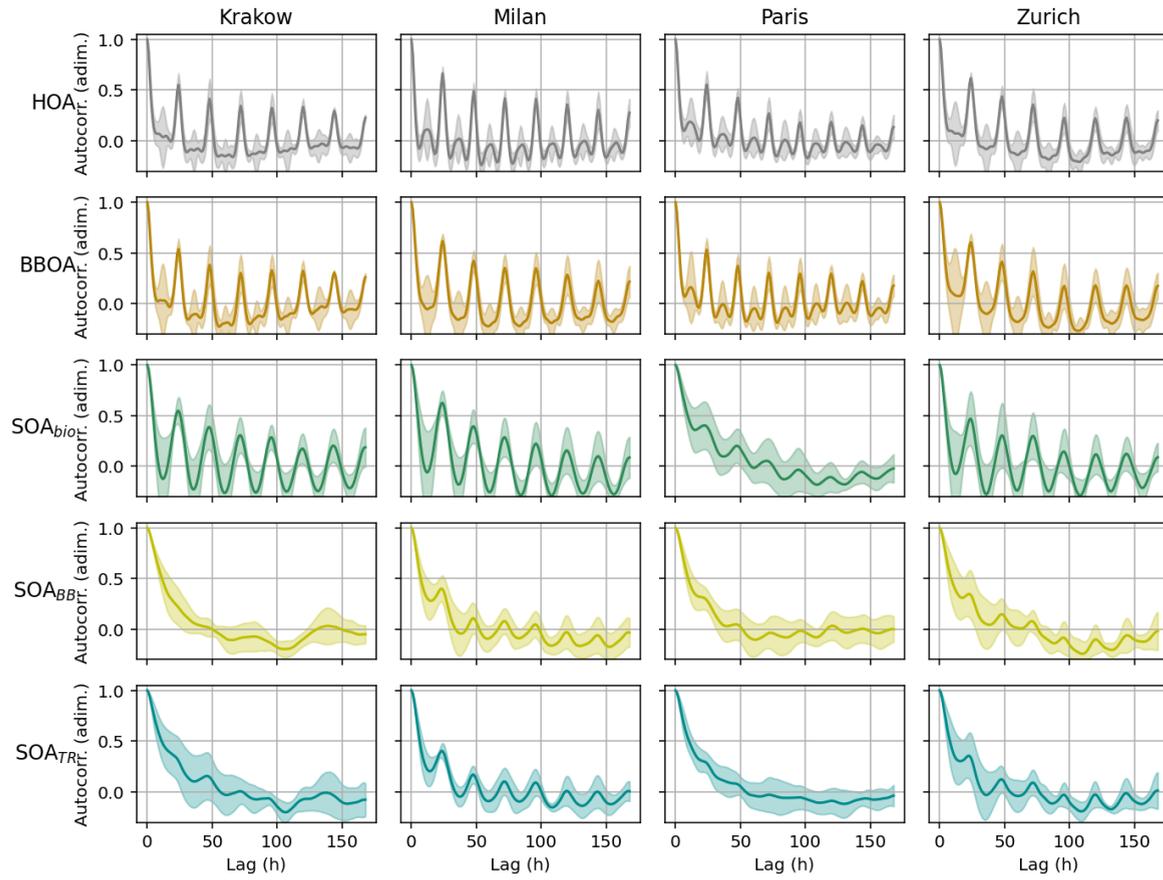
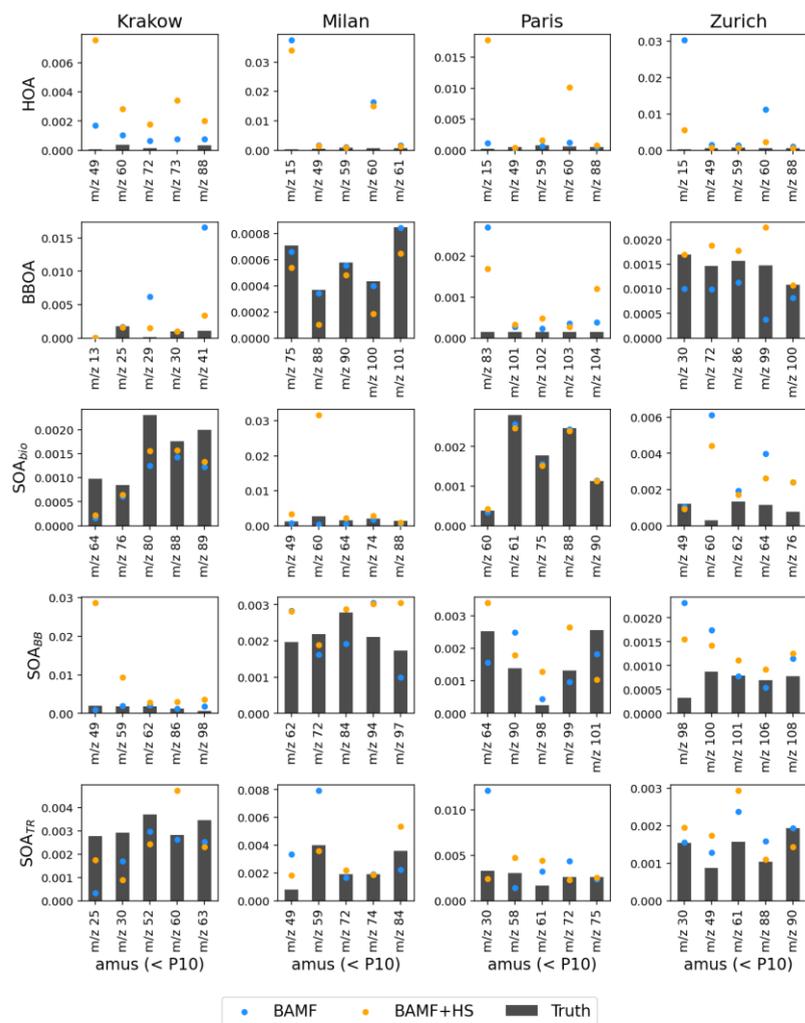


Figure S132. Autocorrelation means (solid lines) and standard deviations (shaded areas) of the 6 datasets per city and source.

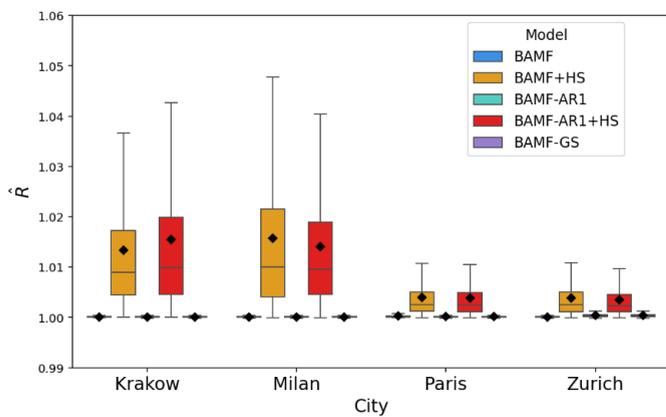


86

87

88

Figure S143. Sparsity introduction results for BAMF and BAMF+HS models in the European synthetic datasets for the 5 less massive m/zs, below the percentile 10. The bars show the truth concentration for the given m/z, and BAMF, BAMF+HS results are shown in markers.



89
90
91 **Figure S154.** \hat{R} for the different cities and models. Model boxplots are in the same order as the legend.
92

93 **B) Supplementary model formulation**

94 **Lasso**

95 The Lasso (Least Absolute Shrinkage and Selection Operator, Rasmussen et al. 2012) is a regularization method that adds an L1 penalty to linear
96 regression, shrinking some coefficients exactly to zero. The L1 penalty is a way to discourage large coefficients in a regression model by adding up
97 the absolute values of all the model's coefficients and including that total in the cost the model tries to minimize. This encourages sparsity in the
98 model, making it useful for variable selection. However, it can over-shrink large coefficients and it treats all coefficients equally, regardless of their
99 signal strength, hence it tends to over-shrink even large signals.

100 **Spike-and-slab**

101 The spike-and-slab prior (Andersen et al. 2014) is a Bayesian approach to variable selection that models each coefficient as coming from a mixture
102 of two distributions: a spike at zero (forcing sparsity) and a slab (a wider distribution allowing nonzero values). It provides strong interpretability
103 and explicit variable inclusion, but is computationally intensive due to its discrete model space. This working scheme provides a binary shrinkage

104 behaviour since the signals are considered either zero or slabs, which enforces 0-like signals, however, the selection of non-zero elements can be
105 too harsh.

106

107 **Horseshoe**

108 The regularized horseshoe prior is an extension of the standard horseshoe prior designed to improve robustness and regularization in sparse Bayesian
109 models. Like the original horseshoe, it combines global shrinkage (for overall sparsity) and local shrinkage (for individual coefficients), using heavy
110 tails to allow large signals while aggressively shrinking noise. What sets the regularized version apart is the addition of a slab component that limits
111 the influence of extremely large coefficients. This slab acts like a soft ceiling, preventing the model from over-trusting outlier variables while
112 maintaining flexibility. This makes the prior more stable in practice, especially in finite data settings or when outliers are present. In comparison to
113 the two above priors, the regularized horseshoe provides a middle ground: it allows for adaptive shrinkage, handles strong signals well, and includes
114 a slab to avoid the horseshoe's instability. It's more robust and scalable than spike-and-slab and more flexible and statistically principled than Lasso.

115

116 **C) Supplementary model evaluation**

117 **C.I. Chemically sparse toy dataset**

118

119 Some alternative autocorrelation formulations to the regular BAMF autocorrelation priors were tried for the toy dataset. BMF, BMF+HS, and BMF-
120 GS – analogs to BAMF, BAMF+HS, BAMF-GS, respectively, but without the autocorrelation term – provided unstable results, with the highest $|\mathbf{Z}$ -
121 $\mathbf{X}|/\sigma$ median and maxima and the poorest \mathbf{X} vs. \mathbf{Z} correlation coefficients, hence, these were discarded for this discussion. Figure S7 shows the time
122 series, profiles and \mathbf{F} components distributions for BAMF, BAMF-AR1, BAMG+GS, BAMF+HS. Both regarding time series and profiles, BAMF
123 seems to be resembling the truth equally or better than BAMF-AR1 and BAMF-GS with the only exception of the time series of Source 3, which is
124 visually better captured by BAMF-GS. However, BAMF-GS profiles are further from the truth and m/z mass closure for this model is much worse
125 than regular BAMF, as also seen by poorer Spearman and Pearson correlation coefficients. Profiles as captured by BAMF-AR1 are not consistently
126 better than those from BAMF either. BAMF appears as the most balanced and accurate model of the three for the presented toy dataset. Regarding
127 \mathbf{F} components distributions, BAMF and BAMF-AR1 and BAMF+HS and BAMF-AR1+HS BAMF-AR1, BAMF-AR1+HS are compared in the
128 right panel of Figure S7. BAMF and BAMF-AR1 show very similar distributions, with the slight differences mentioned before. BAMF+HS and
129 BAMF-AR1+HS present some differences, and even if the shrinkage is also present in BAMF-AR1+HS (although with weaker Gini than

130 BAMF+HS), some components present distributions with a more pronounced multimodal behaviour than the BAMF+HS. Therefore, the BAMF-
131 AR1+HS is shown to be less precise and accurate than BAMF+HS. The BAMF-GS+HS was not tried out since the high concentrations of the
132 components of the non-normalised \mathbf{F} matrix on this model hinder their shrinkage to zero.
133

134 **C.II. Chemically sparse offline synthetic dataset**

135
136 Figure S8 shows the profiles of the \mathbf{F} matrix for BAMF, BAMF-AR1, BAMF-GS (a) and BAMF+HS, BMF+HS, BAMF-AR1+HS (b) models.
137 BAMF, BAMF-AR1, and BAMF-GS were run, initialising \mathbf{F} as a normal distribution to make the sampling more sturdy to avoid initialization
138 failure. The outcomes of these models show general good agreement with the truth, with BAMF+HS being the most accurate model followed by
139 BAMF-AR1+HS. The sum across factors of profile Pearson R^2 with truth for BAMF+HS, BMF+HS, BAMF-AR1+HS, and BAMF-GS were 3.90,
140 3.80, and 3.90, proving the beneficial effect of the horseshoe prior in the description of sparse profiles. The BAMF-AR1+HS model is performing
141 slightly worse than the BAMF+HS in time series (G/G_0 deviation sum of 0.85, 0.91, respectively). Regarding the non-horseshoed models, the sum
142 of the correlation with truth \mathbf{G} for all factors was very similar for BAMF, BAMF-AR1, GS (3.93, 3.92, 3.93, respectively) but the Spearman
143 correlation coefficient with truth \mathbf{F} was better for BAMF (3.75, 3.71, 3.67, respectively), highlighting a slight better performance for BAMF. Hence,
144 with all, the BAMF+HS seems the most accurate model, benefitting both from its autocorrelation and sparsity properties.
145

146 **C.III. Chemically sparse offline real-world dataset**

147 Figure S10 shows the performance of the other models discussed in the Toy dataset and European city datasets sections, BAMF-AR1, BAMF-
148 AR1+HS, BAMF-GS, in comparison to BAMF, used as the base case. All models employed seem to agree, providing overlapping time series and
149 autocorrelations, and similar profiles, with the exception of BAMF-AR1+HS, which is slightly differentiated from the others. However, it does not
150 provide sparsified profiles in the species in which the shrinkage effect was found for BAMF+HS or elsewhere. This is more evidently depicted in
151 Figure S9 (b), in which the distributions of \mathbf{F} components are shown. The BAMF-AR1+HS does not present the expected sparsifying shape except
152 for Ca, K, and some elements in the traffic and salt profiles. In general, the BAMF-AR1+HS distributions are more multimodal, reflecting higher
153 divergence across HMC chains which makes results much less sturdy. Regarding the rest of models, the distributions are very similar ensuring the
154 robustness of all three and stability of the solution. Therefore, all Bayesian models and PMF point to a similar, robust solution for the filters dataset.
155 The horseshoe prior addition to BAMF, though, provides here a useful sparsity introduction in the current dataset which helps purify the profiles

156 from unwanted profile entanglement. However, the implementation of this prior in the BAMF-AR1 model detracts the solution due to HMC
157 chain divergence.
158