Advecting Superspecies

Efficiently Modeling Transport of Organic Aerosol With a Mass-Conserving Dimensionality Reduction Method

Sturm, Patrick Obin; Manders, Astrid; Janssen, Ruud; Segers, Arjo; Wexler, Anthony S.; Lin, Hai Xiang

# Advecting Superspecies: Efficiently Modeling Transport of Organic Aerosol With a Mass-Conserving Dimensionality Reduction Method

**Patrick Obin Sturm[1,2]** (ID)**, Astrid Manders[3], Ruud Janssen[3], Arjo Segers[3]** (ID)**, Anthony S. Wexler[1,4], and Hai Xiang Lin[2,5]** (ID)

[1]Air Quality Research Center, University of California, Davis, Davis, CA, USA, [2]Delft Institute of Applied Mathematics, Delft University of Technology, Delft, The Netherlands, [3]Department of Climate, Air and Sustainability, TNO, Utrecht, The Netherlands, [4]Departments of Mechanical and Aerospace Engineering, Civil and Environmental Engineering, and Land, Air and Water Resources, University of California, Davis, Davis, CA, USA, [5]Institute of Environmental Sciences, Leiden University, Leiden, The Netherlands

**Abstract** The chemical transport model LOTOS-EUROS uses a volatility basis set (VBS) approach to represent the formation of secondary organic aerosol (SOA) in the atmosphere. Inclusion of the VBS approximately doubles the dimensionality of LOTOS-EUROS and slows computation of the advection operator by a factor of two. This complexity limits SOA representation in operational forecasts. We develop a mass-conserving dimensionality reduction method based on matrix factorization to find latent patterns in the VBS tracers that correspond to a smaller set of superspecies. Tracers are reversibly compressed to superspecies before transport, and the superspecies are subsequently decompressed to tracers for process-based SOA modeling. This physically interpretable data-driven method conserves the total concentration and phase of the tracers throughout the process. The superspecies approach is implemented in LOTOS-EUROS and found to accelerate the advection operator by a factor of 1.5–1.8. Concentrations remain numerically stable over model simulation times of 2 weeks, including simulations at higher spatial resolutions than the data-driven models were trained on. The reversible compression of VBS tracers enables detailed, process-based SOA representation in LOTOS-EUROS operational forecasts in a computationally efficient manner. Beyond this case study, the physically consistent data-driven approach developed in this work enforces conservation laws that are essential to other Earth system modeling applications, and generalizes to other processes where computational benefit can be gained from a two-way mapping between detailed process variables and their representation in a reduced-dimensional space.

**Plain Language Summary** The chemical composition of the atmosphere is a complex system involving many physical processes. Computer models can be used to improve our understanding of how these processes interact, as well as simulate hypothetical scenarios to support scientifically-informed climate and air quality policies. However, complicated models with many variables can take a lot of time to run. The LOTOS-EUROS model spends a large fraction of time and computational resources on simulating the transport of chemical species, like particulate matter, by wind. We combine data-driven approaches with domain knowledge to reduce the number of variables while ensuring essential properties are conserved: we model representative combinations of chemical species that are transported all at once, rather than transport each species individually. This leads to faster and cheaper simulations without loss of scientific detail or internal consistency.

## 1. Introduction

Vast amounts of computational resources are required to model phenomena in the Earth sciences. This includes complex models of atmospheric composition that couple a large number of properties and processes (Brasseur & Jacob, 2017). Data-driven approaches, including machine learning (ML), are an emerging set of techniques for decreasing the computational burden of Earth System Models (ESMs) by using more efficient parameterizations, but have documented challenges such as unstable error growth and physical inconsistency when predicted recurrently (Kelp et al., 2018) or when interacting with other processes in the context of larger models (Brenowitz & Bretherton, 2019; Rasp et al., 2018). One approach towards data-driven ML models that can stably interact with

**Methodology:** Patrick Obin Sturm, Astrid Manders, Ruud Janssen, Arjo Segers
**Project Administration:** Astrid Manders, Arjo Segers, Hai Xiang Lin
**Resources:** Astrid Manders, Arjo Segers
**Software:** Patrick Obin Sturm, Arjo Segers
**Supervision:** Astrid Manders, Ruud Janssen, Arjo Segers, Anthony S. Wexler, Hai Xiang Lin
**Validation:** Patrick Obin Sturm, Astrid Manders, Ruud Janssen, Arjo Segers
**Visualization:** Patrick Obin Sturm, Astrid Manders, Ruud Janssen, Hai Xiang Lin
**Writing – original draft:** Patrick Obin Sturm, Anthony S. Wexler
**Writing – review & editing:** Patrick Obin Sturm, Astrid Manders, Ruud Janssen, Anthony S. Wexler, Hai Xiang Lin

other model processes is online training: parameter optimization of neural network surrogates while running the entire model (Kelp et al., 2022; Rasp, 2020).

Other recent efforts have aimed to constrain data-driven approaches using domain knowledge to ensure physically consistent results. One strategy for physically consistent data-driven models reposes the learning targets: rather than estimate important properties or their tendencies, instead estimate fluxes between the properties. The fluxes can then be related to tendencies in a way that balances mass, energy, or atoms (Sturm & Wexler, 2020, 2022; Yuval, O'Gorman, et al., 2021). Custom neural network architectures can also obey conservation laws by incorporating hard constraints in their hidden layers (Beucler et al., 2021), such as flux balances (Sturm & Wexler, 2022): this can also improve the physical interpretability of the inner working of neural networks. Though physical consistency is an important result by itself, imposed constraints do not necessarily improve accuracy of such tools beyond adherence to whichever physical law(s) the constraints enforce. For example, Harder et al. (2022) found the accuracy of a neural network surrogate model of aerosol microphysics was not improved when adding a completion constraint during training, where a chosen variable was reassigned to the sum of all other variables' tendencies to conserve mass. However, constraints can be implemented in ways that add domain knowledge to the data-driven algorithm: Sturm and Wexler (2022) found that by adjusting a feed-forward neural network architecture to include a flux-tendency constraint during training, the overall prediction accuracy of chemical species concentrations improved. Kelp et al. (2020) motivated ML model architectures with built-in assumptions about the physical system as a future research direction. In the approach in Sturm and Wexler (2022) the constraint gives information on the graph relational structure of a chemical mechanism, that is, how different chemical species interact. Recent work toward physically consistent data-driven tools in the Earth sciences, and acknowledgment of their importance (Keller & Evans, 2019; Sturm & Wexler, 2020; Yuval, Pritchard, et al., 2021) has motivated the mass-conserving dimensionality reduction method in this paper.

Within the field of atmospheric chemistry modeling, Kelp et al. (2020) have made progress towards a stable neural network emulating a box model of chemistry and aerosol microphysics processes, through training parameters on the accuracy of multiple future timesteps after predicting in a lower-dimensional latent space. Kelp et al. (2020) pose a future research direction: how the low-dimensional representation of chemical species might interact with other processes, such as advection, in the context of a larger model. The scope of the present study is informed by this direction: we develop and assess a physically consistent data-driven method that compresses the high dimensional set of organic aerosol (OA) tracers to reduce the computational cost of advection in the LOTOS-EUROS chemical transport model (CTM) (Manders et al., 2017). LOTOS-EUROS is a state-of-the-art model that has been compared to the WRF-Chem, CAMx, CMAQ, and EMEP models in several international model intercomparison studies such as AQMEII (Im et al., 2015) and EURODELTA-Trends (Colette et al., 2017) and is part of the European Copernicus Atmospheric Monitoring (CAMS) model ensemble. The advection operator consumes a significant amount of wall time in LOTOS-EUROS, from about 20% of total wall time in sequential runs (only chemistry and sometimes deposition calculations take longer) to over 50% of total wall time in parallel runs using domain decomposition. Wall time of advection can double with the inclusion of organic aerosol tracers (Sturm, 2021). Therefore, the current default in LOTOS-EUROS is to include only one passive tracer for organic aerosol, which is one of the reasons for an underestimation of total particulate matter (Timmermans et al., 2022).

Organic aerosol forms an important contribution to particulate matter (Jimenez et al., 2009). OA can be emitted to the atmosphere as semi-volatile primary organic aerosol (POA) through various direct sources, including vehicle exhaust, wildfire smoke, and residential wood combustion. OA can also be formed in the atmosphere as secondary organic aerosol (SOA) through gas-phase reactions of volatile organic compounds (VOCs), which tend to form less volatile products: intermediate volatility organic compounds (IVOC) and semi-volatile organic compounds (SVOC), referred to together as siVOC. SVOC can partition appreciably to the particle phase under ambient conditions. Both anthropogenic sources, like industrial activity, and biogenic sources, such as forests, emit precursors of SOA. Another source of SOA is the partial evaporation of POA to siVOCs, which in turn react and partition to form SOA (Robinson et al., 2007). This SOA from evaporated and aged POA is often chemically distinct from POA, showing a higher degree of oxidation (Jimenez et al., 2009), and can be tracked separately in models. SOA can form a significant fraction of the total OA concentration (de Gouw et al., 2005; Heald et al., 2005).

Due to the large number of distinct organic species in the atmosphere, organic aerosols are often lumped together into volatility bins according to the magnitude of their saturation vapor pressures (Donahue et al., 2006). This

modeling approach is called the volatility basis set (VBS) and accounts for the tendency of compounds to become less volatile as they are oxidized. The partitioning between gas and particle phase in each volatility bin is governed by its corresponding saturation vapor pressure and the total OA concentration. A 2D-VBS extension has been developed that includes oxygen to carbon ratio along another dimension (Donahue et al., 2011; Jimenez et al., 2009), which can account for fragmentation of larger compounds and estimation of hygroscopicity (Jimenez et al., 2009). A 1D-VBS approach is commonly applied in chemical transport models, including separate basis sets for different classes of OA precursors (Bergström et al., 2012; Hayes et al., 2015; Janssen et al., 2017; Jiang et al., 2019). Use of multiple VBS classes enables distinct properties per class and can give insight into different aerosol systems contributing to total OA. Recent SOA modeling work has concentrated on several topics: (a) further specification of IVOC emissions from specific sources like gasoline and diesel (Jathar et al., 2014; Lu et al., 2020; Ots et al., 2016) and biomass burning (Ciarelli et al., 2017; Jiang et al., 2019; Theodoritsi & Pandis, 2019), (b) effect of aerosol water content on OA partitioning (Pye et al., 2017), (c) the role of SVOC deposition (Knote et al., 2015) and (d) other OA formation pathways, such as reactive uptake of isoprene epoxides (Marais et al., 2016; Nagori et al., 2019; Pye et al., 2013). Hodzic et al. (2016) and Pai et al. (2020) provide a global scale synthesis of some of these ideas.

The inclusion of such detailed, high-dimensional process-based OA models in 3D models is limited by their increased computational burden, for example, to chemical transport models like LOTOS-EUROSv2.2.1 (Manders et al., 2017). Next to an implementation with a single passive organic aerosol tracer, LOTOS-EUROS has an implementation with four VBS classes based on the configuration from Bergström et al. (2012).

Though this approach does not resemble the modern state of the science as discussed in the previous paragraph, it strikes a balance between complexity and level of realism of OA processes: a four-class VBS approach has a higher level of realism than the two-product model (SORGAM) (Odum et al., 1996; Schell et al., 2001) used by other models in air quality forecasts for Europe (Mircea et al., 2019). New developments tend to increase the complexity of the VBS (e.g., by adding specific basis sets for emission sources such as diesel, gasoline, or biomass burning, or by adding more explicit IVOC oxidation). The current VBS module in LOTOS-EUROS v2.2.1 is not used by default, and when included, significantly increases wall time of simulations. The inclusion of VBS tracers adds computation time to other operators in the model relatively more than OA-specific calculations themselves.

Most notably, the high dimensionality caused by adding 58 VBS tracers adds a computational burden to the advection operator in LOTOS-EUROS v2.2.1, which is based on the mixing-ratio conserving scheme in Walcek (2000). When using the VBS module, the number of advected tracers increases from 46 to 104. Model timing experiments in Sturm (2021) found that wall time for the advection operator can double when using the VBS module. Advection is a bulk process and does not perform OA-specific calculations. This motivates dimensionality reduction for a more parsimonious representation of OA in transport processes: we use unsupervised data-driven approaches to find characteristic regimes of VBS tracers, which are used to form lower-dimensional combinations interpreted as superspecies that require fewer transport calculations. Rather than advecting each tracer separately, we instead advect a smaller set of superspecies, which are subsequently mapped back to the OA tracer space after advection. Constraints are applied when compressing to and decompressing from the reduced-dimension space to conserve mass to machine precision. We compare the linear and additive method of non-negative matrix factorization to a nonlinear and more complex neural network autoencoder, and make a model selection after evaluating several configurations based on reconstruction accuracy and physical consistency. Though demonstrated for compression of OA and related compounds during transport to accelerate air quality forecasts over the European continent, the methods developed in this work generalize to other Earth system applications, enabling use of high-dimensional process models whose variables can be reversibly compressed to a physically consistent reduced-dimension representation for use in other processes.

Section 2 outlines the VBS configuration in LOTOS-EUROS and develops four data-driven approaches. These four approaches are tested in Section 3: first, they are trained on the volatility distributions from LOTOS-EUROS model output, then evaluated on reconstruction accuracy of the volatility distributions and physical consistency. One approach from Section 3 is chosen to be implemented in LOTOS-EUROS, with results from various experiments shown in Section 4. More specifically, Section 4 investigates the accuracy of using the superspecies, generalizability of the converged model to other seasons and different spatial resolutions, and corresponding speedup in the 3D model. Section 5 contains a summary of the methods and key results.
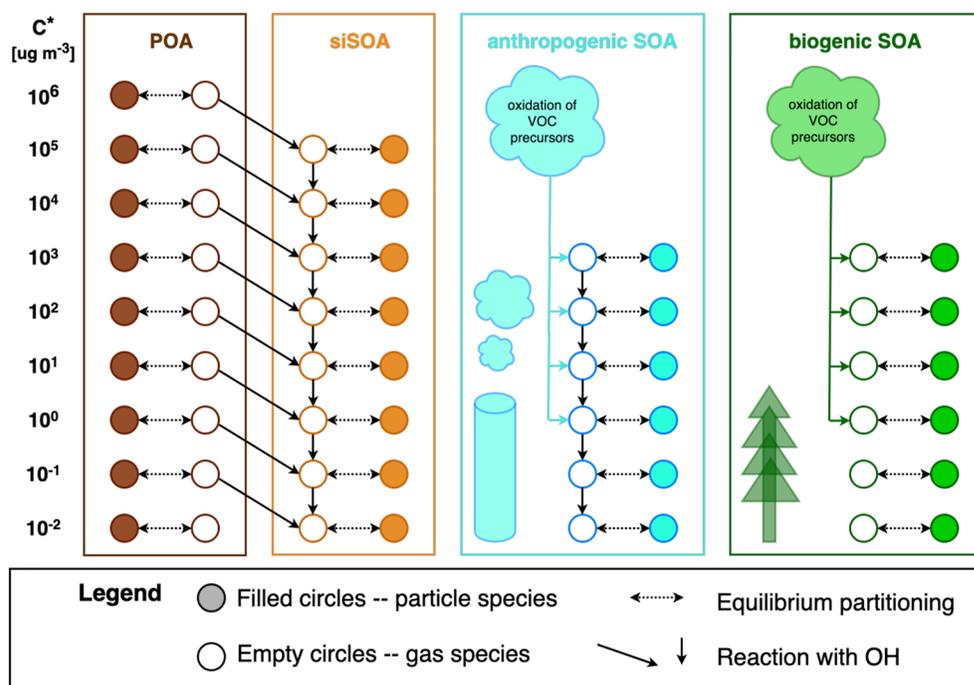
**Figure 1.** Schematic representation of the VBS approach in LOTOS-EUROS v2.2.1, including the 4 VBS classes with 58 tracers, and their thermodynamic and chemical relationships. This diagram was inspired by the schematic in Shrivastava et al. (2008).

## 2. Methods

This section develops four data-driven approaches to reversibly compress VBS-specific tracers. Section 2.1 describes the VBS approach in LOTOS-EUROS v2.2.1. Section 2.2 summarizes several other methods for tracer compression. Section 2.3 details the model configuration used for experiments, as well as the model output used to train the various data-driven approaches. Sections 2.4, 2.5, and 2.6 develop four approaches that are summarized in Section 2.7.

### 2.1. VBS Approach in LOTOS-EUROS

The chemical transport model LOTOS-EUROS v2.2.1 uses a VBS scheme visualized in Figure 1 based on Bergström et al. (2012). This scheme has 4 distinct VBS classes: (a) POA, (b) SOA from siVOCs that are chemically aged after evaporating from semi-volatile POA emissions, (abbreviated as siSOA), and SOA from (c) anthropogenic and (d) biogenic gaseous VOCs, abbreviated as aSOA and bSOA respectively.

Figure 1 provides an overview of the 58 tracers specific to the VBS module. Primary organic material (POM) emissions are modeled using a 9-bin VBS approach: the logarithmically distributed bins represent semi- and intermediate-volatile organics with effective saturation concentrations ranging from $10^{-2}$–$10^6$ μg m$^{-3}$ at 298 K. The reported mass of primary emissions is distributed over the lower 4 volatility bins. As in previous work (Shrivastava et al., 2008), an additional 1.5 times this mass is distributed over the highest 5 volatility bins to represent non-reported intermediate volatility organic compounds (IVOCs). The factor of 1.5 for the VBS in LOTOS-EUROS v2.2.1 is an oversimplification: alternative approaches exist for estimating IVOC emissions from specific sources (e.g., Ciarelli et al., 2017; Jiang et al., 2019; Lu et al., 2020; Ots et al., 2016) and further specification of the VBS in future versions of LOTOS-EUROS could include explicit IVOC emissions per source, adding complexity and underscoring the need for reversible compression for use in transport processes. Only a fraction of the emitted primary material remains in the particle phase: the fraction that evaporates is assumed to be SVOC with effective saturation concentrations on the order of $1 < C^* < 10^3$ μg m$^{-3}$ or IVOC with saturation concentrations on the order of $10^4 < C^* < 10^6$ μg m$^{-3}$, defined at 298 K. The S/IVOCs undergo oxidation by the hydroxyl radical OH and enter the distinct siSOA VBS class. As material moves from the POA VBS to the siSOA

VBS, it also moves to lower volatility bins, as shown in Figure 1. The total siSOA is represented by an 8-bin VBS using effective saturation concentrations from $10^{-2}$–$10^5$ μg m$^{-3}$ (defined at 298 K). Each bin uses two tracers, one aerosol and one gas, to represent the partitioning: this results in 18 tracers for the POA VBS class and 16 tracers for the siSOA VBS class. Formation of SOA from anthropogenic VOCs is represented with a 6-bin VBS class, defined using effective saturation concentrations of $10^{-2}$ to $10^3$ μg m$^{-3}$ at 298 K. This results in 12 tracers (6 in the gas phase and 6 in the particle phase). VOCs such as aromatics, alkenes and alkanes are classified in LOTOS-EUROS as anthropogenic precursors of secondary organic aerosols and upon oxidation are distributed over the 4 highest volatility bins as done by Tsimpidi et al. (2010), linearly interpolating between a low-NOx and high-NOx case as originally suggested by Lane et al. (2008).

An analogous 6-bin VBS class is used to model SOA formation from the biogenic VOCs in LOTOS-EUROS: monoterpene and isoprene. Yields from biogenic gaseous precursors are distributed over the 4 highest volatility bins according to Tsimpidi et al. (2010), with yields calculated by a branching ratio continuously dependent on NOx (Lane et al., 2008). Unlike the anthropogenic VBS class, ageing between bins is turned off for the biogenic VBS in LOTOS-EUROS v2.2.1, as in prior work (Matsui, 2017; Murphy & Pandis, 2009; Tsimpidi et al., 2010, 2014). This is informed by the low sensitivity of biogenic SOA concentration to oxidative ageing (Donahue et al., 2012; Ng et al., 2006), thought to arise from fragmentation effects that balance out functionalization effects on volatility (Murphy et al., 2012). For this reason, material never enters the 2 lowest volatility bins in LOTOS-EUROS v2.2.1, rendering the 4 corresponding tracers effectively inert. However, in LOTOS-EUROS v2.2.1 with the VBS module on, these 4 tracers are still dealt with by the model, contributing to the computational burden on processes such as advection.

## 2.2. Tracer Compression Methods

A method for tracer compression for transport in the GEOS-Chem global CTM is given by Liao et al. (2007), where various oxidation products are lumped together by phase and class, and assumed to behave similarly in transport. The relative compositions from each grid cell's previous time step are used to distribute the lumped tracers back to individual products after transport. This can be thought of as compression to a single lumped superspecies with one degree of freedom (the superspecies concentration) and a fixed composition dictated by the grid cell before the advection operator. Another approach for OA tracers given by Matsui (2017) compresses VBS tracers in a global aerosol model from 106 to 26 (a compression factor of approximately 4) by using fewer volatility bins. This effectively lowers the bin resolution and combines material across a wider range of saturation vapor pressures. Analogously, Matsui (2017) converts between high-resolution and low-resolution bins in a sectional aerosol model for use in processes not directly related to aerosols. An example of tracer compression for advection in a 2D-VBS is given by Zhao et al. (2020) who sum tracers along the O:C axis, resulting in a 1D-VBS for decreased dimensionality in advection.

A partitioning-based compression technique for advection of 1D-VBS tracers could be developed based on partitioning, where the compressed tracers themselves contain all the information needed to decompress to the VBS tracer space without loss of accuracy. This technique advects total concentration for each volatility bin as well as total OA concentration, reducing the 58 phase-specific tracers to 29 combined phase tracers and an additional tracer to keep track of total organic aerosol concentration. After advection, total OA along with the saturation vapor concentration determines the partitioning between phase in each volatility bin. However, this theoretical strategy applied to the VBS tracers would yield a compression factor of only approximately 2 (compressing 58 tracers to 30). This would reduce the total number of advected tracers from 104 to 76. We seek a compression technique that can reduce the number of tracers further, leveraging data-driven approaches optimized on a large amount of representative model output, to reversibly compress VBS distributions with minimal accuracy lost.

## 2.3. Model Configuration and Output

To find latent patterns for a reduced order representation of the 58 VBS tracers, we use LOTOS-EUROS version 2.2.1 (Manders et al., 2017; Manders-Groot et al., 2021) with the optional VBS module. The model is used in its default configuration using 5 levels, the first one being a 25 m surface layer, the second layer reaching the top of the mixing layer, and the other three layers being reservoir layers up to 5 km altitude. The horizontal domain covers 15°W to 35°E and 35–70°N on a lonxlat grid of 0.5 × 0.25°. This grid is termed the MACC (Monitoring

Atmospheric Composition and Change) grid, a predecessor of the current CAMS (Copernicus Atmospheric Monitoring Service). Meteorology is taken from ECMWF IFS 12-hr operational forecasts, using hourly surface values and 3-hr 3D fields interpolated to hourly values. The LOTOS-EUROS advection scheme is based on Walcek (2000). The advection operator does not only refer to bulk horizontal transport by wind, but rather advection in 3 directions: the vertical flux is calculated from the net horizontal flux and continuity. Convection is not implemented as an explicit operator. Instead, the impact of convection is implied by changes in the vertical layer of the model with the first two layers together covering the boundary layer. Other vertical transport is represented by an entrainment and detrainment operator where the vertical structure of the grid is adjusted to mixing layer depth then the pollutant concentrations are linearly interpolated, and a separate vertical diffusion operator. For gas-phase chemistry, a condensed and slighty modified version of CBM-IV is used (Gery et al., 1989). Wet deposition includes in-cloud and below-cloud scavenging as described in Seinfeld and Pandis (2006), deposition of gases is calculated using DEPAC (Zanten et al., 2010), and deposition of particles follows Zhang (2001). The model includes tree-specific biogenic isoprene and terpene emissions as described in Beltman et al. (2013) using a high-resolution tree-species database (Köble & Seufert, 2001) that are combined with land cover data from CORINE2000 (EEA, 2005). Anthropogenic emissions are CAMS emissions for 2015 (CAMS regional air pollutants as delivered in 2018) with a bottom-up estimation for residential wood combustion emissions, providing the best estimate of organic carbon emissions (Denier van der Gon et al., 2015). Wildfire emissions are taken from the MACC global fire assimilation system (Kaiser et al., 2012). Initial and boundary conditions for most species are taken from CAMS near real-time. For organic matter these boundary conditions are not used since they were found to be unrealistically high at some instances. Instead, boundary conditions for OA species were set to zero. With prevailing westerly flow, the assumption of very clean conditions from the western boundary with zero boundary conditions can be justified for most situations and locations not too close to the eastern boundary. In the studied case, boundary conditions act only as a sink.

To generate the model output used in this work, we ran short simulations of 14 days in the last two weeks of February and July 2018 with 5 days of spin-up, the subsequent 5 days for training data-driven models and the last 4 days for evaluation of the converged data-driven models based on their reconstruction error of the volatility distributions. Evaluation of the simulations with observations is outside the scope of the present paper, as the model is regularly evaluated in model validation reports, as well as CAMS ensemble and model evaluations, and peer-reviewed publications, for example, Timmermans et al. (2022). With the first 5 days (February 15 through 19) disregarded as spin-up, 9 days were left for training and testing. With hourly output of surface VBS distributions over 216 hr, and 100 latitudinal grid lines by 140 longitudinal grid lines on the European MACC grid, there are approximately 3 million multi-dimensional data points for each VBS class. The data points range from 12 dimensional from the anthropogenic and biogenic VBS classes to 16- or 18-dimensional for the siSOA and POA VBS classes respectively. Model output from 5 days over February 20 through 24, approximately 1.7 million data points, was used as training data to optimize the parameters of the data-driven models with the objective to compress and reconstruct VBS distributions as accurately as possible. Model output from 4 days over February 25 through 28, approximately 1.3 million data points, was used to evaluate how much reconstruction error each approach introduces: this is detailed in Section 3, which concludes with a selection of the most promising approach.

Section 4 presents the results of implementing the selected approach in LOTOS-EUROS to compress tracers to superspecies before the advection operator and decompress to the VBS tracer distributions after the advection operator. All 3D experiments in Section 4 are run for periods of 2 weeks, more than double the length of the training time horizon. The operator time splitting step in LOTOS-EUROS is chosen dynamically based on wind conditions to satisfy the Courant-Friedrichs-Lewy criterion (Courant et al., 1928, 1967) varying from 1 to 10 min (Manders et al., 2017) with the advection operator called twice in each time step (Manders-Groot et al., 2021). With the advection operator called at a minimum of 12 times an hour for 2-week simulations, the superspecies compression/decompression step is done over 4,000 times for each grid cell. Grid cells interact with each other via transport processes: over the whole MACC domain with 100 longitudinal grid lines, 140 latitudinal grid lines, and 5 levels, superspecies are advected over 280 million times. Section 4 quantifies the effect of advecting superspecies to a baseline run of LOTOS-EUROS advecting all VBS tracers, with model configuration remaining otherwise identical. Also investigated is how well the superspecies optimized on the last 2 weeks of February (winter conditions in Europe) on the MACC grid generalize to (a) the last 2 weeks of July (summer conditions in Europe) with different continental spatial patterns as well as temporal patterns over forested areas and (b) a higher-resolution domain of 0.1°by 0.1°used in CAMS forecasting.

### 2.4. Linear Approach

A linear approach could be used to project the tracer space into a lower dimensional subspace allowing linear combinations of the tracers to be passed to the advection operator. Principal component analysis is a common linear projection method but is mean-centered and can lead to negative values, which are less readily interpretable as concentrations. Non-negative matrix factorization (NMF), also called positive matrix factorization, is an unsupervised data-driven approach chosen in applications where values must remain non-negative, for example, pixel values in image compression (Lee & Seung, 1999) or concentrations in the physical sciences (Paatero & Tapper, 1994). Given a matrix of non-negative data $\mathbf{V} \in \mathbb{R}^{m \times n}$ with $m$ dimensions and $n$ data points, NMF returns two non-negative approximate factors of $\mathbf{V}$ according to an objective function

$$\underset{\mathbf{W},\mathbf{H}}{\mathrm{argmin}} \|\mathbf{V} - \mathbf{WH}\| \qquad s.t. \quad \mathbf{W}, \mathbf{H} \geq 0 \tag{1}$$

where $\mathbf{W} \in \mathbb{R}_{\geq 0}^{m \times r}$ is a mapping from the $m$ dimensional space to a lower dimensional latent space with $r$ features, and $\mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}$ is the latent space representation of each data point. The inequality is interpreted as an element-wise constraint. We use the Frobenius norm in the objective function, which is the default NMF norm in the scikit-learn Python package (Pedregosa et al., 2011). For our application, $m$ is the number of tracers for each class, $n$ the total number of grid cells multiplied by the number of time steps, and $r$ the number of superspecies (a hyperparameter selected in Section 3.1). Each row of $\mathbf{V}$ corresponds to a tracer for each VBS class, and each column the tracer distribution for a given grid cell and time step. $\mathbf{H}$ can be physically interpreted as the concentration of $r$ superspecies representing the tracer concentrations of that VBS class: each column of $\mathbf{H}$ corresponds to the grid cell and time step in $\mathbf{V}$. $\mathbf{W}$ acts as a mapping from the superspecies representation back to the VBS tracer concentrations: a given column of $\mathbf{W}$ can be physically interpreted as the concentration profile of one super-species, with each element representing the relative composition of a VBS tracer in that superspecies. We use NMF to converge on a $\mathbf{W}$ for each VBS class that contains superspecies with characteristic volatility distribution shapes. These superspecies are linearly combined in ways that capture the variation of VBS distributions over all grid cells as well as possible. The coefficients determining the linear combination are the concentrations of each superspecies.

NMF operates on a data matrix, handling batches of observations all at once. For our application, compression of current concentrations of VBS tracers $\vec{v} \in \mathbb{R}^m$ in a given grid cell to a lower dimensional space needs to happen with each new time step. For the purpose of speeding computations, it might be counterproductive to perform the NMF algorithm online in every time step. If $\mathbf{W}$ is optimized using Equation 1 on sufficiently representative training data, it can be used to decompress a set of superspecies $\vec{h} \in \mathbb{R}^r$ to a set of tracers $\vec{v}_{dec} \in \mathbb{R}^m$ approximating $\vec{v}$. However, we still need to obtain the superspecies vector $\vec{h}$. Given a sufficiently representative $\mathbf{W}$, we can use its Moore-Penrose pseudoinverse $\mathbf{W}^+ \in \mathbb{R}^{r \times m}$ to compress a new set of tracers $\vec{v}$ to a corresponding set of new superspecies $\vec{h}$. $\mathbf{W}^+$ may have negative elements for $r > 1$ (more than one superspecies, or degree of freedom), theoretically yielding negative values for superspecies or decompressed tracers. This potential limitation is quantified in Section 3.2. Instead of a Moore-Penrose pseudoinverse, a positive-valued compression matrix $\mathbf{B} \in \mathbb{R}_{\geq 0}^{r \times m}$ can be obtained by similar non-negative matrix factorization methods, using the objective function:

$$\underset{\mathbf{B}}{\mathrm{argmin}} \|\mathbf{H} - \mathbf{BV}\|_F^2 \qquad s.t. \quad \mathbf{B} \geq 0 \tag{2}$$

The full approach to obtain non-negative compression and decompression matrices then becomes

1. Given tracer data $\mathbf{V}$, find $\mathbf{H}$, $\mathbf{W}$ such that $\mathbf{V} - \mathbf{WH}$ is minimized.
2. Given tracer data $\mathbf{V}$, and using $\mathbf{H}$ from the previous step, find $\mathbf{B}$ such that $\mathbf{H} - \mathbf{BV}$ is minimized.
3. Use $\mathbf{B}$ to compress subsequent observations of VBS tracers $\vec{v}$ to a non-negative vector of superspecies $\vec{h}$, and $\mathbf{W}$ to decompress $\vec{h}$ to the original tracer space $\vec{v}_{dec}$.

The compression and decompression matrices $\mathbf{B}$ and $\mathbf{W}$ are optimized for each VBS class, to avoid mixing different classes of OA that have different properties (e.g., molar mass). An important hyperparameter of this approach is $r$, the size of the latent space (number of superspecies). This can be chosen by constructing an elbow plot of error metrics with varying $r$, while also considering compression factor and is done in Section 3.1.

## 2.5. Nonlinear Approach

We investigate whether a more complicated model than the pair of non-negative matrices is appropriate for compressing VBS tracers. Motivated by the recent success of artificial neural networks (NNs) in emulating models of atmospheric composition (Kelp et al., 2020; Schreck et al., 2022; Sturm & Wexler, 2022), we construct a neural network autoencoder that can reversibly compress the VBS tracers to a latent space. Analogously to Section 2.4, the NNs are trained on **V** for each VBS class over the entire domain and training time frame, with the goal of applying a single NN parameterization for each VBS class at all grid cells. Neural networks are connected networks of artificial neurons: each neuron calculates a linear combination of its input, adds a bias scalar, and feeds this result to a (usually non-linear) activation function (Marsland, 2014). Neurons performing this operation on the same input in parallel are designated as a layer within the neural network. Neural networks can have multiple such layers: vector output from neuron layers that are not final output of the NN are called hidden layers. A neural network autoencoder attempts to replicate the identity function via compression, where hidden layers compress the input to the NN to a smaller latent space of size $r$. For our application, the activation function chosen for each neuron is a rectified linear unit that outputs the maximum of its input and zero. This choice of activation function constrains output of both the hidden layer and the NN output to their respective positive half-spaces. In other words, like the non-negative compression/decompression matrices in Section 2.4, this activation function ensures concentrations will not go below zero.

While matrix multiplication to a lower-dimensional space is also part of the linear approach in Section 2.4, the neural network adds complexity in its parameter space via multiple layers with weight parameters, as well as bias and activation functions between layers of neurons. Such complexity obscures physical interpretation: no one layer of the neural network can represent a set of superspecies with distinct compositions as **W** does in NMF. This model should be chosen if it significantly outperforms a linear method using the same size $r$. As the NNs are compared directly to the linear method, one NN per VBS class is chosen.

Training a neural network involves optimizing the coefficients of the linear combination and bias scalar for each perceptron through local minimization methods, often gradient descent. To prevent overfitting of the NNs, dropout layers are used to temporarily remove some neurons during training, and training of NNs is stopped when no further improvement in predictions on a set of validation data (10% of the 5 days training data) after a certain number of passes through the training data is obtained (Li et al., 2020). The neural network models are constructed and trained with the Keras library (Chollet, 2015) using a TensorFlow backend (Abadi et al., 2016).

## 2.6. Physically Consistent Models: Conserving Mass and Phase

Sections 2.4 and 2.5 developed methods to ensure non-negativity of both the compressed superspecies and decompressed tracers. This section refines the linear method to preserve other physical information: concentration and phase.

An advantage of the linear method is that the direction of the decompressed tracer space is invariant to scaling of the superspecies space. In other words, the concentration of superspecies can be adjusted without changing the relative volatility distribution of the decompressed tracers. We can use a scaling factor after compression to ensure that the total concentration of superspecies is equal to the total concentration of the tracers for each VBS class. Similarly, after decompression, we can ensure that the total concentration of decompressed tracers is equal to the total concentrations of superspecies. This ensures that compression and decompression neither add nor remove mass. The scaling factor $s_{com}$ after using **B** to compress tracers $\vec{v}$ to the superspecies vector $\vec{h}$ is

$$s_{com} = \frac{\sum_{i=1}^{m} v_i}{\sum_{j=1}^{r} h_j} \tag{3}$$

After decompression to $\vec{v}_{dec}$ using **W**, the decompressed tracers can be scaled using a factor $s_{dec}$, where

$$s_{dec} = \frac{\sum\limits_{j=1}^{r} h_j}{\sum\limits_{i=1}^{m} v_{dec,i}} \tag{4}$$

Despite conserving total concentration of all tracers, the concentration of total organic aerosol (TOA) may not be conserved due to errors in the mass distribution over volatility bins after decompression. A variation of this method to conserve TOA instead of total concentration, as well as an alternative way to conserve total concentration only using $\mathbf{W}$ from NMF, is explored in Sturm (2021). However, the compromise of conserving TOA versus total concentration is avoidable by adding another cross section: creating compression and decompression matrices $\mathbf{B}$ and $\mathbf{W}$ for each phase as well as VBS class, for example, one transformation for all biogenic gaseous VBS tracers and a separate transformation for all biogenic particle tracers. This phase-specific approach results in eight parameterizations instead of four used in Sections 2.4 and 2.5. The following section gives an overview of all four approaches: these approaches will be tested on their reconstruction accuracy in Section 3.

### 2.7. Four Approaches

The methods developed in Sections 2.4–2.6 lead to the following four approaches.

- Approach 1: NMF/Pseudoinverse linear approach: NMF to find an optimal decompression matrix $\mathbf{W}$, and use its pseudoinverse (with potentially negative elements) $\mathbf{W}^+$ as a compression matrix for each VBS class
- Approach 2: Non-negative matrix factorization: NMF to find an optimal decompression matrix $\mathbf{W}$ and a non-negative compression matrix $\mathbf{B}$ for each VBS class
- Approach 3: Non-negative neural network autoencoder: Create a more complicated neural network with ReLU activation functions in the superspecies and output layers, for each VBS class
- Approach 4: Mass-conserving, non-negative matrix factorization with phase specific superspecies: Create $\mathbf{W}$ and a non-negative compression matrix $\mathbf{B}$, for each phase in each VBS class

Section 3 investigates how well each approach can reconstruct volatility distributions of all four VBS distributions after compression. We select the most promising method in Section 3.3 based on reconstruction accuracy and physical consistency, to be incorporated into a 3D simulation.

## 3. Model Development and Selection

The four approaches developed in Sections 2.4–2.6, and outlined in Section 2.7, were trained on LOTOS-EUROS model output from February 20th through 24th using the model configuration detailed in Section 2.3. This section evaluates the four approaches on their ability to compress and reconstruct the volatility distributions of model output from a different set of days, February 25th through 28th. Section 3.1 uses Approach 1, the simplest approach, to investigate how dimensionality of the latent space $r$ (number of superspecies), inversely related to compression factor, affects reconstruction accuracy. Section 3.2 deals with physical consistency: Section 3.2.1 investigates how Approach 1 can lead to negative concentration values, and motivates the non-negativity constraints in Approaches 2, 3, and 4. Section 3.2.2 demonstrates how Approach 4 conserves mass and phase when mapping tracers to superspecies and back. Finally, Section 3 compares the reconstruction error and physical consistency of all four compression approaches and selects the most promising approach to be implemented in LOTOS-EUROS.

### 3.1. Compression Factor and Accuracy

To obtain a sense of error obtained by a maximum compression factor and the simplest model, we use NMF with a single superspecies ($r = 1$) per VBS class to obtain a decompression matrix (in this case a vector) $\mathbf{W}$ and calculate its pseudoinverse $\mathbf{W}^+$ to be used for compression. This compression strategy is evaluated on reconstruction accuracy of test model output of the entire domain and time period, using average bias and root mean square error (RMSE). While bias is an indicator of the total material that is introduced or removed artificially by compression, RMSE is an absolute metric that indicates how accurately the reconstructed VBS tracers reproduce the volatility distribution. Table 1 shows both reconstruction error metrics for the tracer set of each

**Table 1**
*Test Reconstruction Error Metrics Using the NMF/Pseudoinverse Approach With 1 Superspecies per VBS Class*

|  | Mean ($\mu g\ m^{-3}$) | RMSE ($\mu g\ m^{-3}$) | Bias ($\mu g\ m^{-3}$) | NRMSE (%) | NMB (%) |
|---|---|---|---|---|---|
| aVOC | 0.0043 | 0.0021 | $-3.9 \times 10^{-6}$ | 48.8 | 0.1 |
| bVOC | 0.0262 | 0.0061 | $2.9 \times 10^{-4}$ | 23.3 | 1.1 |
| POA | 0.0558 | 0.0441 | $-0.0021$ | 79.0 | $-3.7$ |
| siSOA | 0.0153 | 0.0205 | $6.4 \times 10^{-5}$ | 134.0 | 0.4 |
| TOA | 0.386 | 0.266 | 0.094 | 68.9 | 24.3 |
| TOM | 1.61 | 0.0978 | $-0.0328$ | 6.1 | $-2.1$ |

class, as well as the reconstruction bias and RMSE's of total organic aerosol concentration (TOA) and total organic material (TOM) from summing across VBS classes. The mean concentrations for each VBS class, as well as TOA and TOM, are included for comparison. We also include normalized root mean square error (NMRSE) and normalized mean bias (NMB) calculated by respectively dividing RMSE and bias by the mean.

Using one superspecies $r = 1$ in Approach 1 leads to high values of RMSE relative to the mean. Moreover, by the use of a single superspecies the tracers pass through a linear transformation of rank 1: the concentration distribution over the volality bins will always have the same shape, with grid cells and different time steps differing only in magnitude, as scaled by the superspecies concentration $h$. This means any spatiotemporal variability of the distribution shape will be lost after passing through a single-dimensional superspecies space. More complexity is needed to capture variation in volatility distribution. This motivates larger matrices that have more degrees of freedom $r$, which comes at the cost of compression factor. Figure 2 visualizes the effect of compression extent on accuracy, using $\mathbf{W}^+$ to convert to superspecies and $\mathbf{W}$ to map back to tracers. Reconstruction accuracy is reported for the set of tracers in each class (both particle and gas) as well as TOA (total organic aerosol, calculated by summing the concentrations of particle tracers across classes).

Figure 2 shows RMSE monotonically decreasing with increasing number of superspecies, with diminishing returns after 3 superspecies. More superspecies to advect will increase the computational burden of the advection operator in LOTOS-EUROS without a substantial improvement in RMSE or bias. In light of the desire to maximize compression factor, the two elbow plots indicate that 3 superspecies strikes a good balance between dimension reduction and accuracy. Using 3 superspecies per class ranges from a compression factor of 4 (the aVOC and bVOC basis sets) to 6 (the POA basis set) with a significant improvement in accuracy from 2 superspecies and minimal improvement in accuracy when using 4 or more superspecies.
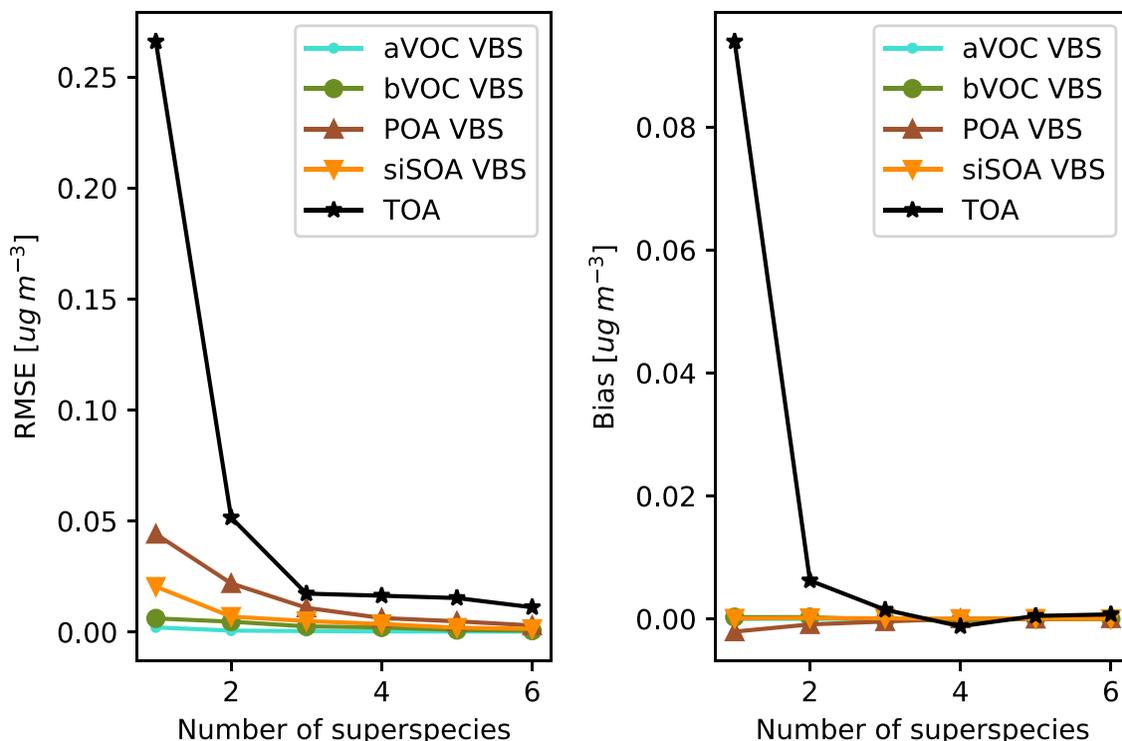


**Figure 2.** Relationship between the number of superspecies and the RMSE and bias for the 4 VBS classes, as well as TOA. There are diminishing returns in accuracy after 3 superspecies per VBS class.

Improved accuracy with number of superspecies comes from the increased degrees of freedom, as each subsequent column of $\mathbf{W}$ adds another basis direction. Each column of $\mathbf{W}$, when normalized, can also be interpreted as a superspecies of unit concentration with elements corresponding to composition of VBS tracers. Each superspecies can also be interpreted as a different regime of organic aerosol, found through a data-driven method. Multiple superspecies can be combined in different amounts, corresponding to their concentrations, to form other distributions.

### 3.2. Physical Consistency of Results

#### 3.2.1. Motivating Non-Negative Constraints

Section 2.4 raised the theoretical possibility of obtaining negative concentrations when using the pseudoinverse $\mathbf{W}^+$ to compress tracers into superspecies. Negative elements in $\mathbf{W}^+$ can lead to negative superspecies. Negative superspecies concentrations are not directly a problem, as the current advection scheme in LOTOS-EUROS v2.2.1 is based on that of Walcek (2000), which is able to handle negative tracer values. However, using the non-negative $\mathbf{W}$ to decompress negative superspecies concentrations back to the tracer space can lead to negative tracer values. Here, we quantify this limitation in practice using 3 superspecies.

Negative concentrations that are extremely small in magnitude can be approximated as zero. This tolerance can of course be set to a threshold, for example, $-1 \times 10^{-8}$ µg m$^{-3}$. However, using the test data of the POA VBS as an example, there are over 4.7 million cases in the test data where a POA VBS tracer is below $-1 \times 10^{-8}$ µg m$^{-3}$, which is more than 19% of the 24 million values in the test data for the POA VBS.

One could choose a more relative, less arbitrary tolerance: for instance, all concentrations that are more negative than the magnitude of the corresponding bias for each VBS. These "significantly negative" concentrations would be negative even after an additive bias correction. For the POA VBS, there were 855,083 such concentrations, about 3.5% of the total test data. Using this relative tolerance, other VBS classes showed even larger proportions of "significantly negative" concentrations: 4.2%, 5.6%, and 7.0% respectively for the siSOA, aSOA, and bSOA VBS classes (for the anthropogenic VBS and siSOA VBS, which had positive biases, the tolerance was chosen to be the negative magnitude of the corresponding bias).

Using the pseudoinverse $\mathbf{W}^+$ for compressing VBS tracers (Approach 1) can result in a number of significantly negative values when using 3 superspecies per VBS class, which motivates the development of non-negative compression strategies. For each VBS class, we find a positive compression matrix $\mathbf{B}$ to replace $\mathbf{W}^+$, according to the objective function and constraints in Equation 2 (Approach 2).

We compare this matrix factorization approach (Approach 2) with a neural network autoencoder (Approach 3) for each VBS class. We construct and train a 5-layer neural network autoencoder with rectified linear unit activation functions in the superspecies and output layers to ensure non-negativity of both superspecies and decompressed VBS tracers. In other hidden layers, a sigmoidal activation function, hyperbolic tangent, is used. In training, a dropout rate of 0.1 is used for every layer except for the superspecies layer. For the autoencoder of each VBS class, the center superspecies layer is chosen to have 3 values: the value of this hyperparameter is chosen for comparison to the linear matrix factorization approach. Section 3.3 compares all four approaches based on how well they reconstruct the VBS tracers after decompression.

#### 3.2.2. Conserving Mass and Phase

Section 2.6 proposed a method for conserving total concentration of the VBS tracers in both the superspecies representation and in subsequent reconstruction to decompressed tracers. Approach 4 applies this method to the cross-sections of VBS class and phase (particle or gas) to ensure that the superspecies transformation does not add or remove mass artificially in the gas and particle phases of every class: this results in conservation of total gas concentration, total aerosol concentration, and concentration of total organic material (TOM). Phase-specific superspecies are composed of entirely gas or entirely particle tracers, conserving information on phase while in the latent space representation.

Phase-specific superspecies require adding another cross-section, halving the number of tracers to be compressed and decompressed by each pair of $\mathbf{B}$ and $\mathbf{W}$, respectively. For this reason, continuing to use 3 superspecies for each phase within each VBS class would reduce the compression factor to slightly over 2.4, not much better than the compression factor of around 2 when using the partitioning-based compression approach. However, using
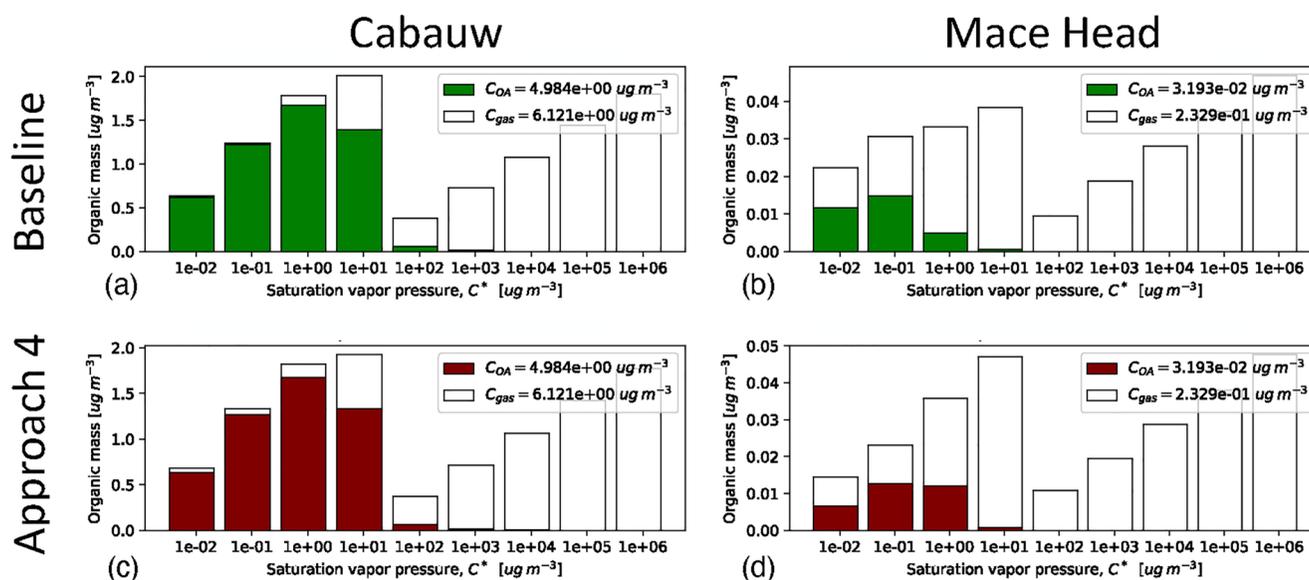
**Figure 3.** Comparison of the volatility basis set distribution for POA near two sites: Cabauw and Mace Head at a snapshot in time on 26 February 2018. The top row in green shows the distributions as modeled by LOTOS-EUROS at Cabauw (a, b) Mace Head. The bottom row in maroon shows the distributions at Cabauw (c) and Mace Head (d) after the non-negative compression/decompression using phase-specific superspecies. Total concentrations are conserved when comparing the legends of the modeled distributions to the reconstructed distributions.

only 1 superspecies per phase per class would fix each corresponding set of tracers to a single shape upon reconstruction, as discussed in Section 3.1. To ensure that this method captures spatiotemporal variability of volatility distributions while maintaining a useful compression factor, we choose to use 2 superspecies per phase per VBS class. This design choice results in a compression factor of approximately 3.6. Its accuracy is compared to the other strategies in the model selection process in Section 3.3.

Figure 3 demonstrates the mass-conserving properties of Approach 4 using representative examples of the primary organic aerosol VBS distribution at two different atmospheric monitoring sites: the Cabauw Experimental Site for Atmospheric Research in the Netherlands, and Mace Head Atmospheric Research Station in Ireland. Mace Head is a more pristine and remote station (O'Dowd et al., 2014). The legend in Figure 3 shows that POA concentration at Cabauw is two orders of magnitude higher than that at Mace Head, 4.984 $\mu g\,m^{-3}$ compared to 0.032 $\mu g\,m^{-3}$.

Figure 3 compares the primary VBS distribution to the reconstructed primary VBS distributions after mapping to phase-specific superspecies and back again using two sites: Cabauw and Mace Head, as representative examples. Comparing the legends of (a) with (c), it can be seen that total POA concentration, as well as total concentration of tracers in the gas phase, is conserved to machine precision after passing through compression. The same holds for the total concentrations at Mace Head, (b) and (d), at orders of magnitude more dilute. With phase information and concentration conserved, the only source of error caused by compression to superspecies is in the shape of the distribution. This reconstruction error is more apparent at Mace Head in Figures 3b and 3d. The reconstructed distribution of Mace Head more closely resembles the constant primary organic emissions profile modeled by LOTOS-EUROS: during training, grid cells with high primary organic emissions are weighted heavily as they tend to have higher aerosol loading. Though the data-driven approaches applied to the primary VBS class are biased to reconstruct the volatility distribution of grid cells with high POA loading, the conservation constraints in Approach 4 ensure that no material will be artificially introduced in more dilute conditions. Though the gas/particle split is not guaranteed to be in equilibrium after reconstruction, the partitioning subroutine (which is not itself a computationally expensive component of the VBS approach) will subsequently determine the gas/particle split.

**Table 2**
*Evaluation RMSE of Selected Approaches*

|  | Approach 1 | Approach 2 | Approach 3 | Approach 4 |
|---|---|---|---|---|
| aVOC VBS | $4.4 \times 10^{-4}$ | 0.0010 | 0.0021 | 0.0011 |
| bVOC VBS | 0.0026 | 0.0078 | 0.0181 | 0.0042 |
| POA | 0.0109 | 0.0285 | 0.0306 | 0.0142 |
| siSOA | 0.0050 | 0.0086 | 0.0094 | 0.0057 |
| TOA | 0.0173 | 0.133 | 0.101 | $6.9 \times 10^{-13}$ |
| TOM | 0.0547 | 0.240 | 0.328 | $1.0 \times 10^{-12}$ |

*Note.* All values reported in $\mu g\ m^{-3}$.

### 3.3. Model Selection

In this section, we compare the four approaches described thus far, and make a judgment about the most promising strategy, evaluated on reconstruction accuracy and physical consistency. The selected approach will be implemented in LOTOS-EUROS v2.2.1 to accelerate the advection operator. The four approaches are restated here, including the number of superspecies used.

- Approach 1: NMF/Pseudoinverse linear approach: NMF to find an optimal decompression matrix $\mathbf{W}$, and use its pseudoinverse (with negative elements) $\mathbf{W}^+$ as a compression matrix using 3 superspecies per VBS class
- Approach 2: Non-negative matrix factorization: NMF to find an optimal decompression matrix $\mathbf{W}$ and a non-negative compression matrix $\mathbf{B}$ using 3 superspecies per VBS class
- Approach 3: Non-negative neural network autoencoder: Create a more complicated neural network with ReLU activation functions in the superspecies and output layers, using 3 superspecies per VBS class
- Approach 4: Mass-conserving, non-negative matrix factorization with phase specific superspecies: Create $\mathbf{W}$, as well as a non-negative compression matrix $\mathbf{B}$ using 2 superspecies per phase per VBS class

Tables 2 and 3 show RMSE and bias of the tracers for each VBS class for the 4 approaches, as well as total organic aerosol (TOA) and total organic material (TOM) concentrations.

Approach 2 uses non-negative $\mathbf{B}$ and $\mathbf{W}$ to linearly combine tracers into three superspecies and shows lower RMSE values than the NN autoencoder in Approach 3, with the exception of TOA concentration. This indicates that matrix factorization is probably suitable for VBS tracer compression. Using the pseudoinverse $\mathbf{W}^+$for compression resulted in lower RMSE for all the VBS classes, but has the critical weakness of producing a significant amount of negative concentrations for superspecies and subsequently reconstructed tracers as explored in Section 3.2.1. Though the phase-specific superspecies approach does not have as low of RMSE for each VBS class as the pseudoinverse approach, it outperforms the other two non-negative approaches. Moreover, it conserves absolute metrics on compression, ensuring that material will stay in each class and each phase, and no material will be added or removed by compression: for this reason, all biases are negligible to machine precision. Preserving information on phase during compression to superspecies has another advantage. This approach can be used in other processes such as dry deposition, which handles particle and gas tracers separately. Because the phase-specific superspecies method (Approach 4) is physically consistent while quite accurate in reconstruction error, and is readily extended to other phase-specific processes, it is chosen for implementation in LOTOS-EUROS v2.2.1.

## 4. Results: Superspecies Implementation in LOTOS-EUROS

The phase-specific, matrix factorization superspecies method (Approach 4) chosen in Section 3.3 was implemented in LOTOS-EUROS v2.2.1. This section explores the accuracy and speedup of replacing VBS tracers with superspecies in advection, as well as the generalizability of the superspecies to different seasonal conditions and spatial resolutions. Additional tracers for superspecies were added to the LOTOS-EUROS tracer list. Subroutines were added to the VBS module to load the parameterizations, as well as perform the compression and decompression operations. When running with the superspecies method, the subroutines are called in the driver program as follows:

1. The initialization subroutine loads offline-optimized $\mathbf{W}$ and $\mathbf{B}$ for each phase and class before the time loop starts.
2. Within the time loop, directly before the call to the advection operator, the compression subroutine is called to map VBS tracers to superspecies concentrations using $\mathbf{B}$, overwriting the current superspecies values. The advection operator skips VBS tracers and advects superspecies instead.

**Table 3**
*Evaluation Bias of Selected Approaches*

|  | Approach 1 | Approach 2 | Approach 3 | Approach 4 |
|---|---|---|---|---|
| aVOC VBS | $2.6 \times 10^{-5}$ | $1.2 \times 10^{-4}$ | $-3.9 \times 10^{-4}$ | $2.8 \times 10^{-20}$ |
| bVOC VBS | $-1.6 \times 10^{-4}$ | $3.8 \times 10^{-4}$ | $-0.0051$ | $-1.6 \times 10^{-16}$ |
| POA | $-4.2 \times 10^{-4}$ | 0.0050 | $-0.0075$ | $-8.8 \times 10^{-18}$ |
| siSOA | $-9.9 \times 10^{-5}$ | $7.7 \times 10^{-4}$ | $-0.0022$ | $1.2 \times 10^{-19}$ |
| TOA | 0.0015 | 0.0657 | $-0.0346$ | $-1.3 \times 10^{-15}$ |
| TOM | $-0.00763$ | 0.108 | $-0.237$ | $-2.1 \times 10^{-15}$ |

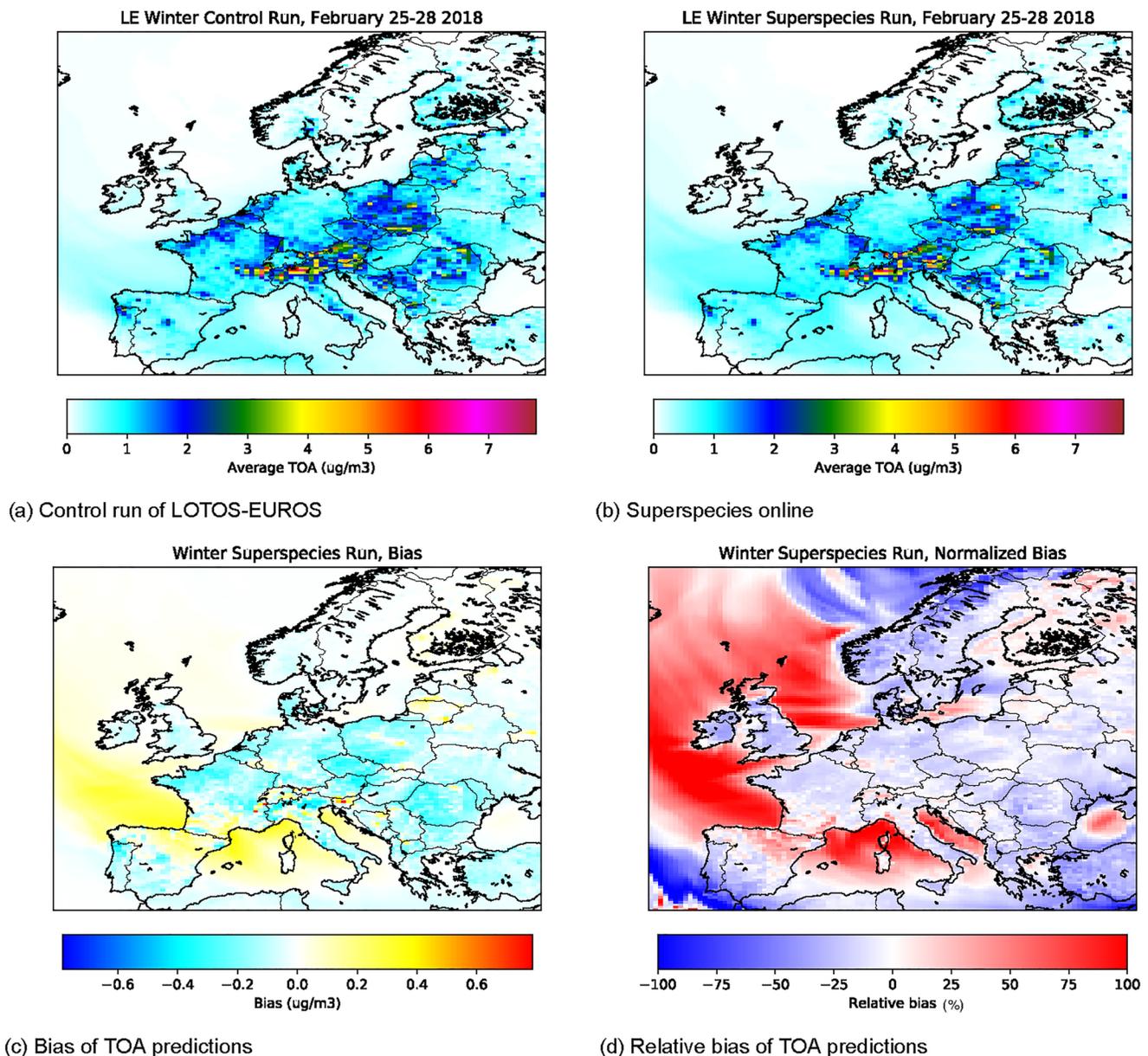*Note.* All values reported in $\mu g\ m^{-3}$.

**Figure 4.** Average TOA for February 25th through 28th 2018, during a 2-week simulation from February 15th through 28th using superspecies matrices optimized offline on winter conditions from February 20th through 24th.

3. Within the time loop, directly after the call to the advection operator, the decompression routine is called to transform superspecies into VBS tracers using **W**, overwriting previous VBS tracer values.

After offline training on data from February 20th through 24th, 2018, the selected superspecies parameterization was loaded into LOTOS-EUROS and used in the advection operator for a run from February 15th through 28th. The results of this run are compared with a control run advecting VBS tracers to directly assess the error from advecting superspecies. Small errors caused by advecting superspecies change subsequent VBS tracer concentrations such that the period of February 20th through 24th differs from the training data set. In that time period, however, meteorological conditions and other processes independent of the VBS and superspecies parameterization are identical to that of the offline training data set. For the sake of comparison, the superspecies run and control run are evaluated on February 25th through 28th, even though the superspecies run has the chance to accumulate error and diverge from the control run from the beginning of the simulation on February 15th.

**Table 4**
*Average TOA Composition in the Control Runs for February and July*

| OA type | February | July |
| --- | --- | --- |
| aSOA | 0.8% | 9.5% |
| bSOA | 4.5% | 34.8% |
| POA | 61.2% | 12.5% |
| siSOA | 33.5% | 43.2% |

Advecting superspecies reproduces the spatial patterns of average TOA across the entire domain. Figure 4 shows average TOA of the control run and the superspecies run, from February 25th through February 28th. This test time period is well into the model run, 10 days after the beginning of the simulation. During this time period and over the entire domain, average bias of TOA of the superspecies run compared to the control run is small and slightly negative, $-0.0095$ µg m$^{-3}$. Small average bias is not in itself indicative of low error, as positive and negative bias cancellations throughout the domain and time period are possible. RMSE, an absolute metric, was larger at 0.217 µg m$^{-3}$. Figure 4 shows total OA, though the VBS classes have partly compensating biases: for example, the positive bias in northern Spain was mainly caused by a positive bias from the siSOA class, of which the corresponding gas-phase species have a longer lifetime than those of the POA class. The north of Spain is less densely populated than other parts of the domain and composition is more affected by long-range transport, as are the ocean parts. Back trajectory analysis of this region revealed both stagnant and long-range trajectories for the averaging period (25–28 February), with long-range transport from more polluted areas in the northeast of the domain. Further analysis revealed that for northern Spain siSOA caused a positive bias for TOA. The condensable gases of the siSOA class have a longer lifetime than those of the POA due to differences in their deposition velocities (arising from different Henry coefficient values). However, the general spatial patterns of total OA across the entire domain are preserved when advecting superspecies.

### 4.1. Seasonal Superspecies

The winter test period from February 25th through 28th directly followed the training test period from February 20th through 24th and had relatively similar conditions to what the superspecies transformation matrices were optimized for. A run in summer from July 20th through August 1st was chosen to assess the robustness of the winter-optimized superspecies to different seasons and weather patterns. Summer conditions differ from winter conditions in Europe for several reasons. One, biogenic precursor gases make up a larger contribution to formation of secondary organic aerosol in the summer, partially due to emissions from forests. Two, average temperatures are higher, affecting the partitioning of the VBS by changing the volatility basis set values $C^*$. The different conditions lead to different modeled compositions of total organic aerosol (TOA). Table 4 compares the modeled average composition of OA for February 25th through 28th to that for July 29th through August 1st.

Though siSOA is on average the largest component of TOA in the run from July 29th through August 1st this is not the full picture, and underscores the importance of bSOA under some conditions. The maximum concentration of surface siSOA over the entire domain over the entire period from July 29th through August 1st was 15.0 µg m$^{-3}$ and 99th percentile 1.3 µg m$^{-3}$, compared to the maximum bSOA concentration of 100.3 µg m$^{-3}$ and 99th percentile 9.4 µg m$^{-3}$. This indicates that although siSOA may dominate in background conditions and when TOA is low, bSOA is the dominant component of TOA in other conditions.

#### 4.1.1. Domain-Wide Assessment

Figure 5 shows average surface TOA, as predicted by the control run (a), the run with superspecies advected (b), and the bias and relative bias of the superspecies run with regards to the control, (c) and (d) respectively. The spatial patterns of TOA are visually different from the winter conditions in Figure 4. Primary organic emissions corresponding to POA are often the largest contributor to winter TOA, and for the time period in Figure 4, TOA is most concentrated in the Po Valley, Czechia, and Poland. The winter superspecies run is able to recreate these large regions of high TOA, as well as other smaller but distinct pockets of TOA, such as Madrid (the most populous city in Spain) and northwestern Portugal, a region with heavy industrial activity. In contrast, summer TOA is concentrated around southern Germany, Switzerland, Austria, and Slovenia. Many places in this region are forested, and contribute to TOA via emission of biogenic precursors of bSOA. The superspecies run shown in (b) is able to capture these spatial patterns, but with a strong bias. For this reason, other regions with high biogenic emissions become visually apparent in (b), such as southern Sweden, Finland Proper, and northwestern Russia, which are all heavily forested. Woodland regions are accounted for in LOTOS-EUROS via land use maps and tree-species emissions (Manders et al., 2017).
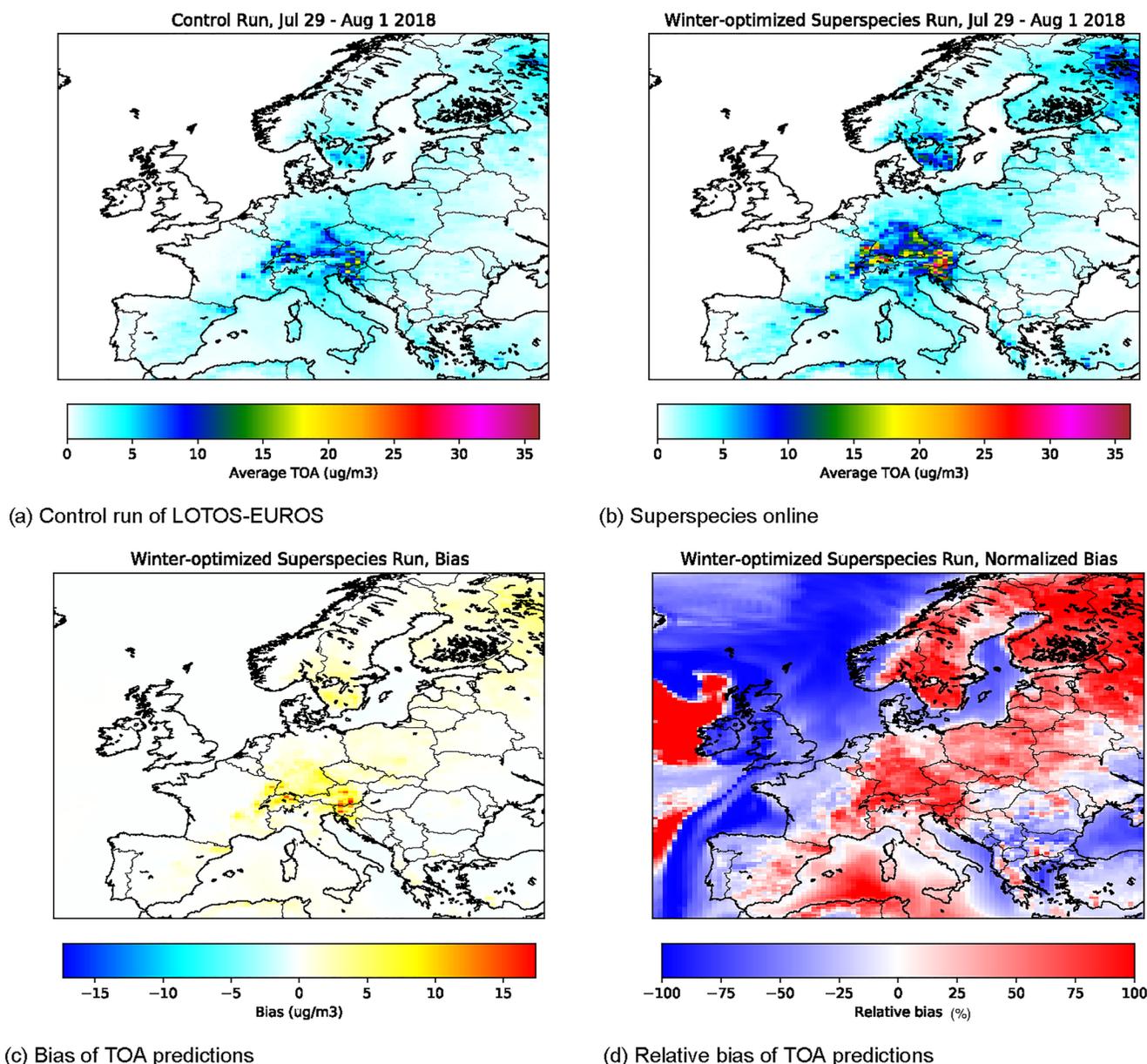
**Figure 5.** Average TOA for July 29th through 1 August 2018, during a 2-week simulation from July 19th through August 1st using superspecies matrices optimized offline on winter conditions from February 20th through 24th.

The superspecies optimized on winter conditions and tested on a 2-week run in July show a large positive bias over the areas with high average TOA, especially heavily forested regions. RMSE for TOA over the whole domain and time period is 2.12 μg m$^{-3}$, with an average bias of 0.321 μg m$^{-3}$. RMSE of the tracers from the biogenic VBS for all times and grid cells is 0.66 μg m$^{-3}$, an order of magnitude higher than tracers from the other VBS classes: the class of tracers with the next highest RMSE value is the siSOA VBS class, at 0.062 μg m$^{-3}$. The average bSOA bias (bias of total biogenic aerosol neglecting gaseous tracers) is 0.068 μg m$^{-3}$, three orders of magnitude smaller than the maximum bSOA bias of 82.9 μg m$^{-3}$. Overestimation of bSOA in the superspecies run under some conditions is likely due to errors in decompression, artificially shifting mass to lower volatility bins. However, the large positive bias in parts of the domain indicate that this tendency to overestimate bSOA only happens in certain conditions: namely, forested regions. The following section analyzes one grid cell in a
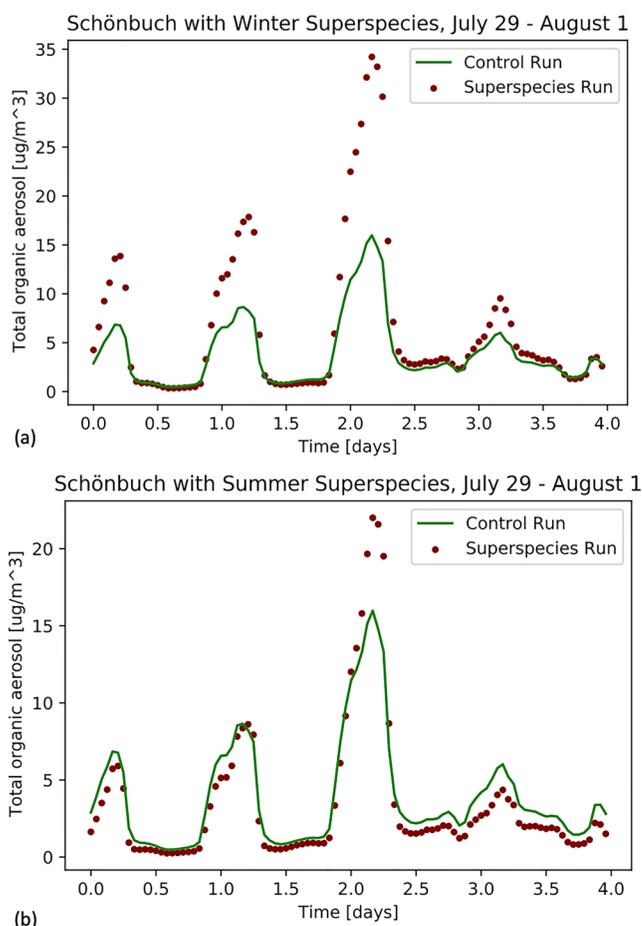
**Figure 6.** Temporal variation of TOA over Schönbuch from July 29th through August 1st using (a) winter-optimized superspecies and (b) summer-optimized species. The maroon points of TOA as predicted with when advecting superspecies are compared to the green line of TOA as modeled by the LE control run used as a baseline.

forested region, and finds additional temporal patterns where bSOA is significantly overestimated, leading to overestimation of TOA.

### 4.1.2. Case Study: Summer Night in a Forest

We choose a single grid cell over a forested area to investigate the superspecies tendency to overestimate bSOA. We study the LOTOS-EUROS grid cell containing the Schönbuch Natural Reserve in southwest Germany, which is 156 square kilometers and 85% forested. Figure 6a shows the temporal variation of TOA in the Schönbuch from July 29th through August 1st. This overestimation systematically occurs at night, with the night of July 30th to July 31st a particularly high TOA event showing the highest bias.

Examining Figure 6a, the peak overestimation occurs at 05:00 on July 31st and overestimates total bSOA with a factor between 2 and 2.5 times that of the control run. The superspecies run has a bSOA concentration of 32.9 $\mu g\ m^{-3}$, which comprises 99% of total OA concentration for that grid cell and time. The control run concentration of bSOA is 14.1 $\mu g\ m^{-3}$, about 95% of TOA for that simulation. By 09:00 on July 31st, both runs return to a total bSOA concentration of less than 3.5 $\mu g\ m^{-3}$. This night episode of high bSOA contains the largest overpredictions for that particular grid cell in the whole time period. However, it is illustrative of a failure mode of the winter-optimized superspecies to capture the total concentration of bSOA, and ultimately TOA due to the importance of bSOA contributions in this example. The spatial patterns and temporal patterns of the superspecies run compared to the control run show that the superspecies are limited in their ability to model conditions over forested areas on summer nights.

Given that winter-optimized superspecies showed limitations in capturing high bSOA events over forested areas at night, we investigate whether superspecies optimized on summer conditions and implemented online reproduce high bSOA conditions with more accuracy. Approach 4 was applied to model output from July 23rd through 28th, 2018, to obtain a superspecies parameterization optimized on summer conditions.

The superspecies approach optimized on summer conditions shows a much lower bias than the winter-optimized superspecies. The temporal behavior of summer-optimized superspecies from July 29th through August 1st after 10 simulated days is shown in Figure 6b. Comparing Figures 6a to 6b, it can be seen that the spatiotemporal pattern of bSOA bias is addressed by using summer-optimized superspecies, which do not show the same nightly overestimation pattern of winter-optimized superspecies. Total bSOA is even slightly underestimated in the day when using summer-optimized superspecies.

Averaged over the entire domain and time period of July 29th through August 1st, the summer-optimized superspecies display a slightly negative average bias for bSOA of −0.023 $\mu g\ m^{-3}$. Small pockets of TOA overestimation (within 10 $\mu g\ m^{-3}$) still occur in the same regions as the winter-optimized superspecies: over highly forested areas. The RMSE over the whole domain of time-averaged TOA was 0.98 $\mu g\ m^{-3}$ when using summer-optimized superspecies, less than half of the RMSE of 2.12 $\mu g\ m^{-3}$ when using winter-optimized superspecies. RMSE of the tracers from the biogenic VBS (both gas and particle phases) for all times and grid cells is reduced by a factor of 2, at 0.32 $\mu g\ m^{-3}$ compared to 0.66 $\mu g\ m^{-3}$. However, in superspecies trained on either season, the biogenic VBS tracers in the summer show significantly higher error than the tracers of the other VBS classes, with the siSOA VBS class having the next highest RMSE value at 0.050 $\mu g\ m^{-3}$. The limitation of winter-optimized superspecies and the subsequent improvement in accuracy when using summer-optimized superspecies indicates that this method might be best applied to different seasons: creating seasonal-specific superspecies results in higher accuracy. Analogously, Kelp et al. (2022) tested neural network surrogate models of atmospheric chemistry optimized online for 3-month seasons against neural networks trained online for a whole year, and concluded that ensembles of ML surrogate models specialized for specific seasons improve accuracy and stability.
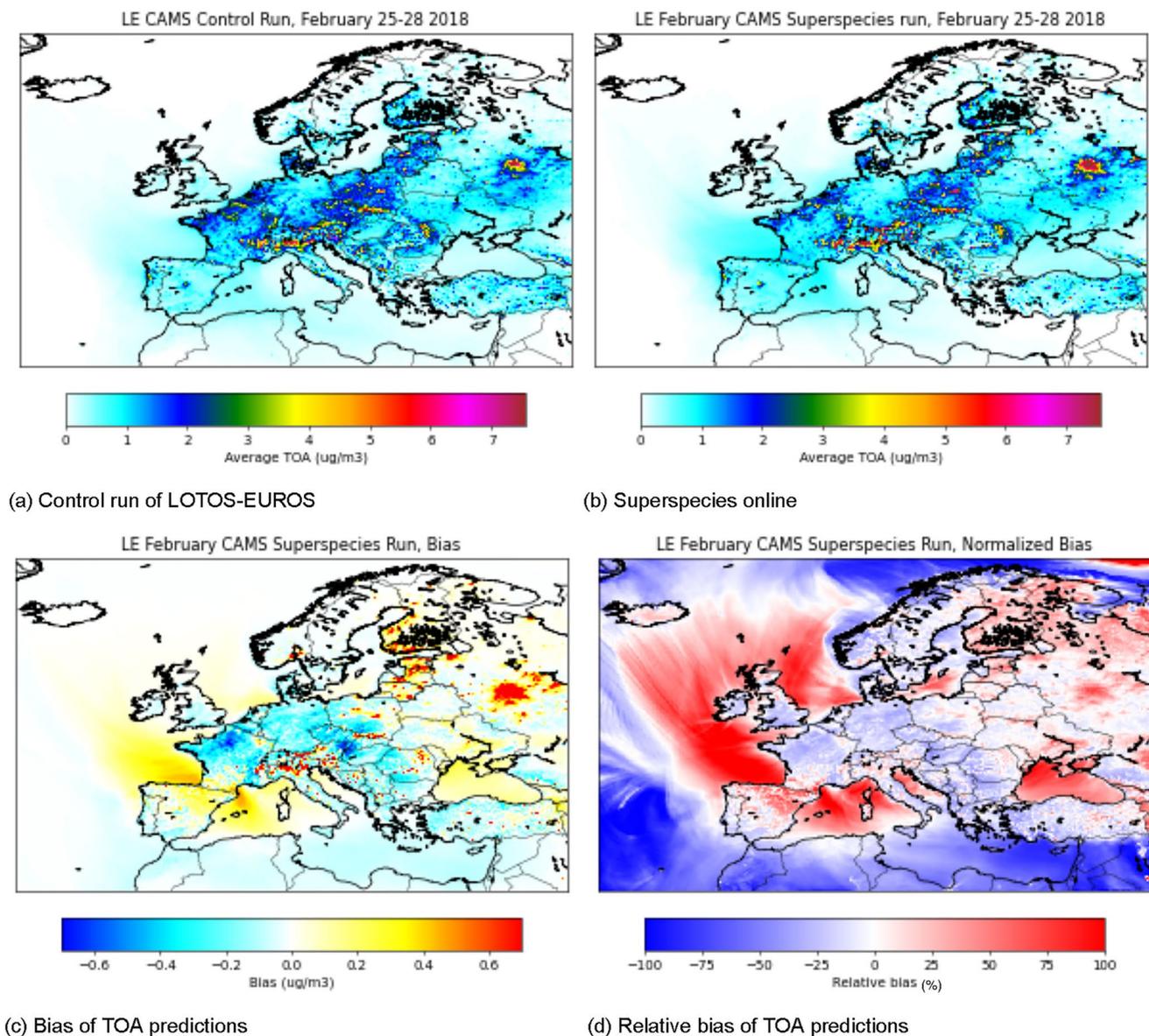
**Figure 7.** Time averaged TOA for the period of February 25th through 28th on the high-resolution domain used in CAMS operational forecasting, from control and superspecies runs, as well as bias and relative bias. The superspecies were optimized on model output from a simulation using the coarse-resolution MACC domain.

### 4.2. Towards Operational Forecasting on Higher-Resolution Domains

LOTOS-EUROS is one model in the ensemble used in the Copernicus Atmospheric Modeling Service (CAMS) operational forecasts, which requires all models to include SOA representation by 2023. The domain used in CAMS operational forecasts has a higher resolution and wider domain than the domain used by MACC: 0.1°by 0.1°for 420 by 700 grid cells compared to the 0.50°by 0.25°used in the MACC domain, and extending past Moscow, Russia. The change of resolution and domain increases the number of grid cells by a factor of 20. One result of this is many more grid cells and computations. Another result is that the operator splitting timestep $\Delta t$ needs to decrease in order to satisfy the Courant-Friedrichs-Lewy criterion as the grid cell distance is smaller. With a smaller operator splitting timestep, the advection operator as well as the compression and decompression steps are called more often. We investigate how the superspecies approach, optimized on model output from February 20th through 24th on the coarse-resolution MACC domain, generalizes to a 2-week run on the extended high-resolution
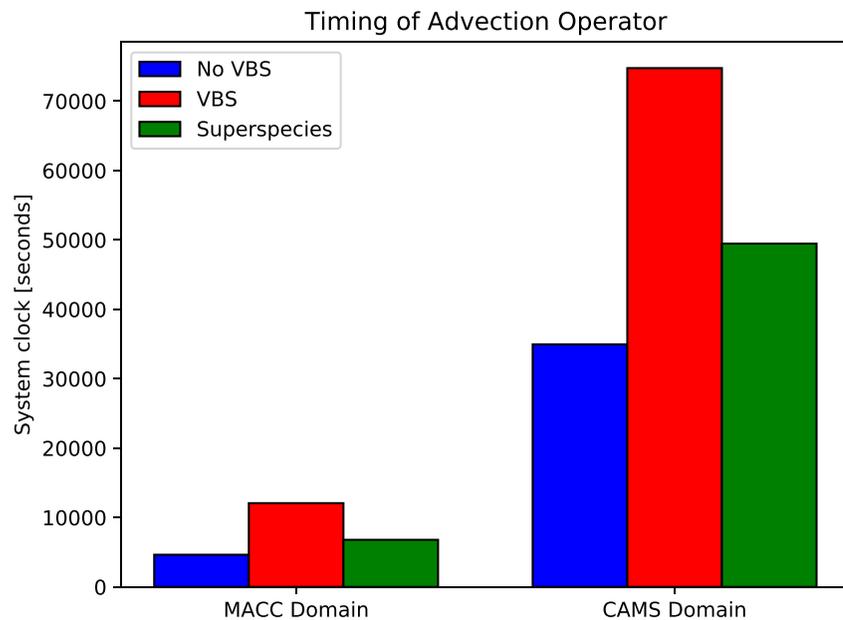
**Figure 8.** Use of the 58 VBS tracers approximately doubles the wall time spent on advection calculations. Advecting superspecies takes 56% and 66% of the time compared to advecting VBS tracers on the MACC and CAMS domains, respectively.

CAMS domain. Figure 7 shows the time-averaged TOA concentration across the entire CAMS domain for the test period of 25th–28th February 2018, chosen for ease of comparison with the winter run on the MACC domain.

The superspecies run has a positive bias for TOA of 0.019 µg m$^{-3}$, with visible overestimation in the area near Moscow, Russia, which is not in the MACC grid used to optimize the compression/decompression matrices. The colorbar limits of Figures 7a–7c were adjusted for visual comparison with Figure 4. For this reason, colors at the upper or lower limits should be interpreted as greater or equal to the limit. Though the maximum grid cell concentration of time-averaged TOA from both the superspecies run and the control run was 28.2 µg m$^{-3}$, 99.85% of the grid cells had a time-averaged TOA under 7.6 µg m$^{-3}$, which was chosen as the upper limit of the colorbar. This means that only 0.15% of the grid cells in Figures 7a and 7b exceed the limit shown in the color-bar. Neglecting the highest 0.15% of average TOA, the spatial patterns of the CAMS control run in Figure 7a are visually very similar to those of the that the CAMS superspecies run in Figure 7b. Both show spatial patterns similar to the simulations performed on the MACC grid for the same time period. The same approach is done for the bias shown in Figure 7c, with very few grid cells in the CAMS simulation exceeding the maximum error of time-averaged TOA on the MACC grid. The maximum absolute error of time-averaged TOA between the superspecies run and the control run was 8.9 µg m$^{-3}$, but 99.2% of all grid cells had an absolute error of less than 0.70 µg m$^{-3}$. Less than 1% of the grid cells in Figure 7c exceed the colorbar limit. The largest instantaneous bias for TOA was 89 µg m$^{-3}$ at a grid cell in northwestern Spain near Ponferrada during a high TOA event on February 25th at 19:00. This grid cell also showed the highest time-averaged TOA concentration of 32.0 µg m$^{-3}$ for the superspecies run, compared to 19.4 µg m$^{-3}$ for the control run. At the highest positive bias of 89 µg m$^{-3}$, TOA concentration as modeled by the superspecies run was 206.4 µg m$^{-3}$ while the control run TOA concentration was 117.4 µg m$^{-3}$. TOA during this event was composed almost wholly of primary material: the superspecies run modeled a POA concentration of 205.9 µg m$^{-3}$ (99.78% of TOA concentration) while the control run POA concentration was 117.1 µg m$^{-3}$ (99.75%). Rather than error compounding and leading to divergence from the control run, the superspecies run restabilized without error accumulation for the rest of the simulation: TOA concentration in the superspecies run converged to that of the control run.

### 4.3. Speed Improvement

The advection operator has an outer for-loop over all tracers that are transported. Using superspecies instead of VBS tracers reduces the number of passes through the outer for-loop. With the superspecies selected in Section 3,

16 superspecies (two gas and two particle superspecies for each of the four VBS classes) are advected rather than the 58 VBS tracers, reducing the total number of advected tracers from 104 to 62. The MACC run on the small domain was run sequentially on one computational node. Figure 8 shows wall time for the advection operator when advecting superspecies rather than VBS tracers was 6,790 s, 56% of the time of (1.8 times faster than) the 12,073 s to advect all tracers in the control run. The high resolution required for CAMS operational forecasts increases the computational intensity of the simulations which were performed using domain decomposition over 24 computing nodes with each node computing a subdomain of 175 by 70 grid cells. Using the VBS on the CAMS domain, advection wall time more than doubled from 34959 to 74,762 s. With superspecies advected instead of VBS tracers, wall time for the advection operator was then reduced to 49,473 s. Advecting superspecies on the CAMS domain took about 66% of the time that advecting all the VBS tracers took, a speedup of approximately 1.5.

The timing results suggest that advection wall time depends linearly on number of tracers, which is expected behavior given the structure of the advection operator: an outer for-loop over all tracers. Compared to a run with no OA, inclusion of 58 VBS tracers increases the total number of advected tracers from 42 to 104 and more than doubles the computation time of the advection operator. Advecting 16 superspecies in place of 58 VBS tracers brings the total number of advected tracers down to 62: the proportion of 62/104 yields an expected 59% speed up, in between the speedup results on the MACC and CAMS domains.

## 5. Conclusions

Modeling of organic aerosol processes via four VBS classes is high-dimensional and computationally expensive in LOTOS-EUROS v2.2.1, slowing the advection operator down by a factor of 2. This work developed data-driven methods to reduce the dimension of VBS tracers to a set of superspecies and reduce the computational burden on the advection operator. These methods were refined to ensure physical consistency, including semi-positive constraints, mass conservation, and information on phase. Multiple approaches were compared in Section 3 and non-negative matrix factorization additionally constrained to conserve mass and phase (Approach 4), after being evaluated on reconstruction accuracy and physical consistency, was selected to be implemented in LOTOS-EUROS v2.2.1 in Section 4. Approach 4 creates 16 phase-specific, class-specific superspecies, a compression factor of 3.6, while preserving phase and conserving total concentration to machine precision. The superspecies parameterization ran stably without runaway error for a model simulation of 2 weeks, exceeding the training time horizon. Higher bias of total OA concentration was shown when the superspecies, optimized to reconstruct winter OA patterns, were used for a 2-week run in the summer. During the summer run, the bias showed a clear spatiotemporal pattern, with biogenic SOA overestimated over forests at night. The superspecies were retrained on model output from summer conditions and implemented in LOTOS-EUROS v2.2.1 to reduce high bias. The resuts of this case study indicate that the superspecies might work best when optimized for season-specific conditions.

We found that the superspecies trained on the coarse-resolution MACC domain performed well when used on the fine-resolution domain used in CAMS operational forecasts for a period of 2 weeks. In an analysis period of 4 days performed at the end of the 2-week CAMS run, over 99% of all grid cells showed an absolute bias of time-averaged TOA within the maximum error of the MACC grid. Evaluating a grid cell that exceeded the maximum average error, we found that high overestimation of total OA concentration occurred at a high OA event, and converged back to the baseline simulation as time progressed rather than displaying continued error growth.

Advecting superspecies reduced the wall time spent on the advection operator: advecting superspecies took 56%–66% of the time that it took to advect VBS tracers. Timing experiments indicate a linear dependence of wall time on number of tracers to advect, an expected relation from the structure of the advection operator, which uses a for-loop over all advected tracers. With linear dependence demonstrated, the design choice of compression factor (number of superspecies) can already give an estimate of theoretical speedup.

The use of physically consistent data-driven methods to find superspecies allows for inclusion of organic aerosol processes without doubling the computational burden on the advection operator. Though this approach has been demonstrated for 2-week forecasts on lower and higher resolutions, more work would have to be done to assess this method on longer timescales, including other seasons or seasonal transitions. This case study has not explored how the superspecies method might interface with other tracer compression methods such as the partitioning-based compression method in Section 2.2, or how the superspecies may perform when used in

other processes. However, preserving information on phase of the superspecies allows for their future use in phase-specific processes such as dry deposition, which can be computationally intensive in LOTOS-EUROS. Though demonstrated on organic aerosol species in a regional CTM as a case study, the focus of this approach on physical consistency and interpretability is relevant to other tracers, processes, and models. As physical consistency and computational efficiency are widely desired aspects of numerical modeling in the physical sciences, this approach could be adapted for use in comprehensive Earth system models with the purpose of providing forecasts of global atmospheric composition, for example, GEOS-CF (Keller et al., 2021). More generally, this approach contributes additional physical consistency to a widely used dimensionality reduction technique (non-negative matrix factorization) that can be used to reversibly map between high and low detail in Earth system models.

## Data Availability Statement

The open source, most current version of LOTOS-EUROS is available online as detailed in Manders et al. (2017). The exact version of LOTOS-EUROS v2.2.1 used to generate the model output in this work, including the superspecies extension, as well as all Python code used for developing the data-driven approaches, analysis of model output, and figure generation, is available at Sturm (2022): https://doi.org/10.5281/zenodo.6601166.

## References

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., et al. (2016). Tensorflow: A system for large-scale machine learning. In *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)* (pp. 265–283).

Beltman, J. B., Hendriks, C., Tum, M., & Schaap, M. (2013). The impact of large scale biomass production on ozone air pollution in Europe. *Atmospheric Environment*, *71*, 352–363. https://doi.org/10.1016/j.atmosenv.2013.02.019

Bergström, R., Denier van der Gon, H. A. C., Prévôt, A. S. H., Yttri, K. E., & Simpson, D. (2012). Modelling of organic aerosols over Europe (2002–2007) using a volatility basis set (VBS) framework: Application of different assumptions regarding the formation of secondary organic aerosol. *Atmospheric Chemistry and Physics*, *12*(18), 8499–8527. https://doi.org/10.5194/acp-12-8499-2012

Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforcing analytic constraints in neural networks emulating physical systems. *Physical Review Letters*, *126*, 098302. https://doi.org/10.1103/PhysRevLett.126.098302

Brasseur, G. P., & Jacob, D. J. (2017). *Modeling of atmospheric chemistry*. Cambridge University Press. https://doi.org/10.1017/9781316544754

Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network parametrization trained by coarse-graining. *Journal of Advances in Modeling Earth Systems*, *11*(8), 2728–2744. https://doi.org/10.1029/2019MS001711

Chollet, F. (2015). *Keras*. GitHub. Retrieved from https://github.com/fchollet/keras

Ciarelli, G., Aksoyoglu, S., El Haddad, I., Bruns, E. A., Crippa, M., Poulain, L., et al. (2017). Modelling winter organic aerosol at the European scale with CAMX: Evaluation and source apportionment with a VBS parameterization based on novel wood burning smog chamber experiments. *Atmospheric Chemistry and Physics*, *17*(12), 7653–7669. https://doi.org/10.5194/acp-17-7653-2017

Colette, A., Andersson, C., Manders, A., Mar, K., Mircea, M., Pay, M.-T., et al. (2017). Eurodelta-trends, a multi-model experiment of air quality hindcast in Europe over 1990–2010. *Geoscientific Model Development*, *10*(9), 3255–3276. https://doi.org/10.5194/gmd-10-3255-2017

Courant, R., Friedrichs, K., & Lewy, H. (1928). Über die partiellen Differenzengleichungen der mathematischen Physik. *Mathematische Annalen*, *100*(1), 32–74. https://doi.org/10.1007/BF01448839

Courant, R., Friedrichs, K., & Lewy, H. (1967). On the partial difference equations of mathematical physics. *IBM Journal of Research and Development*, *11*(2), 215–234. https://doi.org/10.1147/rd.112.0215

de Gouw, J. A., Middlebrook, A. M., Warneke, C., Goldan, P. D., Kuster, W. C., Roberts, J. M., & Bates, T. S. (2005). Budget of organic carbon in a polluted atmosphere: Results from the New England air quality study in 2002. *Journal of Geophysical Research*, *110*(D16), D16305. https://doi.org/10.1029/2004JD005623

Denier van der Gon, H. A. C., Bergström, R., Fountoukis, C., Johansson, C., Pandis, S. N., Simpson, D., & Visschedijk, A. J. H. (2015). Particulate emissions from residential wood combustion in Europe—Revised estimates and an evaluation. *Atmospheric Chemistry and Physics*, *15*(11), 6503–6519. https://doi.org/10.5194/acp-15-6503-2015

Donahue, N. M., Epstein, S. A., Pandis, S. N., & Robinson, A. L. (2011). A two-dimensional volatility basis set: 1. Organic-aerosol mixing thermodynamics. *Atmospheric Chemistry and Physics*, *11*(7), 3303–3318. https://doi.org/10.5194/acp-11-3303-2011

Donahue, N. M., Henry, K. M., Mentel, T. F., Kiendler-Scharr, A., Spindler, C., Bohn, B., et al. (2012). Aging of biogenic secondary organic aerosol via gas-phase oh radical reactions. *Proceedings of the National Academy of Sciences of the United States of America*, *109*(34), 13503–13508. https://doi.org/10.1073/pnas.1115186109

Donahue, N. M., Robinson, A., Stanier, C., & Pandis, S. (2006). Coupled partitioning, dilution, and chemical aging of semivolatile organics. *Environmental Science & Technology*, *40*(8), 2635–2643. https://doi.org/10.1021/es052297c

EEA. (2005). *Image2000 and CLC2000. Products and methods. Corine land cover updating for the year 2000*. Ispra Italy.

Gery, M. W., Whitten, G. Z., Killus, J. P., & Dodge, M. C. (1989). A photochemical kinetics mechanism for urban and regional scale computer modeling. *Journal of Geophysical Research*, *94*(D10), 12925. https://doi.org/10.1029/JD094iD10p12925

Harder, P., Watson-Parris, D., Stier, P., Strassel, D., Gauger, N., & Keuper, J. (2022). Physics-informed learning of aerosol microphysics. *Environmental Data Science*, *1*, E20. https://doi.org/10.1017/eds.2022.22

Hayes, P. L., Carlton, A. G., Baker, K. R., Ahmadov, R., Washenfelder, R. A., Alvarez, S., et al. (2015). Modeling the formation and aging of secondary organic aerosols in Los Angeles during calnex 2010. *Atmospheric Chemistry and Physics*, *15*(10), 5773–5801. https://doi.org/10.5194/acp-15-5773-2015

Heald, C. L., Jacob, D. J., Park, R. J., Russell, L. M., Huebert, B. J., Seinfeld, J. H., et al. (2005). A large organic aerosol source in the free troposphere missing from current models. *Geophysical Research Letters*, *32*(18), L18809. https://doi.org/10.1029/2005GL023831

Hodzic, A., Kasibhatla, P. S., Jo, D. S., Cappa, C. D., Jimenez, J. L., Madronich, S., & Park, R. J. (2016). Rethinking the global secondary organic aerosol (SOA) budget: Stronger production, faster removal, shorter lifetime. *Atmospheric Chemistry and Physics*, *16*(12), 7917–7941. https://doi.org/10.5194/acp-16-7917-2016

Im, U., Bianconi, R., Solazzo, E., Kioutsioukis, I., Badia, A., Balzarini, A., et al. (2015). Evaluation of operational online-coupled regional air quality models over Europe and North America in the context of aqmeii phase 2. Part II: Particulate matter. *Atmospheric Environment*, *115*, 421–441. https://doi.org/10.1016/j.atmosenv.2014.08.072

Janssen, R. H. H., Tsimpidi, A. P., Karydis, V. A., Pozzer, A., Lelieveld, J., Crippa, M., et al. (2017). Influence of local production and vertical transport on the organic aerosol budget over Paris. *Journal of Geophysical Research: Atmospheres*, *122*(15), 8276–8296. https://doi.org/10.1002/2016JD026402

Jathar, S. H., Gordon, T. D., Hennigan, C. J., Pye, H. O. T., Pouliot, G., Adams, P. J., et al. (2014). Unspeciated organic emissions from combustion sources and their influence on the secondary organic aerosol budget in the United States. *Proceedings of the National Academy of Sciences of the United States of America*, *111*(29), 10473–10478. https://doi.org/10.1073/pnas.1323740111

Jiang, J., Aksoyoglu, S., El-Haddad, I., Ciarelli, G., Denier van der Gon, H. A. C., Canonaco, F., et al. (2019). Sources of organic aerosols in Europe: A modeling study using CAMX with modified volatility basis set scheme. *Atmospheric Chemistry and Physics*, *19*(24), 15247–15270. https://doi.org/10.5194/acp-19-15247-2019

Jimenez, J. L., Canagaratna, M. R., Donahue, N. M., Prevot, A. S. H., Zhang, Q., Kroll, J. H., et al. (2009). Evolution of organic aerosols in the atmosphere. *Science*, *326*(5959), 1525–1529. https://doi.org/10.1126/science.1180353

Kaiser, J. W., Heil, A., Andreae, M. O., Benedetti, A., Chubarova, N., Jones, L., et al. (2012). Biomass burning emissions estimated with a global fire assimilation system based on observed fire radiative power. *Biogeosciences*, *9*(1), 527–554. https://doi.org/10.5194/bg-9-527-2012

Keller, C. A., & Evans, M. J. (2019). Application of random forest regression to the calculation of gas-phase chemistry within the GEOS-Chem chemistry model v10. *Geoscientific Model Development*, *12*, 1209–1225. https://doi.org/10.5194/gmd-12-1209-2019

Keller, C. A., Knowland, K. E., Duncan, B. N., Liu, J., Anderson, D. C., Das, S., et al. (2021). Description of the NASA geos composition forecast modeling system geos-cf v1.0. *Journal of Advances in Modeling Earth Systems*, *13*(4), e2020MS002413. https://doi.org/10.1029/2020MS002413

Kelp, M. M., Jacob, D. J., Kutz, J. N., Marshall, J. D., & Tessum, C. W. (2020). Toward stable, general machine-learned models of the atmospheric chemical system. *Journal of Geophysical Research: Atmospheres*, *125*(23), e2020JD032759. https://doi.org/10.1029/2020JD032759

Kelp, M. M., Jacob, D. J., Lin, H., & Sulprizio, M. P. (2022). An online-learned neural network chemical solver for stable long-term global simulations of atmospheric chemistry. *Journal of Advances in Modeling Earth Systems*, *14*(6), e2021MS002926. https://doi.org/10.1029/2021MS002926

Kelp, M. M., Tessum, C. W., & Marshall, J. D. (2018). Orders-of-magnitude speedup in atmospheric chemistry modeling through neural network-based emulation. arXiv preprint arXiv:1808.03874.

Knote, C., Hodzic, A., & Jimenez, J. L. (2015). The effect of dry and wet deposition of condensable vapors on secondary organic aerosols concentrations over the continental us. *Atmospheric Chemistry and Physics*, *15*(1), 1–18. https://doi.org/10.5194/acp-15-1-2015

Köble, R., & Seufert, G. (2001). Novel maps for forest tree species in Europe. In *Proceedings of the 8th European symposium on the physico-chemical behavior of air pollutants: A changing atmosphere* (pp. 17–20).

Lane, T. E., Donahue, N. M., & Pandis, S. N. (2008). Effect of no x on secondary organic aerosol concentrations. *Environmental Science & Technology*, *42*(16), 6022–6027. https://doi.org/10.1021/es703225a

Lee, D. D., & Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, *401*(6755), 788–791. https://doi.org/10.1038/44565

Li, M., Soltanolkotabi, M., & Oymak, S. (2020). Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. In S. Chiappa & R. Calandra (Eds.), *Proceedings of the twenty third international conference on artificial intelligence and statistics*, *PMLR* (Vol. 108, pp. 4313–4324). Retrieved from https://proceedings.mlr.press/v108/li20j.html

Liao, H., Henze, D. K., Seinfeld, J. H., Wu, S., & Mickley, L. J. (2007). Biogenic secondary organic aerosol over the United States: Comparison of climatological simulations with observations. *Journal of Geophysical Research*, *112*(D6), D06201. https://doi.org/10.1029/2006jd007813

Lu, Q., Murphy, B. N., Qin, M., Adams, P. J., Zhao, Y., Pye, H. O. T., et al. (2020). Simulation of organic aerosol formation during the calnex study: Updated mobile emissions and secondary organic aerosol parameterization for intermediate-volatility organic compounds. *Atmospheric Chemistry and Physics*, *20*(7), 4313–4332. https://doi.org/10.5194/acp-20-4313-2020

Manders, A. M. M., Builtjes, P. J. H., Curier, L., Denier van der Gon, H. A. C., Hendriks, C., Jonkers, S., et al. (2017). Curriculum vitae of the lotos–euros (v2.0) chemistry transport model. *Geoscientific Model Development*, *10*(11), 4145–4173. https://doi.org/10.5194/gmd-10-4145-2017

Manders-Groot, A. M. M., Segers, A. J., & Jonkers, S. (2021). *Lotos-euros v2.2.000 reference guide*. TNO Reports.

Marais, E. A., Jacob, D. J., Jimenez, J. L., Campuzano-Jost, P., Day, D. A., Hu, W., et al. (2016). Aqueous-phase mechanism for secondary organic aerosol formation from isoprene: Application to the southeast United States and co-benefit of $SO_2$ emission controls. *Atmospheric Chemistry and Physics*, *16*(3), 1603–1618. https://doi.org/10.5194/acp-16-1603-2016

Marsland, S. (2014). *Machine learning: An algorithmic perspective* (2nd ed.). Chapman & Hall/CRC.

Matsui, H. (2017). Development of a global aerosol model using a two-dimensional sectional method: 1. Model design. *Journal of Advances in Modeling Earth Systems*, *9*(4), 1921–1947. https://doi.org/10.1002/2017ms000936

Mircea, M., Bessagnet, B., D'Isidoro, M., Pirovano, G., Aksoyoglu, S., Ciarelli, G., et al. (2019). Eurodelta III exercise: An evaluation of air quality models' capacity to reproduce the carbonaceous aerosol. *Atmospheric Environment X*, *2*, 100018. https://doi.org/10.1016/j.aeaoa.2019.100018

Murphy, B. N., Donahue, N. M., Fountoukis, C., Dall'Osto, M., O'Dowd, C., Kiendler-Scharr, A., & Pandis, S. N. (2012). Functionalization and fragmentation during ambient organic aerosol aging: Application of the 2-d volatility basis set to field studies. *Atmospheric Chemistry and Physics*, *12*(22), 10797–10816. https://doi.org/10.5194/acp-12-10797-2012

Murphy, B. N., & Pandis, S. N. (2009). Simulating the formation of semivolatile primary and secondary organic aerosol in a regional chemical transport model. *Environmental Science & Technology*, *43*(13), 4722–4728. https://doi.org/10.1021/es803168a

Nagori, J., Janssen, R. H. H., Fry, J. L., Krol, M., Jimenez, J. L., Hu, W., & Vilà-Guerau de Arellano, J. (2019). Biogenic emissions and land–atmosphere interactions as drivers of the daytime evolution of secondary organic aerosol in the southeastern us. *Atmospheric Chemistry and Physics*, *19*(2), 701–729. https://doi.org/10.5194/acp-19-701-2019

Ng, N. L., Kroll, J. H., Keywood, M. D., Bahreini, R., Varutbangkul, V., Flagan, R. C., et al. (2006). Contribution of first-versus second-generation products to secondary organic aerosols formed in the oxidation of biogenic hydrocarbons. *Environmental Science & Technology*, *40*(7), 2283–2297. https://doi.org/10.1021/es052269u

O'Dowd, C., Ceburnis, D., Ovadnevaite, J., Vaishya, A., Rinaldi, M., & Facchini, M. (2014). Do anthropogenic, continental or coastal aerosol sources impact on a marine aerosol signature at mace head? *Atmospheric Chemistry and Physics*, *14*(19), 10687–10704. https://doi.org/10.5194/acp-14-10687-2014

Odum, J. R., Hoffmann, T., Bowman, F., Collins, D., Flagan, R. C., & Seinfeld, J. H. (1996). Gas/Particle Partitioning and Secondary Organic Aerosol Yields. *Environmental Science & Technology*, *30*, 2580–2585.

Ots, R., Young, D. E., Vieno, M., Xu, L., Dunmore, R. E., Allan, J. D., et al. (2016). Simulating secondary organic aerosol from missing diesel-related intermediate-volatility organic compound emissions during the clean air for London (clearflo) campaign. *Atmospheric Chemistry and Physics*, *16*(10), 6453–6473. https://doi.org/10.5194/acp-16-6453-2016

Paatero, P., & Tapper, U. (1994). Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, *5*(2), 111–126. https://doi.org/10.1002/env.3170050203

Pai, S. J., Heald, C. L., Pierce, J. R., Farina, S. C., Marais, E. A., Jimenez, J. L., et al. (2020). An evaluation of global organic aerosol schemes using airborne observations. *Atmospheric Chemistry and Physics*, *20*(5), 2637–2665. https://doi.org/10.5194/acp-20-2637-2020

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*, 2825–2830.

Pye, H. O. T., Murphy, B. N., Xu, L., Ng, N. L., Carlton, A. G., Guo, H., et al. (2017). On the implications of aerosol liquid water and phase separation for organic aerosol mass. *Atmospheric Chemistry and Physics*, *17*(1), 343–369. https://doi.org/10.5194/acp-17-343-2017

Pye, H. O. T., Pinder, R. W., Piletic, I. R., Xie, Y., Capps, S. L., Lin, Y.-H., et al. (2013). Epoxide pathways improve model predictions of isoprene markers and reveal key role of acidity in aerosol formation. *Environmental Science & Technology*, *47*(19), 11056–11064. https://doi.org/10.1021/es402106h

Rasp, S. (2020). Coupled online learning as a way to tackle instabilities and biases in neural network parameterizations: General algorithms and Lorenz 96 case study (v1. 0). *Geoscientific Model Development*, *13*(5), 2185–2196. https://doi.org/10.5194/gmd-13-2185-2020

Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models. *Proceedings of the National Academy of Sciences of the United States of America*, *115*(39), 9684–9689. https://doi.org/10.1073/pnas.1810286115

Robinson, A. L., Donahue, N. M., Shrivastava, M. K., Weitkamp, E. A., Sage, A. M., Grieshop, A. P., et al. (2007). Rethinking organic aerosols: Semivolatile emissions and photochemical aging. *Science*, *315*(5816), 1259–1262. https://doi.org/10.1126/science.1133061

Schell, B., Ackermann, I. J., Hass, H., Binkowski, F. S., & Ebel, A. (2001). Modeling the formation of secondary organic aerosol within a comprehensive air quality model system. *Journal of Geophysical Research*, *106*(D22), 28275–28293. https://doi.org/10.1029/2001JD000384

Schreck, J. S., Becker, C., Gagne, D. J., Lawrence, K., Wang, S., Mouchel-Vallon, C., et al. (2022). Neural network emulation of the formation of organic aerosols based on the explicit gecko-a chemistry model. *Journal of Advances in Modeling Earth Systems*, *14*(10), e2021MS002974. https://doi.org/10.1029/2021MS002974

Seinfeld, J. H., & Pandis, S. N. (2006). *Atmospheric chemistry and physics: From air pollution to climate change*. N.J. J. Wiley.

Shrivastava, M. K., Lane, T. E., Donahue, N. M., Pandis, S. N., & Robinson, A. L. (2008). Effects of gas particle partitioning and aging of primary emissions on urban and regional organic aerosol concentrations. *Journal of Geophysical Research*, *113*(D18), D18301. https://doi.org/10.1029/2007jd009735

Sturm, P. O. (2021). *Advecting superspecies: Reduced order modeling of organic aerosols in lotos-euros using machine learning*. TU Delft Education Repository. Retrieved from http://resolver.tudelft.nl/uuid:2c3be50e-5340-4495-a0b7-1670db9be329

Sturm, P. O. (2022). Code for Sturm et al. advecting superspecies [Software]. Zenodo. https://doi.org/10.5281/zenodo.6601166

Sturm, P. O., & Wexler, A. S. (2020). A mass- and energy-conserving framework for using machine learning to speed computations: A photochemistry example. *Geoscientific Model Development*, *13*(9), 4435–4442. https://doi.org/10.5194/gmd-13-4435-2020

Sturm, P. O., & Wexler, A. S. (2022). Conservation laws in a neural network architecture: Enforcing the atom balance of a Julia-based photochemical model (v0.2.0). *Geoscientific Model Development*, *15*(8), 3417–3431. https://doi.org/10.5194/gmd-15-3417-2022

Theodoritsi, G. N., & Pandis, S. N. (2019). Simulation of the chemical evolution of biomass burning organic aerosol. *Atmospheric Chemistry and Physics*, *19*(8), 5403–5415. https://doi.org/10.5194/acp-19-5403-2019

Timmermans, R., van Pinxteren, D., Kranenburg, R., Hendriks, C., Fomba, K., Herrmann, H., & Schaap, M. (2022). Evaluation of modelled lotos-euros with observational based pm10 source attribution. *Atmospheric Environment*, *14*, 100173. https://doi.org/10.1016/j.aeaoa.2022.100173

Tsimpidi, A. P., Karydis, V. A., Pozzer, A., Pandis, S. N., & Lelieveld, J. (2014). Oracle (v1. 0): Module to simulate the organic aerosol composition and evolution in the atmosphere. *Geoscientific Model Development*, *7*(6), 3153–3172. https://doi.org/10.5194/gmd-7-3153-2014

Tsimpidi, A. P., Karydis, V. A., Zavala, M., Lei, W., Molina, L., Ulbrich, I. M., et al. (2010). Evaluation of the volatility basis-set approach for the simulation of organic aerosol formation in the Mexico City metropolitan area. *Atmospheric Chemistry and Physics*, *10*(2), 525–546. https://doi.org/10.5194/acp-10-525-2010

Walcek, C. J. (2000). Minor flux adjustment near mixing ratio extremes for simplified yet highly accurate monotonic calculation of tracer advection. *Journal of Geophysical Research*, *105*(D7), 9335–9348. https://doi.org/10.1029/1999jd901142

Yuval, J., O'Gorman, P. A., & Hill, C. N. (2021). Use of neural networks for stable, accurate and physically consistent parameterization of subgrid atmospheric processes with good performance at reduced precision. *Geophysical Research Letters*, *48*(6), e2020GL091363. https://doi.org/10.1029/2020gl091363

Yuval, J., Pritchard, M., Gentine, P., Zanna, L., & Fan, J. (2021). Call for papers on machine learning and Earth system modeling. *Eos*, *102*. https://doi.org/10.1029/2021EO160820

Zanten, M. v., Sauter, F., Wichink Kruit, R., Jaarsveld, J. v., Pul, W. v., & Wichink Kruit, R. (2010). Description of the DEPAC module. In *Dry deposition modelling with DEPAC_GCN2010 (Tech. Rep.)*. Bilthoven, The Netherlands: Rijksinstituut voor volksgezondheid en Milieu, RIVM report 680180001. Retrieved from http://www.rivm.nl/Documenten_en_publicaties/Wetenschappelijk/Rapporten/2010/oktober/Description_of_the_DEPAC_module_Dry_deposition_modelling_with_DEPAC_GCN2010

Zhang, L. (2001). A size-segregated particle dry deposition scheme for an atmospheric aerosol module. *Atmospheric Environment*, *35*(3), 549–560. https://doi.org/10.1016/S1352-2310(00)00326-5

Zhao, B., Shrivastava, M., Donahue, N. M., Gordon, H., Schervish, M., Shilling, J. E., et al. (2020). High concentration of ultrafine particles in the amazon free troposphere produced by organic new particle formation. *Proceedings of the National Academy of Sciences of the United States of America*, *117*(41), 25344–25351. https://doi.org/10.1073/pnas.2006716117