



1           **Calibrating a global atmospheric chemistry transport model using Gaussian process**  
2                   **emulation and ground-level concentrations of ozone and carbon monoxide**

3                                   Edmund Ryan<sup>1,2\*</sup>, Oliver Wild<sup>1</sup>

4  
5   <sup>1</sup>Lancaster Environment Centre, Lancaster University, UK.

6   <sup>2</sup>Now at: School of Mathematics, University of Manchester, UK.

7

8

9   \*Corresponding author:

10   School of Mathematics,

11   The University of Manchester,

12   Alan Turing Building,

13   Oxford Road, Manchester.

14   M13 9PL

15   Tel: +44 (0)161 275 5800

16   [edmund.ryan@manchester.ac.uk](mailto:edmund.ryan@manchester.ac.uk)

17

18   Keywords: atmospheric chemistry transport model, model calibration, surface ozone,

19   measurement representativeness, Gaussian Process emulation, Markov Chain Monte Carlo, Just

20   Another Gibbs Sampler.

21

22

23

24



## 1 **Abstract**

2 Atmospheric chemistry transport models are important tools to investigate the local, regional and  
3 global controls on atmospheric composition and air quality. To ensure that these models  
4 represent the atmosphere adequately it is important to compare their outputs with measurements.  
5 However, ground based measurements of atmospheric composition are typically sparsely  
6 distributed and representative of much smaller spatial scales than those resolved in models, and  
7 thus direct comparison incurs uncertainty. In this study, we investigate the feasibility of using  
8 observations of one or more atmospheric constituents to estimate parameters in chemistry  
9 transport models and to explore how these estimates and their uncertainties depend upon  
10 representation errors and the level of spatial coverage of the measurements. We apply Gaussian  
11 process emulation to explore the model parameter space and use monthly averaged ground-level  
12 concentrations of ozone (O<sub>3</sub>) and carbon monoxide (CO) from across Europe and the US. Using  
13 synthetic observations we find that the estimates of parameters with greatest influence on O<sub>3</sub> and  
14 CO are unbiased, and the associated parameter uncertainties are low even at low spatial coverage  
15 or with high representation error. Using reanalysis data, we find that estimates of the most  
16 influential parameter - corresponding to the dry deposition process - are closer to its expected  
17 value using both O<sub>3</sub> and CO data than using O<sub>3</sub> alone. This is remarkable because it shows that  
18 while CO is largely unaffected by dry deposition, the additional constraints it provides are  
19 valuable for achieving unbiased estimates of the dry deposition parameter. In summary, these  
20 findings identify the level of spatial representation error and coverage needed to achieve good  
21 parameter estimates and highlight the benefits of using multiple constraints to calibrate  
22 atmospheric chemistry models.

23



## 1 **Introduction**

2 Changes in atmospheric composition due to human activities make an important contribution to  
3 Earth's changing climate (Stocker et al., 2013) and to outdoor air pollution, which is currently  
4 responsible for about 4.2 million deaths worldwide each year (Cohen et al., 2017). Chemistry  
5 transport models (CTMs) simulate the production, transport, and removal of key atmospheric  
6 constituents, and are important tools for understanding variations in atmospheric composition  
7 across space and time. They permit investigation of future climate and emission scenarios that  
8 fully account for the interactions and feedbacks that characterise physical, chemical and  
9 dynamical processes in the atmosphere. For practical application, CTMs need to reproduce the  
10 magnitude and variation in pollutant concentrations observed at a wide range of measurement  
11 locations. Where biases occur, these can often be reduced by improving process representation  
12 through adjusting model parameters so that the CTM matches the measurements to a sufficient  
13 level of accuracy (e.g. Menut et al., 2014). While estimation of model parameters is common in  
14 many fields of science, it is rarely attempted with atmospheric chemistry models because they  
15 are computationally expensive to run and it is thus burdensome to perform the large number of  
16 model runs required to explore model parameter space. Instead, data assimilation has become a  
17 standard method for ensuring that model states are consistent with measurements, usually  
18 treating model parameters as fixed (Khattatov et al., 2000, Bocquet et al., 2015, van Loon et al.,  
19 2000, Emili et al., 2014).

20 In this study, we explore computationally efficient ways of estimating parameters in  
21 chemistry transport models, focusing on two important tropospheric constituents, ozone (O<sub>3</sub>) and  
22 carbon monoxide (CO). Ozone is a major pollutant that is produced in the troposphere by  
23 oxidation of precursors such as CO and hydrocarbons, which are emitted during combustion



1 processes from vehicular, industrial and residential sources. Ozone is harmful to human health  
2 and has been shown to damage vegetation and reduce crop yields (Goldsmith and Landaw, 1968,  
3 Kampa and Castanas, 2008, Van Dingenen et al., 2009, van Zelm et al., 2008). A recent  
4 assessment of surface O<sub>3</sub> was carried out for the Tropospheric Ozone Assessment Report  
5 (TOAR) based on measurements from an extensive network of 10,000 sites around the world  
6 (Schultz et al., 2017). A simple statistical model of changes in surface O<sub>3</sub> between 2000 and  
7 2014 showed that significant decreases of 28% and 6% have occurred in Eastern North America  
8 and Europe, respectively, but increases of 20% and 45% in south-east and east Asia (Chang et  
9 al., 2017). In recent decades, a similar pattern of decreases in CO in Europe and North America  
10 and increases over parts of Asia has also been observed (Granier et al., 2011). To fully explain  
11 and attribute these changes, a thorough understanding of the processes controlling these  
12 pollutants is needed.

13 To assess the performance of CTMs, it is essential to compare simulations of  
14 tropospheric chemical composition with measurements. A comprehensive evaluation of 15  
15 global models found that they broadly matched measured O<sub>3</sub>, but that modelled O<sub>3</sub> was biased  
16 high in the northern hemisphere and biased low in the southern hemisphere (Young et al., 2018).  
17 The models were unable to capture the long-term trends in tropospheric O<sub>3</sub> observed at different  
18 altitudes. Similar biases were found in an independent study of long-term trends involving three  
19 chemistry climate models (Parrish et al., 2014). While identification of these model biases is  
20 informative, correcting the deficiencies is challenging because it is often unclear why different  
21 models perform well at certain times and for certain places, but poorly elsewhere (Young et al.,  
22 2018). A practical solution is to perform global sensitivity analysis to identify the parameters or  
23 processes that influence the model results most and then to calibrate the model to estimate these



1 parameters and their uncertainties by comparing model predictions with measurements in a  
2 statistically rigorous way. This provides insight into the physical processes causing model biases  
3 that is typically unavailable from simpler approaches.

4         The principal challenge with performing global sensitivity analysis and model calibration  
5 is that they require thousands of model runs, and this is infeasible for a typical global CTM that  
6 may require 12-24 hours to simulate a year on high performance computing facilities. This can  
7 be overcome by replacing the model with a surrogate function such as a Gaussian process  
8 emulator that is computationally much faster to run (Johnson et al., 2018, Ryan et al., 2018, Lee  
9 et al., 2013). Sensitivity analysis and model calibration can then be performed based on  
10 thousands of runs with the emulator rather than the CTM. Since the first application of  
11 emulation methods for model calibration (Kennedy and O'Hagan, 2001), these approaches have  
12 been extended to models with highly multivariate output. Examples include an earth system  
13 model (Wilkinson, 2010), an aerosol model (Johnson et al., 2015), an ice sheet model (Chang et  
14 al., 2016) and a climate model (Salter et al., 2018). In this study, we apply these approaches to  
15 models of tropospheric ozone for the first time to demonstrate the feasibility of parameter  
16 estimation.

17         We identify three issues that need to be addressed for successful atmospheric model  
18 calibration. Firstly, ground-level composition measurements are usually made at a single location  
19 which may not be representative of a wider region at the grid-scale of the model. Global  
20 chemistry transport models typically have a spatial scale of the order of 100 km. Errors  
21 associated with spatial representativeness may be important even for satellite measurements  
22 which provide information at a 10 km scale (Boersma et al., 2016, Schultz et al., 2017). This  
23 representation error is distinct from instrument error, which is often relatively narrow and better



1 understood. The effect of representation errors was explored in simple terrestrial Carbon model  
2 by Hill et al. (2012), who found that as these errors decreased, the accuracy of parameter  
3 estimates improved.

4 Secondly, the spatial coverage of atmospheric composition measurements is typically  
5 relatively poor, and this limits our ability to estimate parameters accurately. Thus, it is important  
6 to explore how the spatial coverage of measurements affects estimates of model parameters and  
7 their associated uncertainties.

8 Thirdly, evaluation of atmospheric chemistry models is typically performed for different  
9 variables independently (e.g., Stevenson et al., 2006, Fiore et al., 2009). However, atmospheric  
10 constituents such as O<sub>3</sub>, CO, NO<sub>x</sub>, and VOC are often closely coupled through interrelated  
11 chemical, physical and dynamical processes. Evaluation of a model with measurements of a  
12 single species neglects the additional process information available from accounting for species  
13 relationships. Lee et al. (2016) highlight the limitation of using a single observational constraint  
14 on modelled aerosol concentrations, finding that this resulted in reduced uncertainty in  
15 concentrations but not in the associated radiative forcing. The benefits of using multiple  
16 constraints have been highlighted previously. For example Miyazaki et al. (2012) used the  
17 Ensemble Kalman Filter and satellite measurements of NO<sub>2</sub>, O<sub>3</sub>, CO and HNO<sub>3</sub> to constrain a  
18 CTM, resulting in a significant reduction in model bias in NO<sub>2</sub> column, O<sub>3</sub> and CO  
19 concentrations simultaneously. Nicely et al. (2016) used aircraft measurements of O<sub>3</sub>, H<sub>2</sub>O and  
20 NO to constrain a photochemical box model, and found estimates of column OH that were 12-  
21 40% higher than those from unconstrained CTMs. They also found that although the CTMs  
22 simulated O<sub>3</sub> well, they underestimated NO<sub>x</sub> by a factor of two, explaining the discrepancy in  
23 column OH.



1           To address these gaps in knowledge, we estimate the probability distributions of eight  
2 parameters from a CTM, given surface O<sub>3</sub> and CO concentrations from the USA and Europe.  
3 We focus on model calibration with a limited number of parameters as a proof of concept, but  
4 show how this could be expanded to a much wider range of parameters in future. To overcome  
5 the excessive computational burden of running the model a large number of times, we replace the  
6 model with a fast surrogate using Gaussian process emulation. After evaluation of the emulator  
7 to ensure that it is an accurate representation of the input-output relationship of the CTM, we  
8 investigate how well model parameters can be estimated from chemical measurement data. We  
9 quantify the impacts of measurement representation error and spatial coverage on the bias and  
10 uncertainty in the estimated model parameters and highlight the extent to which parameter  
11 estimates can be improved using measurements of different variables simultaneously.

## 12 **2. Materials and methods**

### 13 *2.1 Atmospheric Chemical Transport Model*

14 Chemistry transport models simulate the changes in concentration of a range of atmospheric  
15 constituents (e.g. O<sub>3</sub>, CO, NO<sub>x</sub>, CH<sub>4</sub>) with time over a specified three-dimensional domain. They  
16 represent many of the physical and chemical processes involved, usually in a simplified form,  
17 but a detailed understanding is often incomplete. Key processes include the emission of trace  
18 gases into the atmosphere, photochemical reactions that result in chemical transformations,  
19 transport by the winds, convection and turbulence, and removal of trace gases from the  
20 atmosphere through deposition processes. In this study, we apply the Frontier Research System  
21 for Global Change version of the University of California, Irvine chemical transport model, the  
22 FRSGC/UCI CTM (Wild and Prather, 2000; Wild et al, 2004). We focus on eight important  
23 processes affecting tropospheric oxidants that were chosen based on one-at-a-time sensitivity



1 studies with the model (Wild, 2007) and that have been used in previous global sensitivity  
2 analyses of tropospheric ozone burden and methane lifetime (Ryan et al., 2018; Wild et al.,  
3 2020). These processes include the surface emissions of nitrogen oxides (NO<sub>x</sub>), lightning  
4 emissions of NO, biogenic emissions of isoprene, wet and dry deposition of atmospheric  
5 constituents, atmospheric humidity, cloud optical depth and the efficiency of turbulent mixing in  
6 the boundary layer, see Table 1. These do not encompass all sources of uncertainty in the model,  
7 but are broadly representative of major uncertainties across a range of different processes. To  
8 provide a simple and easily interpretable approach to calibration, we define a scaling factor that  
9 spans the range of uncertainty in each process, and these scaling factors form the parameters that  
10 we aim to calibrate. The choice of parameters and uncertainty ranges are described in more detail  
11 in Wild et al. (2020). For this study, we focus on monthly-mean surface O<sub>3</sub> and CO distributions  
12 at the model native grid resolution of 2.8°×2.8° and compare with observations over North  
13 America and Europe for model calibration (Fig. 1). The model uses meteorological driving data  
14 for 2001, a relatively typical meteorological year without strong climate phenomena such as El  
15 Nino (Fiore et al. 2009).

## 16 2.2 Surface O<sub>3</sub> and CO data

17 Ground-based observations of O<sub>3</sub> are relatively abundant in Europe and North America, where  
18 there are ~1800 individual sites that have continuous long-term measurements of O<sub>3</sub> (Chang et  
19 al., 2016, Schultz et al., 2017). Measurements of CO are made at fewer locations, but reliable  
20 long-term data are available from 57 sites that are part of the Global Atmospheric Watch  
21 network (Schultz et al., 2015). To allow more thorough testing of the effects of spatial coverage  
22 over these regions, we use model reanalysis data of surface O<sub>3</sub> and CO from the European Centre  
23 for Medium-Range Weather Forecasts (ECMWF) which has been tuned to match measurements





1 using 4D-Var data assimilation (Flemming et al. 2017). This reanalysis data closely resembles  
2 observed O<sub>3</sub> and CO where measurements are available and has the benefit of complete global  
3 coverage, allowing us to test the importance of measurement coverage directly.

4 Reanalysis data for O<sub>3</sub> and CO are available for 2003–2015, and we average the data by  
5 month across this period to provide a climatological comparison. The control run of the  
6 FRSGC/UCI model matches CO from the reanalysis data reasonably well (Fig. 2), but  
7 overestimates surface O<sub>3</sub>. Overestimation of O<sub>3</sub> in continental regions has been noted in previous  
8 studies and is partly a consequence of rapid photochemical formation from fresh emissions that  
9 is magnified at coarse model resolution (Wild and Prather, 2006). For this exploratory study we  
10 bias-correct the modelled surface O<sub>3</sub> by reducing it by 25%, following the approach taken by  
11 Shindell et al. (2018), so that it matches the reanalysis data (Fig. 2a). This adjustment accounts  
12 for the effect of chemical processes and model resolution which are not explored in this study,  
13 and provides a firmer foundation for investigating the effects of other processes.

### 14 *2.3 Representation error*

15 The “representation error” describes how well measurements made at a single location represent  
16 a wider region at the spatial scale of the model (2.8°×2.8° for this study). The error may be  
17 reduced by averaging measurements made at different stations within a model grid box, although  
18 atmospheric measurements may be too sparse to permit this (Schultz, 2016). The representation  
19 error is sometimes taken as the mean of the spatial standard deviation of different measurements  
20 within a grid-box (Sofen et al. 2016). However, this measure quantifies the spatial variability of  
21 measured O<sub>3</sub> within a grid-box and may not match the representation error

22 To test the effect on parameter estimates of varying this representation error, we use  
23 synthetic data from the control run of the model using parameters set to their nominal default



1 values. Synthetic O<sub>3</sub> and CO data were generated by adding different levels of representation  
2 error for each level of spatial coverage. In mathematical terms:

$$data_i = m_i(x_{control}) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma_i^2) \quad (1)$$

3 where for the *i*th point in space or time, *data<sub>i</sub>* refers to the synthetic data for O<sub>3</sub> or CO,  
4 *m<sub>i</sub>(x<sub>control</sub>)* is the O<sub>3</sub> or CO from the model control run, and *ε<sub>i</sub>* is generated from a Normal  
5 distribution with mean of zero and standard deviation *σ<sub>i</sub>* that is directly proportional to the  
6 magnitude of *m<sub>i</sub>(x<sub>control</sub>)*. In this case, *σ<sub>i</sub>* = *p* × *m<sub>i</sub>(x<sub>control</sub>)* where *p* is a representation error  
7 scaling factor that we varied. We included *p* as one of the parameters to estimate for the  
8 reanalysis data and found values in the range 0.16–0.19. Thus, when using the synthetic data we  
9 set the representation error scaling factor for these variables to *p* = 0.01, 0.1, 0.2 and 0.3.

#### 10 2.4 Global sensitivity analysis

11 Sensitivity analysis was carried out to determine the sensitivity of the simulated surface O<sub>3</sub> and  
12 CO to changes in each of the eight parameters. This allows us to identify which of the  
13 parameters are most important in governing surface O<sub>3</sub> and CO. We use global sensitivity  
14 analysis (GSA), varying each input while averaging over the other inputs. This provides a more  
15 integrated assessment of uncertainty than the traditional one-at-a-time approach varying each  
16 input in turn while fixing the other inputs at nominal values. We use the extended FAST method  
17 (Saltelli et al., 1999), a common and robust approach to GSA in which the sensitivity indices are  
18 quantified by partitioning the total variance in the model output (i.e. modelled surface O<sub>3</sub> or CO)  
19 into different sources of contribution from each input. Like most sensitivity analysis methods,  
20 this approach requires several thousand executions of the model, which would be  
21 computationally expensive for the CTM used here. This is overcome by replacing the CTM with



1 a Gaussian process (GP) emulator. Further details of the implementation of GSA are described  
2 in Ryan et al. et al. (2018).

### 3 *2.5 Gaussian Process Emulation - theory*

4 We replace the CTM with a surrogate model that maps the inputs of the CTM (the eight  
5 parameters listed in Table 1) with its outputs (surface O<sub>3</sub> and CO). We employ a surrogate  
6 model based on Gaussian process (GP) emulation for three reasons. Firstly, due to the attractive  
7 mathematical properties of a GP, the emulator needs very few runs of the computationally  
8 expensive model to train it, typically less than 100. In contrast, methods based on neural  
9 networks can require thousands of training runs. Secondly, a GP emulator is an interpolator and  
10 so predicts the output of the model with no uncertainty at the input points it is trained at.  
11 Thirdly, it gives a complete probability distribution, as a measure of uncertainty, for estimates of  
12 the model output at points it is not trained at.

13 A GP is an extension of the multivariate Gaussian distribution, where instead of a mean  
14 vector  $\mu$  and covariance matrix  $\Sigma$ , mean and covariance functions given by  $E(f(x))$  and  
15  $\text{cov}(f(x), f(x'))$  are used (Rasmussen, 2006). Here,  $f(\cdot): \chi \in \mathbb{R}^q \rightarrow \mathbb{R}^q$  represents the  
16 computationally expensive model and  $\chi$  denotes the input space given by  $x = (x_1, \dots, x_q) \in \chi_1 \times$   
17  $\dots \times \chi_q = \chi \subset \mathbb{R}^q$ , and  $q$  is the number of input variables. GP emulators within a Bayesian  
18 framework were first developed in the 1990s and early 2000s (O'Hagan, 2006, Oakley and  
19 O'Hagan, 2004, Kennedy and O'Hagan, 2000, Currin et al., 1991). The simplest and most  
20 common GP emulator is one where the outputs to be emulated are scalar. Thus, if the  
21 computationally expensive model is given by  $f(\cdot)$ , then the one-dimensional output  $y$  is  
22 calculated by  $y = f(x)$ . This means that if the model output is multidimensional – e.g. a global  
23 map or a time-series – then we need to build a separate emulator for each point in the output



1 space. To build the emulator requires training runs from the expensive model. In general, we  
2 choose  $n$  training inputs, denoted by  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , based on a space filling design such as a  
3 Maximin Latin Hypercube design (Morris and Mitchell, 1995). The number of training points is  
4 based on the rule of thumb  $n = 10 \times q$  (Loeppky et al., 2012).

5 Denoting the scalar outputs by  $y_1 = f(\mathbf{x}_1), y_2 = f(\mathbf{x}_2), \dots, y_n = f(\mathbf{x}_n)$ , we then build  
6 an emulator  $\hat{f}(\cdot)$  given by  $\hat{y} = \hat{f}(x)$ , where  $\hat{y}$  is the estimated output from the emulator. If  $x$   
7 represents one of the training inputs (i.e.  $x = \mathbf{x}_i, 1 \leq i \leq n$ ), then  $\hat{y}$  is equal to the output from  
8  $f(\cdot)$  with no uncertainty (i.e.  $\hat{y} = y$ ). If  $x$  represents an input the emulator is not trained at, then  
9  $\hat{y}$  has a probability distribution represented by a mean function  $m(x)$  and a covariance function  
10  $V(x, x')$ , where  $x'$  is a different input. The mean function is given by:

$$m(x) = h(x)^T \hat{\beta} + t(x)^T \mathbf{A}^{-1} (\mathbf{y} - H \hat{\beta}), \quad (2)$$

11 where  $h(x)^T$  is a  $1 \times (q+1)$  vector given by  $(1, x^T)$ ,  $\hat{\beta}$  is a vector of coefficients determined by  
12  $\hat{\beta} = (H^T \mathbf{A}^{-1} H)^{-1} H^T \mathbf{A}^{-1} \mathbf{y}$ ,  $t(x)^T = (C(x, x_1; \psi), \dots, C(x, x_n; \psi))$ , and  $\mathbf{A}$  is a matrix whose ele-  
13 ments are determined by  $\mathbf{A}_{ij} = C(\mathbf{x}_i, \mathbf{x}_j; \psi)$ ,  $\mathbf{y} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]^T$ ,  $H = [h(\mathbf{x}_1), \dots, h(\mathbf{x}_n)]^T$ .  
14 Here,  $C(x, x'; \psi)$  is a correlation function that represents our prior belief about how the inputs  $x$   
15 and  $x'$  are correlated. A common choice is a Gaussian correlation function which takes the form:  
16  $C(x, x'; \psi) = \exp(-(x - x')^T \mathbf{B} (x - x'))$ , where  $\mathbf{B}$  is a  $p \times p$  matrix with zeros in the off-  
17 diagonals and diagonal elements given by the roughness parameters  $\psi = (\psi_1, \dots, \psi_q)$ . These  
18 give an indication of whether the input-output relationship for each input variable, given the  
19 training data, should be linear. Low values reflect a linear (or smooth) relationship, whereas  
20 high values (e.g.  $> 20$ ) suggest a non-linear (or non-smooth) response surface. For  
21 implementation purposes we express the correlation function as  $C(x, x'; \psi) =$



1  $\sum_{j=1}^{q+1} \exp(-\psi_j(x_j - x'_j)^2)$ , where  $x = (x_1, \dots, x_q)$  and  $x' = (x'_1, \dots, x'_q)$ . The formula for the  
2 covariance function  $V(x, x')$  is given in appendix A.

### 3 *2.6 Gaussian Process Emulation - implementation*

4 Using the Loepky rule we choose  $n=80$  different training inputs for our eight-parameter  
5 calibration study. In total, we emulate two variables (surface O<sub>3</sub> and CO) over 12 months at 272  
6 spatial locations, and so require 6528 different GP emulators. To estimate the model parameters  
7 we evaluate each of the GP emulators tens of thousands of times. Although emulation is  
8 computationally fast, this presents a substantial computational burden, even for more  
9 computationally efficient versions of the emulator (Marrel et al., 2011, Roustant et al., 2012).  
10 We overcome this by computing parts of equation (2) prior to these evaluations. Specifically, we  
11 compute the vectors  $\hat{\beta}$ ,  $m_{LP}$  and  $\psi$  for all points in the output space, where  $m_{LP}$  denotes  
12  $\mathbf{A}^{-1}(\mathbf{y} - H\hat{\beta})$ , the last part of  $m(x)$  from equation (2). We store these three objects as three  
13 matrices  $\hat{\beta}_{ALL}$ ,  $m_{LP,ALL}$  and  $\psi_{ALL}$ . Evaluated at a new input  $x_{new}$ , the mean function of the  
14 emulator (equation 1) can now be expressed as:

$$\begin{aligned} m_i(x_{new}) &= h(x_{new})^T \hat{\beta}_{ALL}[i, :] + t_i(x_{new})^T m_{LP,ALL}[i, :], \\ t_i(x_{new})^T &= (C(x_{new}, x_1; \psi_{ALL}[i, :]), \dots, C(x_{new}, x_n; \psi_{ALL}[i, :])), \end{aligned} \quad (3)$$

15 where  $i$  ( $1 \leq i \leq 6528$ ) denotes the  $i$ th point in the output space. The equivalent formula for  
16  $V(x, x')$  is given in appendix A.

17 To the test the accuracy of GP emulation, we ran each of the 6528 emulators at 20 sets of  
18 parameters which were not used for training the emulators. The estimated O<sub>3</sub> and CO values  
19 from the emulators for all spatial locations and months closely match the simulated O<sub>3</sub> and CO



1 output from the FRSGC/UCI model for these validation runs, with  $R^2 > 0.995$  for each variable,  
2 see Fig. 3.

### 3 *2.7 Parameter Estimation*

4 We estimate the eight model parameters using Bayesian statistics via the software package Just  
5 Another Gibbs Sampler (Plummer, 2003). This uses a Markov Chain Monte Carlo (MCMC)  
6 approach to sample from the multi-dimensional posterior probability distribution of the model  
7 parameters (Berg, 2005). To find the posterior distribution, the MCMC algorithm searches the  
8 parameter space using multiple sets of independent chains. Here, a chain refers to a sequence of  
9 steps in the parameter space that the algorithm takes. A new proposed parameter set in this  
10 search is accepted on two conditions: (1) the set is consistent with the prior probability  
11 distribution, which for our study was a set of Uniform distributions with the lower and upper  
12 bounds given by the defined ranges in Table 1; and (2) the resulting modelled values using the  
13 proposed set of parameters are consistent with measurements, which is assessed using the  
14 following Gaussian likelihood function:

$$L(\theta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(f_i(\theta) - m_i)^2}{\sigma_i^2}\right), \quad (4)$$

15 where  $N$  is the number of measurements used,  $f_i(\theta)$  is the  $i$ th model output ( $1 \leq i \leq N$ ) using the  
16 proposed parameter set  $\theta$ ,  $m_i$  is the measurement corresponding to the  $i$ th model output and  $\sigma_i$  is  
17 the representation error for measurement  $m_i$ .

18 We ran three parallel chains for 10,000 iterations each. After discarding the first half of  
19 these iterations as ‘burn in’, we thinned the chains by a factor of five to reduce within-chain  
20 autocorrelation. Convergence was assessed using the Brooks-Gelman-Rubin diagnostic tool  
21 (Gelman et al., 2013). This produced 3000 independent samples from the posterior distribution



1 for each parameter, which we summarize using their posterior means and 95% credible intervals  
2 (CIs) defined by the 2.5<sup>th</sup> and 97.5<sup>th</sup> percentiles (Gelman et al., 2013). We used the R language  
3 to code up our configuration of the MCMC algorithm.

#### 4 *2.8 Experimental approach*

5 We first perform a global sensitivity analysis to identify the parameters which have the greatest  
6 influence on the two variables we consider. We then perform parameter estimation using  
7 measurement data over the regions of North America and Europe shown in Fig 1 and focus our  
8 analysis on the parameters which have the greatest influence. To provide a demonstration of the  
9 approach we first use “synthetic” measurement data drawn from the control run of the CTM  
10 which was not used to train the emulators, adding increasing levels of noise to represent  
11 measurement representation errors of 1, 10, 20 and 30% ( $p = 0.01, 0.1, 0.2$  and  $0.3$ ), and varying  
12 the spatial coverage of these measurements over the regions considered over a wide range: 2.5, 5,  
13 10, 20, 40 and 100%. We focus on surface O<sub>3</sub> only, surface CO only and then both variables  
14 together. We then use the reanalysis data to represent the measurements, focussing on the effects  
15 of spatial coverage alone, and estimating the representation error  $p$  from this independent dataset.  
16 The 90 different scenarios we consider are summarised in Table 2.

### 17 **3. Results**

#### 18 *3.1 Global sensitivity analysis*

19 Results from global sensitivity analysis reveal that over the continental regions of Europe and  
20 North America considered here, the simulated monthly mean concentrations of surface O<sub>3</sub> are  
21 most sensitive to dry deposition and, to a lesser extent, to isoprene emissions (Fig. 4). This is not  
22 unexpected, given the importance of direct deposition of ozone to the Earth’s surface, and the



1 role of isoprene as a natural source of ozone in continental regions. The simulated surface CO is  
2 most sensitive to isoprene emissions, which represent a source of CO, and to boundary layer  
3 mixing, which influences the transport of CO from polluted emission regions. We thus identify  
4 the scaling parameters corresponding to dry deposition, isoprene emissions and boundary layer  
5 mixing as the most important of the eight considered here to estimate accurately to reduce the  
6 bias in modelled surface O<sub>3</sub> and CO. For completeness, we show the geographical distribution  
7 of sensitivity indices in Figs 5 and 6, which reveal the importance of humidity in governing O<sub>3</sub>  
8 over oceanic regions and highlight the very different responses of surface O<sub>3</sub> and CO to the  
9 major driving processes.

### 10 *3.2 Estimation of scaling parameters using synthetic data*

11 We next use synthetic observation data to calibrate the model and estimate scaling parameters.  
12 For synthetic data we use the model control run with a specified level of representation error  
13 (Table 2), and the default model parameters define the true scaling that we aim to retrieve.  
14 Prescribing surface O<sub>3</sub> with very little error ( $p = 0.01$ ) gives an estimate of the dry deposition  
15 scaling parameter, which has the largest influence on modelled surface O<sub>3</sub>, close to its true value  
16 and the uncertainty is small even when the spatial coverage of measurements is only 2.5% (Fig.  
17 7, column 1). As the representation error is increased to  $p = 0.1$ , the parameter uncertainty is  
18 larger at low spatial coverage but the mean estimate remains unbiased (Fig. 7, column 2). The  
19 uncertainty at all levels of spatial coverage becomes larger as  $p$  increases to 0.2 and 0.3, but the  
20 means remain very close to the true values (Fig. 7, columns 3 and 4). Surface CO is largely  
21 unaffected by dry deposition, and thus provides very little constraint on the scaling parameter.  
22 The effect of prescribing surface CO and O<sub>3</sub> together is very similar to that of using surface O<sub>3</sub>  
23 alone.





1           Using surface CO alone with very little representation error ( $p = 0.01$ ), the mean estimate  
2 of the isoprene emission scaling parameter is equal to the true value with very little uncertainty,  
3 regardless of the spatial coverage (Fig. 8, column 1). When the representation error is increased  
4 to  $p = 0.1$ , the estimate remains very close to the true value, but the uncertainty is substantially  
5 higher at low spatial coverage (2.5% and 5%) than at higher coverage (40% and 100%) (Fig. 8,  
6 column 2). The estimates deviate further from the truth at higher levels of representation error ( $p$   
7 = 0.2 and 0.3) and the uncertainty is greater (Fig. 8, columns 3 and 4). Estimates of the isoprene  
8 scaling parameter are less accurate than those of the dry deposition scaling parameter as the  
9 posterior means are further from the true value of the parameter and the uncertainty intervals are  
10 wider (Fig. 8 vs Fig. 7). As with our findings for dry deposition, the posterior means and the  
11 lengths of the uncertainty intervals for the isoprene scaling parameter remain relatively  
12 unchanged when surface O<sub>3</sub> data is prescribed at the same time.

13           Our findings for the boundary layer mixing scaling parameter follow a similar pattern to  
14 the other two parameters (Fig. 9). In all combinations of representation error and spatial  
15 coverage, we find that the mean estimates are unbiased. Furthermore, we find that the parameter  
16 uncertainty is significantly smaller when the spatial coverage is 10% or higher when  $p = 0.1$ ,  
17 20% or higher when  $p = 0.2$ , and 40% or higher when  $p = 0.3$  (Fig. 9, Table 2). It is clear from  
18 these results that the scalings for these three model parameters can be successfully estimated  
19 from synthetic data with low uncertainty when the representation error is low, and that the  
20 estimates remain good, albeit with higher uncertainty, at higher representation error if the spatial  
21 coverage is relatively good.

### 22   3.3 Estimation of scaling parameters using reanalysis data



1 We consider next the reanalysis data for surface O<sub>3</sub> and CO which are based on assimilated  
2 concentrations from the ECMWF model and are thus independent of the FRSGC/UCI model.  
3 The reanalysis is representative of similar spatial scales to the FRSGC/UCI model, and thus we  
4 ignore the representation error and vary the spatial coverage only. However, we are able to  
5 estimate the representation error factor  $p$  by treating it as a parameter to estimate. With 100%  
6 spatial coverage, this error term is estimated with the MCMC algorithm to be  $p = 0.168 \pm 0.004$   
7 and  $p = 0.191 \pm 0.005$  for surface O<sub>3</sub> and CO, respectively. Although we do not know the true  
8 values of the parameters in this case, the good agreement between the control run of the  
9 FRSGC/UCI model and the reanalysis data suggests that they lie close to their true values.

10 Using the reanalysis data for surface O<sub>3</sub> alone, we find that the posterior means and  
11 uncertainty for the dry deposition parameter are in the upper half of the range defined, indicating  
12 that the real dry deposition flux is greater than that calculated with the FRSGC/UCI model. This  
13 is largely as expected, as the FRSGC/UCI model overestimates surface O<sub>3</sub> at these continental  
14 sites and greater deposition would bring the model into better agreement with the reanalysis. As  
15 the spatial coverage is increased, the estimate of the scaling factor increases to around 1.4 and  
16 the uncertainty is reduced (Fig. 10a). In contrast, using surface O<sub>3</sub> and CO together results in an  
17 estimate closer to 1 and an additional reduction in uncertainty (Fig. 10g). Inclusion of surface  
18 CO measurements, as an additional constraint to surface O<sub>3</sub>, results in an estimate of the dry  
19 deposition parameter closer to that modelled along with a reduction in the associated uncertainty.

20 Using surface CO alone, estimates of the isoprene scaling parameter lie in the central part  
21 of the defined range, whilst estimates of the boundary layer mixing scaling parameter lie in the  
22 upper half of the defined range (Fig 10e,f). For both parameters, increasing the spatial coverage  
23 leads to a reduction in uncertainty. Unlike for dry deposition, inclusion of surface O<sub>3</sub> when



1 estimating either of these parameters results in very little difference in the magnitude of the  
2 estimate or in the associated uncertainty (Fig. 10e vs 10h; Fig. 10f vs 10i).

### 3 **4. Discussion**

#### 4 *4.1 Representation error*

5 Our results show the impact of the size of the representation error on the accuracy of estimated  
6 model parameters. The parametric uncertainty (i.e. the size of the credible intervals in Figs 7-9)  
7 increases at an approximately linear rate as the representation error increases from  $p = 0.01$  to  $p$   
8  $= 0.3$ . This is consistent with Hill et al. (2012) who estimated the parameters and uncertainties of  
9 a simple terrestrial carbon model under varying levels of measurement error.

10 For the reanalysis data, we treat the representation error as a parameter for the MCMC  
11 algorithm to estimate along with the eight model parameters. This is possible because we  
12 assume that the measured value of  $O_3$  is proportional to the simulated value from a forward run  
13 of the FRSGC/UCI model, although such an assumption may not be possible in other situations.  
14 An alternative approach to estimate the representation error would be to carry out an intensive  
15 measurement campaign to determine whether the average  $O_3$  from different measuring stations  
16 within a grid-square is representative of the true average. Satellite products of the terrestrial  
17 biosphere are checked for accuracy using this type of approach (De Kauwe et al., 2011).  
18 Although measurement campaigns at these large spatial and temporal scales would be  
19 challenging and costly, they may not be need to continue for long periods of time since we might  
20 expect representation error to decrease as the temporal scale increases (Schutgens et al., 2016).

#### 21 *4.2 Spatial coverage*



1 We find that as the volume of measurements increase, the estimates of the model parameters are  
2 closer to the truth and the width of the credible intervals decrease. This is particularly clear for  
3 the dry deposition and isoprene emission scaling parameters when using both O<sub>3</sub> and CO  
4 concentrations (Figs 8 and 9). While this highlights the value of good spatial coverage, we note  
5 that the benefits are greatly reduced if the representation error is relatively high. For the  
6 boundary layer mixing parameter, we find little decrease in the credible intervals using synthetic  
7 CO data with the highest representation error ( $p = 0.3$ ), where the spatial coverage is less than  
8 20% (Fig. 9, row 2). In contrast, at the  $p = 0.1$  level, a large decrease in uncertainty is seen  
9 between the 2.5% and 20% coverage. Similar effects are seen, to a lesser extent, for the dry  
10 deposition and isoprene scaling parameters as the spatial coverage increases.

11 Our results using synthetic data show that while the size of the uncertainty intervals vary  
12 substantially depending on the spatial coverage or representation error, the posterior means are  
13 for the most part very close to the true values. Deviation from these typically occurs when the  
14 measurements contain less information either due to low spatial coverage or high representation  
15 error. However, the uncertainty intervals include the true values of the parameters for all the  
16 experimental scenarios considered here, unlike in Hill et al. (2012). This gives strong confidence  
17 in the reliability of the MCMC method used to estimate the parameters.

#### 18 *4.3 Applying multiple constraints*

19 The importance of multiple constraints was most apparent for scenarios involving the  
20 reanalysis data. For the dry deposition scaling parameter, which explains much of the variance  
21 in surface O<sub>3</sub> (Fig. 4), we found that using O<sub>3</sub> data alone results in mean estimates that are in the  
22 upper half of the range of possible values (Fig 10a). However, including CO data brought the  
23 mean estimates into the central part of the range where we would expect the true value to lie



1 (Fig. 10g). This is remarkable given that dry deposition is not an important process for  
2 controlling CO, and highlights the coupling between processes that permits constraints on one  
3 process from one variable to influence those on another. However, it is consistent with previous  
4 studies exploring the uncertainty in estimates of key parameters in an aerosol-chemistry-climate  
5 model (Johnson et al., 2018). For the isoprene emission and boundary layer mixing scaling  
6 parameters, there was little difference in the mean estimates or the size of the uncertainty  
7 intervals when using O<sub>3</sub> and CO together rather than a single constraint. This reveals that the  
8 importance of using multiple constraints is dependent on the process and on the variable  
9 constrained. A judicious choice of these could allow a particular process to be targeted.  
10 Overall, our estimates of the dry deposition and isoprene emission scaling parameters are close  
11 to a priori values from the FRSGC/UCI CTM. In contrast, our estimates of the boundary layer  
12 mixing scaling parameter are substantially larger than those from the model, suggesting that this  
13 process is not represented well in the model.

#### 14 *4.4 Towards constraint with real surface measurements*

15 Our results have demonstrated the feasibility of using measurement data to constrain model  
16 parameters under the right conditions. We have chosen to use synthetic data as they have allowed  
17 us to vary the spatial coverage and to investigate the effects of representation error which is  
18 poorly characterised when using real measurements data. Quantifying this type of error for real  
19 measurements is difficult because measurement sites are relatively sparse and are often  
20 representative of a limited area rather than the larger area typical of a model grid-square.  
21 However, this study has allowed us to estimate the representation error associated with the  
22 reanalysis data, and in the absence of more information these values could be used as a guide  
23 when applying surface measurements as a constraint.



1           The reanalysis data provide a more critical test, as they are independent of the  
2   FRSGC/UCI CTM used here. Although we do not know the true values of the scaling  
3   parameters, we expect them to lie close to those used in the control run given the relatively good  
4   agreement for O<sub>3</sub> and CO concentrations. For the dry deposition parameter, we expect scaling  
5   values to be close to 1, but using surface O<sub>3</sub> reanalysis data alone we found posterior mean  
6   scaling parameters approaching 1.4, with credible intervals that did not include 1 (Fig. 10a).  
7   This likely reflects overestimation of surface O<sub>3</sub> in continental regions in the CTM and may  
8   reflect uncertainties and biases in other processes not considered here, most notably in the  
9   chemical formation and destruction of O<sub>3</sub> and in model transport processes. In the absence of  
10   consideration of the uncertainty in these processes in this feasibility study, the dry deposition  
11   parameter is used as a proxy process to reduce O<sub>3</sub> concentrations. This is an example of  
12   equifinality, where different sets of parameters can result in model predictions that give equally  
13   good agreement with observations (Beven et al., 2001). Applying simultaneous constraints to  
14   CO goes some way to addressing this, but does not remove the problem. Before applying real  
15   surface measurements to constrain the CTM, we propose a more comprehensive assessment of  
16   model uncertainties with a wider range of parameters so that the constraints can more directly  
17   inform process understanding and model development.

## 18   **Conclusion**

19   We have demonstrated the use of surface O<sub>3</sub> and CO concentrations to constrain a global  
20   atmospheric chemical transport model and generate accurate and robust estimates of model  
21   parameters. This would normally be prohibitive for such a model given that thousands of model  
22   runs are required. Our approach is to replace the CTM with a surrogate model using Gaussian  
23   process emulation and then estimate the parameters using the emulator in place of the CTM. In



1 this feasibility study we have shown that surface  $O_3$  has a large sensitivity to dry deposition, and  
2 that surface CO is most sensitive to isoprene emissions and boundary layer mixing processes, as  
3 expected. We find that estimates of the scaling parameters for these processes are dependent on  
4 the spatial coverage and representation error of the surface  $O_3$  and CO data. Our parameter  
5 estimates become less uncertain as coverage increases and as the representation error decreases,  
6 whilst remaining unbiased. Furthermore, we show that using two separate data constraints, in  
7 this case surface  $O_3$  and CO, instead of a single one can result in mean parameter estimates that  
8 are much closer to their likely true values. However, this is dependent on the processes  
9 considered and constraints applied, and while it effective for dry deposition here, we find  
10 relatively little improvement in the estimates or uncertainties for isoprene emission or boundary  
11 layer mixing processes that are also considered here.

12 The approach we adopt here provides a means of constraining atmospheric models with  
13 observations and identifying sources of model error at a process level. Our results suggest that  
14 dry deposition and isoprene emissions are represented relatively well in the FRSGC/UCI CTM  
15 but that boundary layer mixing processes may be somewhat underestimated. However, we have  
16 explored the effect of only eight parameters in this study and consideration of a more complete  
17 set of processes, including those governing photochemistry and dynamics, is needed to generate  
18 more realistic constraints for key pollutants such as  $O_3$ . We aim to expand this study to  
19 investigate a more extensive range of parameters and processes and to constrain with a wider  
20 range of observation data. The emulator-based approach for estimating parameters that we have  
21 successfully demonstrated here can be applied to any model where evaluating the model the  
22 required number of times is too computationally demanding.

23



## 1 Code and data availability

2 The R code used for building and validating the emulators and estimating the posterior  
3 distribution of the model parameters using the Markov Chain Monte Carlo algorithm is available  
4 from the Zenodo data repository via the link: <https://zenodo.org/record/4537614>. The  
5 FRSGC/UCI model output used for training the emulators is available from the CEDA data  
6 repository via the link: <https://catalogue.ceda.ac.uk/uuid/d5afa10e50b44229b079c7c5a036e660>.

## 7 Appendix A

8 The formula for the covariance function  $V(x, x')$  from §2.2 is given by:

$$9 \quad V(x, x') = \sigma^2 [\mathcal{C}(x, x'; \psi) - t(x)^T \mathbf{A}^{-1} t(x) \\ 10 \quad \quad \quad + (h(x)^T + t(x)^T \mathbf{A}^{-1} H)(H^T \mathbf{A}^{-1} H)^{-1} (h(x')^T + t(x')^T \mathbf{A}^{-1} H)^T]$$

11 where,

$$12 \quad \sigma^2 = \frac{\mathbf{y}^T (\mathbf{A}^{-1} - \mathbf{A}^{-1} H (H^T \mathbf{A}^{-1} H)^{-1} H^T \mathbf{A}^{-1}) \mathbf{y}}{n - q - 1}$$

13 To compute the variance or uncertainty of a prediction  $x$  we use the formula for  $V(x, x')$  with  
14  $x' = x$ , which results in  $\mathcal{C}(x, x; \psi) = 1$ . Since we need to evaluate a large number of emulators  
15 for each MCMC iteration step (because we have a separate emulator for every dimension of the  
16 model output), it is more computationally efficient to compute the parts of the above formula  
17 prior to using the emulator. Hence, the above formula can be replaced with:

$$18 \quad V_i(x_{new}, x_{new}) = \sigma_{ALL}^2 [i, 1] \left[ (1 - t_i(x_{new}))^T V_{i,1} t_i(x_{new}) \right. \\ 19 \quad \quad \quad \left. + (h(x_{new})^T + t(x_{new})^T V_{i,2}) V_{i,3} (h(x_{new})^T + t(x_{new})^T V_{i,2})^T \right]$$

20 where:

- 21 •  $i$  ( $1 \leq i \leq r$ ) denoted the  $i$ th point in the  $r$ -dimensional simulator output.
- 22 •  $\sigma_{ALL}^2$  is a  $r \times 1$  vector that stores the values of  $\sigma^2$  for all  $r$  outputs.





- 1 •  $V_{i,1}$  is the  $n \times n$  matrix  $\mathbf{A}^{-1}$  corresponding to the  $i$ th point in the simulator's output. It is  
2 stored as the  $i$ th block of the  $nr \times n$  matrix  $V_1$  defined by:

3 
$$V_1 = \begin{pmatrix} V_{1,1} \\ V_{2,1} \\ \vdots \\ V_{r,1} \end{pmatrix}$$

- 4 •  $V_{i,2}$  is the  $n \times q$  matrix  $\mathbf{A}^{-1}H$  corresponding to the  $i$ th point in the simulator's output. It  
5 is stored as the  $i$ th block of the  $nr \times q$  matrix  $V_2$  defined by:

6 
$$V_2 = \begin{pmatrix} V_{1,2} \\ V_{2,2} \\ \vdots \\ V_{r,2} \end{pmatrix}$$

- 7 •  $V_{i,3}$  is the  $q \times q$  matrix  $(H^T \mathbf{A}^{-1} H)^{-1}$  corresponding to the  $i$ th point in the simulator's  
8 output. It is stored as the  $i$ th block of the  $qr \times q$  matrix  $V_3$  defined by:

9 
$$V_3 = \begin{pmatrix} V_{1,3} \\ V_{2,3} \\ \vdots \\ V_{r,3} \end{pmatrix}$$

#### 10 **Author contributions**

11 ER and OW designed the study. ER carried out the statistical analyses, and OW ran the  
12 FRSGC/UCI model and provided the outputs that were used to train and validate the emulators.  
13 ER wrote the paper with input from OW.

#### 14 **Acknowledgements**

15 This work was supported by the Natural Environment Research Council [grant number  
16 NE/N003411/1]. We thank Karl Hennermann at the ECMWF for making reanalysis data for O<sub>3</sub>



1 and CO from CAMS available. We also thank Lindsay Lee from Sheffield Hallam University for  
2 her feedback and comments on an early version of this paper.

### 3 **References**

- 4 BARET, F., WEISS, M., ALLARD, D., GARRIGUES, S., LEROY, M., JEANJEAN, H., FERNANDES, R.,  
5 MYNENI, R., PRIVETTE, J. & MORISETTE, J. 2005. VALERI: a network of sites and a  
6 methodology for the validation of medium spatial resolution land satellite products.  
7 *Remote Sensing of Environment*, 76, 36-39.
- 8 BERG, B. A. 2005. Introduction to Markov chain Monte Carlo simulations and their statistical  
9 analysis. *Markov Chain Monte Carlo Lect Notes Ser Inst Math Sci Natl Univ Singap*, 7, 1-  
10 52.
- 11 BEVEN, K. and FREER, J., 2001. Equifinality, data assimilation, and uncertainty estimation in  
12 mechanistic modelling of complex environmental systems using the GLUE  
13 methodology. *Journal of hydrology*, 249(1-4), pp.11-29.
- 14 BOCQUET, M., ELBERN, H., ESKES, H., HIRTL, M., ŽABKAR, R., CARMICHAEL, G., FLEMMING, J.,  
15 INNESS, A., PAGOWSKI, M. & PÉREZ CAMAÑO, J. 2015. Data assimilation in atmospheric  
16 chemistry models: current status and future prospects for coupled chemistry  
17 meteorology models. *Atmospheric Chemistry and Physics*, 15, 5325-5358.
- 18 BOERSMA, K., VINKEN, G. & ESKES, H. 2016. Representativeness errors in comparing chemistry  
19 transport and chemistry climate models with satellite UV–Vis tropospheric column  
20 retrievals. *Geoscientific model development*, 9, 875.
- 21 CHANG, K.-L., PETROPAVLOVSKIKH, I., COOPER, O. R., SCHULTZ, M. G. & WANG, T. 2017.  
22 Regional trend analysis of surface ozone observations from monitoring networks in  
23 eastern North America, Europe and East Asia. *Elem Sci Anth*, 5.
- 24 CHANG, W., HARAN, M., APPLGATE, P. & POLLARD, D. 2016. Calibrating an ice sheet model  
25 using high-dimensional binary spatial data. *Journal of the American Statistical*  
26 *Association*, 111, 57-72.
- 27 COHEN, A. J., BRAUER, M., BURNETT, R., ANDERSON, H. R., FROSTAD, J., ESTEP, K., ... & FEIGIN,  
28 V. 2017. Estimates and 25-year trends of the global burden of disease attributable to  
29 ambient air pollution: an analysis of data from the Global Burden of Diseases Study  
30 2015. *The Lancet*, 389(10082), 1907-1918.
- 31 CURRIN, C., MITCHELL, T., MORRIS, M. & YLVIKAKER, D. 1991. Bayesian prediction of  
32 deterministic functions, with applications to the design and analysis of computer  
33 experiments. *Journal of the American Statistical Association*, 86, 953-963.
- 34 DE KAUWE, M. G., DISNEY, M., QUAIFFE, T., LEWIS, P. & WILLIAMS, M. 2011. An assessment of  
35 the MODIS collection 5 leaf area index product for a region of mixed coniferous forest.  
36 *Remote Sensing of Environment*, 115, 767-780.
- 37 EMILI, E., BARRET, B., MASSART, S., LE FLOCHMOEN, E., PIACENTINI, A., EL AMRAOUI, L.,  
38 PANNEKOUCKE, O. & CARIOLLE, D. 2014. Combined assimilation of IASI and MLS  
39 observations to constrain tropospheric and stratospheric ozone in a global chemical  
40 transport model. *Atmospheric Chemistry and Physics*, 14, 177-198.



- 1 FIORE, A. M., DENTENER, F., WILD, O., CUVELIER, C., SCHULTZ, M., HESS, P., TEXTOR, C.,  
2 SCHULTZ, M., DOHERTY, R. & HOROWITZ, L. 2009. Multimodel estimates of  
3 intercontinental source-receptor relationships for ozone pollution. *Journal of*  
4 *Geophysical Research: Atmospheres*, 114.
- 5 FLEMMING, J., BENEDETTI, A., INNESS, A., ENGELEN, R. J., JONES, L., HUIJNEN, V., REMY, S.,  
6 PARRINGTON, M., SUTTIE, M. & BOZZO, A. 2017. The CAMS interim Reanalysis of Carbon  
7 Monoxide, Ozone and Aerosol for 2003-2015. *Atmospheric Chemistry and Physics*, 17,  
8 1945.
- 9 GELMAN, A., CARLIN, J. B., STERN, H. S., DUNSON, D. B., VEHTARI, A. & RUBIN, D. B. 2013.  
10 *Bayesian data analysis*, CRC press.
- 11 GOLDSMITH, J. R. & LANDAW, S. A. 1968. Carbon monoxide and human health. *Science*, 162,  
12 1352-1359.
- 13 GRANIER, C., BESSAGNET, B., BOND, T., D'ANGIOLA, A., VAN DER GON, H. D., FROST, G. J., HEIL,  
14 A., KAISER, J. W., KINNE, S. & KLIMONT, Z. 2011. Evolution of anthropogenic and  
15 biomass burning emissions of air pollutants at global and regional scales during the  
16 1980–2010 period. *Climatic Change*, 109, 163.
- 17 HILL, T. C., RYAN, E. & WILLIAMS, M. 2012. The use of CO<sub>2</sub> flux time series for parameter and  
18 carbon stock estimation in carbon cycle research. *Global Change Biology*, 18, 179-193.
- 19 JOHNSON, J. S., REGAYRE, L. A., YOSHIOKA, M., PRINGLE, K. J., LEE, L. A., SEXTON, D. M.,  
20 ROSTRON, J. W., BOOTH, B. B. & CARSLAW, K. S. 2018. The importance of  
21 comprehensive parameter sampling and multiple observations for robust constraint of  
22 aerosol radiative forcing. *Atmospheric Chemistry and Physics*, 18, 13031-13053.
- 23 JOHNSON, J.S., CUI, Z., LEE, L.A., GOSLING, J.P., BLYTH, A.M. and CARSLAW, K.S., 2015.  
24 Evaluating uncertainty in convective cloud microphysics using statistical  
25 emulation. *Journal of Advances in Modeling Earth Systems*, 7(1), pp.162-187.
- 26 KAMPA, M. & CASTANAS, E. 2008. Human health effects of air pollution. *Environmental*  
27 *pollution*, 151, 362-367.
- 28 KENNEDY, M. C. & O'HAGAN, A. 2000. Predicting the output from a complex computer code  
29 when fast approximations are available. *Biometrika*, 87, 1-13.
- 30 KENNEDY, M. C. & O'HAGAN, A. 2001. Bayesian calibration of computer models. *Journal of the*  
31 *Royal Statistical Society: Series B (Statistical Methodology)*, 63, 425-464.
- 32 KHATTATOV, B. V., LAMARQUE, J. F., LYJAK, L. V., MENARD, R., LEVELT, P., TIE, X., BRASSEUR, G.  
33 P. & GILLE, J. C. 2000. Assimilation of satellite observations of long-lived chemical  
34 species in global chemistry transport models. *Journal of Geophysical Research:*  
35 *Atmospheres*, 105, 29135-29144.
- 36 LANDRIGAN, P. J., FULLER, R., ACOSTA, N. J., ADEYI, O., ARNOLD, R., BALDÉ, A. B., BERTOLLINI,  
37 R., BOSE-O'REILLY, S., BOUFFORD, J. I. & BREYSSE, P. N. 2017. The Lancet Commission on  
38 pollution and health. *The Lancet*.
- 39 LEE, L., PRINGLE, K., REDDINGTON, C., MANN, G., STIER, P., SPRACKLEN, D., PIERCE, J. &  
40 CARSLAW, K. 2013. The magnitude and causes of uncertainty in global model  
41 simulations of cloud condensation nuclei. *Atmos. Chem. Phys*, 13, 8879-8914.
- 42 LEE, L. A., REDDINGTON, C. L. & CARSLAW, K. S. 2016. On the relationship between aerosol  
43 model uncertainty and radiative forcing uncertainty. *Proceedings of the National*  
44 *Academy of Sciences*, 113, 5820-5827.



- 1 LOEPPKY, J. L., SACKS, J. & WELCH, W. J. 2012. Choosing the sample size of a computer  
2 experiment: A practical guide. *Technometrics*.
- 3 MALLEY, C. S., HENZE, D. K., KUYLENSTIERNA, J. C., VALLACK, H. W., DAVILA, Y., ANENBERG, S.  
4 C., TURNER, M. C. & ASHMORE, M. R. 2017. Updated global estimates of respiratory  
5 mortality in adults  $\geq 30$  years of age attributable to long-term ozone exposure.  
6 *Environmental Health Perspectives*, 2017, vol. 125, num. 8, p. 087021.
- 7 MARREL, A., IOOSS, B., JULLIEN, M., LAURENT, B. & VOLKOVA, E. 2011. Global sensitivity  
8 analysis for models with spatially dependent outputs. *Environmetrics*, 22, 383-397.
- 9 MENUT, L., BESSAGNET, B., KHVOROSTYANOV, D., BEEKMANN, M., BLOND, N., COLETTE, A.,  
10 COLL, I., CURCI, G., FORET, G. & HODZIC, A. 2014. CHIMERE 2013: a model for regional  
11 atmospheric composition modelling. *Geoscientific model development*, 6, 981-1028.
- 12 MIYAZAKI, K., ESKES, H., SUDO, K., TAKIGAWA, M., VAN WEELE, M. & BOERSMA, K. 2012.  
13 Simultaneous assimilation of satellite NO<sub>2</sub>, O<sub>3</sub>, CO, and HNO<sub>3</sub> data for the analysis of  
14 tropospheric chemical composition and emissions. *Atmos. Chem. Phys*, 12, 9545-9579.
- 15 MORRIS, M. D. & MITCHELL, T. J. 1995. Exploratory designs for computational experiments.  
16 *Journal of statistical planning and inference*, 43, 381-402.
- 17 NICELY, J. M., ANDERSON, D. C., CANTY, T. P., SALAWITCH, R. J., WOLFE, G. M., APEL, E. C.,  
18 ARNOLD, S. R., ATLAS, E. L., BLAKE, N. J. & BRESCH, J. F. 2016. An observationally  
19 constrained evaluation of the oxidative capacity in the tropical western Pacific  
20 troposphere. *Journal of Geophysical Research: Atmospheres*, 121, 7461-7488.
- 21 O'HAGAN, A. 2006. Bayesian analysis of computer code outputs: a tutorial. *Reliability  
22 Engineering & System Safety*, 91, 1290-1300.
- 23 OAKLEY, J. E. & O'HAGAN, A. 2004. Probabilistic sensitivity analysis of complex models: a  
24 Bayesian approach. *Journal of the Royal Statistical Society: Series B (Statistical  
25 Methodology)*, 66, 751-769.
- 26 PARRISH, D., LAMARQUE, J. F., NAIK, V., HOROWITZ, L., SHINDELL, D., STAEHELIN, J., DERWENT,  
27 R., COOPER, O., TANIMOTO, H. & VOLZ-THOMAS, A. 2014. Long-term changes in lower  
28 tropospheric baseline ozone concentrations: Comparing chemistry-climate models and  
29 observations at northern midlatitudes. *Journal of Geophysical Research: Atmospheres*,  
30 119, 5719-5736.
- 31 PLUMMER, M. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling.  
32 Proceedings of the 3rd international workshop on distributed statistical computing,  
33 2003. Technische Universit at Wien, 125.
- 34 RASMUSSEN, C. E. 2006. Gaussian processes for machine learning.
- 35 RICHARDSON, A. D., WILLIAMS, M., HOLLINGER, D. Y., MOORE, D. J., DAIL, D. B., DAVIDSON, E.  
36 A., SCOTT, N. A., EVANS, R. S., HUGHES, H. & LEE, J. T. 2010. Estimating parameters of a  
37 forest ecosystem C model with measurements of stocks and fluxes as joint constraints.  
38 *Oecologia*, 164, 25-40.
- 39 ROUSTANT, O., GINSBOURGER, D. & DEVILLE, Y. 2012. DiceKriging, DiceOptim: Two R packages  
40 for the analysis of computer experiments by kriging-based metamodeling and  
41 optimization.
- 42 RYAN, E., WILD, O., VOUGARAKIS, A. & LEE, L. 2018. Fast sensitivity analysis methods for  
43 computationally expensive models with multi-dimensional output. *Geoscientific model  
44 development*, 11, 3131-3146.



- 1 SALTELLI, A., TARANTOLA, S. & CHAN, K.-S. 1999. A quantitative model-independent method for  
2 global sensitivity analysis of model output. *Technometrics*, 41, 39-56.
- 3 SALTER, J. M., WILLIAMSON, D. B., SCINOCCA, J. & KHARIN, V. 2018. Uncertainty quantification  
4 for spatio-temporal computer models with calibration-optimal bases. *arXiv preprint*  
5 *arXiv:1801.08184*.
- 6 SCHULTZ, M. G. 2016. Cluster analysis of European surface ozone observations for evaluation of  
7 MACC reanalysis data. *Atmospheric Chemistry and Physics*, 16, 6863.
- 8 SCHULTZ, M. G., AKIMOTO, H., BOTTENHEIM, J., BUCHMANN, B., GALBALLY, I. E., GILGE, S.,  
9 HELMIG, D., KOIDE, H., LEWIS, A. C. & NOVELLI, P. C. 2015. The Global Atmosphere  
10 Watch reactive gases measurement network. *Elem Sci Anth*, 3.
- 11 SCHULTZ, M. G., SCHRÖDER, S., LYAPINA, O., COOPER, O., GALBALLY, I., PETROPAVLOVSKIKH, I.,  
12 VON SCHNEIDEMESSER, E., TANIMOTO, H., ELSHORBANY, Y. & NAJA, M. 2017.  
13 Tropospheric Ozone Assessment Report: Database and metrics data of global surface  
14 ozone observations. *Elem Sci Anth*, 5.
- 15 SCHUTGENS, N. A., GRYSPEERDT, E., WEIGUM, N., TSYRO, S., GOTO, D., SCHULZ, M. & STIER, P.  
16 2016. Will a perfect model agree with perfect observations? The impact of spatial  
17 sampling. *Atmospheric Chemistry and Physics*, 16, 6335-6353.
- 18 SHINDELL, D., FALUVEGI, G., SELTZER, K. and SHINDELL, C., 2018. Quantified, localized health  
19 benefits of accelerated carbon dioxide emissions reductions. *Nature climate*  
20 *change*, 8(4), pp.291-295.
- 21 SOFEN, E., BOWDALO, D., EVANS, M., APADULA, F., BONASONI, P., CUPEIRO, M., ELLUL, R.,  
22 GALBALLY, I., GIRGZDIENE, R. & LUPPO, S. 2016. Gridded global surface ozone metrics  
23 for atmospheric chemistry model evaluation. *Earth System Science Data*, 8, 41.
- 24 STEVENSON, D., DENTENER, F., SCHULTZ, M., ELLINGSEN, K., VAN NOIJE, T., WILD, O., ZENG, G.,  
25 AMANN, M., ATHERTON, C. & BELL, N. 2006. Multimodel ensemble simulations of  
26 present-day and near-future tropospheric ozone. *Journal of Geophysical Research:*  
27 *Atmospheres*, 111.
- 28 STOCKER, T. F., QIN, D., PLATTER, G. K., TIGNOR, M., ALLEN, S. K., BOSCHUNG, J., ... & MIDGLEY,  
29 P. M. 2013. Climate change 2013: The physical science basis. *Contribution of working*  
30 *group I to the fifth assessment report of the intergovernmental panel on climate*  
31 *change*, 1535.
- 32 VAN DINGENEN, R., DENTENER, F. J., RAES, F., KROL, M. C., EMBERSON, L. & COFALA, J. 2009.  
33 The global impact of ozone on agricultural crop yields under current and future air  
34 quality legislation. *Atmospheric Environment*, 43, 604-618.
- 35 VAN LOON, M., BUILTJES, P. J. & SEGERS, A. 2000. Data assimilation of ozone in the atmospheric  
36 transport chemistry model LOTOS. *Environmental Modelling & Software*, 15, 603-609.
- 37 VAN ZELM, R., HUIJBREGTS, M. A., DEN HOLLANDER, H. A., VAN JAARVELD, H. A., SAUTER, F. J.,  
38 STRUIJS, J., VAN WIJNEN, H. J. & VAN DE MEENT, D. 2008. European characterization  
39 factors for human health damage of PM10 and ozone in life cycle impact assessment.  
40 *Atmospheric Environment*, 42, 441-453.
- 41 WILD, O. 2007. Modelling the global tropospheric ozone budget: exploring the variability in  
42 current models. *Atmospheric Chemistry and Physics Discussions*, 7 (1), pp.1995- 2035.  
43 fahal-00302576f



- 1 WILD, O. and PRATHER, M.J., 2006. Global tropospheric ozone modeling: Quantifying errors due  
2 to grid resolution. *Journal of Geophysical Research: Atmospheres*, 111(D11).
- 3 WILD, O., Voulgarakis, A., O'Connor, F., Lamarque, J. F., Ryan, E., & Lee, L. 2020. Global  
4 sensitivity analysis of chemistry-climate model budgets of tropospheric ozone and OH:  
5 Exploring model diversity. *Atmospheric Chemistry and Physics*, 20, 4047-4058.
- 6 WILKINSON, R. D. 2010. Bayesian calibration of expensive multivariate computer experiments.  
7 John Wiley & Sons.
- 8 WILLIAMS, M., RICHARDSON, A., REICHSTEIN, M., STOY, P., PEYLIN, P., VERBEECK, H.,  
9 CARVALHAIS, N., JUNG, M., HOLLINGER, D. & KATTGE, J. 2009. Improving land surface  
10 models with FLUXNET data. *Biogeosciences*, 6, 1341-1359.
- 11 YOUNG, P. J., NAIK, V., FIORE, A. M., GAUDEL, A., GUO, J., LIN, M., NEU, J., PARRISH, D., RIEDER,  
12 H. & SCHNELL, J. 2018. Tropospheric Ozone Assessment Report: Assessment of global-  
13 scale model performance for global and regional ozone distributions, variability, and  
14 trends. *Elem Sci Anth*, 6.

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

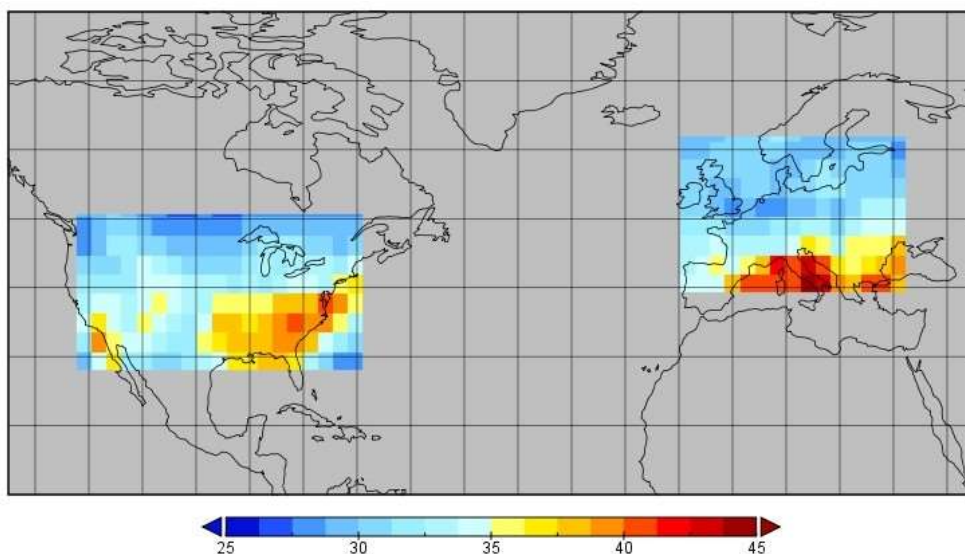
32



1  
2  
3  
4  
5

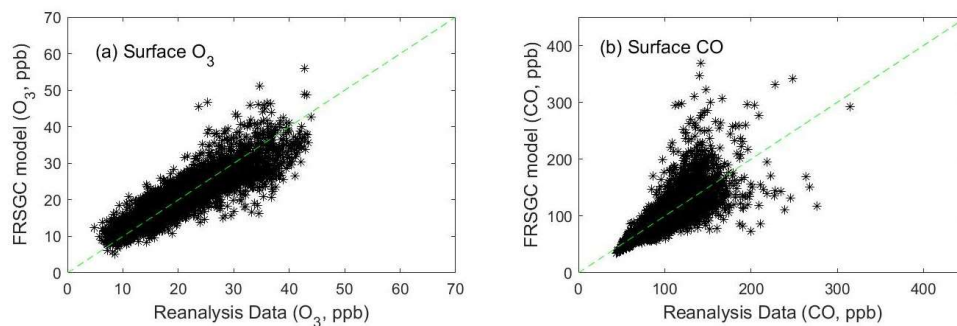
### Figures and Tables for

#### ‘Calibrating a global atmospheric chemistry transport model using Gaussian process emulation and ground-level concentrations of ozone and carbon monoxide’



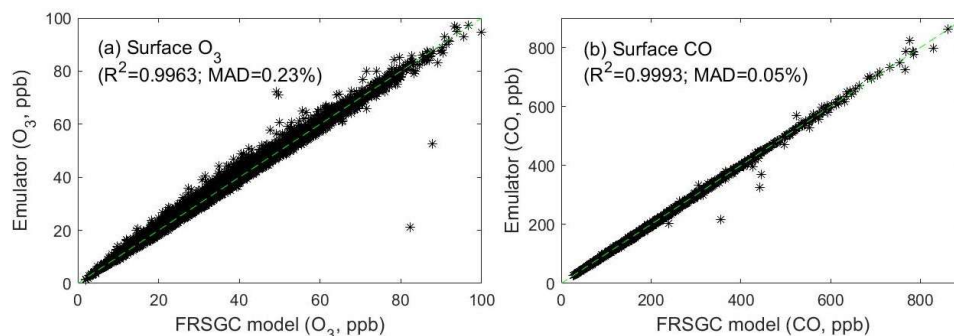
6  
7  
8  
9  
10

**Figure 1.** Annual mean surface ozone mixing ratio (in ppb) from the FRSGC/UCI CTM showing the regions considered here and the 272 grid cells used for model calibration.

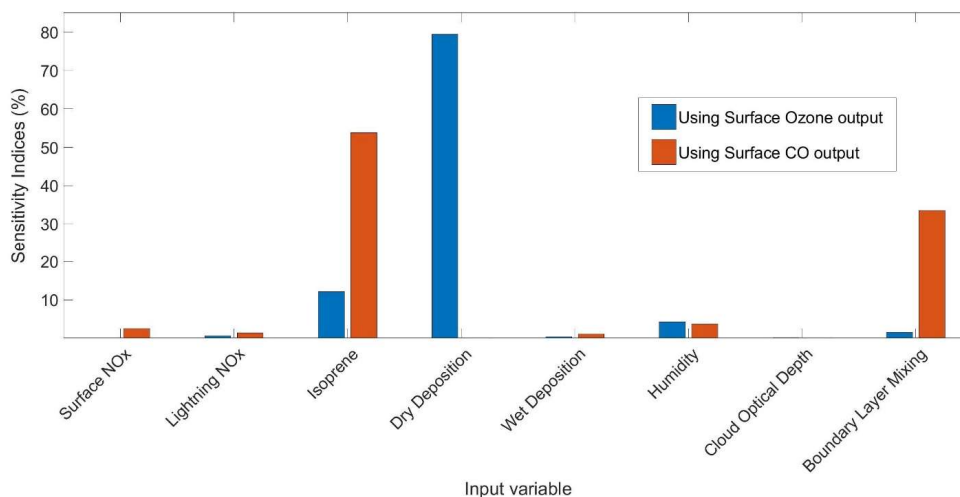


11  
12  
13  
14

**Figure 2.** Monthly mean surface  $O_3$  (panel a) and surface CO (panel b) over Europe and North America simulated with the FRSGC/UCI CTM compared with ECMWF reanalysis data.

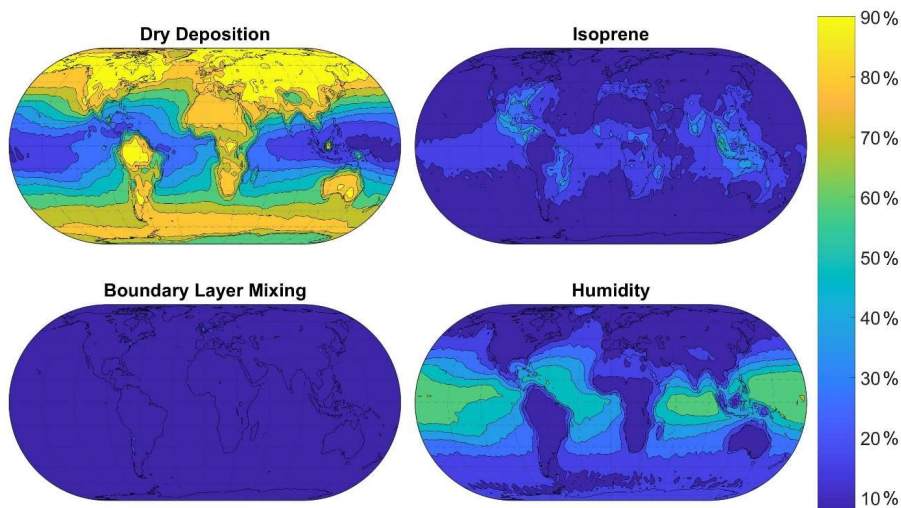


1  
 2 **Figure 3.** Simulated surface O<sub>3</sub> (panel a) and surface CO (panel b) from the FRSGC/UCI CTM versus  
 3 those predicted from the Gaussian process emulators. The simulated and emulated concentrations were  
 4 generated using 20 sets of model parameters that were not used for training the emulators.  
 5

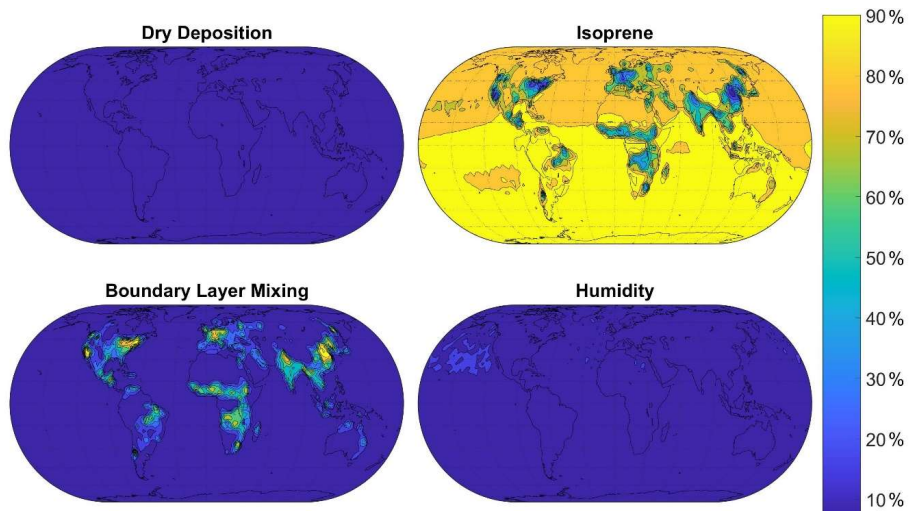


6  
 7 **Figure 4.** Sensitivity indices representing the percentage of the variance in surface O<sub>3</sub> and CO over the  
 8 USA and Europe in the FRSGC/UCI model output due to changes in each parameter.

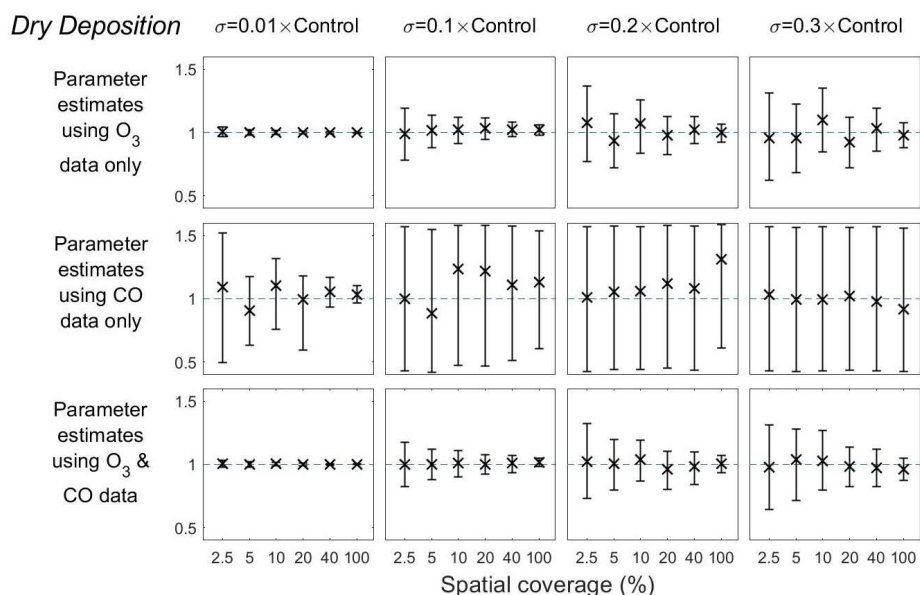




1  
2 **Figure 5.** Sensitivity indices representing the percentage of the variance in surface  $O_3$  in the FRSGC/UCI  
3 model output due to changes in each input parameter. The four parameters displayed here have the  
4 highest sensitivity indices and the largest effect on simulated surface  $O_3$ . Maps of sensitivity indices  
5 corresponding to the other four parameters are shown in Figure S2 of the supplementary material.  
6



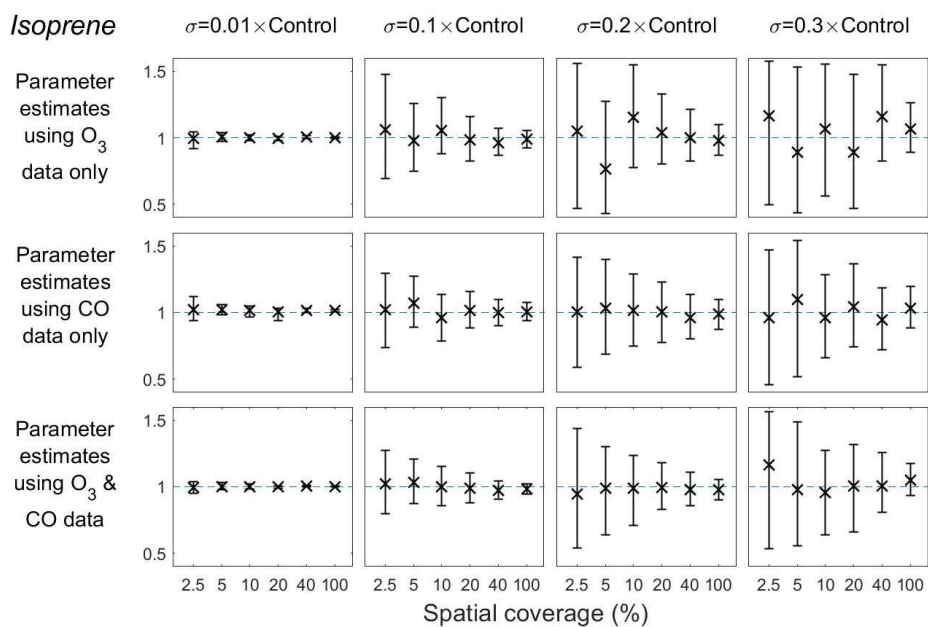
7  
8 **Figure 6.** Sensitivity indices representing the percentage of the variance in surface CO in the  
9 FRSGC/UCI model output due to changes in each input parameter. Maps of sensitivity indices for the  
10 other four parameters are shown in Figure S3 of the supplementary material.



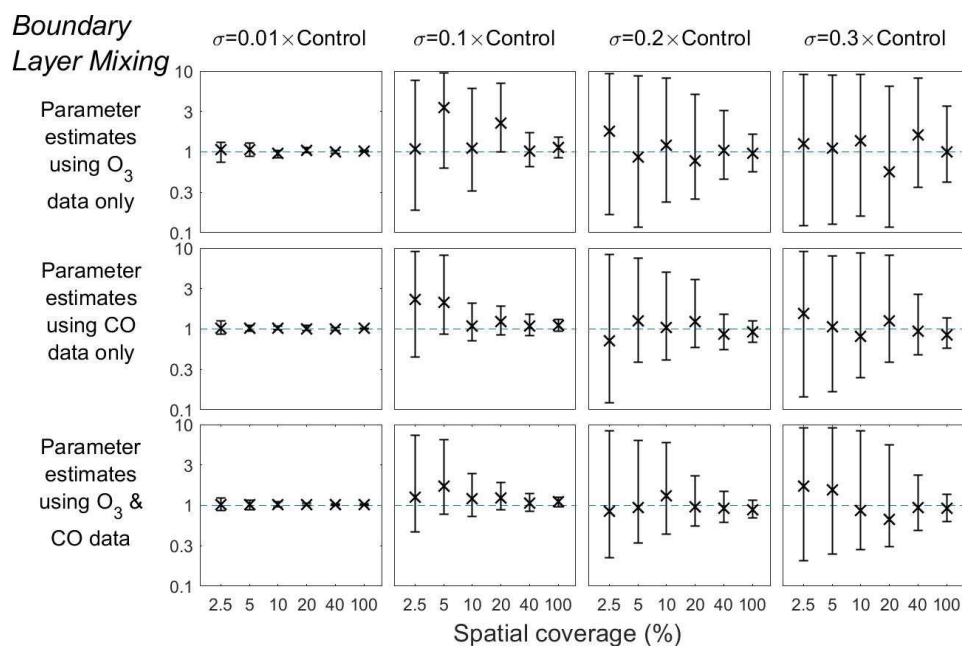
1

2 **Figure 7.** Means and 95% credible intervals of 3000 samples of the **Dry Deposition** scaling parameter  
 3 from posterior distributions using the MCMC algorithm based on synthetic datasets from scenarios 1-72  
 4 (table 1). *Control* refers to the FRSGC control run surface concentration for each output point.

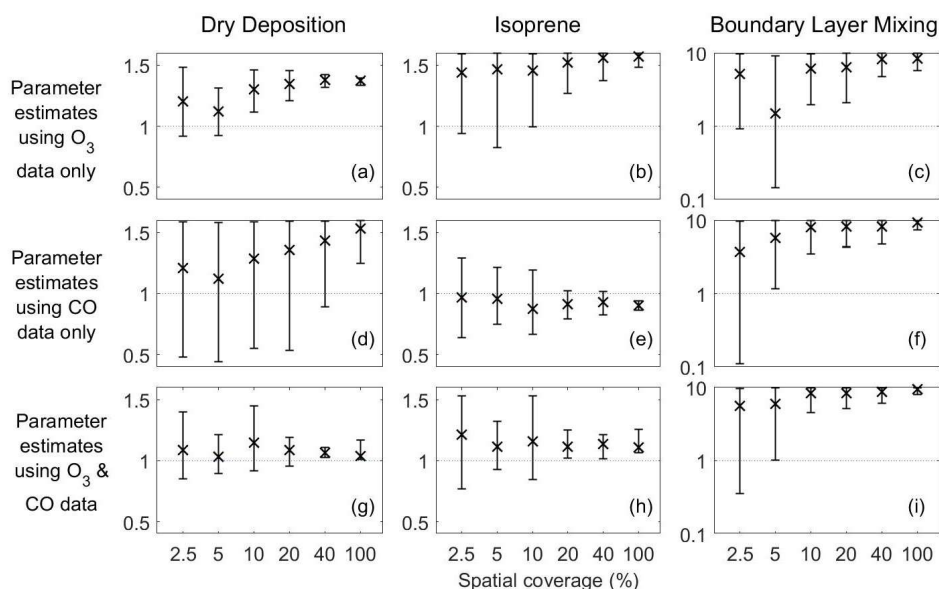
5



1  
 2 **Figure 8.** Means and 95% credible intervals of 3000 samples of the **Isoprene** emission scaling parameter  
 3 from posterior distributions using the MCMC algorithm based on synthetic datasets from scenarios 1-72  
 4 (table 1). *Control* refers to the FRSGC control run surface concentration for each output point.



1  
 2 **Figure 9.** Means and 95% credible intervals of 3000 samples of the **Boundary Layer Mixing** scaling  
 3 parameter from posterior distributions using the MCMC algorithm based on synthetic datasets from  
 4 scenarios 1-72 (table 1). *Control* refers to the FRSGC control run surface concentration at each output  
 5 point. The scaling parameter values are given here on the  $\log_{10}$  scale.



1  
 2 **Figure 10.** Means and 95% credible intervals of 3000 samples of the Dry Deposition, Isoprene and  
 3 Boundary Layer Mixing scaling parameters from posterior distributions using the MCMC algorithm  
 4 based on reanalysis datasets from scenarios 73-90 (table 1). The first and second rows show these  
 5 parameters estimated using one stream of data ( $O_3$  for the first row and CO for the second row), while the  
 6 third row shows estimates using two data streams ( $O_3$  and CO).  
 7

8  
 9  
 10  
 11  
 12  
 13  
 14  
 15  
 16  
 17  
 18  
 19



1 **Table 1.** Model processes and associated scaling parameter ranges used in this study.

| Number | Model process                                       | Control run value | Scaling parameter values |
|--------|-----------------------------------------------------|-------------------|--------------------------|
| 1      | Global surface NO <sub>x</sub> emissions (TgN/year) | 40                | 0.75 – 1.25              |
| 2      | Global lightning NO emissions (TgN/year)            | 5                 | 0.40 – 1.60              |
| 3      | Global isoprene emissions (TgC/year)                | 500               | 0.40 – 1.60              |
| 4      | Dry deposition rates                                | model value       | 0.40 – 1.60              |
| 5      | Wet deposition rates                                | model value       | 0.40 – 1.60              |
| 6      | Humidity                                            | model value       | 0.80 – 1.20              |
| 7      | Cloud optical depth                                 | model value       | 0.33 – 3.00              |
| 8      | Boundary Layer mixing                               | model value       | 0.10 – 10.0              |

2

3 **Table 2.** Summary of the 90 different MCMC scenarios carried out for this study. The scenarios involved  
 4 varying: (i) the type of data (synthetic or reanalysis); (ii) the representation error used for the synthetic  
 5 data ( $p$ ) where  $m_i(x_{control})$  is the control run output of the CTM and  $\sigma_i$  is the amount of statistical noise  
 6 added; (iii) the percentage coverage of grid-squares in the USA and Europe. For the synthetic data the 24  
 7 scenarios correspond to a full factorial combination of four levels of representation error and six levels of  
 8 spatial coverage, while for the reanalysis data the six scenarios correspond to the six levels of spatial  
 9 coverage.

| Scenarios | Dataset                               | Representation error, $p$<br>( $\sigma_i = p \times m_i(x_{control})$ ) | Spatial coverage              |
|-----------|---------------------------------------|-------------------------------------------------------------------------|-------------------------------|
| 1-24      | Synthetic O <sub>3</sub>              | 0.01, 0.1, 0.2, 0.3                                                     | 2.5%, 5%, 10%, 20%, 40%, 100% |
| 25-48     | Synthetic CO                          | 0.01, 0.1, 0.2, 0.3                                                     | 2.5%, 5%, 10%, 20%, 40%, 100% |
| 49-72     | Synthetic O <sub>3</sub> & CO         | 0.01, 0.1, 0.2, 0.3                                                     | 2.5%, 5%, 10%, 20%, 40%, 100% |
| 73-78     | Reanalysis data (O <sub>3</sub> )     | Parameter to be estimated                                               | 2.5%, 5%, 10%, 20%, 40%, 100% |
| 79-84     | Reanalysis data (CO)                  | Parameter to be estimated                                               | 2.5%, 5%, 10%, 20%, 40%, 100% |
| 85-90     | Reanalysis data (O <sub>3</sub> & CO) | Parameter to be estimated                                               | 2.5%, 5%, 10%, 20%, 40%, 100% |

10