

Title

**Multi-stage Ensemble-learning-based Model Fusion for
Surface Ozone Simulations: A Focus on CMIP6 Models**

Authors

Zhe Sun^{1,2*}, Alexander T. Archibald^{1,3*}

Affiliations

¹ Centre for Atmospheric Science, Yusuf Hamied Department of Chemistry,
University of Cambridge, Cambridge CB2 1EW, UK

² Department of Earth Sciences, University of Cambridge, Cambridge CB2
3EQ, UK

³ National Centre for Atmospheric Science, Cambridge CB2 1EW, UK

* Corresponding authors:

Zhe Sun (zs347@cam.ac.uk) and Alexander T. Archibald (ata27@cam.ac.uk)

ABSTRACT

Accurately simulating global surface ozone has long been one of the principal components of chemistry-climate modelling, but divergences in simulation outcomes have been reported as a result of the mechanistic complexity of tropospheric ozone budget. Settling the cross-model discrepancies to achieve higher accuracy thus is a task of priority. Building on the Coupled Model Intercomparison Project Phase 6 (CMIP6), we have transplanted a conventional ensemble learning approach, and also constructed an innovative 2-stage enhanced space-time Bayesian neural network to fuse an ensemble of 57 simulations together with a prescribed ozone dataset, both of which have realised outstanding performances ($R^2 > 0.95$, $RMSE < 2.12$ ppbv). The conventional ensemble learning approach is computationally cheaper and results in higher overall performance, but at the expense of oceanic ozone being overestimated and the learning process being uninterpretable. The Bayesian approach performs better in spatial generalisation and enables perceivable interpretability, but requires heavier computational burdens. Both of these multi-stage learning-based approaches provide frameworks for improving the fidelity of composition-climate model outputs for use in future impact studies.

Keywords

CMIP6; CCM; surface ozone; model ensemble; space-time Bayesian neural network; data fusion

1 INTRODUCTION

Tropospheric ozone (O_3) is a trace-gas, near-term climate forcer with global mean lifetime ~23 days, and also a major air pollutant being of detrimental defects on human and ecosystem health.¹⁻³ Besides warming the atmosphere as a greenhouse gas, ground-level O_3 also reduces crop yields.⁴⁻⁶ Laboratory experiments have confirmed O_3 exposure to cause oxidative stress, inflammatory responses and immunologic diseases.⁷ Epidemiological studies report that short-term exposures to high-level ozone are significantly associated with the exacerbation of asthma⁸ and have increased hospitalisations among children,⁹ while long-term ozone exposure is linked to respiratory diseases like chronic obstructive pulmonary disease, cardiovascular diseases, and even premature deaths.¹⁰⁻¹⁴ Global Burden of Disease (GBD) reported over 0.36 million premature deaths globally in 2019 from exposure to ambient O_3 ;¹⁵ and high O_3 exposure could exacerbate the $PM_{2.5}$ -mortality risk associations.¹⁶ These results underscore the pressing need for research linking population exposure assessment to surface O_3 and its impacts on human health.

Satellite-based observations cannot provide accurate measurements for O_3 at the surface since surface O_3 will be obscured by the climbing O_3 abundance in high-layer atmosphere thus cannot be measured directly from remote-sensing; while the ground-level station-based observation sites are still rather limited in spatial coverage.^{17, 18} The demands for full-coverage surface O_3 concentrations have promoted the application of model simulations, which have been being improved as our understanding of the mechanisms behind tropospheric O_3 has improved.¹⁹⁻
²¹ But model simulations are not perfect, due to imperfections of O_3 chemistry mechanisms built in the models, biases and errors in the underlying emissions, and uncertainties caused by the discretisation and numerical treatment of a non-linear complex system. Archibald et al. have

shown that for future evolution projections of the tropospheric column O_3 , model differences are a leading order term of uncertainty over decadal scales.²¹ There are various types of models used to simulate surface O_3 . Chemical transport models (CTM) perform satisfactorily especially in regional-level simulations;²²⁻²⁵ and are considered to be free of biases in meteorology due to the use of prescribed meteorology. But these models lack important feedbacks from atmospheric composition on to the model meteorology and climate, hence atmospheric composition-climate models (CCM) have been developed; and when coupled with land, sea, and sea-ice modules into earth system models (ESM), it is feasible to simulate multi-decadal or even centennial scale changes in atmosphere.²⁶⁻²⁹

To evaluate and compare the coupled models, a number of research institutes have contributed to the Coupled Model Inter-comparison Project Phase 6 (CMIP6) with a range of experiments conducted by a series of state-of-the-art coupled CCMs and ESMs. The same inputs are used, including emission inventories and land properties.³⁰⁻³³ CMIP6 has endorsed a total of 23 MIPs to answer a wide range of scientific questions in atmospheric chemistry and climate, among which the Aerosols and Chemistry Model Intercomparison Project (AerChemMIP) involves a collection of simulations targeted at reactive gases and aerosols including tropospheric O_3 .³⁴ Large discrepancies have been detected across models; beyond figuring out the mechanistic causes for these differences,^{31, 35} an urgent challenge is how to calibrate and make the maximum use of the simulation ensemble.

Applying frontier machine learning algorithms to assimilate the outputs from multi-source modelling activities like MIPs and observation databases, known as data “assimilation” or data “fusion”, is an important part of environmental research in the big data era. Studies which enhance

the prediction accuracy of ambient air pollution concentrations by ensemble learning have emerged in recent years.³⁶⁻³⁹ However, these studies only used no more than one model simulation integrated with predictor variables contributing to the budget of O₃, without involving fusing multiple simulation ensembles like CMIP6. In addition, the conventional machine- or deep-learning approaches aim purely at brute-force fitting into high accuracy while sacrificing the interpretability of the training processes, so have long been criticised as “black-box” and contradict the nature of mechanism-driven sciences like atmospheric modelling.⁴⁰⁻⁴² Under these circumstances, reaching a performance-interpretability balance for multi-source data fusion following credible observations will be of high value in atmospheric research.

Our current study is an innovative exploration on this issue, emphasising on developing innovative ensemble-learning frameworks to assimilate the multiple CMIP6 model simulation ensembles and TOAR observations to obtain one single surface O₃ dataset capturing the spatiotemporal variabilities as accurate as possible. Fusing a collection of simulation ensembles rather than just using the output from one simulation can give more prominence to the mechanism-driven models so as to avoid brute-force overfitting resulting from external predictor variables, especially when any given model simulation could be largely biased. The primary innovation of this study is in transplanting the conventional ensemble-learning data-assimilation methodology onto multi-source data fusion, and optimising an enhanced 2-stage space-time Bayesian neural network to assimilate the CMIP6 simulation ensemble. The advantages of the conventional approach include a much lower computation burden and higher accuracy in observation-covered regions, while the merits of the innovative Bayesian approach lie in its better spatial generalisability and intuitive perception of spatiotemporal model weighting. In either case,

the multi-model fused surface O₃ concentration can fill in the observational gaps and enable further relevant researches in the long term. As an example we show here using Fourier-series function to fit the temporal surface O₃ variability provides a feasible way to effectively summarise periodic air pollutant concentrations. Detailed evaluations and comparisons on CMIP6 model ensemble, and deeper discussions on model revision insights from deep learning-based calibration processes are beyond the scope of this study.

2 METHODOLOGY AND DATA SOURCES

2.1 CMIP6 simulation ensemble

We collect 14 coupled earth system models having finished the “historical” simulations (1850-2014) of tropospheric O₃ as listed in [Table S1](#), of which 8 models use interactive chemistry schemes. A prescribed O₃ concentration dataset is used for all 4 non-interactive chemistry models (AWI-ESM,⁴³ BCC-CSM2,⁴⁴⁻⁴⁶ IPSL-CM6A,^{47, 48} and MPI-M-ESM1.2⁴⁹⁻⁵²) and 2 CNRM models are not considered for fusion due to the simplified treatment of O₃ chemistry.⁵³⁻⁵⁷ A total of 8 models, including BCC-ESM1,^{58, 59} MPI-ESM1.2-HAM,⁶⁰ MRI-ESM2.0,⁶¹⁻⁶³ NASA-GISS-E2.1,⁶⁴⁻⁶⁶ NCAR-CESM2-WACCM6,^{67, 68} NCC-NorESM,⁶⁹ NOAA-GFDL-ESM4,^{70, 71} and UKESM1-0-LL,^{19, 28, 72-75} consisting of 57 individual simulation experiments (i.e. realisations in terms of CCM simulation labelled as $r_{in}i_{np}n_{fn}$) and 1 prescribed input dataset (from Inputs4MIPs)⁷⁶ are recruited for data fusion. The multiple ensemble members under one model allow for capturing the uncertainties in the chaotic coupled chemistry-climate system; and because of the free-running nature of the simulations, each of the 57 individual simulations is treated separately with no cross-ensemble averaging clustering into each model involved. All simulation outputs are averaged to

monthly time frequency for assimilation with observations. Detailed information of the participant research institutes, design of atmosphere module settings, and experiment labelling rules are illustrated in the Supporting Information.

2.2 Observations

The tropospheric ozone assessment report (TOAR) programme has archived high-quality ground-level O₃ measurements over the period 1990-2014,¹⁷ which are used as “standard” for physical and statistical model evaluation; our study period is thus selected as 1990-2014. To support analyses at the planar spatial resolution of the CCMs involved in this study, TOAR sites are aggregated into 2°×2° latitude-longitude grid as plotted in **Figure S1**, including 585 spatial grids with a total of 5,322 different observational sites; and averaged to monthly temporal interval for the robustness of model-observation evaluation. Such spatiotemporal aggregations can also strengthen the stability of grid-level observation-simulation evaluation, and to some extent abate the statistical compromises by excluding the observation missing records for some certain sites in the early years of the dataset (ca. 1990s). Only spatial grids in which there is at least one observation site are used. Throughout the study, the gridded TOAR observations are used as supervised learning labels.

2.3 Additional auxiliary predictors

Higher prediction accuracy can be achieved when integrating additional features into statistical models.³⁶⁻³⁸ Comprehensively considering the O₃ budget mechanisms, experiences from previous relevant studies, and statistical correlations with surface O₃, we screen out 13 variables as assistant predictors as: CMIP6 simulated concentrations of surface PM_{2.5}, NO₂, higher layers of O₃ (vertical O₃ column), and ambient air temperature obtained from the World Climate Research

Programme (WCRP) Earth System Grid Federation (ESGF) CMIP6 database (<https://esgf-node.llnl.gov/search/cmip6>); emissions of biogenic VOCs, NO_x, CO, black carbon (BC) and organic carbon (OC) together with urbanised land proportions, collected from input datasets for Model Intercomparison Projects (<https://esgf-node.llnl.gov/search/input4mips>); surface elevation downloaded from the Global Multi-resolution Terrain Elevation Data (GMTED);⁷⁷ and gridded urban and rural populations linearly interpolated with corrections towards the actual annual world total populations into year-precision from United Nation's World Population Prospects (UN WPP) Adjusted Population Density and Gridded Population of the World (GPW) operated by NASA Socioeconomic Data and Applications Centre (SEDAC).⁷⁸

2.4 Multi-model Fusion Frameworks

We use “physical model” to refer to the CMIP6 mechanism-driven atmospheric models, and “statistical model” for the data-oriented machine- or deep-learning frameworks to avoid confusion in terminology. No transformations are made for either the observations or model simulations as they follow the Gaussian distribution well with slight temporal imbalance. Following literatures,³⁶⁻³⁸ an adjusted ensemble learning-based multi-model fusion framework is constructed as presented in the upper panel of **Figure 1**. In this approach, raw simulations (i.e. 57 CMIP6 historical simulations and 1 prescribed O₃ dataset, noted as “57+1 ensemble” hereafter) together with the normalised additional predictor variables are first re-gridded onto the 2°×2° TOAR observation grids, following procedures graphically presented in **Figure S2**. Then, all the model simulation ensembles, external predictors, and 6 space-time indices (i.e. 3 Euclidean spherical coordinates in analytic geometry, and 3 helix-shape trigonometricised month sequence t as $[\cos(2\pi tT^{-1}), \sin(2\pi tT^{-1}), t]$ where T is prescribed as 1 year)⁷⁹ are mixed together as inputs for random forest, gradient

boost decision tree, and convolutional neural network regression models separately; and outputs from the 3 algorithms are finally blended by L2-regularisation-based weighting (ridge regression). This approach is entitled as “aggressive” approach because this methodology respects the observations (i.e. labels for supervision) more than the physical models, hence during the process of training, the concentrations in each grid are treated individually so as to compromise the spatiotemporal continuous structure of the original physical model simulations, leading to inexplicability. The aggressive approach involves at least two stages of ensemble: the first CMIP6 multi-model ensemble and second multi-algorithm ensemble, where the random forest regressor essentially is another layer of ensemble learning. The random forest regressor is a large collection of separate decision trees with individual of which generating a single prediction and the final prediction given by averaging all trees, thus the random forest is perceived as an ensemble learning method.⁸⁰

Contrarily, in order to maintain the interpretability of the deep learning processes, we also adopt an enhanced 2-stage space-time Bayesian neural network (BNN) framework as illustrated in the lower panel of **Figure 1**. Space-time indices and additional predictors are put into a 10-layer 1024-node BNN to generate spatiotemporal variant re-scaling factors (k), bias correctors (b) and the randomised noises (σ), under the supervision of TOAR observations to pre-calibrate the raw re-gridded CMIP6 simulations. Then, spatiotemporal variant model weights (α) are estimated by 5-layer 256-node BNN merely from the 6 space-time indices, to finally reach the weighted average ensemble surface O₃ concentration predictions. This approach is named as the “conservative” approach as throughout the process of prediction enhancement, all parameters are clamped by space-time indices with presumed distributions, thus this framework respects the raw

simulations more and might be highly biased on extreme observations. All involved parameters can be thoroughly separated from the framework and presented intuitively by mapping, so that the whole process of assimilation is traceable and interpretable. We construct the two-stage BNN instead of single-stage because the divergences still exist among the calibrated CMIP6 models in the first-stage and hence further mixing is required. Directly using the second-stage BNN will lose the chance to observe the calibration features for individual physical models; and different degrees of initial biases will cast higher weights onto the smaller biased models, possibly leading to undesirable feature monopolisation.

Statistical principles of naïve space-time BNN (i.e. single-stage space-time BNN) are illustrated in details by a recent report.⁷⁹ Mathematically speaking, solutions of the spatiotemporal parameters (i.e. k , b , and α) are not unique, but it is reasonable to assume the observation covered and uncovered regions are of homogeneity in distribution of these parameters, which requires a Bayesian method to replace the single value of parameters with a distribution. The 6 space-time indices can assist in capturing the spatiotemporal autocorrelation of the surface O_3 . 10,000 times of Monte Carlo simulation ensembles are applied to approximate the distribution, so as to guarantee the robustness of BNN estimation, thence the conservative approach involves 3-stage ensemble: first in multi-model ensemble and the latter two in the 2-stage Bayesian parameter generation. For the final predictions based on the optimised distribution parameters trained through the BNN, 69.2% fall into 1 standard deviation (σ) range, 96.2% into 2σ and 99.9% into 3σ , conforming to the regularity of Gaussian distribution and thus justifying our Bayesian model presumption.

To evaluate the performance of 2 approaches, 10-fold cross-validation (CV) assessment is

applied, and 7:3 training-test split is used through the full dataset during 1990-2014. An additional temporal extrapolation test is conducted by manually setting the 1990-2009 TOAR observations with grid-corresponding physical model simulations as training set and 2010-2014 as test set. Three manual cross-validation tests are conducted by splitting the whole dataset into training-testing sets with regional integrity as i) Europe-training for North-America-testing; ii) North-America-training for Europe-testing; and iii) Europe-North-America-training for East-Asia-testing, so as to evaluate the spatial extrapolation capability of the 2 statistical models. Decomposition of model-observation errors follow a previous research.⁸¹ The neural network trainings are accomplished by Adam stochastic optimisation algorithm, setting the initial anchor values from observations and the learning rate as 10^{-4} after centric normalisation.

The complex machine learning frameworks are constructed instead of using simple statistical models owing to their limitations in handling the i) similarities across multiple physical models (i.e. collinearity in statistical term); ii) interaction effects between the input variables; iii) spatiotemporal auto-correlations and discrepancies in calibration parameters; and iv) propensity of overfitting when introducing high-order polynomial terms. Additionally, this cross-disciplinary study closely follows the trends of applying the cutting-edge data sciences onto environmental studies, hence only machine- and deep-learning approaches are transplanted, enhanced and discussed here.

2.5 Other relevant statistics

Fourier-series sinusoid functions theoretically can fit any periodical variables,⁸² so are used to capture the location-specific seasonal periodic variations of surface O_3 in this study to parametrically interpret the final assimilated surface O_3 concentrations by revealing the intra- and

inter-year variability quantitatively with perceivable mapping. Akaike Information Criteria (AIC) is used for statistical model selection, taking the realistic explicability altogether into consideration as listed in Table S2. Given TOAR observations and model outputs are monthly averaged, the final Fourier function is chosen as

$$f(t) = a_0 e^{a_1 t} + (b_0 + b_1 t) \sin\left(\frac{\pi}{6} t + \varphi_1\right) + c_0 \sin\left(\frac{2\pi}{6} t + \varphi_2\right),$$

where t represents the month-sequence; a_0 as starting-point surface O_3 concentration (January 1990); $12a_1$ as annual average change rates; $2b_0$ as the baseline and $24b_1$ as annual change of seasonal variation amplitude (i.e. peak-valley difference); and c_0 as the fine-tuning parameter which can modify the sinusoidal shape, but usually the absolute values are rather small, thus not considered for interpretation. An exponential term for the annual average surface O_3 is applied instead of linear term as the long-term simulations have reported exponential increasing trend of the tropospheric O_3 over centennial scales,³¹ regardless of the fact that the AIC values vote for the linear model.

3 RESULTS

3.1 Raw simulation evaluations

Raw CMIP6 surface O_3 simulations generally perform fairly well across all TOAR covered areas in terms of synchronicity (Figure 2), as the correlations between observations and the 57+1 ensemble averages are 0.74 ± 0.18 (inter-quartile range, IQR: [0.67, 0.87], Range: [-0.58, 0.96]). Overestimations are observed at 4.1 ± 2.0 (IQR: [5.1, 13.1], Range: [-22.2, 31.1]) ppbv across all TOAR covered spatial grids, hence the normalised mean biases (NMB) are high at 9.7 ± 6.3 (IQR: [4.2, 13.5], Range: [-28.1, 48.9]) %. Some regions like west Australia coastline even report

negative correlations (Pearson's $\rho = -0.58$).

The synchronicity and bias for realisation-enssembled model outputs are also evaluated in [Figure S3](#) and [Figure S4](#). NASA-GISS-E2.1 reports negative synchronicity in the USA-Canada border, while NCC-NorESM fails to reproduce the temporal variabilities in most of the studied sites. UKESM1-0-LL predicts closely to the measurements, but underestimates the surface O_3 around the USA-Canada border; while all the rest models present overestimations. Divergences are found between the individual models ([Figure S5](#)), and the high simulation discrepancies are mainly aggregated in the intertropical convergence zone (ITCZ) and eastern China, where the standard deviations exceeded 20% of the ensemble means. The barely satisfactory synchronicities and high overestimation biases indicate that the raw surface O_3 simulation might not be suitable for direct application in health impact studies, verifying the necessity of calibrations, at least statistically.

3.2 Performance of multi-model ensemble fusion

Both aggressive and conservative multi-model fusion perform well in prediction enhancement ([Figure 2](#)). The model-observation correlations are high at 0.98 ± 0.01 (IQR: [0.97, 0.99]) and 0.95 ± 0.08 (IQR: [0.95, 0.98]) for the aggressive and conservative approach, respectively; and NMBs of the aggressive model are 0.29 ± 3.06 (IQR: [-1.22, 1.54]) %, marginally smaller than the conservative model at 0.40 ± 3.57 (IQR: [-1.72, 1.93]) %. The general overestimation issues of the raw CMIP6 simulations have been handled well, but there are still some sporadic high NMBs detected in Asia, Africa, and South America, where the ground-based monitoring sites are rare and spatially scarce.

The full-range fitting R^2 ([Table 1](#)) of the aggressive and conservative approaches are 0.96 and

0.95, respectively, both indicating plausibility of the multi-model fusion with calibration; while the conservative predictions follow more loosely to the observations, especially in the low-concentration ranges (Figure S6), resulting in relatively higher root mean squared error (RMSE) at 2.12 ppbv compared with 1.81 ppbv for the aggressive approach. However, the conservative approach performs better in 1:1 model-observation calibration criteria according to the closer-to-one slope factor ($k_c^{-1} < k_a$, $0.97^{-1} < 1.05$) and closer-to-zero systematic bias ($|b_c| < |b_a|$, $|0.71| < |1.35|$). This is because directly involving additional features (i.e. the aggressive approach) can possibly introduce noise into the calibration, as their association with surface O₃ are not simply linear, especially in higher concentration ranges, so that the 1:1 model-calibration line is deviated.

Both approaches calibrate the physical models effectively, with the conventional aggressive approach performing slightly better than the innovatively established conservative model, which however, is already good. The spatiotemporal stability of the two approaches are also assessed in Table 1, concluding that the aggressive approach performs better in the later years of the dataset, while the conservative approach performs consistently well across the 25-year period. This is because the aggressive approach depends so largely on the observations that defects of observation coverage in early years will compromise the learning effects. However, the aggressive approach performs well across different continents ($R^2 > 0.90$), but the conservative approach performs slightly worse in the southern hemisphere ($R^2 > 0.83$), as a result of insufficient observations. This data sparsity results in the inter model-spread in the raw simulations being, to some extent, retained, as this could not be addressed by the BNN-based weighted linear combination; instead, additional features in the prediction-oriented aggressive approach brute-forcedly correct the large observation-simulation gaps. Both approaches perform well across seasons.

3.3 Extrapolation generalisability

Due to the limitations of lacking systematic observations in China, India, Africa and oceanic regions during 1990-2014, there are no means to verify the simulations in these areas directly; but this problem can be explored indirectly by checking the extrapolation potential on the observation-uncovered locations. Three regional cross-validation tests are graphically summarised in **Figure S7**, all of which reveal better generalisation capability of the conservative approach than aggressive. Neither underfitting nor overfitting issues are detected on the conservative approaches (i.e. CV and test scores are quite close); while underfitting is apparent for the aggressive approach in these regions, mainly reflected by failures in capturing extreme O₃ concentrations. The temporal extrapolation tests of two statistical models reveal high generalisability on the most recent 5-year test sets during 2010-2014 as $R^2 = 0.91$ (CV- $R^2 = 0.88$, test- $R^2 = 0.82$) for the aggressive approach and $R^2 = 0.92$ (CV- $R^2 = 0.89$, test- $R^2 = 0.85$) for the conservative approach. The temporal extrapolation performances are better than spatial generalisation, because the temporal periodic variations of surface O₃ are of a more stable pattern than regional divergences. In a nutshell, the conservative BNN approach wins over towards spatial and temporal generalisability, and we thus regard the conservative BNN results as “standard” for further interpretation.

3.4 Differences between ensemble approaches

Comparisons between the “standard” and aggressive approach outcomes are graphically summarised in **Figure S8**, revealing most of the global regions are of high congruity ($\rho = 0.85 \pm 0.17$, IQR: [0.81, 0.96]), while the divergences mostly occur on the ITCZ and Arabian-African areas ($\rho < 0.02$). Small relative biases have also justified the similarity between the aggressive and conservative approaches, as the NMBs (defined as aggressive – conservative) are 1.38 ± 4.61

(IQR: [-1.59, 3.77]) %. The positive differences mainly aggregate in Africa, Antarctica, Oceania and most of the oceanic basins, while the negative differences cluster in Asia, Europe and America.

The simplest fusion, the arithmetic average, of CMIP6 simulation ensemble would be used as a compromise were there no ground-based observations as used by precedent studies,³¹ which factually could lead to high biases if the real surface O₃ exposure assessment is the main research interest. This study aims to develop innovative approaches to fuse both model simulations and observations, and by comparing with the simplest fusion, advantages of new methods can be highlighted. The conservatively ensembled surface O₃ concentrations are of higher synchronicity ($\rho = 0.97 \pm 0.06$, IQR: [0.97, 0.99]) with the simple ensemble average than the aggressive approach ($\rho = 0.87 \pm 0.14$, IQR: [0.83, 0.96]), as the BNN is essentially an enhanced linear combination of multiple model simulations without substantial changes to the spatiotemporal auto-correlation. The ensemble average exceeds the aggressive fusion by 5.9 ± 9.7 (IQR: [-7.9, 14.3]) %, and the overestimations cluster regularly on land surface, especially the high-population-density regions; but surpass the conservative fusion by 9.6 ± 10.5 (IQR: [0.81, 20.2]) %, with the overestimations mainly detected in the wide-coverage northern-hemisphere without apparent land-ocean distinguishment. In conclusion, the simple ensemble average can lead to overestimations, especially in the northern hemispheric land surface; and the differences also reveal that the aggressive fusion model has modified the spatial auto-correlation of the raw CMIP6 simulation to a larger extent than the conservative approach.

3.5 Bayesian spatiotemporal weights

The differences between the two approaches can also be partially attributed to the different

weighting schemes of the raw individual simulations. The 57+1 ensembles occupy 93.9% weights in the aggressive approach while the additional assistant variables only contribute 6.1%. Generally, for the aggressive approach, 4 among the 58 simulations contribute dominantly by over 10%, as UKESM1-0-LL-r3i1p1f2 (18.6%), the prescribed O₃ (17.4%), NASA-GISS-E2.1-G-r1i1p3f1 (14.7%) and NCAR-CESM2-WACCM6-r1i1p1f1 (14.1%), while 36 ensemble members contribute less than 0.1%, as graphical presented in [Figure S9](#). On the contrary, the conservative approach results in relatively more even weights, where the prescribed O₃ (2.1%), UKESM1-0-LL (1.9%) and NASA-GISS-E2.1 (1.8%).

Besides the physical model weights, the space-time BNN also generates spatiotemporal variant weights, which can reflect the regions of skill for each individual physical model as presented in [Figure S10](#): UKESM1-0-LL and NCAR-CESM2-WACCM6 are weighted higher in northern hemisphere over land, while the prescribed O₃ dataset, NASA-GISS-E2.1, and NOAA-GFDL-ESM4 contribute more in southern hemisphere over land. The temporal variations of the spatial weights are generally small and of regular regional clustering trends, indicating that the physical models have captured the seasonal variability well.

BNN-based multi-model fusion treats the assistant variables independently with the CMIP6 model simulations, so that the weights of these additional features are not at the same level as the physical models like in the aggressive approach. Direct comparisons of the weights of the assistant variables between the two approaches reveal quite similar patterns of using these additional features for model calibration as shown in [Figure S11](#) which indicates that urban-rural populations, ambient air temperature and elevation are important factors. We suggest further work pay more attention to the role of model surface temperature, which is not fixed in these free-

running simulations. High contributions of the space-time indices also indicate that more additional features need to be included for further consideration.

3.6 Long-term surface ozone variations

Spatiotemporal variabilities of the BNN-fused surface O_3 are summarised parametrically using Fourier-series functions (Figure 3). The fitting quality R^2 has reached 0.81 ± 0.12 (IQR: [0.77, 0.87]), where the poor performances ($R^2 < 0.50$) concentrate in ITCZ and the coastlines. The global annual average increasing rate of the surface O_3 is estimated to be 0.23 (95% CI: [0.21, 0.25]) % yr^{-1} , and the highest increasing rates are detected in south Asia, South America, and continental Europe. Decreasing trends are also discovered in eastern China and eastern US. The average intra-year seasonal variation is 13.9 (IQR: [2.1, 49.5]) ppbv, and the highest amplitude differences cluster in eastern US, Africa, Europe, and eastern China. The annual changes of seasonal variations also demonstrate regional variabilities: widening in eastern China by maximum as 1.8 ppbv per year while narrowing in western countries by extreme to -0.8 ppbv per year. The intra-year peak and valley concentrations are generally ascending, as the peaks increase by 8.8 ± 1.1 (IQR: [-6.8, 16.1]) ppbv per year, and the valleys ascend by 0.6 ± 0.8 (IQR: [-7.0, 8.3]) ppbv per year.

4. DISCUSSION

4.1 Multi-model fusion improvement potentials

Decomposition of model-observation errors (Figure S12) can assist in evaluating the optimisation potentials for both the physical and statistical models.⁸³ The overall RMSE for the aggressive approach is 1.81 ppbv, among which the irreducible root-noise is 1.42 ± 0.47 ppbv,

occupying 66.1 ± 16.7 % of the total errors; while the averaged error of the conservative approach is 2.58 ppbv, where the root-noise is 1.87 ± 0.70 ppbv, accounting for 62.2 ± 25.4 %. The noises together with the biases by conservative approach are generally higher than the aggressive approach, while their proportions are close except for the African regions, as listed in [Table S3](#). Most of the unsolvable noises take over more than half of the errors, indicating that both fusion approaches have well approached the realistic observations.

The variances, also known as cross-model divergences, are comparable or even greater than biases for the aggressive approach, while for conservative approach the variances are several folds lower than biases, accounting for less than 10% except for South America (17%). This indicates the conservative fusion model is more robust. The model variances can be statistically perceived as discrepancies of model construction by random draws of the training subset, so that higher model variances represent severe dependences on training inputs, revealing higher sensitivity and lower generalisability.

The current crux of the conservative fusion model falls on the biases, suggesting higher optimisation potentials than the aggressive approach. The biases originate from the inherent systematic biases in physical models, and also the insufficient inclusion of assistant features to enhance the prediction statistically. Comparatively, due to the relatively higher statistical model variances, the aggressive approach shall no longer be the prevalent stream for multi-model fusion, as changes in observation coverage (i.e. labels for supervision in machine learning) will affect the stability of the statistical model substantially.

4.2 Differences in spatial extrapolation

The better spatial generalisation ability of the conservative space-time BNN multi-model

fusion is an advantage over the aggressive approach. Paradoxically, the aggressive approach actually performs well on capturing the extreme values. This shall be attributed to overfitting on the assistant features added directly into the fusion processes, so that the predictions are excessively reliant on these external variables. However, due to the complexity of the mechanisms controlling O₃, the statistical associations between physical models, auxiliary predictors, and the realistic concentrations recognised by the aggressive approach will be superfluous and of localised boundedness so that might be drastically different across regions. Excluding these features from aggressive multi-model fusion alleviates the poor performance in spatial extrapolation, as for each regional cross-validation test, R² rise to 0.81, 0.83, 0.74, and RMSE decline to 3.64, 3.97, 5.95 ppbv for North America, Europe and East Asia, respectively. To put it briefly, the external assistant features can increase the fitting quality in statistical training, but also serve as the limiting factors for model generalisation. This presents an issue towards understanding the processes of aggressive multi-model fusion, as conservative predictions manifested as underfitting by aggressive approach should be ascribed to the overfitting in the additional feature-assisted aggressive pathway. It suggests that conventional ensemble deep-learning approaches respecting the observations as supervision and linking the input variables only statistically rather than respecting the physical and chemistry mechanisms are of rather limited use, hence it is the second reason that the novel conservative multi-model fusion approach by space-time BNN is preferred.

4.3 Cross-approach divergences

Most discrepancies between the two fusion approaches and the simple ensemble average are located in tropics (Figure S8), which is primarily attributable to the lack of observations as training data, and the variations in raw simulations (Figure S5) resulting from the difficulty in

capturing O₃ in this region as a result of complexity in the precursor emissions like biogenic VOCs, soil NO_x, lightning NO_x, etc.³¹ We highlight in particular the need for long-term continuous ground-based measurements of O₃ in the tropics as a research priority.

The differences between the simple ensemble average and the aggressive fusion approach (Figure S8) indicate that the aggressive approach only addressed the systematic overestimations on the land surface; the additional variables lead to a land-ocean contrast (e.g. the population, ambient air temperature, O₃ precursor emissions), which are used as key nodes in the tree-structure regressions, so that the calibrations are only effective over the land rather than the whole global surface. The conservative approach respects the raw simulations more by calibrating uniformly for both lands and oceans, so that the average-conservative differences are more spatially uniform (Figure S8).

4.4 Systematic overestimation

Direct use of the raw CMIP6 surface O₃ simulation ensemble mean, as commonly used in the literature^{2, 31, 35}, causes positive biases around 5-10%, equal to 3.6 ± 4.4 ppbv, with some regions like India high-biased by +40% (+22.7 ppbv), consistent with recent multi-model ensemble studies in this region.⁸⁴ Such large biases have important implications on the use of raw ensemble mean data for work related to public health and pollution control policy studies in these regions, reiterating the necessity of observation-supervised calibration. The systematic overestimations across CMIP6 simulations speculate the major cause as the inadequate vertical stratification in atmospheric module. Essentially speaking, the lowest layers of CMIP6 model simulations are used to approximate the surface O₃, but the layer actually refers to a vertical average. Tropospheric O₃ concentration rises with the altitude,³¹ thus resulting in overestimation. UKESM1-0-LL stratifies

85 vertical layers,¹⁹ which is the most among 8 interactive chemistry CMIP6 models (Table S1), and lowest overestimations are found, with even underestimations observed in quite a few regions (Figure S4). Further experiments by adjusting the vertical stratifications to observe the changes in surface O₃ simulation performances are suggested to rigorously check this speculation.

4.5 Rationality of enhanced space-time BNN

Our enhanced space-time BNN is optimised from the traditional naïve space-time BNN, without additional feature involvement.⁷⁹ The enhancement in part comes from overcoming the inconsistency between the overall and location-specific observation-simulation linear relationships: each simulation cell at different time requires a unique set of k - b parameters for calibration as $y_{l,t}^{obs} = k_{l,t} \cdot y_{l,t}^{mod} + b_{l,t} + \varepsilon_{l,t}$, where the subscripts l and t represent location and time indices, so that using a fixed slope k and intercept b to calibrate all simulation cells is of limited use. However, the calculated sets of parameters are spatially limited to the observations, thus a naïve space-time BNN framework is required for spatial extrapolation onto the full global space.

The BNN generates the space-time variant calibration slopes and intercepts for each CMIP6 model in the pilot attempts, with which the assistant features are significantly correlated Figure S13, indicating these additional factors can contribute to the calibration parameters. For the purpose of increasing the prediction accuracy, the enhanced 2-stage Bayesian neural network regression-based multi-model fusion framework is constructed by firstly incorporating the assistant features into the multi-layer perception structure to generate the calibrated individual simulations, and secondly fusing them up by another naïve space-time BNN without involving any external features.

4.6 Sensitivity Analysis

Considering the cross-realisation variations (0.5 ± 0.1 ppbv) are much lower than the cross-model deviations (4.6 ± 1.7 ppbv, [Figure S5](#)), we conduct an additional sensitivity analysis by firstly averaging the multi-realisation within each model and then putting the 8 realisation-averaged model simulations together with the prescribed O_3 (hereafter noted as 8+1 models) into the aggressive and conservative model as input layer. The results of these new fused data are very similar to the previous calculations, with $R^2 = 0.94$, RMSE = 2.24 ppbv for the aggressive approach, and $R^2 = 0.93$, RMSE = 2.67 ppbv for the conservative approach. It shows that different numbers of realisations for each model will not significantly affect the fusion performance, indicating that the disparity in the number of realisations for a given model (e.g. 21 realisations for NASA-GISS-E2.1 while only a single realisation for NOAA-GFDL-ESM4) is not a significant issue when it comes to model data fusion. It also suggests averaging the multi-realisation ensemble before multi-model fusion takes place will still result in accurate results. This is particularly important if the model-data fusion approach is computationally expensive, as is the case for the conservative approach we have used.

One-dropout sensitivity analysis shows removing one model (with all its realisations) can achieve accuracy R^2 as 0.91 – 0.93 with RMSE 2.49 – 2.82 ppbv with the aggressive approach, and R^2 ranging 0.89 – 0.93 with RMSE 2.97 – 3.46 ppbv by conservative approach; results which are insignificantly lower than using all 8+1 CMIP6 models. However, the multi-model fusion performances are substantially reduced when only 2 models are kept (keeping only one single model will be inappropriate for the basic idea of *multi-model fusion*), as $R^2 = 0.83 - 0.87$, RMSE = 3.68 – 5.14 ppbv with the aggressive approach, and $R^2 = 0.71 - 0.78$, RMSE = 4.79 – 8.02 ppbv

with the conservative approach. The aggressive-conservative performance gap converges when fusing >9 realisations, or >4 realisation-averaged models. It exposes the critical limitation of the conservative approach and that the innovative enhanced space-time BNN will not perform satisfactorily when only a few models are used for fusion, because different models have used different chemistry mechanisms, or simplifications, or have other physical differences,⁸⁵ so that limited numbers of models cannot capture the full variations of the realistic surface O₃ by BNN-based linear-combination. It also further justifies the necessity of the CMIP6 multi-model study from the perspective of raising the signal-noise ratio and enabling more credible surface O₃ datasets (the more models used in the fusion process the better the performance). We keep the aggressively and conservatively-fused outcomes separately as 2 ultimate achievements of this study, instead of mixing them up into a single dataset, because of our aim of maintaining the interpretability of the BNN-fusion processes instead of purely focusing on brute-force fitting.

4.7 Merits and Limitations

Five major merits of our study are highlighted. First, we establish an enhanced 2-stage space-time Bayesian neural network regression-based deep-learning framework to fuse multi-ensemble surface O₃ simulation, which is verified to be of high accuracy and accessible interpretability in spatiotemporal weighting. Second, we verify the better spatial extrapolation generalisability of our newly developed approach than the conventional method; and owing to the commendable spatial and temporal extrapolation potentials, our ensemble learning frameworks can be applied to a wide temporal range of surface O₃ studies. Third, as far as we are aware, our study is the first study to fuse CMIP6 model simulations for surface O₃ over the 25 historical year period of 1990-2014 by machine learning techniques, and such long-term global studies are still rather rare. Fourth, the

fused and calibrated surface O₃ concentration dataset can be used by further researchers for further cross-disciplinary studies. Last but not the least, we innovatively apply Fourier-series functions for the purpose of parametrising and visualising the complex temporal periodical variations of surface O₃. However, our studies are still of several limitations. First, the model evaluation-calibration resolution is coarse as 2°×2°, and some heavily polluted regions like China, India and Africa are still lacking of observations. Second, the additional assistant features to enhance the statistical model prediction are still limited, and more variables shall be considered in further studies. Third, more detailed and deeper discussions concerning the parametric model calibration by 2-stage space-time BNN regression could have been replenished and excavated, but not included in this current paper as it is beyond the scope of this study. We aim to address some of these issues in our further research.

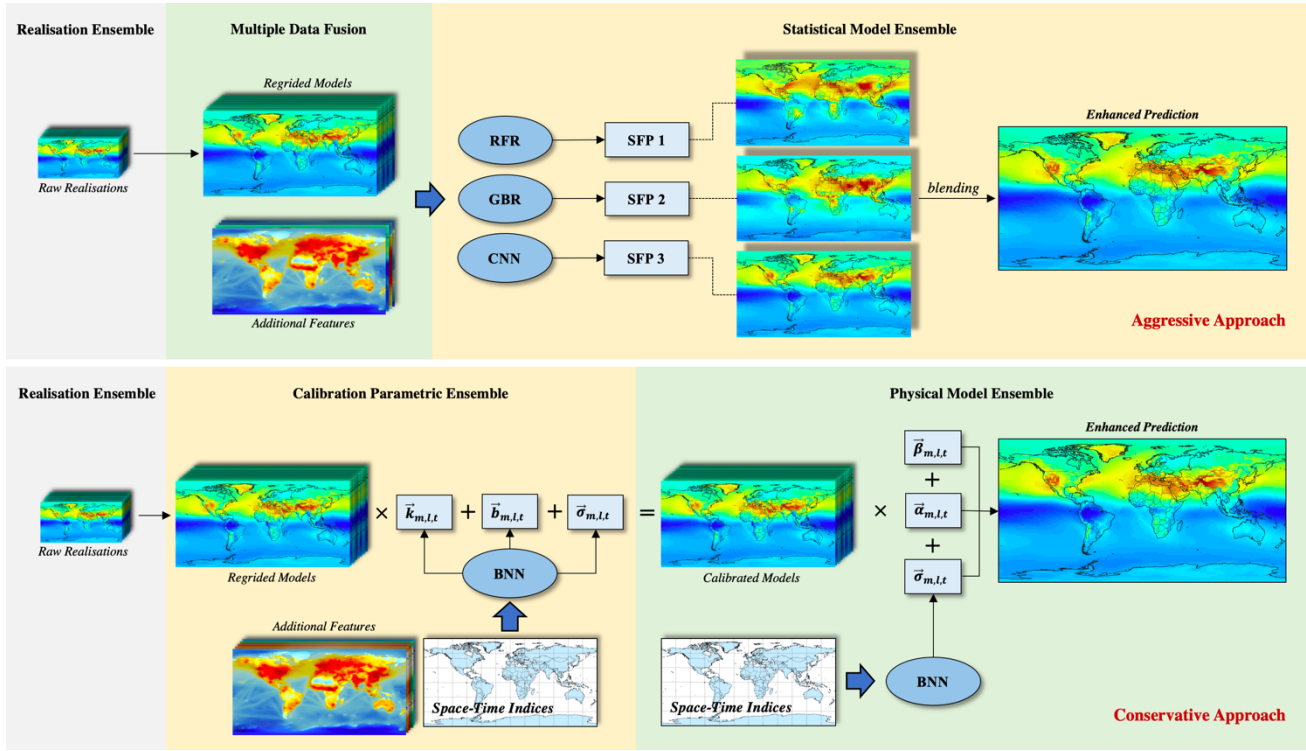


Figure 1 Schematic diagram of machine learning-based multi-model fusion by aggressive and conservative approaches. The stacking of source data layers refers to the collections of datasets with the same level in training models; the ellipses indicate elemental machine learning methodologies; and the rectangles represent the raw outputs from machine learning treatments. A total of 57 physical model simulations and 1 prescribed O₃ concentration dataset (Inputs4MIPs) are considered.

Abbreviations and denotations: RFR, *random forest regression*; GBR, *gradient boosting decision tree regression*; CNN, *convolutional neural network regression*; SFP, *semi-final product*; BNN, *Bayesian neural network regression*; \bar{k} , *re-scaling factor*; \bar{b} , *systematic bias corrector*; $\bar{\alpha}$, *individual model weight*; $\bar{\beta}$, *bias corrector*; \bar{m} , *physical model identifier*; \bar{l} , *location index*; \bar{t} , *temporal index*; $\bar{\sigma}$, *random noise*.

14

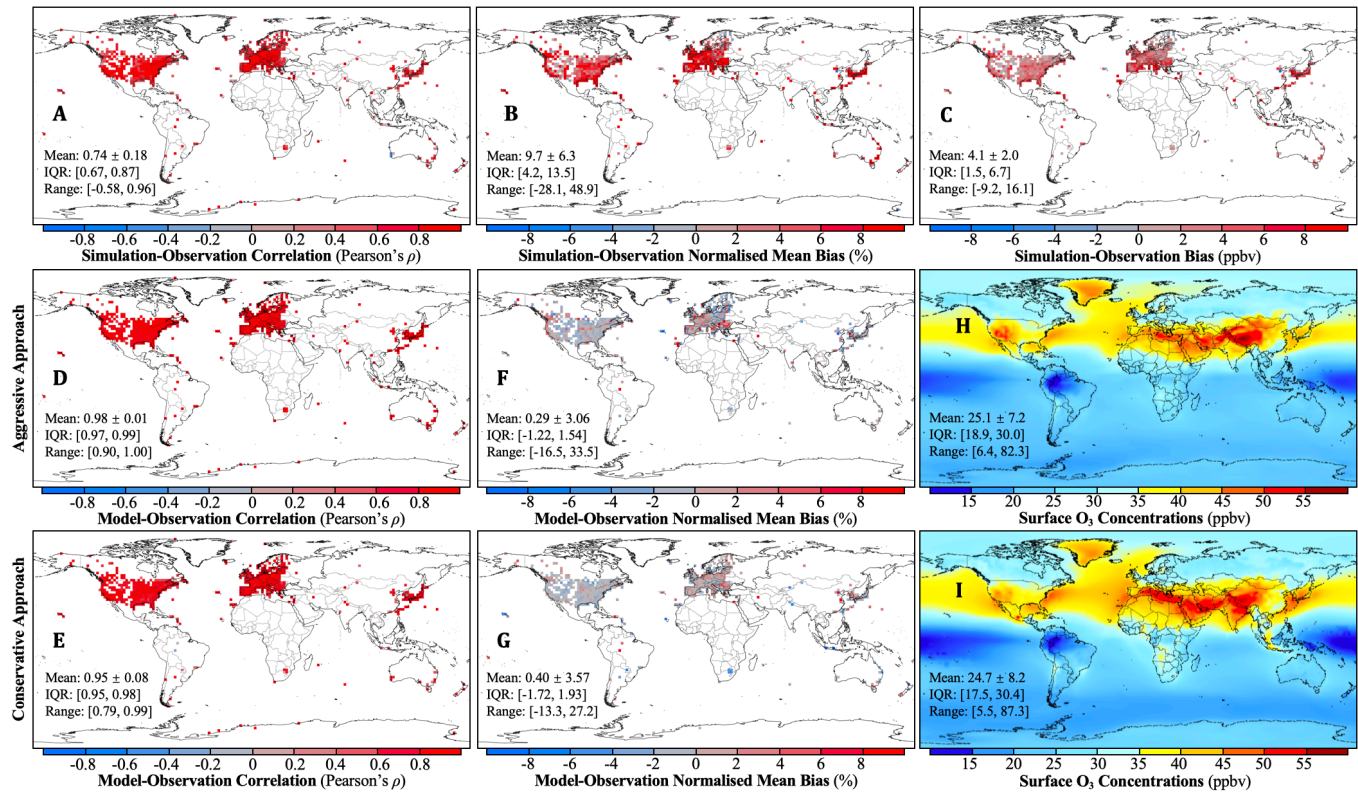
Table 1 Evaluation summary of aggressive and conservative multi-model fusion for surface ozone. The model evaluation metrics include the cross-validation (CV), test and full dataset overall coefficient of determination (R^2), the root mean squared error (RMSE), the normalised mean bias (NMB), and the linear regression slope (k) and intercept (b). Both two statistical models are evaluated separately for each 5-year period, season and continent to assess the spatiotemporal performances.

	Aggressive Approach							Conservative Approach						
	$CV-R^2$	$test-R^2$	$full-R^2$	RMSE	NMB	k	b	$CV-R^2$	$test-R^2$	$full-R^2$	RMSE	NMB	k	b
Period														
1990-1994	0.91	0.90	0.94	2.00	3.41	1.11	-1.62	0.92	0.91	0.93	2.00	0.02	0.98	0.59
1995-1999	0.90	0.90	0.94	1.74	1.71	1.09	-1.26	0.92	0.91	0.92	2.10	0.84	0.97	0.66
2000-2004	0.91	0.91	0.95	1.71	0.88	1.09	-1.16	0.91	0.91	0.93	2.28	0.71	0.97	0.95
2005-2009	0.91	0.91	0.96	1.68	1.11	1.09	-1.17	0.91	0.91	0.91	2.22	0.83	0.97	0.82
2010-2014	0.94	0.93	0.96	1.71	0.88	1.09	-1.16	0.92	0.91	0.94	2.28	0.71	0.97	0.95
Region														
Europe	0.91	0.91	0.94	1.94	2.40	1.12	-1.61	0.92	0.91	0.92	2.02	1.27	0.98	0.37
North America	0.93	0.93	0.96	1.61	1.27	1.08	-1.19	0.91	0.91	0.93	1.96	-0.04	0.97	0.94
Latin America and the Caribbean	0.90	0.87	0.95	1.22	3.12	1.10	-0.89	0.83	0.81	0.83	2.55	3.06	0.92	1.51
Asia	0.92	0.92	0.95	2.14	4.03	1.12	-1.65	0.90	0.90	0.92	2.96	1.85	0.96	0.90
Africa	0.90	0.86	0.90	2.13	2.82	1.19	-2.33	0.82	0.80	0.84	3.69	-3.81	0.93	2.88
Oceania	0.94	0.91	0.96	0.91	0.68	1.08	-0.78	0.83	0.81	0.84	2.13	-1.05	0.88	2.65
Season														
March-May	0.93	0.90	0.97	1.91	0.84	1.13	-0.65	0.94	0.91	0.96	2.06	0.89	0.99	0.97
June-August	0.94	0.92	0.98	1.78	1.12	1.09	-0.86	0.94	0.92	0.95	2.14	0.74	0.97	0.75
September-November	0.93	0.89	0.98	1.75	3.09	1.12	-0.57	0.93	0.90	0.95	2.07	0.10	0.98	0.69
December-February	0.93	0.90	0.98	1.80	3.05	1.14	-0.60	0.93	0.90	0.95	2.19	0.54	0.98	0.51
TOAR	0.94	0.89	0.96	1.81	2.01	1.05	-1.35	0.90	0.88	0.95	2.12	0.57	0.97	0.71

50

51

52



53

54

55

56

57

58

59

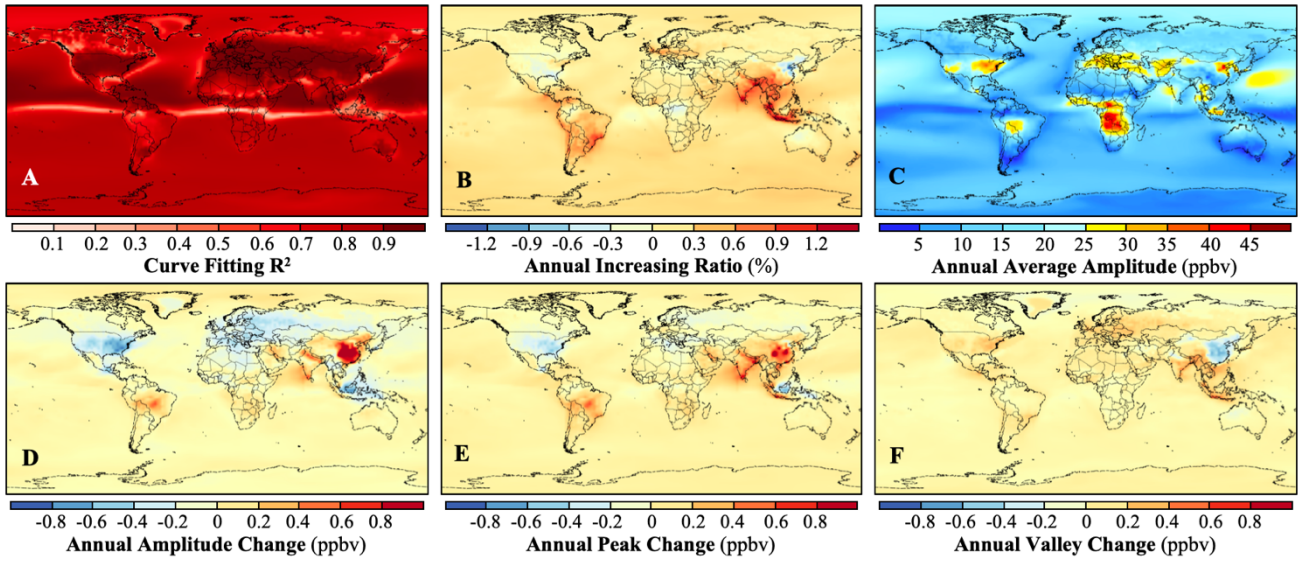
60

61

62

Figure 2 Model-observation evaluation for the raw CMIP6 surface ozone simulation-ensemble and multi-model fusion by both aggressive and conservative approaches. **A-C:** Simulation-observation synchronicity, absolute and relative biases for 57+1 CMIP6 simulation ensemble. Model evaluations are conducted on TOAR observation covered sites across 1990-2014. **D-G:** Evaluations of aggressively and conservatively integrated surface ozone concentrations in terms of the overall model-observation synchronicity and bias. **H-I:** Multi-model and TOAR-observation assimilated historical global surface ozone concentrations by aggressive and conservative approaches. The 25-year average surface ozone concentrations during 1990-2014 are mapped as summary. All spatial resolutions are set as $2^\circ \times 2^\circ$, and the temporal interval is set to month.

33



34

35

36

37

38

39

40

41

42

Figure 3 Spatiotemporal variability parametrisation for CMIP6 multi-model ensemble assimilated surface ozone concentrations during 1990-2014 by the conservative approach. The ensemble-learning predicted concentrations are clustered by month. **A:** Fourier-series function-based curve-fitting quality for grid-specific surface ozone variabilities against temporal sequence, quantified by R^2 . **B:** Annual increasing ratio for yearly average surface ozone concentrations, estimated by $12a_1$. **C:** Annual average intra-year amplitude as the peak-valley differences, estimated by $2b_0$. **D:** Annual average linear change rates of the intra-year amplitudes, estimated by $24b_1$. **E-F:** Averaged annual change rates of peak and valley concentrations, deduced from the fitted second-order Fourier-series function.

73 **Supporting Information.** Further detailed information on CMIP6 AerChemMIP Surface O₃ historical simulation
74 participant research institutes, and annotations on atmospheric module settings. A total of 13 supplementary figures and
75 3 tables.

76

77 **Acknowledgments**

78 The authors are funded by Natural Environment Research Council (NERC) and National Centre for Atmospheric
79 Science (NCAS). We thank Youngsub Matthew Shin (University of Cambridge) for advising the data collection and pre-
30 processing, Ushnish Sengupta (University of Cambridge) and Matt Amos (Lancaster University) for sharing their Python
31 space-time Bayesian neural network core, and Mingtao Xia (University of California, Los Angeles) for scrutinising the
32 code optimisation. We are also grateful to the editor and 3 anonymous reviewers for their insightful revision comments to
33 substantially improve the manuscript.

34

35 **Data and code availability**

36 Core Python codes to construct the first-stage calibration-oriented and second-stage assimilation-targeted Bayesian
37 neural network regressions are available at: <https://github.com/csuen27/BayesNN>, scheduled with regular upgrades every
38 half-year to fit into the latest deep learning frameworks. The CMIP6 simulations with associated metadata can be accessed
39 at: <https://esgf-node.llnl.gov/search/cmip6>. CMIP6 collaborators keep updating the simulation repository, whether adding
40 new ensemble experiments or retracting ones when constructive improvements are to be made, and correspondingly data
41 fusion works will be updated. The up-to-date assimilated surface O₃ concentrations can be shared by the authors for
42 academic use upon request.

43

REFERENCES

1. Stocker, T., *Climate change 2013: the physical science basis: Working Group I contribution to the Fifth assessment report of the Intergovernmental Panel on Climate Change*. Cambridge University Press: 2014.
2. Young, P. J.; Archibald, A. T.; Bowman, K. W.; Lamarque, J. F.; Naik, V.; Stevenson, D. S.; Tilmes, S.; Voulgarakis, A.; Wild, O.; Bergmann, D.; Cameron-Smith, P.; Cionni, I.; Collins, W. J.; Dalsøren, S. B.; Doherty, R. M.; Eyring, V.; Faluvegi, G.; Horowitz, L. W.; Josse, B.; Lee, Y. H.; MacKenzie, I. A.; Nagashima, T.; Plummer, D. A.; Righi, M.; Rumbold, S. T.; Skeie, R. B.; Shindell, D. T.; Strode, S. A.; Sudo, K.; Szopa, S.; Zeng, G., Pre-industrial to end 21st century projections of tropospheric ozone from the Atmospheric Chemistry and Climate Model Intercomparison Project (ACCMIP). *Atmos Chem Phys* **2013**, *13*, (4), 2063-2090.
3. Wilson, S. R.; Madronich, S.; Longstreth, J. D.; Solomon, K. R., Interactive effects of changing stratospheric ozone and climate on tropospheric composition and air quality, and the consequences for human and ecosystem health. *Photochem Photobiol Sci* **2019**, *18*, (3), 775-803.
4. Monks, P. S.; Archibald, A. T.; Colette, A.; Cooper, O.; Coyle, M.; Derwent, R.; Fowler, D.; Granier, C.; Law, K. S.; Mills, G. E.; Stevenson, D. S.; Tarasova, O.; Thouret, V.; von Schneidemesser, E.; Sommariva, R.; Wild, O.; Williams, M. L., Tropospheric ozone and its precursors from the urban to the global scale from air quality to short-lived climate forcer. *Atmos Chem Phys* **2015**, *15*, (15), 8889-8973.
5. Avnery, S.; Mauzerall, D. L.; Liu, J.; Horowitz, L. W., Global crop yield reductions due to surface ozone exposure: 1. Year 2000 crop production losses and economic damage. *Atmos Environ* **2011**, *45*, (13), 2284-2296.
6. Ghude, S. D.; Jena, C.; Chate, D. M.; Beig, G.; Pfister, G. G.; Kumar, R.; Ramanathan, V., Reductions in India's crop yield due to ozone. *Geophys Res Lett* **2014**, *41*, (15), 5685-5691.
7. Wiegman, C. H.; Li, F.; Clarke, C. J.; Jazrawi, E.; Kirkham, P.; Barnes, P. J.; Adcock, I. M.; Chung, K. F., A comprehensive analysis of oxidative stress in the ozone-induced lung inflammation mouse model. *Clin Sci (Lond)* **2014**, *126*, (6), 425-440.
8. McConnell, R.; Berhane, K.; Gilliland, F.; London, S. J.; Islam, T.; Gauderman, W. J.; Avol, E.; Margolis, H. G.; Peters, J. M., Asthma in exercising children exposed to ozone: a cohort study. *Lancet* **2002**, *359*, (9304), 386-391.
9. Sheffield, P. E.; Zhou, J.; Shmool, J. L.; Clougherty, J. E., Ambient ozone exposure and children's acute asthma in New York City: a case-crossover analysis. *Environ Health* **2015**, *14*, (1), 25.
10. Ghude, S. D.; Chate, D. M.; Jena, C.; Beig, G.; Kumar, R.; Barth, M. C.; Pfister, G. G.; Fadnavis, S.; Pithani, P., Premature mortality in India due to PM_{2.5} and ozone exposure. *Geophys Res Lett* **2016**, *43*, (9), 4650-4658.
11. Qiu, X.; Wei, Y.; Wang, Y.; Di, Q.; Sofer, T.; Awad, Y. A.; Schwartz, J., Inverse probability weighted distributed lag effects of short-term exposure to PM_{2.5} and ozone on CVD hospitalisations in New England Medicare participants - Exploring the causal effects. *Environ Res* **2020**, *182*, 109095.
12. Turner, M. C.; Jerrett, M.; Pope, C. A., 3rd; Krewski, D.; Gapstur, S. M.; Diver, W. R.; Beckerman, B. S.; Marshall, J. D.; Su, J.; Crouse, D. L.; Burnett, R. T., Long-Term Ozone Exposure and Mortality in a Large Prospective Study. *Am J Respir Crit Care Med* **2016**, *193*, (10), 1134-42.
13. Di, Q.; Wang, Y.; Zanobetti, A.; Wang, Y.; Koutrakis, P.; Choirat, C.; Dominici, F.; Schwartz, J. D., Air Pollution and Mortality in the Medicare Population. *N Engl J Med* **2017**, *376*, (26), 2513-2522.
14. Bell, M. L.; McDermott, A.; Zeger, S. L.; Samet, J. M.; Dominici, F., Ozone and short-term mortality in 95 US urban communities, 1987-2000. *JAMA* **2004**, *292*, (19), 2372-8.
15. Institute for Health Metrics and Evaluation., GBD Compare Data Visualization. <http://vizhub.healthdata.org/gbd-compare> (accessed on December 18, 2020).
16. Weichenthal, S.; Pinault, L. L.; Burnett, R. T., Impact of oxidant gases on the relationship between outdoor fine particulate air pollution and nonaccidental, cardiovascular, and respiratory mortality. *Scientific Reports* **2017**, *7*, (1), 1-10.
17. Schultz, M. G.; Schröder, S.; Lyapina, O.; Cooper, O.; Galbally, I.; Petropavlovskikh, I.; Von Schneidemesser, E.; Tanimoto, H.; Elshorbany, Y.; Naja, M.; Seguel, R.; Dauert, U.; Eckhardt, P.; Feigenspahn, S.; Fiebig, M.; Hjellbrekke, A.-G.; Hong, Y.-D.; Christian Kjeld, P.; Koide, H.; Lear, G.; Tarasick, D.; Ueno, M.; Wallasch, M.; Baumgardner, D.; Chuang, M.-T.; Gillett, R.; Lee, M.; Molloy, S.; Moolla, R.; Wang, T.; Sharps, K.; Adame, J. A.; Ancellet, G.; Apadula, F.; Artaxo, P.; Barlasina, M.; Bogucka, M.; Bonasoni, P.; Chang, L.; Colomb, A.; Cuevas, E.; Cupeiro, M.; Degorska, A.; Ding, A.; Fröhlich, M.; Frolova, M.; Gadhave, H.; Gheusi, F.; Gilge, S.; Gonzalez, M. Y.; Gros, V.; Hamad, S. H.; Helmig, D.; Henriques, D.; Hermansen, O.; Holla, R.; Huber, J.; Im, U.; Jaffe, D. A.; Komala, N.; Kubistin, D.; Lam, K.-S.; Laurila, T.; Lee, H.; Levy, I.; Mazzoleni, C.; Mazzoleni, L.; McClure-Begley, A.; Mohamad, M.; Murovic, M.; Navarro-Comas, M.; Nicodim, F.; Parrish, D.; Read, K. A.; Reid, N.; Ries, L.; Saxena, P.; Schwab, J. J.; Scorgie, Y.; Senik, I.; Simmonds, P.; Sinha, V.; Skorokhod, A.; Spain, G.; Spangl, W.; Spoor, R.; Springston, S. R.; Steer, K.; Steinbacher, M.; Suharguniyawan, E.; Torre, P.; Trickl, T.; Weili, L.; Weller, R.; Xu, X.; Xue, L.; Zhiqiang, M., Tropospheric Ozone Assessment Report: Database and Metrics Data of Global Surface Ozone Observations. *Elementa-Sci Anthropol* **2017**, *5*, 58-83.

18. Zoogman, P.; Jacob, D. J.; Chance, K.; Worden, H. M.; Edwards, D. P.; Zhang, L., Improved monitoring of surface ozone by joint assimilation of geostationary satellite observations of ozone and CO. *Atmos Environ* **2014**, *84*, 254-261.
19. Archibald, A. T.; O'Connor, F. M.; Abraham, N. L.; Archer-Nicholls, S.; Chipperfield, M. P.; Dalvi, M.; Folberth, G. A.; Dennison, F.; Dhomse, S. S.; Griffiths, P. T.; Hardacre, C.; Hewitt, A. J.; Hill, R.; Johnson, C. E.; Keeble, J.; Köhler, M. O.; Morgenstern, O.; Mulchay, J. P.; Ordóñez, C.; Pope, R. J.; Rumbold, S.; Russo, M. R.; Savage, N.; Sellar, A.; Stringer, M.; Turnock, S.; Wild, O.; Zeng, G., Description and evaluation of the UKCA stratosphere-troposphere chemistry scheme (StratTrop vn 1.0) implemented in UKESM1. *Geosci Model Dev* **2019**, 1223-1266.
20. Wang, T.; Xue, L.; Brimblecombe, P.; Lam, Y. F.; Li, L.; Zhang, L., Ozone pollution in China: A review of concentrations, meteorological influences, chemical precursors, and effects. *Sci Total Environ* **2017**, *575*, 1582-1596.
21. Archibald, A.; Neu, J.; Elshorbany, Y.; Cooper, O.; Young, P.; Akiyoshi, H.; Cox, R.; Coyle, M.; Derwent, R.; Deushi, M., Tropospheric Ozone Assessment Report A critical review of changes in the tropospheric ozone burden and budget from 1850 to 2100. *Elementa-Sci Anthropol* **2020**, *8*, (1), 34-86.
22. Byun, D.; Schere, K. L., Review of the Governing Equations, Computational Algorithms, and Other Components of the Models-3 Community Multiscale Air Quality (CMAQ) Modeling System. *Appl Mech Rev* **2006**, *59*, (2), 51-77.
23. Appel, K. W.; Napelenok, S. L.; Foley, K. M.; Pye, H. O. T.; Hogrefe, C.; Luecken, D. J.; Bash, J. O.; Roselle, S. J.; Pleim, J. E.; Foroutan, H.; Hutzell, W. T.; Pouliot, G. A.; Sarwar, G.; Fahey, K. M.; Gantt, B.; Gilliam, R. C.; Heath, N. K.; Kang, D.; Mathur, R.; Schwede, D. B.; Spero, T. L.; Wong, D. C.; Young, J. O., Description and evaluation of the Community Multiscale Air Quality (CMAQ) modeling system version 5.1. *Geosci Model Dev* **2017**, *10*, (4), 1703-1732.
24. Tesche, T. W.; Morris, R.; Tonnesen, G.; McNally, D.; Boylan, J.; Brewer, P., CMAQ/CAMx annual 2002 performance evaluation over the eastern US. *Atmos Environ* **2006**, *40*, (26), 4906-4919.
25. Mar, K. A.; Ojha, N.; Pozzer, A.; Butler, T. M., Ozone air quality simulations with WRF-Chem (v3.5.1) over Europe: model evaluation and chemical mechanism comparison. *Geosci Model Dev* **2016**, *9*, (10), 3699-3728.
26. Mann, G. W.; Carslaw, K. S.; Spracklen, D. V.; Ridley, D. A.; Manktelow, P. T.; Chipperfield, M. P.; Pickering, S. J.; Johnson, C. E., Description and evaluation of GLOMAP-mode: a modal global aerosol microphysics model for the UKCA composition-climate model. *Geosci Model Dev* **2010**, *3*, (2), 519-551.
27. McLaren, A.; Banks, H.; Durman, C.; Gregory, J.; Johns, T.; Keen, A.; Ridley, J.; Roberts, M.; Lipscomb, W.; Connolley, W., Evaluation of the sea ice simulation in a new coupled atmosphere-ocean climate model (HadGEM1). *Journal of Geophysical Research: Oceans* **2006**, *111*, (C12), 14-30.
28. Sellar, A. A.; Walton, J.; Jones, C. G.; Wood, R.; Abraham, N. L.; Andrejczuk, M.; Andrews, M. B.; Andrews, T.; Archibald, A. T.; Mora, L.; Dyson, H.; Elkington, M.; Ellis, R.; Florek, P.; Good, P.; Gohar, L.; Haddad, S.; Hardiman, S. C.; Hogan, E.; Iwi, A.; Jones, C. D.; Johnson, B.; Kelley, D. I.; Kettleborough, J.; Knight, J. R.; Köhler, M. O.; Kuhlbrodt, T.; Liddicoat, S.; Linova-Pavlova, I.; Mizielski, M. S.; Morgenstern, O.; Mulchay, J.; Neininger, E.; O'Connor, F. M.; Petrie, R.; Ridley, J.; Rioual, J. C.; Roberts, M.; Robertson, E.; Rumbold, S.; Seddon, J.; Shepherd, H.; Shim, S.; Stephens, A.; Teixeira, J. C.; Tang, Y.; Williams, J.; Wiltshire, A.; Griffiths, P. T., Implementation of U.K. Earth System Models for CMIP6. *J Adv Model Earth Syst* **2020**, *12*, (4), e2019MS001946.
29. Young, P. J.; Naik, V.; Fiore, A. M.; Gaudel, A.; Guo, J.; Lin, M. Y.; Neu, J. L.; Parrish, D. D.; Rieder, H. E.; Schnell, J. L.; Tilmes, S.; Wild, O.; Zhang, L.; Ziemke, J.; Brandt, J.; Delcloo, A.; Doherty, R. M.; Geels, C.; Hegglin, M. I.; Hu, L.; Im, U.; Kumar, R.; Luhar, A.; Murray, L.; Plummer, D.; Rodriguez, J.; Saiz-Lopez, A.; Schultz, M. G.; Woodhouse, M. T.; Zeng, G.; Helmig, D.; Lewis, A., Tropospheric Ozone Assessment Report: Assessment of global-scale model performance for global and regional ozone distributions, variability, and trends. *Elementa-Sci Anthropol* **2018**, *6*, (1), 10-50.
30. Eyring, V.; Bony, S.; Meehl, G. A.; Senior, C. A.; Stevens, B.; Stouffer, R. J.; Taylor, K. E., Overview of the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organisation. *Geosci Model Dev* **2016**, *9*, (5), 1937-1958.
31. Griffiths, P. T.; Murray, L. T.; Zeng, G.; Archibald, A. T.; Emmons, L. K.; Galbally, I.; Hassler, B.; Horowitz, L. W.; Keeble, J.; Liu, J.; Moeini, O.; Naik, V.; O'Connor, F. M.; Shin, Y. M.; Tarasick, D.; Tilmes, S.; Turnock, S. T.; Wild, O.; Young, P. J.; Zanis, P., Tropospheric ozone in CMIP6 Simulations. *Atmos Chem Phys* **2021**, *21*, (5), 4187-4218.
32. Stouffer, R. J.; Eyring, V.; Meehl, G. A.; Bony, S.; Senior, C.; Stevens, B.; Taylor, K. E., CMIP5 Scientific Gaps and Recommendations for CMIP6. *B Am Meteorol Soc* **2017**, *98*, (1), 95-105.
33. Feng, L.; Smith, S. J.; Braun, C.; Crippa, M.; Gidden, M. J.; Hoesly, R.; Klimont, Z.; van Marle, M.; van den Berg, M.; van der Werf, G. R., The generation of gridded emissions data for CMIP6. *Geosci Model Dev* **2020**, *13*, (2), 461-482.
34. Collins, W. J.; Lamarque, J.-F.; Schulz, M.; Boucher, O.; Eyring, V.; Hegglin, M. I.; Maycock, A.; Myhre, G.; Prather, M.; Shindell, D.; Smith, S. J., AerChemMIP: quantifying the effects of chemistry and aerosols in CMIP6. *Geosci Model Dev* **2017**, *10*, (2), 585-607.
35. Turnock, S. T.; Allen, R. J.; Andrews, M.; Bauer, S. E.; Deushi, M.; Emmons, L.; Good, P.; Horowitz, L.; John, J. G.; Michou, M.; Nabat, P.; Naik, V.; Neubauer, D.; O'Connor, F. M.; Olivie, D.; Oshima, N.; Schulz, M.; Sellar, A.; Shim, S.; Takemura, T.; Tilmes, S.; Tsigaridis, K.; Wu, T.; Zhang, J., Historical and future changes in air pollutants from CMIP6 models. *Atmos Chem Phys* **2020**, *20*, (23), 14547-14579.
36. Di, Q.; Amini, H.; Shi, L.; Kloog, I.; Silvern, R.; Kelly, J.; Sabath, M. B.; Choirat, C.; Koutrakis, P.; Lyapustin, A., Assessing

- NO₂ Concentration and Model Uncertainty with High Spatiotemporal Resolution across the Contiguous United States Using Ensemble Model Averaging. *Environ Sci Technol* **2019**, *54*, (3), 1372-1384.
37. Di, Q.; Amini, H.; Shi, L.; Kloog, I.; Silvern, R.; Kelly, J.; Sabath, M. B.; Choirat, C.; Koutrakis, P.; Lyapustin, A.; Wang, Y.; Mickley, L. J.; Schwartz, J., An ensemble-based model of PM_{2.5} concentration across the contiguous United States with high spatiotemporal resolution. *Environ Int* **2019**, *130*, 104909.
 38. Di, Q.; Kloog, I.; Koutrakis, P.; Lyapustin, A.; Wang, Y.; Schwartz, J., Assessing PM_{2.5} Exposures with High Spatiotemporal Resolution across the Continental United States. *Environ Sci Technol* **2016**, *50*, (9), 4712-21.
 39. Lyu, B.; Hu, Y.; Zhang, W.; Du, Y.; Luo, B.; Sun, X.; Sun, Z.; Deng, Z.; Wang, X.; Liu, J.; Wang, X.; Russell, A. G., Fusion Method Combining Ground-Level Observations with Chemical Transport Model Predictions Using an Ensemble Deep Learning Framework: Application in China to Estimate Spatiotemporally-Resolved PM_{2.5} Exposure Fields in 2014-2017. *Environ Sci Technol* **2019**, *53*, (13), 7306-7315.
 40. Ren, X.; Mi, Z.; Georgopoulos, P. G., Comparison of Machine Learning and Land Use Regression for fine scale spatiotemporal estimation of ambient air pollution: Modeling ozone concentrations across the contiguous United States. *Environ Int* **2020**, *142*, 105827.
 41. Rudin, C., Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell* **2019**, *1*, (5), 206-215.
 42. Wang, A.; Xu, J.; Tu, R.; Saleh, M.; Hatzopoulou, M., Potential of machine learning for prediction of traffic related air pollution. *Transport Res D-Tr E* **2020**, *88*, 102599.
 43. Danek, C.; Shi, X.; Stepanek, C.; Yang, H.; Barbi, D.; Hegewald, J.; Lohmann, G., AWI AWI-ESM1.1-LR model output prepared for CMIP6 CMIP historical. In Earth System Grid Federation: 2020.
 44. Wu, T.; Lu, Y.; Fang, Y.; Xin, X.; Li, L.; Li, W.; Jie, W.; Zhang, J.; Liu, Y.; Zhang, L.; Zhang, F.; Zhang, Y.; Wu, F.; Li, J.; Chu, M.; Wang, Z.; Shi, X.; Liu, X.; Wei, M.; Huang, A.; Zhang, Y.; Liu, X., The Beijing Climate Center Climate System Model (BCC-CSM): the main progress from CMIP5 to CMIP6. *Geosci Model Dev* **2019**, *12*, (4), 1573-1600.
 45. Wu, T.; Yu, R.; Lu, Y.; Jie, W.; Fang, Y.; Zhang, J.; Zhang, L.; Xin, X.; Li, L.; Wang, Z., BCC-CSM2-HR: A High-Resolution Version of the Beijing Climate Center Climate System Model. *Geosci Model Dev* **2020**, *14*, (5), 2977-3006.
 46. Wu, T.; Chu, M.; Dong, M.; Fang, Y.; Jie, W.; Li, J.; Li, W.; Liu, Q.; Shi, X.; Xin, X.; Yan, J.; Zhang, F.; Zhang, J.; Zhang, L.; Zhang, Y., BCC BCC-CSM2-MR model output prepared for CMIP6 CMIP piControl. In Earth System Grid Federation: 2018.
 47. Boucher, O.; Servonnat, J.; Albright, A. L.; Aumont, O.; Balkanski, Y.; Bastrikov, V.; Bekki, S.; Bonnet, R.; Bony, S.; Bopp, L., Presentation and evaluation of the IPSL-CM6A-LR climate model. *J Adv Model Earth Syst* **2020**, *12*, (7), e2019MS002010.
 48. Boucher, O.; Denvil, S.; Levvasseur, G.; Cozic, A.; Caubel, A.; Foujols, M.-A.; Meurdesoif, Y.; Cadule, P.; Devilliers, M.; Ghattas, J.; Lebas, N.; Lurton, T.; Mellul, L.; Musat, I.; Mignot, J.; Cheruy, F., IPSL IPSL-CM6A-LR model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2018.
 49. Mauritsen, T.; Bader, J.; Becker, T.; Behrens, J.; Bittner, M.; Brokopf, R.; Brovkin, V.; Claussen, M.; Crueger, T.; Esch, M.; Fast, I.; Fiedler, S.; Flaschner, D.; Gayler, V.; Giorgetta, M.; Goll, D. S.; Haak, H.; Hagemann, S.; Hedemann, C.; Hohenegger, C.; Ilyina, T.; Jahns, T.; Jimenez-de-la-Cuesta, D.; Jungclaus, J.; Kleinen, T.; Kloster, S.; Kracher, D.; Kinne, S.; Kleberg, D.; Lasslop, G.; Kornblueh, L.; Marotzke, J.; Matei, D.; Meraner, K.; Mikolajewicz, U.; Modali, K.; Mobis, B.; Muller, W. A.; Nabel, J.; Nam, C. C. W.; Notz, D.; Nyawira, S. S.; Paulsen, H.; Peters, K.; Pincus, R.; Pohlmann, H.; Pongratz, J.; Popp, M.; Raddatz, T. J.; Rast, S.; Redler, R.; Reick, C. H.; Rohrschneider, T.; Schemann, V.; Schmidt, H.; Schnur, R.; Schulzweida, U.; Six, K. D.; Stein, L.; Stemmler, I.; Stevens, B.; von Storch, J. S.; Tian, F.; Voigt, A.; Vrese, P.; Wieners, K. H.; Wilkenskjaeld, S.; Winkler, A.; Roeckner, E., Developments in the MPI-M Earth System Model version 1.2 (MPI-ESM1.2) and Its Response to Increasing CO₂. *J Adv Model Earth Syst* **2019**, *11*, (4), 998-1038.
 50. Müller, W. A.; Jungclaus, J. H.; Mauritsen, T.; Baehr, J.; Bittner, M.; Budich, R.; Bunzel, F.; Esch, M.; Ghosh, R.; Haak, H., A Higher-resolution Version of the Max Planck Institute Earth System Model (MPI-ESM1.2-HR). *J Adv Model Earth Syst* **2018**, *10*, (7), 1383-1413.
 51. Gutjahr, O.; Putrasahan, D.; Lohmann, K.; Jungclaus, J. H.; von Storch, J.-S.; Brüggemann, N.; Haak, H.; Stössel, A., Max Planck Institute Earth System Model (MPI-ESM1.2) for the High-Resolution Model Intercomparison Project (HighResMIP). *Geosci Model Dev* **2019**, *12*, (7), 3241-3281.
 52. von Storch, J.-S.; Putrasahan, D.; Lohmann, K.; Gutjahr, O.; Jungclaus, J.; Bittner, M.; Haak, H.; Wieners, K.-H.; Giorgetta, M.; Reick, C.; Esch, M.; Gayler, V.; de Vrese, P.; Raddatz, T.; Mauritsen, T.; Behrens, J.; Brovkin, V.; Claussen, M.; Crueger, T.; Fast, I.; Fiedler, S.; Hagemann, S.; Hohenegger, C.; Jahns, T.; Kloster, S.; Kinne, S.; Lasslop, G.; Kornblueh, L.; Marotzke, J.; Matei, D.; Meraner, K.; Mikolajewicz, U.; Modali, K.; Müller, W.; Nabel, J.; Notz, D.; Peters, K.; Pincus, R.; Pohlmann, H.; Pongratz, J.; Rast, S.; Schmidt, H.; Schnur, R.; Schulzweida, U.; Six, K.; Stevens, B.; Voigt, A.; Roeckner, E., MPI-M MPI-ESM1.2-XR model output prepared for CMIP6 HighResMIP. In Earth System Grid Federation: 2017.
 53. Voldoire, A.; Saint-Martin, D.; Sénési, S.; Decharme, B.; Alias, A.; Chevallier, M.; Colin, J.; Guérémy, J. F.; Michou, M.; Moine, M. P., Evaluation of CMIP6 deck experiments with CNRM-CM6-1. *J Adv Model Earth Syst* **2019**, *11*, (7), 2177-2213.
 54. Michou, M.; Nabat, P.; Saint-Martin, D.; Bock, J.; Decharme, B.; Mallet, M.; Roehrig, R.; Séférian, R.; Sénési, S.; Voldoire, A., Present-day and historical aerosol and ozone characteristics in CNRM CMIP6 simulations. *J Adv Model Earth Syst* **2020**, *12*,

(1), e2019MS001816.

55. Voldoire, A., CNRM-CERFACS CNRM-CM6-1 model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2018.
56. Séférian, R.; Nabat, P.; Michou, M.; Saint-Martin, D.; Voldoire, A.; Colin, J.; Decharme, B.; Delire, C.; Berthet, S.; Chevallier, M., Evaluation of CNRM Earth System Model, CNRM-ESM2-1: Role of Earth System Processes in Present-Day and Future Climate. *J Adv Model Earth Syst* **2019**, *11*, (12), 4182-4227.
57. Seferian, R., CNRM-CERFACS CNRM-ESM2-1 model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2018.
58. Wu, T.; Zhang, F.; Zhang, J.; Jie, W.; Zhang, Y.; Wu, F.; Li, L.; Yan, J.; Liu, X.; Lu, X.; Tan, H.; Zhang, L.; Wang, J.; Hu, A., Beijing Climate Center Earth System Model version 1 (BCC-ESM1): model description and evaluation of aerosol simulations. *Geosci Model Dev* **2020**, *13*, (3), 977-1005.
59. Zhang, J.; Wu, T.; Shi, X.; Zhang, F.; Li, J.; Chu, M.; Liu, Q.; Yan, J.; Ma, Q.; Wei, M., BCC BCC-ESM1 model output prepared for CMIP6 CMIP piControl. In Earth System Grid Federation: 2018.
60. Neubauer, D.; Ferrachat, S.; Siegenthaler-Le Drian, C.; Stoll, J.; Folini, D. S.; Tegen, I.; Wieners, K.-H.; Mauritsen, T.; Stemmler, I.; Barthel, S.; Bey, I.; Daskalakis, N.; Heinold, B.; Kokkola, H.; Partridge, D.; Rast, S.; Schmidt, H.; Schutgens, N.; Stanelle, T.; Stier, P.; Watson-Parris, D.; Lohmann, U., HAMMOZ-Consortium MPI-ESM1.2-HAM model output prepared for CMIP6 AerChemMIP. In Earth System Grid Federation: 2019.
61. Yukimoto, S.; Kawai, H.; Koshiro, T.; Oshima, N.; Yoshida, K.; Urakawa, S.; Tsujino, H.; Deushi, M.; Tanaka, T.; Hosaka, M.; Yabu, S.; Yoshimura, H.; Shindo, E.; Mizuta, R.; Obata, A.; Adachi, Y.; Ishii, M., The Meteorological Research Institute Earth System Model Version 2.0, MRI-ESM2.0: Description and Basic Evaluation of the Physical Component. *J Meteorol Soc Jpn* **2019**, *97*, (5), 931-965.
62. Yukimoto, S.; Adachi, Y.; Hosaka, M.; Sakami, T.; Yoshimura, H.; Hirabara, M.; Tanaka, T. Y.; Shindo, E.; Tsujino, H.; Deushi, M.; Mizuta, R.; Yabu, S.; Obata, A.; Nakano, H.; Koshiro, T.; Ose, T.; Kitoh, A., A New Global Climate Model of the Meteorological Research Institute: MRI-CGCM3 Model Description and Basic Performance. *J Meteorol Soc Jpn* **2012**, *90A*, (2), 23-64.
63. Yukimoto, S.; Koshiro, T.; Kawai, H.; Oshima, N.; Yoshida, K.; Urakawa, S.; Tsujino, H.; Deushi, M.; Tanaka, T.; Hosaka, M.; Yoshimura, H.; Shindo, E.; Mizuta, R.; Ishii, M.; Obata, A.; Adachi, Y., MRI MRI-ESM2.0 model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2019.
64. Shindell, D. T.; Pechony, O.; Voulgarakis, A.; Faluvegi, G.; Nazarenko, L.; Lamarque, J. F.; Bowman, K.; Milly, G.; Kovari, B.; Ruedy, R.; Schmidt, G. A., Interactive ozone and methane chemistry in GISS-E2 historical and future climate simulations. *Atmos Chem Phys* **2013**, *13*, (5), 2653-2689.
65. NASA-GISS, NASA-GISS GISS-E2.1G model output prepared for CMIP6 ISMIP6. In Earth System Grid Federation: 2018.
66. NASA-GISS, NASA-GISS GISS-E2.1H model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2018.
67. Gettelman, A.; Mills, M. J.; Kinnison, D. E.; Garcia, R. R.; Smith, A. K.; Marsh, D. R.; Tilmes, S.; Vitt, F.; Bardeen, C. G.; McNerny, J.; Liu, H. L.; Solomon, S. C.; Polvani, L. M.; Emmons, L. K.; Lamarque, J. F.; Richter, J. H.; Glanville, A. S.; Bacmeister, J. T.; Phillips, A. S.; Neale, R. B.; Simpson, I. R.; DuVivier, A. K.; Hodzic, A.; Randel, W. J., The Whole Atmosphere Community Climate Model Version 6 (WACCM6). *J Geophys Res-Atmos* **2019**, *124*, (23), 12380-12403.
68. Danabasoglu, G., NCAR CESM2-WACCM model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2019.
69. Seland, Ø.; Bentsen, M.; Olivie, D. J. L.; Toniazzo, T.; Gjermundsen, A.; Graff, L. S.; Debernard, J. B.; Gupta, A. K.; He, Y.; Kirkevåg, A.; Schwinger, J.; Tjiputra, J.; Aas, K. S.; Bethke, I.; Fan, Y.; Griesfeller, J.; Grini, A.; Guo, C.; Ilicak, M.; Karset, I. H. H.; Landgren, O. A.; Liakka, J.; Moseid, K. O.; Nummelin, A.; Spensberger, C.; Tang, H.; Zhang, Z.; Heinze, C.; Iversen, T.; Schulz, M., NCC NorESM2-LM model output prepared for CMIP6 CMIP historical. In Earth System Grid Federation: 2019.
70. Krasting, J.; John, J.; Blanton, C.; McHugh, C.; Nikonov, S.; Radhakrishnan, A.; Rand, K.; Zadeh, N.; Balaji, V.; Durachta, J., NOAA-GFDL GFDL-ESM4 model output prepared for CMIP6 CMIP. In Earth System Grid Federation, 2018.
71. Horowitz, L. W.; Naik, V.; Sentman, L.; Paulot, F.; Blanton, C.; McHugh, C.; Radhakrishnan, A.; Rand, K.; Vahlenkamp, H.; Zadeh, N. T.; Wilson, C.; Ginoux, P.; He, J.; John, J. G.; Lin, M.; Paynter, D. J.; Ploshay, J.; Zhang, A.; Zeng, Y., NOAA-GFDL GFDL-ESM4 model output prepared for CMIP6 AerChemMIP. In Earth System Grid Federation: 2018.
72. Mulcahy, J. P.; Jones, C.; Sellar, A.; Johnson, B.; Boutle, I. A.; Jones, A.; Andrews, T.; Rumbold, S. T.; Mollard, J.; Bellouin, N.; Johnson, C. E.; Williams, K. D.; Grosvenor, D. P.; McCoy, D. T., Improved Aerosol Processes and Effective Radiative Forcing in HadGEM3 and UKESM1. *J Adv Model Earth Syst* **2018**, *10*, (11), 2786-2805.
73. Sellar, A. A.; Jones, C. G.; Mulcahy, J. P.; Tang, Y.; Yool, A.; Wiltshire, A.; O'connor, F. M.; Stringer, M.; Hill, R.; Palmieri, J., UKESM1: Description and evaluation of the UK Earth System Model. *J Adv Model Earth Syst* **2019**, *11*, (12), 4513-4558.
74. Tang, Y.; Rumbold, S.; Ellis, R.; Kelley, D.; Mulcahy, J.; Sellar, A.; Walton, J.; Jones, C., MOHC UKESM1.0-LL model output prepared for CMIP6 CMIP. In Earth System Grid Federation: 2019.
75. Yool, A.; Palmieri, J.; Jones, C.; Sellar, A.; de Mora, L.; Kuhlbrodt, T.; Popova, E.; Mulcahy, J.; Wiltshire, A.; Rumbold, S. T., Spin-up of UK Earth System Model 1 (UKESM1) for CMIP6. *J Adv Model Earth Syst* **2020**, e2019MS001933.

- 19 76. Hegglin, M.; Kinnison, D.; Lamarque, J.-F.; Plummer, D., CCMI ozone in support of CMIP6 - version 1.0. In Earth System
20 Grid Federation: 2016.
- 21 77. Danielson, J. J.; Gesch, D. B., Global multi-resolution terrain elevation data 2010 (GMTED2010). US Department of the
22 Interior, US Geological Survey: 2011.
- 23 78. Lloyd, C. T.; Sorichetta, A.; Tatem, A. J., High resolution global gridded data for use in population studies. *Scientific Data*
24 **2017**, *4*, (1), 1-17.
- 25 79. Sengupta, U.; Amos, M.; Hosking, S.; Rasmussen, C. E.; Juniper, M.; Young, P., Ensembling geophysical models with Bayesian
26 Neural Networks. *Advances in Neural Information Processing Systems* **2020**, *33*.
- 27 80. Sexton, J.; Laake, P., Standard errors for bagged and random forest estimators. *Computational Statistics & Data Analysis* **2009**,
28 *53*, (3), 801-811.
- 29 81. Solazzo, E.; Galmarini, S., Error apportionment for atmospheric chemistry-transport models – a new approach to model
30 evaluation. *Atmos Chem Phys* **2016**, *16*, (10), 6263-6283.
- 31 82. Kovač-Andrić, E.; Brana, J.; Gvozdić, V., Impact of meteorological factors on ozone concentrations modelled by time series
32 analysis and multivariate statistical methods. *Ecological Informatics* **2009**, *4*, (2), 117-122.
- 33 83. Solazzo, E.; Galmarini, S., Error apportionment for atmospheric chemistry-transport models: a new approach to model
34 evaluation. *Atmos Chem Phys* **2016**, *16*, (10), 6263-6283.
- 35 84. Hakim, Z. Q.; Archer-Nicholls, S.; Beig, G.; Folberth, G. A.; Sudo, K.; Abraham, N. L.; Ghude, S.; Henze, D. K.; Archibald, A.
36 T., Evaluation of tropospheric ozone and ozone precursors in simulations from the HTAPII and CCMI model intercomparisons–
37 a focus on the Indian subcontinent. *Atmos Chem Phys* **2019**, *19*, (9), 6437-6458.
- 38 85. Derwent, R. G.; Parrish, D. D.; Archibald, A. T.; Deushi, M.; Bauer, S. E.; Tsigaridis, K.; Shindell, D.; Horowitz, L. W.; Khan,
39 M. A. H.; Shallcross, D. E., Intercomparison of the representations of the atmospheric chemistry of pre-industrial methane and
40 ozone in earth system and other global chemistry-transport models. *Atmos Environ* **2021**, 118248.
- 41