



1 **Recommendations on benchmarks for chemical transport model** 2 **applications in China – Part 2: Ozone and Uncertainty Analysis**

3 Ling Huang¹, Xinxin Zhang¹, Chris Emery², Qing Mu³, Greg Yarwood², Hehe Zhai¹, Zhixu Sun¹,
4 Shuhui Xue¹, Yangjun Wang¹, Joshua S Fu⁴, Li Li^{1*}

5 ¹School of Environmental and Chemical Engineering, Shanghai University, Shanghai, 200444, China

6 ²Ramboll, Novato, California, 94945, USA

7 ³Department of Health and Environmental Sciences, Xi'an Jiaotong-Liverpool University, Suzhou, China

8 ⁴Department of Civil and Environmental Engineering, University of Tennessee, Knoxville, TN 37996, USA

9

10 *Correspondence to:* Li Li (lily@shu.edu.cn)

11 **Abstract**

12 Ground-level ozone (O₃) has emerged as a significant air pollutant in China, attracting increasing attention from
13 both the scientific community and policymakers. Chemical transport models (CTM) serve as crucial tools in
14 addressing O₃ pollution, with frequent applications in predicting O₃ concentrations, identifying source
15 contributions, and formulating effective control strategies. The accuracy and reliability of the simulated O₃
16 concentrations are typically assessed through model performance evaluation (MPE). However, the wide array of
17 CTMs available, variations in input data, model setups, and other factors result in a broad range of simulated O₃
18 concentration differences from observed values, highlighting the necessity for standardized benchmarks in O₃
19 evaluation.

20 Built upon our previous work, this study conducted a thorough literature review of CTM applications simulating
21 O₃ in China from 2006 to 2021. 216 relevant articles out of a total of 667 reviewed were identified to extract
22 quantitative MPE results and key model configurations. From our analysis, two sets of benchmark values for six
23 commonly used MPE metrics are proposed for CTM applications in China, categorized into “goal” benchmarks
24 representing optimal model performance and “criteria” benchmarks representing achievable model performance
25 across a majority of studies. It is recommended that the normalized mean bias (NMB) for hourly O₃ and daily 8-
26 hr maximum O₃ concentrations should ideally fall within ±15% and ±10%, respectively, to meet the “goal”
27 benchmark. If the “criteria” benchmarks are to be met, the NMB should be within ±30% and ±20%, respectively.
28 Moreover, uncertainties in O₃ predictions due to uncertainties in various model inputs were quantified using the
29 decoupled direct method (DDM) in a commonly used CTM. For the simulation period of June 2021, the total
30 uncertainty of simulated O₃ ranged 4-25 µg/m³, with anthropogenic volatile organic compound (AVOC)
31 emissions contributing most to the uncertainty of O₃ in coastal regions and O₃ boundary conditions playing a
32 dominant role in the northwest region. The proposed benchmarks for assessing simulated O₃ concentrations, in
33 conjunction with our previous studies on PM_{2.5} and other criteria air pollutants, represent a comprehensive and
34 systematic effort to establish a model performance framework for CTM applications in China. These benchmarks
35 aim to support the growing modeling community in China by offering a robust set of evaluation metrics and
36 establishing a consistent evaluation methodology relative to the body of prior research, thereby helping to
37 establish the credibility and reliability of their CTM applications. These statistical benchmarks need to be
38 periodically updated as models advance and better inputs become available in the future.

39 **Keywords:** Ozone, chemical transport model, statistical benchmark, uncertainty analysis, China



40 **1 Introduction**

41 Tropospheric ozone (O_3) is a secondary air pollutant generated by complicated photochemical reactions involving
42 nitrogen oxides (NO_x) and volatile organic compounds (VOC) (Seinfeld and Pandis, 2016). Ozone has negative
43 impacts on human health (GBD, 2021), vegetation and ecosystem productions (Ainsworth et al., 2012). Due to
44 rapid economic development and fast industrialization and urbanization over the past several decades, China has
45 experienced heavy haze pollution in winter and severe O_3 pollution in summer, the latter extending into the late-
46 winter haze season (Li et al., 2021). Despite efforts to reduce fine particulate matter ($PM_{2.5}$) and heavy haze days
47 (Wang et al., 2022; Bai et al., 2019; Chu et al., 2020), ground-level O_3 concentrations have continued to increase
48 in recent years (Dang and Liao, 2019; Li et al., 2019; Liu et al., 2019a; Lu et al., 2020; Wang et al., 2017; Yao et
49 al., 2023; Chen et al., 2023; Xu et al., 2023). The challenge in controlling O_3 pollution lies in the significant
50 influences of meteorological conditions on O_3 formation and its nonlinear relationship with precursors (Wang et
51 al., 2022b). In addition, O_3 pollution exhibits strong regional characteristics, necessitating regional-scale control
52 efforts (Yang et al., 2021a).

53 Application of chemical transport models (CTMs) has become increasingly popular in addressing O_3 -related
54 issues in China (Yang and Zhao, 2023), providing insights into the role of local emissions and regional transport
55 (Shen et al., 2022), sectoral contributions (Liu et al., 2020a), policy effectiveness (Liu et al., 2023b), and
56 predictions of future O_3 levels (Yang and Zhao, 2023). Ensuring the representativeness of CTM simulations is
57 crucial, and can benefit from establishing performance standards or benchmarks to help put CTM results in
58 context relative to the existing body of work. While other regions (e.g., the U.S. and Europe) have proposed
59 evaluation criteria for simulated O_3 (Emery et al., 2017), they may not be suitable for China. The increasing
60 prevalence of CTM applications in China necessitates specific CTM benchmarks tailored to this region.

61 This study aims to develop customized CTM benchmarks for O_3 simulations in China, building upon our prior
62 work that proposed evaluation indicators and benchmarks for simulating other criteria air pollutants (Huang et al.,
63 2021; Zhai et al., 2024). A thorough literature review was conducted on O_3 simulations using CTMs from 2006 to
64 2021. Detailed information regarding O_3 performance was extracted and analyzed to recommended model
65 performance evaluation (MPE) metrics and to propose benchmarks tailored to China. Furthermore, uncertainties
66 in O_3 predictions due to various model inputs were quantified using the decoupled direct method of sensitivity
67 analysis (DDM, Cohan and Napelenok, 2011) in a commonly used CTM. The structure of this study is as follows:
68 Section 2 outlines the data source and methodology utilized. Section 3 describes the current status of O_3
69 simulation studies in China and proposes recommended evaluation metrics and associated benchmarks. Section 4
70 delves into discussions on O_3 uncertainties arising from different model inputs and conclusions are given in
71 Section 5.

72 **2 Methodology**

73 **2.1 Data collection**

74 The methodology for data compilation was consistent with our prior studies for other criteria pollutants (Huang et
75 al., 2021; Zhai et al., 2024) and is briefly described here. We considered published O_3 simulations using five
76 CTMs: the Community Multiscale Air Quality (CMAQ, <https://www.epa.gov/cmaq>, accessed on 2024-07-12)
77 model, the Comprehensive Air quality Model with extensions (CAMx, <https://camx.com>, accessed on 2024-07-



78 12), the Goddard Earth Observing System coupled with chemistry (GEOS-Chem, <https://geoschem.github.io>,
79 accessed on 2024-07-12), the Weather Research and Forecasting model coupled with Chemistry (WRF-Chem,
80 <https://www2.acom.ucar.edu/wrf-chem>, accessed on 2024-07-12), and the Nested Air Quality Prediction
81 Modeling System (NAQPMS) (Wang et al., 2014; Ge et al., 2014). We gathered relevant publications using a
82 combination of three keywords in Web of Science: “O₃”, the models’ names (one of the five models), and
83 “China”, with a time range between 2006 and 2021. This process identified a total of 667 records (250 studies for
84 CMAQ, 186 for WRF-Chem, 163 for GEOS-Chem, 36 for CAMx, and 32 for NAQPMS), with subsequent
85 refinement steps to exclude duplicates, non-English publications, conference papers, and journals unrelated to air
86 quality. Through manual selection, which involved identifying studies that provide extractable results (i.e.,
87 studies offering explicit results from model performance evaluations), a final set of 216 studies was chosen for
88 detailed analysis (see Table S1 for a complete list of publications).

89 Detailed information regarding model configurations (e.g., modeling period, spatial resolution, gas-phase
90 chemistry, initial/boundary conditions) and results of 23 MPE metrics (Table S2) were extracted and compiled
91 from those 216 studies. For consistency, we converted O₃ concentrations reported in parts per billion by volume
92 (ppbv) to µg/m³ using a factor of 2.14 (equivalent to 273.15 K at 101.325 kPa) for consistency. Ten regions in
93 China (Table S3), including the Beijing–Tianjin–Hebei (BTH) region, Yangtze River Delta (YRD) region, Pearl
94 River Delta (PRD) region, Sichuan Basin (SCB), North China Plain (NCP), and five other regions (Figure 1),
95 were identified for further analysis.

96 **2.2 Recommended benchmarks for O₃**

97 Among the 23 collected MPE metrics, we derived recommended benchmarks for the six most frequently used
98 metrics (see Table S4 for definitions): mean bias (MB), normalized mean bias (NMB), root mean square error
99 (RMSE), normalized mean error (NME), correlation coefficient (R), and index of agreement (IOA). The
100 derivation of benchmarks follows previous studies by Simon et al. (2012) and Emery et al. (2017). Briefly, each
101 metric’s rank-ordered (from best to worst, for instance, from 1 to 0 for R) distribution was generated to identify
102 the values at the 33rd and 67th percentiles. As highlighted in Emery et al. (2017), these percentiles serve to
103 categorize the entire distribution into three performance categories: studies falling within the 33rd percentile (the
104 “goal”) attain the best performance that current models can be expected to achieve, those between the 33rd and
105 67th percentiles (the “criteria”) attain typical performance achieved by the majority of modeling studies, while
106 those beyond the 67th percentile indicate relatively poor performance for the particular metric under consideration.
107 We present the benchmarks for hourly O₃, maximum daily 8-hr average O₃ (8-hr max O₃), and daily maximum 1-
108 hr O₃ (1-hr max O₃), depending on data availability.

109 **2.3 Uncertain analysis of O₃ simulation**

110 In addition to developing the MPE benchmarks for simulated ozone, we further quantified uncertainties in
111 predicted ozone concentrations using one of the five models (i.e., CMAQ). The CMAQ version 5.3.2
112 (<https://www.epa.gov/cmaq>, accessed on April 17, 2024) was employed to simulate O₃ during June 2021 in
113 China. Base model configurations are the same as our previous study (Sun et al., 2024) and are briefly described
114 here. The modeling domain covers the entirety of China and adjacent Asian regions (Figure 1) with a spatial
115 resolution of 36 km × 36 km grid and 23 vertical layers. Meteorological fields are simulated using the Weather



116 Research and Forecasting model (WRF version 4.0). CB6 and AERO7 were chosen as the gas-phase and aerosol
117 mechanisms, respectively. Emissions data include the 2019 Multi-resolution Emission Inventory for China
118 (MEIC-2019) (<http://www.meicmodel.org>, accessed on June 23, 2022) and the 2010 Emissions Database for
119 Global Atmospheric Research (EDGAR, <http://www.meicmodel.org>, accessed on June 23, 2022). Natural
120 emissions were generated based on the Model of Emissions of Gases and Aerosols from Nature (MEGAN
121 version 3.1, <https://bai.ess.uci.edu/megan>, accessed on June 23, 2022). The CMAQ default O₃ profile (with a
122 uniform O₃ concentration of 29 ppb) was used as the initial and boundary conditions (BCs). A 10-day spin-up run
123 was conducted to mitigate the influence of initial conditions.

124 We followed Dunker et al. (2020) to quantify the uncertainties of predicted O₃ concentrations due to six model
125 inputs: anthropogenic NO_x (ANO_x) and VOC (AVOC) emissions for China, biogenic VOC (BVOC) and soil
126 NO_x (SNO_x) within China; dry deposition velocities for O₃; and BCs for O₃. The uncertainties associated with
127 each of the inputs (Table S5) are based on previous studies addressing emission uncertainties (Cheng et al., 2019),
128 deposition velocities, and BCs (Beddows et al., 2017; Derwent et al., 2018). Like Dunker et al. (2020), these
129 uncertainties were considered independent and lognormally distributed. The CMAQ decoupled direct method
130 (DDM) was used to generate first-order sensitivities of O₃ to each of the inputs (excluding dry deposition). For
131 dry deposition, we conducted two parallel simulations in which the O₃ dry deposition velocities were manually
132 changed by ±10%, and the changes in simulated O₃ concentrations were treated as the O₃ sensitivities to dry
133 deposition velocity:

$$S_{DEP}^{(1)} = \frac{C_{1.1dep_O_3} - C_{0.9dep_O_3}}{2} * 10 \quad \text{Eq. (1)}$$

134 where $S_{DEP}^{(1)}$ is the O₃ sensitivity to dry deposition velocities, and $C_{1.1dep_O_3}$ and $C_{0.9dep_O_3}$ represent the simulated
135 O₃ concentrations as dry deposition velocities are increased and decreased by 10%, respectively. The sensitivities
136 obtained were then combined with their respective uncertainties, enabling us to quantify the contributions to the
137 variance in O₃ concentrations. For example, the O₃ uncertainties due to dry deposition are calculated as:

$$\text{un}(DEP) = \text{var}(DEP) = \left[\frac{\ln(f_{DEP})}{2} * S_{DEP}^{(1)} \right]^2 \quad \text{Eq. (2)}$$

138 where un(DEP) represents the uncertainty of O₃ due to dry deposition at 1σ, and f_{DEP} (=2 from Table S5) is the
139 uncertainty factor for dry deposition and follows an assumption of a lognormal distribution.

140 The contribution of dry deposition to the total uncertainty in O₃ is calculated as follows:

$$\% DEP = \frac{\text{var}(DEP)}{\text{var}(ANO_x) + \text{var}(AVOCs) + \text{var}(BNO_x) + \text{var}(BVOCs) + \text{var}(DEP) + \text{var}(BCs)} \quad \text{Eq. (3)}$$

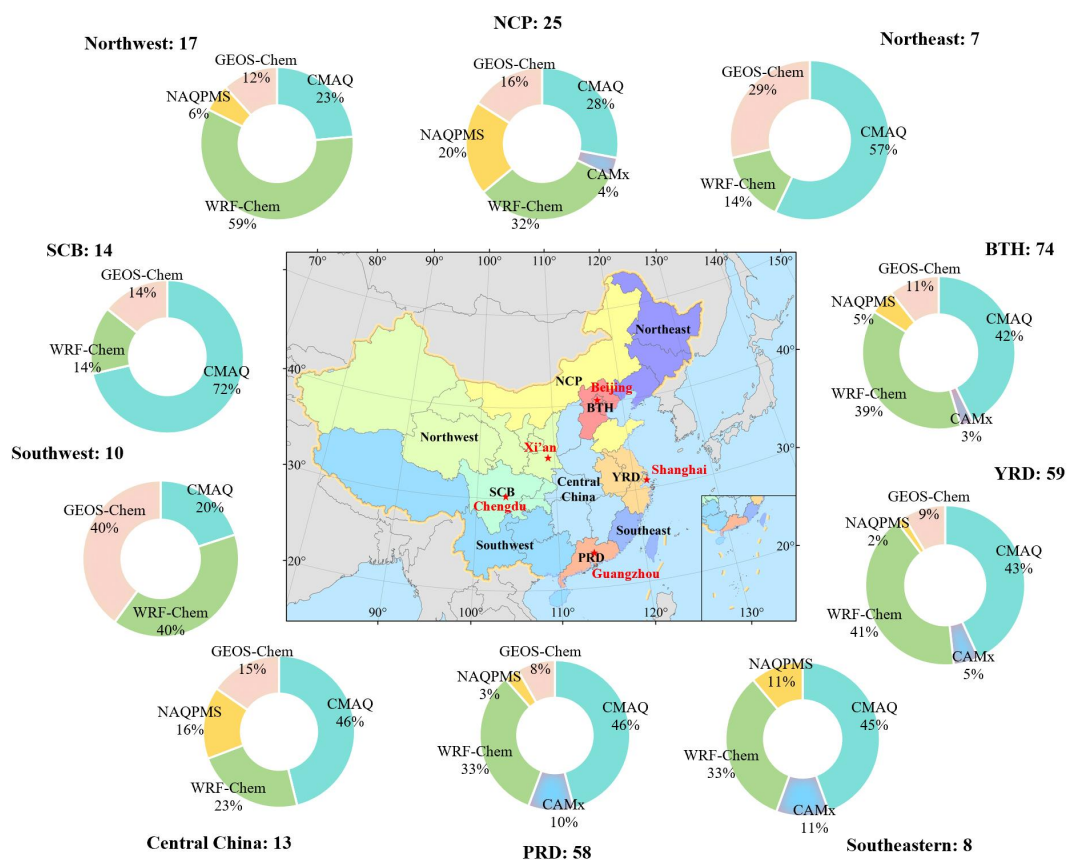
141 3. Results and discussions

142 3.1 General overview of O₃ simulation studies in China

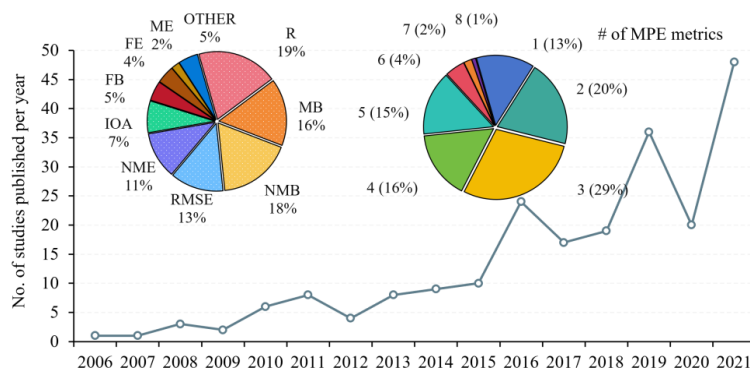
143 In the last decade, there has been a significant increase in research focusing on O₃ in China, as illustrated in
144 Figure 2. The issuance of the Three-Year Action Plan to Win the Blue Sky Defense Battle in 2017
145 (http://www.gov.cn/zhengce/content/2018-07/03/content_5303158.htm, accessed on April 15, 2024) led to a
146 further surge in studies related to O₃, with a noticeable decline in 2020 possibly attributed to the impact of the
147 COVID-19 pandemic. In 2021, there were 48 studies dedicated to addressing O₃-related issues using CTMs,



148 marking a six-fold increase compared to 2011. Similar to PM_{2.5}, BTH (74 studies), YRD (59 studies), and PRD
 149 (58 studies) emerged as the top three most studied regions. Among the various CTMs employed, CMAQ stood
 150 out as the most commonly utilized model (90 studies), followed by WRF-Chem (84 studies). The application of
 151 CAMx (14 studies) and NAQPMS (8 studies) was relatively less frequent. In terms of MPE metrics, R had the
 152 highest frequency of occurrence at 19%, followed by NMB (18%), MB (16%), RMSE (13%), and NME (11%).
 153 Nearly half of the studies incorporated 2 or 3 metrics for evaluating O₃, while less than 7% assessed at least five
 154 different metrics. The three most common types of O₃ concentrations evaluated were hourly O₃ concentration, the
 155 maximum daily 8-hour average O₃ (8-hr max O₃), and the daily maximum 1-hour O₃ (1-hr max O₃). Among all
 156 the articles examined, 77% focused on evaluating hourly O₃, 16% on 8-hr max O₃, and 7% on 1-hr max O₃.



157
 158 **Figure 1** CMAQ modeling domain with definitions of regions used in this study. The surrounding pie charts
 159 display the total number of studies for each region (excluding studies for the entire China) and the percentage of
 160 different CTMs used. Red stars represent the five cities selected in uncertainty analysis.



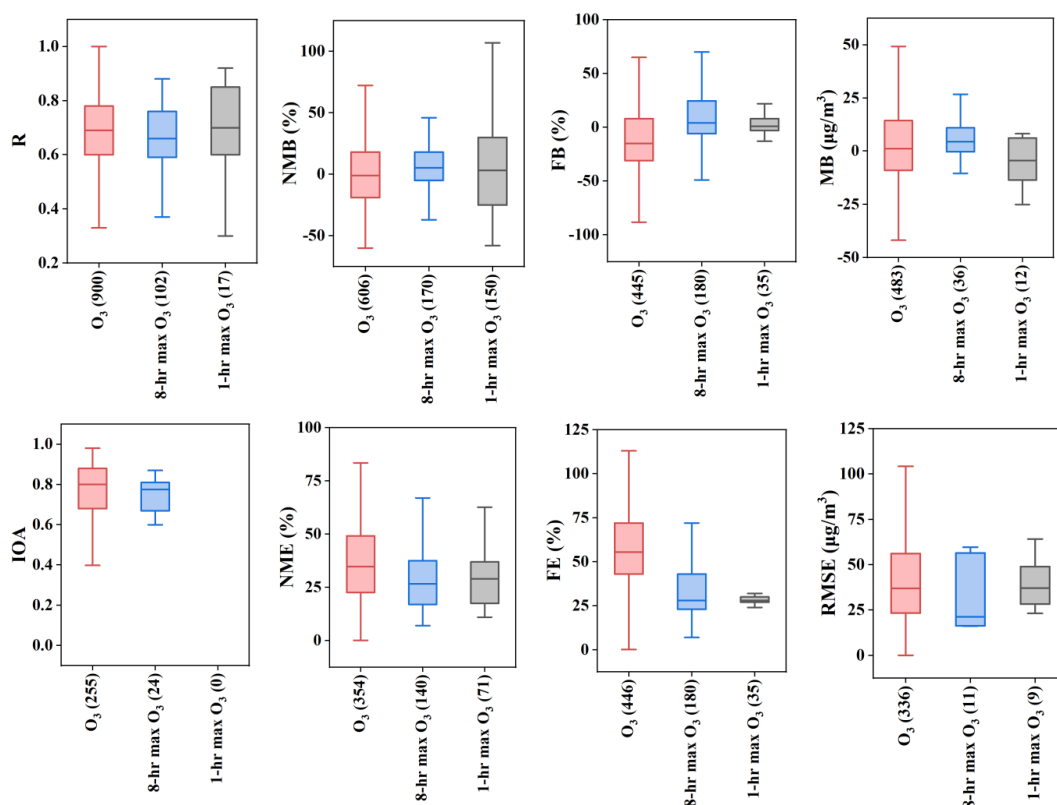
161

162 **Figure 2** Number of O₃ studies published during 2006–2021. Pie charts show the frequency of different MPE
 163 metrics (left) and the number of metrics used in one study (right).

164 **3.2 Quantile distributions of O₃ MPE results**

165 Figure 3 shows the quantile distributions of various evaluation metrics collected in this study. The results are
 166 presented for different types of O₃ concentrations: hourly O₃, 1-hr max O₃, and 8-hr max O₃, whenever data is
 167 available. Previous studies have shown that using maximum O₃ values (i.e. 1-hr max and 8-hr max) instead of
 168 hourly O₃ can lead to differing results within the same study (e.g., Ni et al., 2020; Li et al., 2016). Peak O₃
 169 concentrations typically occur between 12:00 and 18:00. For example in Ni et al. (2018), 8-hr max O₃ showed an
 170 overestimation tendency compared to average hourly O₃, but in another study (Yang et al., 2021b), there was an
 171 opposite trend. Underestimation of peak O₃ concentrations might be offset by overestimation during non-peak
 172 hours and vice versa. Therefore, achieving satisfactory performance in daily averaged O₃ levels does not
 173 necessarily indicate the model's ability to accurately capture high O₃ concentrations.

174 Hourly O₃ exhibited equivalent overestimation and underestimation in terms of MB and NMB, with MB ranging
 175 from as low as -40 µg/m³ to nearly 50 µg/m³ and NMB ranging from less than -50% to more than 70%. However,
 176 fractional bias (FB) indicated more underestimated than overestimated hourly O₃ concentrations. For all three
 177 bias metrics, 8-hr max O₃ exhibited more overestimation than underestimation, suggesting a tendency for models
 178 to overestimate off-peak hours. For 1-hr max O₃, both NMB and FB displayed equivalent overestimation and
 179 underestimation, with NM showing a wider range than FB, likely due to fewer data points. For error metrics, 8-hr
 180 max and 1-hr max O₃ generally performed better than hourly O₃. For instance, the median values of NME were
 181 34.8%, 26.6%, and 29% for hourly O₃, 8-hr max, and 1-hr max O₃, respectively. R and IOA indicate how well the
 182 model captures observed variations, either temporally or spatially. The use of IOA was significantly less than R
 183 and no studies reported IOA values for 1-hr max O₃. For the other two O₃ types, IOA values (median value of 0.8
 184 for O₃ and 0.77 for 8-hr max O₃) were generally higher than R (median value of 0.69 for O₃ and 0.66 for 8-hr
 185 max O₃). Six studies reported both R and IOA values, of which four (Liu and Wang, 2020; Wang et al., 2019; Liu
 186 et al., 2019b; Gao et al., 2017) reported higher IOA values than R.



187 **Figure 3** Quantile distribution of common O₃ performance indicators

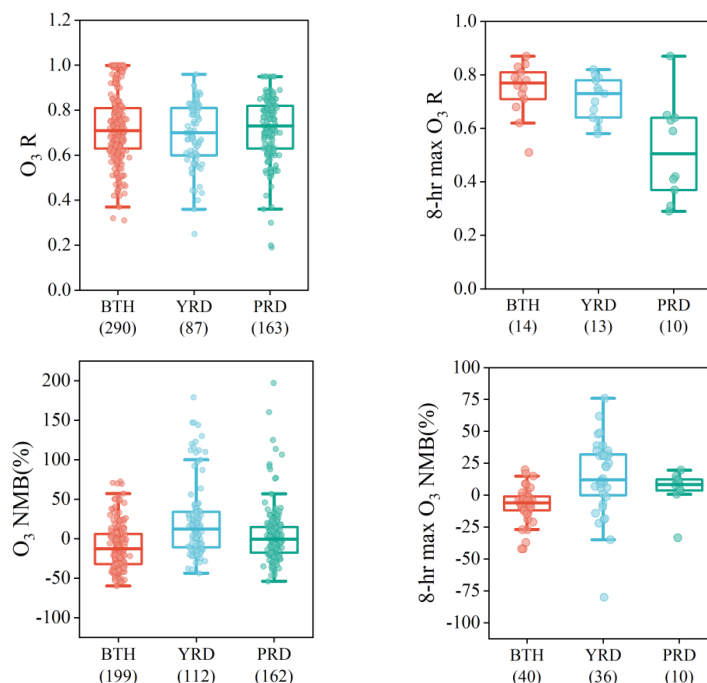
188 **Regional and seasonal differences**

189 Like our previous studies (Huang et al. 2021; Zhai et al. 2024), we discuss the influences of various key factors
 190 on model performance in simulating O₃ concentrations. We first considered whether there were discernible
 191 regional or seasonal differences. Figure 4 presents the distribution of R and NMB values grouped by three key
 192 regions in China: BTH, YRD, and PRD. These regions are the most densely populated and economically
 193 developed urban clusters in China. In terms of hourly O₃, the R values across the three regions display similarity,
 194 with median values around 0.7. For 8-hr max O₃, however, PRD stands out with notably lower R values
 195 compared to BTH and YRD. Regarding NMB values, BTH tends to have more underestimation, while the YRD
 196 and PRD lean towards overestimation. Over the past decade, BTH has consistently recorded the highest O₃ levels
 197 and number of O₃ pollution days among the three regions (Wang et al., 2024). The variations in NMB values
 198 among regions suggest a trend of current models underestimating O₃ levels in areas with more severe O₃
 199 pollution.

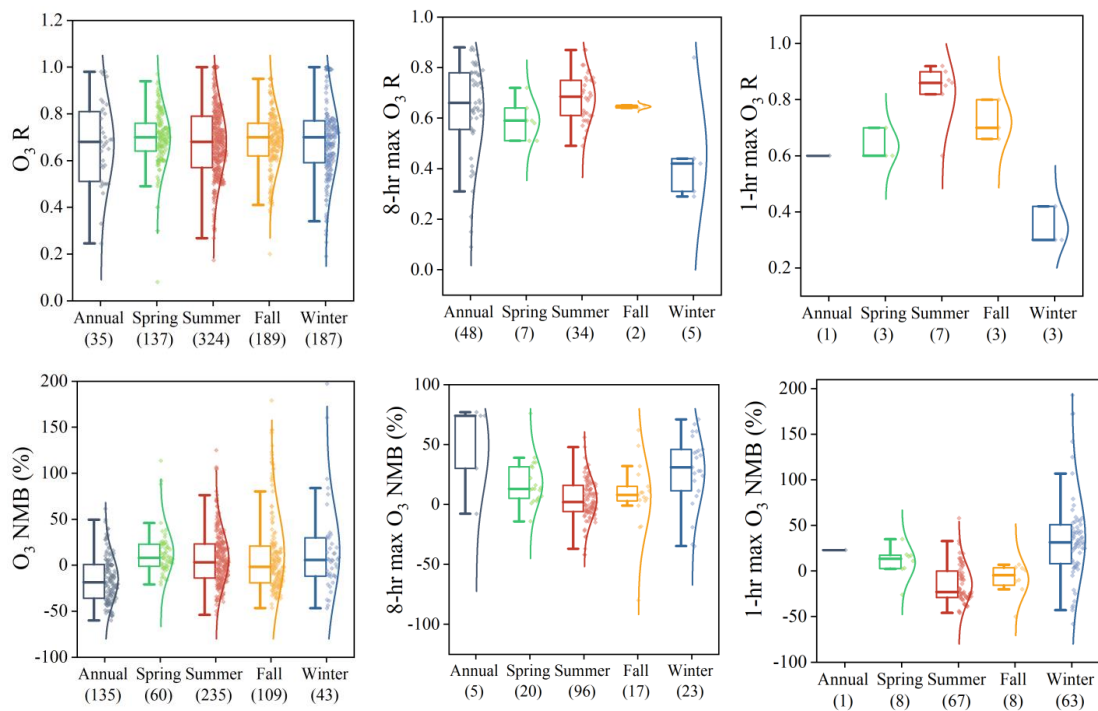
200 In terms of the seasonal variations (Figure 5), the NMB values of hourly O₃ concentrations exhibit similar
 201 patterns across different seasons, showing equivalent overestimation and underestimations. However, when
 202 assessed over the entire year, hourly O₃ concentrations tend to be largely underestimated. The seasonal patterns
 203 of NMB distributions are similar for 8-hr and 1-hr max O₃, with summer O₃ concentrations being more frequently
 204 underestimated compared to other seasons. For instance, in the case of 1-hr max O₃, peak O₃ concentrations are



205 predominantly underestimated (with a median NMB of -23%) while they are overestimated in winter (with a
 206 median NMB of 31.5%).



207 **Figure 4** Quantile distribution of R and NMB of O₃ in BTH, YRD, and PRD



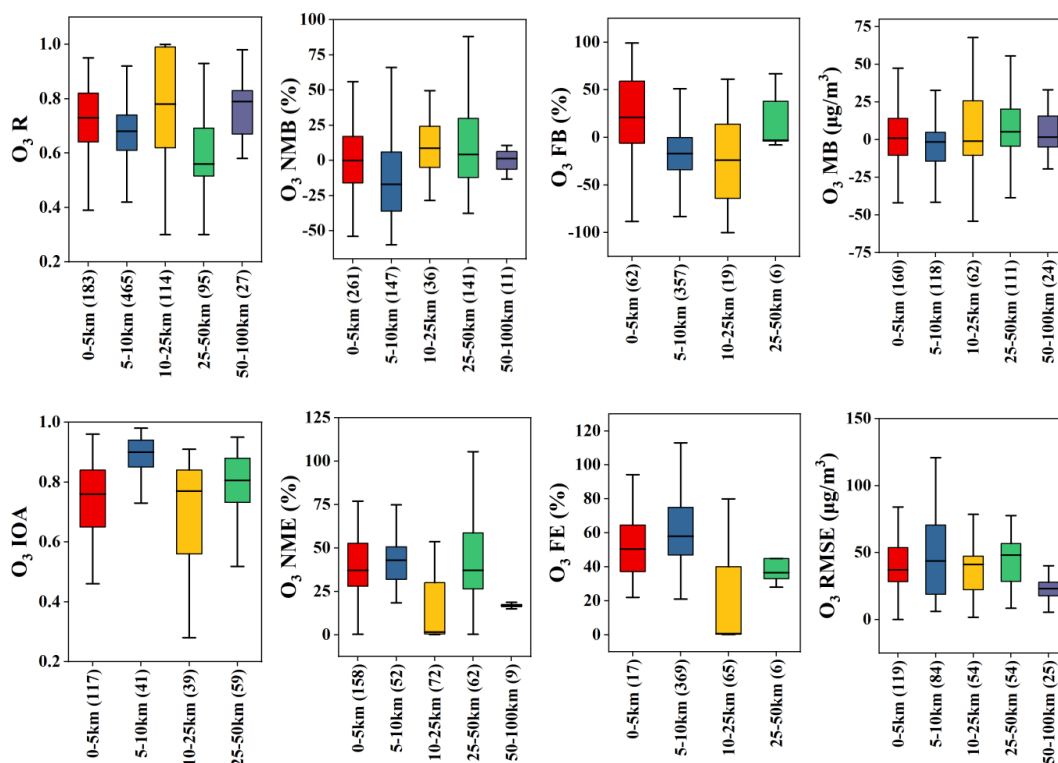
208 **Figure 5** Quantile distribution of O₃ NMB values in different seasons



209 ***Impact of grid spacing***

210 The selection of grid spacing for a CTM application depends on several factors, such as the objective of the study,
211 the geographical scope of the study area, the availability of input data, etc. Generally, a coarse grid spacing (> 50
212 km) is utilized for global simulations (i.e. GEOS-Chem), while a finer grid spacing (< 4km) with nested grids is
213 preferred for regional or city-scale modelling. Coarser grid spacing may result in multiple monitoring stations
214 falling within a single grid cell, potentially smoothing out extreme values observed at specific locations. Among
215 the 216 studies reviewed, 29 different grid resolutions (based on the resolution of the innermost domain) were
216 identified, ranging from 1 km to 200 km. The resolutions were classified into five groups in this study: < 5 km, 5-
217 10 km, 10-25 km, 25-50 km, and 50-100 km (resolutions over 100 km were excluded from the analysis due to
218 limited data points). Figure 6 shows the distribution of eight statistical indicators by different resolutions. Overall,
219 no clear trend was evident to indicate better model performances as grid spacing decreases. For example, the
220 median R value is 0.73 for < 5 km group, surpassing the 5-10 km and 25-50 km groups but falling below the 10-
221 25 km and 50-100 km groups. Studies conducted with a grid spacing of 10-25 km exhibit the best model
222 performance in terms of NME and FE distributions compared to other groups. While most studies assess models
223 within a single domain (usually the innermost domain with the finest resolution), a few studies have conducted
224 multi-domain analyses, where finer spatial resolutions generally have superior results compared to coarse
225 resolutions. Liu et al. (2020b) used WRF-CMAQ to analyze O₃ prediction and health exposure at different spatial
226 resolutions (1, 4, 12, and 36 km). The results showed more than 20% difference in premature mortality due to
227 different model resolutions being used. Nevertheless, reducing grid spacing does not necessarily lead to improved
228 model performance if the input data resolution (i.e., spatial resolution of the emissions) is not correspondingly
229 high or well-matched.

230



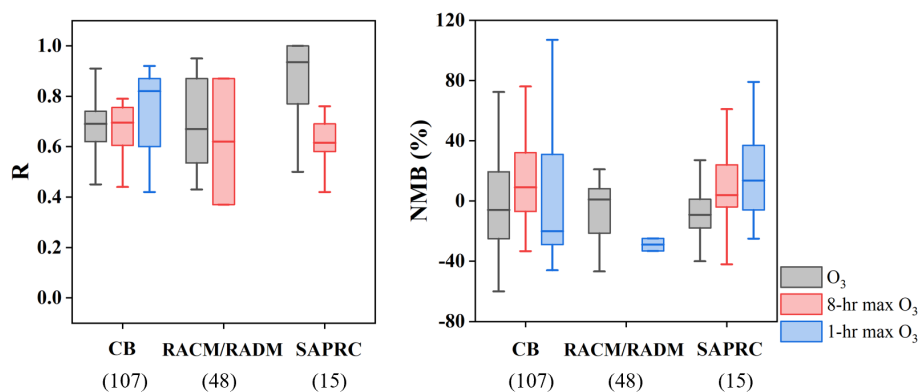
231 **Figure 6** Quantile distribution of O₃ with respect to commonly used assessment indicators at different spatial
 232 resolutions

233 **Choice of gas-phase chemical mechanism**

234 Gas-phase chemical mechanisms play a crucial role in the accurate prediction of atmospheric composition using
 235 CTMs. Some of the commonly used mechanisms include the Carbon Bond mechanism (CB) (Yarwood et al. 1997; Luecken et al., 2019; Appel et al., 2021;
 236 Yarwood and Tuite, 2024), the Statewide Air Pollution Researcher Center (SAPRC) mechanism (Carter, 1996;
 237 Chang et al., 1999; Carter, 2000; Carter, 2010), and the Regional Atmospheric Chemistry Mechanism (RACM)
 238 (Stockwell et al., 1997; Goliff et al., 2013). These mechanisms have undergone rigorous evaluations against
 239 experimental data, showcasing reliable predictive capabilities for O₃ in diverse atmospheric environments. The
 240 CB mechanism is a condensed mechanism in which the carbon bond is treated as a reaction unit, and the carbon
 241 bonds with the same bonding state are treated as a group (Cao et al., 2021). The latest version, CB7, contains 91
 242 gaseous species and 230 reactions ([https://www.tceq.texas.gov/downloads/air-
 243 quality/research/reports/photochemical](https://www.tceq.texas.gov/downloads/air-quality/research/reports/photochemical), accessed on 2024-06-18). In contrast, the SAPRC mechanism categorizes
 244 species based on their reactivity with OH (Carter et al., 2010). The most recent SAPRC22 mechanism includes
 245 162 species and 738 reactions. RACM was developed based on Regional Acid Deposition Model (RADM),
 246 which is an inductive mechanism for treating hydrocarbons with fixed parameterization method and is carried out
 247 according to the reaction rate and activity of different pollutants with ·OH. Compared to the other two
 248 mechanisms, RACM and RACM2 contain detailed chemical processes of radicals, biogenic VOC and less-
 249 reactive VOC able to survive during long distance transport. 119 reactive species and 363 reactions were
 250



251 included in RACM2 describing the oxidation reactions of 21 types of primary VOC in the system (Liu et al.,
252 2023a).
253 Among the 216 studies compiled, nearly half of them used CB mechanism for simulations, approximately a
254 quarter employed RACM/RADM, and only 15 studies utilized SAPRC. Figure 7 compares the distribution of R
255 and NMB grouped by different gas-phase mechanism. In terms of R values, CB tends to perform slightly better
256 than RACM/RADM, with SARPC showing the highest R median value (0.93) for hourly O₃ but the lowest for 8-
257 hr max O₃ among the three mechanisms. Regarding NMB, SAPRC tends to overestimate peak O₃ values
258 compared to the other mechanisms, particularly for 1-hr max O₃, a trend observed in previous studies (Qiao et al.,
259 2019).

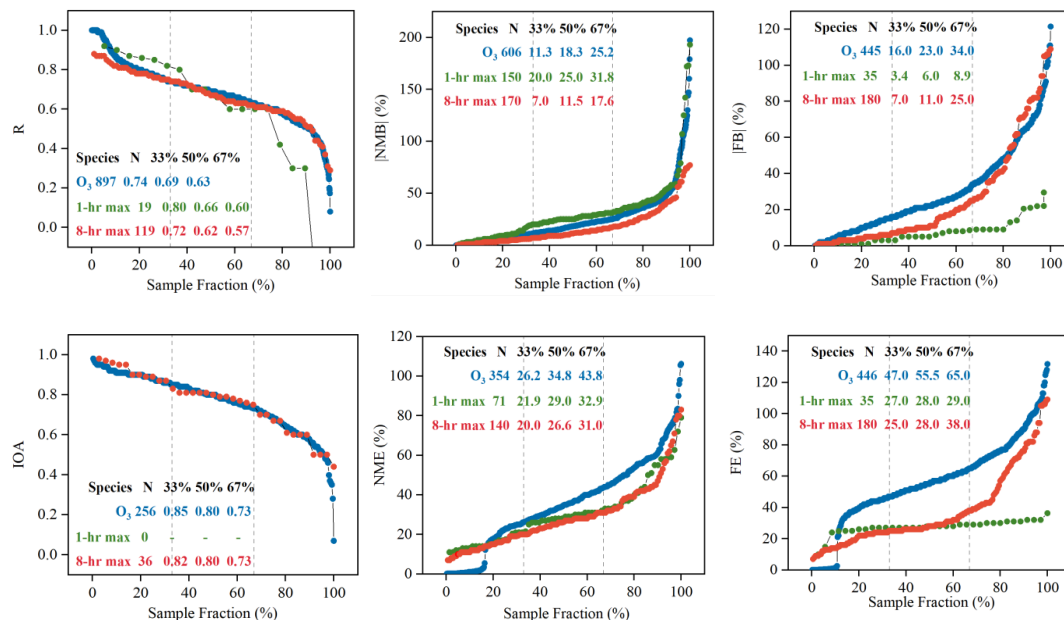


260
261

Figure 7 Quantile distributions of R and NMB by gas-phase chemical mechanism

262 3.3 Recommended benchmarks for O₃ MPE

263 Figure 8 illustrates the ranked distributions of various statistical indicators, including R, IOA, NMB, NME, FB,
264 and FE for hourly O₃, 1-hr max O₃, and 8-hr max O₃. The absolute values of NMB and FB are presented to
265 indicate deviations from zero. In terms of R and IOA, the ranked distributions for hourly O₃ and 8-hr max O₃ are
266 quite similar, with R values ranging from around 0.72 at the 33rd percentile to 0.60 at the 67th percentile. The
267 corresponding IOA values are slightly higher, ranging from ~0.83 at the 33rd percentile to ~0.73 at the 67th
268 percentile. For 1-hr max O₃, the limited number of data points (less than 20) resulted in an R value of 0.80 at the
269 33rd percentile and 0.60 at the 67th percentile, while the IOA distribution was not available due to missing data.
270 For NMB and NME, the results for 8-hr max O₃ show the lowest values, indicating that models perform better in
271 capturing the 8-hr max O₃ concentrations. The 33rd percentile of absolute NMB for 8-hr max O₃ is less than 10%,
272 and the 67th percentile is below 20%. In terms of FB and FE, the ranked distributions for 1-hr max O₃ are flatter
273 compared to the other two O₃ types, likely due to the smaller number of available data points. For both metrics,
274 the 8-hr max O₃ exhibits lower values than O₃. At the 33rd percentile, the absolute FB (FE) is less than 10% (25%)
275 for 8-hr max O₃ and less than 20% (50%) for O₃. At the 67th percentile, the absolute FB (FE) is 25% (38%) for 8-
276 hr max O₃ and 34% (65%) for O₃. In addition, we provide a more detailed ranked distribution in Table S6.



277 **Figure 8** Rank-ordered distributions of R, IOA, NMB, NME, FB, and FE for O₃, 1-hr max O₃ and 8-hr max O₃
 278 speciated components. The number of data points and the 33rd, 50th, and 67th percentile values are also listed.

279 Following Emery et al. (2017) and Huang et al. (2021), we propose recommended statistical indicators and
 280 corresponding benchmarks for evaluating O₃, as detailed in Table 1. The goal values, corresponding to the
 281 threshold at the 33rd percentile, represent the optimal model performance anticipated from current models. The
 282 criteria values, reflecting the threshold at the 67th percentile, represent the performance levels achieved by the
 283 majority of studies. Due to limited data availability, the derivation of benchmarks for certain metrics concerning
 284 1-hr max O₃ remains uncertain. In such cases, benchmarks for IOA and R for hourly O₃ were directly adopted
 285 due to minimal variations among different O₃ types. Similarly, benchmarks proposed for 8-hr max O₃ were
 286 applied to 1-hr max O₃ for FB and FE, given their closer distributions. Our findings indicate that benchmarks
 287 tend to be more stringent for 8-hr max O₃ compared to the other two types, with the exception of IOA where they
 288 remain the same. Based on our results, a value of R greater than 0.70 and 0.55 would meet the goal and criteria
 289 benchmark for 8-hr max O₃. Correspondingly, the goal and criteria values for NMB are 10% and 20%.

290 In contrast to Emery et al. (2017), we provide separate benchmarks for O₃, 8-hr max O₃, and 1-hr max O₃. Emery
 291 et al. (2017) found rather similar results between hourly and 8-hr max O₃ in the U.S and so recommended a single
 292 set of benchmarks for ozone. Out of the 216 studies analyzed, 15 studies evaluated at least two O₃ types. The use
 293 of cutoff for evaluating O₃ is extremely limited in China (only 5 studies applied cutoffs), thereby precluding any
 294 specific recommendation on cutoff values. In addition to the benchmarks for NMB, NME, and R provided by
 295 Emery et al. (2017), we have introduced benchmarks for IOA, FB, and FE, backed by a sufficient number of data
 296 points. The few values marked with an asterisk in Table 1 indicate that our benchmarks are more stringent than
 297 the corresponding values in Emery et al. (2017), implying that achieving our recommended 33rd (or 67th)
 298 percentiles may pose greater challenges.



299 Overall, however, our proposed benchmarks are more lenient than those of Emery et al. (2017), particularly in the
 300 context of hourly O₃. For NME, our suggested goal and criteria for O₃ stand at 30% and 45%, respectively, nearly
 301 double the figures reported by Emery et al. (2017), which recommend 15% for the goal and 25% for the criteria.
 302 The criteria value for R is an exception where our proposed value (0.55 for 8-hr max O₃ and 0.60 for O₃) is
 303 higher than 0.50 in Emery et al. (2017).

304 **Table 1** Recommended benchmarks for evaluating simulated O₃ by CTM applications in China

Metrics	Benchmark level	O ₃	8-hr max O ₃	1-hr max O ₃	Emery et al. (2017)
					1-hr max O ₃ and 8-hr max O ₃
R	Goal	> 0.70	> 0.70	> 0.80*	> 0.75
	Criteria	> 0.60*	> 0.55*	> 0.60*	> 0.50
NMB	Goal	< ±15%	< ±10%	NA	< ±5%
	Criteria	< ±30%	< ±20%	NA	< ±15%
NME	Goal	< 30%	< 20%	< ±20%	< ±15%
	Criteria	< 45%	< 35%	< ±35%	< ±25%
IOA	Goal	> 0.80	> 0.80	< 25%	NA
	Criteria	> 0.70	> 0.70	< 35%	NA
FB	Goal	< ±20%	< ±10%	< ±5%	NA
	Criteria	< ±35%	< ±30%	< ±10%	NA
FE	Goal	< 50%	< 25%	< 25%	NA
	Criteria	< 65%	< 40%	< 30%	NA

305 Note. (1) See descriptions in the main text for bold values. (2) Values with an asterisk indicate that our
 306 benchmarks are stricter than the corresponding values in Emery et al. (2017).

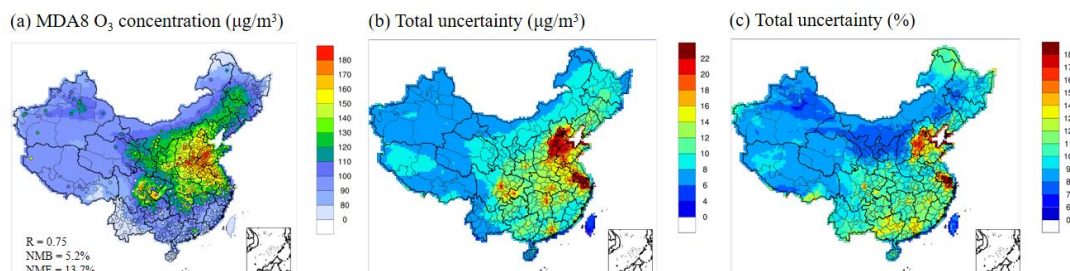
307 4. Uncertainty analysis of O₃ simulation using CMAQ

308 In order to further investigate the uncertainties in simulated O₃ concentrations simulated by CTMs, a base model
 309 simulation was conducted using CMAQ (the most frequently used CTM in China) for June 2021, a typical month
 310 with elevated O₃ in northern and eastern China. The uncertainties due to six model inputs were quantified for this
 311 case: VOC and NO_x emissions in China, differentiation between anthropogenic and biogenic sources, O₃ dry
 312 deposition velocities, and boundary conditions (BCs). The evaluation of the base model results indicates
 313 generally acceptable simulated MDA8 O₃ concentrations when compared to the observations. The results showed
 314 an overall MB of 6.1 µg/m³ and NMB of 5.2% (Figure 9). O₃ underestimation is observed over the BTH region,
 315 while overestimation occurs over the Sichuan Basin. The values of NMB, NME and R meet the goal benchmark
 316 we proposed above.

317 As displayed in Figure 10, the first-order sensitivity of MDA8 O₃ to the six model inputs exhibits substantial
 318 variations in spatial distributions and magnitudes. Higher sensitivity occurs in larger urban areas and is relatively
 319 low in rural areas. The sensitivity to VOC emissions is always positive (i.e., higher VOC leads to higher O₃),
 320 whereas the sensitivity to NO_x emissions could be both positive and negative. High O₃ sensitivity to AVOC
 321 emissions is observed for BTH, northern YRD, PRD, and major metropolitan areas (e.g., Chengdu in Sichuan

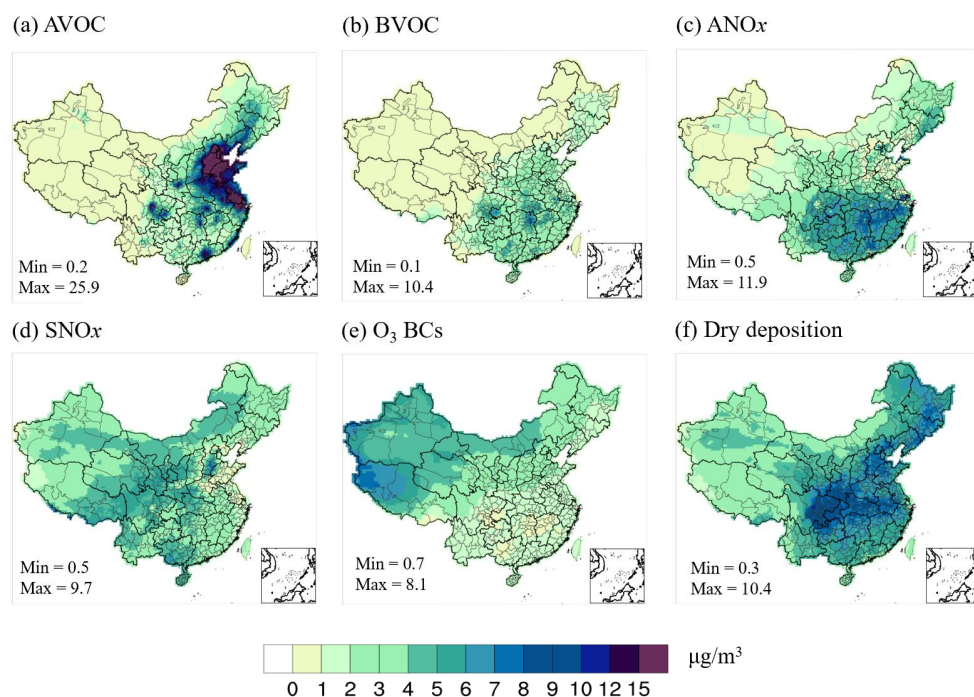


322 province, Xi'an in Shaanxi province), due to NO_x-rich and VOC-limited urban conditions. Conversely,
323 anthropogenic NO_x emissions resulted in negative O₃ sensitivity in the aforementioned regions and positive
324 sensitivity in others where rural conditions are more VOC-rich and NO_x-limited. The sensitivity to biogenic
325 precursor emissions (BVOC and SNO_x) was much lower compared to their anthropogenic counterparts. The
326 sensitivity to O₃ BCs predominantly extends towards the northwest (up to 50 μg/m³), where O₃ precursor
327 emissions are low. The sensitivity to O₃ dry deposition velocity exhibits a uniformly negative distribution (higher
328 deposition rates lead to lower ozone), with higher values in more vegetated areas and an average of -13.7 μg/m³.



329 **Figure 9** Spatial distributions of (a) MDA8 O₃ concentrations (ug/m³), (b) total uncertainties in ug/m³, and (c)
330 total uncertainty in percentage (%). Results are averaged for June 2021.
331

332 When the individual first-order sensitivity coefficient multiplies by the corresponding 1σ uncertainty (Table S5),
333 the contributions to the uncertainty in O₃ predictions can be obtained (Figure 10). Summing up all these
334 uncertainties yields the total uncertainty (Figure 9b). Large ozone uncertainties (> 20 μg/m³) were observed over
335 BTH, central YRD region, and major metropolitan areas (e.g. PRD, Chengdu in Sichuan province). Regions with
336 high uncertainties in O₃ predictions generally align with regions with poorer model performance. In BTH, YRD,
337 and PRD, the total ozone uncertainty due to the six model inputs ranges 11.7~31.8, 7.0~34.6 and 5.0~19.0 μg/m³,
338 respectively, corresponding to a relative percentage of O₃ concentration by 9.2~18.1%, 7.9~25.8%, and
339 7.6~14.6%. It should be noted that our uncertainty estimates represent conservative estimates because the effects
340 of uncertainties in the meteorological inputs and the uncertainties associated with the O₃ chemistry are not
341 included, the latter of which has been shown to have a comparable contribution to the total contributions from
342 emissions, dry deposition, and O₃ BC in the Dallas-Fort Worth region in the U.S. (Dunker et al. 2020).
343 Among the six model inputs, AVOC emissions make the largest contributions (exceeding 15 μg/m³) to the total
344 uncertainty in regions displaying high O₃ sensitivity, such as BTH, northern YRD, PRD, and several metropolitan
345 areas. The large uncertainties, stemming from both the high first-order sensitivities (Figure S1) and a relatively
346 high uncertainty factor (1.68), suggest that in these regions, uncertainties associated with AVOC emission
347 estimates would in more significant biases in simulated O₃ concentrations compared to other areas. O₃
348 uncertainties due to BVOC emissions, ranging 0.1~10.4 μg/m³, are mainly located in southern China, where
349 BVOC emissions are high. A similar spatial pattern is observed for uncertainties in ANO_x emissions, although its
350 contribution is larger (0.5~11.9 μg/m³). While the first-order O₃ sensitivity to SNO_x emissions is minimal (Figure
351 S1), the contribution to O₃ uncertainty from SNO_x emissions is noteworthy (0.5~9.7 μg/m³), given a large
352 uncertainty factor of 2 (Table S5). Uncertainty in O₃ BCs is relatively less important except in the northwest,
353 where it represents the largest contributing factor. Dry deposition serves as an important O₃ sink. Uncertainty
354 contribution from O₃ dry deposition velocities (0.3~10.4 μg/m³) is comparable to that of ANO_x emissions, with a
355 more evenly distributed spatial impact.



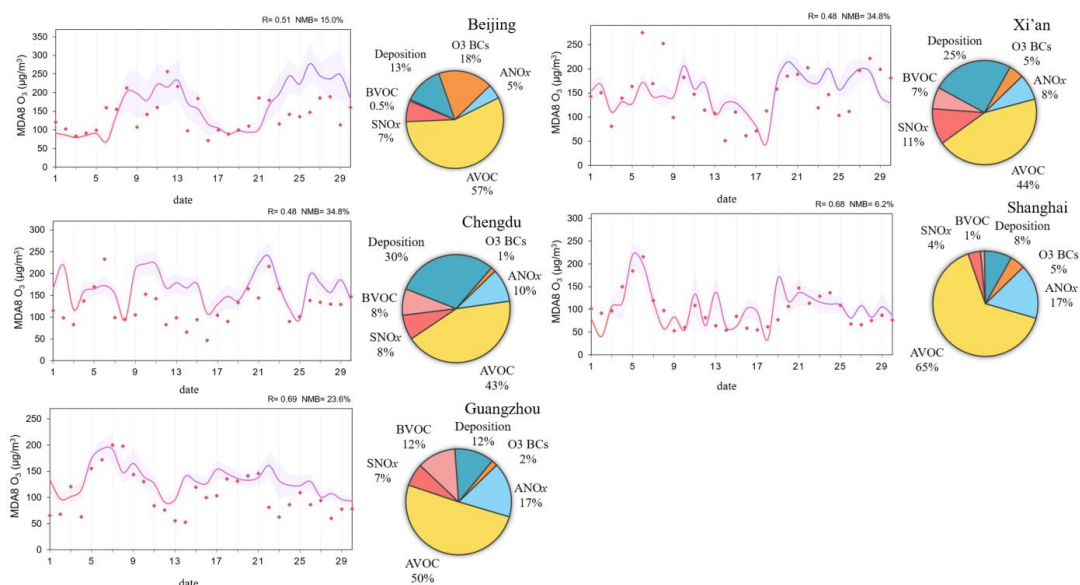
356

357

358 **Figure 10** Contributions to uncertainty in MDA8 O₃ simulation. Contribution of (a) AVOC, (b) BVOC, (c)
359 ANO_x, (d) SNO_x, (e) O₃ BCs, and (f) dry deposition in $\mu\text{g}/\text{m}^3$. Results are averages over all days in June 2021
360 and represent 1σ .

361 Figure 11 compares the observed MDA8 O₃ to the model results with their $\pm 1\sigma$ uncertainty range for five major
362 cities: Beijing, Shanghai, Guangzhou, Chengdu, and Xi'an. In Shanghai, the majority of the observed O₃ fall
363 within the $\pm 1\sigma$ uncertainty range. However, in Beijing, Chengdu, and to a lesser extent in Guangzhou, the model
364 tends to over-predict lower O₃ observations. In Xi'an, the model fails to capture the exceptionally high O₃
365 concentrations (MDA8 O₃ > 250 $\mu\text{g}/\text{m}^3$) on June 6th and 7th. Expanding the uncertainty limits to a $\pm 2\sigma$ range may
366 encompass some of the lower O₃ observations but the current uncertainty estimates do not fully account for all
367 the discrepancies between model results and observations. This discrepancy could be attributed to the coarse grid
368 resolution (36 km) used in this study, which may not adequately resolve the impact of local emission sources.
369 Furthermore, as mentioned earlier, uncertainties related to O₃ chemistry and meteorological inputs were not
370 accounted for and should be quantified in future work.

371 The relative contributions to the total uncertainty are also shown in Figure 11. Across all five cities, uncertainties
372 in the AVOC emissions contribute the most (43%~65%) while the relative importance of other model inputs
373 differs by location. For example, O₃ BCs represent the second largest uncertainty source in Beijing (accounting
374 for 18%) but are negligible in Guangzhou and Chengdu. In Shanghai and Guangzhou, uncertainties in ANO_x
375 emissions (10%~17%) become the second largest contributor. Uncertainties associated with BVOC emissions are
376 minimal in Beijing and Shanghai but noteworthy (7~8%) in Guangzhou and Chengdu. O₃ deposition uncertainty
377 contributes to 8~30% of the total uncertainty, with a higher contribution for cities located in the west.



378

379 **Figure 11** Compared with the average observation results of five urban monitoring points in June 2021, the
380 uncertainty limit of MDA8 O₃ is $\pm 1 \sigma$. The pie chart shows the contribution of each factor to the total uncertainty
381 of the predicted average MDA8 O₃ in June 2021.

382 5. Conclusions

383 Chemical transport models are increasingly being employed to tackle the severe ozone pollution issues in China.
384 This study involved the compilation and analysis of 216 peer-reviewed studies focused on the use of CTMs to
385 simulate O₃ levels in China. Essential model configurations such as study region, simulation season, grid spacing,
386 gas-phase mechanism, and quantitative model performance outcomes were systematically documented. The study
387 presented quantile distributions of common statistical metrics found in the literature and discussed the influence
388 of different model configurations on performance outcomes. Furthermore, we proposed benchmarks for six
389 widely used MPE metrics (R, IOA, NMB, NME, FB, and FE) based on the concepts of "goals" and "standards" to
390 offer guidance to modelers for a more consistent and contextual evaluation of models. Additionally, we utilized
391 CMAQ-DDM to assess the uncertainties in predicted O₃ concentrations resulting from uncertainties in six model
392 inputs. The findings revealed significant variations in spatial distributions and magnitudes of ozone sensitivity to
393 different model inputs, with the most substantial contributions to total uncertainty originating from AVOC
394 emissions in regions with high ozone sensitivity.

395 The proposed benchmarks for assessing simulated O₃ concentrations, in conjunction with previous studies on
396 PM_{2.5} (Huang et al. 2021) and other criteria air pollutants (Zhai et al. 2024), represent a comprehensive and
397 systematic effort to establish a model performance framework for CTM applications in China. These outcomes
398 not only offer valuable guidance to the growing modeling community in China but also support their endeavors
399 in utilizing CTMs to address various research challenges and enhance air quality management.

400

401 **Data availability.** All data is available upon request from the corresponding author.

402 **Acknowledgements.** This work is supported by the Shanghai Technical Service Center of Science and
403 Engineering Computing, Shanghai University.



- 404 **Competing interests.** At least one of the (co-)authors are members of the editorial board of journal ACP.
- 405 **Financial support.** This study was supported by the National Natural Science Foundation of China (Grant No.
- 406 42375103, 42375102).
- 407 **References**
- 408 Ainsworth, E. A., Yendrek, C. R., Sitch, S., Collins, W. J., and Emberson, L. D.: The effects of tropospheric
- 409 ozone on net primary productivity and implications for climate change, *Annual review of plant biology*, 63, 637-
- 410 661, <https://doi.org/10.1146/annurev-arplant-042110-103829>, 2012.
- 411 Appel, K. W., Bash, J. O., Fahey, K. M., Foley, K. M., Gilliam, R. C., Hogrefe, C., Hutzell, W. T., Kang, D.,
- 412 Mathur, R., Murphy, B. N., Napelenok, S. L., Nolte, C. G., Pleim, J. E., Pouliot, G. A., Pye, H. O. T., Ran, L.,
- 413 Roselle, S. J., Sarwar, G., Schwede, D. B., Sidi, F. I., Spero, T. L., and Wong, D. C.: The Community Multiscale
- 414 Air Quality (CMAQ) model versions 5.3 and 5.3.1: system updates and evaluation, *Geosci. Model Dev.*, 14,
- 415 2867-2897, <https://doi.org/10.5194/gmd-14-2867-2021>, 2021.
- 416 Bai, K., Ma, M., Chang, N. B., and Gao, W.: Spatiotemporal trend analysis for fine particulate matter
- 417 concentrations in China using high-resolution satellite-derived and ground-measured PM_{2.5} data, *Journal of*
- 418 *environmental management*, 233, 530-542, <https://doi.org/10.1016/j.jenvman.2018.12.071>, 2019.
- 419 Beddows, A. V., Kitwiroon, N., Williams, M. L., and Beevers, S. D.: Emulation and Sensitivity Analysis of the
- 420 Community Multiscale Air Quality Model for a UK Ozone Pollution Episode, *Environmental Science &*
- 421 *Technology*, 51, 6229-6236, <https://doi.org/10.1021/acs.est.6b05873>, 2017.
- 422 Cao, L., Li, S., and Sun, L.: Study of different Carbon Bond 6 (CB6) mechanisms by using a concentration
- 423 sensitivity analysis, *Atmos. Chem. Phys.*, 21, 12687-12714, <https://doi.org/10.5194/acp-21-12687-2021>, 2021.
- 424 Carter, W. P. L.: Condensed atmospheric photooxidation mechanisms for isoprene, *Atmospheric Environment*,
- 425 30, 4275-4290, [https://doi.org/10.1016/1352-2310\(96\)00088-X](https://doi.org/10.1016/1352-2310(96)00088-X), 1996.
- 426 Carter, W. P. L.: Implementation of the SAPRC-99 chemical mechanism into the models-3 framework, Carter,
- 427 WPL, . 2000.
- 428 Carter, W. P. L.: Development of the SAPRC-07 chemical mechanism, *Atmospheric Environment*, 44, 5324-
- 429 5335, <https://doi.org/10.1016/j.atmosenv.2010.01.026>, 2010.
- 430 Chang, T. Y., Nance, B. I., and Kelly, N. A.: Modeling Smog Chamber Measurements of Vehicle Exhaust
- 431 Reactivities, *Journal of the Air & Waste Management Association* (1995), 49, 57-63,
- 432 <https://doi.org/10.1080/10473289.1999.10463775>, 1999.
- 433 Chen, B., Wang, Y., Huang, J., Zhao, L., Chen, R., Song, Z., and Hu, J.: Estimation of near-surface ozone
- 434 concentration and analysis of main weather situation in China based on machine learning model and Himawari-8
- 435 TOAR data, *Sci. Total Environ.*, 864, 160928, <https://doi.org/10.1016/j.scitotenv.2022.160928>, 2023.
- 436 Cheng, J., Su, J., Cui, T., Li, X., Dong, X., Sun, F., Yang, Y., Tong, D., Zheng, Y., Li, Y., Li, J., Zhang, Q., and
- 437 He, K.: Dominant role of emission reduction in PM_{2.5} air quality improvement in Beijing during 2013–2017:
- 438 a model-based decomposition analysis, *Atmos. Chem. Phys.*, 19, 6125-6146, [https://doi.org/10.5194/acp-19-](https://doi.org/10.5194/acp-19-6125-2019)
- 439 6125-2019, 2019.
- 440 Chu, B., Ma, Q., Liu, J., Ma, J., Zhang, P., Chen, T., Feng, Q., Wang, C., Yang, N., Ma, H., Ma, J., Russell, A. G.,
- 441 and He, H.: Air Pollutant Correlations in China: Secondary Air Pollutant Responses to NO_x and SO₂ Control,
- 442 *Environmental Science & Technology Letters*, 7, 695-700, <https://doi.org/10.1021/acs.estlett.0c00403>, 2020.
- 443 Cohan, D. S. and Napelenok, S. L.: Air Quality Response Modeling for Decision Support, *Atmosphere*, 2, 407-
- 444 425, <https://doi.org/10.3390/atmos2030407>, 2011.



- 445 Dang, R. and Liao, H.: Radiative Forcing and Health Impact of Aerosols and Ozone in China as the Consequence
446 of Clean Air Actions over 2012–2017, *Geophysical Research Letters*, 46, 12511-12519,
447 <https://doi.org/10.1029/2019GL084605>, 2019.
- 448 Derwent, R. G., Parrish, D. D., Galbally, I. E., Stevenson, D. S., Doherty, R. M., Naik, V., and Young, P. J.:
449 Uncertainties in models of tropospheric ozone based on Monte Carlo analysis: Tropospheric ozone burdens,
450 atmospheric lifetimes and surface distributions, *Atmospheric Environment*, 180, 93-102,
451 <https://doi.org/10.1016/j.atmosenv.2018.02.047>, 2018.
- 452 Dunker, A. M., Wilson, G., Bates, J. T., and Yarwood, G.: Chemical Sensitivity Analysis and Uncertainty
453 Analysis of Ozone Production in the Comprehensive Air Quality Model with Extensions Applied to Eastern
454 Texas, *Environmental Science & Technology*, 54, 5391-5399, <https://doi.org/10.1021/acs.est.9b07543>, 2020.
- 455 Emery, C., Liu, Z., Russell, A. G., Odman, M. T., Yarwood, G., and Kumar, N.: Recommendations on statistics
456 and benchmarks to assess photochemical model performance, *Journal of the Air & Waste Management*
457 *Association*, 67, 582-598, <https://doi.org/10.1080/10962247.2016.1265027>, 2017.
- 458 Gao, J., Zhu, B., Xiao, H., Kang, H., Hou, X., Yin, Y., Zhang, L., and Miao, Q.: Diurnal variations and source
459 apportionment of ozone at the summit of Mount Huang, a rural site in Eastern China, *Environmental Pollution*,
460 222, 513-522, <https://doi.org/10.1016/j.envpol.2016.11.031>, 2017.
- 461 Ge, B. Z., Wang, Z. F., Xu, X. B., Wu, J. B., Yu, X. L., and Li, J.: Wet deposition of acidifying substances in
462 different regions of China and the rest of East Asia: Modeling with updated NAQPMS, *ENVIRONMENTAL*
463 *POLLUTION*, 187, 10-21, <https://doi.org/10.1016/j.envpol.2013.12.014>, 2014.
- 464 Global Burden of Disease Collaborative Network. Global Burden of Disease Study 2019 (GBD 2019) Air
465 Pollution Exposure Estimates 1990-2019. Seattle, United States of America: Institute for Health Metrics and
466 Evaluation (IHME), <https://doi.org/10.6069/70JS-NC54>, 2021.
- 467 Goliff, W. S., Stockwell, W. R., and Lawson, C. V.: The regional atmospheric chemistry mechanism, version 2,
468 *Atmospheric Environment*, 68, 174-185, <https://doi.org/10.1016/j.atmosenv.2012.11.038>, 2013.
- 469 Huang, L., Zhu, Y., Zhai, H., Xue, S., Zhu, T., Shao, Y., Liu, Z., Emery, C., Yarwood, G., Wang, Y., Fu, J.,
470 Zhang, K., and Li, L.: Recommendations on benchmarks for numerical air quality model applications in China –
471 Part 1: PM_{2.5} and chemical species, *Atmos. Chem. Phys.*, 21, 2725-2743, [https://doi.org/10.5194/acp-21-2725-](https://doi.org/10.5194/acp-21-2725-2021)
472 2021, 2021.
- 473 Li, K., Jacob, D. J., Liao, H., Shen, L., Zhang, Q., and Bates, K. H.: Anthropogenic drivers of 2013–2017 trends
474 in summer surface ozone in China, *Proceedings of the National Academy of Sciences*, 116, 422-427,
475 <https://doi.org/10.1073/pnas.1812168116>, 2019.
- 476 Li, K., Jacob, D. J., Liao, H., Qiu, Y. L., Shen, L., Zhai, S. X., Bates, K. H., Sulprizio, M. P., Song, S. J., Lu, X.,
477 Zhang, Q., Zheng, B., Zhang, Y. L., Zhang, J. Q., Lee, H. C., and Kuk, S. K.: Ozone pollution in the North China
478 Plain spreading into the late-winter haze season, *Proceedings of the National Academy of Sciences of the United*
479 *States of America*, 118, <https://doi.org/10.1073/pnas.2015797118>, 2021.
- 480 Li, Q., Zhang, L., Wang, T., Tham, Y. J., Ahmadov, R., Xue, L., Zhang, Q., and Zheng, J.: Impacts of
481 heterogeneous uptake of dinitrogen pentoxide and chlorine activation on ozone and reactive nitrogen partitioning:
482 improvement and application of the WRF-Chem model in southern China, *Atmos. Chem. Phys.*, 16, 14875-14890,
483 <https://doi.org/10.5194/acp-16-14875-2016>, 2016.
- 484 Liu, H., Zhang, M., and Han, X.: A review of surface ozone source apportionment in China, *Atmospheric and*
485 *Oceanic Science Letters*, 13, 470-484, <https://doi.org/10.1080/16742834.2020.1768025>, 2020a.
- 486 Liu, H., Zhang, M., Han, X., Li, J., and Chen, L.: Episode analysis of regional contributions to tropospheric
487 ozone in Beijing using a regional air quality model, *Atmospheric Environment*, 199, 299-312,
488 <https://doi.org/10.1016/j.atmosenv.2018.11.044>, 2019a.



- 489 Liu, L., Wu, J., Liu, S., Li, X., Zhou, J., Feng, T., Qian, Y., Cao, J., Tie, X., and Li, G.: Effects of organic coating
490 on the nitrate formation by suppressing the N₂O₅ heterogeneous hydrolysis: a case study during wintertime in
491 Beijing–Tianjin–Hebei (BTH), *Atmos. Chem. Phys.*, 19, 8189–8207, <https://doi.org/10.5194/acp-19-8189-2019>,
492 2019b.
- 493 Liu, T., Wang, C., Wang, Y., Huang, L., Li, J., Xie, F., Zhang, J., and Hu, J.: Impacts of model resolution on
494 predictions of air quality and associated health exposure in Nanjing, China, *Chemosphere*, 249, 126515,
495 <https://doi.org/10.1016/j.chemosphere.2020.126515>, 2020b.
- 496 Liu, Y. and Wang, T.: Worsening urban ozone pollution in China from 2013 to 2017 – Part 1: The complex and
497 varying roles of meteorology, *Atmos. Chem. Phys.*, 20, 6305–6321, <https://doi.org/10.5194/acp-20-6305-2020>,
498 2020.
- 499 Liu, Y., Li, J., Ma, Y., Zhou, M., Tan, Z., Zeng, L., Lu, K., and Zhang, Y.: A review of gas-phase chemical
500 mechanisms commonly used in atmospheric chemistry modelling, *Journal of Environmental Sciences*, 123, 522–
501 534, <https://doi.org/10.1016/j.jes.2022.10.031>, 2023a.
- 502 Liu, Y., Geng, G., Cheng, J., Liu, Y., Xiao, Q., Liu, L., Shi, Q., Tong, D., He, K., and Zhang, Q.: Drivers of
503 Increasing Ozone during the Two Phases of Clean Air Actions in China 2013–2020, *Environmental Science &*
504 *Technology*, 7, 8954–8964, <https://doi.org/10.1021/acs.est.3c00054>, 2023b.
- 505 Lu, X., Zhang, L., Wang, X., Gao, M., Li, K., Zhang, Y., Yue, X., and Zhang, Y.: Rapid Increases in Warm-
506 Season Surface Ozone and Resulting Health Impact in China Since 2013, *Environmental Science & Technology*
507 *Letters*, 7, 240–247, <https://doi.org/10.1021/acs.estlett.0c00171>, 2020.
- 508 Luecken, D. J., Yarwood, G., and Hutzell, W. T.: Multipollutant modeling of ozone, reactive nitrogen and HAPs
509 across the continental US with CMAQ-CB6, *Atmospheric Environment*, 201, 62–72,
510 <https://doi.org/https://doi.org/10.1016/j.atmosenv.2018.11.060>, 2019.
- 511 Ni, R., Lin, J., Yan, Y., and Lin, W.: Foreign and domestic contributions to springtime ozone over China, *Atmos.*
512 *Chem. Phys.*, 18, 11447–11469, <https://doi.org/10.5194/acp-18-11447-2018>, 2018.
- 513 Ni, Z. Z., Luo, K., Gao, Y., Gao, X., Jiang, F., Huang, C., Fan, J. R., Fu, J. S., and Chen, C. H.: Spatial–temporal
514 variations and process analysis of O₃ pollution in Hangzhou during the G20 summit, *Atmos. Chem. Phys.*, 20,
515 5963–5976, <https://doi.org/10.5194/acp-20-5963-2020>, 2020.
- 516 Qiao, X., Guo, H., Wang, P., Tang, Y., Ying, Q., Zhao, X., Deng, W., and Zhang, H.: Fine Particulate Matter and
517 Ozone Pollution in the 18 Cities of the Sichuan Basin in Southwestern China: Model Performance and
518 Characteristics, *Aerosol and Air Quality Research*, 19, 2308–2319, <https://doi.org/10.4209/aaqr.2019.05.0235>,
519 2019.
- 520 Seinfeld, J. H. and Pandis, S. N.: *Atmospheric chemistry and physics : from air pollution to climate change*, 3rd
521 edition, John Wiley & Sons, Inc., ISBN 978-1-119-22117-3, 2016.
- 522 Shen, L., Liu, J., Zhao, T., Xu, X., Han, H., Wang, H., and Shu, Z.: Atmospheric transport drives regional
523 interactions of ozone pollution in China, *Sci. Total Environ.*, 830, 154634,
524 <https://doi.org/10.1016/j.scitotenv.2022.154634>, 2022.
- 525 Simon, H., Baker, K. R., and Phillips, S.: Compilation and interpretation of photochemical model performance
526 statistics published between 2006 and 2012, *Atmospheric Environment*, 61, 124–139,
527 <https://doi.org/10.1016/j.atmosenv.2012.07.012>, 2012.
- 528 Stockwell, W. R., Kirchner, F., Kuhn, M., and Seefeld, S.: A new mechanism for regional atmospheric chemistry
529 modeling, *Journal of Geophysical Research: Atmospheres*, 102, 25847–25879,
530 <https://doi.org/10.1029/97JD00849>, 1997.



- 531 Sun, Z., Tan, J., Wang, F., Li, R., Zhang, X., Liao, J., Wang, Y., Huang, L., Zhang, K., Fu, J. S., and Li, L.:
532 Regional background ozone estimation for China through data fusion of observation and simulation, *Sci. Total*
533 *Environ.*, 912, 169411, <https://doi.org/10.1016/j.scitotenv.2023.169411>, 2024.
- 534 Wang, Z., Li, J., Wang, Z., Yang, W., Tang, X., Ge, B., Yan, P., Zhu, L., Chen, X., Chen, H., Wand, W., Li, J.,
535 Liu, B., Wang, X., Wand, W., Zhao, Y., Lu, N., and Su, D.: Modeling study of regional severe hazes over mid-
536 eastern China in January 2013 and its implications on pollution prevention and control, *Science China Earth*
537 *Sciences*, 57, 3-13, <https://doi.org/10.1007/s11430-013-4793-0>, 2014.
- 538 Wang, B., Sun, M., Si, L., and Niu, Z.: Spatio-temporal variation of O₃ concentration and exposure risk
539 assessment in key regions of China, 2015–2021, *Atmospheric Pollution Research*, 15, 101941,
540 <https://doi.org/10.1016/j.apr.2023.101941>, 2024.
- 541 Wang, M. Y., Yim, S. H. L., Wong, D. C., and Ho, K. F.: Source contributions of surface ozone in China using
542 an adjoint sensitivity analysis, *Sci. Total Environ.*, 662, 385-392, <https://doi.org/10.1016/j.scitotenv.2019.01.116>,
543 2019.
- 544 Wang, T., Xue, L., Feng, Z., Dai, J., Zhang, Y., and Tan, Y.: Ground-level ozone pollution in China: a synthesis
545 of recent findings on influencing factors and impacts, *Environmental Research Letters*, 17, 063003,
546 <https://doi.org/10.1088/1748-9326/ac69fe>, 2022.
- 547 Wang, W.-N., Cheng, T.-H., Gu, X.-F., Chen, H., Guo, H., Wang, Y., Bao, F.-W., Shi, S.-Y., Xu, B.-R., Zuo, X.,
548 Meng, C., and Zhang, X.-C.: Assessing Spatial and Temporal Patterns of Observed Ground-level Ozone in China,
549 *Scientific Reports*, 7, 3651, <https://doi.org/10.1038/s41598-017-03929-w>, 2017.
- 550 Xu, T., Zhang, C., Liu, C., and Hu, Q.: Variability of PM_{2.5} and O₃ concentrations and their driving forces over
551 Chinese megacities during 2018-2020, *Journal of Environmental Sciences*, 124, 1-10,
552 <https://doi.org/10.1016/j.jes.2021.10.014>, 2023.
- 553 Yang, J. and Zhao, Y.: Performance and application of air quality models on ozone simulation in China – A
554 review, *Atmospheric Environment*, 293, 119446, <https://doi.org/10.1016/j.atmosenv.2022.119446>, 2023.
- 555 Yang, L., Xie, D., Yuan, Z., Huang, Z., Wu, H., Han, J., Liu, L., and Jia, W.: Quantification of Regional Ozone
556 Pollution Characteristics and Its Temporal Evolution: Insights from Identification of the Impacts of
557 Meteorological Conditions and Emissions, *Atmosphere*, 12, 279, <https://doi.org/10.3390/atmos12020279>, 2021a.
- 558 Yang, Y., Zhao, Y., Zhang, L., Zhang, J., Huang, X., Zhao, X., Zhang, Y., Xi, M., and Lu, Y.: Improvement of
559 the satellite-derived NO_x emissions on air quality modeling and its effect on ozone and secondary inorganic
560 aerosol formation in the Yangtze River Delta, China, *Atmos. Chem. Phys.*, 21, 1191-1209,
561 <https://doi.org/10.5194/acp-21-1191-2021>, 2021b.
- 562 Yao, Y., Ma, K., He, C., Zhang, Y., Lin, Y., Fang, F., Li, S., and He, H.: Urban Surface Ozone Concentration in
563 Mainland China during 2015-2020: Spatial Clustering and Temporal Dynamics, *International journal of*
564 *environmental research and public health*, 20, <https://doi.org/10.3390/ijerph20053810>, 2023.
- 565 Yarwood, G., Jung, J., Whitten, G. Z., Heo, G., Mellberg, J., and Estes, M.: UPDATES TO THE CARBON
566 BOND MECHANISM FOR VERSION 6 (CB6), <https://doi.org/10.1093/bioinformatics/btp533>, 1997.
- 567 Yarwood, G. and Tuite, K.: Representing Ozone Formation from Volatile Chemical Products (VCP) in Carbon
568 Bond (CB) Chemical Mechanisms, *Atmosphere*, 15, 178, <https://doi.org/10.3390/atmos15020178>, 2024.
- 569 Zhai, H., Huang, L., Emery, C., Zhang, X., Wang, Y., Yarwood, G., Fu, J. S., and Li, L.: Recommendations on
570 benchmarks for photochemical air quality model applications in China — NO₂, SO₂, CO and PM₁₀, *Atmospheric*
571 *Environment*, 319, 120290, <https://doi.org/10.1016/j.atmosenv.2023.120290>, 2024.