



# Ensemble Forecasts of Air Quality in Eastern China

## Part 2. Evaluation of the MarcoPolo-Panda Prediction System, Version 1.

Anna Katinka Petersen<sup>1</sup>, Guy P Brasseur<sup>1,2</sup>, Idir Bouarar<sup>1</sup>, Johannes Flemming<sup>3</sup>, Michael Gauss<sup>4</sup>, Fei Jiang<sup>5</sup>, Rostislav Kouznetsov<sup>6</sup>, Richard Kranenburg<sup>7</sup>, Bas Mijling<sup>8</sup>, Vincent-Henri Peuch<sup>3,9</sup>, Matthieu Pommier<sup>4</sup>, Arjo Segers<sup>7</sup>, Mikhail Sofiev<sup>6</sup>, Renske Timmermans<sup>7</sup>, Ronald van der A<sup>8,9</sup>, Stacy Walters<sup>7</sup>, Ying Xie<sup>10</sup>, Jianming Xu<sup>10</sup>, Guangqiang Zhou<sup>10</sup>

<sup>1</sup> Max Planck Institute for Meteorology, Hamburg, Germany

<sup>2</sup> National Center for Atmospheric Research, Boulder, CO, USA

<sup>3</sup> European Centre for Middle Range Weather Forecasts, Reading, UK.

<sup>4</sup> Norwegian Meteorological Institute, Oslo, Norway

<sup>5</sup> Nanjing University, Nanjing, China

<sup>6</sup> Finnish Meteorological Institute, Helsinki, Finland.

<sup>7</sup> TNO, Utrecht, The Netherlands

<sup>8</sup> Royal Netherlands Meteorological Institute (KNMI), De Bilt, The Netherlands

<sup>9</sup> Nanjing University of Information Science and Technology, Nanjing, China

<sup>10</sup> Shanghai Meteorological Service, Shanghai, China

### Abstract:

An operational multi-model forecasting system for air quality has been developed to provide air quality services for urban areas of China. The initial forecasting system included seven state-of-the-art computational models developed and executed in Europe and China (CHIMERE, IFS, EMEP MSC-W, WRF-Chem-MPIM, WRF-Chem-SMS, LOTOS-EUROS and SILAMtest). Several other models joined the prediction system recently, but are not considered in the present analysis. In addition to the individual models, a simple multi-model ensemble was constructed by deriving statistical quantities such as the median and the mean of the predicted concentrations.

The prediction system provides daily forecasts and observational data of surface ozone, nitrogen dioxides and particulate matter for the 37 largest urban agglomerations in China (population higher than 3 million in 2010). These individual forecasts as well as the multi-model ensemble predictions for the next 72 hours are displayed as hourly outputs on a publicly accessible web site ([www.marcopolo-panda.eu](http://www.marcopolo-panda.eu)).

In this paper, the performance of the predictions system (individual models and the multi-model ensemble) for the first operational year (April 2016 until June 2017) has been analysed through statistical indicators using the surface observational data reported at Chinese national monitoring stations. This evaluation aims to investigate a) the seasonal behavior, b) the geographical distribution and c) diurnal variations of the ensemble and model skills. Statistical indicators show that the ensemble product usually provides the best performance compared to the individual model forecasts. The ensemble product is robust even if occasionally some individual model results are missing.

Overall and in spite of some discrepancies, the air quality forecasting system is well suited for the prediction of air pollution events and has the ability to provide alert warning (binary prediction) of air pollution events if bias corrections are applied to improve the ozone predictions.



## 47 1. Introduction

48

49 With the rapid development of its economy, China has been experiencing repeated intense air  
50 pollution episodes (e.g. *Guo et al., 2014, Huang et al., 2014, Wang et al., 2014*) with a wide range  
51 of health effects (*Kampa and Castanas 2008; Wu et al., 2012; Hamra et al. 2015; Boynard et al.,*  
52 *2014; WHO, 2018*) and serious consequences on ecosystems (*Fowler et al., 2008, Ashmore, 2005;*  
53 *Leisner et al., 2012; Sinha et al., 2015*) and on climate (*Sitch et al. 2007; Brasseur et al., 1999;*  
54 *Akimoto, 2003*). High concentrations of particulate matter often cover a large area of eastern China  
55 during winter when air remains stagnant for several days and chemical compounds emitted by  
56 power plants, industrial complexes, traffic and domestic infrastructures remain trapped near the  
57 surface (e.g. *Wang et al., 2014; Zhao et al., 2013*). During summer, photochemical processes  
58 convert nitrogen oxides (NO<sub>x</sub>) and volatile organic compounds (VOCs) into tropospheric ozone  
59 (O<sub>3</sub>) (e.g. *Xu et al., 2008, Sun et al., 2016*).

60

61 Long-term solutions to mitigate air pollution require a fundamental transformation of the energy  
62 system, which may require decades to be fully implemented. Short-term actions to avoid severe air  
63 pollution episodes, however, can be put in place immediately if such episodes can be reliably  
64 predicted a few days prior to their occurrence. Comprehensive air quality models that capture  
65 meteorological, chemical and physical processes in the troposphere and predict the fate of air  
66 pollutants are key tools to forecast the likelihood of air pollution episodes and hence to inform the  
67 authorities.

68

69 Within the EU projects MarcoPolo and Panda, that include European as well as Chinese partner  
70 organizations, an operational multi-model forecasting system for air quality including a number of  
71 different chemical transport models has been developed, and is providing daily forecasts of ozone,  
72 nitrogen oxides, and particulate matter for the 37 largest urban areas of China (population higher  
73 than 3 million in 2010). These individual forecasts as well as the mean and median concentrations  
74 for the next 3 days are posted on a dedicated website ([www.marcopolo-panda.eu/forecast](http://www.marcopolo-panda.eu/forecast)) together  
75 with the hourly observational data from local measurements reported by the Chinese monitoring  
76 network of the China National Environmental Monitoring Centre (CNEMC) (data available at  
77 [www.pm25.in](http://www.pm25.in)). This operational air quality analysis and forecasting system is presented in detail in  
78 a companion paper (*Brasseur et al, 2018*), where the individual models contributing to the  
79 MarcoPolo-Panda prediction system are described, and details about the individual models and their  
80 individual settings are provided. Information about selected parametrization options for the physical  
81 processes, including boundary layer, radiation, convection and surface processes, and about the  
82 emissions adopted in MarcoPolo-Panda prediction system are also provided.

83

84 In the present study, we evaluate the prediction system of the MarcoPolo and Panda projects that  
85 have been in operation for more than one year. We concentrate on the period April 2016 to June  
86 2017 and analyse the model forecasts (7 individual models and the ensemble median) and  
87 observational data for 34 cities (covered by most of the models, depending on the extent of the  
88 domains, for two models only 31 and 32 cities).

89

90 We evaluate the performance of the individual models involved in the present study, and to examine  
91 the performance of the overall forecasting system by comparing the predicted surface  
92 concentrations to values reported by the Chinese air pollution monitoring network. Section 2 of the  
93 paper provides a brief description of the forecasting system, while Section 3 investigates the  
94 performance of the system using different statistical indicators including the mean bias (BIAS), the  
95 root mean square error (RMSE), the modified normalised bias (MNBIAS), the fractional gross error  
96 (FGE) and the correlation coefficient. We derive in particular (a) statistical indicators for each



97 model over the time of the year (on a monthly basis) in order to analyse seasonal characteristics, (b)  
98 the geographical distribution of the statistical indicators for the ensemble median in order to derive  
99 regional characteristics and issues, (c) the statistical indicators of all models and of the ensemble  
100 median over the time of the day (considering all model-observation pairs of all cities and for the  
101 whole time period) and for a specific city (Beijing) together with the diurnal variation of the  
102 pollutants during the whole time period. In Section 4, we assess the impacts of missing forecasts  
103 from one or more models on the production of the ensemble. As the prediction system intends to  
104 provide warning of air pollution episodes to the general public, the system performance has been  
105 evaluated regarding its ability to predict the exceedence of air quality thresholds (binary prediction  
106 of pollution events). This analysis is presented in Section 5. We show that the application of bias  
107 correction to the models improves the forecasting skills of binary ozone predictions. We conclude  
108 with a summary and outlook in Section 6.

109  
110

## 111 2. Description of the Analysis and Forecasting System

112 Within the EU projects MarcoPolo and Panda, a number of chemistry transport models have been  
113 applied to provide daily air quality forecasts for a selection of 37 large Chinese agglomerations  
114 (population over 3 million, 2010 census). Initially, seven models, CHIMERE (Royal Netherlands  
115 Meteorological Institute (KNMI)), IFS (European Centre for Medium Range Weather Forecast  
116 (ECMWF)), WRF-chem-SMS (Shanghai Meteorological Service (SMS)), SILAMtest (Finish  
117 Meteorological Institute (FMI)), WRF-chem-MPIM (Max Planck Institute for Meteorology  
118 (MPIM) in Hamburg), EMEP MSC-W (hereafter referred to as ‘EMEP’, Norwegian Meteorological  
119 Institute (MET Norway)) and LOTOS-EUROS (The Netherlands Organisation for Applied  
120 Scientific Research (TNO)) were providing daily forecasts every day at 0:00 UTC for the next 72  
121 hours (three days) for NO<sub>2</sub>, O<sub>3</sub>, PM<sub>10</sub> and PM<sub>2.5</sub> (see Figure 1). WRF-CMAQ and WRMS-  
122 CMAQ, both used by Chinese institutions (Nanjing University and SMS), have joined recently the  
123 prediction system, but are not considered in the present analysis.

124

125 We should note that the models considered in the present study may have significantly evolved  
126 since the present analysis was performed. This is the case, for example, of the SILAM model  
127 developed by the Finish Meteorological Institute, whose configuration was still in a test mode, and  
128 is therefore referred to as SILAMtest.

129

130 The individual models are executed independently on the computing systems available in each  
131 partner institution. The surface concentrations of the key chemical species are extracted locally  
132 from the model outputs and forwarded to a central database operated by the Royal Netherlands  
133 Meteorological Institute (KNMI).

134

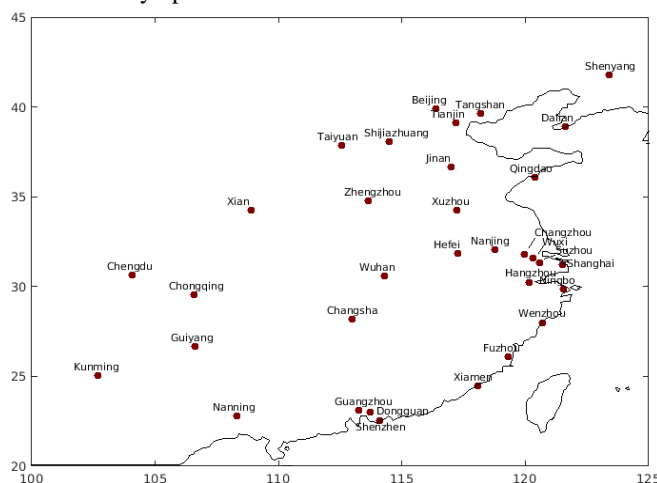
135 Hourly predictions of surface concentrations (expressed in  $\mu\text{g}/\text{m}^3$ ), are provided by the models as  
136 grid values, which are bi-linearly interpolated to city center coordinates. The average for the data  
137 provided by the urban network (usually around 5-12 stations), is posted together with the  
138 corresponding standard deviation and the number of contributing stations. In the present analysis,  
139 we consider only the model simulations corresponding to 34 cities, since the cities of Ürümqi (most  
140 western, only covered by three models), Changchun and Harbin (most northern cities), are located  
141 outside of the domains covered by most individual models, which are indicated in the companion  
142 paper (*Brasseur et al., 2018*).

143

144 In addition to the forecasts provided by the individual participating models, a multi-model ensemble  
145 was constructed from which the median and the mean were derived. To process the ensemble



146 median, all seven individual models are first interpolated to a common horizontal grid. For each  
147 grid point, the ensemble model is calculated as the median value of the individual model forecasts.  
148 The median is relatively insensitive to outliers in the forecasts. The method is also less vulnerable to  
149 occasionally missing data from individual models, as the minimum number of model results needed  
150 to calculate a meaningful ensemble mean or median is almost always available. This will be  
151 discussed in detail in Section 4. The multi-model approach also provides more accurate forecasts  
152 and thus reduces the underlying uncertainties (as will be shown in the following section). More  
153 advanced methods, e.g. based on individual model skills, are discussed in the literature (e.g.  
154 *Galmarini et al, 2013*). They are significantly more costly from a computational point of view and  
155 therefore not well suited for daily operations.



156

*Figure 1: Map of the 34 cities/urban clusters (population over 3 million (2010 census)) with available data (observational and model ensembles), used in this evaluation.*

### 157 3. Evaluation of the performance of the system

158

159 The evaluation of the performance of a forecasting system is a necessary step for assessing the  
160 quality of the predictions and demonstrating its usefulness. It also provides important information  
161 that can lead to the improvement of the forecasting system and to further model development. The  
162 comparison between model output and in situ measurements is not straightforward because of the  
163 different nature of the respective quantities: air quality models provide volume averaged quantities  
164 over each model grid cell and time averages over the modeling time step. Observations are available  
165 at fixed measurement sites and at a fixed time. Further, they are influenced by local processes that  
166 are not necessarily well captured by relatively coarse models. Thus, the representativeness of the  
167 observational site is not always guaranteed.

168

169 The MarcoPolo-Panda forecasting and analysis system uses the surface observations available at the  
170 web site [www.pm25.in](http://www.pm25.in) for 37 Chinese cities. For a given city, the observational data considered for  
171 the evaluation of the model consist of an average of the measurements made at the different stations  
172 of the urban network, usually 5 – 12 stations, which are aggregated to one value for the whole city.  
173 The model fields are bilinearly interpolated to the city center coordinates.

174

175 The mean bias



176

177

$$BIAS = \frac{1}{N} \sum_i (m_i - o_i),$$

178

179 where  $m_i$  and  $o_i$  are the model forecast value and the observation value, and  $N$  the number of  
 180 model-observation pairs, the root mean square error

181

182

$$RMSE = \sqrt{\frac{1}{N} \sum_i (m_i - o_i)^2},$$

183

184 the modified normalized bias

185

186

$$MNBIAS = \frac{2}{N} \sum_i \frac{(m_i - o_i)}{(m_i + o_i)},$$

187

188 the fractional gross error

189

$$FGE = \frac{2}{N} \sum_i \frac{|m_i - o_i|}{|m_i + o_i|}$$

190

191 and the correlation coefficient between the model forecast and observed values

192

$$R = \frac{\frac{1}{N} \sum_i (m_i - \bar{m})(o_i - \bar{o})}{\sigma_m \sigma_o}$$

193

194 are used to measure the system performance. Here  $\bar{m}$  and  $\bar{o}$  are the mean values of the model  
 195 forecast and observed values, and  $\sigma_m$  and  $\sigma_o$  are the corresponding standard deviations.

196

197 The evaluation presented here aims to investigate a) the statistical indicators for each model over  
 198 the time of the year (on a monthly basis) so that the seasonal features can be characterized and  
 199 related issues of individual models can be identified (Section 3.1); b) the geographical distribution  
 200 of the statistical indicators of the ensemble median to highlight regional characteristics and related  
 201 issues (Section 3.2); c) statistical indicators of all models and the ensemble median over the time of  
 202 the day (considering all model-observation pairs of all cities and for the whole time period) and for  
 203 a specific city (Beijing) together with the diurnal variation of the pollution species over the whole  
 204 time period (Section 3.3).

205

206

### 207 3.1 Evaluation of the Seasonal Behavior of the Models

208

209 We start our evaluation of the multi-model prediction system by examining the seasonal behavior of  
 210 the predicted concentrations of key chemical species. The statistical indicators mentioned above  
 211 have been calculated separately for each month from April 2016 to June 2017 and for the entire  
 212 period during which the forecasting system was operational. Due to storage issues, only the  
 213 predictions for the first 24 hours (0-23h) were saved while the predictions from 24h-72h were not  
 214 retained and not analyzed in this work.

215

216



217 Figure 2 shows the RMSE, BIAS, MNBIAS and FGE of NO<sub>2</sub> (left panel) and O<sub>3</sub> (right panel) for  
218 each of the seven individual models included in the system and for the model ensemble median, for  
219 each individual month between April 2016 and June 2017. The same results are also provided for  
220 the whole period (“all”). It can be seen, that there is a wide spread of the results produced by the  
221 seven models. The individual models are continuously improving during the first months because  
222 many changes have been applied by the different modeling groups in order to improve their  
223 individual predictions. In the case of NO<sub>2</sub>, most individual models slightly overestimate the  
224 concentrations compared to observations. In the EMEP model, it may be explained by the larger  
225 nitric oxide emissions used in comparison with the other models (Brasseur et al., 2018). This  
226 results in a positive BIAS and MNBIAS for most models and the ensemble median. The RMSE of  
227 the model ensemble is highest in July/August/September 2016 and remains relatively constant after  
228 October 2016. It can be seen, that the median of the model ensemble has the lowest RMSE for NO<sub>2</sub>,  
229 the smallest BIAS and MNBIAS (slightly positive) and the lowest FGE. This demonstrates the  
230 advantage of adopting a model ensemble rather than the prediction provided by individual models.

231  
232 Most models underestimate O<sub>3</sub> (likely as a result of the overestimated NO<sub>2</sub> because the O<sub>3</sub>  
233 production is not NO<sub>x</sub>-limited) during the whole period under consideration. For O<sub>3</sub>, the CHIMERE  
234 model shows slightly better performance (lowest RMSE) than the model ensemble median. The  
235 median BIAS for O<sub>3</sub> is relatively constant (slightly negative). For this particular species, the model  
236 ensemble median does not provide the best results regarding the BIAS. In fact, in this case, the  
237 model LOTOS-EUROS gives the best performance for ozone. Interestingly, this particular model  
238 has the largest negative BIAS for NO<sub>2</sub>. The median BIAS of O<sub>3</sub> remains relatively constant during  
239 the period, while the MNBIAS exhibits higher negative values during the winter months, as a result  
240 of the relative low O<sub>3</sub> concentrations during winter time.

241  
242 As stated above, the MarcoPolo-Panda prediction system has the tendency to overestimate surface  
243 NO<sub>2</sub>, which leads to O<sub>3</sub> titration especially during night time. The emission injection height is also  
244 a relevant factor here since it can largely influence the results in the planetary boundary layer.  
245 During night-time, emissions from stacks may be take place above the mixing layer and explain  
246 model-data discrepancies since the models often assume that the injection of primary pollutants  
247 takes place in the first layer above the surface.

248  
249  
250 Anthropogenic emissions of primary pollutants are changing extremely rapidly in China. The  
251 adopted emissions inventories usually reflect to the situation a few years before the period during  
252 which the model simulations were performed. Since the recent NO<sub>x</sub> emissions have decreased  
253 significantly in some urban areas of China in response to measures taken by the local authorities (*F.*  
254 *Liu et al., 2017*), the anthropogenic emissions used for the current forecasts may be overestimated  
255 in some areas. Some models use reduced NO<sub>x</sub> and SO<sub>x</sub> anthropogenic emissions (for details see  
256 *Brasseur et al., 2018*), however, daytime concentrations of ozone are generally underestimated in  
257 most models, even when the level of NO<sub>2</sub> is in reasonable agreement with the observational values.  
258 The discrepancy could be caused by an underestimation of the emissions of some VOCs, especially  
259 in the center of urban areas where ozone is often VOC-limited.

260  
261 For PM<sub>10</sub> and PM<sub>2.5</sub>, the model ensemble median shows the best performance compared to all  
262 individual models during the time period under consideration (see Figure 3). For PM<sub>10</sub>, there is an  
263 overall slight underestimation by all models except by CHIMERE and hence, by the median of the  
264 model ensemble. For PM<sub>2.5</sub>, the BIAS is relatively constant (apart in the WRF-Chem-SMS model  
265 which exhibits a lot of variation in the BIAS of PM<sub>10</sub> and PM<sub>2.5</sub>). In this case, the BIAS is slightly  
266 overestimated, but close to zero.



267  
268 Figure 4 shows the temporal correlation coefficients for NO<sub>2</sub>, O<sub>3</sub>, PM10 and PM2.5 for each  
269 individual month, and for the whole time period. It can be seen, that there is a wide spread between  
270 the individual models: the calculated correlations range from 0.2 to 0.7 for NO<sub>2</sub>, PM10 and PM2.5  
271 and from 0.3 to 0.8 for O<sub>3</sub>. The model ensemble median and CHIMERE are characterized by high  
272 correlation coefficients in the case of NO<sub>2</sub>, O<sub>3</sub> and PM2.5. For PM10, the model ensemble median  
273 and the LOTOS-EUROS model provide the highest correlation coefficients. In general, the model  
274 ensemble median gives the best performance.

275  
276 The correlation coefficient of O<sub>3</sub> for the ensemble median remains relatively unchanged during the  
277 whole time period, and ranges between 0.6 and 0.8. Considering the whole time period, it is of the  
278 order of 0.75, with CHIMERE providing a slightly higher correlation coefficient for the whole time  
279 period, and also for each individual months. All models exhibit small correlation coefficients in  
280 March 2017. High correlation coefficients are found during the early summer months (June/July).  
281 For PM10 and PM2.5 the correlation coefficients exhibit more variability, starting with very low  
282 correlation for all models and for the ensemble during April and May 2016, high correlation from  
283 June 2016 to March 2017, and again low correlation during April and May 2017. These differences  
284 may be due to missing sources of biomass burning or dust or to individual model tunings. For the  
285 entire time period, the correlation coefficient of the ensemble mean is higher than for each  
286 individual models (~0.58 for PM10 and ~0.78 for PM2.5). The correlation between the model  
287 ensemble and the observations is therefore relatively satisfactory.

### 289 3.2 Evaluation of the Geographical Distribution

290 The statistical indicators, described above for all contributing cities, have also been calculated for  
291 the individual cities. The purpose here is to assess regional characteristics and to identify model  
292 issues. Figure 5 shows the statistical indicators (RMSE, BIAS and correlation coefficient) for O<sub>3</sub>,  
293 NO<sub>2</sub> and PM2.5 of the Ensemble Median for each city during the time period under consideration  
294 (April 2016 until June 2017). In the upper most left panel, the BIAS of ozone for each city is  
295 shown. It can be seen, that the ensemble median is underestimating the ozone concentrations in the  
296 north and northeastern regions of China, while no significant bias compared to the observations is  
297 found in cities in the southern part of the country. RMSE in the northern/northeastern cities are  
298 higher (around 40 µg m<sup>-3</sup>) than in southern and western cities (around 20-30 µg m<sup>-3</sup>).

299  
300 The temporal correlation coefficients for ozone calculated for each city over the whole period under  
301 consideration are slightly higher in the northern part of the country and slightly smaller in the  
302 southern regions. This indicates that the day-to-day variability is well simulated, even though the  
303 models are slightly underestimating the ozone pollution in the north. NO<sub>2</sub> concentrations (see the  
304 middle panels of Figure 5) are overestimated in some cities and underestimated in other cities.  
305 There is, however, no systematic geographical characterization of the bias. When considering  
306 individual cities, it can be seen that the NO<sub>2</sub> concentrations are slightly overestimated in most urban  
307 areas including Beijing, Shanghai, Chengdu, Wuhan and Changsha. The RMSE for NO<sub>2</sub> in the  
308 middle panel of Figure 5 is very uniform (around 20 µg m<sup>-3</sup>) in the whole country. The correlation  
309 coefficients of NO<sub>2</sub> (between 0.5 and 0.7) are smaller than those of O<sub>3</sub>, as NO<sub>2</sub> exhibits more  
310 temporal variability than O<sub>3</sub>. In the case of PM2.5, (see upper most right panel), the concentrations  
311 are well simulated in the northern and southern parts of China, but there are a few city clusters in  
312 the middle of the domain (Chengdu, Chongqing, Wuhan and Changsha) in which the PM2.5  
313 concentrations are overestimated by more than 50 µg m<sup>-3</sup>. These cities also show overestimation  
314 of NO<sub>2</sub>. The overestimation of PM2.5 may therefore be related to the errors in precursor emissions,  
315 e.g. NO<sub>x</sub>, SO<sub>2</sub>. The RMSE of PM2.5 is smaller in the southern part of the domain and along the



316 coastline of China, while the model results are less satisfactory in the city clusters located in the  
317 central part of the domain, with very high RMSE of 60-80 $\mu\text{g m}^{-3}$  in three cities. The correlation  
318 coefficients for the individual cities are relatively constant around 0.7 with few cities characterized  
319 by lower correlation coefficients (mostly in the central part of the domain).

320

### 321 3.3 Evaluation of the diurnal variation

322 We now examine the ability of the models to reproduce the diurnal variations of the chemical  
323 species' concentrations. We first provide a general view based on all observations in China and then  
324 examine the particular situation in the city of Beijing.  
325

#### 326 3.3.a Analysis based on all observations in China

327 The RMSE, BIAS, MNBIAS, and FGE of O<sub>3</sub>, NO<sub>2</sub>, PM<sub>10</sub> and PM<sub>2.5</sub> for the seven models and the  
328 ensemble median for all available observations in China are displayed over the forecasting time (0-  
329 23h) (Figure 6 and 7). Due to storage limitations, only the predictions for the first 24 hours (0-23h)  
330 were saved while the predictions for the 24h-72h period performed by all models were not retained.  
331 Unfortunately, this does not allow the investigation of a day to day degradation of the statistical  
332 indicators (from day1 to day3). Only the diurnal behavior of the statistical indicators can be  
333 assessed, which provides important hints for possible model issues.

334

335 It can be seen in the left panels of Figure 6 that the statistical indicators of NO<sub>2</sub> for the ensemble  
336 median is relatively stable over the time of the day, with slightly higher RMSE and higher  
337 BIAS/MNBIAS during the night time hours. For the individual models, the variability of the RMSE  
338 is somewhat higher during daytime, while some models exhibit very high RMSE and BIAS during  
339 the night time hours. Most models show a positive BIAS of NO<sub>2</sub> during the night, but a few of them  
340 exhibit a negative bias; this results in a relatively small BIAS for the ensemble median, showing  
341 good results with respect to the BIAS throughout the day.

342

343 In the case of ozone, the statistical indicators exhibit a variation over the time of the day. The  
344 RMSE is smallest between 7:00 and 9:00 local time, after which it increases until 18:00 in the  
345 evening to become constant at about 30  $\mu\text{g m}^{-3}$  during the night.

346

347 An examination of the BIAS and MNBIAS for O<sub>3</sub> over the day shows that O<sub>3</sub> is underestimated by  
348 nearly all models, apart from WRF-Chem-SMS. This might result from the slight overestimation of  
349 NO<sub>2</sub> concentrations by most models. Especially during nighttime when the height of the boundary  
350 layer is low, near surface NO<sub>2</sub> concentrations are high, and ozone is underestimated by 50% – 100%  
351 by most models. In the first hours of the day, only SILAMtest, WRF-Chem-SMS and LOTOS-  
352 EUROS exhibit slightly positive O<sub>3</sub> BIAS. The same models produce a negative BIAS for NO<sub>2</sub>  
353 during the first hours of the day.

354

355 Figure 7 shows that the BIAS and MNBIAS of both PM<sub>10</sub> and PM<sub>2.5</sub> stay relatively constant over  
356 the time of the day. PM<sub>10</sub> is slightly underestimated by the ensemble median (-5 to -10%), while  
357 PM<sub>2.5</sub> is slightly overestimated (10 to 25%). In most cases, the models overestimate the PM<sub>2.5</sub>  
358 observations, while for PM<sub>10</sub> there are stronger differences between the individual models.

359

360 For PM<sub>10</sub> and PM<sub>2.5</sub>, the ensemble median exhibits a better performance than the individual  
361 models: the RMSE BIAS, MNBIAS and FGE of the ensemble are on average lower than the





362 corresponding statistical parameters of the individual models. This demonstrates again the  
363 advantage of using the ensemble median for the prediction of PM10 and PM2.5.

364

365 Figure 8 presents the diurnal variation of the concentrations of O<sub>3</sub>, NO<sub>2</sub>, O<sub>3</sub> + NO<sub>2</sub> and PM2.5 from  
366 the individual models (and the ensemble median) and from the observations at a specific location  
367 (Beijing). The RMSE and the BIAS are also provided during the whole period under consideration.

368

369 It can be seen that the ensemble median (black line) underestimates the O<sub>3</sub> observations (red line)  
370 throughout the day, especially during the nighttime hours and in the late afternoon. Only WRF-  
371 Chem-SMS reproduces the amplitude of the O<sub>3</sub> diurnal cycle, but it also underestimates the O<sub>3</sub>  
372 concentrations after 18:00 when the height of the boundary layer is rapidly decreasing. All models  
373 and the ensemble median reproduce the diurnal cycle with a maximum in the late afternoon, but this  
374 maximum produced by the model appears about 2 hours earlier than observed. When considering  
375 the RMSE, the models produce the best results during the morning, and with increasing O<sub>3</sub>  
376 concentrations as the day progresses, the RMSE is also increasing. The negative BIAS is increasing  
377 for all models and for the model ensemble throughout the day.

378

### 379 3.3.b Analysis for the specific case of Beijing

380

381 In Beijing, the diurnal variation of the NO<sub>2</sub> concentrations is overestimated by the individual  
382 models as also reflected by the ensemble median. During the nighttime, for example, the observed  
383 concentrations are about 20-30 µg m<sup>-3</sup> lower than the concentrations associated with the ensemble  
384 median. The individual models and the ensemble median show a much stronger diurnal behavior  
385 than the observations. Atmospheric measurements suggest that the concentrations of NO<sub>2</sub> are  
386 relatively constant over the time of the day. This might be due to applied temporal profiles of the  
387 anthropogenic emissions or issues in the vertical mixing of the individual models. Also, the models  
388 with their spatial resolution may not capture the details seen in the observations by the ground  
389 network. The RMSE of all models and for the ensemble median is highest in late afternoon and  
390 during the night. The MarcoPolo-Panda prediction system has thus a tendency to overestimate  
391 surface NO<sub>2</sub>, which leads to an overestimation of the O<sub>3</sub> titration especially at night.

392

393 To further analyze the chemical coupling between ozone and NO<sub>2</sub>, we have added at each time step  
394 the mixing ratios of O<sub>3</sub> and NO<sub>2</sub>. The resulting variable, called Ox and expressed here in ppbv, has  
395 the advantage of not being affected by the fast interchange (null cycle) and the resulting partitioning  
396 between ozone and NO<sub>2</sub> produced by reactions NO + O<sub>3</sub>, NO<sub>2</sub> + hv and O + O<sub>2</sub> + M. If only these  
397 three rapid photochemical reactions are considered, Ox is a conserved quantity. In other words,  
398 even when a more comprehensive chemical scheme is adopted, the diurnal cycle of Ox should be  
399 considerably less pronounced than the diurnal cycle of NO<sub>2</sub> and O<sub>3</sub>.

400

401 In fact, in the model forecasts, the sum of O<sub>3</sub> and NO<sub>2</sub>, is nearly constant during the day, but  
402 exhibits nevertheless some diurnal variation, which appears to be weaker than in the observation.  
403 The calculated O<sub>x</sub> is slightly too high at night and too low during daytime, suggesting an  
404 overestimation in photochemical activity by the majority of the models. The partitioning of O<sub>x</sub> into  
405 NO<sub>2</sub> and O<sub>3</sub> is not well reproduced despite the simple chemistry that determines this partitioning:  
406 NO<sub>2</sub> is generally too high and O<sub>3</sub> too low, especially in the afternoon and early night. The simple  
407 partitioning approach does not seem to work properly under high NO<sub>x</sub> loading. As a result, the  
408 diurnal cycle of O<sub>3</sub> is not well reproduced by the forecasting ensemble and high ozone events are  
409 generally underestimated. This issue is discussed in more detail in the companion paper by  
410 *Brasseur et al., 2018*.



411

412 The observed diurnal variation of PM<sub>2.5</sub> is not well reproduced by the models and by the ensemble  
413 median. The calculated variability in Beijing is substantially higher than suggested by the  
414 observations (which are characterized by relatively constant concentrations throughout the day).  
415 The models show a maximum in PM<sub>2.5</sub> concentrations around 8-9 a.m., and a second maximum  
416 during nighttime hours. This morning maximum is not present in the observations. The model  
417 ensemble is overestimating the observations in the morning and underestimating them in the early  
418 afternoon, resulting in a diurnal variability of the BIAS, shown in the lowest panel. Again, this  
419 might be related to the adopted diurnal profiles of the anthropogenic emission sources or might be  
420 due to errors in the formulation of vertical mixing in the PBL.

421

422

423

#### 424 4. The impact of missing model data on the ensemble performance

425 To assess the impact on the ensemble forecast of occasionally missing results from one or several  
426 models, we compare the following ensembles during a given test period (1-30 May 2017),  
427 separately for O<sub>3</sub>, NO<sub>2</sub> and PM<sub>2.5</sub>: This approach has already been adopted by *Marécal et al., 2015*,  
428 to evaluate European air quality predictions. We consider the following cases:

429

430 - “MEDIAN 7”, the median provided by the operational ensemble method, which includes all seven  
431 models;

432 - “MEDIAN 5”, the median built on five individual models, excluding the “best” and the “worst”  
433 models;

434 - “MEDIAN 3”, the median built on three individual models, excluding the two “best” and the  
435 “two” worst models;

436 - “BEST”, the model with the highest performance;

437 - “WORST”, the model with the lowest performance.

438

439 Since the relative performance of individual models varies in time and space, the criterion to order  
440 the seven individual models from “worst to best” is provided by the value of their respective RMSE  
441 over the test period. For ozone, the criterion is measured by the RMSE over the 30 days between  
442 12:00 and 18:00 LST (ozone peak time) (this criterion is based on the fact that the “best” model  
443 refers to the best forecast of daytime ozone levels). RMSE is seen as the most objective criterion  
444 since MB and MNMB can include compensating effects.

445

446 Figure 9 shows the statistical indicators for May 2017 as a function of the forecasting time (0-23h)  
447 of the ensemble median based on all 7 models (MEDIAN7, shown in red), 5 models (MEDIAN5,  
448 shown in blue), and 3 models (MEDIAN3, shown in black). The results are also shown for the  
449 “best” and the “worst” model (BEST (magenta) and WORST (light blue)). For all three species, the  
450 ensemble median based on 7 models is of highest quality (based on the statistical indicators used in  
451 this analysis), and generally surpasses the results provided by the “best” model. When only 5  
452 models (excluding the best and the worst) are available to calculate the ensemble, all statistical  
453 indicators show only very small differences with the more inclusive MEDIAN7 case based on seven  
454 models. Reducing the ensemble calculation further to three models (MEDIAN3), the statistical  
455 scores degrade slightly compared to the MEDIAN7 and MEDIAN5 for all three species, but remain  
456 higher or at least similar to the score of the “best” model (BEST).

457

458 It is interesting to note that the “best” model (BEST) is not the same model for the different months  
459 that are investigated, nor the same model for all species. For example, in August 2016, the “best”



460 model for O<sub>3</sub> and PM<sub>2.5</sub> is IFS, while LOTOS-EUROS shows the best performance for NO<sub>2</sub>. In  
 461 May 2017, the best model for PM<sub>2.5</sub> is LOTOS-EUROS and the worst model is IFS, but the results  
 462 remain the same: the ensemble product performs better than (or at a similar level as) the best model.  
 463 Since the “BEST” model can change depending on time period and species, the ensemble product is  
 464 particularly valuable for the sustained quality of the forecasting system. This study shows therefore  
 465 that using the ensemble product (median) of models, even if occasionally based on fewer models, is  
 466 more useful than using a single model, even if the performance of this individual model is high. The  
 467 ensemble product is still robust compared to the observations if the output of some contributing  
 468 models is occasionally missing. It also shows that an ensemble product remains valuable even if  
 469 only few models are available for the production of the forecast.

470  
471

## 472 5. Performance of the Forecasting System for Alert Warnings

473 The prediction system has been designed to support the development of policies and the calculation  
 474 of air quality indexes. One of the applications of the system is to provide alerts to the general public  
 475 when acute air pollution episodes are expected. Thus, the performance of the forecast system has  
 476 been tested regarding the likelihood to predict air pollution events. We will refer to this type of  
 477 forecast as binary prediction of events (*Brasseur and Jacob, 2017*).

478

479 A model prediction of a specific event such as an air pollution episode at a given location (e.g.  
 480 concentration of pollutants exceeding a regulatory threshold) is evaluated by considering a binary  
 481 variable and by distinguishing between four possible situations: (1) the event is predicted and  
 482 observed, (2) the event is not predicted and not observed, (3) the event is predicted but not  
 483 observed, (4) the event is not predicted but is observed. Cases (1) and (2) are regarded as successful  
 484 predictions (hits), while (3) and (4) are considered to be failures (misses). The skill of the model for  
 485 binary prediction (event or no event) is measured by the fractions of observed events that are  
 486 correctly predicted (probability of detection (POD)). The fraction of predicted events, that did not  
 487 occur is measured by the false alarm rate (FAR)).

488

489 We have calculated the POD and the FAR for the ensemble median for the cities of Beijing,  
 490 Shanghai and Guangzhou between April 2016 and June 2017, specifically for ozone (based on the 8  
 491 hour and the daily maximum value), NO<sub>2</sub> and PM<sub>2.5</sub>. The air quality indexes are calculated for 1)  
 492 1-hour ozone, 2) 8-hour ozone concentrations 3) 24-hour mean NO<sub>2</sub> concentrations, 4) 1-hour NO<sub>2</sub>  
 493 concentrations and 5) 24-hour PM<sub>2.5</sub> concentrations. The definitions breakpoints for the individual  
 494 air quality indexes (AQI) are shown in Table 1 and Table 2; they are based on current definitions of  
 495 AQI from the Chinese government.

496

497

498 **Table 1:** Chinese AQI categories

499

Index values	AQI levels	AQI categories
0-50	1	Good
51-100	2	Moderate
101-150	3	Lightly polluted
151-200	4	Moderately polluted
201-300	5	Heavily polluted
>300	6	Severely polluted



500  
501  
502  
503

**Table 2:** Individual AQI for 1-hour and 8-hour Ozone, 24-hour and 1-hour NO<sub>2</sub> and 24-hour PM<sub>2.5</sub>

IAQI	1-hour O <sub>3</sub> [μg m <sup>-3</sup> ]	8-hour O <sub>3</sub> [μg m <sup>-3</sup> ]	24-hour NO <sub>2</sub> [μg m <sup>-3</sup> ]	1-hour NO <sub>2</sub> [μg m <sup>-3</sup> ]	24-hour PM <sub>2.5</sub> [μg m <sup>-3</sup> ]
0	0	0	0	0	0
50	160	100	40	100	35
100	200	160	80	200	75
150	300	215	180	700	115
200	400	265	280	1200	150
300	800	800	565	2340	250
400	1000	Use hourly	750	3090	350
500	1200	Use hourly	940	3840	500

504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514

In order to highlight the presence of thresholds violated during the time period under consideration, Figure 10-12 show the time series for the period April 2016 – July 2017 of the 1) daily maximum ozone concentrations, 2) 8-hour moving average of ozone, 3) the 24-hour mean NO<sub>2</sub> concentrations, 4) the daily maximum NO<sub>2</sub> concentrations and 5) the 24-hour mean PM<sub>2.5</sub> concentrations for Beijing (Figure 10), Shanghai (Figure 11) and Guangzhou (Figure 12) derived from the model and from the observations at each location. Pink lines indicate the thresholds for the air quality indexes for moderate (line), lightly polluted (dashed line) and moderately polluted (dotted line) conditions for each pollutant.

515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526

In Beijing and Shanghai, the daily maximum ozone concentrations exceeded the thresholds of 160 (moderate) and 200 (lightly polluted) within the considered time period only during the months of April to September 2016. During the months of October 2016 to March 2017, the ozone concentrations remained below the threshold of 160, highlighting fair air quality conditions with regard to ozone in wintertime. In Beijing, the ensemble median has a probability of detection of air pollution events for moderate 1-hour ozone AQI of 0.44 (55 out of 126 events of 1-hour ozone breaking the threshold of 160 μg m<sup>-3</sup> have been detected). The False Alarm Rate (FAR) is 0.05 (the model ensemble predicted 58 events where ozone exceeds the threshold of 160 μg m<sup>-3</sup>, where 3 out of these 58 events were false alarm (observations below the threshold). Lightly polluted events (1-hour ozone exceeding 200 μg m<sup>-3</sup>) were correctly predicted only 14 times, while the observations exceeded the threshold 79 times. The FAR for lightly polluted ozone events is 0.12 (2 out of 16).

527  
528  
529  
530  
531  
532  
533  
534

For moderately polluted ozone events (1-hour ozone exceeding 300 μg m<sup>-3</sup>), the POD is 0, the model ensemble was not able to predict the 4 observed events (FAR is not applicable, (0 out of 0)). Looking at the 8-hour ozone predictions for Beijing, the model ensemble is very similar, with a POD of 0.45 (864 out of the 1921 observed events have been predicted correctly) and a FAR of 0.06 (56 counts are false alarm out of 920 events). For lightly polluted ozone conditions, the POD is 0.18 (118 out of 657 observed events) with a FAR = 0.06 (7 out of 125 are false alarm). For moderately polluted conditions, the model ensemble predicted 7 out of 150 observed events correctly with a FAR of 0.22 (2 out of 9 alarms are false).

535  
536  
537

For Shanghai, the PODs for ozone predictions are lower than in Beijing: for moderate air quality conditions, the POD is 0.16 (15 out of 92 observed events are predicted correctly) with a FAR of 0



538 (no false alarm) for 1-hour ozone predictions, and  $POD = 0.21$  (488 out of 2346 observed events)  
539 with a  $FAR$  of 0.01 (7 false alarms relative to 495 counts) for 8-hour ozone predictions. For lightly  
540 polluted conditions, the  $POD$  is decreasing:  $POD = 0.08$  (3 correct predictions out of 38 observed  
541 events) with  $FAR$  of 0 (no false alarm, 3 correct predictions) for 1-hour ozone, and  $POD = 0.07$  (27  
542 out of 398 observed) with a  $FAR$  of 0.10 (3 false alarms out of 30) for 8-hour ozone. For  
543 moderately polluted conditions (1-hour ozone exceeding  $300 \mu\text{g m}^{-3}$  or 8-hour ozone exceeding  $215$   
544  $\mu\text{g m}^{-3}$ ), the  $POD$  for 1-hour ozone is not applicable (no predicted, no observed events), and for 8-  
545 hour ozone  $POD = 0$  (0 predicted out of the 29 observed),  $FAR = 1$  (2 false alarms out of 2  
546 predicted, but not observed).

547

548 In Guangzhou, there is no clear difference between ozone conditions in summer or wintertime  
549 during the considered time period. Ozone observations regularly exceed the threshold of 160  
550 (moderate) and  $200 \mu\text{g m}^{-3}$  (lightly polluted) during the whole time period, and 5 times 1-hour  
551 ozone is exceeding the threshold of  $300 \mu\text{g m}^{-3}$ .

552 The  $POD$  of 1-hour ozone in Guangzhou is 0.16 (15 correct predictions out of 94 observed) with  
553  $FAR = 0.21$  (4 false alarms out of 19 predicted) for moderate conditions, and  $POD = 0.03$  (1  
554 predicted out of 36 observed) with  $FAR = 0$  (0 out of 1 predicted) for lightly polluted conditions,  
555 and  $POD = 0$  (0 predicted out of 5 observed events) for moderately polluted ozone conditions. For  
556 8-hour ozone, the  $POD$  is 0.31 (315 correct predicted out of 1032 observed) with  $FAR = 0.28$  (122  
557 false alarms of 437 predicted events) for moderate conditions,  $POD = 0.06$  (12 out of 217 observed)  
558 with  $FAR = 0$  (no false alarm out of 12 predicted events) for lightly polluted ozone conditions, and  
559  $POD = 0$  (0 out of 47 observed events) for moderately polluted ozone conditions.

560

561 In general, the ability of the model ensemble to predict correctly ozone air pollution events is best  
562 for light ozone pollution, while it fails to predict correctly the ozone pollution events for moderately  
563 polluted situations. This is mostly a result of the model ensemble being too low compared to the  
564 observations. The predictions can be improved by applying a bias correction to the ozone  
565 predictions. This is investigated in the following Section 5.1.

566

567 The  $\text{NO}_2$  predictions of the ensemble median are in general too high compared to the observation,  
568 especially in Beijing and Shanghai. Especially, in summertime (June/July/August/September), the  
569 model predictions are sometimes twice as high as the observations, which might be a result of  
570 uncertainties in the emissions. In all three cities under consideration, the  $\text{NO}_2$  concentrations are  
571 only exceeding the thresholds of  $40 \mu\text{g m}^{-3}$  for 24-hour  $\text{NO}_2$  (100 for 1-hour  $\text{NO}_2$ ) and  $80 \mu\text{g m}^{-3}$  for  
572 24-hour  $\text{NO}_2$  ( $200 \mu\text{g m}^{-3}$  for 1-hour  $\text{NO}_2$ ) during the considered period (moderate and lightly  
573 polluted conditions for  $\text{NO}_2$ ). During wintertime (November/December/January), the observations  
574 are slightly higher than in summer and the ensemble system is in better agreement with the  
575 observations.

576

577 In Beijing, the  $POD$  for 24-hour  $\text{NO}_2$  is 1 (214 of 214 observed events are predicted) for moderate  
578 conditions with a  $FAR$  of 0.46 (180 false alarms relative to 394 predicted events). This indicates  
579 that  $\text{NO}_2$  is generally overestimated by the model ensemble. For lightly polluted events, the  $POD$  is  
580 0.79 (27 predicted out of 34 observed events) with  $FAR = 0.70$  (63 false alarms out of 90  
581 predicted). For the 1-hour  $\text{NO}_2$ , the  $POD$  for moderate conditions is 0.61 (36 out of 59 observed  
582 events) with  $FAR = 0.80$  (141 false alarms out of 177 predicted). For lightly polluted conditions, no  
583 events have been observed nor predicted for 1-hour  $\text{NO}_2$  in Beijing during the considered period. In  
584 Beijing, the threshold for moderately polluted  $\text{NO}_2$  conditions has not been exceeded neither by 1-  
585 hour  $\text{NO}_2$  nor by 24h-  $\text{NO}_2$  during the considered period.

586



587 In Shanghai, the numbers are very similar to those in Beijing: POD for 24-hour NO<sub>2</sub> is 1 (208 of  
588 208 observed events are predicted) for moderate conditions with a FAR of 0.42 (152 false alarms of  
589 360 predicted events). There is also a general overestimation by the model ensemble compared to  
590 the observations. For lightly polluted conditions, the POD for 24-hour NO<sub>2</sub> is 0.67 (10 out of 15  
591 observed) and a FAR of 0.86 (60 false alarms of 70 predicted), which is a clear result of the  
592 overestimated NO<sub>2</sub>. For the 1-hour NO<sub>2</sub>, the POD is 0.91 (48 predicted out of 53 observed) with a  
593 FAR of 0.70 (111 false alarms out of 159 predicted) for moderate conditions. The thresholds for  
594 lightly polluted and moderately polluted conditions for 1-hour NO<sub>2</sub> have not been exceeded in  
595 Shanghai during the considered period, but there was 1 false alarm (1 out of 1) for lightly polluted  
596 conditions.

597

598 In Guangzhou, the model ensemble and the observations for NO<sub>2</sub> are in better agreement. There is  
599 slight overestimation of the NO<sub>2</sub> concentrations from May to September 2016, and in May 2017,  
600 but in general, there is a good agreement between the model time series and the observations. The  
601 POD for 24h-NO<sub>2</sub> exceeding the threshold for moderate conditions is 0.94 (208 predicted out of 222  
602 observed) with a FAR of 0.35 (110 false alarms of 318 predicted events), for lightly polluted  
603 conditions POD is 0.56 (15 predicted out of 27 observed) with 32 false alarms out of 47 predicted  
604 events (FAR = 0.69). Stronger polluted events have not been observed nor predicted for NO<sub>2</sub> in  
605 Guangzhou. For the 1-hour NO<sub>2</sub>, 58 events have been predicted out of 76 observed for moderate  
606 conditions (POD = 0.76, FAR = 0.63 (97 false alarms out of 155 predicted). For lightly polluted  
607 conditions, there was 1 false alarm (1 out of 1), with neither observed nor correctly predicted  
608 events.

609 The thresholds for moderately polluted conditions for 24-hour NO<sub>2</sub> and 1-hour NO<sub>2</sub> have not been  
610 exceeded in Guangzhou during the considered period, no events have been predicted nor observed.

611

612 The predictions of PM<sub>2.5</sub> (24-hour PM<sub>2.5</sub>) of the model ensemble are in very good agreement with  
613 the observations in all three cities during the considered period.

614

615 In Beijing, the POD for the prediction of moderate condition for 24-h PM<sub>2.5</sub> is 0.95 (268 correctly  
616 predicted events out of 283 observed) with a FAR of 0.19 (61 false alarms out of 329 predicted  
617 events). For lightly polluted conditions, the POD is 0.76 (111 correct predicted events of 146  
618 observed events) with a FAR of 0.28 (43 false alarms for 154 predicted events). Moderately  
619 polluted PM<sub>2.5</sub> events have been correctly predicted 33 times out of 64 observed events (POD =  
620 0.52) with a FAR of 0.35 (18 false alarms out of 51 predicted events).

621

622 In Shanghai, 191 moderate condition-events for PM<sub>2.5</sub> have been correctly predicted out of 220  
623 observed events (POD = 0.87, FAR = 0.19), with 46 false alarms out of the 237 predicted events.  
624 For lightly polluted events, the POD is 0.84 (32 out of 38 observed events) with a FAR of 0.47 (28  
625 false alarms of 60 predicted events). For moderately polluted conditions of PM<sub>2.5</sub>, the POD is 0.50  
626 (3 correctly predicted events out of 6 observed) with a relatively high FAR (0.67, 6 false alarms out  
627 of 9 predicted).

628

629 In Guangzhou, the POD for moderate conditions of PM<sub>2.5</sub> is 0.85 (149 correctly predicted out of  
630 175 observed) with 65 false alarms out of 214 predicted events (FAR = 0.30). Lightly polluted  
631 events have been observed only 7 times, the ensemble median predicted 4 of them correctly (POD =  
632 0.57), but with a very high false alarm rate (16 false alarms out of 20 predicted events, FAR =  
633 0.80), this indicates a slight overestimation of the PM<sub>2.5</sub> concentrations of the models compared to  
634 the observations. In Guangzhou, no moderately polluted events of PM<sub>2.5</sub> have been observed nor  
635 predicted during the considered period.

636



637 Only in Beijing, and only with regard to 24-hour PM<sub>2.5</sub>, heavily polluted conditions have been  
638 observed and predicted during the considered period in the winter months 2016/2017: The POD is  
639 0.5 (18 correct predicted out of 36 observed events) with a FAR of 0.28 (7 false alarms out of 25).  
640

641 These investigations show, that the model ensemble is well suited to be used in air quality  
642 predictions of PM<sub>2.5</sub>. For ozone, due to biases of the model ensemble compared to observations,  
643 the model ensemble is not able to predict ozone pollution in an appropriate way. Although the FAR  
644 is very low for ozone predictions, the POD of model ensemble is not very high. In the following  
645 Section, we apply bias correction to improve the predictions for ozone pollution events.  
646

## 647 5.1 Bias Correction for Ozone Predictions

648 Bias corrections can be applied to improve the predictions of an individual model or a model  
649 ensemble. In our case, we have calculated the summertime bias of the time series of the hourly  
650 ozone concentrations from the model ensemble with respect to the hourly observations, and  
651 subtracted the bias from the hourly time series. For predictions of ozone air pollution, the  
652 summertime is an appropriate season to consider since the ozone thresholds are exceeded only  
653 during this season. As the bias between the observations and the model might not be the same for  
654 each month, and our goal is to obtain the best improvement in the ozone predictions for  
655 summertime, we have subtracted the mean summertime bias (mean of the bias of  
656 June/July/August/September 2016) from the original time series. The daily maximum ozone values  
657 and the 8-hour moving average for the corrected time series have then been calculated. The  
658 resulting, POD and FAR for 1-hour ozone and 8-hour ozone under different air quality conditions  
659 are shown in Table 3. This table shows that, for bias-corrected predictions, the POD in all three  
660 cities is larger than for the non-corrected time series, especially in the case of moderate and lightly  
661 polluted conditions of ozone. Thus, the predictions of air pollution events are significantly  
662 improved when the bias correction is applied in the case of ozone. Only for the predictions of  
663 moderately polluted conditions of ozone, the POD is not changing. The FAR is also slightly  
664 decreasing for all cities, but the improvement is small.  
665

666 In Beijing, the POD air pollution events represented by a moderate AQI for 1-hour ozone increased  
667 from 0.44 for Beijing (55 out of 126 observed events) before bias correction to 0.69 (87 out of 126  
668 events) after bias correction. The False Alarm Rate (FAR) also increased from 0.05 (3 false alarms  
669 out of these 58 events) to 0.10 (10 false alarms out of 97 predicted events). Lightly polluted events  
670 (1-hour ozone exceeding  $200 \mu\text{g m}^{-3}$ ) have been predicted correctly 31 times (14 times without the  
671 corrections), while the observations exceeded the threshold 79 times. The FAR for lightly polluted  
672 ozone events also slightly increased from 0.125 (2 out of 16) to 0.2 (8 false alarms out of 40).  
673

674 For moderately polluted ozone events (1-hour ozone exceeding  $300 \mu\text{g m}^{-3}$ ), the POD for the bias-  
675 corrected prediction is still 0. The model ensemble was not able to predict the 4 observed events  
676 (FAR is not applicable, (0 out of 0)).  
677

678 Looking at the 8-hour ozone predictions for Beijing, the POD of 0.45 (864 out of the 1921 observed  
679 events have been predicted correctly) increased to 0.76 (1452 out of 1921) after bias corrections,  
680 and the FAR from 0.06 (56 counts are false alarm out of 920) to 0.23 (424 false alarms out of 1876  
681 predictions) for moderate ozone pollution. For lightly polluted ozone conditions, the POD increased  
682 to 0.44 (291 out of 657) and FAR = 0.22 (81 false alarms of 372 predicted) for the bias corrected  
683 predictions compared to POD = 0.18 (118 out of 657 observed events) with a FAR = 0.06 (7 out of  
684 125 are false alarm). For moderately polluted conditions, the model ensemble with bias corrected



685 predicted 27 (instead of only 7) out of 150 observed events correctly with a FAR of 0.28 (13 false  
686 alarms of 47 predictions) compared to FAR of 0.22 (2 out of 9 are false alarm).

687

688 For Shanghai, for moderate air quality conditions of ozone, the POD increased from 0.16 to 0.51  
689 (47 (15 for non-corrected) out of 92 observed events are predicted correctly); the FAR increased  
690 from 0 (no false alarm) to 0.10 (5 false alarms out of 52) for 1-hour Ozone predictions. For 8-hour  
691 ozone predictions, the POD increased from 0.21 to 0.66 (1554 (non-corrected: 488) out of 2346  
692 observed events), the FAR increased from 0.01 (7 false alarms of 495 predicted events) to 0.32  
693 (726 false alarms of 2280 counts) for 8-hour ozone predictions. For lightly polluted ozone  
694 conditions, the POD increased from 0.08 (3 correct predictions out of 38 observed) with FAR of 0  
695 (no false alarm, 3 correct predictions) to  $POD = 0.34$  (13 out of 38) with  $FAR = 0.07$  (1 false alarm  
696 of 14 predicted events) for 1-hour ozone, and for 8-hour ozone, the POD increased from 0.07 to  
697 0.27 (109 (non-corrected: 27) out of 398 observed) and the FAR increased from 0.10 (3 false alarms  
698 out of 30) to 0.13 (16 false alarms in 125 predicted events). For moderately polluted ozone  
699 conditions, the POD for 1-hour ozone is not applicable for both non-corrected and bias-corrected  
700 predictions (no predicted, no observed events), but for the bias-corrected prediction, one false alarm  
701 is observed ( $FAR = 1$ , 1 false alarm in 1 predicted event), and for 8-hour ozone POD increased  
702 from 0 to 0.10 (3 (non-corrected: 0) predicted out of the 29 observed), the FAR decreased from 1 (2  
703 false alarms out of 2 predicted, but not observed) to 0.8 (12 false alarms of 15 predicted events).

704

705 In Guangzhou, the predictions are not as accurate as in Beijing and Shanghai, and the bias  
706 corrections result only in slight improvements of the ozone forecasts for Guangzhou. The POD of 1-  
707 hour ozone in Guangzhou increased from 0.16 to 0.32 (30 (non-corrected: 15) correct predictions  
708 out of 94 observed) and the FAR slightly increased from 0.21 (4 false alarms out of 19 predicted) to  
709 0.33 (15 false alarms out of 45 predicted events) for moderate conditions. For lightly polluted ozone  
710 conditions, the POD increased from 0.03 to 0.14 (5 (non corrected: 1) predicted out of 36 observed)  
711 and the FAR increased from 0 (0 out of 1 predicted) to 0.29 (2 false alarms of 7 predicted events).  
712 For moderately polluted ozone predictions, the POD and FAR did not change with bias corrections  
713 ( $POD = 0$  (0 predicted out of 5 observed events), FAR not applicable).

714

715 For 8-hour ozone of moderate conditions, the POD increased from 0.31 to 0.49 (508 (non-corrected:  
716 315) correct predicted out of 1032 observed) and the FAR increased from 0.28 (122 false alarms of  
717 437 predicted events) to 0.37 (296 false alarms for 804 predictions). For lightly polluted ozone  
718 conditions the POD increased from 0.06 to 0.13 (29 (non-corrected: 12) out of 217 observed) and  
719 the FAR increased from 0 (no false alarm out of 12 predicted events) to 0.19 (7 false alarms for 36  
720 predicted events). For moderately polluted ozone conditions, the POD and FAR did not change with  
721 bias corrections ( $POD = 0$  (0 out of 47 observed events), FAR not applicable).

722

723 Figure 13 a–c shows the time series of the model ensemble, the bias corrected time series of the  
724 model ensemble and the observations. For the daily maximum ozone, the bias correction results in a  
725 better agreement with the observations, which also results in better event predictions. For 8-hour  
726 ozone, there is better agreement during summertime, while during the wintertime, the bias-corrected  
727 ozone time series are too high compared to the observations (both correcting for the bias derived  
728 from the total time series, or only from the summertime time series). This shows (as we have seen  
729 in Section 3.1), that the bias is not the same during the whole year, and also that the diurnal cycle of  
730 ozone is not well captured by the model ensemble. While the bias corrected daily maximum ozone  
731 is in better agreement with the observations, the 8-hour bias corrected moving average is too high  
732 during winter time (with very low ozone concentrations). As the ozone is too low in winter to  
733 exceed the lowest threshold (moderate conditions) for air quality index calculations, this is not  
734 affecting the quality of the event prediction. A more sophisticated bias-correction (bias correction





735 with diurnal and annual variation included) could be applied to further improve the predictions,  
736 provided that a longer time series (more than one year of data) is available. The statistical bias  
737 correction can then be used for the improvement of future predictions.  
738  
739

## 740 6. Conclusions and Future Developments

741  
742 In this paper, we evaluate the forecasting system developed and implemented as part of the EU  
743 Panda and MarcoPolo projects after a little more than one year of operation. The forecasting system  
744 is based on an ensemble of seven state-of-the-art chemistry-transport models (CHIMERE, EMEP,  
745 IFS, LOTOS-EUROS, WRF-Chem-MPIM, WRF-Chem-SMS, SILAMtest). Each model is  
746 executed on a computer platform hosted by individual institutes in China and Europe. Input for  
747 meteorological forcing, emissions and boundary conditions have been carefully chosen and adopted  
748 for the specific situation of China, but vary from model to model. The forecasting system provides  
749 every day hourly forecasts for 3 days ahead for four major chemical pollutants (O<sub>3</sub>, NO<sub>2</sub>, PM<sub>10</sub> and  
750 PM<sub>2.5</sub>) together with hourly observational data provided by the Chinese observational network  
751 ([www.pm25.in](http://www.pm25.in)).  
752

753 The models, whose predictions are strongly influenced by the adopted weather forecast, reproduce  
754 in general the regional features and capture many air pollution events. In most cases, the model  
755 ensemble reproduces satisfactorily the day-to-day variability of the concentrations of the primary  
756 and secondary air pollutants and in particular, predicts the occurrence of pollution events a few days  
757 before they occur. Overall, and in spite of some discrepancies, the air quality forecasting system is  
758 well suited for the prediction of air pollution events and has the ability to be used for alert warning  
759 (binary prediction) of the general public, specifically if bias corrections are applied to improve the  
760 ozone forecasts.  
761

762 In most cases, the ensemble approach provides more accurate forecasts and reduces the  
763 uncertainties in comparison with the individual models results. The calculation of the median of all  
764 models is also relatively insensitive to model outliers, and is computationally efficient. Using the  
765 ensemble median based on all models provides the best performance for all species, as the relative  
766 performance of any individual model may vary in time, space and species. We showed, that the  
767 ensemble product, even if occasionally based on fewer models, is more useful than a single model  
768 of good quality, and that the ensemble product is still robust compared to the observations if data  
769 from some contributing models are occasionally missing.  
770

771 Despite the fact that the prediction system is in its development phase and that the resources  
772 available to improve the system are limited, the MarcoPolo and Panda forecasting system can be  
773 viewed as already quite successful. The inter-comparison presented in the companion paper by  
774 *Brasseur et al., 2018* and the present evaluation were performed to diagnose differences between  
775 models, identify problems and contribute to individual model improvements. Specifically, the  
776 underestimation of ozone under high NO<sub>x</sub> conditions and the resulting errors in the diurnal cycle of  
777 ozone need to be addressed in an effort to improve the model forecasts in China. Although major  
778 efforts are ongoing to improve emission inventories for China, the remaining uncertainties,  
779 especially in regard to local emissions, may partly explain the differences between models and  
780 observations. This is subject of further investigation. Furthermore, data assimilation of satellite and  
781 in situ observations should significantly improve the performance of the forecasting system. Finally,  
782 a more advanced approach to extract observations provided by the Chinese network is expected to  
783 improve the model-data comparison.



784 **Data Availability**

785

786 The models described here are used operationally by the participating research and service  
787 organizations involved in the present study. The data produced by the multi-model forecasting  
788 system are available from the Royal Dutch Meteorological Institute (KNMI).

789

790

791 **Acknowledgements**

792

793 The model inter-comparison presented in the present study has been conducted during a workshop  
794 organized in May 2017 by the Shanghai Meteorological Service (SMS) in China. The authors thank  
795 Dr. Jianming Xu for hosting this meeting and providing support to the participants. The ensemble of  
796 models described here has been produced under the Panda and MarcoPolo projects supported by the  
797 European Commission within the Framework Program 7 (FP7) under grant agreements n°606719  
798 and n°606953. The National Center for Atmospheric Research (NCAR) is sponsored by the US  
799 National Science Foundation.

800



801 **Table 3:** POD and FAR for Beijing, Shanghai and Guangzhou

	Probability of Detection (POD)			False Alarm Rate (FAR)		
	AQI 2 (moderate)	AQI 3 (lightly poll.)	AQI 4 (moderately poll.)	AQI 2 (moderate)	AQI 3 (lightly poll.)	AQI 4 (moderately poll.)
<b>Beijing</b>						
1-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.44 (55/126)	0.18 (14/79)	0 (0/4)	0.05 (3/58)	0.12 (2/16)	NaN (0/0)
Bias corrected 1-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.69 (87/126)	0.41 (32/79)	0 (0/4)	0.10 (10/97)	0.20 (8/40)	NaN (0/0)
8-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.45 (864/1921)	0.18 (118/657)	0.05 (7/150)	0.06 (56/920)	0.06 (7/125)	0.22 (2/9)
Bias corrected 8-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.76 (1452/1921)	0.44 (291/657)	0.23 (34/150)	0.23 (424/1876)	0.21 (81/372)	0.28 (13/47)
24-hour NO <sub>2</sub> [µg m <sup>-3</sup> ]	1 (214/214)	0.79 (27/34)	NaN (0/0)	0.46 (180/394)	0.70 (63/90)	NaN (0/0)
1-hour NO <sub>2</sub> [µg m <sup>-3</sup> ]	0.61 (36/59)	NaN (0/0)	NaN (0/0)	0.80 (141/177)	NaN (0/0)	NaN (0/0)
24-hour PM <sub>2.5</sub> [µg m <sup>-3</sup> ]	0.95 (268/283)	0.76 (111/146)	0.52 (33/64)	0.19 (61/329)	0.28 (43/154)	0.35 (18/51)
<b>Shanghai</b>						
1-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.16 (15/92)	0.08 (3/38)	NaN (0/0)	0 (0/15)	0 (0/3)	NaN (0/0)
Bias corrected 1-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.51 (47/92)	0.34 (13/38)	NaN (0/0)	0.10 (5/52)	0.07 (1/14)	1 (1/1)
8-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.21 (488/2346)	0.07 (27/398)	0 (0/29)	0.01 (7/495)	0.10 (3/30)	1 (2/2)
Bias corrected 8-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.66 (1554/2346)	0.27 (109/398)	0.10 (3/29)	0.32 (726/2280)	0.13 (16/125)	0.80 (12/15)
24-hour NO <sub>2</sub> [µg m <sup>-3</sup> ]	1 (208/208)	0.67 (10/15)	NaN (0/0)	0.42 (152/360)	0.86 (60/70)	NaN (0/0)
1-hour NO <sub>2</sub> [µg m <sup>-3</sup> ]	0.91 (48/53)	NaN (0/0)	NaN (0/0)	0.70 (111/159)	1 (1/1)	NaN (0/0)
24-hour PM <sub>2.5</sub> [µg m <sup>-3</sup> ]	0.87 (191/220)	0.84 (32/38)	0.50 (3/6)	0.19 (46/237)	0.47 (28/60)	0.67 (6/9)
<b>Guangzhou</b>						
1-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.16 (15/94)	0.03 (1/36)	0 (0/5)	0.21 (4/19)	0 (0/1)	NaN (0/0)
Bias corrected 1-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.32 (30/94)	0.14 (5/36)	0 (0/5)	0.33 (15/45)	0.29 (2/7)	NaN (0/0)
8-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.31 (315/1032)	0.06 (12/217)	0 (0/47)	0.28 (122/437)	0 (0/12)	NaN (0/0)
Bias corrected 8-hour O <sub>3</sub> [µg m <sup>-3</sup> ]	0.49 (508/1032)	0.13 (29/217)	0 (0/47)	0.37 (296/804)	0.19 (7/36)	NaN (0/0)
24-hour NO <sub>2</sub> [µg m <sup>-3</sup> ]	0.94 (208/222)	0.56 (15/27)	NaN (0/0)	0.35 (110/318)	0.68 (32/47)	NaN (0/0)
1-hour NO <sub>2</sub> [µg m <sup>-3</sup> ]	0.76 (58/76)	NaN (0/0)	NaN (0/0)	0.63 (97/155)	1 (1/1)	NaN (0/0)
24-hour PM <sub>2.5</sub> [µg m <sup>-3</sup> ]	0.85 (149/175)	0.57 (4/7)	NaN (0/0)	0.30 (65/214)	0.80 (16/20)	NaN (0/0)



802

803

804

**Table 4:** POD and FAR for PM<sub>2.5</sub> for Beijing under heavily polluted conditions.

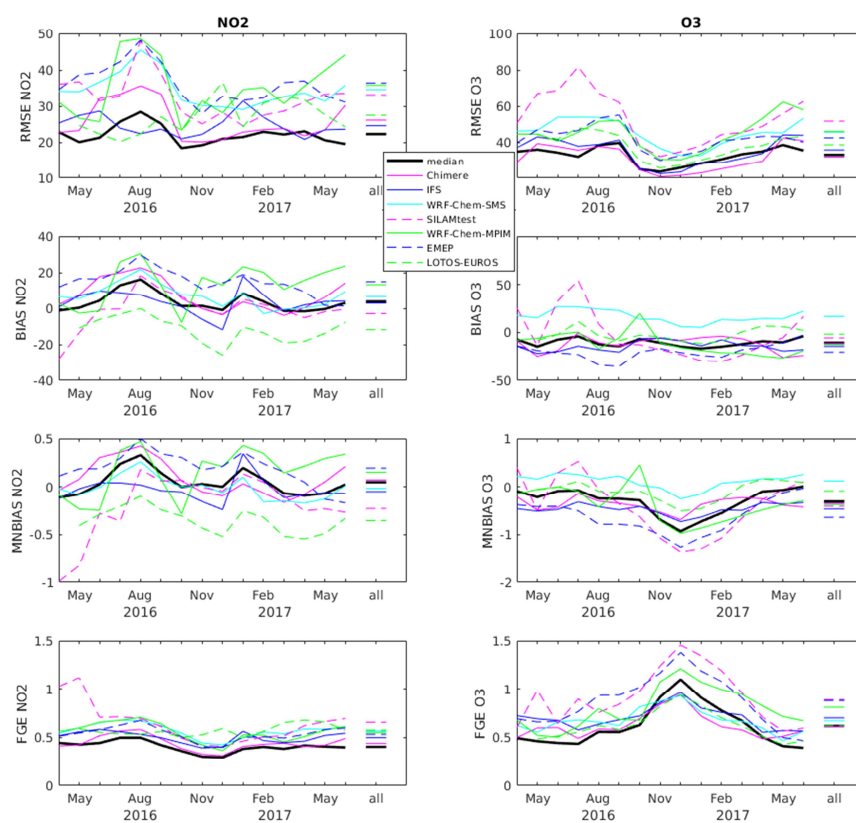
Beijing AQI heavily polluted	POD	FAR
24-hour PM <sub>2.5</sub> [ $\mu\text{g m}^{-3}$ ]	0.50 (18/36)	0.28 (7/25)

805

806

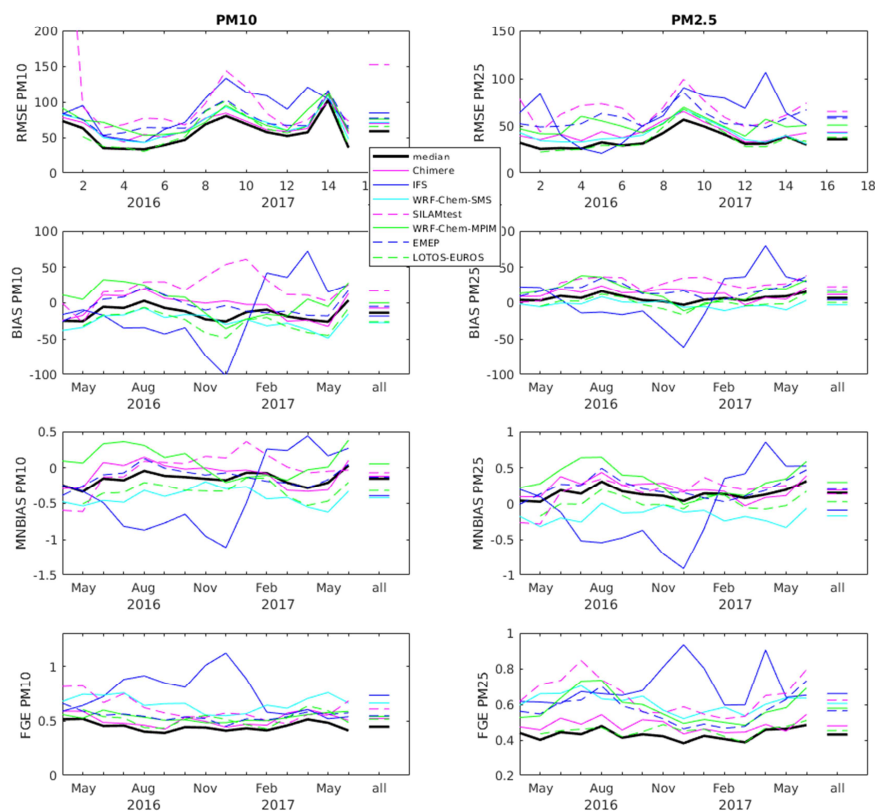


807

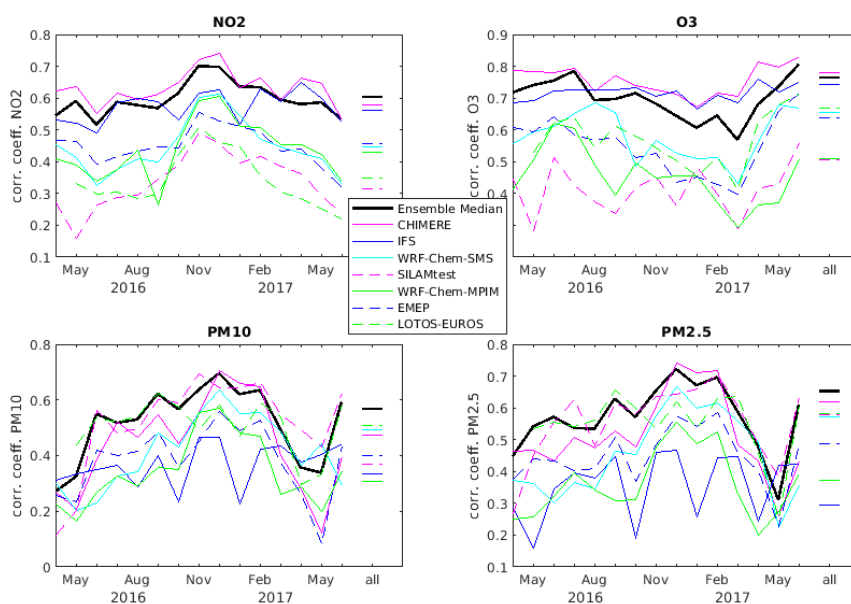


808  
809  
810  
811  
812  
813

Figure 2: RMSE, BIAS, MNBIAS and FGE of NO<sub>2</sub> and O<sub>3</sub> for each month and for the entire time period (April 2016 – June 2017, lines on the right side of each panel).

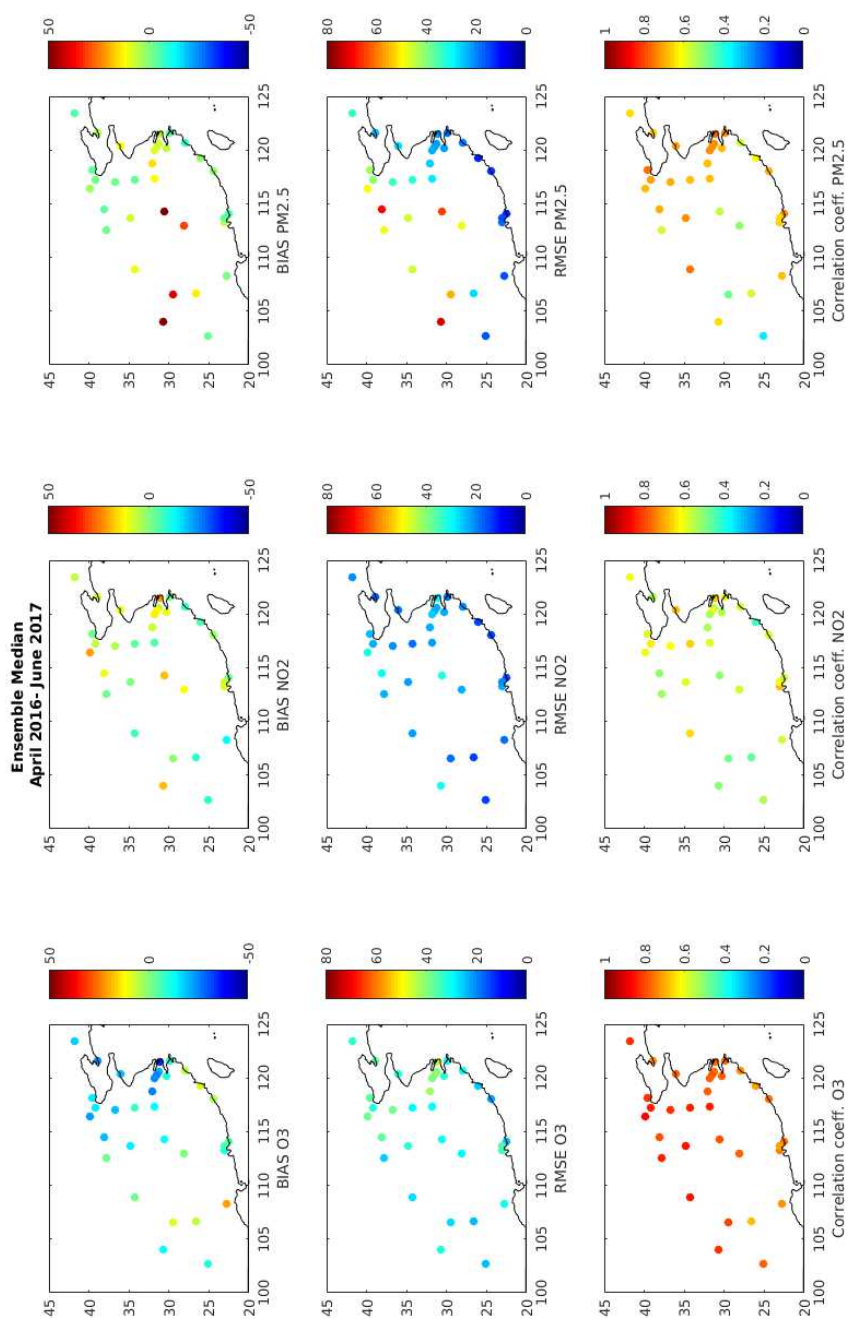


814  
815 *Figure 3: RMSE, BIAS, MNBIAS and FGE of PM10 and PM2.5 for each month and for the entire*  
816 *time period (April 2016 – June 2017, lines on the right side of each panel).*  
817



818  
819 *Figure 4: Correlation coefficients based on hourly concentrations of  $NO_2$ ,  $O_3$ ,  $PM_{10}$  and  $PM_{2.5}$  for*  
820 *each month and for the entire time period between April 2016 and June 2017 (lines on the right*  
821 *side of each panel).*

822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842

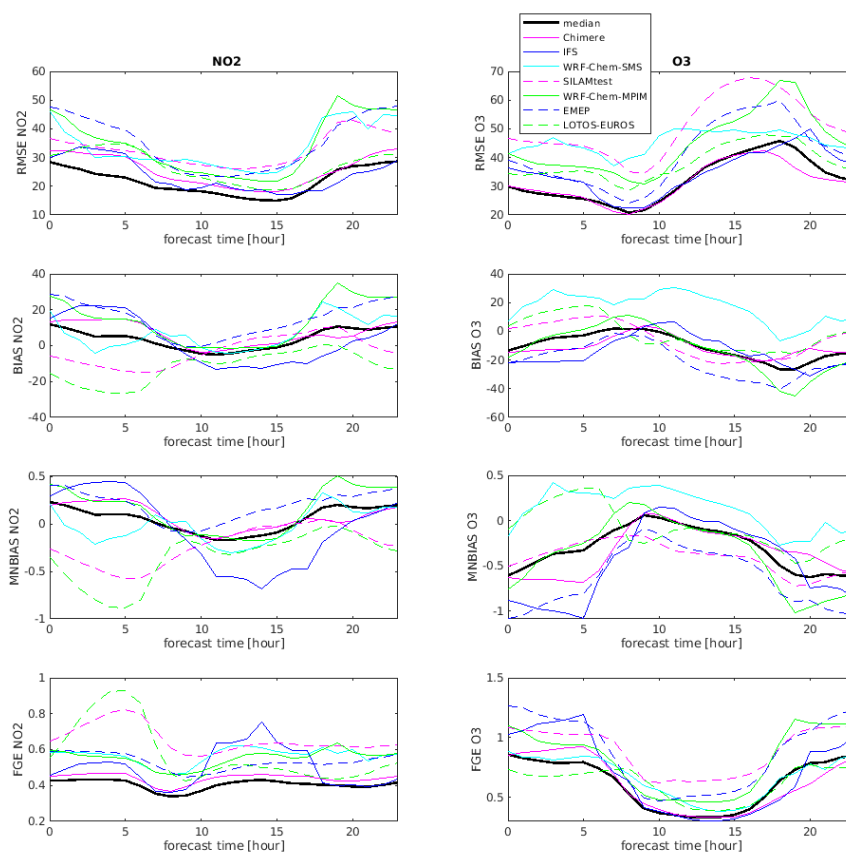


887

888

889 *Figure 5: Map of the BIAS, RMSE and temporal correlation coefficient of O<sub>3</sub>, NO<sub>2</sub> and PM<sub>2.5</sub> for*  
890 *the whole time period (April 2016 until June 2017) for each city.*



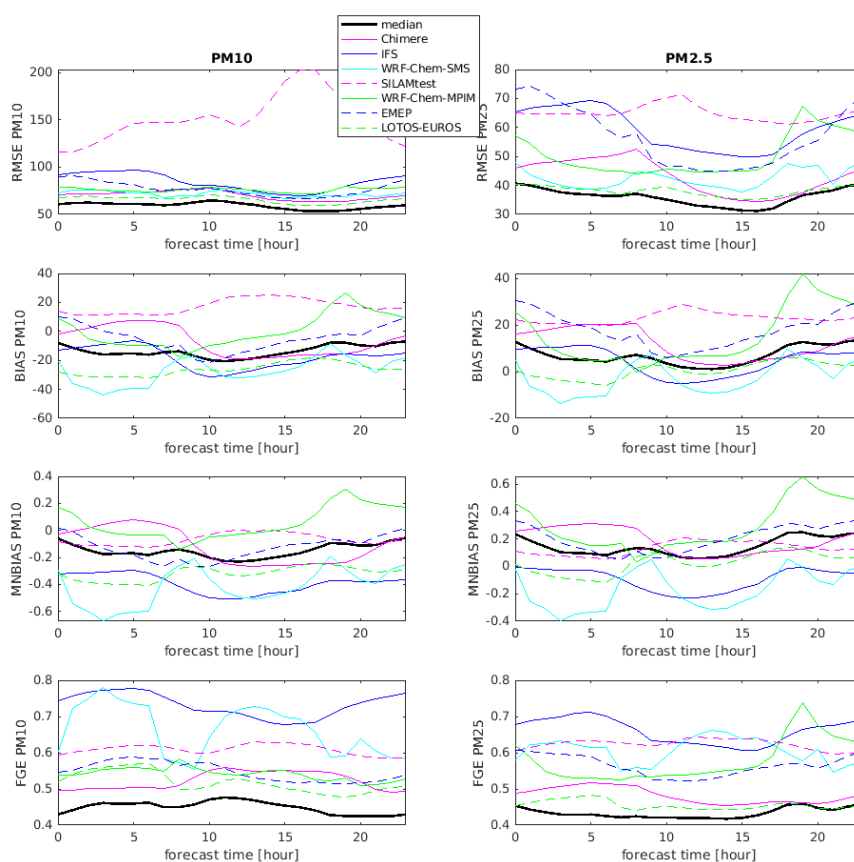


891  
892 *Figure 6: RMSE, BIAS, MNBIAS and FGE of NO<sub>2</sub> and O<sub>3</sub> over the forecasting time (time of the*  
893 *day).*

894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910



911  
912

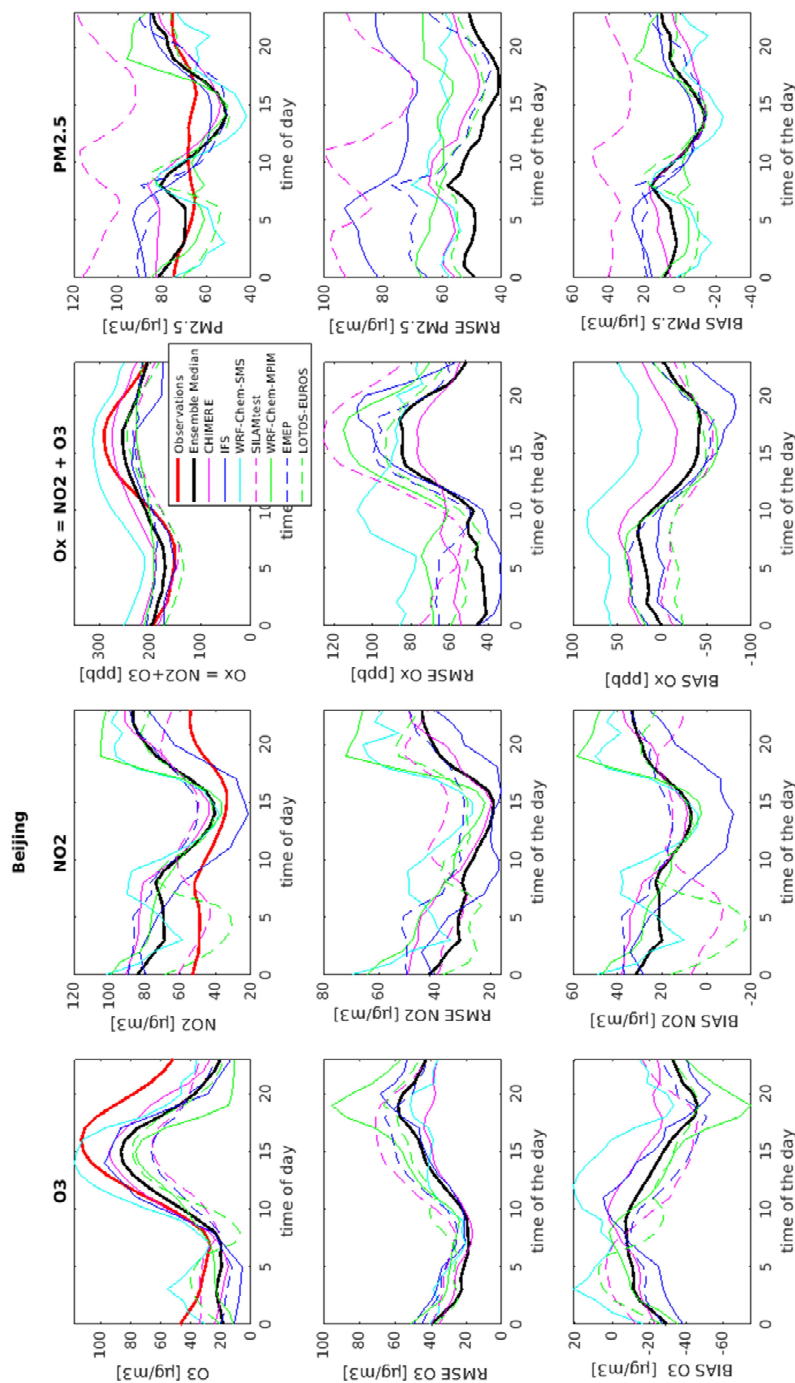


913  
914 *Figure 7: RMSE, BIAS, MNBIAS and FGE of PM10 and PM2.5 over the forecasting time (time of*  
915 *the day).*

916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930



931

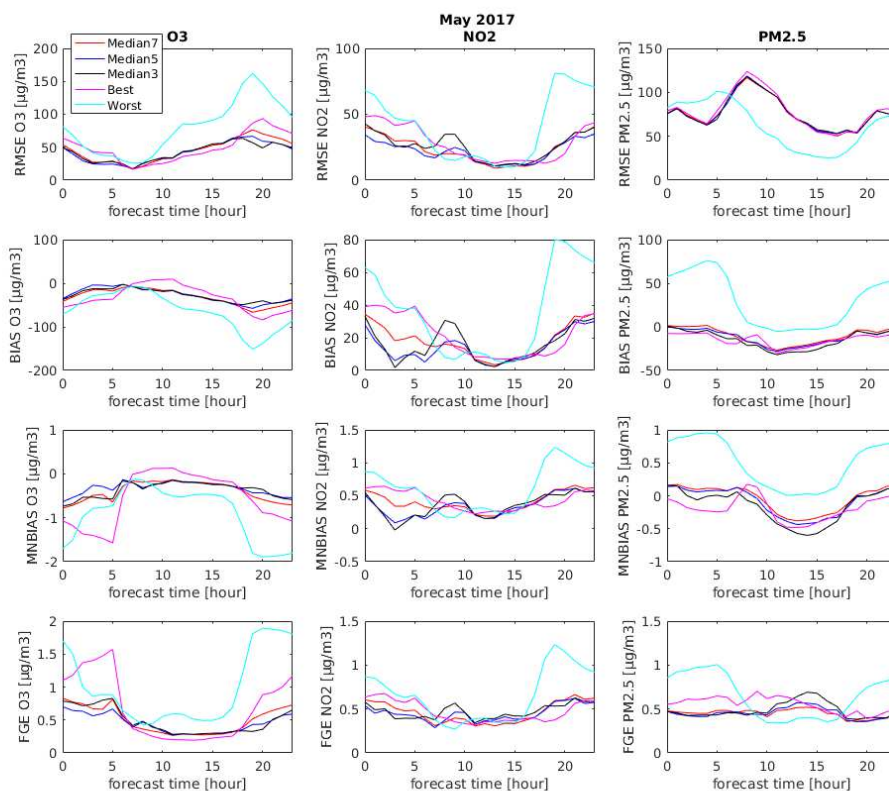


932  
 933  
 934

Figure 8: Diurnal variations of the concentrations and of the RMSE and BIAS of  $O_3$ ,  $NO_2$ ,  $O_X$  and  $PM_{2.5}$  for Beijing for the whole time period (April 2016 – June 2017).

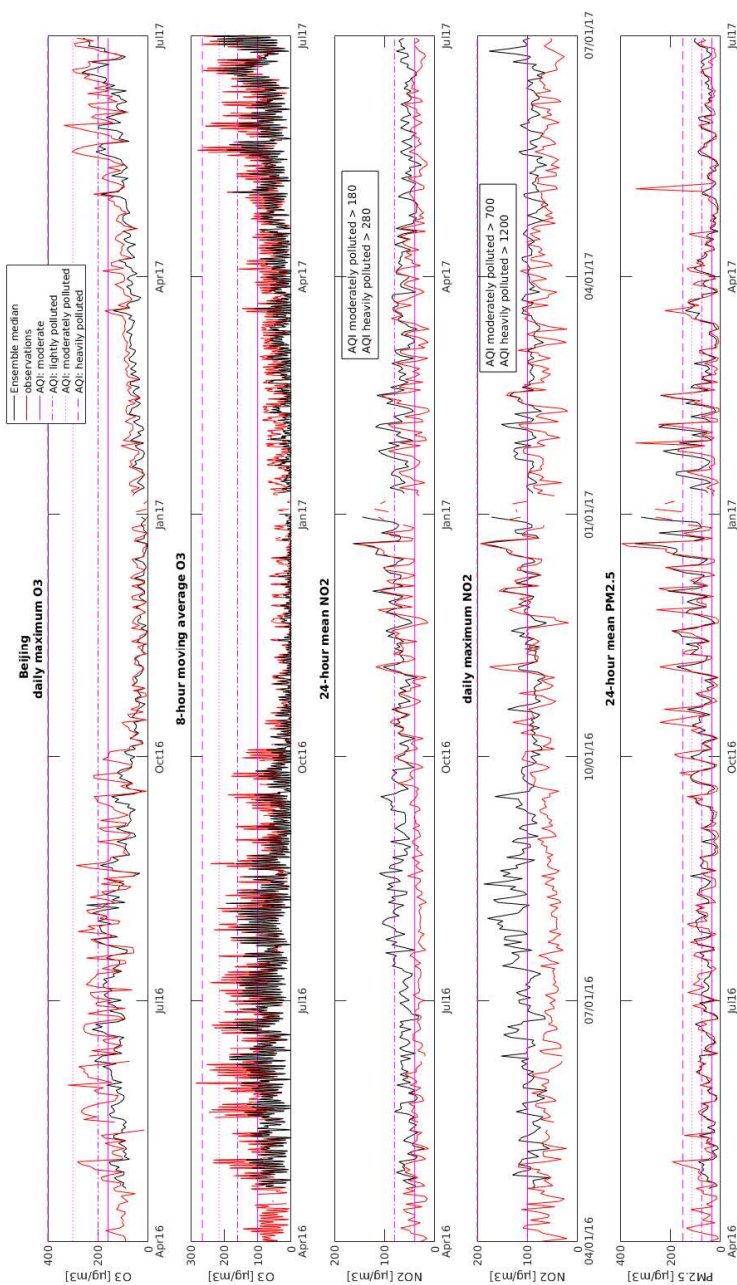


935  
936



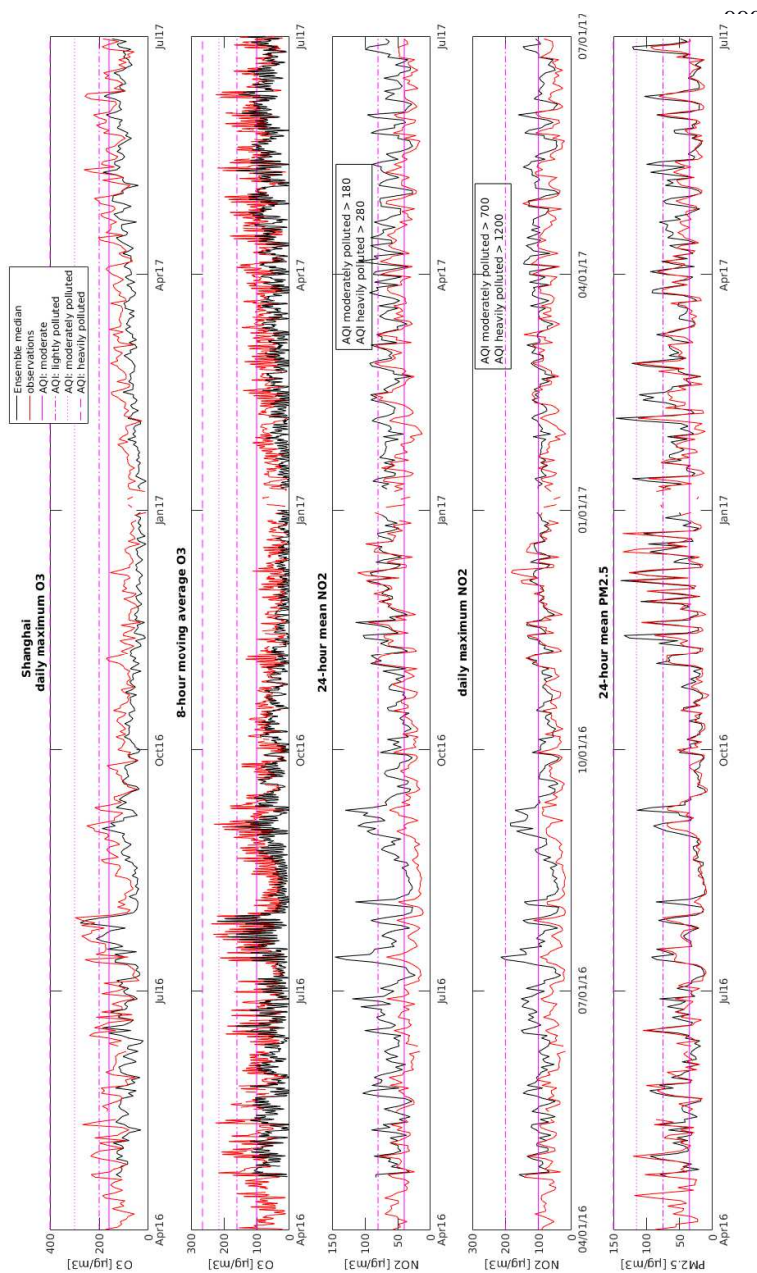
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951

Figure 9: RMSE, BIAS, MNBIAS and FGE of  $O_3$ ,  $NO_2$  and  $PM_{2.5}$  over the forecasting time (time of the day) for the Median7, Median5, Median3 and the best and the worst model.



996

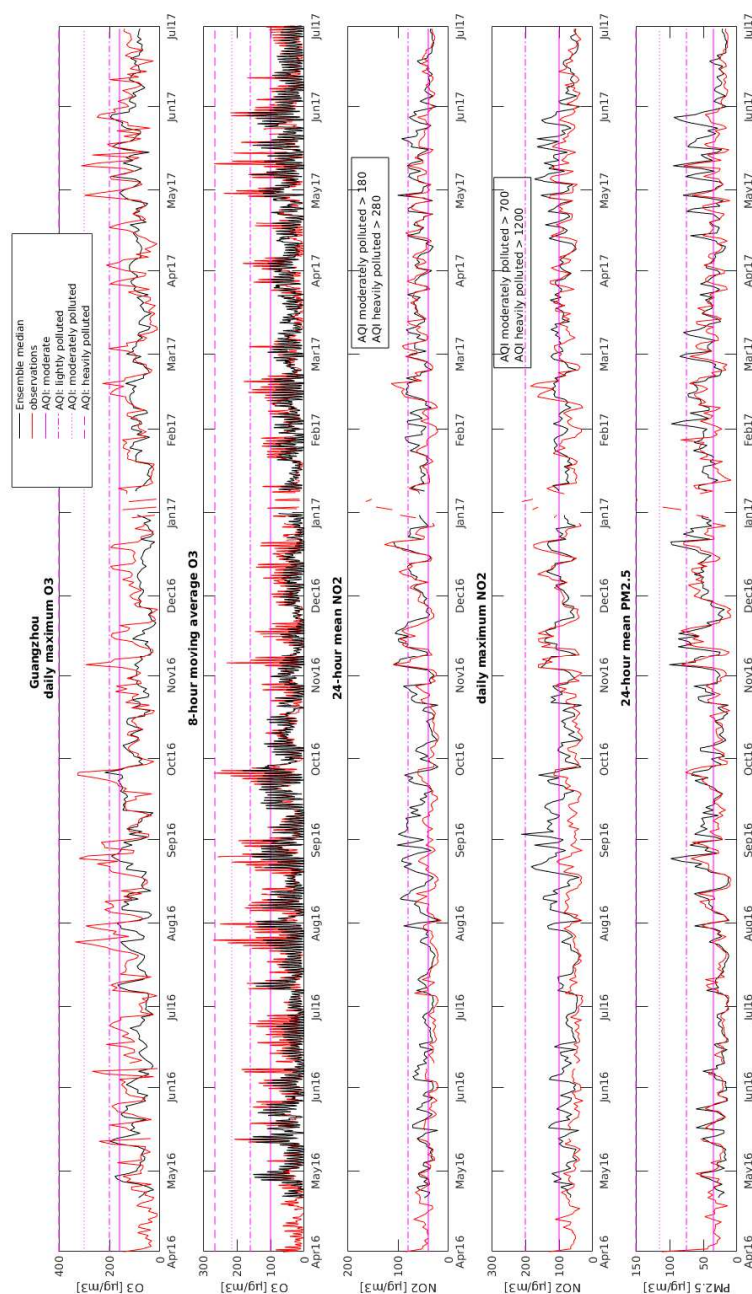
997 *Figure 10: Timeseries of daily maximum O<sub>3</sub>, 8-hour moving average O<sub>3</sub>, 24-hour mean NO<sub>2</sub>, daily*  
 998 *maximum NO<sub>2</sub> and 24-hour mean PM<sub>2.5</sub> for Beijing from April 2016 until June 2017.*



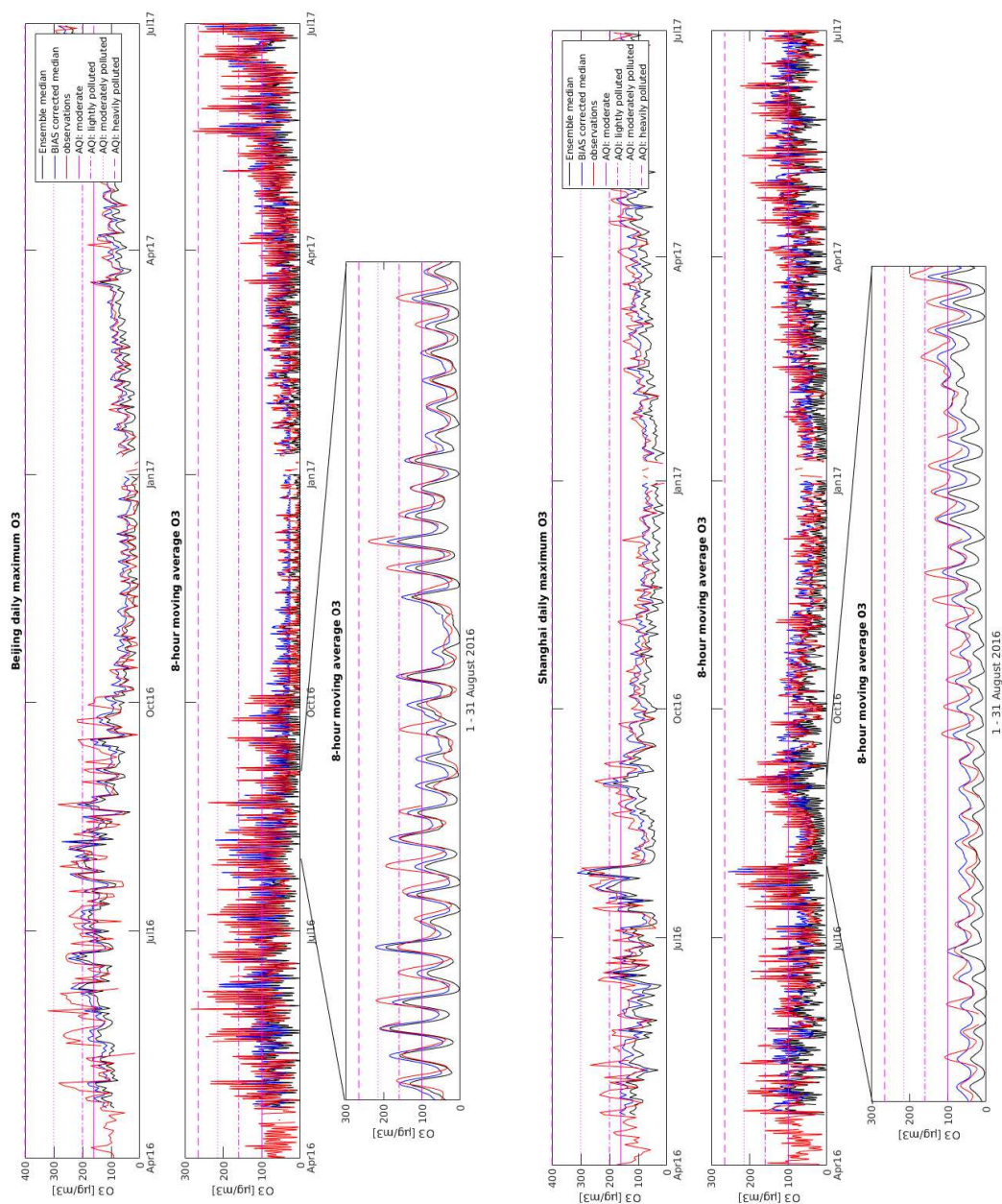
1042

1043  
 1044  
 1045  
 1046

Figure 11: Timeseries of daily maximum  $O_3$ , 8-hour moving average  $O_3$ , 24-hour mean  $NO_2$ , daily maximum  $NO_2$  and 24-hour mean  $PM_{2.5}$  for Shanghai from April 2016 until June 2017.

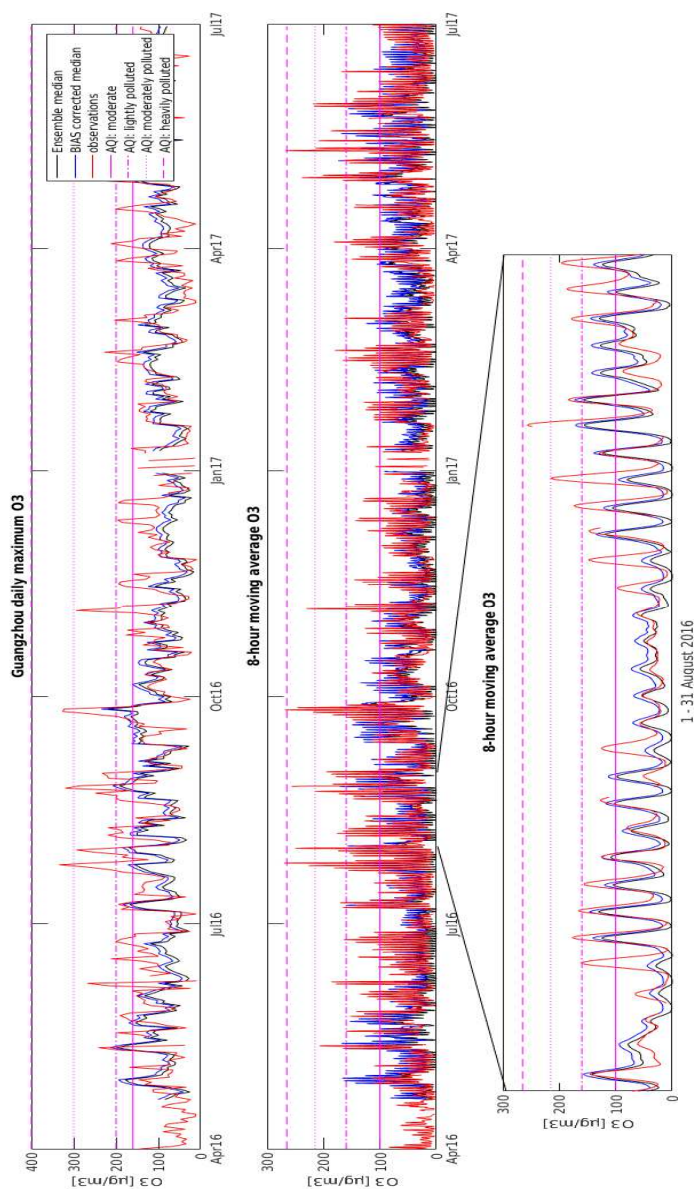


1091 *Figure 12: Calculated (ensemble median) and observed timeseries of daily maximum O<sub>3</sub>, 8-hour*  
 1092 *moving average O<sub>3</sub>, 24-hour mean NO<sub>2</sub>, daily maximum NO<sub>2</sub> and 24-hour mean PM<sub>2.5</sub> for*  
 1093 *Guangzhou from April 2016 until June 2017.*  
 1094



1095 *Figure 13 a and b: Timeseries of calculated (ensemble*  
 1096 *median) and observed daily maximum and 8-hour moving average O<sub>3</sub> for Beijing and Shanghai*  
 1097 *together with the bias corrected calculated timeseries.*  
 1098  
 1099





1141  
 1142 *Figure 13 c: Timeseries of calculated (ensemble median) and observed daily maximum and 8-hour*  
 1143 *moving average O<sub>3</sub> for Guangzhou together with the bias corrected calculated timeseries.*  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149



1150

1151 **References**

1152

1153

1154 Akimoto, H., Global air quality and pollution, *Science*, 302(5651):1716-9, 2003.

1155

1156 Ashmore, M. R. (2005), Assessing the future global impacts of ozone on vegetation. *Plant, Cell &*  
1157 *Environment*, 28: 949–964. doi:10.1111/j.1365-3040.2005.01341.x

1158

1159 Boynard, A., Clerbaux, C., Clarisse, L., Safieddine, S., Pommier, M., Van Damme, M., Bauduin, S.,  
1160 Oudot, C., Hadji-Lazaro, J., Hurtmans, D., Coheur, P.-F, First simultaneous space measurements of  
1161 atmospheric pollutants in the boundary layer from IASI: A case study in the North China Plain,  
1162 *Geophys. Res. Lett.*, 41, 645–651, doi:10.1002/2013GL058333, 2014.

1163

1164 Brasseur, G.P., Xie, Y., Petersen, A.K., Bouarar, I., Flemming, J., Gauss, M., Jiang, F., Kouznetsov,  
1165 R., Kranenburg, R., Mijling, B., Peuch, V.-H., Pommier, M., Segers, A., Sofiev, M., Timmermans,  
1166 R., Van der A, R., Walters, S., Xu, J., Zhou, G., Ensemble Forecasts of Air Quality in Eastern  
1167 China, Part 1. Model Description and Implementation, submitted to *Geosci. Model Dev*, 2018.

1168

1169 Brasseur, G. P. and D. J. Jacob, *Modeling of Atmospheric Chemistry*, Cambridge University Press,  
1170 2017.

1171

1172 Brasseur, G. P., J. Orlando, and G. Tyndall, *Atmospheric Chemistry and Global Change*, Oxford  
1173 University Press, New York, 1999.

1174

1175 Fei, Liu, Beirle, S., Zhang, Q., Van der A, R., Zheng, B., Tong, D., He, K., NOx emission trends  
1176 over Chinese cities estimated from OMI observations during 2005 to 2015, *Atmos. Chem. Phys.*, 17,  
1177 9261–9275, 2017, <https://doi.org/10.5194/acp-17-9261-2017>, 2017.

1178

1179 Fowler, D., Amann, M., Anderson, F., Ashmore, M., Cox, P., Depledge, M., Derwent, D., Grennfelt, P,  
1180 Hewitt, N., Hov, O., Jenkin, M., Kelly, F., Liss, PS, Pilling, M, Pyle, J, Slingo, J and Stevenson, D  
1181 (2008) *Ground-level ozone in the 21st century: Future trends, impacts and policy implications*.  
1182 Royal Society Science Policy Report, 15 (08).

1183

1184 Galmarini, S., Kioutsioukis, I., and Solazzo, E.: E pluribus unum\*: ensemble air quality predictions,  
1185 *Atmos. Chem. Phys.*, 13, 7153–7182, doi:10.5194/acp-13-7153-2013, 2013.

1186

1187 Guo, S., Hu, M., Zamora, M. L., Peng, J., Shang, D., Zheng, J., Du, Z., Wu, Z., Shao, M., Zeng, L.,  
1188 Molina, M. J. and R. Zhang, Elucidating severe urban haze formation in China, *Proc. Natl. Acad.*  
1189 *Sci. USA*, 111(49): 17373–17378., 2014.

1190

1191 Hamra, G. B., Laden, F., Cohen, A. J., Raaschou-Nielsen, O., Brauer, M., and D. Loomis, Lung  
1192 Cancer and Exposure to Nitrogen Dioxide and Traffic: A Systematic Review and Meta-Analysis,  
1193 *Environ Health Perspect*, 123 | 11, DOI:10.1289/ehp.1408882, 2015.

1194

1195 Huang, K., Zhuang, G., Wang, Q., Fu, J. S., Lin, Y., Liu, T., Han, L., and Deng, C.: Extreme haze  
1196 pollution in Beijing during January 2013: chemical characteristics, formation mechanism and role  
1197 of fog processing, *Atmos. Chem. Phys. Discuss.*, 14, 7517-7556, doi:10.5194/acpd-14-7517-2014,  
1198 2014.

1199



- 1200 Huang, R.-J., Y. Zhang, et al. (2014). High secondary aerosol contribution to particulate pollution  
1201 during haze events in China. *Nature* 514(7521): 218-222.  
1202
- 1203 Kampa, M., and Castanas, E., Human health effects of air pollution, *Environmental Pollution*,  
1204 151:362–367, DOI: 10.1016/j.envpol.2007.06.012, 2008.  
1205
- 1206 Leisner, C. P. and Ainsworth, E. A.: Quantifying the effects of ozone on plant reproductive growth  
1207 and development, *Global Change Biol.*, 18, 606–616, 2012.  
1208
- 1209 Marécal, V., Peuch, V.-H., Andersson, C., Andersson, S., Arteta, J., Beekmann, M., Benedictow, A.,  
1210 Bergstroem, R., Bessagnet, B., Cansado, A., Cheroux, F., Colette, A., Coman, A., Curier, R.L.,  
1211 Denier van der Gon, H.A.C., Drouin, A., Elbern, H., Emili, E., Engelen, R.J., Eskes, H.J., Foret, G.,  
1212 Friese, E., Gauss, M., Giannaros, C., Guth, J., Joly, M., Jaumouille, E., Josse, B., Kadygrov, N.,  
1213 Kaiser, J.W., Krajsek, K., Kuenen, J., Kumar, U., Liora, N., Lopez, E., Malherbe, L., Martinez, I.,  
1214 Melas, D., Meleux, F., Menut, L., Moinat, P., Morales, T., Parmentier, J., Piacentini, A., Plu,  
1215 M., Poupkou, A., Queguiner, S., Robertson, L., Rouil, L., Schaap, M., Segers, A., Sofiev, M.,  
1216 Thomas, M., Timmermans, R., Valdebenito, A., van Velthoven, P., van Versendaal, R., Vira, J.,  
1217 Ung, A., 2015. A regional air quality forecasting system over Europe: the MACC-II daily ensemble  
1218 production. *Geosci. Model Dev.* 8, 2777 e 2813. <http://dx.doi.org/10.5194/gmd-8-2777-2015>.  
1219
- 1220 Sinha, B., Singh Sangwan, K., Maurya, Y., Kumar, V., Sarkar, C., Chandra, B. P., and Sinha, V.:  
1221 Assessment of crop yield losses in Punjab and Haryana using 2 years of continuous in situ ozone  
1222 measurements, *Atmos. Chem. Phys.*, 15, 9555-9576, doi:10.5194/acp-15-9555-2015, 2015.  
1223
- 1224 Sitch, S., Cox, P.M., Collins, W.J., Huntingford, C., Indirect radiative forcing of climate change  
1225 through  
1226 ozone effects on the land-carbon sink. *Nature*. 448: 791-794, 2007.  
1227
- 1228 Sun, L., Xue, L., Wang, T., Gao, J., Ding, A., Cooper, O. R., Lin, M., Xu, P., Wang, Z., Wang, X.,  
1229 Wen, L., Zhu, Y., Chen, T., Yang, L., Wang, Y., Chen, J., and Wang, W.: Significant increase of  
1230 summertime ozone at Mount Tai in Central Eastern China, *Atmos. Chem. Phys.*, 16, 10637-10650,  
1231 <https://doi.org/10.5194/acp-16-10637-2016>, 2016.  
1232
- 1233 Wang, Y., Yao, L., Wang, L., Ji, D., Tang, G., Zhang, J., Sun, Y., Hu, B., Xin, J., Mechanism for  
1234 the formation of the January 2013 heavy haze pollution episode over central and eastern China, *Sci.*  
1235 *China Earth Sci.*, 57: 14. doi:10.1007/s11430-013-4773-4, 2014.  
1236
- 1237 WHO, <http://www.who.int/airpollution/data/cities/en/>, 2018.  
1238
- 1239 Wu, Q., Wang, Z., Chen, H., Zhou, W., Wenig, M., An evaluation of air quality modeling over the  
1240 Pearl River Delta during November 2006, *Meteorol. Atmos. Phys.*, 116: 113. doi:10.1007/s00703-  
1241 011-0179-z, 2012.  
1242
- 1243 Xu, J., Zhang, Y., Fu, J.S., Zheng, S., Wang, W., Process analysis of typical summertime ozone  
1244 episodes over the Beijing area, *Science of The Total Environment*, 399 (1–3), 147-157,  
1245 <http://dx.doi.org/10.1016/j.scitotenv.2008.02.013>, 2008.  
1246
- 1247 Zhao, X. J., Zhao, P. S., Xu, J., Meng, W., Pu, W. W., Dong, F., He, D., and Shi, Q. F.: Analysis of  
1248 a winter regional haze event and its formation mechanism in the North China Plain, *Atmos. Chem.*  
1249 *Phys.*, 13, 5685-5696, doi:10.5194/acp-13-5685-2013, 2013.



1250

1251 Zhang, B., Wang, Y., and Hao, J., Simulating aerosol–radiation–cloud feedbacks on meteorology  
1252 and air quality over eastern China under severe haze conditions in winter, *Atmos. Chem. Phys.*, *15*,  
1253 2387–2404, doi:10.5194/acp-15-2387-2015, 2015.