

Authors Response: ACP-minor revision of acp-2015-118 – “Use of North American and European air quality networks to evaluate global chemistry-climate modeling of surface ozone”

The authors thank the editor for his useful comments and suggestions. We have adopted all of the editorial and explanatory suggestions, but did not make those changes where we felt that they represented new work or were part of the authors' judgment (as in parts of the conclusion section). Below we list the comments made by the editor in black with our response in blue. Text added to the manuscript is in “**bold**” with all modifications in the tracked-changes as the end of this document.

Abstract: Having read the abstract I somehow miss at the end a strong concluding statement. You now end with still stating a result from the model evaluation but think that especially your abstract would benefit from a closing statement that makes a clear point what can be concluded from the presented model evaluation.

Yes, there should be a stronger concluding statement, thanks. We have added the following the last sentence to the abstract: “**We conclude that the skill of the models evaluated here provides confidence in their projections of future surface ozone.**”

Introduction; reading over again the statement that explains what controls surface ozone levels I was surprised to see that the term deposition was missing there. I recall that I mentioned this in my response on the initial discussions on this paper and that you also then responded indicating that you indeed were also recognizing that this is an additional relevant process for surface ozone. I see that later on in the introduction you mention the importance of surface deposition and land use changes but think that you should already include this in the listing of the important features that determine surface ozone.

We definitely agree that surface deposition is a large part of the tropospheric ozone budget. However, in the introduction, we discuss the factors driving *future* surface ozone changes. Although surface ozone deposition will be altered by future CO₂ levels and land use change, these were not included in the study here. These topics need to be noted in the introduction, and we felt that their inclusion at the end of the paragraph covered this topic sufficiently.

Description of the models, 2.2: although most readers should be familiar with the notations, to deal with acronym slang I suggest to properly introduce the terms CCMs and CTMs.

Thanks, we thought that these definitions in Section 1.0 were adequate, but have added the expansion of the acronyms in 2.2.

Results: “**The shape of the diurnal cycle of O₃ is driven primarily by sunlight, meteorology (e.g., temperature and variations in boundary layer mixing), and the daily cycle of**

precursor emissions”, also here add the deposition term. But then does this statement **“The hour of the maximum phase h occurs when these factors align, usually in midafternoon”** still holds?

Agreed, this sentence has been changed to include deposition.

“The underestimate of the summertime diurnal amplitude H by most ACCMIP models suggests that they either underestimate net daytime production or have too little nighttime loss of O₃.” I would add here that they might underestimate entrainment of free troposphere air masses or have too little nighttime loss of O₃ associated with chemical destruction or surface deposition. You discuss hereafter that the representation of the nocturnal inversion layer (I would rather call it inversion layer than PBL) is the reason for the misrepresentation of the early morning increase in O₃. However, the growth of the PBL in the early morning is one of the essential processes in 1) mixing in residual and free troposphere air masses enhanced in O₃ and this growth of the PBL also determines the efficiency of mixing of the emitted species with the depositing species and, because of that, to some extent the efficiency of the chemical production of the O₃. This could be analyzed by comparing the turbulent transport tendencies versus the chemical tendencies, diagnostics that are generally unfortunately not available in global model simulations.

To address the editor’s comment we have made the following two additions:

- 1) The above bolded sentence now reads: **“...net daytime production or have too little nighttime loss of O₃ or its precursors, through either in situ chemical loss or dry deposition.”**
- 2) We have added the following sentence: **“The mismatch of the slope in the early morning, during which the boundary layer grows rapidly, may also be caused by the models underestimating entrainment of free troposphere air.”**

This sentence reads a little weird due to the combination of all abbreviations and propose the next small addition: **“This reduces daytime production and could partly explain why ‘the models’ G and H consistently underestimate H more than others, however A overestimates H.”**

Thanks, have fixed.

Regarding the statement: **“The boundary layer schemes may be responsible for these underestimates, however, Menut et al., (2013) notes that at least for one model, increasing its vertical resolution results in very small surface O₃ changes”**, it suggests that only increasing the vertical resolution might be the way to more properly resolve the role of boundary layer dynamics and turbulent transport in properly resolving the diurnal cycles in O₃. However, if the models would all fail on properly simulating the surface energy partitioning, which might be likely with the models e.g., not considering the urban tile/surrounding areas differences in

energy balance, then increasing vertical resolution would definitely not help in improving the models performance on PBL dynamics.

This is an important scientific point, but a discussion of surface energy partitioning goes beyond the scope of this paper. Our goal was to highlight that the boundary layer schemes are wide ranging yet nearly all the ACCMIP models underestimate the amplitude of the diurnal cycle.

Page 9, lines 14-19:” **The lack of hourly variation of emissions may account for the overall underestimates of H by the ACCMIP models, since NO emissions can be lost heterogeneously at night, less effectively than those during the morning and afternoon peaks in traffic. In addition, if the early morning peak in transport NO_x emission was included, the modeled morning rise in O₃ would most likely be augmented, thus yielding larger values of H.**” Reading through this statement I fully appreciate the observation about the potential importance of missing temporal variability in the emissions that might be relevant for explaining some of the issues on the representation of the diurnal cycle in O₃. But also recognizing the fact that you deal with large-scale models that do not generally distinguish the urban tile and simply apply grid average emissions, you would not consider the representation of diurnal exchange processes which combines the role of anthropogenic and biogenic emissions but also dry deposition as a function of the large contrasts in nocturnal and daytime turbulent exchange. Those emissions are now generally assumed to be directly into the models surface layer (or above when emission height profiles are considered) whereas most emissions, with the suppressed nocturnal and early morning mixing conditions might happen in the “urban” of vegetation canopy (BVOCs, NO_x) and where consideration of such subtle features of the actual location of emissions and role of turbulent transport might result in very different diurnal cycles in the effective emissions into the surface layer. Another point of discussion is if the emissions are included in the form of NO or NO₂?

The points raised by the editor are certainly valid, but the text represents our best call as to the cause of the discrepancy and the role diurnal variations in emissions might play. This is our best judgment, and, as the editor points out, it could prove to be wrong with further study. As such, we choose to leave our conclusions as is.

Page 9; lines 28-29: “..and ‘which’ will be ‘done/conducted’ after submission of this manuscript.”

Yes, changed to: “... and **which** will be **done** after...”

Page 10, “In northern mid-latitudes, processes that drive the shape of the annual cycle are similar to those of the diurnal cycle (i.e., sunlight, temperature, and precursor emissions)”. Again making my point about dry deposition, what about the role of seasonal cycles in dry deposition due to large differences in biomass and stomatal uptake?

In the second sentence of the second paragraph of Sect. 3.2, we have added stomatal uptake as follows: “**Dry deposition through stomatal uptake**, large-scale meteorological...”

Page 11, line 5: I would suggest “The amplitude M is controlled by both meteorology and photochemistry and dry deposition”

We have chosen not to include this comment, believing that the correction above covers it.

Having read section 3.2 I realize that an essential aspect on the analysis of the 24 hour versus the MDA8 values might be the actual reference heights that are used for the models and the observations. I guess that you use the model simulated O₃ concentrations at the middle of the surface layer, which is let's typically about 30-40m height. What is the measurement height of most of the observations? If this is much closer to the surface, you might get large differences between the measurements and the models for stable nocturnal conditions. So, this difference would especially impact the 24 hour comparison whereas it would not largely affect the comparison of the daytime only data. This is pointing again at the challenge that also for analysis of these large-scale model output on atmospheric chemistry we need to integrate some of the knowledge on nocturnal exchange phenomena.

While there may (likely) be a mismatch in this, it is not clear that it is simply height. The ACCMIP models reported “surface (2m) ozone” and how they did that could be either from very thin surface layers (our CTM layer 1 is ~20 m in thickness when running the full grid) or interpolation to the surface. A bigger issue is the variations in surface elevation across a 200km x 200km grid. The measurements are not necessarily in entirely flat terrain and probably involve nighttime mesoscale circulation. This is clearly a known problem for matching these models with measurements and why we developed an optimal method for using the measurements to generate averages over the grid squares (Schnell et al, 2014).

The authors agree that a mismatch between the height at which the observations are made and what the models report as “surface” ozone may result in large measurement-model differences under stable nocturnal conditions. It is clear from Fig. 1 that the ACCMIP models have higher biases as night and thus have higher 24-hour average biases compared to MDA8. In this paragraph, we are not discussing the reasons for this since we feel we extensively covered that in the section on diurnal cycles. Instead, this paragraph aims to highlight the difference between MDA8 and 24-hour monthly average data by demonstrate the diagnostics represent different values as well as highlight that some quantities (e.g. percentiles) cannot be calculated without daily data. As such, the authors thank the editor for the comment, but have chosen to not make any modifications.

Page 15: **“but above the 60th ‘percentile’ (where UCI and observed)”**

We have made this addition.

Page 19: **“It remains unclear whether such errors result from chemical or physical processes.”** If you would indeed agree on my statements on the role of surface deposition and

boundary layer dynamics than I would suggest to modify this statement (but also previous ones) on the role chemical versus non-chemical processes in O₃ temporal variability to chemical and biogeophysical processes (to consider the role of the biosphere and PBL dynamics).

We agree with the editor's statements on the role of surface deposition and boundary layer dynamics, however, the "physical" ensures that we are not just talking about "chemical". The statement here is about the reproducibility of the scales of extreme air quality episodes and thus we do not feel additional qualifications are needed.

Section 4.1: on this discussion on what kind of output would be required to more optimally diagnose the temporal variability, having worked myself with global chemistry and climate models I used generally the output at an output frequency on the order of 7/13/23 hours or so to sample at least in a month each time of the day to construct from this a monthly mean diurnal cycle. Also as mentioned, to differentiate between the role of chemistry versus surface deposition and vertical and horizontal transport, the output of process tendencies would be optimal but am aware that this is simply too much to get from these global model simulations. You could however consider to get this extra diagnostics for some of the particular locations where such more detailed analysis could be insightful, e.g. getting the process tendencies for strong contrasting regions of the tropospheric O₃ budget.

The authors thank the editor for this comment, it is true and indeed points to the incredible difficulty of diagnosing a diurnally varying species in a CCM. As usual, however, the specification of the ACCMIP was brief and we were lucky to get the hourly output. We would have been overjoyed to get the diurnal process tendencies (our CTM does generate these, but not as diurnally resolved), and if revisited, we would hope to design a shorter simulation (e.g., one summer) where we acquire a fuller set of diagnostics. We decided here to not comment on the CMIP experiments, but are participating in the CMIP6 to see if we can get more.

1 **Use of North American and European air quality networks**
2 **to evaluate global chemistry-climate modeling of surface**
3 **ozone**

4
5 **J. L. Schnell¹, M. J. Prather¹, B. Josse², V. Naik³, L. W. Horowitz⁴, P. Cameron-**
6 **Smith⁵, D. Bergmann⁵, G. Zeng⁶, D. A. Plummer⁷, K. Sudo^{8,9}, T. Nagashima¹⁰, D.**
7 **T. Shindell¹¹, G. Faluvegi¹², and S. A. Strode^{13,14}**

8 [1]{Department of Earth System Science, University of California, Irvine, CA, USA}

9 [2]{GAME/CNRM, Météo-France, CNRS – Centre National de Recherches Météorologiques,
10 Toulouse, France}

11 [3]{UCAR/NOAA Geophysical Fluid Dynamics Laboratory, National Oceanic and
12 Atmospheric Administration, Princeton, NJ, USA}

13 [4]{Geophysical Fluid Dynamics Laboratory, National Oceanic and Atmospheric
14 Administration, Princeton, NJ, USA}

15 [5]{Lawrence Livermore National Laboratory, Livermore, CA, USA}

16 [6]{National Institute of Water and Atmospheric Research, Lauder, New Zealand}

17 [7]{Canadian Centre for Climate Modeling and Analysis, Environment Canada, Victoria,
18 British Columbia, Canada}

19 [8]{Department of Earth and Environmental Science, Graduate School of Environmental
20 Studies, Nagoya University, Nagoya, Japan}

21 [9]{Department of Environmental Geochemical Cycle Research, Japan Agency for Marine-
22 Earth Science and Technology, Yokohama, Japan}

23 [10]{Center for Regional Environmental Research, National Institute for Environmental
24 Studies, Tsukuba, Japan}

25 [11]{Nicholas School of the Environment, Duke University, Durham, NC, USA}

26 [12]{NASA Goddard Institute for Space Studies, and Columbia Earth Institute, Columbia
27 University, New York, NY, USA}

1 [13]{NASA Goddard Space Flight Center, Greenbelt, MD, USA}

2 [14]{Universities Space Research Association, Columbia, MD, USA}

3 Correspondence to: J. L. Schnell (jschnell@uci.edu)

4

5 **Abstract**

6 We test the current generation of global chemistry-climate models in their ability to simulate
7 observed, present-day surface ozone. Models are evaluated against hourly surface ozone from
8 4,217 stations in North America and Europe that are averaged over $1^\circ \times 1^\circ$ grid cells, allowing
9 commensurate model-measurement comparison. Models are generally biased high during all
10 hours of the day and in all regions. Most models simulate the shape of regional summertime
11 diurnal and annual cycles well, correctly matching the timing of hourly ($\sim 15:00$) and monthly
12 (mid-June) peak surface ozone abundance. The amplitude of these cycles is less successfully
13 matched. The observed summertime diurnal range (~ 25 ppb) is underestimated in all regions
14 by about 7 ppb, and the observed seasonal range (~ 21 ppb) is underestimated by about 5 ppb
15 except in the most polluted regions where it is overestimated by about 5 ppb. The models
16 generally match the pattern of the observed summertime ozone enhancement, but they
17 overestimate its magnitude in most regions. Most models capture the observed distribution of
18 extreme episode sizes, correctly showing that about 80% of individual extreme events occur
19 in large-scale, multi-day episodes of more than 100 grid cells. The models also match the
20 observed linear relationship between episode size and a measure of episode intensity, which
21 shows increases in ozone abundance by up to 6 ppb for larger-sized episodes. We conclude
22 that the skill of the models evaluated here provides confidence in their projections of future
23 surface ozone.

24

25 **1 Introduction**

26 We test simulated present-day surface ozone in global chemistry-climate models on temporal
27 scales from diurnal to multi-year variability and on statistics from median geographic patterns
28 to the timing and size of extreme air quality episodes. The tests use gridded hourly surface
29 ozone abundances based on a decade of observations from 4,217 air quality monitoring sites
30 in North America and Europe. Chemistry-climate models provide a valuable means for
31 projecting future air quality in a changing climate (Kirtman et al., 2013), but recent

1 assessments have lacked commensurate observational comparisons to establish their
2 credibility in reproducing current cycles in surface ozone over polluted regions (Young et al.,
3 2013). Model-measurement comparisons to date have identified model faults; yet, they often
4 have been limited to monthly statistics, biased to picking clean-air sites over limited parts of
5 the continents (Fiore et al., 2009; Doherty et al., 2013), and avoided evaluating diurnal cycles
6 and the patterns of major pollution episodes (Schnell et al., 2014, henceforth S2014).

7 The factors driving future surface ozone (O_3) changes include: (1) local-to-regional emissions,
8 (2) global-scale emissions of air pollution transported across continents and oceans, (3) global
9 emissions and physical climate change that alters the hemispheric-scale abundances of
10 tropospheric O_3 , and (4) climatic shifts in the meteorology that creates the worst pollution
11 episodes. Factors 1, 2, and 3 have been studied extensively with global chemical transport
12 models (CTMs) and chemistry-climate models (CCMs), and there is some agreement on
13 model projections given an emissions scenario (e.g., Prather et al., 2003; Reidmiller et al.,
14 2009; HTAP, 2010; Wild et al., 2012; Doherty et al., 2013; Young et al., 2013). The
15 importance of (4), however, lies in the recognition that air quality extremes (AQX), the worst
16 pollution episodes in a decade, are triggered by meteorological conditions. Air quality
17 absolute exceedances are known to occur in multi-day, spatially-extensive episodes over the
18 US (Logan, 1989; Seinfeld et al., 1991), but it was not until the regular gridding of all station
19 data over North America and Europe and the statistical definition of extremes in S2014 that
20 the extent, coherence, and decadal variability of the episodes became clear. If climate change
21 increases the duration and/or extent of the worst decadal AQX episodes, then the overall
22 health impact of poor air quality may be worse than expected based on precursor emission
23 changes alone (Fiore et al., 2012). A warming climate appears to increase the number of
24 stagnation days (Horton et al., 2014) and may decrease the frequency of ventilating mid-
25 latitude cyclones (e.g., Mickley et al., 2004), but it is unclear how these meteorological
26 indices relate to surface O_3 or particulate matter, especially with respect to the worst AQX
27 episodes as identified in S2014.

28 The models in the Atmospheric Chemistry and Climate Model Intercomparison Project
29 (ACCMIP; Lamarque et al., 2013) were used in the recent assessment of the
30 Intergovernmental Panel on Climate Change (IPCC; Kirtman et al., 2013) and represent the
31 most advanced attempt to simulate global surface O_3 in a future climate. However, in order to
32 place any confidence in their projections, their ability to simulate the observed, present-day

1 surface O₃ climatology must be evaluated. In this paper we present the first such model-
2 measurement comparisons, specifically addressing (4) by applying the methodologies from
3 S2014 to the current generation of CCMs in an effort to quantify their ability to simulate the
4 decadal statistics of the AQX episodes. Due to the complexity and nonlinearity of the
5 underlying processes, accurately simulating surface O₃ over both clean and polluted
6 environments is a formidable task for global models with resolutions of 100 km at best. For
7 example, it has been shown that choices in the parameterization of surface deposition can
8 shift modeled surface O₃ levels by ten ppb or more (Val Martin et al., 2014). Moreover, there
9 are new, phenologically-based land-surface models for interactions between atmospheric
10 chemistry and the biosphere (Büeker et al., 2012) that have yet to be fully implemented in
11 global models. In any case, the history of land-use change - both recent and future - is
12 expected to impact surface O₃ abundances (Ganzeveld et al., 2010). Thus, we recognize that
13 this model-measurement comparison is just one of the first steps in evaluating global model
14 simulations of surface O₃ pollution. A summary of the observational and model datasets as
15 well as a brief overview of the methods developed in S2014, and used here, is presented in
16 Sect. 2. Model-measurement comparisons are presented in Sect. 3 with concluding remarks
17 and further discussion in Sect. 4.

18

19 **2 Data and Methods**

20 **2.1 Observations of surface O₃**

21 We use 10 years (2000-2009) of hourly surface O₃ measurements from air quality networks in
22 North America (NA) and Europe (EU). Following S2014, in NA we use 1,633 stations from
23 the US Environmental Protection Agency's (EPA) Air Quality System (AQS), but also
24 increase the spatial coverage in NA by including 92 stations from the US EPA's Clean Air
25 Status and Trends Network (CASTNet) and 207 stations from Environment Canada's
26 National Air Pollution Surveillance Program (NAPS). The datasets used for EU remain the
27 same as S2014: 2,123 stations from the European Environment Agency's air quality database
28 (AirBase) and 162 stations from the European Monitoring and Evaluation Programme
29 (EMEP; Hjellbrekke et al., 2013). Table 1 provides a summary of the observational datasets.

30 A major advance by S2014 was the generation of average surface O₃ abundance in a grid cell
31 from observational products, one that could be directly compared to gridded model output.

1 The station measurements are used to generate a $1^\circ \times 1^\circ$ hourly grid cell average surface O_3
2 product over NA and EU using the interpolation scheme described in S2014. The
3 interpolation is similar to an inverse distance-weighted (IDW) interpolation, but additionally
4 incorporates a declustering technique employed to reduce data redundancy, similar to that of
5 Kriging (Wackernagel, 2003). The method also avoids disproportionately representing
6 stations that often are preferentially placed in the most polluted urban environments. S2014
7 first derived the maximum daily 8 h averages (MDA8) of the individual stations and then
8 interpolated onto the $1^\circ \times 1^\circ$ grid, while here we interpolate the hourly measurements and
9 subsequently derive the MDA8 at each grid cell. Differences between the two methods are
10 small (e.g., some missing station data, different 8 h periods for nearby stations), but the new
11 approach allows modeled diurnal cycles to be analyzed. The effects of (i) the new hourly $1^\circ \times$
12 1° cells being used to calculate MDA8 and (ii) the addition of CASTNet and NAPS stations
13 on the decadal 25th, 50th, and 95th percentiles at each grid cell in NA are shown in Fig. S1.
14 Overall, the difference (this work minus S2014) is about -0.6 parts per billion (ppb) O_3 for
15 each of the three percentiles. These decreases are most likely a result of deriving MDA8 from
16 the interpolated hourly abundances rather than first deriving each station's MDA8 and then
17 interpolating. Other notable changes are: the northeast edge of the domain (-5 ppb) for all
18 three percentiles due to the generally lower O_3 abundances of Canadian NAPS stations; and
19 Wyoming and Colorado at the 25th percentile (+5 ppb) possibly from CASTNet stations
20 reflecting either cumulative production of O_3 as polluted air reaches them or else more
21 prevalent stratospheric influx.

22 **2.2 Description of Models (ACCMIP + UCI CTM)**

23 The Atmospheric Chemistry & Climate Model Intercomparison Project (ACCMIP) consists
24 of 16 global models (12 [Chemistry Climate Models \(CCMs\)](#), 2 [Chemical Transport Models](#)
25 [\(CTMs\)](#), and 2 Chemistry-General Circulation Models (CGCMs)) and was designed with the
26 intent to better understand the relationships between atmospheric chemistry and climate
27 change (Lamarque et al., 2013). We focus on the *acchist* experiment, designed to test the
28 models' ability to reproduce the observed climatology of quantities specifically relevant to
29 chemistry modeling (Lamarque et al., 2013). We use the eight ACCMIP models (6 CCMs, 1
30 CTM, and 1 CGCM) with archived hourly surface O_3 , incorporating the years from each
31 model most closely aligned with observations. Most models provide 10 years of data, starting
32 in either model year 2000 or 2001. In any case, all ACCMIP simulations are climatologically

1 representative of the average 2000s with respect to meteorology and emissions. Table 2
2 provides a brief summary and the references of the models used in this study. Detailed
3 descriptions of the ACCMIP models can be found in Lamarque et al. (2013) and references
4 therein.

5 We also include a hindcast simulation over the same period as the observations from the
6 University of California, Irvine Chemical Transport Model (UCI CTM) performed at T42L60
7 resolution (Holmes et al., 2013) to both compare our model with the current generation
8 models and to highlight differences between model simulations using free-running and
9 hindcast meteorological conditions. The UCI CTM had many updates since the $1^\circ \times 1^\circ \times L40$
10 version (Tang and Prather, 2010) used by S2014, but calculates similar, not unexpectedly
11 high-biased patterns of surface O_3 .

12 For commensurate comparison of the models and measurements, we regrid the modeled
13 hourly O_3 abundances (typically at 2° to 3° resolution) to the same $1^\circ \times 1^\circ$ cells as the
14 observations using first-order conservative mapping (i.e., proportion of overlapping grid cell
15 areas). Modeled hourly abundances are adjusted by 1 h per 15° longitude to be consistent
16 with the local time of the observations. Our two major domains are: NA bounded by $25^\circ N$ -
17 $49^\circ N$ and $125^\circ W$ - $67^\circ W$; and EU bounded by $36^\circ N$ - $71^\circ N$ and $11^\circ W$ - $34^\circ E$. A further masking
18 drops coastal grid cells for which the quality of prediction index, $Q^P < 2/3$ (the number of
19 independent stations at an effective distance of 100 km used to calculate the grid-cell values),
20 see S2014 and Fig. S2 in the Supplement. Supplementary Table S1 provides the latitudes and
21 longitudes used in the final masking for both domains. Because of their differing chemical
22 regimes, some of our analyses split the NA domain into Western (WNA) and Eastern (ENA)
23 regions at $96^\circ W$, and EU into Southern (SEU) and Northern (NEU) regions at $53^\circ N$.

24 **2.3 Air quality extremes (AQX)**

25 We define air quality extreme (AQX) events on a daily basis using local (i.e., grid-cell)
26 climatologies to identify the 10 times N worst days (i.e., highest MDA8) in an N-year period
27 (i.e., the ~ 97.3 percentile; e.g., the 100 worst days in a decade). The space-time
28 connectedness of the AQX events into episodes is defined using a hierarchal clustering
29 algorithm described in S2014. Because AQX episodes span across the regions, statistics for
30 these analyses are done only on the two major domains NA and EU. The total size of an
31 AQX episode (S , units = km^2 -days) is calculated by integrating the areal extent of an episode

1 (km²) through time (days). For a given set of episodes, the mean size \bar{S} is calculated as a
2 weighted geometric mean, with the weights equal to the AQX episode sizes (Eq. 6 in S2014).
3 Because the lower native resolutions of the models typically map onto 4 to 8 contiguous 1° x
4 1° grid cells, the modeled episode sizes have artificial minimums, however, S2014
5 demonstrated that this has little effect on the resultant episode size distributions.

6

7 **3 Results**

8 **3.1 Diurnal cycles**

9 We test the models' abilities to reproduce the observed shape (i.e., phase and amplitude) of
10 the diurnal cycle, averaged over summer (JJA) and winter (DJF) months. For each of the four
11 regions, average hourly values (local solar time) are calculated as the area-weighted mean of
12 all grid cells' O₃ abundances. We calculate the phase (h , hour of peak O₃ abundance, with $h =$
13 0.0 corresponding to 00:00 local time) and peak-to-peak amplitude (H , ppb difference from
14 minimum to maximum) of the diurnal cycle using a cosine fit with a period of 24 hours.
15 Although the diurnal cycle could be more accurately represented by a higher-order fit, this
16 simple method provides objective and continuous measures of h and H for each dataset,
17 avoiding subjective, ambiguous results in cases of flat and/or multiple maxima.

18 Figure 1a-h shows the diurnal cycle of the observations and models averaged over JJA (top
19 row) and DJF (second row) in WNA, ENA, SEU, and NEU (columns from left to right). A
20 triangle for each dataset is plotted as $(x, y) = (h, H)$. The large number of data points ($\sim 10^6$ x
21 24 h per model) provides a smooth and robust estimate of each dataset's diurnal cycle. The
22 color scheme and model abbreviations in the legend of Fig. 1 are common to all similar
23 figures and text throughout. The Taylor diagrams (Taylor, 2001) in Supplementary Fig. S3a-
24 h show an alternate, commonly-used summary of the results in terms of the correlation
25 coefficient (R), the normalized standard deviation (NSD), and centered root-mean-square
26 difference (RMSD). Figures 1 and S3 show very similar quantities (e.g., model-measurement
27 discrepancies in h and H roughly correspond to R and NSD, respectively); however, we
28 consider the representation in Fig. 1 to be more useful. The panels of Fig. S3 correspond to
29 panels in Fig. 1 in terms of region and variable. Summary statistics on diurnal cycles, annual
30 cycles, and AQX events for ENA are presented in Table 3, with all regions and additional
31 statistics provided in Supplementary Tables S2-S4.

1 The shape of the diurnal cycle of O₃ is driven primarily by sunlight, meteorology (e.g.,
2 temperature and variations in boundary layer mixing), [surface deposition](#), and the daily cycle
3 of precursor emissions. The hour of the maximum phase h occurs when these factors align,
4 usually in midafternoon. Indeed, for seven of eight region-seasons in Fig. 1a-h, the observed
5 value of h ranges from 14.8 to 15.5 hours. For DJF in NEU, where photochemical O₃
6 formation is negligible, there is no obvious diurnal cycle in observations and the double
7 minimum may simply reflect the titration of O₃ from the morning and afternoon peaks in
8 transport NO_x emissions. In this case there is little information from the diurnal cycle except
9 that the amplitude H is small. The ACCMIP models, but not the UCI CTM, mostly show h
10 within ± 1 hour, generally later than observed (Tables 3 and S2).

11 Although the ACCMIP models' diurnal phase closely matches the observed, the peak-to-peak
12 amplitude H is less successfully simulated. For JJA the observed H is 27, 29, 24 and 14 ppb
13 in WNA, ENA, SEU, and NEU, respectively; while for DJF, H is 10, 9, 5, and 0.2 ppb. We
14 characterize the three largest H 's as high-photochemical region-seasons (JJA in WNA, ENA
15 and SEU), and the remaining five as low-photochemical. In this sense JJA in NEU is closer
16 to DJF in ENA in terms of near-surface O₃ production. The ACCMIP models generally
17 underestimate H by about 7 ppb in the high-three region-seasons, but cluster around H for the
18 low-five. Model A is the only ACCMIP model to overestimate H in any of the high-three,
19 possibly a result of its large total VOC (volatile organic compounds, excluding methane)
20 emissions (55% larger than the average of the other 7 models). The 24 h mean bias (MB, see
21 Tables 3 and S2) for the ACCMIP models is typically positive in all 8 region-seasons (up to
22 28 ppb), but with some models (e.g., C and E in JJA, E in DJF) showing little or no mean
23 bias, even though they underestimate H in JJA by about 25% like all ACCMIP models.

24 The underestimate of the summertime diurnal amplitude H by most ACCMIP models
25 suggests that they either underestimate net daytime production or have too little nighttime loss
26 of O₃: [or its precursors, through either in situ chemical loss or dry deposition](#). From the
27 derivative of the diurnal cycles in Fig. 1a-d, there are two periods of model-observation
28 discrepancy: in the early morning (~06:00) models underestimate the observed slope; and in
29 the early evening (~19:00) they overestimate it. The models generally match the observed
30 slope to within $\pm 1\%$ h⁻¹ during midday and throughout the night. Thus the model error is to
31 underestimate net O₃ production in the early morning and overestimate it in early evening,
32 which may be caused by the lack of a diurnal emission cycle in these global models. [The](#)

1 mismatch of the slope in the early morning, during which the boundary layer grows rapidly,
2 may be caused by the models underestimating entrainment of free troposphere air. We find
3 no clear evidence that modeling errors in the nocturnal planetary boundary layer (Lin et al.,
4 2008) or missing near-surface processes affect the diurnal cycle on a regional average.

5 Underestimated daytime production could result from limited representation of VOC
6 chemistry, since discrepancies are largest in summer when VOCs play a larger role. Indeed,
7 model A, which simulates the most chemical species of all the ACCMIP models in addition
8 having the largest VOC emissions, is one of the few models to consistently overestimate H .
9 To the contrary, however, C and E are two of the better performing models despite their
10 comparatively simple representation of VOC chemistry (C – only isoprene, E – none). The
11 only models to include the small and relatively uncertain fractional yield of HNO_3 from the
12 reaction of HO_2 and NO are A, G, and H (Lamarque et al., 2013). This reduces daytime
13 production and could partly explain why the models G and H consistently underestimate H
14 more than others, however model A overestimates H .

15 The ACCMIP models reproduce the phase of the observed diurnal cycle in both seasons
16 despite not accounting for hourly variation in emissions. The weekly, emission-driven cycles
17 in MDA8 O_3 were diagnosed by S2014, but we do not apply that diagnostic here because the
18 models did not include such variability in emissions. The lack of hourly variation of
19 emissions may account for the overall underestimates of H by the ACCMIP models, since NO
20 emissions can be lost heterogeneously at night, less effectively than those during the morning
21 and afternoon peaks in traffic. In addition, if the early morning peak in transport NO_x
22 emission was included, the modeled morning rise in O_3 would most likely be augmented, thus
23 yielding larger values of H . The ACCMIP models use a wide range of boundary layer mixing
24 schemes but consistently underestimate H . The boundary layer schemes may be responsible
25 for these underestimates, however, Menut et al., (2013) notes that at least for one model,
26 increasing its vertical resolution results in very small surface O_3 changes.

27 The UCI CTM's values of h and H show that it drastically overestimates net daytime O_3
28 production, especially during early morning hours. Its values of h are about 2.4 hours earlier
29 than observed in the high-three region-seasons, and in contrast to the ACCMIP models its H
30 values are too large by 10s of ppb. This diagnostic identifies a serious problem with the UCI
31 CTM diurnal cycle over polluted regions that needs to be investigated (e.g., missing
32 heterogeneous loss of NO_2 at night, capped boundary layer in the morning) and which will be

1 [done](#) after submission of this manuscript. S2014 found that the UCI CTM accurately hindcast
2 the summertime probability distribution of MDA8 O₃, the occurrence of AQX events, and the
3 size of these episodes, albeit with high bias of about +29 ppb in JJA over both NA and EU.
4 This new diurnal diagnostic has clearly identified model errors and pathways to improve our
5 model as well as models like G, which gravely underpredicts the amplitude of the diurnal
6 cycle. The tests shown here emphasize a large-scale average over different photochemical
7 regimes in the four regions, and thus individual model developers may wish to analyze the
8 observations for smaller regions using the datasets generated here, which are available by
9 request from the corresponding author.

10 **3.2 Annual cycle**

11 We test the models' abilities to reproduce the observed phase and amplitude of the annual
12 cycle over the four regions. Average monthly values for each region are calculated as the
13 area-weighted mean of all encompassed cells' MDA8 O₃ abundance, reflecting the EPA air
14 quality metric (www.epa.gov/air/criteria.html). Similar to the diurnal cycle, we derive the
15 phase (m , month of peak O₃ abundance, with $m = 0.0$ corresponding to 1 Jan) and peak-to-
16 peak amplitude (M , ppb difference from minimum to maximum) using a cosine fit assuming
17 12 equally spaced monthly means. Figure 1i-l shows the annual cycle of the observations and
18 models over our 4 regions with triangles plotted for each model and dataset as $(x, y) = (m, M)$.
19 The filled gray curve shows ± 1 standard deviation of each monthly mean based on 10 years of
20 observations. This interannual variability is quite narrow, much less than the spread across
21 models. As for the diurnal cycle, the Taylor diagrams in Fig. S3i-l show an alternate
22 presentation of the annual cycle results with summary statistics given in Tables 3 and S3.

23 In northern mid-latitudes, processes that drive the shape of the annual cycle are similar to
24 those of the diurnal cycle (i.e., sunlight, temperature, and precursor emissions) but occur on
25 continental to hemispheric scales. [LargeDry deposition through stomatal uptake, large-scale](#)
26 meteorological conditions including stratosphere-troposphere exchange and the position of
27 the jet stream (Barnes and Fiore, 2013) ~~can~~ also play important roles. These surface
28 observations show the same well-known cycle that has been seen in the northern hemisphere
29 mid-latitude troposphere from ozone sondes and clean-air remote sites (Logan, 1999; Fiore et
30 al., 2009): lowest values in late fall (ND), increasing through winter (JFM) followed by a
31 broad flat peak over spring-summer (AMJJA). The lower reactivity region NEU peaks in
32 April and declines until January, indicating meteorologically driven increases through the

1 winter (e.g., stratospheric influx). The observations show a phase $m = 5.6, 5.3, 5.5,$ and 4.3
2 month-of-year for WNA, ENA, SEU, and NEU, respectively; and corresponding amplitudes
3 $M = 22, 21, 26,$ and 17 ppb. By fitting a cosine curve to each grid cell's time series, we find
4 that in terms of specific locations, the earliest m occur in Canada, Florida, and NEU while the
5 latest m occur in California, south-central NA, and SEU (not shown). Most ACCMIP models
6 have m within ± 1 month of the observations, generally earlier in NEU, later in ENA and SEU,
7 and split in WNA. Models C and G have difficulty producing the observed seasonal cycles,
8 and their derived phases are not meaningful.

9 The amplitude M is controlled by both meteorology and photochemistry. For the very large
10 regional values of M , it is clearly chemical, occurring in regions with large O_3 precursor
11 emissions: California, ~ 40 ppb; the Great Lakes region ~ 30 ppb; and northern Italy, ~ 45 ppb
12 (not shown). The smallest values of M (~ 15 ppb) are found in northwest and southeast NA,
13 and NEU. The ACCMIP models generally underestimate M by about 5 ppb in WNA, SEU,
14 and NEU, while they overestimate it by about 5 ppb in ENA. The low values of M for C and
15 G suggest they are either overestimating net production of O_3 in winter or underestimating it
16 in summer, however their wintertime biases (see Fig. 1e-h, Tables 3 and S2) indicate that
17 wintertime production or representation of wintertime physical climate could be causing the
18 low M values.

19 The annual cycles here are constructed using the MDA8 O_3 derived from hourly data. Many
20 models, including 8 other ACCMIP models not analyzed here, do not report hourly surface O_3
21 but only monthly means (i.e., the average of all hours within a month). We chose MDA8
22 values to conform to the US EPA primary air quality standards and statistics, but if we used
23 monthly averages then more models could be evaluated. Unfortunately, without at least daily
24 diagnostics (e.g., daily mean or maximum value) analysis of percentile patterns and AQX
25 events and episodes (see Sects. 3.3 - 3.7) are precluded. Further, we tested the difference in
26 annual cycles diagnosed both ways and found that the bias of a model can differ and thus
27 these two diagnostics cannot be mixed. For example, the ACCMIP ensemble mean bias for
28 JJA using MDA8 averages is 2, 11, 11, and 8 ppb in WNA, ENA, SEU and NEU,
29 respectively; however, the corresponding bias using 24 h averages is consistently larger at 6,
30 14, 13, and 9 ppb. This result was expected since the ACCMIP model ensemble generally has
31 the largest biases outside of MDA8 hours. These conclusions are generally true for all

1 seasons and models, as illustrated in Supplementary Fig. S4, which shows the mean bias
2 (model minus observed) of MDA8 minus 24 h average for each model, season, and region.
3 For the UCI model, excess production in the diurnal cycle is also evident in the annual cycle,
4 overestimating M in all regions, most in ENA (+44 ppb) and least in NEU (+9 ppb). In
5 addition, the month of peak abundance is always later than observed, sometimes by more than
6 1 month. Not unexpectedly, the bias in M using 24 h averages is significantly less than that
7 using MDA8 (e.g., +30 ppb vs. +44 ppb in ENA) because largest errors occur near midday.
8 We conclude that using 24 h averages to construct the annual cycle is basically a different,
9 almost independent diagnostic than that constructed from the daily MDA8 O_3 , and further it
10 would predict different health impacts if used to project summertime surface O_3 in a future
11 climate.

12 **3.3 AQX events**

13 Next, we test the models' ability to reproduce the annual cycle of the individual AQX events,
14 identified for each grid cell as the 100 days with the highest MDA8 in the decade (40 in 4
15 years for A, 50 in 5 years for G). Figure 1m-p shows the annual cycle of AQX events for the
16 observations and models over our 4 regions. The filled gray curve shows ± 1 standard
17 deviation for each month based on 10 years of observations. The interannual variability is
18 much larger than that seen in the observed MDA8 cycle with most models falling in its range
19 in SEU and NEU, but not in WNA or ENA. An alternate presentation as Taylor diagrams is
20 shown in Fig. S3, and the summary statistics are given in Tables 3 and S4. The month of
21 maximum AQX events for most models is within ± 1 month of that observed in each region
22 (m_{AQX} in Tables 3 and S4). Based on S2014, we expect the annual cycle of AQX events to be
23 highly correlated with that of MDA8, as the observations show correlations R_{MDA8} (i.e., AQX
24 vs. MDA8) of 0.81 to 0.87 for all regions. For the ACCMIP models this correlation is not as
25 good, but they still show $R_{MDA8} > 0.70$ (Tables 3 and S4). Models whose monthly MDA8
26 correlates well with observed MDA8 also have monthly AQX events that correlate well with
27 observed. Nevertheless, matching the AQX events annual cycle is more difficult than
28 matching the cycle of MDA8 (Tables 3, S3, S4, and Fig. S3) because AQX events are driven
29 by meteorological extremes which are not necessarily represented in these climatological
30 simulations.

1 The UCI CTM also reproduces the annual AQX events well, and since it is a hindcast, we can
2 extend the analysis to how well it identifies each AQX event on an exact-match basis ('model
3 skill' by S2014). For a climatological model that exactly matches the annual cycle (i.e.,
4 matching the number of AQX events in each month) but is synoptically random in each
5 month, a skill score of ~8% is expected; but the UCI hindcast correctly identifies 28%, 33%,
6 33%, and 21% of AQX individual cell events in WNA, ENA, SEU, and NEU, respectively.

7 **3.4 Mapping O₃ percentiles and enhancements**

8 We can define baseline levels of O₃ from observations as the statistically lowest percentiles
9 (NRC, 2009). Baseline levels are independent of attribution to specific emissions or policy
10 relevance implied by US EPA's use of the term background. We can expect, or possibly
11 assume, that baseline levels are not influenced by recent, locally-emitted or produced
12 pollution (HTAP, 2010). To estimate the daytime enhancement in summertime O₃,
13 presumably caused by continental emissions, we first want to define a baseline level for each
14 grid cell as a lower percentile of the daily surface O₃. We seek a percentile that represents the
15 cleanest air possible over the summer season (even if it is never realized during the summer),
16 and one that does not change across years. We use MDA8 rather than 24 h average data to
17 prevent nighttime values from determining the baseline. We calculate percentiles for each
18 cell on an annual basis and then derive regional area-weighted averages of the percentiles.
19 The resulting percentiles by region (Fig. 2) show that the year-to-year variability is small
20 below the 40th percentile, but the largest pollution years are evident at and above the 50th
21 percentile. Thus, we select the 30th percentile as each grid cell's baseline level, which
22 corresponds roughly to the lower levels of spring-fall days. One might argue choosing, for
23 example, the 10th percentile of JJA to estimate summertime enhancement, however, this
24 assumes JJA in all models is the peak of the annual cycle and still sees clean air. We define
25 O₃ enhancement (E_X , unit = ppb) here as the difference between the 30th percentile and any
26 larger value, where subscripts will describe the reference value.

27 To estimate the summertime O₃ enhancement from local to continental-scale pollution, we
28 assume that the 92 days of JJA are the highest O₃ values of the year, pick their median value
29 (87th percentile), and subtract from it the spring-fall baseline (30th percentile). Maps of the
30 summer enhancement E_{JJA} (i.e., 87th minus 30th percentile) in NA and EU in observations
31 and models are shown in Fig. 3. While O₃ levels for the 87th, 30th, and other percentiles vary

1 considerably from cell-to-cell (see S2014), the maps of observed E_{JJA} show mostly large-scale
2 structures.

3 Many models (A, B, D, E, F, H, I) have similar patterns of E_{JJA} over NA, with large
4 enhancements (30 to 50 ppb) from the Mississippi through the Ohio River valley to the
5 Northeast, whereas the observations show such a pattern but with smaller enhancements (25
6 to 30 ppb). Model A greatly overestimates E_{JJA} in the most polluted areas (e.g., California,
7 northeast NA, south and central EU) as well as coastal areas near the Gulf of Mexico. The
8 extremely large bias near the Gulf of Mexico is unique to model A, presumably resulting from
9 natural JJA emission sources such as lightning NO_x , wildfires, or biogenic VOCs since the area
10 is not known for large anthropogenic sources. Two models (C, G) are unusually uniform
11 across NA (except California). Surprisingly, this sorting of the models does not hold for EU.
12 For example, there must be some clue as to why model B greatly overestimates E_{JJA} over NA
13 but underestimates it over EU. Such behavior from model C (uniform E_{JJA}) may be expected
14 since the tropospheric VOC chemistry is highly simplified. The uniform pattern of E_{JJA} is
15 also somewhat evident in EU for model E, which has even simpler VOC chemistry compared
16 to C, although this may be due to biases in the representation of physical climate rather than
17 chemistry.

18 The E_{JJA} diagnostic provides an excellent geographically resolved test for CCM development.
19 It also provides a useful measure of O_3 regional pollution changes in a future climate with
20 shifting O_3 baselines due to hemispheric-scale changes in methane, water vapor, temperature,
21 and stratospheric influx. Over each of our four regions, we calculate the average summertime
22 enhancement \bar{E}_{JJA} (see Tables 3 and S3), expecting to find the values and model-measurement
23 differences similar to those found in the seasonal amplitude M . Indeed, this is true, albeit E_{JJA}
24 is generally smaller than M . In addition, the spatial pattern of the values and model-
25 measurement differences are also consistent between E_{JJA} and M (not shown).

26 **3.5 AQX episode size**

27 We examine the models' ability to simulate the observed distribution of AQX episode sizes
28 over the decade 2000-2009. Our hierarchical clustering analysis identifies connected-cell,
29 multi-day AQX episodes of size S (given here in units of $10^4 \text{ km}^2\text{-days}$). We do not split the
30 NA and EU domains here because episodes span across regions. Figure 4a-b shows the
31 distribution of episode sizes in the observations and each model as the complementary

1 cumulative distribution (CCD, %), i.e., the fraction of total AQX area-day events occurring in
2 episodes of size S or larger.

3 For NA observations, the fraction of AQX area-weighted events that occur in episodes with S
4 $> 100 \times 10^4 \text{ km}^2\text{-days}$ (CCD_{100}) is 79%; and those with $S > 1000 \times 10^4 \text{ km}^2\text{-days}$ (CCD_{1000}) is
5 38%. For EU observations, most AQX events also occur in large episodes: $CCD_{100} = 80\%$
6 and $CCD_{1000} = 35\%$. Model C is aberrant in having extremely large episodes (e.g., NA
7 $CCD_{100} = 93\%$), which fall mostly in the spring rather than summer months (see Fig. 1m-p).
8 This may result from the model's simplified chemistry or unrealistic, widespread stratospheric
9 intrusion of O_3 . In any case, this model's summertime high ozone events are obscured.
10 Model A, with much more complex chemistry, however, shows significantly smaller
11 episodes. For CCD_{100} , the other models (B, D-I) are close to the observed: 73-85% for NA,
12 and 71-85% for EU. For CCD_{1000} , however, this model spread diverges substantially, 13-69%
13 for NA, and 15-70% for EU. In general, models A, B, E, and G do not produce the larger
14 episodes and thus their physical climate may lack the synoptically correlated persistent
15 stagnation episodes. The UCI CTM, using observed meteorology, captures the shape of the
16 observed CCDs extremely well compared to the free-running climate of the ACCMIP models.
17 Integrating over all episodes, we calculate the weighted geometric mean size \bar{S} (see S2014).
18 Observations have mean episode sizes \bar{S} of 415 ($10^4 \text{ km}^2\text{-days}$) and 444 in NA and EU,
19 respectively. Models C, D, F, H, and I are biased high in \bar{S} , while models A, B, E, and G are
20 biased low for both NA and EU (Fig. 4a-b, Tables 3 and S4).

21 **3.6 Non-stationarity and possible trends**

22 One problem with diagnosing decadal AQX size statistics is that they can be biased if more
23 AQX events occur at one end of the decade due to a trend in O_3 precursor emissions. A
24 greater density of events in one summer generally means larger episodes. A linear fit of
25 annually derived O_3 percentiles calculated over years 2000-2008 (2009 was excluded due to
26 lack of NO_x and VOC emission data, see below) for each of the 4 regions (Fig. S5) shows
27 clearly decreasing surface O_3 abundances at the higher percentiles (see also Fig. 2),
28 presumably through reductions in NO_x and VOC emissions (Hudman et al., 2009; Xing et al.,
29 2014). To test if these trends are emissions-driven or artifacts of the meteorological time
30 slice, we analyze the UCI CTM results (dashed lines, Fig. S5), which are forced by observed
31 meteorology but have constant anthropogenic pollution emissions over the time period. We

1 also obtain total NO_x and VOC emissions from version 4.2 of the Emission Database for
2 Global Atmospheric Research (EDGAR, EC-JRC/PBL, 2009) for years 2000-2008 (2009 was
3 unavailable at time of publication) and calculate their trends over the period. Over WNA and
4 ENA, meteorology seems to be driving the small positive trends at lower O₃ percentiles
5 (where UCI and observed trends roughly agree), but above the 60th percentile (where UCI
6 and observed trends diverge) emissions reductions are the most likely cause. In SEU and
7 NEU the trends are less conclusive for either meteorology or emission based, but most EU
8 NO_x reductions occurred prior to 2000 (Xing et al., 2015). Koumoutsaris and Bey (2012)
9 compare GEOS-Chem hindcasts with NA and EU trends at a limited number of stations from
10 CASTNet and EMEP (~40 in each domain) and find similar trends. They also attribute the
11 negative trends at high percentiles to reduced precursor emissions, however they attribute the
12 positive trends at low percentiles to changing background O₃ as opposed to changing
13 meteorology posited here.

14 In an effort to correct the AQX decadal statistics for changes in O₃ precursors, we searched
15 for correlations on a cell-by-cell basis between high-percentile MDA8 O₃ vs. NO_x emissions
16 on an annual basis for years 2000-2008. No simple linear relation emerged, and we could
17 find no satisfactory way to “correct” the observations for this regionally varying, monotonic,
18 but non-linear, decline in NO_x and VOC emissions that did not corrupt the data. The post-
19 CMIP5 plans for the Chemistry-Climate Model Initiative (CCMI) include hindcast
20 simulations with time-dependent emissions that will allow for the simulation of the observed
21 O₃ non-stationarity.

22 One option for analyzing extremes in a non-stationary decadal data set is to define AQX
23 events annually on a 10-per-year basis. This approach greatly dampens the observed episode
24 mean size and across-year standard deviation from 415 ± 307 (100 per decade) to 249 ± 67
25 (10 per year) in NA and from 444 ± 720 to 355 ± 48 in EU. Moreover, it gives a false
26 positive impression of the severity of air pollution in extreme years. Thus, we maintain our
27 primary analysis with AQX defined as 100-per-decade. In parallel with Fig. 4a-b, we show
28 the CCDs using a 10-per-year basis for AQX in Supplementary Fig. S6a-b.

29 **3.7 Severity of pollution in largest episodes**

30 As a measure of O₃ produced during AQX events/episodes, we map out the enhancement at
31 the AQX threshold level E_{AQX} (~97.3 percentile) as shown in Fig. S7 (parallel to Fig. 3, also

1 relative to the local 30th percentile). We also calculate the average AQX enhancement \bar{E}_{AQX}
2 over our regions (Tables 3 and S4). For ENA, the ACCMIP modeled range of \bar{E}_{AQX} is 29-52
3 ppb, spanning the observed of 35 ppb (Table 3). This average result is encouraging for the
4 ACCMIP models except that, as for E_{JJA} (Fig. 3), the pattern match is not as good (Tables 3,
5 S3 and S4).

6 Of the 100 AQX events in each cell, many will lie above the local AQX threshold value. We
7 expect that larger, longer-duration episodes accumulate more O_3 , and thus these super
8 episodes might have O_3 enhancements (relative to the 30th percentile) well above the AQX
9 threshold enhancement, E_{AQX} . For each AQX event, we calculate an enhancement (ppb) as
10 the MDA8 value of that AQX event minus the local 30th percentile value. For each episode
11 of size S , we calculate the area-weighted average enhancement E_S . Figure 4c-d plots the
12 observed density distribution of all E_S , quantized every 2.5 ppb for E_S and every decade in 10^4
13 km^2 -days for S . These plots show large variability in the observed E_S frequency (gray pixels),
14 but a consistent picture of the mean enhancements as a function of S (open circles). For
15 episode sizes of 0.3 (i.e., 0.1 to 0.99), 3 and 30, E_S is almost constant (~ 32 ppb for both NA
16 and EU), but for sizes 300 and 3000 it increases almost linearly per decade. We calculate this
17 slope $\Delta\bar{E}_S$ as the average of the 30-to-300 increase (1 decade in S) plus half of 30-to-3000
18 increase (2 decades), getting values of 2.9 (NA) and 1.7 (EU) ppb increase per decade of
19 episode size. Similar results are seen for the 10-per-year AQX definition (Fig. S6c-d), with
20 $\Delta\bar{E}_S$ of 2.7 (NA) and 3.3 (EU). The slope $\Delta\bar{E}_S$ is not simply an expected result from our
21 statistical sorting since in NA we find that compared to the observations, model C has slope
22 that is a factor of about 4 smaller, while A has a slope nearly a factor of 4 larger, and F has a
23 negative slope.

24 The models generally produce the shape of E_S vs. S , although most models (except A and I,
25 see Fig. 4 caption) underestimate the enhancement for all sizes. The obvious discrepancies
26 are for NA episodes, where many models predict that the largest enhancements occur in the
27 smallest episodes ($S = 0.3$). This anomaly does not occur for EU episodes. These small
28 episodes are uncommon, representing only a small fraction of events (see Sect. 3.5), and we
29 find them mostly along the coasts at the edge of the mask. We understand them to be the
30 effect of very polluted air masses being advected to the neighboring ocean cells which are
31 typically low- O_3 regions with very low 30th percentile baselines, resulting in large
32 enhancements from the highly polluted air. The observations are interpolated and not capable

1 of following a pollution shift offshore. Thus the models are probably correct, but the method
2 of masking and station interpolation makes this discrepancy a systematic feature. The lack of
3 such a feature in EU can be understood by the lack of such sharp coastal gradients. Overall,
4 most models agree with the observations, showing that the super-episodes have the largest O₃
5 enhancements.

6

7 **4 Conclusions and Discussion**

8 Confidence in modeled projections of future air quality is based fundamentally on our ability
9 to accurately simulate the present-day observed climatology of surface O₃ and particulate
10 matter over North America (NA) and Europe (EU) where dense, long-term, reliable
11 measurements are available. In this work we evaluate the surface O₃ climatologies from 8
12 global models (6 CCMs, 1 CTM, and 1 CGCM) that reported hourly surface O₃ as part of the
13 ACCMIP. In addition we test the UCI CTM simulation as an exact hindcast of the 2000-
14 2009 decade of observations used here. Our tests follow the unique approach of S2014 in
15 which over 4,000 heterogeneously spaced air quality stations are used to calculate the hourly
16 O₃ averaged over 1° x 1° grid cells that can then be compared unambiguously with the
17 modeled grid. Diagnostics include the hourly diurnal cycle, monthly seasonal cycles, and
18 sizes and intensity of air quality extreme (AQX) episodes. For the most part, the models are
19 biased high during all hours of the day, all months of the year, and in all regions.

20 Averaged over large regions, the ACCMIP models simulate the shape of the observed
21 summertime diurnal cycle well, with the hour of maximum within ±1 hours of observed
22 (~15:00). The observed peak-to-peak amplitude (25 to 29 ppb over the more polluted
23 regions) is not as well matched and typically underestimated by about 7 ppb. The UCI CTM
24 hindcast, which performed well in the S2014 tests except for a uniform high bias, clearly fails
25 these new diurnal tests and indicates model error in the morning boundary layer chemistry. In
26 general, the ACCMIP models simulate the observed regional annual cycle of monthly mean
27 MDA8 O₃. They match the month of maximum to within ±1 months of observed (mid-June),
28 although two models are in error with almost no annual cycle and no clear maximum. The
29 other models overestimate the peak-to-peak amplitude of the observed cycle by about 5 ppb
30 (20%) in the most polluted region (Eastern North America) while underestimating it by about
31 5 ppb in the other three regions. Model skill in matching the annual cycle of AQX events is
32 fair but not good. This annual cycle has much larger interannual variability than that of

1 MDA8 O₃, and many models shift the month of maximum AQX events to later in the summer
2 than is observed.

3 Measures of the enhancement in surface O₃ driven by pollution are derived from the statistics
4 of the decade of daily gridded MDA8 values. For our measure of summertime enhancement
5 (87th minus 30th percentile), the models generally replicate the observed spatial structures but
6 overestimate the magnitude in the most polluted regions. Two models are surprisingly
7 uniform across both continents and fail to highlight areas with the largest emissions of O₃
8 precursors. Typically, modeled high biases appear in the upper percentiles, not the 30th
9 percentile, which appears to be a good measure of the baseline O₃ across the decade.

10 About 80% of the AQX events in NA and EU occur in large, connected, multi-day episodes
11 consisting of 100 grid cells or more. This result is closely matched by all but two models,
12 with C producing much larger episodes and A, much smaller ones. It remains unclear
13 whether such errors result from chemical or physical processes. The observations show that
14 super-sized episodes of 100 cells or more have successively greater O₃ levels as they become
15 bigger, with the 100-times-larger episodes having 4-6 ppb greater O₃. Most, but not all of the
16 models match this increase. It is likely that larger, longer-lasting episodes allow for greater
17 accumulation of O₃ from neighboring pollution sources.

18 **4.1 What are the best air quality diagnostics for model development?**

19 For testing and identifying the model strengths and weaknesses and improving simulations of
20 air quality, modelers save a large number of diagnostics during the model development
21 process. This typical model development process is far less limiting than the experiment
22 analyzed here, which is based on the voluntary contributions of many models and many
23 terabytes of diagnostics imposed in the ACCMIP. We would still recommend saving the
24 diagnostic of hourly surface O₃ over a decade or more of simulation from which all of the
25 primary diagnostics here can be readily derived and compared with the observations. To
26 segue from the surface O₃ over NA and EU to the sondes and remote sites, a monthly
27 averaged 3-D O₃ would probably suffice. Hourly data observed at coastal or mountain sites
28 likely includes a diurnal meteorology that is not represented in the global models, even at a
29 resolution of 0.5° x 0.5°. Furthermore, the 24 h and MDA8 averages show different biases,
30 and should not be treated as the same diagnostic. There may be inventive ways to avoid the
31 massive hourly data sets by storing the diurnal cycle as a monthly mean and calculating

1 MDA8 inline or just storing the maximum daily O₃ value, which would then require similar
2 analysis of the observations.

3 The open questions are what model simulations are practical and which would be most useful
4 to identify model errors. The ACCMIP simulations forced by a decade of 2000s climate-
5 model sea surface temperatures are useful in comparing decadal statistics, but the UCI CTM
6 hindcast provides unique tests on the ability to simulate specific events and years. Even if the
7 observed sea surface temperatures were used, the synoptic extreme events would not likely
8 coincide with the observed, so a hindcast meteorology based on reanalysis for forecast fields
9 provides an important test of the model.

10 The surface O₃ data here is based on an interpolation algorithm that was optimized for the 50-
11 100 km scale averages. Thus, the supplied grid-cell averaged data could be regenerated at
12 0.5° resolution, but if one wants 10 km cell averages for regional models then the parameters
13 in the current algorithm would need to be revised and re-optimized. The surface O₃ data set
14 will be expanded to include more than 2 decades (1993-2015) and thus longer simulations
15 would be desirable to investigate interannual variability.

16 **4.2 What are the most important tests for these chemistry-climate models,** 17 **assuming that hindcasts and detailed emission data are not being used?**

18 Another major question is what emissions to use. With ACCMIP the choice of a single year
19 of representative emissions for the decade was the optimal choice. The downward trending
20 emissions in NA and EU over the 2000-2009 decade, however, created a non-stationary data
21 set. Going to a longer data set, 1993-2015, will make the comparison between models and
22 measurements more awkward. Model developers will need to take some account of this non-
23 stationarity, possibly as a sensitivity study using two different emissions sets representative of
24 the early and late periods of observations, when not tracking emission changes each year.

25 An emissions problem not resolved here is whether the modeled diurnal cycle over heavily
26 polluted regions in summer would be affected by imposing a more accurate diurnal and
27 weekly cycle in emissions. This is probably beyond what can be imposed in a MIP, but
28 should be part of the individual model development as a sensitivity assessment.

29 The four-region decadal average statistics here provide a fairly broad view of the models'
30 ability to predict the buildup of O₃ and extreme events in polluted regions. Clear examples of
31 model error are identified. The general agreement of the diurnal cycle between models and

1 measurements still needs to be tested with diurnal emissions. Going beyond the mean
2 regional cycles, the ability to test models at the grid cell level provides clear geographic
3 coverage, identifying patterns of the discrepancy that are sometimes disturbing, as shown in
4 Fig. 3, but not developed further in this paper. The next study of the CMIP-generated surface
5 O₃ needs to evaluate this.

6 **4.3 What tests provide the best confidence in model prediction of future air** 7 **quality?**

8 Accurate projections of future air quality rely on our ability to predict the changes in both
9 baseline level and pollution buildup in response to both specified future climatic conditions
10 and a change in local-to-global emissions. Both the baseline and the amount of O₃ produced
11 from pollution are likely to change and need to be assessed separately. For that purpose, we
12 find that the maps of summertime (87th percentile) and baseline (30th percentile) and their
13 difference are one of the more important tests of a model's simulation of the present-day. The
14 annual cycle of monthly means is also in some way a measure of the summertime
15 enhancement, but not as useful as the percentiles. One key measure of future change would
16 be in the size and intensity of extreme episodes. The intensity needs to be assessed relative to
17 the baseline, but the size of the episodes clearly relates to their intensity and would be
18 independent of shifts in baseline. Thus the AQX statistics based on the daily MDA8 values
19 here are an important model test.

20

21 **Acknowledgements**

22 Research at UCI was supported by NASA grants NNX09AJ47G, NNX13AL12G,
23 NNX15AE35G, and DOE award DE-SC0007021. JLS was supported by the National
24 Science Foundation's Graduate Research Fellowship Program (DGE-1321846). The work of
25 DB and PC was funded by the U.S. Dept. of Energy (BER), performed under the auspices of
26 LLNL under Contract DE-AC52-07NA27344, and used the supercomputing resources of
27 NERSC under contract No. DE-AC02-05CH11231. GZ acknowledges the use of New
28 Zealand's national HPC facilities that are provided by the NZ eScience Infrastructure and
29 funded jointly by NeSI's collaborator institutions and through the Ministry of Business,
30 Innovation & Employment's Research Infrastructure Programme. The simulations with
31 MIROC-CHEM was supported by the Global Environment Research Fund (S-7) by the

1 Ministry of the Environment Japan and completed with the supercomputer (NEC SX-8R) at
2 the National Institute for Environmental Studies (NIES). We are grateful to the US
3 Environmental Protection Agency's (EPA) Air Quality System (AQS) and Clean Air Status
4 and Trends Network (CASTNet), Environment Canada's National Air Pollution Surveillance
5 Program (NAPS), the European Monitoring and Evaluation Programme (EMEP), and the
6 European Environment Agency's (EEA) air quality database (AirBase) for providing the
7 observational datasets used in this study. We are also grateful to the British Atmospheric
8 Data Centre (BADDC), which is part of the NERC National Centre for Atmospheric Science
9 (NCAS), for collecting and archiving the ACCMIP data.

10

11 **References**

12 Barnes, E. A., and Fiore, A. M.: Surface ozone variability and the jet position: Implications
13 for projecting future air quality, *Geophys. Res. Lett.*, 40, 2839-2844, doi:10.1002/grl.50411,
14 2013.

15 Büeker, P., Morrissey, T., Briolat, A., Falk, R., Simpson, D., Tuovinen, J.-P., Alonso, R.,
16 Barth, S., Baumgarten, M., Grulke, N., Karlsson, P. E., King, J., Lagergren, F., Matyssek, R.,
17 Nunn, A., Ogaya, R., Peñuelas, J., Rhea, L., Schaub, M., Uddling, J., Werner, W., and
18 Emberson, L. D.: DO₃SE modelling of soil moisture to determine ozone flux to forest trees,
19 *Atmos. Chem. Phys.*, 12, 5537-5562, doi:10.5194/acp-12-5537-2012, 2012.

20 Cameron-Smith, P., Lamarque, J. F., Connell, P., Chuang, C., and Vitt, F.: Toward an Earth
21 system model: atmospheric chemistry, coupling, and petascale computing, *J. Phys.: Conf.*
22 *Ser.*, 46, 343-350, doi:10.1088/1742-6596/46/1/048, 2006.

23 Doherty, R. M., Wild, O., Shindell, D. T., Zeng, G., MacKenzie, I. A., Collins, W. J., Fiore,
24 A. M., Stevenson, D. S., Dentener, F. J., Schultz, M. G., Hess, P., Derwent, R. G., and
25 Keating, T. J.: Impacts of climate change on surface ozone and intercontinental ozone
26 pollution: A multi-model study, *J. Geophys. Res.-Atmos.*, 118, 3744-3763,
27 doi:10.1002/jgrd.50266, 2013.

28 Donner, L. J., Wyman, B. L., Hemler, R. S., Horowitz, L. W., Ming, Y., Zhao, M., Golaz, J.-
29 C., Ginoux, P., Lin, S. J., Schwarzkopf, M. D., Austin, J., Alaka, G., Cooke, W. F., Delworth,
30 T. L., Freidenreich, S. M., Gordon, C. T., Griffies, S. M., Held, I. M., Hurlin, W. J., Klein, S.
31 A., Knutson, T. R., Langenhorst, A. R., Lee, H.-C., Lin, Y., Magi, B. I., Malyshev, S. L.,

1 Milly, P. C. D., Naik, V., Nath, M. J., Pincus, R., Ploshay, J. J., Ramaswamy, V., Seman, C.
2 J., Shevliakova, E., Sirutis, J. J., Stern, W. F., Stouffer, R. J., Wilson, R. J., Winton, M.,
3 Wittenberg, A. T., and Zeng, F.: The Dynamical Core, Physical Parameterizations, and Basic
4 Simulation Characteristics of the Atmospheric Component AM3 of the GFDL Global
5 Coupled Model CM3, *J. Climate*, 24, 3484-3519, doi:10.1175/2011jcli3955.1, 2011.

6 European Commission, Joint Research Centre (JRC)/Netherlands Environmental Assessment
7 Agency (PBL), EC-JRC/PBL, Emission Database for Global Atmospheric Research
8 (EDGAR), release version 4.0., <http://edgar.jrc.ec.europa.eu> (last access: 23 August 2014),
9 2009.

10 Fiore, A. M., Dentener, F. J., Wild, O., Cuvelier, C., Schultz, M. G., Hess, P., Textor, C.,
11 Schulz, M., Doherty, R. M., Horowitz, L. W., MacKenzie, I. A., Sanderson, M. G., Shindell,
12 D. T., Stevenson, D. S., Szopa, S., Van Dingenen, R., Zeng, G., Atherton, C., Bergmann, D.,
13 Bey, I., Carmichael, G., Collins, W. J., Duncan, B. N., Faluvegi, G., Folberth, G., Gauss, M.,
14 Gong, S., Hauglustaine, D., Holloway, T., Isaksen, I. S. A., Jacob, D. J., Jonson, J. E.,
15 Kaminski, J. W., Keating, T. J., Lupu, A., Marmer, E., Montanaro, V., Park, R. J., Pitari, G.,
16 Pringle, K. J., Pyle, J. A., Schroeder, S., Vivanco, M. G., Wind, P., Wojeik, G., Wu, S., and
17 Zuber, A.: Multimodel estimates of intercontinental source-receptor relationships for ozone
18 pollution, *J. Geophys. Res.-Atmos*, 114, D04301, doi:10.1029/2008jd010816, 2009.

19 Ganzeveld, L., Bouwman, L., Stehfest, E., van Vuuren, D. P., Eickhout, B., and Lelieveld, J.:
20 Impact of future land use and land cover changes on atmospheric chemistry-climate
21 interactions, *J. Geophys. Res.*, 115, D23301, doi:10.1029/2010JD014041, 2010.

22 Hjellbrekke, A.-G., Solberg, S., and Fjærraa, A. M.: Ozone measurements 2011, EMEP/CCC-
23 Report 3/2013, 0-7726, Tech. Rep., Norwegian Institute for Air Research, Norway, available
24 at: <http://www.nilu.no/projects/CCC/reports/cccr3-2013.pdf> (last access: 25 July 2013), 2013.

25 Holmes, C. D., Prather, M. J., Sovde, O. A., and Myhre, G.: Future methane, hydroxyl, and
26 their uncertainties: key climate and emission parameters for future predictions, *Atmos. Chem.*
27 *Phys.*, 13, 285-302, doi:10.5194/acp-13-285-2013, 2013.

28 Horton, D. E., Skinner, C. B., Singh, D., and Diffenbaugh, N. S.: Occurrence and persistence
29 of future atmospheric stagnation events, *Nature Clim. Change*, 4, 698-703,
30 doi:10.1038/nclimate2272, 2014.

1 HTAP: Hemispheric Transport Of Air Pollution 2010, Part A: Ozone And Particulate Matter,
2 United Nations, Geneva, Switzerland, 2010.

3 Hudman, R. C., Murray, L. T., Jacob, D. J., Turquety, S., Wu, S., Millet, D. B., Avery, M.,
4 Goldstein, A. H., and Holloway, J.: North American influence on tropospheric ozone and the
5 effects of recent emission reductions: Constraints from ICARTT observations, *J. Geophys.*
6 *Res.-Atmos*, 114, D07302, doi:10.1029/2008jd010126, 2009.

7 Josse, B., Simon, P., and Peuch, V. H.: Radon global simulations with the multiscale
8 chemistry and transport model MOCAGE, *Tellus B*, 56, 339-356, doi:10.1111/j.1600-
9 0889.2004.00112.x, 2004.

10 Kirtman, B., Power, S., Adedoyin, A. J., Boer, G., Bojariu, R., Camilloni, I., Doblus-Reyes,
11 F., Fiore, A., Kimoto, M., Meehl, G., Prather, M., Sarr, A., Schaer, C., Sutton, R.,
12 Oldenborgh, G. J. v., Vecchi, G., and Wang, H.-J.: Near-term Climate Change: Projections
13 and Predictability, in *Climate Change 2013: The Physical Science Basis*, chapter 11, IPCC
14 WGI Contribution to the Fifth Assessment Report, 2013.

15 Koch, D., Schmidt, G. A., and Field, C. V.: Sulfur, sea salt, and radionuclide aerosols in GISS
16 ModelE, *J. Geophys. Res.-Atmos*, 111, D06206, doi:10.1029/2004jd005550, 2006.

17 Koumoutsaris, S., and Bey, I.: Can a global model reproduce observed trends in summertime
18 surface ozone levels?, *Atmos. Chem. Phys.*, 12, 6983-6998, doi:10.5194/acp-12-6983-2012,
19 2012.

20 Lamarque, J. F., Shindell, D. T., Josse, B., Young, P. J., Cionni, I., Eyring, V., Bergmann, D.,
21 Cameron-Smith, P., Collins, W. J., Doherty, R., Dalsoren, S., Faluvegi, G., Folberth, G.,
22 Ghan, S. J., Horowitz, L. W., Lee, Y. H., MacKenzie, I. A., Nagashima, T., Naik, V.,
23 Plummer, D., Righi, M., Rumbold, S. T., Schulz, M., Skeie, R. B., Stevenson, D. S., Strode,
24 S., Sudo, K., Szopa, S., Voulgarakis, A., and Zeng, G.: The Atmospheric Chemistry and
25 Climate Model Intercomparison Project (ACCMIP): overview and description of models,
26 simulations and climate diagnostics, *Geosci. Model Dev.*, 6, 179-206, doi:10.5194/gmd-6-
27 179-2013, 2013.

28 Lin, J. T., Youn, D., Liang, X. Z., and Wuebbles, D. J.: Global model simulation of
29 summertime US ozone diurnal cycle and its sensitivity to PBL mixing, spatial resolution, and
30 emissions, *Atmos. Environ.*, 42, 8470-8483, doi:10.1016/j.atmosenv.2008.08.012, 2008.

1 Logan, J. A.: Ozone in rural-areas of the united-states, *J. Geophys. Res.-Atmos*, 94, 8511-
2 8532, doi:10.1029/JD094iD06p08511, 1989.

3 Menut, L., Bessagnet, B., Colette, A., and Khvorostiyarov, D.: On the impact of the vertical
4 resolution on chemistry-transport modelling, *Atmos. Environ.*, 67, 370-384,
5 doi:10.1016/j.atmosenv.2012.11.026, 2013.

6 Mickley, L., Jacob, D., Field, B., and Rind, D.: Effects of future climate change on regional
7 air pollution episodes in the United States, *Geophys. Res. Lett.*, 31, L24103,
8 doi:10.1029/2004GL021216, 2004.

9 Naik, V., Horowitz, L. W., Fiore, A. M., Ginoux, P., Mao, J., Aghedo, A. M., and Levy, H.,
10 II: Impact of preindustrial to present-day changes in short-lived pollutant emissions on
11 atmospheric composition and climate forcing, *J. Geophys. Res.-Atmos.*, 118, 8086-8110,
12 doi:10.1002/jgrd.50608, 2013.

13 National Research Council (US). Committee on the Significance of International Transport of
14 Air Pollutants. Global Sources of Local Pollution: An Assessment of Long-Range Transport
15 of Key Air Pollutants to and from the United States. National Academies Press, 2009.

16 Oman, L. D., Ziemke, J. R., Douglass, A. R., Waugh, D. W., Lang, C., Rodriguez, J. M., and
17 Nielsen, J. E.: The response of tropical tropospheric ozone to ENSO, *Geophys. Res. Lett.*, 38,
18 L13706, doi:10.1029/2011gl047865, 2011.

19 Prather, M., Gauss, M., Bernsten, T., Isaksen, I., Sundet, J., Bey, I., Brasseur, G., Dentener,
20 F., Derwent, R., Stevenson, D., Grenfell, L., Hauglustaine, D., Horowitz, L., Jacob, D.,
21 Mickley, L., Lawrence, M., von Kuhlmann, R., Muller, J. F., Pitari, G., Rogers, H., Johnson,
22 M., Pyle, J., Law, K., van Weele, M., and Wild, O.: Fresh air in the 21st century?, *Geophys.*
23 *Res. Lett.*, 30, No. 2, 1100, doi:10.1029/2002gl016285, 2003.

24 Reidmiller, D. R., Fiore, A. M., Jaffe, D. A., Bergmann, D., Cuvelier, C., Dentener, F. J.,
25 Duncan, B. N., Folberth, G., Gauss, M., Gong, S., Hess, P., Jonson, J. E., Keating, T., Lupu,
26 A., Marmer, E., Park, R., Schultz, M. G., Shindell, D. T., Szopa, S., Vivanco, M. G., Wild,
27 O., and Zuber, A.: The influence of foreign vs. North American emissions on surface ozone in
28 the US, *Atmos. Chem. Phys.*, 9, 5027-5042, doi:10.5194/acp-9-5027-2009, 2009.

29 Schnell, J. L., Holmes, C. D., Jangam, A., and Prather, M. J.: Skill in forecasting extreme
30 ozone pollution episodes with a global atmospheric chemistry model, *Atmos. Chem. Phys.*,
31 14, 7721-7739, doi:10.5194/acp-14-7721-2014, 2014.

1 Scinocca, J. F., McFarlane, N. A., Lazare, M., Li, J., and Plummer, D.: Technical Note: The
2 CCCma third generation AGCM and its extension into the middle atmosphere, *Atmos. Chem.*
3 *Phys.*, 8, 7055-7074, doi:10.5194/acp-8-7055-2008, 2008.

4 Seinfeld, J. H., Atkinson, R., Berglund R. L., Chameides, W. L., Elston, J. C., Fehsenfeld, F.,
5 Finlayson-Pitts, B. J., Harriss, R. C., Kolb, C. E., Liou, P. J., Logan, J. A., Prather, M. J.,
6 Russell, A., and Steigerwald, B.: Rethinking the ozone problem in urban and regional air
7 pollution, National Academy Press, Washington, D.C., 1991.

8 Shindell, D. T., Pechony, O., Voulgarakis, A., Faluvegi, G., Nazarenko, L., Lamarque, J. F.,
9 Bowman, K., Milly, G., Kovari, B., Ruedy, R., and Schmidt, G. A.: Interactive ozone and
10 methane chemistry in GISS-E2 historical and future climate simulations, *Atmos. Chem.*
11 *Phys.*, 13, 2653-2689, doi:10.5194/acp-13-2653-2013, 2013.

12 Silva, R. A., West, J. J., Zhang, Y. Q., Anenberg, S. C., Lamarque, J. F., Shindell, D. T.,
13 Collins, W. J., Dalsoren, S., Faluvegi, G., Folberth, G., Horowitz, L. W., Nagashima, T.,
14 Naik, V., Rumbold, S., Skeie, R., Sudo, K., Takemura, T., Bergmann, D., Cameron-Smith, P.,
15 Cionni, I., Doherty, R. M., Eyring, V., Josse, B., MacKenzie, I. A., Plummer, D., Righi, M.,
16 Stevenson, D. S., Strode, S., Szopa, S., and Zeng, G.: Global premature mortality due to
17 anthropogenic outdoor air pollution and the contribution of past climate change, *Environ. Res.*
18 *Lett.*, 8, 034005, doi:10.1088/1748-9326/8/3/034005, 2013.

19 Tang, Q., and Prather, M. J.: Correlating tropospheric column ozone with tropopause folds:
20 the Aura-OMI satellite data, *Atmos. Chem. Phys.*, 10, 9681-9688, doi:10.5194/acp-10-9681-
21 2010, 2010.

22 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J.*
23 *Geophys. Res.-Atmos.*, 106, 7183-7192, doi:10.1029/2000jd900719, 2001.

24 Teysedre, H., Michou, M., Clark, H. L., Josse, B., Karcher, F., Olivie, D., Peuch, V. H.,
25 Saint-Martin, D., Cariolle, D., Attie, J. L., Nedelec, P., Ricaud, P., Thouret, V., van der A, R.
26 J., Volz-Thomas, A., and Cheroux, F.: A new tropospheric and stratospheric Chemistry and
27 Transport Model MOCAGE-Climat for multi-year studies: evaluation of the present-day
28 climatology and sensitivity to surface processes, *Atmos. Chem. Phys.*, 7, 5815-5860,
29 doi:10.5194/acp-7-5815-2007, 2007.

1 Val Martin, M., Heald, C. L., and Arnold, S. R.: Coupling dry deposition to vegetation
2 phenology in the Community Earth System Model: Implications for the simulation of surface
3 O₃, *Geophys. Res. Lett.*, 41, 2988-2996, doi:10.1002/2014GL059651, 2014.

4 Wackernagel, H.: *Multivariate Geostatistics: An introduction with applications*, 3rd ed.,
5 Springer, Berlin, 387 pp., 2003.

6 Watanabe, S., Hajima, T., Sudo, K., Nagashima, T., Takemura, T., Okajima, H., Nozawa, T.,
7 Kawase, H., Abe, M., Yokohata, T., Ise, T., Sato, H., Kato, E., Takata, K., Emori, S., and
8 Kawamiya, M.: MIROC-ESM 2010: model description and basic results of CMIP5-20c3m
9 experiments, *Geosci. Model Dev.*, 4, 845-872, doi:10.5194/gmd-4-845-2011, 2011.

10 Wild, O., Fiore, A. M., Shindell, D. T., Doherty, R. M., Collins, W. J., Dentener, F. J.,
11 Schultz, M. G., Gong, S., MacKenzie, I. A., Zeng, G., Hess, P., Duncan, B. N., Bergmann, D.
12 J., Szopa, S., Jonson, J. E., Keating, T. J., and Zuber, A.: Modelling future changes in surface
13 ozone: a parameterized approach, *Atmos. Chem. Phys.*, 12, 2037-2054, doi:10.5194/acp-12-
14 2037-2012, 2012.

15 Xing, J., Mathur, R., Pleim, J., Hogrefe, C., Gan, C.-M., Wong, D. C., Wei, C., Gilliam, R.,
16 and Pouliot, G.: Observations and modeling of air quality trends over 1990-2010 across the
17 Northern Hemisphere: China, the United States and Europe, *Atmos. Chem. Phys.*, 15, 2723-
18 2747, doi:10.5194/acp-15-2723-2015, 2015.

19 Young, P. J., Archibald, A. T., Bowman, K. W., Lamarque, J. F., Naik, V., Stevenson, D. S.,
20 Tilmes, S., Voulgarakis, A., Wild, O., Bergmann, D., Cameron-Smith, P., Cionni, I., Collins,
21 W. J., Dalsoren, S. B., Doherty, R. M., Eyring, V., Faluvegi, G., Horowitz, L. W., Josse, B.,
22 Lee, Y. H., MacKenzie, I. A., Nagashima, T., Plummer, D. A., Righi, M., Rumbold, S. T.,
23 Skeie, R. B., Shindell, D. T., Strode, S. A., Sudo, K., Szopa, S., and Zeng, G.: Pre-industrial
24 to end 21st century projections of tropospheric ozone from the Atmospheric Chemistry and
25 Climate Model Intercomparison Project (ACCMIP), *Atmos. Chem. Phys.*, 13, 2063-2090,
26 doi:10.5194/acp-13-2063-2013, 2013.

27 Zeng, G., Pyle, J. A., and Young, P. J.: Impact of climate change on tropospheric ozone and
28 its global budgets, *Atmos. Chem. Phys.*, 8, 369-387, doi:10.5194/acp-8-369-2008, 2008.

29 Zeng, G., Morgenstern, O., Braesicke, P., and Pyle, J. A.: Impact of stratospheric ozone
30 recovery on tropospheric ozone and its budget, *Geophys. Res. Lett.*, 37, L09805,
31 doi:10.1029/2010gl042812, 2010.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Table 1. Observational datasets (2000 to 2009).			
<i>Domain</i>	<i>Surface ozone network</i>	<i>No. stations</i>	<i>URL or reference</i>
North America (NA)	US EPA Air Quality System (AQS)	1633	http://www.epa.gov/ttn/airs/aqsdatamart
	US EPA Clean Air Status and Trends Network (CASTNet)	92	http://epa.gov/castnet/javaweb/index.html
	Environment Canada's National Air Pollution Surveillance Program (NAPS)	207	http://maps-cartes.ec.gc.ca/rnspa-naps/data.aspx?lang=en
Europe (EU)	European Monitoring and Evaluation Programme (EMEP)	162	<i>Hjellbrekke et al. (2013)</i>
	European Environment Agency's air quality database (AirBase)	2123	www.eea.europa.eu/data-and-maps/data/airbase-the-european-air-quality-database-8

Table 2. Model summary.

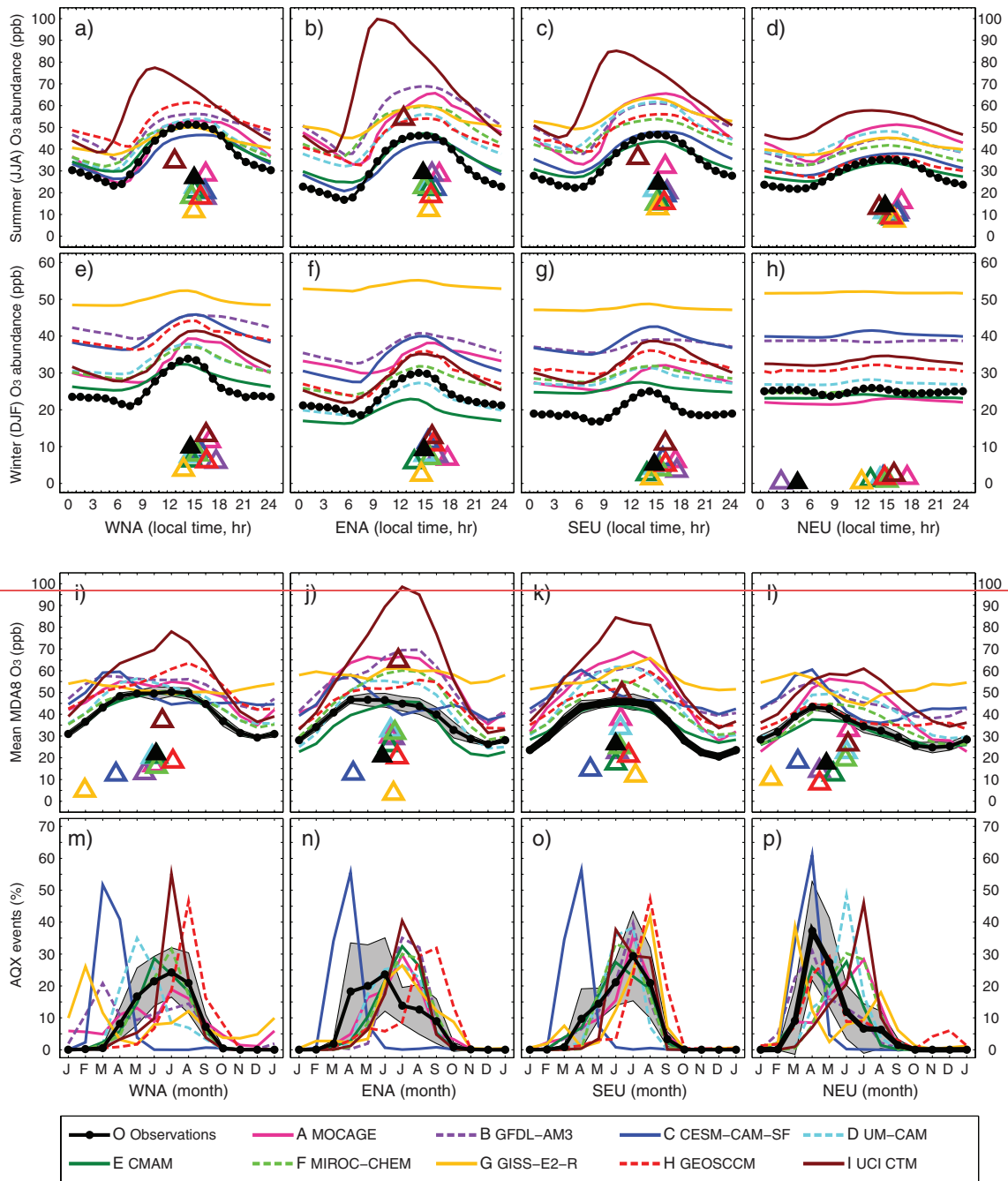
<i>(abbreviation) Model</i>	<i>Modeling Center</i>	<i>Member^a</i>	<i>Resolution (lat. x lon)</i>	<i>No. years</i>	<i>Reference(s)</i>
(A) MOCAGE	MeteoFrance	r2i1p1, v2	2° x 2°	4	<i>Josse et al. (2004)</i> <i>Teyssèdre et al. (2007)</i>
(B) GFDL-AM3	GFDL	r1i1p1, v2	2° x 2.5°	10	<i>Donner et al. (2011)</i> <i>Naik et al. (2013)</i>
(C) CESM-CAM-SF	LLNL-NCAR	r1i1p1, v4	~1.9° x 2.5°	10	<i>Cameron-Smith et al. (2006)</i> <i>Lamarque et al. (2013)</i>
(D) UM-CAM	NIWA	r1i1p1, v2	2.5° x 3.75°	10	<i>Zeng et al. (2008, 2010)</i>
(E) CMAM	CCCma	r1i1p1, v2	~3.7° x 3.75°	10	<i>Scinocca et al. (2008)</i>
(F) MIROC-CHEM	JAMSTEC-NU- NIES	r1i1p1, v2	~2.8° x 2.8125°	10	<i>Watanabe et al. (2011)</i>
(G) GISS-E2-R	GISS	r1i1p3, v1	2° x 2.5°	5	<i>Koch et al. (2006)</i> <i>Shindell et al. (2013)</i>
(H) GEOSCCM	NASA-GSFC	r1i1p1, v1	2° x 2.5°	10	<i>Oman et al. (2011)</i>
(I) UCI CTM	UCI	--	~2.8° x 2.8125°	10	<i>Holmes et al. (2013)</i>

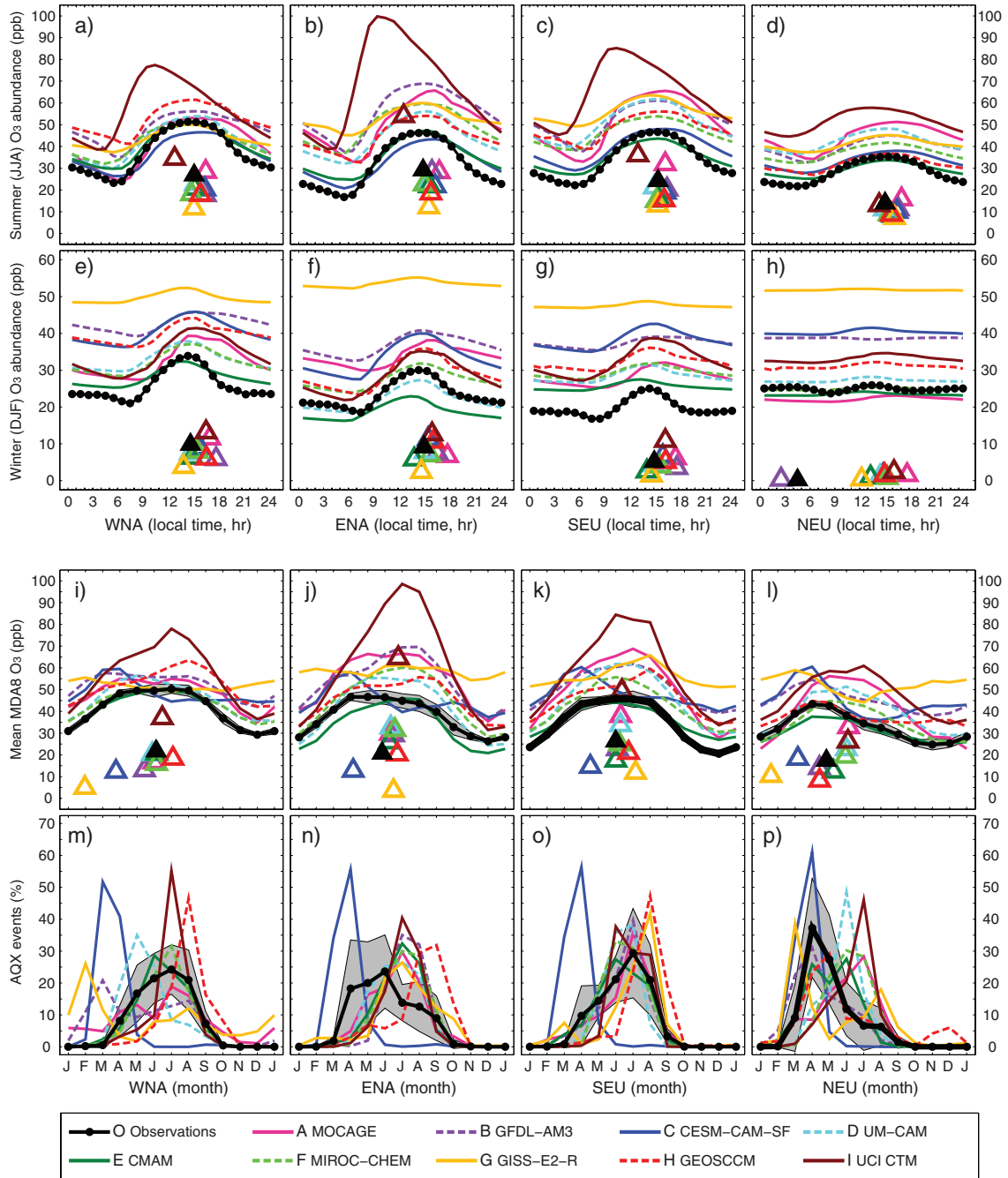
^aThe format r<N>i<M>p<L>, vX distinguishes among closely related simulations by a single model where the set of integers (N, M, L, X) formatted as shown (e.g., r2i1p1, v2) define each model simulation's realization number (N), initialization method (M), perturbed physics version (L), and version of publication-level dataset (X).

1
2
3
4
5
6
7
8
9

Table 3. Example summary statistics for the observations (OBS), the ACCMIP models (A-H), and the UCI CTM (I) for Eastern North America’s (ENA) summer (JJA) and winter (DJF) diurnal cycles, annual cycle of MDA8, annual cycle of AQX events, and North America’s (NA, combined Western North America (WNA) and ENA) AQX episodes (100 AQX events per decade case).

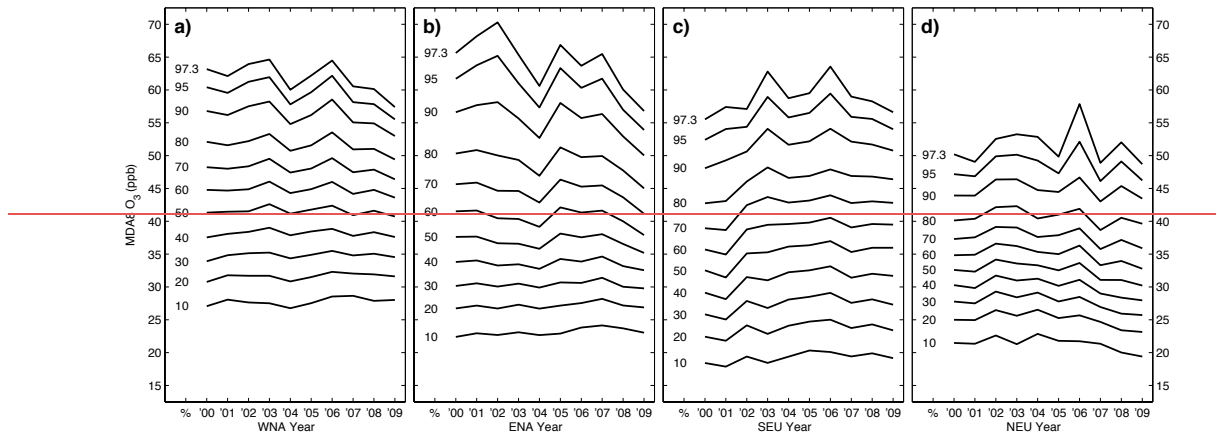
Data	Metric, description (unit)	OBS	A	B	C	D	E	F	G	H	I
JJA diurnal cycle	h , maximum phase (hour)	15.0	17.0	16.1	16.5	15.5	15.8	15.2	15.7	16.0	12.7
	H , peak-to-peak amplitude (ppb)	29.1	28.3	28.4	21.8	22.7	21.8	22.6	12.1	18.5	54.0
	MB, mean bias (ppb)	-	19.0	24.4	1.1	12.2	3.5	17.9	21.1	12.9	37.0
DJF diurnal cycle	h , maximum phase (hour)	15.1	18.0	16.7	15.7	15.3	14.0	15.9	14.8	16.3	16.1
	H , peak-to-peak amplitude (ppb)	9.1	6.7	7.5	11.3	7.8	5.8	6.9	2.4	10.6	12.6
	MB, mean bias (ppb)	-	10.2	13.2	9.8	-1.5	-4.6	4.0	30.1	5.5	4.8
MDA8 annual cycle	m , maximum phase (month)	5.3	5.8	6.0	3.7	5.8	5.7	6.1	6.0	6.2	6.3
	M , peak-to-peak amplitude (ppb)	20.7	29.8	29.1	12.8	32.7	25.9	31.5	3.5	20.3	64.6
	MB, mean bias (ppb)	-	16.9	16.6	6.8	4.2	-4.2	8.1	20.1	8.0	24.8
	\bar{E}_{JJA} , 87th – 30th percentile (ppb)	22.8	33.0	27.5	19.4	27.0	22.4	28.3	19.1	21.9	56.0
	R_{E-JJA} , spatial correlation of E_{JJA} maps	1.00	0.70	0.81	0.52	0.69	0.69	0.34	0.27	0.69	0.71
AQX event annual cycle	m_{AQX} , maximum phase (month)	5.5	6.2	6.8	3.2	6.2	6.4	6.6	6.8	7.7	6.6
	R_{MDA8} , correlation of AQX and MDA8 cycles	0.84	0.76	0.78	0.88	0.78	0.82	0.80	0.78	0.70	0.83
	\bar{E}_{AQX} , AQX threshold – 30th percentile (ppb)	34.7	53.8	39.9	29.1	36.1	30.4	41.1	32.1	31.5	82.3
	R_{E-AQX} , spatial correlation of E_{AQX} maps	1.00	0.70	0.78	0.28	0.63	0.53	0.44	0.60	0.74	0.68
NA AQX episodes	\bar{S} , weighted geometric mean AQX episode size (10^4 km ² -days)	415	128	229	1426	461	290	522	243	774	463
	CCD_{100} , fraction of AQX events’ areas in AQX episodes > 100 x 10^4 km ² -days (%)	79.0	56.1	73.7	92.6	85.3	76.1	80.3	73.0	83.0	80.2
	CCD_{1000} , fraction of AQX events’ areas in AQX episodes > 1000 x 10^4 km ² -days (%)	38.0	9.7	12.8	69.2	30.8	19.2	43.6	12.7	48.7	37.5
	$\Delta\bar{E}_S$, average increase in E_S for AQX episodes of size S (ppb-dec ⁻¹)	2.9	9.9	4.6	0.8	2.3	2.9	-0.1	3.5	2.9	6.0



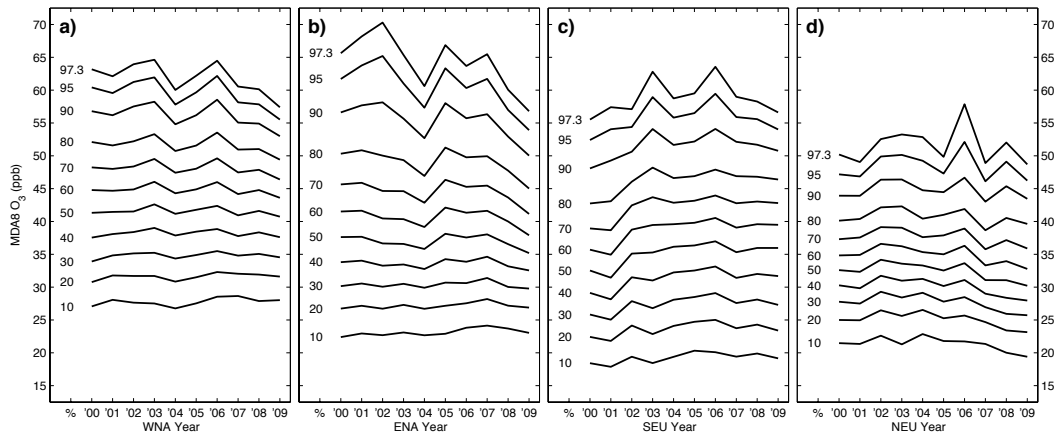


1
2 **Figure 1.** (a-h) Diurnal cycles of hourly O₃ abundances (ppb) for the observations (O), ACCMIP models (A-H),
3 and UCI CTM (I) averaged over (a-d) summer (JJA) and (e-h) winter (DJF) months in (a, e) WNA, (b, f) ENA,
4 (c, g) SEU, and (d, h) NEU. Triangles show the observation's and models' cosine fit derived values of the hour
5 of maximum phase h and peak-to-peak amplitude H plotted as $(x, y) = (h, H)$ for each season, region,
6 observation, and model. (i-p) Annual cycles of (i-l) MDA8 O₃ and (m-p) AQX events in (i, m) WNA, (j, n)
7 ENA, (k, o) SEU, and (l, p) NEU. The filled gray curve shows $\pm 1\sigma$ for each month (calculated across years) for
8 the observations. Triangles show the observations' and models' cosine fit derived values of the MDA8 cycle

- 1 month of maximum phase m and peak-to-peak amplitude M plotted as $(x, y) = (m, M)$ for each region,
- 2 observation, and model.



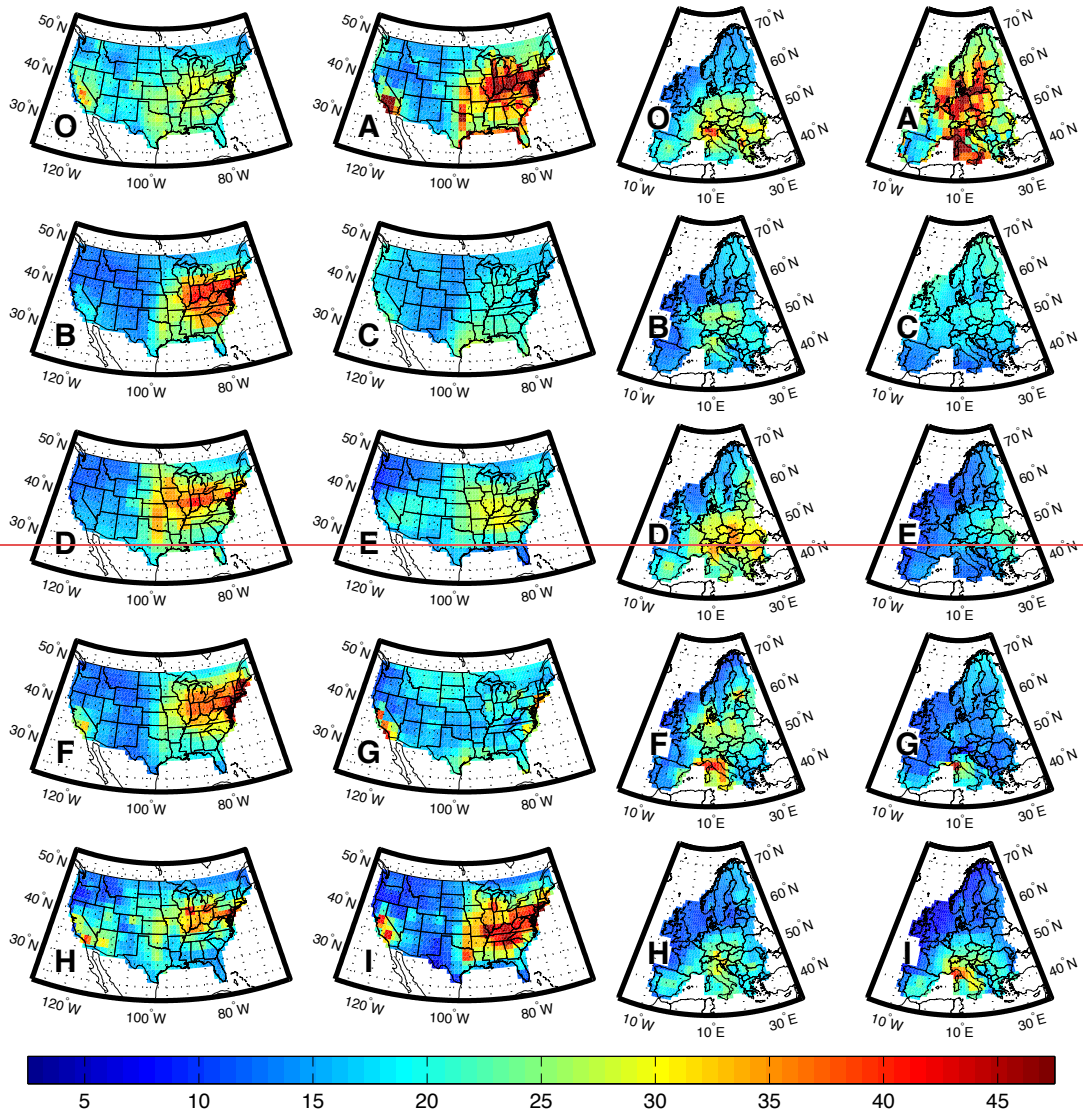
3



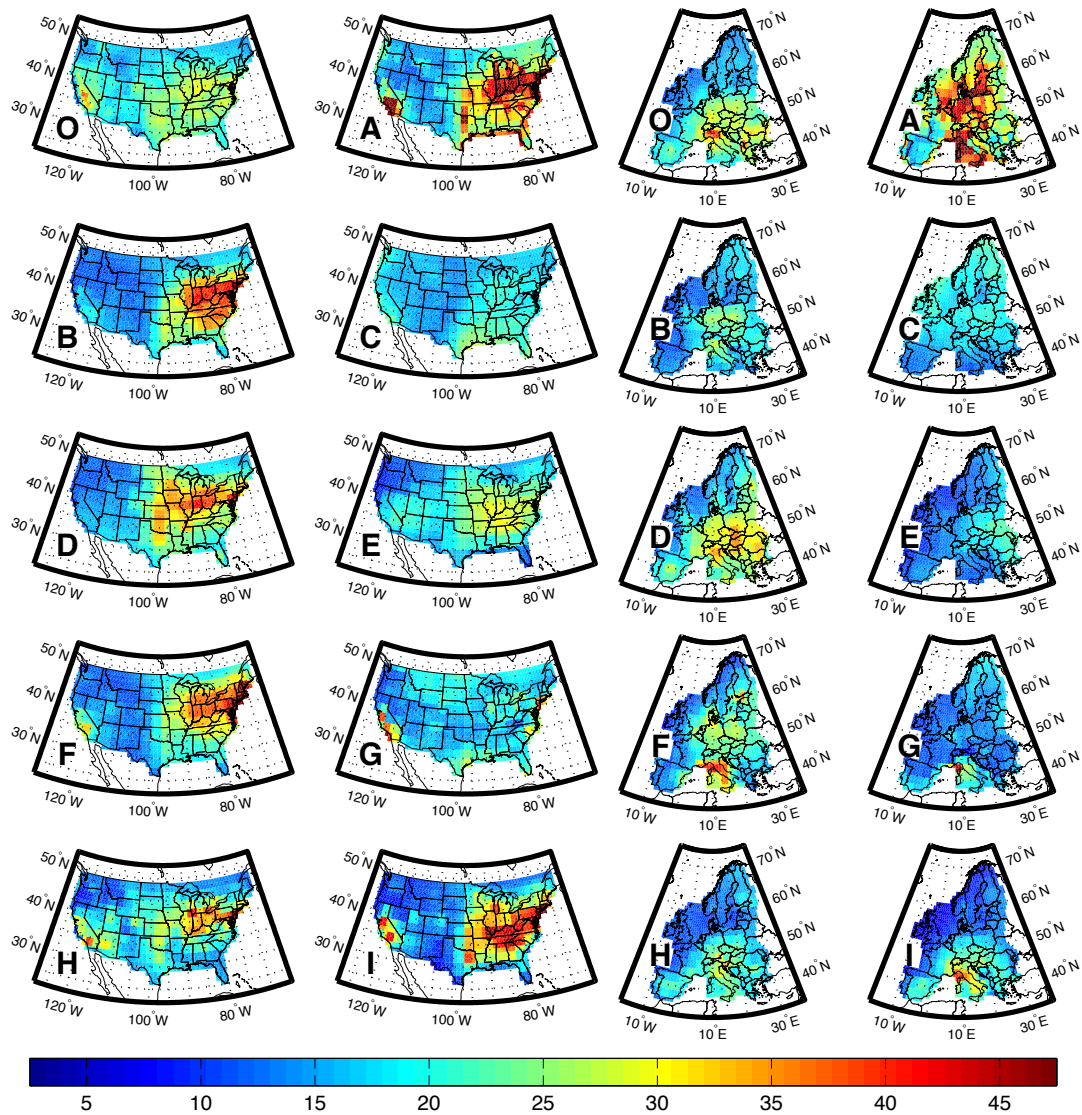
4

5 **Figure 2.** Values of MDA8 O₃ (ppb) for years 2000 to 2009 corresponding to the 10th, 20th, ..., 90th, 95th, and
 6 97.3 (i.e., AQX threshold) percentiles in (a) WNA, (b) ENA, (c) SEU, and (d) NEU. The percentile for each
 7 line is shown at the beginning of the curves in each panel.

8

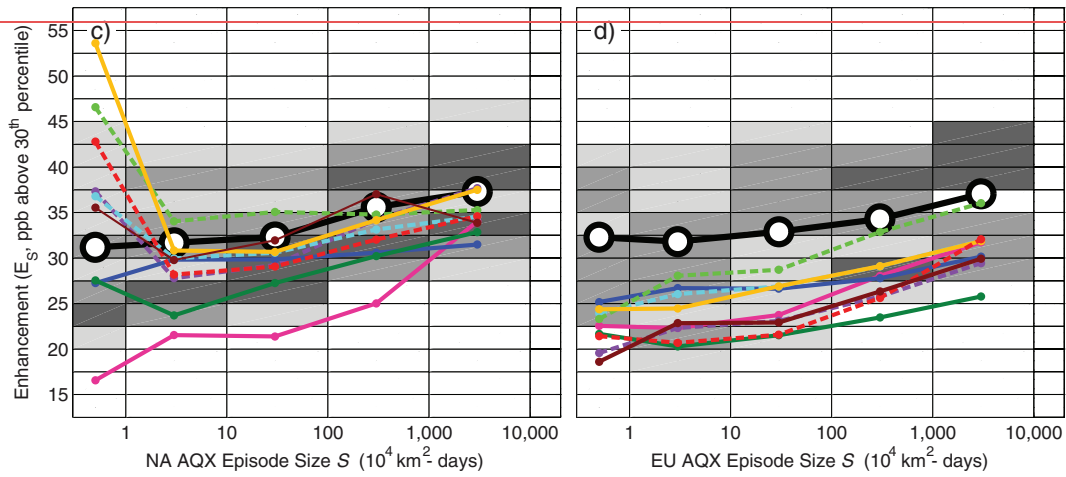
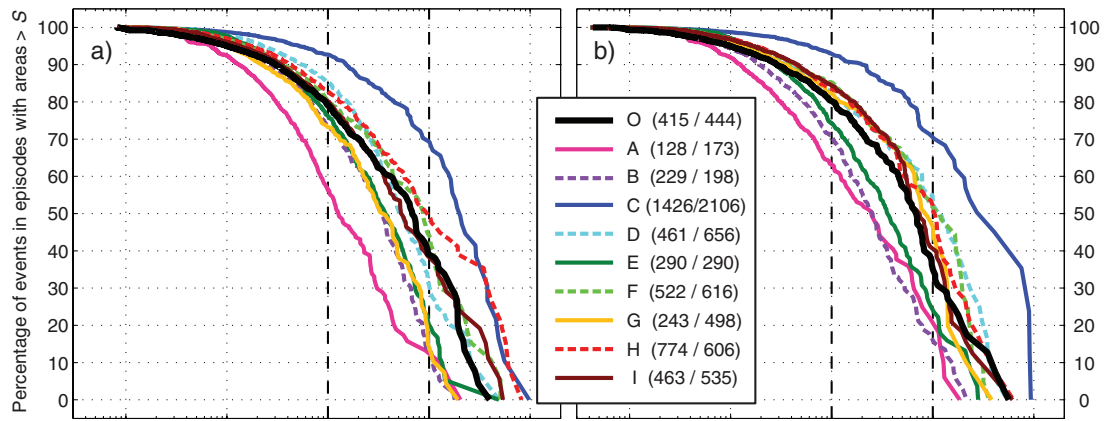


1

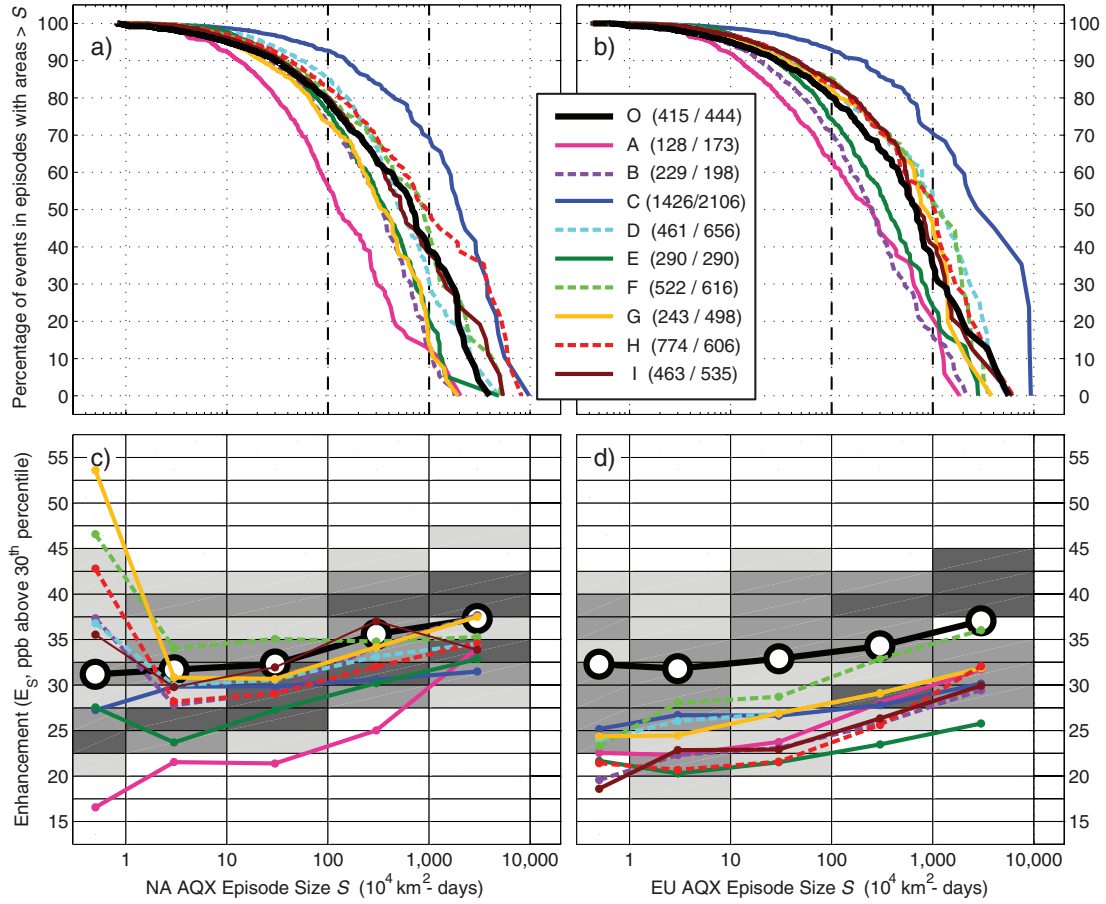


1
 2 **Figure 3.** Summertime O_3 enhancement $E_{JJA} =$ difference between the 87th and 30th percentile of the gridded
 3 surface MDA8 O_3 (ppb) over (left two columns) NA and (right two columns) EU for the observations (O),
 4 ACCMIP models (A-H), and UCI CTM (I). The values of model I are scaled by 0.5 so the same color scale can
 5 be used.

6
 7
 8
 9
 10
 11
 12



1



1
2 **Figure 4.** (a-b) Complementary cumulative distribution (CCD) of the percentage of total areal extent of all
3 individual AQX events (100-per-decade case) as a function of AQX episode size, (S , $10^4 \text{ km}^2\text{-days}$) for the
4 observations (O), ACCMIP models (A-H), and UCI CTM (I) in (a) NA and (b) EU. Dashed vertical lines show
5 the graphical representations of CCD_{100} and CCD_{1000} . Mean episode size \bar{S} for each dataset and domain is given
6 in the legend as (NA/EU). (c-d) Density scatterplot of the observations enhancement of AQX episodes E_S versus
7 their size S (E_S binned at 2.5 ppb increments from <15 ppb to >55 ppb, S binned at each log-decade) in (c) NA
8 and (d) EU. The gray scale represents the relative percentage of AQX episodes in each $(x, y) = (S, E_S)$ bin and
9 includes percent ranges of $\leq 5\%$ (white), 5-10%, 10-15%, and $>15\%$ (darkest gray) where the size bins (i.e.,
10 columns) are normalized to sum to 100%. The overlain curves show the observation's and each model's area-
11 weighted mean enhancement E_S for each size bin. The values of E_S in each size bin for models A and I have
12 been scaled by 0.5 since they are largely outside the range of the others.

13
14