Response to reviews on Representativeness errors in comparing chemistry transport and chemistry climate models with satellite UV/Vis tropospheric column retrievals.

San Francisco, 18 December 2015

Dear editor,

Below please find a point-by-point response to the reviewers' comments on our manuscript. The changes made to the manuscript in response to these comments are evident in our response.

Below the list of responses, we added the revised manuscript with changes in blue. Removed parts are indicated by strikethrough font.

We hope that with this revision, our manuscript is accepted for publication in GMD.

Kind regards,

Folkert Boersma

Anonymous Referee #1

This manuscript presents a methodology for properly comparing models for atmospheric composition with Level-2 retrievals from satellites in the UV/VIS. It does so by identifying three potential issues, testing the potential introduced errors for these issues and recommend procedures to avoid the errors. The manuscript is generally well written (apart from Section 2; see detailed comments below) and the scientific methods are sound. However, the presented concepts and methods are not really new and the manuscript is in that sense a bit disappointing. Sampling issues and the use of averaging kernels have been discussed in the literature for some time now. Nevertheless, I support the publishing of this manuscript as a model assessment methods paper, subject to minor revision, because it provides a good summary of the issues and guidelines, which will hopefully encourage scientists in the field to make better use of these concepts.

The concepts and methods used in our study indeed exist for some 20 years now. But whereas the data assimilation community may generally well aware of the issues at hand, we feel that the modeling community using tropospheric column retrievals of NO2, HCHCO, and SO2 from UV/Vis measurements is still struggling with representativeness errors, and generally lacks guidelines on how to perform a proper model-column comparison. This is partly because the UV/Vis retrievals do not always provide averaging kernels in their data products. Another reason is that the use of UV/Vis retrievals is relatively recent compared with the use of temperature and O3 profiles from satellite instruments in formal data assimilation systems, and new users lack a body of literature to guide them how to actually do that. Our paper indeed has the potential to encourage proper use of the UV/Vis satellite retrievals in the community.

Detailed comments:

Section 1: I am missing some important earlier work, such as the studies by Rodgers and Connor, JGR, 2003 and Migliorini et al., Monthly Weather Review, 2008. These could maybe be added in lines 16/17 of page 7827?

We added these relevant references in section 1.

P7827, line 11: The authors go a bit fast here. While the statement is correct, I had to read it several times to understand the logic. I would advise to add a sentence that explains that the contribution of the prior profile to the final solution increases with decreasing sensitivity of the measurement.

Thanks for the suggestion. We now added a sentence along these lines.

P7827, line 22-25: Assuming xa=0 makes the problem non-Gaussian, because xa cannot be negative. How does this assumption affect the discussion?

The DOAS retrieval is not an optimal estimation technique based on Bayes theorem. So this is not a problem. We refer to Eskes an Boersma [2003] for an extensive discussion of the retrieval theory for DOAS columns.

P7828, line 23: 'whereas satellite measurements provide "snapshots" at a particular local time': this is not a description of spatial error but of temporal error. Please remove. We removed "at a particular local time" because this may be confusing.

Section 2.3 is a slightly odd section. In section 1 the authors say that section 2 will introduce the issues, while section 2.3 is actually used to describe the proposed solution. There is therefore overlap with text later in the manuscript, which does not help the reader. I suggest to either remove section 2.3 or make it part of text later in the document (e.g., Section 3.2). Also, the end of section 2.2 describes conclusions, which should not be part of this introduction. Thanks for this excellent suggestion. We have moved section 2.3 to 3.2 and streamlined the text as necessary.

P7830, line 3: Should you not do the opposite in this case, i.e., average over all the grid boxes that fall within the satellite foot print?

That is in principle possible, but at the cost of no longer evaluating the model on its own grid, or a regular grid at all. Such evaluations are sometimes done on a daily basis (e.g. Huijnen et al., ACP, 2010). Ultimately, they still re-grid the satellite footprints to a coarser grid for final seasonal mean comparisons (maps, regional averages). Re-gridding the model-retrieval pairs from the irregular, daily varying satellite grid to a coarser common grid, brings the total number of manipulations on three (1. model-to-satellite, 2. model-at-satellite back to model grid, and 3. satellite to model grid), whereas the recipe in Eq. (3) only requires one such manipulation.

Section 2.4 is slightly strange as well. Of course, one has to be aware of clouds and their impact on the retrieval. But wasn't the whole point of recommending the use of averaging kernels exactly to deal with these kind of issues. Maybe, the authors could mention the impact of clouds in a shortened paragraph earlier in the manuscript and then present the use of averaging kernels later in the manuscript both in the context of vertical sensitivity in general and in the context of clouds.

Thanks for the suggestion. We have now merged the section on the influence of clouds with section 2.1. In that section we also state that it is in principle possible to use the cloudy retrievals as long as the kernels are being applied.

P7832, last paragraph: Treating systematic errors as random is a strange assumption considering the fact that the authors acknowledge later on that the systematic model errors can be as large as 50 %. This needs some further clarification.

We admit that this was formulated inadequately. It is better to interpret Eq. (6) as providing an overall envelop of the model-retrieval comparison error, where the individual (model, retrieval, representativeness) error contributions may have substantial systematic components. While some of these systematic contributions may cancel out in an area-averaged, monthly mean model-retrieval comparison such as presented in Section 6, robust systematic error will persist. We therefore removed the sentence "Here and in Sect. 6 we treat them as random errors, an approach usually followed in retrieval and data assimilation studies (Rodgers, 2000)".

P7833, lines 10 – 12: I am not sure if I understand this argument. Averaging will remove the random component from the error budget, so systematic errors will actually become more dominant. Or do I miss something?

In general, a priori data that is always and everywhere biased in the same way will then indeed become dominant. But a priori data can also suffer from systematic errors that are variable in space and time. Our argument has its origin in the use of a priori data in the retrieval (e.g. surface albedo) that may be systematically biased high for pixel 1, but biased low for pixel 2. Averaging retrievals over a large area, will then tend to cancel the effect of systematic errors in the a priori data, as long as the mean systematic error (in the albedo) is zero.

Section 3.2: Maybe the text of Section 2.3 could be used here. Done.

Section 4.1: It is unclear what the relationship is between the model profile used for the AMF and the a priori profile (set to zero under the weak absorbing assumption) used in the retrieval. Both prior assumptions play a role in the definition of the end product, but it is not clear how this is accounted for in the averaging kernel.

The a priori profile \mathbf{x}_a used in the air mass factor calculation is typically taken from a CTM. The a priori profile is used in the calculation of the AMF, and also in the averaging kernel, as discussed already in section 2.1, right below Eq. (1). The model that provides \mathbf{x}_a may or may not be the same as the model under evaluation. In our study, the a priori profiles \mathbf{x}_a in the OMI NO2 retrieval are from the TM4 model, and the models that are being evaluated are TM5 and GEOS-Chem. We decided to not go into much detail here, since this has been extensively described in many retrieval papers including our own (Eskes and Boersma [2004], Boersma et al., [2003; 2007; 2011]).

Page 7837, eq. 8: This was already described in Section 2.3

Good point. We removed Eq. (8) and instead refer to Eq. (3) now in the text.

Page 7838, lines 2 – 4: Does this mean that the DOAS error estimates are not correct? If they would be, they should be taken into account. I am not convinced by the argument that the method drives the statistical interpretation of the results.

The DOAS error, calculated from formal error propagation, is our best estimate of the uncertainty in a single retrieval. We stress that here we calculate an area average of individual retrievals that is the best spatial representation of the retrieved column over a large, grid-cell area. We take the individual errors (and their spatial correlation) into account when computing a superobservation error as in Eq. (9) in the manuscript.

Page 7841, line 13: Are the TM5 cloud fractions indeed simulated by the model or do they come from ERA-Interim. If the former is the case, it would be worth mentioning how these cloud fractions are simulated. If it is the latter, this should be mentioned as well. The TM5 cloud information comes from ERA-Interim, as stated already in Section 4.2 and in 5.2. The details of the cloud fraction simulation are described in Dee et al. [2011].

Section 5.3: Although the presented comparisons do indicate a better agreement between model and observations, it is not rock-solid proof. There remains the possibility that the use of averaging kernels is masking/compensating other errors in the model. This is probably worth mentioning for completeness.

It depends on what is understood as 'proof'. Our point is not that the model is doing better if the kernel is applied, but rather that if the model is sampled as if the satellite would observe the modeled state, the model and retrieval can be properly compared according to the retrieval's capabilities. Such a comparison allows for a better-constrained evaluation of the model. Of course some model errors may still go undetected following such an evaluation, but the evaluation is no longer biased because of the difference in vertical sensitivity of the retrievals vs. the model.

Appendix D is not referenced in the text. However, it contains interesting content that could maybe be used in Section 4.1? See also my earlier comment about Section 4.1. Thanks for this nice suggestion. We don't think it would be useful to include it in Section 4.1 that is about the retrieval settings. Still, we think referencing Appendix D is a good idea, and we do that now in Section 5.3, because this is where the actual application of the kernel is discussed.

Anonymous Referee #2

The authors have described an approach to quantity three types of errors that can arise when trying to compare modelled NO2 fields with UV/Vis retrievals of column NO2. They have focused on horizontal and vertical representative errors, and errors associate with cloud cover. These are all issues that the community is aware of and is struggling to address using various approaches. In that context, the manuscript is not innovative, but I believe that it will be useful to the community. By providing a coherent approach for dealing with these errors, the manuscript will help reduce misuse of the NO2 column data, and it may even spur new approaches for mitigating the errors. The manuscript is well written and I recommend it for publication after minor revisions to address my comments below. We thank the reviewer for the constructive comments.

Comments

1) Page 7828, line 25, and Page 7840, line 8 (title of Section 5.2): It is not clear to me what is the temporal component here. It seems to me that this error, which is the focus of Section 5.2, is really an issue with the representativeness of the cloud cover. Referring to this as temporal (or meteorological) representativeness error is vague and confusing. Why not just call it representativeness errors in cloud cover?

On page 7828, we removed the adjective 'meteorological' to clarify that this category of errors is addressing temporal aspects that not just relate to cloud cover. Temporal representativeness errors arise from models not properly representing the temporal evolution associated with

changing cloud cover or with changing emissions such as the weekend effect, so just labelling them as 'representativeness errors in cloud cover' would not be a complete description. The title of section 5.2 is in our opinion appropriate as it is.

2) Page 7832, lines 21-23: Here the authors state that they treat all of the errors as random errors. However, in Section 5.1 and Appendix B they account for a correlation in the errors, so clearly they are accounting for some systematic errors. This confusing and needs to be better explained.

We removed the sentence that we treat the retrieval, modelling, and representativeness errors as random errors in response to Rev#1. In section 5.1, we discuss how to calculate representativeness errors, taking into account error correlation, so that combined observation and (spatial) representativeness may be obtained. In many data assimilation studies, representativeness errors are included in the observation errors, as already stated in our original manuscript, and now also stated right after introducing Eq. (3).

3) Page 7833, lines 8-12: I do not agree with the statement that the vertical transport errors will be smaller, in an average sense, when aggregated over a month. That is, for example, unlikely to be the case in the tropics or over eastern North America in summer, when convective transport is strong. Indeed, the monthly mean differences in the vertical distribution of NO2 in the lower troposphere between GEOS-Chem andTM5 shown in Figure 7a are not that different from those shown in Figure 6a for Feb 18th.

They will be smaller in a monthly average sense, because not all days in the month will have strong convective activity, especially in a region like North America. Furthermore, Figures 6 and 7 do not necessarily indicate that GEOS-Chem suffers from vertical transport errors but rather that GEOS-Chem and TM4 simulate vertical transport of NO2 differently.

4) Page 7836, line 17: Are there additional references besides Lin (2012) that should be included here? What about Lamsal et al. (2010, J. Geophys. Res., D05302, doi:10.1029/2009JD013351)? We have now included the Lamsal et al. [2010] reference, also in a number of other occasions where we thought it appropriate.

5) Page 7845, lines 2-6: Not all of the NO2 products provide averaging kernels. It would be helpful if the authors could explain what one should do to mitigate the vertical representativeness errors in the case when averaging kernels are not provided. There is not much one can do except but go back to the NO2 product developers and make a request to make a available the averaging kernels or the scattering weights for each pixel in the product. Using the averaging kernel of another product would not help, unless all the a priori data used in calculating the averaging kernels (albedo, cloud parameters, atmospheric profile) is exactly the same, which is highly unlikely to ever happen.

6) Figure 7 caption: Explain that the numbers given in each panel, e.g. 5.75 for GEOS-Chem, are the integrated column abundances. Done, also for Figure 6.

7) Figure 10 caption: Should "100% x (B/C - 1)" be "100% x (B/A - 1)"? Why not make these differences relative to experiment A, which is believed to be better, so that for C vs A is it 100% x (C/A - 1) and for B vs A it is 100% x (B/A - 1)?

Thanks for spotting this. The experiments were indicated relative to experiment A. We adapted the caption accordingly.

Representativeness errors in comparing chemistry transport and chemistry climate models with satellite UV/Vis tropospheric column retrievals

4

5 K. F. Boersma^{1,2}, G. C. M. Vinken³, and H. J. Eskes¹.

6 [1]{Royal Netherlands Meteorological Institute, De Bilt, The Netherlands}

7 [2]{Wageningen University, Meteorology and Air Quality department, Wageningen, The8 Netherlands}

9 [3] {Eindhoven University of Technology, Eindhoven, The Netherlands}

10 Correspondence to: K. F. Boersma (boersma@knmi.nl)

11

12 Abstract

13 UV/Vis satellite retrievals of trace gas columns of nitrogen dioxide (NO₂), sulphur dioxide (SO₂), and formaldehyde (HCHO) are useful to test and improve models of atmospheric 14 15 composition, for data assimilation, air quality hindcasting and forecasting, and to provide top-16 down constraints on emissions. However, because models and satellite measurements do not 17 represent the exact same geophysical quantities, the process of confronting model fields with 18 satellite measurements is complicated by representativeness errors, which degrade the quality 19 of the comparison beyond contributions from modelling and measurement errors alone. Here 20 we discuss three types of representativeness errors that arise from the act of carrying out a 21 model-satellite comparison: (1) horizontal representativeness errors due to imperfect 22 collocation of the model grid cell and an ensemble of satellite pixels called superobservation, 23 (2) temporal representativeness errors originating mostly from differences in cloud cover 24 between the modelled and observed state, and (3) vertical representativeness errors because of 25 reduced satellite sensitivity towards the surface accompanied with necessary retrieval 26 assumptions on the state of the atmosphere. To minimize the impact of these 27 representativeness errors, we recommend that models and satellite measurements be sampled 28 as consistently as possible, and our paper provides a number of recipes to do so. A practical 29 confrontation of tropospheric NO₂ columns simulated by the TM5 chemistry transport model

(CTM) with Ozone Monitoring Instrument (OMI) tropospheric NO₂ retrievals suggests that 1 2 horizontal representativeness errors, while unavoidable, are limited to within 5-10% in most 3 cases and of random nature. These errors should be included along with the individual 4 retrieval errors in the overall superobservation error. Temporal sampling errors from 5 mismatches in cloud cover, and, consequently, in photolysis rates, are on the order of 10% for NO₂ and HCHO, and systematic, but partly avoidable. In the case of air pollution applications 6 7 where sensitivity down to the ground is required, we recommend that models should be 8 sampled on the same mostly cloud-free days as the satellite retrievals. The most relevant 9 representativeness error is associated with the vertical sensitivity of Ultraviolet-visible (UV/Vis) satellite retrievals. Simple vertical integration of modelled profiles leads to 10 11 systematically different model columns compared to application of the appropriate averaging kernel. In comparing OMI NO₂ to GEOS-Chem NO₂ simulations, these systematic 12 13 differences are as large as 15-20% in Summer, but, again, avoidable.

14 **1** Introduction

15 Chemistry transport models (CTMs) are increasingly being evaluated with satellite column 16 retrievals from UV/Vis solar backscatter satellite instruments. Satellite retrievals of trace gas 17 concentrations constitute a rich source of information on key tropospheric species such as 18 nitrogen dioxide (NO₂), sulphur dioxide (SO₂), and formaldehyde (HCHO) that is beginning 19 to be exploited on an ever-larger scale. Ultraviolet-visible (UV/Vis) satellite observations are 20 being used to:

- evaluate the capability of models to simulate atmospheric concentrations of various
 species (e.g. Uno et al. [2007], Herron-Thorpe et al. [2010], Huijnen et al. [2010^a]),
- drive data assimilation experiments aimed at improving estimates of the atmospheric
 state (e.g. Wang et al. [2011], Inness et al. [2013], Miyazaki et al. [2014]),
- provide constraints on uncertain model inputs such as emission inventories through
 inverse modelling (e.g, Wang et al. [2007], Müller et al. [2008], Mijling and van der A
 [2012], Barkley et al. [2013]), and to identify new emissions sources (for instance
 newly built power plants [Zhang et al., 2009]),
- test processes influencing the lifetime of crucial chemical species (e.g. Schaub et al.,
 2007; Lamsal et al., 2010; Beirle et al., 2011; Stavrakou et al., 2013), or, more

broadly, the chemical regime of the atmosphere (e.g. Martin et al., 2004; Duncan et
 al., 2010)

3 When comparing model simulations to satellite measurements, both modelling errors and 4 measurement errors are usually taken into account. Measurement errors are often reasonably 5 well characterized (e.g. Boersma et al., 2004; De Smedt et al., 2008; Lee et al., 2009), but 6 modelling errors are more difficult to establish, because of the large number of uncertain 7 model processes, uncertain boundary (e.g. emissions) and initial conditions, and unresolved or 8 misrepresented aspects of atmospheric physics and chemistry. Modelling errors are best 9 characterized by comparing model simulations to observations. Unfortunately, the observations available for such comparisons are mostly limited in vertical range and regional 10 11 coverage such as in the case of ground-based networks, or they are merely sporadic in space and time, such as for aircraft campaigns. Satellite data records are based on robust retrieval 12 methods, provide global coverage, and cover decadal time spans. Satellite data has recently 13 14 been successfully used for dedicated modelling error studies (e.g. Lin et al., 2012; Stavrakou 15 et al., 2013).

16 When using satellite data, modellers need to be aware that most UV/Vis-retrievals generally contain little information on the vertical distribution of a species (the exception is 17 18 stratospheric ozone profile retrieval in the far UV of the spectrum, but this species will not be 19 considered in this study). Here we focus on the application of tropospheric UV/Vis retrievals, 20 and we limit ourselves to retrievals of tropospheric species NO₂, SO₂, and HCHO for comparison with models. These species are all relatively short-lived and their retrievals are 21 22 generally based on differential optical absorption spectroscopy (DOAS, Platt and Stutz [2008]). DOAS retrievals in the UV/Vis match relevant absorption cross section spectra to the 23 24 solar backscatter spectrum measured by the satellite instrument in order to infer the column integral (slant column density, expressed in molecules cm⁻²) of a species along the effective 25 26 atmospheric photon path. The subsequent retrieval step requires the conversion of the slant column density into a vertical column density, and this conversion depends on knowledge 27 (assumptions) of the state of the atmosphere, e.g. on the presence of clouds and aerosols, the 28 vertical distribution of the species, and surface properties. When these assumptions are very 29 30 different from the atmospheric state modelled by a chemistry transport model (CTM), this will lead to inflated differences between modelled (by, say, CTM 1) and retrieved columns 31 32 (aided by CTM 2). Such differences, however, can be avoided or in any case minimized, if the

1 user of satellite data accounts for the representativeness and averaging kernels of the satellite 2 data while interpreting model simulations. Representativeness here is defined as the context in which the satellite measurement holds, i.e. the horizontal coverage, the temporal 3 4 representativeness, and the vertical information content of the retrieval. It is the goal of this 5 study to provide guidelines on how users can take the representativeness of the UV/Vis column retrievals into account when comparing CTM simulations to satellite retrievals, and 6 7 by how much the model-retrieval differences would inflate if aspects of representativeness are 8 neglected.

9 In Sect. 2, we introduce the definitions and terminology for sources of error in the comparison of models and observations, and relate these to what is common practice in the data 10 11 assimilation community. In doing so, we follow the notation proposed by Ide et al. [1997], also used in relevant work by Rodgers and Connor [2003] and Migliorini et al. [2008]. 12 13 Section 3 will give an overview of the common features shared by various UV/Vis retrievals 14 with an emphasis on the assumptions made in the retrieval approach that are relevant to 15 modellers and other data users, and it provides a recipe for constructing an appropriate observation operator. Section 4 introduces the TM5 and GEOS-Chem models that we will 16 17 evaluate to demonstrate the nature and magnitude of representativeness errors. In Sect. 5, we discuss the error budgets associated with a confrontation of CTM simulations with satellite 18 19 measurements, and, in particular, how the representativeness errors contribute to that budget. 20 Section 6 presents the result of a practical assessment of representativeness errors made when 21 comparing global CTM simulations of tropospheric NO₂ to satellite measurements from the Ozone Monitoring Instrument, and provides recommendations on how to minimize these. 22

23 2 Comparing models and UV/Vis satellite measurements

24 2.1 UV-VIS satellite retrievals

Over the last two decades, tropospheric NO₂, SO₂, and HCHO columns have been retrieved from measurements by the GOME, SCIAMACHY, OMI and GOME-2 (on Metop-A and Metop-B) satellite sensors. The retrievals generally use a two-step approach, based on the DOAS-technique. In step 1, the reflectance spectra measured by the satellite instruments are modelled with a fitting routine that accounts for the spectral signatures from trace gas absorption, inelastic scattering, and (broadband) Rayleigh, Mie, and surface scattering. For each of the above species, spectral regions are selected where the absorption structures are 1 most distinct, and spectral interference from other species is minimal. The species' slant 2 column density is then calculated from the inferred absorption in combination with 3 knowledge of the species' absorption cross-section. Before converting the slant column 4 densities into tropospheric vertical columns, background corrections may be required to 5 account for the fact that a portion of the slant column has originated from the species' 6 absorption of light in the stratosphere.

7 In step 2 of the retrieval, the tropospheric slant column densities are converted into vertical 8 column estimates, using a radiative transfer (forward) model and forward model parameters, 9 that influence the retrieval. For DOAS UV/Vis retrievals, forward model parameters typically include the sensor viewing geometry, and best estimates of the surface albedo, terrain height, 10 cloud and aerosol properties (or an effective representation thereof), as well as the a priori 11 vertical distribution of the species (\mathbf{x}_a) of interest. The radiative transfer calculations are 12 expressed as so-called air mass factors, defined as the (forward) modelled ratio of slant (N_S) 13 and vertical columns (N_V), given the set of forward model parameters: $M = N_S/N_V$. 14 15 Tropospheric air mass factors have been shown to be very sensitive to choices for surface 16 albedo, for cloud correction, and for a priori vertical distribution, and, consequently, air mass 17 factor uncertainties are large, and dominate the retrieval error budget for tropospheric 18 columns (e.g. Boersma et al. [2004], Millet et al. [2006], Lee et al. [2009]).

19 Data users need to be aware of the important role played by clouds in UV/Vis-retrievals. With 20 the exception of elevated plumes resulting from volcanoes, lightning, and aircraft, most 21 tropospheric NO₂, SO₂, and HCHO generally resides in the lower atmosphere, close to their 22 surface sources. Clouds thus typically obscure the absorbing species from (satellite) view, 23 leading retrieval groups to advise against the use of their satellite data when taken under cloudy conditions. Trace gas retrievals under cloudy situations suffer from larger errors (e.g. 24 25 Schaub et al. [2006]), because the detectable fraction corresponds to the column above the cloud, leaving a so-called 'ghost column' below the cloud to be added somehow. Because 26 27 ghost columns are generally taken from climatology or a CTM, they do not contribute to the 28 measured information in any way, so that inclusion of columns under cloudy situations compromises a model - satellite comparison, unless the averaging kernels are taken into 29 account [Schaub et al., 2006]. This does not mean that satellite measurements taken under 30 cloudy conditions should not be used at all. In data assimilation systems, cloudy 31 32 measurements still provide valuable information on the abundance and vertical information of trace gases above the cloud, for instance for constraints on e.g. lightning-produced NO₂
 [Boersma et al., 2005] and in recent cloud-slicing techniques [Choi et al., 2014; Belmonte Rivas et al., 2015]. In cloud covered situations, it is essential to take the averaging kernels
 into account.

5 DOAS UV/Vis nadir retrievals are characterized by a vertical sensitivity that generally reduces with increasing atmospheric pressure, and require an a priori vertical profile of the 6 species \mathbf{x}_{a} to interpret the slant column (e.g. Palmer et al. [2001]; Richter et al. [2006]). 7 Because Rayleigh scattering of sunlight is more effective in the UV, fewer photons reach the 8 9 lower atmosphere in the spectral range where SO₂ has distinct absorption spectral features 10 (300-330 nm), compared to the spectral windows for HCHO (340-360 nm) or NO₂ (400-500 11 nm). This implies that the measurement sensitivity to species in the lower atmosphere is lowest for SO₂, followed by HCHO, and highest for NO₂. The contribution of the a priori 12 13 profile to the retrieved column increases with decreasing sensitivity of the measurement. Uncertainty in the species a priori vertical profile thus propagates stronger for SO₂ (up to 22% 14 15 error [Lee et al., 2009]), and somewhat less for NO₂ (10-15% error, e.g. Hains et al. [2010]; Vinken et al. [2014]). 16

17 This a-priori profile error contribution to model-satellite comparisons can be eliminated by 18 application of the averaging kernel to the model output (Eskes and Boersma [2003]; Boersma 19 et al. [2004]; Rodgers and Connor [2003]). The averaging kernel for UV/Vis retrievals 20 describes the relationship between the true column and the estimated, or retrieved column \hat{y}_o 21 where the hat denotes that the retrieval represents an estimated value of the true column:

22

$$\hat{y}_o = \mathbf{A} \cdot \mathbf{x}_{\text{true}} \tag{1}$$

with **A** the averaging kernel whose discretized elements can be described as $A_l = \frac{m_l}{M(x_a)}$, with m_l the scattering weights [Palmer et al., 2001], or box air mass factors for layer *l* (see Eskes and Boersma [2003], and Boersma et al. [2004] for more detail). Note that the retrieval problem has been linearised around $\mathbf{x_a} = 0$, related to the weak absorber character of the species, which implies that the a-priori state does not explicitly appear in Eq. (1).

28 **2.2** Model evaluation with UV/Vis satellite retrievals

A comparison between satellite measurements \hat{y}_o (e.g. the retrieved tropospheric NO₂ columns within a model grid cell), and the model state \mathbf{x}_m (e.g. the modelled vertical NO₂ distribution in the troposphere), in the form of measurement-minus-model departures (d) is
expressed as:

3

$$d = \hat{y}_o - \mathbf{H}\mathbf{x}_\mathbf{m} \tag{2}$$

4 with **H** the observation operator that describes the relation between the observed data and the 5 modelled state. Apart from the observation errors (σ_0 in the following) and the modelling errors (σ_m), we also need to take into account representativeness errors (σ_r) associated with 6 7 the fact that model simulations and satellite measurements provide different representations of 8 a geophysical quantity. We generalise the representativeness errors as the errors introduced in 9 a satellite-to-model evaluation by an incorrect description of the relation between the grid-cell mean concentrations and the satellite retrieval(s), i.e. we can think of them as errors in the 10 11 observation operator **H**. In data assimilation, representativeness errors are normally included in the observation errors (e.g. Jones et al., [2003]; Miyazaki et al., [2012]). 12

Substantial representativeness error may arise when the observation operator **H** is simplified and the model is not sampled in a manner fully consistent with the satellite observation. We can identify three types of representativeness errors associated with model-satellite comparisons:

- Spatial representativeness errors. Such errors will arise because models provide a
 spatially smoothed representation of the atmospheric state, whereas satellite
 measurements provide 'snapshots' at a particular local time, and often resolve
 variability at scales (pixels) smaller than the model grid cell.
- Temporal and meteorological representativeness errors. In applications focusing on
 clear-sky situations such as emission estimates, failure to sample the model for the
 same clear-sky conditions and overpass time as the satellite measurements, will lead to
 systematic sampling errors.
- Vertical representativeness errors. Because the sensitivity of the UV/Vis satellite
 retrievals is altitude-dependent (Palmer et al., [2001]), UV/Vis retrievals should be
 regarded as estimates of the state weighted by the averaging kernel (Eskes and
 Boersma [2003]). Neglecting the averaging kernel or vertical sensitivity of the
 retrieval in the comparison will inevitably introduce additional representativeness
 errors to the comparison in Eq. (2).

- 1 To minimize these representativeness errors in comparing CTMs and satellite measurements,
- 2 we recommend to follow the recipe given in Section 3.2 the box inset on this page.
- 3 This recipe on how to compare a CTM with satellite observations is a set of mathematical

Recipe for minimizing representativeness errors

1. The first step in comparing satellite observations to model simulations is to ensure that the satellite measurements are spatially representative for the area of the model grid cell. This is achieved by calculating the weighted average of all individual retrievals \hat{y}_t^{ρ} within the superobservation model grid cell over the entire area covered by all (valid) retrievals, where the weight is given by the pixel area w_t (in km²):

$$\hat{y}_{\theta} = \frac{\sum_{t} w_{t} \hat{y}_{\theta}^{\theta}}{\sum_{t} w_{t}} \tag{B1}.$$

If the model grid cell happens to be smaller than the satellite pixel, Eq. (B1) will reduce to $\hat{y}_{o} = \hat{y}_{\pm}^{o}$ for grid cells that are completely overlapped by a single satellite pixel (w_{t} =1).

2. The second step is to sample the CTM field sequence $\mathbf{x}_{\mathbf{m}}[t]$, here expressed as a discrete series of periodic fields with *t* an integer, when model time *t* is closest to the satellite overpass time t_{o} .

$$\mathbf{x}_{\mathbf{m}} = \mathbf{x}_{\mathbf{m}}[t]\delta[t], \text{ with } \delta[t] = \begin{cases} 0, \ t \neq t_{\sigma} \\ 1, \ t = t_{\sigma} \end{cases}$$
(B2)

The model sequence is sometimes also sampled with somewhat looser criteria, by requiring that the absolute model-satellite time difference stays within 1-2 hours (e.g. Martin et al. [2003]).

3. The third step is to apply the averaging kernel on the model vertical distribution $\mathbf{x}_{\mathbf{m}}$ to obtain the model estimate $\hat{y}_{\mathbf{m}}$ that can be directly compared to the observed state \hat{y}_{σ} :

$$\hat{\mathbf{y}}_{m} = \mathbf{A} \cdot \mathbf{x}_{m} = \sum_{l=1}^{L} A_{l} S_{l} \mathbf{x}_{m,l} \tag{B3}$$

where S_t are the components at the *l*-th vertical layer of an operator that executes a mass-conserving vertical interpolation or integration followed by a conversion to sub-columns (molec.em⁻²) in case the model vertical distribution $x_{m,t}$ is not yet given in those units. The product of the mathematical expressions (B2) and (B3) forms the observation operator **H** in Eq. (2), which describes the relation between the superobservation and the modelled state.

operations on satellite and model data. This is particularly relevant for short-lived species that have a high spatial and diurnal variability such as NO₂, SO₂, and HCHO (e.g. Boersma et al. [2008], Vrekoussis et al. [2009], Barkley et al. [2013]). Details of the approach may differ (e.g. spatial interpolation of the model state to the location of the pixel, averaging over different model times close to the satellite measurement time, replacing the a priori profile with the model profile in the retrieval), as long as the general principle of consistent sampling is observed. We advise against a comparison of the original satellite column (retrieved with a priori profile \mathbf{x}_a) to the model column \hat{x}_m because in that case differences between the a priori and modelled vertical profiles would inflate the overall error *d*, see Sect. 6 and recommendations in Sect. 2.3 of Boersma et al. [2004], and Duncan et al. [2014].

4 2.3 The role of clouds in UV/Vis retrievals

5 Data users need to be aware of the important role played by clouds in UV/Vis-retrievals. With the exception of elevated plumes resulting from volcanoes, lightning, and aircraft, most 6 7 tropospheric NO₂, SO₂, and HCHO generally resides in the lower atmosphere, close to their 8 surface sources. Clouds thus typically obscure the absorbing species from (satellite) view, leading retrieval groups to advise against the use of their satellite data when these are taken 9 under cloudy conditions. Trace gas retrievals under cloudy situations suffer from larger errors 10 (e.g. Schaub et al. [2006]), because the detectable column then corresponds to the column 11 12 above the cloud, leaving a so-called 'ghost column' below the cloud to be added. Because 13 ghost columns are generally taken from climatology or a CTM, they do not contribute to the 14 measured information in any way, so that inclusion of columns under cloudy situations compromises a model satellite comparison, unless the averaging kernels are taken into 15 16 account [Schaub et al., 2006]. This does not mean that satellite measurements taken under cloudy conditions should not be 17 18 used at all. In data assimilation systems, for above-cloud constraints on e.g. lightning-

19 produced NO₂ [Boersma et al., 2005], but also in recent cloud-slicing techniques [Choi et al.,

20 2014; Belmonte-Rivas et al., 2015], cloudy measurements provide valuable information on

21 the abundance and vertical information of trace gases above the cloud. In cloud-covered

22 situations, it is essential to take the averaging kernels into account.

23 **3** Theoretical model evaluation error budget

3.1 Sources of errors in evaluating CTMs with UV/Vis retrievals

A comparison between model simulations and satellite retrievals begins with a comparison of
 their theoretical capabilities. A model-satellite comparison will be influenced by:

- 1. modelling errors σ_m , related to an incomplete knowledge and description of the atmospheric state \mathbf{x}_m ,
- 29 2. retrieval errors σ_o , because of instrument noise and uncertainty in the (external) 30 forward model parameters, and

1 3. representativeness errors σ_r , arising from fundamental differences between the-2 atmospheric sampling by models and satellites, i.e. errors in the observation operator 3 H.

4 Assuming that these error terms are independent, the error analysis for a satellite-model 5 column difference $(\hat{y}_o - \mathbf{H}\mathbf{x}_m)$ can be written as:

6

$$\sigma = (\sigma_o^2 + \sigma_m^2 + \sigma_r^2)^{1/2}$$
(3)

7 with σ_o^2 the best estimate for the (relative) column retrieval errors, σ_m^2 for the (relative) 8 modelling error, and σ_r^2 the contribution to the error arising from the act of carrying out the 9 comparison itself (i.e. from errors in the observation operator). Some studies (e.g. Jones et al. 10 [2003]) include representativeness errors in the observation errors. Below we will show that 11 representativeness errors may contribute substantially to the overall error in satellite-model 12 confrontations.

13 The retrieval, modelling, and representativeness errors will all have systematic and random components. Here and in Sect. 6 we treat them as random errors, an approach usually 14 followed in retrieval and data assimilation studies [Rodgers, 2000]. In principle, one would 15 16 like to distinguish between the random and systematic contributions, but in practice this is very complicated, because many systematic contributions to retrieval and model errors are 17 18 only weakly correlated in space and time. Examples of subtle systematic retrieval effects are errors in individual albedo values with a small spatial correlation length but with 100% 19 20 correlation in time (for instance because residual cloud effects in the albedo climatology are 21 strongly variable from one location to the other [Kleipool et al., 2008]). When averaged over 22 a larger region such as the spatial extent of a coarse model grid cell, the impact of such errors tends to reduce. Likewise, models will suffer from systematic errors in for instance the 23 description of vertical transport. In particular circumstances, such as strong, small-scale 24 convective activity, such errors tend to be acute, but in an average sense, such as comparisons 25 26 aggregated over a month and a region, we may expect these errors to be smaller.

27 3.2 Recipe for minimizing representativeness errors

1. The first step in comparing satellite observations to model simulations is to ensure that the satellite measurements are spatially representative for the area of the model grid cell. This is achieved by calculating the weighted average of all individual retrievals \hat{y}_i^o within the 1 superobservation model grid cell over the entire area covered by all (valid) retrievals, where 2 the weight is given by the pixel area w_i (in km²):

 $\hat{y}_o = \frac{\sum_i w_i \hat{y}_i^o}{\sum_i w_i} \tag{4}.$

4 If the model grid cell happens to be smaller than the satellite pixel, Eq. (4) will reduce

3

5 to $\hat{y}_o = \hat{y}_1^o$ for grid cells that are completely overlapped by a single satellite pixel ($w_l = 1$).

6 2. The second step is to sample the CTM field sequence $\mathbf{x}_{\mathbf{m}}[t]$, here expressed as a discrete 7 series of periodic fields with *t* an integer, when model time *t* is closest to the satellite overpass 8 time t_o .

9
$$\mathbf{x}_{\mathbf{m}} = \mathbf{x}_{\mathbf{m}}[t]\delta[t], \text{ with } \delta[t] = \begin{cases} 0, \ t \neq t_o \\ 1, \ t = t_o \end{cases}$$
(5)

The model sequence is sometimes also sampled with somewhat looser criteria, by requiring
that the absolute model-satellite time difference stays within 1-2 hours (e.g. Martin et al.
[2003]).

13 3. The third step is to apply the averaging kernel on the model vertical distribution $\mathbf{x_m}$ to 14 obtain the model estimate \hat{y}_m that can be directly compared to the observed state \hat{y}_o :

15
$$\hat{y}_m = \mathbf{A} \mathbf{x}_m = \sum_{l=1}^L A_l S_l x_{m,l}$$
(6)

16 where S_l are the components at the *l*-th vertical layer of an operator that executes a mass-

17 conserving vertical interpolation or integration followed by a conversion to sub-columns

18 (molec. cm⁻²) in case the model vertical distribution $x_{m,l}$ is not yet given in those units. The

19 product of the mathematical expressions (5) and (6) forms the observation operator **H** in Eq.

20 (2), which describes the relation between the superobservation and the modelled state.

3.3 Representativeness errors in evaluating CTMs with UV/Vis retrievals

The total representativeness error σ_r is composed of horizontal representation errors, (temporal model) sampling errors, and vertical smoothing errors, and these three contributions may be assumed to be largely uncorrelated:

25
$$\sigma_r = (\sigma_h^2 + \sigma_t^2 + \sigma_v^2)^{1/2}$$
(7)

For an appropriate comparison between model simulations and satellite retrievals, it is important to sample the CTM as closely as possible to the satellite's sampling of the atmosphere (see Sect. 3.2). These may seem like trivial conditions for comparison, yet one or
 more of these conditions are often violated.

3 4 Data used in this study

4 4.1 Satellite data

5 In this study, we use tropospheric NO_2 retrievals from the Dutch OMI NO_2 (DOMINO) algorithm v2.0 [Boersma et al., 2011]. These retrievals proceed along the lines discussed 6 7 above, with spectral fitting of NO₂ in the 405-465 nm window [van Geffen et al., 2014], data 8 assimilation of the NO₂ slant columns in the TM4 chemistry transport model [Williams et al., 9 2009] to estimate the stratospheric background [Dirksen et al., 2011], and final conversion of the tropospheric slant columns with air mass factors based on radiative transfer calculations 10 11 with the DAK model. In the DOMINO algorithm, altitude-dependent AMFs are interpolated from pre-calculated look-up tables using the best available information on the satellite 12 viewing geometry, surface albedo [Kleipool et al., 2008], terrain height (3km-resolution 13 elevation data provided with Aura data). Subsequently, the local altitude-dependent AMFs are 14 15 combined with the predicted local vertical NO₂ distributions (from TM4), to produce the (tropospheric) air mass factors. The air mass factor step also includes a correction for the 16 17 temperature-dependency of the NO₂ absorption cross-section [Boersma et al., 2004], because only the 220 K cross section is used in the spectral fit. The DOMINO v2.0 data have been 18 19 evaluated in a number of validation exercises (e.g. Irie et al. [2012], Ma et al. [2013], Lin et al. [2014]), showing their quality and use, although a number of relevant improvements is 20 21 planned and currently being implemented (Maasakkers [2013], van Geffen et al. [2014]). 22 DOMINO v2.0 has been used in many applications and model studies (e.g. Stavrakou et al. [2013], Castellanos et al. [2014], McLinden et al. [2014], Verstraeten et al. [2015]), which 23 24 makes the data product well-suited for evaluating satellite-to-model comparisons and the 25 errors associated with such comparisons, which is the purpose of this study.

Chemistry-transport models (CTMs) are the central tools to simulate tropospheric concentrations of NO₂, SO₂, and HCHO, and to help interpret and use satellite measurements of these species. For the short-lived species studied here, previous studies indicate modelling biases of $\pm 20-30\%$ for NO₂ (e.g. van Noije et al. [2006]), and 20-50% for HCHO (e.g. Dufour et al. [2009]; Williams et al. [2012]) over regions with substantial pollution.

1 4.2 TM5

2 We use the TM5, the global 3-D CTM version 3.0 [Huijnen et al., 2010] with a grid of 3° longitude \times 2° latitudes \times 34 vertical layers, and a model top at 0.1 hPa [Krol et al., 2005]. 3 The TM5 model is used in many studies for atmospheric chemistry (e.g. Williams et al. 4 5 [2014]), aerosol haze (e.g. von Hardenberg et al. [2012]), data assimilation, and inversion 6 applications (e.g. Hooghiemstra et al. [2012]; Krol et al. [2013]). The model is driven by 7 ERA-Interim meteorological reanalysis data from the European Centre for Medium Range 8 Weather Forecats (ECMWF) [Dee et al., 2011] and the base time step is 1 hour. In the version 9 used here, TM5 operates with Carbon Bond Mechanism 4 chemistry [Gery et al., 1989] to describe the production of ozone, hydrogen oxide radicals (HO_x=OH+HO₂) and oxidation of 10 nitrogen oxides (NO_x=NO+NO₂), SO₂, and volatile organic compounds (VOCs), with 40 11 species, 64 gas-phase, and 16 photolysis reactions. In TM5, SO₂ is oxidized in clouds and on 12 13 aerosols, and nighttime hydrolysis of N₂O₅ into nitric acid (HNO₃) is parameterized with a 14 global mean uptake coefficient of 0.02 following recommendations by Evans and Jacob 15 [2005]. NO_x emissions are from the RETRO inventory for the anthropogenic sectors (Regional Emission inventory in ASia – REAS for Asia) with a total of 33 Tg N/yr, 9 Tg N/yr 16 17 from soil, 5 TgN/yr from biomass burning (from the Global Fire Emissions Database v2 (GFED2) van der Werf et al. [2006]), and 6 TgN/yr for lightning. Global anthropogenic SO₂ 18 19 emissions are taken from the AeroCom project at 108 Tg SO₂/yr [Dentener et al., 2006]. Biogenic VOC emissions, including the important HCHO and its precursor isoprene, are from 20 the ORCHIDEE database [Lathière et al., 2006], and are 10 TgC/yr for HCHO and 565 Tg 21 22 C_5H_8 /yr for isoprene. We simulated the year 2006 with a one-year spin-up.

TM5 simulations of NO₂ and HCHO have been evaluated by Huijnen et al. [2010] and Williams et al. [2012]. These studies indicate that tropospheric NO₂ columns in TM5 are 20-30% low compared to DOMINO v2.0 columns, but the model captures the seasonality, and shows realistic vertical distributions of NO₂ relative to INTEX-B aircraft measurements. TM5 captures the seasonality of HCHO tropospheric columns but also overestimates these columns by 0-50%, partly because of inadequate photolysis rates in the model [Williams et al., 2012].

29 **4.3 GEOS-Chem**

We also use the GEOS-Chem model, v9-02i, with a grid of 2.5° longitude $\times 2^{\circ}$ latitude $\times 47$ vertical layers, and the model top at 80 km. The GEOS-Chem model is a CTM in use by a

large community of scientists for a wide range of applications including, shipping NO_x 1 plume-in-grid chemistry [Vinken et al., 2011], and estimating isoprene and ammonia 2 3 emissions (e.g. Millet et al. [2008]; Paulot et al., [2014]). GEOS-Chem is driven by GEOS-5 4 meteorological fields from NASA GMAO, with a time step of 30 minutes. As TM5, GEOS-5 Chem uses a condensed O₃-NO_x-HO_x-VOC-aerosol chemistry scheme (described in Mao et al. [2010] and references therein). The standard chemistry scheme has 66 species, and 236 6 7 chemical reactions. GEOS-Chem takes into account heterogeneous chemistry on aerosol and 8 cloud particles [Mao et al., 2010], including the uptake of N₂O₅ on aerosols leading to 9 nighttime HNO₃ formation following the parameterization by Evans and Jacob [2005]. Anthropogenic NOx emissions are from the global EDGAR 3.2FT2000 inventory [Olivier 10 11 and Berdowski, 2000], but these are replaced by regional inventories over various continents. Other NO_x emissions in GEOS-Chem include soil, lightning, biomass burning, biofuel, 12 13 aircraft and ship, resulting in a global total source of 51.5 TgN/yr for 2006 (similar to TM5 14 with 53 TgN/yr for the same year). A two-year spin-up was performed (2004-2005), and GEOS-Chem output was stored for the year 2006. For more details on the GEOS-Chem 15 16 simulation, see Vinken et al. [2014].

GEOS-Chem simulations of tropospheric NO₂ columns have been evaluated before by Lamsal et al. [2010] and Lin [2012], who found, similar to the TM5 evaluation discussed above, that the model underestimates tropospheric NO₂ by 20-35% (over China). Zhang et al. [2012], in a study targeting nitrogen deposition over the United States, found excellent agreement between the modelled and OMI-observed spatial distribution of tropospheric NO₂, but also underestimates of 10% in the northeastern US, and 40% locally in southern California, were also evident.

24 5 Representativeness errors

25 **5.1** Horizontal representativeness errors

If the complete spatial extent of a model grid cell is covered with valid retrievals, a good comparison is straightforward because a spatially fully representative area average can be calculated. For partly covered cases, the difficulty lies in estimating the magnitude of the (horizontal representativeness) errors associated with limited coverage of a model grid cell. One way to calculate a representative grid cell average is by averaging all valid satellite observations that were taken within the boundaries of the grid cell within a given model time

step, as in Eq. (3), with w_i the fractional grid cell coverage defined as A_{pixel}/A_{cell} with A_{pixel} the 1 2 area (in km²) covered by the fraction of the satellite pixel that falls within the boundaries of the model grid cell with area A_{cell} (in km²). In this manner, one obtains a 'superobservation' 3 4 that may be considered as representative for the grid cell average (Dirksen et al. [2011]; 5 Miyazaki et al. [2013]). In some model-satellite confrontations, the number of satellite retrievals is thinned out to 1 per grid cell, but we advise against such an approach in view of 6 7 the strong sub-grid variations and the considerable errors in individual measurements. In 8 many global applications, the spatial resolution of the model is coarser than the resolution of 9 the satellite observations. In that case, we recommend computing a set of superobservations defined as the grid cell average trace gas column \overline{N} : 10

11
$$\overline{N} = \frac{\sum_{i=1}^{n} w_i N_i}{\sum_{i=1}^{n} w_i}$$
(5)

with w_t the fractional grid cell coverage of retrieval N_t , defined as A_{nixel}/A_{cell} with A_{nixel} the area 12 (in km²) covered by the fraction of the satellite pixel that falls within the boundaries of the 13 model grid cell with area A_{cell} (in km²), and *n* the total number of valid retrievals within the 14 eell. We caution against applying additional weighting by the individual retrieval errors in Eq. 15 16 (4). Because, by nature of the DOAS approach, retrieval errors are largest for large column 17 values (see e.g. Boersma et al. [2004]), error weighting would skew the average to the lower values in the distribution. The measurement error for superobservations can be calculated 18 from area-weighting the individual pixel errors $\sigma_{o,i}$ to provide an area-weighted average 19 (statistical) retrieval error σ , and by accounting for a partial correlation in the errors between 20 21 pixels as in Eskes et al. [2003] (see Appendix B for a derivation):

$$\sigma_o = \sigma \sqrt{\frac{1-c}{n} + c} \tag{8}$$

with the second term on the right side representing the error correlation (*c*) between the *n* retrievals. Miyazaki et al. [2012] propose c=0.15, based on the consideration that errors in clouds, albedo, a priori profile, and aerosol in retrievals are typically correlated in space, but they acknowledge that the exact number is difficult to estimate.

22

Some studies take a different approach than the superobservations proposed in Eq. (4) and interpolate the model simulations to the centre of a satellite pixel, but the difficulty with this approach is the questionable spatial representativeness of the interpolated model value, especially if the model grid cells cover a larger area than the satellite pixels.

1 Both individual pixel errors and representativeness errors contribute to the total error in the 2 superobservation. Following Miyazaki et al. [2012], we calculate the horizontal 3 representativeness error σ_r as a function of the total fractional coverage achieved by all valid pixels by random reduction of the number of retrievals used to calculate the mean grid cell 4 5 value. For homogeneous scenes with little variability of NO₂, SO₂, or HCHO, such errors will obviously be small. But for grid cells covering strong inhomogeneous sources of air pollution, 6 7 such as megacities or coal plants, we may expect the area average to depend strongly on the 8 spatial sampling. Figure 1 illustrates the horizontal representativeness error as a function of 9 total fractional coverage for one polluted model grid cell, here taken over the eastern United States (greater New York City), at two resolutions, i.e. $3^{\circ} \times 2^{\circ}$ (typical for a global CTM) and 10 $0.5^{\circ} \times 0.5^{\circ}$ (regional CTM). To calculate the horizontal representativeness error, we randomly 11 reduced the number of pixels n in Eq. (4) first by 1, then by 2, and so on, until there was only 12 13 one pixel left, to obtain new estimates \hat{y}'_{0} . We repeated this 100 times and interpret the root 14 mean squared difference with the original \hat{y}_0 as the horizontal representativeness error, which 15 is zero in situations of full coverage. Complete coverage of the grid cell is typically achieved by more than 100 OMI pixels in the case of $3^{\circ} \times 2^{\circ}$ resolution grid cells, and by ± 5 pixels¹ for 16 $0.5^{\circ} \times 0.5^{\circ}$. The horizontal representativeness errors appear higher for the $0.5^{\circ} \times 0.5^{\circ}$ than for 17 the $3^{\circ} \times 2^{\circ}$ grid cell, due to the smaller sample (*n*=5) size and the strong spatial gradients over 18 19 the central New York area for the higher resolution model. For models with higher spatial resolution $(0.5^{\circ} \times 0.5^{\circ})$, there is less tolerance for reduced area coverage over strongly 20 21 inhomogeneous areas such as central New York, as indicated by the steeper 22 representativeness error increase with reduced cover (blue dashed line in Figure 1). This 23 reflects the more heterogeneous distribution of polluted NO₂ column values for the high-24 resolution model with a small sample (5 pixels) than for the coarse resolution with a large sample (> 100 pixels). The $3^{\circ} \times 2^{\circ}$ case with complete area coverage by OMI NO₂ pixels (on 25 17 July 2006) illustrates the potential for horizontal representativeness errors. For a fractional 26 27 coverage of 0.5, the horizontal representativeness error increases to 10-15%, which is still 28 considerably smaller than the 20-30% errors in the satellite measurements themselves. For 29 fractional coverage of 0.1 however, the representativeness error increases to 35%, a level that exceeds the theoretical NO₂ retrieval error (Boersma et al., 2011) and NO₂ validation errors 30

¹ Because OMI pixel sizes vary with viewing zenith angle (largest pixels at the edge of the swath), the exact number of pixels covering a model grid cell depends on which part of the OMI swath covers the grid cell.

(e.g. Irie et al., 2012). However, by averaging over multiple days, the representativeness error
 can be reduced further, depending on the day-to-day variability of the columns. Table S1
 (Supplementary Material) shows the statistics of a comparison between monthly mean
 observed and simulated columns over the greater eastern United States in July 2006, for
 different degrees of fractional coverage required.

6 In data assimilation systems, any fractional coverage may be used as long as the horizontal 7 representativeness error is well described and accounted for along with the observation error. 8 This can be achieved by adding in quadrature the measurement error and representativeness 9 error $\sqrt{\sigma_{N,o}^2 + \sigma_{N,r}^2}$ to represent the overall superobservation error.





Figure 1. Relative horizontal representativeness errors as a function of the covered fraction of one model grid cell in the case of OMI tropospheric NO₂ columns for polluted area(s) (mean column 5×10^{15} molec.cm⁻²). The black line indicates the error as a function of the fractional coverage for a $3^{\circ} \times 2^{\circ}$ grid cell over the area of New York City on one day (17 July 2006, 114

1 OMI pixels). The blue asterisks indicate the mean error as a function of fractional coverage 2 for various $0.5^{\circ} \times 0.5^{\circ}$ grid cells on 17 July 2006.

Temporal representativeness errors related to clouds

3

5.2

4 In the case where UV/Vis satellite retrievals of the tropospheric column are used for air pollution applications (taken under cloud-free situations, see e.g. Schaub et al. [2006], Millet 5 6 al. [2006], Geddes et al. [2012]), both measurements and models should be sampled under similar clear-sky situations. As long as the model appropriately simulates the effects of clouds 7 8 on photolysis rates, this ensures that measurement and model represent the trace gas 9 concentrations under similar photochemical regimes. Failure to sample the model on clear-sky days only, will introduce a bias in the modelled average. Short-lived trace gases may have a 10 longer lifetime against photochemical loss in situations with overhead clouds (assuming they 11 12 are represented well in models), when actinic fluxes and temperatures are lower and chemistry slower than in clear-sky situations. For trace gases whose emissions reflect distinct 13 14 anthropogenic patterns, it is also necessary to sample the model according to the observations, in order to properly weigh well-documented weekend (e.g. Beirle et al., 2003; Boersma et al., 15 2009) and national holiday reductions (Lin et al., 2011) when calculating the model average. 16

17 We first evaluate the TM5 model's ability to simulate the effective cloud cover as observed 18 by OMI at 13:30 hrs local time. Cloud cover (and cloud optical thickness) data in TM5 are 19 hourly interpolated from 3-hourly pre-processed ECMWF fields [Huijnen et al., 2010]. Since the OMI cloud retrieval reports effective cloud fractions, based on the assumption that clouds 20 21 are optically thick (optical thickness of 40, with a corresponding cloud albedo of 0.8) 22 [Acarreta et al., 2004; Stammes et al., 2008], we converted the TM5 geometrical cloud cover into an effective fraction comparable to the OMI observations. To do so, we used the 23 24 maximum-random overlap assumption [Morcrette and Jakob, 2000] to compute the total geometrical cloud cover and total cloud optical thickness from the vertically resolved cloud 25 26 cover and optical thickness in TM5. We used the modelled relationship between the total cloud optical thickness for a liquid water cloud and its spherical cloud albedo in Buriez et al. 27 28 [2005] to calculate the effective cloud albedo associated with each grid cell's cloud cover. Finally, we weighted the total geometric cloud cover with the ratio of the effective cloud 29 30 albedo to 0.8, the value assumed for all clouds in the OMI retrieval (Acarreta et al. 2004]; Stammes et al. [2008]). For more details we refer to the Appendix C. 31

Figure 2 shows monthly mean effective cloud fractions as retrieved from OMI and simulated 1 2 with TM5 for February and August 2006. The model was sampled within 30 minutes of the 3 OMI overpass time of 13:30 hrs, and model and satellite were matched in space and time for 4 further analysis. We see that TM5 captures the spatial patterns observed by OMI, with low 5 cloud fractions over the subtropics, and high cloud fractions over the tropical ITCZ and the middle-to-high latitudes (> 40°). Largest differences occur at the edges of areas flagged as 6 7 snow-covered in the OMI retrieval (February 2006), and over areas where TM5 predicts cloud 8 optical thickness to exceed 40, such as over the tropics, where ice clouds often occur (and the 9 relationship for water clouds from Buriez et al. [2005] is less valid).

10 To evaluate the simulated effective cloud fractions, we report the correlation coefficient, 11 mean bias, and root mean square error relative to the OMI-observed cloud fractions over Europe for February and August 2006. Figure 2 shows significant positive correlation 12 13 between TM5 and OMI effective cloud fractions over Europe both in February (r=0.70, 14 n=3379) and August (r=0.75, n=4665). The mean bias between TM5 and OMI is -0.08 in 15 February and +0.02 in August, and the root mean square error is 0.23 in February and 0.20 in August. The agreement between TM5 and OMI, while far from perfect, suggests that TM5 16 has some success in simulating the contrast between 'cloud-free' ($f_{OM} < 0.2$) and 'cloudy sky' 17 $(f_{OM}>0.2)$ situations, i.e. the likelihood that OMI reports a clear-sky scene, while TM5 18 19 simulates a cloudy sky, and vice versa is <20% and <14%, respectively.





2 Figure 2. Monthly average effective cloud fraction observed from OMI (upper panels) and 3 simulated by TM5 based on ECMWF meteorological fields (middle panels) in February (left 4 column) and August 2006 (right). Cloud fractions have been selected only for those days and locations that had a successful OMI O₂-O₂ retrieval. Grey areas indicate less than 3 successful 5 6 coincidences. Bottom panels: scatterplot of daily pairs of OMI (x-axis) and TM5 cloud 7 fractions (y-axis) in February 2006 (left) and August 2006 (right) over Europe (10° W-30° E; 8 35°-60° N). The colours indicate the number of times a particular grid cell has been filled, 9 where light blue corresponds to $2\times$, green $3\times$, yellow $4\times$, orange $5\times$, red $6\times$, and magenta to 10 7× or more. TM5 effective cloud fractions can be expressed as $-0.10 + 1.06 f_{OMI}$ (February) 11 and $-0.01 + 1.07 f_{OMI}$ (August).

Figure 3 shows a box and whisker plot for OMI and TM5 effective cloud fractions over 1 2 Europe in February and August 2006. The figure indicates that for OMI measurements of effective cloud fractions smaller than 0.2, TM5 reproduces similar small effective cloud 3 fractions (February median OMI: 0.09, TM5: 0.06; August median OMI: 0.05, TM5: 0.04). 4 5 For days and locations when OMI observes effective cloud fractions larger than 0.2 (February: 0.59, August: 0.47), TM5 simulates comparable high effective cloud fractions 6 7 (January: 0.49, July: 0.45), providing some confidence in the TM5 model, driven by ECMWF 8 meteorological fields, to capture the observed effective cloud fractions.



9

10 Figure 3. Box and whisker plots for OMI (black) and TM5 (red) effective cloud fractions 11 over Europe in February 2006 (left panel) and August 2006 (right panel). The two left boxes 12 of each panel indicate the clear sky situations when the OMI cloud fraction < 0.2. The centre line of the boxes indicate the median cloud fraction, the upper and lower edges indicate the 13 25th and 75th percentiles and the lower and upper whiskers represent the minimum and 14 maximum value in the sample. For February 2006, the sample consisted of 3379 pairs (737 15 clear sky, 2642 cloudy), and for August, the sample size was 4665 (1991 clear sky, 2674 16 17 cloudy).

Figure 4(a) shows a comparison of average TM5 tropospheric NO₂ columns simulated under 18 clear-sky and cloudy situations over Europe in February and August 2006. TM5 was sampled 19 20 for polluted situations (cells with monthly mean NO₂ columns in excess of 1.0×10^{15} molec.cm⁻²) between 12:00-15:00 hrs local time, on days with clear skies and on days with 21 22 cloud-cover. Under clear-sky situations, TM5 simulates tropospheric NO₂ columns that are on 23 average 15-20% lower than under cloudy circumstances, in line with in situ observations 24 reported by Boersma et al. [2009] and Geddes et al. [2012] over Israeli and Canadian cities, respectively. Both in February and August, the clear-sky mean NO₂ column is 12% below the 25

1 28-day monthly mean in February and 31-day monthly mean in August. Although we cannot 2 rule out that other effects than enhanced photochemical loss may have contributed to lower 3 NO_2 columns over the polluted grid cells (e.g. increased ventilation or deposition) on clear-4 sky days, a comparison of NO_2 columns for all European grid cells showed that the 5 geometrical mean of the local clear-sky to cloudy column ratios was 0.74 in February and 6 0.89 in August, suggesting that reduced clear-sky NO_2 columns presented in Fig. 4 show a 7 robust effect.

8 The results for August 2006 indicate that clear-sky sampling of the model is also relevant for 9 HCHO in the growing season (Fig. 4(b)). Average HCHO columns are 12% higher under 10 clear-sky situations than on cloudy days and the clear-sky mean HCHO column is 8% higher 11 than the all-sky monthly mean (August 2006). In winter, HCHO concentrations are generally 12 low over Europe and differences between clear and cloudy sky are well below the detection 13 limit of UV-Vis satellite sensors.



14

Figure 4.(a) monthly mean tropospheric NO₂ columns simulated by TM5 for polluted grid cells (with all-sky monthly means > 1.0×10^{15} molec.cm⁻², *n*=18 in February, *n*=17 in August). The blue bars represent the average of the tropospheric NO₂ column sampled on days when the OMI cloud fraction was smaller than 0.2. Light blue: average for columns sampled when OMI cloud fraction > 0.2. (b): monthly mean TM5 HCHO columns for clearsky and cloudy situations (*n*=18 in February, for August: all-sky monthly mean > 7.5×10^{15} molec.cm⁻², *n*=12).

Exclusive sampling of the model on clear-sky days is important, because photolysis rates J[NO₂] in the lower troposphere are significantly higher on those days and can be simulated well by TM5 [Williams et al., 2012], so that NO₂ columns will be systematically lower. The differences between HCHO columns sampled on clear-sky and cloudy days are somewhat smaller than for NO₂ columns because both the formation and destruction of HCHO are driven by photochemistry. Nevertheless, the stronger summertime production of HCHO from the (OH-driven) oxidation of methane and especially isoprene outpaces the increased loss of HCHO through photolysis and oxidation [Fried et al., 1997] on clear-sky days compared to cloudy days, in line with observations (e.g. Munger et al. [1995], Cerquiera et al. [2003]).

7 To estimate the magnitude of the temporal representativeness errors arising from the 8 particular choice of model sampling, we evaluated the satellite-model comparison results for 9 different sampling strategies. Again, we use the averaged ratio of satellite measurements to model simulations (\hat{y}_o/\hat{x}_m) , and the spatio-temporal correlation coefficient, as appropriate 10 indicators of representativeness errors. Since the model - measurement bias may well be due 11 to unrelated systematic errors in either the CTM (emissions, chemistry) or the satellite 12 retrievals, we are not concerned with the absolute value of the measurement-to-model ratio, 13 14 but we are interested in the sensitivity of the ratio to various sampling strategies. We tested 15 four strategies for comparing tropospheric NO₂ over large polluted regions: (A) both OMI 16 (for OMI effective cloud-fraction) and TM5 (TM5 effective cloud fraction) collocated and 17 sampled for mostly clear-sky scenes only at the OMI overpass time of 13:30 hrs, (B) OMI and 18 TM5 collocated and co-sampled for situations with OMI effective cloud radiance fractions < 0.5^2 , (C) OMI sampled for situations with OMI effective cloud radiance fractions < 0.5, but 19 TM5 more loosely sampled for OMI effective cloud fractions < 0.6, and (D) OMI sampled for 20 21 situations with OMI effective cloud radiance fractions < 0.5, but TM5 sampled for all days in 22 the month (i.e. no temporal collocation except for appropriate overpass time). Strategy (A) is 23 considered to be optimal, but to our knowledge has not been applied in studies to date. 24 Strategy (B) has been followed in numerous studies, and relies on the assumption that CTMs 25 capture the observed cloud cover well. In spite of its erroneous co-sampling with the satellite 26 measurements, strategy (D) has also been used frequently, and therefore we tested its impact 27 on the temporal representativeness errors. Finally, strategy (C) holds middle ground between 28 (B) and (D). Figure 5 shows that the model-to-measurement ratio shows substantial dependence on the comparison strategy, especially in Winter. The differences between 29 strategies (A) and (B) are negligible, but with strategy (D) the OMI/TM5 ratio drops more 30

² The cloud radiance fraction is defined as the relative contribution of top-of-atmosphere radiance received by the cloud part of the pixel. A cloud radiance fraction of 0.5 corresponds to a geometric cloud fraction of ± 0.2 .

1 than 25% below the values obtained by strategies (A) and (B). These strategies also 2 demonstrate that strategy (D) leads to a reduced capacity of the model to explain the observed 3 variability in the NO₂ spatial patterns, with R^2 dropping almost 10% (from 0.64 to 0.55 in 4 Winter and from 0.66 to 0.59 in Summer).



Figure 5. Impact of sampling strategy on monthly averaged OMI:TM5 ratio of tropospheric NO2 columns (black dots) and on spatial correlation coefficient (\mathbb{R}^2 , blue dots) over the eastern United States (30°-44° N, 90°-72° W). Left panel: ratio and \mathbb{R}^2 for February 2006 (*n*=28). Right panel: August 2006 (*n*=32). Grid-cells were selected in the comparison when the covered fraction exceeded 0.5. The dashed black line shows the normalized OMI:TM5 ratio for strategy (A), and the dashed grey line shows the \mathbb{R}^2 for strategy (A) as a guide to the eye.

5

Analyses for other regions showed similar results as in Fig. 5. These results imply that for applications of satellite data such as emission estimates or model evaluations, substantial systematic errors may occur in the final estimate, if sampling strategies such as (D) are used. We therefore strongly discourage the use of such comparison strategies, as they lead to considerable temporal representativeness errors, and, thus, systematic underestimations in measurement:model ratios.

5.3 Vertical representativeness errors

Here we evaluate the representativeness errors introduced in a satellite-model comparison if 2 the averaging kernel is not accounted for. To illustrate the way the kernels work, Figure 6 3 4 shows GEOS-Chem NO₂ vertical profiles with and without the averaging kernel applied over the Beijing grid cell on clear-sky days with excellent spatial coverage (18 February and 23 5 6 August 2006). On both days, application of the kernel leads to a higher value for the model 7 column, reflecting the relatively larger amounts of NO₂ aloft in GEOS-Chem simulations 8 compared to the a priori TM4 NO₂ profiles. The lower panels show that on two other clear-9 sky days (17 February and 31 August 2006) the kernel has only little effect on the GEOS-Chem tropospheric NO₂ column. On these days, the TM4 a priori and GEOS-Chem NO₂ 10 profiles show similar, less pronounced vertical distributions. Nevertheless, in Figure 7 we see 11 that, on average, for February and August 2006, the OMI averaging kernels result in increases 12 13 in GEOS-Chem NO₂ columns over Beijing of 15% (February) and 8% (August), and a closer 14 agreement with OMI NO₂ retrievals. This result can be understood from the stronger vertical 15 mixing in the GEOS-Chem model compared to TM4, rather than from differences in NO_x emissions or chemistry between models (NO2 amounts are quite similar between TM4 and 16 GEOS-Chem over Beijing in 2006). 17



1

Figure 6. Vertical averaging kernel (black dashed line) and NO₂ profiles simulated by GEOS-Chem (blue), TM4 (red, a priori profiles in OMI NO₂ retrieval), and GEOS-Chem convolved with the averaging kernel (purple) following Eq. (6); (a) 18 February 2006, (b) 23 August 2006 over the Beijing grid cell (centered on 40°N, 116.25°E), (c) 17 February 2006, and (d) 31 August 2006. The numbers given in blue, purple, and red indicate the tropospheric vertical NO₂ columns in GEOS-Chem and TM4.

8 The above finding does not have general validity in the sense that applying the kernel on any 9 other model will also result in a tropospheric column increase. Applying the kernels to NO₂ 10 profiles from a model with weaker vertical mixing than TM4 (rather than generally stronger 11 vertical mixing as in the case of GEOS-Chem) is likely to reduce those columns. Figure S1 in the Supplementary Information shows as much for the North Sea grid cell in February 2006,
 when GEOS-Chem exceeds TM4 NO₂ concentrations below 900 hPa, and for Siberia in
 August 2006, when GEOS-Chem simulates a substantially enhanced tropospheric NO₂
 column compared to TM4.



5

Figure 7. Monthly mean averaging kernel (black dashed line) and NO₂ profiles simulated by
GEOS-Chem (blue), TM4 (red, a priori profiles in OMI NO₂ retrieval), and GEOS-Chem
convolved with the averaging kernel (purple) following Eq. (6); left panel: February 2006,
right panel: August 2006 over the Beijing grid cell (centered on 40°N, 116.25°E). The
numbers given in blue, purple, and red indicate the tropospheric vertical NO₂ columns in
GEOS-Chem and TM4.

We next compare the monthly averaged GEOS-Chem tropospheric NO₂ column fields for February and August 2006 with and without the kernels applied. Figure 8 shows that applying the kernel leads to substantial increases of up to 2×10^{15} molec. cm⁻² in the columns for the polluted source regions in the northern hemisphere (eastern USA, Europe, and China). At the periphery of these regions in wintertime, and over regions with possible biomass burning in summer, we see that the smoothed columns can be lower than the original columns, indicating

- 1 that the GEOS-Chem vertical NO₂ profile is more skewed towards the surface than the TM4 a
- 2 priori in those situations, as confirmed by the profiles shown in Figure S1.

3



Figure 8. Difference between monthly mean GEOS-Chem with AK (Eq. (6)) and GEOSChem tropospheric NO₂ columns without AK for February 2006 (upper panel) and August
2006 (lower panel). Only grid cells with more than 3 days of better than 40% coverage of
clear-sky pixels have been selected.

8 Here we evaluate the level of agreement between the original GEOS-Chem and OMI NO₂ 9 columns, compared to the level of agreement between the kernel-based GEOS-Chem and 10 OMI NO₂ column for the polluted source regions in the northern hemisphere, as the differences provide a measure of the representativeness errors that can be avoided by using 11 12 the averaging kernel. Figure 9 shows the agreement between OMI and the GEOS-Chem NO₂ 13 columns with and without kernel over Europe in February and August 2006. The upper panels indicate that the spatial correlation between the model and OMI tropospheric columns 14 improves when the kernel is applied on the model NO₂ profiles, especially in Summer when 15 16 differences between the TM4 a priori and GEOS-Chem NO₂ profile shapes are strong. Application of the kernel also results in geometric mean OMI:GEOS-Chem ratios with 17 smaller uncertainty intervals at values of 1.15^{1.82}_{0.73} (February) and 1.24^{1.78}_{0.86} (August) compared 18 to $1.13_{0.70}^{1.82}$ and $1.42_{0.94}^{2.14}$. We find similar results over the eastern United States and China (see 19 20 Table 1 in Sect. 7). Figure 9(d) further supports the notion that application of the kernel

allows for a better-constrained evaluation of the model, as witnessed by the more peaked and 1 2 narrower histogram of satellite; model ratios. We conclude that sampling the model according to the averaging kernel is especially relevant in Summer, and improves the satellite-model 3 4 evaluation by removing differences between (TM4 apriori and GEOS-chem) profile shapes 5 contributing to the discrepancies [Boersma et al., 2004]. Neglecting the kernels for GEOS-Chem would lead to up to 15% stronger discrepancies between OMI and GEOS-Chem, and 6 7 this portion could be wrongfully attributed in a model evaluation to e.g. too low NO_x emissions, or too fast NO₂ removal by chemistry or deposition. Appendix D presents an 8 9 alternative to the application of the averaging kernel by providing a recipe to replace the a priori profile used in the retrieval by the profile from the CTM under evaluation. Such a 10 11 recipe results in a modified retrieval that can be directly compared with the CTM under 12 evaluation.



13

Figure 9. Comparison between monthly average OMI and GEOS-Chem tropospheric NO₂ columns over Europe in February 2006 (left panels) and August 2006 (right panels); (a) scatter diagram of monthly average GEOS-Chem with AK (black circles) and GEOS-Chem without AK (grey circles) vs. OMI tropospheric NO₂ columns for February 2006. The black

and grey lines indicate the geometric mean of the OMI:GEOS-Chem ratio; (b) as for (a) but
for August 2006; (c) histogram of per-grid cell OMI-to-GEOS-Chem with AK tropospheric
NO₂ column ratios (black bars) and OMI-to-GEOS-Chem without AK ratios (grey bars) for
February 2006; (d) as (c) but for August 2006. Only grid cells with more than 3 days of better
than 40% coverage of clear-sky pixels have been selected.

6 6 **Combined representativeness errors**

To obtain an estimate of typical, overall representativeness errors in model evaluations with UV/Vis satellite measurements, we define three types of model evaluations, executed with increasing degree of detail. We again evaluate tropospheric NO_2 from the GEOS-Chem model here (with OMI NO_2 retrievals), as this model is sufficiently different from the TM4 model used to provide the a priori profiles in the OMI retrievals. The three types of evaluations can be characterised as advanced, common, and naïve:

- (A) advanced evaluation: accounting for sufficient spatial coverage and appropriate
 temporal representativeness, and also taking into account vertical representativeness,
- 15 (B) common evaluation: as (A) but without taking into account vertical sensitivity,
- 16 (C) naïve evaluation: no consideration of potential representativeness errors whatsoever,

17 For evaluation (C), the model monthly average was based on samples from all days of the 18 month (on OMI overpass time), irrespective of cloud coverage, and no kernel was applied (in 19 other words a 31-day, all-sky, without AK monthly mean). We first evaluate the (avoidable) 20 representativeness errors by comparing local OMI:GEOS-Chem ratios evaluated with approaches (A) vs. (C), and approaches (A) vs. (B). Figure 10 shows the relative difference in 21 22 the local OMI:GEOS-Chem ratios for February and August 2006. We see that the systematic, 23 avoidable errors in the OMI:GEOS-Chem ratio are largest with evaluation approach (C). The 24 blue colours in the upper panel of Figure 10(a) indicate that, in winter, sampling the model on 25 all (including cloudy sky) days leads to too low (by 15-20%) OMI/GEOS-Chem ratios reflecting the too high GEOS-Chem NO₂ values resulting from temporal representativeness 26 27 errors (cloudy-sky sampling cf. Figure 4).

The similarity between the panels of Figure 10(b) shows that appropriate sampling is not as important in Summer, a season with ample clear-sky days, and, consequently, a smaller sampling error. Figure 10(b) suggests that application of the averaging kernel when sampling the model is the most important step, with the red colours indicating that failure to apply the averaging kernel leads to OMI/GEOS-Chem NO₂ ratios that are too high by up to 30%. We 1 conclude that appropriate clear-sky sampling is mainly important in winter, but vertical 2 smoothing is less relevant in that season. The reverse holds in Summer: with sufficient clear-3 sky days available, application of the averaging kernel becomes essential, reflecting the fact 4 NO₂ vertical distributions are especially different between (the TM4 and GEOS-Chem) 5 models in that season.



8 Figure 10. Relative difference between local monthly mean OMI:GEOS-Chem NO₂ column

10 (B) and (A) (lower panel), and (b) August 2006. Relative difference defined as $100\% \times$

ratio's for (a) February 2006 between method (C) and (A) (upper panel) and between method

11 ((C/A)-1), and $100\% \times$ ((B/A)-1).

9

Table 1 summarizes the results of the OMI/GEOS-Chem comparisons for the three specific 1 regions of the United States, Europe, and China following the different evaluation 2 approaches. In all cases, the spatial correlation between model and measurements within the 3 regions is highest for evaluation approach (A), and generally lowest for approach (C). 4 5 Wintertime OMI:GEOS-Chem ratios are too low by 15-20% with approach (C) and too high by 5-10% in Summer. Using the common approach (B), OMI/GEOS-Chem ratios are 6 7 primarily biased in Summer, by +15-20% for Europe and the United States, and by -5% for 8 China. The results in Table 1 and Figure 9 also indicate that the spread of local OMI/GEOS-9 Chem ratios is $\pm 30\%$ for approach (A), smaller than for approaches (B) and (C) with spreads 10 of $\pm 35\%$, corroborating the fact that using the kernel results in a better-defined comparison 11 between satellite measurements and model simulations.

12 We summarize the contribution of the model sampling errors to the overall representativeness 13 errors for the evaluation of GEOS-Chem simulations with OMI NO₂ in Table 2. The table 14 should not be interpreted as a general recommendation for all applications, but rather as a 15 recommendation for air pollution applications such as model evaluation and inversions to estimate emissions. For instance, for data assimilation and studies of the higher atmosphere, 16 17 retrievals under cloudy situations can still be used, and the main recommendation there is to 18 apply the averaging kernel. The table shows that naïve comparison strategies (C) that do not 19 account for appropriate temporal or vertical sampling will result in a largely systematic 20 representativeness error of up to 25%. Following the motivated recommendations discussed 21 above however (i.e. comparison strategy (A)) would eliminate temporal and vertical 22 representativeness errors and limit the overall comparison error to not more than 5-10% from 23 imperfect horizontal sampling.

24

25 7 Discussion and conclusions

Evaluations of chemistry-transport model simulations with UV/Vis satellite retrievals of short-lived gases, notably NO_2 and HCHO, are strongly influenced by the exact comparison strategy. The characteristics of these satellite retrievals –with ground pixels typically smaller than model grid cells, clear-sky sampling needed for air pollution applications, and reduced vertical sensitivity towards the lower troposphere- require that models and retrievals are sampled as consistently as possible. This pertains to consistent sampling in space (horizontally and vertically) and in time (day-of-week, clear-sky day, time-of-day). Of these

aspects, appropriate horizontal sampling is a relatively minor, but unavoidable concern. In 1 most model-to-satellite comparisons, we recommend using the concept of the 2 3 superobservation, which has the distinct advantage of providing a grid cell average observed 4 state along with a realistic measurement plus horizontal representativeness error. Depending 5 on the model resolution and the satellite instrument resolution, users can impose a minimum fractional coverage (of the model grid cell area) by the ensemble of pixels to reduce 6 7 horizontal representativeness errors down to levels where the measurement contribution 8 becomes the dominant term in the superobservation error budget.

Recommendations on and error estimates of the fractional coverage requirement depend on 9 the exact method of comparing model simulations and satellite retrievals and on the spatial 10 11 variability of the species of interest. Generally speaking, fractional coverage requirements may be rather loose for comparisons over regions with little spatial variability in gas 12 13 concentration, for coarse-resolution model simulations, and for temporal averages over 14 multiple days (e.g. monthly means). In contrast, total fractional coverage requirements need to 15 be strict for comparisons over regions with strong variability in gas concentrations (i.e. SO₂) and NO₂ source regions), on a high spatial resolution with (regional) CTMs. 16

17 In these situations we recommend limiting horizontal representativeness errors to within 18 $\pm 10\%$ because representativeness errors are then still considerably smaller than the satellite 19 observation error $\sigma_{\overline{N},o}$.

20 A faithful comparison between satellite measurements and model simulations requires that 21 models need to be sampled appropriately in time. Sampling models irrespective of 22 photochemical regime (such as when calculating a 31-day monthly mean without collocating 23 the model with individual measurements) gives rise to systematic temporal representativeness 24 errors on the order of +12% for NO₂ and -8% for HCHO. Such errors should (and can) be 25 avoided, as they may misdirect interpretation of model-satellite differences, for instance by 26 misinforming inversion studies by requiring changes in the rates of emissions, or chemical 27 reactions to better match the observations. Our comparison of OMI O₂-O₂ and co-sampled 28 TM5 cloud information indicated that a strict requirement on the TM5 model to simulate a 29 clear-sky scene along with a mostly clear-sky OMI superobservation has little effect over 30 omitting such a filter. In the case of TM5, driven by ECMWF ERA Interim meteorological fields, the model shows good correlation with OMI-observed cloud fractions, with little 31 32 probability (<15%) of simulating false positives or negatives.

Larger systematic errors in model-satellite ratios will be introduced when model profiles are 1 2 not sampled according to the averaging kernel associated with most UV/Vis satellite products. 3 While the exact magnitude effect depends on the model under evaluation and on the a priori 4 profiles and other assumptions used in the retrievals, our analysis showed that for a comparison between OMI and GEOS-Chem NO₂, application of the averaging kernels results 5 in up to 20% lower satellite-to-model ratios, and more coherent values of these ratios within 6 7 relevant regions such as the eastern United States and Europe. The effect of applying the 8 kernel is most relevant in Summer, when the vertical distribution of species like NO₂ and 9 HCHO is variable, and differences between the model profiles and the profiles used in the 10 retrieval are most prominent. We strongly recommend using averaging kernels in satellite-11 model evaluations. Use of the averaging kernel allows for a better satellite-to-model 12 comparison, by ensuring that the model is sampled in a manner consistent with the satellite 13 retrievals because identical assumptions are made on vertical sensitivity, and differences 14 between the model and satellite a priori vertical distribution cancel. Here we focused on an 15 evaluation of tropospheric NO₂ simulations from the GEOS-Chem model with retrievals of 16 tropospheric NO₂ columns with substantial vertical sensitivity down to the lower troposphere. However, application of the averaging kernel will be even more relevant for model 17 evaluations of HCHO and SO₂, since these retrievals are less sensitive to the lower 18 19 troposphere. Recently, retrieval scientists have also made averaging kernel information 20 available along with the HCHO and SO₂ data products (e.g. González Abad et al. [2014]; 21 Theys et al. [2013]).

For future evaluations of chemistry transport models and data assimilation with UV/Vis satellite retrievals (of NO₂, HCHO, CHO-CHO, or SO₂), we advocate the use of the recommendations laid out in this paper, especially with respect to the required clear-sky sampling and appropriate vertical smoothing.

26

1 Appendix A

2 Calculating horizontal representativeness errors

3 The horizontal representativeness error of an ensemble of satellite measurements
4 (superobservation) for a model grid cell of any size can be calculated as follows:

5 (1) First compute the distance from one corner coordinate to the two adjacent (not 6 opposite) corners to obtain estimates for the 'base' and the 'height' of the pixel³. Then 7 approximate the pixel as a parallelogram, to calculate the pixel area A_i as base × 8 height.

9 (2) Calculate the fractional coverage f_{cov} of all valid satellite pixels in the model grid cell 10 as the ratio of the area covered by all *n* valid pixels to the complete area covered by 11 the grid cell A_{cell} :

12

$$f_{cov} = \frac{\sum_{i=1}^{n} A_i}{A_{cell}}$$
(A1)

(3) Given the fractional coverage f_{cov}, the horizontal representativeness error can be read
off from Figure 1 for models with 3° × 2° and 0.5° × 0.5° resolution. Figure 1(b) of
Miyazaki et al. [2012] provides a similar figure for a model resolution of 2.5° × 2.5°.
For example, a 0.6 fractional coverage for a 3° × 2° model grid cell corresponds to a
horizontal representativeness error of ~10%. 0.6 coverage for a 0.5° × 0.5° model
corresponds to a representativeness error of ~15%.

Note that the recipe laid out above provides a horizontal representativeness error that is at the high end of the possible range. The variability in the complete ensemble of pixels will often be much smaller than the variability in the ensemble of pixels from Figure 1 (over New York City) or Figure 1(b) from Miyazaki et al. [2012] (which excluded situations with small NO₂ columns).

24 Appendix B

³ A simple distance calculation between two latitude, longitude pairs (lat1, lon1) and (lat2, lon2) is provided by the following Fortran90 pseudo-code:

real, intent(in) :: lat1, lon1 ! coordinates of pixel 1
real, intent(in) :: lat2, lon2 ! coordinates of pixel 2
real, parameter :: dtkm = 111.32 ! at equator 1deg equals 111.32 km

deg_to_rad = acos(-1.)/180.
angle1 = 0.5 * (lat1-lat2)*deg_to_rad
angle2 = 0.5 * (lon1-lon2)*deg_to_rad
arg = (sin(angle1))2 + cos(lat1*deg_to_rad)*cos(lat2*deg_to_rad) + (sin(angle2))2
y = dtkm * 2. * asin(sqrt(arg))*180./acos(-1.)

1 Derivation of the superobservation error

If the retrieval errors within a superobservation grid cell have some degree of correlation, we cannot simply take the area-weighted average retrieval error σ (calculated as $\frac{\sum_{i=1}^{n} w_i \sigma_i}{\sum_{i=1}^{n} w_i}$) as representative for the superobservation error. The error expectation value for the ensemble of pixels composing a superobservation is written as:

6

$$\langle \varepsilon_N^2 \rangle = \sum_{ij} w_i w_j \langle \epsilon_i \epsilon_j \rangle \tag{B1}$$

7 with ϵ_i the individual retrieval error in pixel *i*, the area weights now normalized ($\sum_i w_i = 1$) 8 to facilitate notation. Now, for a partly correlated error between pixels *i* and *j*, we write:

9
$$\langle \epsilon_i \epsilon_j \rangle = \begin{cases} \sigma_i^2 \text{ for } i = j \\ c \sigma_i \sigma_j \text{ for } i \neq j \end{cases}$$
 (B2)

10 so that the superobservation error σ_N^2 can be written as follows:

11
$$\sigma_N^2 = \langle \varepsilon_N^2 \rangle = (1 - c) \sum_i w_i^2 \sigma_i^2 + c (\sum_i w_i \sigma_i)^2$$
(B3)

12 For $\sigma_i = \sigma$, and $w_i = \frac{1}{n}$ this reduces to Eq. (6): $\sigma_N = \sigma \sqrt{\frac{1-c}{n}} + c$.

13 Appendix C

14 Calculating CTM-simulated effective cloud fractions

We can express the modelled cloud properties into a quantity that is comparable to the 15 16 effective cloud fraction provided by the OMI O2-O2 cloud retrieval, and defined as the 17 radiometric equivalent fraction of a viewing scene covered by a Lambertian reflector with an 18 albedo of 0.8 (corresponding to a cloud with an optical thickness of ~40) [Stammes et al., 19 2008]. Some data products use cloud information retrieved with different approaches, but 20 many UV/Vis trace gas retrievals use the effective cloud fraction approach. The TM5 cloud 21 information (geometric cloud cover, and cloud optical thickness) was converted into an 22 effective cloud fraction in a two-step approach. In the first step the maximum-random overlap 23 assumption is used to calculate the one column-representative geometrical cloud cover $f_{tm5,geo}$ following practical guidelines for similar model evaluations with MODIS clouds by 24 25 Quaas [2011]. The maximum-random overlap assumption implies maximum overlap for 26 cloud cover in adjacent layers (one cloud layer is exactly on top of the other), and random overlap for (layers of) cloud cover f_l separated by at least one clear-sky layer: 27

$$f_{tm5,geo} = \prod_{l=1}^{L} \frac{1 - \max(f_{l,l-1})}{1 - \min(f_{l-1}, 1 - \epsilon_f)}$$
(C1)

where ϵ_f (here 0.001) is the threshold value for which a layer is considered to be cloud-free. In the second step the albedo of the cloud is determined based on the cloud optical thickness and the sensitivity of cloud spherical albedo to cloud optical thickness modelled by Buriez et al. [2005] for a liquid water cloud⁴. The final step to obtain the effective (OMI equivalent) TM5 cloud fraction $f_{tm5,eff}$ from the geometrical cloud fraction and the obtained cloud albedo a_c proceeds as:

8

25

1

$$f_{tm5,eff} = f_{tm5,geo} \frac{a_c}{0.8} \tag{C2}$$

9 Appendix D

10 Alternatives to the application of the averaging kernel

In a satellite-model comparison, the vertical sensitivity needs to be taken into account, and this can be done in alternatively by replacing the a priori profile \mathbf{x}_{a} from the CTM used in the retrieval, by the profile \mathbf{x}_{m} from the CTM used by the modeller, i.e. by re-calculating the air mass factors as follows:

15
$$M'(\mathbf{x}_{\mathbf{m}}) = M(\mathbf{x}_{\mathbf{a}}) \frac{\sum_{l=1}^{L} A_l x_{m,l}}{\sum_{l=1}^{L} x_{m,l}}$$
(D1)

with $M(\mathbf{x_a})$ the original tropospheric air mass factor used in the retrieval, and A_l the elements of the averaging kernel. The new air mass factors $(M'(\mathbf{x_m}))$ need to be applied on the retrieved slant column densities (instead of $M(\mathbf{x_a})$), to generate modified columns \hat{y}'_o . These modified columns can be directly compared to the model column \hat{x}_m , without the need to explicitly apply the averaging kernel. Such approaches have been shown to improve the consistency of the comparison considerably -- for instance by 10-20% in the case of tropospheric NO₂, see Lamsal et al. [2010], Vinken et al. [2014] and Lamsal et al. [2014].

The modified averaging kernels associated with \hat{y}'_o retrieved with the new a priori profiles \mathbf{x}_m become:

 $\mathbf{A}' = \frac{M(\mathbf{x}_a)}{M'(\mathbf{x}_m)} \mathbf{A} \tag{D2}$

⁴ A 6th order polynomial fitted in close approximation to the relationship between cloud albedo a_c and cloud optical thickness τ_c in Figure 2 of Buriez et al. [2005] was used for the conversion: $a_c = \sum_{i=0}^{6} c_i \tau^i$ with $c_0=0.00808$, $c_1=0.11153$, $c_2=-0.09734$, $c_3=0.00052$, $c_4=-0.0000154$, $c_5=0.00000029$, and $c_6=-0.0000000013$.

1 Acknowledgements

K. Folkert Boersma acknowledges funding from Vidi Project 'Attributing the sources of
tropospheric ozone from space', NWO Grant 864.09.001 and from the FP7 Project Quality
Assurance for Essential Climate Variables (QA4ECV), n° 607405. We appreciate the
constructive suggestions from reviewers that helped improve this paper.

6

7 Author contributions

K. F. B. performed the research, drafted the manuscript, prepared the figures and developed
the analysis methods. G. C. M. V. contributed to the development of the averaging kernel
code. H. J. E. aided in drafting the manuscript and methods, and supported the development
of the analysis methods and interpretation. All authors contributed to discussions of the results
and preparation of the manuscript.

1 References

- 2 Acarreta, J. R., J. F. De Haan, J. F., and Stammes, P.: Cloud pressure retrieval using the O₂-
- 3 O₂ absorption band at 477 nm, J. Geophys. Res., 109, D05204, doi:10.1029/2003JD003915,
- 4 2004.
- 5 Barkley, M. P., De Smedt, I., Van Roozendael, M., Kurosu, T. P., Chance, K., Arneth, A.,
- 6 Hagberg, D., Guenther, A., Paulot, F., Marais, E., Mao, J.: Top-down isoprene emissions over
- 7 tropical South America inferred from SCIAMACHY and OMI formaldehyde columns, J.
- 8 Geophys. Res., 118, 6849–6868, doi:10.1002/jgrd.50552, 2013.
- 9 Beirle, S., Platt, U., Wenig, M., and Wagner, T.: Weekly cycle of NO₂ by GOME
- measurements: a signature of anthropogenic sources, Atmos. Chem. Phys., 3, 2225-2232,
 doi:10.5194/acp-3-2225-2003, 2003.
- Beirle, S., Salzmann, M., Lawrence, M. G., and Wagner, T.: Sensitivity of satellite
 observations for freshly produced lightning NO_x, Atmos. Chem. Phys., 9, 1077-1094,
 doi:10.5194/acp-9-1077-2009, 2009.
- Beirle, S., Boersma, K.F., Platt, U., Lawrence, M.G., and Wagner, T.: Megacity emissions
 and lifetimes of nitrogen oxides probed from space. Science, 333(6050), 1737-1739, doi:
 10.1126/science.1207824, 2011.
- Boersma, K. F., Eskes, H. J., and Brinksma, E. J., Error analysis for tropospheric NO₂
 retrieval from space, *J. Geophys. Res.*, 109, D04311, doi:10.1029/2003JD003962, 2004.
- Boersma, K. F., Eskes, H. J., Meijer, E. W., and Kelder, H. M.: Estimates of lightning NO_x
 production from GOME satellite observations, Atmos. Chem. Phys., 5, 2311-2331,
 doi:10.5194/acp-5-2311-2005, 2005.
- Boersma, K. F., Jacob, D. J., Eskes, H. J., Pinder, R. W., Wang, J., and van der A, R. J.:
 Intercomparison of SCIAMACHY and OMI tropospheric NO₂ columns: Observing the
 diurnal evolution of chemistry and emissions from space, J. Geophys. Res., 113, D16S26,
 doi:10.1029/2007JD008816, 2008.
- Boersma, K. F., Jacob, D. J., Trainic, M., Rudich, Y., De Smedt, I., Dirksen, R., and
 Eskes, H. J.: Validation of urban NO₂ concentrations and their diurnal and seasonal variations
- 29 observed from the SCIAMACHY and OMI sensors using in situ surface measurements in
- 30 Israeli cities, Atmos. Chem. Phys., 9, 3867-3879, doi:10.5194/acp-9-3867-2009, 2009.

- Boersma, K. F., Eskes, H. J., Dirksen, R. J., van der A, R. J., Veefkind, J. P., Stammes, P.,
 Huijnen, V., Kleipool, Q. L., Sneep, M., Claas, J., Leitão, J., Richter, A., Zhou, Y., and
 Brunner, D.: An improved tropospheric NO₂ column retrieval algorithm for the Ozone
 Monitoring Instrument, Atmos. Meas. Tech., 4, 1905-1928, doi:10.5194/amt-4-1905-2011,
 2011.
- Buriez, J. C., F. Parol, C. Cornet, M. Doutriaux-Boucher, An improved derivation of the topof-atmosphere albedo from POLDER/ADEOS-2: Narrowband albedos, J. Geophys. Res., 110,
 doi:10:1029/2004JD005243, 2005.
- 9 Cerquiera, M. A., Pio, C. A., Gomes, P. A., Matos, J. S., and Nunes, T. V.: Volatile organic
- compounds in rural atmospheres in central Portugal, The Science of the Total Environment,
 313, 49-60, doi:10.1016/S0048-9697(03)00250-X, 2003.
- Curier, R. L., Kranenburg, R., Segers, A. J. S., Timmermans, R. M. A., and Schaap, M.:
 Synergistic use of OMI NO₂ tropospheric columns and LOTOS-EUROS to evaluate the NO_x
 emissions trends across Europe, Remote Sensing of the Environment, 149(2014), 58-69,
 dx.doi.org/10.1016/j.rse.2014.03.032, 2014.
- De Smedt, I., Van Roozendael, M., Stavrakou, T., Müller, J.-F., Lerot, C., Theys, N., Valks,
 P., Hao, N., and van der A, R.: Improved retrieval of global tropospheric formaldehyde
 columns from GOME-2/ MetOp-A addressing noise reduction and instrumental degradation
 issues, Atmos. Meas. Tech., 5, 2933-2949, doi:10.5194/amt-5-2933-2012, 2012.
- Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., et al.: The ERA-Interim
 reanalysis: Configuration and performance of the data assimilation system, Q.J. Roy.
 Meteorol. Soc., 137, 553-597. doi: 10.1002/qj.828, 2011.
- Dirksen, R. J., Boersma, K. F., Eskes, H. J., Ionov, D. V., Bucsela, E. J., Levelt, P. F., and
 Kelder, H. M.: Evaluation of stratospheric NO₂ retrieved from the Ozone Monitoring
 Instrument: Intercomparison, diurnal cycle, and trending, J. Geophys. Res., 116, D08305,
 doi:10.1029/2010JD014943, 2011.
- 27 Dentener, F., Kinne, S., Bond, T., Boucher, O., Cofala, J., Generoso, S., Ginoux, P., Gong, S.,
- 28 Hoelzemann, J. J., Ito, A., Marelli, L., Penner, J. E., Putaud, J.-P., Textor, C., Schulz, M., van
- 29 der Werf, G. R., and Wilson, J.: Emissions of primary aerosol and precursor gases in the years
- 30 2000 and 1750 prescribed data-sets for AeroCom, Atmos. Chem. Phys., 6, 4321-4344,
- 31 doi:10.5194/acp-6-4321-2006, 2006.

- 1 Dufour, G., Wittrock, F., Camredon, M., Beekmann, M., Richter, A., Aumont, B., and
- 2 Burrows, J. P.: SCIAMACHY formaldehyde observations: constraint for isoprene emission
- 3 estimates over Europe?, Atmos. Chem. Phys., 9, 1647-1664, doi:10.5194/acp-9-1647-2009,
- 4 2009.
- 5 Duncan, B. N., Yoshida, Y., Olson, J. R., Sillman, S., Martin, R. V., Lamsal, L., Hu, Y.
- 6 Pickering, K. E., Retscher, C., Allen, D. J., Crawford, J. H.: Application of OMI observations
- 7 to a space-based indicator of NO_x and VOC controls on surface ozone formation:
- 8 Atmospheric Environment, 44(18), 2213-2223, doi:10.1016/j.atmosenv.2010.03.010, 2010.
- 9 Duncan, B. N., Prados, A. I., Lamsal, L. N., Liu, Y., Streets, D. G., Gupta, P., Hilsenrath, E.,
- 10 Kahn, R. A., Nielsen, J. E., Beyersdorf, A. J., Burton, S. P., Fiore, A. M., Fishman, J., Henze,
- 11 D. K., Hostetlere, C. A., Krotkov, N. A., Lee, P., Lin, M., Pawson, S. Pfister, G., Pickering,
- 12 K. E., Pierce, R. B., Yoshida, Y., and Ziemba, L.D.: Satellite data of atmospheric pollution for
- 13 U.S. air quality applications: Examples of applications, summary of data end-use resources,
- 14 answers to FAQs, and common mistakes to avoid, Atmospheric Environment, 94, 647-662,
- 15 doi:10.1016/j.atmosenv.2014.05.061, 2014.
- 16 Eskes, H. J., van Velthoven, P. F. J., Valks, P. J. M., and Kelder, H. M.: Assimilation of
- 17 GOME total-ozone satellite observations in a three-dimensional tracer-transport model, Q. J.
- 18 Roy. Meteorol. Soc., 129, 1663-1681, doi:10.1256/qj.02.14, 2003.
- $19 \qquad \text{Evans, M. J., and Jacob, D. J.: Impact of new laboratory studies of N_2O_5 hydrolysis on global}$
- 20 model budgets of tropospheric nitrogen oxides, ozone, and OH, Geophys. Res. Lett., 32,
- 21 L09813, doi:10.1029/2005GL022469, 2005.
- Fried, A., McKeen, S., Sewell, S., Harder, J., Henry, B., Goldan, P., Kuster, W., Williams, E.,
 Baumann, K., Shetter, R., and Cantrell, C.: Photochemistry of formaldehyde during the 1993
- 24 Tropospheric OH Photochemistry Experiment, J. Geophys. Res., 102(D5), 6283-6296, 1997.
- Gery, M. W., Whitten, G. Z., Killus, J. P., & Dodge, M. C.: A photochemical kinetics
 mechanism for urban and regional scale computer modeling. J. Geophys. Res., 94(D10),
 12925-12956, 1989.
- 28 González Abad, G., Liu, X., Chance, K., Wang, H., Kurosu, T. P., and Suleiman, R.: Updated
- 29 SAO OMI formaldehyde retrieval, Atmos. Meas. Tech. Discuss., 7, 1-31, doi:10.5194/amtd-
- 30 7-1-2014, 2014.

- Herron-Thorpe, Lamb, B. K., Mount, G. H., Vaughan, J. K.: Evaluation of a regional air
 quality forecast model for tropospheric NO2 columns using the OMI/Aura satellite
 tropospheric NO2 product, Atmos. Chem. Phys., 10, 8839–8854, doi:10.5194/acp-10-8839 2010, 2010.
- Holmes, C. D., Jacob, D. J., Corbitt, E. S., Mao, J., Yang, X., Talbot, R., and Slemr, F.:
 Global atmospheric model for mercury including oxidation by bromine atoms, Atmos. Chem.
- 7 Phys., 10, 12037-12057, doi:10.5194/acp-10-12037-2010, 2010.
- 8 Hooghiemstra, P. B., M. C. Krol, P. Bergamaschi, A. T. J. deLaat, G. R. van derWerf, P. C.
- 9 Novelli, M. N. Deeter, I. Aben, and T. Röckmann, Comparing optimized CO emission
- 10 estimates using MOPITT or NOAA surface network observations, J. Geophys. Res., 117,
- 11 D06309, doi:10.1029/2011JD017043, 2012.
- 12 Huijnen, V., Eskes, H. J., Poupkou, A., Elbern, H., Boersma, K. F., Foret, G., Sofiev, M.,
- 13 Valdebenito, A., Flemming, J., Stein, O., Gross, A., Robertson, L., D'Isidoro, M.,
- 14 Kioutsioukis, I., Friese, E., Amstrup, B., Bergstrom, R., Strunk, A., Vira, J., Zyryanov, D.,
- 15 Maurizi, A., Melas, D., Peuch, V.-H., and Zerefos, C.: Comparison of OMI NO₂ tropospheric
- 16 columns with an ensemble of global and European regional air quality models, Atmos. Chem.
- 17 Phys., 10, 3273-3296, doi:10.5194/acp-10-3273-2010, 2010^a.
- Huijnen, V., Williams, J. E., van Weele, M., van Noije, T. P. C., Krol, M. C., Dentener, F.,
 Segers, A., Houweling, S., Peters, W., de Laat, A. T. J., Boersma, K. F., Bergamaschi, P.,
 van Velthoven, P. F. J., Le Sager, P., Eskes, H. J., Alkemade, F., Scheele, M. P., Nédélec, P.,
 and Pätz, H.-W.: The global chemistry transport model TM5: description and evaluation of
 the tropospheric chemistry version 3.0, Geosci. Model Dev. Discuss., 3, 1009-1087,
 doi:10.5194/gmdd-3-1009-2010, 2010^b.
- Ide, K., Courtier, P., Ghil, M., Lorenc, A. C., Unified Notation for Data Assimilation:
 Operational, Sequential, and Variational, J. Met. Soc. Japan, Special Issue on "Data
 Assimilation in Meteorology and Oceanography: Theory and Practice, Vol. 75, No. 1B, 181189, 1997.
- Inness, A., Baier, F., Benedetti, A., Bouarar, I., Chabrillat, S., Clark, H., et al. (2013). The
 MACC reanalysis: An 8 yr data set of atmospheric composition. *Atmospheric Chemistry and Physica* 12, 4072, 4100, doi:10.5104/com 12, 4072, 2012, 2012.
- 30 *Physics, 13*, 4073-4109. doi:10.5194/acp-13-4073-2013, 2013.

- Irie, H., Boersma, K. F., Kanaya, Y., Takashima, H., Pan, X., and Wang, Z. F.: Quantitative
 bias estimates for tropospheric NO₂ columns retrieved from SCIAMACHY, OMI, and
 GOME-2 using a common standard for East Asia, Atmos. Meas. Tech., 5, 2403-2411,
 doi:10.5194/amt-5-2403-2012, 2012.
- Jones, D. B. A., Bowman, K. W., Palmer, P. I., Worden, J. D., Jacob, D. J., Hoffman, R. N.,
 Bey, I., and Yantosca, R. M., Potential of observations from the Tropospheric Emission
 Spectrometer to constrain continental sources of carbon monoxide, J. Geophys. Res.,
 108(D24), 4789, doi:10.1029/2003JD003702, 2003.
- 9 Kleipool, Q. L., Dobber, M. R., de Haan, J., & Levelt, P. F., Earth surface reflectance 10 climatology from 3 years of OMI data. *J. Geophys. Res.*, *113*(D18), D18308, 2008.
- Krol, M., Houweling, S., Bregman, B., van den Broek, M., Segers, A., van Velthoven, P.,
 Peters, W., Dentener, F., and Bergamaschi, P.: The two-way nested global chemistrytransport zoom model TM5: algorithm and applications, Atmos. Chem. Phys., 5, 417-432,
 doi:10.5194/acp-5-417-2005, 2005.
- Krol, M., Peters, W., Hooghiemstra, P., George, M., Clerbaux, C., Hurtmans, D.,
 McInerney, D., Sedano, F., Bergamaschi, P., El Hajj, M., Kaiser, J. W., Fisher, D.,
 Yershov, V., and Muller, J.-P.: How much CO was emitted by the 2010 fires around
 Moscow?, Atmos. Chem. Phys., 13, 4737-4747, doi:10.5194/acp-13-4737-2013, 2013.
- Lamsal, L. N., R. V. Martin, R. V., van Donkelaar, A., Celarier, E. A., Bucsela, E. J.,
 Boersma, K. F., Dirksen, R., Luo, C., and Wang, Y.: Indirect validation of tropospheric
 nitrogen dioxide retrieved from the OMI satellite instrument: Insight into the seasonal
 variation of nitrogen oxides at northern midlatitudes, J. Geophys. Res., 115, D05302,
 doi:10.1029/2009JD013351, 2010.
- Lamsal, L. N., Krotkov, N. A., Celarier, E. A., Swartz, W. H., Pickering, K. E., Bucsela, E. J.,
 Gleason, J. F., Martin, R. V., Philip, S., Irie, H., Cede, A., Herman, J., Weinheimer, A.,
 Szykman, J. J., and Knepp, T. N.: Evaluation of OMI operational standard NO₂ column
 retrievals using in situ and surface-based NO₂ observations, Atmos. Chem. Phys., 14, 11587-
- 28 11609, doi:10.5194/acp-14-11587-2014, 2014.
- 29 Lathière, J., Hauglustaine, D. A., Friend, A. D., De Noblet-Ducoudré, N., Viovy, N., and
- 30 Folberth, G. A.: Impact of climate variability and land use changes on global biogenic volatile

- organic compound emissions, Atmos. Chem. Phys., 6, 2129-2146, doi:10.5194/acp-6-2129 2006, 2006.
- Lee, C., Martin, R. V., van Donkelaar, A., O'Byrne, G., Krotkov, N., Richter, A., Huey, L. G.,
 and Holloway, J. S., Retrieval of vertical columns of sulfur dioxide from SCIAMACHY and
 OMI: Air mass factor algorithm development, validation, and error analysis, J. Geophys. Res.,
- 6 114, D22303, doi:10.1029/2009JD012123, 2009.
- Lin, J.-T., and McElroy, M. B.: Detection from space of a reduction in anthropogenic
 emissions of nitrogen oxides during the Chinese economic downturn, Atmos. Chem. Phys.,
- 9 11, 8171-8188, doi:10.5194/acp-11-8171-2011, 2011.
- 10 Lin, J.-T., Liu, Z., Zhang, Q., Liu, H., Mao, J., and Zhuang, G.: Modeling uncertainties for
- 11 tropospheric nitrogen dioxide columns affecting satellite-based inverse modeling of nitrogen
- 12 oxides emissions, Atmos. Chem. Phys., 12, 12255-12275, doi:10.5194/acp-12-12255-2012,
 13 2012^a.
- Lin, J.-T.: Satellite constraint for emissions of nitrogen oxides from anthropogenic, lightning
 and soil sources over East China on a high-resolution grid, Atmos. Chem. Phys., 12, 28812898, doi:10.5194/acp-12-2881-2012, 2012^b.
- Ma, J. Z., Beirle, S., Jin, J. L., Shaiganfar, R., Yan, P., and Wagner, T.: Tropospheric NO₂
 vertical column densities over Beijing: results of the first three years of ground-based MAXDOAS measurements (2008–2011) and satellite validation, Atmos. Chem. Phys., 13, 15471567, doi:10.5194/acp-13-1547-2013, 2013.
- Maasakkers, J. D.: Vital improvements to the retrieval of tropospheric NO₂ columns from the
 Ozone Monitoring Instrument, M.Sc Report, R-1882-A, July 2013, Eindhoven University of
 Technology, 2013.
- 24 Mao, J., Jacob, D. J., Evans, M. J., Olson, J. R., Ren, X., Brune, W. H., Clair, J. M. St., 25 Spencer, K. M., Beaver, M. R., Crounse, J. D., Wennberg, P. O., Cubison, M. J., 26 Jimenez, J. L., Fried, A., Weibring, P., Walega, J. G., Hall, S. R., Weinheimer, A. J., Cohen, R. C., Chen, G., Crawford, J. H., McNaughton, C., Clarke, A. D., Jaeglé, L., 27 28 Fisher, J. A., Yantosca, R. M., Le Sager, P., and Carouge, C.: Chemistry of hydrogen oxide 29 radicals (HO_x) in the Arctic troposphere in spring, Atmos. Chem. Phys., 10, 5823-5838,
- 30 doi:10.5194/acp-10-5823-2010, 2010.

- Martin, R. V., Fiore, A. M., and van Donkelaar, A.: Space-based diagnosis of surface ozone
 sensitivity to anthropogenic emissions, Geophys. Res. Lett., 31, L06120,
 doi:10.1029/2004GL019416, 2004.
- McLinden, C. A., Fioletov, V., Boersma, K. F., Kharol, S. K., Krotkov, N., Lamsal, L.,
 Makar, P. A., Martin, R. V., Veefkind, J. P., and Yang, K.: Improved satellite retrievals of
- 6 NO₂ and SO₂ over the Canadian oil sands and comparisons with surface measurements,
- 7 Atmos. Chem. Phys., 14, 3637-3656, doi:10.5194/acp-14-3637-2014, 2014.
- 8 Migliorini, S., Piccolo, C., and Rodgers, C. D.: Use of the Information Content in Satellite
- 9 Measurements for an Efficient Interface to Data Assimilation, Mon. Wea. Rev., 136, 2633-
- 10 2650, doi:10.1175/2007MWR2236.1, 2008.
- 11 Millet, D. B., Jacob, D. J., Boersma, K. F., Fu, T. M., Kurosu, T. P., Chance, K., Heald, C. L.,
- 12 and Guenther, A.: Spatial distribution of isoprene emissions from North America derived
- 13 from formaldehyde column measurements by the OMI satellite sensor, J. Geophys. Res., 113,
- 14 D02307, doi:10.1029/2007JD008950, 2008.
- 15 Miyazaki, K., H. J. Eskes, K. Sudo, M. Takigawa, M. van Weele, and K. F. Boersma,
- 16 Simultaneous assimilation of satellite NO2, O3, CO, and HNO3 data for the analysis of
- 17 satellite NO2, O3, CO, and HNO3 data for the analysis of tropospheric chemical composition
- 18 and emissions, Atmos. Chem. Phys., 12, 9545-9579, doi:10.5194/acp-12-9545-2012, 2012.
- Mijling, G., and van der A, R.J.: Using daily satellite observations to estimate emissions of
 short-lived air pollutants on a mesoscopic scale, J. Geophys. Res., 117, D17302,
 doi:10.1029/2012JD017817, 2012.
- Morcrette, J. J., & Jakob, C., The response of the ECMWF model to changes in the cloud overlap assumption. *Monthly Weather Review*, *128*(6), 1707-1732, 2000.
- Müller, J.-F., Stavrakou, T., Wallens, S., De Smedt, I., Van Roozendael, M., Potosnak, M. J.,
 Rinne, J., Munger, B., Goldstein, A., and Guenther, A. B.: Global isoprene emissions
- estimated using MEGAN, ECMWF analyses and a detailed canopy environment model,
- 27 Atmos. Chem. Phys., 8, 1329-1341, doi:10.5194/acp-8-1329-2008, 2008.
- 28 Olivier, J. G. J. and Berdowski, J. J. M.: Global emissions sources and sinks, in The Climate
- 29 System, Lisse, Netherlands, 33–78, 2001.

- Platt, U., and Stutz, J.: Differential Optical Absorption Spectroscopy: Principles and
 Applications, Springer, 2008.
- 3 Quaas, J., Final versions of CALIPSO-PARASOL observational analysis product and of
- 4 MODIS simulator, Deliverable D1.2: Final version of MODIS simulator,
- 5 http://www.euclipse.eu/downloads/D1.2_euclipse_modissimulator.pdf, 2011.
- 6 Richter, A., Wittrock, F., & Burrows, J. P., SO₂ measurements with SCIAMACHY, In Proc.
- 7 Atmospheric Science ConferenceF (pp. 8-12), 2006.
- 8 Rodgers, C. D., Inverse methods for atmospheric sounding: Theory and Practice, Series on
- 9 Atmospheric, Oceanic and Planetary Physics–Vol. 2., Singapore, World Scientific (2000).

10 Rodgers, C. D., and B. J. Connor, Intercomparison of remote sounding instruments, J.

11 Geophys. Res., 108(D3), 4116, doi:10.1029/2002JD002299, 2003.

- 12 Schaub, D., Boersma, K. F., Kaiser, J. W., Weiss, A. K., Folini, D., Eskes, H. J., and
- 13 Buchmann, B.: Comparison of GOME tropospheric NO₂ columns with NO₂ profiles deduced
- 14 from ground-based in situ measurements, Atmos. Chem. Phys., 6, 3211-3229,
 15 doi:10.5194/acp-6-3211-2006, 2006.
- Schaub, D., Brunner, D., Boersma, K. F., Keller, J., Folini, D., Buchmann, B.,
 Berresheim, H., and Staehelin, J.: SCIAMACHY tropospheric NO₂ over Switzerland:
 estimates of NO_x lifetimes and impact of the complex Alpine topography on the retrieval,
- 19 Atmos. Chem. Phys., 7, 5971-5987, doi:10.5194/acp-7-5971-2007, 2007.
- Stammes, P., Sneep, M., de Haan, J. F., Veefkind, J. P., Wang, P., and Levelt, P. F.: Effective
 cloud fractions from the Ozone Monitoring Instrument: Theoretical framework and
 validation, J. Geophys. Res., 113, D16S38, doi:10.1029/2007JD008820, 2008.
- 23 Stavrakou, T., Müller, J.-F., Boersma, K. F., van der A, R. J., Kurokawa, J., Ohara, T., and
- 24 Zhang, Q.: Key chemical NO_x sink uncertainties and how they influence top-down emissions
- of nitrogen oxides, Atmos. Chem. Phys., 13, 9057-9082, doi:10.5194/acp-13-9057-2013,
 2013.
- Theys, N., Campion, R., Clarisse, L., Brenot, H., van Gent, J., Dils, B., Corradini, S.,
 Merucci, L., Coheur, P.-F., Van Roozendael, M., Hurtmans, D., Clerbaux, C., Tait, S., and
 Europeiric F.: Valeenia SO, fluxes derived from actallite date: a survey using OML COME 2
- 29 Ferrucci, F.: Volcanic SO₂ fluxes derived from satellite data: a survey using OMI, GOME-2,

- IASI and MODIS, Atmos. Chem. Phys., 13, 5945-5968, doi:10.5194/acp-13-5945-2013,
 2013.
- 3 Uno, I., He, Y., Ohara, T., Yamaji, K., Kurokawa, J.-I.: Systematic analysis of interannual
- and seasonal variations of model-simulated tropospheric NO₂ in Asia and comparison with
 GOME-satellite data. Atmos. Chem. Phys., 7, 1671-1681, 2007.
- 6 van der Werf, G. R., Randerson, J. T., Giglio, L., Collatz, G. J., Kasibhatla, P. S., and
- 7 Arellano Jr., A. F.: Interannual variability in global biomass burning emissions from 1997 to
- 8 2004, Atmos. Chem. Phys., 6, 3423–3441, doi:10.5194/acp-6-3423-2006, 2006.
- 9 van Geffen, J. H. G. M., Boersma, K. F., Van Roozendael, M., Hendrick, F., Mahieu, E.,
- 11 from OMI in the 405–465 nm window, Atmos. Meas. Tech., 8, 1685-1699, doi:10.5194/amt-

De Smedt, I., Sneep, M., and Veefkind, J. P.: Improved spectral fitting of nitrogen dioxide

- 12 8-1685-2015, 2015.

10

- 13 Vinken, G. C. M., Boersma, K. F., Jacob, D. J., and Meijer, E. W.: Accounting for non-linear
- 14 chemistry of ship plumes in the GEOS-Chem global chemistry transport model, Atmos.
- 15 Chem. Phys., 11, 11707-11722, doi:10.5194/acp-11-11707-2011, 2011.
- Vinken, G. C. M., Boersma, K. F., van Donkelaar, A., and Zhang, L.: Constraints on ship
 NOx emissions in Europe using GEOS-Chem and OMI satellite NO₂ observations, Atmos.
- 18 Chem. Phys., 14, 1353-1369, doi:10.5194/acp-14-1353-2014, 2014.
- von Hardenberg, J., Vozella, L., Tomasi, C., Vitale, V., Lupi, A., Mazzola, M.,
 van Noije, T. P. C., Strunk, A., and Provenzale, A.: Aerosol optical depth over the Arctic: a
 comparison of ECHAM-HAM and TM5 with ground-based, satellite and reanalysis data,
 Atmos. Chem. Phys., 12, 6953-6967, doi:10.5194/acp-12-6953-2012, 2012.
- Vrekoussis, M., Wittrock, F., Richter, A., and Burrows, J. P.: Temporal and spatial variability
 of glyoxal as observed from space, Atmos. Chem. Phys., 9, 4485-4504, doi:10.5194/acp-94485-2009, 2009.
- Wang, Y., McElroy, M. B., Martin, R. V., Streets, D. G., Zhang, Q., and Fu, T.-M.: Seasonal
 variability of NO_x emissions over east China constrained by satellite observations:
 Implications for combustion and microbial sources, J. Geophys. Res., 112, D06301,
 doi:10.1029/2006JD007538, 2007.

- 1 Wang, X., Mallet, V, Berroir, J.-P., and Herlin, I.: Assimilation of OMI NO₂ retrievals into a
- 2 regional chemistry-transport model for improving air quality forecasts over Europe, Atm.
- 3 Environm., 45 (2011), 485-492, doi:10.1016/j.atmosenv.2010.09.028, 2011.
- 4 Williams, J. E., Scheele, M. P., van Velthoven, P. F. J., Cammas, J.-P., Thouret, V. Galy-
- 5 Lacaux, C., and Volz-Thomas, A.: The influence of biogenic emissions from Africa on 6 tropical tropospheric ozone during 2006: a global modeling study, Atmos. Chem. Phys., 9,
- 7 5729-5749, doi:10.5194/acp-9-5729-2009, 2009.
- Williams, J. E., Strunk, A., Huijnen, V., and van Weele, M.: The application of the Modified
 Band Approach for the calculation of on-line photodissociation rate constants in TM5:
 implications for oxidative capacity, Geosci. Model Dev., 5, 15-35, doi:10.5194/gmd-5-152012, 2012.
- Williams, J. E., Le Bras, G., Kukui, A., Ziereis, H., and Brenninkmeijer, C. A. M.: The impact of the chemical production of methyl nitrate from the NO + CH_3O_2 reaction on the global distributions of alkyl nitrates, nitrogen oxides and tropospheric ozone: a global modelling study, Atmos. Chem. Phys., 14, 2363-2382, doi:10.5194/acp-14-2363-2014, 2014.
- Zhang, L., Jacob, D. J., Knipping, E. M., Kumar, N., Munger, J. W., Carouge, C. C.,
 van Donkelaar, A., Wang, Y. X., and Chen, D.: Nitrogen deposition to the United States:
 distribution, sources, and processes, Atmos. Chem. Phys., 12, 4539-4554, doi:10.5194/acp12-4539-2012, 2012.
- Zhang, Q., Streets, D. G., and He, K.: Satellite observations of recent power plant
 construction in Inner Mongolia, China, Geophys. Res. Lett., 36, L15809,
 doi:10.1029/2009GL038984, 2009.
- 23

- 1 Table 1. Summary of tropospheric NO₂ GEOS-Chem model evaluations following recipes
- 2 (A), (B), and (C) with OMI NO₂ retrievals for February and August 2006. *n* refers to the
- 3 number of grid cells used in the comparison.

	\mathbb{R}^2			Geometric mean			
Model evaluation	(A)	(B)	(C)	(A)	(B)	(C)	п
Europe Feb 2006	0.66	0.63	0.54	$1.15^{1.88}_{0.70}$	$1.13^{1.89}_{0.67}$	$0.92^{1.59}_{0.54}$	120
Europe Aug 2006	0.66	0.57	0.59	$1.24^{1.80}_{0.86}$	$1.42_{0.91}^{2.21}$	$1.40^{2.13}_{0.92}$	137
US Feb 2006	0.82	0.79	0.75	$1.12^{1.29}_{0.97}$	$1.08^{1.34}_{0.87}$	$0.95^{1.25}_{0.72}$	41
US Aug 2006	0.83	0.61	0.67	$0.75_{0.59}^{0.95}$	$0.91^{1.21}_{0.68}$	$0.90^{1.16}_{0.70}$	42
China Feb 2006	0.58	0.57	0.54	$1.00^{1.37}_{0.72}$	$0.99^{1.63}_{0.59}$	$0.86^{1.08}_{0.69}$	35
China Aug 2006	0.61	0.58	0.58	$1.13^{1.39}_{0.92}$	$1.07^{1.36}_{0.84}$	$1.06^{1.37}_{0.82}$	44

4

5 **Table 2.** Overview of magnitude and nature of various model sampling errors, their 6 contribution to the overall comparison error budget, and ways to avoid them. Based on the 7 GEOS-Chem evaluation with OMI NO₂ retrievals for February and August 2006. Note that 8 these recommendations hold for air pollution applications of UV/Vis satellite retrievals such 9 as model evaluation and top-down emission estimates.

	Relative error	Type of error	Recommendation
Horizontal sampling	<5-10%	Inevitable and random	Require at least 40% coverage of model grid cell.
Temporal sampling	10%	Avoidable and systematic	Sample model grid cells exclusively on clear-sky days
Vertical sampling	20%	Avoidable and systematic	Apply averaging kernel on model vertical distribution
Overall representativeness error	10%-25%		Follow recommendations listed above to keep $\sigma_r < \sigma_o$

10