# Author's response to the general comments from referee III:

Thank you for revealing your valuable criticism regarding the manuscript. Below, please find our responses to your specific comments, along with the implemented changes to our manuscript. All page and line numbers as well as figure numbering refer to the *revised* manuscript. Note specifically that the figure numbering has changed during the review process.

## MAIN COMMENTS:

**1 ) Main comment from referee:**

a) The paper is hard to read. Often it requires re-reading a paragraph a number of times, to understand. It also has to do with the structure of the paper. It would help to define the main methodology of the data analysis and to have this as the main thread throughout the paper.

I think I understand the methodology, but I am still not sure.
Let me explain my interpretation of the method:
(i) Four dimensional look-up tables (LUT) are created of co-located data, so differences between data sets are stratified according to $q\_a$, U, SST, and wvpa.
(ii) The mean difference between HOAPS and insitu data are interpreted as biases.
(iii) The variances of the difference are used for the triple co-location method, resulting in error estimates.
(iv) This results in LUT's of biases and random error estimates.
(v) In applications (e.g. global maps of mean and random error of $q\_a$) the observations of $q\_a$, U, SST and wvpa point to the table and provide errors of each observation. These can be averaged to obtain the desired map.

b) I feel that it would be helpful to describe upfront that this is the general methodology and follow it throughout the paper. So this would lead to 3 main sections in the paper: (i) Description of the methodology, (ii) Results of the methodology, i.e statistics on the LUT data, and (iii) Application to HOAPS evaporation. In case I am completely wrong on the interpretation of the paper, there is even more reason to be clear about the methodology.

c) Another question is: what is the main result of the paper? If my interpretation is correct, then the 4-dimensional table of error estimates is the main result, because it would allow a user to make estimates of anything he/she is interested in (e.g. monthly averages, daily averages, or El-Nino years). So it is worth thinking about communicating this 4D table to the users. Most of the current paper is about applying the methodology, but these are in fact just examples.

**Author's reponse:**

Related to a) We have restructured the paper to be more clear about the methodology. The introduction has been rearranged and shortened (see also specific comment #1 below) for a much clearer understanding of the motivation, the benefit, and the structure of the paper. We have furthermore appended a flowchart to this document, which guides you through the individual steps of data processing, intermediate data products, and resulting HOAPS uncertainty measures.

Regarding your methodology interpretation (the flowchart assists):
(i) Correct.
(ii) Yes, differences of paired collocations are considered as biases. Depending on U, $q_a$, SST, and wvpa, these single biases are assigned to one of the $20^4$ bins. Once all collocations have been

assigned to a bin, bin-averaged systematic uncertainties are computed based on the *absolute* differences of all assigned biases. That is, we consider the upper end of bias considerations (see comment #10 on this).

(iii) Exactly. Although the variances of differences, applied to the triplets, only help to decompose the *random* uncertainty estimates to end up with HOAPS-related (that is, retrieva-related) random uncertainty estimates.

(iv) The LUTs of systematic and (uncorrected) random uncertainties already result from i). iii) helps to decompose the random uncertainty components, isolate the retrieval-related part, and finanly correct the random uncertainty LUTs.

(v) Yes, this is correct.

Related to b) Regarding your proposed three main sections: We believe that this has already been done in a similar manner and is reflected in the numbering of the Section: Methodology (Sect. 3) and Results/Applications (Sect. 4). We do not want to dedicate „Results of the methodology" an own section, as we think that it belongs to the methodological part of the manuscript (Sects. 3.3, 3.4).

However, we agree that the submitted manuscript was not structured clearly enough. This shortcoming has been improved (see comment related to a) above).

Related to c) We agree that *one* of the main outcomes of the manuscript is the benefit of the multi-dimensional bias analyses. Particularly because the approach can be easily transferred to other satellite retrievals and potentially also other remotely sensed parameters. The approach itself should be stressed more clearly in the conclusion. Communicating our specific LUTs to the users is not helpful, though, as they are tailored to HOAPS-3.3 (due to the double collocations described in Sect. 3.1). The results of applying an updated version of the LUTs to instantaneous HOAPS data are implemented in the most recent HOAPS 4.0 Version (Andersson et al. (2017)) in form of systematic and random uncertainties.

We believe that the application of the uncertainty characterization approach is equally important, as it leads to uncertainty estimates for a widely used data record. On the one hand, none of the remaining LHF-related satellite climatologies are equipped with such estimates. On the other hand, Sect. 4 demonstrates a variety of different approaches for illustrating the uncertainties and allows for identifying regions where uncertainties in the satellite retrieval are an issue and need to be accounted for.

The focus of this paper is therefore twofold: 1) describing the method and 2) applying the method to arrive at HOAPS-3.3 uncertainty estimates.

**Changes in the manuscript:**

Related to a and b) The whole manuscript has been revised for a clearer reading experience. This specifically targets Sects. 1 and 3. The last paragraph of Sect. 1 now guides the reader through the manuscript step by step. Section 3.3. and 3.4 have been swapped to be consistent with the sequence of analyses.

Related to c) The benefit of the multi-dimensional bias analyses for uncertainty characterizations has been highlighted more clearly in Sect. 5 (P.20, L.21ff).

---

**2 ) Main comment from referee:**

Estimation of biases is non-trivial. In fact this is very important because, as the authors point out, for long term averages the systematic errors dominate.

My concern is two-fold:
a) I have the feeling that it is assumed that DWD-ICOADS data is bias-free? If this is the basis for the bias estimates, then it deserves more discussion also in view of what has been published in literature.
b) Fig. 1 is used as an example to illustrate the estimation of biases. However, it is likely that artificial biases occur in binned scatter plots of noisy data if correlated variables are used on abscissa and ordinate. This applies to Fig. 1a where hair(HOAPS) is used on both vertical and horizontal axes. It also applies to hair versus wind because these variables are correlated due to the physics of the mixing (more wind brings hair closer to the surface value). To check, one could e.g. bin the differences of Fig. 1a in classes of hair(insitu). Also hair(insitu) is noisy because it has large representativeness errors (point observation, whereas HOAPS has a large footprint).

c) Finally, if one can be confident about the bias estimation, then it should also be trivial to apply a bias correction to HOAPS. This would just leave the uncertainty in C_E which is a parametrization constant used for satellite as well as in-situ data. Please discuss.

**Author's reponse:**

Related to a) It is correct that we assume the DWD-ICOADS data base to be bias free (see last paragraph in Sect. 2.2 of submitted manuscript). Our filtering procedure ensures that only high-quality in situ data is used for collocation analysis. Systematic effects of known origin are thought to have been removed or at least minimized within the quality checking procedure at the Marine Climate Data Center of DWD. Other systematic uncertainties like differing sensor heights and cool skin effects have been eliminated prior to our analysis due to sensor height corrections using in situ platform meta data (U) and cool skin corrections ($q_s$). We are aware of the fact that no ground "truth" exists, but are confident that our extensive data base is the best ground "reference" available. Freeman et al. (2016) present a great overview of the variety of ICOADS applications, which also include the calibration and validation of satellite data (e.g. Bentamy et al. (2003), Bentamy et al. (2013), Jackson et al. (2009), Jackson and Wick (2010)).
It should be kept in mind that our systematic uncertainty estimates represent the upper limit of a more simple bias estimation. Assuming a bias free ground reference therefore does not violate our conclusions, although a small contribution to the systematic uncertainties may be caused by the in situ reference.
One could argue that our uncertainty estimates in regions of poor in situ data coverage are questionable. However, as picked up in Sect. 3.2, we overcome the regional dependency by characterizing uncertainties as a function of ambient atmospheric conditions. Poor in situ data densities are therefore of secondary importance, as their ambient atmospheric conditions may be similar in regions with considerably more match ups.

Related to b) Thank you for the suggestion to investigate the one-dimensional patterns of dq as a function of the in situ source. We exemplarily performed this analysis for 2001 with approximately 1.8 million match ups. We compared the magnitudes of the mean 5-percentiles, which (in case of HOAPS) are illustrated as black squares in Fig. 2. For U, $q_s$, and $q_a$, our results indicate that in 80% of all match ups (i.e., excluding the margins), relative differences between HOAPS and in situ mean 5-percentiles range between ± 6-10%, which we consider as negligible. We presume that a two- instead of one-sided regression approach would lead to even more robust 5-percentile means. Towards the margins of the distributions, relative differences become larger. We believe that this does not have a noteworthy impact on the four-dimensional analyses, as the biases in one-dimensional space may become smaller or even cancel out when the remaining three atmospheric state parameters are considered concurrently.
Independent of this, biases as a function of in situ LHF-related parameters cannot be investigated in four-dimensional space, as vertically integrated water vapour ("wvpa"), an important indicator for

the prevalent atmospheric condition, is not available from in situ measurements. This would lead to an undesirable simplification of our uncertainty analysis approach. Additionally, our match up data base only lasts until 2008. In consequence, no uncertainties could be assigned to pixel level HOAPS data from 2008 onwards, if the multi-dimensional bias approach was based on in situ data.

Related to c) Regarding a bias correction of HOAPS data: Our approach aims at characterizing uncertainties inherent to HOAPS. This allows users to implement this information into their analyses and arrive at appropriate conclusions. We have further emphasized the benefit of such estimates in the revised version of Sect. 1. The focus is therefore not put on bias correction with respect to DWD-ICOADS. A sustainable consequence of large uncertainties should in fact point at the need of modifying the retrieval algorithm instead of bias correcting the data. It is our impression that a bias correction is feasible, if a constant bias (in terms of dependent variables, region, or time) is present relative to a fiducial reference. Such a reference is not available at present.

**Changes in the manuscript:**

regarding a) Freeman et al. (2016) is picked up in context of describing the in situ data base (P.6, L.18f). Furthermore, some further references are given regarding our assumption of bias-free ICOADS measurements (P.8, L.3ff).

regarding b) We briefly mention the artificial biases due to correlating variables and conclude that two-sided regression analyses could reduce these spurious biases (P.9, L.17ff).

---

<u>**SPECIFIC COMMENTS:**</u>

**1 ) Comment from referee:**
Section 1: Although well written, the introduction is rather long and contains sometimes fluffy language. No reference is made to an earlier study by Kinzel et al. (2016). What is new compared to earlier work? Reference is made to other data sets and to studies that provide error estimates. However, nothing is said about published error estimation methods.
**Author's reponse:** We agree that the introduction of the submitted manuscript is too long. We have restructured Sect. 1 following your suggestions and believe that this essentially improved the manuscript (see also "main comment 1" on this). Furthermore, Kinzel et al. (2016) has been included in the revised manuscript to point at the  random uncertainty decomposition approach. We agree it is important to distinguish between earlier work and new aspects of this manuscript. This also includes a statement regarding earlier error estimation methods and those present in our manuscript.
**Changes in the manuscript:** The first 24 lines have been considerably shortened.
The whole introduction has been restructured to be clearer about the motivation and benefit of our study. It now clearly differentiates between earlier approaches (that is, mostly intercomparison studies, P.3, L.25ff) and the novelty of our uncertainty characterization (that is, uncertainty estimates that are exclusively related to a specific data set, in particular HOAPS, P.4, L.1f.; P.4, L.14ff). At the same time, we highlight the new aspects of our approach (e.g. four-dimensional LUTs, P.4, L.20f). In this regard, Kinzel et al. (2016) has been included into the revised manuscript and is put into context (P.4, L.10f, L.21f). The advantage of multi-dimensional LUTs has been included into the abstract (P.1, L.7f).

**2 ) Comment from referee:** Page 5, Line 32: The sentence with "The latter depends" suggests that it refers to $q_a$ in the sentence before, but what it intends to say is that the COARE algorithm needs stability and that specific assumptions are made. Please rephrase.
**Author's reponse:** The wording was chosen on purpose, as we wanted to point out that the

saturation vapour pressure and hence $q_a$ depends on the surface air temperature. However, we agree that this may be confusing and the focus should be put on the stability calculation.

**Changes in the manuscript:** We changed the wording to "It includes atmospheric stability calculations, which necessitate surface air temperatures as input. These are estimated by assuming..." (P.6, L.7f).


**3 ) Comment from referee:** Page 6, Line 20-24: The non-correction of $q_a$ for measuring height is confusing. Why not using the real measuring height in the bulk formula? Perhaps it is possible to say in one sentence what the results are of the height difference effects as estimated by Kent et al. (2014).

**Author's reponse:** Prytherch et al. (2014) and Kinzel et al. (2016) point at the disadvantages related to $q_a$ height corrections. We agree that a statement regarding the height correction effect is useful. Kent et al. (2014) quantify the height correction effect to be 0.11 g $kg^{-1}$ for the time period 1971-2006, owing to the continuously increasing measurement platform heights. However, this effect is masked by bias corrections associated with measurement techniques, which are thought to be 2-3 times larger.

**Changes in the manuscript:** Results by Kent et al. (2014) regarding the height difference effects are briefly mentioned (P.7, L.3ff).


**4 ) Comment from referee:** Page 7, Line 13: Cool skin corrections are applied to in situ observation but not to HOAPS-3.3 SST (AVHRR based). This makes sense in priciple because AVHRR measures the skin temperature. However, there must be a calibration procedure of AVHRR, which is probably against bulk SST data. So, what does calibrated AVHRR data represent, bulk or skin SST?

**Author's reponse:** Thank you for bringing this up. Indeed, AVHRR was calibrated against bulk SST. Formally, this would necessitate a cold skin correction. However, compared to OISST (Reynolds et al. (2007)), AVHRR has a cold bias of unknown origin, which is in the order of the skin correction. We therefore refrained from performing the correction and consider the AVHRR SST as a skin SST. Note that this cold bias problem is overcome in HOAPS 4.0 (Andersson et al. (2017), which is based on OISST. For the HOAPS 4.0 retrieval (Andersson et al. (2017)), OISST is corrected for the cold skin effect.


**5 ) Comment from referee:** Pages 4-5 section 2.1: It would be informative to mention pixel size of the microwave sensors.

**Author's reponse:** Yes, we agree.

**Changes in the manuscript:** Pixel sizes have been included into Sect. 2.1 of the revised manuscript (P.5, L.11-13).


**6 ) Comment from referee:** Page 8, Line 11: The sentence "Figure 1a overestimates ...." is confusing. Formally it is correct, but, after reading the first time it suggests that the biases range from 7-12 g/kg and that the plot is for the inner tropics.

**Author's reponse:** Indeed, this may be misunderstood.

**Changes in the manuscript:** The wording has been changed in the revised manuscript to: "For $q_a$ values between 7-12g $kg^{-1}$ , HOAPS-3.3 overestimates near-surface specific humidities (see Figure 2a). Overestimations are also observed in the inner tropics, where $q_a$ is in the order of 20 g $kg^{-1}$" (P.9, L6f).


**7 ) Comment from referee:** Page 8, Line 17: The expression "over-(under-)estimated" is perhaps better than "over-(under-)represented"

**Author's reponse:** Thank you for this suggestion.

**Changes in the manuscript:** "over-(under-)represented has been replaced by "over-(under-)estimated" (P.9, L.12).

**8 ) Comment from referee:** Scatter plots in Fig.1: In all the plots except (c) the variables on the vertical axis are correlated with the variable of the horizontal axis. This is most obvious for Fig. (a) where hair-HOAPS is used in both abscissa and ordinate. In such cases the binning according to one axis can show biases that are not necessarily real. Whether this is really the case can be easily demonstrated by making the same plot but now with hair-insitu on the horizontal axis. Similarly unrealistic bias may be seen in (b) and (d) because wind and wvpa are derived with from the same satellite channels and therefore correlate with hair-HOAPS. Please discuss.

**Author's reponse:** We assume that "hair-insitu" means "$q_a$(in situ)" and not the mathematical difference between HOAPS and in situ $q_a$? We are aware of the correlation between the individual variables. The aspect of correlating variables is an important remark, which we thoroughly discuss in context of the "main comment 2" (part b), see further above). In fact, this is fundamental for our multi-dimensional approach: characterizing systematic and random uncertainty estimates of U, $q_s$, and $q_a$ as a function of atmospheric state parameters, which (as we believe) have an impact on the parameters themselves. Specifically regarding Figure 2d): wvpa is not available from in situ measurements, which is why a bias dependency on in situ wvpa cannot be investigated.

**Changes in the manuscript:** See "main comment 2" (part b) further above.

**9 ) Comment from referee:** Page 9, Line 21: Please specify what "even stronger winds" are.

**Author's reponse:** "stronger wind" mean wind speeds exceeding 20 m s$^{-1}$.

**Changes in the manuscript:** The wording has been changed in the revised manuscript (P.10, L.9).

**10 ) Comment from referee:** Page 9, Lines 24-26: This paragraph is hard to read. After reading, a number of times times, I think I understand. Is it not better to say: "Our goal is to document the upper bound of the bias and therefore we take the absolute value of the possible systematic error in CE"?

**Author's reponse:** We agree that this paragraph is somewhat confusing and out of place. It has been moved further up into the appropriate context.

**Changes in the manuscript:** The wording has been modified and has been moved further up into the appropriate context (P.10, L.17ff).

**11 ) Comment from referee:** Page 10, line 15 and page 11, line 7: I suggest to replace "Next to" by "In addition to"

**Author's reponse:** Thank you for this suggestion.

**Changes in the manuscript:** The wording has been changed in the revised manuscript (P.12, L.12).

**12 ) Comment from referee:** Page 11, section 3.5: This section is hard to read. If I understand correctly, it addresses the question: Does it matter for the averages that the satellites sample the ocean at particular times of the day only, given that a diurnal cycle may be present? The authors investigate by looking at buoy data and by comparing averages that cover the full diurnal cycle with samples at satellite overpass times only. Part of the confusion is because it mentions spatial sampling, but I don't think this section covers that? Please simplify for clarity.

**Author's reponse:** Exactly. For the monthly mean HOAPS product (HOAPS-G), sampling uncertainties need to be quantified because of the diurnal cycle of the geopyhsical parameters. Due to the sun-synchronous satellite overflights, diurnal cycles or frontal passages are likely to be missed. This will affect the monthly mean averages. We agree that the the term "spatial sampling" is misleading, as we only cover the temporal sampling issue.

**Changes in the manuscript:** The aspect of "spatial sampling uncertainties" has been removed from the revised manuscript to avoid confusion.

**13 ) Comment from referee:** page 13, Lines 9-11: I am not sure that it is helpful here to refer to Fig. 1a, because it is showing the combination of E_ins($q_a$_a) and E_retr($q_a$), which is different from

Fig. 2a. The authors point this out but instead of clarifying something it confuses.
**Author's reponse:** We agree that this may be confusing.
**Changes in the manuscript:** This section has been shortened to become more clear. (P. 14, L.18ff)

**14 ) Comment from referee:** Page 13, 23: Suggestion: replace "merely" by "only"
**Author's reponse:** Thank you.
**Changes in the manuscript:** "merely" has been replaced by "only". (P.14, L.31).

**15 ) Comment from referee:** Page 13, Line 24: What is meant by "local minimum in that region for q_a"? E_retr(q_a) has a maximum over the warm pool.
**Author's reponse:** This is a mistake in our manuscript, thank you for pointing this out. We wanted to point at the $q_a$ random retrieval uncertainty, not $q_a$ itself.
**Changes in the manuscript:** the wording has been changed (P.14, L.31f).

**16 ) Comment from referee:** Page 13, Line 29: In the sentence "Respective values partly exceed 50 W/m$^2$", what is meant by "respective" and "partly"? Do the authors mean: "In these areas, values are found in excess of 50 W/m$^2$"?
**Author's reponse:** Yes, this is correct.
**Changes in the manuscript:** The wording has been changed in the revised manuscript. (P.15, L.3).

**17 ) Comment from referee:** Page 14, Line 33: "direct eddy covariance" is not wind speed.
**Author's reponse:** Sorry for not being correct here. The wind stresses are based on inertial-dissipation methods. Together with eddy covariance based LHF, the turbulent fluxes of a variety of satellite, reanalysis, and combined products are evaluated.
**Changes in the manuscript:** The wording has been changed. (P.15, L.31ff).

**18 ) Comment from referee:** Page 16, Lines 1-2: This is an interesting example, where it is explained that $q_a$ retrievals may be in error because of dry air advection. However, it is not clear how the systematic error analysis picks up the area of the Agulhas current. The systematic error estimation is entirely driven by U, $q_a$ and SST and wvpa (if I understand correctly).
**Author's reponse:** It is correct that the systematic uncertainty estimation is entirely driven by combinations of ambient U, $q_a$, SST, and wvpa. Our multi-dimensional bias approach does not point at specific regions. This implies that we cannot be 100% certain that the observed uncertainties over the Agulhas region are exclusively associated with local retrieval issues. In general, match ups over a region contribute to the look up tables (LUTs), which implies these regions are somewhat mirrored in the LUTs. However, they are not explicitly resolved.
The following serves to explain how the LUTs pick up the Agulhas region: Figure 1 (left) indicates that numerous collocations between buoys/ships and satellite exist in this area, which is characterized by a unique combination of ambient U, $q_a$, SST, and wvpa. In case of the mentioned dry cold air outbreaks from the South, $q_a$ will be anomalously low and hence $q_s$-$q_a$ and LHF anomalously large. According to Santorelli et al. (2011), satellite retrievals seem to encounter difficulties with these dry cold air outbreaks, which implies that they will not capture $q_a$ correctly. This would for example be seen when investigating $dq_a$. That is, differences between satellite and in situ $q_a$ would be negative, which directly impacts our four-dimensional uncertainty analysis. In conclusion, repetitive retrieval issues over a specific regions will be manifested in the LUTs and will eventually be seen in systematic uncertainty maps. At the same time, underestimated $q_a$ along the Agulhas Current contribute to an increase in the random uncertainty component of the LUTs.

**19 ) Comment from referee:** Section 4.6 and Fig. 4: Here both systematic and random errors are discussed region by region and climatologically versus January/July. Earlier in the paper it was concluded that the random errors were small compared to the systematic errors. However in Fig. 4 the random errors are larger than the systematic errors. Furthermore I would expect that the

climatological data (I assume averaged over the entire period) has much more data than the January or July data and therefore much smaller random errors.

**Author's reponse:** As mentioned in Sect. 4.6, the error bars in Fig. 5 point at *instantaneous* random uncertainties (such as those shown in Fig. 3). The idea is to show the maximum uncertainty to be expected for a specific region and season on an *instantaneous* basis. This approach allows for illustrating random uncertainties, as they often even exceed the systematic counterpart for pixel-level data, as is seen when comparing Fig. 3 to Fig. 4. Fig. 3 shows averaged *instantaneous* random uncertainties as a function of region and time. If properly scaled according to the considered period of time, they decease with increasing time period and become insignificant at monthly or multi-annual (that is, climatological) time scale. Keeping this in mind, this also answers the question as to the smaller random uncertainties for multi-annual mean (1988-2012) compared to seasonal means (1988-2012): The difference in error bar magnitudes is not related to averaging periods, as these are averaged *instantaneous* random uncertainties as a function of region and time. We agree, however, that this is not clearly stated in the manuscript.

**Changes in the manuscript:** The wording as been modified in the revised manuscript to clarify what is shown in Fig. 5. (P.18, L.29-33) This also targets the caption of Fig. 5.

**20 ) Comment from referee:** Page 18, Line 31: Please replace "outperforms" by "exceeds"
**Author's reponse:** Thank you for this suggestion.
**Changes in the manuscript:** This section has been considerably shortened. The phrase is no longer included in the revised manuscript.

---

### Cited:

**Andersson,** A. Graw, K., Schröder, M., Fennig, K., Liman, J., Bakan, S., Hollmann, R.and Klepp, C.: Hamburg Ocean Atmosphere Parameters and Fluxes from Satellite Data - HOAPS 4.0, Satellite Application Facility on Climate Monitoring (CM SAF), doi:10.5676/EUM\_SAF\_CM/HOAPS/V002, 2017.

**Freeman**, E., Woodruff, S. D., Worley, S. J., Lubker, S. J., Kent, E. C., Angel, W. E., Berry, D. I., Brohan, P., Eastman, R., Gates, L., Gloeden, W., Ji, Z., Lawrimore, J., Rayner, N. A., Rosenhagen, G., and Smith, S. R.: ICOADS Release 3.0: a major update to the historical marine climate record, Int. J. Climatol., p. in press, doi:10.1002/joc.4775, 2016.

**Kinzel,** J., Fennig, K., Schröder, M., Andersson, A., Bumke, K., and Hollmann, R.: Decomposition of Random Errors Inherent to HOAPS-3.2 Near-Surface Humidity Estimates Using Multiple Triple Collocation Analysis, J. Atmos. Oceanic Technol., 33, 1455–1471, doi:10.1175/JTECH-D-15-0122.1, 2016.

**Reynolds**, R. W., Smith, T. M., Liu, C., Chelton, D. B., Casey, K., and Schlax, M. G.: Daily High-Resolution-Blended Analyses for Sea Surface Temperature, J. Climate, 20, 5473–5496, doi:10.1175/2007JCLI1824.1, 2007.