

Response to reviewer #2 (David Whiteman):

The authors are very grateful to Dave Whiteman for the numerous comments and suggestions which really sharpened the interpretation and presentation of our results. Corrections prompted by the referee appear in red in the revised manuscript.

Answers to the “General comments” section

The authors are comparing lidar measurements of mixing ratio with point measurements in the near field as well as integrated total column water versus GPS. Both of these methods are potentially influenced by the lidar overlap function, yet this issue is not addressed. Also, the lidar profile is extended upward beyond a certain altitude through the use of a radiosonde profile and model. Some quantitative discussion of how much the non-lidar portions of the profile influence the entire column is needed.

>> More quantitative discussion of how much the non lidar portion of the profile influence the entire column is given in the text according to the specific comments (see below).

The authors make use of measurements using a nitrogen filter in both the water vapor and nitrogen channels to monitor drift in measurements. Exactly how this is implemented is not clear from the current discussion. Presumably the water vapor filter is removed to perform this test, in which case the drift measured does not include any changes in the water vapor filter. Also, since the PMT cathodes presumably have different quantum efficiency at the water vapor and nitrogen wavelengths, it seems that using the channel ratios to correct the drift must involve some calibration normalization procedure which is not explained.

>> We detailed about the N2 calibration at the beginning of section 5.2. As hypothesized by the reviewer the procedure implies the removal of filters, the drift of which is therefore not measured with this procedure. Note that both the water vapor and the usual nitrogen filters are removed. An adjustment to external water vapor measurements is needed to achieve an absolute H2O calibration. The N2 calibrations are easy to perform and are used to monitor and correct for the variations in the detection system (changes in the optical alignments, fluctuations in PMTs quantum efficiency and gain). They are capable of detecting both short-term and long-term fluctuations.

There are some references that are not well chosen to support the statements that are being made. Finally, there are numerous awkward uses of the English language which sometimes make for confusing sentences. I point some out in the specific comments. The authors are encouraged to address these and the other awkward uses of English for improved clarity of the discussion.

>> References have been added where the reviewer suggested.

>> English has been checked and made as clear as possible.

Answers to “Specific comments” section

1) Line 17. The term “a unique nitrogen filter” is used where something like “a single nitrogen filter” is likely intended.

>> Done.

2) Line 30. The statement made is that the ability of the Raman lidar to make water vapor measurements on small spatial and temporal scales “is of prime interest for weather forecasting and climate studies”. The citation given is Held and Soden, 2000. Do Held and Soden explicitly mention Raman lidar for providing such measurements? The way this sentence is structured, that would seem to be implied.

>> For sure, the citation was not intended to suggest that. The beginning of the introduction was slightly reorganized and this sentence was removed.

3) Paragraph starting with Line 46. A discussion of independent and dependent calibrations is given. But the authors neglect to discuss the work of Venable et al, which demonstrated a scanning lamp technique using, Avila, 2004 water vapor cross section values, that agreed with corrected radiosonde to better than 5%. The Avila 2004 cross section calculations were based on a full ab initio simulation of the water vapor molecule unlike the earlier model used in Sherlock et al, which only considered trace scattering as I recall.

>> We added this discussion at the end of the paragraph, including the reference to Venable et al. 2011.

4) Line 58. *Convolutd* → *Convolved*

>> Changed

5) Line 59. Statement is made that the absolute accuracy of the independent method is not better than 10% and is not referenced. I suspect the authors are referring to the uncertainty of 10% given for the Raman cross section ratio in the Penny and Lapp 1976 paper. The Avila cross sections agree to within 8% of these numbers (and are thus within the P&L uncertainty estimate) but the Venable et al. calibration work shows agreement with carefully performed radiosonde-based calibration to better than 5%. These two results would tend to indicate that the absolute uncertainty of the integrated q1 cross section based on the Avila, 2004 work is better than 10%. In fact, Fernandez-Sanchez, the lead of the Spanish group that performed the Avila, 2004 cross section measurements, has stated (private communication) that the uncertainty of the integrated water vapor q1 Raman cross section from the Avila work is approximately 5%.

>> The following text has been added: "More recent determination of water vapor Raman cross sections (Avila et al., 2004) achieved an absolute accuracy better than 10%."

6) Lines 78-82. There are various ways that the technique described in Leblanc and McDermid, 2008 can fail to detect sources of calibration variation. Those failure modes are discussed in the reference below which should be cited. It is for the reasons in the cited paper that NDACC did not, in fact, adopt the technique described in Leblanc and McDermid, 2008 as the standard calibration technique within NDACC, despite what has been published. At the NDACC Raman water vapor lidar workshop held in Greenbelt, MD in 2010, which I chaired, we recommended the use of one or more lamps for assessing instrument performance but neither the hybrid technique, nor any other technique, was selected as a reference calibration technique for NDACC. Whiteman, D. N., D. Venable, E. Landulfo, Comments on "Accuracy of Raman lidar water vapor calibration and its applicability to long-term measurements", *Applied Optics* Vol. 50, Iss. 15, pp. 2170–2176 (2011)

>> The discussion lines 78-82 has been changed to:

"Lately, Leblanc and Mc-Dermid (2008) developed a stability control procedure with calibrated spectral lamps, reaching a standard deviation of 2% of the calibration factor over more than a year. Venable et al. (2011) proposed an improvement of this calibration technique which consists in scanning a known light source over the telescope aperture instead of a stationary lamp system. Advantages and drawbacks of these methods are discussed in Whiteman et al. (2011)."

7) Page 5. The discussion here seems to imply that there was no lidar sub-system for automating the bore-site alignment. It would be helpful for the authors to address the question of how many of their troubles may have been reduced or eliminated if they had had an auto-alignment system that could hold the alignment to within, say, 10-20 microradians. This is the performance that we see in using a variation of the Licel bore site system that is commercially available. Such questions could perhaps be addressed through the authors' use of the Zemax program.

>> The resolution of the alignment for our system is about 40 μ rad. And it can be expected that an auto-alignment with a 10-20 μ rad resolution would improve the stability of our system.

8) *Figure 3. These fiber optic mode patterns are consistent with those presented in the Whiteman et al, 2011 paper cited in #6 above. This reference would be helpful to include here.*

>> Reference has been included.

9) *Figure 4. This figure is quite small and there are various notations in the figure that are not described. The figure could be expanded and much greater detail used in explaining the figure and the optimization procedure. But the authors describe standard optical optimization considerations, the details of which I am not sure are of benefit to the reader. So I suggest that Figure 4 be eliminated and the optimization discussion reduced to speak in more general terms about what was tested and what was finally determined by the optimization.*

>> The discussion has been reduced and the title of the section was changed to “Correction of the vignetting problem in the detecting system”. However, we decided to keep the figure because we think it is useful to visualize the three optical configurations discussed in section 4 (cf. figures 6, 7, and 8). The figure was enlarged and more information is described in the captions and in the text. The new optical configurations are referred to as option 1 and 2. Both are optimized based on standard optical considerations but take different opto-mechanical constraints into account.

10) *Line 255. The term “adjustment tolerance” is used. Please be more specific in describing what adjustment you are referring to: x,y,x displacement? angular displacement?*

>> We completed the term “adjustment tolerance” with “along the optical axis”.

11) *Sections 4.1, 4.2 and 4.3 all show significant improvements in the calibration stability using the optimized optical configuration. It would be helpful for the authors to clearly state what aspects of the optimized optical design lead to this improvement. From my reading of the paper, this was not clear. Was it the removal of vignetting? Use of a larger portion of the photocathode? Reduced wandering of the spot on the photocathode due to displacements of the input beam?*

>> The improvements in calibration stability using the new optical configurations (option 1 and option 2) are explained by the combined removal of vignetting on L1 and a more careful alignment of the optical beam in the detection system avoiding also vignetting on other optical elements and on the PMTs. We removed the first sentence of section 4 and clarified the aim of the experiments at the end of the paragraph:

“The experiments are repeated three times for the different options of the detection system layout: initial (configuration of the Demevap campaign), option 1 and option 2 (Figure 4). The configuration of the Demevap campaign includes vignetting which, combined with beam wandering on the PMT surface, is a major source of instability.”

12) *Section 4.3. If I understand the experiment described, the authors translated the device shown in figure 5 along the optical axis to simulate the range dependence of the spot size that impinges on the fiber optic. In an actual lidar implementation, the angular information in the focusing lidar beam will also change as the laser beam propagates away from the telescope. It does not seem that this effect is simulated in the experiment that the authors performed. Please clarify this and comment on whether this has any effect on the results presented.*

>> The angular variation expected at the input of the fiber is of order 1 mrad (diameter of fiber/focal length of telescope). We did not simulate this variation source and agree that it may have some impact. We added a sentence in the text to mention it: “...moving away the fiber from the tube which simulate a variation of spot size (but not of the angular variation associated with NA variations in the actual lidar implementation which might also have some impact).”

13) Line 363. The authors refer to “photon detection threshold”. I believe that this is the quantity that is usually expressed as “discrimination level” which should be provided with a unit, which I assume is mV in this case. Did the authors use the standard technique of determining the pulse height distributions of the PMTs to set these levels? If not, how did they do it?

>> We replaced “photon detection threshold” with “discrimination level”. The levels were not given in mV but we removed this technical detail. We changed the sentence into: “The discrimination levels have been determined using the standard technique of pulse height distributions”.

14) Line 364. The authors state that after making the adjustment discussed in #13 that “the detected photon rate matched properly the expected Poisson Law”. More detail is needed here for the reader (at least this reader) to understand what is meant by this statement.

>> The Raman signals are assumed to follow a Poisson distribution which states that the mean is equal to the variance. We clarified it in the text:

“We also checked that the detected photon rates matched properly the expected Poisson law (mean equal to variance).”

15) Line 366. This is another example of where the presence of an auto-alignment system could have benefited the measurements shown here. At some point in the manuscript, please comment on the feasibility of having such a device in the lidar system and which of the problems would be reduced or eliminated if such a system were in place.

>> We agree with this suggestion, and added a sentence in the conclusion: “Auto-alignment of the transmitter and the receiver would further improve the short-term stability.” But first, we think it will be necessary to better control the thermal deformations of the optical bench.

16) Equations (1) and (2). The authors earlier describe the Sherlock technique of convolving the instrument transmission function with the Raman cross section. The authors also state that they are accounting for the temperature dependence of the Raman cross sections in their calculations, but the older versions of the equations given here do not capture those temperature dependent details. Please use equations that account for the variation of Raman cross sections with temperature to agree with your earlier discussions.

>> We added the temperature-dependent functions $F(T)$ in equation (2), consistently with Whiteman 2003b.

17) Discussion between lines 380 – 385. The authors refer to uncertainties of laboratory components (10%), differential transmission (2-5%) and Raman cross sections (10%). As mentioned earlier, Venable et al., achieved calibration agreement between their independent technique and that of corrected radiosondes to better than 5%. Thus, the numbers that the authors use here seem too large. Also, for the purposes of climate studies, it is very important to distinguish between random and systematic uncertainties. Let us assume that an independent calibration is performed using a Raman cross section ratio of X . If at a later date the cross section ratio is found to be Y , earlier data may be reprocessed simply by multiplying by the ratio of X and Y . In other words, systematic uncertainties, once they are discovered and quantified can (and need to) be corrected. Such correction is much more difficult to do when using a dependent calibration technique since issues of atmospheric variability may complicate comparisons. The ability to reprocess data when higher accuracy cross section information becomes available is a strong argument in favor of independent calibration techniques in developing time series of water vapor mixing ratio.

>> We agree and changed the text: “The uncertainties associated with these terms are typically 5-10%. Venable et al. (2011) achieved calibration agreement between their independent technique and that of corrected radiosondes to better than 5%.”

18) Line 394. Reference is made to an a priori calibration coefficient. What is this a priori calibration?

>> By a priori calibration we meant a calibration coefficient obtained from previous measurements. We removed the sentence because it is misleading. The goal of the calibration is to determine the coefficient from new measurements (even if in practice this consists in two factors composed of an a priori coefficient and a correction factor).

19) Lines 395 – 400. Discussion is given concerning the comparison of lidar measurements in two ways: 1) comparing a lidar range cell to a point sensor and 2) comparing IPW derived in part from the lidar profile with that of GPS. There are several issues that need discussion in these comparisons including that of the influence of the lidar overlap function on these calculations. The authors should not assume that the overlap functions for the water vapor and nitrogen channels are identical. The degree to which the channel overlap functions cancel in the ratio of H₂O and N₂ channels needs to be demonstrated within some level of uncertainty. See also #24.

>> Our laser beam is sent coaxially with the receiving telescope, so we think that we don't have significant differential overlap variations. However, we estimated the overlap functions as the ratio of H₂O and N₂ Raman channel signals with measurements made with a common N₂ filter (in N₂ calibration configuration) over all measurements made during the Saint-Mandé campaign. We found that the ratio reaches unity at about 150 m and below that altitude the overlap effect varies in the range of +/- 3%. The determination of the overlap ratio is not very accurate at the lowest elevations because the measurements are very noisy. We decided not to correct this effect. So a residual bias might affect the PTU calibration results which can be seen on figure 12 (up). This description has been included in the manuscript.

20) Lines 400 – 405. Discussion is given of using the radiosonde profile to extend the lidar profile in order to calculate the total column water vapor. Additional discussion is needed that quantifies the influence of errors in the radiosonde profile (due either to wet/dry bias, time lag, collocation issues, etc) on the derived total column water. See also #24.

>> The radiosondes used are MODEM M10. Recent studies show that they are consistent in quality with Vaisala RS92 radiosondes (e.g. Ingleby, 2016). This information has been included in the manuscript.

21) Section 5.2. N₂ calibration procedure. This procedure seems quite similar to what Vaughan et al implemented and what we discussed in our 1992 paper. If so, please provide those references here. Later in the paper (around line 512) the point is made that dispersion of PMT cathode quantum efficiency is a factor in applying this technique. Those points should be made here when the technique is being introduced. At some point, authors should consider how that dispersion might be affecting their results, such as whether some kind of normalization is needed to make the calibration corrections that are mentioned below.

>> References were added.

>> Regarding the dispersion of PMT cathode QE we included the following discussion in section 5.2: "It should be noted that the ratio is computed with measurements at the same wavelength whereas in operations different filters are used in both channels, and measurements are made at different wavelengths (386.7 nm and 407.6 nm). This difference is corrected with the final calibration factor based on comparison with external WVMR measurements."

22) Lines 420-421. Authors state that the slope obtained for the upper layer during Demevap is not well determined due to low SNR. Then they speculate that "this drift" may be caused by aging optics or deposition of dust. Presumably the authors refer to some other slopes in the table, but not the ones for the upper layer during Demevap, which are shown to be 0 within the uncertainty. Please reconcile this.

>> The slope obtained for the upper layer during Demevap was not well determined because of a few outlier points associated with low SNR. We removed values with standard errors > 0.3 and fitted

again the slopes for all three layers. The results for all three layers are now consistent. Figure 11 and Table 2 were modified accordingly.

Also, some of the slopes shown are as large as 8.7%/month which seems a very large drift. I would find it difficult to believe that aging of components could explain such a change unless, for example, some electronic component was in the process of failing. Was that the case? Regarding the possibility that dust deposition could cause such a drift, I would think that such dust would need to be a brown carbon-like substance that demonstrates considerable dispersion of UV absorption. Were the components subjected to outbreaks of combustion aerosols that could have caused brown carbon to be deposited on the optics? And was that deposition continuous so as to explain the drift observed? In summary regarding Table 2, the large slopes shown for some cases do not seem well explained in the current discussion. Are there other effects that could explain the drifts in addition to aging of components and deposition of dust that absorbs differentially in the UV?

>> Inspecting carefully the time series of calibration factors shows that the values do not align well on a straight line and thus that one should be cautious about interpreting the drift as a linear one. For example, during the Saint-Mandé campaign, some PTU calibration values in April are as low as in July. We thus think that the large temporal fluctuations in the calibration factors are mainly explained by displacements or changes in the spatial patterns of the optical beams on the PMTs due to re-alignments of the transmitted beam or interventions on the system (e.g. we made tests with different interference filters; but only those using the same Barr Filters as described in the manuscript are shown here). In general, a better control of the optical alignments during the Saint-Mandé campaign explains the smaller overall drift and short-term fluctuations.

23) Line 424. Column 5 is referred to where, I suspect, column 4 was intended.

>> Changed.

24) Section 5.3.

1. The authors discuss using a PTU sensor, that is earlier stated to be at a height of 15m above the ground, to compare with a range cell of the lidar that is centered at a height of 67.5m. It was no doubt a practical choice to not use lidar data at a height of 15 m above the ground since it is very difficult to get reliable lidar measurements so close to the ground. But this non-collocation is certainly an issue for calibration and needs to be discussed.

>> During the Saint-Mandé campaign, a second PTU sensor was mounted on another building at a height of 25 m above the ground. Comparisons between the two sensors at 15 and 25 m gave us an estimate of the gradient of humidity in the surface layer of 0.1 g/kg per 10 m. Extrapolated to the height of the lidar range bin the resulting bias would be about 5% but we think that the vertical humidity gradient above the first 10 m is smaller and the real bias would be less. We included a discussion of this point in section 5.1.

The authors could present a high temporal resolution time series of comparisons to justify the technique to give the reader confidence that the lidar measurements are consistently reliable at such short range. Our experience is that there is considerable random variation in such near range measurements under the best of circumstances and, if the alignment is permitted to wander, significant systematic uncertainties can be introduced as well. So, I think the authors need to convince the reader that such point comparisons in the near field are valid to be making beyond what is shown in Figure 12.

>> A time series comparing lidar and PTU measurements was included in the new manuscript as new Figure 12 in section 5.3.

Then they should discuss what systematic differences that might be expected due to the fact that the PTU sensor and lidar are not sampling the same atmosphere.

>> A discussion of this point is included in the manuscript (see answer to comment above).

The profile comparison shown in Figure 9 is somewhat difficult to see but would seem to indicate a significant difference between the lidar and PTU values of mixing ratio near the surface.

>> There was indeed a difference in the values of mixing ratio because the lidar profile was calibrated with an a priori constant (see answer to question 18). This is corrected since we now show a calibrated lidar profile and also add a second plot in log-log scale to better show the profiles in the surface layer.

Please also clarify the range resolution of the lidar measurements in this discussion.

>> This point is clarified at the end of section 5.1: “The integration time of the lidar profile is 5 min and the vertical resolution varies with the altitude between 7.5 m and 240 m (Bosser et al., 2007)”

During the Dimevap campaign, the authors mention using 2 point sensors at different ranges and compare with slant range lidar measurements. Even though it may be described in more detail in earlier publications, some additional description of this technique would be helpful here to contrast with what was done during the Saint-Mande campaign.

>> This discussion has been removed because the results reported in Table 3 have been updated after a new reprocessing of the data. The results from both campaigns are slightly better and more consistent after the N2 calibration.

2. Regarding the GPS comparisons, the authors use a composite water vapor profile derived from 3 sources of information (lidar, sonde, model) to compare with GPS. Please quantify, on average, what fraction of the IPW is sampled by the lidar and what fraction is quantified by these other sources. Then discuss reasonable estimates of the uncertainties of the nonlidar data sources and, from this, provide an estimate of the total uncertainty in IPW due to the use of three data sources for the profile. Since most of the IPW is measured by lidar, it is possible that this technique can provide a useful comparison to GPS for the purposes of lidar calibration, but the burden is on the authors to convince the reader that the technique is a valid way to derive a lidar calibration. The current description is not sufficient to accomplish this task.

We used the lidar profile up to 5 km, then the radiosonde profile from Trappes up to 10km. The fraction of the lidar ZWD represents 90 to 95% of the total ZWD. The remainder is covered by the sonde measurements while the model component is always below 0.1 mm (< 0.1%). The model is mainly useful for the dry part of the refractivity profile which is used in a ray-tracing algorithm. But since it has negligible contribution to the total ZWD, we removed it from the discussion. Now, assuming that the accuracy of the sonde measurements in the layer between 5 and 10 km is about 10-20%, the accuracy of the correction is about 1-2%. This is the expected accuracy of the GPS ZWD calibration technique. This discussion is included in section 5.3.

25) Line 446. Authors state that the drift in H2O calibration is consistent with the N2 calibration results and that the drift can thus be corrected. But the authors do not describe how they make the correction although they may be using the same technique described in Whiteman et al., 1992. But, since the N2 calibration technique involves measurements in the H2O channel at a different wavelength than for measurements of water vapor, the question of dispersion of quantum efficiency of the photocathode must be addressed. Also, as stated earlier, the authors did not really describe how this technique is employed and they should add such description. But, assuming that there is some QE difference to account for in the water vapor channel, how are the N2 calibration results actually used to adjust the drift in water vapor calibration? It would seem that some normalization of results is needed. If that is the case, how is the normalization done and what assumptions are made in the normalization technique?

>> The correction based on N2 calibration measurements is made by applying a linear function of time fitted to the N2 calibration results. This procedure is applied here because the data cover a limited period of time. It is different from the “operational” procedure described in Whiteman et al (1992) which applies a different correction coefficient each night. This description is added in the text (also in answer to a question from referee n°1):

“These coefficients have been corrected by linear function of time fitted from the N2 calibration results. We did not use the night to night coefficients correction presented by Whiteman et al. (1992) because our N2 calibration coefficients are noisy and we did not want to add dispersion to the H2O calibration coefficients.”

Regarding the normalization necessary because of the photocathode dispersion, this difference is corrected with the final calibration factor based on comparison with external WVMR measurements (see also answer to question 21).

26) Figure 12. Upper panel shows a significant difference between the calibration constants derived from GPS and PTU comparisons (perhaps as large as 10%) as well as a consistent drift of both calibrations. Lower panel shows the corrected times series where the drift has been essentially removed but also where now the relative calibration constants are both near 1.0. How has the N2 calibration been used to correct the difference in calibration values shown in the upper panels? This is not clear and would seem to require an assumption about what calibration source is “better” than the other or, at least, that some portion of one time series comparison is “better” than another. Please explain how this was done since it is really not clear at this point.

>> The systematic difference in calibration constants derived from GPS and PTU (observed on both campaigns) can be explained partly by the differential overlap function which is not corrected and partly by the bias in the PTU calibration due to non-collocation of lidar and PTU measurements (see answer to question 24). This remark has been included in the manuscript.

>> The correction using the N2 calibration results is explained above.

27) Line 499. Authors state that the main problem of the lidar calibration drift is the thermal expansion of the optical bench. This seems quite unusual an explanation to explain a long term drift and would require that the mean temperature of the lidar system is somehow changing consistently over extended periods. Was that the case? What was the range of temperatures experienced by the lidar system? I can understand how thermal expansion could perhaps explain drift within a single session but do not see how it could result in drift over extended periods of time since presumably the temperature of the lidar system will return to some mean value (determined by the lidar enclosure temperature) between experiments. Is the lidar system open to the atmosphere during measurement sessions and was the external temperature consistently changing in one direction during the experimental campaign? Still do the materials of the optical table have such a large coefficient of thermal expansion to explain the calibration drift seen? Such things can be modeled with the Zemax software. Did the authors perform such a modeling experiment with Zemax that showed the drift could be explained by the range of temperatures experienced by the lidar system?

>> See answer to question 22. The discussion starting line 499 was revised accordingly.

28) Line 500. Thermalization → Athermalization

>> Done.

29) Lines 511-513. This is where the authors consider the issue of dispersion of PMT QE in the N2 calibration technique. These points need to be made earlier and discussed in much more detail as requested above.

>> Done.