## Responses to referee comments and changes to the manuscript

*Quality control and error assessment of the Aeolus L2B wind results*
*from the Joint Aeolus Tropical Atlantic Campaign*
*(https://doi.org/10.5194/amt-2022-223)*

We appreciate the referees' very insightful and helpful remarks on our manuscript. The responses to the individual comments and the corresponding changes that were made to the manuscript are presented in the following.

### Response to Referee Comment 1

General comment #1.1:

*The paper has a bit of a gift wrap structure. The methods and goals are only rather vaguely touched in the beginning of the introduction, and become clear more and more only as one reads on. I recommend to give a more specific outline for what the reader is to expect.*

*Somewhat related, I find the naming / separation of sections 3 and 4 a bit confusing. Both sections are structured pretty much identically, Sec. 3 showing results for the model evaluation against ECMWF, Sec. 4 showing results for the AVATAR-T evaluation. The methods are indeed introduced in Sec. 3 as the name suggests, but it is somewhat of a results / methods hybrid.*

*I can see the benefit of introducing the methods on the examples with ECMWF, but I think a dedicated methods-only section that gives a summary (perhaps even a schematic figure?) of all methods that are going to be used could help the reader to have a clear picture on what to expect. As mentioned earlier, right now the exact methodology reveals itself only over long stretches of the manuscript.*

Response to General Comment #1.1:

Thank you for this comment. We agree that the goal and structure of the paper should be outlined more clearly. Therefore, we have extended the introduction to better describe the objective of the paper and the methodical approach used to achieve this (new text is underlined):

> *[...] Therefore, a more detailed treatment of different QC schemes and how they affect the resulting statistics is necessary for comparable validation results and for a more objective assessment of the Aeolus wind data quality. Moreover, it is an important aspect with regards to the operational data assimilation in NWP centres and allows for a more rigorous error characterization of the Aeolus winds.*

*This paper aims to raise the awareness to the influence of the chosen QC schemes on the validation results, particularly when using the L2B EE. It also demonstrates the usefulness of specific statistical tools for the purpose of outlier removal and the assessment of normality, which are necessary to retrieve the Aeolus wind errors in accordance with the MRD. The presented methods are applied in the context of the AVATAR-T validation campaign in 2021 by comparisons against the ECMWF model background winds and the 2-µm DWL wind data. The specifics of the campaign and available datasets are outlined in Sect. 2 together with a description of the L2B Rayleigh-clear and Mie-cloudy EE and their temporal evolution over the past three years. <u>The model comparison in Sect. 3 serves as an example to introduce the reader to the detailed treatment of the Aeolus wind data in terms of QC and error assessment. In particular, the modified Z-score (Sect. 3.2) and normal quantile plots (Sect. 3.4) are discussed as powerful tools for removing gross errors and assessing the normality of the wind error distribution, respectively.</u> In addition, the impact of the QC settings on the results from the model comparison is elaborated (Sect. 3.5). In Sect. 4, the statistical methods are then applied to the comparison of Aeolus wind observations against 2-µm DWL data. The paper concludes with a summary and outlook to future studies of the L2B wind error characteristics in Sect. 5.*

Regarding the structure and content of sections 3 and 4, we think that the statistical methods can be best illustrated by a concrete example rather than by an abstract explanation in a dedicated methods section. Most of the presented statistical parameters and tools refer to the Aeolus wind error, which is defined as the difference between the wind speed measured by Aeolus and that determined by the model or instrument employed for the validation. Consequently, introducing, for instance, the normal quantile plot in section 3.4 or the combined bar and line graph in section 3.5, necessitates a wind reference to demonstrate the relevance of these plots for assessing the wind error characteristics and the influence of the used QC scheme. For the purpose of keeping the text as concise as possible, despite the complexity of the topic, we decided to use the model comparison as an example to guide the reader through the steps of a careful statistical analysis of the Aeolus wind error. In section 4, these guidelines are then transferred to the validation against the 2-µm DWL. We think that with the extended introduction, as described above, this approach will become clear to the reader and facilitate the understanding of the statistical methods.

Also, we have renamed section 2 from "Datasets and methods" to "Datasets and the L2B estimated error", as this header better describes the content. In this manner, it becomes clear that the methods are introduced in section 3 on the example of the model comparison.

<u>Specific comment #1.1:</u>

*L25: I find the term "biased gross errors" a bit unusual and unclear.*

Response to Specific comment #1.1:

The term was changed to "positively biased wind results".

Specific comment #1.2:

*L68: Is a Gaussian distribution actually appropriate for wind retrievals? Since wind speed is a bound variable, shouldn't the error distribution get more and more skewed as the retrievals get closer to the limit? I appreciate this is a common assumption for many variables that has also some practical reasons, but since such a strong focus is put on forcing the data into a normal distribution, some words on that might be helpful.*

Response to Specific comment #1.2:

Our study follows the error definitions that are formulated in the Mission Requirements Document (MRD) (ESA, 2016), which assume Gaussian distributions for the Rayleigh and Mie wind error with respect to other wind observations or the model background. This is justified by the fact that the wind error is dominated by Poisson-distributed photon noise on the detectors, particularly for the Rayleigh channel. For the Mie channel, deficiencies in the signal analysis (Mie Core algorithm) give rise to gross errors which additionally contribute to the error distribution. The assumption in the MRD that the Mie gross errors are uniformly distributed did not prove correct, as pointed out in Sect. 3.3. Nevertheless, after sorting out the gross errors, the assumption of a Gaussian distribution is also valid for the Mie channel. When the wind retrieval gets closer to the limit, e.g., at low signal levels, the random error will increase with the noise accordingly, but without skewing the wind error distribution.

We have emphasized these considerations in the revised manuscript in the context of the error definitions (Sect. 3.1):

> *The definitions given here are in line with the those stated in the Aeolus MRD, which assumes that the wind error with respect to other wind observations or the model background can be described by a Gaussian distribution whose centre and width represent the accuracy and precision of the Aeolus winds. This is justified by the fact that the wind error is dominated by Poisson-distributed photon noise on the detectors, particularly for the Rayleigh channel. For the Mie channel, deficiencies in the signal analysis (Mie Core algorithm) give rise to gross errors which additionally contribute to the error distribution, as will be pointed out in Sect. 3.3. Nevertheless, after sorting out the gross errors, the assumption of a Gaussian distribution is also valid for the Mie channel.*

Specific comment #1.3:

*L91-93: This statement is distracting and unnecessary here. Also, I have the impression that "Rayleigh-clear" and "Mie-cloudy" are used more often than not later on, so perhaps just delete this statement altogether (or stick with the simple notation consistently, which I'd actually prefer).*

Response to Specific comment #1.3:

We removed the statement from the text.

Specific comment #1.4:

*L103: Where does this strong signal decrease over time actually come from?*

Response to Specific comment #1.4:

The root cause analysis of the decline in the atmospheric return signal is still ongoing. Laser-induced contamination, laser-induced damage and bulk darkening of the instrument's optics are the most probable causes, in addition to clipping losses at the instrument field stop. We have added a short comment on the potential root causes for the signal loss to Sect. 2:

> *The root cause analysis of the decline in the atmospheric return signal was still ongoing as of the writing of this paper. Laser-induced contamination, laser-induced damage, and bulk darkening of the instrument's optics are the most probable causes, in addition to clipping losses at the instrument field stop.*

Specific comment #1.5:

*L121/L321: "Whereas" is a rather unusual conjunction to start a sentence with and sounds a bit awkward to me. Perhaps better use "while"?*

Response to Specific comment #1.5:

The text was changed accordingly.

Specific comment #1.6:

*L297: "the distribution is far from normal" sounds a bit funny to me. Perhaps better "far from Gaussian"?*

Response to Specific comment #1.6:

The text was changed accordingly.

Specific comment #1.7:

*L298: The value of 3.5 seems purely empirical. What was the decision criterion of Iglewicz and Hoaglin (1993), and is it likely to make it a good choice for your study as well?*

Response to Specific comment #1.7:

The motivation of *Iglewicz and Hoaglin (1993)* for using a threshold value of 3.5 to detect outliers was based on a simulation study where the portion of outliers (as identified by the modified Z-score) in a random normal distribution was determined in dependence on the sample size and threshold value. The results were based on 10,000 replications for each sample size and showed that the portion of observations that are labelled as outliers does not vary much with the sample size if a threshold value of 3.5 is chosen. The so derived outlier-labelling rule as they call it, however, still contains a certain degree of arbitrariness. Nevertheless, we found in our studies that threshold values ranging from 3.0 to 3.5 ensure wind error distributions with a high degree of normality, i.e. small residuals in the normal quantile plots, as stated in l. 404ff.:

> *For the model comparison of Rayleigh-clear winds, a Z-score limit ranging from 3.0 to 3.5 was found to yield a high degree of normality, which is in accordance with the recommendation by Iglewicz and Hoaglin (1993).*

We then decided to use a threshold value of 3.5 to allow more data points to pass the QC, thereby providing more robust statistics. Also, we have checked the variability of the statistical parameters within the threshold range from 3.0 to 3.5. The outcome is discussed in l. 477ff.:

> *For the model comparison discussed above, the statistical parameters change by less than 7% if the Z-score limit is reduced from 3.5 to 3.0. The largest influence is found for the Rayleigh standard deviation which decreases to 7.3 $m \cdot s^{-1}$ (compared to 7.8 $m \cdot s^{-1}$ for a Z-score limit of 3.5), as the portion of outliers accounts for 4.5% (compared to 3.2%).*

Specific comment #1.8:

*L335: remove the comma after "Gaussian distribution"*

Response to Specific comment #1.8:

Done.

Specific comment #1.9:

*L345: ECMWF are considered as absolute "truth" in the presented analyses. I'm wondering how likely it is that some of the supposed "gross errors" are actual rare extreme wind occurrences that were not modelled properly but captured correctly by the observations?*

<u>Response to Specific comment #1.9:</u>

It is true that, in addition to representativity errors, model deficiencies in terms of parametrization and resolution can cause large discrepancies between the model background winds and the Aeolus wind observations. As stated in l. 601ff., such model errors are mainly located in highly-convective regions, e.g. in the tropics, and can exceed 10 m·s$^{-1}$:

> *Moreover, discrepancies between the 2-µm DWL and model background wind data can result from model deficiencies that are caused by imperfect parametrization or too low resolution. Errors of the model background, i.e. before the assimilation of Aeolus winds, are found to be especially large in convective areas in the tropics, exceeding even 10 m·s$^{-1}$ on several occasions (Rennie et al., 2021).*

However, it is very unlikely that the modified Z-score filter will sort out these wind observations, as the threshold value is sufficiently relaxed. Taking the scaled MAD values from Table 2 before applying the Z-score filter ($k$ = 5.1 m·s$^{-1}$ for Mie, $k$ = 7.2 m·s$^{-1}$ for Rayleigh) and a Z-score threshold of 3.5, only wind results that deviate by more than ≈18 m·s$^{-1}$ (Mie) or ≈25 m·s$^{-1}$ (Rayleigh) from the model are considered as outliers (see also Fig. 2c,d). Hence, the risk of misinterpreting extreme wind occurrences as gross errors is very low.

The following paragraph was added to section 3.2 to clarify this point:

> *Note that the Z-score filter is sufficiently relaxed to avoid rejection of wind observations that might be due to extreme wind occurrences. Taking the scaled MAD values for the Mie and Rayleigh winds before applying the filter (5.1 and 7.2 m·s$^{-1}$, respectively) and considering that the median of the wind speed differences is close to zero, only wind results that deviate by more than ≈18 m·s$^{-1}$ (Mie) or ≈25 m·s$^{-1}$ (Rayleigh) from the model are considered as outliers. This is well above typical model errors that are caused by deficiencies in terms of parametrization and resolution, and can in extreme cases, like in highly-convection regions, exceed 10 m·s$^{-1}$ (Rennie et al., 2021).*

<u>Specific comment #1.10:</u>

*L352: Perhaps change to "it was \*assumed\* that"?*

<u>Response to Specific comment #1.10:</u>

The text was changed accordingly.

<u>Specific comment #1.11:</u>

*L386: Is 10 ms$^{-1}$ a commonly used threshold? Perhaps add a rationale / reference?*

Response to Specific comment #1.11:

As summarized in the introduction (l. 70ff.), the Rayleigh EE threshold typically ranges between 7 and 12.5 m·s$^{-1}$, with 8 m·s$^{-1}$ being the most commonly used value in many validation studies. We chose 10 m·s$^{-1}$ as an appropriate threshold given the fact that the EE of the Rayleigh winds had already significantly increased towards the AVATAR-T campaign (median of 4.5 to 5 m·s$^{-1}$, see Fig. 1) compared to the start of the FM-B period due to the signal loss and increased noise levels accordingly. The threshold also leads to portion of identified outliers (5.5%) that is comparable to that when the QC is solely based on the modified Z-score (3.2%), thereby providing better comparability between the two cases shown in Fig. 5.

The choice of the threshold value was specified in the revised manuscript as follows:

*Here, a threshold of 10 m·s$^{-1}$ was chosen, which is slightly larger than the most commonly used value of 8 m·s$^{-1}$, but regarded as an appropriate threshold given the increased noise levels during the campaign compared to the start of the FM-B and elevated EE of the Rayleigh winds accordingly (see Fig. 1).*

Specific comment #1.12:

*L429/Fig 6.: The graph contains a whole lot of information which is hard to grasp as a whole. Perhaps add a more high-level description of what the purpose of the plot is before describing the axes and lines, etc. specifically?*

Response to Specific comment #1.12:

Agreed. We have added the following sentence to section 3.5:

*The diagram intends to give an overview of the most relevant statistical parameters (μ, σ, k, fraction of valid winds that is considered in the statistics) in dependence on the chosen EE threshold.*

Specific comment #1.13:

*L499: Perhaps change to "1:1 line"?*

Response to Specific comment #1.13:

The term was changed accordingly.

Specific comment #1.14:

*L573: There are no orange bars. I assume this should be "black bars"?*

Response to Specific comment #1.14:

The text was changed accordingly.

Specific comment #1.15:

*L599: Replace "On the contrary" with "In contrast"?*

Response to Specific comment #1.15:

Done.

Specific comment #1.16:

*L639: The presented graph is a great way to summarize and visualize a lot of information, but I, personally, find it a bit exaggerated to call it "developing a new graph". I think it just draws the attention a bit away from the actually interesting part, which is the proposed systematic approach for selecting an EE threshold in a somewhat more objective manner.*
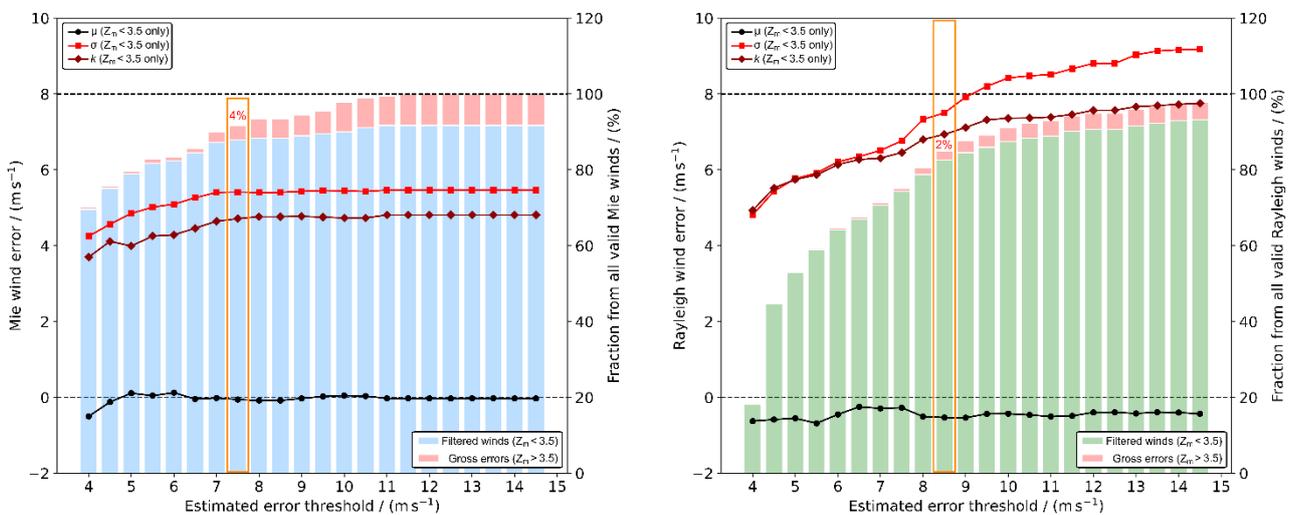
Response to Specific comment #1.16:

The text was changed to "a combined bar and line graph <u>can be used</u> to illustrate […]".

Specific comment #1.17:

*L654-657: I'm wondering how the statistics against ECMWF predictions would look like when using only the exact same locations as those in the AVATAR-T campaign? That is, how much of the mentioned deviations are just related to spatial representativeness?*

Response to Specific comment #1.17:

Very good point indeed. For answering this question, we produced another combined line and bar graph depicting the statistics for the model comparison after restriction to those Aeolus wind results that were also validated by the 2-µm DWL (see below).



Same as Fig. 6, but after restriction to those Aeolus wind results that were also validated by the 2-µm DWL in Fig. 8.

In contrast to Fig. 6, this plot is thus based on only the first five underflights when the DWL was still operable, i.e. with a number of data points that is comparable to the DWL comparison. Also, for the sake of clarity, it does not include the statistical parameters ($\mu$, $\sigma$, $k$) without applying the modified Z-score which are shown as dashed lines in Figs. 6 and 8. When comparing the new figure to Figs. 6 and 8, it becomes obvious that the statistical results are quite comparable to those of the model comparison of the entire dataset for the Mie winds. The mean bias is close to zero, while the standard deviation and scaled MAD converge toward 5.5 m·s$^{-1}$ and 4.8 m·s$^{-1}$, respectively, as the EE threshold is relaxed beyond 10 m·s$^{-1}$. These values are slightly higher than those obtained for the entire model dataset ($\sigma = 5.3$ m·s$^{-1}$, $k = 4.3$ m·s$^{-1}$) which might be due to the different spatial representativeness of the model and the 2-µm DWL or simply because of the restriction of the dataset. As for the Rayleigh-clear winds, the discrepancies are larger with a more negative mean bias of -0.4 m·s$^{-1}$ and significantly larger standard deviation and scaled MAD ($\sigma = 9.1$ m·s$^{-1}$ compared to $\sigma = 7.8$ m·s$^{-1}$ and $k = 7.7$ m·s$^{-1}$ compared to $k = 6.9$ m·s$^{-1}$) if no QC based on the EE is applied. When using the two-step QC approach with an EE threshold of 8.5 m·s$^{-1}$, the values are reduced to $\sigma = 7.5$ m·s$^{-1}$ and $k = 6.9$ m·s$^{-1}$, which is closer to the values obtained from the comparison against the 2-µm DWL (Fig. 8) ($\sigma = 8.2$ m·s$^{-1}$, $k = 7.2$ m·s$^{-1}$), but still lower. This result suggests that the spatial representativeness indeed contributes to the deviations in the validation results when comparing Figs. 6 and 8 in addition to the restricted overlap of the Aeolus data with the DWL winds. These aspects are covered in the discussion of the different overlap regions of the 2-µm DWL data coverage with the Mie and Rayleigh winds. Other factors that are mentioned in the text, e.g. the model errors, additionally contribute to the observed discrepancies.

Specific comment #1.18:

*L670-- I find the concluding remarks straying a bit off-topic.*

Response to Specific comment #1.18:

We have split the last section into two sections named "Discussion and summary" (sect. 5) and "Conclusion and outlook" (sect. 6) to better distinguish between the resume of the article and some remarks about the relevance of the results with a view to the Aeolus data assimilation and forthcoming studies on the Aeolus error characteristics. We think that the final remarks are important to highlight the impact of the paper, as it was also requested by the other two referees.

**Response to Referee Comment 2**

General comment #2.1:

*The abstract appears to be a bit too verbose and I would suggest to make it more focused on the key results and their implications.*

Response to General comment #2.1:

We agree that the abstract can be formulated more concise and have shortened it a bit accordingly:

*Since the start of the European Space Agency's Aeolus mission in 2018, various studies were dedicated to the evaluation of its wind data quality, and particularly to the determination of the systematic and random errors of the Rayleigh-clear and Mie-cloudy wind results provided in the Aeolus Level-2B (L2B) product. The quality control (QC) schemes applied in the analyses mostly rely on the estimated error (EE), reported in the L2B data, using different and often subjectively chosen thresholds for rejecting data outliers, thus hampering the comparability of different validation studies. This work gives insight into the calculation of the EE for the two receiver channels and reveals its limitations as a measure of the actual wind error due to its spatial and temporal variability. It is demonstrated that a precise error assessment of the Aeolus winds necessitates a careful statistical analysis, including a rigorous screening for gross errors to be compliant with the error definitions formulated in the Aeolus mission requirements. To this end, the modified Z-score and normal quantile plots are shown to be useful statistical tools for effectively eliminating gross errors and for evaluating the normality of the wind error distribution in dependence on the applied QC scheme, respectively. The influence of different QC approaches and thresholds on key statistical parameters is discussed in the context of the Joint Aeolus Tropical Atlantic Campaign (JATAC), which was conducted in Cabo Verde in September 2021. Aeolus winds are compared against model background data from the European Centre for Medium-range Weather Forecasts (ECMWF) before assimilation of Aeolus winds and against wind data measured with the 2-µm heterodyne-detection Doppler wind lidar (DWL) onboard the Falcon aircraft. The two studies make evident that the error distribution of the Mie-cloudy winds is strongly skewed with a preponderance of positively biased wind results distorting the statistics if not filtered out properly. Effective outlier removal is accomplished by applying a two-step QC based on the EE and the modified Z-score, thereby ensuring an error distribution with a high degree of normality while retaining a large portion of wind results from the original dataset. After utilization of the described QC approach, the systematic errors of the L2B Rayleigh-clear and Mie-cloudy winds are determined to be below $0.3\ m{\cdot}s^{-1}$ with respect to both the ECMWF model background and the 2-µm DWL. Differences in the random errors relative to the two reference datasets (Mie vs. model: $5.3\ m{\cdot}s^{-1}$, Mie vs. DWL: $4.1\ m{\cdot}s^{-1}$; Rayleigh vs. model: $7.8\ m{\cdot}s^{-1}$; Rayleigh vs. DWL: $8.2\ m{\cdot}s^{-1}$) are elaborated in the text.*

General comment #2.2:

*The discussion of the potential implications of the results presented is mostly missing. It would be useful to develop this aspect in the abstract and in the concluding section.*

Response to General comment #2.2:

Thank you for this comment which concurs with the remark from referee #3 (General comment #3.1). To address this point, we have extended the introduction to better describe the objective of the paper and the methodical approach used to achieve this (new text is underlined):

> *[...] Therefore, a more detailed treatment of different QC schemes and how they affect the resulting statistics is necessary for comparable validation results and for a more objective assessment of the Aeolus wind data quality. <u>Moreover, it is an important aspect with regards to the operational data assimilation in NWP centres and allows for a more rigorous error characterization of the Aeolus winds.</u>*
>
> *This paper aims to raise the awareness to the influence of the chosen QC schemes on the validation results, particularly when using the L2B EE. It also demonstrates the usefulness of specific statistical tools for the purpose of outlier removal and the assessment of normality, which are necessary to retrieve the Aeolus wind errors in accordance with the MRD. The presented methods are applied in the context of the AVATAR-T validation campaign in 2021 by comparisons against the ECMWF model background winds and the 2-µm DWL wind data. The specifics of the campaign and available datasets are outlined in Sect. 2 together with a description of the L2B Rayleigh-clear and Mie-cloudy EE and their temporal evolution over the past three years. <u>The model comparison in Sect. 3 serves as an example to introduce the reader to the detailed treatment of the Aeolus wind data in terms of QC and error assessment. In particular, the modified Z-score (Sect. 3.2) and normal quantile plots (Sect. 3.4) are discussed as powerful tools for removing gross errors and assessing the normality of the wind error distribution, respectively.</u> In addition, the impact of the QC settings on the results from the model comparison is elaborated (Sect. 3.5). In Sect. 4, the statistical methods are then applied to the comparison of Aeolus wind observations against 2-µm DWL data. The paper concludes with a summary and outlook to future studies of the L2B wind error characteristics in Sect. 5.*

Moreover, we have split the last section into two sections named "Discussion and summary" (Sect. 5) and "Conclusion and outlook" (Sect. 6). The latter includes a new paragraph to highlight the relevance of the presented results not only in terms of the validation of Aeolus wind data, but also with regards to its assimilation in NWP centres, along with a short description of the QC scheme that is used in the Aeolus data assimilation at the ECMWF.

*This work is intended to provide a guideline on how to perform a rigorous QC when working with Aeolus wind data. The presented results have demonstrated that a careful QC scheme is crucial for rejecting gross errors and, in turn, for providing an accurate estimation of the wind data quality. The shown statistical methods form the basis for a standardization and objectification of the Aeolus wind validation and will be applied in forthcoming studies involving DLR's wind lidar instruments. Furthermore, apart from the better comparability among different validation studies, the investigation fosters the analysis of the individual channel error characteristics and stimulates the refinement of the QC schemes that are currently used in the assimilation of Aeolus wind data into operational models. Both aspects are important to further improve the impact of the Aeolus products for NWP centres around the world.*

*In this context, it should be noted that the operational assimilation of Aeolus wind data at the ECMWF involves a multi-step QC scheme which also largely relies on the imperfect L2B EE. It comprises a first-guess check, which rejects observations with very large (O-B) departures ($5\sigma$), followed by the so-called variational QC (VarQC) method (Andersson and Järvinen, 1998). The VarQC assumes that the distribution of the normalized wind error, i.e., the (O-B) wind error divided by the assigned observation error, takes the form of a Gaussian function including an offset. The assigned observation error is proportional to the EE and additionally considers a representativeness error of 2 $m \cdot s^{-1}$ for the Mie winds (Rennie et al., 2021). Finally, there is a blacklist in the ECMWF assimilation which removes Rayleigh winds below 850 hPa pressure altitude as well as Rayleigh-clear and Mie-cloudy winds with EE larger than 12 and ~5 $m \cdot s^{-1}$, respectively. The multi-step approach ensures effective removal of the largest gross errors, but the VarQC assumption does not well represent the Aeolus normalized wind error distribution, especially for the Mie winds. In this regard, the use of the modified Z-score may help to improve the performance of the QC in the Aeolus data assimilation.*

Specific comment #2.1:

*The title of the last section should probably be Discussion and summary instead of Summary and conclusions.*

Response to Specific comment #2.1:

As described in the Response to General comment #2.2, we have split the last section into two sections named "Discussion and summary" (sect. 5) and "Conclusion and outlook" (sect. 6). The former section wraps up the results from the manuscript in accordance with its title.

Specific comment #2.2:

*L.573: should it be "Black bars"?*

Response to Specific comment #2.2:

Yes, we have changed the text accordingly.

**Response to Referee Comment 3**

General comment #3.1:

*I agree with the comment of RC-2 regarding the mostly missing discussion on the implications of the improved QC. While, in my mind, the paper presents results with high scientific value (understanding the instrument characteristics, improving the outcomes of validation and the outcomes of DA efforts), the importance of the improved QC are not strongly highlighted either in terms of motivation, or in terms of impact on the scientific or applications efforts (e.g. operational DA).*

Response to General comment #3.1:

Thank you for this advice which concurs with the comment from referee #2 (General comment #2.2). We have revised the last part of the introduction to elaborate on the motivation of our study and the used methodology (new text is underlined):

*[...] Therefore, a more detailed treatment of different QC schemes and how they affect the resulting statistics is necessary for comparable validation results and for a more objective assessment of the Aeolus wind data quality. Moreover, it is an important aspect with regards to the operational data assimilation in NWP centres and allows for a more rigorous error characterization of the Aeolus winds.*

*This paper aims to raise the awareness to the influence of the chosen QC schemes on the validation results, particularly when using the L2B EE. It also demonstrates the usefulness of specific statistical tools for the purpose of outlier removal and the assessment of normality, which are necessary to retrieve the Aeolus wind errors in accordance with the MRD. The presented methods are applied in the context of the AVATAR-T validation campaign in 2021 by comparisons against the ECMWF model background winds and the 2-μm DWL wind data. The specifics of the campaign and available datasets are outlined in Sect. 2 together with a description of the L2B Rayleigh-clear and Mie-cloudy EE and their temporal evolution over the past three years. The model comparison in Sect. 3 serves as an example to introduce the reader to the detailed treatment of the Aeolus wind data in terms of QC and error assessment. In particular, the modified Z-score (Sect. 3.2) and normal quantile plots (Sect. 3.4) are discussed as powerful tools for removing gross errors and assessing the normality of the wind error distribution, respectively. In addition, the impact of the QC settings on the results from the model comparison is elaborated (Sect. 3.5). In Sect. 4, the statistical methods are then applied to the comparison of Aeolus wind observations against 2-μm DWL data. The paper concludes with a summary and outlook to future studies of the L2B wind error characteristics in Sect. 5.*

Regarding the implications for the Aeolus wind data assimilation, it should be pointed out that the DA at the ECMWF involves a complicated QC which relies on a multi-step procedure to ensure effective outlier removal. To start with, there is a first-guess check of the (O-B) wind error with respect to the model background, rejecting Aeolus winds with departures greater than $5\sqrt{\sigma_O^2 + \sigma_B^2}$, where $\sigma_O$ is the assigned observation error (1.4 times the EE for the Rayleigh; 1.25 times the EE plus 2 m·s$^{-1}$ representativeness error for the Mie). $\sigma_B$ denotes the background error which is derived from Ensemble of Data Assimilation (EDA) statistics and can vary from about $\sim\sigma_O$ for Mie-cloudy to $\ll\sigma_O$ for Rayleigh-clear winds, so that the first-guess check typically discards wind data with deviations larger than $\sim 5\,\sigma_O$ to $5\sqrt{2}\,\sigma_O$. The first QC step is followed by the so-called variational QC (VarQC) method (Andersson and Järvinen, 1998) which assumes that the distribution of the normalised wind error, i.e., the (O-B) error divided by $\sigma_O$, takes the form of a Gaussian plus flat distribution (Gaussian function including an offset) to account for gross errors. The VarQC applies a weighting to the Aeolus observations depending on the normalized wind error. However, since for the current settings, observations are only down-weighted significantly for departures larger than 4.71 $\sigma_O$, in most cases the VarQC has only a small filtering effect in addition to the first-guess check. Finally, there is a blacklist in the ECMWF assimilation which removes Rayleigh-clear winds below 850 hPa pressure altitude as well as Rayleigh-clear and Mie-cloudy winds with EE larger than 12 and $\sim 5$ m·s$^{-1}$, respectively. Hence, this multi-step QC scheme used in the ECMWF DA has some similarities to the two-step QC approach described in the paper, as it combines a rather relaxed EE threshold with a filter that rejects winds with large departures from the model wind. However, as this multi-step QC scheme also largely relies on the imperfect EE and does not directly aim at a Gaussian wind error distribution, there is probably some room for improvement with regard to the DA in NWP.

Due to the complexity of the QC scheme used at the ECMWF, we refrained from elaborating on it in detail in the text. However, we have split the last section into two new sections "Discussion and summary" and "Conclusions and outlook", and have added a short paragraph to the latter including a brief description of the ECMWF QC to highlight the relevance of the presented results not only in terms of the validation of Aeolus wind data, but also with regards to its assimilation in NWP centres:

> *This work is intended to provide a guideline on how to perform a rigorous QC when working with Aeolus wind data. The presented results have demonstrated that a careful QC scheme is crucial for rejecting gross errors and, in turn, for providing an accurate estimation of the wind data quality. The shown statistical methods form the basis for a standardization and objectification of the Aeolus wind validation and will be applied in forthcoming studies involving DLR's wind lidar instruments.*

*Furthermore, apart from the better comparability among different validation studies, the investigation fosters the analysis of the individual channel error characteristics and stimulates the refinement of the QC schemes that are currently used in the assimilation of Aeolus wind data into operational models. Both aspects are important to further improve the impact of the Aeolus products for NWP centres around the world.*

*In this context, it should be noted that the operational assimilation of Aeolus wind data at the ECMWF involves a multi-step QC scheme which also largely relies on the imperfect L2B EE. It comprises a first-guess check, which rejects observations with very large (O-B) departures ($5\sigma$), followed by the so-called variational QC (VarQC) method (Andersson and Järvinen, 1998). The VarQC assumes that the distribution of the normalized wind error, i.e., the (O-B) wind error divided by the assigned observation error, takes the form of a Gaussian function including an offset. The assigned observation error is proportional to the EE and additionally considers a representativeness error of 2 $m \cdot s^{-1}$ for the Mie winds (Rennie et al., 2021). Finally, there is a blacklist in the ECMWF assimilation which removes Rayleigh winds below 850 hPa pressure altitude as well as Rayleigh-clear and Mie-cloudy winds with EE larger than 12 and ~5 $m \cdot s^{-1}$, respectively. The multi-step approach ensures effective removal of the largest gross errors, but the VarQC assumption does not well represent the Aeolus normalized wind error distribution, especially for the Mie winds. In this regard, the use of the modified Z-score may help to improve the performance of the QC in the Aeolus data assimilation.*

Specific comment #3.1:

*When using the ECMWF winds as the truth against which the Aeolus winds are compared, up to a 12-hour ECMWF forecast is used so that the model data are independent from the Aeolus winds (they have not been assimilated yet). However, Aeolus data have been assimilated in the previous model runs (cycles). What is the possible impact of the fact that the previous cycles have already assimilated the Aeolus winds?*

Response to Specific comment #3.1:

That's a very good question which also concerns the assessment of the impact of Aeolus wind retrievals, e.g. on ECMWF global weather forecasts, as discussed by Rennie et al. (2021). In this work, it is stated that since the start of the operational assimilation of Aeolus winds at ECMWF

> *"[...] the background departures are no longer independent of past Aeolus winds; it is unclear if this affected the Aeolus error estimates but there is no obvious discontinuity in the time series, so it probably did not." (Rennie et al, 2021).*

Given this statement, we assume that the influence of the assimilation of Aeolus winds in previous model runs on the (O-B) statistics is only minor. A more detailed investigation would be required to

verify this assumption, but goes beyond the scope of this paper. Aside from this, the fact that even model background is not entirely independent of the Aeolus wind observations emphasizes the relevance of the performed Cal/Val activities using ground-based and airborne instruments such as the 2-µm DWL.

Specific comment #3.2:

*The text on P. 23 that describes Fig. 10 has several inconsistencies with the figure - e.g. Figure 10 does not have orange bars; it is said that the 1-step QC statistics are given in a black inset while it is gray.*

Response to Specific comment #3.2:

We have corrected the text as follows:

> *Finally, the PDFs of the Mie and Rayleigh wind errors are presented in Fig. 10, indicating those wind results that are filtered out by the EE threshold (red bars) as well as those that are additionally filtered out by the modified Z-score (<u>black</u> bars). The statistical results that are provided in the boxes refer to the different subsets without QC (red), one-step QC using solely the EE threshold (<u>grey</u>) and two-step QC additionally applying the modified Z-score filter (blue/green).*

***Additional changes and corrections***

Changes to the manuscript #4.1:

We have extended the Acknowledgements section to thank Michael Rennie (ECMWF) and Jos de Kloe (KNMI) for providing insights to the calculation and temporal evolution of the Aeolus L2B estimated error as well as to the QC schemes that are used in the Aeolus data assimilation at ECMWF.

Changes to the manuscript #4.2:

The reference to the paper by Witschas et al. (2022b), which was published in AMTD, was updated.