

We thank the referees for their time reviewing our manuscript and their useful comments and feedback. Based on the reviewers' feedback, we have made several changes which we feel significantly improves the manuscript.

Below, reviewer comments are in **bold** while our responses are in regular type. Attached we have also provided a 'track changes' version of the manuscript, with added text in **blue** and deleted and/or moved text in **red**.

Comment by Anonymous Referee #1

General Comments

Overall, the paper is well written, and makes some important points regarding the limitations of simple performance metrics and the need for more intensive investigations of measurement error. I agree with the conclusions of the paper in principle, but I think that the paper could better support these conclusions through its examples.

Response 1: Thank you very much for the positive comments, and we appreciate the suggestions to better support our arguments. In the following paragraphs we have addressed this through a more detailed explanation of the limitations and advantages of common metrics, as well as revised figures to better illustrate and support the conclusions made.

As written, besides the contrived example of Figure 2, I don't see a clear case where differences and deficiencies in the measurements are not at least hinted at through relatively worse simple performance metric values (R-squared, RMSE). As shown in several examples of Section 4, the simple performance metrics do have utility in allowing comparisons between alternative measurement devices or techniques; for the most part, the sensors showing difficulties in the B-A and REU plots also showed relatively worse performance metric values. This is partly due to how these results are presented; they mostly compare data from common collocation experiments at the same location and covering the same time period, and therefore represent situations where it would be more appropriate to compare simple performance metrics (this is well stated in lines 394-396). The exception is Figure 6, but in that case, it isn't clear that the B-A or REU plots show any more ability to anticipate the poor observed performance at the roadside site than the simple performance metrics; rather, this is a general issue of relying on single-site collocation studies for characterization. More directly relevant to the topic of the paper would be to show attempts to compare different collocation datasets using simple metrics only, and to illustrate the shortcomings

of that approach; these shortcomings can then be addressed through the approaches you suggest. Perhaps such an example might be constructed from the existing data you present in the paper. For example, a collocation dataset could be divided across time by taking data collected in different seasons (if possible) and treating these as separate collocation experiments. In different seasons, the same sensor could have different performance metrics due to the differences in concentrations and variability in environmental conditions between seasons. These differences and their effects on errors likely would be much more apparent in B-A or REU plots (e.g., the collocation data would span different sections on the horizontal axis). Therefore, the information on error characteristics from each collocation analyzed via B-A or REU plots would tend to complement each other, as opposed to the simple performance metrics which might seemingly contradict each other. This is just a thought; while in general I agree with the logic underlying the arguments being presented here, I don't think that the examples, as they are currently presented, do a strong enough job of backing up these arguments.

Response 2: We thank the reviewer for this very useful comment. We agree with the reviewer that simple performance metrics do have utility. However, we argue that although these metrics are useful as a quick assessment or sanity check of performance, approaches such as the BA/REU plots enable a potential end user to view the nature of the errors and thus assess how these errors will impact any end use application. We appreciate that we potentially did not make this clear enough in the initial manuscript, and have therefore expanded the example given in the introduction (lines 48-53) and we have added the following text to "2.1 Performance indices, error structure and uncertainty" (lines 125-139):

When evaluating multiple sensors during a collocation experiment, single metrics can be a useful way to globally compare instruments/sensors. However, these metrics do little to communicate the nature of the measurement errors and the impacts these will have in any end use application, in part because they reduce the error down to a single value (Tian et al., 2016). Even more if a specific concentration range is of paramount interest to the end-user, these metrics are not capable of characterising the weight of noise and/or the bias effect. The R^2 shows globally the data set linearity and gives an idea of the measurement noise. However, it is unable to distinguish whether a specific range of concentrations is more or less linear (or more or less noisy) than another. Similarly, the RMSE is also a very useful metric and perhaps more complete than R^2 , as it considers both noise and bias (although they need to be explicitly decomposed from RMSE). Nevertheless, the RMSE is an average measure (of noise and bias) over the entire dataset under analysis. Using combinations of simple metrics increases the information communicated, but does not necessarily make it easy to assess how the errors will likely impact a particular measurement application. Visualising the absolute and relative

measurement errors across the concentration range (unreachable by global metrics) enables end users to view the errors, and any features (non-linearities, step changes, etc.) that would impact the measurement but that global metrics (and in some cases time-series and/or regression plots) are incapable of showing.

Complementary to the text added to the manuscript we have taken the reviewers suggestion of better using the QUANT dataset to support our arguments. In order to explicitly demonstrate the advantages of the BA and REU plots we have updated figure 5 using data from sensors with different error characteristics. We have also updated figure 6, as suggested by the reviewer, to present data from the same sensor but during different periods, in addition to data from an identical sensor at a different location. The revised Figures 5 and 6, along with the associated edited text, are shown below. We feel these plots much better support the arguments made in the paper, and want to thank the reviewer again for suggesting this. We have also updated all the figures to show the density of data points, as we feel this further increases the information communicated:

Figure 5 shows two colocated measurements from two different LCS devices: one measuring NO₂ (a-panels) and the other O₃ (b-panels). Both measurements are compared with colocated reference measurements at an urban background site in the city of Manchester. Unlike the modelled instruments in Sect. 4.1, the combination of error sources is unknown in this case and we can thus only assess the LCS measurement performance through comparison with the reference measurements using metrics and visual tools.

Single value metrics indicate an acceptable performance for both measurements: high linearity (both R^2 are higher than 0.8) and relatively low errors (RMSE ~ 5ppb). However, the plots present the data in a variety of ways that enable the user to identify patterns in the measurement errors that would be less obvious if only global metrics were used. For example, the NO₂ sensor (LCS1, a-panels) has a non-linear response that is almost imperceptible from the regression plot but stands out in the B-A plot. Furthermore (despite the high R^2 and relatively low RMSE), the REU plot shows high relative errors that do not meet the Class 2 DQO for the measured concentration range. Regarding the O₃ sensor (LCS2, b-panels), the B-A plot shows two high density measurement clusters, one with positive absolute errors (over-measuring) and a larger one with negative errors (under-measuring). These are the result of a step change in the correction algorithm applied by the manufacturer and could easily have been missed if only summary metrics and a regression plot were used, especially if the density of the data points was not coloured.

It is worth noting that these plots do not directly identify the source of the proportional bias, with sensor response to the target compound or another covarying compound possible, but provides information on how much it impacts the data.

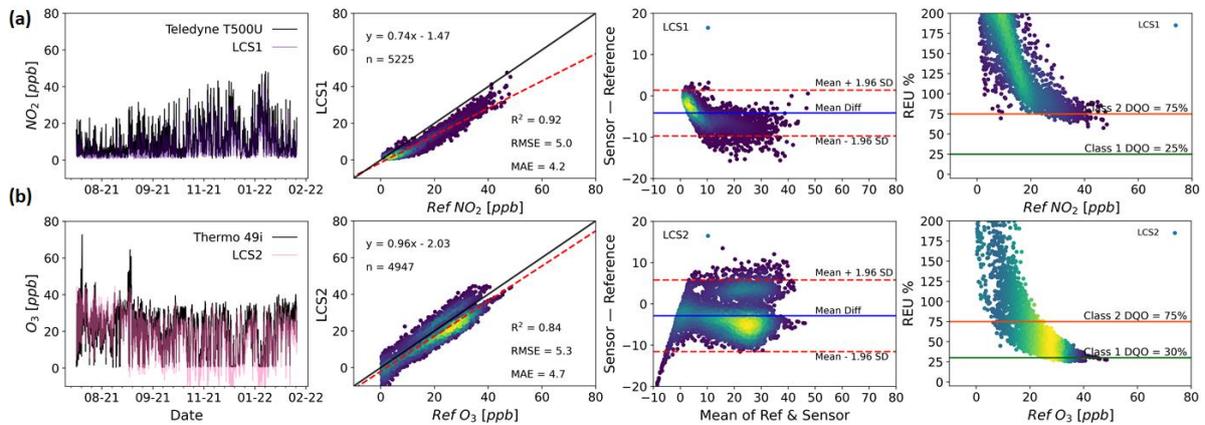
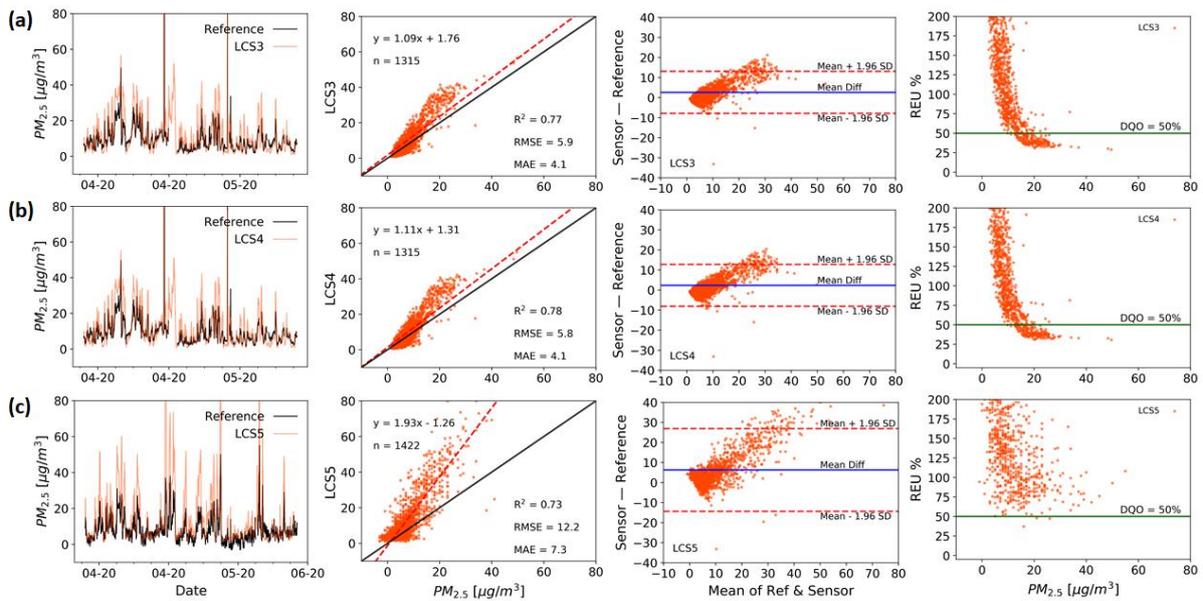


Figure 5. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels; NO₂ Class 1 DQO = 25% & Class 2 DQO = 75% ; O₃ Class 1 DQO = 30% & Class 2 DQO = 75%) for NO₂ (a-panels) and O₃ (b-panels) measurements by two LCS systems of different brands in the same location and time span (Manchester Supersite, July 2021 to February 2022. Time res 1 h). All but the time-series plots, have coloured by data density.

We have also replaced the “old” figure 6 for the one below, in which we present two of the same sensors as previously, but now LCS3 is shown for two different periods: panel a, from Apr to May 2020; panel b, Oct to Nov 2020 (exactly 6 months after the initial period). For LCS4 (panel c) the period is also Apr to May 2020:



Old figure 6.

Figure 6 shows three out-of-the-box PM_{2.5} measurements made by two devices (LCS3 & LCS4) from the same brand in spring (LCS3: a-panels; LCS4: c-panels) and in autumn (b-panels, only LCS3). The collocation shown correspond to two different sites: an urban background site (LCS3, a and c-panels) and a roadside site (LCS4, c-panels).

As the regression and the B-A plots show, all LCS measurements in Fig. 6 have a proportional bias compared with the reference, with the LCS over predicting the reference values. The device at the urban background site (LCS3) show a dissimilar performance in spring and autumn, indicating that the errors this device suffers are differently influenced by local conditions in the two seasons (all the duplicates at the urban background show the same pattern). While for LCS3 during spring the error have a more linear behaviour, in autumn a non-linear pattern is clearly observed in the regression and B-A plots. Despite the utility that single metrics can have in certain circumstances, the non-linear pattern goes completely unnoticed by them: while for the two different seasons RMSE and the MAE are almost constant the R² indicates a higher linearity for autumn.

A number of duplicates were deployed at both sites showing a very similar performance in terms of the single metric values but also in regard to the more visual tools (not shown here). This internal consistency is highly desirable, especially when LCS's are to be deployed in networks, as although mean absolute measurement error may be high, differences between identical devices are likely to be interpretable.

Having prior knowledge of the nature of the measurement errors allows informed experimental design prior to data collection. This is key if an end user is to maximise the power of a dataset, and the information it provides, to answer a specific question. For example, if an end-user wanted to identify pollution hotspots within a relatively small geographical area, then using a dense network of sensor devices that possess errors large and variable enough to make quantitative comparisons with limit values difficult (possibly due to an interference from a physical parameter like relative humidity) but show internal consistency could be a viable option. Providing the hotspot signal is large enough relative to any random error magnitude.

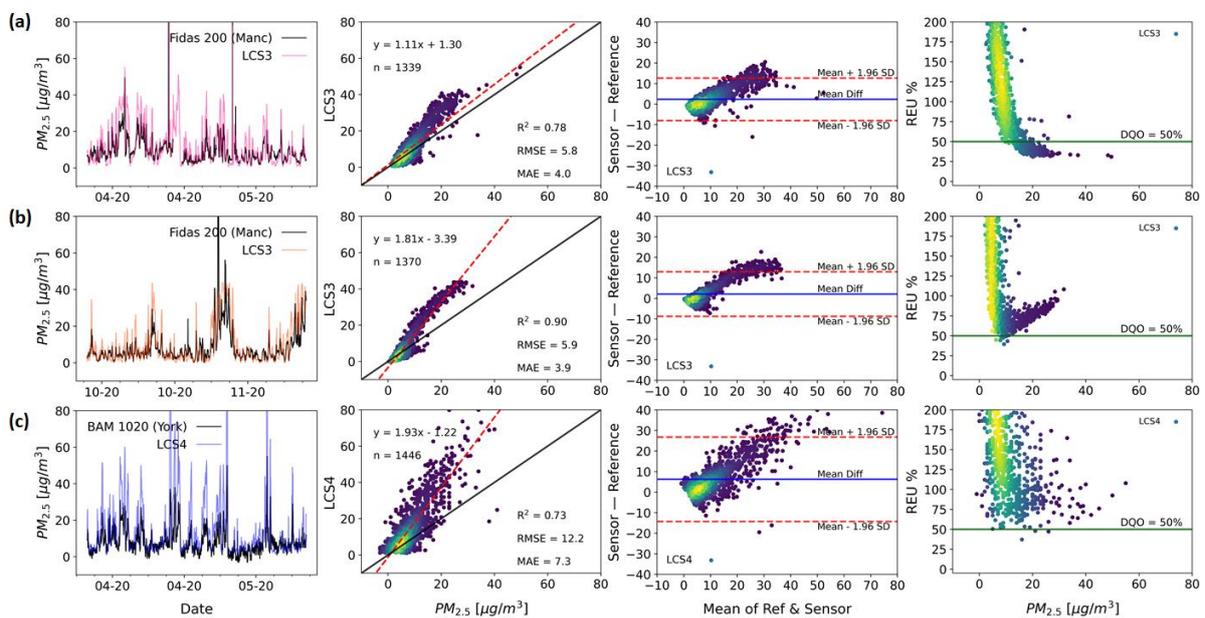


Figure 6. Two LCS systems (LCS3 & LCS4, same brand) measuring PM_{2.5} (Time res 1 h). While LCS3 is shown for the same location (Manchester) but unfolded in two different seasons (a-panels: Apr to May 2020; b-panels: Oct to Nov 2020), LCS4 is at a different location (c-panels: York, Apr to May 2020). Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels; DQOPM_{2.5} = 50%) are used to characterise the device's error structure. All but the time-series plots have been coloured by data density.

Specific Comments

Line 75: One of the commas here seems misplaced.

Response 3: Corrected.

Lines 78-79: Might be better stated as “a linear additive model is often assumed”.

Response 4: Corrected.

Figure 1: REU should be defined before it is used in this figure.

Response 5: Corrected.

Line 103: Should be “data are communicated”.

Response 6: Corrected.

Line 136: Remove “And” at start of sentence.

Response 7: Corrected.

Line 174: Suggest replacing “data” with “data set”.

Response 8: Corrected.

Lines 306-310: This is background information, better included as part of the introduction, where it can be integrated with similar statements already there.

Response 9: We have preferred keeping that sentence as originally set in “4.2 Real-world instruments”, but we have added a paragraph to the introduction expanding on this important point, where now it can be read (lines 72-78):

The covariance of many of the physical and chemical parameters of the atmosphere, makes accurately identifying particular sources of measurement interference or error very difficult in the real world. Unfortunately, specific laboratory experiments for the characterization of errors is complex and very expensive, resulting in many sources of error being essentially unknown for many measurement devices. The use of imperfect error correction algorithms that are not available to the end-user (e.g. in many LCS devices) makes error identification and quantification even more complex. For this reason, colocation experiments in relevant environments are often the best option to assess the applicability of a given measurement method for its intended purpose.

Lines 342-344: This is an important point, often used as justification for the use of LCS for applications like hotspot identification. I wonder if the authors could comment more on this, either here or elsewhere. My prompting question would be: what kinds of analysis approaches could be used to verify the ability of LCS to qualitatively identify meaningful differences between measurements, even in situations where relative uncertainties are too high to make reliable quantitative comparisons? Alternatively, is such a distinction (qualitative versus quantitative analysis) meaningful here, or is this “qualitative analysis” merely a quantitative analysis performed under higher relative uncertainty.

Response 10: The discussion of analysis approaches for specific applications is beyond the scope of this work, and would likely be best supported through a number of case studies. In this work we focus on tools that aid the interpretation of performance data in order to inform measurement strategies. We agree with the reviewer that this is an important point, and we are of the opinion that, when it comes to air quality measurements, qualitative analysis is merely quantitative analysis performed under higher uncertainty. All the devices discussed in this work report values for target pollutants, and as such are quantitative. Understanding the impact of likely measurement errors on the power of the data to answer specific questions (e.g. hotspot identification) is important for all devices, not just LCS. Especially as criteria pollutant limit values continue to decrease based on revised health evidence. We therefore argue that more emphasis should be placed on informed experimental design when making the measurements than on analysis methods that attempt to extract signals from data with uncharacterised errors. In the case that LCS devices show high levels of internal consistency, an informed experimental design should be able to take advantage of this to minimise the impact of measurement errors on the information gathered from the measurements. We thank the reviewer for highlighting that we have not said this explicitly in the text, and have added the text below to the manuscript:

To what have been said in the original lines:

“This internal consistency is highly desirable, especially when LCS’s are to be deployed in networks, as although mean absolute measurement error may be high, differences between identical devices are likely to be interpretable.”

We have added the following to the text (lines 440-446):

Having prior knowledge of the nature of the measurement errors allows informed experimental design prior to data collection. This is key if an end user is to maximise the power of a dataset, and the information it provides, to answer a specific question. For example, if an end-user wanted to identify pollution hotspots within a relatively small geographical area, then using a dense network of sensor devices that posses errors large and variable enough to make quantitative comparisons with limit values difficult (possibly due to an interference from a physical parameter like relative humidity) but show internal consistency could be a viable option. Providing the hotspot signal is large enough relative to any random error magnitude.

Line 367: Second “at” is superfluous.

Response 11: Corrected.

Line 370: “deviations” should be “deviation”.

Response 12: Corrected.

Line 372: “appears” should be “appear”.

Response 13: Corrected.

Lines 376-378: This is another important point. Since air quality regulations are based on these reference instruments, the traceability of LCS to these reference instruments has been a major focus of work. However, we must acknowledge that these references themselves are imperfect. Is it thus inappropriate to hold LCS to certain performance standards which the reference instruments themselves may not meet (especially if improperly operated)? On the other hand, what is the alternative to ensuring data quality? I think that, as you suggest, comparing different reference instruments among themselves should be done more frequently, and these intercomparisons more widely used as a benchmark against which the performance of LCS can be judged (instead of establishing arbitrary performance metric targets, especially if these targets are not connected in some way to the different conditions under which the sensors are expected to operate). However, there is of course the practical question of the cost and feasibility of doing this at the necessary scale. Generally speaking, this is a major point which could be explored further by the authors either here or elsewhere.

Response 14: We thank the reviewer for this comment. We also feel that this is an important point and more attention needs to be placed on measurement uncertainty across the field of air pollution measurements, not just low cost sensors. Especially as limit values continue to fall. A more detailed discussion of this point and its implications is outside the scope of this work, but this is something we plan to expand on in the future and are in the process of collocating reference instruments for this purpose.

Line 443: “data is” should be “data are”.

Response 15: Corrected.

Lines 450-453: The meaning of this sentence is unclear; consider breaking it into several simpler sentences.

Response 16: Corrected.

Previously it was said:

“If end users are to be able to make well informed decisions about device applicability to a particular task, then an argument can be made for information on the scale of the error corrections made to a

reported measurement to be made available, ideally alongside and a demonstration of its benefits in a relevant environment.”

Now it can be read (lines 571-574):

If end users are to be able to make well informed decisions about device applicability then information on the scale of the measurement errors, and the impact of corrections made to minimise these, should be made available. Exemplar case studies in a range of relevant environments would also be highly valuable.

Comments by Anonymous Referee #2

This is a well written paper on air sensor uncertainty. Uncertainty in air sensors is a very important topic in the field. While the authors lay out a number of issues with current uncertainty methods their method seems to make only minor improvements on current methods. This paper is still helpful as it provides another way to visualize similar information in different ways which may speak more clearly to some people. I have a number of specific comments below that I hope the authors will address to improve the strength of the paper.

Response 1: We thank the reviewer for their comments. We would like to take this opportunity to clarify that our intention is not to replace the commonly used performance evaluation methods (R^2 , RMSE, MAE, etc.). On the contrary, we think that they are useful tools, but like any tool they have certain deficiencies, and the end user needs to be aware of this. We have added more clarification on this point in a response to the first reviewer (please see Response #2 to reviewer 1).

Yes, these plots you are proposing may be more helpful than just R^2 , MAE, and RMSE but typically I'm seeing those metrics reported along with slope and intercept (and often a scatter plot). This seems like a false comparison you talk about repeatedly in the paper. Slope, intercept, and R^2 seem to provide much of the same info as BA or REU plots just in a different form.

Response 2: We agree with the reviewer that using combinations of simple metrics in conjunction with a regression plot provides significantly more information than any single value metric. However, we argue that viewing the errors directly, as in the BA and REU plots, provides a clearer picture of the nature of the errors and thus how they would likely impact any application of the measurement device. Often in air pollution measurement applications there are specific target concentrations where the data is of most interest, for example around a legal limit value. Single value metrics give a global picture of a data set, but do not describe the error distribution in specific ranges or concentration intervals. The use of visualisations such as B-A and REU is complementary to the aforementioned metrics, with the added value that the user is now more aware of how the data looks in an absolute and/or relative error space, allowing them to distinguish some characteristics of interest.

In order to clarify this we have expanded the original text and now in lines 228-235 it can be read:

On the other hand, the use of visualisations such as B-A and REU is complementary to the aforementioned metrics, with the added value that the user is now more aware of how the data looks like in an absolute and/or relative error space, allowing them to distinguish some characteristics of interest. These visualizations are indeed more laborious and the interpretation can be challenging for non-experts, but they provide additional insights into the nature of the errors, not attainable by one

or more combined performance metrics: while B-A plots shows the noise (dispersion of the data) and the bias effect (tendency of the data) in an absolute scale, the REU can be explicitly decomposed in the noise and bias components (see Yatkin et al., 2022).

We admit that we did not use the best examples from our dataset to support these points, and have updated Figures 5 and 6 to better highlight the strengths of these approaches over just using global metrics (see Response #2 to Reviewer 1).

The BA plot seems to be just a less intuitive form of a scatterplot but maybe I'm missing how to interpret it in a helpful way? I see that there is value though in visualizing things in different ways since people see things differently.

Response 3: We agree with the reviewer that the scatter or correlation plot and BA plot are similar to a point, and much of the same information can be extracted. However, we argue that the BA plot is better placed to evaluate the agreement between two different methods for measuring the same variable than a scatter plot. As both methods being compared are in theory measuring the same parameter, but with different measurement errors, then it is to be expected that the two measurements should have good correlation when sampling over a wide range of parameter values. A high correlation (high R^2), however, does not necessarily imply good agreement between the two measurements. It is also not always possible in atmospheric collocation studies to guarantee that a sufficiently wide range of parameter values will be observed.

In order to clarify the information that a Bland-Altman plot is capable of provide we have re-written and expanded the ideas originally set in lines 253-262:

In contrast to the regression plot -where the measured values from the two measurements (e.g. LCS vs Ref) are plotted against each other- the Bland-Altman plots essentially display the difference between measurements (abscissa) as a function of the average measurement (ordinate), enabling more information on the nature of the error to be communicated. This direct visualisation of the absolute error acknowledges that the true value is unknown and that both measurements have errors. The B-A plot enables the easy identification of any systematic bias between the measurements or possible outliers, and is the reason B-A plots are extensively used in analytical chemistry and biomedicine to evaluate agreement between measurement methods (Doğan, 2018). The mean difference between the measurements, (represented by the blue line in the figures), is the estimated bias between the two observations. The spread of error values around this average line indicates if the error shows purely random fluctuations around this mean, or if it has structure across the observed concentration range.

We therefore feel that the BA plot is better placed to show particular features or characteristics of the error than a scatter plot. It is hoped that the updated Figure 5 (see Response #2 to Reviewer 1) and the accompanying text now illustrates this more clearly.

In the end the plots you've made reveal very little about temperature, RH, and other pollutant interferent biases. Is there any way to modify the plots you are proposing to make them more helpful in addressing the issues you've brought up about interferences and error?

Response 4: We thank the reviewer for this comment. Yes, it would be possible to plot the absolute or relative error against variables other than those in the BA and REU plots in order to investigate correlations of error with other variables (e.g. relative humidity). Although these would no longer be BA or REU plots in the established definition, they could prove insightful in demonstrating an error correlation. Unfortunately, in the situation where the error sources are not known or fully understood (e.g. LCS) these plots would only be able to show correlations between error and variables, not diagnosing error causes. This is because in the real atmosphere there are a vast number of covarying physical and chemical parameters, making it very difficult to prove causation from correlation. We strongly feel that the best way in which to identify interference biases is through controlled laboratory experiments, where confounding variables can be controlled. However, as we mention in the text, this can also be very difficult due to real atmospheric complexity, and plots such as those discussed could be insightful in at least confirming the apparent presence or not of a bias correlation. Although a thorough demonstration of this is outside of the scope of this paper, and the focus of a future piece of work, we have added the following text to the manuscript discussion (lines 540-546, see the added text in blue font):

A major challenge comes from complex errors, such as interferences from other compounds or with environmental factors, that vary temporally and/or spatially. Similar graphical techniques to those presented above can be used to identify the existence of such relationships, but correcting for them remains a challenge. *For example, the correlation between measurement errors and relative humidity could be explored by replacing the abscissa with measured relative humidity in both the B-A and REU plots. This would visualise the relationship between absolute and relative errors with relative humidity, but would not be able to confirm causality. The complex and covarying nature of the atmosphere means that the best way to identify a device error source is through controlled laboratory experiments, where confounding variables can be controlled, although these experiments are often difficult and expensive to perform in a relevant way.*

Are all the DQOs for REU just percentages? There is no absolute target? (e.g. 25% or 5 ppb)

Response 5: Yes, the DQO (Data Quality Objective) is a defined percentage (defined by of the European Air Quality Directive 2008/50/EC) but for regulatory purposes it needs to be evaluated at a fixed concentration (called Limit Value, see GDE (2010) for its complete definition). However, as the main message is focused on end-user needs and not necessarily for regulatory purposes, we have preferred minimising the use of this as a performance target and instead the DQO is used just as a reference line in the REU plots.

Did you consider how your estimation of uncertainty compares to the method in this recent paper?

<https://amt.copernicus.org/articles/14/7369/2021/>

Response 6: We thank the reviewer for this suggestion. This particular work considers uncertainty from a prognostic perspective, in contrast to the diagnostic uncertainty discussed in our paper. To make this distinction more explicit, we have added the following text in the introduction (lines 69-71):

Also, when the term “uncertainty” is used here, it is referring to “diagnosis uncertainty”, in contrast with “prognosis uncertainty” (see Sayer et al., 2020 for more details).

You may want to add a reference to the EPA performance targets. They recommend slope, intercept, R2, RMSE, along with precision metrics, and making plots looking at error vs. T/RH/Dewpoint. Base testing is colocations with enhanced lab testing.

Response 7: Thanks for the suggestions. Yes, we are aware of the EPA performance targets but for simplicity we have tried to minimise the use of performance targets in the manuscript to not to deviate our main message which is focused on the questions end-users may want to answer with the provided data.

Can you include equations for RMSE, MAE, and REU.

Response 8: These metrics are widely used and we feel that including these are unnecessary. In addition the REU derivation is not trivial, and would require a significant addition to the paper, which we feel would detract from the paper. We would therefore prefer not to add the equations, as they are defined elsewhere, but leave the decision to the editor. Instead, we have added the references below to where to find the equations and definitions, and have also added a Zenodo and a github link (lines 626-627) to an open source repository of python and R code to generate the plots, with example data, enabling readers to take advantage of these tools.

For the RMSE and MAE equations (line 221): [see the equation definitions in Cordero et al., 2018](#)

For the REU equations (line 106): [as defined by the GDE \(2010\)](#)

Figures 1 and 2: These are really nice illustrations!

Response 9: Many thanks! We are glad you liked it.

Figure 1. can you define the acronyms (e.g. REU) that you haven't defined in the text yet.

Response 10: Corrected.

Line 353 do you need "roadside side" or just "roadside"

Response 11: Corrected.

Line 367: "at during" only need "during"

Response 12: Corrected.

1 **Air pollution measurement errors: Is your data fit for** 2 **purpose?**

3 Sebastian Diez¹, Stuart E. Lacy¹, Thomas J. Bannan², Michael Flynn², Tom Gardiner³, David
4 Harrison⁴, Nicholas Marsden², Nick Martin³, Katie Read^{1,5}, Pete M. Edwards¹

5 ¹Wolfson Atmospheric Chemistry Laboratories, University of York, York, YO10 5DD, UK

6 ²Department of Earth and Environmental Science, Centre for Atmospheric Science, School of Natural Sciences,
7 The University of Manchester, Manchester, M13 9PL, UK

8 ³National Physical Laboratory, Teddington TW11 0LW, UK

9 ⁴Bureau Veritas UK, London, E1 8HG, UK

10 ⁵National Centre for Atmospheric Science, University of York, York, YO10 5DD, UK

11 *Correspondence:*

12 Sebastian Diez (sebastian.diez@york.ac.uk); Pete Edwards (pete.edwards@york.ac.uk)

13 **Abstract.** When making measurements of air quality, having a reliable estimate of the measurement uncertainty
14 is key to assessing the information content that an instrument is capable of providing, and thus its usefulness in a
15 particular application. This is especially important given the widespread emergence of Low Cost Sensors (LCS)
16 to measure air quality. To do this, end users need to clearly identify the data requirements a priori and design
17 quantifiable success criteria by which to judge the data. All measurements suffer from errors, with the degree to
18 which these impact the accuracy of the final data often determined by our ability to identify and correct for them.
19 The advent of LCS has provided a challenge in that many error sources show high spatial and temporal variability,
20 making laboratory derived corrections difficult. Characterising LCS performance thus currently depends primarily
21 on colocation studies with reference instruments, which are very expensive and do not offer a definitive solution
22 but rather a glimpse of LCS performance in specific conditions over a limited period of time. Despite the
23 limitations, colocation studies do provide useful information on measurement device error structure, but the results
24 are non-trivial to interpret and often difficult to extrapolate to future device performance. A problem that obscures
25 much of the information content of these colocation performance assessments is the exacerbated use of global
26 performance metrics (R^2 , RMSE, MAE, etc.). Colocation studies are complex and time-consuming, and it is easy
27 to fall into the temptation to only use these metrics when trying to define the most appropriate sensor technology
28 to subsequently use. But the use of these metrics can be limited, and even misleading, restricting our understanding
29 of the error structure and therefore the measurements' information content. In this work, the nature of common
30 air pollution measurement errors is investigated, and the implications these have on traditional metrics and other
31 empirical, potentially more insightful, approaches to assess measurement performance. With this insight we
32 demonstrate the impact these errors can have on measurements, using a selection of LCS deployed alongside
33 reference measurements as part of the QUANT project, and discuss the implications this has on device end-use.

34

35 1. Introduction

36 The measurement of air pollutants is central to our ability to both devise and assess the effectiveness of policies
37 to improve air quality and reduce human exposure (Molina & Molina, 2004). The emergence of low-cost sensor
38 (LCS) based technologies means a growing number of measurement devices are now available for this purpose
39 (Morawska et al., 2018), ranging from small low-cost devices that can be carried on an individual's person all the
40 way through to large, expensive reference and research-grade instrumentation. A key question that needs to be
41 asked when choosing a particular measurement technology is whether the data provided is fit for purpose
42 (Andrewes et al., 2021; Lewis & Edwards, 2016). In order to answer this, the user must first clearly define the
43 question that is to be asked of the data, and thus the information required. For example, a measurement to
44 characterize “rush hour” concentrations, or to determine if the concentration of a pollutant exceeded an 8 h average
45 legal threshold value at a particular location would demand a very different set of data requirements than a
46 measurement to determine if a change in policy had modified the average pollutant concentration trend in a
47 neighbourhood. ~~Considerations such as measurement time resolution and ability to capture spatial variability~~
48 ~~would be important for such examples (Feinberg et al., 2019).~~ Would the R^2 or RMSE or any other global single-
49 value metric be enough to decide between the different device’s options? Considerations such as the origin of the
50 performance data, type of experiment (laboratory or colocation) (Jiao et al., 2016), the test location (Feenstra et
51 al., 2019) and period (i.e. duration, season, etc.), the LCS and reference measurement method (Giordano et al.,
52 2021), measurement time resolution and ability to capture spatial variability (Feinberg et al., 2019) would be
53 important factors to consider for such examples. The measurement uncertainty is also of critical consideration, as
54 this ultimately determines the information content of the data, and hence how it can be used (Tian et al., 2016).

55 All measurements have an associated uncertainty, and even in highly controlled laboratory assessments, the true
56 value is not known, with any measurement error defined relative to our best estimate of the range of possible true
57 values. However, quantifying and representing error and uncertainty is a challenge for a wide range of analytical
58 fields, and often what these concepts represent is not the same to all practitioners. This results in a spectrum of
59 definitions that take into account the way truth, error, and uncertainty are conceived (Grégis, 2019; Kirkham et
60 al., 2018; Mari et al., 2021). For atmospheric measurements assessing uncertainty is complex and non-trivial.
61 Firstly, given the “true” value can never be known, an agreed reference is needed. Secondly, the constantly
62 changing atmospheric composition means that repeat measurements cannot be made and the traditional methods
63 for determining the random uncertainty are not applicable. And finally, a major challenge arises from the multiple
64 sources of error both internal and external to the sensor that can affect a measurement. Signal responses from a
65 non-target chemical or physical parameter or electromagnetic interference are examples of an almost limitless
66 number of potential sources of measurement error. In this work, we will follow the definitions given by the
67 International Vocabulary of Metrology (JCGM, 2012) for measurement error (“measured quantity value minus a
68 reference quantity value”) and for measurement uncertainty (“non-negative parameter characterising the
69 dispersion of the quantity values being attributed to a measurand, based on the information used”). Also, when
70 the term “uncertainty” is used here, it is referring to “diagnosis uncertainty”, in contrast with “prognosis
71 uncertainty” (see Sayer et al., 2020 for more details).

72 The covariance of many of the physical and chemical parameters of the atmosphere, makes accurately identifying
73 particular sources of measurement interference or error very difficult in the real world. Unfortunately, specific
74 laboratory experiments for the characterization of errors is complex and very expensive, resulting in many sources
75 of error being essentially unknown for many measurement devices. The use of imperfect error correction
76 algorithms that are not available to the end-user (e.g. in many LCS devices) makes error identification and
77 quantification even more complex. For this reason, colocation experiments in relevant environments are often the
78 best option to assess the applicability of a given measurement method for its intended purpose.

79 The mentioned difficulties in defining and quantifying uncertainty across the full range of end-use applications of
80 a measurement device, means that often the quoted measurement uncertainty is not applicable, or in some cases
81 not provided or provided in an ambiguous manner. This makes assessing the applicability of a measurement device
82 to a particular task difficult for users. In this work, we investigate the nature of common air pollution measurement
83 errors, and the implications these have on traditional goodness-of-fit metrics and other, potentially more insightful
84 approaches to assess measurement uncertainty. We then use this insight to demonstrate the impact these errors
85 can have on measurements, using a selection of LCS deployed alongside reference measurements as part of the
86 UK Clean Air program funded QUANT (Quantification of Utility of Atmospheric Network Technologies) project,
87 a 2-year colocation study of 26 commercial LCS devices (56 gases measurements and 56 PM measurements) at
88 multiple urban, background and roadside locations in the UK. After analysing some of the real-life uncertainty
89 characteristics we discuss the implications this has on data use.

90 2. Error characterization

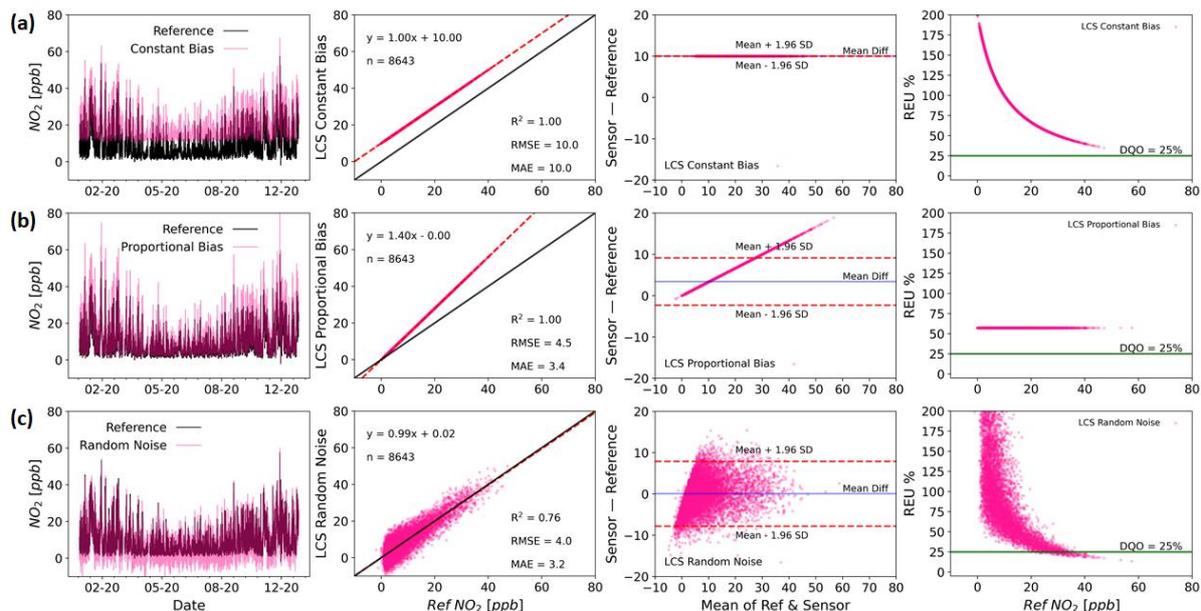
91 When characterising measurement error, in the absence of evidence to the contrary ~~it is often assumed~~ a linear
92 additive model ~~is often assumed~~. Once the analytical form of the model is defined, its parameters aim to capture
93 the error characteristics, and in the case of linear models (Eq. (1)), these are typically separated into three types
94 (Tian et al., 2016): (i) proportional bias or scale error (b_1), (ii) constant bias or displacement error (b_0) and (iii)
95 random error (ε) (Tian et al., 2016). Any measurement (y_i , e.g from the LCS) can therefore be thought of as a
96 combination of the reference value (x_i) and the three error types, such that:

$$97 \quad y_i = b_1 x_i + b_0 + \varepsilon \quad (1)$$

98 As the simplest approximation, this linear relationship for the error characteristics is often used to correct for
99 observed deviations between measurements and the agreed reference. It is worth to note, however, that this
100 equation assumes time-independent error contributions and that the three error components are not correlated,
101 which is often not the case on both counts (e.g. responses to non-target compounds). The parameter values
102 determined for Eq. (1) are also generally only applicable for individual instruments, potentially in specific
103 environments, unless the transferability of these parameters between devices has been explicitly demonstrated.

104 Figure 1 shows examples of how pure constant bias (a-panels), pure proportional bias (b-panels), and pure random
105 noise (c-panels) would look like in time-series, regression, Bland-Altman (B-A) (Altman & Bland, 1983) and
106 Relative Expanded Uncertainty (REU, as defined by the GDE (2010)) plots. In each of these ideal cases, the error
107 plots enable the practitioner to view the error characteristics in slightly different ways, allowing the impacts of the
108 observed measurement uncertainty to be placed into the context of the data requirements. In this work, we will

109 refer to them as “error types” (in contrast to “error sources”), which is the way they are distilled by the linear error
 110 model.



111
 112 **Figure 1. Time series (left panels), regression (middle-left panels), B-A Bland-Altman (middle-right panels) and REU**
 113 **(right panels, DQO for $NO_2 = 25\%$) plots for arbitrary examples of pure constant bias (Slope = 1, Intercept = 1, $SD_\epsilon =$**
 114 **0; a-panels), pure proportional bias (Slope = 1.4, Intercept = 0, $SD_\epsilon = 0$; b-panels) and pure random noise (Slope = 1,**
 115 **Intercept = 0, $SD_\epsilon = 4$; c-panels) simulated errors.**

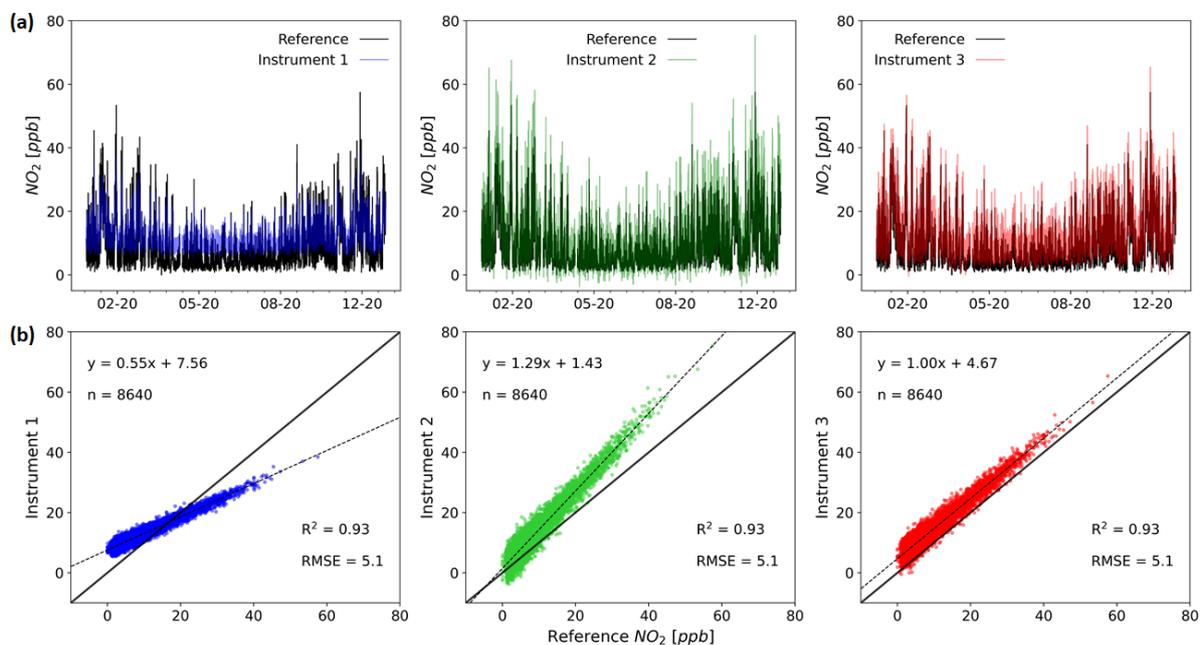
116

117 2.1 Performance indices, error structure and uncertainty

118 A major challenge faced by end-users of measurement devices characterised using collocation studies is the non-
 119 trivial question of how the comparisons themselves are performed and how the data are is communicated. Often
 120 single value performance metrics, such as the coefficient of determination (R^2) or root mean squared error
 121 (RMSE), are calculated between the assessed method (e.g. LCS) and an agreed reference, and the user is expected
 122 to infer an expected device performance or uncertainty for a measurement in their application (Duvall et al., 2016;
 123 Malings et al., 2019). ~~These metrics contain useful information about the measurement, but they are unable to~~
 124 ~~fully describe the error characteristics, in part because they reduce the error down to a single value (Tian et al.,~~
 125 ~~2016).~~ When evaluating multiple sensors during a collocation experiment, single metrics can be a useful way to
 126 globally compare instruments/sensors. However, these metrics do little to communicate the nature of the
 127 measurement errors and the impacts these will have in any end use application, in part because they reduce the
 128 error down to a single value (Tian et al., 2016). Even more if a specific concentration range is of paramount
 129 interest to the end-user, these metrics are not capable of characterising the weight of noise and/or the bias effect.
 130 The R^2 shows globally the data set linearity and gives an idea of the measurement noise. However, it is unable to
 131 distinguish whether a specific range of concentrations is more or less linear (or more or less noisy) than another.
 132 Similarly, the RMSE is also a very useful metric and perhaps more complete than R^2 , as it considers both noise
 133 and bias (although they need to be explicitly decomposed from RMSE). Nevertheless, the RMSE is an average

134 measure (of noise and bias) over the entire dataset under analysis. Using combinations of simple metrics increases
 135 the information communicated, but does not necessarily make it easy to assess how the errors will likely impact
 136 a particular measurement application. Visualising the absolute and relative measurement errors across the
 137 concentration range (unreachable by global metrics) enables end users to view the errors, and any features (non-
 138 linearities, step changes, etc.) that would impact the measurement but that global metrics (and in some cases time-
 139 series and/or regression plots) are incapable of showing.

140 Unfortunately, the widespread use of a small number of metrics as the sole method to assess measurement
 141 uncertainty, without a thorough consideration of the nature of the measurement errors, means measurement
 142 devices are often chosen that are unable to provide data that is fit for purpose. In addition, unconscious about
 143 potential flaws, users (e.g. researchers) could communicate findings or guide decision making based on results
 144 that may not justify the conclusions drawn from the data. Figure 2 shows three simulated measurements compared
 145 with the true values. Despite the measurements having identical R^2 and RMSE values, the time series and
 146 regression plots show that the error characteristics are significantly different, and would impact how the data from
 147 such a device could viably be used.



148
 149 **Figure 2. Time series (a-panels) and regression plots (b-panels) for three hypothetical instruments and a reference (1**
 150 **year of data). The most used metrics for evaluating the performance of LCS (R^2 and RMSE) are identical for the**
 151 **systems shown, even when the errors have very different characteristics (time res 1 h).**

152 There are multiple performance metrics that can be used for the assessment of measurement errors and uncertainty.
 153 Tian et al (2016) present an excellent summary of some of the major pitfalls of performance metrics and promote
 154 an approach of error modelling as a more reliable method of uncertainty quantification. These modelling
 155 approaches, however, rely on the assumption of statistical stationarity, whereby the statistical properties of the
 156 error are constant in the temporal and spatial domains. The presence of unknown or poorly characterised sources
 157 of error, for example, due to interferences from other atmospheric constituents or drifts in sensor behaviour, makes
 158 this assumption difficult to satisfy, especially when the dependencies of these errors show high spatial and

159 temporal variability. Thus, if field collocation studies are the primary method for performance assessment, as is
160 the case for LCS, only through a detailed assessment of the measurement errors across a wide range of conditions
161 and timescales can the uncertainty of the measurement be realistically estimated.

162 **2.2 Dealing with errors: established techniques vs Low-Cost Sensors**

163 Different approaches are available to the user to minimise the impact of errors, generally by making corrections
164 to the sensor data. For example, in the case of many atmospheric gas analysers, if the error is dominated by a
165 proportional bias, a multi-point calibration can be performed using standard additions of the target gas.
166 Displacement errors can be quantified, and then corrected for, by sampling a gas stream that contains zero target
167 gas. **And** Random errors can be reduced by applying a smoothing filter (e.g moving average filter, time-averaging
168 the data, etc.), at the cost of losing some information (Brown et al., 2008). These approaches work well for simple
169 error sources that, ideally, do not change significantly over timescales from days to months. Unfortunately, more
170 complex error sources can manifest in such a way that they contribute across all three error types, and also vary
171 temporally and spatially. For example, an interference from another gas-phase compound could in part manifest
172 itself as a displacement error, based on the instrument response to its background value, and in part as a
173 proportional bias if its concentration correlates with the target compounds, with any short-term deviations from
174 perfect correlation contributing to the random error component. In this case, time-averaging combined with
175 periodic calibrations and zeros would not necessarily minimise the error, and the user would need to employ
176 different tactics. One option would be to independently measure the interferent concentration, albeit with
177 associated uncertainty, and then use this to derive a correction. This is feasible if a simple and cost-effective
178 method exists for quantifying the interferent and its influence on the result is understood, but can make it very
179 difficult to separate out error sources, and can become increasingly complex if this measurement also suffers from
180 other interferences.

181 For many measurement devices, in particular for LCS based instruments, a major challenge is that the sources and
182 nature of all the errors are unknown or difficult to quantify across all possible end-use applications, meaning
183 estimates of measurement uncertainty are difficult. In the case of most established research and reference-grade
184 measurement techniques, comprehensive laboratory and field experiments have been used to explore the nature
185 of the measurement errors (Gerboles et al., 2003; Zucco et al., 2003). Calibrations have then been developed,
186 where traceable standards are sampled and measurement bias, both constant and proportional, can be corrected
187 for. Interferences from variables such as temperature, humidity, or other gases, have also been identified and then
188 either a solution engineered to minimise their effect or robust data corrections derived. Unfortunately, these
189 approaches have been shown not to perform well in the assessment of LCS measurement errors, due to the
190 presence of multiple, potentially unknown, sensor interferences from other atmospheric constituents (Thompson
191 & Ellison, 2005). These significant sensitivities to constituents such as water vapour and other gases mean
192 laboratory-based calibrations of LCS become exceedingly complex, and expensive, as they attempt to simulate
193 the true atmospheric complexity, often resulting in observed errors being very different to real-world sampling
194 (Rai et al., 2017; Williams, 2020). This has resulted in collocation calibration becoming the accepted method for
195 characterising LCS measurement uncertainties (De Vito et al., 2020; Masson et al., 2015; Mead et al., 2013;
196 Popoola et al., 2016; Sun et al., 2017), where sensor devices are run alongside traditional reference measurement
197 systems for a period of time, and statistical corrections derived to minimise the error between the two. As the true

198 value of a pollutant concentration cannot be known, this collocation approach assumes all the error is in the low-
199 cost measurement. Although this assumption may often be approximately valid (i.e. reference error variance \ll
200 LCS error variance), no measurement is absent of uncertainty and this can be transferred from one measurement
201 to another, obscuring attempts to identify its sources and characteristics. A further consideration when the fast
202 time-response aspect of LCS data is important, is that reference measurement uncertainties are generally
203 characterised at significantly lower reported measurement frequencies (typically 1 hr). This means that a high
204 time-resolution (e.g. 1 min) reference uncertainty must be characterised in order to accurately estimate the LCS
205 uncertainty (requiring specific experiments and additional costs). If a lower time-resolution reference data set is
206 used as a proxy, then the natural variability timescales of the target compound should be known and any impact
207 of this on the reported uncertainty caveated.

208 Another challenge with this approach is that, unlike targeted laboratory studies, real-world collocation studies at a
209 single location, and for a limited time period, are not able to expose the measurement devices to the full range of
210 potential sampling conditions. As many error sources are variable, both spatially and temporally, using data
211 generated under a limited set of conditions to predict the uncertainty on future measurements is risky. Deploying
212 a statistical model makes the tacit assumption that all factors affecting the target variable are captured by the
213 model (and the data set used to build the model). This is very often an unrealistic demand, and in the complex
214 multifaceted system that is atmospheric chemistry, this is extremely unlikely to be tenable, resulting in a clear
215 potential for overfitting to the training dataset. Ultimately, however, these collocation comparisons with
216 instruments with a well-quantified uncertainty need to be able to communicate a usable estimate of the information
217 content of the data to end-users, so that devices can be chosen that are fit for a particular measurement purpose.

218 3. Methods

219 In this work, we explore measurement errors, and their impacts, using the most common single value metrics: the
220 Coefficient of Determination or R^2 , the Root Mean Squared Error or RMSE and the Mean Absolute Error or MAE
221 (see the equation definitions in Cordero et al., 2018), ~~along with two additional widely used approaches to~~
222 ~~visualise the error distribution across a dataset:~~. To visualise the error distribution across a dataset we have also
223 employed two additional widely used approaches: the Bland-Altman plots (B-A) and Relative Expanded
224 Uncertainty (REU).

225 The performance metrics provide a single value irrespective of the size of the dataset, and might appear convenient
226 for users when comparing across devices or datasets, but can encourage over-reliance on the metric, often at the
227 expense of looking at the data in more detail or bringing an awareness of the likely physical processes driving the
228 error sources. On the other hand, the use of visualisations such as B-A and REU is complementary to the
229 aforementioned metrics, with the added value that the user is now more aware of how the data looks like in an
230 absolute and/or relative error space, allowing them to distinguish some characteristics of interest. These
231 visualizations ~~The B-A and the REU plots~~ are indeed more laborious techniques and the interpretation can be
232 challenging for non-experts, but they provide additional insights into the nature of the errors, not attainable by
233 one or more combined performance metrics: while B-A plots shows the noise (dispersion of the data) and the bias
234 effect (tendency of the data) in an absolute scale, the REU can be explicitly decomposed in the noise and bias
235 components (see Yatkin et al., 2022).

236 In order to understand how the different tools used here show different characteristics of the error structure, some
237 errors commonly found in LCS are examined through simulation studies. Subsequently, two real world case
238 studies are presented: (i) LCS duplicates for NO₂ and PM_{2.5} belonging to the QUANT project located in two sites
239 -the Manchester Natural Environment Research Council (NERC) measurement Supersite, and the York Fishergate
240 Automatic Urban and Rural Network (AURN) roadside site- and (ii) a set of duplicate reference instruments (only
241 at Manchester Supersite). Table S1 shows the research grade instrumentation used for this study.

242 3.1 Visualisation tools

243 An ideal performance metric should be able to deliver not only a performance index but also an idea of the
244 uncertainty distribution (Chai & Draxler, 2014). This is difficult to deliver through a simple numerical value, and
245 easy to interpret visualisations of the data are often much more useful for conveying multiple aspects of data
246 performance. Figure 2 shows the two most common data visualisation tools, the time-series plot and the regression
247 plot. In the time series plot the instrument under analysis and the agreed reference are plotted together as a function
248 of time. This allows a user to visually assess tendencies of over or under prediction, differences in the base line
249 or other issues, but can be readily over interpreted and does not allow for easy quantification of the observed
250 errors. In the regression plot the data from the instrument under analysis is plotted against the agreed reference
251 data. This allows for the correlation between the two methods to be more readily interpreted, in particular any
252 deviations from linearity, but gives little detail on the nature of the errors themselves.

253 In contrast to the regression plot -where the measured values from the two measurements (e.g. LCS vs Ref) are
254 plotted against each other- the Bland-Altman plots essentially display the difference between measurements
255 (abscissa) as a function of the average measurement (ordinate), enabling more information on the nature of the
256 error to be communicated. This direct visualisation of the absolute error acknowledges that the true value is
257 unknown and that both measurements have errors. The B-A plot enables the easy identification of any systematic
258 bias between the measurements or possible outliers, and is the reason B-A plots are extensively used in analytical
259 chemistry and biomedicine to evaluate agreement between measurement methods (Doğan, 2018). The mean
260 difference between the measurements, (represented by the blue line in the figures), is the estimated bias between
261 the two observations. The spread of error values around this average line indicates if the error shows purely
262 random fluctuations around this mean, or if it has structure across the observed concentration range.

~~263 In contrast to the regression plot, Bland Altman (B-A) plots essentially display the difference between
264 measurements, enabling more information on the nature of the error to be communicated. B-A plots (Altman &
265 Bland, 1983) are extensively used in analytical chemistry and biomedicine to evaluate the differences between
266 two measurement techniques (Doğan, 2018). The B-A is a scatter plot, in which the abscissa represents the average
267 of these measures (e.g. LCS and a reference measurement), acknowledging that the true value is unknown and that
268 both measurements have errors, and the ordinate shows the difference between the two paired measurements.~~

269 In the case where all the error is assumed to be in one of the measurements, e.g. comparing a LCS to a reference
270 grade measurement, there is an argument that the B-A abscissa could be the agreed reference value instead of the
271 average of two measurements. However, in this work we use the average of the two values as per the traditional

272 B-A analysis. To illustrate the B-A interpretation, from the error model (Eq. (1)) we can derive the following
273 expression:

$$274 \quad y_i - x_i = x_i (b_1 - 1) + b_0 + \varepsilon \quad (2)$$

275 From Eq. (2) it can be seen that if $b_1 \neq 1$ or if the error term (ε) variance is non-constant (e.g. heteroscedasticity)
276 the difference will not be normally distributed. The B-A plot (with x_i as the reference instrument results) allows a
277 quick visual assessment of the error distribution without the need to calculate the model parameters. In the case
278 the differences are normally distributed, the so-called “agreement interval” (usually defined as $\pm 2\sigma$ around the
279 mean) will hold 95% of the data points. Even though the estimated limits of agreement will be biased if the
280 differences are not normally distributed, it can still be a valuable indicator of agreement between the two
281 measurements.

282 If the ultimate goal of studying measurement errors is to diagnose the measurement uncertainty in a particular
283 target measurement range, then visualising the uncertainty in pollutant concentration space can be very
284 informative. The REU (GDE, 2010) provides a relative measure of the uncertainty interval about the measurement
285 within which the true value can be confidently asserted to lie. The abscissa in an REU plot represents the agreed
286 reference pollutant concentration, whose error is taken into account, something not considered by the other metrics
287 or visualisations discussed. The REU is regularly used to assess measurement compliance with the Data Quality
288 Objective (DQO) of the European Air Quality Directive 2008/50/EC, and is mandatory for the demonstration of
289 equivalence of methods other than the EU reference methods. For LCS the REU is widely used as a performance
290 indicator (Bagkis et al., 2021; Bigi et al., 2018; Castell et al., 2017; Cordero et al., 2018; Spinelle et al., 2015).
291 However, the evaluation of this metric is perceived as arduous and cumbersome and it is not included in the
292 majority of sensor studies (Karagulian et al., 2019). There is now a new published European Technical
293 Specification (TS) for evaluating the LCS performance for gaseous pollutants (CEN/TS 17660-1:2021). It
294 categorises the devices in 3 classes according to the DQO (Class 1 for “indicative measurements”, Class 2 for
295 “objective estimations”, and Class 3 for non-regulatory purposes, e.g. research, education, citizen science, etc.).
296 In the following sections, we use these established methods for assessing measurement uncertainty, alongside
297 simple time series and regression plots, to explore different error sources and their implications for air pollution
298 measurements.

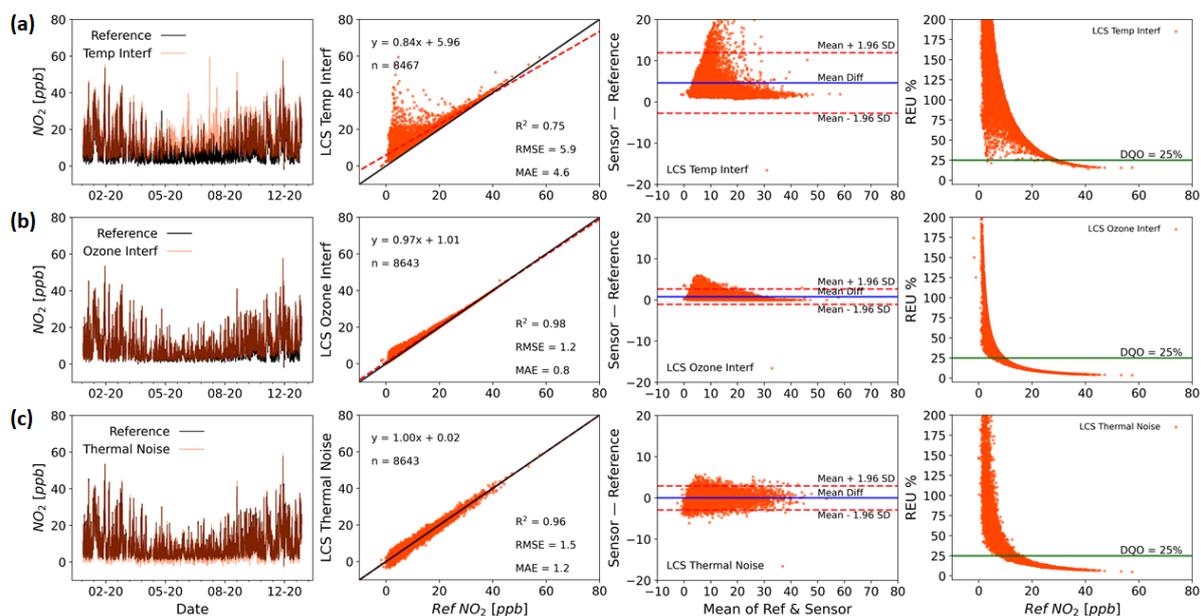
299 **4. Case studies**

300 **4.1 Simulated instruments**

301 In order to investigate the impact of different origins of measurement error on measurement performance, a set of
302 simulated datasets have been created. These data are derived using real-world reference data as the true values,
303 with the subsequent addition of errors of different origins to generate the simulated measurement data. Error
304 origins were chosen for which examples have been described in the LCS literature. Performance metrics along
305 with visualisation methods are then used to assess measurement performance.

306 As the complexity of the error increases, the impact of the assumption of statistical stationarity can become more
307 difficult to satisfy, with the magnitude of the errors becoming less uniform across the observed concentration, and

308 hence spatial, or time domains. Figure 3 shows examples of modelled sources of errors on NO₂ measurements:
 309 temperature interference (correction model taken from (Popoola et al., 2016), a-panels), a non-target gas (ozone)
 310 interference (correction model taken from (Peters et al., 2021), b-panels) and thermal electrical noise (white noise,
 311 c-panels).

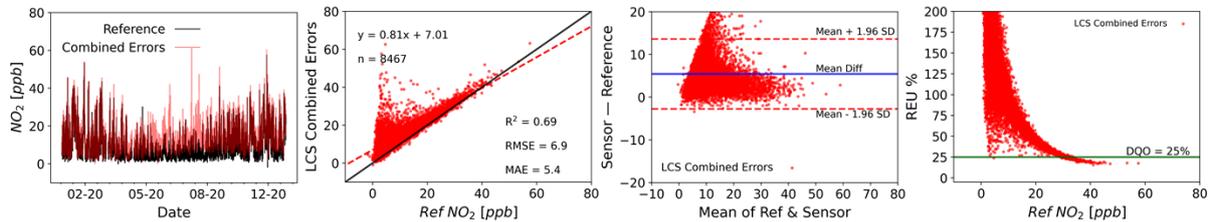


312
 313 **Figure 3. Time series (left panels), regression plots (middle-left panels, including R², RMSE & MAE), Bland-Altman**
 314 **plots (middle-right panels) and REU (right panels, DQO for NO₂ = 25%) for temperature (a-panels), ozone (b-panels)**
 315 **and thermal electrical noise (c-panels) modelled interferences on NO₂ measurements (time res 1 h).**

316 The above simulations show examples of how individual sources of error can impact measurement performance.
 317 Figure S1 shows some more examples, this time for different drift effects (baseline drift, temperature interference
 318 drift and instrument sensitivity drift). This set of error origins is not exhaustive, with countless others potentially
 319 impacting the measurement, such as those coming from (i) hardware (sensor-production variability, sampling,
 320 thermal effects due to materials expansion, drift due to ageing, RTC lag, Analog-to-Digital conversion,
 321 electromagnetic interference, etc.), (ii) software (signal sampling frequency, signal-to-concentration conversion,
 322 concept drift, etc.), (iii) sensor technology/measurement method (selectivity, sensitivity, environmental
 323 interferences, etc.) and (iv) local effects (spatio-temporal variation of concentrations, turbulence, sampling issues
 324 etc.).

325 Each error source impacts the uncertainty of the measurement, which in turn impacts its ability to provide useful
 326 information for a particular task. For example, the form of the temperature interference shown in Fig. 3 (a-panels)
 327 results in the largest errors being seen at the lower NO₂ values. This is because NO₂ concentrations are generally
 328 lowest during the day, due to photolytic loss when temperatures are highest. Thus this device would be better
 329 suited to an end-user intending to assess daily peak NO₂ concentration compared with the daytime hourly exposure
 330 values, providing the environment the device was deployed in showed a similar relationship between temperature
 331 and true NO₂ as that used here. The O₃ interference shown in Fig. 3 (b-panels) is similar, due again to a general
 332 anti-correlation observed between ambient O₃ and NO₂ concentrations. This type of interference can often be

333 interpreted incorrectly as a proportional bias, and a slope correction applied to the data. However, this type of
 334 correction will ultimately fail as O₃ concentrations are dependent on a range of factors, such as hydrocarbon
 335 concentrations and solar radiation, and as these change the O₃ concentration relative to the NO₂ concentration will
 336 change. To further complicate matters, multiple error sources can act simultaneously, meaning that the majority
 337 of measurements will contain multiple sources of error. Figure 4 shows a simple linear combination of the
 338 modelled errors shown in Fig 3, and the impact this has on the performance metrics.



339
 340 **Figure 4. Time series (left panel), regression plot (middle-left panel, including R², RMSE & MAE), Bland-Altman**
 341 **plot (middle-right panel) and REU (right panel, DQO for NO₂ = 25%) for a linear combination of temperature, ozone**
 342 **and thermal electrical noise modelled interferences (time res 1 h).**

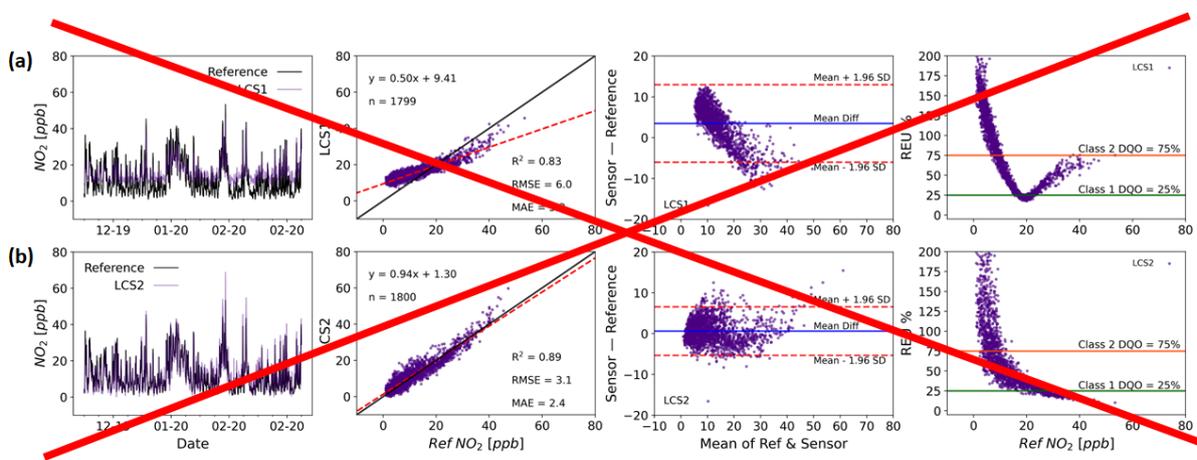
343 As the simulations show, the nature of the errors determine the observed effect on the measurement performance.
 344 In an ideal situation, like those shown in figures 3 and 4, the error sources would be well characterised, allowing
 345 the error to be modelled and approaches such as calibrations (for bias) and smoothing (for random errors)
 346 employed to minimise the total uncertainty. Unfortunately, in scenarios where sources of error and their
 347 characteristics are not known, modelling the error becomes more difficult and a more empirical approach to
 348 assessing the measurement performance and uncertainty may be required. The growing use of LCS represents a
 349 particular challenge in this regard. The susceptibility of LCS to multiple, often unknown or poorly characterised,
 350 error sources means that in order to determine if a particular LCS is able to provide data with the required level
 351 of uncertainty for a given task, a relevant uncertainty assessment is required. The following section explores the
 352 uncertainty characteristics of several LCS, with unknown error sources, deployed alongside reference
 353 instrumentation in UK urban environments as part of the QUANT study.

354 4.2 Real-world instruments

355 The difficulty in generating representative laboratory error characterisation data means for many measurement
 356 devices the error sources are essentially unknown. This, combined with the use of imperfect algorithms that are
 357 not available to the end-user (i.e. “black-box” models) to minimise errors, means that, collocation data is often the
 358 best option available to end-users in order to assess the applicability of a measurement method for their desired
 359 purpose. This is particularly the case for LCS air pollution measurement devices. In this section, we show
 360 collocation data collected as part of the UK Clean Air program funded QUANT project, and use the tools described
 361 above to investigate the impact of the observed errors on end-use.

362 ~~Figure 5 shows two colocated NO₂ measurements, from two different LCS devices using only their out of box~~
 363 ~~calibrations (i.e. no collocation data from that site was used to improve performance), compared with colocated~~
 364 ~~reference measurements at an urban background site in the city of Manchester. Unlike the modelled instruments~~
 365 ~~in Sect. 4.1, the combination of error sources is unknown in this case and we can thus only assess the LCS~~

366 measurement performance through comparison with the reference measurements using metrics and visual tools.
 367 There are obvious differences in the performance of both LCS instruments shown in Fig. 5. LCS1 (a panels)
 368 shows an appreciable difference in the time series baseline, which can be interpreted from both the regression (b1
 369 <1) and the B-A plots as a proportional bias. This bias also impacts the REU plot, with a minimum in the region
 370 where the regression best fit line crosses the 1:1 line (~ 17 ppb). It is worth noting that these plots do not directly
 371 identify the source of the proportional bias, with sensor response to the target compound or another covarying
 372 compound possible, but provides information on how much it impacts the data. For LCS2 (Fig. 5, b panels) any
 373 proportional bias is significantly smaller, with the B-A plot showing a much more symmetrical distribution of
 374 points around the central line across the observed mixing ratio range, although this is not a normal distribution as
 375 evidenced by the heteroscedastic nature of the differences, indicating the cause is not entirely random in nature.
 376 The lack of a large proportional bias also results in the REU plot showing a continued reduction in relative
 377 uncertainty as the true NO₂ concentration increases. Interestingly, both LCS's also show an additional bias at the
 378 highest NO₂ values observed. This does not significantly impact the REU, due to its relative nature, but can be
 379 seen in the regression and B-A plots. Correcting for the observed proportional bias in LCS1 and LCS2 improves
 380 the observed performance by providing the errors with a more symmetrical distribution (LCS1* and LCS2* shown
 381 in Fig. S2).



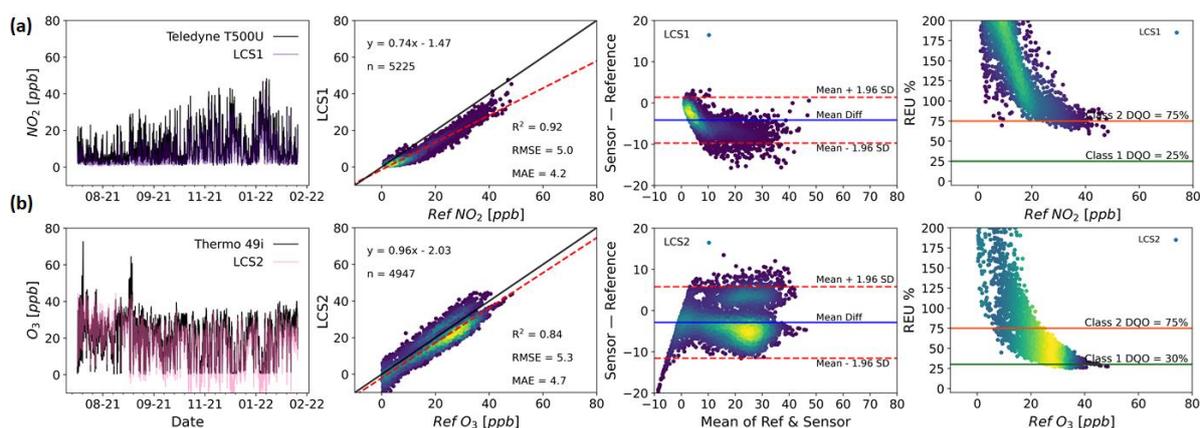
382
 383 **Figure 5. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and**
 384 **REU (right panels; NO₂-Class 1 DQO = 25% & Class 2 DQO = 75%) for NO₂ measurements by two LCS systems of**
 385 **different brands (a and b panels) in the same location (Manchester Supersite, December 2019 to February 2020. Time**
 386 **res 1 h).**

387 Figure 5 shows two colocated measurements from two different LCS devices: one measuring NO₂ (a-panels) and
 388 the other O₃ (b-panels). Both measurements are compared with colocated reference measurements at an urban
 389 background site in the city of Manchester. Unlike the modelled instruments in Sect. 4.1, the combination of error
 390 sources is unknown in this case, and we can thus only assess the LCS measurement performance through
 391 comparison with the reference measurements using global metrics and visual tools.

392 Single value metrics indicate an acceptable performance for both measurements: high linearity (both R² are higher
 393 than 0.8) and relatively low errors (RMSE \sim 5ppb). However, the plots present the data in a variety of ways that
 394 enable the user to identify patterns in the measurement errors that would be less obvious if only global metrics

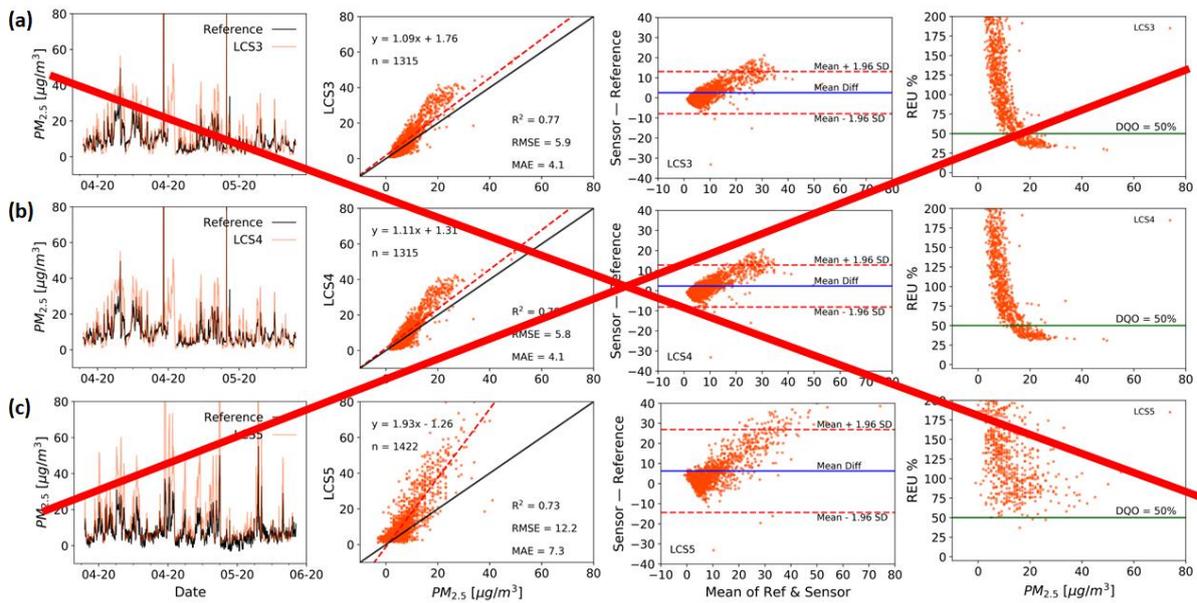
395 were used. For example, the NO₂ sensor (LCS1, a-panels) has a non-linear response that is almost imperceptible
 396 from the regression plot but stands out in the B-A plot. Furthermore (despite the high R² and relatively low
 397 RMSE), the REU plot shows high relative errors that do not meet the Class 2 DQO for the measured concentration
 398 range. Regarding the O₃ sensor (LCS2, b-panels), the B-A plot shows two high density measurement clusters, one
 399 with positive absolute errors (over-measuring) and a larger one with negative errors (under-measuring). These are
 400 the result of a step change in the correction algorithm applied by the manufacturer and could easily have been
 401 missed if only summary metrics and a regression plot were used, especially if the density of the data points was
 402 not coloured.

403 It is worth noting that these plots do not directly identify the source of the proportional bias, with sensor response
 404 to the target compound or another covarying compound possible, but provides information on how much it impacts
 405 the data.



406
 407 **Figure 5. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and**
 408 **REU (right panels; NO₂ Class 1 DQO = 25% & Class 2 DQO = 75% ; O₃ Class 1 DQO = 30% & Class 2 DQO = 75%)**
 409 **for NO₂ (a-panels) and O₃ (b-panels) measurements by two LCS systems of different brands in the same location and**
 410 **time span (Manchester Supersite, July 2021 to February 2022. Time res 1 h). All but the time-series plots, have coloured**
 411 **by data density.**

412 ~~Figure 6 shows three out of the box PM_{2.5} measurements made by three devices from the same brand in spring,~~
 413 ~~located at two sites: the first two at an urban background (LCS3 & LCS4, a and b panels) and the third at a roadside~~
 414 ~~(LCS5, c panels). As the regression and the B-A plots show, all LCS measurements in Fig. 6 have a proportional~~
 415 ~~bias compared with the reference, with the LCS over predicting the reference values. Both LCS's at the urban~~
 416 ~~background site show very similar performance, indicating that the devices are similarly affected by errors. This~~
 417 ~~internal consistency is highly desirable, especially when LCS's are to be deployed in networks, as although mean~~
 418 ~~absolute measurement error may be high, differences between identical devices are likely to be interpretable.~~



419
420
421
422
423

Figure 6. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels, DQO for $PM_{2.5} = 50\%$) for $PM_{2.5}$ measurements by three LCS systems of the same brand (panels a, b and c) in different locations: an urban background (Manchester Supersite, panels a and b) and a roadside site (York, panel c) (April & May 2020, time res 1 h).

424
425
426
427

Figure 6 shows three out-of-the-box $PM_{2.5}$ measurements made by two devices (LCS3 & LCS4) from the same brand in spring (LCS3: a-panels; LCS4: c-panels) and in autumn (b-panels, only LCS3). The collocation shown correspond to two different sites: an urban background site (LCS3, a and c-panels) and a roadside site (LCS4, c-panels).

428
429
430
431
432
433
434
435

As the regression and the B-A plots show, all LCS measurements in Fig. 6 have a proportional bias compared with the reference, with the LCS over predicting the reference values. The device at the urban background site (LCS3) show a dissimilar performance in spring and autumn, indicating that the errors this device suffers are differently influenced by local conditions in the two seasons (all the duplicates at the urban background show the same pattern). While for LCS3 during spring the error have a more linear behaviour, in autumn a non-linear pattern is clearly observed in the regression and B-A plots. Despite the utility that single metrics can have in certain circumstances, the non-linear pattern goes completely unnoticed by them: while for the two different seasons RMSE and the MAE are almost constant the R^2 indicates a higher linearity for autumn.

436
437
438
439

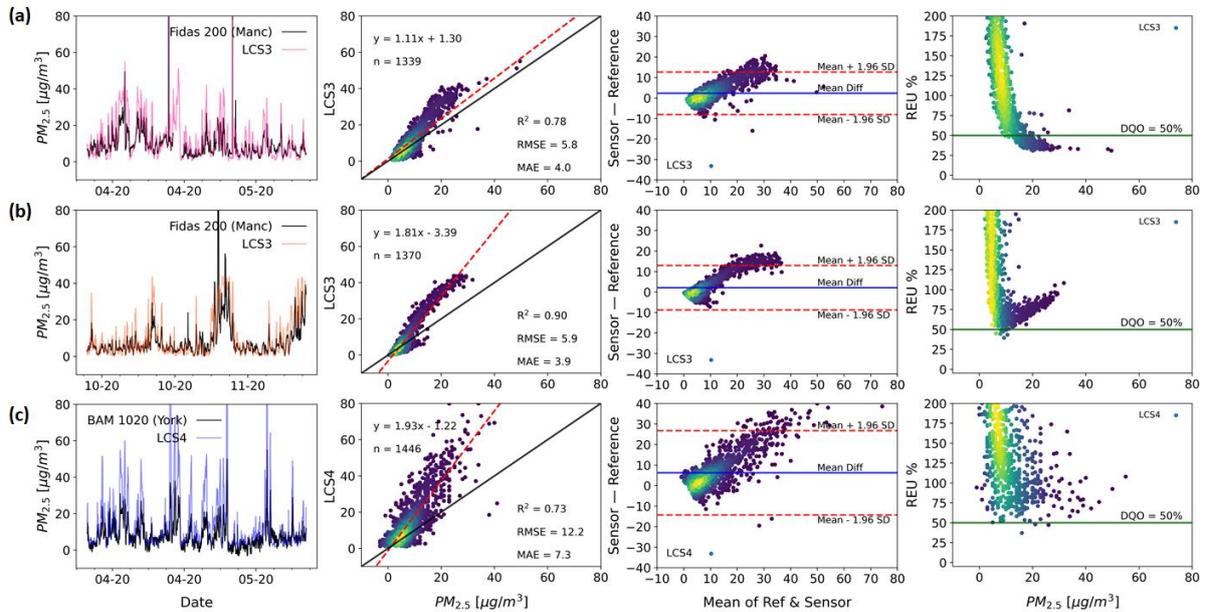
A number of duplicates were deployed at both sites showing a very similar performance in terms of the single metric values but also in regard to the more visual tools (not shown here). This internal consistency is highly desirable, especially when LCS's are to be deployed in networks, as although mean absolute measurement error may be high, differences between identical devices are likely to be interpretable.

440
441
442
443
444

Having prior knowledge of the nature of the measurement errors allows informed experimental design prior to data collection. This is key if an end user is to maximise the power of a dataset, and the information it provides, to answer a specific question. For example, if an end-user wanted to identify pollution hotspots within a relatively small geographical area, then using a dense network of sensor devices that posses errors large and variable enough to make quantitative comparisons with limit values difficult (possibly due to an interference from a physical

445 parameter like relative humidity) but show internal consistency could be a viable option. Providing the hotspot
 446 signal is large enough relative to any random error magnitude.

447



448

449 **Figure 6. Two LCS systems (LCS3 & LCS4, same brand) measuring PM_{2.5} (Time res 1 h). While LCS3 is shown for**
 450 **the same location (Manchester) but unfolded in two different seasons (a-panels: Apr to May 2020; b-panels: Oct to Nov**
 451 **2020), LCS4 is at a different location (c-panels: York, Apr to May 2020). Time series (left panels), regression plots**
 452 **(middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels; DQO_{PM_{2.5}} = 50%) are used to**
 453 **characterise the device's error structure. All but the time-series plots have been coloured by data density.**

454

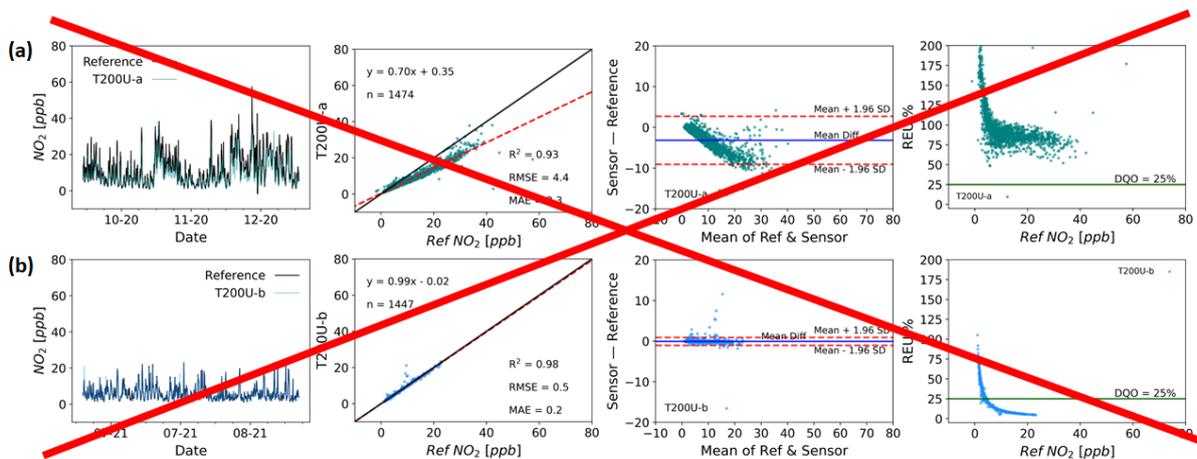
455 The LCS data from the roadside location (LCS4) show significantly lower precision than those at the urban
 456 background site, as seen in the B-A plot. This could be caused by differences in particle properties and size
 457 distributions between the two sites (Gramsch et al., 2021), and by the high frequency variation of transport
 458 emissions close to the roadside **site side** and turbulence effects (Baldauf et al., 2009; Makar et al., 2021). Duplicate
 459 measurements show that all sensors of this type responded similarly in this roadside environment (not shown
 460 here), supporting the high internal consistency of this device, but indicating a spatial heterogeneity in some key
 461 error sources. It is also worth noting that the gold standard instruments at the two sites are not “reference method”
 462 but “reference equivalent methods” (GDE, 2010), each using a different measurement technique: while an optical
 463 spectrometer (Palas Fidas 200) is used in Manchester, the York instrument uses a Beta attenuation method (Met
 464 One BAM 1020), which could also potentially lead to some of the observed differences. The increased apparent
 465 random variability for LCS4, combined with the proportional bias, results in significantly higher measurement
 466 uncertainty across the observed range, as can be seen by the REU plots, with LCS4 never reaching an acceptable
 467 DQO level (50% for PM_{2.5}). **As with the NO₂ sensors (Fig. 5).** If the observed proportional bias is corrected the
 468 linearly bias-corrected sensors (Fig. S3) show a much improved comparison with the reference measurement,
 469 specially **LCS3* in autumn and LCS4*.** The error distribution for the LCS3 (autumn) shown by the B-A plot is

470 greatly narrowed (~3x times) and now the sensor is accomplishing the DQO below 10 $\mu\text{g m}^{-3}$ as the REU plot
 471 indicates. For LCS4 In this case the B-A plot shows an error characteristic more dominated by random errors, and
 472 the REU plots shows a significant reduction of the relative uncertainty, with the REU at 10 $\mu\text{g m}^{-3}$ reducing from
 473 ~75 to ~50%.

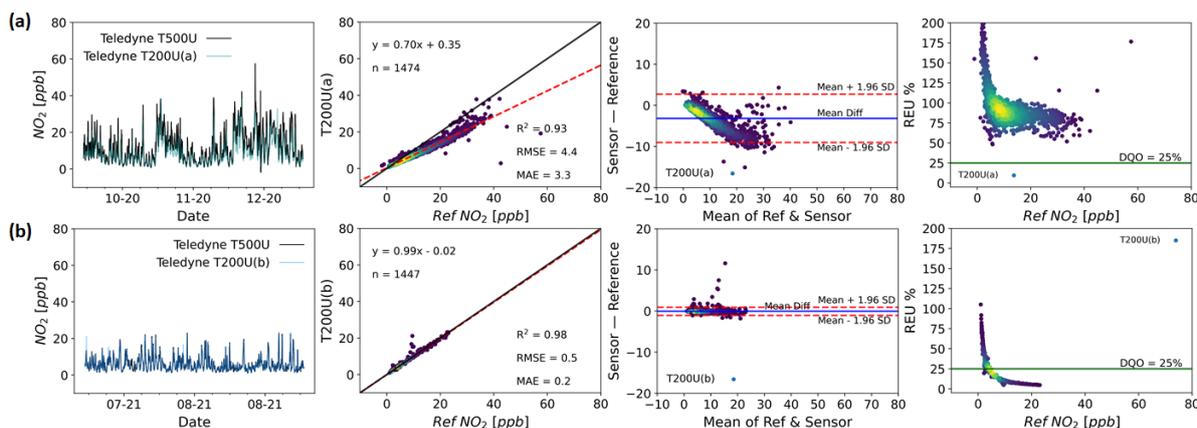
474 As a comparison for the LCS data shown above, Fig. 7 shows two identical NO₂ reference grade instruments,
 475 Teledyne T200U (Chemiluminescence method) at the Manchester urban background site (panels a and b) during
 476 two different time periods, with a Teledyne T500U (CAPS detection method) used as the “ground truth”
 477 instrument. Instrument “a” manifests a significant proportional bias, in contrast to instrument “b”, but both show
 478 differences that could be non-negligible depending on the application. The deviations observed in instrument “a”
 479 was due to the cell pressure being above specification by ~20%, unnoticed while the instrument was in operation.
 480 This demonstrates the importance of checking instrument parameters regularly in the field even if the data appears
 481 reasonable.

482 As the LCS error structure is determined relative to the performance of a reference measurement, if the reference
 483 instrument suffers from significant errors this will affect the outcomes of the performance assessment, due to the
 484 assumption that all the errors reside with the LCS. As Fig. 7 shows, however, this assumption is not necessarily
 485 always valid and potentially argues that reference instruments used in colocation studies should be subject to
 486 further error characterisation, including possible colocation with other reference instruments. As a similar
 487 comparison of reference instruments, Fig. S3 shows two ozone research grade instruments (a Thermo 49i and a
 488 2B).

489 It is worth noting that even when using reference, or reference equivalent, grade instrumentation, inherent
 490 measurement errors mean that relative uncertainty, as shown in the REU plot, increases asymptotically at lower
 491 values. This is not unexpected, but is potentially important as ambient target concentration recommendations
 492 continue to fall based on updated health evidence (World Health Organization, 2021).



493



494

495 **Figure 7. Time series (left panel), regression plots (middle-left panel), Bland-Altman plots (middle-right panel) and**
 496 **REU (right panel, DQO for NO₂ = 25%) for two identical (Teledyne T200U) reference NO₂ instruments (panels a and**
 497 **b) colocated at the Manchester Supersite (1h time res). The first instrument between October & November 2020 and**
 498 **the second between July & August 2021. All but the time-series plots have been coloured by data density.**

499 **5. Discussion**

500 The widespread use of collocation studies to assess measurement device performance, means many examples exist
 501 in the LCS literature where different devices are compared using summary metrics for field or laboratory studies
 502 (Broday, 2017; Duvall et al., 2016; Hofman et al., 2022; Karagulian et al., 2019; Mueller et al., 2017; Rai et al.,
 503 2017; van Zoest et al., 2019). Although these comparisons do provide useful information, they can be misleading
 504 for end users wanting to compare the performance of different devices, as they are often carried out under different
 505 conditions and do not present the data or experimental design in full. Even in the case where comparisons have
 506 been done under identical conditions, the data still needs to be treated with caution, as inevitable differences
 507 between assessment environment and proposed application environment, as well as any changes to
 508 instrument/sensor design or data processing, mean that past performance does not guarantee future performance.

509 All measurement devices suffer from measurement errors, many of which are potentially significant depending
 510 on the application, with devices and their error susceptibility covering a broad spectrum. As evidenced by Fig. 7,
 511 reference instruments are not immune from this phenomena, with the proportional bias of one of the NO_x
 512 instruments clearly affecting its measurements resulting in the absolute error increasing with concentration. As
 513 the requirements on measurement devices continue to increase, driven in part by new evidence supporting the
 514 reduction of air pollutant target values, the devices currently being used for a particular application could no longer
 515 be fit-for-purpose in the situation where the limit value has decreased to the point where it is small relative to the
 516 device's uncertainty.

517 Single value performance metrics, such as R² and RMSE, can seem convenient when comparing multiple co-
 518 located devices as they facilitate decision making when a threshold criterion is defined. However, these scalar
 519 values hide important information about the scale and / or distribution of the errors within a dataset; graphical
 520 summaries of the measurements themselves can offer significantly more insight into the impact of measurement
 521 errors on device performance and ultimate capabilities. Of particular use in air pollution measurements is the
 522 ability to see how the errors manifest themselves in relation to our best estimate of the true pollutant concentration,

523 as often applications have specific target pollutant concentration ranges of interest. For example, the two NO_2
524 LCS devices shown in Fig. 5 have ~~similar R^2 values of 0.83 and 0.89, but one is suffering from a strong~~
525 ~~proportional bias that impacts on measurements either side of the 18ppb crossing point~~ considerably high R^2
526 values (0.92 and 0.84) and relatively low RMSE and MAE, but one suffering of non-linear errors (LCS1) and the
527 other with data coming from two different calibration states (LCS2).

528 Errors, or combinations of errors, frequently result in varying magnitude of the observed measurement
529 inaccuracies across the concentration space observed, and it is often useful to assess both the absolute and relative
530 effects of the errors. By getting a more complete picture of the device performance, decisions can be made on the
531 effectiveness of simple corrections, such as correcting for an apparent proportional bias using an assumption of a
532 linear error model. Ultimately end users need to identify the data requirements a priori and design quantifiable
533 success criteria by which to judge the data. For example, rather than just wanting to measure the 8-hour average
534 NO_2 , be more specific and require that this needs to be accurate to within 5 ppb, have demonstrated approximately
535 normally distributed errors in a representative environment for the period of interest, and no statistical evidence
536 of deviation from a linear correlation with the reference measurement over the target concentration range for the
537 period of interest.

538 A major challenge comes from complex errors, such as interferences from other compounds or with environmental
539 factors, that vary temporally and/or spatially. Similar graphical techniques to those presented above can be used
540 to identify the existence of such relationships, but correcting for them remains a challenge. For example, the
541 correlation between measurement errors and relative humidity could be explored by replacing the abscissa with
542 measured relative humidity in both the B-A and REU plots. This would visualise the relationship between absolute
543 and relative errors with relative humidity, but would not be able to confirm causality. The complex and covarying
544 nature of the atmosphere means that the best way to identify a device error source is through controlled laboratory
545 experiments, where confounding variables can be controlled, although these experiments are often difficult and
546 expensive to perform in a relevant way.

547 This brings into question the power of colocation studies, as they can ultimately never be performed under the
548 exact conditions for every intended application. The $\text{PM}_{2.5}$ sensors shown in Fig. 6 demonstrate this, as if a
549 colocation dataset generated at the urban background site was used to inform a decision about the applicability of
550 these devices to a roadside monitoring task, then an overly optimistic assessment of the scale of the errors to be
551 expected would be likely. It is therefore always desirable that colocation studies are as relevant as possible to the
552 desired application, and this is even more paramount in the case where the error sources are poorly specified. For
553 this reason, complete meta-data on the range of conditions over which a study was conducted is key information
554 in judging its applicability to different users.

555 Although there is no strict definition on what makes a device a LCS, we often make the categorization based on
556 the hardware used. Standard reference measurement instruments are generally based on well-characterised
557 techniques developed and improved over years, based primarily on the progressive refinement of hardware (e.g.
558 materials used for the detection elements, electronic circuits to filter noise, refinement of production methods,
559 etc.). Although LCS sensor technologies are improving, it is interesting that many of the significant improvements
560 that have been made to LCS performance have been through software, rather than hardware advances. As more

561 colocation data ~~are~~ ~~is~~ generated in different environments, many LCS manufacturers have been able to develop
562 data correction algorithms that minimise the scale of the errors that are present on the LCS hardware. This can
563 greatly improve the performance of LCS devices, and has been a large factor in the improvements seen in these
564 devices over recent years. These algorithms are, however, inevitably imperfect and can suffer from concept drift
565 (De Vito et al., 2020), caused by the lack of available colocation data over a full spectrum of atmospheric
566 complexity. Furthermore, any kind of statistical model introduces a new error source that can work in conjunction
567 with the pre-existing measurement errors to drastically change the observed error characteristics, making it much
568 more difficult for users to interpret and extrapolate from colocation study performance to intended application. ~~If~~
569 ~~end users are to be able to make well informed decisions about device applicability to a particular task, then an~~
570 ~~argument can be made for information on the scale of the error corrections made to a reported measurement to be~~
571 ~~made available, ideally alongside and a demonstration of its benefits in a relevant environment.~~ If end users are
572 to be able to make well informed decisions about device applicability then information on the scale of the
573 measurement errors, and the impact of corrections made to minimise these, should be made available. Exemplar
574 case studies in a range of relevant environments would also be highly valuable. Unfortunately, this colocation
575 data ~~are~~ ~~is~~ costly to generate, meaning relevant data often does not exist, and when it does is often not
576 communicated in such a way that enables the user to make a fully informed decision.

577 6. Conclusions

578 In situ measurements of air pollutants are central to our ability to identify and mitigate poor air quality.
579 Measurement applications are wide ranging, from assessing legal compliance to quantifying the impact of an
580 intervention. The range of available measurement tools for key pollutants is also increasingly broad, with
581 instrument price tags spreading several orders of magnitude. In order for a measurement device to be of use for a
582 particular application it must be fit-for-purpose, with cost, useability and data quality all needing to be considered.
583 Understanding measurement uncertainty is key in choosing the correct tool for the job, but in order for this to be
584 assessed the job needs to be fully specified a priori. The specific data requirements of each measurement
585 application need to be understood and a measurement solution chosen that is capable of providing data with
586 sufficient information content.

587 In order to aid end users in extrapolating from colocation study performance to potential performance in a specific
588 application, performance metrics are often used. Although single value performance metrics do convey some
589 useful information about the agreement between the data from the measurement device being assessed and the
590 reference data, they can often be misleading in their evaluation of performance. This dictates a more rigorous and
591 empirical approach to data uncertainty assessment in order to determine if a measurement is fit for purpose. The
592 ability to assess device performance across the observed concentration range, as in the B-A and REU plots, enables
593 an end-user to make an informed decision about the capabilities of a measurement device in the target
594 concentration range. These visual tools also help identify any simple corrections that can be applied to improve
595 performance. In contrast, if an end-user was only provided with a single value metric, such as R^2 or RMSE then
596 it would be significantly more difficult to understand the likely implications of the measurement uncertainties.

597 All measurement devices suffer from errors, which result in deviations between the reported and true values.
598 These errors can come from a multitude of sources, with the scale of the deviation from the true value being

599 dependent on the nature of the error. Although a known measurement uncertainty for all applications would be
600 ideal for end users to be able to assess measurement device suitability for purpose, in many cases, especially for
601 LCS, this is not possible due to the presence of poorly characterised, or sometimes unknown, error sources. In the
602 absence of this, useful information on likely measurement performance can be obtained using colocation data
603 compared with a measurement with a quantified uncertainty. It is important that such a colocation study is carried
604 out in an environment as similar as possible to the application environment, as the unknown nature of many error
605 sources means their magnitude can change significantly between different locations and/or seasons (e.g. Fig. 6).
606 Ideally, depending on the measurement task, the user could use the colocation data to model the error causes and
607 use this to develop strategies to minimise final measurement uncertainty. Unfortunately, relevant colocation study
608 **are is** often not available, and to generate the data would be prohibitively costly, which limits the user's ability to
609 make a realistic assessment of likely uncertainties. The presence of, often complex, error minimisation post
610 processing or calibration algorithms further complicates things. This additional uncertainty is most likely to bias
611 any performance prediction if the end user is unaware of the purpose or scale of the data corrections, and their
612 applicability to the target environmental conditions. Ideally, long term colocation data sets demonstrating the
613 performance of measurement hardware and software, in a range of relevant locations, over multiple seasons, and
614 carried out by impartial bodies would be available to inform measurement solution decisions.

615 In order for end users to take full advantage of the ever increasing range of air pollution measurement devices
616 available, the questions being asked of the data must be consummate with the information content of the data.
617 Ultimately this information content is determined by the measurement uncertainty. Thus, providing end users with
618 as accurate an estimate as possible of the likely measurement uncertainty, in any specific application, is essential
619 if end users are to be able to make informed decisions. Similarly, end users must specify the data uncertainty
620 requirements for each specific task if the correct tool for the job is to be identified. This requirement for air quality
621 management strategies to acknowledge the capabilities of available devices, both in the setting and monitoring of
622 limits, will only become increasingly important as target levels continue to decrease.

623 **Supplementary**

624 The supplement related to this article is available online at:

625 **Code and data availability**

626 [The code and data for this study can be found on Zenodo: https://zenodo.org/record/6518027#.YnKbH9PMJhE.](https://zenodo.org/record/6518027#.YnKbH9PMJhE)

627 [The live code can be found on GitHub: https://github.com/wacl-york/quant-air-pollution-measurement-errors.](https://github.com/wacl-york/quant-air-pollution-measurement-errors)

628 **Author contributions**

629 PE: Funding acquisition; Supervision. SD and PE: Project administration; Formal analysis. SD, PE & SL:
630 Conceptualization; Methodology; Investigation. SD & SL: Visualisation; Software. KR, NM, MF: Resources. SD,
631 SL, KR, NM, MF: Data curation. SD, PE, SL, TB, NM, TG & DH: Writing – review & editing.

632 **Competing interests**

633 The authors declare that they have no conflict of interest.

634 **Acknowledgements**

635 This work was funded as part of the UKRI Strategic Priorities Fund Clean Air program (NERC NE/T00195X/1),
636 with support from Defra. We would also thank the OSCA team (Integrated Research Observation System for
637 Clean Air) at the Manchester Air Quality Supersite (MAQS), for help in data collection for the regulatory-grade
638 instruments. The secondary research grade instruments used here (Thermo ozone 49i, 2B Technologies 202 ozone,
639 and Teledyne T200U NO_x) are provided through the Atmospheric Measurement and Observation Facility
640 (AMOF) and the calibrations were carried out in the COZI Laboratory, a facility housed at the Wolfson
641 Atmospheric Chemistry Laboratories (WACL). Both funded through the National Centre for Atmospheric Science
642 (NCAS). Special thanks to Elena Martin Arenos, Chris Anthony, Killian Murphy, Stuart Young, Steve Andrews
643 and Jenny Hudson-Bell from WACL for the help and support to the project. Also thanks to Stuart Murray and
644 Chris Rhodes from the Department of Chemistry Workshop for their technical assistance and advice. Thanks to
645 Andrew Gillah and Michael Golightly from the York Council who assisted with site access.

646

647 **References**

- 648 Altman, D. G. and Bland, J. M.: Measurement in Medicine: The Analysis of Method Comparison Studies, *J. R.*
649 *Stat. Soc. Ser. Stat.*, 32, 307–317, <https://doi.org/10.2307/2987937>, 1983.
- 650 Andrewes, P., Bullock, S., Turnbull, R., and Coolbear, T.: Chemical instrumental analysis versus human
651 evaluation to measure sensory properties of dairy products: What is fit for purpose?, *Int. Dairy J.*, 121,
652 105098, <https://doi.org/10.1016/j.idairyj.2021.105098>, 2021.
- 653 Bagkis, E., Kassandra, T., Karteris, M., Karteris, A., and Karatzas, K.: Analyzing and Improving the Performance
654 of a Particulate Matter Low Cost Air Quality Monitoring Device, *Atmosphere*, 12, 251,
655 <https://doi.org/10.3390/atmos12020251>, 2021.
- 656 Baldauf, R., Watkins, N., Heist, D., Bailey, C., Rowley, P., and Shores, R.: Near-road air quality monitoring:
657 Factors affecting network design and interpretation of data, *Air Qual. Atmosphere Health*, 2, 1–9,
658 <https://doi.org/10.1007/s11869-009-0028-0>, 2009.
- 659 Bigi, A., Mueller, M., Grange, S. K., Ghermandi, G., and Hueglin, C.: Performance of NO,
660 NO₂ low cost sensors and three calibration approaches within a real world
661 application, *Atmospheric Meas. Tech.*, 11, 3717–3735, <https://doi.org/10.5194/amt-11-3717-2018>, 2018.
- 662 Broday, D. M.: Wireless Distributed Environmental Sensor Networks for Air Pollution Measurement—The
663 Promise and the Current Reality, *Sensors*, 17, 2263, <https://doi.org/10.3390/s17102263>, 2017.
- 664 Brown, R. J. C., Hood, D., and Brown, A. S.: On the Optimum Sampling Time for the Measurement of Pollutants
665 in Ambient Air, *J. Autom. Methods Manag. Chem.*, 2008, 814715, <https://doi.org/10.1155/2008/814715>,
666 2008.

667 Castell, N., Dauge, F. R., Schneider, P., Vogt, M., Lerner, U., Fishbain, B., Broday, D., and Bartonova, A.: Can
668 commercial low-cost sensor platforms contribute to air quality monitoring and exposure estimates?,
669 *Environ. Int.*, 99, 293–302, <https://doi.org/10.1016/j.envint.2016.12.007>, 2017.

670 CEN/TS 17660-1: 2021, Air quality – Performance evaluation of air quality sensor systems – Part 1 Gaseous
671 pollutants in ambient air.
672 [https://standards.cen.eu/dyn/www/f?p=204:110:0:::FSP_PROJECT,FSP_LANG_ID:60880,25&cs=1B](https://standards.cen.eu/dyn/www/f?p=204:110:0:::FSP_PROJECT,FSP_LANG_ID:60880,25&cs=1B6992D14C0BCD6D6333E555D297F1306)
673 [6992D14C0BCD6D6333E555D297F1306](https://standards.cen.eu/dyn/www/f?p=204:110:0:::FSP_PROJECT,FSP_LANG_ID:60880,25&cs=1B6992D14C0BCD6D6333E555D297F1306) (accessed on 15 January 2022). 2021.

674 Chai, T. and Draxler, R. R.: Root mean square error (RMSE) or mean absolute error (MAE)? – Arguments against
675 avoiding RMSE in the literature, *Geosci. Model Dev.*, 7, 1247–1250, [https://doi.org/10.5194/gmd-7-](https://doi.org/10.5194/gmd-7-1247-2014)
676 [1247-2014](https://doi.org/10.5194/gmd-7-1247-2014), 2014.

677 Cordero, J. M., Borge, R., and Narros, A.: Using statistical methods to carry out in field calibrations of low cost
678 air quality sensors, *Sens. Actuators B Chem.*, 267, 245–254, <https://doi.org/10.1016/j.snb.2018.04.021>,
679 2018.

680 De Vito, S., Esposito, E., Castell, N., Schneider, P., and Bartonova, A.: On the robustness of field calibration for
681 smart air quality monitors, *Sens. Actuators B Chem.*, 310, 127869,
682 <https://doi.org/10.1016/j.snb.2020.127869>, 2020.

683 Doğan, N. Ö.: Bland-Altman analysis: A paradigm to understand correlation and agreement, *Turk. J. Emerg. Med.*,
684 18, 139–141, <https://doi.org/10.1016/j.tjem.2018.09.001>, 2018.

685 Duvall, R. M., Long, R. W., Beaver, M. R., Kronmiller, K. G., Wheeler, M. L., and Szykman, J. J.: Performance
686 Evaluation and Community Application of Low-Cost Sensors for Ozone and Nitrogen Dioxide, *Sensors*,
687 16, 1698, <https://doi.org/10.3390/s16101698>, 2016.

688 Feenstra, B., Papapostolou, V., Hasheminassab, S., Zhang, H., Boghossian, B. D., Cocker, D., and Polidori, A.:
689 Performance evaluation of twelve low-cost PM_{2.5} sensors at an ambient air monitoring site, *Atmos.*
690 *Environ.*, 216, 116946, <https://doi.org/10.1016/j.atmosenv.2019.116946>, 2019.

691 Feinberg, S. N., Williams, R., Hagler, G., Low, J., Smith, L., Brown, R., Garver, D., Davis, M., Morton, M.,
692 Schaefer, J., and Campbell, J.: Examining spatiotemporal variability of urban particulate matter and
693 application of high-time resolution data from a network of low-cost air pollution sensors, *Atmos.*
694 *Environ.*, 213, 579–584, <https://doi.org/10.1016/j.atmosenv.2019.06.026>, 2019.

695 GDE. Guidance for the Demonstration of Equivalence of Ambient Air Monitoring Methods.
696 <https://ec.europa.eu/environment/air/quality/legislation/pdf/equivalence.pdf> (accessed on 20 Dec 2021).
697 2010.

698 Gerboles, M., Lagler, F., Rembges, D., and Brun, C.: Assessment of uncertainty of NO₂ measurements by the
699 chemiluminescence method and discussion of the quality objective of the NO₂ European Directive, *J.*
700 *Environ. Monit.*, 5, 529–540, <https://doi.org/10.1039/B302358C>, 2003.

701 Giordano, M. R., Malings, C., Pandis, S. N., Presto, A. A., McNeill, V. F., Westervelt, D. M., Beekmann, M., and
702 Subramanian, R.: From low-cost sensors to high-quality data: A summary of challenges and best
703 practices for effectively calibrating low-cost particulate matter mass sensors, *J. Aerosol Sci.*, 158,
704 105833, <https://doi.org/10.1016/j.jaerosci.2021.105833>, 2021.

705 Gramsch, E., Oyola, P., Reyes, F., Vásquez, Y., Rubio, M. A., Soto, C., Pérez, P., Moreno, F., and Gutiérrez, N.:
706 Influence of Particle Composition and Size on the Accuracy of Low Cost PM Sensors: Findings From
707 Field Campaigns, *Front. Environ. Sci.*, 9, 2021.

708 Grégis, F.: On the meaning of measurement uncertainty, *Measurement*, 133, 41–46,
709 <https://doi.org/10.1016/j.measurement.2018.09.073>, 2019.

710 Hofman, J., Nikolaou, M., Shantharam, S. P., Stroobants, C., Weijs, S., and La Manna, V. P.: Distant calibration
711 of low-cost PM and NO₂ sensors; evidence from multiple sensor testbeds, *Atmospheric Pollut. Res.*, 13,
712 101246, <https://doi.org/10.1016/j.apr.2021.101246>, 2022.

713 JCGM. International vocabulary of metrology - Basic and general concepts and associated terms.
714 [https://www.bipm.org/documents/20126/2071204/JCGM_200_2012.pdf/f0e1ad45-d337-bbeb-53a6-](https://www.bipm.org/documents/20126/2071204/JCGM_200_2012.pdf/f0e1ad45-d337-bbeb-53a6-15fe649d0ff1?version=1.15&t=1641292389029&download=true)
715 [15fe649d0ff1?version=1.15&t=1641292389029&download=true](https://www.bipm.org/documents/20126/2071204/JCGM_200_2012.pdf/f0e1ad45-d337-bbeb-53a6-15fe649d0ff1?version=1.15&t=1641292389029&download=true) (accessed on 20 Dec 2021). 2012.

716 Jiao, W., Hagler, G., Williams, R., Sharpe, R., Brown, R., Garver, D., Judge, R., Caudill, M., Rickard, J., Davis,
717 M., Weinstock, L., Zimmer-Dauphinee, S., and Buckley, K.: Community Air Sensor Network
718 (CAIRSENSE) project: evaluation of low-cost sensor performance in a suburban environment in the
719 southeastern United States, *Atmospheric Meas. Tech.*, 9, 5281–5292, [https://doi.org/10.5194/amt-9-](https://doi.org/10.5194/amt-9-5281-2016)
720 [5281-2016](https://doi.org/10.5194/amt-9-5281-2016), 2016.

721 Karagulian, F., Barbieri, M., Kotsev, A., Spinelle, L., Gerboles, M., Lagler, F., Redon, N., Crunaire, S., and
722 Borowiak, A.: Review of the Performance of Low-Cost Sensors for Air Quality Monitoring, *Atmosphere*,
723 10, 506, <https://doi.org/10.3390/atmos10090506>, 2019.

724 Kirkham, H., Riepnicks, A., Albu, M., and Laverty, D.: The nature of measurement, and the true value of a
725 measured quantity, in: 2018 IEEE International Instrumentation and Measurement Technology
726 Conference (I2MTC), 2018 IEEE International Instrumentation and Measurement Technology
727 Conference (I2MTC), 1–6, <https://doi.org/10.1109/I2MTC.2018.8409771>, 2018.

728 Lewis, A. and Edwards, P.: Validate personal air-pollution sensors, *Nat. News*, 535, 29,
729 <https://doi.org/10.1038/535029a>, 2016.

730 Makar, P. A., Stroud, C., Akingunola, A., Zhang, J., Ren, S., Cheung, P., and Zheng, Q.: Vehicle-induced
731 turbulence and atmospheric pollution, *Atmospheric Chem. Phys.*, 21, 12291–12316,
732 <https://doi.org/10.5194/acp-21-12291-2021>, 2021.

733 Malings, C., Tanzer, R., Hauryliuk, A., Kumar, S. P. N., Zimmerman, N., Kara, L. B., Presto, A. A., and R.
734 Subramanian: Development of a general calibration model and long-term performance evaluation of low-

735 cost sensors for air pollutant gas monitoring, *Atmospheric Meas. Tech.*, 12, 903–920,
736 <https://doi.org/10.5194/amt-12-903-2019>, 2019.

737 Mari, L., Wilson, M., and Maul, A.: Measurement across the Sciences: Developing a Shared Concept System for
738 Measurement, Springer International Publishing, Cham, <https://doi.org/10.1007/978-3-030-65558-7>,
739 2021.

740 Masson, N., Piedrahita, R., and Hannigan, M.: Approach for quantification of metal oxide type semiconductor gas
741 sensors used for ambient air quality monitoring, *Sens. Actuators B Chem.*, 208, 339–345,
742 <https://doi.org/10.1016/j.snb.2014.11.032>, 2015.

743 Mead, M. I., Popoola, O. A. M., Stewart, G. B., Landshoff, P., Calleja, M., Hayes, M., Baldovi, J. J., McLeod, M.
744 W., Hodgson, T. F., Dicks, J., Lewis, A., Cohen, J., Baron, R., Saffell, J. R., and Jones, R. L.: The use of
745 electrochemical sensors for monitoring urban air quality in low-cost, high-density networks, *Atmos.*
746 *Environ.*, 70, 186–203, <https://doi.org/10.1016/j.atmosenv.2012.11.060>, 2013.

747 Molina, M. J. and Molina, L. T.: Megacities and Atmospheric Pollution, *J. Air Waste Manag. Assoc.*, 54, 644–
748 680, <https://doi.org/10.1080/10473289.2004.10470936>, 2004.

749 Morawska, L., Thai, P., Liu, X., Asumadu-Sakyi, A., Ayoko, G., Bartonova, A., Bedini, A., Chai, F., Christensen,
750 B., Dunbabin, M., Gao, J., Hagler, G., Jayaratne, R., Kumar, P., Lau, A., Louie, P., Mazaheri, M., Ning,
751 Z., Motta, N., Mullins, B., Rahman, M., Ristovski, Z., Shafiei, M., Tjondronegoro, D., Westerdahl, D.,
752 and Williams, R.: Applications of low-cost sensing technologies for air quality monitoring and exposure
753 assessment: How far have they gone?, *Environ. Int.*, 116, 286–299,
754 <https://doi.org/10.1016/j.envint.2018.04.018>, 2018.

755 Mueller, M., Meyer, J., and Hueglin, C.: Design of an ozone and nitrogen dioxide sensor unit and its long-term
756 operation within a sensor network in the city of Zurich, *Atmospheric Meas. Tech.*, 10, 3783–3799,
757 <https://doi.org/10.5194/amt-10-3783-2017>, 2017.

758 Peters, D. R., Popoola, O. A. M., Jones, R. L., Martin, N. A., Mills, J., Fonseca, E. R., Stidworthy, A., Forsyth,
759 E., Carruthers, D., Dupuy-Todd, M., Douglas, F., Moore, K., Shah, R. U., Padilla, L. E., and Alvarez, R.
760 A.: Evaluating uncertainty in sensor networks for urban air pollution insights, *Gases/In Situ*
761 *Measurement/Validation and Intercomparisons*, <https://doi.org/10.5194/amt-2021-210>, 2021.

762 Popoola, O. A. M., Stewart, G. B., Mead, M. I., and Jones, R. L.: Development of a baseline-temperature
763 correction methodology for electrochemical sensors and its implications for long-term stability, *Atmos.*
764 *Environ.*, 147, 330–343, <https://doi.org/10.1016/j.atmosenv.2016.10.024>, 2016.

765 Rai, A. C., Kumar, P., Pilla, F., Skouloudis, A. N., Di Sabatino, S., Ratti, C., Yasar, A., and Rickerby, D.: End-
766 user perspective of low-cost sensors for outdoor air pollution monitoring, *Sci. Total Environ.*, 607–608,
767 691–705, <https://doi.org/10.1016/j.scitotenv.2017.06.266>, 2017.

768 Sayer, A. M., Govaerts, Y., Kolmonen, P., Lipponen, A., Luffarelli, M., Mielonen, T., Patadia, F., Popp, T., Povey,
769 A. C., Stebel, K., and Witek, M. L.: A review and framework for the evaluation of pixel-level uncertainty
770 estimates in satellite aerosol remote sensing, *Atmospheric Meas. Tech.*, 13, 373–404,
771 <https://doi.org/10.5194/amt-13-373-2020>, 2020.

772 Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Field calibration of a cluster of
773 low-cost available sensors for air quality monitoring. Part A: Ozone and nitrogen dioxide, *Sens.*
774 *Actuators B Chem.*, 215, 249–257, <https://doi.org/10.1016/j.snb.2015.03.031>, 2015.

775 Sun, L., Westerdahl, D., and Ning, Z.: Development and Evaluation of A Novel and Cost-Effective Approach for
776 Low-Cost NO₂ Sensor Drift Correction, *Sensors*, 17, E1916, <https://doi.org/10.3390/s17081916>, 2017.

777 Thompson, M. and Ellison, S. L. R.: A review of interference effects and their correction in chemical analysis
778 with special reference to uncertainty, *Accreditation Qual. Assur.*, 10, 82–97,
779 <https://doi.org/10.1007/s00769-004-0871-5>, 2005.

780 Tian, Y., Nearing, G. S., Peters-Lidard, C. D., Harrison, K. W., and Tang, L.: Performance Metrics, Error
781 Modeling, and Uncertainty Quantification, *Mon. Weather Rev.*, 144, 607–613,
782 <https://doi.org/10.1175/MWR-D-15-0087.1>, 2016.

783 Williams, D. E.: Electrochemical sensors for environmental gas analysis, *Curr. Opin. Electrochem.*, 22, 145–153,
784 <https://doi.org/10.1016/j.coelec.2020.06.006>, 2020.

785 World Health Organization: WHO global air quality guidelines: particulate matter (PM_{2.5} and PM₁₀), ozone,
786 nitrogen dioxide, sulfur dioxide and carbon monoxide: executive summary, World Health Organization,
787 Geneva, 10 pp., 2021.

788 van Zoest, V., Osei, F. B., Stein, A., and Hoek, G.: Calibration of low-cost NO₂ sensors in an urban air quality
789 network, *Atmos. Environ.*, 210, 66–75, <https://doi.org/10.1016/j.atmosenv.2019.04.048>, 2019.

790 Yatkin, S., Gerboles, M., Borowiak, A., Davila, S., Spinelle, L., Bartonova, A., Dauge, F., Schneider, P., Van
791 Poppel, M., Peters, J., Matheussen, C., and Signorini, M.: Modified Target Diagram to check
792 compliance of low-cost sensors with the Data Quality Objectives of the European air quality directive,
793 *Atmos. Environ.*, 273, 118967, <https://doi.org/10.1016/j.atmosenv.2022.118967>, 2022.

794 Zucco, M., Curci, S., Castrofino, G., and Sassi, M. P.: A comprehensive analysis of the uncertainty of a
795 commercial ozone photometer, *Meas. Sci. Technol.*, 14, 1683–1689, <https://doi.org/10.1088/0957-0233/14/9/320>, 2003.

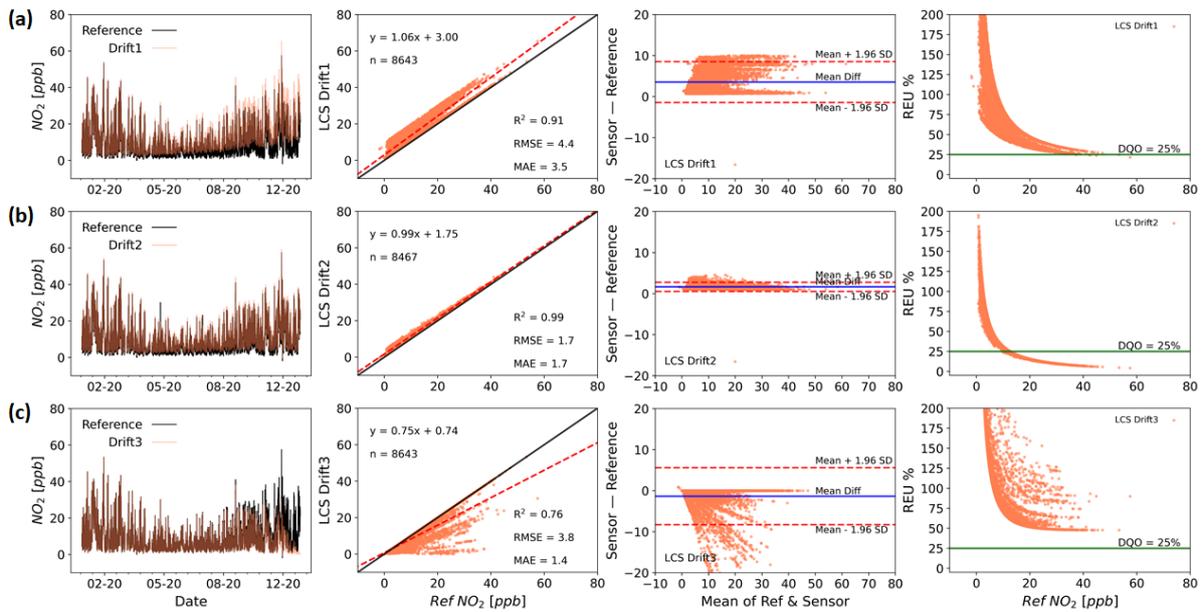
1 **Table S1. Research grade instrumentation used for this study.**

Analyte	Manchester			York
NO ₂	*Teledyne T500U (CAPS)	**Teledyne T200U (Chemiluminescence)	**Teledyne T200U (Chemiluminescence)	-----
O ₃	*Thermo 49i (UV photometry)	**Thermo 49i (UV photometry)	**2B (UV photometry)	-----
PM _{2.5}	-----	-----	-----	*Met One BAM 1020 (Beta attenuation)

2 *Instruments permanently deployed at the site.

3 **Instruments temporarily deployed as part of the QUANT study.

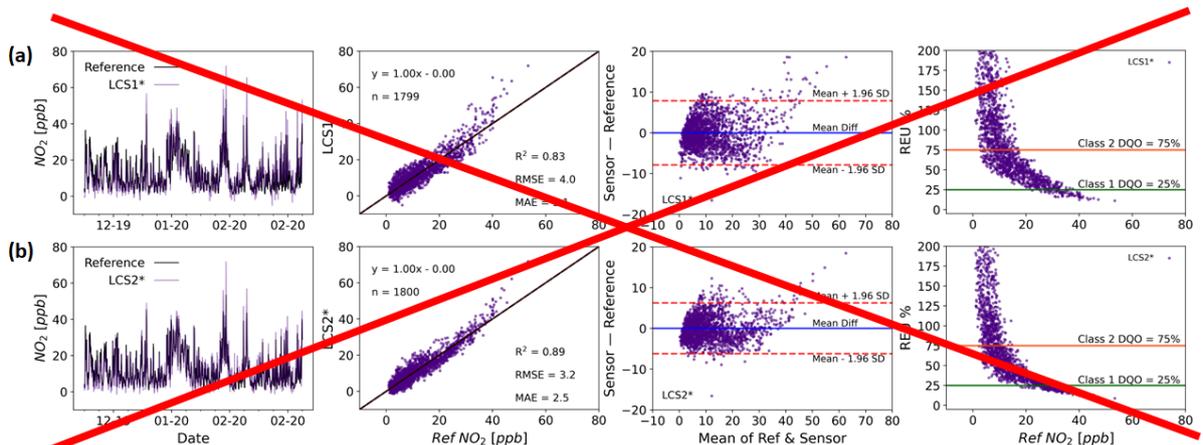
4



5

6 **Figure S1. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and**
 7 **REU (right panels, DQO for NO₂ = 25%) for baseline drift (a-panels), temperature interference drift (b-panels), and**
 8 **instrument sensitivity drift (c-panels) simulated errors.**

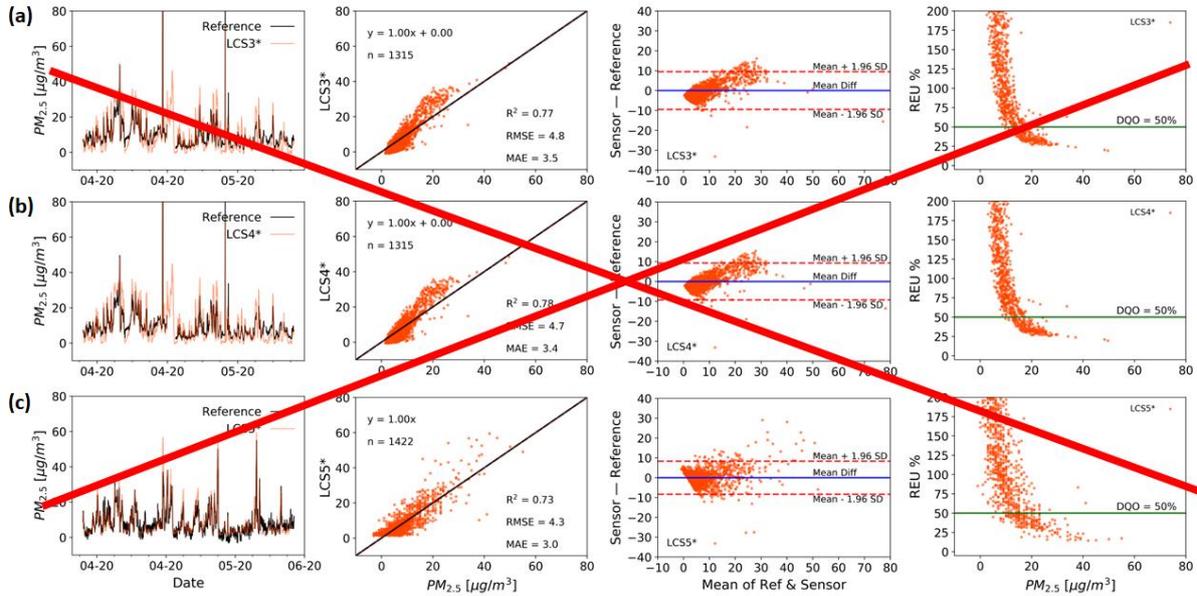
9



10

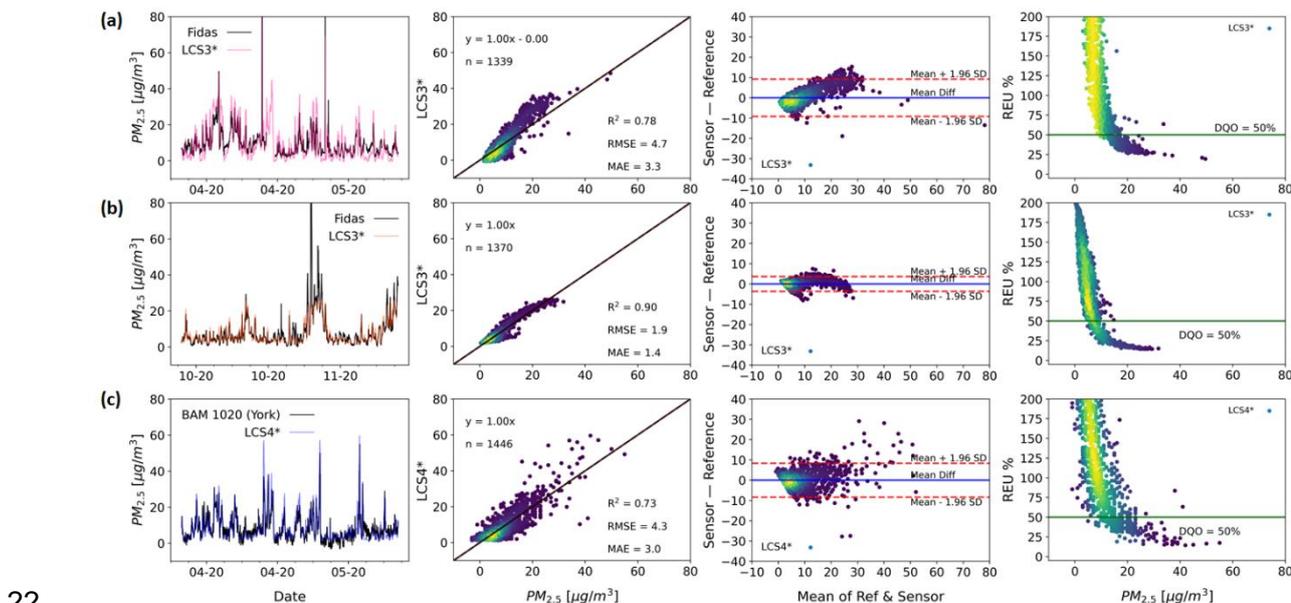
11 **Figure S2. Time series (left panels), regression plots (middle-left panels), Bland-Altman-plots (middle-right panels) and**
 12 **REU (right panels; NO₂ Class 1 DQO = 25% & Class 2 DQO = 75%) for NO₂ measurements by two LCS bias-corrected**
 13 **systems of different brands (panels a and b) in the same location (Manchester Supersite, December 2019 to February**
 14 **2020, 1hr time res):**

15



16 **Figure S3. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and**
 17 **REU (right panels, DQO for PM_{2.5} = 50%) for PM_{2.5} measurements by three LCS bias-corrected systems of the same**
 18 **brand (panels a, b and c) in different locations (April & May 2020, 1hr time res): an urban background (Manchester**
 19 **Supersite, a and b panels) and a roadside site (York, c panels).**

20



22 **Figure S2. Two bias corrected LCS systems (LCS3 & LCS4, same brand) measuring PM_{2.5} (Time res 1 h). While LCS3**
 23 **is shown for the same location (Manchester) but unfolded in two different seasons (a-panels: Apr to May 2020; b-**
 24 **panels: Oct to Nov 2020), LCS4 is at a different location (c-panels: York, Apr to May 2020). Time series (left panels),**
 25 **regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels; DQO_{PM_{2.5}} =**
 26

27 50%) are used to characterise the device's error structure. All but the time-series plots have been coloured
28 by density.
29

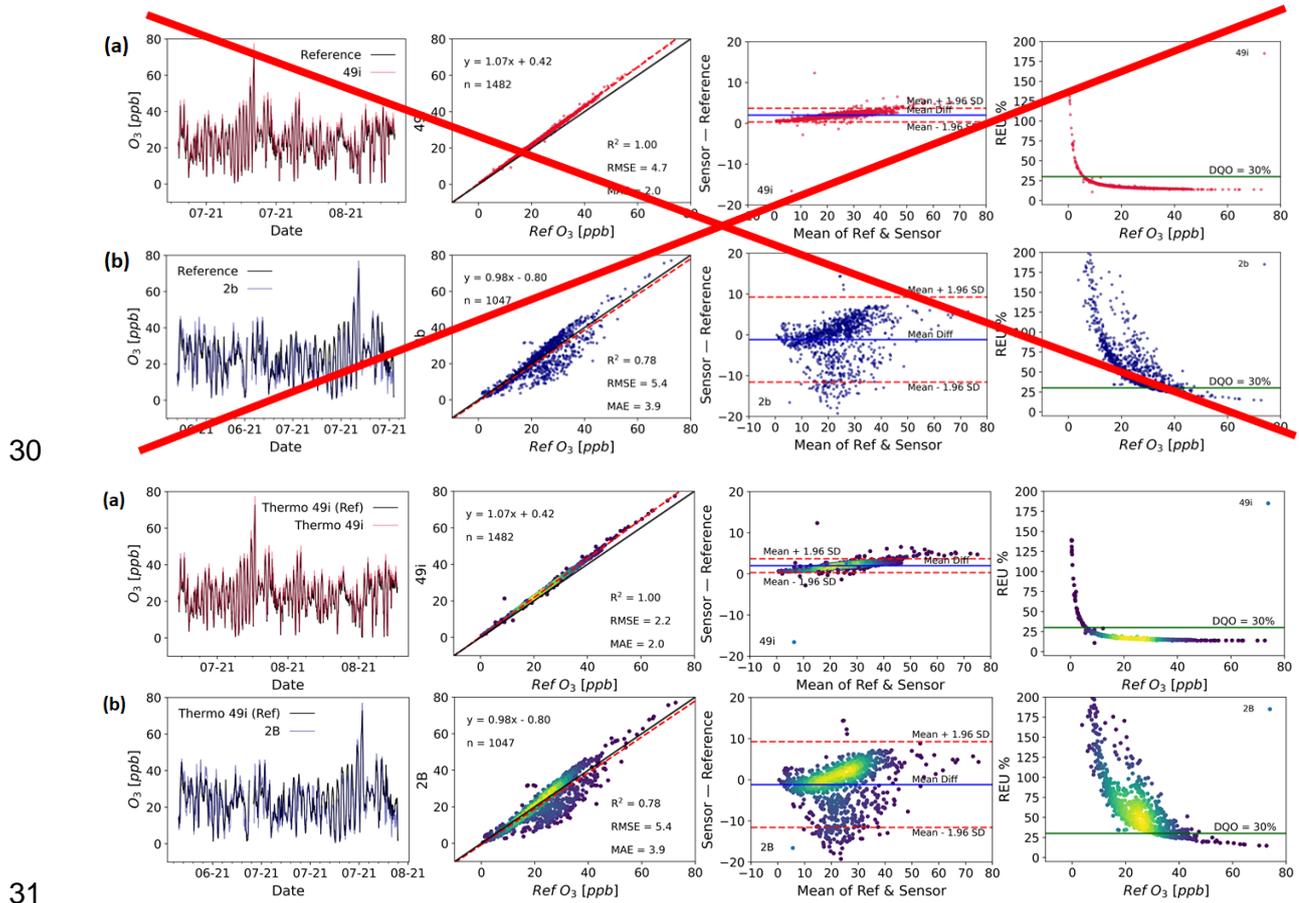


Figure S3. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels, DQO for O₃ = 30%) for two ozone research grade instruments (1hr time res): a Thermo 49i (a-panels, July & August 2021) and a 2B (b-panels, June and July 2021). All but the time-series plots have been coloured by data density.