

Response to reviewers - Long-term Evaluation of Commercial Air Quality Sensors: An Overview from the QUANT Study

Sebastian Diez, Stuart Lacy, Hugh Coe, Josefina Urquiza, Max Priestman, Michael Flynn, Nicholas Marsden, Nicholas A. Martin, Stefan Gillott, Thomas Bannan, and Pete Edwards.

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 improve the manuscript.

Below, you will see:

- reviewer comments in **bold**
- our responses are in regular type (Calibri font).
- cited text from manuscript and supplementary in Times New Roman font.

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**.

Reviewer#1

This study is a comprehensive, long-term study of a wide range of air sensors currently available in the market. Because of the collocation in time and space during the comparison, environmental variables are lessened thereby focusing the comparison on the performance of the air sensors relative to one another and in comparison, to a “reference” monitor. This adds great value to the study. The sensors were deployed and collocated in a real-world environment that they will most likely be used, so the conditions at which the sensors are compared were not biased because of the environment (i.e., compared to if it were collocated in a “pristine” environment, or under laboratory conditions). This study includes both gas and PM sensors which adds to the novelty of the study.

This paper is recommended for publication in AMT with revisions as outlined in this discussion.

Response: Thank you for your encouraging feedback and recognition of our study's comprehensive and practical approach to evaluating air sensors.

GENERAL COMMENTS

- The paper can benefit from tightening up language and being more succinct and concise in its statements.

Response: We appreciate the feedback on the need for clearer and more concise language throughout our manuscript. To address this, we have carefully reviewed our text, focusing on simplifying complex sentences, removing redundant phrases, and ensuring that each statement directly contributes to our argument or findings.

- A glossary of terminologies for commonly used yet widely misused or confused terms in the field (e.g., “sensor” vs “sensor systems” vs “sensing unit”, “manufacturer” vs “company”, “model”/“unit”/“type”/“device”) may be useful to the reader, also to help guide the authors in using consistent terminology all throughout the paper.

Response: We acknowledge the importance of clear and consistent terminology in our manuscript. We have chosen to complement the initial footnote instead of creating a glossary of terms, as we think this is a more practical use of these definitions. Please refer to section “1. Introduction” to check the usage of “sensor” and “sensor systems”. As for the terms “manufacturer” and “company”, please refer to the section “2. QUANT study design”. This approach offers immediate contextual explanations without diverting the reader’s attention from the main content. It also accommodates the diverse backgrounds of our audience by providing specific clarifications tailored to the context of this specific study. These footnotes now can be read:

¹ The term “sensor systems” refers to sensors housed within a protective case, which includes a sampling and power system, electronic hardware and software for data acquisition, analog-to-digital conversion, data processing and their transfer (Karagulian et al., 2019). Unless specified otherwise, the term “sensor” will be used as a synonym of “sensor systems”. Other alternative names for “sensor systems” used here are “sensor devices” (or “devices”), “sensor units” (or “units”).

² In a narrower sense, “sensor” typically denotes the specific component within a sensor system that detects and responds to environmental inputs, producing a corresponding output signal. To distinguish this from the broader use of “sensor” as equivalent to “sensor system” in our text, we will utilise alternative terms such as “detector”, “sensing element”, or “OEM” (original equipment manufacturer) when referring specifically to this component, thereby preventing confusion.

³ Throughout this article, the terms “manufacturers” and “company” are used interchangeably to refer to entities that produce, and/or sell sensor systems or devices. This usage reflects the industry practice of referring to businesses involved in the production and distribution of technology products without distinguishing between their roles in manufacturing or sales.

As for the term “type” it has been removed from the text and replaced with “model and brand” for more specificity.

- A major issue is the description and reliance on usual metrics like R2 in comparing two instrumental methods. With a goal of prediction and calibration in mind, R2 is an appropriate statistical metric; however, in plainly comparing the correlation (specifically: the concordance or agreement between two measurements), the authors are recommended to use more appropriate statistical metrics that measure concordance such as the concordance correlation coefficient. The paper can also greatly benefit a wider audience if the authors expanded on the statistical discussion and provide a separate discussion of the statistical metrics used, thus serving as a technical guidance that outlines metrics that can be used in such an intercomparison or calibration exercise.

Response: Thank you for your insightful comments. We acknowledge the limitations of relying solely on metrics like R2 for comparing instrumental methods. To address this, we have previously published a paper discussing the limitations of single-value metrics and advocated for a combined approach using visualisation tools and a range of metrics for a more nuanced analysis. For an in-depth discussion on this topic, we direct readers to Diez et al., 2022. Furthermore, the lead author of this manuscript has contributed to a chapter on performance metrics in an upcoming World Meteorological Organization report (to be published in May 2024), which delves into the advantages, limitations, and best practices regarding performance metrics in greater detail.

In this overview paper, we aim to provide readers with a holistic view of sensor performance by employing a variety of metrics and visualisation techniques. While we recognize the value of the Concordance Correlation Coefficient and its relevance, we also note its limitations. No single metric can fully capture all aspects of sensor performance, prompting our choice of widely recognized metrics like R2, RMSE, and MAE. This choice aims to facilitate comparisons with previous studies and standards in scientific literature.

In response to your suggestion, we have expanded our discussion on the statistical metrics used, adding a section in the supplementary materials (section “S5. Performance Metrics”) that serves as a summary discussion of the topic. Additionally, we have included a summary of recent standardisation efforts (Table S6), which we believe will be beneficial for end-users seeking guidance on metric selection for intercomparison or calibration exercises.

SPECIFIC COMMENTS

- **For the title: A more specific term than “Evaluation” can be used. Suggestions: “Intercomparison”; “Precision Analysis”**

Response: Thank you for your suggestion. We've chosen "Evaluation" as it most accurately encompasses the study's scope, which goes beyond comparison and precision analysis to include a comprehensive assessment of devices across various end-use scenarios. This term reflects our analysis's breadth, covering precision, accuracy, long-term performance, and adaptability to environmental conditions, thereby providing a holistic view of the sensors' capabilities for potential users.

- **The abstract is missing some key findings and results, e.g. how many air sensors and reference sensors were quantified, statistical metrics used to quantify the performance of the sensors, etc. For example, Lines 89-90 can be added to the abstract.**

Response: Thank you for your valuable feedback. In response, we have enriched the abstract to reflect the number of sensors and companies, as well as the range of metrics and visualisation tools utilised in our study.

- **A separate section on the methodology summarizing performance metrics used, and explaining each under a subheader, e.g., “Bias” as a subheader and explaining R2, RMSE, etc. under this heading would be useful to the reader and also makes this a good reference paper for intercomparison studies.**

Response: Thank you for your suggestion. As mentioned in response to an earlier comment, we have expanded our discussion on the metrics used. Please to section “S5. Performance Metrics” in the supplementary.

- **Section 3.3 explores reference instrumentation. Authors need to make what is meant by “reference”, e.g. a reference method designated by an authority (EU, US EPA, etc.) or a self-defined or agreed-upon reference method.**

Response: Thank you for your suggestion. The text has been complemented in order to clarify this point:

For an overview of reference and equivalent-to-reference instrumentation, as defined in the European Union Air Quality Directive 2008/50/EC (hereafter referred to as EU AQ Directive), at each site, please refer to Section S2 (Table S1). For details on the quality assurance procedures applied to the reference instruments, see Table S2.

Table S1 has also been updated to reflect this.

- **In one section, the manufacturers / models of air sensors were referred to; however, in Figure 10, it was anonymized. What was the rationale? Can it be consistently anonymized or named? And if not, explain why and make sure that the transition is clearly explained within each section.**

Response: Thank you for highlighting this aspect. Our initial decision to anonymize the names of companies in certain figures, including the original version of Figure 10, aimed to center the discussion on the generalizable features of sensor technologies rather than on specific findings

tied to individual manufacturers. This aimed to minimise potential biases and promote a comprehensive understanding of the technology in question. Nevertheless, acknowledging the significance of transparency and in response to constructive feedback, we have updated Figure 10 to reflect the names of the companies. Furthermore, we have elaborated on this rationale within the “3. Results and discussion” section to ensure a clear and consistent explanation of our approach to anonymization versus explicit naming:

To highlight broad implications and insights into sensor technology, rather than focusing on the performance of specific manufacturers, figures illustrating brand-specific features have been anonymized. This is intended to prevent potential bias and encourage a holistic view of the data, ensuring interpretations remain focused on general trends rather than isolated examples.

In addition, we are preparing a series of articles that will delve into more granular aspects of our study, as outlined in the conclusion section. Moreover, the dataset has been made publicly available, enabling comprehensive scrutiny of sensor performance by the broader community. A data descriptor paper detailing the QUANT dataset has also been submitted and is currently under review; once published, it will provide users with full access to and understanding of the dataset, further enhancing transparency and facilitating research in this field.

- It is useful from a consumer perspective to mention which devices are available readily as-is (without add-ons) and/or which ones require customization from the manufacturer’s end. This can possibly be added to the summary in Table 1 and/or Supplementary S1, with a short reference (a sentence or two) in the main text.

Response: Thank you for your suggestion. The text description of tables 1 and 2 has been complemented in order to clarify this point:

Table 1. Main QUANT devices description. The 20 units, all commercially available and ready for use as-is, offered 56 gas and 56 PM measurements in total. For a detailed description of the devices see Section S31 in the Supp.

Table 2. The 23 WPS devices deployed at the Manchester supersite, all commercially available and ready for use as-is, provided 63 gases and 62 PM measurements in total. For a detailed description of the devices see the Section S43 in the Supp.

- Might be useful to add in the conclusions / recommendations section for future researchers: quantify inter-location variability.

Response: Thank you for your suggestion to include more recommendations in our conclusions section. In addition to our final paragraph which includes future work we intend to do to address research needs we have also added the sentence below which highlights the end-user need for impartial performance data, which researchers are in a unique position to address.

Ultimately, this work shows that sensor performance can be highly variable between different devices and end-users need to be provided with impartial performance data on characteristics such as accuracy, inter-device precision, long-term drift and calibration transferability in order to decide on the right measurement tool for their specific application.

- Might be useful to explain and emphasize (including in the abstract) why correction with satellite data was not explored in this study.

Response: We appreciate the reviewer's suggestion regarding the exploration of correction with satellite data. However, after careful consideration, we have decided not to include an explanation for the omission of satellite data correction in our study, both in the manuscript and the abstract. This decision is based on the focused scope of our research, which is the direct evaluation of commercial air quality sensors in urban environments. Our primary aim is to assess these sensors' performance and applicability in settings where ground-level monitoring provides the most immediate and relevant data for urban air quality assessment. Including satellite data, which typically offers broader spatial coverage but lacks the fine-scale resolution required for our analysis, would not align with the specific objectives of our study. Furthermore, we aim to maintain a concise and focused narrative that is directly relevant to our core findings and methodology. We believe this approach will serve our audience best, keeping the manuscript clear and streamlined.

- Can employ the terms “inter-device” and “inter-location” for succinctness of ideas.

Response: Thank you for your suggestion. We have reviewed our manuscript and incorporated these terms where appropriate to more precisely describe the variability in sensor performance across different units and sites.

Line by line comments

- Line 45: “cross-sensitivity” seems to be a term usually used in the medical context. It might be beneficial to define this term in this context, and differentiate it from “interference”. Levy Zamora (2022) used “cross-sensitivity” in the title of their article and Bitner (2022) defined it, so it might be helpful if these two comes as the first articles cited in this instance).

Response: We appreciate the reviewer's input and have included Bitner (2022) for further clarification on “cross-sensitivities”. However, to keep our overview concise and focused, we believe the addition of this reference, along with Cross et al. (2017) and Pang et al. (2018), adequately informs readers about sensor challenges without overextending on definitions.

- Lines 48-49: citations for temperature and humidity might be combined since they are usually explained in cited references in combination.

Response: Thank you for the input. We have made the change accordingly.

- Lines 63-68: This paragraph could benefit from differentiating “calibration” from “correction” and how these two terms are sometimes interchangeably used (albeit incorrectly). A reference to an article that explains this difference will also be helpful. The Liang (2022) paper cited explains some of these nuances, including mathematical equations for calibrations, but does not fully differentiate “correction” from “calibration”.

Response: Thank you for your valuable feedback. To clarify the distinction between "calibration" and "correction", we have refined the text to include a direct reference to VIM (2012). Now the text reads:

The calibration of any instrument used to measure atmospheric composition is fundamental to guarantee their accuracy (Alam et al., 2020; Long et al., 2021; Wu et al., 2022). Using out-of-the-box sensor data without fit-for-purpose calibration can produce misleading results (Liang & Daniels, 2022). An effective calibration **not only** involves identifying but also **compensating for estimated and correcting** systematic effects ~~errors~~ in the sensor readings, a process defined as a correction (for a detailed definition and differentiation of calibration and correction see JCGM, 2012).

- Lines 67-68: True for gases. Mention examples of acceptable calibration method(s) for PM?

Response: Thank you for your suggestion. We have modified the original text in this way:

For standard air pollution measurement techniques, calibration is often performed in a controlled laboratory environment (Liang, 2021), ~~or by sampling gas from a certified standard cylinder in the field. For PM, particles of known density and size are used, controlling the airflow conditions.~~ For example, for gases, a known concentration is sampled from a certified standard. Similarly, for PM, particles of known density and size are generated. Both gases and PM calibration are conducted under controlled airflow condition

- Lines 81-83: Another reason is that there are a lot of sensors/sensor systems with different configurations commercially available, and also individual sensing units are sold and can be “DIY”-ed—the market is diluted with many options and many different iterations of the same underlying technology with marginal differences.

Response: Thank you for your suggestion. We have modified the original text in this way:

This is largely due to the significant variability in both the number of sensors and the variety of applications tested, **compounded by the proliferation of commercially available sensors/sensor systems with different configurations.** ~~as well as the availability of highly accurate measurement instrumentation and/or regulatory networks to those outside of the atmospheric measurement academic field.~~ Furthermore, the access to highly accurate measurement instrumentation and/or regulatory networks remains limited for those outside of the atmospheric measurement academic field (e.g. Lewis and Edwards (2016) and Popoola et al. (2018)).

- Line 94: Clarify or add examples of “data products” e.g., APIs, mobile apps, etc.

Response: Thank you for your suggestion. The revised sentence now reads:

Furthermore, we tested multiple manufacturers' data products, **such as out-of-the-box data versus locally calibrated data**, for a significant number of these sensors to understand the implications of local calibration.

Section “2.3 Data collection, co-located reference data and data products” has also been updated (see also response to line 158 comment, page 10).

- Lines 105-106 and 116: Useful to add a subsection that describes the UK urban environment including seasonality, sources of pollution (transportation? Household commercial products use?) in the three locations (London, Manchester and York).

Response: We appreciate the suggestion and recognize the importance of seasonality and pollution sources on sensor performance. However, a detailed exploration of seasonal variations at the 3 sites extends beyond this study's scope. The need for such analysis, considering the uncertainties around the UK's environmental characteristics, motivated the

initiation of the [Integrated Research Observation System for Clean Air \(OSCA\)](#), the measurement component of which was underway during the QUANT study. OSCA's forthcoming outputs are expected to provide in-depth environmental insights, and these data will be used in future work for a more in depth study of sensor interferences etc. Nonetheless, we've updated the supplementary material to include a more detailed description on urban environment details and anthropogenic activities, for each site. Please refer to "S1. Co-location sites".

- Line 106: "replicates" or "units" are more appropriate terms than "duplicates" if you are talking of the units of the same model

- Line 109: define what is meant by "near real-time" in this context.

Response to the last two comments: Thank you for your suggestion. The text has been updated accordingly, and now reads:

Four ~~units duplicates~~ of five different commercial sensor devices (Table 1) were purchased in Sept 2019 for inclusion in the study, with the selection criteria being: market penetration and/or previous performance reported in the literature, ability to measure pollutants of interest (e.g. NO₂, NO, O₃, and PM_{2.5}), and capacity to run continuously reporting high time resolution data (1-15 min data) ideally in near real-time (*i.e.*, [available within minutes of measurement](#)) with data accessible via an API.

- Line 113: Were the units tested together before deploying separately? Clarify.

Response: Thank you for your comment. All the units were first deployed together in Manchester as stated in the original lines 122 and 123:

Initially, all the sensors were deployed in Manchester for approximately 3 months (mid-Dec 2019 to mid-Mar 2020) before being split up amongst the three sites (Fig. 1).

- Line 121: A sentence or two succinctly describing the sites will also be useful in this line. Then you can refer to the Supplementary.

Response: Thank you for your suggestion. We've slightly modified the current description for readers interested in more detailed information, as we have added the official web address describing each of the mentioned sites. We have also made a more thorough description of each of the sites at the supplementary (please see the response to an earlier comment regarding Lines 105-106 and 116).

- Lines 122-126: Consider moving up before Lines 113-121.

Response: Thank you for your suggestion. We moved the text according to the suggestion.

- Line 125: "inter-device consistency" may also be rewritten as "precision".

Response: Thank you for your suggestion. We acknowledge that "precision" could effectively communicate the aspect of measurement variability among devices. However, we believe "inter-device consistency" encompasses not only the precision aspect but also the reliability and stability of device measurements across various conditions and over time. Thus, we have opted for the initial term to convey the broader scope of our analysis more accurately.

- Lines 134-136: “vendors were invited to contribute multiple sensor devices throughout the WPS study”. How does a “sensor device” differ from a “sensor” or “sensor system”? Does this mean that manufacturers can contribute different sensor models? Also, does vendor = company = manufacturer? Note consistent terminology all throughout the manuscript (might be useful to have a glossary or footnote, like that for “sensor” and “sensor system” on page 2).

Response: We appreciate the suggestion to clarify terminology. Following the initial recommendation in “General comments” (page 2), we have carefully defined and differentiated these terms as footnotes in the “Introduction” and in “QUANT study design” sections. As for the word “vendor” (it appeared two times in the original text), it was replaced by the word “manufacturer”.

- Lines 139-141: Table 1. Does AQMesh AQM, Kunak AP, and SCS Prax have all of the sensors listed (from NO to PM10) in one unit? This table might benefit from a clarification (can be added to the Table caption). Also, as mentioned in a previous comment, add in the description if these are consumer-ready (eg already sold in the market as that unit), or customizations available from the manufacturer.

Response: Thank you for your suggestion. As stated in “Specific comments” (page 5), the text description for tables 1 and 2 has been adapted to help clarify this point.

- Lines 148-149: I understand that PurpleAir does not have a mobile data connection, only WiFi, but WiFi was not good in the location so you opted to download the data from the device memory instead. The text can be enhanced by better explaining the issue as described. (i.e., differentiating from WiFi and mobile data connectivity)

Response: Thank you for your suggestion. We updated the text, and now reads:

~~PurpleAir units were exempt from this due to a lack of mobile data connection and poor internet signal at the sites; instead, readings were locally collected and manually uploaded.~~ Unlike other brands that utilize mobile data connections, PurpleAir sensors rely on WiFi for data transmission. Due to poor internet signal at the sites, we locally collected and manually uploaded readings for these units.

- Line 150: and harmonize? In the methodology section, it might be useful to mention that temporal and spatial scales of the sensor systems was important to match, thus aggregation and harmonization was necessary. How was incomplete data treated? Were there imputed data? Might be useful to add it in the supplementary.

Response: Thank you for your suggestion. We've clarified in our manuscript that “standardize” includes formatting, aggregating, and ensuring data compatibility across devices, without altering data characteristics. This revision is now reflected:

Minor pre-processing was applied at this stage, including temporal harmonisation to ensure that all measurements had a minimum sampling period of 1-minute, ensuring consistency in measurement units and labels, and coercing into the same format to allow for full compatibility across sensor units.

On the other hand, incomplete data didn't receive special consideration as was originally stated. We've now expanded this in the text to avoid confusion:

No additional modifications to the original measurements were applied; missing values were kept as missing and no additional flags were created based on the measurements beyond those provided by the manufacturers. ~~No outlier checks or data modifications were applied at this stage.~~

- Line 150: by data format, do you mean datetime / time and date?

Response: Thank you for your inquiry. As per the response to the previous question, the phrase “data format” has been replaced by an explanation of the specific pre-processing steps applied.

- Line 158: “calibrated data products”: is this referring to API? Measurements? As with my previous comment – clarify what “data products” mean.

Response: Thank you for your comment. As described in our manuscript, this term refers to the various versions of data provided by manufacturers, reflecting different stages of calibration and adjustment based on colocated reference data. We believe this description, along with the supplementary material, offers a comprehensive insight into what is encompassed by this term. However, we have decided to complement the text, for clarity:

However, those who did were expected to create and submit calibrated data products, subsequently named as “out-of-box” (initial data product), “cal1” (first calibrated product), and “cal2” (second calibrated product). This differentiation highlighted the varying degrees of engagement and application of the reference data by different manufacturers. Figures S2 and S3 (section S3 and S4 respectively) show a time-line of the different data products.

- Lines 160-166: What is cal1? cal2? Clearly define / describe these in the text and/or supplementary. This section may benefit from a subsection explaining / describing these.

Response: see previous response.

- Lines 170-174: Is this a caveat / weakness of using these statistical metrics used herein (R2, MAE, etc)? What is the alternative? I suggest concordance (agreement) metrics, such as the Concordance Correlation Coefficient: See Lin, Biometrics (1989): <https://doi.org/10.2307/2532051>. The reader might also benefit from a separate subsection and/or supplementary section describing the metrics or including a glossary of the metrics used.

Response: Thank you for your comment. While we agree that the Concordance Correlation Coefficient is a valuable metric for assessing inter-rater agreement, it still suffers from the same limitation as described in the text, wherein over-reliance upon a single metric can obscure the full picture. Instead, we opt for a more holistic assessment comprising multiple facets of a sensor’s performance. Also, we’ve focused on more commonly used metrics within both the scientific community and technical guidelines for sensor evaluation. This choice aims to provide a comprehensive assessment of sensor performance and facilitate comparisons with existing guidelines and research findings. We’ve updated the text to better convey this:

~~Furthermore, the overreliance on global performance metrics, such as R^2 (i.e., the Coefficient of Determination), RMSE (i.e., the Root Mean Squared Error), and MAE (i.e., the Mean Absolute Error) is an important issue when assessing sensors. While these metrics provide a general understanding of sensor~~

~~performance, they can be limiting or even misleading, restricting a comprehensive understanding of the error structure and the measurement information content (Diez et al., 2022).~~ Furthermore, the overreliance on global performance metrics is a significant concern in sensor assessment. The Coefficient of Determination (R^2), Root Mean Squared Error (RMSE), and Mean Absolute Error (MAE) are among the most popular single-value metrics for evaluating sensor performance, alongside others (e.g., the bias, the slope and intercept of the regression fit). However, while single-value metrics offer an overview of performance, they can be limiting or misleading. They condense vast amounts of data into a single value, simplifying complexity at the expense of a nuanced understanding of error structures and information content (Diez et al., 2022), potentially overlooking critical aspects of sensor performance (Chai & Draxler, 2014). Visualisation tools (such as Regression plots, Target plots, and Relative Expanded Uncertainty plots) complement these metrics, allowing end users to identify relevant features, which could be beyond the scope of global metrics. For further discussion on metrics and visualisation tools for performance evaluation, readers are directed to Diez et al. (2022).

As for the suggestion section describing the used metrics, please refer to response to previous suggestions (page 3).

- Line 183: “multiple devices of the same type” when you mean “type” do you mean similar underlying principles of measurement? Model? Be consistent in terminology. Also, it might be useful to cite an example of which devices you are considering a same “type”, e.g. AQM and Clarity—are these of the same “type” as described?

Response: Thank you for your suggestion. As stated in a previous response from “General comments” (see page 2), the term “type” has been replaced with “model and brand” for more specificity.

- Lines 190-196. See also deSouza (2023): An analysis of degradation in low-cost particulate mater sensors <https://doi.org/10.1039/D2EA00142J>

Response: Thank you for your suggestion.

- Lines 202-204: Good point.

Response: Thanks!

- Line: 217: 75% inclusion criteria is common—but perhaps not for readers not familiar with this data type. Readers might benefit from a citation, explanation in methodology or supplementary. Suggested section to add it in: Section 2.1, lines 150-151.

Response: Thank you for your suggestion. To address the comment and enhance clarity, we have specifically mentioned the data inclusion criterion of 75% at the end of the last paragraph in the “3. Results and discussion” section:

The following sections aim to provide an overview of the data and provide initial findings, with a focus on those that are most relevant to end-users of these technologies. All metrics and plots presented here are based on 1-hour averaged data. Unless otherwise specified, a data inclusion criterion of 75% was uniformly applied across our analyses to ensure the reliability and representativeness of the results. This threshold

aligns with the EU AQ Directive, which mandates this proportion when aggregating air quality data and calculating statistical parameters.

- Line 225: Clarify: did you mean closer together spatially / physical location? How “close” is close?

Response: Thank you for seeking clarification. “Closer together” refers to the clustering of sensor performance data points in the plot, indicating improved measurement consistency among sensors from the same manufacturer after calibration. We have updated the text to avoid this confusion:

Secondly, it can help to improve within-manufacturer precision ~~by grouping sensor systems from the same company closer together.~~, as evidenced by sensor systems from the same company grouping more closely as the right plot in Fig. 4 shows.

- Line 232: could benefit more from a further explanation of the bias-variance tradeoff.

Response: Thank you for the suggestion. We have updated the text to be more descriptive.

For the out-of-the-box data, these regions are noticeably larger than in the calibrated results for most manufacturers, suggesting that colocation calibration has helped to tailor the response of each device to the specific site conditions. ~~This is reinforced by the cRMSE component reducing by a greater extent than the MBE; in the terminology of machine learning, the calibration has helped reduce the variance portion of the bias-variance trade-off.~~ This observation suggests that colocation calibration effectively improves each device's response to particular site conditions. This improvement is underscored by the more substantial reduction in the cRMSE component compared to the MBE. The cRMSE, representing the portion of error that persists after bias removal, essentially measures errors attributable to variance within the data space. In the context of out-of-the-box data, this “data space” spans all potential deployment locations used by manufacturers for initial calibration model training (i.e., before shipping the sensors for the QUANT study), thus exhibiting high variability. However, applying site-specific calibration significantly narrows this variability, leveraging local training data to minimise variance.

- Section [3.3. and Supplementary S4. Cite the authorities that consider the instruments mentioned as reference. e.g., are they considered “reference” because they are listed in an EU directive or US EPA documentation? If so, cite these. If not, provide a rationale or a citation as to how these instruments were categorized as “reference”.

Response: Thank you for your suggestion. This point was already addressed in the “Specific comments” (refer to page 3).

- Lines 241-242: Expound on the significance and advantages of REU as opposed to the other metrics.

Response: Thank you for the suggestion. We have updated the text accordingly:

For applications where it is important to understand how calibrations impact lower or higher percentiles, considering other metrics or visual tools would be advisable. An example of this is the absolute and Relative

Expanded Uncertainty (REU, defined by the Technical Specification CEN/TS 17660-1:202). Unlike the more commonly used metrics such as R^2 , RMSE, and MAE, which measure performance of the entire dataset, the REU offers a unique “point by point” evaluation, enabling its representation in various graphical forms, such as time series or concentration space (for the REU mathematical derivation, refer to section “S5. Performance Metrics”). The REU approach, also incorporates the uncertainty of the reference method into its assessment, highlighting the intrinsic uncertainty present in all measurements, including those from reference instruments. This consideration of reference uncertainty is crucial for a holistic understanding of sensor performance and calibration effectiveness. For a comprehensive discussion on this, refer to Diez et al. (2022).

- Lines 268-269: clarify / reiterate the said minimum requirements in this text.

Response: Thanks for this request. The modified text now reads:

The REU demonstrates that, under these circumstances, an instrument designated as a reference does not meet the minimum requirements ($REU \leq 15\%$) set out by the Data Quality Objectives (DQOs) of the EU ~~AQ Air Quality Directive 2008/50/EC~~.

- Lines 283-284: R2 can many times be subjective, e.g. how can we say that an R2 of 0.87 means “it does not fully agree” and a slope of 0.80 is considered a pronounced bias? This is where the definition of concordance and using concordance metrics might be useful. If two measurements are concordant (in agreement), then slope is expected to be unity (=1). Also, a high R2 does not necessarily mean agreement between the two instruments. Also, clarify what is meant by “limiting the linearity”. The authors are cautioned against using R2 in quantifying agreement between two instruments that are being compared.

Response: We appreciate the feedback and understand the nuances involved in interpreting these metrics within the specific context of air quality studies. The interpretation of any single-value metric, like R^2 , RMSE, MAE—including the CCC—inevitably involves a degree of subjectivity, relying on the analyst's expertise to discern their significance within the context of the study. These statistical measures compress a vast amount of data into a singular value, potentially obscuring the broader picture. In our study we opted for a holistic assessment trying to encompass multiple facets of a sensor's performance, integrating both quantitative metrics and visual analyses to offer a comprehensive evaluation rather than placing sole emphasis on any single metric

Regarding the interpretation of an R^2 value of 0.87, we welcome the correction that R^2 doesn't directly measure the degree of agreement between 2 sensors and have updated the text to be more precise by referring to the linearity. We have described this as a “strong association”, which albeit a subjective term, is justified by a Pearson's $R > 0.90$, which is in-line with conventional usage (for example <https://pubs.acs.org/doi/10.1021/acsearthspacechem.8b00079>).

In response to the comment regarding the slope of 0.80, it appears there may have been a misunderstanding. Our original manuscript stated that such a slope “is not considered a very pronounced bias”. We value this opportunity to clarify our intent and have refined our wording for greater clarity and to avoid any potential ambiguity. The modified text reads:

To illustrate these differences in practice, Fig. 6 compares these two equivalent-to-reference PM_{2.5} measurements obtained with a BAM (AURN York site, located on a busy avenue), and a FIDAS unit specifically installed for QUANT. During this specific period, they ~~show a strong linear association do not fully agree~~ ($R^2 = 0.87$). ~~Although the bias is not extremely pronounced (slope=0.80), the FIDAS measurements are, on average, systematically lower compared to BAM. Despite a not very pronounced bias (slope=0.80), the dispersion of points around the best fit line is noticeable, limiting the linearity of the FIDAS compared to the BAM.~~

- Lines 288-289: Specify criterion stipulated by EU DQOs

Response: Thanks for this request. The modified text now reads:

In the hypothetical case that the BAM were to be considered the reference method (arbitrarily chosen for this example as it is the current instrument at the AURN York site) when assessing the FIDAS under these test conditions, it would only meet the criterion stipulated by the EU DQOs for indicative measurements ($REU \leq 25\%$ for PM_{2.5}), but not for fixed (i.e., reference) measurements ($REU \leq 50\%$ for PM_{2.5}).

- Line 308: “reasonably consistent” – reasonably is subjective and qualitative. Suggest dropping the word, or provide a percentage of the time that the RMSE is consistent (e.g. provide a qualitative measure).

Response: Thank you for highlighting this. In response to your suggestion, we have clarified the statement and the revised text now reads:

The RMSE remains reasonably consistent ($\text{range } 2.27 \text{ to } 3.47 \text{ ppb}$) between the devices across the periods and locations

- Line 312: “local conditions”: give or name some examples. Do you mean weather conditions? Traffic? etc.

Response: Thanks for this request. The modified text now reads:

The precise cause of this change is not immediately evident and will be the focus of a follow-up study, but could be due to changes in local conditions (e.g., weather, emissions, etc.) impacting sensor calibration and/or differences in actual PM_{2.5} sources and particle characteristics at the sites (Raheja et al., 2022).

- Line 334: quickly define (add a phrase) that describes “overfitting”

Response: Thank you for your suggestion. Accordingly, the text has been revised to include a brief definition of overfitting:

Furthermore, it is important to consider that excessive post-processing may lead to overfitting—a situation where a model excessively conforms to specific patterns in the training data, resulting in poor performance on new, unseen data (Aula et al., 2022).

- Line 336: “linear correction” – “linear regression” might be the appropriate term.

Response: Thank you for your suggestion. In the context of our study, “linear correction” was deliberately chosen to describe the application of a simple linear adjustment to sensor data (e.g., zero and span correction).

- Line 347: RMSE also showed seasonality.

Response: Thank you for your comment. You're correct that RMSE exhibits some seasonality. Our analysis primarily aimed to emphasise the behaviour of its orthogonal components: MBE and cRMSE (i.e., RMSE is as a function or consequence the latter). We have slightly modified the figure caption in order to acknowledge this point:

Figure 9. Seasonal variation of error (as RMSE, red line) of one of the systems belonging to the Main QUANT...

- Lines 345-351: Add more explanations about the seasonality. Add recommendations.

Response: Thank you for your valuable feedback. We believe that the existing sections of this overview paper, particularly 3.5 and 3.6, already offer a detailed discussion on the temporal nature of sensor errors, including seasonality aspects (lines 345-351), and outline practical recommendations for improving sensor data quality, highlighting methods like NO₂ bias correction using diffusion tubes (lines 388-396). Recognizing the importance of a deeper investigation, our manuscript also outlines future studies (lines 463-466) dedicated to a more thorough exploration of seasonality effects and the development of detailed recommendations. This is part of our broader commitment to improving the understanding and application of sensor performance, with successive studies planned to delve into these aspects further.

- Line 355: Clarify what a “1-day slide” means – it can be added in the supplementary or a quick description in the figure caption.

Response: Thank you for your suggestion. We have updated the description in the figure caption, and now reads:

Figure 9. Error (as RMSE, red line) of one of the systems belonging to the Main QUANT, decomposed into cRMSE (in blue) and MBE (in yellow) estimated based on a 40-day (aligning with the sample size recommendation by the CEN/TS 17660-1:2021 standard for on-field tests) moving window approach with a 1-day slide (i.e., advancing the calculation 1 day at a time) (1-day slide) moving window. Panel a) is for O₃ measurements, and panel b) is for NO₂ (April 2020-Oct 2022). Panel c) is also for NO₂, this time showing the effect of a linear correction using diffusion tubes (see next section for more details).

- Lines 363-366: This can be said more succinctly. Also, what sort of information can be provided? Be specific, based on your results so far—what sort of information can you recommend be provided?

Response: Thank you for your feedback. In response, we have revised this section to make it more concise and to explicitly outline the type of information that should be provided. The revised text now reads:

~~In order to realise the potential of air pollution sensor technologies, end-users need to be provided with the information required to critically assess the strengths and weaknesses of potential candidate sensor devices, ideally in an easy to access and interpret manner.~~ To realise the potential of air pollution sensor

technologies, end users need to align their specific measurement needs with the capabilities of available devices. Achieving this necessitates access to unbiased performance data, such as long-term stability and accuracy across varying conditions, ideally in an easy-to-access and interpret manner.

- Line 373: Inconsistency in the use of the term “sensor”, “system” and “sensor system”.

Response: See previous response on “General comments” (see page 2).

- Lines 377-379: This discussion can benefit from a more detailed explanation of the tiers (Classes) assigned and what was the basis of the assignment to different classes. The figure caption for Figure 10 offers an explanation, which should be repeated and explained in more detail in-text.

Response: Thank you for your suggestion. While we value the importance of detailed explanations regarding the classification into different tiers, we aimed to provide a broad overview within the scope of this manuscript, directing readers to specific documents (i.e., EU AQ Directive and CEN/TS 17660-1:2021) for comprehensive details on class assignments and the rationale behind them. This approach ensures that our manuscript remains focused on providing an accessible overview of the QUANT study, while still making in-depth information readily available to those interested in delving deeper. We have however modified this section, and now reads:

Both REU and DC are key criteria within the EU scheme (EU 2008/50/EC) for evaluating the performance of measurement methods, and are complemented by the CEN/TS 17660-1:2021 specifically for sensors. The latter ~~This~~ document defines three different sensor system tiers. Class 1 NO₂ sensors, bounded by the green rectangle (REU < 25% and DC > 90%), offer higher accuracy than Class 2 sensors (REU < 75% and DC > 50%), delimited ~~highlighted~~ by the red rectangle (Class 3 sensors have no set requirements). Presenting the REU and DC data like in Fig. 10 ~~this~~ helps users anticipate the performance of sensor systems —under the assumption that all sensors from the same brand will behave similarly in equivalent environmental conditions— providing more insight into selecting the appropriate instrument for a given project or study.

- Line 388 onwards can benefit from a separate subheader / subsection.

Response: Thank you for your suggestion. After careful consideration, we would prefer to maintain the current structure of Section 3.6. As the intention of this overview paper is to summarise the study and some of the key findings as well as point potentially interested readers towards our unique dataset, we are keen to avoid it being overly exhaustive and long. We also feel that our manuscript has thematic coherence, as Section 3.6 transitions from discussing sensor performance nuances to practical implications for end-use applications. The seamless flow into "Informing end-use applications" is intentional, reflecting our comprehensive approach to presenting both technical analysis and its practical implications in a unified narrative.

- Lines 390-391: Specify an example of “simpler methods”

Response: Thank you for highlighting the need for clarification. We intended to convey that, depending on the application, there might be other feasible alternative methods for bias, rather

than “simpler” from a technical point of view. To clarify this in our manuscript, we have refined the text as follows:

~~Depending on the application, simpler methods could also be available to reduce the magnitude of the changing bias, and thus significantly improve the accuracy of an individual sensor system, but also that of broader sensor networks. For the case shown in Fig.9b, one possible way to do this would be using supporting observations of NO₂ made via diffusion tubes.~~ Depending on the application and available options, users can access alternative methods to reduce bias, thus enhancing the accuracy of sensor systems and networks. For example, “Indicative methods”, as defined by the EU AQ Directive, such as diffusion tubes (e.g., NO_x, SO₂, VOCs, etc.), can be an option. Specifically, our study leverages diffusion tube data for NO₂, illustrating one effective approach to bias correction using supporting observations, as exemplified in Fig. 9b.

- Line 393: explain more by what instrumental method is NO₂ measured/detected from these diffusion tubes. Cite a reference as well.

Response: Thank you for raising this question. We would like to clarify that a detailed explanation of the instrumental method used for NO₂ detection via diffusion tubes was already included in the supplementary material of our original submission, specifically in the section “S6. NO₂ Diffusion Tubes” (renamed after revisions as “S7...”)

- Line 411: “change points”: does this mean inflection points? Periods of biggest slopes? Peak concentration periods? Consider changing verbiage.

Response: Thank you for your suggestion. To clarify, we've added a brief description of change points within the statistical field of change detection to our manuscript:

An example of this from the QUANT dataset is the use of sensor devices to successfully identify change points in a pollutant’s concentration profile. ~~These are points in time where the parameters governing the data generation process are identified to change, commonly the mean or variance, and can arise from human-made or natural phenomena (Aminikhanghahi and Cook, 2017).~~

- Lines 414-415: Is it applied in this paper? If so, this line should be described in the methodology and explained further.

Response: Thank you for your comment. The mention of change point analysis was intended to illustrate the potential applications of the QUANT dataset, rather than to detail a methodology applied within this specific study. To clarify and prevent any misunderstanding, we have revised our manuscript, and now reads:

Determining when a specific pollutant has changed its temporal nature is a challenging task as there are a large number of confounding factors that influence ~~atmospheric concentrations a pollutant’s concentration at a specific point in time~~, including but not limited to seasonal factors, environmental conditions (both natural and arising from human behaviour), and meteorological factors. ~~This challenge has lead to several “deweathering” techniques being proposed in the literature (Carslaw et al., 2007; Grange and Carslaw, 2019; Ropkins et al., 2022).~~ While change point detection is highlighted here as a promising application of

sensor data, it represents just one of many potential methodologies that could be explored with the QUANT dataset.

We have also updated the text explaining the methodology applied to the QUANT dataset:

A state-space based deweathering model was applied to NO₂ concentrations measured from the sensor systems that had remained in Manchester throughout 2020 to remove these confounding factors, with the overarching objective to identify whether the well-documented reduction in ambient NO₂ concentrations due to changes in travel patterns associated with COVID-19 restrictions could be observed in the low-cost sensor systems. To provide a quantifiable measure of whether a meaningful reduction had occurred, the Bayesian online change-point detection (Adams & MacKay, 2007) was applied. Of the 8 devices that measured NO₂, clear change points corresponding to the introduction of a lockdown were identified in 2 (Fig.11), demonstrating the potential of these devices to identify long-term trends with appropriate processing, even with only 3 months of training data.

- Line 421: Expound what “unsupervised analysis” means in this context. General verbiage related to machine learning may sometimes be unnecessary to use in this text, and can be avoided, because fundamental/rudimentary statistical metrics (as opposed to complex “black-box” machine learning algorithms) are used.

Response: Thank you for your comment and suggestion. In this context, "unsupervised analysis" refers to the application of statistical techniques without explicit guidance or labelled data. For this case, it means without directly comparing the modelled output (estimated change-point) to the actual measured outcome (e.g., date of Covid lockdown). Acknowledging this can cause misleading interpretations, this term has been removed with the text.

- Line 422: consistency in terminology. Do the authors mean “sensor system” when they mention “devices”?

Response: See previous response on “General comments”

- Line 433: “..use of these devices has been primarily limited...” I would disagree, because consumers and many users still use these devices (sensor systems) and they aren’t necessarily limited by accuracy concerns, e.g. many users are willing to accept a large margin of error for awareness purposes.

Response: Thank you for your perspective. We acknowledge that despite concerns over data quality, there is a significant user base that utilises these sensor systems for various purposes, including general air quality awareness. We had intended the limitations mentioned in line 433 to be contextualised by the preceding statement about the potential of low-cost sensors to enhance air pollution management and understanding, but have edited the sentence to provide clarity. This now reads:

Large-scale uptake in the use of these devices [for air quality management](#) has, [however](#), been primarily limited by concerns over data quality and a general lack of a realistic characterisation of the measurement uncertainties making it difficult to design end uses that make the most of the data information content.

- Line 439: suggested addition: (limitations in) technical ability in post-processing of data

Response: Thank you for the suggestion. To clarify this point, we have modified the text in this way:

A challenge with the use of sensor-based devices is that many of the end-use communities do not have access to extensive reference-grade air pollution measurement capability (Lewis & Edwards, 2016), or in many cases, expertise in making atmospheric measurements [or the technical ability for data post-processing](#).

- Lines 460-461: Will this be done by the authors in a future study, or is this a call/recommendation for other researchers?

Response: Thank you for your inquiry. The future studies mentioned in the concluding paragraphs of our manuscript (lines 460-469) are currently being undertaken by our team and will be detailed in forthcoming publications.

- Line 466: suggestion for a future study: explore different VOC-NO_x regimes (see Wennberg, ES&T Air: <https://doi.org/10.1021/acs.air.3c00055>)

Response: Thanks for your suggestion!

TECHNICAL CORRECTIONS

Grammatical, Typographical, Figure and Formating comments throughout the text

- Note the usage of “data” as a plural noun, e.g. “data were” rather than “data was”

Response: Thank you for pointing this out. This is now corrected.

- “Manufacturer” rather than “Company” might be a more descriptive noun for the intended usage.

Response: Thank you for your suggestion. In a previous comment, we've already addressed this concern.

- “co-location” vs “collocation”? Stay consistent.

Response: Thank you for your observation. We've addressed this by ensuring the usage of "co-location" throughout the text.

- Many links in the “References” section of the supplementary point to a Zotero page that is meant for Google docs, thus rendering the links inaccessible

Response: Thank you for bringing this to our attention. We have revised the links in the “References” section to ensure accessibility and functionality.

Figure comments

In general, labeling the figure panels with letters (e.g. (a), (b), (c), (d)) allows for easier and clearer reference in text and in figure captions. (e.g. Line 327 mentions the “top row” in Figure 8)

Response: While we appreciate the recommendation, we believe that maintaining the current status is appropriate to avoid potential visual clutter. We will leave the decision to the editor's discretion.

- Figure 1. Good visual—a nice representation of the timeline of events.

Response: Thank you for the compliment; we appreciate your positive feedback. In the reviewed draft, we have also slightly enhanced Figure 1 by explicitly adding important dates to the study's timeline, as well as the names of the companies and the number of systems involved.

- Figure 7. Which sensors are being compared here? Why the anonymity compared to the other section(s)? Also, the readers may benefit from a colorblind-friendly and more contrasting color palette. “Class 1” and “Class 2” sensors are not actually described until page 15 (line 377 onwards) – it might be useful to refer to this section (i.e. Section 3.6) in the figure caption or the accompanying text (paragraphs) that describes this figure, and mention that it will be thoroughly explained in that section.

Response: We sincerely appreciate the suggestion for clarification. Our decision to anonymize sensors in specific figures intentionally focuses our analysis on evaluating broader sensor technology rather than individual brands. This approach aims to prevent biased interpretations, encouraging a general understanding of technological capabilities and limitations. We elaborate on our reasoning for anonymization in the “3. Results and discussion” section (see the earlier response to this point).

Following Copernicus guidelines, we ensured Figure 7 is accessible to all readers, including those with CVD, by adopting a colorblind-friendly palette (using Python's seaborn library “colorblind” option). We further validated the figure's colours via the Color Blindness Simulator (<https://www.color-blindness.com/coblis-color-blindness-simulator/>), as per Copernicus guidelines.

Regarding the figure's initial oversight in referring to “Class 1” and “Class 2” sensors for PM_{2.5}, we acknowledge that the CEN/TS 17660-1:2021 standard applies only to gases. This error has been corrected to include the Data Quality Objectives (DQOs) of the EU AQ Directive, serving here only as a reference. Consequently, the figure caption has been updated to read:

Figure 7. Regression (top) and REU (bottom) plots showing data from four PM_{2.5} sensors (same manufacturer) over 2 time periods: Apr-Jun 2022 and Aug-Oct 2022. The four devices were in separate locations in the first period, but all deployed in Manchester in the second. Only for reference, we have included the PM_{2.5} DQOs as outlined by the EU AQ Directive (for “fixed” PM_{2.5} measurements, REU < 25%; for “indicative” PM_{2.5} measurements, REU < 50%) as horizontal dashed lines.

- Figures 8 and 10. Explain the colorations, e.g. is it meant to be a heat map? What do the specific colors mean? Figure 10 may also benefit from higher contrasting – difficult to see the contrast especially in the lower left panel, and when the plots are printed. Dashing is also difficult to see—might benefit from greater color contrast.

Response: Thank you for your constructive feedback. The colour gradients in Figure 10 are indeed representative of a heat map, where darker colours indicate higher densities of sensor readings within the specified REU and DC values. We have modified this figure adjusting the

contrast. Additionally, the dashed lines have been thickened and their transparency reduced. As for Figure 8, the colour gradient indicates data point density, with darker colours representing lower densities and brighter colours highlighting higher densities.

Line by Line

- **Line 80: Suggestion: “academia” or “academic research” instead of “academic arena”.**

Response: Thank you for the suggestion. We have decided to keep the term “academic arena” to broadly encompass the variety of scholarly activities related to this topic.

- **Line 104: Suggestion: reword “transparent”. Suggested synonyms: open, comprehensive (this changes the meaning a bit)**

Response: Thank you for the suggestion. We have decided to keep “transparent” as it precisely conveys our intended meaning, in that all methodologies and assessment criteria are open. Much of the performance data used by manufacturers to advertise sensor devices is not transparent and thus is difficult to extrapolate to end-user applications

- **Line 118: Typo: “influenced”**

Response: Thank you for catching that typo. It has been corrected to “influenced”.

- **Line 155: Suggestion: reword “ratified” to “validated”**

Response: Thank you for your suggestion. We have opted to maintain “ratified” as it is the specific terminology used by the National Physical Laboratory (NPL) in this context.

- **Line 160: is “time-line” the correct term? Perhaps “comparison” or “matrix” would be more apt for Figure S1; Figure S2 is a scatter plot or a bivariate plot.**

Response: Thank you for the suggestion. We confirm that “time-line” is the correct term, as both figures Figure S1 & S2 (now renamed as Figure S2 & S3) illustrate chronological sequences.

- **Line 162: change “to use this data” to “to use these data”**

Response: Thank you for pointing this out. We have corrected the phrase to “to use these data” to adhere to the grammatical convention.

- **Line 206: “Mean Bias Error” rather than “Mean Error Bias”**

Response: Thank you for your correction, the term has been updated.

- **Line 209: enclose “out-of-box” in quotation marks; typically “out-of-the-box”**

Response: Thank you for your suggestion. The term “out-of-box” is used in this context as an abbreviated form of “out-of-the-box”, facilitating its encoding within our documentation, data processing (see “Data collection”) and our metadata.

- **Line 231: semi-colon after “MBE”, comma after “machine learning”**

Response: Thank you for the suggestions. We have implemented the suggested changes.

- **Line 257: actual “metrological” as in measurements and units, or “meteorological” as in RH and Temp?**

Response: Thank you for your inquiry. The term “metrological” is indeed correct in this context, reflecting the focus of our discussion on sensor data uncertainty.

- **Line 271: “...hypothetical scenario where it...” does “it” refer to T200U? T500?**

Response: Thank you for your comment. The reference to “it” pertains to the T200U, as contextually established in the preceding sentences.

- **Line 275: “All of this” to “all of these”**

Response: Thank you for the suggestion. Upon review, we find the original phrasing “All of this” accurately encompasses the list of required actions as a collective process, and therefore we would prefer to retain the original wording.

- **Line 276: Add comma after “monitoring”**

Response: Thank you for your attention to detail. The change has been made.

- **Line 278: “equivalent-to-reference” – consistency in hyphenation**

Response: Thank you for your comment. The use of “equivalent to reference” (without hyphenation) within quotation marks is deliberate to signify direct terminology as specified in the EU Air Quality Directive. This precise phrasing is retained to reflect the source accurately.

- **Line 282: “obtained with a BAM at the AURN York site, located on a busy avenue” – delete parentheses**

Response: Thank you for the comment. We believe the current use of parentheses enhances the reader's understanding. Therefore, we have chosen to retain it as is.

- **Line 289: Omit “of course”**

Response: Thank you for your suggestion. We have removed “of course” from the text.

- **Line 299: capitalize “FIDAS”**

Response: Thank you for your suggestion. The term “Fidas” is presented in a manner consistent with certain source materials (including the instrument manufacturer website) and common usage within our document. Thus, we have decided to keep it in the text.

- **Lines 300-301: the choice of the reference measurement**

Response: Thanks for the suggestion. “reference method” aligns with our use of PM instruments using different measurement methods/techniques. We'd therefore like to keep the original wording for consistency.

- **Line 309: Paraphrase “saw its slope change”. Suggested: ...”a slope change from 0.69 to 0.86 was observed...”**

Response: Thank you for your suggestion. We have made the adjustment as suggested.

- **Line 310: change “when” to “while”**

Response: Thank you for your suggestion. We agree with your recommendation and have updated the manuscript accordingly.

- **Line 321: “despite” might not be the correct conjunction here.**

Response: Thank you for your comment. After review, we believe it accurately conveys the intended contrast, so we have decided to retain it.

- **Line 331: change “akin to this later” to “akin to the later”**

Response: Thank you for your suggestion. The correction has been applied as suggested.

- **Line 268: Redundant. Change “with a measurement instrument” to “with an instrument”**

Response: Thank you for your recommendation. We have updated the text.

- **Lines 368-369: Can be paraphrased to be more succinct.**

Response: Thank you for your suggestion. We have opted to keep the original phrasing, as it precisely communicates the critical concept of uncertainty in measurement instruments and their implications.

- **Line 383: “4-system” rather than “4 systems companies”**

Response: Thank you for your suggestion. We have decided to retain the original wording as it accurately reflects our analysis of data from the sensor systems provided by four distinct companies. Each plot in Figure 10 represents the aggregated data from all operative sensors of each of the shown companies, making “4 systems companies” the most precise description of our evaluation.

- **Line 386: add “dashed”, i.e., green dashed rectangle**

Response: Thank you for your suggestion. It was corrected.

- **Line 398-399. “high time-resolution” (note hyphen placement)**

Response: Thank you for your suggestion. It was corrected.

- **Line 400: subscript on NO₂**

Response: Thank you for your suggestion. It was corrected.

- **Line 400: Is “DEFRA” all capitalized, or is it “Defra” as mentioned in the acknowledgement (Line 489)?**

Response: Thank you for your pointing this out. We have chosen to retain “Defra”.

- **Line 407: Consider using a different word from “digestible”**

Response: Thank you for your feedback on this term. We replaced “digestible” by “accessible”.

- **Line 433: change “uptake” to “uptick”**

Response: Thank you for your suggestion. We believe that “uptake” is more appropriate in this context as it is a well-established term commonly used to describe the widespread adoption or acceptance of new technologies or practices. Therefore, we have decided to retain it.

- Line 452: “high level” seems unnecessary.

Response: Thank you for your observation. We believe this term is necessary to accurately convey the depth of the dataset analysis conducted.

- Line 455: “accuracy with respect to reference methods”

Response: Thank you for your suggestion. We believe that the current wording effectively conveys the ideas.

- Line 471: Lacks the link to supplementary information (online version link is accessible).

Response: Thank you for your observation. Including the link to supplementary information is indeed part of the editorial process.

Reviewer#2

General Comments

Overall this paper provides a good overview of the QUANT study and some salient results. A few clarifications are needed, as outlined below.

Response: Thank you for your positive feedback on the paper and for acknowledging the overview it provides of the QUANT study along with its key findings.

Section 2.3 should describe any harmonization of the data from the sensors' reporting frequencies to a standard frequency, i.e., what was the common time frequency for which the measurements were averaged for analysis and comparison with the reference? Or was this done differently for the native reporting frequencies of each instrument? Finally, in the available QUANT dataset, are the measurements reported at the initial sampling frequency or at the down-averaged frequency (or both)?

Response: Thank you for your inquiries regarding the handling of sensor data frequencies.

-In regards to the data harmonization, we have updated the methodology text ("2.3 Sensor deployment and data collection"), and now reads: Minor pre-processing was applied at this stage, [including temporal harmonisation to ensure that all measurements had a minimum sampling period of 1-minute, ensuring consistency in measurement units and labels, and coercing into the same format to allow for full compatibility across sensor units.](#)

-as for the data collection frequencies, we have added the following text in order to clarify this (see "2.3 Sensor deployment and data collection"):

[For an overview of the sensor measurands and their corresponding data time resolutions as provided by the companies participating in the Main QUANT study and the WPS, please see Seccion S3 and S4 \(Table S4 and S5\) respectively.](#)

-Regarding the analysis showcased in this overview, we processed the sensor data into hourly averages. We have added the following text to clarify this (see the "Results" section):

[All metrics and plots presented here are based on 1-hour averaged data.](#)

-The QUANT dataset reports data at 1-min time resolution. We have recently submitted a detailed manuscript that delves into the QUANT database (still under review). For more details, please refer to the response to this reviewer's last General Comment response.

In Section 3.1, results are only presented for the PM2.5 data. I would suggest that information on the inter-sensor precision for all measurands should be provided, maybe as part of the supplemental materials, since this is a basic feature of the different sensors which can inform all the other results presented later.

Response: While the primary aim of this manuscript is to serve as an overview —introducing the methodology used in the QUANT study, showcasing the data's potential, and highlighting broader findings— to align with this feedback we have added NO2 and O3 inter-sensor precision plots to the supplemental materials. It's important to clarify that subsequent publications will delve into detailed analyses, where more specific findings will be explored extensively.

Since one of the goals of this paper is to introduce the QUANT dataset as a public resource for long-term performance assessment, it may be worth adding a section which details the dataset itself, or expanding the “Data Availability” section to do this. Some points to consider for this section would be the size of the dataset, the parameters included, what quality controls are applied (especially to the reference data), and any licensing of the dataset or policies associated with its use. Currently, the link provided in the “Data Availability” section does not seem to be working; presumably this will be active by the time of publication.

Response: Thank you for your valuable suggestions. Our manuscript is primarily an overview intended to introduce the QUANT study's methodology, showcase the collected data's potential, and present general findings, rather than a detailed dataset description. The dataset's complexity, including multiple calibration products for each measured species for certain devices, made a comprehensive description challenging within this paper's scope. However, to thoroughly address the dataset specifics, we recently submitted a detailed data descriptor manuscript, providing extensive details on the collection, processing, accessibility, and structure of the QUANT dataset, including variables, reporting frequencies, and QA/QC measures. This manuscript is currently under review, and we believe it will greatly aid in understanding and using the QUANT dataset upon publication. Complementarily, we have updated the “Data Availability” text, and now reads:

The QUANT dataset, accessible at the Centre for Environmental Data Analysis (CEDA) (Lacy et al., 2023; <https://catalogue.ceda.ac.uk/uuid/ae1df3ef736f4248927984b7aa079d2e>), is the most extensive collection to date assessing air pollution sensors' performance in UK urban settings. It encompasses gas and PM sensor data recorded in the native reporting frequency of each device. The reference data from the three monitoring sites can be found at:

- MAQS: <https://data.ceda.ac.uk/badc/osca/data/manchester>;
- LAQS: <https://www.londonair.org.uk/london/asp/datadownload.asp>);
- YoFi: https://uk-air.defra.gov.uk/data/data_selector.

A comprehensive data descriptor manuscript, detailing the QUANT dataset's collection methods, processing protocols, accessibility features, and overall structure—including variables, data reporting frequencies, and QA/QC practices—has been submitted for publication. At the time of this writing, the manuscript is still under review.

A GitHub repository at <https://github.com/wacl-york/quant-air-pollution-measurement-errors> provides access to Python and R scripts designed for generating diagnostic visuals and metrics related to the QUANT study, along with sample analyses using the QUANT dataset.

Specific Comments

Line 19: suggest clarification that this technology is providing the first steps for regions without pre-existing monitoring.

Response: We appreciate your suggestion. We have revised it as follows:

In times of growing concern about the impacts of air pollution across the globe, lower-cost sensor technology is giving the first steps in helping to enhance our understanding and ability to manage air quality issues, [particularly in regions without established monitoring networks](#).

Line 34: “end-users” should be “end-user”.

Response: Thank you for your suggestion. Correction made.

Line 35: “capabilities the” should be “capabilities, the”.

Response: Thank you for your suggestion. The wording was corrected.

Line 54: “helping mitigating” should be “helping to mitigate”.

Response: Thank you for your correction. It was rephrased accordingly.

Line 61: suggest removing “of”.

Response: Thank you for your suggestion. This was removed.

Line 90: “extensive” is repeated.

Response: Thank you for your suggestion. The sentence was corrected and now it reads: “...alongside extensive reference measurements, to generate the data for a [comprehensive extensive](#) in-depth performance assessment.”

Line 118: “inlfuenced” should be “influenced”.

Response: Thank you for your suggestion. Correction made.

Line 128: “Quant” should be “QUANT”.

Response: Thank you for your suggestion. It was corrected.

Line 142: “Polludrone: Poll” should be “Poll: Polludrone”.

Response: Thank you for your suggestion. The correction was made.

Figure 2: Suggest using the same colors for the different sensors between the left and right panels.

Response: Thank you for noticing this. This was corrected.

Figure 3: Suggest moving this figure and associated discussion to the next section, since it is an assessment of performance against a reference rather than an assessment of inter-sensor consistency.

Response: Thank you for your feedback regarding Figure 3. Although it assesses performance against a reference, it also reveals inter-device precision through the dispersion of points for sensors of the same brand. Its strategic location before the section “3.2 Device accuracy and collocation calibrations” provides a transition to discussions on accuracy, underscoring not only reference comparison but also variability among devices of the same make—essential for understanding sensor consistency and reliability. Thus, we would like to maintain Figure 3 in its

current position, but will move it if the reviewer and editor insist. We have slightly adjusted the preceding text for clarity, as follows:

~~In addition to highlighting which devices are most accurate, Fig. 3 also provides an additional perspective of inter-device precision. In addition to showcasing inter-device precision, Fig. 3 also serves as a transition to accuracy evaluation (the focus of the subsequent section).~~

Figure 8: These seems to be a switch between the use of uncalibrated and calibrated data between the left and right panels as well. It is not clear what these calibrations are based on, and the application of the calibration might be a contributing factor to the difference in performance, together with the move between sites. It may be more illustrative to present a comparison at both sites with either the calibrated or the uncalibrated data only.

Response: We'd like to clarify that the "out-of-the-box" and "calibrated" data products are associated with specific periods (as summarised in Figures S2 and S3). Calibrations were performed by the companies using data from Manchester, during the first co-location period (Dec 2019 - Feb 2020). At the end of this period, the brands ceased providing "out-of-the-box" data and began supplying data adjusted for the co-location data from Manchester. A few days later, one quarter of the instruments were moved to London. Thus, while the data are labelled as "calibrated", it does not imply (in this case) that they have been corrected to local conditions in London.

Regarding the transition between uncalibrated and calibrated data across the panels, it's important to note that we lack access to the specific calibration methods used by the manufacturers. This limits our ability to comprehensively detail the foundation of these calibrations.

As for the observed differences in performance between sites, and the potential influence of the calibration approach employed, we acknowledge that both aspects can be significant. Although London and Manchester are classified as "urban background" sites, one might expect comparable sensor performance, disparate calibration methods—applied as manufacturers assimilate local reference data—may lead to divergent outcomes. This is exemplified by "Sensor A", which, upon relocation to London, exhibits a shift in bias while maintaining response linearity. In contrast, "Sensor B" shows a notable degradation in linearity. We suspect that the distinct calibration methodologies each company employs markedly influence these performance variances. Yet, beyond this speculation, the point we aim to highlight with this figure is the potential for end-users to implement simple corrections. Specifically, "Sensor A" appears amenable to linear correction, whereas for "Sensor B", such an approach may not yield significant benefits.

Concerning the suggestion to present a comparison at both sites using exclusively calibrated or uncalibrated data, we are limited by the nature of the data products available during the periods in question (as detailed in Figures S2 and S3).

The original text has been reworded in order to convey these points and now reads:

~~The primary distinction between both systems' behaviour lies in the fact that the sensor located in the top row, even after being relocated to London, maintains a linear response (albeit slightly more degraded than that observed in Manchester, as the R^2 and RMSE show). In contrast, in the second system (bottom row),~~

~~the response is notably noisier as the Standard Error (SE)—which is the dispersion of the data around the best line fit line, i.e., the remaining error after bias correction. In scenarios akin to this latter, where there is a high variance in the residuals, a linear correction will not provide a significant improvement. While more sophisticated corrections could be applied, these will be limited by domain knowledge of the end-user, and potentially by other complex data sources that might be available. However, it is important to remember that additional post-processing could increase the risk of overfitting (Aula et al., 2022). On the other hand, for cases like the top plots, users might benefit from trying to correct them using simple linear correction (e.g. using reference instruments if available) or other approaches that could provide means for zero and span correction. A straightforward and cost effective example could be the use of diffusion tubes for the case of NO₂, as discussed in Section 3.6. The primary distinction between both systems' behaviour lies in the fact that the sensor located in the top row (Sensor A), even after being relocated to London, maintains a linear response (albeit slightly more degraded than that observed in Manchester, as indicated by the R² and RMSE). In contrast, Sensor B's response becomes significantly noisier upon relocation to London, as highlighted by the Standard Error (SE) —which represents the remaining error after applying a perfect bias correction. Despite both systems utilising identical sensing elements, the variance in residuals between them may stem from the distinct calibration approaches applied by the respective companies.~~

For cases resembling Sensor A, users might find it beneficial to implement simple linear correction methods (e.g., using reference instruments if available) or explore other strategies for zero and span correction. A practical and cost-effective approach, for example, is using diffusion tubes for NO₂ measurements, as discussed in Section 3.6. Conversely, in scenarios characterised by high variance in residuals, such as those observed with Sensor B, a-posteriori attempts to apply a simple linear correction are unlikely to result in significant improvement. While more sophisticated corrections are theoretically feasible, their effectiveness is limited by the end-user's domain knowledge and the availability of additional complex data sources. Furthermore, it is important to consider that excessive post-processing may lead to overfitting — a situation where a model excessively conforms to specific patterns in the training data, resulting in poor performance on new, unseen data (Aula et al., 2022).

Lines 329-331: Sentence may be incomplete.

Response: Thank you for noticing this. We have adapted the text. Please see the previous response text.

Lines 373-374: The meaning of this is unclear; does this mean that results from 2 systems were combined (e.g., to increase coverage)? Or were coverage and REU assessed separately for each device and then data from both devices combined to create the density plots of Figure 10?

Response: Thank you for your inquiry. Each sensor device was independently assessed in terms of Data Coverage (DC) and Relative Expanded Uncertainty (REU). After this, we aggregated the data to create the density plots for all units of a unique brand, thus illustrating the collective behaviour of NO₂ sensors from the same company in relation to DC and REU. Recognizing that the original text may not have clearly conveyed this, we have revised it as follows:

~~Figure 10 shows the REU (y-axis) and Data Coverage (DC, x-axis) of companies measuring NO₂ with more than 2 systems running to avoid ambiguity in the results. Using multiple systems, not only avoids ambiguity in results but also enhances the robustness of the data collected.~~ Figure 10 illustrates the collective behaviour of NO₂ sensors from each of the four companies with more than two working systems, showcasing their REU (y-axis) versus Data Coverage (DC, x-axis). Both parameters were calculated for each sensor system using a 40-day moving window approach and then aggregated by brand, ensuring a comprehensive analysis. This methodology leverages overlapping data from multiple sensors to provide a robust representation of company-wide sensor performance and aims to prevent biased interpretations.

Line 374: “systems, not” should be “systems not”.

Response: Thank you for your suggestion. The correction was made.

Line 423-424: Please explain further the use of the reference data as a prior in this method.

Response: Thank you for your comment. We have removed the mention of the prior as it distracted from the overall results, and instead have provided references to several papers that explain the general deweathering strategy.

Lines 435-436: Consider changing one of “developments” or “developing”.

Response: Thank you for your comment. We have revised the sentence to eliminate the redundancy and improve the readability of the text. The modified sentence now reads:

~~Developments in the field of air pollution sensor technology are also developing rapidly, with advances in both the measurement technology and particularly in the data post-processing and calibration.~~ Advances are occurring rapidly, in both the measurement technology and particularly in the data post-processing and calibration.

Lines 460-469: I would suggest adding a sentence earlier in the document (and perhaps in the abstract) noting that further analysis will be left for future publications. I was expecting at several points a more comprehensive presentation of results across all pollutants and for all phases of the study, while only particular aspects of the results were highlighted. This is alright, but I think it needs to be more clearly stated up-front that this is not a comprehensive presentation of the study results. I would also suggest, as a topic for future work, examining the manufacturer-supplied calibrations in more detail, seeing where these improved upon the raw and where they perhaps did not, and how robust the calibrations are to environment changes and movement of the sensors to new sites. This is briefly presented in several figures, e.g., Figure 8, but a more comprehensive assessment across all sensors and pollutants could be made.

Response: We appreciate your feedback and have taken steps to clarify the scope and intent of our study both in the abstract and in the introduction of our document. In the abstract, we have added the following sentence:

While more comprehensive analyses are reserved for future detailed publications, the results shown here highlight the significant variation between systems, the incidence of corrections made by manufacturers, the effects of relocation to different environments, and the long-term behaviour of the systems.

Similarly, we have modified the introduction, and now reads:

This comprehensive approach offers unprecedented insights into the operational capabilities and limitations of these sensors in real-world conditions. Significantly, some of the insights gathered during QUANT have contributed to the development of the Publicly Available Specification (PAS 4023, 2023), which provides guidelines for the selection, deployment, maintenance, and quality assurance of air quality sensor systems. [While this manuscript serves as an initial overview, detailed analyses of the measured pollutants and study phases, offering a more comprehensive perspective on sensor performance, are planned for future publications.](#)

We appreciate your suggestion concerning the detailed examination of manufacturer-supplied calibrations. Indeed, this aspect is being considered in our current efforts and we anticipate publishing these findings in the near future.

Supplemental Information, Lines 4-6: Indicate which of these channels and/or data products were considered for this study. Also report the sampling frequency for this sensor.

Response: Thank you for your feedback on this point. We wish to clarify that the original supplementary text did indeed specify that the PA sensors provided data at a 2-minute resolution. Furthermore, in response to concerns about channels and data products, we have added a note to the manuscript for clarity, which states:

[*Note: For this study, only Channel A and the data product “cf_atm” were included in the analysis and shown in the plots.](#)

To ensure a comprehensive understanding of the sensor data utilised, we collected all data products offered by each company, preserving their native resolution. To enhance clarity, we have now included two additional tables in the supplementary materials—one for the QUANT study (table S4) and another for the WPS study (table S5). These tables summarise the data products collected during QUANT and their native resolution.

Supplemental Information, Line 55: “y” should be “and”.

Response: Thank you for your suggestion. The correction was applied.

Reviewer#3

GENERAL COMMENTS

This article provides an important contribution to the advancement of studies associated with air quality sensors. It provides a good overview and information about the QUANT study and some important results. Discussions associated with data quality add value in an important way to alert to errors and possible corrections associated with time and space. Shows the importance of using reference sensors in calibrations detailing correct use and necessary considerations.

I recommend this publication. However, as this is an important study that can be reused or used as a basis for others, I think it is important to go into more detail especially methodologically so that it can be continued and used as the authors suggest at the end.

Response: Thank you for recognizing the contribution of our article and for your supportive remarks on the overview and insights provided by the QUANT study.

SPECIFIC COMMENTS

Section 2.1

As spatial analyzes are carried out, I consider the spatial description of the areas of the article to be important such as distances and spatial layout. A spatial image would enhance spatial visualization and discussion. This arrangement is important in analyzing spatial differences and environmental conditions that influence the data.

Response: Thank you for this suggestion. Recognizing this, we have expanded the description of the study areas to include more detailed information on distances and spatial layouts. Additionally, to further enhance spatial visualisation and support the discussion, we have incorporated satellite images taken from Google Maps into our manuscript.

Section 2.3

Line 135. The sensors were implemented according to the manufacturer's specifications. Was any standardization found in the logistics or studied at this stage? I think it's important to describe this stage perhaps in supplementary material. The layout of the sensors, whether it was completely open or needed some protection, ground height, proximity to the reference, obstacles, necessary infrastructure, etc. These are all factors that influence the data and are still the subject of much discussion when it comes to implementing the sensors.

Response: Thank you for highlighting this point. In response to your suggestions, we have adapted and renamed one of the sub-sections in the methodology, in order to describe these important points. Please refer to:

2.3 Sensor deployment and data collection, ~~co-located reference data and data products~~

Section about treatments, analysis, and metrics

It would be important to include a section describing the data analysis treatments and statistical metrics that were used for these specific analyses.

Response: Thank you for your detailed suggestion. We've created a new section in the supplementary ("S5. Performance Metrics") that succinctly describe the statistical metrics employed in our analysis.

It would be important to include: data standardizations such as sensor frequencies for comparison with the reference, if there was a change in frequency, how the amount of valid data for the calculations was considered; a description of the calibrations or validations applied; statistical metrics used in analyzes such as RMSE, REU, etc., a simple description would add a lot to the article; Another point would be the pollutants used (PM2.5, NO2) and because these if there are analyzes for the others, it would be interesting to mention.

Response: Thank you for your detailed suggestion. Following previous reviewers' comments, we have taken steps to address these aspects in the manuscript:

-for the updated text on data standardisation and sensor frequencies, please refer to the "2.3 Data collection" section.

-in regards to the amount of valid data used for metrics and plots, please see the added text in the "Results" section.

-we've expanded our description on the calibration processes applied to the sensors in the newly created section "2.4 Data products and co-located reference data".

-as for statistical metrics, we've created a new section in the supplementary ("S5. Performance Metrics") that describes the statistical metrics employed in our analysis.

SECTION 3

Why were some analyzes used PM2.5 and others NO2? Would there be any explanation?

Response: Thank you for your inquiry. The primary air pollution issues in the UK are PM2.5 and NO2 exceedances, and as such in this overview paper we aimed to showcase the NO2 and PM2.5 measurements over those of other pollutants due to their relevance. The choice of our use of NO2 or PM2.5 for any particular example shown is either to highlight a specific facet of the data, such as the potential use of NO2 diffusion tubes to reduce NO2 sensor bias, or is arbitrary in order to avoid focussing more on one pollutant over another. We have added the following text to the results to clarify this:

[The majority of examples presented here focus on PM_{2.5} and NO₂ measurements, due to both a larger dataset available for these pollutants and their critical role in addressing the exceedances that predominantly impact UK air quality.](#)

From section 3.4 onwards, sensors are no longer specified from which manufacturer. For example, in Figure 7, which sensors are being compared? Is it just from one manufacturer? Or multiple manufacturers? Is only one sensor from each manufacturer considered or multiple?

Response: Thank you for your question. To emphasise the broader implications and insights into sensor technology, we chose to anonymize figures illustrating brand-specific features. This aims to mitigate potential bias and foster a broader view of the technology performance, focusing on general trends rather than the performance of individual manufacturers. We've provided our reasoning in the "3. Results and discussion" section for clarity.

As for Figure 7, the original caption specifies that the comparison involves sensors from a single manufacturer, though we have anonymized the details to align with our overarching goal of emphasising generalizable findings.

In figure 8, what would sensors A and B be, are they from the same manufacturer or different?

Response: Thank you for your inquiry. Sensors A and B represent two distinct systems from different manufacturers. We have adapted the manuscript accordingly to clarify this point for our readers:

A second example of **inter-location** performance ~~changing between locations~~ is presented in Fig. 8, showing NO₂ data from two sensor systems **(from two different manufacturers, identified as Systems A and B)** ~~(different brands, one shown on top of the other)~~ before (left plots) and after (right plots) they were moved from Manchester to London in March 2020.

Figure 10. Is the analysis for NO₂? If yes, specify in the legend. The companies are unidentified, wouldn't it be possible to associate them with each one?

Response: Thank you for your input. The specified corrections have been made. We've also identified the companies in Figure 10. As we responded to an earlier comment, we initially chose the anonymize companies to focus discussions on broad technological features over specific manufacturer data.

1 Long-term Evaluation of Commercial Air Quality Sensors: An 2 Overview from the QUANT Study

3 Sebastian Diez^{1,2}, Stuart Lacy², Hugh Coe³, Josefina Urquiza^{4,5}, Max Priestman⁶, Michael
4 Flynn³, Nicholas Marsden³, Nicholas A. Martin⁷, Stefan Gillott⁶, Thomas Bannan³, Pete
5 Edwards²

6 ¹Centro de Investigación en Tecnologías para la Sociedad, Universidad del Desarrollo, Santiago, Chile, CP 7550000

7 ²Wolfson Atmospheric Chemistry Laboratories, University of York, York, YO10 5DD, UK

8 ³Department of Earth and Environmental Science, Centre for Atmospheric Science, School of Natural Sciences, The
9 University of Manchester, Manchester, M13 9PL, UK

10 ⁴Grupo de Estudios de la Atmósfera y el Ambiente (GEAA), Universidad Tecnológica Nacional, Facultad Regional
11 Mendoza (UTN-FRM), Cnel. Rodríguez 273, Mendoza, 5501, Argentina

12 ⁵Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET) Argentina

13 ⁶MRC Centre for Environment and Health, Environmental Research Group, Imperial College, London, W12 0BZ,
14 UK

15 ⁷National Physical Laboratory, Teddington TW11 0LW, UK

16 *Correspondence:* Sebastian Diez (sebastian.diez@udd.cl); Pete Edwards (pete.edwards@york.ac.uk)

17 **Abstract.** In times of growing concern about the impacts of air pollution across the globe, lower-cost sensor
18 technology is giving the first steps in helping to enhance our understanding and ability to manage air quality issues,
19 particularly in regions without established monitoring networks. While the benefits of greater spatial coverage and
20 real-time measurements that these systems offer are evident, challenges still need to be addressed regarding sensor
21 reliability and data quality. Given the limitations imposed by intellectual property, commercial implementations are
22 often “black boxes”, which represents an extra challenge as it limits end-users' understanding of the data production
23 process. In this paper we present an overview of the QUANT (Quantification of Utility of Atmospheric Network
24 Technologies) study, a comprehensive 3-year assessment across a range of urban environments in the United
25 Kingdom, evaluating 43 sensor devices, including 119 gas sensors and 118 particulate matter sensors, from multiple
26 companies. QUANT stands out as one of the most comprehensive studies of commercial air quality sensor systems
27 carried out to date, encompassing a wide variety of companies in a single evaluation and including two generations
28 of sensor technologies. Integrated into an extensive data set open to the public, it was designed to provide a long-term
29 evaluation of the precision, accuracy, and stability of commercially available sensor systems. To attain a nuanced
30 understanding of sensor performance, we have complemented commonly used single-value metrics (e.g., Coefficient
31 of Determination (R^2), Root Mean Square Error (RMSE), Mean Absolute Error (MAE)) with visual tools. These
32 include Regression plots, Relative Expanded Uncertainty (REU) plots, and Target plots, enhancing our analysis
33 beyond traditional metrics. This overview discusses the assessment methodology, and key findings showcasing the
34 significance of the study. While more comprehensive analyses are reserved for future detailed publications, the results

35 shown here highlight the significant variation between systems, the incidence of corrections made by manufacturers,
36 the effects of relocation to different environments, and the long-term behaviour of the systems. Additionally, the
37 importance of accounting for uncertainties associated with reference instruments in sensor evaluations is emphasised.
38 Practical considerations in the application of these sensors in real-world scenarios are also discussed, and potential
39 solutions to end-user data challenges are presented. Offering key information about the sensor systems' capabilities,
40 the QUANT study will serve as a valuable resource for those seeking to implement commercial solutions as
41 complementary tools to tackle air pollution.

42 **Keywords:** air pollution, commercial sensor systems, QUANT, long-term evaluation.

43 1. Introduction

44 Emerging lower-cost sensor systems¹ offer a promising alternative to the more expensive and complex monitoring
45 equipment traditionally used for measuring air pollutants such as PM_{2.5}, NO₂, and O₃ (Okure et al., 2022). These
46 innovative devices hold the potential to expand spatial coverage (Malings et al., 2020) and deliver real-time air
47 pollution measurements (Tanzer-Gruener et al., 2020). However, concerns regarding the variable quality of the data
48 they provide still hinder their acceptance as reliable measurement technologies (Karagulian et al., 2019; Zamora et
49 al., 2020).

50 Sensors² face key challenges such as cross-sensitivities (Bittner et al., 2022; Cross et al., 2017; Levy Zamora et al.,
51 2022; Pang et al., 2018), internal consistency (Feenstra et al., 2019; Ripoll et al., 2019), signal drift (A. Miech et al.,
52 2023; Li et al., 2021; Sayahi et al., 2019), long term performance (Bulot et al., 2019; Liu et al., 2020) and data coverage
53 (Brown & Martin, 2023; Duvall et al., 2021; Feinberg et al., 2018). Additionally, environmental factors such as
54 temperature and humidity (Bittner et al., 2022; Farquhar et al., 2021; ~~and humidity~~ Crilley et al., 2018; Williams,
55 2020) can significantly influence sensor signals.

56 In recent years, manufacturers of both sensing elements (Han et al., 2021; Nazemi et al., 2019) and sensor systems
57 have made significant technological advances (Chojer et al., 2020). For example, there are now commercial and non-
58 commercial systems equipped with multiple detectors to measure distinct pollutants (Buehler et al., 2021; Hagan et
59 al., 2019; Pang et al., 2021) helping to mitigate the effects of cross-interferences. Additionally, enhancements in
60 electrochemical OEMs have been demonstrated in terms of their specificity (Baron & Saffell, 2017; Ouyang, 2020).

61 However, the complex nature of their responses, coupled with their dependence on local conditions means sensor
62 performance can be inconsistent (Bi et al., 2020). This complicates the comparison of results or anticipating sensor
63 future performance across different studies. Moreover, assessments of sensor performance found in the academic

¹ The term “sensor systems” refers to sensors housed within a protective case, which includes a sampling and power system, electronic hardware and software for data acquisition, analog-to-digital conversion, data processing and their transfer (Karagulian et al., 2019). Unless specified otherwise, the term “sensor” will be used as a synonym of “sensor systems”. Other alternative names for “sensor systems” used here are “sensor devices” (or “devices”), “sensor units” (or “units”).

² In a narrower sense, “sensor” typically denotes the specific component within a sensor system that detects and responds to environmental inputs, producing a corresponding output signal. To distinguish this from the broader use of “sensor” as equivalent to “sensor system” in our text, we will utilise alternative terms such as “detector”, “sensing element”, or “OEM” (original equipment manufacturer) when referring specifically to this component, thereby preventing confusion.

64 literature often rely on a range of protocols (e.g., CEN (2021) and Duvall et al. (2021)) and data quality metrics (e.g.,
65 Spinelle et al. (2017) and Zimmerman et al. (2018)), with many studies limited to a single-site co-location and/or
66 short-term evaluations that do not fully account for broader environmental variations (Karagulian et al., 2019).

67 The calibration of any instrument used to measure atmospheric composition is fundamental to guarantee their accuracy
68 (Alam et al., 2020; Long et al., 2021; Wu et al., 2022). Using out-of-the-box sensor data without fit-for-purpose
69 calibration can produce misleading results (Liang & Daniels, 2022). An effective calibration **not only** involves
70 identifying but also **compensating for estimated and correcting systematic effects errors** in the sensor readings, a
71 **process defined as a correction (for a detailed definition and differentiation of calibration and correction see JCGM,**
72 **2012).** For standard air pollution measurement techniques, calibration is often performed in a controlled laboratory
73 environment (Liang, 2021), ~~or by sampling gas from a certified standard cylinder in the field. For PM, particles of~~
74 ~~known density and size are used, controlling the airflow conditions.~~ **For example, for gases, a known concentration is**
75 **sampled from a certified standard. Similarly, for PM, particles of known density and size are generated. Both gases**
76 **and PM calibration are conducted under controlled airflow conditions**

77 Yet, the aforementioned challenges with lower-cost sensor-based devices suggest that such calibrations may not
78 always accurately reflect real-world conditions (Giordano et al., 2021). A frequent approach involves co-locating
79 sensors alongside regulatory instruments in their intended deployment areas and/or conditions and using data-driven
80 methods to match the reference data (Liang & Daniels, 2022). Numerous studies have investigated the effectiveness
81 of calibration methods for sensors e.g. (Bigi et al., 2018; Bittner et al., 2022; Malings et al., 2020; Spinelle et al., 2017;
82 Zimmerman et al., 2018), including selecting appropriate reference instruments (Kelly et al., 2017), the need for
83 regular calibration to maintain accuracy (Gamboa et al., 2023), the necessity of rigorous calibration protocols to ensure
84 consistency (Kang et al., 2022), and transferability (Nowack et al., 2021) of results. Ultimately, the reliability and
85 associated uncertainty of any applied calibration will influence the final sensor data quality.

86 For end-users to make informed decisions on the applicability of air pollution sensors, a realistic understanding of the
87 expected performance in their chosen application is necessary (Rai et al., 2017). Despite this, there has been relatively
88 little progress in clarifying the performance of sensors for air pollution measurements outside of the academic arena.
89 This is largely due to the significant variability in both the number of sensors and the variety of applications tested,
90 **compounded by the proliferation of commercially available sensors/sensor systems with different configurations. as**
91 ~~well as the availability of highly accurate measurement instrumentation and/or regulatory networks to those outside~~
92 ~~of the atmospheric measurement academic field.~~ **Furthermore, the access to highly accurate measurement**
93 **instrumentation and/or regulatory networks remains limited for those outside of the atmospheric measurement**
94 **academic field** (e.g. Lewis and Edwards (2016) and Popoola et al. (2018)). From a UK clean air perspective, this
95 ambiguity represents a major problem. The lack of a consistent message undermines the exploitation of these devices'
96 unique strengths, notably their capability to form spatially dense networks with rapid time resolution. Consequently,
97 there is potential for a mismatch in users' expectations of what sensor systems can deliver and their actual operating
98 characteristics, eroding trust and reliability.

99 In this work, as part of the UK Clean Air program funded QUANT project, we deployed a variety of sensor
100 technologies (43 commercial devices, 119 gas and 118 PM measurements) at 3 representative UK urban sites —
101 Manchester, London and York— alongside extensive reference measurements, to generate the data for an
102 **comprehensive extensive** in-depth performance assessment. This project aims to not only evaluate the performance
103 of sensor devices in a UK urban climatological context but also provide critical information for the successful

104 application of these technologies in various environmental settings. To our knowledge, QUANT is the most extensive
105 and longest-running evaluation of commercial sensor systems globally to date. Furthermore, we tested multiple
106 manufacturers' data products, such as out-of-the-box data versus locally calibrated data, for a significant number of
107 these sensors to understand the implications of local calibration. This comprehensive approach offers unprecedented
108 insights into the operational capabilities and limitations of these sensors in real-world conditions. Significantly, some
109 of the insights gathered during QUANT have contributed to the development of the Publicly Available Specification
110 (PAS 4023, 2023), which provides guidelines for the selection, deployment, maintenance, and quality assurance of
111 air quality sensor systems. While this manuscript serves as an initial overview, detailed analyses of the measured
112 pollutants and study phases, offering a more comprehensive perspective on sensor performance, are planned for future
113 publications.

114 In the following sections, we delve into the methodology and provide an overview of the QUANT dataset, as well as
115 a discussion of some of the key findings and potential considerations for end-users.

116 2. QUANT study design

117 To capture the variability of UK urban environments, identical units were installed at three carefully selected field
118 sites. Two of these sites are highly instrumented urban background measurement supersites: the London Air Quality
119 Supersite (LAQS; for more details, refer here: https://uk-air.defra.gov.uk/networks/site-info?site_id=HPI) and the
120 Manchester Air Quality Supersite (MAQS; for more details, see: <http://www.cas.manchester.ac.uk/restools/firs/>),
121 located in densely populated urban areas with unique air quality challenges. The third site is a roadside monitoring
122 site in York, which is part of the Automatic Urban and Rural Network (AURN; click here for more details: [https://uk-
123 air.defra.gov.uk/networks/site-
124 info?uka_id=UKA00524&search=View+Site+Information&action=site&provider=archive](https://uk-air.defra.gov.uk/networks/site-info?uka_id=UKA00524&search=View+Site+Information&action=site&provider=archive)), representing a urban
125 environment more influenced by traffic. This selection strategy ensures that the QUANT study's findings reflect the
126 dynamics of urban air quality across different UK settings, while providing comprehensive reference measurements.
127 Further details about each site can be found in Section S1 in the Supp.

128 2.1 Main study

129 The Main QUANT assessment study aimed to perform a transparent long-term (19 Dec 2019 - 31 Oct 2022) evaluation
130 of commercially available sensor technologies for outdoor air pollution monitoring in UK urban environments. Four
131 units ~~duplicates~~ of five different commercial sensor devices (Table 1) were purchased in Sept 2019 for inclusion in
132 the study, with the selection criteria being: market penetration and/or previous performance reported in the literature,
133 ability to measure pollutants of interest (e.g. NO₂, NO, O₃, and PM_{2.5}), and capacity to run continuously reporting
134 high time resolution data (1-15 min data) ideally in near real-time (i.e., available within minutes of measurement)
135 with data accessible via an API.

136 **Table 1. Main QUANT devices description. The 20 units, all commercially available and ready for use as-is, offered 56 gas
137 and 56 PM measurements in total. For a detailed description of the devices see Section S3† in the Supp.**

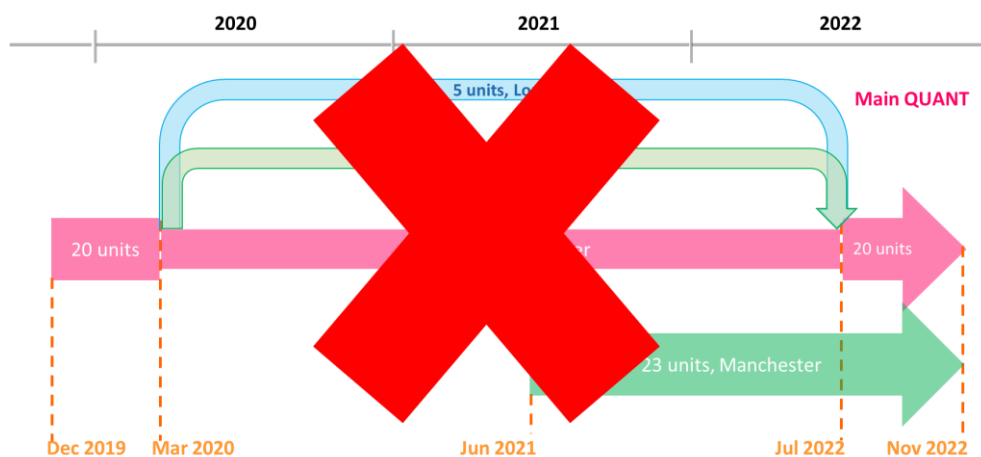
Product*	Measurements	Cost (£)**
----------	--------------	------------

(# units)	Company ³	NO	NO ₂	O ₃	CO	CO ₂	PM ₁	PM _{2.5}	PM ₁₀	
AQY (4)	Aeroqual	-	✓	✓	-	-	-	✓	✓	~4.7K
AQM (4)	AQMesh	✓	✓	✓	-	✓	✓	✓	✓	~8.6K
Ari (4)	QuantAQ	✓	✓	✓	✓	✓	✓	✓	✓	~8.6K
PA (4)	PurpleAir	-	-	-	-	-	✓	✓	✓	~0.3K
Zep (4)	Earthsense	✓	✓	✓	-	-	✓	✓	✓	~7K

*AQY: Aeroqual; AQM: AQMesh; Ari: Arisense; PA: PurpleAir; Zep: Zephyr. **Cost (Sep 2019) per unit including UK taxes and associated contractual costs (i.e., communication, data access, sensor replacement, etc.).

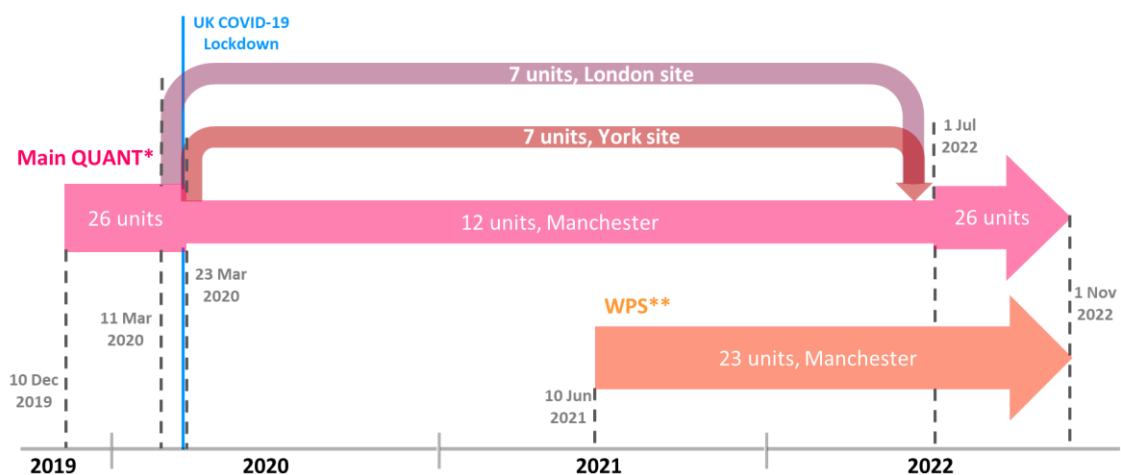
138 ~~To capture the variability of UK urban environments, identical units were installed at three carefully selected field~~
139 ~~sites. Two of these sites are highly instrumented urban background measurement supersites: the London Air Quality~~
140 ~~Supersite (LAQS) and the Manchester Air Quality Supersite (MAQS), located in densely populated urban areas with~~
141 ~~unique air quality challenges. The third site is a roadside monitoring site in York, which is part of the Automatic~~
142 ~~Urban and Rural Network (AURN, <https://uk-air.defra.gov.uk/data/>), representing a urban environment more~~
143 ~~influenced by traffic. This selection strategy ensures that the QUANT study's findings reflect the dynamics of urban~~
144 ~~air quality across different UK settings, while providing comprehensive reference measurements. Further details about~~
145 ~~each site can be found in Section S3 in the Supp., and the available reference instrumentation in Section S4.~~

146 Initially, all the sensors were deployed in Manchester for approximately 3 months (mid-Dec 2019 to mid-Mar 2020)
147 before being split up amongst the three sites (Fig. 1). At least one unit per brand was re-deployed to the other two
148 sites (mid-March 2020 to early-July 2022) leaving two devices per company in Manchester to assess inter-device
149 consistency. In the final 4 months of the study, all the sensor systems were relocated back to Manchester (early July
150 2022 to the end of October 2022).



151

³ Throughout this article, the terms “manufacturers” and “company” are used interchangeably to refer to entities that produce, and/or sell sensor systems or devices. This usage reflects the industry practice of referring to businesses involved in the production and distribution of technology products without distinguishing between their roles in manufacturing or sales.



*: Aeroqual (x4), AQMesh (x4), Zephyr (x4), QuantAQ (x4), PurpleAir (x10)

** : AQMesh (x3), Bosch (x2), Clarity (x3), Kunak (x3), Oizom (x2), QuantAQ (x3), South Coast Science (x2), Respirer Living Sciences (x2), Vortex (x3)

152

153 **Figure 1. Main QUANT Quant and Wider Participation Study (WPS) timeline.**

154 **2.2 Wider Participation Study**

155 The Wider Participation Study (WPS) was a no-cost complementary extension of the QUANT assessment, specifically
 156 designed to foster innovation within the air pollution sensors domain. This segment of the study took place entirely at
 157 the MAQS from 10th June 2021 to 31st October 2022 (Fig. 1). It included a wider array of commercial platforms (9
 158 different sensor systems brands), and offered manufacturers the opportunity to engage in a free-of-charge impartial
 159 evaluation process. Although participation criteria matched those of the Main QUANT study, a key distinction lay in
 160 the voluntary nature of participation: ~~manufacturers vendors~~ were invited to contribute multiple sensor devices
 161 throughout the WPS study (see Table 2). Participants were able to demonstrate their systems' performance against
 162 collocated high-resolution (1-minute) reference data at a state-of-the-art measurement site such as the Manchester
 163 supersite.

164 **Table 2. The 23 WPS devices deployed at the Manchester supersite, all commercially available and ready for use as-is,**
 165 **provided 63 gases and 62 PM measurements in total. For a detailed description of the devices see the Section S43 in the**
 166 **Supp.**

Product* (# units)	Company	Measurements							
		NO	NO ₂	O ₃	CO	CO ₂	PM ₁	PM _{2.5}	PM ₁₀
Mod (3)	QuantAQ	-	-	-	-	-	✓	✓	✓
AQM (3)	AQMesh	✓	✓	✓	✓	✓	✓	✓	✓
Atm (2)	RLS**	-	-	-	-	-	✓	✓	✓
IMB (2)	Bosch	-	✓	✓	-	-	-	✓	✓
Poll (2)	Oizom	✓	✓	✓	✓	✓	-	✓	✓
AP (3)	Kunak	✓	✓	✓	✓	✓	✓	✓	✓
SA (3)	Vortex IoT	-	✓	✓	-	-	-	✓	✓
NS (3)	Clarity	-	✓	-	-	-	✓	✓	✓

Prax (2) SCS*** ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓

167 *Mod: Modulair; AQM: AQMesh; Atm: Atmos, Poll: Polludrone ~~Polludrone: Poll~~; AP: Kunak Air Pro; SA: Silax Air, NS: Node-
168 S, Prax: Praxis. **RLS: Respirer Living Sciences. ***SCS: South Coast Science.

169 2.3 Sensor deployment and data collection, ~~co-located reference data and data products~~

170 All sensor devices were installed at the measurement sites as per manufacturer recommendations, adhering strictly to
171 manufacturers' guidelines for electrical setup, mounting, cleaning, and maintenance guaranteed proper installation.
172 Since all deployed systems were designed for outdoor use, no additional protective measures were necessary. Each of
173 the systems were mounted on poles acquired specifically for the project or on rails at the co-location sites, without the
174 need for special protections. Following the manufacturer's suggestions, sensors were positioned within 3 metres of
175 the reference instruments' inlets. Custom electrical setups were developed for each sensor type, incorporating local
176 energy sources and weather-resistant safety features, alongside security measures to deter vandalism and ensure
177 uninterrupted operation. Routine maintenance was conducted monthly, although the COVID-19 pandemic
178 necessitated longer intervals between visits. Despite these obstacles, efforts to maintain sensor security and
179 functionality continued unabated, employing both physical safeguards and remote monitoring to preserve data
180 integrity.

181 In addition to the device supplier's own cloud storage (accessed on-demand via each supplier's web portals), an
182 automated daily scraping of each company's API was performed to save data onto a secure server at the University
183 of York to ensure data integrity. ~~PurpleAir units were exempt from this due to a lack of mobile data connection and
184 poor internet signal at the sites; instead, readings were locally collected and manually uploaded.~~ Unlike other brands
185 that utilise mobile data connections, PurpleAir sensors rely on WiFi for data transmission. Due to poor internet signal
186 at the sites, we locally collected and manually uploaded readings for these units. Minor pre-processing was applied at
187 this stage, including temporal harmonisation to ensure that all measurements had a minimum sampling period of 1-
188 minute, ensuring consistency in measurement units and labels, and coercing into the same format to allow for full
189 compatibility across sensor units. No additional modifications to the original measurements were applied; missing
190 values were kept as missing and no additional flags were created based on the measurements beyond those provided
191 by the manufacturers. ~~No outlier checks or data modifications were applied at this stage.~~ For an overview of the sensor
192 measurands and their corresponding data time resolutions as provided by the companies participating in the Main
193 QUANT study and the WPS, please see Seccion S3 and S4 (Table S4 and S5) respectively.

194 2.4 Data products and co-located reference data

195 In addition to providing an independent assessment of sensor performance, QUANT also aimed to contribute to device
196 manufacturers to help advance the field of air pollution sensors. During QUANT, device calibrations were performed
197 solely at the discretion of the manufacturers without any intervention from our team, thus limiting the involvement of
198 manufacturers in the provision of standard sensor outputs and unit maintenance as would be required by any standard
199 customer. This approach enabled manufacturers to independently assess and benchmark their sensors' performance,
200 using provided reference data to potentially develop calibrated data products. It's noteworthy that not all manufacturers
201 chose to utilise these data for corrections or enhancements. However, those who did were expected to create and
202 submit calibrated data products, subsequently named as "out-of-box" (initial data product), "cal1" (first calibrated
203 product), and "cal2" (second calibrated product). This differentiation highlighted the varying degrees of engagement

204 and application of the reference data by different manufacturers. Figures S2 and S3 (section S3 and S4 respectively)
205 show a time-line of the different data products.

206 To this end, three separate 1-month periods of reference data, spaced every 6 months, were shared with each supplier,
207 provisional data soon after each period, and ratified data when available. ~~For an overview of reference instrumentation
208 at each site refer Table S1, and for details on the quality assurance procedures applied to the reference instruments
209 see Table S2.~~ All reference data were embargoed until it was released to all manufacturers simultaneously to ensure
210 consistency across manufacturers. For an overview of reference and equivalent-to-reference instrumentation, as
211 defined in the European Union Air Quality Directive 2008/50/EC (hereafter referred to as EU AQ Directive), at each
212 site, please refer to Section S2 (Table S1). For details on the quality assurance procedures applied to the reference
213 instruments, see Table S2. To see the dates and periods of the shared reference data refer to Table S3. ~~Access to
214 collocated reference data allowed the companies to assess sensors' performance and, if they chose, to generate and
215 provide additional calibrated data products. These products are distinct data versions provided by manufacturers
216 throughout QUANT, before and/or after sharing reference data—for instance, “out of box”, “cal1”, “cal2”, etc.
217 Figures S1 and S2 show a time line of the different data products. To see the dates and periods of the shared reference
218 data refer to Table S3. All reference data was embargoed until it was released to all manufacturers simultaneously to
219 ensure consistency across manufacturers. Not every manufacturer opted to use this data to apply corrections or
220 improve calibrations, but if they chose to do so, the updated measurements were treated as a separate data product.
221 Device calibrations were performed solely at the discretion of the manufacturers without any intervention from our
222 team, thus limiting the involvement of vendors/manufacturers in the provision of standard sensor outputs and unit
223 maintenance as would be required by any standard customer.~~

224 3. Results and discussion

225 A key challenge in sensor performance evaluation is the high spatial and temporal variability errors that impact the
226 accuracy of their readings, making the application of laboratory corrections more challenging. ~~Furthermore, the
227 overreliance on global performance metrics, such as R^2 (i.e., the Coefficient of Determination), RMSE (i.e., the Root
228 Mean Squared Error), and MAE (i.e., the Mean Absolute Error) is an important issue when assessing sensors. While
229 these metrics provide a general understanding of sensor performance, they can be limiting or even misleading,
230 restricting a comprehensive understanding of the error structure and the measurement information content (Diez et
231 al., 2022).~~ Furthermore, the overreliance on global performance metrics is a significant concern in sensor assessment.
232 The Coefficient of Determination (R^2), Root Mean Squared Error (RMSE), and Mean Absolute Error (MAE) are
233 among the most popular single-value metrics for evaluating sensor performance, alongside others (e.g., the bias, the
234 slope and intercept of the regression fit). However, while single-value metrics offer an overview of performance, they
235 can be limiting or misleading. They condense vast amounts of data into a single value, simplifying complexity at the
236 expense of a nuanced understanding of error structures and information content (Diez et al., 2022), potentially
237 overlooking critical aspects of sensor performance (Chai & Draxler, 2014). Visualisation tools (such as Regression
238 plots, Target plots, and Relative Expanded Uncertainty plots) complement these metrics, allowing end users to identify
239 relevant features, which could be beyond the scope of global metrics. For additional details on the metrics utilised in
240 this study, including some of their limitations and advantages refer to section “S5. Performance Metrics”. This section
241 also provides a summary of current guidelines and standardisation initiatives, which may offer a foundation for end-
242 users to select appropriate metrics for their own analyses (refer to table S6). For further discussion on metrics and
243 visualisation tools for performance evaluation, readers are directed to Diez et al. (2022).

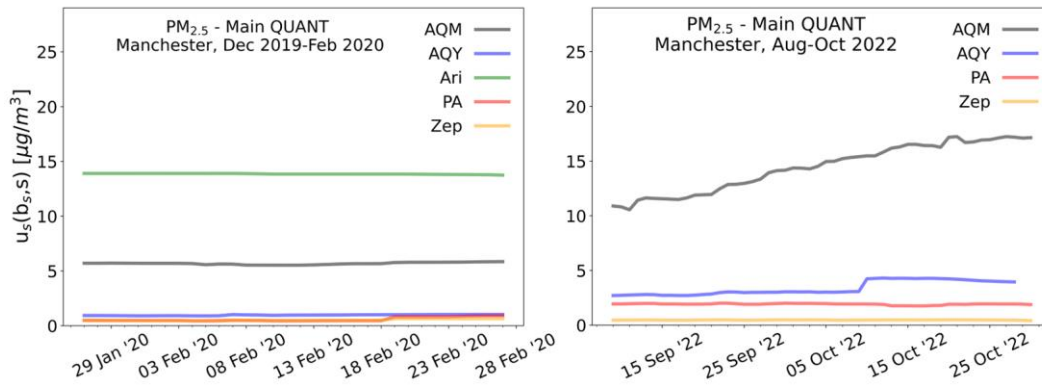
244 In response to these challenges, the QUANT assessment represents the most extensive independent appraisal of air
245 pollution sensors in UK urban atmospheres. As the results presented here illustrate, QUANT is dedicated to examining
246 sensor performance through multiple complementary ~~perspectives and~~ metrics and visualisation tools, aiming to
247 integrate these to accurately reflect the complexity of this dataset. This methodology promotes a nuanced
248 understanding of sensor performance, extending beyond the limitations of conventional global single-value metrics.

249 Furthermore, by providing open access to the dataset, we encourage stakeholders to explore and utilise the data
250 according to their unique needs and contexts, as detailed in the “Data Availability” section. In addition, we have
251 developed a publicly accessible analysis platform (<https://shiny.york.ac.uk/quant/>), designed for straightforward
252 offline analysis of the QUANT dataset. This platform enables users to interactively visualise the data through various
253 representations, such as time series, regression plots, and Bland-Altman plots. It also offers statistical parameters
254 (including regression equation, R^2 , and RMSE) for analysing different pollutants, selecting specific sensors or
255 manufacturers, and comparing across various co-location timeframes.

256 The following sections aim to provide an overview of the data and provide initial findings, with a focus on those that
257 are most relevant to end-users of these technologies. The majority of examples presented here focus on $PM_{2.5}$ and
258 NO_2 measurements, due to both a larger dataset available for these pollutants and their critical role in addressing the
259 exceedances that predominantly impact UK air quality. All metrics and plots presented here are based on 1-hour
260 averaged data. Unless otherwise specified, a data inclusion criterion of 75% was uniformly applied across our analyses
261 to ensure the reliability and representativeness of the results. This threshold aligns with the EU AQ Directive, which
262 mandates this proportion when aggregating air quality data and calculating statistical parameters. To highlight broad
263 implications and insights into sensor technology, rather than focusing on the performance of specific manufacturers,
264 figures illustrating brand-specific features have been anonymized. This is intended to prevent potential bias and
265 encourage a holistic view of the data, ensuring interpretations remain focused on general trends rather than isolated
266 examples.

267 3.1 Inter-device precision

268 Inter-device precision refers to the consistency of measurements across multiple identical devices (i.e., same brand
269 and model) ~~of the same type~~, an important characteristic to ensure the reliability of sensor outputs over time (Moreno-
270 Rangel et al., 2018). During QUANT, all the devices were collocated for the first 3 months and the final 3 months of
271 the deployment to assess inter-device precision and its changes over time. Fig. 2 shows the inter-device precision (as
272 defined by the CEN/TS 17660-1:2021, i.e., the “between sensor system uncertainty” metric: $u_s(b_s, s)$) of $PM_{2.5}$
273 measurements during these periods. For an overview of NO_2 and O_3 inter-device precision, see the “S6.
274 Complementary plots” section in the supplementary (figures S4 and S5). While most of the companies display a
275 certain level of inter-device precision stability in each period (except for one, with a seemingly upward trend in the
276 final period), there are evident long-term changes. Notably, out of the four manufacturers assessed in the final period
277 (each having 3 devices running simultaneously), three experienced a decline in their inter-device precision compared
278 to two years earlier. This is likely due to both hardware degradation but also drift in the calibration, which at this point
279 had been applied between 16 and 34 months prior (depending on the manufacturer). For extended periods,
280 inconsistencies among devices from the same manufacturer might emerge, leading to varying readings under similar
281 conditions. Consequently, data collected from different devices may not be directly comparable, which could result
282 in inaccuracies or misinterpretations when analysing air quality trends or making decisions.

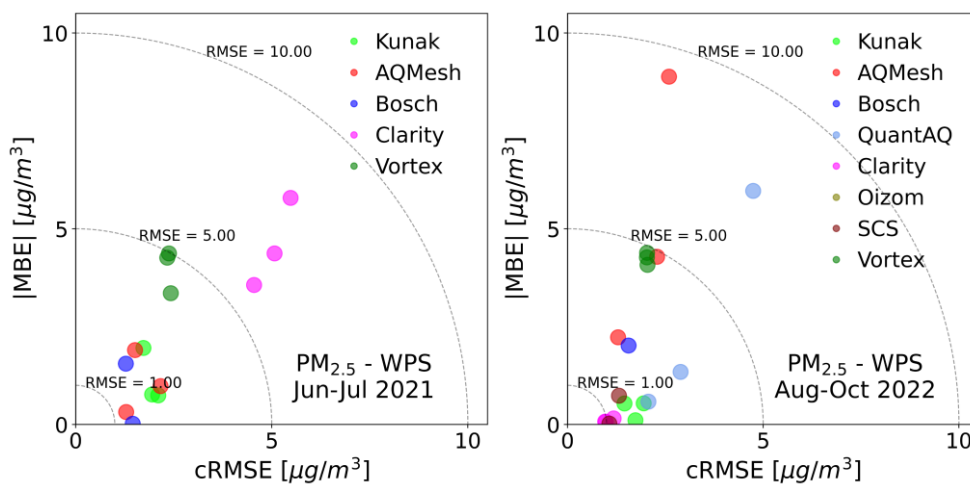


283
 284 **Figure 2. The inter-device precision of PM_{2.5} measurements from “identical” devices across the 5 companies participating**
 285 **in QUANT is assessed using the “between sensor system uncertainty” metric (defined by the CEN/TS 17660-1:2021 as $u(b_s, s)$).**
 286 **Each line represents this metric as a composite of all sensors per brand (excluding units with less than 75% data) within**
 287 **a 40-day sliding window.**

288 It is worth noting that the inter-device precision provides no information on the accuracy of the sensor measurements;
 289 a batch of devices may provide a highly consistent, but also highly inaccurate measurement of the target pollutant.

290 The “target plot” (as shown in Fig. 3) is a tool commonly used to depict the bias/variance decomposition of an
 291 instrument’s error relative to a reference (for more details see Jolliff et al. (2009)). The mean bias error (MBE) is used
 292 to characterise accuracy and precision is quantified by the centered Root Mean Squared Error (cRMSE, e.g. Kim et
 293 al. (2022) also called unbiased Root Mean Squared Error (uRMSE, e.g. Guimarães et al. (2018)). Fig. 3 visualises the
 294 performance of a set of PM_{2.5} sensors of the WPS deployment for the first 2 months (out-of-box data) and the last 3
 295 months of colocation (manufacturer-supplied calibrations). ~~In addition to highlighting which devices are most~~
 296 ~~accurate, Fig. 3 also provides an additional perspective of inter device precision.~~ In addition to showcasing inter-
 297 device precision, Fig. 3 also serves as a transition to accuracy evaluation (the focus of the subsequent section).

298

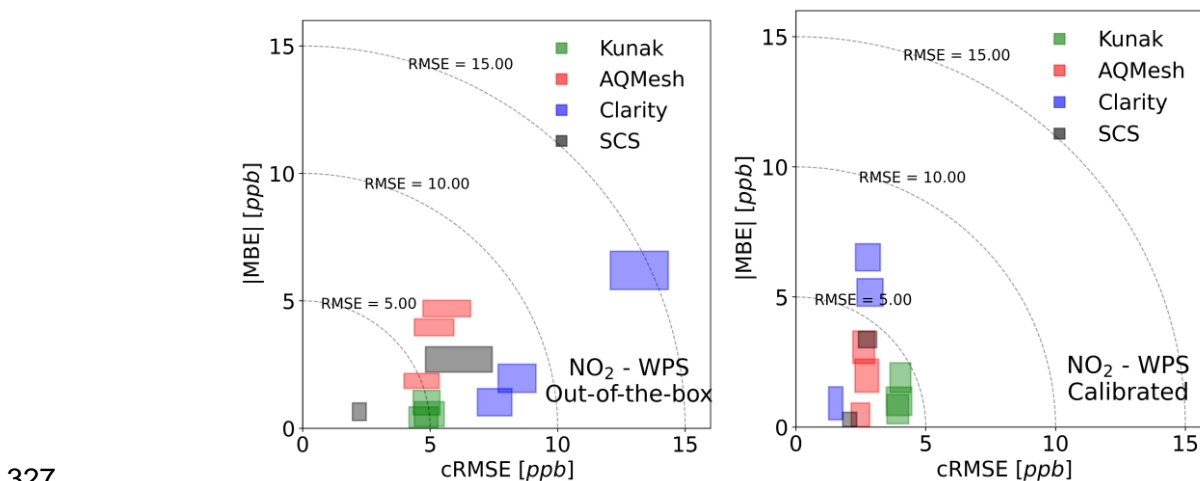


299
 300 **Figure 3. Target diagrams for the WPS PM_{2.5} measurements during the initial co-location period (Jun-Jul 2021, left) and**
 301 **final co-location period (Aug-Oct 2022, right). The error (RMSE) for each instrument is decomposed into the MBE (y-axis)**
 302 **and cRMSE (x-axis). Each point represents an individual sensor device, with duplicate devices having the same colour.**

303 Since only units with more than 75% of the data were considered, the plot on the right shows fewer units than the plot on
304 the left.

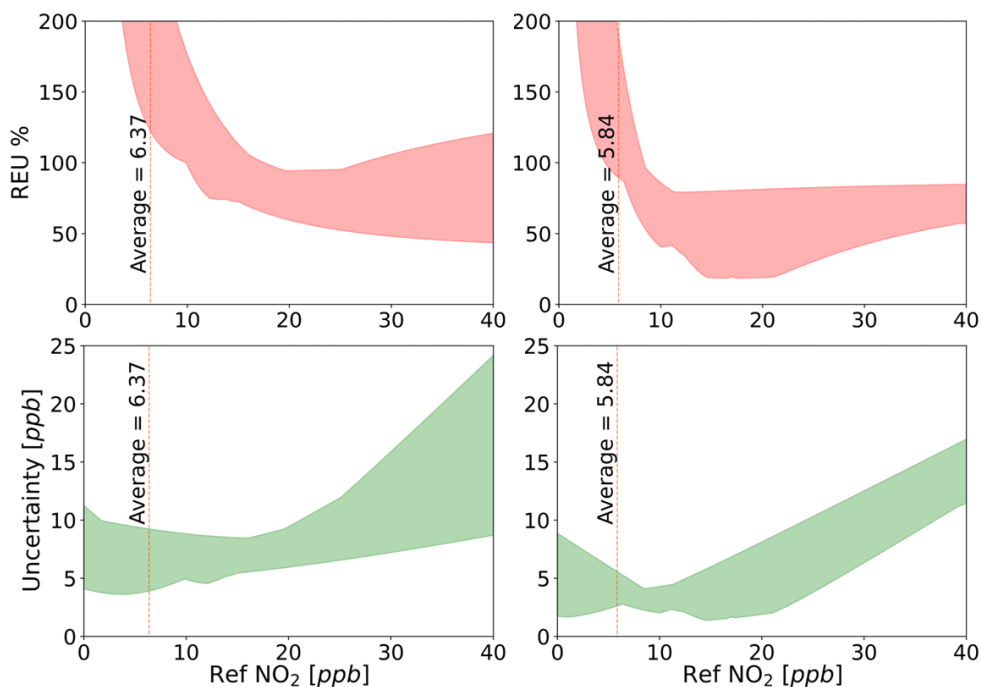
305 3.2 Device accuracy and co-location calibrations

306 Sensor measurement accuracy denotes how close a sensor's readings are to reference values (Wang et al., 2015).
307 Characterising this feature is imperative for establishing sensor reliability and making informed decisions based on
308 its data. Fig. 4 shows that co-location calibration can greatly impact observed NO₂ sensor performance in a number
309 of ways. Firstly, measurement bias is often, but not always, reduced following calibration, as evidenced by a general
310 trend for devices to migrate towards the origin (RMSE = 0 ppb). Secondly, it can help to improve within-manufacturer
311 precision ~~by grouping sensor systems from the same company closer together~~, as evidenced by sensor systems from
312 the same company grouping more closely as the right plot in Fig. 4 shows. The figure also highlights a fundamental
313 challenge with evaluating sensor systems: the measured performance can vary dramatically over time —and space—
314 as the surrounding environmental conditions change. To quantify this, 95% Confidence Intervals (CIs) were estimated
315 for each device using bootstrap simulation and are visualised as a shaded region. For the out-of-the-box data, these
316 regions are noticeably larger than in the calibrated results for most manufacturers, suggesting that colocation
317 calibration has helped to tailor the response of each device to the specific site conditions. ~~This is reinforced by the~~
318 ~~eRMSE component reducing by a greater extent than the MBE; in the terminology of machine learning, the calibration~~
319 ~~has helped reduce the variance portion of the bias variance trade-off.~~ This observation suggests that colocation
320 calibration effectively improves each device's response to particular site conditions. This improvement is underscored
321 by the more substantial reduction in the cRMSE component compared to the MBE. The cRMSE, representing the
322 portion of error that persists after bias removal, essentially measures errors attributable to variance within the data
323 space. In the context of out-of-the-box data, this “data space” spans all potential deployment locations used by
324 manufacturers for initial calibration model training (i.e., before shipping the sensors for the QUANT study), thus
325 exhibiting high variability. However, applying site-specific calibration significantly narrows this variability,
326 leveraging local training data to minimise variance.



327
328 **Figure 4. Effect of colocation calibration on NO₂ sensor accuracy. The accuracy is quantified using RMSE, which is**
329 **decomposed into MBE (y-axis) and cRMSE (x-axis). 95% confidence regions were estimated using bootstrap sampling. The**
330 **left panel displays results from the period Jun - Jul 2021 (‘out-of-the-box’ data), while the right-hand panel summarises**
331 **Aug 2021 when calibrations were applied for all the WPS manufacturers.**

332 However, it is important to note a limitation of Target Plots: they primarily focus on sensor behaviour around the
 333 mean. Therefore, the collective improvement evidenced by Fig. 4 might be only partial. For applications where it is
 334 important to understand how calibrations impact lower or higher percentiles, considering other metrics or visual tools
 335 would be advisable. An example of this is the absolute and Relative Expanded Uncertainty (REU, defined by the
 336 Technical Specification CEN/TS 17660-1:202). Unlike the more commonly used metrics such as R^2 , RMSE, and
 337 MAE, which measure performance of the entire dataset, the REU offers a unique “point by point” evaluation, enabling
 338 its representation in various graphical forms, such as time series or concentration space (for the REU mathematical
 339 derivation, refer to section “S5. Performance Metrics”). The REU approach also incorporates the uncertainty of the
 340 reference method into its assessment, highlighting the intrinsic uncertainty present in all measurements, including
 341 those from reference instruments. This consideration of reference uncertainty is crucial for a holistic understanding of
 342 sensor performance and calibration effectiveness. For a comprehensive discussion on this, refer to Diez et al. (2022).
 343 Fig. 5 illustrates how NO₂ calibrations might not only improve collective performance around the mean (as indicated
 344 by the dotted red line in Fig. 5 and previously displayed in the target plot) but across the entire concentration range.



345
 346 **Figure 5. The top plots display the REU (%) across the concentration range, while the bottom plots depict the Absolute**
 347 **Uncertainty (ppb) —both before (left plots) and after (right plots) calibrating NO₂ WPS systems. The shaded areas**
 348 **represent the collective variability evolution (all sensors from all companies) of both metrics. These plots were constructed**
 349 **using the minimum and maximum value of the REU and the Absolute Uncertainty for the entire concentration range.**

350 However, a note of caution when interpreting results from observational studies such as these is that it is impossible
 351 to ascertain a direct causal relationship between calibration and sensor performance as there are numerous other
 352 confounding factors at play (Diez et al., 2022). Notably these two data products are being assessed over different
 353 periods when many other factors will have changed, for example, the local meteorological conditions as well as
 354 human-made factors such as reduced traffic levels following the COVID-19 lockdown that commenced in March
 355 2020.

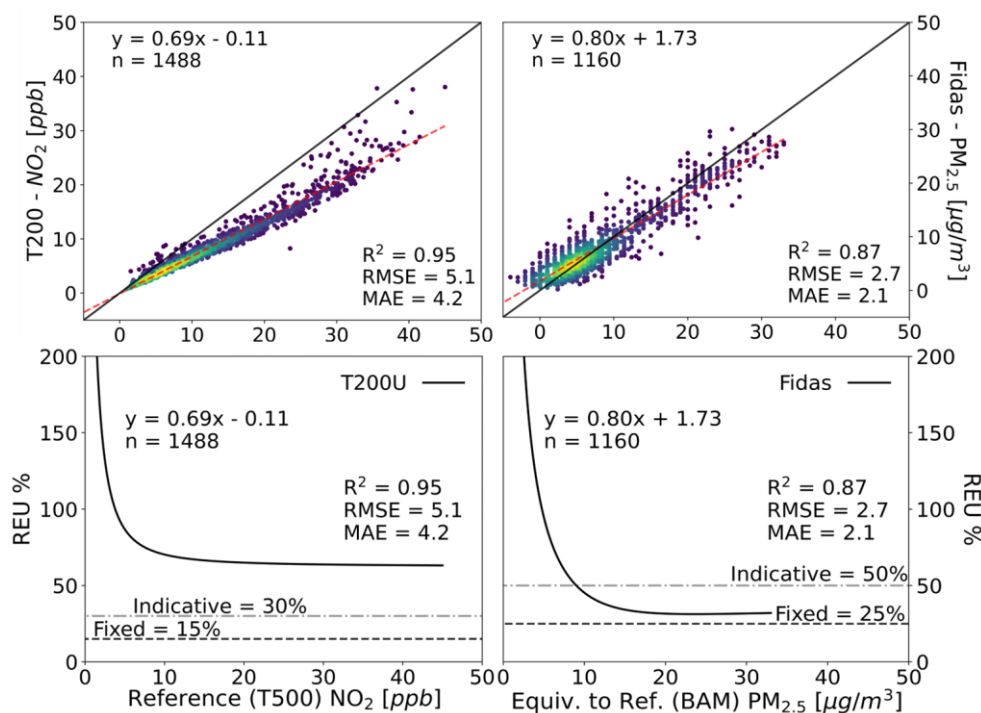
356 3.3 Reference instrumentation is key

357 A common assumption when evaluating the performance of sensors is that the metrological characteristics of the
358 sensor predominantly influence discrepancies detected in co-locations. While this presumption can often be justified
359 due to both devices' (sensor and the reference method) relative scales of measurement errors, it is not always the case.
360 Since every measurement is subject to uncertainties, it is crucial to consider those associated with the reference when
361 deriving the calibration factors of placement.

362 Fig. 6 (left plots) displays the performance of a NO₂ reference instrument (Teledyne T200U) specifically installed for
363 QUANT, located next to the usual instrument at the Manchester supersite (Teledyne T500). Although they use
364 different analytical techniques (chemiluminescence for the T200U and Cavity Attenuated Phase Shift Spectroscopy
365 for the T500), their measurements are highly correlated ($R^2 \sim 0.95$). However, it's possible to identify a proportional
366 bias (slope=0.69), attributed to retaining the initial calibration (conducted in York) without subsequent adjustments,
367 a situation exacerbated by an unnoticed mechanical failure of one of the instrument's components. The REU
368 demonstrates that, under these circumstances, an instrument designated as a reference does not meet the minimum
369 requirements ($REU \leq 15\%$ for NO₂ reference measurements) set out by the Data Quality Objectives (DQOs) of
370 the EU AQ Air Quality Directive 2008/50/EC. Figure S63 shows a unique sensor evaluated against both the T500 and
371 the T200U. The comparison against the T200U yields better results, suggesting that, in a hypothetical scenario where
372 it was the only instrument at the site, this could lead to misleading conclusions. This situation reinforces the idea that
373 instruments should not only be adequately characterised but also undergo rigorous quality assurance and data quality
374 control programs, as well as receive appropriate maintenance (Pinder et al., 2019). All of this must be performed
375 before and during the use of any instrument.

376 For PM monitoring, the current EU reference method is the gravimetric technique (CEN EN 12341, 2023), which is
377 a non-continuous monitoring method that requires weighing the sampled filters and off-line processing of the results.
378 Techniques that have proven to be equivalent to the reference method (called "equivalent to reference" in the EU AQ
379 Air Quality Directive) are very often used in practice. In the UK context, the Beta Attenuated Monitor (BAM) and
380 FIDAS (optical aerosol spectrometer) are equivalent-to-reference methods commonly used as part of the Urban
381 AURN Network (Allan et al., 2022). To illustrate these differences in practice, Fig. 6 compares these two equivalent-
382 to-reference PM_{2.5} measurements obtained with a BAM (AURN York site, located on a busy avenue), and a FIDAS
383 unit specifically installed for QUANT. During this specific period, they show a strong linear association ~~do not fully~~
384 ~~agree~~ ($R^2 = 0.87$). Although the bias is not extremely pronounced (slope=0.80), the FIDAS measurements are, on
385 average, systematically lower compared to BAM. ~~Despite a not very pronounced bias (slope=0.80), the dispersion of~~
386 ~~points around the best fit line is noticeable, limiting the linearity of the FIDAS compared to the BAM.~~

387 In the hypothetical case that the BAM were to be considered the reference method (arbitrarily chosen for this example
388 as it is the current instrument at the AURN York site) when assessing the FIDAS under these test conditions, it would
389 only meet the criterion stipulated by the EU DQOs for indicative measurements ($REU \leq 25\%$ for PM_{2.5}), but not
390 for fixed (i.e., reference) measurements ($REU \leq 50\%$ for PM_{2.5}). ~~Of course,~~ This example is primarily intended to
391 illustrate the magnitude of differences between both methods for this particular application, and by no means does
392 this observation imply that the FIDAS measurements are inherently problematic.

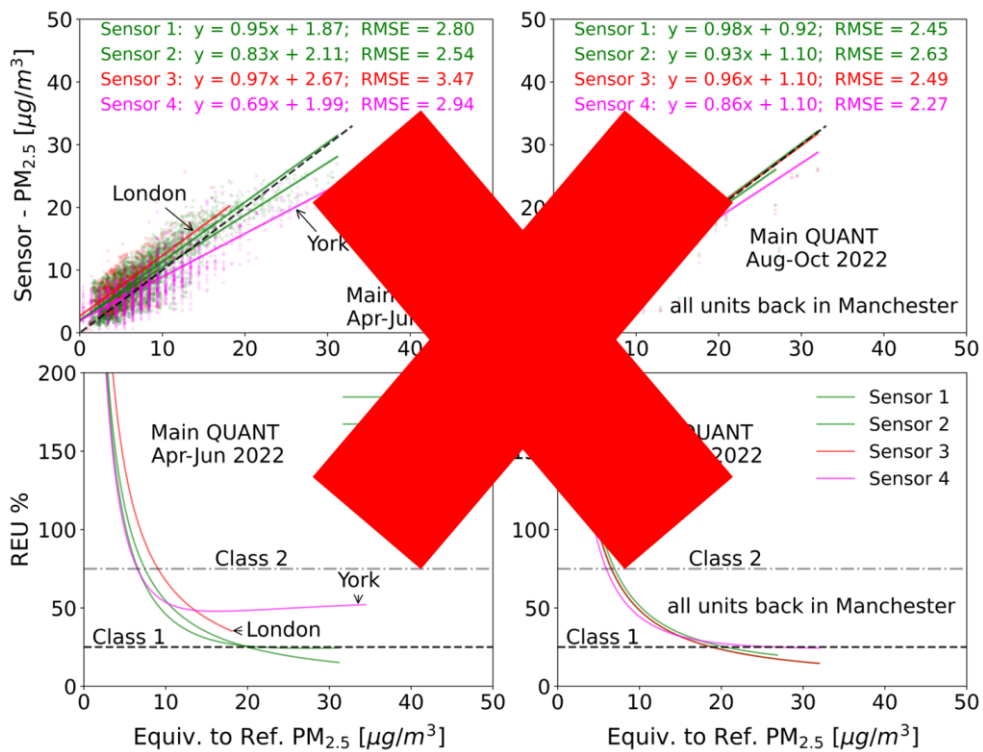


393
 394 **Figure 6. The left plots depict the comparison between the Teledyne T200U (chemiluminescence analyzer) and the reference**
 395 **method (Teledyne T500 CAPS analyzer) at the Manchester supersite. The plots to the right illustrate PM_{2.5} measurements**
 396 **in York, taken with a FIDAS instrument (optical aerosol spectrometer) and a BAM 1020 (beta attenuation monitor), both**
 397 **equivalent-to-reference methods. While the top plots show the regression (including some typical single-value metrics),**
 398 **those on the bottom present the REU alongside the DQOs defined by the EU AQ Directive [European Directive 2008/50/EC](#).**

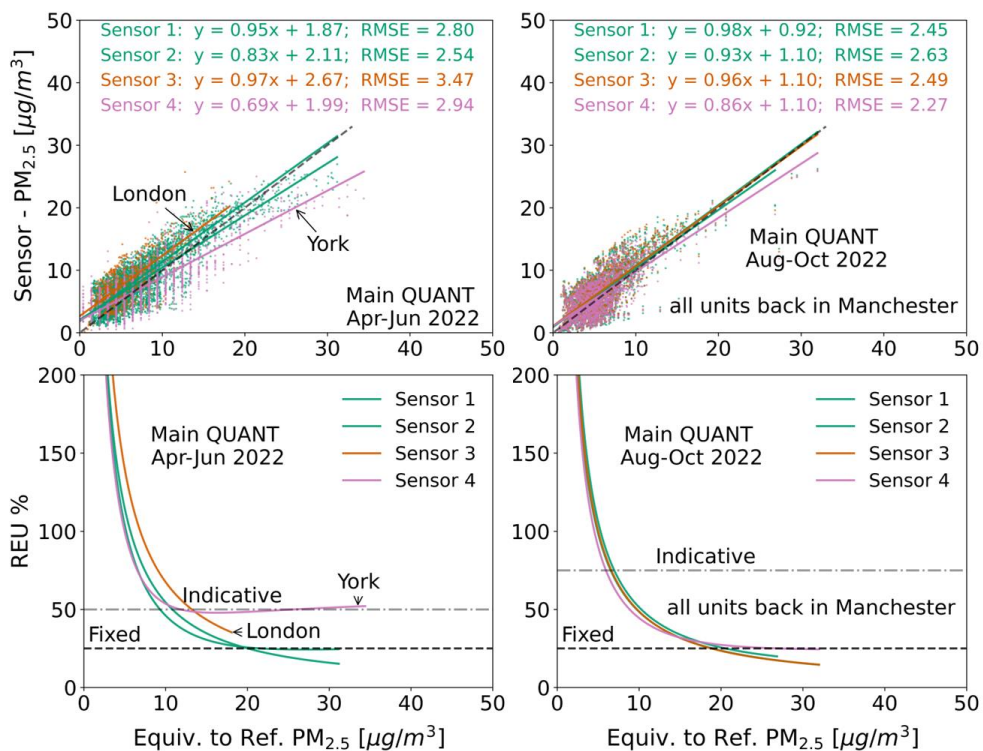
399 Although these two instruments (BAM and Fidas) show a greater concordance between themselves than with sensors
 400 (for the comparison of two sensor systems against the BAM and the Fidas, refer to Fig. S74), the choice of the
 401 measurement method can have a considerable impact on evaluations of this type. This underscores the importance of
 402 adequately characterising the uncertainties of the reference monitor when evaluating sensors.

403 3.4 Inter-location performance ~~Systems performance after location transfer~~

404 An extreme example of sensor performance varying due to environmental conditions is when sensors are moved
 405 between locations, as their apparent performance may vary drastically. Fig. 7 displays the REU and regression plots
 406 for four of the same PM_{2.5} sensor system in two periods: April-June 2022 when the devices were working across the
 407 3 sites (York, Manchester and London), and August-October 2022 when they were all reunited in Manchester. The
 408 RMSE remains reasonably consistent (range 2.27 to 3.47 ppb) between the devices across the periods and locations.
 409 However, for the device that moved from York to Manchester, a change in slope from 0.69 to 0.86 was observed.
 410 Because this device's slope is consistent with the other units while running in Manchester, this is likely due to the
 411 different sensor responses in the specific environments. The precise cause of this change is not immediately evident
 412 and will be the focus of a follow-up study, but could be due to changes in local conditions (e.g., weather, emissions,
 413 etc.) impacting sensor calibration and/or differences in actual PM_{2.5} sources and particle characteristics at the sites
 414 (Raheja et al., 2022).



415

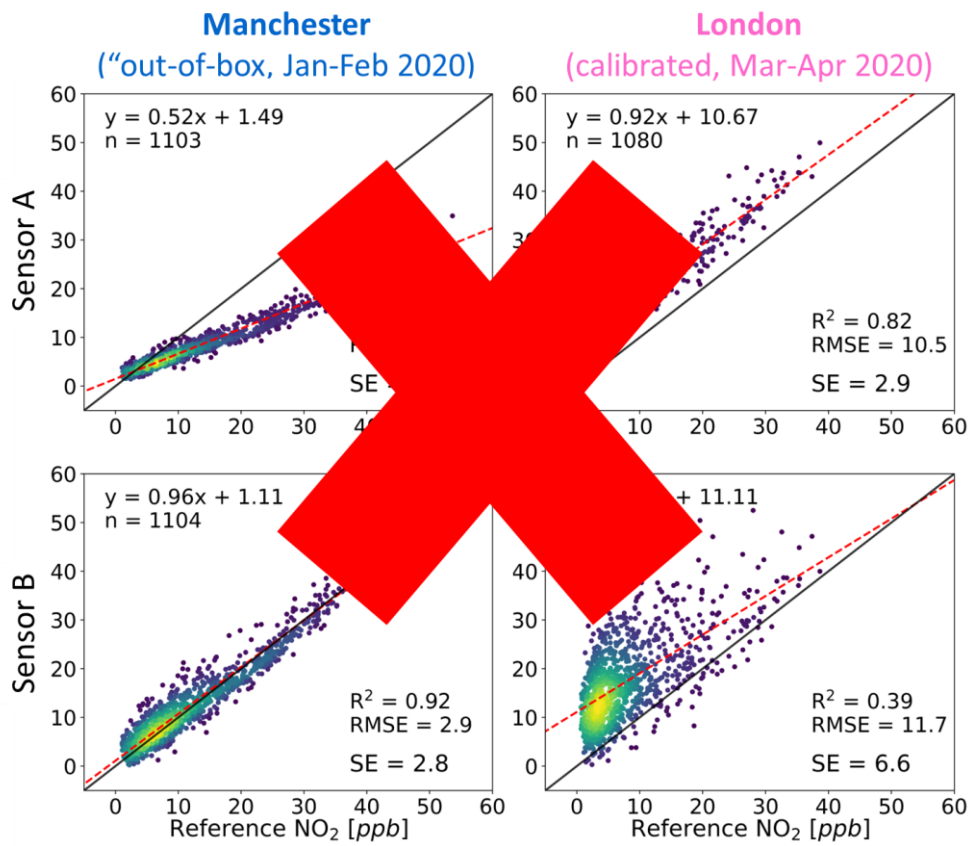


416

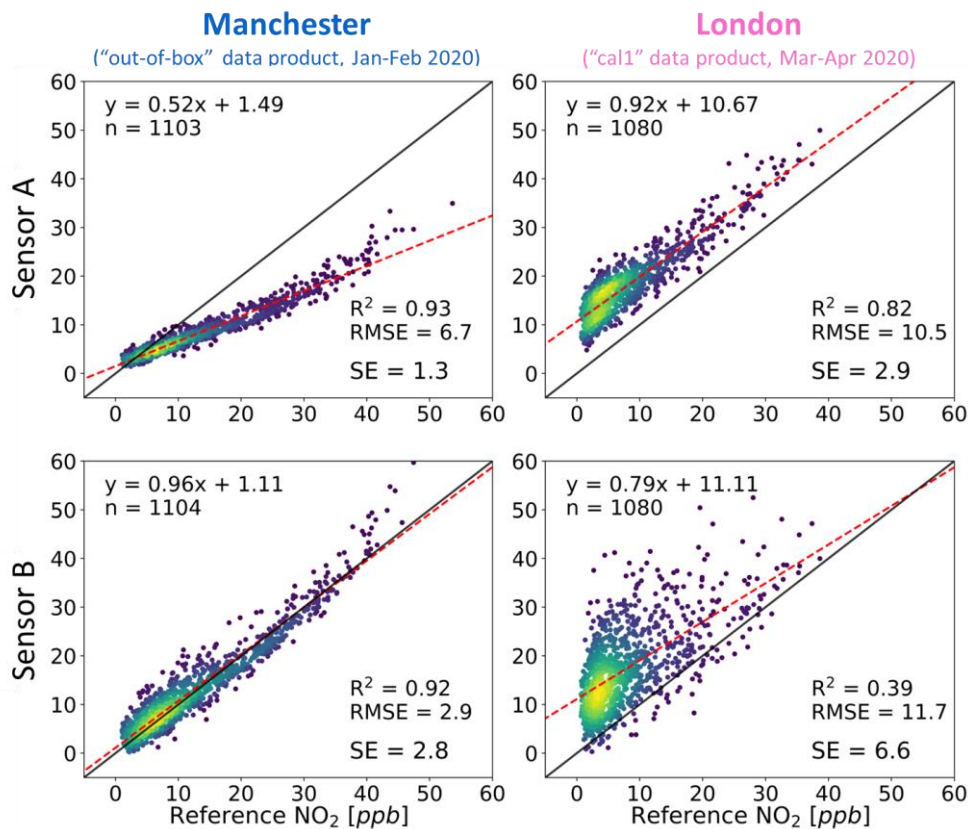
417 **Figure 7. Regression (top) and REU (bottom) plots showing data from four PM_{2.5} sensors (same manufacturer) over 2 time**
 418 **periods: Apr-Jun 2022 and Aug-Oct 2022. The four devices were in separate locations in the first period, but all deployed**
 419 **in Manchester in the second. The horizontal dashed lines represent a reference for the PM_{2.5} DQOs as defined by the EU**
 420 **AQ Directive (for “fixed” PM_{2.5} measurements, REU < 25%; for “indicative” PM_{2.5} measurements, REU < 50%). Readers**
 421 **are encouraged to consult the specified standard for further details.**

422 A second example of inter-location performance ~~changing between locations~~ is presented in Fig. 8, showing NO₂ data
 423 from two sensor systems (from two different manufacturers, identified as Systems A and B) ~~(different brands, one~~

424 ~~shown on top of the other~~ before (left plots) and after (right plots) they were moved from Manchester to London in
425 March 2020. Both sensors saw a reduction in agreement with the reference instrument at the London site compared
426 to Manchester, despite both these sites being classified as urban-background with reference instrument performance
427 regularly audited by the UK National Physical Laboratory.



428



429
 430 **Figure 8. Comparison of NO₂ measurements for two systems (A and B) that were moved between Manchester (left plots)**
 431 **and London (right plots). The Manchester deployment was from January–February 2020, and the London data were**
 432 **recorded from April–May 2020.**

433 **Figure 8. Comparative analysis of NO₂ measurements from two systems (A and B), across two urban settings. The left plots**
 434 **display Manchester “out-of-box” data product (January to February 2020), while the right plots show London “cal1” data**
 435 **product (April to May 2020). This “cal1” label does not indicate corrections specific to London’s conditions but denotes a**
 436 **data product from a specific period (as detailed in Figures S2 and S3). The colour gradient represents the density of data**
 437 **points, with darker shades indicating lower densities and brighter shades signifying higher densities.**

438 **The primary distinction between both systems’ behaviour lies in the fact that the sensor located in the top row, even**
 439 **after being relocated to London, maintains a linear response (albeit slightly more degraded than that observed in**
 440 **Manchester, as the R² and RMSE show). In contrast, in the second system (bottom row), the response is notably**
 441 **noisier as the Standard Error (SE) — which is the dispersion of the data around the best line fit line, i.e., the remaining**
 442 **error after bias correction. In scenarios akin to this latter, where there is a high variance in the residuals, a linear**
 443 **correction will not provide a significant improvement. While more sophisticated corrections could be applied, these**
 444 **will be limited by domain knowledge of the end-user, and potentially by other complex data sources that might be**
 445 **available. However, it is important to remember that additional post-processing could increase the risk of overfitting**
 446 **(Aula et al., 2022). On the other hand, for cases like the top plots, users might benefit from trying to correct them**
 447 **using simple linear correction (e.g. using reference instruments if available) or other approaches that could provide**
 448 **means for zero and span correction. A straightforward and cost-effective example could be the use of diffusion tubes**
 449 **for the case of NO₂, as discussed in Section 3.6. The primary distinction between both systems’ behaviour lies in the**
 450 **fact that the sensor located in the top row (Sensor A), even after being relocated to London, maintains a linear response**
 451 **(albeit slightly more degraded than that observed in Manchester, as indicated by the R² and RMSE). In contrast, Sensor**

452 B's response becomes significantly noisier upon relocation to London, as highlighted by the Standard Error (SE) —
453 which represents the remaining error after applying a perfect bias correction. Despite both systems utilising identical
454 sensing elements, the variance in residuals between them may stem from the distinct calibration approaches applied
455 by the respective companies.

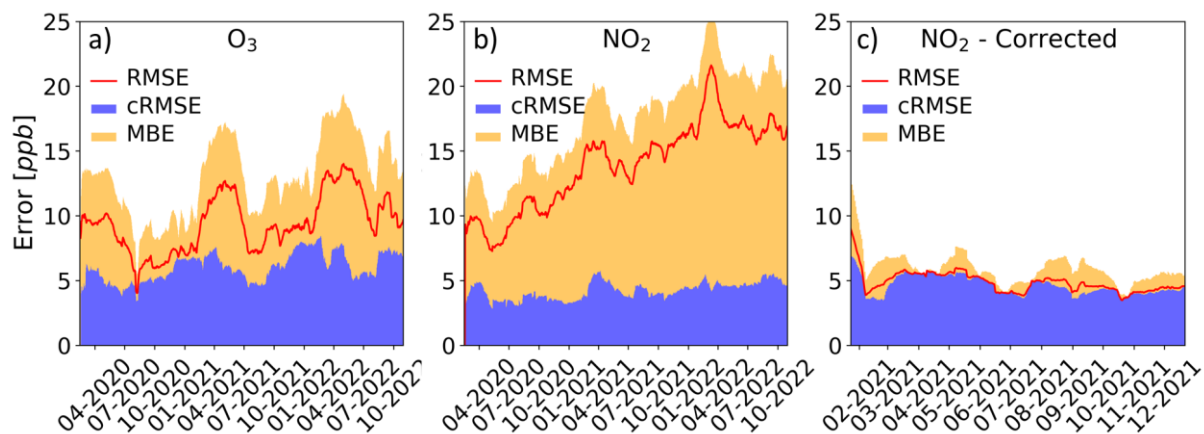
456 For cases resembling Sensor A, users might find it beneficial to implement simple linear correction methods (e.g.,
457 using reference instruments if available) or explore other strategies for zero and span correction. A practical and cost-
458 effective approach, for example, is using diffusion tubes for NO₂ measurements, as discussed in Section 3.6.
459 Conversely, in scenarios characterised by high variance in residuals, such as those observed with Sensor B, a-
460 posteriori attempts to apply a simple linear correction are unlikely to result in significant improvement. While more
461 sophisticated corrections are theoretically feasible, their effectiveness is limited by the end-user's domain knowledge
462 and the availability of additional complex data sources. Furthermore, it is important to consider that excessive post-
463 processing may lead to overfitting—a situation where a model excessively conforms to specific patterns in the training
464 data, resulting in poor performance on new, unseen data (Aula et al., 2022).

465 3.5 Long-term stability

466 The long-term stability of sensor response is also an important facet of its performance, especially for certain use
467 cases such as multi-year network deployments. There can be multiple causes of long-term changes to sensor response,
468 for example, particles settling inside the sampling chamber in optical-based sensors(e.g. Hofman et al. (2022)), or the
469 gradually changing composition of electrochemical cells (e.g. Williams (2020)). How these changes manifest
470 themselves in the data must be identified if ways to account for them are to be implemented.

471 Fig. 9 shows the temporal nature of the O₃ and NO₂ errors (MBE, cRMSE and RMSE) from a sensor system between
472 February 2020 and October 2022. The O₃ shows (Fig. 9a) a gradual increase in the overall measurement error, largely
473 due to an increase in the MBE. It also shows a distinct seasonality MBE, increasing by a factor of 3-4 between March
474 and July compared to the August-February period. The cRMSE component shows fluctuations during the study but
475 only has a small increasing trend. The NO₂ system (Fig. 9b) demonstrates a consistently increasing overall error, with
476 a less pronounced seasonal influence. The bias contributes greatly to the total error (see Section 3.6 for NO₂ sensor
477 correction, Fig. 9c).

478



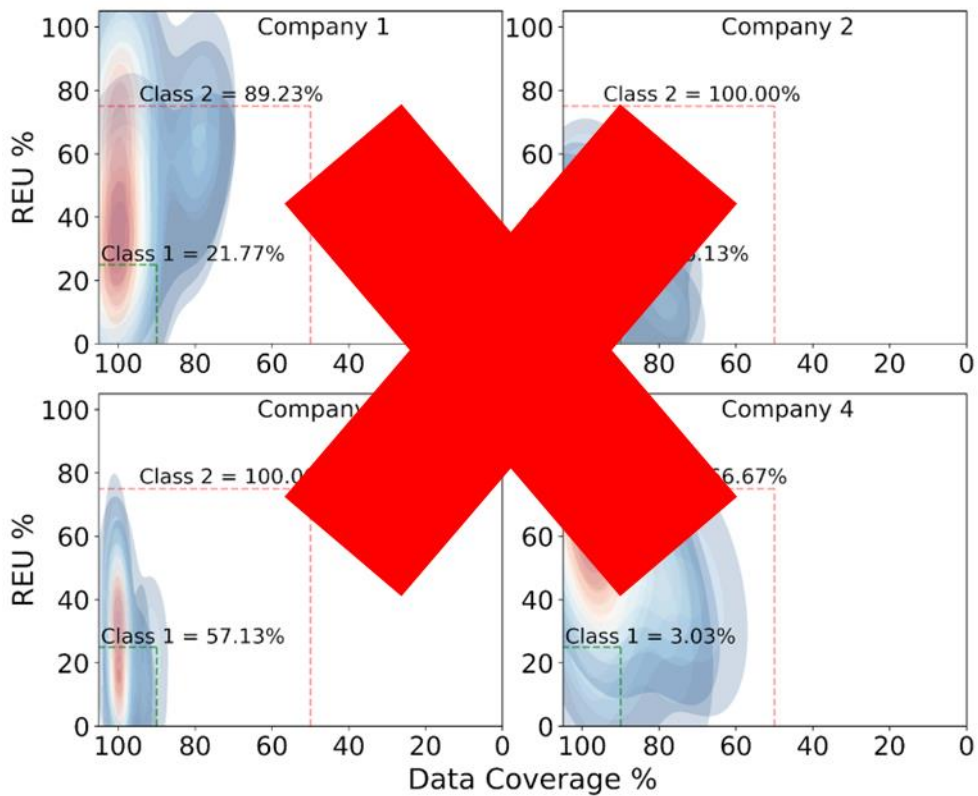
479

480 **Figure 9. Seasonal variation of error (as RMSE, red line) of one of the systems belonging to the Main QUANT, decomposed**
481 **into cRMSE (in blue) and MBE (in yellow) estimated based on a 40-day (aligning with the sample size recommendation by**
482 **the CEN/TS 17660-1:2021 standard for on-field tests) moving window approach with a 1-day slide (i.e., advancing the**
483 **calculation 1 day at a time) (1-day slide) moving window.** Panel a) is for O₃ measurements, and panel b) is for NO₂ (April
484 2020-Oct 2022). Panel c) is also for NO₂, this time showing the effect of a linear correction using diffusion tubes (see next
485 section for more details).

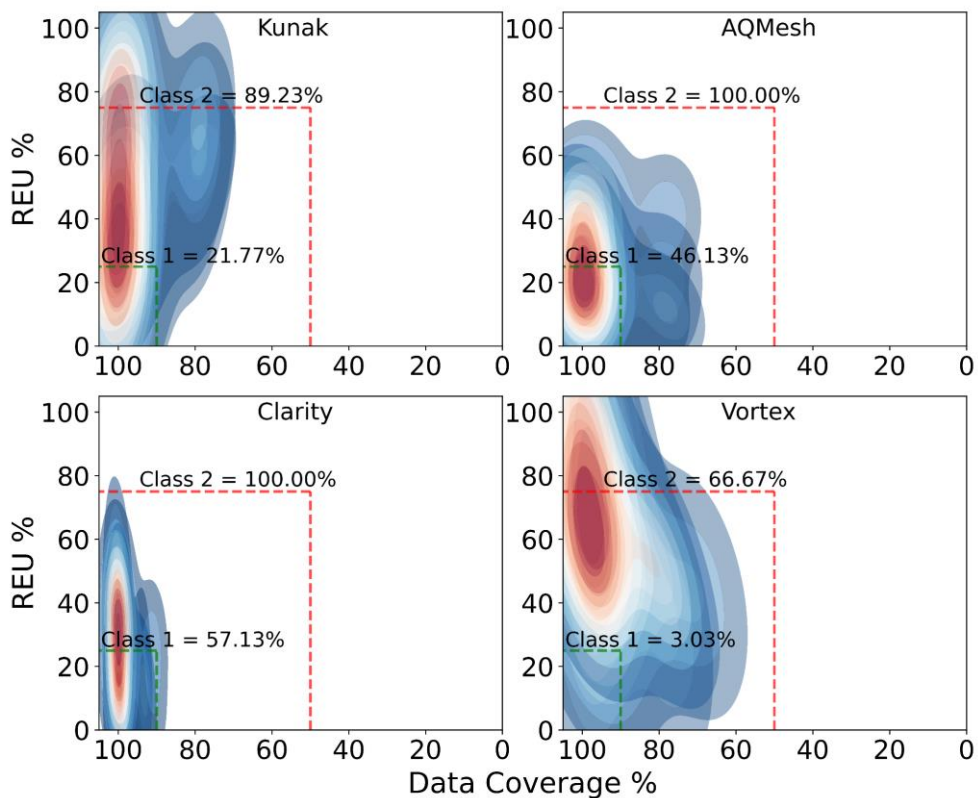
486 3.6 Informing end-use applications

487 Ultimately, for any air pollution monitoring application, the requirements of the task should dictate the measurement
488 technology options available. For example, if the requirement for a particular measurement is to assess legal
489 compliance, then lower measurement uncertainty must be a key consideration as the reported values need to be
490 compared to a limit value. In contrast, if an application aimed to look at long-term trends in pollutants, then absolute
491 accuracy may not be as important as the long-term stability of sensor response. ~~In order to realise the potential of air~~
492 ~~pollution sensor technologies, end users need to be provided with the information required to critically assess the~~
493 ~~strengths and weaknesses of potential candidate sensor devices, ideally in an easy-to-access and interpret manner.~~ To
494 realise the potential of air pollution sensor technologies, end users need to align their specific measurement needs
495 with the capabilities of available devices. Achieving this necessitates access to unbiased performance data, such as
496 long-term stability and accuracy across varying conditions, ideally in an easy-to-access and interpret manner.

497 Understanding the uncertainty associated with a ~~measurement~~ instrument is essential for recognizing its capabilities
498 and limitations. Accurate instruments are crucial, especially in areas like public health decision-making, where
499 inaccurate data can have profound implications (Molina Rueda et al., 2023). Furthermore, instruments that operate
500 autonomously ensure consistent, uninterrupted data collection, making them more efficient and cost-effective in terms
501 of maintenance and calibration. ~~Figure 10 shows the REU (y-axis) and Data Coverage (DC, x-axis) of companies~~
502 ~~measuring NO₂ with more than 2 systems running to avoid ambiguity in the results. Using multiple systems, not only~~
503 ~~avoids ambiguity in results but also enhances the robustness of the data collected.~~ Figure 10 illustrates the collective
504 behaviour of NO₂ sensors from each of the four companies with more than two working systems, showcasing their
505 REU (y-axis) versus Data Coverage (DC, x-axis). Both parameters were calculated for each sensor system using a 40-
506 day moving window approach and then aggregated by brand, ensuring a comprehensive analysis. This methodology
507 leverages overlapping data from multiple sensors to provide a robust representation of company-wide sensor
508 performance and aims to prevent biased interpretations. Both REU and DC are key criteria within the EU scheme
509 (EU 2008/50/EC) for evaluating the performance of measurement methods, and are complemented by the CEN/TS
510 17660-1:2021 specifically for sensors. ~~The latter~~ This document defines three different sensor system tiers. Class 1
511 NO₂ sensors, bounded by the green rectangle (REU < 25% and DC > 90%), offer higher accuracy than Class 2 sensors
512 (REU < 75% and DC > 50%), delimited ~~highlighted~~ by the red rectangle (Class 3 sensors have no set requirements).
513 Presenting the REU and DC ~~data~~ like in Fig. 10 ~~this~~ helps users anticipate the performance of sensor systems —under
514 the assumption that all sensors from the same brand will behave similarly in equivalent environmental conditions—
515 providing more insight into selecting the appropriate instrument for a given project or study.



516



517

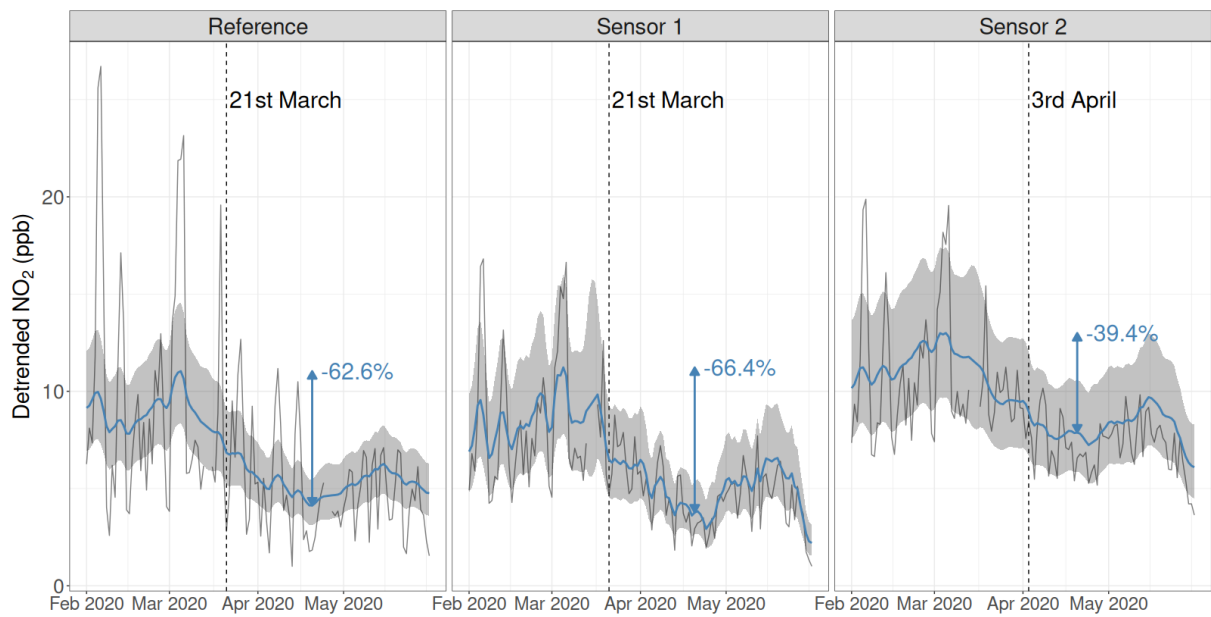
518 **Figure 10. The REU vs. Data Coverage (DC) for 4 systems-companies was evaluated during the WPS for the period Nov**
 519 **2021-Oct 2022 (after all companies had at least one calibrated product). Both the REU and the DC were estimated based**
 520 **on a 40-day size (which is the number of days used by CEN/TS 17660-1:2021 for on-field tests) moving window (1-day slide).**
 521 **While the green rectangle represents the DQOs for Class 1 sensors, the red one limits the DQOs for Class 2 sensors (Class**
 522 **3 sensors have no requirements).**

523 **Figure 10. REU vs. DC for 4 sensor system companies measuring NO₂, with more than two units working simultaneously**
524 **during the WPS (period Nov 2021-Oct 2022, after all companies provided at least one calibrated product). Each heat map**
525 **plot (cooler colours for lower densities and warmer colours for higher densities) aggregates the REU and DC from sensors**
526 **of the same brand working concurrently. The calculation of these two parameters employ a 40-day (aligning with the sample**
527 **size recommendation by the CEN/TS 17660-1:2021 standard for on-field tests) moving window approach with a 1-day slide**
528 **(i.e., advancing the calculation 1 day at a time). The green dashed rectangle limits the Data Quality Objectives (DQOs) for**
529 **Class 1 sensors, and the red dashed rectangle outlines the DQOs for Class 2 sensors.**

530 Depending on the nature of the sensor data uncertainty, methods can be implemented to improve certain aspects of
531 the data quality for a particular application. One such example is the use of distributed networks to estimate sensor
532 measurement errors, such as that described by (J. Kim et al., 2018). ~~Depending on the application, simpler methods~~
533 ~~could also be available to reduce the magnitude of the changing bias, and thus significantly improve the accuracy of~~
534 ~~an individual sensor system, but also that of broader sensor networks. For the case shown in Fig.9b, one possible way~~
535 ~~to do this would be using supporting observations of NO₂ made via diffusion tubes.~~ Depending on the application and
536 available options, users can access alternative methods to reduce bias, thus enhancing the accuracy of sensor systems
537 and networks. For example, “Indicative methods”, as defined by the EU AQ Directive, such as diffusion tubes (e.g.,
538 NO_x, SO₂, VOCs, etc.), can be an option. Specifically, our study leverages diffusion tube data for NO₂, illustrating
539 one effective approach to bias correction using supporting observations, as exemplified in Fig. 9b. These
540 measurements are widely used to monitor NO₂ concentrations in UK urban environments, due to their lower cost (~£5
541 per tube) and ease of deployment, but only provide average concentrations over periods of weeks to months
542 (Butterfield et al., 2021). During QUANT, NO₂ diffusion tubes were deployed at the 3 colocation sites (see Section
543 S7 at the Supp. for more details). Combining these measurements offers the possibility of quantifying the average
544 sensor bias, thus reducing the error on the sensor measurement whilst maintaining the benefits of its high time-
545 resolution observations. It is important to note that while bias correction has been applied to the sensor data, the NO₂
546 diffusion tube concentrations used for comparison purposes must also be adjusted (e.g. following Defra ~~DEFRA~~
547 (2022)). Fig. 9c shows the accuracy of the same NO₂ sensor data shown in Fig. 9b but applies a monthly offset
548 calculated as the difference between its monthly average measurement and that from the diffusion tube (see Figure
549 S85). This shows a dramatic reduction in overall error largely driven by its bias correction. What remains largely
550 resulting from the cRMSE, i.e. the error variance that might arise from limitations from the sensing technology itself
551 and/or the conversion algorithms used to transform the raw signals into the concentration output. To validate the
552 efficacy and reliability of this bias correction method, further long-term studies are warranted.

553 The development and communication of methods that improve sensor data quality, ideally in **accessible digestible**
554 case studies, would likely increase the successful application of sensor devices for local air quality management. There
555 is also a need for similar case studies showcasing the successful application of sensor devices for particular monitoring
556 tasks. An example of this from the QUANT dataset is the use of sensor devices to successfully identify change points
557 in a pollutant’s concentration profile. **These are points in time where the parameters governing the data generation**
558 **process are identified to change, commonly the mean or variance, and can arise from human-made or natural**
559 **phenomena (Aminikhanghahi and Cook, 2017).** Determining when a specific pollutant has changed its temporal
560 nature is a challenging task as there are a large number of confounding factors that influence **atmospheric**
561 **concentrations** ~~a pollutant’s concentration at a specific point in time~~, including but not limited to seasonal factors,
562 environmental conditions (both natural and arising from human behaviour), and meteorological factors. **This challenge**
563 **has lead to several “deweathering” techniques being proposed in the literature (Carslaw et al., 2007; Grange and**

564 Carslaw, 2019; Ropkins et al., 2022). While change point detection is highlighted here as a promising application of
565 sensor data, it represents just one of many potential methodologies that could be explored with the QUANT dataset.



566
567 **Figure 11. NO₂ measurements (black solid line) and detrended estimates (blue solid line with 95% confidence interval in**
568 **the shaded grey region) from the reference instrument (left panel) and 2 sensor systems (middle and right panels) from**
569 **Manchester in 2020. Vertical dashed lines and their corresponding dates indicate identified change points, which**
570 **correspond to the introduction of the first national lockdown due to COVID-19 on the 23rd of March 2020. The percentage**
571 **in blue represents the relative peak-trough decrease from 5th March to 20th April.**

572 ~~A novel statistical approach to smoothing air quality measurements was applied, accounting for these external factors~~
573 ~~(Lacy & Moller). This method was applied to NO₂ concentrations determined from the sensor systems that had~~
574 ~~remained in Manchester throughout 2020, aiming to identify whether the well-documented reduction in ambient NO₂~~
575 ~~concentrations could be observed due to changes in travel patterns associated with COVID-19 restrictions. To provide~~
576 ~~an objective quantification of whether a change point had occurred, the Bayesian online change point detection~~
577 ~~(Adams & MacKay, 2007) was applied. Of the 8 devices that measured NO₂, clear changepoints corresponding to the~~
578 ~~introduction of a lockdown were identified in 2 (Fig.11). While this is an unsupervised analysis, it demonstrates the~~
579 ~~potential of these devices to identify long-term trends with appropriate processing, even with only having had 3~~
580 ~~months of training data to fit the model to. This is especially aided by the given algorithm's ability to use reference~~
581 ~~data as a prior allowing sensor systems to fine-tune the model.~~

582 A state-space based deweathering model was applied to NO₂ concentrations measured from the sensor systems that
583 had remained in Manchester throughout 2020 to remove these confounding factors, with the overarching objective to
584 identify whether the well-documented reduction in ambient NO₂ concentrations due to changes in travel patterns
585 associated with COVID-19 restrictions could be observed in the low-cost sensor systems. To provide a quantifiable
586 measure of whether a meaningful reduction had occurred, the Bayesian online change-point detection (Adams &
587 MacKay, 2007) was applied. Of the 8 devices that measured NO₂, clear change points corresponding to the
588 introduction of a lockdown were identified in 2 (Fig.11), demonstrating the potential of these devices to identify long-
589 term trends with appropriate processing, even with only 3 months of training data.

590 4. Conclusions

591 Lower-cost air pollution sensor technologies have significant potential to improve our understanding and ability to
592 manage air pollution issues. Large-scale uptake in the use of these devices [for air quality management](#) has, [however](#),
593 been primarily limited by concerns over data quality and a general lack of a realistic characterisation of the
594 measurement uncertainties making it difficult to design end uses that make the most of the data information content.
595 ~~Developments in the field of air pollution sensor technology are also developing rapidly, with advances in both the~~
596 ~~measurement technology and particularly in the data post-processing and calibration.~~ [Advances are occurring rapidly,](#)
597 [in both the measurement technology and particularly in the data post-processing and calibration.](#) A challenge with the
598 use of sensor-based devices is that many of the end-use communities do not have access to extensive reference-grade
599 air pollution measurement capability (Lewis & Edwards, 2016), or in many cases, expertise in making atmospheric
600 measurements [or the technical ability for data post-processing](#). For this reason, reliable information on expected sensor
601 performance needs to be available to aid effective end-use applications. Large-scale independent assessments of air
602 sensor technologies are non-trivial and costly, however, making it difficult for end users to find relevant performance
603 information on current sensor technologies. The QUANT assessment is a multi-year study across multiple locations,
604 that aims to provide relevant information on the strengths and weaknesses of commercial air pollution sensors in UK
605 urban environments.

606 The QUANT sensor systems were installed at two highly instrumented urban background measurement sites, in
607 Manchester and London, and one roadside monitoring station in York. The study design ensured that multiple devices
608 were collocated to assess inter-device precision, and devices were also moved between locations and able to test
609 additional calibration data products to assess and enable developments in sensor performance under realistic end-use
610 scenarios. A wider participation component of the Main QUANT assessment was also run at the Manchester site to
611 expand the market representation of devices included in the study, and also to assess recent developments in the field.

612 A high-level analysis of the dataset has highlighted multiple facets of air pollution sensor performance that will help
613 inform their future usage. Inter-device precision has been shown to vary, both between different devices [of the same](#)
614 [brand and model types](#) and over different periods of time, with the most accurate devices generally showing the highest
615 levels of inter-device precision. The accuracy of the reported data for a particular device can be impacted by a variety
616 of factors, from the calibrations applied to its location or seasonality. This has important implications for the way
617 sensor-based technologies are deployed and supports the case made by others (Bittner et al., 2022; Farquhar et al.,
618 2021; Crilley et al., 2018; Williams, 2020; Bi et al., 2020) that practical methods to monitor sensor bias will be crucial
619 in uses where data accuracy is paramount. [Ultimately, this work shows that sensor performance can be highly variable](#)
620 [between different devices and end-users need to be provided with impartial performance data on characteristics such](#)
621 [as accuracy, inter-device precision, long-term drift and calibration transferability in order to decide on the right](#)
622 [measurement tool for their specific application.](#)

623 In addition to these findings, this overview lays the groundwork for more detailed research to be presented in future
624 publications. Subsequent analyses will focus on providing a more nuanced understanding of the uncertainty in air
625 pollution sensor measurements, thus equipping end-users with better insights into the capability of sensor data. Future
626 studies will delve into specific aspects of air pollution sensor performance: 1) a comprehensive performance
627 evaluation of PM_{2.5} data, assessing their accuracy and reliability under different environmental conditions; 2) an in-
628 depth analysis of NO₂ measurements, examining their sensitivity and response in various urban environments; and 3)
629 a detailed investigation into the detection limits of these sensor technologies, targeting their optimised application in

630 low concentration scenarios. These focused studies are basic steps needed to further advance our understanding of
631 sensors' capabilities and limitations, ensuring informed and effective application in air quality monitoring.

632 **Supplementary**

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

634 **Data availability**

635 ~~The data for this study can be found at the Centre for Environmental Data Analysis (CEDA): Lacy et al. (2023):~~
636 ~~Quantification of Utility of Atmospheric Network Technologies: (QUANT): Low cost air quality measurements from~~
637 ~~52 commercial devices at three UK urban monitoring sites. NERC EDS Centre for Environmental Data Analysis, date~~
638 ~~of citation (<https://catalogue.ceda.ac.uk/uuid/ae1df3ef736f4248927984b7aa079d2e>).~~

639 The QUANT dataset, accessible at the Centre for Environmental Data Analysis (CEDA) (Lacy et al., 2023;
640 <https://catalogue.ceda.ac.uk/uuid/ae1df3ef736f4248927984b7aa079d2e>), is the most extensive collection to date
641 assessing air pollution sensors' performance in UK urban settings. It encompasses gas and PM sensor data recorded
642 in the native reporting frequency of each device. The reference data from the three monitoring sites can be found at:

- 643 • MAQS: <https://data.ceda.ac.uk/badc/osca/data/manchester>;
- 644 • LAQS: <https://www.londonair.org.uk/london/asp/datadownload.asp>;
- 645 • YoFi: https://uk-air.defra.gov.uk/data/data_selector.

646 A comprehensive data descriptor manuscript, detailing the QUANT dataset's collection methods, processing
647 protocols, accessibility features, and overall structure—including variables, data reporting frequencies, and QA/QC
648 practices—has been submitted for publication. At the time of this writing, the manuscript is still under review.

649 A GitHub repository at <https://github.com/wacl-york/quant-air-pollution-measurement-errors> provides access to
650 Python and R scripts designed for generating diagnostic visuals and metrics related to the QUANT study, along with
651 sample analyses using the QUANT dataset.

652 **Author contributions**

653 The initial draft of the manuscript was created by SD, PME, and SL. The research was conceptualised, designed, and
654 conducted by PME and SD. Methodological framework and conceptualization were developed by SD, PME, and SL.
655 Data analysis was primarily conducted by SD and SL. The software tools for data visualisation and analysis were
656 developed by SD and enhanced by SL. MF, MP and NM supplied the reference data critical for the study. TB, HC,
657 DH, SG, NAM and JU made substantive revisions to the manuscript, enriching the final submission.

658 **Competing interests**

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

660 **Acknowledgements**

661 This work was funded as part of the UKRI Strategic Priorities Fund Clean Air program (NERC NE/T00195X/1), with
662 support from Defra. We would also like to thank the OSCA team (Integrated Research Observation System for Clean
663 Air, NERC NE/T001984/1, NE/T001917/1) at the MAQS, for their assistance in data collection for the regulatory-

664 grade instruments. The authors wish to acknowledge Dr. Katie Read and the Atmospheric Measurement and
665 Observation Facility (AMOF), a Natural Environment Research Council (UKRI-NERC) funded facility, for providing
666 the Teledyne 200U used in this study and for their expertise on its deployment. Special thanks are due to Dr David
667 Green (Imperial College London) for granting access and sharing the data from LAQS (NERC NE/T001909/1).
668 Special thanks to Chris Anthony, Killian Murphy, Steve Andrews and Jenny Hudson-Bell from WACL for the help
669 and support to the project. Our acknowledgment would be incomplete without mentioning Stuart Murray and Chris
670 Rhodes from the Department of Chemistry Workshop for their technical assistance and advice. Further, we
671 acknowledge Andrew Gillah, Jordan Walters, Liz Bates and Michael Golightly from the City of York Council, who
672 were instrumental in facilitating site access and regularly checking on instrument status. We acknowledge the use of
673 ChatGPT to improve the writing style of this article.

674 **References**

- 675 Adams, R. P. and MacKay, D. J. C.: Bayesian Online Changepoint Detection,
676 <https://doi.org/10.48550/arXiv.0710.3742>, 19 October 2007.
- 677 Alam, M. S., Crilley, L. R., Lee, J. D., Kramer, L. J., Pfrang, C., Vázquez-Moreno, M., Ródenas, M.,
678 Muñoz, A., and Bloss, W. J.: Interference from alkenes in chemiluminescent NO_x measurements,
679 *Atmospheric Meas. Tech.*, 13, 5977–5991, <https://doi.org/10.5194/amt-13-5977-2020>, 2020.
- 680 Allan, J., Harrison, R., and Maggs, R.: Measurement Uncertainty for PM_{2.5} in the Context of the UK
681 National Network, 2022.
- 682 [Aminikhangahi, S. and Cook, D. J.: A survey of methods for time series change point detection, *Knowl.*](#)
683 [*Inf. Syst.*, 51, 339–367, <https://doi.org/10.1007/s10115-016-0987-z>, 2017.](#)
- 684 A. Miech, J., Stanton, L., Gao, M., Micalizzi, P., Uebelherr, J., Herckes, P., and P. Fraser, M.: In situ drift
685 correction for a low-cost NO₂ sensor network, *Environ. Sci. Atmospheres*, 3, 894–904,
686 <https://doi.org/10.1039/D2EA00145D>, 2023.
- 687 Aula, K., Lagerspetz, E., Nurmi, P., and Tarkoma, S.: Evaluation of Low-cost Air Quality Sensor
688 Calibration Models, *ACM Trans. Sens. Netw.*, 18, 72:1-72:32, <https://doi.org/10.1145/3512889>, 2022.
- 689 Baron, R. and Saffell, J.: Amperometric Gas Sensors as a Low Cost Emerging Technology Platform for Air
690 Quality Monitoring Applications: A Review, *ACS Sens.*, 2, 1553–1566,
691 <https://doi.org/10.1021/acssensors.7b00620>, 2017.
- 692 Bi, J., Wildani, A., Chang, H. H., and Liu, Y.: Incorporating Low-Cost Sensor Measurements into High-
693 Resolution PM_{2.5} Modeling at a Large Spatial Scale, *Environ. Sci. Technol.*, 54, 2152–2162,
694 <https://doi.org/10.1021/acs.est.9b06046>, 2020.
- 695 Bigi, A., Mueller, M., Grange, S. K., Ghermandi, G., and Hueglin, C.: Performance of NO, NO₂ low cost
696 sensors and three calibration approaches within a real world application, *Atmospheric Meas. Tech.*, 11,

697 3717–3735, <https://doi.org/10.5194/amt-11-3717-2018>, 2018.

698 Bittner, A. S., Cross, E. S., Hagan, D. H., Malings, C., Lipsky, E., and Grieshop, A. P.: Performance
699 characterization of low-cost air quality sensors for off-grid deployment in rural Malawi, *Atmospheric*
700 *Meas. Tech.*, 15, 3353–3376, <https://doi.org/10.5194/amt-15-3353-2022>, 2022.

701 Brown, R. J. C. and Martin, N. A.: How standardizing ‘low-cost’ air quality monitors will help measure
702 pollution, *Nat. Rev. Phys.*, 5, 139–140, <https://doi.org/10.1038/s42254-023-00561-8>, 2023.

703 Buehler, C., Xiong, F., Zamora, M. L., Skog, K. M., Kohrman-Glaser, J., Colton, S., McNamara, M., Ryan,
704 K., Redlich, C., Bartos, M., Wong, B., Kerkez, B., Koehler, K., and Gentner, D. R.: Stationary and portable
705 multipollutant monitors for high-spatiotemporal-resolution air quality studies including online calibration,
706 *Atmospheric Meas. Tech.*, 14, 995–1013, <https://doi.org/10.5194/amt-14-995-2021>, 2021.

707 Bulot, F. M. J., Johnston, S. J., Basford, P. J., Easton, N. H. C., Apetroaie-Cristea, M., Foster, G. L.,
708 Morris, A. K. R., Cox, S. J., and Loxham, M.: Long-term field comparison of multiple low-cost particulate
709 matter sensors in an outdoor urban environment, *Sci. Rep.*, 9, 7497, [https://doi.org/10.1038/s41598-019-](https://doi.org/10.1038/s41598-019-43716-3)
710 [43716-3](https://doi.org/10.1038/s41598-019-43716-3), 2019.

711 Butterfield, D., Martin, N. A., Coppin, G., and Fryer, D. E.: Equivalence of UK nitrogen dioxide diffusion
712 tube data to the EU reference method, *Atmos. Environ.*, 262, 118614,
713 <https://doi.org/10.1016/j.atmosenv.2021.118614>, 2021.

714 [Carslaw, D. C., Beevers, S. D., and Tate, J. E.: Modelling and assessing trends in traffic-related emissions](#)
715 [using a generalised additive modelling approach, *Atmos. Environ.*, 41, 5289–5299,](#)
716 <https://doi.org/10.1016/j.atmosenv.2007.02.032>, 2007.

717 CEN: CEN/TS 17660-1 Air quality - Performance evaluation of air quality sensor systems - Part 1:
718 Gaseous pollutants in ambient air, 2021.

719 CEN EN 12341: Ambient air - Standard gravimetric measurement method for the determination of the
720 PM10 or PM2,5 mass concentration of suspended particulate matter, 2023.

721 Chojer, H., Branco, P. T. B. S., Martins, F. G., Alvim-Ferraz, M. C. M., and Sousa, S. I. V.: Development
722 of low-cost indoor air quality monitoring devices: Recent advancements, *Sci. Total Environ.*, 727, 138385,
723 <https://doi.org/10.1016/j.scitotenv.2020.138385>, 2020.

724 Crilley, L. R., Shaw, M., Pound, R., Kramer, L. J., Price, R., Young, S., Lewis, A. C., and Pope, F. D.:
725 Evaluation of a low-cost optical particle counter (Alphasense OPC-N2) for ambient air monitoring,
726 *Atmospheric Meas. Tech.*, 11, 709–720, <https://doi.org/10.5194/amt-11-709-2018>, 2018.

727 Cross, E. S., Williams, L. R., Lewis, D. K., Magoon, G. R., Onasch, T. B., Kaminsky, M. L., Worsnop, D.
728 R., and Jayne, J. T.: Use of electrochemical sensors for measurement of air pollution: correcting

729 interference response and validating measurements, *Atmospheric Meas. Tech.*, 10, 3575–3588,
730 <https://doi.org/10.5194/amt-10-3575-2017>, 2017.

731 DEFRA: Technical Guidance (TG22). Local Air Quality Management, 2022.

732 Diez, S., Lacy, S. E., Bannan, T. J., Flynn, M., Gardiner, T., Harrison, D., Marsden, N., Martin, N. A.,
733 Read, K., and Edwards, P. M.: Air pollution measurement errors: is your data fit for purpose?, *Atmospheric*
734 *Meas. Tech.*, 15, 4091–4105, <https://doi.org/10.5194/amt-15-4091-2022>, 2022.

735 [Diez, S., Lacy, S., Read, K., Pete, E., and Josefina, U.: QUANT: A Three-Year, Multi-City Air Quality](#)
736 [Dataset of Commercial Air Sensors and Reference Data for Performance Evaluation,](#)
737 <https://doi.org/10.5281/zenodo.10775692>, 2024.

738 Duvall, R. M., Clements, A. L., Hagler, G., Kamal, A., Kilaru, V., Goodman, L., Frederick, S., Barkjohn,
739 K. K., Greene, D., and Dye, T.: Performance Testing Protocols, Metrics, and Target Values for Fine
740 Particulate Matter Air Sensors, 2021.

741 Farquhar, A. K., Henshaw, G. S., and Williams, D. E.: Understanding and Correcting Unwanted Influences
742 on the Signal from Electrochemical Gas Sensors, *ACS Sens.*, 6, 1295–1304,
743 <https://doi.org/10.1021/acssensors.0c02589>, 2021.

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

747 Feinberg, S., Williams, R., Hagler, G. S. W., Rickard, J., Brown, R., Garver, D., Harshfield, G., Stauffer,
748 P., Mattson, E., Judge, R., and Garvey, S.: Long-term evaluation of air sensor technology under ambient
749 conditions in Denver, Colorado, *Atmospheric Meas. Tech.*, 11, 4605–4615, [https://doi.org/10.5194/amt-11-](https://doi.org/10.5194/amt-11-4605-2018)
750 [4605-2018](https://doi.org/10.5194/amt-11-4605-2018), 2018.

751 Gamboa, V. S., Kinast, É. J., and Pires, M.: System for performance evaluation and calibration of low-cost
752 gas sensors applied to air quality monitoring, *Atmospheric Pollut. Res.*, 14, 101645,
753 <https://doi.org/10.1016/j.apr.2022.101645>, 2023.

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

758 [Grange, S. K. and Carslaw, D. C.: Using meteorological normalisation to detect interventions in air quality](#)
759 [time series, *Sci. Total Environ.*, 653, 578–588, <https://doi.org/10.1016/j.scitotenv.2018.10.344>, 2019.](#)

760 Guimarães, U. S., Narvaes, I. da S., Galo, M. de L. B. T., da Silva, A. de Q., and Camargo, P. de O.:

761 Radargrammetric approaches to the flat relief of the amazon coast using COSMO-SkyMed and TerraSAR-
762 X datasets, *ISPRS J. Photogramm. Remote Sens.*, 145, 284–296,
763 <https://doi.org/10.1016/j.isprsjprs.2018.09.001>, 2018.

764 Hagan, D. H., Gani, S., Bhandari, S., Patel, K., Habib, G., Apte, J. S., Hildebrandt Ruiz, L., and Kroll, J.
765 H.: Inferring Aerosol Sources from Low-Cost Air Quality Sensor Measurements: A Case Study in Delhi,
766 India, *Environ. Sci. Technol. Lett.*, 6, 467–472, <https://doi.org/10.1021/acs.estlett.9b00393>, 2019.

767 Han, J., Liu, X., Jiang, M., Wang, Z., and Xu, M.: A novel light scattering method with size analysis and
768 correction for on-line measurement of particulate matter concentration, *J. Hazard. Mater.*, 401, 123721,
769 <https://doi.org/10.1016/j.jhazmat.2020.123721>, 2021.

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

773 [JCGM: The international vocabulary of metrology—basic and general concepts and associated terms](#)
774 [\(VIM\), 3rd edn. JCGM 200:2012, 2012.](#)

775 Jolliff, J. K., Kindle, J. C., Shulman, I., Penta, B., Friedrichs, M. A. M., Helber, R., and Arnone, R. A.:
776 Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment, *J. Mar. Syst.*, 76, 64–82,
777 <https://doi.org/10.1016/j.jmarsys.2008.05.014>, 2009.

778 Kang, Y., Aye, L., Ngo, T. D., and Zhou, J.: Performance evaluation of low-cost air quality sensors: A
779 review, *Sci. Total Environ.*, 818, 151769, <https://doi.org/10.1016/j.scitotenv.2021.151769>, 2022.

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

783 Kelly, K. E., Whitaker, J., Petty, A., Widmer, C., Dybwad, A., Sleeth, D., Martin, R., and Butterfield, A.:
784 Ambient and laboratory evaluation of a low-cost particulate matter sensor, *Environ. Pollut.*, 221, 491–500,
785 <https://doi.org/10.1016/j.envpol.2016.12.039>, 2017.

786 Kim, H., Müller, M., Henne, S., and Hüglin, C.: Long-term behavior and stability of calibration models for
787 NO and NO₂ low-cost sensors, *Atmospheric Meas. Tech.*, 15, 2979–2992, [https://doi.org/10.5194/amt-15-](https://doi.org/10.5194/amt-15-2979-2022)
788 [2979-2022](https://doi.org/10.5194/amt-15-2979-2022), 2022.

789 Kim, J., Shusterman, A. A., Lieschke, K. J., Newman, C., and Cohen, R. C.: The BErkeley Atmospheric
790 CO₂ Observation Network: field calibration and evaluation of low-cost air quality sensors, *Atmospheric*
791 *Meas. Tech.*, 11, 1937–1946, <https://doi.org/10.5194/amt-11-1937-2018>, 2018.

792 [Lacy, S., Diez, S., and Edwards, P.: Quantification of Utility of Atmospheric Network Technologies:](#)

793 (QUANT): Low-cost air quality measurements from 52 commercial devices at three UK urban monitoring
794 sites., 2023.

795 Levy Zamora, M., Buehler, C., Lei, H., Datta, A., Xiong, F., Gentner, D. R., and Koehler, K.: Evaluating
796 the Performance of Using Low-Cost Sensors to Calibrate for Cross-Sensitivities in a Multipollutant
797 Network, *ACS EST Eng.*, 2, 780–793, <https://doi.org/10.1021/acsestengg.1c00367>, 2022.

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

800 Li, J., Hauryliuk, A., Malings, C., Eilenberg, S. R., Subramanian, R., and Presto, A. A.: Characterizing the
801 Aging of Alphasense NO₂ Sensors in Long-Term Field Deployments, *ACS Sens.*, 6, 2952–2959,
802 <https://doi.org/10.1021/acssensors.1c00729>, 2021.

803 Liang, L.: Calibrating low-cost sensors for ambient air monitoring: Techniques, trends, and challenges,
804 *Environ. Res.*, 197, 111163, <https://doi.org/10.1016/j.envres.2021.111163>, 2021.

805 Liang, L. and Daniels, J.: What Influences Low-cost Sensor Data Calibration? - A Systematic Assessment
806 of Algorithms, Duration, and Predictor Selection, *Aerosol Air Qual. Res.*, 22, 220076,
807 <https://doi.org/10.4209/aaqr.220076>, 2022.

808 Liu, X., Jayaratne, R., Thai, P., Kuhn, T., Zing, I., Christensen, B., Lamont, R., Dunbabin, M., Zhu, S.,
809 Gao, J., Wainwright, D., Neale, D., Kan, R., Kirkwood, J., and Morawska, L.: Low-cost sensors as an
810 alternative for long-term air quality monitoring, *Environ. Res.*, 185, 109438,
811 <https://doi.org/10.1016/j.envres.2020.109438>, 2020.

812 Long, R. W., Whitehill, A., Habel, A., Urbanski, S., Halliday, H., Colón, M., Kaushik, S., and Landis, M.
813 S.: Comparison of ozone measurement methods in biomass burning smoke: an evaluation under field and
814 laboratory conditions, *Atmospheric Meas. Tech.*, 14, 1783–1800, [https://doi.org/10.5194/amt-14-1783-](https://doi.org/10.5194/amt-14-1783-2021)
815 2021, 2021.

816 Malings, C., Tanzer, R., Hauryliuk, A., Saha, P. K., Robinson, A. L., Presto, A. A., and Subramanian, R.:
817 Fine particle mass monitoring with low-cost sensors: Corrections and long-term performance evaluation,
818 *Aerosol Sci. Technol.*, 54, 160–174, <https://doi.org/10.1080/02786826.2019.1623863>, 2020.

819 Molina Rueda, E., Carter, E., L'Orange, C., Quinn, C., and Volckens, J.: Size-Resolved Field Performance
820 of Low-Cost Sensors for Particulate Matter Air Pollution, *Environ. Sci. Technol. Lett.*, 10, 247–253,
821 <https://doi.org/10.1021/acs.estlett.3c00030>, 2023.

822 Moreno-Rangel, A., Sharpe, T., Musau, F., and McGill, G.: Field evaluation of a low-cost indoor air
823 quality monitor to quantify exposure to pollutants in residential environments, *J. Sens. Sens. Syst.*, 7, 373–
824 388, <https://doi.org/10.5194/jsss-7-373-2018>, 2018.

825 Nazemi, H., Joseph, A., Park, J., and Emadi, A.: Advanced Micro- and Nano-Gas Sensor Technology: A
826 Review, *Sensors*, 19, 1285, <https://doi.org/10.3390/s19061285>, 2019.

827 Nowack, P., Konstantinovskiy, L., Gardiner, H., and Cant, J.: Machine learning calibration of low-cost
828 NO₂ and PM₁₀ sensors: non-linear algorithms and their impact on site transferability, *Atmospheric Meas.*
829 *Tech.*, 14, 5637–5655, <https://doi.org/10.5194/amt-14-5637-2021>, 2021.

830 Okure, D., Ssematimba, J., Sserunjogi, R., Gracia, N. L., Soppelsa, M. E., and Bainomugisha, E.:
831 Characterization of Ambient Air Quality in Selected Urban Areas in Uganda Using Low-Cost Sensing and
832 Measurement Technologies, *Environ. Sci. Technol.*, 56, 3324–3339,
833 <https://doi.org/10.1021/acs.est.1c01443>, 2022.

834 Ouyang, B.: First-Principles Algorithm for Air Quality Electrochemical Gas Sensors, *ACS Sens.*, 5, 2742–
835 2746, <https://doi.org/10.1021/acssensors.0c01129>, 2020.

836 Pang, X., Shaw, M. D., Gillot, S., and Lewis, A. C.: The impacts of water vapour and co-pollutants on the
837 performance of electrochemical gas sensors used for air quality monitoring, *Sens. Actuators B Chem.*, 266,
838 674–684, <https://doi.org/10.1016/j.snb.2018.03.144>, 2018.

839 Pang, X., Chen, L., Shi, K., Wu, F., Chen, J., Fang, S., Wang, J., and Xu, M.: A lightweight low-cost and
840 multipollutant sensor package for aerial observations of air pollutants in atmospheric boundary layer, *Sci.*
841 *Total Environ.*, 764, 142828, <https://doi.org/10.1016/j.scitotenv.2020.142828>, 2021.

842 PAS 4023: Selection, deployment, and quality control of low-cost air quality sensor systems in outdoor
843 ambient air – Code of practice, 2023.

844 Pinder, R. W., Klopp, J. M., Kleiman, G., Hagler, G. S. W., Awe, Y., and Terry, S.: Opportunities and
845 challenges for filling the air quality data gap in low- and middle-income countries, *Atmos. Environ.*, 215,
846 116794, <https://doi.org/10.1016/j.atmosenv.2019.06.032>, 2019.

847 Raheja, G., Sabi, K., Sonla, H., Gbedjangni, E. K., McFarlane, C. M., Hodoli, C. G., and Westervelt, D.
848 M.: A Network of Field-Calibrated Low-Cost Sensor Measurements of PM_{2.5} in Lomé, Togo, Over One to
849 Two Years, *ACS Earth Space Chem.*, 6, 1011–1021, <https://doi.org/10.1021/acsearthspacechem.1c00391>,
850 2022.

851 Rai, A. C., Kumar, P., Pilla, F., Skouloudis, A. N., Di Sabatino, S., Ratti, C., Yasar, A., and Rickerby, D.:
852 End-user perspective of low-cost sensors for outdoor air pollution monitoring, *Sci. Total Environ.*, 607–
853 608, 691–705, <https://doi.org/10.1016/j.scitotenv.2017.06.266>, 2017.

854 Ripoll, A., Viana, M., Padrosa, M., Querol, X., Minutolo, A., Hou, K. M., Barcelo-Ordinas, J. M., and
855 Garcia-Vidal, J.: Testing the performance of sensors for ozone pollution monitoring in a citizen science
856 approach, *Sci. Total Environ.*, 651, 1166–1179, <https://doi.org/10.1016/j.scitotenv.2018.09.257>, 2019.

857 Ropkins, K., Walker, A., Philips, I., Rushton, C., Clark, T., and Tate, J.: Change Detection of Air Quality
858 Time-Series Using the R Package Aqeval, <https://doi.org/10.2139/ssrn.4267722>, 4 November 2022.

859 Sayahi, T., Butterfield, A., and Kelly, K. E.: Long-term field evaluation of the Plantower PMS low-cost
860 particulate matter sensors, *Environ. Pollut.*, 245, 932–940, <https://doi.org/10.1016/j.envpol.2018.11.065>,
861 2019.

862 Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., and Bonavitacola, F.: Field calibration of a
863 cluster of low-cost commercially available sensors for air quality monitoring. Part B: NO, CO and CO₂,
864 *Sens. Actuators B Chem.*, 238, 706–715, <https://doi.org/10.1016/j.snb.2016.07.036>, 2017.

865 Tanzer-Gruener, R., Li, J., Eilenberg, S. R., Robinson, A. L., and Presto, A. A.: Impacts of Modifiable
866 Factors on Ambient Air Pollution: A Case Study of COVID-19 Shutdowns, *Environ. Sci. Technol. Lett.*, 7,
867 554–559, <https://doi.org/10.1021/acs.estlett.0c00365>, 2020.

868 Wang, Y., Li, J., Jing, H., Zhang, Q., Jiang, J., and Biswas, P.: Laboratory Evaluation and Calibration of
869 Three Low-Cost Particle Sensors for Particulate Matter Measurement, *Aerosol Sci. Technol.*, 49, 1063–
870 1077, <https://doi.org/10.1080/02786826.2015.1100710>, 2015.

871 Williams, D. E.: Electrochemical sensors for environmental gas analysis, *Curr. Opin. Electrochem.*, 22,
872 145–153, <https://doi.org/10.1016/j.coelec.2020.06.006>, 2020.

873 Wu, T. Y., Horender, S., Tancev, G., and Vasilatou, K.: Evaluation of aerosol-spectrometer based PM_{2.5}
874 and PM₁₀ mass concentration measurement using ambient-like model aerosols in the laboratory,
875 *Measurement*, 201, 111761, <https://doi.org/10.1016/j.measurement.2022.111761>, 2022.

1 **S1. Co-location sites**

2 For the main QUANT deployment, 3 field sites were chosen: Manchester, London, and York, all providing
3 extensive reference measurements across a range of chemical environments representative of UK urban
4 atmospheres. On the other hand, only the Manchester site was used for the WPS colocation.

5 ~~The Manchester Air Quality Supersite (MAQS, 53° 26' 39.2"N, 2° 12' 51.9"W) is one of the largest air quality
6 research facilities in the UK, and also because it is located in the south of the city of Manchester (the second
7 largest metropolitan area in the UK, with approx. 3.3 million inh.) in an urban background environment (avg.
8 temp. in winter of about 4-5 °C and RH ~87 %, avg. temp. in summer around 16-17 °C and RH ~88 %. MAQS
9 reference instrumentation details can be found in the Section S4. All the data provided by MAQS was 1-min time
10 resolution.~~

11 The Manchester Air Quality Supersite (MAQS, 53° 26' 39.2"N, 2° 12' 51.9"W) stands as one of the largest air
12 quality research facilities in the UK. Situated in an urban background setting approximately four kilometres south
13 of Manchester city center — the UK's second-largest metropolitan area with around 3.3 million residents —
14 MAQS benefits from a strategic location on the University of Manchester's Fallowfield Campus. This location is
15 notably distanced from direct traffic emissions, surrounded by student accommodations, university administrative
16 buildings, and sports facilities. The campus's vicinity to shops, bars, and restaurants introduces a range of human
17 activities, including varying levels of foot traffic and associated vehicular movement. Additionally, the presence
18 of these commercial and recreational spaces, alongside residential buildings, contributes to the area's ambient air
19 quality through emissions from heating and cooking, among other sources. For a visual representation of MAQS's
20 surroundings, please refer to Figure S1 (panel a). The site experiences an average winter temperature of
21 approximately 4-5°C with relative humidity around 87%, and an average summer temperature of about 16-17°C
22 with relative humidity near 88%. Detailed information on MAQS's reference instrumentation and the
23 methodologies employed for air quality measurements can be found in section S2. Data from MAQS are provided
24 with a 1-minute time resolution, facilitating a granular temporal analysis of air quality metrics.

25 ~~The London Air Quality Supersite (LAQS) is an urban background monitoring site located at Honor Oak Park
26 (51° 26' 58.9"N 0° 02' 14.6"W) in Greater London, the third biggest European urban conglomeration with approx.
27 14.8 million inh. (avg. temp. in winter ~5 °C and RH ~84 %, avg. temp. in summer ~17 °C and RH ~72 %). All
28 gas data provided by LAQS is 1-min time resolution and 15-min for PM.~~

29 The London Air Quality Supersite (LAQS, 51° 26' 58.9"N 0° 02' 14.6"W) serves as an urban background
30 monitoring site, nestled within Honor Oak Park in Greater London. Situated 9 km southeast of the city center of
31 the third-largest European urban conglomeration, LAQS offers a unique window into the air quality challenges of
32 an area inhabited by approximately 14.8 million people. Nestled within the serene King's College sports grounds,
33 is surrounded by middle-class neighbourhoods, abundant parks, and green spaces. This tranquil setting, is
34 distanced from major roads and pollution sources, provides a representative snapshot of the ambient air quality
35 typical of residential London. LAQS's surroundings are marked by a low level of commercial activity, with local
36 shops and restaurants contributing minimally to the area's overall noise and bustle. Figure S1 (panel b) offers an
37 aerial view of LAQS, illustrating the overall urban layout. The area is characterised by a temperate climate,

38 experiencing average winter temperatures of around 5°C with RH of approx. 84%, and milder summers with
39 temperatures averaging 17°C and RH of around 72%. Gas measurements at LAQS are conducted with a 1-minute
40 time resolution, while PM data are collected at a 15-minute resolution (see section S2 for more details).

~~41 The York Fishergate roadside site (YoFi), located in the city of York (~210,000 inh., avg. temp. in winter of ~4°C
42 and RH ~87 %, avg. temp. in summer around 15 °C and RH ~80 %). This site is a self-contained air quality
43 monitoring station located very close to the city centre on a traffic island (53° 57' 06.9"N 1° 04' 33.1"W)
44 surrounded by a residential/commercial area. This site was chosen to evaluate the LCS responses to a greater
45 pollutant variability typical of traffic-related sites (in contrast with urban background monitoring stations as in the
46 case of MAQS and LAQS). While PM and NO_x data from YoFi are 1-hr time resolution, the O₃ data is 1-min
47 (deployed on the 15th of May 2020, specifically as part of the QUANT study).~~

48 The York Fishergate roadside site (YoFi, 53° 57' 06.9"N, 1° 04' 33.1"W), in the historic city of York, which is
49 home to approximately 210,000 inhabitants (avg. temp. in winter of ~4°C and RH ~87 %, avg. temp. in summer
50 around 15 °C and RH ~80 %). Situated just about 1 km from the city center on a traffic island, YoFi stands amidst
51 a predominantly residential area that also encompasses commercial and light industrial elements. Unique to its
52 location, the site is sandwiched between two lanes of Fishergate Road, a major avenue that bifurcates to facilitate
53 traffic flow into and out of the city's southern part. Directly across from YoFi, a primary school adds to the daily
54 human activity around the site, while the nearby River Ouse, located merely 300 metres to the west, contributes
55 to the area's environmental characteristics. A vibrant commercial zone, featuring pubs and restaurants, is found
56 just 100 metres to the north. Moreover, the site is flanked by Walmgate Stray, an expanse of recreational fields,
57 located about 300 metres to the southeast, offering a green respite amidst the urban setting. Additional details can
58 be visualised in Figure S1 (panel c), providing an aerial perspective of the site's key features and its urban context.
59 This self-contained air quality monitoring station was specifically selected for the QUANT study to assess sensors'
60 responses to the greater pollutant variability typical of traffic-related sites, contrasting with the urban background
61 settings of MAQS and LAQS. YoFi provides data on PM and NO_x with a 1-hour time resolution. Additionally,
62 in a targeted effort to enhance our understanding of air quality dynamics, O₃ measurements (deployed on the 15th
63 of May 2020, specifically as part of the QUANT study), utilising a 1-minute time resolution to offer detailed
64 insights into temporal variations (refer to section S2 for more details).



66 **Figure S1: Aerial views of the air quality monitoring sites: a) MAQS, b) LAQS, and c) YoFi, captured from Google**
 67 **Earth. These images illustrate the diverse urban settings of each site, emphasising aspects such as their proximity to**
 68 **traffic sources, presence of green spaces, and the general urban layout. Image credits: Google Earth.**

69 **S2. Reference instrumentation, QA/QC, and data-sharing periods**

70 Table S1 summarises the reference instrumentation at each site, Table S2 describes some of the QA/QC processes
 71 at the supersites, and Table S3 shows the data periods shared with the suppliers.

72 **Table S1. Research grade instrumentation used for the QUANT study.**

Analyte	Manchester	London	York
NO	Thermo 42i-y (Chem)	Teledyne T200U (Chem)	Teledyne T200UP (Chem)
NO ₂	*Teledyne T500U (CAPS)	*Teledyne T500U (CAPS)	
O ₃	*Thermo 49i (UV)	*Teledyne 400E (UV)	*2B 205 (UV)
PM	*Palas FIDAS200 (OAS)	*Palas FIDAS200 (OAS)	*Met One BAM 1020 (BA)

73 *Equivalent to reference (as defined in the [European Air Quality Directive 2008/50/EC](#))

74 Acronyms: Chem: Chemiluminescence; CAPS: Cavity Attenuated Phase Shift Spectroscopy; UV: Ultraviolet; OAS:
 75 Optical aerosol spectrometer; BA: Beta attenuation.

76 **Table S2. Summary of Quality Assurance processes in MAQS and LAQS**

Instrument	Frequency	*Process
NO _y	At least monthly	Zero and span checks using standard cylinder and scrubber. Corrections to zero and span values.
NO ₂	Daily	Automatic zero and span checks using internal NO ₂ diffusion tube and scrubber. Zero corrections, span monitored.
O ₃	Daily	Automatic zero and span checks using internal O ₃ lamp and scrubber. Corrections to zero, span monitored.
CO	Every three hours & monthly	Zero checks every three hours and span checks monthly using onsite cylinder. Adjustments to zero and span values.
CO ₂ and CH ₄	Regular	Stability checks using onsite cylinder, no corrections made.
*PM	Semiannual	Sizing response verified with Mono dust, flow rate checked with Gilibrator.

77 *Checked with external standards by NPL every 6 months. These external standards are also used to provide a certification of the on-site
 78 standard cylinders. Final corrections to the data are provided by using the audit data to define the concentration of the on-site standards, with
 79 zero and span values interpolated between the calibration points.

80 **Sizing and flow checked every 6-month NPL audit process.1

81 **Table S3. Reference data is shared with the sensor manufacturers.**

QUANT main study			Wider Participation Study		
Reference dataset	Period	Released	Reference dataset	Period	Released
1	10-12-2019 - 17-02-2020	15-04-2020	1	17-06-2021 - 16-07-2021	23-07-2021
2	18-02-2020 - 17-08-2020	27-10-2020	2	01-12-2021 - 31-12-2021	26-01-2022
3	18-08-2020 - 17-02-2021	15-04-2021	3	01-05-2022 - 31-05-2022	15-06-2022

82

83 S3. QUANT main study devices

84 In this section, a brief description of the QUANT main study systems' components is offered.

85 PurpleAir (PA) (<https://www2.purpleair.com>) devices (PA-II-SD model, firmware v4.11) reports particulate
86 matter (PM₁, PM_{2.5}, and PM₁₀), and it was chosen for its penetration around the world. Two identical Plantower
87 PMS5003 (Plantower) sensors (channels A and B) are found in each PA. It offers two data products (2-min avg.
88 time): the “cf_atm” (for outdoor applications) and the “cf_1” (for indoor or controlled environment applications).
89 The PMS behaves like a nephelometer rather than an optical particle counter to measure the light scattered by the
90 PM (Ouimette et al., 2022) and is composed of a laser, a photodiode, a fan, and a microprocessor control unit.
91 They also measure temperature (Temp), relative humidity (RH), and atmospheric pressure (Pres) (Bosch). The
92 data can be communicated via Wi-Fi or stored locally (microSD card), which was the preferred way during the
93 colocation. No calibrated products are offered by the company.

94 **Note: For this study, only Channel A and the data product “cf_atm” were included in the analysis and shown in*
95 *the plots.*

96 AQMesh (<https://www.aqmesh.com>) reports NO₂, NO, O₃ using electrochemical (EC) sensors (Alphasense), CO₂
97 with a non-dispersive infrared sensor (NDIR, Alphasense), PM₁, PM_{2.5}, and PM₁₀ through a light-scattering sensor
98 (Nephelometer, Environmental Instr.) with 1-minute time resolution (algorithm v5.1 for gases and v3.0 for PM).
99 This instrument also registers Temp, RH, and Pres (Solid-State sensors) (Zauli-Sajani et al., 2022) and the
100 sampling mechanism employs a pump. The collected data is sent to the company server via a cellular network and
101 post-processed (Temp, RH, and cross-interference correction) in the cloud by a proprietary algorithm. Finally, the
102 data is released to the final user via secure web login or through its Application Programming Interface (API).
103 Although the first 4 months of the deployment the data had a 15-min resolution, since then the provided resolution
104 is 1-min average.

105 AQY (v.1.0) is also a multi-species device (<https://www.aeroqual.com>) and measures O₃, NO₂, PM_{2.5}, PM₁₀,
106 Temp, and RH. This is the only device system that does not use Alphasense sensors for gases. While O₃ is
107 quantified using a metal oxide sensor (WO₃-based, Aeroqual Ltd), the NO₂ is measured by an EC sensor
108 (Membrapore type O₃/M5, Aeroqual Ltd) (Weissert et al., 2019). For PM it uses a light scattering method (Nova)

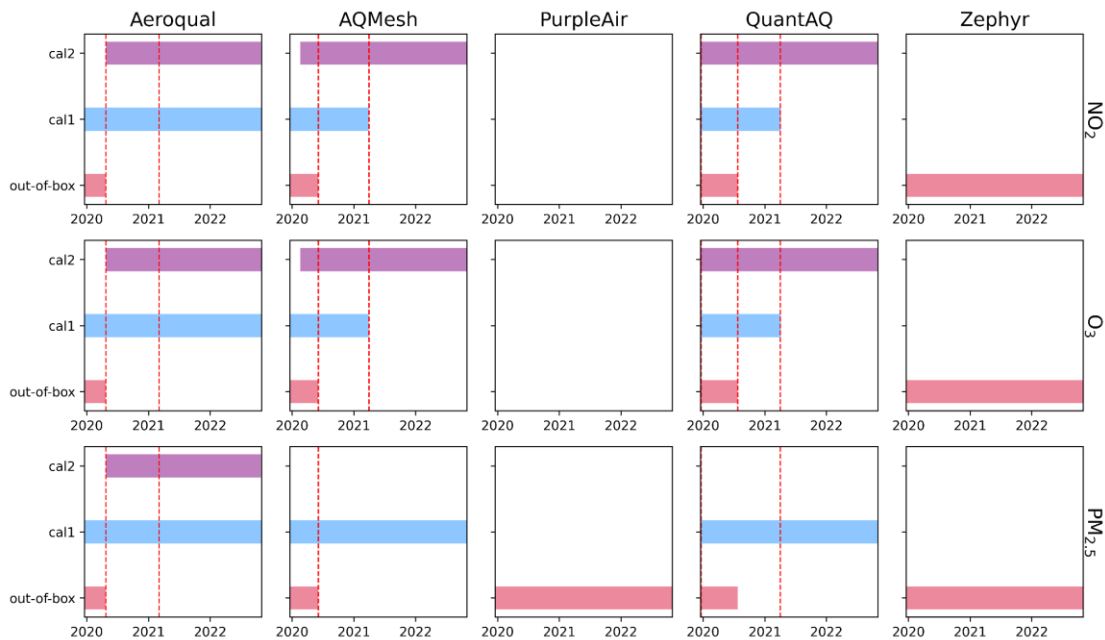
109 to convert size and particle count to a mass fraction and behaves like a nephelometer (Myklebust et al., 2022).
 110 These LCS devices send their data (1-min time resolution) to the Aeroqual server via cellular (WiFi could also be
 111 used for this purpose) or stored locally (microSD card). The non-local data access is through a web portal or via
 112 API.

113 Zephyr units (<https://www.earthsense.co.uk>) measure PM (Nephelometer, Plantower), Temp & RH (Sensirion),
 114 and Press (Bosch) (the sample uptake uses a fan). As most of the commercial units tested here, it used Alphasense
 115 EC sensors (the “A series”, a smaller version than the B series) for gases (NO, NO₂, and O₃). These devices send
 116 their raw data to the server via a cellular network, where they pre-process the raw signals. We have secure access
 117 to the measurements with a time resolution of 1-min per species through the website or via its API.

118 ARIsense v200 devices (<https://quant-aq.com>) measure NO, NO₂, O₃, CO (EC, Alphasense), CO₂ (NDIR,
 119 Alphasense), Temp & RH (Sensirion), and Press (Bosch) (Cross et al., 2017). Of all the devices tested, this is the
 120 only one that uses an Optical Particle Counter (OPC) for PM (Particles Plus). Communication is carried out
 121 through a cellular network and the data products are accessed through a web portal or API (1-minute time
 122 resolution). According to the company policy, only the gas data products are subjected to calibrations (if
 123 colocation data is available).

124 **Table S4. Summary of sensor measurements and the time resolution data provided by participating companies in the Main**
 125 **QUANT study.**

System	Measurands	Time Resol.
PA	PM ₁ , PM _{2.5} , PM ₁₀	2min
AQM	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃ , CO ₂	1min/15min
AQY	PM _{2.5} , PM ₁₀ , NO ₂ , O ₃	1min
Zep	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃	1min
Ari	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃ , CO; CO ₂	1min



126

127 **Figure S2.** Data product for each of the participating companies during Main QUANT. The top panels are for
 128 **NO₂**, the middle panels for **O₃** and the bottom panels for **PM_{2.5}**. The y-axis represents the different products: “out-
 129 of-box”, cal1 and cal2. The x-axis shows the dates for which each company provided the mentioned products.

130 **S4. WPS devices**

131 A short description of the WPS devices’ components is shown in this section

132 Modulair-PM instruments (<https://quant-aq.com>) employ two different techniques to obtain PM mass
 133 concentration (it samples the air using a fan), an OPC (Alphasense, OPC-N3) and a nephelometer (Plantower,
 134 PMS5003). This system provides 1-min time resolution data for PM₁, PM_{2.5}, and PM₁₀, plus size-resolved particle
 135 number concentration (range 350 nm to 40 μm) (Meyer et al., 2022; Westgate and Ng, 2022). Temp, RH, and
 136 Press are also measured, but no data was found about the sensing elements it uses. The post-processed data can
 137 be accessed locally (microSD card) or through its server (cellular network comm) via its web portal or API.

138 AQMesh (see earlier description).

139 The Atmos device (<http://urbansciences.in/>) reports PM₁, PM_{2.5}, PM₁₀ (Plantower, PMS7003) plus Temp and RH
 140 (Adafruit), employing a fan as a means to sample the air. The system transmits the data (1-min time resolution)
 141 to a cloud server (only via Wi-Fi) and also stores it locally (Puttaswamy et al., 2022). The data can be accessed
 142 via a web dashboard or API. Unfortunately, and due to the meteorological conditions at the Manchester supersite
 143 these co-located devices only survived for about 2 months.

144 The IMB instrument (<https://www.bosch-mobility-solutions.com>) measures NO₂, O₃ PM_{2.5} and PM₁₀,
 145 (Alphasense sensors), plus Press, RH an Temp (no details were found about the brand and model). The raw data
 146 is transmitted to their cloud using cellular connectivity (3G or LTE). The final data is 1-min resolution (accessed
 147 only via API).

148 Polludrone (<https://oizom.com>) uses Alphasense sensors for gas measurements (B4 series for NO, NO₂, O₃. No
 149 data available about CO, CO₂ and SO₂) and a Wuhan Cubic PM3006S for PM (PM_{2.5} and PM₁₀) (Oizom -
 150 Polludrone Smart, 2023). It also registers RH and Temp, but no data was found in regards to sensor model/brand.
 151 The sampling mechanism uses a fan and data transmission is wireless. The final product (time res is 10-min) can
 152 be obtained through the Oizom webpage and/or via API.

153 Kunak Air Pro (<https://www.kunak.es/>) uses a fan for sampling and all sensors are from Alphasense (EC, B series
 154 for CO, NO, NO₂ and O₃; an NDIR sensor for CO₂; and an OPC-N3 for PM₁, PM_{2.5}, and PM₁₀) (Hofman et al.,
 155 2022). It also provides Temp, RH, and Press (no data was found in regards to environmental sensor model/brand).
 156 The raw data is transmitted via a multi-band network, and the final data (time res is 5-min) can be accessed through
 157 their website or via API.

158 The Silax Air (<https://vortexiot.com>) system measures NO₂, O₃, PM₁₀ and PM_{2.5}. Their webpage mentions that for
 159 PM an optical scattering sensor is used and EC sensors for the gases. Further details weren't found. The raw data
 160 is transmitted via 4G or WiFi and the final user accesses the final product (5-min time res) through API or website.

161 The Node-S system (<https://www.clarity.io>) holds a nephelometer (Plantower PMS6003) to measure 3 PM size
 162 cuts (PM₁, PM_{2.5}, PM₁₀) (Liu et al., 2022) and EC sensors for NO₂ (Alphasense) (Miech et al., 2021). The air is
 163 dragged into the system by a fan and a Bosch sensor is used for press, RH, and temp. The data is communicated
 164 to Clarity's cloud via cellular signal (4G) and the final product is ~3-min time res (something unusual for sensor
 165 systems). Access to the final data is via the web portal or through API.

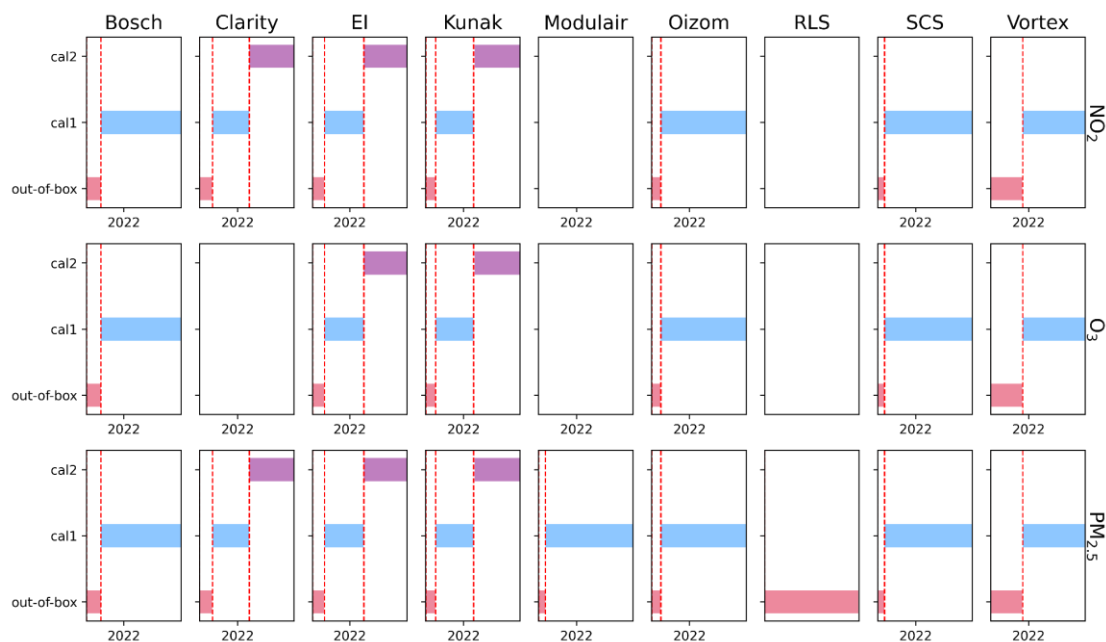
166 Praxis/Urban (<https://www.southcoastscience.com>) system employs EC sensors for NO, NO₂, O₃ (Alphasense, A
 167 series), an NDIR for CO₂ (Alphasense), and particle counter (Alphasense, OPC-N3) for PM₁, PM₁₀ and PM_{2.5}.
 168 The Temp/RH is Sensirion and the Press sensor is TDK. The raw data is communicated to the company server
 169 using 4G and the user can access it and post-processed data through an API (1-min time res).

170 **Table S5. Summary of sensor measurements and the time resolution data provided by participating companies in the WPS**
 171 **study.**

System	Measurands	Time Resol.
Mod	PM ₁ , PM _{2.5} , PM ₁₀	1min
AQM	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃ , CO; CO ₂	15min
Atm	PM ₁ , PM _{2.5} , PM ₁₀	2min
IMB	PM ₁ , PM _{2.5} , PM ₁₀ , NO ₂ , O ₃	1min
Poll	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃	10min
AP	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃ , CO; CO ₂	5min

SA	PM ₁ , PM _{2.5} , PM ₁₀ , NO ₂ , O ₃	5min
NS	PM ₁ , PM _{2.5} , PM ₁₀ , NO ₂	~5min
Prax	PM ₁ , PM _{2.5} , PM ₁₀ , NO, NO ₂ , O ₃ , CO; CO ₂	1min

172



173

174 **Figure S32.** Data product for each of the participating companies in the WPS. The top panels are for NO₂, the
 175 middle panels for O₃ and the bottom panels for PM_{2.5}. The y-axis represents the different products: “out-of-box”,
 176 cal1 and cal2. The x-axis shows the dates for which each company provided the mentioned products.

177 **S5. Performance Metrics**

178 In the assessment of sensor measurement error, it is standard practice to employ a linear additive model, described
 179 by the following equation:

180
$$y_i = b_1 x_i + b_0 + \varepsilon_i \tag{1}$$

181 In this model, the dependent variable “y” represents the sensor measurements, while the independent variable “x”
 182 denotes the reference measurements. The coefficient b₁ corresponds to the slope of the regression line (the
 183 response sensitivity of the sensor relative to the reference) and b₀ is the ordinate at the origin (the sensor's output
 184 when the reference measurement is zero). ε_i, assumed to have a mean of zero and a standard deviation of σ_ε,
 185 captures the portion of “y” that cannot be explained by “x”. For a sensor to perfectly match the reference
 186 measurements (i.e., y = x), b₁ would equal one, with both b₀ and ε_i being zero.

187 **Coefficient of Determination (R²)**

188 R^2 is an adimensional metric that quantifies the proportion of variance in the sensor measurements (“y”) that can
189 be explained by its linear relationship with the reference measurements (“x”):

$$190 \quad R^2 = \frac{\sum_{i=1}^n (x_i - \hat{y})^2}{\sum_{i=1}^n (y_i - \hat{y})^2} \quad (2)$$

191 As a bounded metric, R^2 varies between zero and one ($0 \leq R^2 \leq 1$), where a value closer
192 to one indicates a stronger linear association between the sensor and reference
193 data. Despite being one of the most widely used metrics in sensor evaluation, as
194 highlighted by Karagulian et al. (2019), R^2 comes with limitations that warrant careful consideration.
195 Notably, R^2 does not account for bias in the data; a regression line diverging from the ideal 1:1 relationship
196 between “x” and “y” does not affect its value. Additionally, R^2 is influenced by the dynamic range of the
197 measurements, which can skew its interpretation. Given these nuances, it is prudent to report R^2 alongside
198 complementary metrics that can offer a more rounded view of sensor performance. For a more in-depth analysis
199 of the limitations and proper use of R^2 , readers are directed to the discussion in Legates and McCabe Jr. (1999).

200 *Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE)*

201 MAE and RMSE (both dimensional metrics, expressed in the same units as the measured variable), also stand as
202 very popular metrics for performance evaluation, as they offer insights into the accuracy of sensors, presenting a
203 fuller picture than the R^2 alone. These metrics can be estimated as follows:

$$204 \quad MAE = \frac{1}{n} \sum_{i=1}^n |y_i - x_i| \quad (3)$$

$$205 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - x_i)^2} \quad (4)$$

206
207 However, both MAE and RMSE quantify average errors. MAE does so by calculating the average magnitude of
208 errors without directionality, utilising absolute differences, while RMSE gauges the standard deviation of these
209 differences, highlighting the squared differences between sensor readings and reference grade measurements.
210 Although MAE and RMSE are both valued for their measure of accuracy, they bear distinct implications in
211 practice. MAE treats all errors equally, allocating proportional weight across the board. Conversely, RMSE
212 disproportionately penalises larger errors due to its squaring of difference values, an aspect noted by (Willmott
213 and Matsuura, 2005). This characteristic makes RMSE particularly sensitive to outliers, shaping its utility in
214 identifying and rectifying significant deviations.

215 *Mean Bias Error (MBE)*

216 The MBE quantifies the average bias in sensor measurements relative to reference values. Expressed in the same
217 units as the variable being measured, MBE reflects the systematic error, offering a straightforward indication of a
218 sensor's tendency to overestimate or underestimate the reference:

$$219 \quad MBE = \frac{1}{n} \sum_{i=1}^n (y_i - x_i) \quad (5)$$

220 A zero value of MBE indicates no consistent over- or underestimation, while positive or negative values signal
 221 systematic bias in measurement. This simplicity in interpretation makes MBE particularly valuable for initial
 222 assessments of sensor accuracy and for guiding calibration efforts to correct for systematic bias. However, the
 223 MBE does not capture the precision of the measurements. For this reason, MBE is most effective when used in
 224 conjunction with other metrics, such as RMSE and MAE, to gain a comprehensive understanding of sensor
 225 performance, encompassing both systematic and random errors.

226 *Relative Expanded Uncertainty (REU)*

227 In contrast to single-value metrics such as R², RMSE, and MAE, which assess data sets as a whole, REU offers a
 228 “point by point” metric. This allows for graphical representations (like the REU in the concentration space or as
 229 a time series), offering detailed insights into measurement performance variability. The REU’s mathematical
 230 framework is outlined in the “Guidance for the Demonstration of Equivalence of Ambient Air Monitoring
 231 Methods” (European Commission, 2010), as follows:

$$232 \quad U(y_i) = \sqrt{\frac{RSS}{n-2} + u^2(x_i) + (y_i - b_0 - b_1 x_i)^2} \quad (6)$$

$$233 \quad REU(y_i) = \frac{k \cdot U(y_i)}{\hat{x}} \quad (7)$$

$$234 \quad RSS = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 \quad (8)$$

235 here, U(y_i) represents the measurement uncertainty [concentration units]; REU(y_i) denotes the REU [percentage];
 236 u(x_i) is the random uncertainty of the reference monitor [concentration units]; “n” stand for the number of
 237 collocated data points considered; RSS is the Residual Sum of Squares; k is the coverage factor (set at 2 for a 95%
 238 confidence level).

239 A distinctive feature of REU is its incorporation of the uncertainty associated with the reference method (i.e.,
 240 u(x_i)). This aspect recognizes that all measurements, including those from reference methods, are subject to
 241 inherent uncertainties. While calculating REU is more complex than traditional metrics, it's essential to
 242 acknowledge that, like any metric, REU is based on specific assumptions and considerations. These factors must
 243 be thoughtfully evaluated when interpreting data to ensure that conclusions are firmly rooted in the context of the
 244 study.

245 *Current guidance and normalisation efforts*

246 Table S6 summarises the key metrics addressed in some of the most recent guidance documents and technical
 247 standards. These metrics have been categorised under various labels: linearity, bias, error, uncertainty, data
 248 coverage, and inter-sensor precision. Each of these guidelines and regulations has its own set of procedures,
 249 protocols, and thresholds. Therefore, it is advisable for readers to consult the original documents for a detailed
 250 understanding of these specificities.

251 **Table S6. Summary of field evaluation metrics for sensors according to different guidelines and technical standards.**

Feature	EPA ^{1&2}	CEN ³	ASTM ^{4&5}
---------	------------------------	------------------	-------------------------

<i>Pollutants covered</i>	PM _{2.5} & O ₃	NO ₂ , O ₃ , CO, SO ₂ & Bencene	PM _{2.5} , PM ₁₀ NO ₂ , O ₃ , CO & SO ₂
<i>Linearity</i>	R ²	----	R ²
<i>Bias</i>	Slope	Slope	Slope
	Intercept	Intercept	Intercept
<i>Error</i>	----	----	MAE
	RMSE	----	RMSE
	NRMSE	----	NRMSE
<i>Uncertainty</i>	----	REU	----
<i>Data coverage</i>	Data	Data	Data
	completeness	Capture	Capture Rate
<i>Inter-sensor precision</i>	SD	u _(bs,s)	S _{r,f}
	CV	----	----

252 [References in the table:](#)

253 ¹EPA/600/R-20/279 Performance Testing Protocols, Metrics, and Target Values for Ozone Air Sensors.

254 ²EPA/600/R-20/280 Performance Testing Protocols, Metrics, and Target Values for Fine Particulate Matter
255 Air Sensors.

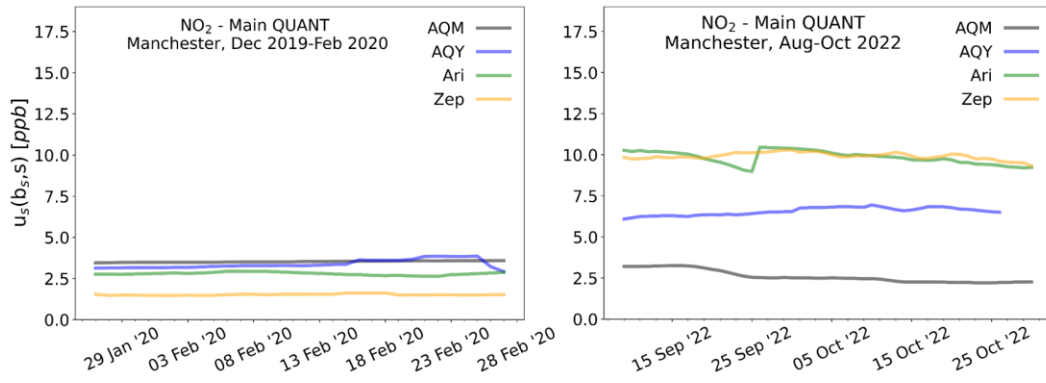
256 ³CEN/TS 17660-1: Air quality - Performance evaluation of air quality sensor systems - Part 1 Gaseous
257 pollutants in ambient air.

258 ⁴ASTM D8406-22: Standard Practice for Performance Evaluation of Ambient Outdoor Air Quality Sensors
259 and Sensor-based Instruments for Portable and Fixed-point Measurement.

260 ⁵ASTM WK74812: Standard Specification for Ambient Outdoor Air Quality Sensors and Sensor-based
261 Instruments for Portable and Fixed-Point Measurement.

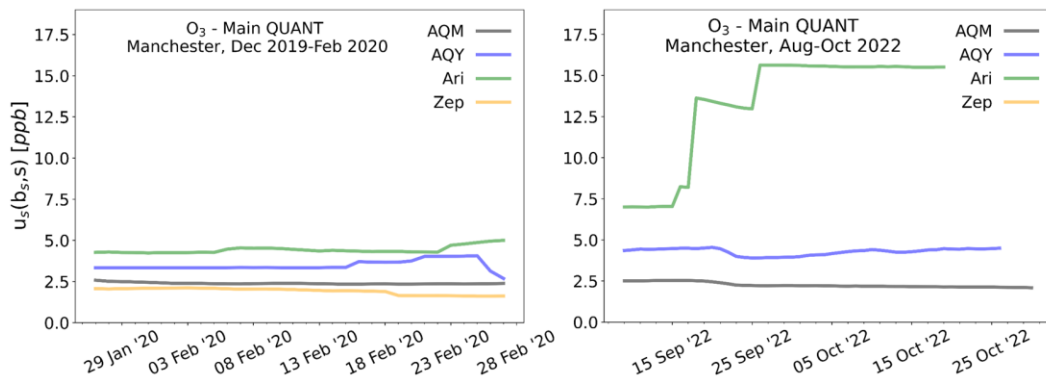
262 Acronyms: EPA: U.S. Environmental Protection Agency; CEN: European Committee for Standardization;
263 ASTM: American Society for Testing and Material. CV: Coefficient of Variation; SD: Standard Deviation
264 (see the definition in the EPA Performance Testing Protocols); u_(bs,s): Between sensor system uncertainty
265 (see the definition in the CEN TS 17660-1); S_{r,f}: field reproducibility standard deviation (see the definition
266 in the ASTM protocols).

267 **S6. Complementary plots**



268

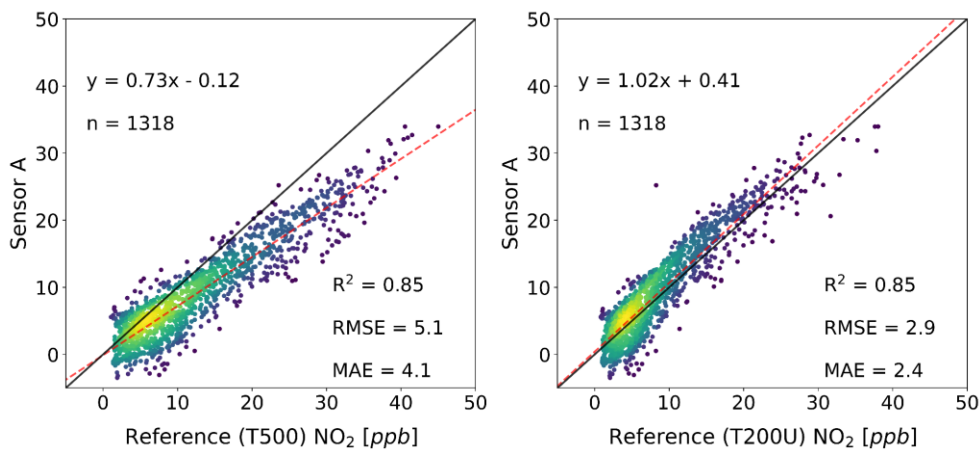
269 Figure S4. Inter-device precision of NO₂ measurements from “identical” devices across the 4 companies
 270 participating in QUANT is assessed using the “between sensor system uncertainty” metric (defined by the CEN/TS
 271 17660-1:2021 as $u(bs, s)$). Each line represents this metric as a composite of all sensors per brand (excluding units
 272 with less than 75% data) within a 40-day sliding window.



273

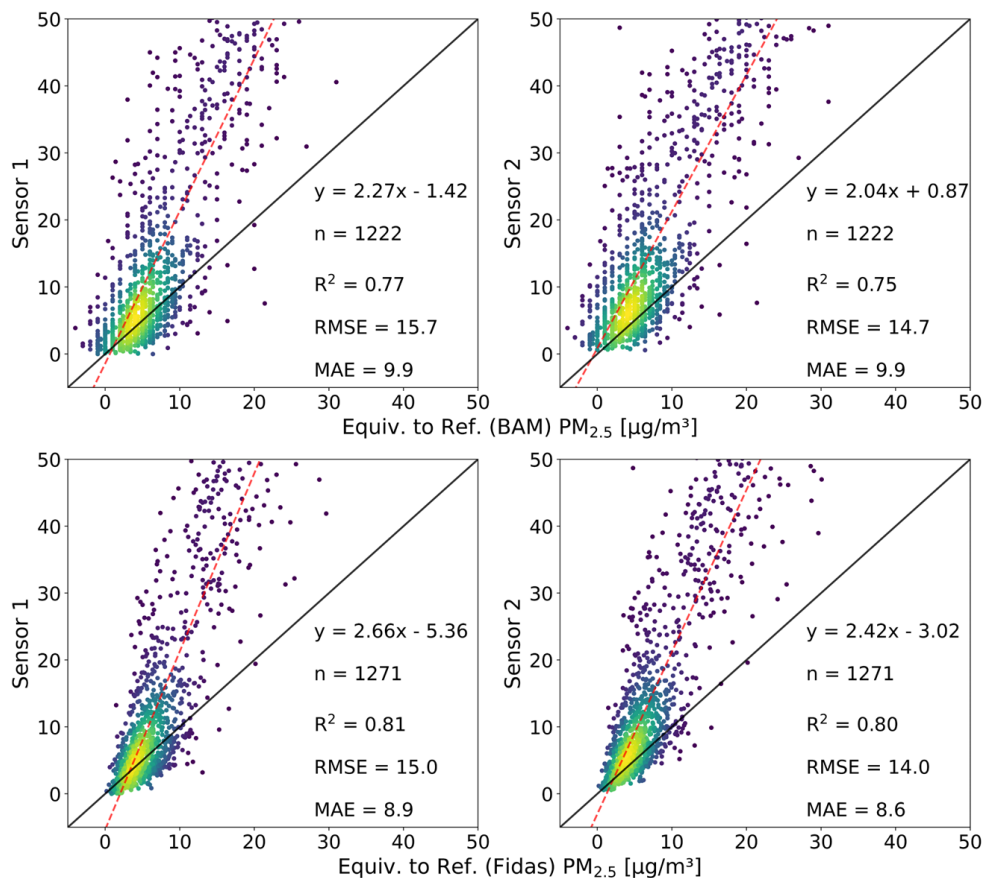
274 Figure S5. The inter-device precision of O₂ measurements from “identical” devices across the 4 companies
 275 participating in QUANT is assessed using the “between sensor system uncertainty” metric (defined by the CEN/TS
 276 17660-1:2021 as $u(bs, s)$). Each line represents this metric as a composite of all sensors per brand (excluding units
 277 with less than 75% data) within a 40-day sliding window.

278 **S5. Sensor performance estimated using different reference methods**



279

280 **Figure S63.** Comparative analysis of “Sensor A” performance against two reference instruments for NO₂
 281 measurements. The left plot shows the correlation with the Teledyne T500 (Cavity Attenuated Phase Shift
 282 Spectroscopy), while the right plot is against the Teledyne T200U (chemiluminescence) and specifically installed at
 283 the Manchester supersite for the QUANT study. The dashed red line represents the line of best fit for the sensor
 284 data against each reference, indicating a closer agreement with the T200U (slope=1.02) compared to the T500
 285 (slope=0.73).

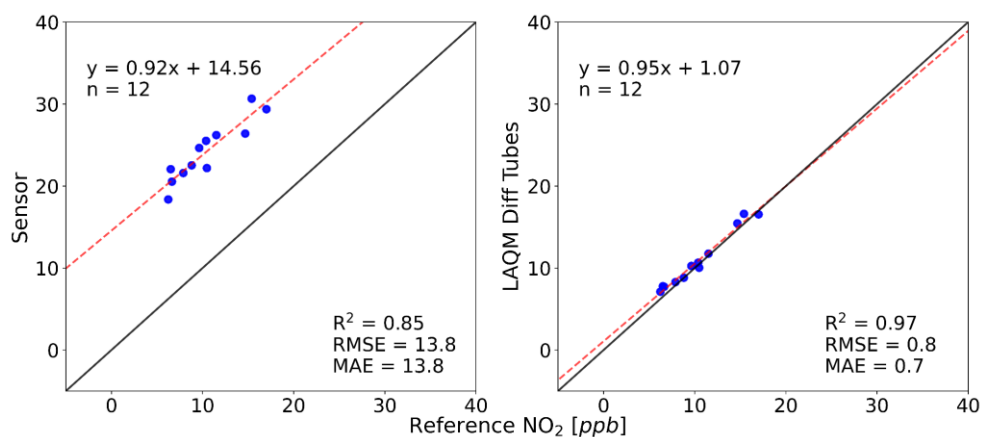


286
 287 **Figure S74.** Comparative regression analysis and performance metrics of two distinct PM_{2.5} sensor systems
 288 benchmarked against a BAM for the top plots and a Fidas for the bottom plots. Each plot demonstrates the
 289 correlation and agreement between the sensor readings and the two equivalent-to-reference instruments in a
 290 roadside site located in York.

291 **S76. NO₂ Diffusion tubes**

292 A diffusion tube co-location study was carried out between November 2020 and November 2021 at the MAQS,
 293 LAQS and York sites, using two types of diffusion tubes: the conventional (also known as LAQM, for Local Air
 294 Quality Management) and UUNN (for UK Urban NO₂ Network). LAQM tubes have an open end and capture
 295 NO₂ which is converted to nitrite when reacting with triethanolamine (TEA) for subsequent analysis. On the other
 296 hand, UUNN tubes, similar in the sampling process to LAQM, include an amorphous polyethylene filter at the
 297 open end to further mitigate the effect of wind on NO₂ measurements. For more details refer to (Butterfield et al.,
 298 2021). Both types of tubes (conventional and UUNN) were installed in duplicates, either in shelters (to limit the
 299 incidence of wind) or directly exposed without protection in mounting blocks. Figure S5 illustrates the

300 performance comparison of traditional diffusion tubes and a sensor system in Manchester. The data from these
301 diffusion tubes have been used to correct the sensor shown here and explained in detail in Section 3.6 (Figures 9b
302 and 9c).



303
304 **Figure S85.** The left plot displays the correlation between an air quality sensor's readings and those from a reference
305 monitor for NO₂, while the right plot demonstrates the LAQM diffusion tube performance. The LAQM plot shows
306 a tighter correlation with the 1:1 line, indicating a higher accuracy in measuring NO₂ concentrations for the period
307 Nov 2020 - Nov 2021 at the Manchester supersite (blue dots represent monthly averages).

308 References

- 309 Butterfield, D., Martin, N. A., Coppin, G., and Fryer, D. E.: Equivalence of UK nitrogen dioxide diffusion tube
310 data to the EU reference method, *Atmos. Environ.*, 262, 118614,
311 <https://doi.org/10.1016/j.atmosenv.2021.118614>, 2021.
- 312 Cross, E. S., Williams, L. R., Lewis, D. K., Magoon, G. R., Onasch, T. B., Kaminsky, M. L., Worsnop, D. R.,
313 and Jayne, J. T.: Use of electrochemical sensors for measurement of air pollution: correcting
314 interference response and validating measurements, *Atmospheric Meas. Tech.*, 10, 3575–3588,
315 <https://doi.org/10.5194/amt-10-3575-2017>, 2017.
- 316 European Commission: Guide to the demonstration of equivalence of ambient air monitoring methods, Report
317 by an EC Working, Group on Guidance. European Commission, 2010.
- 318 Hofman, J., Peters, J., Stroobants, C., Elst, E., Baeyens, B., Van Laer, J., Spruyt, M., Van Essche, W., Delbare,
319 E., Roels, B., Cochez, A., Gillijns, E., and Van Poppel, M.: Air Quality Sensor Networks for Evidence-
320 Based Policy Making: Best Practices for Actionable Insights, *Atmosphere*, 13, 944,
321 <https://doi.org/10.3390/atmos13060944>, 2022.
- 322 Karagulian, F., Barbieri, M., Kotsev, A., Spinelle, L., Gerboles, M., Lagler, F., Redon, N., Crunaire, S., and
323 Borowiak, A.: Review of the Performance of Low-Cost Sensors for Air Quality Monitoring,
324 *Atmosphere*, 10, 506, <https://doi.org/10.3390/atmos10090506>, 2019.

325 Legates, D. R. and McCabe Jr., G. J.: Evaluating the use of “goodness-of-fit” Measures in hydrologic and
326 hydroclimatic model validation, *Water Resour. Res.*, 35, 233–241,
327 <https://doi.org/10.1029/1998WR900018>, 1999.

328 Liu, G., Moore, K., Su, W.-C., Delclos, G. L., Gimeno Ruiz de Porras, D., Yu, B., Tian, H., Luo, B., Lin, S.,
329 Lewis, G. T., Craft, E., and Zhang, K.: Chemical explosion, COVID-19, and environmental justice:
330 Insights from low-cost air quality sensors, *Sci. Total Environ.*, 849, 157881,
331 <https://doi.org/10.1016/j.scitotenv.2022.157881>, 2022.

332 Meyer, M., Afshar-Mohajer, N., Cross, E., and Mudgett, P.: Feasibility of using Low-Cost COTS Sensors for
333 Particulate Monitoring in Space Missions, 2022.

334 Miech, J. A., Stanton, L., Gao, M., Micalizzi, P., Uebelherr, J., Herckes, P., and Fraser, M. P.: Calibration of
335 Low-Cost NO₂ Sensors through Environmental Factor Correction, *Toxics*, 9, 281,
336 <https://doi.org/10.3390/toxics9110281>, 2021.

337 Myklebust, H., Aarhaug, T. A., and Tranell, G.: Use of a Distributed Micro-sensor System for Monitoring the
338 Indoor Particulate Matter Concentration in the Atmosphere of Ferroalloy Production Plants, *JOM*, 74,
339 4787–4797, <https://doi.org/10.1007/s11837-022-05487-7>, 2022.

340 Oizom - Polludrone Smart: <http://www.aqmd.gov/aq-spec/sensordetail/oizom---polludrone-smart>, last
341 access: 27 January 2023.

342 Ouimette, J. R., Malm, W. C., Schichtel, B. A., Sheridan, P. J., Andrews, E., Ogren, J. A., and Arnott, W. P.:
343 Evaluating the PurpleAir monitor as an aerosol light scattering instrument, *Atmospheric Meas. Tech.*,
344 15, 655–676, <https://doi.org/10.5194/amt-15-655-2022>, 2022.

345 Puttaswamy, N., Sreekanth, V., Pillarisetti, A., Upadhya, A. R., Saidam, S., Veerappan, B., Mukhopadhyay, K.,
346 Sambandam, S., Sutaria, R., and Balakrishnan, K.: Indoor and Ambient Air Pollution in Chennai, India
347 during COVID-19 Lockdown: An Affordable Sensors Study, *Aerosol Air Qual. Res.*, 22, 210170,
348 <https://doi.org/10.4209/aaqr.210170>, 2022.

349 Weissert, L. F., Alberti, K., Miskell, G., Pattinson, W., Salmond, J. A., Henshaw, G., and Williams, D. E.: Low-
350 cost sensors and microscale land use regression: Data fusion to resolve air quality variations with high
351 spatial and temporal resolution, *Atmos. Environ.*, 213, 285–295,
352 <https://doi.org/10.1016/j.atmosenv.2019.06.019>, 2019.

353 Westgate, S. and Ng, N. L.: Using in-situ CO₂, PM₁, PM_{2.5}, and PM₁₀ measurements to assess air change
354 rates and indoor aerosol dynamics, *Build. Environ.*, 224, 109559,

355 <https://doi.org/10.1016/j.buildenv.2022.109559>, 2022.

356 Willmott, C. J. and Matsuura, K.: Advantages of the mean absolute error (MAE) over the root mean square error
357 (RMSE) in assessing average model performance, *Clim. Res.*, 30, 79–82,
358 <https://doi.org/10.3354/cr030079>, 2005.

359 Zauli-Sajani, S., Marchesi, S., Boselli, G., Broglia, E., Angella, A., Maestri, E., Marmioli, N., and Colacci, A.:
360 Effectiveness of a Protocol to Reduce Children’s Exposure to Particulate Matter and NO₂ in Schools
361 during Alert Days, *Int. J. Environ. Res. Public. Health*, 19, 11019,
362 <https://doi.org/10.3390/ijerph191711019>, 2022.