# 1 Air pollution measurement errors: Is your data fit for

# 2 purpose?

3 Sebastian Diez[1], Stuart Lacy[1], Thomas J. Bannan[2], Michael Flynn[2], Tom Gardiner[3], David

4 Harrison[4], Nicholas Marsden[2], Nick Martin[3], Katie Read[1,5], Pete M. Edwards[1]

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

6 [2]Department of Earth and Environmental Science, Centre for Atmospheric Science, School of Natural Sciences,

7 The University of Manchester, Manchester, M13 9PL, UK

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

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

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

11 *Correspondence:*

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

13 **Abstract.** When making measurements of air quality, having a reliable estimate of the measurement uncertainty

14 is key to assessing the information content that an instrument is capable of providing, and thus its usefulness in a

15 particular application. This is especially important given the widespread emergence of Low Cost Sensors (LCS)

16 to measure air quality. To do this, end users need to clearly identify the data requirements a priori and design

17 quantifiable success criteria by which to judge the data. All measurements suffer from errors, with the degree to

18 which these impact the accuracy of the final data often determined by our ability to identify and correct for them.

19 The advent of LCS has provided a challenge in that many error sources show high spatial and temporal variability,

20 making laboratory derived corrections difficult. Characterizing LCS performance thus currently depends primarily

21 on colocation studies with reference instruments, which are very expensive and do not offer a definitive solution

22 but rather a glimpse of LCS performance in specific conditions over a limited period of time. Despite the

23 limitations, colocation studies do provide useful information on measurement device error structure, but the results

24 are non-trivial to interpret and often difficult to extrapolate to future device performance. A problem that obscures

25 much of the information content of these colocation performance assessments is the exacerbated use of global

26 performance metrics ($R^2$, RMSE, MAE, etc.). Colocation studies are complex and time-consuming, and it is easy

27 to fall into the temptation to only use these metrics when trying to define the most appropriate sensor technology

28 to subsequently use. But the use of these metrics can be limited, and even misleading, restricting our understanding

29 of the error structure and therefore the measurements' information content. In this work, the nature of common

30 air pollution measurement errors is investigated, and the implications these have on traditional metrics and other

31 empirical, potentially more insightful, approaches to assess measurement performance. With this insight we

32 demonstrate the impact these errors can have on measurements, using a selection of LCS deployed alongside

33 reference measurements as part of the QUANT project, and discuss the implications this has on device end-use.

Atmospheric
Measurement
Techniques

Discussions

34    **Keywords:** air quality measurements, errors, uncertainty, information content, low-cost sensors.

35    **1. Introduction**

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

51    All measurements have an associated uncertainty, and even in highly controlled laboratory assessments, the true
52    value is not known, with any measurement error defined relative to our best estimate of the range of possible true
53    values. However, quantifying and representing error and uncertainty is a challenge for a wide range of analytical
54    fields, and often what these concepts represent is not the same to all practitioners. This results in a spectrum of
55    definitions that take into account the way truth, error, and uncertainty are conceived (Cross et al., 2017; Grégis,
56    2019; Kirkham et al., 2018; Mari et al., 2021). For atmospheric measurements assessing uncertainty is complex
57    and non-trivial. Firstly, given the "true" value can never be known, an agreed reference is needed. Secondly, the
58    constantly changing atmospheric composition means that repeat measurements cannot be made and the traditional
59    methods for determining the random uncertainty are not applicable. And finally, a major challenge arises from the
60    multiple sources of error both internal and external to the sensor that can affect a measurement. Signal responses
61    from a non-target chemical or physical parameter or electromagnetic interference are examples of an almost
62    limitless number of potential sources of measurement error. In this work, we will follow the definitions given by
63    the International Vocabulary of Metrology (JCGM, 2012) for measurement error ("measured quantity value minus
64    a reference quantity value") and for measurement uncertainty ("non-negative parameter characterizing the
65    dispersion of the quantity values being attributed to a measurand, based on the information used").

66    The difficulties in defining and quantifying uncertainty across the full range of end-use applications of a
67    measurement device, means that often the quoted measurement uncertainty is not applicable, or in some cases not
68    provided or provided in an ambiguous manner. This makes assessing the applicability of a measurement device
69    to a particular task difficult for end-users. In this work, we investigate the nature of common air pollution
70    measurement errors, and the implications these have on traditional goodness-of-fit metrics and other, potentially
71    more insightful approaches to assess measurement uncertainty. We then use this insight to demonstrate the impact
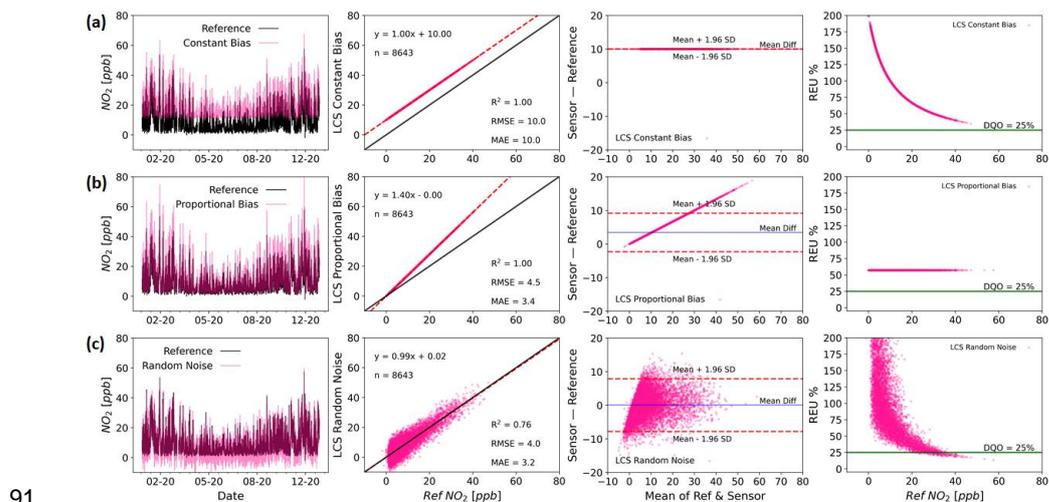
72    these errors can have on measurements, using a selection of LCS deployed alongside reference measurements as

73    part of the UK Clean Air program funded QUANT (Quantification of Utility of Atmospheric Network

74    Technologies) project, a 2 year colocation study of 26 commercial LCS devices (56 gases measurements and 56

75    PM measurements) at multiple urban, background and roadside, locations in the UK. After analysing some of the

76    real-life uncertainty characteristics we discuss the implications this has on data use.

**2. Error characterization**

78    When characterizing measurement error, in the absence of evidence to the contrary it is often assumed a linear

79    additive model. Once the analytical form of the model is defined, its parameters aim to capture the error

80    characteristics, and in the case of linear models (Eq. (1)), these are typically separated into three types (Tian et al.,

81    2016): (i) proportional bias or scale error ($b_1$), (ii) constant bias or displacement error ($b_0$) and (iii) random error

82    ($\varepsilon$) (Tian et al., 2016). Any measurement ($y_i$, e.g from the LCS) can therefore be thought of as a combination of

83    the reference value ($x_i$) and the three error types, such that:

84    $$y_i = b_1 x_i + b_0 + \varepsilon \qquad (1)$$

85    As the simplest approximation, this linear relationship for the error characteristics is often used to correct for

86    observed deviations between measurements and the agreed reference. It is worth to note, however, that this

87    equation assumes time-independent error contributions and that the three error components are not correlated,

88    which is often not the case on both counts (e.g. responses to non-target compounds). The parameter values

89    determined for Eq. (1) are also generally only applicable for individual instruments, potentially in specific

90    environments, unless the transferability of these parameters between devices has been explicitly demonstrated.
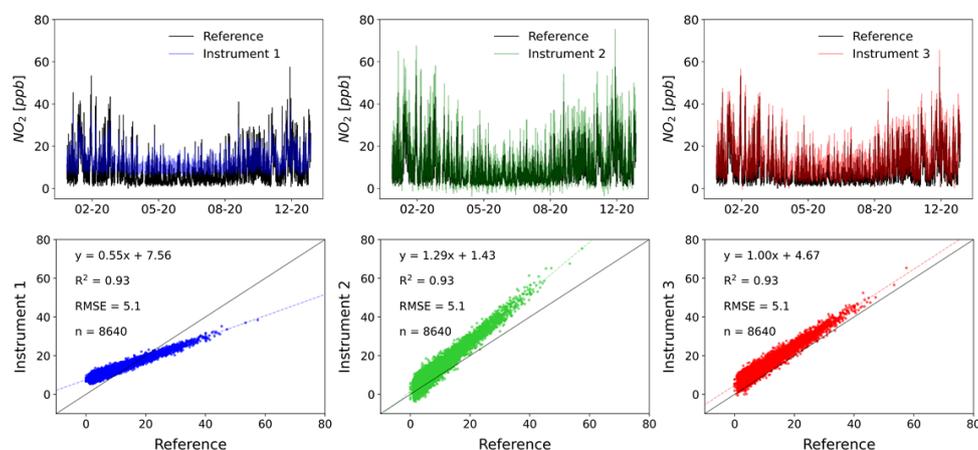


91

**Figure 1. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and REU (right panels, DQO for NO₂ = 25%) for arbitrary examples of pure constant bias (Slope = 1, Intercept = 1, SD$_\varepsilon$ = 0; a-panels), pure proportional bias (Slope = 1.4, Intercept = 0, SD$_\varepsilon$ = 0; b-panels) and pure random noise (Slope = 1, Intercept = 0, SD$_\varepsilon$ = 4; c-panels) simulated errors.**

96   Figure 1 shows examples of how pure constant bias (a-panels), pure proportional bias (b-panels), and pure random

97   noise (c-panels) would look like. In each of these ideal cases, the error plots enable the practitioner to view the

98   error characteristics in slightly different ways, allowing the impacts of the observed measurement uncertainty to

99   be placed into the context of the data requirements. In this work, we will refer to them as "error types" (in contrast

100  to "error sources"), which is the way they are distilled by the linear error model.

### 2.1 Performance indices, error structure and uncertainty

102  A major challenge faced by end-users of measurement devices characterised using colocation studies is the non-

103  trivial question of how the comparisons themselves are performed and how the data is communicated. Often single

104  value performance metrics, such as the coefficient of determination ($R^2$) or root mean squared error (RMSE), are

105  calculated between the assessed method (e.g. LCS) and an agreed reference, and the user is expected to infer an

106  expected device performance or uncertainty for a measurement in their application (Cross et al., 2017; Duvall et

107  al., 2016; Malings et al., 2019). These metrics contain useful information about the measurement, but they are

108  unable to fully describe the error characteristics, in part because they reduce the error down to a single value (Tian

109  et al., 2016). Unfortunately, the widespread use of a small number of metrics as the sole method to assess

110  measurement uncertainty, without a thorough consideration of the nature of the measurement errors, means

111  measurement devices are often chosen that are unable to provide data that is fit for purpose. In addition,

112  unconscious about potential flaws, users (e.g. researchers) could communicate findings or guide decision making

113  based on results that may not justify the conclusions drawn from the data. Figure 2 shows three simulated

114  measurements compared with the true values. Despite the measurements having identical $R^2$ and RMSE values,

115  the time series and regression plots show that the error characteristics are significantly different, and would impact

116  how the data from such a device could viably be used.



117

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

121  There are multiple performance metrics that can be used for the assessment of measurement errors and uncertainty.

122  Tian et al (2016) present an excellent summary of some of the major pitfalls of performance metrics and promote

123    an approach of error modelling as a more reliable method of uncertainty quantification. These modelling

124    approaches, however, rely on the assumption of statistical stationarity, whereby the statistical properties of the

125    error are constant in the temporal and spatial domains. The presence of unknown or poorly characterised sources

126    of error, for example, due to interferences from other atmospheric constituents or drifts in sensor behaviour, makes

127    this assumption difficult to satisfy, especially when the dependencies of these errors show high spatial and

128    temporal variability. Thus, if field colocation studies are the primary method for performance assessment, as is

129    the case for LCS, only through a detailed assessment of the measurement errors across a wide range of conditions

130    and timescales can the uncertainty of the measurement be realistically estimated.

131    **2.2 Dealing with errors: established techniques vs Low-Cost Sensors**

132    Different approaches are available to the user to minimise the impact of errors, generally by making corrections

133    to the sensor data. For example, in the case of many atmospheric gas analysers, if the error is dominated by a

134    proportional bias, a multi-point calibration can be performed using standard additions of the target gas.

135    Displacement errors can be quantified, and then corrected for, by sampling a gas stream that contains zero target

136    gas. And random errors can be reduced by applying a smoothing filter (e.g moving average filter, time-averaging

137    the data, etc.), at the cost of losing some information (Brown et al., 2008). These approaches work well for simple

138    error sources that, ideally, do not change significantly over timescales from days to months. Unfortunately, more

139    complex error sources can manifest in such a way that they contribute across all three error types, and also vary

140    temporally and spatially. For example, an interference from another gas-phase compound could in part manifest

141    itself as a displacement error, based on the instrument response to its background value, and in part as a

142    proportional bias if its concentration correlates with the target compounds, with any short-term deviations from

143    perfect correlation contributing to the random error component. In this case, time-averaging combined with

144    periodic calibrations and zeros would not necessarily minimise the error, and the user would need to employ

145    different tactics. One option would be to independently measure the interferent concentration, albeit with

146    associated uncertainty, and then use this to derive a correction. This is feasible if a simple and cost-effective

147    method exists for quantifying the interferent and its influence on the result is understood, but can make it very

148    difficult to separate out error sources, and can become increasingly complex if this measurement also suffers from

149    other interferences.

150    For many measurement devices, in particular for LCS based instruments, a major challenge is that the sources and

151    nature of all the errors are unknown or difficult to quantify across all possible end-use applications, meaning

152    estimates of measurement uncertainty are difficult. In the case of most established research and reference-grade

153    measurement techniques, comprehensive laboratory and field experiments have been used to explore the nature

154    of the measurement errors (Gerboles et al., 2003; Zucco et al., 2003). Calibrations have then been developed,

155    where traceable standards are sampled and measurement bias, both constant and proportional, can be corrected

156    for. Interferences from variables such as temperature, humidity, or other gases, have also been identified and then

157    either a solution engineered to minimise their effect or robust data corrections derived. Unfortunately, these

158    approaches have been shown not to perform well in the assessment of LCS measurement errors, due to the

159    presence of multiple, potentially unknown, sensor interferences from other atmospheric constituents (Thompson

160    & Ellison, 2005). These significant sensitivities to constituents such as water vapour and other gases mean

161    laboratory-based calibrations of LCS become exceedingly complex, and expensive, as they attempt to simulate

162 the true atmospheric complexity, often resulting in observed errors being very different to real-world sampling
163 (Rai et al., 2017; Williams, 2020). This has resulted in colocation calibration becoming the accepted method for
164 characterizing LCS measurement uncertainties (De Vito et al., 2020; Masson et al., 2015; Mead et al., 2013;
165 Popoola et al., 2016; Sun et al., 2017), where sensor devices are run alongside traditional reference measurement
166 systems for a period of time, and statistical corrections derived to minimise the error between the two. As the true
167 value of a pollutant concentration cannot be known, this colocation approach assumes all the error is in the low-
168 cost measurement. Although this assumption may often be approximately valid (i.e. reference error variance <<
169 LCS error variance), no measurement is absent of uncertainty and this can be transferred from one measurement
170 to another, obscuring attempts to identify its sources and characteristics. A further consideration when the fast
171 time-response aspect of LCS data is important, is that reference measurement uncertainties are generally
172 characterised at significantly lower reported measurement frequencies (typically 1 hr). This means that a high
173 time-resolution (e.g. 1 min) reference uncertainty must be characterized in order to accurately estimate the LCS
174 uncertainty (requiring specific experiments and additional costs). If a lower time-resolution reference data is used
175 as a proxy, then the natural variability timescales of the target compound should be known and any impact of this
176 on the reported uncertainty caveated.

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

187 **3. Methods**

188 In this work, we explore measurement errors, and their impacts, using the most common single value metrics ($R^2$,
189 RMSE and MAE), along with two additional widely used approaches to visualize the error across a dataset: the
190 Bland-Altman plots (B-A) and Relative Expanded Uncertainty (REU).

191 The performance metrics provide a single value irrespective of the size of the dataset, and might appear convenient
192 for users when comparing across devices or datasets, but can encourage over-reliance on the metric, often at the
193 expense of looking at the data in more detail or bringing an awareness of the likely physical processes driving the
194 error sources. On the other hand, the B-A and the REU plots are more laborious techniques and the interpretation
195 can be challenging for non-experts, but they provide additional insights into the nature of the errors, not attainable
196 by one or more combined performance metrics.

197 In order to understand how the different tools used here show different characteristics of the error structure, some
198 errors commonly found in LCS are examined through simulation studies. Subsequently, two real world case

199    studies are presented: (i) LCS duplicates for $NO_2$ and $PM_{2.5}$ belonging to the QUANT project located in two sites

200    -the Manchester Natural Environment Research Council (NERC) measurement Supersite, and the York Fishergate

201    Automatic Urban and Rural Network (AURN) roadside site- and (ii) a set of duplicate reference instruments (only

202    at Manchester Supersite). Table S1 shows the research grade instrumentation used for this study.

203    **3.1 Visualization tools**

204    An ideal performance metric should be able to deliver not only a performance index but also an idea of the

205    uncertainty distribution (Chai & Draxler, 2014). This is difficult to deliver through a simple numerical value, and

206    easy to interpret visualisations of the data are often much more useful for conveying multiple aspects of data

207    performance. Figure 2 shows the two most common data visualisation tools, the time-series plot and the regression

208    plot. In the time series plot the instrument under analysis and the agreed reference are plotted together as a function

209    of time. This allows a user to visually assess tendencies of over or under prediction, differences in the base line

210    or other issues, but can be readily over interpreted and does not allow for easy quantification of the observed

211    errors. In the regression plot the data from the instrument under analysis is plotted against the agreed reference

212    data. This allows for the correlation between the two methods to be more readily interpreted, in particular any

213    deviations from linearity, but gives little detail on the nature of the errors themselves.

214    In contrast to the regression plot, Bland-Altman (B-A) plots essentially display the difference between

215    measurements, enabling more information on the nature of the error to be communicated. B-A plots (Altman &

216    Bland, 1983) are extensively used in analytical chemistry and biomedicine to evaluate the differences between

217    two measurement techniques (Doğan, 2018). The B-A is a scatter plot, in which the abscissa represents the average

218    of these measures (e.g LCS and a reference measurement), acknowledging that the true value is unknown and that

219    both measurements have errors, and the ordinate shows the difference between the two paired measurements. In

220    the case where all the error is assumed to be in one of the measurements, e.g. comparing a LCS to a reference

221    grade measurement, there is an argument that the B-A abscissa could be the agreed reference value instead of the

222    average of two measurements. However, in this work we use the average of the two values as per the traditional

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

224    expression:

225    $$y_i - x_i = x_i (b_1 - 1) + b_0 + \varepsilon \qquad (2)$$

226    From Eq. (2) it can be seen that if $b_1 \neq 1$ or if the error term ($\varepsilon$) variance is non-constant (e.g. heteroscedasticity)

227    the difference will not be normally distributed. The B-A plot (with $x_i$ as the reference instrument results) allows a

228    quick visual assessment of the error distribution without the need to calculate the model parameters. In the case

229    the differences are normally distributed, the so-called "agreement interval" (usually defined as $\pm 2\sigma$ around the

230    mean) will hold 95% of the data points. Even though the estimated limits of agreement will be biased if the

231    differences are not normally distributed, it can still be a valuable indicator of agreement between the two

232    measurements.

233    If the ultimate goal of studying measurement errors is to diagnose the measurement uncertainty in a particular

234    target measurement range, then visualising the uncertainty in pollutant concentration space can be very

235    informative. The REU (GDE, 2010) provides a relative measure of the uncertainty interval about the measurement

236 within which the true value can be confidently asserted to lie. The abscissa in an REU plot represents the agreed

237 reference pollutant concentration, whose error is taken into account, something not considered by the other metrics

238 or visualisations discussed. The REU is regularly used to assess measurement compliance with the Data Quality

239 Objective (DQO) of the European Air Quality Directive 2008/50/EC, and is mandatory for the demonstration of

240 equivalence of methods other than the EU reference methods. For LCS the REU is widely used as a performance

241 indicator (Bagkis et al., 2021; Bigi et al., 2018; Castell et al., 2017; Cordero et al., 2018; Spinelle et al., 2015).

242 However, the evaluation of this metric is perceived as arduous and cumbersome and it is not included in the

243 majority of sensor studies (Karagulian et al., 2019). There is now a new published European Technical

244 Specification (TS) for evaluating the LCS performance for gaseous pollutants (CEN/TS 17660-1:2021). It

245 categorizes the devices in 3 classes according to the DQO (Class 1 for "indicative measurements", Class 2 for

246 "objective estimations", and Class 3 for non-regulatory purposes, e.g. research, education, citizen science, etc.).
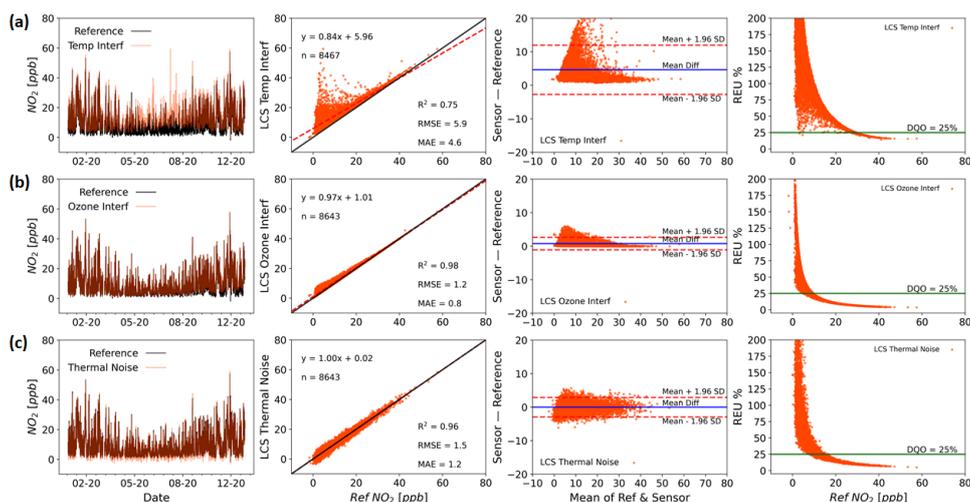
247 In the following sections, we use these established methods for assessing measurement uncertainty, alongside

248 simple time series and regression plots, to explore different error sources and their implications for air pollution

249 measurements.

**4. Case studies**

**4.1 Simulated instruments**

252 In order to investigate the impact of different origins of measurement error on measurement performance, a set of

253 simulated datasets have been created. These data are derived using real-world reference data as the true values,

254 with the subsequent addition of errors of different origins to generate the simulated measurement data. Error

255 origins were chosen for which examples have been described in the LCS literature. Performance metrics along

256 with visualization methods are then used to assess measurement performance.
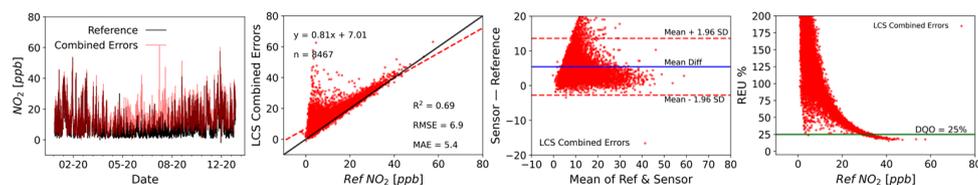
257 As the complexity of the error increases, the impact of the assumption of statistical stationarity can become more

258 difficult to satisfy, with the magnitude of the errors becoming less uniform across the observed concentration, and

259 hence spatial, or time domains. Figure 3 shows examples of modelled sources of errors on $NO_2$ measurements:

260 temperature interference (correction model taken from (Popoola et al., 2016), a-panels), a non-target gas (ozone)

261 interference (correction model taken from (Peters et al., 2021), b-panels) and thermal electrical noise (white noise,

262 c-panels).

Atmospheric
Measurement
Techniques
Discussions
Open Access
EGU

**Figure 3. Time series (left panels), regression plots (middle-left panels, including R$^2$, RMSE & MAE), Bland-Altman plots (middle-right panels) and REU (right panels, DQO for NO$_2$ = 25%) for temperature (a-panels), ozone (b- panels) and thermal electrical noise (c-panels) modelled interferences on NO$_2$ measurements (time res 1 h).**

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

Each error source impacts the uncertainty of the measurement, which in turn impacts its ability to provide useful information for a particular task. For example, the form of the temperature interference shown in Fig. 3 (a-panels) results in the largest errors being seen at the lower NO$_2$ values. This is because NO$_2$ concentrations are generally lowest during the day, due to photolytic loss when temperatures are highest. Thus this device would be better suited to an end-user intending to assess daily peak NO$_2$ concentration compared with the daytime hourly exposure values, providing the environment the device was deployed in showed a similar relationship between temperature and true NO$_2$ as that used here. The O$_3$ interference shown in Fig. 3 (b-panels) is similar, due again to a general anti-correlation observed between ambient O$_3$ and NO$_2$ concentrations. This type of interference can often be interpreted incorrectly as a proportional bias, and a slope correction applied to the data. However, this type of correction will ultimately fail as O$_3$ concentrations are dependent on a range of factors, such as hydrocarbon concentrations and solar radiation, and as these change the O$_3$ concentration relative to the NO$_2$ concentration will change. To further complicate matters, multiple error sources can act simultaneously, meaning that the majority of measurements will contain multiple sources of error. Figure 4 shows a simple linear combination of the modelled errors shown in Fig 3, and the impact this has on the performance metrics.

9

**Figure 4. Time series (left panel), regression plot (middle-left panel, including R$^2$, RMSE & MAE), Bland-Altman plot (middle-right panel) and REU (right panel, DQO for NO$_2$ = 25%) for a linear combination of temperature, ozone and thermal electrical noise modelled interferences (time res 1 h).**

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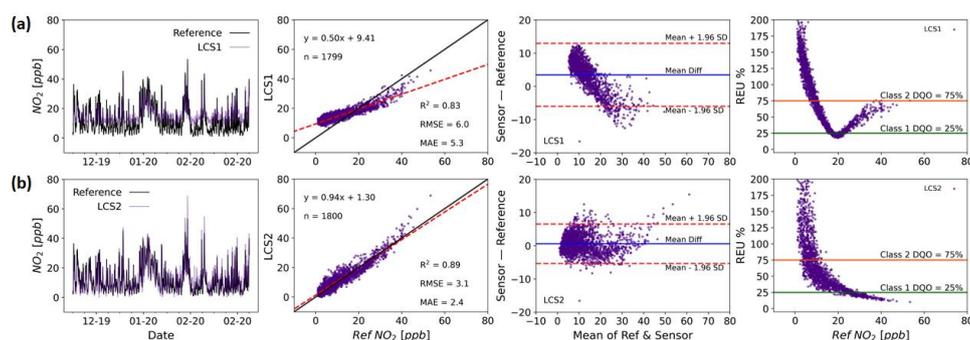
## 4.2 Real-world instruments

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

Figure 5 shows two colocated NO$_2$ measurements, from two different LCS devices using only their out-of-box calibrations (i.e. no colocation data from that site was used to improve performance), compared with colocated reference measurements at an urban background site in the city of Manchester. Unlike the modelled instruments in Sect. 4.1, the combination of error sources is unknown in this case and we can thus only assess the LCS measurement performance through comparison with the reference measurements using metrics and visual tools. There are obvious differences in the performance of both LCS instruments shown in Fig. 5. LCS1 (a-panels) shows an appreciable difference in the time-series baseline, which can be interpreted from both the regression (b$_1$ <1) and the B-A plots as a proportional bias. This bias also impacts the REU plot, with a minimum in the region where the regression best fit line crosses the 1:1 line (~17ppb). It is worth noting that these plots do not directly identify the source of the proportional bias, with sensor response to the target compound or another covarying compound possible, but provides information on how much it impacts the data. For LCS2 (Fig. 5, b-panels) any

Atmospheric
Measurement
Techniques
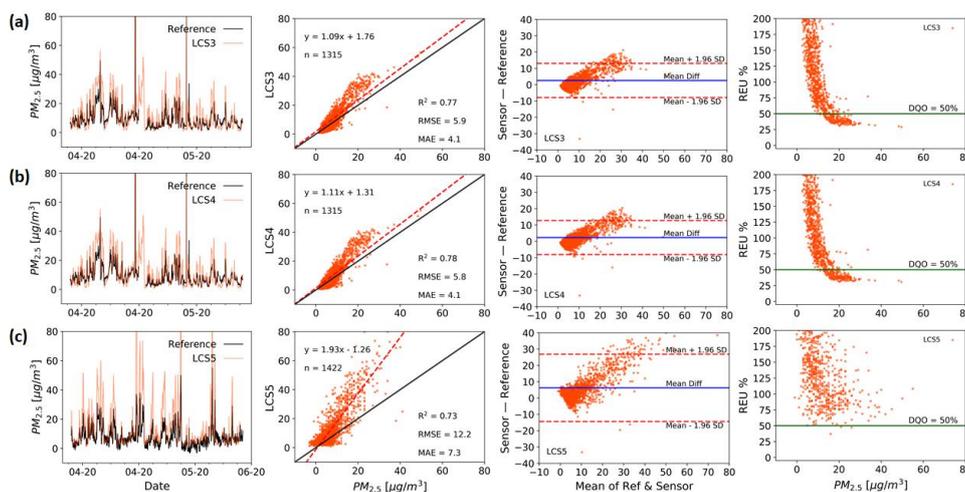Discussions

Open Access

EGU

324 proportional bias is significantly smaller, with the B-A plot showing a much more symmetrical distribution of

325 points around the central line across the observed mixing ratio range, although this is not a normal distribution as

326 evidenced by the heteroscedastic nature of the differences, indicating the cause is not entirely random in nature.

327 The lack of a large proportional bias also results in the REU plot showing a continued reduction in relative

328 uncertainty as the true $NO_2$ concentration increases. Interestingly, both LCS's also show an additional bias at the

329 highest $NO_2$ values observed. This does not significantly impact the REU, due to its relative nature, but can be

330 seen in the regression and B-A plots. Correcting for the observed proportional bias in LCS1 and LCS2 improves

331 the observed performance by providing the errors with a more symmetrical distribution (LCS1* and LCS2* shown

332 in Fig. S2).



333

334 **Figure 5. Time series (left panels), regression plots (middle-left panels), Bland-Altman plots (middle-right panels) and**

335 **REU (right panels; $NO_2$ Class 1 DQO = 25% & Class 2 DQO = 75%) for $NO_2$ measurements by two LCS systems of**

336 **different brands ( a and b panels) in the same location (Manchester Supersite, December 2019 to February 2020. Time**

337 **res 1 h).**

338 Figure 6 shows three out-of-the-box $PM_{2.5}$ measurements made by three devices from the same brand in spring,

339 located at two sites: the first two at an urban background (LCS3 & LCS4, a and b panels) and the third at a roadside

340 (LCS5, c panels). As the regression and the B-A plots show, all LCS measurements in Fig. 6 have a proportional

341 bias compared with the reference, with the LCS over predicting the reference values. Both LCS's at the urban

342 background site show very similar performance, indicating that the devices are similarly affected by errors. This

343 internal consistency is highly desirable, especially when LCS's are to be deployed in networks, as although mean

344 absolute measurement error may be high, differences between identical devices are likely to be interpretable.

Atmospheric
Measurement
Techniques
Discussions

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

The LCS data from the roadside location (LCS5) show significantly lower precision than those at the urban background site, as seen in the B-A plot. This could be caused by differences in particle properties and size distributions between the two sites (Gramsch et al., 2021), and by the high frequency variation of transport emissions close to the roadside side and turbulence effects (Baldauf et al., 2009; Makar et al., 2021). Duplicate measurements show that all sensors of this type responded similarly in this roadside environment (not shown here), supporting the high internal consistency of this device, but indicating a spatial heterogeneity in some key error sources. It is also worth noting that the gold standard instruments at the two sites are not "reference method" but "reference equivalent methods" (GDE, 2010), each using a different measurement technique: while an optical spectrometer (Palas Fidas 200) is used in Manchester, the York instrument uses a Beta attenuation method (Met One BAM 1020), which could also potentially lead to some of the observed differences. The increased apparent random variability for LCS5, combined with the proportional bias, results in significantly higher measurement uncertainty across the observed range, as can be seen by the REU plots, with LCS5 never reaching an acceptable DQO level (50% for PM$_{2.5}$). As with the NO$_2$ sensors (Fig. 5), if the observed proportional bias is corrected the linearly bias-corrected sensors (Fig. S3) show a much improved comparison with the reference measurement, specially LCS5*. In this case the B-A plot now shows an error characteristic more dominated by random errors, and the REU plots shows a significant reduction, with the REU at 10 ugm$^{-3}$ reducing from ~75 to ~50%.
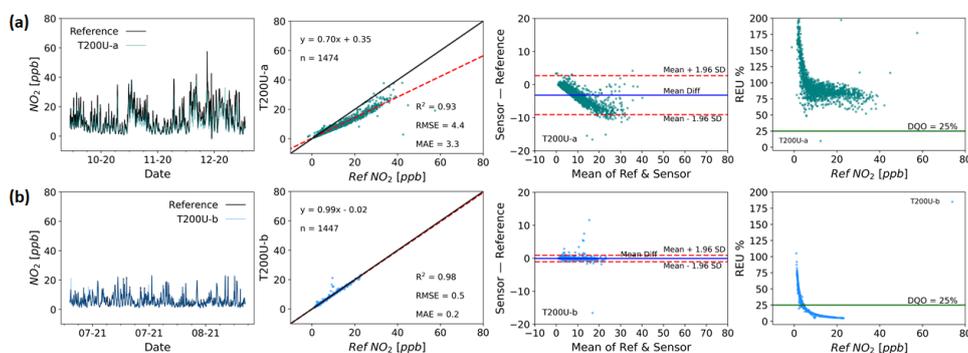
As a comparison for the LCS data shown above, Fig. 7 shows two identical NO$_2$ reference grade instruments, Teledyne T200U (Chemiluminescence method) at the Manchester urban background site (panels a and b) at during two different time periods, with a Teledyne T500U (CAPS detection method) used as the "ground truth" instrument. Instrument "a" manifests a significant proportional bias, in contrast to instrument "b", but both show differences that could be non-negligible depending on the application. The deviations observed in instrument "a" was due to the cell pressure being above specification by ~20%, unnoticed while the instrument was in operation.

372  This demonstrates the importance of checking instrument parameters regularly in the field even if the data appears

373  reasonable.

374  As the LCS error structure is determined relative to the performance of a reference measurement, if the reference

375  instrument suffers from significant errors this will affect the outcomes of the performance assessment, due to the

376  assumption that all the errors reside with the LCS. As Fig. 7 shows, however, this assumption is not necessarily

377  always valid and potentially argues that reference instruments used in colocation studies should be subject to

378  further error characterisation, including possible colocation with other reference instruments. As a similar

379  comparison of reference instruments, Fig. S4 shows two ozone research grade instruments (a Thermo 49i and a

380  2B).

381  It is worth noting that even when using reference, or reference equivalent, grade instrumentation, inherent

382  measurement errors mean that relative uncertainty, as shown in the REU plot, increases asymptotically at lower

383  values. This is not unexpected, but is potentially important as ambient target concentration recommendations

384  continue to fall based on updated health evidence (World Health Organization, 2021).

385



386  **Figure 7. Time series (left panel), regression plots (middle-left panel), Bland-Altman plots (middle-right panel) and**

387  **REU (right panel, DQO for NO₂ = 25%) for two identical (Teledyne T200U) reference NO₂ instruments (panels a and**

388  **b) colocated at the Manchester Supersite (1h time res). The first instrument between October & November 2020 and**

389  **the second between July & August 2021.**

390  **5. Discussion**

391  The widespread use of colocation studies to assess measurement device performance, means many examples exist

392  in the LCS literature where different devices are compared using summary metrics for field or laboratory studies

393  (Broday, 2017; Duvall et al., 2016; Hofman et al., 2022; Karagulian et al., 2019; Mueller et al., 2017; Rai et al.,

394  2017; van Zoest et al., 2019). Although these comparisons do provide useful information, they can be misleading

395  for end users wanting to compare the performance of different devices, as they are often carried out under different

396  conditions and do not present the data or experimental design in full. Even in the case where comparisons have

397  been done under identical conditions, the data still needs to be treated with caution, as inevitable differences

398  between assessment environment and proposed application environment, as well as any changes to

399  instrument/sensor design or data processing, mean that past performance does not guarantee future performance.

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

408 Single value performance metrics, such as $R^2$ and RMSE, can seem convenient when comparing multiple co-
409 located devices as they facilitate decision making when a threshold criterion is defined. However, these scalar
410 values hide important information about the scale and / or distribution of the errors within a dataset; graphical
411 summaries of the measurements themselves can offer significantly more insight into the impact of measurement
412 errors on device performance and ultimate capabilities. Of particular use in air pollution measurements is the
413 ability to see how the errors manifest themselves in relation to our best estimate of the true pollutant concentration,
414 as often applications have specific target pollutant concentration ranges of interest. For example the two $NO_2$ LCS
415 devices shown in Fig. 5 have similar $R^2$ values of 0.83 and 0.89, but one is suffering from a strong proportional
416 bias that impacts on measurements either side of the 18ppb crossing-point. Errors, or combinations of errors,
417 frequently result in varying magnitude of the observed measurement inaccuracies across the concentration space
418 observed, and it is often useful to assess both the absolute and relative effects of the errors. By getting a more
419 complete picture of the device performance, decisions can be made on the effectiveness of simple corrections,
420 such as correcting for an apparent proportional bias using an assumption of a linear error model. Ultimately end
421 users need to identify the data requirements a priori and design quantifiable success criteria by which to judge the
422 data. For example, rather than just wanting to measure the 8 hour average $NO_2$, be more specific and require that
423 this needs to be accurate to within 5 ppb, have demonstrated approximately normally distributed errors in a
424 representative environment for the period of interest, and no statistical evidence of deviation from a linear
425 correlation with the reference measurement over the target concentration range for the period of interest.

426 A major challenge comes from complex errors, such as interferences from other compounds or with environmental
427 factors, that vary temporally and/or spatially. Similar graphical techniques to those presented above can be used
428 to identify the existence of such relationships, but correcting for them remains a challenge. This brings into
429 question the power of colocation studies, as they can ultimately never be performed under the exact conditions
430 for every intended application. The $PM_{2.5}$ sensors shown in Fig. 6 demonstrate this, as if a colocation dataset
431 generated at the urban background site was used to inform a decision about the applicability of these devices to a
432 roadside monitoring task, then an overly optimistic assessment of the scale of the errors to be expected would be
433 likely. It is therefore always desirable that colocation studies are as relevant as possible to the desired application,
434 and this is even more paramount in the case where the error sources are poorly specified. For this reason, complete
435 meta-data on the range of conditions over which a study was conducted is key information in judging its
436 applicability to different users.

437 Although there is no strict definition on what makes a device a LCS, we often make the categorization based on
438 the hardware used. Standard reference measurement instruments are generally based on well-characterized

439 techniques developed and improved over years, based primarily on the progressive refinement of hardware (e.g.
440 materials used for the detection elements, electronic circuits to filter noise, refinement of production methods,
441 etc.). Although LCS sensor technologies are improving, it is interesting that many of the significant improvements
442 that have been made to LCS performance have been through software, rather than hardware advances. As more
443 colocation data is generated in different environments, many LCS manufacturers have been able to develop data
444 correction algorithms that minimise the scale of the errors that are present on the LCS hardware. This can greatly
445 improve the performance of LCS devices, and has been a large factor in the improvements seen in these devices
446 over recent years. These algorithms are, however, inevitably imperfect and can suffer from concept drift (De Vito
447 et al., 2020), caused by the lack of available colocation data over a full spectrum of atmospheric complexity.
448 Furthermore, any kind of statistical model introduces a new error source that can work in conjunction with the
449 pre-existing measurement errors to drastically change the observed error characteristics, making it much more
450 difficult for users to interpret and extrapolate from colocation study performance to intended application. If end
451 users are to be able to make well informed decisions about device applicability to a particular task, then an
452 argument can be made for information on the scale of the error corrections made to a reported measurement to be
453 made available, ideally alongside and a demonstration of its benefits in a relevant environment. Unfortunately,
454 this colocation data is costly to generate, meaning relevant data often does not exist, and when it does is often not
455 communicated in such a way that enables the user to make a fully informed decision.

456 **6. Conclusions**

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

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

Atmospheric
Measurement
Techniques
Discussions

Open Access
EGU

476 All measurement devices suffer from errors, which result in deviations between the reported and true values.
477 These errors can come from a multitude of sources, with the scale of the deviation from the true value being
478 dependent on the nature of the error. Although a known measurement uncertainty for all applications would be
479 ideal for end users to be able to assess measurement device suitability for purpose, in many cases, especially for
480 LCS, this is not possible due to the presence of poorly characterised, or sometimes unknown, error sources. In the
481 absence of this, useful information on likely measurement performance can be obtained using colocation data
482 compared with a measurement with a quantified uncertainty. It is important that such a colocation study is carried
483 out in an environment as similar as possible to the application environment, as the unknown nature of many error
484 sources means their magnitude can change significantly between different locations and/or seasons (e.g. Fig. 6).
485 Ideally, depending on the measurement task, the user could use the colocation data to model the error causes and
486 use this to develop strategies to minimise final measurement uncertainty. Unfortunately, relevant colocation study
487 data is often not available, and to generate the data would be prohibitively costly, which limits the user's ability
488 to make a realistic assessment of likely uncertainties. The presence of, often complex, error minimisation post
489 processing or calibration algorithms further complicates things. This additional uncertainty is most likely to bias
490 any performance prediction if the end user is unaware of the purpose or scale of the data corrections, and their
491 applicability to the target environmental conditions. Ideally, long term colocation data sets demonstrating the
492 performance of measurement hardware and software, in a range of relevant locations, over multiple seasons, and
493 carried out by impartial bodies would be available to inform measurement solution decisions.

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

502 **Supplementary**

503 The supplement related to this article is available online.

504 **Code and data availability**

505 The QUANT intercomparison study data will be made publically available once the study has finished, in early
506 2023. In the meantime, data can be obtained upon request from the corresponding authors.

507 **Author contributions**

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

511 **Competing interests**

Atmospheric
Measurement
Techniques

Discussions

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

525

526    **References**

527    Altman, D. G., & Bland, J. M. (1983). Measurement in Medicine: The Analysis of Method Comparison Studies.

528          *Journal of the Royal Statistical Society. Series D (The Statistician)*, *32*(3), 307–317.

529          https://doi.org/10.2307/2987937

530    Andrewes, P., Bullock, S., Turnbull, R., & Coolbear, T. (2021). Chemical instrumental analysis versus human

531          evaluation to measure sensory properties of dairy products: What is fit for purpose? *International Dairy*

532          *Journal*, *121*, 105098. https://doi.org/10.1016/j.idairyj.2021.105098

533    Bagkis, E., Kassandros, T., Karteris, M., Karteris, A., & Karatzas, K. (2021). Analyzing and Improving the

534          Performance of a Particulate Matter Low Cost Air Quality Monitoring Device. *Atmosphere*, *12*(2), 251.

535          https://doi.org/10.3390/atmos12020251

536    Baldauf, R., Watkins, N., Heist, D., Bailey, C., Rowley, P., & Shores, R. (2009). Near-road air quality

537          monitoring: Factors affecting network design and interpretation of data. *Air Quality, Atmosphere &*

538          *Health*, *2*(1), 1–9. https://doi.org/10.1007/s11869-009-0028-0

539    Bigi, A., Mueller, M., Grange, S. K., Ghermandi, G., & Hueglin, C. (2018). Performance of NO,

540          $NO_2$ low cost sensors and three calibration approaches within a real world

541          application. *Atmospheric Measurement Techniques*, *11*(6), 3717–3735. https://doi.org/10.5194/amt-11-

542          3717-2018

543   Broday, D. M. (2017). Wireless Distributed Environmental Sensor Networks for Air Pollution Measurement—

544          The Promise and the Current Reality. *Sensors (Basel, Switzerland)*, *17*(10), 2263.

545          https://doi.org/10.3390/s17102263

546   Brown, R. J. C., Hood, D., & Brown, A. S. (2008). On the Optimum Sampling Time for the Measurement of

547          Pollutants in Ambient Air. *Journal of Automated Methods and Management in Chemistry*, *2008*,

548          814715. https://doi.org/10.1155/2008/814715

549   Castell, N., Dauge, F. R., Schneider, P., Vogt, M., Lerner, U., Fishbain, B., Broday, D., & Bartonova, A. (2017).

550          Can commercial low-cost sensor platforms contribute to air quality monitoring and exposure estimates?

551          *Environment International*, *99*, 293–302. https://doi.org/10.1016/j.envint.2016.12.007

552   CEN/TS 17660-1: 2021, Air quality – Performance evaluation of air quality sensor systems – Part 1 Gaseous

553          pollutants in ambient air.

554   Chai, T., & Draxler, R. R. (2014). Root mean square error (RMSE) or mean absolute error (MAE)? –

555          Arguments against avoiding RMSE in the literature. *Geoscientific Model Development*, *7*(3), 1247–

556          1250. https://doi.org/10.5194/gmd-7-1247-2014

557   Cordero, J. M., Borge, R., & Narros, A. (2018). Using statistical methods to carry out in field calibrations of low

558          cost air quality sensors. *Sensors and Actuators B: Chemical*, *267*, 245–254.

559          https://doi.org/10.1016/j.snb.2018.04.021

560   Cross, E. S., Williams, L. R., Lewis, D. K., Magoon, G. R., Onasch, T. B., Kaminsky, M. L., Worsnop, D. R., &

561          Jayne, J. T. (2017). Use of electrochemical sensors for measurement of air pollution: Correcting

562          interference response and validating measurements. *Atmospheric Measurement Techniques*, *10*(9),

563          3575–3588. https://doi.org/10.5194/amt-10-3575-2017

564   De Vito, S., Esposito, E., Castell, N., Schneider, P., & Bartonova, A. (2020). On the robustness of field

565          calibration for smart air quality monitors. *Sensors and Actuators B: Chemical*, *310*, 127869.

566          https://doi.org/10.1016/j.snb.2020.127869

567   Doğan, N. Ö. (2018). Bland-Altman analysis: A paradigm to understand correlation and agreement. *Turkish

568          Journal of Emergency Medicine*, *18*(4), 139–141. https://doi.org/10.1016/j.tjem.2018.09.001

569   Duvall, R. M., Long, R. W., Beaver, M. R., Kronmiller, K. G., Wheeler, M. L., & Szykman, J. J. (2016).

570          Performance Evaluation and Community Application of Low-Cost Sensors for Ozone and Nitrogen

571          Dioxide. *Sensors*, *16*(10), 1698. https://doi.org/10.3390/s16101698

572   Feinberg, S. N., Williams, R., Hagler, G., Low, J., Smith, L., Brown, R., Garver, D., Davis, M., Morton, M.,

573    Schaefer, J., & Campbell, J. (2019). Examining spatiotemporal variability of urban particulate matter

574        and application of high-time resolution data from a network of low-cost air pollution sensors.

575        *Atmospheric Environment*, *213*, 579–584. https://doi.org/10.1016/j.atmosenv.2019.06.026

576    GDE. (2010). Guidance for the Demonstration of Equivalence of Ambient Air Monitoring Methods. Report by

577        an EC Working Group.

578    Gerboles, M., Lagler, F., Rembges, D., & Brun, C. (2003). Assessment of uncertainty of NO2 measurements by

579        the chemiluminescence method and discussion of the quality objective of the NO2 European Directive.

580        *Journal of Environmental Monitoring*, *5*(4), 529–540. https://doi.org/10.1039/B302358C

581    Gramsch, E., Oyola, P., Reyes, F., Vásquez, Y., Rubio, M. A., Soto, C., Pérez, P., Moreno, F., & Gutiérrez, N.

582        (2021). Influence of Particle Composition and Size on the Accuracy of Low Cost PM Sensors:

583        Findings From Field Campaigns. *Frontiers in Environmental Science*, *9*.

584        https://www.frontiersin.org/article/10.3389/fenvs.2021.751267

585    Grégis, F. (2019). On the meaning of measurement uncertainty. *Measurement*, *133*, 41–46.

586        https://doi.org/10.1016/j.measurement.2018.09.073

587    Hofman, J., Nikolaou, M., Shantharam, S. P., Stroobants, C., Weijs, S., & La Manna, V. P. (2022). Distant

588        calibration of low-cost PM and NO2 sensors; evidence from multiple sensor testbeds. *Atmospheric

589        Pollution Research*, *13*(1), 101246. https://doi.org/10.1016/j.apr.2021.101246

590    JCGM. (2012). International vocabulary of metrology - Basic and general concepts and associated terms, Paris:

591        BIPM.

592    Karagulian, F., Barbiere, M., Kotsev, A., Spinelle, L., Gerboles, M., Lagler, F., Redon, N., Crunaire, S., &

593        Borowiak, A. (2019). Review of the Performance of Low-Cost Sensors for Air Quality Monitoring.

594        *Atmosphere*, *10*(9), 506. https://doi.org/10.3390/atmos10090506

595    Kirkham, H., Riepnieks, A., Albu, M., & Laverty, D. (2018). The nature of measurement, and the true value of a

596        measured quantity. *2018 IEEE International Instrumentation and Measurement Technology

597        Conference (I2MTC)*, 1–6. https://doi.org/10.1109/I2MTC.2018.8409771

598    Lewis, A., & Edwards, P. (2016). Validate personal air-pollution sensors. *Nature News*, *535*(7610), 29.

599        https://doi.org/10.1038/535029a

600    Makar, P. A., Stroud, C., Akingunola, A., Zhang, J., Ren, S., Cheung, P., & Zheng, Q. (2021). Vehicle-induced

601        turbulence and atmospheric pollution. *Atmospheric Chemistry and Physics*, *21*(16), 12291–12316.

602        https://doi.org/10.5194/acp-21-12291-2021

603   Malings, C., Tanzer, R., Hauryliuk, A., Kumar, S. P. N., Zimmerman, N., Kara, L. B., Presto, A. A., & R.

604       Subramanian. (2019). Development of a general calibration model and long-term performance

605       evaluation of low-cost sensors for air pollutant gas monitoring. *Atmospheric Measurement Techniques*,

606       *12*(2), 903–920. https://doi.org/10.5194/amt-12-903-2019

607   Mari, L., Wilson, M., & Maul, A. (2021). *Measurement across the Sciences: Developing a Shared Concept*

608       *System for Measurement*. Springer International Publishing. https://doi.org/10.1007/978-3-030-65558-7

609   Masson, N., Piedrahita, R., & Hannigan, M. (2015). Approach for quantification of metal oxide type

610       semiconductor gas sensors used for ambient air quality monitoring. *Sensors and Actuators B:*

611       *Chemical*, *208*, 339–345. https://doi.org/10.1016/j.snb.2014.11.032

612   Matejka, J., & Fitzmaurice, G. (2017). Same Stats, Different Graphs: Generating Datasets with Varied

613       Appearance and Identical Statistics through Simulated Annealing. *Proceedings of the 2017 CHI*

614       *Conference on Human Factors in Computing Systems*, 1290–1294.

615       https://doi.org/10.1145/3025453.3025912

616   Mead, M. I., Popoola, O. A. M., Stewart, G. B., Landshoff, P., Calleja, M., Hayes, M., Baldovi, J. J., McLeod,

617       M. W., Hodgson, T. F., Dicks, J., Lewis, A., Cohen, J., Baron, R., Saffell, J. R., & Jones, R. L. (2013).

618       The use of electrochemical sensors for monitoring urban air quality in low-cost, high-density networks.

619       *Atmospheric Environment*, *70*, 186–203. https://doi.org/10.1016/j.atmosenv.2012.11.060

620   Molina, M. J., & Molina, L. T. (2004). Megacities and Atmospheric Pollution. *Journal of the Air & Waste*

621       *Management Association*, *54*(6), 644–680. https://doi.org/10.1080/10473289.2004.10470936

622   Morawska, L., Thai, P., Liu, X., Asumadu-Sakyi, A., Ayoko, G., Bartonova, A., Bedini, A., Chai, F.,

623       Christensen, B., Dunbabin, M., Gao, J., Hagler, G., Jayaratne, R., Kumar, P., Lau, A., Louie, P.,

624       Mazaheri, M., Ning, Z., Motta, N., … Williams, R. (2018). Applications of low-cost sensing

625       technologies for air quality monitoring and exposure assessment: How far have they gone?

626       *Environment International*, *116*, 286–299. https://doi.org/10.1016/j.envint.2018.04.018

627   Mueller, M., Meyer, J., & Hueglin, C. (2017). Design of an ozone and nitrogen dioxide sensor unit and its long-

628       term operation within a sensor network in the city of Zurich. *Atmospheric Measurement Techniques*,

629       *10*(10), 3783–3799. https://doi.org/10.5194/amt-10-3783-2017

630   Peters, D. R., Popoola, O. A. M., Jones, R. L., Martin, N. A., Mills, J., Fonseca, E. R., Stidworthy, A., Forsyth,

631       E., Carruthers, D., Dupuy-Todd, M., Douglas, F., Moore, K., Shah, R. U., Padilla, L. E., & Alvarez, R.

632       A. (2021). *Evaluating uncertainty in sensor networks for urban air pollution insights* [Preprint].

633    Gases/In Situ Measurement/Validation and Intercomparisons. https://doi.org/10.5194/amt-2021-210

634    Popoola, O. A. M., Stewart, G. B., Mead, M. I., & Jones, R. L. (2016). Development of a baseline-temperature

635    correction methodology for electrochemical sensors and its implications for long-term stability.

636    *Atmospheric Environment*, *147*, 330–343. https://doi.org/10.1016/j.atmosenv.2016.10.024

637    Rai, A. C., Kumar, P., Pilla, F., Skouloudis, A. N., Di Sabatino, S., Ratti, C., Yasar, A., & Rickerby, D. (2017).

638    End-user perspective of low-cost sensors for outdoor air pollution monitoring. *Science of The Total*

639    *Environment*, *607–608*, 691–705. https://doi.org/10.1016/j.scitotenv.2017.06.266

640    Spinelle, L., Gerboles, M., Villani, M. G., Aleixandre, M., & Bonavitacola, F. (2015). Field calibration of a

641    cluster of low-cost available sensors for air quality monitoring. Part A: Ozone and nitrogen dioxide.

642    *Sensors and Actuators B: Chemical*, *215*, 249–257. https://doi.org/10.1016/j.snb.2015.03.031

643    Sun, L., Westerdahl, D., & Ning, Z. (2017). Development and Evaluation of A Novel and Cost-Effective

644    Approach for Low-Cost $NO_2$ Sensor Drift Correction. *Sensors (Basel, Switzerland)*, *17*(8), E1916.

645    https://doi.org/10.3390/s17081916

646    Thompson, M., & Ellison, S. L. R. (2005). A review of interference effects and their correction in chemical

647    analysis with special reference to uncertainty. *Accreditation and Quality Assurance*, *10*(3), 82–97.

648    https://doi.org/10.1007/s00769-004-0871-5

649    Tian, Y., Nearing, G. S., Peters-Lidard, C. D., Harrison, K. W., & Tang, L. (2016). Performance Metrics, Error

650    Modeling, and Uncertainty Quantification. *Monthly Weather Review*, *144*(2), 607–613.

651    https://doi.org/10.1175/MWR-D-15-0087.1

652    van Zoest, V., Osei, F. B., Stein, A., & Hoek, G. (2019). Calibration of low-cost NO2 sensors in an urban air

653    quality network. *Atmospheric Environment*, *210*, 66–75.

654    https://doi.org/10.1016/j.atmosenv.2019.04.048

655    Williams, D. E. (2020). Electrochemical sensors for environmental gas analysis. *Current Opinion in*

656    *Electrochemistry*, *22*, 145–153. https://doi.org/10.1016/j.coelec.2020.06.006

657    World Health Organization. (2021). *WHO global air quality guidelines: Particulate matter (PM2.5 and PM10),*

658    *ozone, nitrogen dioxide, sulfur dioxide and carbon monoxide: executive summary*. World Health

659    Organization. https://apps.who.int/iris/handle/10665/345334

660    Zucco, M., Curci, S., Castrofino, G., & Sassi, M. P. (2003). A comprehensive analysis of the uncertainty of a

661    commercial ozone photometer. *Measurement Science and Technology*, *14*(9), 1683–1689.

662    https://doi.org/10.1088/0957-0233/14/9/320