# Artificial intelligence based software facilitates spirometry quality control in asthma and COPD clinical trials

Eva Topole[1], Sonia Biondaro[1], Isabella Montagna[1], Sandrine Corre[1], Massimo Corradi[2], Sanja Stanojevic[3], Brian Graham[4], Nilakash Das [ORCID][5,6], Kevin Ray[6] and Marko Topalovic[6]

[1]Global Clinical Development, Chiesi Farmaceutici, S.p.A., Parma, Italy. [2]Department of Medicine and Surgery, University of Parma, Parma, Italy. [3]Department of Community Health and Epidemiology, Dalhousie University, Halifax, NS, Canada. [4]Division of Respirology, Critical Care and Sleep Medicine, University of Saskatchewan, Saskatoon, SK, Canada. [5]Laboratory of Respiratory Diseases and Thoracic Surgery, Department of Chronic Diseases Metabolism and Ageing, KU Leuven, Leuven, Belgium. [6]ArtiQ NV, Leuven, Belgium.

Corresponding author: Marko Topalovic (marko.topalovic@artiq.eu)

## Abstract

*Rationale* Acquiring high-quality spirometry data in clinical trials is important, particularly when using forced expiratory volume in 1 s or forced vital capacity as primary end-points. In addition to quantitative criteria, the American Thoracic Society (ATS)/European Respiratory Society (ERS) standards include subjective evaluation which introduces inter-rater variability and potential mistakes. We explored the value of artificial intelligence (AI)-based software (ArtiQ.QC) to assess spirometry quality and compared it to traditional over-reading control.
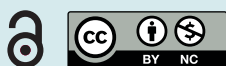
*Methods* A random sample of 2000 sessions (8258 curves) was selected from Chiesi COPD and asthma trials (n=1000 per disease). Acceptability using the 2005 ATS/ERS standards was determined by over-reader review and by ArtiQ.QC. Additionally, three respiratory physicians jointly reviewed a subset of curves (n=150).

*Results* The majority of curves (n=7267, 88%) were of good quality. The AI agreed with over-readers in 91% of cases, with 97% sensitivity and 93% positive predictive value. Performance was significantly better in the asthma group. In the revised subset, n=50 curves were repeated to assess intra-rater reliability ($\kappa=0.83$, 0.86 and 0.80 for each of the three reviewers). All reviewers agreed on 63% of 100 unique tests ($\kappa=0.5$). When reviewers set the consensus (gold standard), individual agreement with it was 88%, 94% and 70%. The agreement between AI and "gold-standard" was 73%; over-reader agreement was 46%.

*Conclusion* AI-based software can be used to measure spirometry data quality with comparable accuracy as experts. The assessment is a subjective exercise, with intra- and inter-rater variability even when the criteria are defined very precisely and objectively. By providing consistent results and immediate feedback to the sites, AI may benefit clinical trial conduct and variability reduction.

## Introduction

Measures of lung function, typically using spirometry, are critical clinical outcomes used in respiratory clinical trials [1, 2]. The efficacy of treatments targeting the lungs is often assessed through two main spirometry derived parameters: the forced expiratory volume in 1 s ($FEV_1$) and the forced vital capacity (FVC) [3–5]. Since both of these indices are dependent on a maximal forced expiratory effort, spirometry requires skilled technicians to ensure the subjects' compliance and optimal performance of the test. Therefore, in the context of clinical trials, ensuring and assessing the quality of spirometry is of paramount importance to guarantee that the measured $FEV_1$ and FVC values are sufficiently reliable to assess the efficacy of a therapeutic.

To ensure valid data and reduce variability in the clinical trials, it is normal to use central spirometry that incorporates standardised instrumentation, site training and centralised review (over-reading) of all spirometry sessions [6, 7]. The over-reading is performed by independent reviewers who are registered respiratory experts trained to review spirometry in a clinical trial environment where the objective is to ensure consistent results across all patients, sites and time points. These specialists typically use the criteria defined by an international task force appointed by the American Thoracic Society (ATS) and European Respiratory Society (ERS) to evaluate spirometry data quality, with an average turnaround time of 24–48 h for feedback to the investigational sites [8, 9]. This delay between subjects performing the spirometry session and the investigational site receiving feedback can delay study-specific decisions, such as subject randomisation and treatment changes, and may lead to repeat visits with additional spirometry sessions to be performed by the subjects, increasing the burden on participants.

The ATS/ERS spirometry quality standards define both quantitative and qualitative criteria to ensure high quality standards for spirometry measurements. Quantitative criteria (*e.g.* expiration duration longer than a minimum time period, or back-extrapolated volume below a certain threshold) can be evaluated relatively easily by software. However, the qualitative components (*e.g.* artefact detection) require subjective assessment by the technicians acquiring the data [10–12]. Technician training, experience and other factors can lead to heterogeneity of the assessment [13]. Agreement between two reviewers assessing the same data has been measured as low as 52% in previous studies [14–17]. It is important to reduce this inter-rater variability, especially in a clinical trial setting where spirometry data are being used for measuring the therapeutic effect of a pharmaceutical intervention.

In this investigation, we explored whether an artificial intelligence (AI)-based software (ArtiQ.QC) can perform spirometry data quality assessment with an accuracy at least as high as manual over-reading, bringing additional benefits of a shorter turnaround time and improved consistency.

## Materials and methods
### Study objectives
The primary objective of the study was to determine whether an AI-based software (ArtiQ.QC) can be used to perform spirometry data quality assessment with similar accuracy as expert over-readers. Secondary objectives were to assess the level of inter-rater variability and intra-rater reliability in spirometry data quality assessment; determine to what extent a single reviewer's judgment can be considered as a gold-standard comparator; and determine whether AI-based software (ArtiQ.QC) may outperform the over-reader's quality assessment.

### Data
Spirometry data from past clinical trials sponsored by Chiesi in the therapeutic areas of asthma and COPD were used for analysis [18–20]. These studies were registered with ClinicalTrials.gov (identifier numbers NCT02579850, NCT02676076 and NCT02676089). Using a function for random selection of data, a completely randomly chosen sample of 2000 sessions was selected for analysis, with 1000 sessions each from asthma and COPD subjects. A total of 8258 spirometry curves (n=4085 COPD, n=4173 asthma) were used. All data selected for analysis were collected and assessed using the ATS/ERS 2005 spirometry standards [9]. Each curve was labelled by an expert over-reader at the time of data collection according to the over-reader guidelines for each clinical trial, with possible outcomes of "acceptable" or "unacceptable". Over-reader guidelines were aligned with the ATS/ERS 2005 standards, and stated that if a curve was "usable" according to the ATS/ERS 2005 standards (*i.e.* have an acceptable start of test and are free from artefacts, such as a cough) then the over-reading label should be acceptable. The acceptability evaluation of the whole session is not reported within this article.

### Study design
In this noninterventional retrospective study, the AI-based software (ArtiQ.QC; ArtiQ NV, Leuven, Belgium) was used to assess the quality of spirometry data using the ATS/ERS 2005 standards [9]. ArtiQ. QC is AI-based software which can apply the quantitative and qualitative criteria defined in the ATS/ERS 2005 standards to assess spirometry data quality. The qualitative criteria (*i.e.* artefact detection) are assessed using a previously published deep-learning based approach [21]. ArtiQ.QC is a fully validated software, and in this study was used by Chiesi for testing in a clinical trial context.

An overview of the study design is shown in figure 1. Firstly, ArtiQ.QC was used to assess data quality in all 8258 curves, and the performance of ArtiQ.QC was measured (accuracy, sensitivity, specificity, positive predictive value (PPV) and negative predictive value (NPV)) using the original over-reading labels as comparators. Secondly, ArtiQ.QC was recalibrated (without altering the AI) on 80% of the data selected
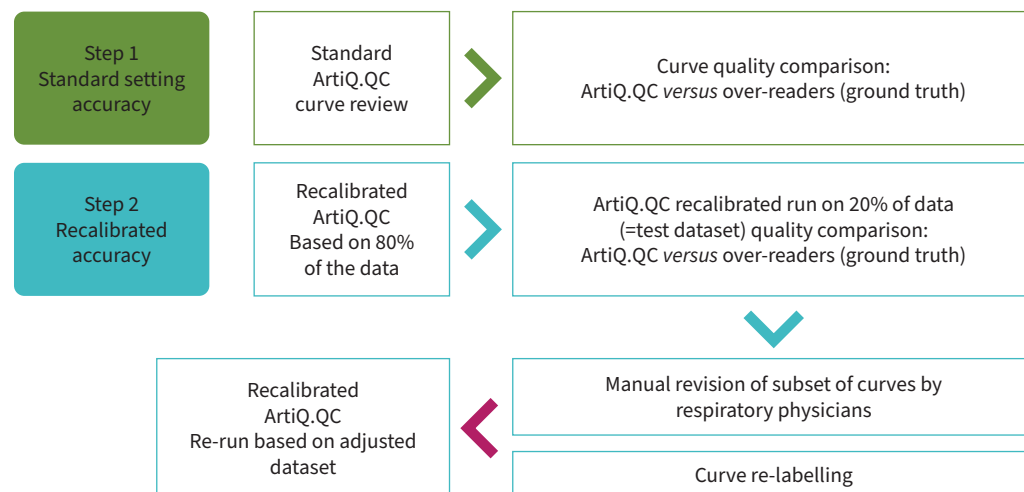
**FIGURE 1** Study design and flow.

for analysis (n=1600 sessions, equally split between asthma and COPD) and the performance of the recalibrated version determined on the remaining 20% of the data. The original over-reading labels were again used as comparators in this second performance evaluation to evaluate the extent to which ArtiQ.QC performance could be improved.

Subsequently, a selection of curves was manually assessed by three reviewers to determine gold-standard labels (n=150 curves total, with n=100 unique curves, equally split between diseases; full breakdown of curves selected for review is provided in the supplementary material). Initially, all curves were assessed separately by each reviewer. All curves where there was disagreement between the reviewers' applied labels were jointly assessed to eventually agree on a gold-standard label for each curve. Reviewers were blinded to the original over-reader and ArtiQ.QC labels, and used the criteria defined in the ATS/ERS 2005 standards to assess the spirometry data quality. Inter-rater variability and intra-rater reliability were measured based on reviewer labels. Finally, ArtiQ.QC and original over-reader labels were both compared to the gold-standard labels provided by the joint consensus of reviewers.

### Ensuring a fair comparison of ArtiQ.QC to over-readers

Individual spirometry curves in the dataset used for analysis were labelled as either acceptable or unacceptable by the over-readers. ArtiQ.QC generates labels for individual spirometry curves in accordance with the ATS/ERS 2005 standards, which permits "curves that have an acceptable start of test and are free from artefact, such as a cough" to be "usable curves" for the measurement of $FEV_1$ and FVC. For the purpose of comparing ArtiQ.QC and over-reading labels, an ArtiQ.QC label of "usable" was converted to "acceptable", consistent with the rules followed by over-readers in the data source studies.

### Analysis

A comparison of proportions test was used to evaluate whether significant differences in the baseline dataset characteristics existed. For performance evaluation of ArtiQ.QC, a 2×2 confusion matrix of labels assigned by ArtiQ.QC or original over-reading at a curve level was created. Accuracy, sensitivity, specificity, PPV and NPV were measured from the confusion matrix. McNemar's test was used to test whether recalibration significantly improved performance of ArtiQ.QC. Cohen's κ was used to assess intra-rater reliability, and Fleiss's κ (the multi-rater generalisation of Cohen's κ) was used to assess inter-rater variability from the manual review. Accuracy, sensitivity, specificity, PPV and NPV were calculated for ArtiQ.QC and original over-reading labels using the expert reviewer labels as comparator. Performance of ArtiQ.QC and original over-reading was compared using McNemar's test.

### Results

### Dataset characteristics

The vast majority of curves (n=7267, 88%) were labelled as acceptable by the original over-reader assessment. More curves were labelled acceptable in the asthma subgroup (3818 out of 4173, 91.5%) compared to the COPD subgroup (3391 out of 4085, 83%; p<0.0001 comparison of proportions test).

**TABLE 1** Performance evaluation of baseline ArtiQ.QC using original over-reader labels as the comparator in the full dataset

|  | Overall | Asthma subgroup | COPD subgroup |
|---|---|---|---|
| Curves n | 8258 | 4173 | 4085 |
| Accuracy | 87 | 90 | 84 |
| Sensitivity | 93 | 93 | 94 |
| Specificity | 35 | 52 | 25 |
| PPV | 92 | 96 | 87 |
| NPV | 41 | 40 | 43 |

Data are presented as %, unless otherwise stated. PPV: positive predictive value; NPV: negative predictive value.

### Performance of AI software using original over-reader labels as comparator

A summary of the performance of the AI-based software using original over-reader labels as the comparator is shown in table 1. From all n=8258 curves, the AI software agreed with over-readers in 87% of cases, with 93% sensitivity and 92% PPV. Performance was significantly better in the asthma subgroup (accuracy 90%, n=4173 curves) than the COPD subgroup (accuracy 84%, n=4085 curves; $p < 0.0001$ comparison of proportions test).

### Performance of recalibrated AI software using original over-reader labels as comparator

Tables 2 and 3 summarise the performance evaluation results of the recalibrated and baseline AI software, respectively, using original over-reader labels as a comparator. Following recalibration of the AI software in 80% of the full dataset, 67 curves in the 20% test set changed label compared to the quality label assigned with the standard AI software. This resulted in a significant improvement in performance in the test dataset ($p < 0.0001$, McNemar's test). The parameters adjusted in the recalibration were the artefact probability and hesitation volume thresholds (full details are presented in the supplementary material). Overall recalibrated accuracy was 91%, with sensitivity 97% and PPV 93% (1659 curves), compared to baseline accuracy of 87%, sensitivity 92% and PPV 93%. The greatest improvement in performance was in the asthma subgroup (increase in accuracy from 89% to 94% in 841 curves following recalibration; $p < 0.001$ McNemar's test). In addition, a nonsignificant improvement in COPD performance was seen (accuracy improved from 84% to 86% in 818 curves; $p = 0.13$ McNemar's test).

### Manual review: reviewer reliability and variability

A set of 150 curves were additionally manually reviewed by three reviewers. 50 of these curves were repeated in a random fashion in order to assess intra-rater reliability. Cohen's $\kappa$ values for the three reviewers were 0.83, 0.86 and 0.80, indicating that although reviewers mainly provided consistent labels when presented with the same data twice, some disagreement with their own evaluation was present.

From the 100 unique curves, all reviewers agreed after their separate review sessions in 63% of curves (inter-rater Fleiss's $\kappa$=0.5), indicating disagreement between reviewers of a magnitude consistent with previous studies [14–17]. Agreement between reviewers was higher in the subset of curves where the AI software and original over-reading provided the same label (18 out of 25, 75%; Fleiss's $\kappa$=0.63) than in the subset where the AI software and original over-reader labels differed (45 out of 75, 60%; Fleiss's

**TABLE 2** Performance evaluation of baseline ArtiQ.QC using original over-reader labels as the comparator in the 20% test dataset

|  | Overall | Asthma subgroup | COPD subgroup |
|---|---|---|---|
| Curves n | 1659 | 841 | 818 |
| Accuracy | 87 | 89 | 84 |
| Sensitivity | 92 | 92 | 94 |
| Specificity | 24 | 43 | 15 |
| PPV | 93 | 96 | 87 |
| NPV | 30 | 29 | 32 |

Data are presented as %, unless otherwise stated. PPV: positive predictive value; NPV: negative predictive value.

**TABLE 3** Performance evaluation of recalibrated ArtiQ.QC using original over-reader labels as the comparator in the 20% test dataset

|  | Overall | Asthma subgroup | COPD subgroup |
|---|---|---|---|
| Curves n | 1659 | 841 | 818 |
| Accuracy | 91 | 94 | 86 |
| Sensitivity | 97 | 98 | 98 |
| Specificity | 22 | 40 | 14 |
| PPV | 93 | 96 | 87 |
| NPV | 58 | 59 | 57 |

Data are presented as %, unless otherwise stated. PPV: positive predictive value; NPV: negative predictive value.

κ=0.45). Despite using a set of objective criteria defined by the ATS/ERS 2005 standards, this result indicates that inter-rater disagreement exists between expert reviewers.

### Performance of recalibrated AI software, original over-readers and each reviewer using gold-standard labels as comparator

The 37 curves where disagreement existed between reviewers were jointly reviewed to derive labels for all 100 curves that all reviewers agreed on (gold-standard labels). The percentage agreement with gold-standard labels and Cohen's κ values for each reviewer, the reviewers' majority opinion, AI-based software and the original over-reader labels are shown in table 4. Note that the lowest possible value for percentage agreement with gold-standard labels for any individual reviewer is 63%, since 63 out of 100 curves had the same label from the separate reviews. A percentage agreement of 89% for the majority opinion of reviewers from their separate review sessions indicates that the gold-standard label was opposite to the reviewer majority opinion in 11 curves. Unsurprisingly, the agreement of individual reviewers with the gold-standard labels is high, since it is the reviewers' joint opinion that defines the gold standard. The agreement between AI-based software and gold-standard label (73%) is similar to reviewer 3 (70%). This agreement is far superior to the original over-reader labels (46%) in this subset of curves, although the high inter-rater variability in this subset may not generalise to the whole dataset.

### Discussion

This study demonstrates that AI-based software can be used to measure spirometry data quality with similar accuracy to expert over-readers. Additionally, this analysis confirmed that assessment of spirometry quality is a subjective exercise, with intra- and inter-rater variability even when the criteria are defined very precisely and objectively. When using the joint opinion of three expert reviewers as a gold standard, AI-based software has a much higher agreement with the gold-standard labels than original over-reader labels. Additionally, by providing consistent results and immediate feedback to the sites, using ArtiQ.QC may benefit clinical trial conduct and reduce the variability in outcomes, although these benefits have not been considered in this retrospective study.

Assessing the performance of AI software using the original over-reader labels as a comparator assumes that the original over-reader labels are 100% accurate. However, the existence of inter- and intra-rater variability observed in this and other studies [1–5] suggests that this assumption is unlikely to be correct.

**TABLE 4** Percentage agreement and Cohen's κ value for agreement with the gold-standard label, which was defined using the joint opinion of the three reviewers, from 100 spirometry curves

|  | Agreement with gold standard | Cohen's κ |
|---|---|---|
| Reviewer 1 | 88% | 0.73 |
| Reviewer 2 | 94% | 0.85 |
| Reviewer 3 | 70% | 0.42 |
| Majority opinion of reviewers | 89% | 0.75 |
| AI software | 73% | 0.47 |
| Original over-reader | 46% | −0.08 |

AI: artificial intelligence.

Even when comparing three very experienced reviewers, and the exact criteria from the 2005 ATS/ERS standards were used to derive quality scores for each curve, disagreement between reviewers was still present in over a third of all 100 unique curves reviewed in this study.

The inter- and intra-rater variability present when spirometry data are reviewed manually means that a single reviewer's opinion is not a reliable comparator for assessing the accuracy or reliability of an automated tool. While the inter-rater agreement measured in this study was consistent with previous reports, this agreement may be influenced by experience, training and expertise of the reviewers [22–24]. However, this is reflective of the real-world scenario, where usually only a single reviewer assesses spirometry data quality. The joint opinion of the three reviewers may be considered a viable gold standard, but this approach is not realistic for implementation in clinical trials. Therefore, for AI-based software to be useful its performance simply needs to match that of another manual reviewer. With percentage agreement of 73% with the gold-standard labels in the subset of curves manually reviewed, the AI software performance is similar to reviewer 3 in this study. Compared to the original over-reader labels (percentage agreement with the gold standard of 46%), it may be the case that AI software such as ArtiQ. QC could outperform over-readers in assessing the quality of spirometry data in clinical trials.

### Limitations

This study was conducted using legacy spirometry data acquired during the two clinical trials (one in asthma and one in COPD). The over-reading services for the two studies were provided by two different vendors, potentially causing a difference in the level of agreement between AI algorithm and over-reading for each study. Additionally, the algorithm was initially developed and trained on a diverse set of curves (from healthy subjects, different respiratory and other internal diseases), potentially leading to a lower accuracy than with the algorithm that would have been developed and trained on a set of curves within specific therapeutic area. Further work is required to explore the generalisability of these results to other disease areas, paediatric population and non-clinical trial settings [25, 26].

This study compares the performance value of the software that is trained to mimic experts and experts themselves, once the complete spirometry session is performed. However, the largest practical impact on the running clinical trials will be if the software were used after each spirometry blow. This would immediately reject initial bad blows and help subjects and technicians to secure the acceptable quality of spirometry by the end of the session.

This study was conducted using data acquired using the 2005 ATS/ERS standards, and with over-reading performed according to those standards. New ATS/ERS standards were published in 2019, which offer more guidance in quality assessment and change the way in which spirometry data must be acquired explicitly recommending a maximal inspiration following forced expiration [8]. This study serves as proof of concept that the AI methodology can be tailored to follow the 2005 spirometry standards to assess spirometry acceptability. Although the AI algorithms will be updated to comply with the revisions in the 2019 spirometry standards, the proven underlying AI methodology remains unchanged. With the new standards and AI software working in parallel, inter- and intra-rater variability in spirometry data quality assessment may decrease; this is the subject of ongoing work.

### Conclusion

With this retrospective analysis on a subset of asthma and COPD spirometry data acquired according to the 2005 ATS/ERS standards, we confirmed that AI software can be used with high accuracy to evaluate the quality of spirometry data in clinical trials. AI software such as ArtiQ.QC may offer an alternative approach, or assistance, to manual over-reading with the main advantage that AI software can provide immediate feedback allowing a real time evaluation with the subject still at the site. Additional advantages that AI software brings include increased consistency and repeatability of spirometry data quality assessment because the AI model is more objective than manual over-reading. Further evaluation and testing of this technology are needed to understand its practical implications and accuracy according to the most recent 2019 ATS/ERS standards and in disease areas other than asthma and COPD. Prospective real-time testing would also allow to explore and evaluate additional potential benefits of this technology on reducing patient burden and enhancing quality of data in clinical trials. Next to this, such AI algorithms may contribute to increased reliability and enhanced validity of epidemiological studies where spirometry data is collected.

## References

1   European Medicines Agency (EMA). Guideline on Clinical Investigation of Medicinal Products in the Treatment of Chronic Obstructive Pulmonary Disease. Amsterdam, EMA, 2012.

2   European Medicines Agency (EMA). Guideline on the Clinical Investigation of Medicinal Products for the Treatment of Asthma. Amsterdam, EMA, 2015.

3   Busse WW, Morgan WJ, Taggart V, et al. Asthma outcomes workshop: overview. J Allergy Clin Immunol 2012; 129: S1–S8.

4   Glaab T, Vogelmeier C, Buhl R. Outcome measures in chronic obstructive pulmonary disease (COPD): strengths and limitations. Respir Res 2010; 11: 79.

5   Cazzola M, MacNee W, Martinez FJ, et al. Outcomes for COPD pharmacological trials: from lung function to biomarkers. Eur Respir J 2008; 31: 416–469.

6   Pérez-Padilla R, Vázquez-García JC, Márquez MN, et al. Spirometry quality-control strategies in a multinational study of the prevalence of chronic obstructive pulmonary disease. Respir Care 2008; 53: 1019–1026.

7   Malmstrom K, Peszek I, Botto A, et al. Quality assurance of asthma clinical trials. Control Clin Trials 2002; 23: 143–156.

8   Graham BL, Steenbruggen I, Miller MR, et al. Standardization of spirometry 2019 update. An official American Thoracic Society and European Respiratory Society technical statement. Am J Respir Crit Care Med 2019; 200: e70–e88.

9   Miller MR, Hankinson J, Brusasco V, et al. Standardisation of spirometry. Eur Respir J 2005; 26: 319–338.

10  Müller-Brandes C, Krämer U, Gappa M, et al. LUNOKID: can numerical American Thoracic Society/European Respiratory Society quality criteria replace visual inspection of spirometry? Eur Respir J 2014; 43: 1347–1356.

11  Townsend MC. The American Thoracic Society/European Respiratory Society 2019 spirometry statement and occupational spirometry testing in the United States. Am J Respir Crit Care Med 2020; 201: 1010–1011.

12  Beeckman-Wagner L-AF, Freeland D. Spirometry Quality Assurance; Common Errors and Their Impact on Test Results. DHHS (NIOSH) publication 2012-116. Cincinnati, NIOSH, 2012.

13  Borg BM, Hartley MF, Bailey MJ, et al. Adherence to acceptability and repeatability criteria for spirometry in complex lung function laboratories. Respir Care 2012; 57: 2032–2038.

14  Velickovski F, Ceccaroni L, Marti R, et al. Automated spirometry quality assurance: supervised learning from multiple experts. IEEE J Biomed Heal Inform 2018; 22: 276–284.

15  Hankinson JL, Eschenbacher B, Townsend M, et al. Use of forced vital capacity and forced expiratory volume in 1 second quality criteria for determining a valid test. Eur Respir J 2015; 45: 1283–1292.

16  Tan WC, Bourbeau J, O'Donnell D, et al. Quality assurance of spirometry in a population-based study – predictors of good outcome in spirometry testing. COPD 2014; 11: 143–151.

17  Eaton T, Withy S, Garrett JE, et al. Spirometry in primary care practice: the importance of quality assurance and the impact of spirometry workshops. Chest 1999; 116: 416–423.

18  Virchow JC, Kuna P, Paggiaro P, et al. Single inhaler extrafine triple therapy in uncontrolled asthma (TRIMARAN and TRIGGER): two double-blind, parallel-group, randomised, controlled phase 3 trials. Lancet 2019; 394: 1737–1749.

19  Kots M, Georges G, Guasconi A, et al. S26 Effect of single-inhaler extrafine beclometasone/formoterol/glycopyrronium pMDI (BDP/FF/GB) compared with two-inhaler fluticasone furoate/vilanterol DPI+ tiotropium DPI (FLF/VIL+ TIO) triple therapy on health-related quality of life (HRQoL) in patients with COPD: The TRISTAR study. Thorax 2021; 76: A20.

20  Papi A, Vestbo J, Fabbri L, et al. Extrafine inhaled triple therapy versus dual bronchodilator therapy in chronic obstructive pulmonary disease (TRIBUTE): a double-blind, parallel group, randomised controlled trial. Lancet 2018; 391: 1076–1084.

21 Das N, Verstraete K, Stanojevic S, *et al.* Deep-learning algorithm helps to standardise ATS/ERS spirometric acceptability and usability criteria. *Eur Respir J* 2020; 56: 2000603.

22 Enright P, Vollmer WM, Lamprecht B, *et al.* Quality of spirometry tests performed by 9893 adults in 14 countries: the BOLD study. *Respir Med* 2011; 105: 1507–1515.

23 Enright PL, Beck KC, Sherrill DL. Repeatability of spirometry in 18,000 adult patients. *Am J Respir Crit Care Med* 2004; 169: 235–238.

24 Parsons R, Schembri D, Hancock K, *et al.* Effects of the spirometry learning module on the knowledge, confidence, and experience of spirometry operators. *NPJ Prim Care Respir Med* 2019; 29: 30.

25 Bonell C, Oakley A, Hargreaves J, *et al.* Assessment of generalisability in trials of health interventions: suggested framework and systematic review. *BMJ* 2006; 333: 346–349.

26 Futoma J, Simons M, Panch T, *et al.* The myth of generalisability in clinical research and machine learning in health care. *Lancet Digit Heal* 2020; 2: e489–e492.