

Using Virtual-Reality Simulation to Assess Performance in Endobronchial Ultrasound

Lars Konge^a Jouke Annema^d Paul Clementsen^b Valentina Minddal^b
Peter Vilmann^c Charlotte Ringsted^e

^aCentre for Clinical Education, University of Copenhagen and the Capital Region of Denmark, Copenhagen,

^bDepartment of Pulmonology, Gentofte Hospital, University of Copenhagen, Hellerup, and ^cDepartment of Surgical Gastroenterology, Copenhagen University Hospital Herlev, Herlev, Denmark; ^dDepartment of Pulmonology, Leiden University Medical Center, Leiden, and Department of Pulmonology, Academic Medical Center, Amsterdam, The Netherlands; ^eDepartment of Anesthesia, University of Toronto, Toronto, Ont., Canada

Key Words

Endobronchial ultrasound · Transbronchial needle aspiration · Virtual-reality simulator

Abstract

Background: For optimal treatment of patients with non-small cell lung carcinoma, it is essential to have physicians with competence in endobronchial ultrasound-guided transbronchial needle aspiration (EBUS-TBNA). EBUS training and certification requirements are under discussion and the establishment of basic competence should be based on an objective assessment of performance. **Objectives:** The aims of this study were to design an evidence-based and credible EBUS certification based on a virtual-reality (VR) EBUS simulator test. **Methods:** Twenty-two respiratory physicians were divided into 3 groups: experienced EBUS operators (group 1, n = 6), untrained novices (group 2, n = 8) and simulator-trained novices (group 3, n = 8). Each physician performed two standardized simulated EBUS-TBNA procedures. Simulator metrics with discriminatory ability were identified and reliability was explored. Finally, the contrasting-groups method was used to establish a pass/fail standard, and the consequences of this standard were explored.

Results: Successfully sampled lymph nodes and procedure time were the only simulator metrics that showed statistically significant differences of $p = 0.047$ and $p = 0.002$, respectively. The resulting quality score (QS, i.e. sampled lymph nodes per minute) showed an acceptable reliability and a generalizability coefficient of 0.67. Reliability of 0.8 could be obtained by testing in 4 procedures. Median QS was 0.24 (range 0.21–0.26) and 0.098 (range 0.04–0.21) for groups 1 and 2, respectively ($p = 0.001$). The resulting pass/fail standard was 0.19. Group 3 had a median posttraining QS of 0.11 (range 0–0.17). None of them met the pass/fail standard. **Conclusions:** With careful design of standardized tests, a credible standard setting and appropriate transfer studies, VR simulators could be an important first line in credentialing before proceeding to supervised performance on patients.

Copyright © 2013 S. Karger AG, Basel

Introduction

Correct TNM (tumor-node-metastasis) staging is essential for optimal diagnosis and staging of patients with non-small cell lung carcinoma. Integrated computed/

Table 1. The 5-step approach used to set pass/fail standards based on performance metrics from a VR simulator

Step 1 Identify simulator metrics with discriminative ability (table 2)
Step 2 Create an aggregate score that combines these metrics into a single QS
Step 3 Test the reliability of the QS under different circumstances (fig. 1)
Step 4 Perform a standard-setting procedure to set the pass/fail standard (fig. 2)
Step 5 Check the consequences of this standard (pass/fail ratio of different groups) (fig. 3)
Steps 1, 3, 4 and 5 are referenced with the associated tables and figures

positron emission tomography increases the accuracy of staging but tissue sampling is necessary to determine the final N-classification in the case of enlarged lymph nodes and/or positive findings [1]. Surgical mediastinoscopy has been the gold standard for many years, but studies on endoscopic ultrasound-guided needle aspiration indicate that mediastinal tissue staging should start with an endosonographic assessment and only proceed to mediastinoscopy in the case of benign findings and high clinical suspicion of malignancy in order to achieve optimal nodal staging [2]. Endosonographic techniques require specific technical skills, and studies on endobronchial ultrasound-guided transbronchial needle aspiration (EBUS-TBNA) have shown that the diagnostic yield is highly influenced by the competence of the operator [3, 4]. Earlier guidelines proposed that basic competency was ensured after performing 40 or 50 procedures, respectively [5, 6], but the latest EBUS guidelines from the British Thoracic Society acknowledge that all individuals learn at a different pace, which makes these numbers arbitrary [7]. Focus should rather be on monitoring the performance of the trainees and the outcome. However, necessary guidance and sometimes direct intervention by the supervisor for the sake of the patient makes it more difficult to assess the real competence of the trainee. Furthermore, performance is influenced by the difficulty of the procedure; restless and coughing patients can be challenging even for experienced operators. Finally, direct observation of a trainee by his/her supervisor is prone to

observer bias such as subjectivity, false impressions, the three 'isms' (ageism, racism and sexism), rumor, grudge and misinterpretation [8].

Using virtual-reality (VR) simulators could overcome all three problems. They allow the trainee to perform independently without risk to patients, provide a highly controlled, standardized measurement environment and give access to a multitude of metrics that can be used as unbiased performance data [9]. However, even high-fidelity simulators are never identical to real life and the clinical relevance of the simulator metrics needs to be explored. When creating a test set-up, it is also important to contemplate that the performances of trainees are highly variable, especially in the first, cognitive stage of the classic motor skills learning model by Fitts and Posner [10]. Extending the test by including multiple scenarios improves its reliability [11], but care must be taken to ensure that it is still feasible to administer. Finally, when setting a pass/fail standard, it is important to consider the limited difficulty of the simulated procedure and a possible ceiling effect when using simulator metrics; trainees should not be able to pass the test after only a brief training session.

The aims of this study were to explore the validity of EBUS simulator metrics and the reliability of different test set-ups, to create a credible pass/fail standard and to investigate if a limited amount of simulator training would allow novices to pass the simulator test.

The research questions were:

- Which simulator metrics can discriminate between novices and experienced operators?
- How many procedures need to be performed to achieve an acceptable reliability?
- What is a credible pass/fail standard of the test?
- What is the pass rate for novices after approximately 2 h of simulator training?

Materials and Methods

A 5-step approach was used to answer the research questions (table 1). Three groups of physicians were included in the study. Group 1 was a convenience sample of experienced EBUS operators from four different hospitals who had each performed more than 200 EBUS-TBNA procedures (4 males and 2 females). Groups 2 and 3 were respiratory physicians from Denmark and the Netherlands (8 males and 8 females). They had all attended a 1-day theoretical EBUS course in January 2012, and had only participated in 2–5 EBUS procedures with real patients before their simulator session. All physicians were experienced in flexible bronchoscopy and mediastinal anatomy. These physicians were paired and then randomized to either group 2 or group 3 by an impartial nurse or secretary using sets of envelopes.

Table 2. The metrics automatically generated by the simulator software after the performance of simulated EBUS-TBNA

Simulator metric	Group 1: (n = 6) experienced operators	Group 2: (n = 8) novices	p value
Successful samples performed, n	1.9±0.20	1.4±0.50	0.047*
Total procedure time, minutes and seconds	8 min 28 s ± 47 s	16 min 1 s ± 4 min 43 s	0.002*
Percentage of time with ultrasound visualization	66±5	62±7	0.29
Amount of lidocaine used, mg	301±194	376±174	0.46
Hemodynamic complications, n	0.33±0.44	0.63±0.44	0.18
Attempted samples with the scope in a nonoptimal position, n	2.2±0.75	1.9±1.7	0.70
Blood vessels punctured, n	0.75±0.27	0.38±0.35	0.053
Risk of scope damage, n ¹	0.58±0.38	0.94±0.32	0.081

Means and standard deviations for a group of experienced EBUS operators (group 1) and a group of untrained novices (group 2) are shown along with the p values. * Shows significant differences.

¹ The number of times the biopsy needle was introduced or retracted while the tip of the scope was flexed.

Training and testing were done at individual sessions in Copenhagen, Leiden or Amsterdam by the same researcher (L.K.) using identical equipment placed in rooms assigned for the purpose where each participant could perform without disturbances. Before the test, the participants in groups 1 and 2 were allowed only a warm-up time in the simulator of a maximum of 5 min, whereas the participants in group 3 all completed a standardized training program consisting of 5 cases focusing on mediastinal anatomy, landmark recognition using ultrasound and performance of ultrasound-guided biopsies. Participants in group 3 were tested immediately after completing the training program.

The test consisted of EBUS-TBNA procedures on simulator patient cases. Patient 1 was a 73-year-old male with an intrapulmonary mass in the left lower lobe, and enlarged lymph node stations: left lower paratracheal (station 4L), subcarinal (station 7) and left interlobar (station 11L). Patient 2 was a 72-year-old male with squamous cell carcinoma of the right upper lobe and enlarged lymph node stations: right upper paratracheal (station 2R), right lower paratracheal (station 4R), subcarinal (station 7) and right hilar (station 10R). During the procedures, the participants had to introduce the scope, demonstrate anatomical knowledge by identifying six landmarks (station 4 left, station 7, station 10 or 11 left, station 10 or 11 right, the azygos vein and station 4 right), and perform two biopsies from stations 7 and 4R, respectively. The 2 simulator patient cases in the test were not included in the training program undertaken by group 3.

The VR simulator used in the study was the GI Bronch Mentor™ (Symbionix, Cleveland, Ohio, USA) that consisted of a proxy bronchoscope, a proxy syringe for applying local anesthesia, a proxy EBUS biopsy needle, an interface tracking the motions of the equipment and a monitor displaying the computer-generated endoscopic and ultrasound images. After each procedure, the software automatically stores the following metrics: the number of successful samples performed, the total procedure time, the percentage of time with ultrasound visualization (= wall contact), the amount of lidocaine used, the number of hemodynamic complications (desaturation or hypertension), the number of attempted samples with the scope in a nonoptimal position, the number of

blood vessels punctured and the number of times the biopsy needle was introduced or retracted while the tip of the scope was flexed (= risk of scope damage).

All simulator metrics were exported into a statistics program for analysis (see below), and the metrics that could discriminate between group 1 (experienced operators) and group 2 (untrained novices) were combined to a single, aggregated score: the quality score (QS). The reliability of the QS was tested, and the effect of changing the number of procedures in the test was explored. An examinee-based method, the contrasting-groups method [12], was used for setting the standard. The score distributions of groups 1 and 2 were plotted and the passing score was set at the intersection of the distributions, as we considered false-positives (novices passing the test) and false-negatives (experienced EBUS operators failing the test) to be of equal weight. Finally, the consequences of the pass/fail standard for the participants in group 3 (simulator-trained novices) were investigated.

Statistical Analysis

Metrics with discriminatory ability were identified using a 2-way mixed analysis of variance (ANOVA) with procedures as the repeated-measures variable and experience (either experienced operator or novice) as the between-group variable. Generalizability theory was used to explore the reliability of the resulting aggregated score and the effect of changing the number of procedures in the test. First, the estimated variance components were calculated in a 1-facet, balanced G-study with test persons (p) crossed with simulated EBUS procedures (e), i.e. $p \times e$. These components were then used in the context of a D-study design to estimate generalizability coefficients for different number of EBUS procedures performed. The aim was a generalizability coefficient >0.8.

Statistical analysis was performed using a statistical software package (PASW, version 18.0; SPSS Inc., Chicago Ill., USA). Differences were considered to be statistically significant when the p value was <0.05. The G-study and D-study were performed using the GENOVA software version 3.1 [13].

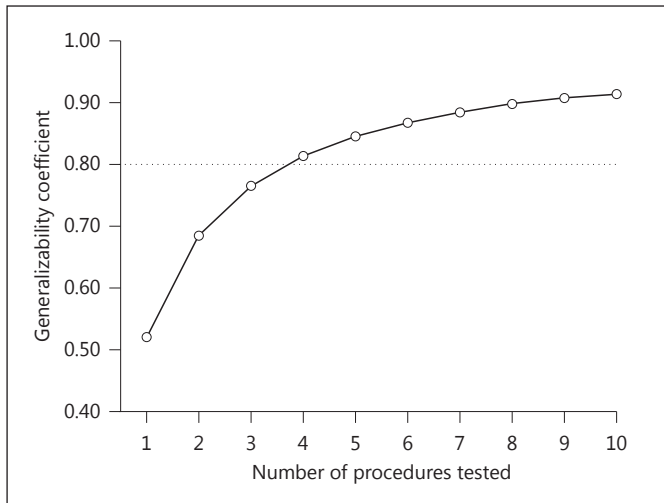


Fig. 1. The effect on the generalizability coefficient of assessing more procedures (the reliability of different test set-ups). A generalizability coefficient of 0.8 is generally accepted as necessary for high-stakes assessment.

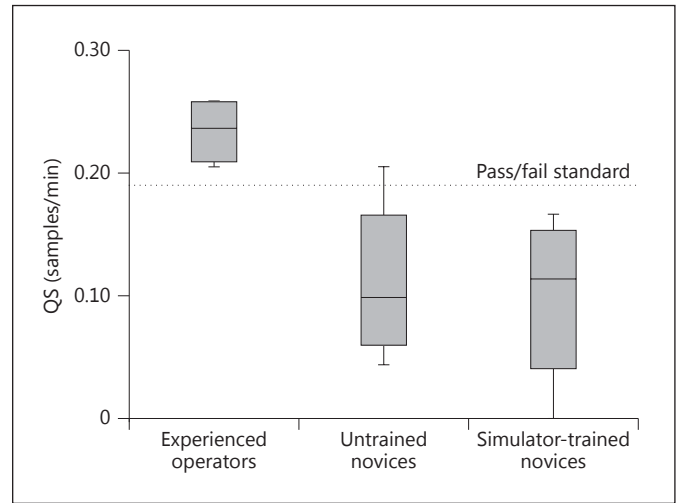


Fig. 3. QSs (samples/min) of experienced operators, untrained novices and simulator-trained novices. Box plot shows outliers, minimum, first quartile, median, third quartile and maximum.

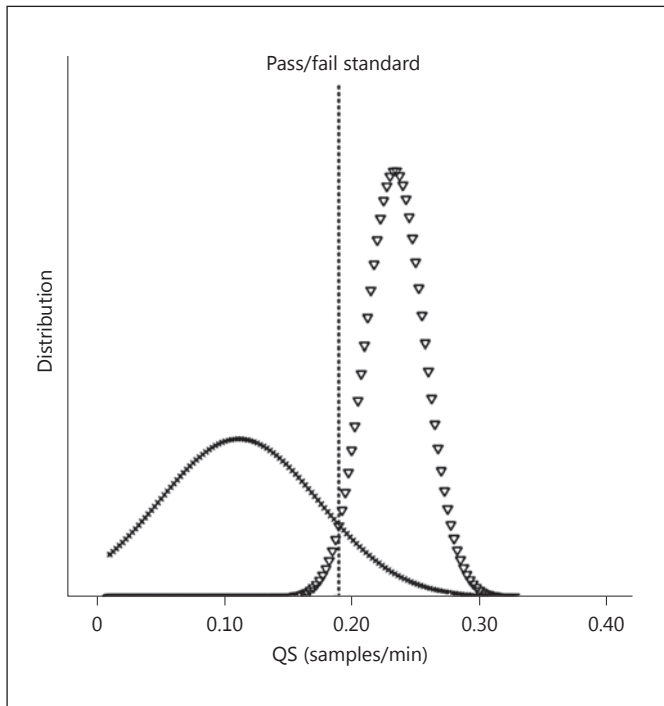


Fig. 2. Establishing a pass/fail standard using the contrasting-groups method: the distribution of scores of untrained novices (crosses) and experienced operators (triangles), respectively, are plotted using the means and standard deviations. The intersection marks the pass/fail standard of 0.19 samples/min.

Results

Six experienced EBUS operators joined group 1, and there were 8 respiratory physicians in both groups 2 and 3. The first physician was tested on January 5, 2012, and the last test was performed on August 6, 2012. Table 2 shows the performance of groups 1 and 2. Procedure time and successfully sampled lymph nodes were the only simulator metrics that showed statistically significant differences. These two metrics were combined to a QS (sampled lymph nodes per minute) that showed an acceptable reliability and a generalizability coefficient of 0.67. Figure 1 shows that a reliability of 0.8 could be obtained by testing in 4 procedures. The median QS of the untrained novices was 0.098 (range 0.04–0.21) and that of the experienced operators was 0.24 (range 0.21–0.26) ($p = 0.001$). Figure 2 shows these score distributions, and the intersection marks the pass/fail standard, QS = 0.19. The mean, effective simulator practice time of group 3 was 1 h and 46 min (SD 17 min). Figure 3 shows the consequences of the established pass/fail standard: all the experienced operators passed the test as well as a single outlier in the untrained novice group. The novices who had a brief simulator training session did not achieve a significantly higher QS compared to novices without training, and none of them managed to pass the test.

Discussion

A recent survey among pulmonary fellowship directors in the USA found that only 30% of programs had a formal protocol in place to evaluate EBUS competency [14]. We argue that high-fidelity VR simulators like the EBUS simulator used in this study could be used for initial training and assessment of basic competence before proceeding to supervised practice on patients. The simulators automatically deliver a vast amount of nonbiased performance metrics immediately after each simulated procedure, and it is easy to use these data to assess the competency of the trainee. However, care must be taken to identify which metrics are clinically relevant, i.e. by selecting metrics where experienced operators significantly outperform novices. Seemingly meaningful metrics can lack discriminatory ability and therefore be useless. An example from a study on a VR bronchoscopy simulator is the metric 'percentage of bronchial segments visualized' that was used to reflect the completeness of the examination but found that complete novices with limited anatomical knowledge navigated the bronchial tree like a rat in a maze and thereby achieved the same scores as experts [15].

The experienced operators in this study managed to sample more lymph nodes than novices, which was not surprising as performance of EBUS-TBNA was a core competence in their everyday clinical practice. The experienced group also spent less time on the procedures and performed more consistently, both of which are recognized features of the final stage of learning in motor skills learning theory, i.e. the autonomous stage in the Fitts and Posner model [10]. Another trademark of skilled operators is their ability to avoid errors, but surprisingly, group 1 in this study did not make significantly fewer errors than group 2 (table 1). A possible explanation is that the experienced operators seemed more determined to obtain good samples thereby also risking errors, i.e. sampling with the scope in a nonoptimal position or puncturing a blood vessel. Novices typically terminated the procedure after the planned two attempts, not realizing that one or both had failed, and a test solely focusing on avoiding errors would encourage this approach. The optimal percentage of time with ultrasound visualization and the correct amount of lidocaine used during an EBUS procedure are difficult or impossible to define, and were not significantly different for the 2 groups; these metrics should not be used to assess competence.

Even though only two metrics had discriminatory abilities, the combined QS was able to discriminate between experienced operators and untrained novices, thereby in-

dicating that the simulator test possessed construct validity [16]. However, the test should also be reliable, i.e. it should be possible to replicate or reproduce the data [17]. The individual variance in performance among the participating physicians in this study meant that applying the test for a single procedure would result in a generalizability coefficient of only 0.52 (fig. 1). Four procedures were needed to ensure generalizability coefficients >0.8 which is generally accepted for higher-stakes examinations [18]. This finding accords well with clinical studies on diagnostic bronchoscopy and mediastinal endoscopic ultrasound which found that 3 and 4 procedures, respectively, were needed for the reliable assessment of competence [19, 20]. A test of EBUS-TBNA performance consisting of 4 procedures would take approximately 1 h to administer and a few minutes to evaluate (using the simulator metrics). Hence, it would be feasible to apply for new trainees.

A reliable and valid assessment method can be used to provide important feedback to trainees (formative assessment), but a formal test created for certification (summative assessment) also needs a pass/fail score. There is no single correct or best standard-setting method, but an appropriate method has to be utilized and a credible standard has to produce reasonable outcomes [21]. We found a meaningful pass/fail score with appropriate consequences (fig. 2, 3) using the contrasting-groups method that has also been successfully used for setting standards for a simulation-based test of flexible bronchoscopy [22]. However, if a simulated task is easy to master (i.e. has a steep learning curve) there will be a considerable ceiling effect allowing trainees to pass the test after a brief training session. A study on flexible bronchoscopy simulation showed no differences between novices that had tried the simulator 5 times and experts [23], and in another study, 8 out of 10 medical students passed a simulator test after 1–3 h of training [24]. EBUS-TBNA is a difficult procedure to master and the fact that approximately 2 h of training did not allow the trainees to reach expert performance indicates that the VR simulator mimics the difficulties of real-life patients, and that the QS is an appropriate measure of competence.

A limitation to this study is that the simulator metrics could not explore the completeness of the diagnostic examination. The participants were asked to demonstrate anatomical knowledge by identifying six anatomical landmarks, but the simulator could not record the quality of this investigation. An earlier study on EBUS simulation created additional clinically relevant measurements based on direct observation, but this approach requires the participation of a skilled EBUS supervisor and reintroduces

possible bias into the assessment [25]. Bias can be reduced, but not removed, by using a validated assessment instrument and trained expert raters [26]. A test based solely on simulator metrics does not require resources from busy attending physicians, but it is important to realize that the simulator only assesses the technical issues of the procedure. The ability to provide the correct indication for the procedure, interpret the radiological examinations and communicate with the patient still need to be assessed under direct observation during clinical practice. Another limitation is the limited number of participating respiratory physicians (n = 22). We propose incorporating simulation-based training and testing in international, standardized EBUS teaching curricula in the future and continue to monitor experiences. Appropriate transfer studies exploring the impact of simulator training on actual patient management are necessary.

Conclusion

Great caution has to be applied when using simulator metrics to assess procedural competence. We only found two metrics with discriminative ability, procedure time

and successfully sampled lymph nodes. However, these data could be used to create a reliable test by testing in 4 procedures, and set a credible pass/fail standard. No novices passed the test after a brief training session. With proper and careful design of standardized tests, we believe that VR simulators could be an important first line in credentialing procedural skills before proceeding to supervised performance on patients.

Acknowledgments

We thank Drs Mark Krasnik, Mette Siemsen and Klaus Richter for their participation as experienced EBUS operators.

Financial Disclosure and Conflicts of Interest

The authors have no conflicts of interest or financial disclosure. The trial was investigator-initiated. Design and execution of the study were performed independently from the simulator company. Funds for the study were supplied by the institutions of the authors which also owned all three simulators used in the study.

References

- 1 Fischer BM, Mortensen J, Hansen H et al: Multimodality approach to mediastinal staging in non-small cell lung cancer. Faults and benefits of PET-CT: a randomised trial. *Thorax* 2011;66:294–300.
- 2 Tournoy KG, Keller SM, Annema JT: Mediastinal staging of lung cancer: novel concepts. *Lancet Oncol* 2012;13:e221–e229.
- 3 Steinfors DP, Hew MJ, Irving LB: Bronchoscopic evaluation of the mediastinum using endobronchial ultrasound – a description of the first 216 cases performed at an Australian tertiary hospital. *Intern Med J* 2011;41:815–824.
- 4 Kemp SV, El Batrawy SH, Harrison RN et al: Learning curves for endobronchial ultrasound using cusum analysis. *Thorax* 2010;65:534–538.
- 5 Bolliger CT, Mathur PN, Beamis JF et al: European Respiratory Society/American Thoracic Society statement on interventional pulmonology. *Eur Respir J* 2002;19:356–373.
- 6 Ernst A, Silvestri GA, Johnstone D: Interventional pulmonary procedures: guidelines from the American College of Chest Physicians. *Chest* 2003;123:1693–1717.
- 7 Du Rand IA, Barber PV, Goldring J et al: Summary of the British Thoracic Society guidelines for advanced diagnostic and therapeutic flexible bronchoscopy in adults. *Thorax* 2011;66:1014–1015.
- 8 McGaghie WC, Butter J, Kaye M: Observational assessment; in Downing SM, Yudkowsky R (eds): *Assessment in Health Professions Education*, ed 1. New York, Routledge, 2009, pp 185–215.
- 9 McGaghie WC, Issenberg SB: Simulations in assessment; in Downing SM, Yudkowsky R (eds): *Assessment in Health Professions Education*, ed 1. New York, Routledge, 2009, pp 245–268.
- 10 Magill RA: *The Stages of Learning*. Motor Learning and Control, ed 8. New York, McGraw-Hill, 2007, pp 263–289.
- 11 Streiner DL, Norman GR: Reliability; in Streiner DL, Norman GR (eds): *Health Measurement Scales – A Practical Guide to Their Development and Use*, ed 4. Oxford, Oxford University Press, 2008, pp 167–210.
- 12 Livingston S, Zieky M: *Passing Scores: A Manual for Setting Standards of Performance on Educational and Occupational Tests*. Princeton, Educational Testing Service, 1982.
- 13 Crick JE, Brennan RL: Freeware computer-program, University of Iowa, USA. Download: http://www.uiowa.edu/~casma/computer_programs.htm.
- 14 Tanner NT, Pastis NJ, Silvestri GA: Training for linear endobronchial ultrasound among US pulmonary/critical care fellowships: a survey of fellowship directors. *Chest* 2013;143:423–428.
- 15 Konge L, Arendrup H, von BC, Ringsted C: Using performance in multiple simulated scenarios to assess bronchoscopy skills. *Respiration* 2011;81:483–490.
- 16 Ringsted C, Hodges B, Scherpier A: ‘The Research Compass’: an introduction to research in medical education: AMEE Guide No. 56. *Med Teach* 2011;33:695–709.
- 17 Axelson RD, Kreiter CD: Reliability; in Downing SM, Yudkowsky R (eds): *Assessment in Health Professions Education*. New York, Routledge, 2009, pp 57–73.
- 18 Downing SM: Reliability: on the reproducibility of assessment data. *Med Educ* 2004;38:1006–1012.
- 19 Konge L, Larsen KR, Clementsen P, Arendrup H, von BC, Ringsted C: Reliable and valid assessment of clinical bronchoscopy performance. *Respiration* 2012;83:53–60.

- 20 Konge L, Vilmann P, Clementsen P, Annema JT, Ringsted C: Reliable and valid assessment of competence in endoscopic ultrasonography and fine-needle aspiration for mediastinal staging of non-small cell lung cancer. *Endoscopy* 2012;44:928–933.
- 21 Norcini J, Guille R: Combining tests and setting standards; in Norman GR, Van der Vleuten CMP, Newble DI (eds): *International Handbook of Research in Medical Education*. Amsterdam, Kluwer Academic Publishers, 2002.
- 22 Konge L, Clementsen P, Larsen KR, Arendrup H, Buchwald C, Ringsted C: Establishing pass/fail criteria for bronchoscopy performance. *Respiration* 2012;83:140–146.
- 23 Moorthy K, Smith S, Brown T, Bann S, Darzi A: Evaluation of virtual reality bronchoscopy as a learning and assessment tool. *Respiration* 2003;70:195–199.
- 24 Krogh CL, Konge L, Bjurström J, Ringsted C: Training on a new, portable, simple simulator transfers to performance of complex bronchoscopy procedures. *Clin Respir J* 2012, Epub ahead of print.
- 25 Stather DR, Maceachern P, Rimmer K, Hergott CA, Tremblay A: Validation of an endobronchial ultrasound simulator: differentiating operator skill level. *Respiration* 2011;81:325–332.
- 26 Davoudi M, Colt HG, Osann KE, Lamb CR, Mullon JJ: Endobronchial ultrasound skills and tasks assessment tool: assessing the validity evidence for a test of endobronchial ultrasound-guided transbronchial needle aspiration operator skill. *Am J Respir Crit Care Med* 2012;186:773–779.