



Archived by Flinders University

This is the peer reviewed version of the following article:  
Sam, A. H., Field, S. M., Collares, C. F., van der Vleuten, C. P. M., Wass, V. J., Melville, C., Harris, J., & Meeran, K. (2018).  
Very-short-answer questions: reliability, discrimination and acceptability. *Medical Education*, 52(4), 447–455. <https://doi.org/10.1111/medu.13504>

which has been published in final form at  
<https://doi.org/10.1111/medu.13504>

This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for self-archiving.

© 2018 John Wiley & Sons Ltd and The Association for the Study of Medical Education.



medical education

www.mededuc.com

**Very short answer questions: reliability, discrimination and acceptability**

Journal:	<i>Medical Education</i>
Manuscript ID	MED-2017-0755.R1
Manuscript Type:	Research Papers
Keywords:	Testing/Assessment

SCHOLARONE™  
Manuscripts

view

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

**1 Very short answer questions: reliability, discrimination and acceptability**

- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25

For Review

1  
2  
3 1 **Abstract**  
4

5 2  
6

7 3 **Background:**  
8

9 4 Single best answer questions (SBAs) have been widely used to test knowledge because  
10 5 they are easy to mark and demonstrate high reliability. However, SBAs are open to the  
11 6 criticism that they are subject to cueing. We used a novel assessment tool that facilitates  
12 7 efficient marking of open-ended very short answer questions (VSAs). We compared VSAs  
13 8 with SBAs with regards to reliability, discrimination and student performance, as well as and  
14 9 evaluating the acceptability of VSAs.  
15  
16  
17  
18

19 10  
20

21 11 **Method:**  
22

23 12 Medical students were randomised to sit a 60-question assessment in either VSA then SBA  
24 13 format (group 1, n=155) or the reverse (group 2, n=144). The VSAs were delivered on a  
25 14 tablet, which was computer-marked and subsequently reviewed by two examiners. Standard  
26 15 error of measurement (SEM) across the ability spectrum was estimated using item response  
27 16 theory.  
28  
29  
30

31 17  
32

33 18 **Results:**  
34

35 19 Review of the machine-marked questions took on average one minute 36 seconds per  
36 20 question for all students. The VSAs had a high reliability (alpha: 0.91), significantly lower  
37 21 SEM than the SBAs ( $p < 0.001$ ) and higher mean item-total point biserial correlations  
38 22 ( $p < 0.001$ ). The VSA scores were significantly lower than the SBA scores ( $p <$   
39 23  $0.001$ ). The difference in scores between VSAs and SBAs was attenuated in group 2.  
40 24 Although 80.4% of students found the VSAs more difficult, 69.2% found them more  
41 25 authentic.  
42  
43  
44  
45

46 26  
47

48 27 **Conclusions:**  
49

50 28 VSAs demonstrated high reliability and discrimination and were perceived as more  
51 29 authentic. SBAs were associated with significant cueing. Our results suggest a higher  
52 30 degree of validity for VSAs.  
53  
54  
55

56 31  
57  
58  
59  
60

## 1 Introduction

Multiple choice single best answer (SBA) questions are widely used in undergraduate and postgraduate medical assessment programmes worldwide, as they tend to have high levels of reliability and can be machine-marked efficiently and accurately. However, there have long been concerns that multiple choice questions may not provide a true reflection of knowledge as they rely on answer recognition rather than recall.<sup>1, 2</sup> Students tend to perform better in multiple choice questions than open-ended, free response short answer questions.<sup>2-7</sup> This may be due to the effect of cueing when candidates are presented with a list of options. A core principle of the validity of an assessment is the extent to which the test measures the competency it is supposed to measure.<sup>8</sup> Thus, if the aim of the assessment is to examine the student's ability to synthesise or generate rather than recognise a correct answer, short answer questions may provide a greater validity.<sup>9</sup>

It is well recognised that assessment drives learning,<sup>10, 11</sup> and that students prepare differently for exams in different formats.<sup>10, 12-14</sup> Indeed recognition and recall require different learning operations and their distinction has long been recognized in cognitive psychology.<sup>15</sup> Learning exam technique for multiple choice questions is a recognised phenomenon, which may lead to exam success at the expense of a deeper understanding of the subject being tested.<sup>16-18</sup> Furthermore, students demonstrate greater long-term information retention after studying for or completing short answer question rather than multiple choice questions.<sup>19-22</sup>

A meta-analysis investigating construct equivalence of multiple-choice and constructed response (e.g. short answer) items showed a higher mean correlation between the two formats when stem-equivalent items were used.<sup>23</sup> However, studies of stem-equivalent items include assessment of topics that lend themselves to multiple-choice questions. In our experience, question-writers can find it challenging to think of sufficient plausible distractors for some topics. The content of the assessment may also be skewed as core knowledge becomes too easy, causing question writers to resort to the testing of obscure material.<sup>24</sup> Short answer questions can offer greater flexibility for question-writers by allowing them to focus on common and relevant themes rather than academic minutiae.<sup>24, 25</sup>

Despite the potential advantages of short answer questions, their use in large-scale assessments has been limited by feasibility as, historically, they have not been amenable to being machine-marked<sup>4, 26</sup> and thus have been unable to sample the curriculum efficiently due to resource limitation. We have previously used Microsoft Excel to mark very short answer (VSA) questions<sup>9</sup>. In the current study, we used an online assessment management system, which allows questions to be posed in a VSA format on an electronic platform

1 requiring one-to-four word answers. Assessment item types can be categorised according  
2 to the degree of constraint on the respondent's options for answering ranging from fully  
3 constrained/selected responses (multiple choice questions) to fully constructed responses  
4 (essays).<sup>27</sup> Very short answer questions fall in the 'intermediate constraint' items which can  
5 be marked efficiently and accurately by computer using new information technologies.

6 The utility of an assessment method can be evaluated by its reliability, validity, acceptability,  
7 educational impact and costs.<sup>28</sup> The purpose of this study was to compare VSA and SBA  
8 questions to assess the utility of VSA questions as an assessment method. Although it is  
9 difficult to link assessment formats with learning behaviour directly,<sup>13</sup> we assessed student  
10 opinions and potential educational impact of VSA questions with a post-test student survey.

## 11 12 **Methods**

### 13 ***Participants and assessments***

14 This study was approved by the Medical Education Ethics Committee at Imperial College  
15 London. Ethical approval was granted to invite all medical students in year 3 at Imperial  
16 College School of Medicine by an email from the faculty to sit a formative exam. All students  
17 had been on one surgical and two medical attachments in year 3. There were no other  
18 inclusion or exclusion criteria.

19 Medical students sat a formative examination consisting of 60 questions under exam  
20 conditions. The clinical vignettes were constructed to allow them to be posed in both SBA  
21 and VSA formats without changing the content (Figure 1). During the construction of the  
22 VSA questions we were able to generate a list of acceptable answers for each item within  
23 five minutes. Content validity was ensured by blueprinting against the year 3 curriculum at  
24 Imperial College School of Medicine to ensure broad sampling of relevant topics and close  
25 alignment with the syllabus. Items were written with short case descriptions and tested a  
26 range of cognitive processes, including clinical reasoning, decision-making and knowledge  
27 recall. They were independently reviewed to minimise construction errors.

28 Examinees were randomly assigned to two groups. In group 1, the examinees were  
29 presented with 60 questions posed in VSA format (a 90 minute exam), immediately followed  
30 by the same 60 questions posed in SBA format with five options (a 60 minute exam). In  
31 group 2, the first test consisted of SBA questions and the second of VSA questions. The  
32 students were given more time for the VSA questions because of the time taken to type the  
33 answers. Students had to complete the first test before commencing the second one, and  
34 could not return to the previous test. VSA questions were posed using a new online

1 examination management software (*Practique, Fry*) where the students provided answers on  
2 an iPad. The SBA format was delivered using a traditional paper-based system with a  
3 machine-marked scoring card (*Multiquest, Speedwell*). Following completion of both test  
4 formats, the students were invited to complete a feedback form to evaluate student opinions  
5 on the VSA questions.

### 6 **Marking**

7 Answers to the SBA questions were machine-marked with the individual student and  
8 question performance exported for statistical analysis. The students' answers to the VSA  
9 questions captured by the iPad App (*Practique, Fry*) were sent to a server over an encrypted  
10 connection. At the end of the exam the server applied an automated matching algorithm  
11 using the Levenshtein distance to match each answer against pre-approved acceptable  
12 answers for each question. Subsequently two markers reviewed all the non-exact matches  
13 and match failures to see if any of the non-exact matches should be disallowed, or any of the  
14 match failures should be allowed. Grouping similar answers in blocks facilitated this  
15 verification process. The system applied the examiner marking judgments to all the identical  
16 answers, ensuring consistency and saving the marker time. The system also learns the new  
17 marking judgments for each question and adds this to the pre-approved answer list for each  
18 question to improve the automatic marking for the next time that particular question is used.

19 Figure 2 shows an example of the marking system. The green answers have been  
20 automatically marked as correct based on the pre-approved answers. The yellow answers  
21 have been marked as correct based on their similarity to pre-approved answers. Answers  
22 marked as incorrect are shown in red, however during the review process this can be over-  
23 ridden, with all identical answers automatically given the same mark. The time taken for the  
24 examiners to review each item was recorded to give a measure of acceptability.

### 25 **Analysis**

26 Statistical analyses were performed with *IBM SPSS Statistics* version 24.0 and *Prism*  
27 version 5.0c software (GraphPad Software). The Kolmogorov-Smirnov test showed a  
28 normal distribution of all variables. Pearson's correlation coefficient was used to assess the  
29 correlation between the raw scores of the two formats. The difference in sex distribution  
30 between groups was tested using the chi-square test. Differences in raw scores, item-total  
31 correlations and standard error of measurement (SEM) within and between groups were  
32 analyzed using mixed design analysis of variance (ANOVA) with effect sizes expressed as  
33 partial eta squared ( $\eta_p^2$ ). Differential item functioning (DIF) for sex was assessed with  
34 *Xcalibre* 4.2 (Assessment Systems) for all items using the significance of the z-test based on  
35 the Mantel-Haenszel coefficient.

1 Students' responses to the same questions in the two formats were compared, with 'positive  
2 cueing' defined as percentage of questions answered correctly in the SBA format and  
3 incorrectly in the VSA format. 'Negative cueing', occurred where distractors caused students  
4 who knew the correct answer in the VSA question to answer the SBA question incorrectly.<sup>6</sup>

5 Three-parameter logistic model analysis was carried out to include a specific parameter to  
6 estimate the probability of correct answers due to guessing. Analysis was carried out using  
7 the R package *mirt*.

8 In the post-test survey, students were asked to rate the following four statements on a 5-  
9 point Likert scale (strongly disagree, disagree, neutral, agree, strongly agree).

- 10 1. Questions in the single best answer format are easier than the very short answer format.
- 11 2. Very short answer questions are a better representation of how I would be expected to  
12 answer questions in clinical practice.
- 13 3. Having examinations in very short answer format would change my learning and revision  
14 strategy.
- 15 4. Using very short answer questions in assessments would help improve my preparation  
16 for clinical practice.

17 The questionnaire also provided a space for students to write comments about the use of  
18 VSA questions in medical school assessments. The *NVivo* software was used to identify  
19 themes in the students' free text responses.

## 20 21 **Results**

22 Of 340 students in year 3 at Imperial College School of Medicine, 302 students (89%) sat the  
23 formative examination, with 299 students (99%) completing both parts. The sex distribution  
24 was similar in both groups (female/male ratio in group 1: 72:83, and group 2: 63:81,  $p=0.64$ ).  
25 Among the tests in both groups, only one item (a VSA question group 1) had significant  
26 differential item functioning (DIF) for sex, with bias against males ( $p = 0.04$ ).

27 There was a significant positive correlation between the two formats ( $p<0.001$ ,  $r=0.83$ ). Due  
28 to the high reliability and the identical content of the tests, correction for attenuation was not  
29 performed, as this was likely to overestimate the correlation.

## 30 **Acceptability**

31 The VSA questions were reviewed by two examiners. Based on the preloaded acceptable  
32 answers, the system was able to identify 80.2% of correct answers prior to review, which  
33 was instrumental in allowing efficient marking. The remainder of correct answers were  
34 marked during the review process. Of answers marked 'correct' by the system, 0.2% were



1 deemed to be 'incorrect' on review (due to a spelling error significantly changing the  
2 meaning of the answer). The total time taken to review the machine-marked answers to all  
3 60 VSA questions for all 299 students was 95 minutes 51 seconds. The average time spent  
4 by the examiners to review the answers to each question for all 299 students was 1 minute  
5 and 36 seconds (standard deviation=1 minute 2 seconds). The marking system allowed for  
6 multiple correct answers aside from trivial spelling or terminology differences, and in 8%  
7 (5/60) of questions at least one student offered an alternative answer to the question that  
8 was judged to be correct following review by examiners.

### 9 ***Effect of cueing***

10 The raw scores for each test are shown in Table 1. The raw scores for VSA and SBA  
11 questions in group 1 were 52.4% and 68.2% respectively. In group 2, the raw scores for VSA  
12 and SBA questions were 65.7% and 69.7% respectively. There was a significant difference  
13 in raw scores between item type with a large effect size [ $F(1,297)=384,339$ ,  $p<0.001$ ,  
14  $\eta_p^2=0.56$ ]. There were also significant differences in the interaction between item type and  
15 group [ $F(1,297)=136,343$ ,  $p<0.001$ ,  $\eta_p^2=0.32$ ,  $p<0.001$ ] and between the groups  
16 [ $F(1,297)=19,854$ ,  $p<0.001$ ,  $\eta_p^2=0.06$ ]. The difference in VSA and SBA scores between the  
17 two groups is likely due to students seeing the answer options in the SBA questions before  
18 answering the VSA questions and the associated cueing effect. Positive cueing, with  
19 students answering the SBA question correctly and the equivalent VSA question incorrectly,  
20 was seen in 19.2% of items in group 1 and 7.5% of items in group 2. Negative cueing, with  
21 students answering the VSA question correctly and the equivalent SBA question incorrectly,  
22 was in seen in 3.5% of items in group 1 and 3.5% of items in group 2. On an item level,  
23 positive cueing occurred for every item (100%) and negative cueing occurred in 50 out of 60  
24 (83.3%) items. There were four items where the students performed better in the VSA than  
25 the SBA questions, possibly because the answer options distracted the students. Notably in  
26 20% of items, over 30% of students in group 1 could only get the answer correct in its SBA  
27 question format. For example, whilst 13% of students in group 1 were able to generate a  
28 diagnosis of erythema multiforme in the VSA question, 68% were able to select the correct  
29 answer in the SBA question.

### 30 ***Reliability and Discrimination***

31 Table 1 shows reliability (Cronbach's alpha) and standard error of measurement (SEM)  
32 values for the VSA and SBA tests in both groups. SBA tests had Cronbach's alpha values of  
33 0.84 and 0.85 in groups 1 and 2 respectively. VSA tests had a Cronbach's alpha of 0.91 in  
34 both groups 1 and 2.

35 We also used the three-parameter logistic model to estimate the standard error of  
36 measurement (SEM) for the two tests in both groups across the ability spectrum (Figure 3).

1 VSA questions had significantly lower standard error of measurement [ $F(1,297)=213,782$ ,  
2  $p<0.001$ ,  $\eta_p^2=0.42$ ]. The effects of the group and the interaction between item format and  
3 group were not significant ( $p=0.23$  and  $0.97$  respectively).

4 Figure 4 shows individual estimates for information according to theta ability estimates and  
5 item type and group combination. VSA questions had significantly higher information with a  
6 large effect size [ $F(1,297)=311,998$ ,  $p<0.001$ ,  $\eta_p^2=0.51$ ]. There was no significant interaction  
7 between item format and group [ $F(1,297)=3,584$ ,  $p=0.06$ ,  $\eta_p^2=0.01$ ] and the group effect size  
8 was small [ $F(1,297)=4,742$ ,  $p=0.03$ ,  $\eta_p^2=0.02$ ].

9 Mean item-total score point-biserial correlations for VSA questions were 0.36 and 0.35 in  
10 groups 1 and 2 respectively. Mean item-total score point-biserial correlations for SBA  
11 questions were 0.26 and 0.27 in groups 1 and 2 respectively. VSA questions had  
12 significantly higher item-total correlations [ $F(1,118)=89,235$ ,  $p<0.001$ ,  $\eta_p^2=0.43$ ]. The  
13 effects of the group and the interaction between format and group were not significant  
14 ( $p=0.81$  and  $0.30$  respectively).

### 15 ***Potential impact on learning behaviour***

16 Figure 5 shows the percentage of students selecting each point on the Likert scale for each  
17 statement. 80.4% of students agreed/strongly agreed that SBA questions were easier than  
18 VSA questions. With regards to authenticity, 69.2% agreed/strongly agreed that VSA  
19 questions were more representative of how they would be expected to answer questions in  
20 clinical practice. Almost half the cohort (49.3%) agreed/strongly agreed that having VSA  
21 questions in summative examinations would change their revision and learning strategy and  
22 would help improve their preparation for clinical practice.

23 Thirty one percent of students who thought SBA questions were easier, commented this was  
24 due to the presence of options. Fifty six percent of free text comments regarding authenticity  
25 were related to how in practice they will be expected to recall information without options.  
26 Seventy percent of comments on a change in revision strategy indicated more emphasis on  
27 thoroughness and 9% would spend more time learning spelling.

### 29 **Discussion**

30 Our results indicate that VSA questions have advantages over SBA questions in terms of  
31 reliability and discrimination. The cueing effect associated with SBA questions was apparent  
32 in the analysis of the scores for the VSA questions and SBA questions in the two groups.  
33 The difference in scores between the two formats was attenuated when the students saw the

1 SBA questions first, likely due to the cueing effect of the options. Therefore, VSA questions  
2 have higher validity when testing ability to arrive at a correct answer without cueing or  
3 guessing.

4 Another potential limitation of SBA questions is the implication that there is one best answer  
5 for any question, therefore discouraging question writing in areas of medicine where there  
6 may be multiple defensible answers.<sup>1, 3</sup> Very short answer questions allow students to offer  
7 alternative answers that may be as good or second-option alternatives. Allowing students to  
8 demonstrate the scope of their knowledge improves the validity of the assessment.

9 Almost 70% of students thought VSA questions better represent how they would be  
10 expected to answer questions in real life clinical practice, and this suggests they have  
11 greater authenticity. More authentic examinations promote both deeper learning methods  
12 and increase student motivation.<sup>29</sup> Furthermore, a major component of learning from  
13 assessment is the quality of the feedback provided.<sup>8, 11</sup> Very short answer questions can  
14 offer an opportunity to provide more specific and detailed feedback based on the diverse  
15 range of incorrect answers proposed by the students. These can therefore be addressed or  
16 targeted by curriculum developers.

17 A major limitation to the widespread use of VSA questions would have been their  
18 acceptability in terms of resources. Introduction of any novel assessment method should be  
19 weighed against the extra time and resources incurred. Machine-marked VSA questions  
20 should be reviewed by subject matter experts. With the assessment software used in this  
21 study, the total time taken to review the machine-marked answers to all 60 VSA questions  
22 for 299 students was under two hours. Questions that took longer to mark had a higher  
23 number of permutations of the correct answers making it hard to create a comprehensive  
24 answer key (for example: ultrasound left leg, Doppler USS LL, Left leg US). However, every  
25 time an unforeseen answer is marked as 'correct' by the examiner, the system will learn the  
26 variation. Therefore, in future uses of the question, the answers will already be present in the  
27 system making marking more efficient. With advances in computational linguistics and  
28 machine-learning across educational fields,<sup>30, 31</sup> it is likely that speed and reliability of  
29 marking short answer questions will continue to improve.

30 Limitations of our study include the sample size and inclusion of students from one centre.  
31 Whilst the students were randomly assigned to the two groups and there were no differences  
32 in sex distribution, the possibility of differences in other characteristics cannot be excluded.  
33 For example we did not have access to data on the cultural background of the students for  
34 inclusion in the differential item functioning. Only one VSA question showed differential item  
35 functioning with a bias against males. Interestingly this item was based on the presentation

1 of urinary tract infection in a young female. Another limitation of our study is the lack of data  
2 on clinical experience and anticipated specialty of those who agreed VSA questions were  
3 more representative of real-life practice. Furthermore, there was no external validation  
4 measure to assess and compare students' competence. Future studies should investigate  
5 the utility of VSA questions in multi-centre studies involving a larger number of students. It  
6 would also be interesting to examine the effects of cueing in candidates with varying levels  
7 of expertise.

8 Given their high correlation with short answer questions, SBA questions are widely used as  
9 an efficient assessment tool across medical education as a proxy for assessing applied  
10 knowledge. However, high correlations between assessments do not necessarily imply that  
11 the same cognitive facility is being tested. Indeed answer generation rather than recognition  
12 is tested in short answer questions.<sup>2, 4, 32</sup> Very short answer questions have efficiency and  
13 acceptability that is approaching that of SBA questions. Furthermore, compared with the  
14 SBA format, VSA questions demonstrate higher reliability, discrimination and authenticity.  
15 The results of this study demonstrate the utility of VSA questions as a an assessment  
16 instrument, which has the potential to improve existing assessment programs.

#### 17 18 **Conflict of interests**

19 The authors and their institutions do not have a financial interest in the software (*Practique*,  
20 *www.fry-it.com*) used in the study.

#### 21 22 **Acknowledgement**

23 We thank Mr Kean Schupke for his help with the provision of the *Practique* software.

#### 24 25 **References**

- 26 [1] Elstein AS. Beyond multiple-choice questions and essays: the need for a new way to  
27 assess clinical competence. *Acad Med*. 1993;**68**:244-249.
- 28 [2] Veloski JJ, Rabinowitz HK, Robeson MR, Young PR. Patients don't present with five  
29 choices: an alternative to multiple-choice tests in assessing physicians' competence. *Acad*  
30 *Med*. 1999;**74**:539-546.
- 31 [3] Shaibah HS, van der Vleuten CP. The validity of multiple choice practical  
32 examinations as an alternative to traditional free response examination formats in gross  
33 anatomy. *Anat Sci Educ*. 2013;**6**:149-156.
- 34 [4] Newble DI, Baxter A, Elmslie RG. A comparison of multiple-choice tests and free-  
35 response tests in examinations of clinical competence. *Med Educ*. 1979;**13**(4):263-268.
- 36 [5] Desjardins I, Touchie C, Pugh D, Wood TJ, Humphrey-Murto S. The impact of cueing  
37 on written examinations of clinical decision making: a case study. *Med Educ*. 2014;**48**(3):  
38 255-261.
- 39 [6] Schuwirth LWT, van der Vleuten CPM, Donkers HJLM. A closer look at cueing  
40 effects in multiple-choice questions. *Med Educ*. 1996;**30**(1):44-49.

- 1  
2  
3 1 [7] Schuwirth LWT, van der Vleuten CPM, Stoffers HEJH, Peperkamp AGW. Computerized  
4 2 long-menu questions as an alternative to open-ended questions in computerized  
5 3 assessment. *Med Educ*. 1996;**30**:50-55.  
6 4 [8] van der Vleuten CPM, Schuwirth LWT. Assessing professional competence:  
7 5 from methods to programmes. *Med Educ*. 2005;**39** (3):309-317.  
8 6 [9] Sam AH, Hameed S, Harris J, Meeran K. Validity of very short answer versus single  
9 7 best answer questions for undergraduate assessment. *BMC Med Educ*. 2016;**16**:266.  
10 8 [10] Epstein RM. Assessment in medical education. *N Engl J Med*. 2007;**356**:387-396.  
11 9 [11] Wass V, Van der Vleuten C, Shatzer J, Jones R. Assessment of clinical competence.  
12 10 *Lancet*. 2001;**357**:945-949.  
13 11 [12] Cilliers FJ, Schuwirth LW, van der Vleuten CP. A model of the pre-assessment  
14 12 learning effects of assessment is operational in an undergraduate clinical context. *BMC Med*  
15 13 *Educ*. 2012;**12**:9.  
16 14 [13] Al-Kadri HM, Al-Moamary MS, Roberts C, Van der Vleuten CP. Exploring  
17 15 assessment factors contributing to students' study strategies: literature review. *Med Teach*.  
18 16 2012;**34 Suppl 1**:S42-50.  
19 17 [14] Newble DI, Jaeger K. The effect of assessments and examinations on the learning of  
20 18 medical students. *Med Educ*. 1983;**17**(3):165-171.  
21 19 [15] Eagle M, Leiter E. Recall and Recognition in Intentional and Incidental Learning. *J*  
22 20 *Exp Psychol*. 1964;**68**:58-63.  
23 21 [16] McCoubrie P. Improving the fairness of multiple-choice questions: a literature review.  
24 22 *Med Teach*. 2004;**26**:709-712.  
25 23 [17] Newble DI, Entwistle NJ. Learning styles and approaches: implications for medical  
26 24 education. *Med Educ*. 1986;**20**(3):162-175.  
27 25 [18] Willing S, Ostapczuk M, Musch J. Do sequentially-presented answer options prevent  
28 26 the use of testwiseness cues on continuing medical education tests? *Adv Health Sci Educ*  
29 27 *Theory Pract*. 2015;**20**:247-263.  
30 28 [19] McConnell MM, St-Onge C, Young ME. The benefits of testing for learning on later  
31 29 performance. *Adv Health Sci Educ Theory Pract*. 2015;**20**:305-320.  
32 30 [20] Larsen DP, Butler AC, Roediger HL, III. Test-enhanced learning in medical  
33 31 education. *Med Educ*. 2008;**42**(10):959-966.  
34 32 [21] Wood T. Assessment not only drives learning, it may also help learning. *Med Educ*.  
35 33 2009;**43**(1):5-6.  
36 34 [22] McDaniel MA, Roediger HL, 3rd, McDermott KB. Generalizing test-enhanced  
37 35 learning from the laboratory to the classroom. *Psychon Bull Rev*. 2007;**14**:200-206.  
38 36 [23] Rodriguez MC. Construct Equivalence of Multiple-Choice and Constructed-Response  
39 37 Items: A Random Effects Synthesis of Correlations. *Journal of Educational Measurement*.  
40 38 2003;**40**:163-184.  
41 39 [24] Fenderson BA, Damjanov I, Robeson MR, Veloski JJ, Rubin E. The virtues of  
42 40 extended matching and uncued tests as alternatives to multiple choice questions. *Hum*  
43 41 *Pathol*. 1997;**28**:526-532.  
44 42 [25] Damjanov I, Fenderson BA, Veloski JJ, Rubin E. Testing of medical students with  
45 43 open-ended, uncued questions. *Hum Pathol*. 1995;**26**:362-365.  
46 44 [26] Case SM, Swanson DB. Extended-matching items: A practical alternative to  
47 45 free-response questions. *Teaching and Learning in Medicine*. 1993;**5**:107-115.  
48 46 [27] Scalise K, Gifford B. Computer-Based Assessment in E-Learning: A Framework for  
49 47 Constructing "Intermediate Constraint" Questions and Tasks for Technology Platforms. *The*  
50 48 *Journal of Technology, Learning, and Assessment*. 2006;**4**:1-44.  
51 49 [28] Van Der Vleuten CP. The assessment of professional competence: Developments,  
52 50 research and practical implications. *Adv Health Sci Educ Theory Pract*. 1996;**1**:41-67.  
53 51 [29] Gulikers JTM, Bastiaens TJ, Kirschner PA. A five-dimensional framework for  
54 52 authentic assessment. *Educational Technology Research and Development*. 2004;**52**:67-86.  
55 53 [30] Burrows S, Gurevych I, Stein B. The Eras and Trends of Automatic Short Answer  
56 54 Grading. *International Journal of Artificial Intelligence in Education*. 2015;**25**:60-117.  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1 [31] Pulman SG, Sukkariéh JZ. Automatic short answer marking. *Proceedings of the*  
2 *second workshop on Building Educational Applications, Ann Arbour, Michigan.* 2005:9-16.  
3 [32] Ozuru Y, Briner S, Kurby CA, McNamara DS. Comparing comprehension measured  
4 by multiple-choice and open-ended questions. *Can J Exp Psychol.* 2013;**67**:215-227.

5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

For Review

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Single best answer (SBA) question**

A 24 year old woman has two months of lethargy, dizziness, weight loss and nausea. She has type 1 diabetes and reports erratic blood sugars and one episode of loss of consciousness. She has hyperpigmentation in her palmar creases and her oral mucosa. Her temperature is 36.8°C, pulse rate 101 bpm, BP 78/61 mmHg (standing), respiratory rate 16 breaths per minute and oxygen saturation 99% breathing air. Her capillary blood glucose is 3.2 mmol/litre.

Investigations:		
Sodium	129 mmol/L	(135–146)
Potassium	5.4 mmol/L	(3.4–5.0)
Urea	7.7 mmol/L	(2.5–7.8)
Creatinine	67 µmol/L	(50–95)

What is the most likely diagnosis?

- A. Addison's disease
- B. Congenital adrenal hyperplasia
- C. Cushing's disease
- D. Hypothyroidism
- E. SIADH

**Very short answer (VSA) question**

A 24 year old woman has two months of lethargy, dizziness, weight loss and nausea. She has type 1 diabetes and reports erratic blood sugars and one episode of loss of consciousness. She has hyperpigmentation in her palmar creases and her oral mucosa. Her temperature is 36.8°C, pulse rate 101 bpm, BP 78/61 mmHg (standing), respiratory rate 16 breaths per minute and oxygen saturation 99% breathing air. Her capillary blood glucose is 3.2 mmol/litre.

Investigations:		
Sodium	129 mmol/L	(135–146)
Potassium	5.4 mmol/L	(3.4–5.0)
Urea	7.7 mmol/L	(2.5–7.8)
Creatinine	67 µmol/L	(50–95)

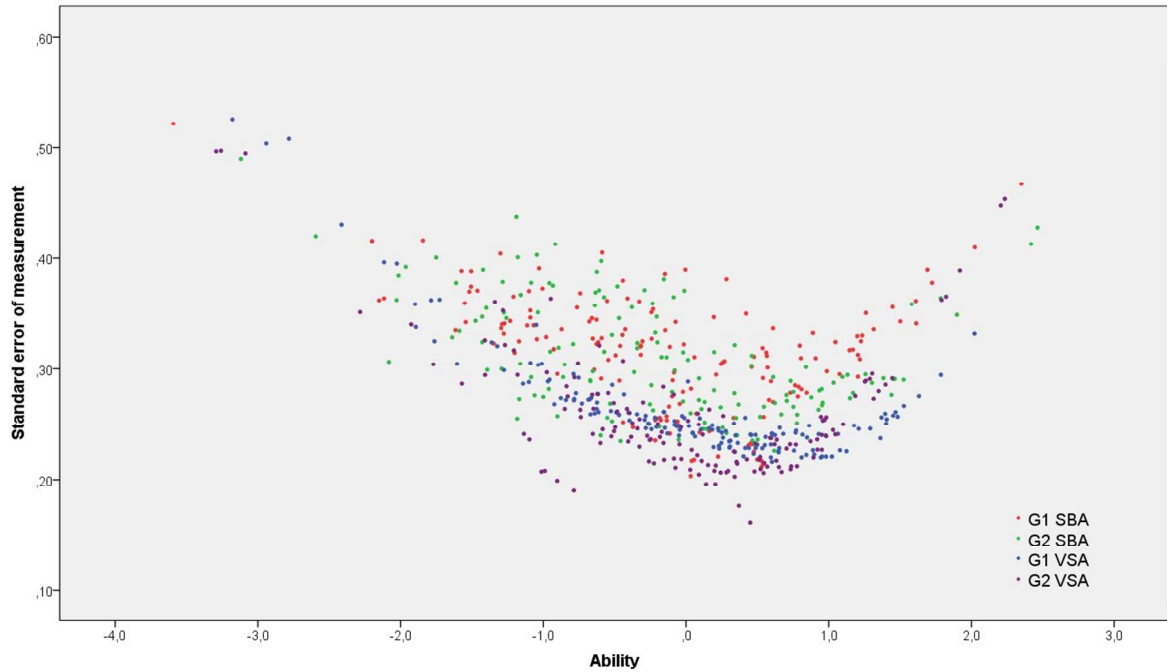
What is the most likely diagnosis?

- Correct answers:
- Addison's disease
  - Addison's
  - Adrenal insufficiency
  - Primary adrenal insufficiency
  - Hypoadrenalism

**Figure 1.** Example of a question written in both single best answer (SBA) format with five options (left) and very short answer (VSA) format with the acceptable variations of the correct answer that would automatically receive a mark (right).

Addison's disease	1.00
Primary adrenal insufficiency	1.00
Addisons disease	1.00
SIADH	0.00
Primary hypoadrenalism	1.00

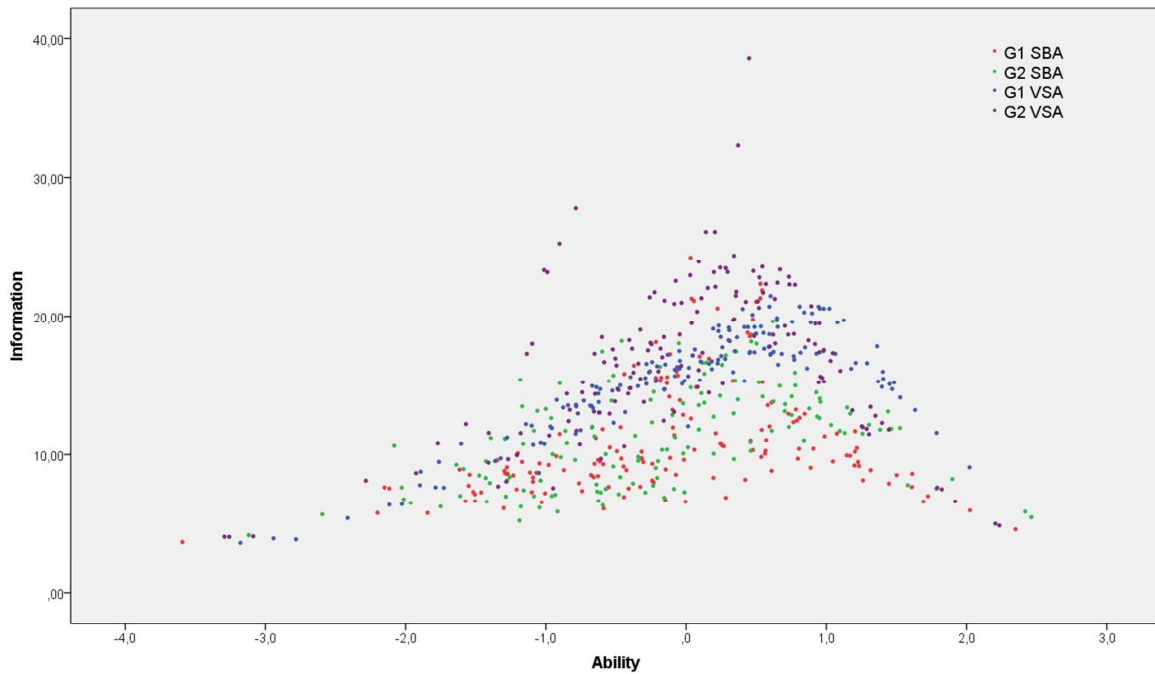
**Figure 2.** Example of marking of a very short answer question by computer. The green answers have been automatically assigned a mark (1.00) as they were included on the list of acceptable answers. The yellow answers have been marked as correct based on their similarity to the acceptable answers. Answers marked by computer as incorrect are shown in red, however during the review process this can be over-ridden (e.g. 'primary hypoadrenalism'), with all identical answers automatically receiving a mark.



1

2 **Figure 3.** Individual estimates for the standard error of measurement according to theta  
 3 ability estimates, and item type (single best answer: SBA, very short answer: VSA) and  
 4 group (G) combination.

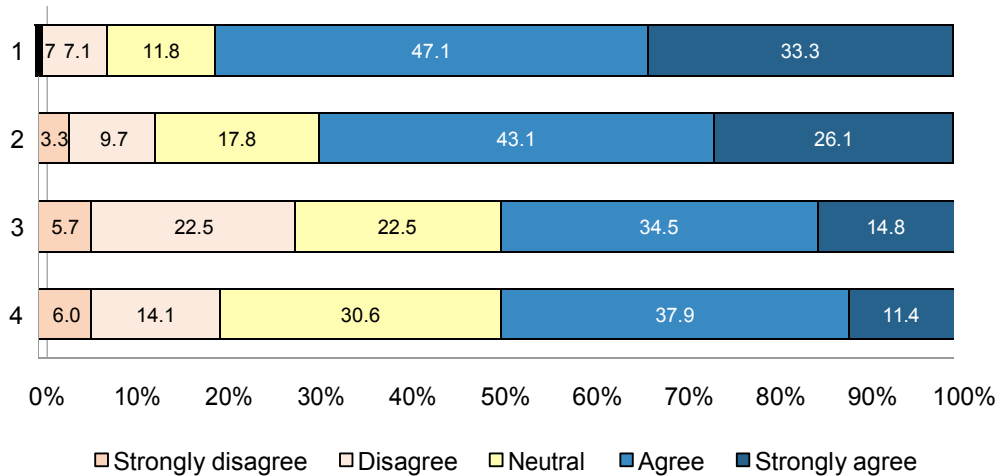
5



6

7 **Figure 4.** Individual estimates for information according to theta ability estimates, and item  
 8 type (single best answer: SBA, very short answer: VSA) and group (G) combination.





**Figure 5.** Students were asked to rate their agreement with each statement on a 5-point Likert scale (strongly disagree, disagree, neutral, agree, strongly agree). For each of the following four statements, the percentage of students selecting each point on the Likert scale is shown.

1. Questions in the single best answer format are easier than the very short answer format.
2. Very short answer questions are a better representation of how I would be expected to answer questions in clinical practice.
3. Having examinations in very short answer format would change my learning and revision strategy.
4. Using very short answer questions in assessments would help improve my preparation for clinical practice.

	Group 1 (n=155)		Group 2 (n=144)	
	VSA	SBA	SBA	VSA
Mean (%)	52.4	68.2	69.7	65.7
SD %	17.4	12.5	12.9	16.5
Cronbach's alpha	0.91	0.84	0.85	0.91
SEM	5.235	5.03	4.970	5.09
Mean item-total score point-biserial correlation	0.36	0.26	0.27	0.35

**Table 1.** The mean raw scores, standard deviation (SD), Cronbach's alpha, standard error of measurement (SEM) and mean item-total score point-biserial correlations for the very short answer (VSA) and single best answer (SBA) tests in groups 1 and 2.