

A modified electronic key feature examination for undergraduate medical students: validation threats and opportunities

MARTIN R. FISCHER¹, VERONIKA KOPP¹, MATTHIAS HOLZER¹, FRANZ RUDERICH² & JANA JÜNGER²

¹Klinikum der Universität München, Medizinische Klinik-Innenstadt, München, Germany;

²Med. Klinik II der Universität Heidelberg, Heidelberg, Germany

SUMMARY *The purpose of our study was the development and validation of a modified electronic key feature exam of clinical decision-making skills for undergraduate medical students. Therefore, the reliability of the test (15 items), the item difficulty level, the item-total correlations and correlations to other measures of knowledge (40 item MC-test and 580 items of German MC-National Licensing Exam, Part II) were calculated. Based on the guidelines provided by the Medical Council of Canada, a modified electronic key feature exam for internal medicine consisting of 15 key features (KFs) was developed for fifth year German medical students. Long menu (LM) and short menu (SM) question formats were used. Acceptance was assessed through a questionnaire. Thirty-seven students from four medical schools voluntarily participated in the study. The reliability of the key feature exam was 0.65 (Cronbach's alpha). The items' difficulty level scores were between 0.3 and 0.8 and the item-total correlations between 0.0 and 0.4. Correlations between the results of the KF exam and the other measures of knowledge were intermediate (r between 0.44 and 0.47) as well as the learners' level of acceptance. The modified electronic KF examination is a feasible and reliable evaluation tool that may be implemented for the assessment of clinical undergraduate training.*

Introduction

Decision-making skills are important for undergraduate medical students (Bordage, 1994; Hatala & Norman, 2002). The purpose of this study was to develop a feasible, reliable, and valid electronic examination for decision-making skills and to evaluate students' acceptance of this test.

Key feature approach

One approach used for assessing the acquisition of decision-making skills is the key feature approach (Bordage *et al.*, 1995; Page & Bordage, 1995; Page *et al.*, 1995). A *key feature* is defined "as a critical step in the resolution of a problem. Two corollaries were added to the general definition of a key feature: (1) it focuses on a step in which examinees are most likely to make errors in the resolution of the problem, and (2) it is a difficult aspect of the identification and management of the problem in practice" (Page & Bordage, 1995). Problems embedding a key feature (KF), referred to as *key feature problems*, consist of a brief clinical stem followed by one or more questions, the KFs.

Hatala & Norman (2002) introduced the KF assessment method in an undergraduate setting when evaluating the clinical decision-making skills of internal medicine clerks. They reported a Cronbach's Alpha of 0.49 with a 15 KF exam in a paper and pencil format.

Question format

The KF approach recommends two types of questions, which were used in the study from Page *et al.* (1995): short-answer "write-in" and the "short menu" (SM) format. In the short menu, the examinees have to select their responses from prepared lists of options. These lists can range between 2 and 45 options. Thereby, it seems important that all keyed responses, synonyms, and incorrect responses including common misconceptions are provided in the list in order to reduce cueing effects (Schuwirth *et al.*, 1996a). In the write-in format, the students have to supply their responses as free text entries. As an alternative to the write-in format, there is the "long-menu" (LM) format. Long menus are alphabetically ordered long lists of possible answers (over 500). Used in a paper and pencil exam, this format is very time consuming and error-prone (Case *et al.*, 1994). By using a computerized assessment tool, these difficulties were overcome (Schuwirth *et al.*, 1996b). In our study the students selected an answer from the list by typing it into a dialogue box. The computer then searches through the LM-list for "hits". The alternatives found are reported back to the examinee immediately so he can check whether the retrieved option is the desired one. It was shown that here are only negligible differences concerning performance, cueing and reliability between the LM format and the write-in format (Schuwirth *et al.*, 1996b).

However, an electronic format of a KF examination including LM questions has not been studied yet. We selected an undergraduate setting because of the importance of clinical decision-making skills in an early stage of medical education and the need for efficient tools to assess it (Groves *et al.*, 2003).

Furthermore, the students' performance and their opinion concerning the electronic KF exam were assessed, as the latter is an often neglected research question (Ogilvie *et al.*, 1999).

Correspondence: Martin R. Fischer, Klinikum der Universität München, Medizinische Klinik-Innenstadt, Ziemssenstr. 1, 80336 München, Germany. Email: fischer.martin@med.uni-muenchen.de

Methods

Participants

Thirty-seven students (22 female, 15 male) participated in the study. On average, the students were about 26 years old (ranging between 24 and 41). Participation was voluntary and participants were financially reimbursed. The study was conducted in July 2003 at four German universities (Freiburg, Heidelberg, Munich (LMU), and Ulm).

Test procedure

After a short introduction to the test procedure and the test-tool, the students had to process the MC-test. Subsequently, the examinees had to answer 15 KF problems, each consisting of 3 to 5 KFs leading to a total of 60 KFs. At the end, they had to fill in a questionnaire (Table 1). About three weeks later, all participants had to take the second part of the German national medical licensing examination (NBE), thus ensuring a high level of knowledge activation.

Instruments

Assessment tool. As a testing-tool, the computer-based training system CASUS was used (Fischer, 2000). The client-server architecture of the computer system uses a standard web-browser as the user interface. Interactivity is handled by servlets and Java script. Each participant was assigned to a unique login and password. Each KF could only be answered once, as the answer was mostly contained in the next KF. Backward navigation was only possible to review information, not for editing. Thus, items were dependent and this test format is therefore technically difficult to apply hardly applicable in a paper and pencil format. All answers were centrally recorded and scored from the server.

Key feature problems

Test length. In order to get close to the recommended examination reliability of 0.8 for summative assessment tools, Page & Bordage (1995) calculated that a KF exam should contain 40 KF problems and could be completed in 4.1 h. As this seemed not feasible for undergraduate students due to logistical limitations, we reduced the examination to 15 KF problems like Hatala & Norman (2002). Concerning the time to answer one KF, we acted in accordance with the time limits given in the NBE, which is 90 seconds per question.

Domain definition and examination blueprint. As content domain we selected general internal medicine. In the second NBE, 15.5% or 90 out of a total of 580 MC-questions concern internal medicine as a core subject of the clinical curriculum.

The Swiss one-dimensional blueprint for internal medicine was used (IAWF, 1999) for weighing of subdomains in internal medicine (Table 2). The content relation

Table 1. Procedure of the study.

Procedure	Technical Introduction	MC-test (40 questions)	15 KF problems	Acceptance questionnaire
Time	10 min.	60 min.	90 min.	10 min.

of our KF problems to these subdomains is fitting these recommendations.

Key features

Fifteen KF problems were written at the University departments of internal medicine in Heidelberg and Munich and were reviewed and adapted by four physicians with long-term clinical expertise in general internal medicine from both institutions who were not involved as authors. Each problem contained between 3 and 5 KFs. As a total, the students had to answer 60 KFs (30 in the SM format, 30 in the LM format).

Scoring keys

Since differential weighting of responses does not improve score reliability (Page & Bordage, 1995), each KF question was scored using a dichotomous scoring system, with the problem score being the average of the questions scores. As one KF problem was interpreted as one item, the maximum score was 15 points.

Multiple choice test

As mentioned before, the written part of the second German NBE consists exclusively of A-type MC-questions (one answer has to be selected from five alternatives). A multiple choice test, containing 40 MC-questions from former second NBEs on internal medicine, was composed to assess factual knowledge. The questions from internal medicine together with their psychometric characteristics were provided by the German National Institution for Medical and Pharmaceutical Examinations Institute (IMPP). It was estimated that 40 questions were needed to reach an examination reliability of 0.7 (Cronbach's alpha).

Second national medical licensing examination

The German NBE part II consists of a written and an oral part. In the written part, the students had to answer 580 MC-questions; 90 of them related to internal medicine. A proportion of about 15% of the NBE questions had been used previously. For the rest of the questions the psychometric characteristics are not known as piloting is impossible for legal reasons.

Table 2. Subjects of the KF problems and their weight in the exam in percent (Swiss Blueprint as reference in parenthesis).

Subject	%
Rheumatology/immunology	13 (13)
General internal medicine/neurology/ psychosomatic medicine/alcoholism/miscellaneous	14 (20)
Cardiology/angiology/hypertension	14 (15)
Endocrinology/nephrology	13 (11)
Gastroenterology	13 (10)
Oncology/infectiology	13 (15)
Pneumology/emergency medicine	20 (16)

Questionnaire

The questionnaire consisted of 2 essay questions and 22 scaled items concerning the examinees' acceptance of the KFs, the computerized test format, and its usability. A Likert scale from 1 (total disagreement) to 5 (total agreement) was used.

Question format

The SM and the LM question format were used. In the SM format, the students had to tick the right answer(s) out of the given alternatives (between 5 and 27). In the LM format, the students had to choose one right answer out of a long menu. The list was generated by filtering out the 500 most commonly used answer terms for the NBE internal medicine questions. Beyond, the right answers with synonyms and the relevant distractors of all 15 KFs were added to the list, leading to about 700 terms in our study.

Analysis. Following Bordage's advice (personal communication, 2003), the reliability of the test was calculated by taking each KF problem as one item. The Cronbach's alpha was computed as the equivalent of a repeated measures ANOVA (Hatala & Norman, 2002).

Validity is the sine qua non of assessment, but also the most difficult criterium to assess (Downing, 2003). We focused on the item difficulty level and the item-total correlations, and on the relationship to other variables as a sort of criterium related validity. As external criteria other measures of knowledge (MC-test and the second NBE) were used. In order to compare the results of these different measures of knowledge, an analysis of variance

with repeated measures and *t*-test analyses were calculated. For the analysis of the questionnaire, descriptive statistics were used.

Results

Feasibility

The feasibility of the online KF exam was impaired only by temporary uncontrolled network traffic interfering with the performance of the server in one instance. Nevertheless, this examinee could finish the test in time. Other major technical problems did not occur.

Reliability

The KF problems had a reliability score of 0.65 (Cronbach's alpha). Extrapolated to 25 KF-problems, a reliability of 0.75 can be expected.

The reliability of the 40 item MC-test was 0.71 (Cronbach's alpha), 0.97 (Cronbach's alpha) in the NBE, and 0.80 (Cronbach's alpha) in the 90-item part concerning internal medicine, respectively.

Item difficulty level and item-total correlations

An item difficulty level between 0.2 and 0.8 is recommended (Bortz & Döring, 2002) to differentiate between high and low achieving students and the item-total correlations should reach positive scores. For the 15 KF-items the item difficulty level was between 0.3 and 0.8. Seven items reached middle correlations, between 0.3 and 0.6; five items were only slightly correlated with the total test score. Three items did not correlate with the other items. In Table 3 and Figure 1

Table 3. Item difficulty level and item-total correlations (1 = item difficulty level; 2 = item total correlation).

	KF1	KF2	KF3	KF4	KF5	KF6	KF7	KF8	KF9	KF10	KF11	KF12	KF13	KF14	KF15
1	0.42	0.52	0.64	0.58	0.60	0.32	0.80	0.61	0.54	0.83	0.68	0.39	0.34	0.47	0.67
2	0.42	0.24	0.19	0.06	0.19	0.21	0.14	0.25	0.38	0.42	0.37	0.14	0.29	0.19	0.56

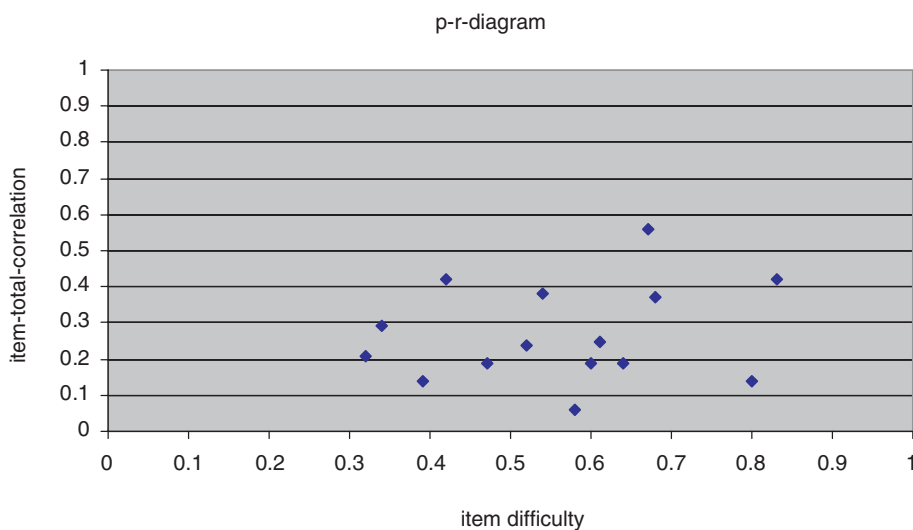


Figure 1. The p-r-diagram of the item-total correlation and the item difficulty.

the item difficulty and the item-total correlations for each item are shown.

Performance of the students in the KF exam and other measures of knowledge

In the KF exam, the students scored 8.4 points (SD = 1.32) or 56% on average. In the other measures of knowledge the students reached scores between 77 and 79% (Table 4). The differences between the KF exam and the other measures of knowledge (the MC-test, the NBE and the part of the NBE concerning internal medicine) are significant ($t(36) = 12.1$; $p < 0.01$; $t(35) = 15.0$; $p < 0.01$; $t(36) = 16.5$; $p < 0.01$).

Students reached on average 44% ($M = 13.3$ (SD = 3.56)) in the SM subscale and 68% ($M = 20.5$ (SD = 2.78)) in the LM subscale, respectively. The divergence between these two subscales is highly significant ($t(36) = 12.6$; $p < 0.01$).

Correlations between the key feature examination and other measures of knowledge

The correlations between the KF examination and the other measures of knowledge were intermediate (Table 5).

Students' acceptance of the computerized test-setting

The students had to evaluate the KFs concerning their problem-orientation, authenticity and interdisciplinarity. The results from the questionnaire concerning acceptance are displayed in Table 6.

Besides the closed questions, there were two open questions asking for positive and negative statements. The answers to these questions were not limited in space. Answers were categorized and summarized as follows:

Positive aspects. The most positive statements (7) were made in a general manner, like "Overall, I liked it". Also, seven statements emphasized the case- and problem-orientation or authenticity of the problems and their clinical relevance. Five liked the efficiency given through working with the computer, because the transfer from the booklets to coding sheets is omitted. Four students especially liked doing tests on a computer and three students wrote that they liked being given the right answer together with the next question. Four students liked the fact they could learn from the given answers provided with the next question.

Negative aspects. Most criticism was aimed at the LM question format. The students complained that they had problems finding the right answer term and that the list did not contain all possible answers they were looking for (10 statements). Besides general statements saying that they did not like the multiple response question format, they mentioned that there were too many possible alternatives for answering (five statements).

Some criticism was related to the system's user interface and to the system itself. Eight students complained that the font size and the size of the pictures were too small.

A further point of criticism concerned the KF problems: first of all that there is no possibility of changing the given

Table 4. Students' performance in the KF exam and other measures of knowledge.

	Key features	MC-test	Written NBE	NBE – internal medicine
Max. score	15	40	580	90
M (SD)	8.4 (1.32)	30.8 (4.5)	459.7 (40.5)	70.8 (7.0)
Percent	56%	77%	79%	78%

Table 5. Correlations between the KF exam and other measures of knowledge and students' performance in the tests.

	MC-test	National boards (580 questions)	NBE (90 questions internal medicine)
KF exam	0.47**	0.44**	0.46**

**The correlation was significant on a 0.01 level (two-tailed test).

answers at a later point in time (4 statements). Further criticism concerned details of the KFs' formulation and text length.

Discussion

The purpose of the study was to develop a feasible, reliable and valid electronic examination for decision-making skills used for the medical undergraduate assessment. Based on the KF approach, a 90 minute, 15 KF exam was generated. Technically the online format of the test seemed to be feasible for broader use. Our experiences speak in favour of a separated assessment intranet to be protected against unpredictable bandwidth. We were able to achieve a surprisingly high reliability of 0.65 (Cronbach's alpha) compared to Hatala and Norman (2002), who developed a written test of clinical decision-making with KFs used to evaluate a clinical clerkship in internal medicine. Within this 2-hour exam with 15 problems, they achieved an overall test reliability of 0.49.

Referring to item difficulty, only acceptable scores between 0.3 and 0.8 were achieved. Concerning item-total correlation, the most items could reach at least small correlations. Three KF problems did not contribute to the differentiation between high and low achieving students. When analysing the KFs within these problems it became clear that they were either too easy or contained cueing answers. However, the content validation process showed a high relevance for these specific questions. The analysis of the item-total correlation and the difficulty level in respect to content relevance will lead to a change of some questions for the next version of our KF-test. These predominantly positive results contribute to the validity of the exam. The correlations between the KF exam's results and the MCQ-tests are moderate ranging between 0.44 and 0.47. This could be related to differential validity and supports the intention of measuring two different kinds of knowledge: factual knowledge evaluated by the MCQ-tests and decision making skills measured by the KF-exam. Although, the KFs

Table 6. Results of the questionnaire ($n = 37$).

	M	SD	Min	Max
Getting to know this way of assessment (case- and computer-based) was interesting to me.	4.31	0.75	2	5
Overall the examination was fun.	3.78	0.90	2	5
I felt the difficulty level of the assessment was appropriate.	3.53	0.85	2	5
The key features are problem-oriented.	4.06	0.67	3	5
The key features are interdisciplinary.	3.47	0.94	1	5
The key features are close to problems from clinical practice.	3.92	0.87	2	5
I enjoyed working on the short cases.	2.97	1.40	1	5
Working on the short key features is a useful way of assessing my knowledge.	3.25	1.44	1	5
I wish to have an examination with key features in the future curriculum.	2.80	1.43	1	5
The selection of answers form a long list is a good compromise between MC-answers and free text answers.	2.56	1.46	1	5
The time frame for the MC-questions was appropriate.	4.78	0.49	3	5
The time frame for the key features was appropriate.	4.22	1.12	1	5
The online format of the examination was appealing to me.	3.52	1.01	1	5
Working on the key features and questions on the computer was more strenuous than compared to paper and pencil examinations.	3.08	1.27	1	5
Assessment should be preferably done online in the future.	2.83	1.30	1	5
This test examination was a valuable preparation for me for the second national board exam.	2.97	0.88	1	5
The planned time schedule of the examination ran without any problems.	3.31	1.19	1	5
The software ran smoothly without technical problems.	3.29	1.30	1	5
The screen design was appropriate for the conduction of an online examination.	3.67	1.12	1	5
The images were of sufficient quality	3.64	1.31	1	5
The text was legible.	4.25	1.05	1	5
My preferred answer was mentioned in the long-menu list.	2.67	0.93	1	5

went through repeated review loops, a more careful content validation process is necessary (Bordage *et al.*, 1995). A further limitation, which should be considered in further studies, is the small sample of voluntary students with a potential selection bias. Similarly, the computer literacy of participants should additionally be objectively assessed.

A further question is the students' performance in the KF examination. The results demonstrated that they reached an average score of 56%. In this respect, the difficulty of the KFs was appropriate for students passing the second NBE. There were neither ceiling nor floor effects. But compared with their results in different MC-tests, where they reached nearly 80%, this result is surprisingly low.

By taking a closer look at the results of the KF exam, it has to be determined that the students performed very differently depending on the question format: in the KFs in LM format, they reached nearly 70%, whereas in the SM format, they achieved only 44%. They performed significantly worse in the KFs with SM format than in all other tests. This is an interesting finding, because we expected that the students would perform in the SM KFs as well as in LM KFs. That the students performed best in the MC-tests, is not astonishing, as they are used to this kind of testing and can be assumed to be well prepared for these exams. But why do they perform better in answering KFs in the LM format than in the SM format?

One possible answer is that the LM did not sufficiently contain all possible distractors and common misconceptions, so that the students rethought and chose one of the right ones. This possible explanation is supported by the results from the questionnaire. One question aimed at if the students could always find the answer they were looking for in the long list. The mean score of the answers to this question was only moderate. Furthermore, 10 students complained, that they had problems finding the right answer in the long list or that the list was not complete. These findings could be interpreted as a lack of distractors.

Another possible explanation for this is the following: The question format had an influence on the way of generating and formulating as well as on the content of a question. It is possible that the questions posed in the LM format were easier for the students to answer than the questions in the SM format. Due to the wish of KF authors, we did not reinforce the recommendations of Bordage for the selection of answer-types.

Both explanations would have influenced the results. This brings out the importance of a correct and complete list and indicates how difficult it is to generate good KFs. Regardless of these possible positive influences, the students performed lower in the KF than in the MCQ exam. Possibly, the students' decision-making skills were not sufficient and need improvement in the future. It also has to be noted, that the

participants could not prepare for the KF exam format. In contrast, they excessively trained for the MC-questions of the NBE. Thus, the novelty effect of the KF exam may contribute to performance differences.

The last question concerned the students' opinion about the KF exam and its realization on the computer. Although the students perceived the KF problems as authentic, problem-based and interdisciplinary, and felt that the KF exam in itself was a reasonable approach, they were quite sceptical concerning the implementation of such examinations.

Concerning the use of computerized exams, there was only a moderate acceptance, though the questions concerning usability were responded to positively. These results differ a little from those reported in the literature, where residents and medical students like the use of computer administered examinations (Butzin *et al.*, 1984; Legler & Realini, 1994; Ogilvie *et al.*, 1999). Furthermore, the students found working on the computer a little bit more strenuous than working with paper and pencil. To find a solution for this problem, more investigations are necessary. The system's user interface was adequate when concerning its use in examinations. Nevertheless, there are points of criticism, which are very helpful for a further development and should be implemented in future versions.

In the light of an increasing workload related to valid assessment strategies, online tools with relevant contents and formats are a key factor for successful curriculum reform. Therefore, further studies are needed to learn more about the characteristics of online KF exams.

To sum up, the developed modified electronic KF examination promises to be a reliable assessment tool for faculty wide summative evaluations when the number of items is adequately high. Further studies for a broader validation of this assessment format and its technical feasibility are needed with special respect to the LM answer format.

Practice points

- A 15-item electronic key feature examination is feasible for undergraduate context and has a reliability of 0.65 (Cronbachs alpha), potentially useful for high-stakes student evaluations.
- Key feature problems correlate only moderately with measures of factual knowledge (classical MC-examinations).
- Acceptance of online KF exams is intermediate and could be improved by the curricular integration of electronic case-studies.

Acknowledgements

We are indebted to Matthias Angstwurm, Andreas Eigler, Roland Gärtner and Jochen Schopohl for case review; to Dietmar Neuman and Dr Christian Götz for statistical

support; to Reiner Singer, Martin Adler, and Benjamin Körner for technical support and Sibyl Hermann, Magnus Müller, Claudia Seydi for organizational support. This study was financially supported by the German Ministry for Education and Research (BMBF) within the CASEPORT-Project (FKZ 08 NM 111).

Notes on contributors

MARTIN R. FISCHER

VERONIKA KOPP

MATTHIAS HOLZER

FRANZ RUDERICH

JANA JÜNGER

References

- BORDAGE, G. (1994) Elaborated knowledge: a key to successful diagnostic thinking, *Academic Medicine*, 69, pp. 883–885.
- BORDAGE, G., BRAILOVSKY, C., CARRETIER, H. & PAGE, G. (1995) Content validation of key features on an national examination of clinical decision-making skills, *Academic Medicine*, 70, pp. 276–281.
- BORTZ, J. & DÖRING, N. (2002) *Forschungsmethoden und Evaluation*, 3rd edn (Berlin, Springer).
- BUTZIN, D.W., FRIEDMAN, C.P. & BROWNLEE, R.C. (1984) A pilot study of microcomputer testing in pediatrics, *Medical Education*, 18, pp. 339–342.
- CASE, S.M., SWANSON, D.B. & RIPKEY, D.R. (1994) Comparison of items in five-option and extended-matching formats for assessment of diagnostic skills, *Academic Medicine*, 69, pp. S1–3.
- DOWNING, S.M. (2003) On the meaningful interpretation of assessment data, *Medical Education*, 37, pp. 830–837.
- FISCHER, M.R. (2000) CASUS – An authoring and learning tool supporting diagnostic reasoning, in: C. DAETWYLER (Ed.) *Use of Computers in Medical Education (Part II)*. *Zeitschrift für Hochschuldidaktik*, 1, pp. 87–98.
- GROVES, M., O'ROURKE, P. & ALEXANDER, H. (2003) Clinical reasoning: the relative contribution of identification, interpretation and hypothesis errors to misdiagnosis, *Medical Teacher*, 25, pp. 621–625.
- HATALA, R. & NORMAN, G. (2002) Adapting the key features examination for a clinical clerkship, *Medical Education*, 36, pp. 160–165.
- IAWF (1999) *Kompetent prüfen, Handbuch zur Planung, Durchführung und Auswertung von Facharztprüfungen* (Bern).
- LEGLER, J.D. & REALINI, J.P. (1994) Computerized student testing as a learning tool in a family practice clerkship, *Family Medicine*, 26, pp. 14–17.
- OGILVIE, R.W., TRUSK, T.C. & BLUE, A.V. (1999) Students' attitudes towards computer testing in a basic science course, *Medical Education*, 33, pp. 828–831.
- PAGE, G. & BORDAGE, G. (1995) The medical council of Canada's key features project: a more valid written examination of clinical decision-making skills, *Academic Medicine*, 70, pp. 104–110.
- PAGE, G., BORDAGE, G. & ALLEN, T. (1995) Developing key-feature problems and examinations to assess clinical decision-making skills, *Academic Medicine*, 70, pp. 194–201.
- SCHUWIRTH, L., VAN DER VLEUTEN, C.P. & DONKERS, H.H. (1996a) A closer look at cueing effects in multiple-choice questions, *Medical Education*, 30, pp. 44–49.
- SCHUWIRTH, L., VAN DER VLEUTEN, C.P., STOFFERS H.E., PEPEKAMP, A.G. (1996b) Computerized long-menu questions as an alternative to open-ended questions in computerized assessment, *Medical Education*, 30, pp. 50–55.

AUTHOR QUERIES

Journal id: CMTE-107830

Query number	Query
--------------	-------

- 1 Please supply notes on contributors in usual journal style.