



The design and evaluation of a computerized adaptive test on mobile devices

Evangelos Triantafyllou *, Elissavet Georgiadou, Anastasios A. Economides

University of Macedonia, Egnatia 156 Thessaloniki 54006, Greece

Received 3 September 2006; received in revised form 13 November 2006; accepted 8 December 2006

Abstract

The use of computerized adaptive testing (CAT) has expanded rapidly over recent years mainly due to the advances in communication and information technology. Availability of advanced mobile technologies provides several benefits to e-learning by creating an additional channel of access with mobile devices such as PDAs and mobile phones. This paper describes the design issues that were considered for the development and the implementation of a CAT on mobile devices, the computerized adaptive test on mobile devices (CAT-MD). Throughout the development of the system, formative evaluation was an integral part of the design methodology. The recommendations, suggestions and the results of the formative evaluation were used to improve the system in order to make the assessment procedure more effective and efficient. These results of the formative evaluation are also presented here.

© 2006 Elsevier Ltd. All rights reserved.

Keywords: Computerized adaptive testing; Mobile learning; Formative evaluation

1. Introduction

Recent advances in computer technology and psychometric theories have accelerated the change of test format from conventional paper-and-pencil tests to computerized adaptive testing (CAT). A CAT is a method for administering tests that adapts to the examinee's ability level. Moreover, in recent years, much attention has been paid to mobile computing. Availability of advanced mobile technologies have started to extend e-learning by adding an extra channel of access for mobile users with mobile devices such as hand phones, personal digital assistants (PDAs) or pocket PCs. Although mobile computing has become an important and interesting research issue, little research has been done regarding the implementation of CAT using mobile devices, and this is the focus of our research.

Our research is an attempt to examine the design and development issues that may be important in the implementation of a CAT on different mobile products such as mobile phones and PDAs. As a case study

* Corresponding author.

E-mail addresses: vtrianta@csd.auth.gr (E. Triantafyllou), elisag@otenet.gr (E. Georgiadou), economid@uom.gr (A.A. Economides).

an educational assessment prototype was developed, called computerized adaptive test on mobile devices (CAT-MD). A formative evaluation of the CAT-MD was carried out in order to investigate the effectiveness and efficiency of the system and also to assess its usability and appeal.

This article describes the design, the development and the formative evaluation of the CAT-MD. It first discusses the design issues that were considered for the development of the computerized adaptive test. Next, it describes the architecture and the implementation of the computerized adaptive test on a PDA and, finally, it presents the description and the results of the formative evaluation.

2. Computerized adaptive test

Testing is one of the most widely used tools in higher education. The main goal of testing is to measure student knowledge level in one or more concepts or subjects, i.e. in pieces of knowledge that can be assessed. Since education was established as an institution, different methods of assessment were used in different contexts, such as class presentations, essays, projects, practicum, etc. However, the most common tools of measuring performance are the oral test and the ‘paper-and-pencil’ test. Given that the computer has been an educational tool over the last few decades and its use has spread rapidly in all levels of education and training, the use of computer-based tests (CBTs) has increased significantly over the last few years. CBTs became feasible for licensure, certification and admission.

The most common type of CBT is the linear one that is a fixed-length computerized assessment that presents the same number of items to each examinee in a specified order and the score usually depends on the number of items answered correctly. A linear test consists of a full range of easy and difficult test items that are either randomly selected from a larger pool or are the same for all examinees. Evidently the type of CBT described here imitates a ‘paper-and-pencil’ test that is presented in a digital format and pays little or no attention to the ability of each individual examinee.

By contrast, in computerized adaptive testing (CAT), a special case of computer-based testing, each examinee takes a unique test that is tailored to his/her ability level. As an alternative to giving each examinee the same fixed test, CAT item selection adapts to the ability level of individual examinees. After each response the ability estimate is updated and the next item is selected such that it has optimal properties according to the new estimate (van der Linden & Glas, 2003). The CAT first presents an item of moderate difficulty in order to initially assess each individual’s level. During the test, each answer is scored immediately and if the examinee answers correctly, then the test statistically estimates her/his ability as higher and then presents an item that matches this higher ability. If the next item is again answered correctly, it re-estimates the ability as higher still and presents the next item to match the new ability estimate. The opposite occurs if the item is answered incorrectly. The computer continuously re-evaluates the ability of the examinee until the accuracy of the estimate reaches a statistically acceptable level or when some limit is reached; such as a maximum number of test items. The score is determined from the level of the difficulty, and as a result, while all examinees may answer the same percentage of questions correctly, the high-ability ones will get a better score as they correctly answer more difficult items.

Regardless of some disadvantages reported in the literature – for example, high cost of development, item calibration, item exposure (Boyd, 2003; Eggen, 2001), the effect of a flawed item (Abdullah, 2003), or the use of CAT for summative assessment (Lilley & Barker, 2002, 2003) – CAT has several advantages. Testing on demand can be facilitated so as an examinee can take the test whenever and wherever s/he is ready. Multiple media can be used to create innovative item formats and more realistic testing environments. Other possible advantages are flexibility of test management; immediate availability of scores; increased test security; increased motivation etc. However, the main advantage of CAT over any other computerized based test is efficiency. Since fewer questions are needed to achieve a statistically acceptable level of accuracy, significantly less time is needed to administer a CAT compared to a fixed-length computerized based test (Linacre, 2000; Rudner, 1998).

Since the mid-80s when the first CAT systems became operational, i.e. the armed services vocational aptitude battery for the US department of defense account (van der Linden & Glas, 2003) using adaptive techniques to administer multiple choice items, much research and many technical challenges have made new assessment tools possible. Current research in CAT is not limited to educational admissions, but focuses

on applications that address self-assessment, training, employment, professional teacher development for schools, industry, military, the assessment of non-cognitive skills etc.

3. Mobile learning

In recent years the use of different mobile products such as mobile phones and Personal Digital Assistant (PDA) devices has increased rapidly. Moreover, much attention has been paid to mobile computing within the Information Technology industry. The availability of advanced mobile technologies, such as high bandwidth infrastructure, wireless technologies, and handheld devices, has started to extend e-learning towards mobile learning (Sharples, 2000). Mobile learning (m-learning) intersects mobile computing with e-learning; it combines individualized (or personal) learning with anytime/anywhere learning. The advantages of m-learning include: flexibility, low cost, small size, ease of use and timely application (Jones & Jo, 2004).

M-learning is often thought of as a form of e-learning, but Georgiev, Georgieva, and Smrikarov (2004) believe it would be more correctly defined as a part, or sub-level, of e-learning. They believe m-learning is a new stage in the progress of e-learning and that it resides within its boundaries. M-learning is not only wireless or Internet based e-learning but should include the anytime/anyplace concept without permanent connection to physical networks (Jones & Jo, 2004). The place independence of mobile devices provides several benefits for the m-learning environment, such as allowing students and instructors to utilize their spare time while travelling in a train or bus to finish their homework or lesson preparation. Furthermore, m-learning has the potential to change the way students behave and interact with each other (Luvai, in press).

The introduction of mobile devices into the learning process can complement e-learning by creating an additional channel of assessment with mobile devices such as PDAs, mobile phones, portable computers. Due to their convenient size and reasonable computing power, mobile devices have emerged as a potential platform for computer-based testing. Although there is previous research evaluating the impact of computer-based testing (using desktop systems), there is little research on the impact of using mobile devices to administer tests since handheld and wireless technologies are relatively new and continually evolving. Most of the literature is focused on practical topics concerned with the development and implementation of new educational applications that can be delivered with these devices (Chen, Myers, & Yaron, 2000; Cook, 2000). However, Segall, Doolen, and Porter (2005) conducted research in order to compare the usability effectiveness, efficiency, and satisfaction of a PDA-based quiz application to that of standard paper-and-pencil quizzes in a university course. Their study showed that the PDA-based quiz is more efficient since students completed it in less time than the paper-and-pencil quiz. Furthermore, in their study no differences in effectiveness and satisfaction were found between the two quiz types. However, the authors concluded that PDAs are an attractive test administration option for schools and universities.

Research and development on CAT is mostly targeted towards PC based users. Although mobile computing has become an important and interesting research issue, little research has been done on the implementation of CAT using mobile devices and this is the focus of the research presented in this paper. The current study is an attempt to examine the design and development issues, which may be important in the implementation of a CAT using mobile devices such as mobile phones and PDAs. As a case study an educational assessment prototype was developed, called computerized adaptive test on mobile devices (CAT-MD), to support the assessment procedure of the subject “Physics” which is typically offered to second grade students in senior high school in Greece.

4. System architecture

The prototype CAT-MD uses the item response theory (IRT) as an underlying psychometric theory, which is the basis for many adaptive assessment systems and depicts the relationship between examinees and items through mathematical models (Hambleton, Swamination, & Rogers, 1991; Lord, 1980; Wainer, 1990). Psychometric theory is the psychological theory or technique of mental measurement, which forms the basis for understanding general testing theory and methods. The central element of IRT is mathematical functions that calculate the probability of a specific examinee answering a particular item correctly. IRT is used to estimate the student’s knowledge level, in order to determine the next item to be posed, and to decide when to

finish the test. There are four main components needed for developing IRT-based CAT: the item pool, the item selection procedure, the ability estimation and the stopping rule (Dodd, De Ayala, & Koch, 1995). The following sections describe these components of the CAT-MD system.

4.1. Item pool

The most important element of a CAT is the item pool, which is a collection of test items that includes a full range of levels of proficiency, and from which varying sets of items are presented to the examinees. The success of any CAT program is largely dependent on the quality of the item pool that can be conceptualized according to two basic criteria: (a) the total number of items in the pool must be sufficient to supply informative items throughout a testing session, and (b) the items in the pool must have characteristics that provide adequate information at the proficiency levels that are of greatest interest to the test developer. This criterion mainly suggests that at all important levels of proficiency there are sufficient numbers of items whose difficulty parameters provide valuable information. Therefore, a high-quality item pool will include sufficient numbers of useful items that allow efficient, informative testing at important levels of proficiency (Wise, 1997).

The item parameters included in the pool are dependent upon the Item Response Theory (IRT) model selected to model the data and to measure the examinees' ability levels. In IRT-based CATs, the difficulty of an item describes where the item functions along the ability scale. For example, an easy item functions among the low-ability examinees and a hard item functions among the high-ability examinees; thus, difficulty is a location index.

An ideal item pool needs many items, best spread evenly over the possible range of difficulty. In our approach CAT-MD includes a database that contains 80 items related to the chapter "electricity" from the "physics" subject. For every item, the item pool includes the item's text, details on the correct answer and the difficulty level. The difficulty level varies from "very easy" to "very hard" and the values used cover the range between -2 and $+2$.

4.2. Item selection

In IRT theory, the item selection procedure is the process of selecting an item from the item pool to be administered to the examinee. A reasonable assumption is that each examinee responding to a test item possesses some amount of the underlying ability. Thus, one can consider each examinee to have a numerical value, a score that places him or her somewhere on the ability scale. This ability score will be denoted by the Greek letter theta, θ . At each ability level, there will be a certain probability that an examinee with that ability will give a correct answer to the item. This probability will be denoted by $P(\theta)$. In the case of a typical test item, this probability will be low for examinees of low-ability and high for examinees of high ability (Baker, 2001).

Two common classes of IRT models are determined by the way items' responses are scored. Items with only two response options (correct or incorrect) are modelled with the dichotomous IRT models. Items with more than two response options can be modelled with polytomous IRT models (Boyd, 2003). Our prototype CAT-MD includes multiple choice items and true false items. Since, these are examples of items that can be scored dichotomously, CAT-MD is based on a dichotomous IRT model.

The main aspect of IRT is the Item Characteristic Curve (ICC) (Baker, 2001). ICC is an exponential function, which expresses the probability of a learner with certain skill level correctly answering a question of a certain difficulty level. ICC is a cumulative distribution function with a discrete probability. The models most commonly used as ICC functions are the family of logistics models of one (1PL), two (2PL) and three parameters (3PL).

The 1-parameter logistic (1PL), or Rasch model is the simplest IRT model. The Danish mathematician Georg Rasch first published the 1-parameter logistic model in 1960s and as its name implies, it assumes that only a single item parameter is required to represent the item response process. This item parameter is termed difficulty and the equation for this model is given by:

$$P(\theta) = \frac{1}{1 + e^{-1(\theta-b)}}, \quad (1)$$

where, e is the constant 2.718, b is the difficulty parameter and θ is an ability level.

In CAT-MD, as each student answers a question, his or her response is evaluated as being either correct or incorrect. In the event of a correct response, the probability $P(\theta)$ is estimated applying the formula shown in Eq. (1). Otherwise, the function $Q(\theta) = 1 - P(\theta)$ is used.

The item information function (IIF) is also considered as an important value in the item response theory's item selection process. It gives information about the item to be presented to the learner in an adaptive assessment. Before selecting a question appropriate to the learner, the IIF for all the questions in the assessment should be calculated and the question with highest value of IIF is presented to the learner. This provides more information about the learner's ability and is given by the equation:

$$I_i(\theta) = P_i(\theta)(1 - P_i(\theta)), \quad (2)$$

where $P_i(\theta)$ is the probability of a correct response to item i conditioned on ability θ (Baker, 2001; Lord, 1980).

4.3. Ability estimation

After each item is administered and scored, an interim estimate of the examinee's ability (θ) is calculated and used by the item selection procedure to select the next item. The most commonly used estimation procedure is maximum likelihood estimation (MLE) (Lord, 1980). Similar to the item parameter estimation, this procedure is an iterative process. It begins with some a priori value for the ability of the examinee. In CAT-MD, it begins with $\theta = 1$. The estimation calculation approach is the modification of the Newton–Raphson iterative method for solving equations method outlined by Lord. The estimation equation used is shown below:

$$\theta_{n+1} = \theta_n + \frac{\sum_{i=1}^n S_i(\theta_n)}{\sum_{i=1}^n I_i(\theta_n)}, \quad (3)$$

where

$$S_i(\theta) = [u_i - P_i(\theta)] \frac{P_i'(\theta)}{P_i(\theta)[1 - P_i(\theta)]}, \quad (4)$$

where

θ is the skill level after n questions, and
 $u_i = 1$ if the response is correct and $u_i = 0$ for the incorrect response.

4.4. Stopping rule

One important characteristic of CAT is the test termination criterion. The termination criterion is generally based on the accuracy with which the examinees' ability has been assessed. In most CATs, the termination of the test may be based on the number of items administered, the precision of measurement or a combination of both (Boyd, 2003). Measurement precision is usually assessed based on error associated with a given ability. The standard error associated with a given ability is calculated by summing the values of the item information functions (IIF) at the candidate's ability level to obtain the test information. Test information, $TI(\theta)$, is given by the equation:

$$TI(\theta) = \sum_{i=1}^N I_i(\theta). \quad (5)$$

Next, the standard error is calculated by using the equation:

$$SE(\theta) = \frac{1}{\sqrt{TI(\theta)}}. \quad (6)$$

After each administration of an item, the standard error associated with a given ability is calculated to determine whether a new item must be selected or whether the administration of the test can be terminated. It is common in practice, to design CATs so that the standard errors are about 0.33 or smaller (Rudner, 1998). In CAT-MD the test terminates for each examinee when the standard error associated with a given ability (θ) is less than 0.30 or when the maximum number (that is 20) of items has been administered.

5. System implementation

To date, the basic architecture of the system has been implemented. The prototype software has been developed using macromedia flash as it offers competitive advantages. It is a lightweight, cross-platform runtime that can be used not just for enterprise applications, but also for communications, and mobile applications. According to macromedia company, 98% of all Internet enabled computers and 30 million mobile devices use the Flash technology (<www.macromedia.com>). Currently, many manufacturers license macromedia flash on their branded consumer electronics devices, such as mobile phones, portable media players, PDAs, and other devices. These licensees include leading mobile device manufacturers such as nokia, samsung, motorola, and sony ericsson.

Fig. 1 presents two screenshots of the implementation of CAT-MD on a mobile phone and on a PDA. Moreover, the CAT-MD is portable to any device that has the macromedia standalone-flash player installed. In addition, if a macromedia plug-in for the web browser (Internet Explorer, Mozilla, etc.) is installed, the CAT-MD can also be accessed as a flash shockwave film.



Fig. 1. Interface of CAT-MD (a) The CAT-MD on HP iPAQ (PDA). (b) The CAT-MD on motorola MPx220 (mobile).

6. Formative evaluation

Throughout the development of CAT-MD, formative evaluation was an integral part of the design methodology. Formative evaluation is the judgment of the strengths and weaknesses of the system in its developing stages, for the purpose of revising the system. The goals of the evaluation were to test the satisfaction of the learners with the use of CAT-MD, and to examine the effectiveness and efficiency of the assessment procedure. In the CAT-MD case, three types of formative evaluation were used: expert review, one-to-one evaluation and small group evaluation. Tessmer (1993) defines these evaluations types as follows:

- Expert review – experts review the system with or without the evaluator present. The experts can be content experts, designers or instructors.
- One-to-one evaluation – one learner at a time reviews the system with the evaluator and comments upon it.
- Small group evaluation – the evaluator tries out the system with a group of learners and records their performance and comments.

As discussed above, CAT-MD is portable to any device that has the macromedia standalone-flash player installed. However, we evaluated the PDA version (HP iPAQ 1950) as it is a more useful device in an educational setting. The current interface of CAT-MD and its functionalities are the result of revisions based on the analysis of the data collected during the formative evaluation. Next, we will discuss the process of the formative evaluation i.e. subjects, procedure, and results.

6.1. Expert review and one-to-one evaluation

Four experts acted as evaluators in the expert review phase: a teaching/training expert, an instructional design expert, a subject-matter expert, and an educational technologist. In this phase, a semi-structured interview aimed at determining the reactions of experts and a debriefing session were used. During the semi-structured interview, although the subjects were prompted with questions, the main aim was to get their subjective reactions on the clarity, completeness and ease of use of the prototype CAT-MD. The debriefing session was used to conclude the evaluation with general questions about the assessment procedure and the design of the prototype CAT-MD and to prompt subjects' suggestions for improvement. In addition, the experts were asked to evaluate the prototype's effectiveness as a computer-assisted assessment tool in an academic context.

All the experts declared satisfaction with the use of the PDA as an additional tool of the assessment procedure and they agreed with its user friendliness. They indicated that PDAs are attractive devices because they are lightweight, extremely portable, and feature touch-sensitive displays. Furthermore, they felt that the introduction of mobiles devices into the learning process adds extra value to the use of computers in educational settings. They pointed out that the mobility eliminates the need for a specialized computer lab. PDAs can be used anywhere, including a traditional classroom.

Furthermore, they stated that computerized assessment procedure is a promising alternative to the traditional paper-and-pencil assessment procedure. Experts stated that they were satisfied with the use of CAT-MD and they indicated that such assessment can also be implemented within other subjects and especially those that use labs, such as physics, chemistry etc. Furthermore, they agreed that the CAT-MD interface was clear and straightforward. In addition, they observed that the adaptation based on examinees' knowledge estimate was easy to understand.

The expert review phase was followed by the one-to-one evaluation. Ten subjects participated in the one-to-one evaluation. They were second grade senior high school students studying the chapter "electricity" from the "physics" subject. In this phase, a semi-structured interview and a debriefing session were used as well. The interview covered the following topics: (1) the use of the PDA, (2) the review of the computerized adaptive testing procedure, and (3) the evaluation of the CAT-MD system.

All the students agreed on the user friendliness of the PDA. They did not have any previous experience with a similar tool and they felt content using it. Furthermore, they stated that they wanted to use the PDA again. With regard to the adaptive testing algorithm, they agreed that the structure of the CAT was clear and easy to understand. They also felt that the CAT was accurate, exact and reliable. Concerning the CAT-MD system,

their comments were also satisfactory. They agreed on its effectiveness and user friendliness. Furthermore, most of the students felt that the CAT-MD was enjoyable mainly due to the use of the PDA.

Although most comments were positive, the students pointed out some weaknesses in the software. These comments from one-to-one evaluation were compared and processed together with experts' suggestions. Two of the more significant recommendations that were implemented during the revision phase are as follows:

- In the initial implementation of CAT-MD, the examinee selected his/her response and then the test did not forward automatically to the next item. This was in order to give the student the possibility to rethink and alter his/her response. To move to the next item, the examinee had to press the “next” button. Although the students received a short introduction on how to use the CAT-MD system, most of them expected that the test would administer the next item automatically after their response. Most of the students and experts suggested that this would be improved by altering the state of the “next” button in order to indicate that the selection has been finished and that the examinee can proceed to the next phase (Fig. 2).
- The process of displaying items, evaluating responses and selecting the next item to be administered is repeated until the accuracy of the estimate reaches a statistically acceptable level or twenty (20) items have been administered. In the prototype CAT-MD, the system did not display any message about the correctness of the student's response. Most of the students indicated that they prefer to know if their answer was correct or not immediately after the processing of each item. Furthermore, they pointed out as a weakness that they did not know the total number of the items that they had already answered. After the revision, these suggestions were implemented (Fig. 3).

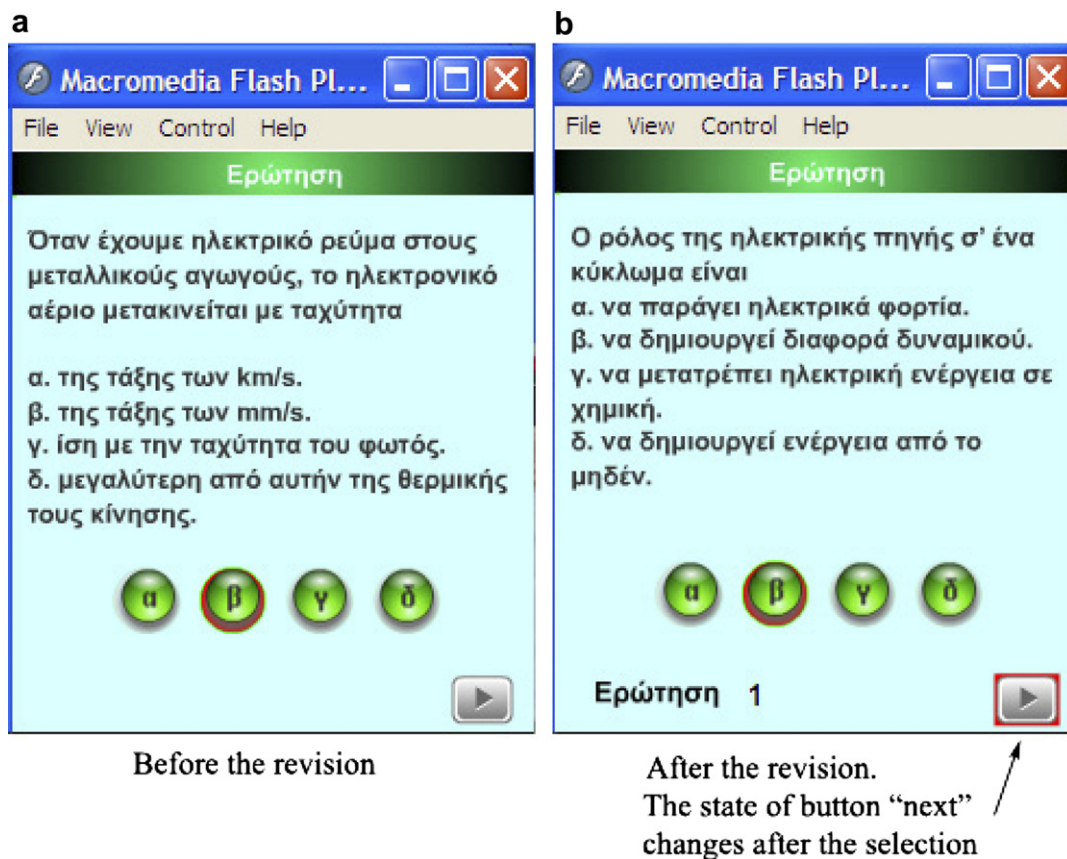


Fig. 2. System screen before an after the revision.

6.2. Small group evaluation

The final stage of formative evaluation consisted of the small group evaluation in a “real world” environment. Twelve subjects took part in the small group evaluation. Students at the second grade of senior high school were asked to volunteer for the small group evaluation. The subjects were selected carefully from the volunteers in order to represent the target population. There were seven males and five females. While some of the students had little experience with computers, none had received any test on a computer and none had ever used a PDA. The evaluation was conducted by collecting data about the assessment procedure from a variety of sources, using a variety of data gathering methods and tools.

At the beginning of the evaluation, all students were administered a paper-and-pencil test in order to compare testing time between CAT-MD and paper-and-pencil test. Twenty (20) multiple choice items and true false items were included in the paper-and-pencil test. These items were not included in the item pool of the CAT-MD to avoid the items’ exposure. However, their difficulty level was similar to those included in the item pool in order to have comparable results. After this procedure the students received a short introduction on how to use the CAT-MD system. During the assessment procedure, the students used comment logs in order to note specific strengths or weaknesses of the system. When all the students had finished using the CAT-MD, an attitude questionnaire was used to determine their experience of using the system. Finally, a debriefing session was used to assess the subjective satisfaction of subjects on the instructional and interface design of CAT-MD.

One of the major purposes of the small group evaluation was to investigate the efficiency and the effectiveness of CAT-MD. Effectiveness and efficiency are interrelated evaluation goals. The efficiency of the assessment is related to the time required for learners to conclude the test. The majority of the students

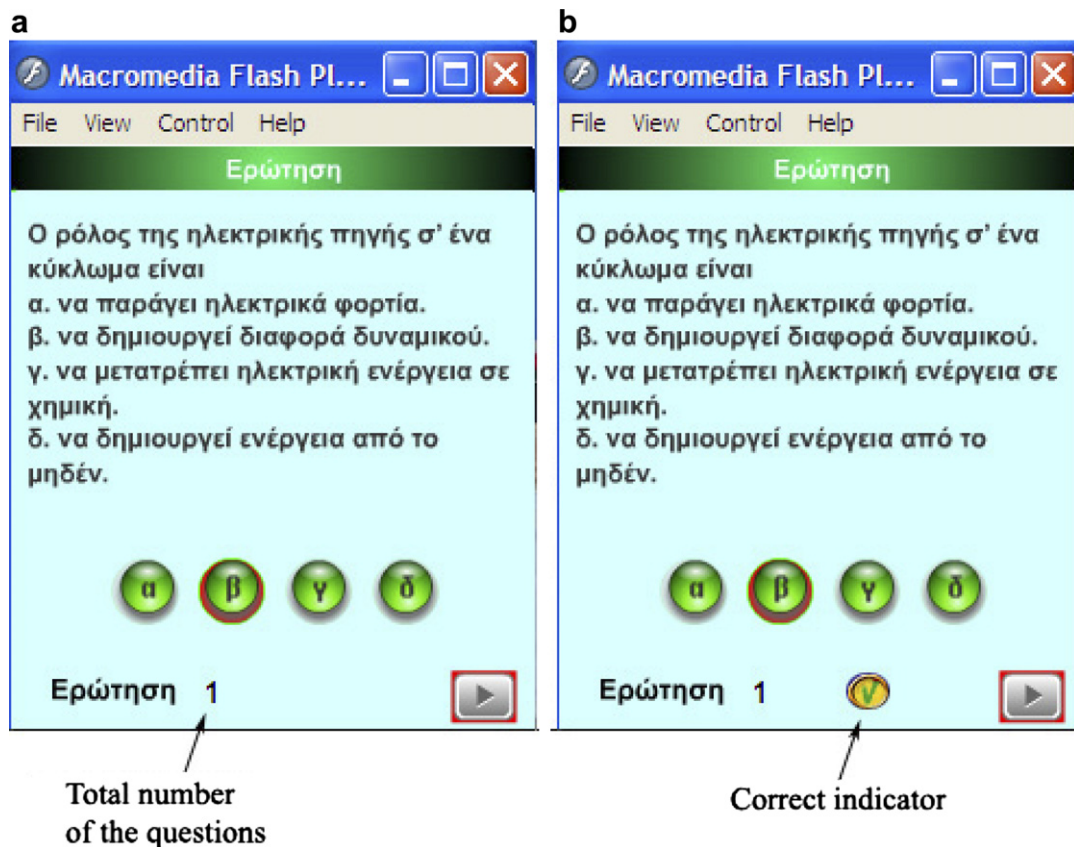


Fig. 3. System screen after the revision.

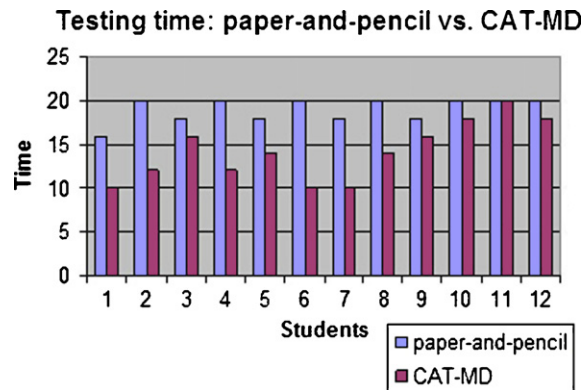


Fig. 4. Paper-and-pencil vs. CAT-MD testing time results.

completed the assessment procedure with CAT-MD in significantly less time than the paper-and-pencil test since, in general, fewer items were needed with CAT-MD to estimate their ability (Fig. 4). According to students' comments they did not waste time attempting items that were too hard or trivially easy. So, regarding the efficiency, the instruction was considered successful.

Furthermore, small group evaluation can offer useful implementation information concerning the usability and the appeal of the system. Usability is a measure of the ease with which a system can be learned and used. The usability and the appeal of the system were investigated in this study through the attitude questionnaires and the debriefing session. An analysis of the data collected showed that overall the students were satisfied with CAT-MD. Ten out of twelve students indicated that they were pleased with the design of the interface. However, two of the students noticed that the images included in some items were not so clear due to small resolution and size.

Moreover, the majority of the students felt that the test was clear and secure. They indicated that it was easy to understand the adaptive assessment procedure. They were satisfied with the test adaptation based on their knowledge estimate. Furthermore, the majority of students felt that CAT-MD offers more accurate results as they were graded not just on the basis of correct answers, but also on the difficulty level of the items answered. In addition, they expressed satisfaction with the fact that the score was available immediately, providing instantaneous feedback. To summarise, positive comments referred to the ease of taking the mobile test, the speed of the mobile test, and the relative absence of stress during the exam (Table 1).

Most of the students were satisfied with the system's mobility. They indicated that the use of the PDA was very interesting and attractive. The majority of the students were enthusiastic with the implementation of the CAT on the PDA and they stated that they wanted to use CAT-MD again (Table 1). Eleven out of twelve students answered affirmatively to the question of whether they would like to use this tool for other subjects. In addition, they pointed out that CAT-MD can be used for any subject and especially for those subjects that use labs, such as physics, chemistry etc.

Table 1
Students' answers on the attitude questionnaire

Attitude questionnaire	Not at all (%)	A little (%)	Average (%)	Much (%)	Very much (%)
Test is clear and secure	0	0	0	25	75
Easy to understand the adaptive assessment procedure	0	0	8	17	75
Satisfaction with the test adaptation based on student knowledge estimate	0	0	17	17	66
Accurate results	0	0	0	25	75
Instantaneous feedback	0	8	8	17	67
Satisfaction with system's mobility	0	8	17	17	58
The use of the PDA is very interesting and attractive	0	0	0	17	83

7. Conclusion

This article describes the design and development of the CAT-MD (Computerized Adaptive Test on Mobile Devices), a prototype CAT on mobile devices such as PDAs. The design process was driven by continuous formative evaluation. The positive comments of the users and the experts are summarized below:

- CAT-MD is an effective and efficient assessment tool. The assessment procedure with CAT-MD was completed in significantly less time than the paper-and-pencil test since, in general, fewer items were needed with CAT-MD to estimate their ability. According to students' comments they did not waste time attempting items that were too hard or trivially easy.
- CAT-MD was accurate, exact and reliable. The students felt that CAT-MD offers more accurate results as they were graded not just on the basis of correct answers, but also on the difficulty level of the items answered.
- CAT-MD can also be implemented within other subjects and especially those that use labs, such as physics, chemistry. Furthermore, the mobility eliminates the need for a specialized computer lab and it can be used anywhere.

Moreover, the recommendations from experts and the suggestions that resulted from one-to-one and small group evaluation are as follows:

- The system needs to clearly indicate when the next item is ready to be administered.
- The system should provide immediate feedback on the correctness of the examinee's answer.
- Students should be able to access information at any time with regards to the total numbers of items administered.
- Images on PDAs should be carefully selected taking into consideration the small resolution and size of the device.

These recommendations result from the formative evaluation of CAT-MD. Although they do not have universal value because the design and development of a system always depend on the target population and the subject-matter, they could be seen as points worth considering for designers of computerized adaptive testing on mobile devices. The next step of our research is summative evaluation in order to assess the effectiveness of the system with reference to validity.

Acknowledgement

The work presented in this paper is partially funded by the General Secretariat for Research and Technology, Hellenic Republic, through the e-learning, EL-51, FlexLearn project.

References

- Abdullah, S. C. (2003). *Student modelling by adaptive testing - a knowledge-based approach*. Unpublished PhD Thesis, University of Kent at Canterbury, June 2003.
- Baker, Frank (2001). *The basics of item response theory*. ERIC clearinghouse on assessment and evaluation. College Park, MD: University of Maryland.
- Boyd, A. M. (2003). *Strategies for controlling testlet exposure rates in computerized adaptive testing systems*. Unpublished PhD Thesis, The University of Texas at Austin, May 2003.
- Chen, F., Myers, B., & Yaron, D. (2000). *Using handheld devices for tests in classes* (Tech. Report CMU-CS-00-152). Pittsburgh: Carnegie Mellon University, School of Computer Science.
- Cook, R. P. (2000). The national classroom project an experience report. *Proceedings of the 30th ASEE/ISEE frontiers in education conference*. Champaign, IL: Stripes Publishing, pp. T1E/1–6.
- Dodd, B. G., De Ayala, R. J., & Koch, W. R. (1995). Computerized adaptive testing with polytomous items. *Applied Psychological Measurement*, 19(1), 5–22.
- Eggen, T. J. H. M. (2001). *Overexposure and underexposure of items in computerized adaptive testing*. Measurement and Research Department Reports 2001-1, Citogroep Arnhem.

- Georgiev, T., Georgieva, E., & Smrikarov, A. (2004). M-learning – a New Stage of e-learning, International conference on computer systems and technologies – CompSysTech'2004, Rousse, Bulgaria, 17–18 June, 2004. <<http://ecet.ecs.ru.acad.bg/cst04/Docs/sIV/428.pdf>>.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. California: Sage.
- Jones V., & Jo H. J. (2004). Ubiquitous learning environment: an adaptive teaching system using ubiquitous technology. In: *Proceedings of the 21st ASCILITE conference*, Perth, Western Australia, 5–8 December, 2004.
- Lilley, M., & Barker, T. (2002). The development and evaluation of a computer-adaptive testing application for english language. *6th Computer assisted assessment conference*, July 2002, Loughborough.
- Lilley, M., & Barker, T. (2003). An evaluation of a computer adaptive test in a UK university context. *7th Computer assisted assessment conference*, 8th and 9th July, 2003, Loughborough.
- Linacre, J.M. (2000). Computer-adaptive testing: a methodology whose time has come. MESA memorandum no. 69. Published in Sunhee Chae, Unson Kang, Eunhwa Jeon and J. M. Linacre. Development of computerised middle school achievement test, Seoul, South Korea: Komesa Press (in Korean).
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. New Jersey: Lawrence Erlbaum.
- Luvai, F. M. (in press). Mobile learning: a framework and evaluation. *Computers and education*.
- Rudner, L. M. (1998). An online, interactive, computer adaptive testing tutorial. 11/98. <<http://EdRes.org/scripts/cat>>.
- Segall, N., Doolen, T., & Porter, D. (2005). A usability comparison of PDA-based quizzes and paper-and-pencil quizzes. *Computer and Education*, 45(2005), 417–432.
- Sharples, M. (2000). The design of personal mobile technologies for lifelong learning. *Computers and Education*, 34, 177–193.
- Tessmer, M. (1993). *Panning and conducting formative evaluations*. Kogan Page.
- van der Linden, W. J., & Glas, C. A. W. (2003). Preface. In W. J. van der Linden & C. A. W. Glas (Eds.), *Computerised adaptive testing: theory and practice* (pp. vi–xii). Dordrecht, Boston, London: Kluwer Academic Publishers.
- Wainer, H. (1990). *Computerized adaptive testing (a primer)*. New Jersey: Lawrence Erlbaum.
- Wise, S. L. (1997, March). Overview of practical issues in a CAT program. *Paper presented at the annual meeting of the National Council on Measurement in Education*, Chicago IL.