

E-assessment: screen design and question difficulty

Changing the way we see test-items in a computer-based environment

Matt Haigh: mrh47@cam.ac.uk



Lessons with a penguin

James is a newly-qualified teacher who is enthusiastic about using technology in the classroom. He has persuaded his head of department to purchase a very impressive software package for assessing his students' science knowledge. The students are particularly excited by the animated penguin that helps them navigate around the features of the software and gives them feedback on their answers.

After several lessons he looks at the students' assessment results. Following a promising rise in their scores, he notices a sudden downturn in performance; even his brightest students are performing badly on the assessments. He checks the software and observes his students more closely in the following lesson. His observations draw a blank and, following a string of incorrect answers, he asks one of his students what is going on:

"Well sir, you see, when you get a question right, the penguin just gives you a little smile, but when you get a question wrong, he rolls off his iceberg and makes a big splash in the water which is really funny, so I'm kind of getting the questions wrong so I can see the penguin get wet!"



Background

Aim

The aim of this study was to investigate how aspects of screen design in a computer based assessment affect the difficulty of the test items.

Research questions

1. What are the effects of changing the student-item interaction on measures of item difficulty with a computer-based test item?
2. What insights are given by the test users for any potential sources of difficulty in computer-based test items?

Context and sample

The research was undertaken with 112 15-year-old students in seven English secondary schools. The computer-based assessment was based on GCSE Science qualifications and was administered online.

Research design

Mixed methods

The study used a mixed methods design, with an initial quantitative phase to examine measures of item difficulty, followed by a qualitative phase to collect the perceptions of difficulty from the students taking the assessment.

Stage 1 – Quantitative

- Two parallel forms of the online test were created; in the second of these forms a specific aspect of screen design was modified.
- Students taking the test were randomly assigned one of the two parallel forms.
- The test consisted of five items common to both tests for the purposes of checking that the samples were matched with regards to their ability.
- Ten further items were modified for each parallel form.
- Student responses were analysed to determine the measures of item difficulty.
- Background data were collected on prior attainment and ICT competence which have been shown to be key factors in computer-based test performance.

Stage 2 - Qualitative

- Two focus groups (each with four students) were held following the online-test.
- Students were presented with the parallel forms of each item.
- Responses were collected on features that students perceived to make an item easier or more difficult.
- Responses were analysed for common themes to identify insights into how the screen design altered their perception of an item's difficulty.

Results

Analysis and implications

Measures of difficulty

No significant differences were reported in the measures of item difficulty in any of the ten items modified in the test. This implies that students are able to demonstrate their subject knowledge through a variety of screen formats and interactions.

Student perceptions

In many cases students clearly articulate perceived differences in the difficulty of items. On some they are able to reach consensus on what makes an item more difficult, and on others they list factors relating to both versions of the item that have an influence on difficulty.

Implications for practice

The quantitative data indicates that students are able to handle a range of interactions for providing evidence of their abilities which is reassuring for the development of tests in different formats. However the qualitative data indicate that some forms are perceived differently by students in terms of their difficulty and that practitioners should take note of student input when designing computer-based assessments.

Further research

It may be that particular sub-groups of students (e.g. those with poor ICT skills or of lower ability) may be more susceptible to changes in the presentation of items on screen and this is worth investigating further. The study was undertaken with a limited number of screen-modifications and in a specific context. Further research could be undertaken with a broader range of item types and in a wider set of contexts to identify whether the findings were more generalisable.

Parallel test forms – Spot the difference...

Key to results:

	Description of the modification to screen design in the parallel forms of the test		The difficulty measure of the question (item facility): 0 represents a very hard question 1 represents a very easy question
	Indicates whether the difference in question difficulty is significant (at 5% level)		Student comments selected from the qualitative phase of the study

A subset of five out of the ten modified items in the study are shown below

Question 6
Which one of the following is a valid argument for using nuclear power stations?
 for maximum efficiency, they have to be sited on the coast
 they have high decommissioning costs
 they use a renewable energy source
 they do not produce gases that pollute the atmosphere

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Presence / absence of colour image

0.45

0.44

"The one with the picture is easier because you can see what it's about."

Question 6
Which one of the following is a valid argument for using nuclear power stations?
 for maximum efficiency, they have to be sited on the coast
 they have high decommissioning costs
 they use a renewable energy source
 they do not produce gases that pollute the atmosphere

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Question 9
When James exercises his breathing rate gets faster.
Drag the correct words below to complete the sentence
His breathing rate gets faster so that his muscles can receive more quickly, the muscles also need to remove more
carbon dioxide nitrogen protein
vitamins oxygen

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Dragging response vs drop-down menu response

0.82

0.87

"The dragging is easier because you can see all the words like actually on the page... instead of just clicking and looking down at the choices"

Question 9
When James exercises his breathing rate gets faster.
Drag the correct words below to complete the sentence
His breathing rate gets faster so that his muscles can receive [Select...] more quickly, the muscles also need to remove more [Select...]
Select...
nitrogen
protein
vitamins

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Question 10
Join the boxes to show the metals present in each alloy.
Click on each dot to start each line

alumini
brass
copper
nickel
steel
titanium
water

Lead & tin
Mercury
Copper & zinc
Nickel & cobalt

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Connecting points response vs dragging response

0.60

0.61

"The drawing is a bit more fun than just dragging it"

"The drawing lines one looks messy with all the lines going different places which is more difficult"

Question 10
Drag the boxes into the table to show the metals present in each alloy.

Alloy	Metals Present
alumini	
brass	
solder	
steel	

Lead & tin
Mercury
Copper & zinc
Nickel & cobalt

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Question 13
Cracking is a process that takes place at an oil refinery.
Which two sentences below about cracking are correct?
This is the TWO correct sentences

Cracking converts small molecules into large molecules
 Cracking needs a catalyst and a high temperature
 Cracking separates crude oil fractions
 Cracking is used at an oil refinery to make more petrol
 Cracking works because different fractions have different boiling points

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Multiple choice tick-boxes vs multiple choice whole-sentence selection

0.23

0.28

"[Highlighting] could make it easier when you're going over your answers at the end"

Question 13
Cracking is a process that takes place at an oil refinery.
Which two sentences below about cracking are correct?
Select the TWO correct sentences

Selecting needs a catalyst and a high temperature
Cracking separates crude oil fractions
Cracking is used at an oil refinery to make more petrol
Cracking works because different fractions have different boiling points

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Question 14
Hennetta is testing three fuels using the apparatus shown. The table shows her results.
Which fuel gives the most energy per gram?

Explain your answer:

apparatus results

Fuel	Mass of fuel	Volume of water	Initial temperature	Final temperature	Temperature change	Heat lost	Heat gained
gasoline	0.8	0.9	21.2	22.1	0.9	10.8	10.8
diesel	20	22	19	20	1	21	21
PDO	40	42	28	29	1	42	42

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Full stimulus on screen vs stimulus accessed through tabbed panels

0.50

0.61

"Most people would know that [you need to click on the tabs], but not if you are in a rush in your exam... on the last few questions you might not know the tab was there."

Question 14
Hennetta is testing three fuels using the apparatus shown. The table shows her results.
Which fuel gives the most energy per gram?

apparatus results

Explain your answer:

Fuel	Mass of fuel	Volume of water	Initial temperature	Final temperature	Temperature change	Heat lost	Heat gained
gasoline	0.8	0.9	21.2	22.1	0.9	10.8	10.8
diesel	20	22	19	20	1	21	21
PDO	40	42	28	29	1	42	42

Q1 Q2 Q3 Q4 Q5 Q6 Q7 Q8 Q9 Q10 Q11 Q12 Q13 Q14 Q15

< Previous Q Next Q > Finish Test

Qualitative analysis – Word clouds

Student responses to the focus groups were divided into two categories: reasons given for features making an item easier and reasons associated with features that make an item more difficult. The word clouds below are based on a word-frequency analysis where font size directly corresponds to the word-frequency.

