

# Computer delivered listening tests: a sad necessity or an opportunity?

*John Field, CRELLA Institute, University of Bedfordshire*



# Technological improvements in

- Digitised sound quality
- Transmission of sound via the internet (cf phone transmission)
- Standardisation of headphone quality
- Availability and quality of listening centre facilities

# Benefits

- Flexibility
- Numbers
- Minimal supervision
- Costs
- Reduction in listening anxiety?
- Greater test taker control over content?
- Individualisation of the tasks to suit test taker's level?
- Greater construct validity??

# Some little discussed questions

- Can CBT pave the way for greater use of visual material?
- What are the constraints and opportunities for multi-level tests?
- Are some formats more appropriate than others?
- How best can a test deliver items on-screen?
- What type and length of recording is best suited to computer delivery?
- Is single or double play most appropriate?
- Are there problems with computerised scoring?
- Can computer delivered materials lend themselves to diagnostic testing and remedial practice?



# Benefits of visual support

In real-world conditions...

- A listener is assisted by paralinguistic cues from a speaker
  - Facial expressions
  - Gestures
  - Mouth movements (known to be integrated into the listening process: see e.g. McGurk, )
- There is a physical context for a dialogue
- There are probably powerpoint slides in a presentation

# Visuals in CBT listening



# A helpful / motivational graphic?

**PERIODIC TABLE OF THE ELEMENTS**

1 <b>H</b> 1.0079																	2 <b>He</b> 4.0026
3 <b>Li</b> 6.941	4 <b>Be</b> 9.012											5 <b>B</b> 10.811	6 <b>C</b> 12.011	7 <b>N</b> 14.007	8 <b>O</b> 16.00	9 <b>F</b> 18.998	10 <b>Ne</b> 20.179
11 <b>Na</b> 22.99	12 <b>Mg</b> 24.30											13 <b>Al</b> 26.98	14 <b>Si</b> 28.09	15 <b>P</b> 30.974	16 <b>S</b> 32.06	17 <b>Cl</b> 35.453	18 <b>Ar</b> 39.948
19 <b>K</b> 39.09	20 <b>Ca</b> 40.08	21 <b>Sc</b> 44.96	22 <b>Ti</b> 47.88	23 <b>V</b> 50.94	24 <b>Cr</b> 52.00	25 <b>Mn</b> 54.938	26 <b>Fe</b> 55.845	27 <b>Co</b> 58.933	28 <b>Ni</b> 58.69	29 <b>Cu</b> 63.546	30 <b>Zn</b> 65.38	31 <b>Ga</b> 69.723	32 <b>Ge</b> 72.63	33 <b>As</b> 74.922	34 <b>Se</b> 78.96	35 <b>Br</b> 79.904	36 <b>Kr</b> 83.80
37 <b>Rb</b> 85.468	38 <b>Sr</b> 87.62	39 <b>Y</b> 88.906	40 <b>Zr</b> 91.224	41 <b>Nb</b> 92.906	42 <b>Mo</b> 95.94	43 <b>Tc</b> (98)	44 <b>Ru</b> 101.07	45 <b>Rh</b> 102.905	46 <b>Pd</b> 106.42	47 <b>Ag</b> 107.868	48 <b>Cd</b> 112.411	49 <b>In</b> 114.818	50 <b>Sn</b> 118.710	51 <b>Sb</b> 121.757	52 <b>Te</b> 127.603	53 <b>I</b> 126.905	54 <b>Xe</b> 131.29
55 <b>Cs</b> 132.905	56 <b>Ba</b> 137.327	57 <b>*La</b> 138.905	72 <b>Hf</b> 178.49	73 <b>Ta</b> 180.948	74 <b>W</b> 183.84	75 <b>Re</b> 186.207	76 <b>Os</b> 190.23	77 <b>Ir</b> 192.222	78 <b>Pt</b> 195.084	79 <b>Au</b> 196.967	80 <b>Hg</b> 200.59	81 <b>Tl</b> 204.38	82 <b>Pb</b> 207.2	83 <b>Bi</b> 208.98	84 <b>Po</b> (209)	85 <b>At</b> (210)	86 <b>Rn</b> (222)
87 <b>Fr</b> (223)	88 <b>Ra</b> (226)	89 <b>*Ac</b> (227)	103 <b>Rf</b> (261)	104 <b>Db</b> (262)	105 <b>Sg</b> (263)	106 <b>Bh</b> (264)	107 <b>Hs</b> (265)	108 <b>Mt</b> (266)	109 <b>Uu</b> (269)	110 <b>Uu</b> (271)	111 <b>Uu</b> (272)	(Not yet named)					
*Lanthanide Series		58 <b>Ce</b> 140.12	59 <b>Pr</b> 140.91	60 <b>Nd</b> 144.24	61 <b>Pm</b> (145)	62 <b>Sm</b> 150.4	63 <b>Eu</b> 151.96	64 <b>Gd</b> 157.25	65 <b>Tb</b> 158.93	66 <b>Dy</b> 162.50	67 <b>Ho</b> 164.93	68 <b>Er</b> 167.26	69 <b>Tm</b> 168.93	70 <b>Yb</b> 173.04	71 <b>Lu</b> 174.97		
†Actinide Series		90 <b>Th</b> 232.04	91 <b>Pa</b> 231.04	92 <b>U</b> 238.03	93 <b>Np</b> 237.05	94 <b>Pu</b> (244)	95 <b>Am</b> (243)	96 <b>Cm</b> (247)	97 <b>Bk</b> (247)	98 <b>Cf</b> (251)	99 <b>Es</b> (252)	100 <b>Fm</b> (257)	101 <b>Md</b> (258)	102 <b>No</b> (259)	103 <b>Lr</b> (260)		

## Accompanying text

- *Um, actually, that reminds me of a good example of all this ... element 43. See on the table, the symbols for elements 42 and 44? Well, in early versions of the table, there was no symbol for an element with 43 protons because an element had yet been discovered with 43 protons. So the periodic table had a gap between elements 42 and 44. And, then, uh, in 1925 a team of chemists led by a scientist named Ida Tacke claimed that they had found element 43. They had been, uh, using a relatively new technology called x-ray spectroscopy—and they were using this to examine an ore sample—and they claimed they'd found an element with 43 protons. And they named it masurium.*

# A helpful / motivational graphic?

**PERIODIC TABLE OF THE ELEMENTS**

1 <b>H</b> 1.0079																	2 <b>He</b> 4.0026
3 <b>Li</b> 6.941	4 <b>Be</b> 9.012											5 <b>B</b> 10.811	6 <b>C</b> 12.011	7 <b>N</b> 14.007	8 <b>O</b> 16.00	9 <b>F</b> 18.998	10 <b>Ne</b> 20.179
11 <b>Na</b> 22.99	12 <b>Mg</b> 24.30											13 <b>Al</b> 26.98	14 <b>Si</b> 28.08	15 <b>P</b> 30.974	16 <b>S</b> 32.06	17 <b>Cl</b> 35.453	18 <b>Ar</b> 39.948
19 <b>K</b> 39.09	20 <b>Ca</b> 40.08	21 <b>Sc</b> 44.96	22 <b>Ti</b> 47.88	23 <b>V</b> 50.94	24 <b>Cr</b> 51.996	25 <b>Mn</b> 54.938	26 <b>Fe</b> 55.845	27 <b>Co</b> 58.933	28 <b>Ni</b> 58.69	29 <b>Cu</b> 63.546	30 <b>Zn</b> 65.38	31 <b>Ga</b> 69.72	32 <b>Ge</b> 72.63	33 <b>As</b> 74.922	34 <b>Se</b> 78.96	35 <b>Br</b> 79.904	36 <b>Kr</b> 83.80
37 <b>Rb</b> 85.47	38 <b>Sr</b> 87.62	39 <b>Y</b> 88.91	40 <b>Zr</b> 91.22	41 <b>Nb</b> 92.91	42 <b>Mo</b> 95.94	43 <b>Tc</b> 98.906	44 <b>Ru</b> 101.07	45 <b>Rh</b> 102.91	46 <b>Pd</b> 106.42	47 <b>Ag</b> 107.87	48 <b>Cd</b> 112.41	49 <b>In</b> 114.82	50 <b>Sn</b> 118.71	51 <b>Sb</b> 121.76	52 <b>Te</b> 127.60	53 <b>I</b> 126.91	54 <b>Xe</b> 131.29
55 <b>Cs</b> 132.91	56 <b>Ba</b> 137.33	57 <b>*La</b> 138.91	72 <b>Hf</b> 178.49	73 <b>Ta</b> 180.95	74 <b>W</b> 183.85	75 <b>Re</b> 186.21	76 <b>Os</b> 190.2	77 <b>Ir</b> 192.22	78 <b>Pt</b> 195.08	79 <b>Au</b> 196.97	80 <b>Hg</b> 200.59	81 <b>Tl</b> 204.38	82 <b>Pb</b> 207.2	83 <b>Bi</b> 208.98	84 <b>Po</b> (209)	85 <b>At</b> (210)	86 <b>Rn</b> (222)
87 <b>Fr</b> (223)	88 <b>Ra</b> (226)	89 <b>*Ac</b> (227)	104 <b>Rf</b> (261)	105 <b>Db</b> (262)	106 <b>Sg</b> (263)	107 <b>Bh</b> (264)	108 <b>Hs</b> (265)	109 <b>Mt</b> (266)	110 <b>Uu</b> (269)	111 <b>Uu</b> (271)	112 <b>Uu</b> (273)	(Not yet named)					
*Lanthanide Series		58 <b>Ce</b> 140.12	59 <b>Pr</b> 140.91	60 <b>Nd</b> 144.24	61 <b>Pm</b> (145)	62 <b>Sm</b> 150.4	63 <b>Eu</b> 151.97	64 <b>Gd</b> 157.25	65 <b>Tb</b> 158.93	66 <b>Dy</b> 162.50	67 <b>Ho</b> 164.93	68 <b>Er</b> 167.26	69 <b>Tm</b> 168.93	70 <b>Yb</b> 173.04	71 <b>Lu</b> 174.97		
†Actinide Series		90 <b>Th</b> 232.04	91 <b>Pa</b> 231.04	92 <b>U</b> 238.03	93 <b>Np</b> 237.05	94 <b>Pu</b> (244)	95 <b>Am</b> (243)	96 <b>Cm</b> (247)	97 <b>Bk</b> (247)	98 <b>Cf</b> (251)	99 <b>Es</b> (252)	100 <b>Fm</b> (257)	101 <b>Md</b> (258)	102 <b>No</b> (259)	103 <b>Lr</b> (260)		

# Still images

- Ginther (2002) cites evidence that test takers benefit from 'content' images which complement the information provided by the recording but not from others that just provide background context.
- Drawing on various findings, Ockey concludes (2007: 533) that **listeners are not usually helped or motivated by still visual stimuli.**

# But for video the picture is not clear...

- The research evidence that visual information enhances L2 listening comprehension under test conditions is not as clear as might be assumed.
- Some studies (Ginther, 2002, Wagner, 2010) report clear benefits for video material over audio.
- Others (Coniam, 2001) suggest that video makes no difference
- Others (Suvarov, 2009) suggest that performance declines .
- An interpretation of the last finding is that visual information actually serves as a distraction, and that audio material focuses a listener's attention better. Or it may be that individuals differ in the use they make of the evidence provided by video, with some benefiting and others not.

# Tests of academic listening

- It is anomalous that CB tests of academic listening do not aim to achieve greater cognitive validity by reproducing the real-world circumstances of listening to a lecture
- They have the opportunity of accompanying the words of the lecturer with on-screen Powerpoint slides – not simply visuals but also the text that would normally be available as a digest of the main points being made.
- There has been little research into the nature of the cognitive interaction that takes place in the mind of the academic listener between PP slides and the voice of the speaker. But see some early evidence in an eye-tracking study by Suvarov (2015).



# Computer delivery of multi-level tests

- Should material be in order of difficulty or presented randomly?
- In increasing order of difficulty
- How many exponents of each level do we need to ensure reliable grading of test takers?
- I'd suggest around 8
- How do we ensure that a test taker spends most time on items at an appropriate level?
- Limit time spent on items at lower levels to ensure that higher-level candidates don't get demotivated or stop treating the test seriously.

# Computational alternatives

- A. **A cut off point** where an individual's performance ceases to be consistent (within a margin of error)
- B. **A seeding solution**. According to the score on A<sub>1</sub> items, the individual is moved on either to A<sub>2</sub> or to B<sub>2</sub>. At a subsequent stage, the individual is redirected back to B<sub>1</sub> or on to C levels.
- C. **An adaptive algorithm** enabling the test taker to skip sections that are too easy and finally to focus on a section at the most appropriate level.



## Formats that work...

- Multiple-choice
- Visual multiple-choice
- Multiple matching
- True/false/not mentioned
- Identify three main points (not in same order as recording)
- **NOT gap fill**
- Tick box
- Drag and drop
- Drag and drop
- Tick box
- Tick box
- ? Drag and drop at lower levels?

# Scoring: test designers beware

Test providers are quite often reliant on standard platforms. The platform may not be capable of

- Scoring in terms of location on the screen
  - ?? Drag and drop??
- Scoring according to whether pieces of text have been placed in order
  - ?? Identify speaker's points and place in order??
- Scoring according to whether ticks appear in correct place in a table
  - ?? Multiple matching??

# Delivering items on-line

- So a tester has found a riveting 3 minute recording about the life of the housefly.
- She has found 6 major points to target with 4-option multiple-choice items
- Ergo: 30 lines of text. The test taker therefore has to scroll down while listening
- Problem: This adds to the divided attention demands of needing to read while listening. It also adds to item-order pressures: 'Have I missed Q<sub>1</sub> and should I now be listening out for Q<sub>2</sub>?' (Field, 2012)

# Solutions

- The CBT mode lends itself to extensive use of short clips of 15-40 seconds followed by a single question or by two questions.
- These can be used, even at higher proficiency levels, so long as they target higher-level processes (identify main point or speaker's attitude, infer information not made explicit, etc).
- Though some longer (3 min.) recordings are obviously necessary at higher levels, MCQs should be avoided. Response formats can be simplified by the use of multiple-matching, drag & drop etc.



# Options

- **Single play.** Fails to compensate for the lack of the visual evidence or the unnatural cognitive demands of most test formats.
- **Optional double play.** Test taker can opt to play the recording twice.
- **Optional replay.** Test taker can revert to any part of the recording that has not been understood
- **Double play.** Recording heard twice.

The last three can be claimed to reduce the extent of listening anxiety experienced by the test taker.

# Issues raised by optional replay

- The test designer has to determine maximum length for the test, allowing adequate time for replays.
- Should there be additional marks for a test taker who completes the test within a shorter time limit (i.e. with minimal replay)?
- Test takers, especially weaker ones, need to make decisions about a) how accurately the recording has been understood; b) whether sufficient time remains for a second hearing. This adds an important construct-irrelevant element to the cognitive demands of the test.
- Contrary to what has just been suggested, the stress associated with decision-making might particularly affect those suffering from listening anxiety.

# A double play solution that fits CBT well

- Research by Sherman (1997) demonstrated that performance in listening tests improves when test items are delivered between the two plays of a recording.
- It is much easier to adopt this procedure in CB tests than it is in paper based ones ,where answer sheets have to be handed out.
- Advantages: It ensures that, at least during the first play, test takers listen to a recording in a way that resembles real world listening : i.e. without
  - a) the prior knowledge gained from scanning *written* information in the form of pre-set test items
  - b) the test-wise strategies that the availability of this information encourages.See Field (2015) for a full discussion.

## Double play to reduce scrolling down

- Test takers hear a standard 3-4 min. recording all the way through.
- On a second play, the recording is divided into three parts. Before each part, two or three items appear on the screen. Test takers are allowed time to answer them, before the next part begins and a new set of questions appears...



# Progress tests

- In a listening centre, test takers listen to the recording once , and attempt to answer questions.
- On a second play, they are provided with a transcript. They highlight sections which they find difficult.
- This feeds into small-scale practice exercises that focus
  - a. on perceptual features that cause problems of word / phrase recognition
  - b. on syntactic and connective features that cause problems of parsing or understanding

# Transcript

- *I've had quite a lot of experience of overseas travel + so I do know that it's important to get to the airport with plenty of time + and I guess that's what I should have done this morning*

# Transcript

- *I've had quite a lot of experience of overseas travel + so I do know that it's important to get to the airport with plenty of time + and I guess that's what I should've done this morning*

# Practice materials

- *I should've done*
- *He might've done*
- *They must've done*
- *I needn't've done*
  
- *I should've thought*
- *I should've remembered*
- *I should've paid*
- *I should've apologised*

# References

- Coniam, D. (2001). The use of audio or video comprehension as an assessment instrument in the certification of English language teachers: a case study. *System*, vol. 29: 1-14.
- Field, J. (2012a) The cognitive validity of the lecture based question in the IELTS listening paper. In L. Taylor & C. Weir et al. (eds.) *IELTS Collected Papers 2: Research in Reading and Listening Assessment*. Cambridge: Cambridge University Press: pp. 391-453.
- Field, J. (2015) The effects of single and double play upon listening test outcomes and cognitive processing. London: British Council ARAGS Reports, 2015  
<http://www.britishcouncil.org/exam/aptis/research/publications>
- Ginther, A. (2002). Context and content visuals and performance on listening comprehension stimuli. *Language Testing*, vol. 19: 133-167.
- Ockey, G.J. (2007) Construct implications of including still image or video in computer-based listening tests. *Language Testing*, 24: 517-537
- Sherman, J. (1997) The effect of question preview in listening comprehension tests. *Language Testing*, 14, 185-213.
- Suvarov (2015). Interacting with visuals in L2 listening tests: an eye-tracking study. London: British Council ARAGS Reports, 2015. Available at:
- <https://www.britishcouncil.org/exam/aptis/research/publications/interacting>

## Solution 1: Control recording length

- A1-B1: clips of around 15 seconds
- B1+: recordings of 30 seconds
- B2 / C levels: recordings of around 2.5 – 3 minutes
- Advantages:
- Higher level test takers don't lose motivation through spending a lot of time on easy material
- Possible to recognise a cut-off points where test takers no longer perform consistently