



ICAR News

Issue Two
08 | 16

Editorial

by Philipp Doebler

Welcome to the second issue of the ICAR News. Since the last newsletter in the fall of 2015 the ICAR Project has grown substantially. ICAR now has more than 300 registered users out of which more than a 100 are actively using ICAR items in their own research.

We are more than happy to see the rapid growth of the number of ICAR members, and we believe that this growth is indicative for the large demand for public domain measures. The ICAR community is hence in the focus of this issue with three contributions from outside of the ICAR core team. The common denominator of these three contributions is that each article describes innovative software for cognitive ability assessment.

David Stillwell showcases the new version of the **Concerto platform**, an R-based solution for the web based presentation of test material. In contrast to many existing solutions, the concerto platform has been developed with adaptive testing and automatic on-the-fly item generation in mind.

ICAR Project homepage

The ICAR homepage is found at <https://icar-project.com>.

ICAR News Editor Team

The editors of ICAR News currently include several of the lead investigators of the ICAR Project: David Condon, Philipp Doebler, Heinz Holling, William Revelle, John Rust, and Luning Sun.

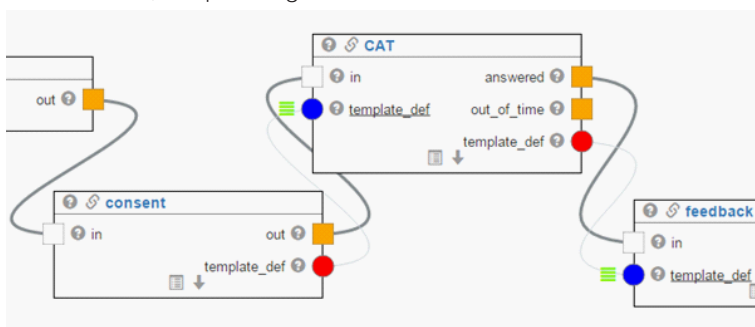
Diego Blum explains an item generator for **figural analogies**. The generator has kindly been made available in the form of an R package and is available on the ICAR website to registered users.

Finally, Peter Harrison, Tom Collins and Daniel Müllensiefen describe their state of the art approach to the assessment of **melodic discrimination ability**. The Goldsmiths Musical Sophistication Index (Gold-MSI) combines computerised adaptive testing and automatic item generation.

Concerto v5

by David Stillwell

Cambridge Psychometrics Centre are pleased to announce that **Concerto, the Open-source Online R-based Adaptive Testing Platform, has a new version (v5) in public beta**. Concerto allows users to create various online assessments, from simple surveys to complex IRT-based adaptive tests. Whereas Concerto's previous version was code-based and squarely aimed at programmers, v5 has the goal of making it possible to create **90% of tests without coding**. To do this, v5 has a new **flowchart interface** where test creators can link up nodes like "Demographics Page," "Questionnaire Page," and "CAT Page" in order to build a test in a modular format. Each time you add a node, a user interface takes you step by step through the process of adding information like a title, or uploading items to an item bank.



Our new website is www.concertoplatform.com and Concerto can be downloaded from <https://github.com/campsyh/concerto-platform/>. ICAR users who want to make use of the new Concerto platform hosted by the ICAR server (available at <http://concerto5.icar-project.com/admin>) can contact the ICAR admin (admin@icar-project.com) to apply for a trial account. We will be actively adding new tutorial videos over the next couple of months before launch, and are very interested to hear feedback or suggestions to d.stillwell@jbs.cam.ac.uk.

The Psychometrics Centre

The Psychometrics Centre is a centre of excellence within the University of Cambridge dedicated to research, teaching and product development in both pure and applied psychological assessment in the online environment. Active in Cambridge since 2005, it has seen significant growth in the past three years as a consequence of the explosion of activity in on-line communication and social networks. Visit the Psychometrics Centres homepage at www.psychometrics.cam.ac.uk.

about the author



Dr David Stillwell is a Lecturer in Big Data Analytics and Quantitative Social Science at the Judge Business School of the University of Cambridge. He is also Deputy Director of the The Psychometrics Centre at the University of Cambridge.

Automatic generation of figural analogy items

by Diego Blum

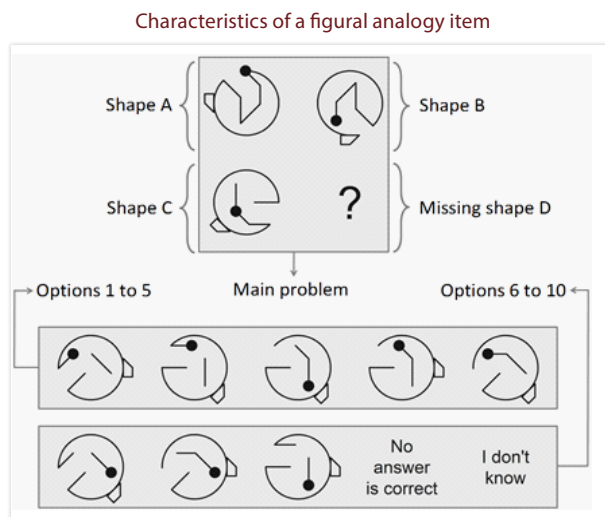
The research team headed by Prof. Dr. Heinz Holling at the University of Münster is currently working on several projects involving **Automatic Item Generation** (AIG). One of them is based on my research experience which I used for the creation of the **Item Maker** (IMak). IMak is an R-package available online from the ICAR homepage and thereby accessible to every researcher willing to perform studies on a self-generated scale. The functions 'build_fa' and 'plot_fa' are used to construct and display items respectively.

How to obtain IMak

IMak is available to registered ICAR users, and you can obtain it by clicking here: https://icar-project.com/projects/icar-project/wiki/Item_types

IMak in current research

Initially, items of this kind were manually constructed for a research project. See the **Intelligence** article here <http://dx.doi.org/10.1016/j.intell.2016.03.001>



Each item consists of a main problem and a set of options, only one of which is the right answer. The main problem comprises of three shapes arranged in a 2x2 array. The respondent is required to pick a shape from the options to complete the array. These shapes are related to each other by **analogy**, meaning that shape A is to shape B as shape C is to the missing shape D (A:B::C:D). Alternatively, the analogical relation could be A:C::B:D. During item construction, a cognitive operation leading to task solution (i.e., a rule) is matched with item parts and applied. Shape A, which is the top-left shape of the 2x2 array, is the one from where the reasoning always begins.

While the rule(s) affect(s) the analogical relation, the initial position of each part of the shape A affects how this relation is going to look like. This is a key difference between structural (i.e., conceptual) relations

and visual appearance. In AIG language, it is possible to say that the structure is given by the so-called **radicals**, which are the rules in this case, while the visual state of shape A is given by **incidentals**. Both radical and incidental arguments can be manipulated with function 'build_fa', but incidental arguments are allowed to be left random whereas at least one radical argument should always be affected. Moreover, the amount of radical arguments to be used equals the amount of rules applied, and a maximum of four different radical arguments can be manipulated at the same time inside the package.

Examples of items

Item with a subtraction rule plotted with default values.

Item with a reflection rule and some shape form changes plus German language.

Item with four rules and changes in plot mode as well as shape form plus Spanish language.

For users who are only interested in simple item generation, or if they are about to use IMak for the first time, it should be remembered that two pieces of information are the most important:

1. Items made with the 'build_fa' function of the IMak package require at least one rule to work with, and the following general rules can be manipulated: main shape rotation, main shape reflection, trapezium rotation, line segment subtraction and dot movement. Thus, **radical arguments** can be used which are named in a similar way as the aforesaid general rules. The beginner should play a little with radical arguments by specifying some rules of his choice and even combining them in a single function. It should be remembered that the following rule combinations are **not** allowed: main shape rotations combined with

about the author



Diego Blum

Lic. Diego Blum has worked in Buenos Aires as a Research Associate for several years and is now working on a PhD project funded by the German Academic Exchange Service (DAAD) at the University of Münster in cognitive ability assessment. He thanks for technical support by Philipp Doebler and Ehsan Masoudi. Anyone interested in item generation with IMak is free to consult Diego at blumworx@gmail.com

Website Link:
<http://www.uni-muenster.de/PsyIFP/AEHolling/de/blum/personen/>

reflection, and rules that have a general rule in common combined with each other. When it comes to rotations, it is strongly recommended to work with two numerical values that are not so distant from one another (see the examples inside the package).

2. It is useless to create an object with the 'build_fa' function if this object is not then plotted. The easiest way to plot is by simply using the 'plot_fa' function. It is not recommended to save the plot in a directory by means of the R interface buttons; instead, a directory can be given as another argument so as to let the function perform this procedure precisely and straightforward.

Assessing melodic discrimination abilities with computerised adaptive testing and automatic item generation

by Peter Harrison, Tom Collins and Daniel Müllensiefen

The ability to process and understand music is a universal human faculty. However, individuals can vary immensely in their ability to process musical materials, and many tests have been developed over the past century to assess these differences. Historically, these tests have typically been used to measure musical 'aptitude', with the aim of selecting the most able children for music education and instrument tuition. However, musical ability tests are increasingly being used in psychological and neuroscientific research, investigating how musical abilities affect music cognition, and how individual differences in musical abilities relate to individual differences in other cognitive abilities.

At the **Music, Mind & Brain group** at Goldsmiths, University of London, we have aimed to develop a musical test battery and corresponding self-report instrument called the **Goldsmiths Musical Sophistication Index** (Gold-MSI). This instrument is intended not to rely on formal musical training, not to be biased towards a particular musical style, and to be representative of a

large range of active musical behaviours in the general population.

It can be difficult to test for musical listening abilities effectively. If the test is to measure a broad range of abilities, as is common in the general population, it must contain items spanning a wide range of difficulty levels. This usually requires the test to be relatively long. Moreover, if listening tests work by aural presentation (and do not use musical notation), then they tend to provide only a small number of response options for each question. This contributes noise to test scores, because it is easy to guess the correct answer by chance. In order to compensate for this noise, test length needs to be increased further. However, musical listening tests are already tiring for test-takers, because test items usually take a long time to administer and demand a lot of concentration. This makes it difficult to increase test length further without increasing fatigue and hence diminishing reliability.

One possible way of addressing this problem is through computerised adaptive testing. Computerised adaptive tests (CATs) continuously tailor their difficulty to the estimated ability level of the test-taker as the test progresses (e.g. Figure 1). This can greatly increase test efficiency, as participants no longer have to take items that are far too easy or too difficult for their ability level. As a result, CATs typically can achieve the same reliability as traditional fixed-length tests even when test length is reduced by 50-70%.

Unfortunately, however, CATs can be very expensive to construct. CATs are usually built within Item Response Theory (IRT), a powerful psychometric framework for modelling ability tests. Each test item is modelled in terms of several psychometric parameters, such as difficulty, discrimination, and chance success rate. These parameters need to be estimated before the test can be used, but this can be a very expensive process, especially in the case of CATs, which typically require very large item banks. One recently developed technique for making CATs more efficient to construct is automatic item generation (AIG). In AIG, item parameters are not estimated for each item individually, but are instead predicted on the basis of structural characteristics of the items. This relies on a good understanding of the cognitive processes behind test-taking and a clear conceptualisation of how individual differences in abilities can contribute to task performance. AIG can make CAT construction more efficient, because much less response data is required to calibrate an AIG model than would be required to estimate every item parameter individually.

We have applied these techniques to the construction of a test of melodic discrimination abilities. In this test, test-takers have to discriminate between three similar versions of the same melody, each of which is transposed slightly higher in pitch. In every

Music Testing and Bias

Many traditional measures of musical abilities rely heavily on skills that are typically acquired through formal musical training. This is problematic because formal musical training tends to be biased towards Western art music and typically relies heavily on knowledge of Western musical notation. While musical notation is certainly important for some musical styles and some modes of musical engagement, it is by no means important for all of them. For example, it is very possible to be an effective DJ, music journalist, or music producer without being able to read music.

See the Gold-MSI in Action

The test is implemented on the Concerto platform and an example implementation is available at <http://concerto.icar-project.com/v4/?wid=5&tid=1>.

trial, two of the melodies possess exactly the same interval content (i.e. the pitch relationships between notes are unchanged), whereas one melody has different interval content. The test-taker's task is to identify the 'odd-one-out', which can either come first, second, or third (Figure 2).

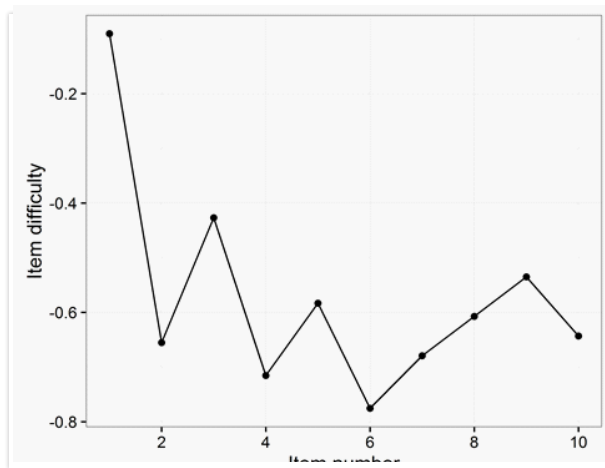


Figure 1

Simulation of the CAT item selection process for a low-ability test-taker on the melodic discrimination test. Over the course of the test, the CAT algorithm homes in on the test-taker's true ability level of -0.6.

We identified **four key cognitive processes** that underlie this test: perceptual encoding, memory retention, similarity comparison, and decision-making. In perceptual encoding, cognitive representations of the melodies are derived from the audio signal. These melody representations are retained in working memory, to allow the melodies to be compared. The comparisons take the form of similarity judgements between pairs of melodies. Finally, a decision-making process combines information from the similarity judgements to determine the final response. This cognitive model forms the basis of our AIG system. Melodies for test items are automatically generated using an automatic composition algorithm, and the difficulty of these items is then predicted on the basis of our cognitive understanding of the task. Memory encoding difficulty is varied by manipulating the length of the melodies; longer melodies place higher demands on working memory, and so are harder to retain. Difficulty of similarity comparison is varied by manipulating the degree and type of differences between the melodies; more similar melodies are harder to discriminate.

We calibrated our AIG system using response data from 425 participants, each of whom took a 10-minute online melodic discrimination test. The AIG system was then used to construct a CAT with an item bank of 1,200 items. We then investigated how the test-retest reliability (Figure 3) and standard error of the ability estimates (Figure 4) of this CAT varied for different test lengths,

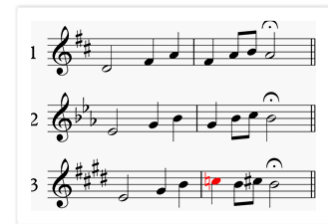


Figure 2

Example of one trial of the melodic discrimination test. Here the third melody is the 'odd-one-out'.

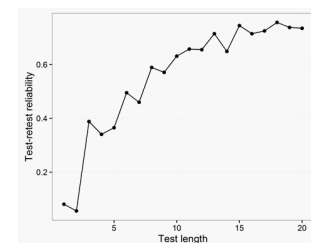


Figure 3

Test-retest reliability as a function of test length for the melodic discrimination CAT.

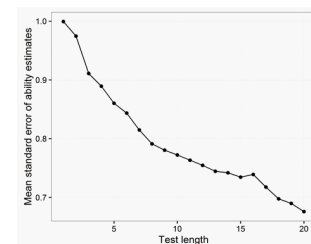


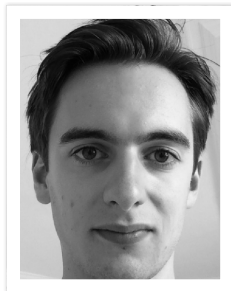
Figure 4

Standard error of ability estimates as a function of test length for the melodic discrimination CAT.

using a nationally representative sample group of 42 online test-takers. As Figure 3 demonstrates, a peak test-retest reliability of 0.75 is reached with about 15 items. On average, pre-existing melodic discrimination tests reach a similar reliability with about 30 items. Our approach therefore allows us to reduce test length by approximately half without compromising reliability. Test reliability should still be higher under controlled laboratory conditions.

In the future, we hope to improve the melodic discrimination test further by extending the range of item difficulties available and improving item difficulty predictions. We are also developing a complementary beat perception test using similar techniques of computerised adaptive testing and automatic item generation. So far, our results suggest that these psychometric techniques have exciting potential for the future of musical ability testing.

about the authors



Peter Harrison

Peter Harrison is a PhD student at Queen Mary, University of London.

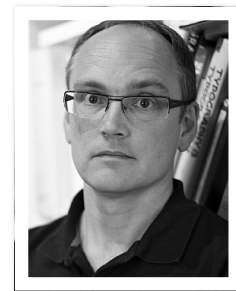
✉ [Contact Author via Email:](mailto:pharr011@gold.ac.uk)
pharr011@gold.ac.uk



Tom Collins

Tom Collins is a Visiting Assistant Professor in the Department of Psychology, Lehigh University, Pennsylvania, USA. His research interests include the development of Web-based music software and its effect on student education and work, pattern discovery in music and other domains, modelling musical expectancy, and automatic identification of high-level music-theoretic concepts.

🌐 [Author's website:](http://www.tomcollinsresearch.net/)
http://www.tomcollinsresearch.net/



Daniel Müllensiefen

Daniel is a music psychologist and member of the Music, Mind and Brain research group at Goldsmiths. His research interests include musicality and individual differences in musical abilities, the psychometrics of music, memory for music and melodies, involuntary musical imagery (i.e. 'earworms'), the perception of musical similarity, statistical models of cognition, music in advertising, corpus-based musicology, cognitive biases in musical judgement and cognitive issues related to music copyright.

🌐 [Author's website:](http://www.doc.gold.ac.uk/~mas03dm/)
http://www.doc.gold.ac.uk/~mas03dm/