

The International Cognitive Ability Resource: Development and initial validation of a public-domain measure

David M. Condon*, William Revelle

Northwestern University, Evanston, IL

Abstract

For all of its versatility and sophistication, the extant toolkit of cognitive ability measures lacks a public-domain method for large-scale, remote data collection. While the lack of copyright protection for such a measure poses a theoretical threat to test validity, the effective magnitude of this threat is unknown and can be offset by the use of modern test-development techniques. To the extent that validity can be maintained, the benefits of a public-domain resource are considerable for researchers, including: cost savings; greater control over test content; and the potential for more nuanced understanding of the correlational structure between constructs. The International Cognitive Ability Resource was developed to evaluate the prospects for such a public-domain measure and the psychometric properties of the first four item types were evaluated based on administrations to both an offline

*Correspondence concerning this article should be addressed to David M. Condon, Department of Psychology, Northwestern University, Evanston IL, 60208. Telephone number: 847-491-4515. Email: davidcondon2009@u.northwestern.edu. With thanks to Melissa Mitchell.

Please do not cite without permission.

university sample and a large online sample. Concurrent and discriminative validity analyses suggest that the public-domain status of these item types did not compromise their validity despite administration to 97,000 participants. Further development and validation of extant and additional item types is recommended.

Keywords: cognitive ability, intelligence, online assessment, psychometric validation, public-domain measures, spatial reasoning, matrix reasoning

1. Introduction

The domain of cognitive ability assessment is now populated with dozens, possibly hundreds, of proprietary measures (Camara et al., 2000; Carroll, 1993; Cattell, 1943; Eliot and Smith, 1983; Goldstein and Beers, 2004; Murphy et al., 2011). While many of these are no longer maintained or administered, the variety of tests in active use remains quite broad, providing those who want to assess cognitive abilities with a large menu of options. In spite of this diversity, however, assessment challenges persist for researchers attempting to evaluate the structure and correlates of cognitive ability. We argue that it is possible to address these challenges through the use of well-established test development techniques and report on the development and validation of an item pool which demonstrates the utility of a public-domain measure of cognitive ability for basic intelligence research. We conclude by imploring other researchers to contribute to the on-going development, aggregation and maintenance of many more item types as part of a broader, public-domain tool – the International Cognitive Ability Resource (“ICAR”).

17 **2. The Case For A Public Domain Measure**

18 To be clear, the science of intelligence has historically been well-served
19 by commercial measures. Royalty income streams (or their prospect) have
20 encouraged the development of testing “products” and have funded their on-
21 going production, distribution and maintenance for decades. These assess-
22 ments are broadly marketed for use in educational, counseling and industrial
23 contexts and their administration and interpretation is a core service for
24 many applied psychologists. Their proprietary nature is fundamental to the
25 perpetuation of these royalty streams and to the privileged status of trained
26 psychologists. For industrial and clinical settings, copyright-protected com-
27 mercial measures offer clear benefits.

28 However, the needs of primary researchers often differ from those of com-
29 mercial test users. These differences relate to issues of score interpretation,
30 test content and administrative flexibility. In the case of score interpretation,
31 researchers are considerably less concerned about the nature and quality of
32 interpretative feedback. Unlike test-takers in selection and clinical settings,
33 research participants are typically motivated by monetary rewards, course
34 credit or, perhaps, a casual desire for informal feedback about their perfor-
35 mance. This does not imply that researchers are less interested in quality
36 norming data – it is often critical for evaluating the degree to which a sample
37 is representative of a broader population. It simply means that, while many
38 commercial testing companies have attempted to differentiate their products
39 by providing materials for individual score interpretation, these materials
40 have relatively little value for administration in research contexts.

41 The motivation among commercial testing companies to provide useful

42 interpretative feedback is directly related to test content however, and the
43 nature of test content is of critical importance for intelligence researchers.
44 The typical rationale for cognitive ability assessment in research settings is
45 to evaluate the relationship between constructs and a broad range of other
46 attributes. As such, the variety and depth of a test's content are very mean-
47 ingful criteria for intelligence researchers – ones which are somewhat incom-
48 patible with the provision of meaningful interpretative feedback for each type
49 of content. In other words, the ideal circumstance for many researchers would
50 include the ability to choose from a variety of broadly-assessed cognitive abil-
51 ity constructs (or perhaps to choose a single measure which includes the as-
52 sessment of a broad variety of constructs). While this ideal can sometimes
53 be achieved through the administration of multiple commercial measures,
54 this is rarely practical due to issues of cost and/or a lack of administrative
55 flexibility.

56 The cost of administering commercial tests in research settings varies
57 considerably across measures. While published rates are typically high, many
58 companies allow for the qualified use of their copyright-protected materials
59 at reduced rates or free-of-charge in research settings (e.g., the ETS Kit
60 of Factor-Referenced Cognitive Tests (Ekstrom et al., 1976)). Variability
61 in administration and scoring procedures is similarly high across measures.
62 A small number of extant tests allow for brief, electronic assessment with
63 automated scoring conducted within the framework of proprietary software,
64 though none of these measures allow for customization of test content. The
65 most commonly-used batteries are more arduous to administer, requiring
66 one-to-one administration for over an hour followed by an additional 10 to

67 20 minutes for scoring (Camara et al., 2000). All too often, the result of
68 the combination of challenges posed by these constraints is the omission of
69 cognitive ability assessment in psychological research.

70 Several authors have suggested that the pace of scientific progress is di-
71 minished by reliance on proprietary measures (Gambardella and Hall, 2006;
72 Goldberg, 1999; Liao et al., 2008). While it is difficult to evaluate this claim
73 empirically in the context of intelligence research, the circumstances sur-
74 rounding development of the International Personality Item Pool (“IPIP”)
75 (Goldberg, 1999; Goldberg et al., 2006) provide a useful analogy. Prior to
76 the development of the IPIP, personality researchers were forced to choose
77 between validated but restrictive proprietary measures and a disorganized
78 collection of narrow-bandwidth public-domain scales (these having been de-
79 veloped by researchers who were either unwilling to deal with copyright issues
80 or whose needs were not met by the content of proprietary options). In the
81 decade ending in 2012, at least 500 journal articles and book chapters using
82 IPIP measures were published (Goldberg, 2012).

83 In fact, most of the arguments set forth in Goldberg’s (1999) proposal
84 for public-domain measures are directly applicable here. His primary point
85 was that unrestricted use of public-domain instruments would make it less
86 costly and difficult for researchers to administer scales which are flexible
87 and widely-used. Secondary benefits would include a collaborative medium
88 through which researchers could contribute to test development, refinement,
89 and validation. The research community as a whole would benefit from an
90 improved means of empirically comparing hypotheses across many diverse
91 criteria.

92 Critics of the IPIP proposal expressed concern that a lack of copyright
93 protection would impair the validity of personality measures (Goldberg et al.,
94 2006). This argument would seem even more germane for tests of cogni-
95 tive ability given the “maximal performance/typical behavior” distinction
96 between intelligence and personality measures. The widely-shared presump-
97 tion is that copyright restrictions on proprietary tests maintain validity by
98 enhancing test security. Testing materials are, in theory, only disseminated
99 to authorized users who have purchased licensed access and further dissemi-
100 nation is discouraged by the enforcement of intellectual property laws. Un-
101 fortunately, it is difficult to ascertain the extent to which test validity would
102 be compromised in the general population without these safeguards. Con-
103 cerns about disclosure have been called into question with several prominent
104 standardized tests (Field, 2012). There is also debate about the efficacy of in-
105 tellectual property laws for protection against the unauthorized distribution
106 of testing materials via the internet (Field, 2012; Kaufmann, 2009; McCaffrey
107 and Lynch, 2009). Further evaluation of the relationship between copyright-
108 protection and test validity seems warranted by these concerns, particularly
109 for research applications where individual outcomes are less consequential.

110 Fortunately, copyright protection is not a prerequisite for test validity.
111 Modern item-generation techniques (Arendasy et al., 2006; Dennis et al.,
112 2002) present an alternate strategy that is less dependent on test security.
113 Automatic item-generation makes use of algorithms which dictate the param-
114 eters of new items with predictable difficulty and in many alternate forms.
115 These techniques allow for the creation of item types where the universe of
116 *possible* items is very large. This, in turn, reduces the threat to validity that

117 results from item disclosure. It can even be used to enhance test validity un-
118 der administration paradigms that expose participants to sample items prior
119 to testing and use alternate forms during assessment as this methodology
120 reduces the effects of differential test familiarity across participants.

121 While automatic item-generation techniques represent the optimal method
122 for developing public-domain cognitive ability items, this approach is often
123 considerably more complicated than traditional development methods and it
124 may be some time before a sizable number of automatically-generated item
125 types is available for use in the public domain. For item types developed by
126 traditional means, the maintenance of test validity depends on implementa-
127 tion of the more practical protocols used by commercial measures (i.e., those
128 which do not invoke the credible threat of legal action). A public domain
129 resource should set forth clear expectations for researchers regarding appro-
130 priate and ethical usage and make use of “warnings for nonprofessionals”
131 (Goldberg et al., 2006). Sample test items should be made easily available
132 to the general public to further discourage wholesale distribution of testing
133 materials. Given the current barriers to enforcement for intellectual property
134 holders, these steps are arguably commensurate with protocols in place for
135 copyright-protected commercial measures.

136 To the extent that traditional and automatic item-generation methods
137 maintain adequate validity, there are many applications in which a non-
138 proprietary measure would be useful. The most demanding of these applica-
139 tions would involve distributed, un-proctored assessments *in situ*, presumably
140 conducted via online administration. Validity concerns would be most acute
141 in these situations as there would be no safeguards against the use of external

142 resources, including those available on the internet.

143 The remainder of this paper is dedicated to the evaluation of a public-
144 domain measure developed for use under precisely these circumstances. This
145 measure, the International Cognitive Ability Resource (“ICAR”), has been
146 developed in stages over several years and further development is on-going.
147 The first four item types (described below) were initially designed to provide
148 an estimation of general cognitive ability for participants completing person-
149 ality surveys at SAPA-Project.org, previously test.personality-project.org.

150 The primary goals when developing these initial item types were to: (1)
151 briefly assess a small number of cognitive ability domains which were rela-
152 tively distinct from one another (though considerable overlap between scores
153 on the various types was anticipated); (2) avoid the use of “timed” items in
154 light of potential technical issues resulting from telemetric assessment (Wilt
155 et al., 2011); and (3) avoid item content that could be readily referenced else-
156 where given the intended use of un-proctored online administrations. The
157 studies described below were conducted to evaluate the degree to which these
158 goals of item development were achieved.

159 The first study evaluated the item characteristics, reliability and struc-
160 tural properties of a 60-item ICAR measure. The second study evaluated
161 the validity of the ICAR items when administered online in the context of
162 self-reported achievement test scores and university majors. The third study
163 evaluated the construct validity of the ICAR items when administered offline,
164 using a brief commercial measure of cognitive ability.

165 **3. Study 1**

166 We investigated the structural properties of the initial version of the In-
167 ternational Cognitive Ability Resource based on internet administration to a
168 large international sample. This investigation was based on 60 items repre-
169 senting four item types developed in various stages since 2006 (and does not
170 include deprecated items or item types currently under development). We
171 hypothesized that the factor structure would demonstrate four distinct but
172 highly correlated factors, with each type of item represented by a separate
173 factor. This implied that, while individual items might demonstrate moder-
174 ate or strong cross-loadings, the primary loadings would be consistent among
175 items of each type.

176 *3.1. Method*

177 *3.1.1. Participants*

178 Participants were 96,958 individuals (66% female) from 199 countries who
179 completed an online survey at SAPA-project.org (previously test.personality-
180 project.org) between August 18, 2010 and May 20, 2013 in exchange for
181 customized feedback about their personalities. All data were self-reported.
182 The mean self-reported age was 26 years ($sd = 10.6$, median = 22) with a
183 range from 14 to 90 years. Educational attainment levels for the partici-
184 pants are given in Table 1. Most participants were current university or sec-
185 ondary school students, although a wide range of educational attainment lev-
186 els were represented. Among the 75,740 participants from the United States
187 (78.1%), 67.5% identified themselves as White/Caucasian, 10.3% as African-
188 American, 8.5% as Hispanic-American, 4.8% as Asian-American, 1.1% as

189 Native-American, and 6.3% as multi-ethnic (the remaining 1.5% did not
190 specify). Participants from outside the United States were not prompted
191 for information regarding race/ethnicity.

192 3.1.2. Measures

193 Four item types from the International Cognitive Ability Resource were
194 administered, including: 9 Letter and Number Series items, 11 Matrix Rea-
195 soning items, 16 Verbal Reasoning items and 24 Three-Dimensional Rotation
196 items. A 16 item subset of the measure, hereafter referred to as the *ICAR*
197 *Sample Test*, is included as Appendix A in the Supplemental Materials.¹
198 Letter and Number Series items prompt participants with short digit or let-
199 ter sequences and ask them to identify the next position in the sequence
200 from among six choices. Matrix Reasoning items contain stimuli that are
201 similar to those used in Raven’s Progressive Matrices. The stimuli are 3x3
202 arrays of geometric shapes with one of the nine shapes missing. Partici-
203 pants are instructed to identify which of six geometric shapes presented as
204 response choices will best complete the stimuli. The Verbal Reasoning items
205 include a variety of logic, vocabulary and general knowledge questions. The
206 Three-Dimensional Rotation items present participants with cube renderings
207 and ask participants to identify which of the response choices is a possible
208 rotation of the target stimuli. None of the items were timed in these admin-

¹In addition to the sample items available in Appendix A, the remaining ICAR items can be accessed through ICAR-Project.org. A sample data set based on the items listed in Appendix A is also available (‘iqitems’) through the *psych* package (Revelle, 2013) in the R computing environment (R Core Team, 2013).

209 istrations as untimed administration was expected to provide more stringent
210 and conservative evaluation of the items' utility when given online (there
211 are no specific reasons precluding timed administrations of the ICAR items,
212 whether online or offline).

213 Participants were administered 12 to 16 item subsets of the 60 ICAR
214 items using the Synthetic Aperture Personality Assessment (“SAPA”) tech-
215 nique (Revelle et al., 2010), a variant of matrix sampling procedures discussed
216 by Lord (1955). The number of items administered to each participant varied
217 over the course of the sampling period and was independent of participant
218 characteristics. The number of administrations for each item varied con-
219 siderably (median = 21,764) as did the number of pairwise administrations
220 between any two items in the set (median = 2,610). This variability reflected
221 the introduction of newly developed items over time and the fact that item
222 sets include unequal numbers of items. The minimum number of pairwise
223 administrations among items (422) provided sufficiently high stability in the
224 covariance matrix for the structural analyses described below (Kenny, 2012).

225 *3.1.3. Analyses*

226 Internal consistency measures were assessed by using the Pearson correla-
227 tions between ICAR items to calculate α , ω_h , and ω_{total} reliability coefficients
228 (Revelle, 2013; Revelle and Zinbarg, 2009; Zinbarg et al., 2005). The use of
229 tetrachoric correlations for reliability analyses is discouraged on the grounds
230 that it typically over-estimates both alpha and omega (Revelle and Condon,
231 2012).

232 Two latent variable exploratory factor analyses (“EFA”) were conducted
233 to evaluate the structure of the ICAR items. The first of these included all

234 60 items (9 Letter and Number Series items, 11 Matrix Reasoning items,
235 16 Verbal Reasoning items and 24 Three-Dimensional Rotation items). A
236 second EFA was required to address questions regarding the structural im-
237 pact of including disproportionate numbers of items by type. This was done
238 by using only the subset of participants ($n = 4,574$) who were administered
239 the 16 item *ICAR Sample Test*. This subset included four items each from
240 the four ICAR item types. These items were selected as a representative set
241 on the basis of their difficulty relative to the full set of 60 items and their
242 factor loadings relative to other items of the same type. Note that the factor
243 analysis of this 16 item subset was not independent from that conducted on
244 the full 60 item set. EFA results were then used to evaluate the omega hier-
245 archical general factor saturation (Revelle and Zinbarg, 2009; Zinbarg et al.,
246 2006) of the 16 item *ICAR Sample Test*.

247 Both of these exploratory factor analyses were based on the Pearson cor-
248 relations between scored responses using Ordinary Least Squares (“OLS”) re-
249 gression models with oblique rotation (Revelle, 2013). The factoring method
250 used here minimizes the χ^2 value rather than minimizing the sum of the
251 squared residual values (as is done by default with most statistical software).
252 Note that in cases where the number of administrations is consistent across
253 items, as with the 16 item *ICAR Sample Test*, these methods are identical.
254 The methods differ in cases where the number of pairwise administrations
255 between items varies because the squared residuals are weighted by sample
256 size rather than assumed to be equivalent across variables. Goodness-of-fit
257 was evaluated using the Root Mean Square of the Residual, the Root Mean
258 Squared Error of Approximation (Hu and Bentler, 1999), and the Tucker

259 Lewis Index of factoring reliability (Kenny, 2012; Tucker and Lewis, 1973).

260 Analyses based on two-parameter Item Response Theory (Baker, 1985;
261 Embretson, 1996; Revelle, 2013) were used to evaluate the unidimensional
262 relationships between items on several levels, including (1) all 60 items, (2)
263 each of the four item types independently, and (3) for the 16 item *ICAR*
264 *Sample Test*. In these cases, the tetrachoric correlations between items were
265 used. These procedures allow for estimation of the correlations between items
266 as if they had been measured continuously (Uebersax, 2000).

267 3.2. Results

268 Descriptive statistics for all 60 ICAR items are given in Table 2. Mean
269 values indicate the proportion of participants who provided the correct re-
270 sponse for an item relative to the total number of participants who were
271 administered that item. The Three-Dimensional Rotation items had the
272 lowest proportion of correct responses ($m = 0.19$, $sd = 0.08$), followed by
273 Matrix Reasoning ($m = 0.52$, $sd = 0.15$), then Letter and Number Series (m
274 $= 0.59$, $sd = 0.13$), and Verbal Reasoning ($m = 0.64$, $sd = 0.22$). Internal
275 consistencies for the ICAR item types are given in Table 3. These values
276 are based on the composite correlations between items as individual partici-
277 pants completed only a subset of the items (as is typical when using SAPA
278 sampling procedures).

279 Results from the first exploratory factor analysis using all 60 items sug-
280 gested factor solutions of three to five factors based on inspection of the scree
281 plots in Figure 1. The fit statistics were similar for each of these solutions.
282 The four factor model was slightly superior in fit (RMSEA = 0.058, RMSR
283 = 0.05) and reliability (TLI = 0.71) to the three factor model (RMSEA =

284 0.059, RMSR = 0.05, TLI = 0.7) and was slightly inferior to the five factor
285 model (RMSEA = 0.055, RMSR = 0.05, TLI = 0.73). Factor loadings and
286 the correlations between factors for each of these solutions are included in
287 the supplementary materials (see Supplementary Tables 1 to 6).

288 The second EFA, based on a balanced number of items by type, demon-
289 strated very good fit for the four-factor solution (RMSEA = 0.014, RMSR
290 = 0.01, TLI = 0.99). Factor loadings by item for the four-factor solution
291 are shown in Table 4. Each of the item types was represented by a different
292 factor and the cross-loadings were small. Correlations between factors (Table
293 5) ranged from 0.41 to 0.70.

294 General factor saturation for the 16 item *ICAR Sample Test* is depicted
295 in Figures 2 and 3. Figure 2 shows the primary factor loadings for each
296 item consistent with the values presented in Table 4 and also shows the
297 general factor loading for each of the second-order factors. Figure 3 shows
298 the general factor loading for each item and the residual loading of each item
299 to its primary second-order factor after removing the general factor.

300 The results of IRT analyses for the 16 item *ICAR Sample Test* are pre-
301 sented in Table 6 as well as Figures 4 and 5. Table 6 provides item information
302 across levels of the latent trait and summary information for the test as a
303 whole. The item information functions are depicted graphically in Figure 4.
304 Figure 5 depicts the test information function for the *ICAR Sample Test* as
305 well as reliability in the vertical axis on the right (reliability in this context
306 is calculated as one minus the reciprocal of the test information). The re-
307 sults of IRT analyses for the full 60 item set and for each of the item types
308 independently are available in the supplementary materials (Supplementary

309 Tables 7 to 11). The pattern of results was similar to those for the *ICAR*
310 *Sample Test* in terms of the relationships between item types and the spread
311 of item difficulties across levels of the latent trait, though the reliability was
312 higher for the full 60 item set across the range of difficulties (Supplementary
313 Figure 1).

314 3.3. Discussion

315 A key finding from Study 1 relates to the broad range of means and
316 standard deviations for the ICAR items as these values demonstrated that
317 the un-proctored and untimed administration of cognitive ability items online
318 does not lead to uniformly high scores with insufficient variance. To the
319 contrary, all of the Three-Dimensional Rotation items and more than half
320 of all 60 items were answered incorrectly more often than correctly and the
321 weighted mean for all items was only 0.53. This point was further supported
322 by the IRT analyses in that the item information functions demonstrate a
323 relatively wide range of item difficulties.

324 Internal consistency was good for the Three-Dimensional Rotation item
325 type, adequate for the Letter and Number Series and the Verbal Reason-
326 ing item types, and marginally adequate for the Matrix Reasoning item
327 type. This suggests that the 11 Matrix Reasoning items were not uni-
328 formly measuring a singular latent construct whereas performance on the
329 Three-Dimensional Rotation items was highly consistent. For the compos-
330 ites based on both 16 and 60 items however, internal consistencies were ad-
331 equate ($\alpha=0.81$; $\omega_{total}=0.83$) and good ($\alpha=0.93$; $\omega_{total}=0.94$), respectively.
332 While higher reliabilities reflect the greater number of items in the ICAR60,
333 it should be noted that the general factor saturation was slightly higher for

334 the shorter 16-item measure (ICAR16 $\omega_h=0.66$; ICAR60 $\omega_h=0.61$). When
335 considered as a function of test information, reliability was generally ade-
336 quate across a wide range of latent trait levels, and particularly good within
337 approximately ± 1.5 standardized units from the mean item difficulty. All of
338 the factor analyses demonstrated evidence of both a positive manifold among
339 items and high general factor saturation for each of the item types. In the
340 four factor solution for the 16 item scale, the Verbal Reasoning and the Letter
341 and Number Series factors showed particularly high ‘*g*’ loadings (0.8).

342 **4. Study 2**

343 Following the evidence for reliable variability in ICAR scores in Study
344 1, it was the goal of Study 2 to evaluate the validity of these scores when
345 using the same administration procedures. While online administration pro-
346 tocols precluded validation against copyrighted commercial measures, it was
347 possible to evaluate the extent to which ICAR scores correlated with (1) self-
348 reported achievement test scores and (2) published rank orderings of mean
349 scores by university major. In the latter case, ICAR scores were expected
350 to demonstrate group discriminant validity by correlating highly with the
351 rank orderings of mean scores by university major as previously described by
352 the Educational Testing Service (Educational Testing Service, 2010) and the
353 College Board (College Board, 2012).

354 In the former case, ICAR scores were expected to reflect a similar rela-
355 tionship with achievement test scores as extant measures of cognitive ability.
356 Using data from the National Longitudinal Study of Youth 1979, Frey and
357 Detterman (2004) reported simple correlations between the SAT and the

358 Armed Services Vocational Aptitude Battery ($r = 0.82$, $n = 917$) and sev-
359 eral additional IQ measures ($rs = 0.53 - 0.82$) with smaller samples ($ns =$
360 $15 - 79$). In a follow-up study with a university sample, Frey and Detterman
361 (2004) evaluated the correlation between combined SAT scores and Raven's
362 Progressive Matrices scores, finding an uncorrected correlation of 0.48 ($p <$
363 $.001$) and a correlation after correcting for restriction of range of 0.72. Similar
364 analyses with ACT composite scores (Koenig et al., 2008) showed a correla-
365 tion of 0.77 ($p < .001$) with the ASVAB, an uncorrected correlation with the
366 Raven's Advanced Progressive Matrices of 0.61 ($p < .001$), and a correlation
367 corrected for range restriction with the Raven's APM of 0.75.

368 Given the breadth and duration of assessment for the ASVAB, the SAT
369 and the ACT, positive correlations of a lesser magnitude were expected be-
370 tween the ICAR scores and the achievement tests than were previously re-
371 ported with the ASVAB. Correlations between the Raven's APM and the
372 achievement test scores were expected to be more similar to the correlations
373 between the achievement test scores and the ICAR scores, though it was not
374 possible to estimate the extent to which the correlations would be affected
375 by methodological differences (i.e., the un-proctored online administration of
376 relatively few ICAR items and the use of self-reported, rather than indepen-
377 dently verified, achievement test scores as described in the Methods section
378 below).

379 *4.1. Method*

380 *4.1.1. Participants*

381 The 34,229 participants in Study 2 were a subset of those used for Study 1,
382 chosen on the basis of age and level of educational attainment. Participants

383 were 18 to 22 years old ($m = 19.9$, $s.d. = 1.3$, median = 20). Approximately
384 91% of participants had begun but not yet attained an undergraduate de-
385 gree; the remaining 9% had attained an undergraduate degree. Among the
386 26,911 participants from the United States, 67.1% identified themselves as
387 White/Caucasian, 9.8% as Hispanic-American, 8.4% as African-American,
388 6.0% as Asian-American, 1.0% as Native-American, and 6.3% as multi-ethnic
389 (the remaining 1.5% did not specify).

390 4.1.2. Measures

391 Both the sampling method and the ICAR items used in Study 2 were
392 identical to the procedures described in Study 1, though the total item ad-
393 ministrations (median = 7,659) and pairwise administrations (median = 906)
394 were notably fewer given that the participants in Study 2 were a sub-sample of
395 those in Study 1. Study 2 also used self-report data for three additional vari-
396 ables collected through SAPA-project.org: (1) participants' academic major
397 on the university level, (2) their achievement test scores, and (3) participants'
398 scale scores based on randomly administered items from the Intellect scale of
399 the "100-Item Set of IPIP Big-Five Factor Markers" (Goldberg, 2012). For
400 university major, participants were allowed to select only one option from
401 147 choices, including "undecided" ($n = 3,460$) and several categories of
402 "other" based on academic disciplines. For the achievement test scores, par-
403 ticipants were given the option of reporting 0, 1, or multiple types of scores,
404 including: SAT Critical Reading ($n = 7,404$); SAT Mathematics ($n = 7,453$);
405 and the ACT ($n = 12,254$). Intellect scale scores were calculated using IRT
406 procedures, assuming unidimensionality for the Intellect items only (items
407 assessing Openness were omitted). Based on composites of the Pearson cor-

408 relations between items without imputation of missing values, the Intellect
409 scale had an α of 0.74, an ω_h of 0.60, and an ω_{total} of 0.80. The median
410 number of pairwise administrations for these items was 4,475.

411 *4.1.3. Analyses*

412 Two distinct methods were used to calculate the correlations between the
413 achievement test scores and the ICAR scores in order to evaluate the effects
414 of two different corrections. The first method used ICAR scale scores based
415 on composites of the tetrachoric correlations between ICAR items (compos-
416 ites are used because each participant was administered 16 or fewer items).
417 The correlations between these scale scores and the achievement test scores
418 were then corrected for reliability. The α reliability coefficients reported in
419 Study 1 were used for the ICAR scores. For the achievement test scores,
420 the need to correct for reliability was necessitated by the use of self-reported
421 scores. Several researchers have demonstrated the reduced reliability of self-
422 reported scores in relation to official test records (Cassady, 2001; Cole and
423 Gonyea, 2009; Kuncel et al., 2005; Mayer et al., 2006), citing participants'
424 desire to misrepresent their performance and/or memory errors as the most
425 likely causes. Despite these concerns, the reported correlations between self-
426 reported and actual scores suggest that the rank-ordering of scores is main-
427 tained, regardless of the magnitude of differences (Cole and Gonyea, 2009;
428 Kuncel et al., 2005; Mayer et al., 2006). Reported correlations between self-
429 reported and actual scores have ranged from 0.74 to 0.86 for the SAT -
430 Critical Reading section, 0.82 to 0.88 for the SAT - Mathematics, and 0.82
431 to 0.89 for the SAT - Combined (Cole and Gonyea, 2009; Kuncel et al., 2005;
432 Mayer et al., 2006). Higher correlations were found by Cole and Gonyea

433 (2009) for the ACT Composite (0.95). The Study 2 sample approximated
434 the samples on which these reported correlations were based in that (1) par-
435 ticipants were reminded about the anonymity of their responses and (2) the
436 age range of participants was limited to 18 to 22 years. The weighted mean
437 values from these findings (SAT-CR = 0.86; SAT-M = 0.88; SAT-Combined
438 = 0.88; ACT = 0.95) were used as reliability coefficients for the achievement
439 test scores when correcting correlations between the achievement tests and
440 other measures (ICAR scores and the IPIP-100 Intellect scores).

441 The second method for calculating correlations between ICAR scores and
442 achievement test scores used IRT-based (2PL) scoring (Revelle, 2013). Scale
443 scores for each item type and the full test were calculated for each partici-
444 pant, and these scale scores were then correlated with the achievement test
445 scores. In this case, corrections were made to address the potential for an
446 incidental selection effect due to optional reporting of achievement test scores
447 (Cassady, 2001; Frucot and Cook, 1994). 52.5% of participants in Study 2 did
448 not report any achievement test scores; 10.1% reported scores for all three
449 (SAT - CR, SAT - M, and ACT). These circumstances would result in an
450 incidental selection effect if the correlations between self-reported achieve-
451 ment test scores and the ICAR measures were affected by the influence of
452 a third variable on one or both measures (Sackett and Yang, 2000). The
453 so-called “third” variable in this study likely represented a composite of la-
454 tent factors which are neither ergodic nor quantifiable but which resulted
455 in group differences between those who reported their scores and those who
456 did not. If the magnitude of differences in achievement test scores between
457 groups were non-trivial, the effect on the overall correlations would also be

458 non-trivial given the proportion of participants not reporting. The need
459 for correction procedures in this circumstance was elaborated by both Pear-
460 son (1903) and Thorndike (1949), though the methods employed here were
461 developed in the econometrics literature and are infrequently used by psy-
462 chologists (Sackett and Yang, 2000). Clark and Houle (2012) and Cuddeback
463 et al. (2004) provide useful illustrations of these procedures. The two-step
464 method of the “Heckman correction” (Greene, 2008; Heckman, 1976, 1979;
465 Toomet and Henningsen, 2008) was used to evaluate and correct for selection
466 effects where warranted using IPIP-100 Intellect scores.

467 In addition to these analyses of the relationship between ICAR scores
468 and achievement test scores, the Study 2 sample was used to evaluate the
469 correlations between the ICAR items and the published rank orderings of
470 mean scores by university major. This was done using IRT-based ICAR
471 scores when grouped by academic major on the university level. These were
472 evaluated relative to similar data sets published by the Educational Testing
473 Service (Educational Testing Service, 2010) and the College Board (College
474 Board, 2012) for the GRE and SAT, respectively. GRE scores were based on
475 group means for 287 “intended graduate major” choices offered to fourth-year
476 university students and non-enrolled graduates who took the GRE between
477 July 1, 2005 and June 30, 2008 ($N = 569,000$). These 287 groups were
478 consolidated with weighting for sample size in order to match the 147 uni-
479 versity major choices offered with the ICAR. Of these 147 majors, only the
480 91 with $n > 20$ were used. SAT scores were based on group means for 38
481 “intended college major” choices offered to college-bound seniors in the high
482 school graduating class of 2012 ($N = 1,411,595$). In this case, the 147 uni-

483 versity major choices offered with the ICAR were consolidated to match 29
484 of the choices offered with the SAT. The 9 incompatible major choices col-
485 lectively represented only 1.3% of the SAT test-takers. The omitted majors
486 were: Construction Trades; Mechanic and Repair Technologies/Technician;
487 Military Technologies and Applied Sciences; Multi/Interdisciplinary Stud-
488 ies; Precision Production; Security and Protective Services; Theology and
489 Religious Vocations; Other; and Undecided.

490 *4.2. Results*

491 Descriptive statistics for the self-reported achievement test scores are
492 shown in Table 7. Correlations between self-reported achievement test scores
493 and ICAR scale scores calculated using composites of the tetrachoric corre-
494 lations are shown in Table 8, with uncorrected correlations shown below the
495 diagonal and the correlations corrected for reliability shown above the diag-
496 onal. Reliabilities for each measure are given on the diagonal. Correlations
497 between composites which were not independent have been omitted. Cor-
498 rected correlations between the achievement test scores and both the 16 and
499 60 item ICAR composites ranged from 0.52 - 0.59 ($ses \leq 0.016$).²

500 Table 9 presents the correlations between the self-reported achievement
501 test scores and the IRT-based ICAR scores, with the uncorrected correlations
502 below the diagonal and the correlations corrected for incidental selection

²The standard error of the composite scores are a function of both the number of items and the number of participants who took each pair of items (Revelle and Brown, 2013). Estimates of the standard errors can be identified through the use of bootstrapping procedures to derive estimates of the confidence intervals of the correlations (Revelle, 2013). In this case, the confidence intervals were estimated based on 100 sampling iterations.

503 effects above the diagonal. Correlations between non-independent scores
504 were omitted. Scores for the ICAR measures were based on a mean of 2 to 4
505 responses for each of the item types (mean number of LN items administered
506 = 3.2, $sd = 1.3$; MR items $m = 2.8$, $sd = 1.1$; R3D items $m = 2.0$, $sd =$
507 1.5 ; VR items $m = 4.3$, $sd = 2.2$) and 12 to 16 items for the ICAR60 scores
508 ($m = 12.4$, $sd = 3.8$). Corrected correlations between the achievement test
509 scores and ICAR60 ranged from 0.44 to 0.47 ($ses \leq 0.016$).

510 Tables 10 and 11 contain group-level correlations using mean scores for
511 university major. Table 10 shows the correlations between the published
512 norms for the SAT, the mean self-reported SAT scores for each major in the
513 Study 2 sample, and the mean IRT-based ICAR scores for each major in the
514 Study 2 sample. The correlation between mean ICAR scores by major and
515 mean combined SAT scores by major in the published norms was 0.75 ($se =$
516 0.147). Table 11 shows the correlations between the published norms for the
517 GRE by major and the IRT-based ICAR scores for the corresponding majors
518 in the Study 2 sample (self-reported GRE scores were not collected). The
519 correlation between mean ICAR scores by major and mean combined GRE
520 scores by major in the published norms was 0.86 ($se = 0.092$).

521 4.3. Discussion

522 After correcting for the “reliability” of self-reported scores, the 16 item
523 *ICAR Sample Test* correlated 0.59 with combined SAT scores and 0.52 with
524 the ACT composite. Correlations based on the IRT-based ICAR scores were
525 lower though these scores were calculated using even fewer items; correlations
526 were 0.47 and 0.44 with combined SAT scores and ACT composite scores
527 respectively based on an average of 12.4 ICAR60 items answered per partic-

528 ipant. As expected, these correlations were smaller than those reported for
529 longer cognitive ability measures such as the ASVAB and the Raven's APM
530 (Frey and Detterman, 2004; Koenig et al., 2008).

531 The ICAR items demonstrated strong group discriminant validity on the
532 basis of university majors. This indicates that the rank ordering of mean
533 ICAR scores by major is strongly correlated with the rank ordering of mean
534 SAT scores and mean GRE scores. Consistent with the individual-level cor-
535 relations, the group-level correlations were higher between the ICAR subtests
536 and the mathematics subtests of the SAT and the GRE relative to the verbal
537 subtests.

538 **5. Study 3**

539 The goal of the third study was to evaluate the construct validity of the
540 ICAR items against a commercial measure of cognitive ability. Due to the
541 copyrights associated with commercial measures, these analyses were based
542 on administration to an offline sample of university students rather than an
543 online administration.

544 *5.1. Method*

545 *5.1.1. Participants*

546 Participants in Study 3 were 137 college students (76 female) enrolled at
547 a selective private university in the midwestern United States. Students par-
548 ticipated in exchange for credit in an introductory psychology course. The
549 mean age of participants in this sample was 19.7 years ($sd = 1.2$, median =
550 20) with a range from 17 to 25 years. Within the sample, 67.2% reported

551 being first-year students, 14.6% second-year students, 8.0% third-year stu-
552 dents and the remaining 10.2% were in their fourth year or beyond. With
553 regards to ethnicity, 56.2% identified themselves as White/Caucasian, 26.3%
554 as Asian-American, 4.4% as African-American, 4.4% as Hispanic-American,
555 and 7.3% as multi-ethnic (the remaining 1.5% did not specify).

556 5.1.2. Measures

557 Participants in the university sample were administered the 16 item *ICAR*
558 *Sample Test*. The presentation order of these 16 items was randomized across
559 participants. Participants were also administered the *Shipley-2*, which is a
560 2009 revision and restandardization of the *Shipley Institute of Living Scale*
561 (Shipley et al., 2009, 2010). The *Shipley-2* is a brief measure of cognitive
562 functioning and impairment that most participants completed in 15 to 25
563 minutes. While the *Shipley-2* is a timed test, the majority of participants
564 stopped working before using all of the allotted time. The *Shipley-2* has
565 two administration options. Composite A ($n = 69$) includes a vocabulary
566 scale designed to assess crystallized skills and an abstraction scale designed
567 to assess fluid reasoning skills (Shipley et al., 2009). Composite B ($n = 68$)
568 includes the same vocabulary scale and a spatial measure of fluid reasoning
569 called the “Block Patterns” scale (Shipley et al., 2009). All three scales in-
570 cluded several items of low difficulty with little or no variance in this sample.
571 After removal of items without variance, internal consistencies were low for
572 the Abstraction scale (10 of 25 items removed, $\alpha = 0.37$; $\omega_{total} = 0.51$) and
573 the Vocabulary scale (7 of 40 items removed, $\alpha = 0.61$; $\omega_{total} = 0.66$). The
574 Block Patterns scale had fewer items without variance (3 of 26) and adequate
575 consistency ($\alpha = 0.83$, $\omega_{total} = 0.88$). Internal consistencies were calculated

576 using Pearson correlations between items.

577 5.1.3. Analyses

578 Correlations were evaluated between scores on the *ICAR Sample Test* and
579 a brief commercial measure of cognitive ability, the *Shibley-2*. Two types of
580 corrections were relevant to these correlations; one for the restriction of range
581 among scores and a second for reliability. The prospect of range restriction
582 was expected on the grounds that participants in the sample were students at
583 a highly selective university. The presence of restricted range was evaluated
584 by looking for reduced variance in the sample relative to populations with
585 similar characteristics. In this case, the university sample was evaluated
586 relative to the online sample. Where present, the appropriate method for
587 correcting this type of range restriction uses the following equation (case 2c
588 from Sackett and Yang, 2000) (Bryant and Gokhale, 1972; Alexander, 1990):

$$\hat{\rho}_{xy} = r_{xy}(s_x/S_x)(s_y/S_y) \pm \sqrt{[1 - (s_x/S_x)^2][1 - (s_y/S_y)^2]} \quad (1)$$

589 where s_x and s_y are the standard deviations in the restricted sample, S_x
590 and S_y are the standard deviations in the unrestricted sample and the \pm
591 sign is conditional on the direction of the relationship between the selection
592 effect and each of the variables, x and y . When correcting for reliability, the
593 published reliabilities (Shibley et al., 2010) were used for each of the *Shibley-*
594 *2* composites (0.925 for Composite A and 0.93 for Composite B) instead of
595 the reliabilities within the sample due to the large number of items with little
596 or no variance.

597 5.2. Results

598 The need to correct for restriction of range was indicated by lower stan-
599 dard deviations of scores on all of the subtests and composites for the *Shiple-*
600 *2* and the *ICAR Sample Test*. Table 12 shows the standard deviation of scores
601 for the participants in Study 3 (the “restricted” sample) and the reference
602 scores (the “unrestricted” samples).

603 Correlations between the ICAR scores and *Shiple-2* scores are given in
604 Table 13, including the uncorrected correlations, the correlations corrected
605 for range restriction and the correlations corrected for reliability and range re-
606 striction. The range and reliability corrected correlations between the *ICAR*
607 *Sample Test* and the *Shiple-2* composites were nearly identical at 0.81 and
608 0.82 ($se = 0.10$).

609 5.3. Discussion

610 Correlations between the ICAR scores and the *Shiple-2* were comparable
611 to those between the *Shiple-2* and other measures of cognitive ability. The
612 correlations after correcting for reliability and restricted range between the
613 16 item *ICAR Sample Test* and *Shiple-2* composite A and B were 0.82
614 and 0.81, respectively. Correlations between *Shiple-2* composite A and B
615 were 0.64 and 0.60 with the *Wonderlic Personnel Test*, 0.77 and 0.72 with
616 the Full-Scale IQ scores for the *Wechsler Abbreviated Scale of Intelligence*
617 in an adult sample, and 0.86 and 0.85 with the Full-Scale IQ scores for the
618 *Wechsler Adult Intelligence Scale* (Shiple et al., 2010).

619 **6. General Discussion**

620 Reliability and validity data from these studies suggest that a public-
621 domain measure of cognitive ability is a viable option. More specifically, they
622 demonstrate that brief, un-proctored, and untimed administrations of items
623 from the International Cognitive Ability Resource are moderately-to-strongly
624 correlated with measures of cognitive ability and achievement. While this
625 method of administration is inherently less precise and exhaustive than many
626 traditional assessment methods, it offers many benefits. Online assessment
627 allows for test administration at any time of day, in any geographic location,
628 and over any type of internet-enabled electronic device. These administra-
629 tions can be conducted either with or without direct interaction with the
630 research team. Measures constructed with public-domain item types like
631 those described here can be easily customized for test length and content
632 as needed to match the research topic under evaluation. All of this can be
633 accomplished without the cost, licensing, training, and software needed to
634 administer the various types of copyright-protected commercial measures.

635 These data also suggest that there are many ways in which the ICAR
636 can be improved. With regard to the existing item types, more - and more
637 difficult - items are needed for all of the item types except perhaps the Three-
638 Dimensional Rotation items. While the development of additional Letter and
639 Number Series items can be accomplished formulaically, item development
640 procedures for the Verbal Reasoning items is complicated by the need for
641 items to be resistant to basic internet word searches. The Matrix Reasoning
642 items require further structural analyses before further item development as
643 these items demonstrated less unidimensionality than the other three item

644 types. This may be appropriate if they are to be used as a measure of
645 general cognitive ability, but it remains important to identify the ways in
646 which these items assess subtly different constructs. This last point relates
647 to the additional need for analyses of differential item functioning for all of
648 the item types and the test as a whole.

649 The inclusion of many more item types in the ICAR is also needed as is
650 more extensive validation of new and existing item types. The most useful
651 additions in the near term would include item types which assess constructs
652 distinct from the four item types described here. Several such item types
653 are in various stages of development and piloting by the authors and their
654 collaborators. These item types should be augmented with extant, public-
655 domain item types when feasible.

656 **7. Conclusion**

657 Public-domain measures of cognitive ability have considerable potential.
658 We propose that the International Cognitive Ability Resource provides a
659 viable foundation for collaborators who are interested in contributing ex-
660 tant or newly-developed public-domain tools. To the extent that these tools
661 are well-suited for online administration, they will be particularly useful for
662 large-scale cognitive ability assessment and/or use in research contexts be-
663 yond the confines of traditional testing environments. As more item types
664 become available, the concurrent administration of ICAR item types will
665 become increasingly valuable for researchers studying the structure of cog-
666 nitive abilities on both the broad, higher-order levels (e.g., spatial and verbal
667 abilities) as well as the relatively narrow (e.g., more closely related abilities

668 such as two- and three-dimensional rotation). The extent to which a public-
669 domain resource like the ICAR fulfills this potential ultimately depends on
670 the researchers for whom it offers the highest utility. We entreat these poten-
671 tial users to consider contributing to its on-going development, improvement,
672 validation and maintenance.

673 **References**

- 674 Alexander, R. A. (1990). Correction formulas for correlations restricted by
675 selection on an unmeasured variable. *Journal of Educational Measurement*,
676 27(2):187–189.
- 677 Arendasy, M., Sommer, M., Gittler, G., and Hergovich, A. (2006). Auto-
678 matic generation of quantitative reasoning items. *Journal of Individual*
679 *Differences*, 27(1):2–14.
- 680 Baker, F. B. (1985). *The basics of item response theory*. Heinemann Educa-
681 tional Books, Portsmouth, NH.
- 682 Bryant, N. D. and Gokhale, S. (1972). Correcting correlations for restrictions
683 in range due to selection on an unmeasured variable. *Educational and*
684 *Psychological Measurement*, 32(2):305–310.
- 685 Camara, W. J., Nathan, J. S., and Puente, A. E. (2000). Psychological test
686 usage: Implications in professional psychology. *Professional Psychology:*
687 *Research and Practice*, 31(2):141–154.
- 688 Carroll, J. B. (1993). *Human cognitive abilities: A survey of factor-analytic*
689 *studies*. Cambridge University Press, Cambridge, UK.
- 690 Cassady, J. C. (2001). Self-reported GPA and SAT: A methodological note.
691 *Practical Assessment, Research & Evaluation*, 7(12).
- 692 Cattell, R. B. (1943). The measurement of adult intelligence. *Psychological*
693 *Bulletin*, 40(3):153–193.

- 694 Clark, S. J. and Houle, B. (2012). Evaluation of Heckman selection model
695 method for correcting estimates of HIV prevalence from sample surveys.
696 *Center for Statistics and the Social Sciences*, Working Paper no. 120:1–18.
- 697 Cole, J. S. and Gonyea, R. M. (2009). Accuracy of self-reported SAT and
698 ACT test scores: Implications for research. *Research in Higher Education*,
699 51(4):305–319.
- 700 College Board (2012). *2012 college-bound seniors total group profile report*.
701 New York: The College Board. Retrieved September 13, 2013, from
702 [http://media.collegeboard.com/digitalServices/pdf/research/TotalGroup-](http://media.collegeboard.com/digitalServices/pdf/research/TotalGroup-2012.pdf)
703 [2012.pdf](http://media.collegeboard.com/digitalServices/pdf/research/TotalGroup-2012.pdf).
- 704 Cuddeback, G., Wilson, E., Orme, J. G., and Combs-Orme, T. (2004). De-
705 tecting and statistically correcting sample selection bias. *Journal of Social*
706 *Service Research*, 30(3):19–33.
- 707 Dennis, I., Handley, S., Bradon, P., Evans, J., and Newstead, S. (2002). Ap-
708 proaches to modeling item-generative tests. In Irvine, S. H. and Kyllonen,
709 P. C., editors, *Item generation for test development*, pages 53–71. Lawrence
710 Erlbaum Associates, Mahwah, New Jersey.
- 711 Educational Testing Service (2010). Table of GRE scores by intended grad-
712 uate major field.
- 713 Ekstrom, R. B., French, J. W., Harman, H. H., and Dermen, D. (1976).
714 *Manual for kit of factor-referenced cognitive tests*. Educational Testing
715 Service, Princeton, NJ.

- 716 Eliot, J. and Smith, I. M. (1983). *An International Directory of Spatial Tests*.
717 NFER-NELSON Publishing Company Ltd., Great Britain.
- 718 Embretson, S. E. (1996). The new rules of measurement. *Psychological*
719 *Assessment*, 8(4):341–349.
- 720 Field, T. G. (2012). *Standardized tests: Recouping development*
721 *costs and preserving integrity*. Retrieved September 13, 2013, from
722 http://papers.ssrn.com/sol3/papers.cfm?abstract_id=1989584.
- 723 Frey, M. C. and Detterman, D. K. (2004). Scholastic assessment or g? The
724 relationship between the scholastic assessment test and general cognitive
725 ability. *Psychological Science*, 15(6):373–378.
- 726 Frucot, V. and Cook, G. (1994). Further research on the accuracy of stu-
727 dents’ self-reported grade point averages, SAT scores, and course grades.
728 *Perceptual and Motor Skills*, 79(2):743–746.
- 729 Gambardella, A. and Hall, B. H. (2006). Proprietary versus public domain
730 licensing of software and research products. *Research Policy*, 35(6):875–
731 892.
- 732 Goldberg, L. R. (1999). A broad-bandwidth, public-domain, personality in-
733 ventory measuring the lower-level facets of several Five-Factor Models. In
734 Mervielde, I., Deary, I., De Fruyt, F., and Ostendorf, F., editors, *Person-*
735 *ality Psychology in Europe*, pages 1–7. Tilburg University Press, Tilburg,
736 The Netherlands.
- 737 Goldberg, L. R. (2012). *International Personality Item Pool: A scien-*
738 *tific collaboratory for the development of advanced measures of personality*

- 739 *traits and other individual differences*. Retrieved November 16, 2012, from
740 <http://ipip.ori.org/>.
- 741 Goldberg, L. R., Johnson, J. A., Eber, H. W., Hogan, R., Ashton, M. C.,
742 Cloninger, C. R., and Gough, H. G. (2006). The international personality
743 item pool and the future of public-domain personality measures. *Journal*
744 *of Research in Personality*, 40(1):84–96.
- 745 Goldstein, G. and Beers, S. R., editors (2004). *Comprehensive Handbook of*
746 *Psychological Assessment, Volume 1: Intellectual and Neuropsychological*
747 *Assessment*. John Wiley & Sons, Inc., Hoboken, NJ.
- 748 Greene, W. H. (2008). *Econometric Analysis*. Pearson Prentice Hall, Upper
749 Saddle River, NJ, 6th edition edition.
- 750 Heckman, J. J. (1976). The common structure of statistical models of trun-
751 cation, sample selection and limited dependent variables and a simple es-
752 timator for such models. In Berg, S. V., editor, *Annals of Economic and*
753 *Social Measurement, Volume 5, number 4*, pages 475–492. NBER, Cam-
754 bridge, MA.
- 755 Heckman, J. J. (1979). Sample Selection Bias as a Specification Error. *Econo-*
756 *metrica*, 47(1):153–161.
- 757 Hu, L. and Bentler, P. M. (1999). Cutoff criteria for fit indexes in covariance
758 structure analysis: Conventional criteria versus new alternatives. *Struc-*
759 *tural Equation Modeling*, 6(1):1–55.
- 760 Kaufmann, P. (2009). Protecting raw data and psychological tests from

- 761 wrongful disclosure: A primer on the law and other persuasive strategies.
762 *The Clinical Neuropsychologist*, 23(7):1130–1159.
- 763 Kenny, D. A. (2012). *Measuring model fit*. Retrieved November 7, 2012, from
764 <http://www.davidakenny.net/cm/fit.htm>.
- 765 Koenig, K. A., Frey, M. C., and Detterman, D. K. (2008). ACT and general
766 cognitive ability. *Intelligence*, 36(2):153–160.
- 767 Kuncel, N. R., Crede, M., and Thomas, L. L. (2005). The validity of self-
768 reported grade point averages, class ranks, and test scores: A meta-analysis
769 and review of the literature. *Review of Educational Research*, 75(1):63–82.
- 770 Liao, H.-Y., Armstrong, P. I., and Rounds, J. (2008). Development and initial
771 validation of public domain Basic Interest Markers. *Journal of Vocational*
772 *Behavior*, 73(1):159–183.
- 773 Lord, F. M. (1955). Sampling fluctuations resulting from the sampling of
774 test items. *Psychometrika*, 20(1):1–22.
- 775 Mayer, R. E., Stull, A. T., Campbell, J., Almeroth, K., Bimber, B., Chun, D.,
776 and Knight, A. (2006). Overestimation bias in self-reported SAT scores.
777 *Educational Psychology Review*, 19(4):443–454.
- 778 McCaffrey, R. J. and Lynch, J. K. (2009). Test security in the new millen-
779 nium: Is this really psychology's problem? *Emerging Fields*, 21(2):27.
- 780 Murphy, L. L., Geisinger, K. F., Carlson, J. F., and Spies, R. A. (2011).
781 *Tests in Print VIII*. An Index to Tests, Test Reviews, and the Literature on
782 Specific Tests. Buros Institute of Mental Measurements, Lincoln, Nebraska.

- 783 Pearson, K. (1903). Mathematical contributions to the theory of evolution.
784 XI. On the influence of natural selection on the variability and correla-
785 tion of organs. *Philosophical Transactions of the Royal Society of London.*
786 *Series A, Containing Papers of a Mathematical or Physical Character,*
787 200:1–66.
- 788 R Core Team (2013). *R: A Language and Environment for Statistical Com-*
789 *puting.* R Foundation for Statistical Computing, Vienna, Austria. ISBN
790 3-900051-07-0.
- 791 Revelle, W. (2013). *psych: Procedures for psychological, psychometric, and*
792 *personality research.* Northwestern University, Evanston, Illinois. R pack-
793 age version 1.3.9.13.
- 794 Revelle, W. and Brown, A. (2013). Standard errors for SAPA correlations.
795 In *Society for Multivariate Experimental Psychology*, pages 1–1, St. Peters-
796 burg, FL.
- 797 Revelle, W. and Condon, D. M. (2012). *Estimating ability for two*
798 *samples.* Retrieved November 1, 2013, from [http://www.personality-](http://www.personality-project.org/revelle/publications/EstimatingAbility.pdf)
799 [project.org/revelle/publications/EstimatingAbility.pdf](http://www.personality-project.org/revelle/publications/EstimatingAbility.pdf).
- 800 Revelle, W., Wilt, J., and Rosenthal, A. (2010). Individual differences in
801 cognition: New methods for examining the personality-cognition link. In
802 Gruszka, A., Matthews, G., and Szymura, B., editors, *Handbook of Individ-*
803 *ual Differences in Cognition: Attention, Memory and Executive Control,*
804 chapter 2, pages 27–49. Springer, New York.

- 805 Revelle, W. and Zinbarg, R. E. (2009). Coefficients alpha, beta, omega, and
806 the glb: Comments on Sijtsma. *Psychometrika*, 74(1):145–154.
- 807 Sackett, P. R. and Yang, H. (2000). Correction for range restriction: An
808 expanded typology. *Journal of Applied Psychology*, 85(1):112–118.
- 809 Shipley, W. C., Gruber, C., Martin, T., and Klein, A. M. (2010). *Shipley*
810 *Institute of Living Scale, 2nd edition*. Western Psychological Services, Los
811 Angeles, CA.
- 812 Shipley, W. C., Gruber, C. P., Martin, T. A., and Klein, A. M. (2009).
813 *Shipley-2*. Western Psychological Services, Los Angeles, CA.
- 814 Thorndike, R. L. (1949). *Personnel selection: Test and measurement tech-*
815 *niques*. John Wiley & Sons Inc, London.
- 816 Toomet, O. and Henningsen, A. (2008). Sample selection models in R: Pack-
817 age sampleSelection. *Journal of statistical software*, 27(7):1–23.
- 818 Tucker, L. R. and Lewis, C. (1973). A reliability coefficient for maximum
819 likelihood factor analysis. *Psychometrika*, 38(1):1–10.
- 820 Uebersax, J. S. (2000). *Estimating a latent trait model by factor anal-*
821 *ysis of tetrachoric correlations*. Retrieved September 13, 2013, from
822 <http://www.john-uebersax.com/stat/irt.htm>.
- 823 Wilt, J., Condon, D. M., and Revelle, W. (2011). Telemetrics and online
824 data collection: Collecting data at a distance. In Laursen, B., Little, T. D.,
825 and Card, N. A., editors, *Handbook of Developmental Research Methods*,
826 chapter 10, pages 163–180. Guilford Press.

827 Zinbarg, R. E., Revelle, W., Yovel, I., and Li, W. (2005). Cronbach's α ,
828 Revelle's β , and McDonald's ω_h : Their relations with each other and two
829 alternative conceptualizations of reliability. *Psychometrika*, 70(1):123–133.

830 Zinbarg, R. E., Yovel, I., Revelle, W., and McDonald, R. P. (2006). Es-
831 timating generalizability to a latent variable common to all of a scale's
832 indicators: A comparison of estimators for omega hierarchical. *Applied*
833 *Psychological Measurement*, 30(2):121–144.

Parallel Analysis Scree Plots

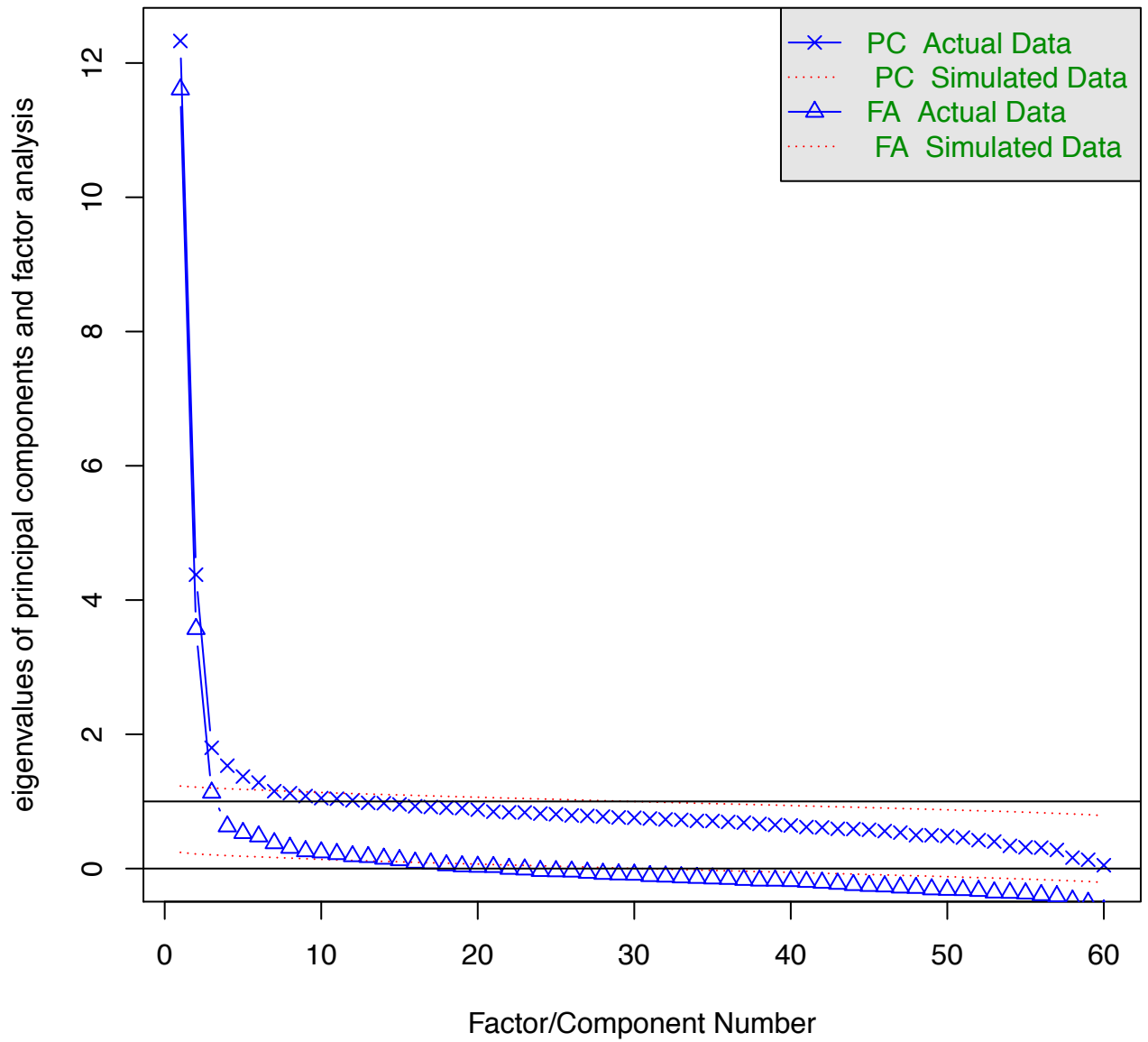


Figure 1: Scree plots based on all 60 ICAR items

Figure 2: Omega hierarchical for the *ICAR Sample Test*

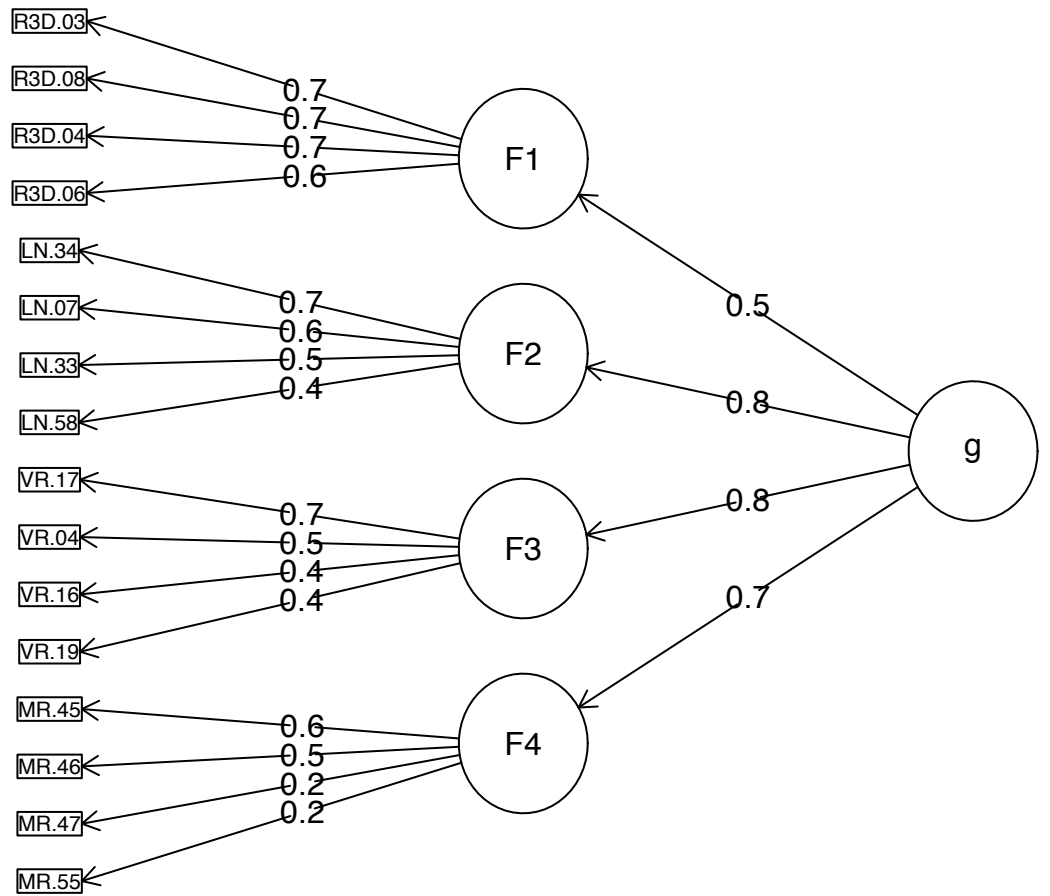
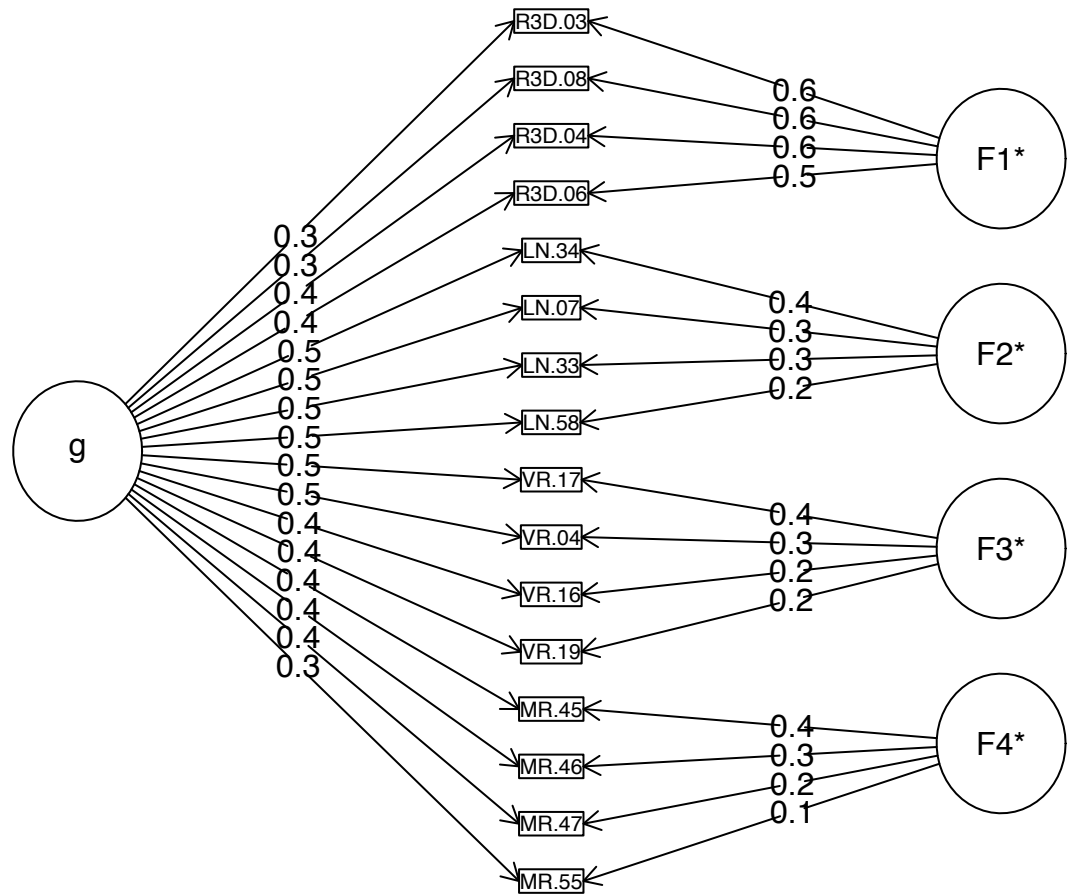


Figure 3: Omega with Schmid-Leiman transformation for the *ICAR Sample Test*



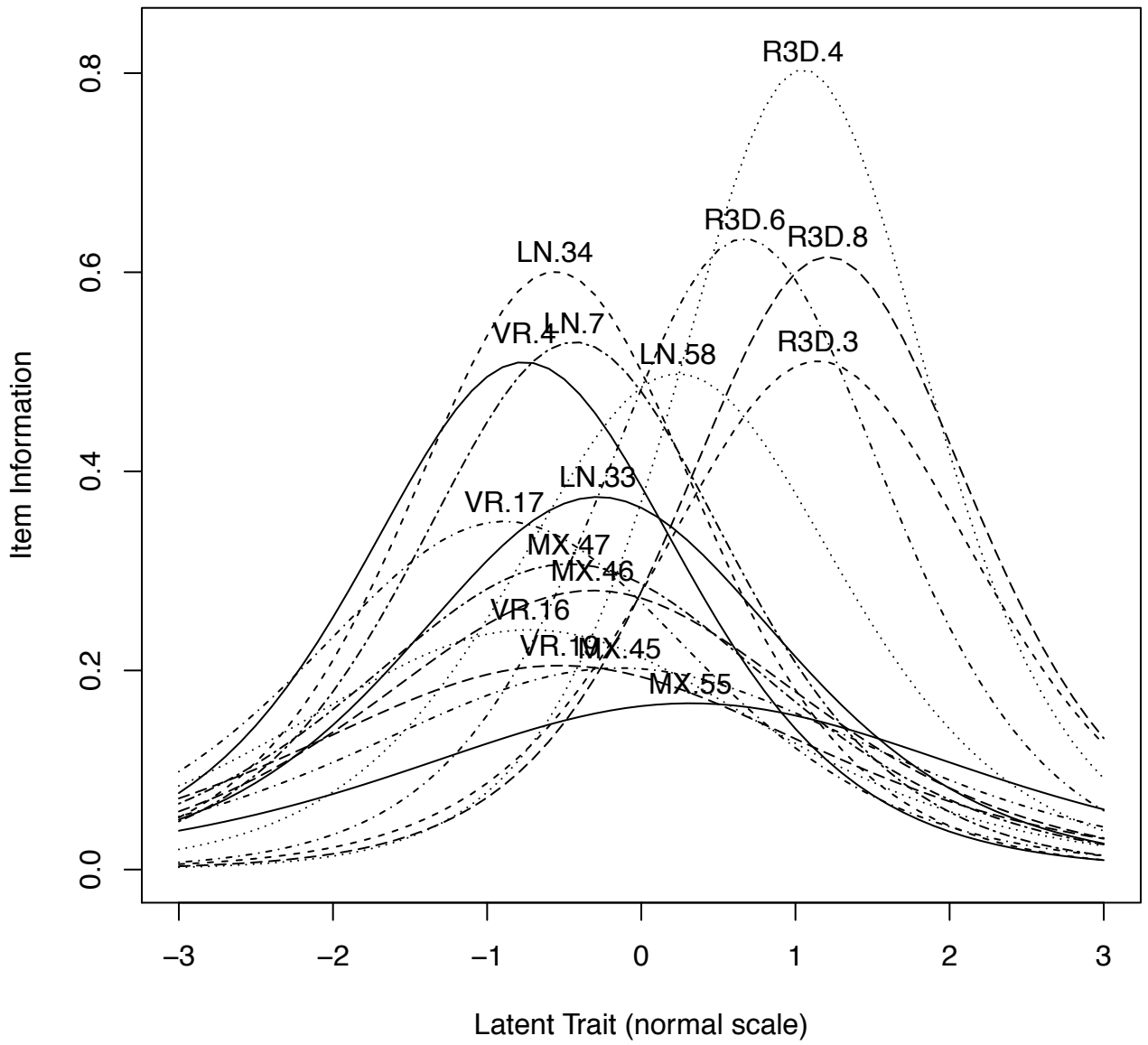


Figure 4: Item Information Functions for the 16 item *ICAR Sample Test*

ICAR Sample Test Test Information Function

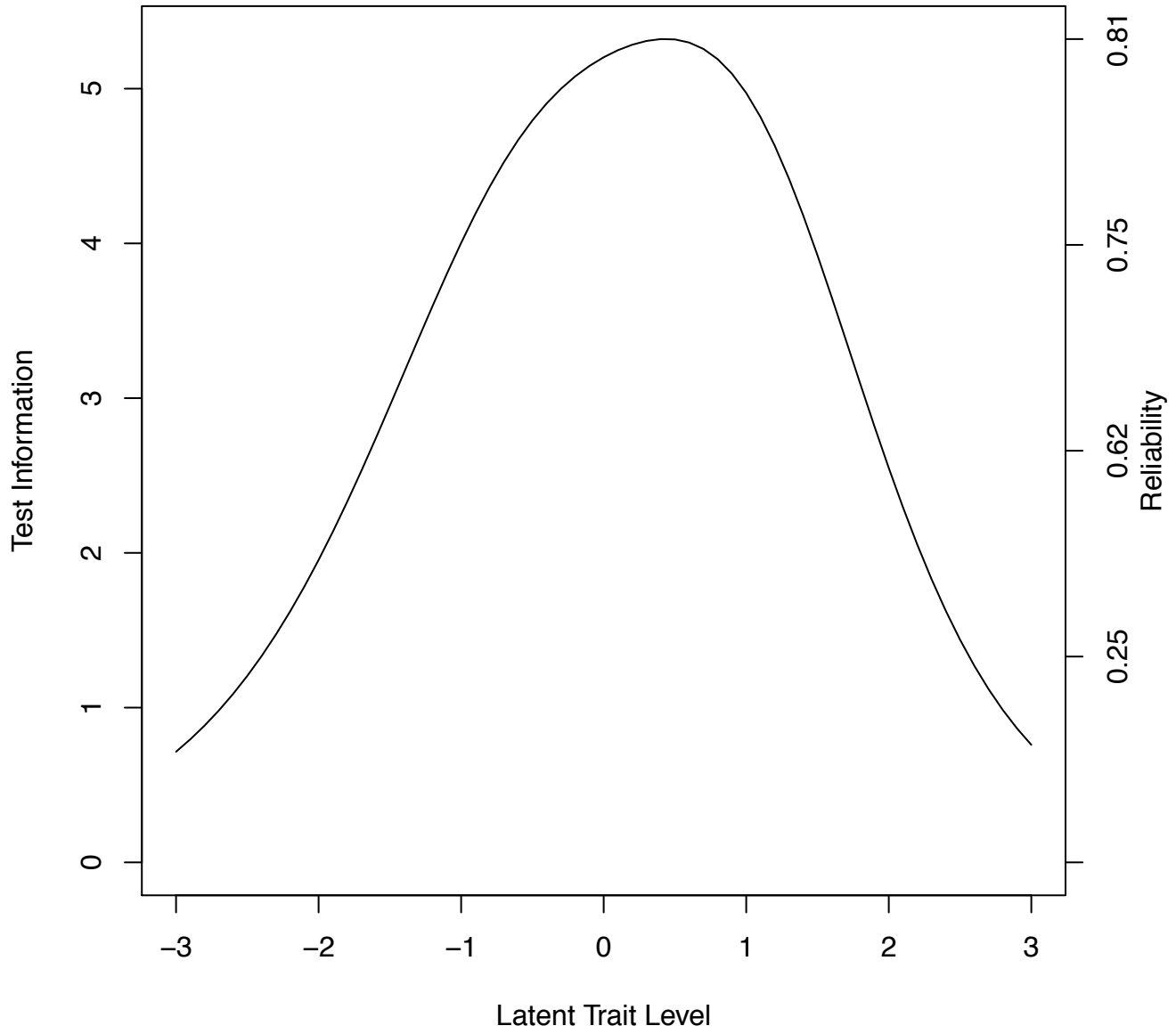


Figure 5: Test Information Function for the 16 item *ICAR Sample Test*

Table 1: Study 1 Participants by Educational Attainment

Educational attainment	% of total	Mean age	Median age
Less than 12 years	14.5%	17.3	17
High school graduate	6.2%	23.7	18
Currently in college/university	51.4%	24.2	21
Some college/university, but did not graduate	5.0%	33.2	30
College/university degree	11.7%	33.2	30
Currently in graduate or professional school	4.4%	30.0	27
Graduate or professional school degree	6.9%	38.6	36

Table 2: Descriptive statistics for the ICAR items administered in Study 1

Item	<i>n</i>	<i>mean</i>	<i>sd</i>	Item	<i>n</i>	<i>mean</i>	<i>sd</i>
LN.01	31,239	0.79	0.41	R3D.11	7,165	0.09	0.29
LN.03	31,173	0.59	0.49	R3D.12	7,168	0.13	0.34
LN.05	31,486	0.75	0.43	R3D.13	7,291	0.10	0.30
LN.06	34,097	0.46	0.50	R3D.14	7,185	0.14	0.35
<i>LN.07</i>	<i>36,346</i>	<i>0.62</i>	<i>0.49</i>	R3D.15	7,115	0.22	0.42
<i>LN.33</i>	<i>39,384</i>	<i>0.59</i>	<i>0.49</i>	R3D.16	7,241	0.30	0.46
<i>LN.34</i>	<i>36,655</i>	<i>0.62</i>	<i>0.48</i>	R3D.17	7,085	0.15	0.36
LN.35	34,372	0.47	0.50	R3D.18	6,988	0.13	0.34
<i>LN.58</i>	<i>39,047</i>	<i>0.42</i>	<i>0.49</i>	R3D.19	7,103	0.16	0.37
MR.43	29,812	0.77	0.42	R3D.20	7,203	0.39	0.49
MR.44	17,389	0.66	0.47	R3D.21	7,133	0.08	0.28
<i>MR.45</i>	<i>24,689</i>	<i>0.52</i>	<i>0.50</i>	R3D.22	7,369	0.30	0.46
<i>MR.46</i>	<i>34,952</i>	<i>0.60</i>	<i>0.49</i>	R3D.23	7,210	0.19	0.39
<i>MR.47</i>	<i>34,467</i>	<i>0.62</i>	<i>0.48</i>	R3D.24	7,000	0.19	0.39
MR.48	17,450	0.53	0.50	<i>VR.04</i>	<i>29,975</i>	<i>0.67</i>	<i>0.47</i>
MR.50	19,155	0.28	0.45	VR.09	25,402	0.70	0.46
MR.53	29,548	0.61	0.49	VR.11	26,644	0.86	0.35
MR.54	19,246	0.39	0.49	VR.13	24,147	0.24	0.43
<i>MR.55</i>	<i>24,430</i>	<i>0.36</i>	<i>0.48</i>	VR.14	26,100	0.74	0.44
MR.56	19,380	0.40	0.49	<i>VR.16</i>	<i>31,727</i>	<i>0.69</i>	<i>0.46</i>
R3D.01	7,537	0.08	0.28	<i>VR.17</i>	<i>31,552</i>	<i>0.73</i>	<i>0.44</i>
R3D.02	7,473	0.16	0.37	VR.18	26,474	0.96	0.20
<i>R3D.03</i>	<i>12,701</i>	<i>0.17</i>	<i>0.37</i>	<i>VR.19</i>	<i>30,556</i>	<i>0.61</i>	<i>0.49</i>
<i>R3D.04</i>	<i>12,959</i>	<i>0.21</i>	<i>0.41</i>	VR.23	24,928	0.27	0.44
R3D.05	7,526	0.24	0.43	VR.26	13,108	0.38	0.49
<i>R3D.06</i>	<i>12,894</i>	<i>0.29</i>	<i>0.46</i>	VR.31	26,272	0.90	0.30
R3D.07	7,745	0.12	0.33	VR.32	25,419	0.55	0.50
<i>R3D.08</i>	<i>12,973</i>	<i>0.17</i>	<i>0.37</i>	VR.36	25,076	0.40	0.49
R3D.09	7,244	0.28	0.45	VR.39	26,433	0.91	0.28
R3D.10	7,350	0.14	0.35	VR.42	25,108	0.66	0.47

Note: “LN” denotes Letter and Number Series, “MR” is Matrix Reasoning, “R3D” is Three-Dimensional Rotation, and “VR” is Verbal Reasoning. Italicized items denote those included in the 16-Item *ICAR Sample Test*.

Table 3: Alpha and omega for the ICAR item types

	α	ω_h	ω_t	items
ICAR60	0.93	0.61	0.94	60
LN items	0.77	0.66	0.80	9
MR items	0.68	0.58	0.71	11
R3D items	0.93	0.78	0.94	24
VR items	0.76	0.64	0.77	16
ICAR16	0.81	0.66	0.83	16

Note: ω_h = omega hierarchical, ω_t = omega total. Values are based on composites of Pearson correlations between items.

Table 4: Four-factor item loadings for the *ICAR Sample Test*

Item	Factor 1	Factor 2	Factor 3	Factor 4
R3D.03	0.69	-0.02	-0.04	0.01
R3D.08	0.67	-0.04	-0.01	0.02
R3D.04	0.66	0.03	0.01	0.00
R3D.06	0.59	0.06	0.07	-0.02
LN.34	-0.01	0.68	-0.01	-0.02
LN.07	-0.03	0.60	-0.01	0.05
LN.33	0.04	0.52	0.01	0.00
LN.58	0.08	0.43	0.07	0.01
VR.17	-0.04	0.00	0.65	-0.02
VR.04	0.06	-0.01	0.51	0.05
VR.16	0.02	0.05	0.41	0.00
VR.19	0.03	0.02	0.38	0.06
MR.45	-0.02	-0.01	0.01	0.56
MR.46	0.02	0.02	0.01	0.50
MR.47	0.05	0.18	0.10	0.24
MR.55	0.14	0.09	-0.04	0.21

Table 5: Correlations between factors for the *ICAR Sample Test*

	R3D Factor	LN Factor	VR Factor	MR Factor
R3D Factor	1.00			
LN Factor	0.44	1.00		
VR Factor	0.70	0.45	1.00	
MR Factor	0.63	0.41	0.59	1.00

Note: R3D = Three-Dimensional Rotation, LN = Letter and Number Series, VR = Verbal Reasoning, MR = Matrix Reasoning

Table 6: Item and test information for the 16 item *ICAR Sample Test*

Item	Latent Trait Level (normal scale)						
	-3	-2	-1	0	1	2	3
VR.04	0.07	0.23	0.49	0.42	0.16	0.04	0.01
VR.16	0.08	0.17	0.25	0.23	0.13	0.06	0.02
VR.17	0.09	0.27	0.46	0.34	0.13	0.04	0.01
VR.19	0.07	0.14	0.24	0.25	0.16	0.07	0.03
LN.07	0.06	0.18	0.38	0.39	0.19	0.06	0.02
LN.33	0.05	0.15	0.32	0.37	0.21	0.08	0.02
LN.34	0.05	0.20	0.46	0.45	0.19	0.05	0.01
LN.58	0.03	0.09	0.26	0.43	0.32	0.13	0.04
MR.45	0.05	0.11	0.17	0.20	0.16	0.09	0.04
MR.46	0.06	0.13	0.22	0.24	0.17	0.08	0.04
MR.47	0.06	0.16	0.31	0.32	0.18	0.07	0.02
MR.55	0.04	0.07	0.11	0.14	0.13	0.10	0.06
R3D.03	0.00	0.01	0.06	0.27	0.64	0.47	0.14
R3D.04	0.00	0.01	0.07	0.35	0.83	0.45	0.10
R3D.06	0.00	0.03	0.14	0.53	0.73	0.26	0.05
R3D.08	0.00	0.01	0.06	0.26	0.64	0.48	0.14
TIF	0.72	1.95	4.00	5.20	4.97	2.55	0.76
SEM	1.18	0.72	0.50	0.44	0.45	0.63	1.15
Reliability	NA	0.49	0.75	0.81	0.80	0.61	NA

Table 7: Self-reported achievement test scores and national norms

	Study 2			published	
	self-reported			norms	
	n	mean	s.d.	mean	s.d.
SAT - Critical Reading	7,404	609	120	496	114
SAT - Math	7,453	611	121	514	117
ACT	12,254	25.4	5.0	21.1	5.2

Note: SAT norms are from the 2012 *Total Group Profile Report*. ACT norms are from the 2011 *ACT Profile Report*.

Table 8: Correlations between self-reported achievement test scores and ICAR composite scales

	ICAR composite scale scores									
	SAT-CR	SAT-M	SAT-CR+M	ACT	ICAR60	LN	MR	R3D	VR	ICAR16
SAT-CR ¹	<i>0.86</i>	0.83		0.69	0.52	0.41	0.37	0.39	0.68	0.52
SAT-M ²	0.72	<i>0.88</i>		0.66	0.60	0.50	0.47	0.49	0.67	0.59
SAT-CR+M ³			<i>0.89</i>	0.71	0.59	0.48	0.44	0.47	0.72	0.59
ACT ⁴	0.62	0.60	0.65	<i>0.95</i>	0.52	0.39	0.35	0.44	0.61	0.52
ICAR60 ⁵	0.46	0.54	0.54	0.49	<i>0.93</i>					
LN ⁵	0.33	0.41	0.40	0.33		<i>0.77</i>	0.84	0.59	0.90	
MR ⁵	0.28	0.36	0.34	0.28		0.61	<i>0.68</i>	0.67	0.81	
R3D ⁵	0.35	0.44	0.43	0.41		0.50	0.53	<i>0.93</i>	0.58	
VR ⁵	0.55	0.55	0.59	0.52		0.69	0.58	0.49	<i>0.76</i>	
ICAR16 ⁵	0.43	0.50	0.50	0.46						<i>0.81</i>

Note: Uncorrected correlations below the diagonal, correlations corrected for reliability above the diagonal. Reliability values shown on the diagonal.

¹ $n = 7,404$

² $n = 7,453$

³ $n = 7,348$

⁴ $n = 12,254$

⁵ Composite scales formed based on item correlations across the full sample ($n = 34,229$).

Table 9: Correlations between self-reported achievement test scores and IRT-based ICAR scores

	ICAR IRT-based scores								
	SAT-CR	SAT-M	SAT-CR+M	ACT	ICAR60	LN	MR	R3D	VR
SAT-CR ¹					0.44	0.37	0.35	0.37	0.44
SAT-M ²	0.72				0.44	0.33	0.29	0.35	0.39
SAT-CR+M ³	0.93	0.93			0.47	0.37	0.33	0.38	0.45
ACT ⁴	0.62	0.60	0.65		0.44	0.35	0.32	0.38	0.43
ICAR60 ⁵	0.36	0.42	0.42	0.39					
LN ⁵	0.24	0.28	0.28	0.24					
MR ⁵	0.18	0.22	0.21	0.18		0.30			
R3D ⁵	0.25	0.32	0.30	0.28		0.26	0.23		
VR ⁵	0.35	0.36	0.38	0.36		0.36	0.26	0.22	

Note: IRT scores for ICAR measures based on 2 to 4 responses per participant for each item type (LN, MR, R3D, VR) and 12 to 16 responses for ICAR60. Uncorrected correlations are below the diagonal, correlations corrected for incidental selection are above the diagonal.

¹ $n = 7,404$

² $n = 7,453$

³ $n = 7,348$

⁴ $n = 12,254$

⁵ $n = 34,229$

Table 10: Correlations between mean SAT norms, mean SAT scores in Study 2 and mean IRT-based ICAR scores when ranked by university major

	College Board Norms			Study 2 Self-Reported			Study 2 IRT-based			
	SAT-CR	SAT-M	SAT-CR+M	SAT-CR	SAT-M	SAT-CR+M	ICAR60	LN	MR	R3D
SAT-M norms	0.66									
SAT-CR+M norms	0.91	0.91								
SAT-CR study 2	0.79	0.61	0.77							
SAT-M study 2	0.56	0.80	0.74	0.81						
SAT-CR+M study 2	0.71	0.74	0.80	0.95	0.95					
ICAR60 study 2	0.53	0.84	0.75	0.60	0.77	0.72				
LN study 2	0.41	0.80	0.66	0.49	0.76	0.66	0.96			
MR study 2	0.22	0.66	0.48	0.23	0.52	0.39	0.83	0.78		
R3D study 2	0.42	0.80	0.67	0.50	0.71	0.64	0.94	0.92	0.82	
VR study 2	0.69	0.79	0.81	0.76	0.80	0.82	0.91	0.82	0.64	0.76

Note: $n = 29$.

Table 11: Correlations between mean GRE norms and mean IRT-based ICAR scores when ranked by university major

	ETS Norms			Study 2 IRT-based			
	GREV	GREQ	GREVQ	ICAR60	LN	MR	R3D
GREQ norms	0.23						
GREVQ norms	0.63	0.90					
ICAR60 study 2	0.54	0.78	0.86				
LN study 2	0.41	0.72	0.76	0.93			
MR study 2	0.42	0.71	0.75	0.86	0.81		
R3D study 2	0.44	0.80	0.83	0.92	0.86	0.75	
VR study 2	0.67	0.63	0.80	0.92	0.80	0.79	0.77

Note: $n = 91$.

Table 12: Standard deviations of scores for the unrestricted samples and Study 3

Sample	Block Patterns	Abstraction	<i>Shipley-2</i>			<i>ICAR</i>
			Vocabulary	Composite A	Composite B	<i>Sample Test</i>
Unrestricted	15.0	15.0	15.0	15.0	15.0	1.86
Study 3	11.1	9.8	6.8	6.8	8.9	1.48

Note: Unrestricted standard deviations based on the published norms for the *Shipley-2* and the Study 1 sample for the *ICAR Sample Test*.

Table 13: Correlations between the *ICAR Sample Test* and the *Shipley-2*

ICAR16	Block Patterns ¹	Abstraction ²	Vocabulary ³	Composite A ²	Composite B ¹
Uncorrected	0.40	0.44	0.15	0.41	0.41
Range corrected	0.64	0.69	0.59	0.68	0.68
Range & reliability corrected				0.82	0.81

¹ $n = 68$

² $n = 69$

³ $n = 137$