

A peer-reviewed version of this preprint was published in PeerJ on 22 December 2015.

[View the peer-reviewed version](https://doi.org/10.7717/peerj.1460) (peerj.com/articles/1460), which is the preferred citable publication unless you specifically need to cite this preprint.

Piper BJ, Mueller ST, Geerken AR, Dixon KL, Kroliczak G, Olsen RHJ, Miller JK. 2015. Reliability and validity of neurobehavioral function on the Psychology Experimental Building Language test battery in young adults. PeerJ 3:e1460 <https://doi.org/10.7717/peerj.1460>

Reliability and validity of neurobehavioral function on the Psychology Experimental Building Language test battery in young-adults

Brian J Piper, Shane T Mueller, Alexander R Geergen, Kyle L Dixon, Gregory Kroliczak, Reid HJ Olsen, Jeremy K Miller

Background. The Psychology Experiment Building Language (PEBL) software consists of over one-hundred computerized tests based on classic cognitive neuropsychology and behavioral neurology measures. Although the PEBL tests are becoming more widely utilized, there is currently very limited information about the psychometric properties of these measures. **Methods.** Study I examined inter-relationships among ten PEBL tests including indices of motor-function (Pursuit Rotor and Dexterity), attention (Test of Attentional Vigilance and Time-Wall), working memory (Digit Span Forward), and executive-function (PEBL Trail Making Test, Berg/Wisconsin Card Sorting Test, Iowa Gambling Test, and Mental Rotation) in a normative sample (N = 189, ages 18-22). Study II evaluated test-retest reliability with a two-week interval between administrations in a separate sample (N = 79, ages 18-22). **Results.** Moderate intra-test, but low inter-test, correlations were observed and ceiling/floor effects were uncommon. Sex differences were identified on the Pursuit Rotor (Cohen's $d = 0.89$) and Mental Rotation ($d = 0.31$) tests. The correlation between the test and retest was high for tests of motor learning (Pursuit Rotor time on target $r = .86$) and attention (Test of Attentional Vigilance response time $r = .79$), intermediate for memory (digit span $r = .63$) but lower for the executive function indices (Wisconsin/Berg Card Sorting Test perseverative errors = .45, Tower of London moves = .15). Significant practice effects were identified on several indices of executive function. **Conclusions.** These results are broadly supportive of the reliability and validity of individual PEBL tests in this sample. These findings indicate that the freely downloadable, open-source, PEBL battery <http://pebl.sourceforge.net> is a versatile research tool to study individual differences in neurocognitive performance.

1 **Reliability and Validity of Neurobehavioral Function on the Psychology Experimental**
2 **Building Language Test Battery in Young-Adults**

3
4 Brian J. Piper^{abc}, Shane T. Mueller^d, Alexander R. Geergen^a, Kyle L. Dixon^{ae}, Gregory
5 Króliczak^f, Reid H. J. Olsen^c, & Jeremy K. Miller^a

6 ^aDepartment of Psychology, Willamette University, Salem, OR 97301 USA

7 ^bDepartment of Psychology, Bowdoin College, Brunswick, ME 04011 USA

8 ^cDepartment of Behavioral Neuroscience, Oregon Health and Science University, Portland, OR
9 97239 USA

10 ^dDepartment of Cognitive & Learning Sciences, Michigan Technological University, Houghton,
11 MI 49931 USA

12 ^eDepartment of Psychology, University of New Mexico, Albuquerque, NM 87131 USA

13 ^fInstitute of Psychology, Adam Mickiewicz University of Poznań, Poznań, Poland

14
15
16
17
18
19
20
21
22
23
24
25
26
27 Corresponding Author:
28 Brian J. Piper, Ph.D.
29 255 Main Street
30 Department of Psychology & Neuroscience Program
31 Bowdoin College
32 Brunswick, ME 04001
33 psy391@gmail.com; bpiper@bowdoin.edu

Abstract

34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58

Background. The Psychology Experiment Building Language (PEBL) software consists of over one-hundred computerized tests based on classic cognitive neuropsychology and behavioral neurology measures. Although the PEBL tests are becoming more widely utilized, there is currently very limited information about the psychometric properties of these measures.

Methods. Study I examined inter-relationships among ten PEBL tests including indices of motor-function (Pursuit Rotor and Dexterity), attention (Test of Attentional Vigilance and Time-Wall), working memory (Digit Span Forward), and executive-function (PEBL Trail Making Test, Berg/Wisconsin Card Sorting Test, Iowa Gambling Test, and Mental Rotation) in a normative sample (N = 189, ages 18-22). Study II evaluated test-retest reliability with a two-week interval between administrations in a separate sample (N = 79, ages 18-22).

Results. Moderate intra-test, but low inter-test, correlations were observed and ceiling/floor effects were uncommon. Sex differences were identified on the Pursuit Rotor (Cohen's $d = 0.89$) and Mental Rotation ($d = 0.31$) tests. The correlation between the test and retest was high for tests of motor learning (Pursuit Rotor time on target $r = .86$) and attention (Test of Attentional Vigilance response time $r = .79$), intermediate for memory (digit span $r = .63$) but lower for the executive function indices (Wisconsin/Berg Card Sorting Test perseverative errors = .45, Tower of London moves = .15). Significant practice effects were identified on several indices of executive function.

Conclusions. These results are broadly supportive of the reliability and validity of individual PEBL tests in this sample. These findings indicate that the freely downloadable, open-source, PEBL battery <http://pebl.sourceforge.net> is a versatile research tool to study individual differences in neurocognitive performance.

59 INTRODUCTION

60 A large collection of classic tests from the behavioral neurology and cognitive
61 psychology fields have been computerized and made available (<http://pebl.sf.net>). This
62 Psychology Experiment Building Language (PEBL) (Mueller, 2010, 2014a, 2014b; Mueller &
63 Piper, 2014) has been downloaded over 168 thousand times with 73% of downloads by
64 institutions located outside of the United States, and used in scores of published manuscripts
65 (e.g. Barrett & Gonzalez-Lima, 2013; Danckert et al., 2012; Fox et al., 2013; Gonzalez-Giraldo
66 et al. 2014, 2015a, 2015b; Piper, 2010; Piper et al. 2012; Premkumar et al., 2013; Wardle et al.,
67 2012; Supplementary Table 1). The growth in PEBL use is likely due to three factors. First,
68 PEBL is free while other similar programs (Robbins et al., 1994) have costs that preclude use by
69 all but the largest laboratories and are beyond the capacities of the majority of investigators in
70 developing countries. Second, PEBL is open-source software and therefore the computational
71 operations are more transparent than may be found with proprietary measures. Third, the
72 distributors of some commercial tests restrict test availability to those who have completed
73 specific coursework whereas PEBL is available to anyone with an internet connection. This
74 investigation reports on the use of ten PEBL measures including convergent and divergent
75 validity (Study I) and test-retest reliability (Study II). A brief history of the more commonly
76 utilized of these tests is provided below.

77 *Digit Span*

78 The origins of digit span, an extremely simple test in which strings of numbers of
79 increasing length are presented and must be repeated back to the experimenter, are ambiguous
80 but procedures that are analogous to what are frequently employed today date back at least as far
81 as the pioneering developmental studies of Alfred Binet (Richardson, 2007). Although digit span

82 is frequently described as an index of working memory, the importance of attention for optimal
83 performance should not be underestimated (Lezak et al. 2012).

84 *Rotary Pursuit*

85 The rotary pursuit test measures motor performance by using a stylus to track a target that
86 moves clockwise at a fixed rate (Ammons, Alprin, & Ammons, 1955). Procedural learning
87 deficits using the rotary pursuit have been shown among patients with Huntington's (Schmidtke
88 et al., 2002). As a result of the wide-spread use of the rotary pursuit in experimental psychology
89 laboratories, a computerized version was developed. Unfortunately, this version could only
90 generate linear target paths due to technical limitations at that time (Willingham, Hollier, &
91 Joseph, 1995). The PEBL pursuit rotor is a more faithful version of the original rotary pursuit
92 distributed by Lafayette instruments. Importantly, prior computer experience does account for a
93 small portion of the variance in time on target (Piper, 2010).

94 *Wisconsin Card Sorting Test*

95 The Wisconsin Card Sorting Test is a classic neuropsychological measure of cognitive
96 flexibility and was originally developed at the University of Wisconsin following WWII by Esta
97 Berg and David Grant (Grant & Berg, 1948). In the original version, participants sorted physical
98 cards into piles and determined the underlying classification principle by trial and error. Once
99 consistent correct matching was achieved, the principle would be changed. The subsequent
100 development of a computerized version of this complex task made for both more efficient use of
101 the participants time and automated scoring (Lezak et al. 2012). Another key discovery was that
102 64 cards could be used instead of 128 (Axelrod, Woodard, & Henry, 1992; Fox et al. 2013).

103 *Trail Making Test*

104 The Trail Making Test is another of the oldest and most commonly employed
105 neurobehavioral measures (Lezak et al., 2012). The Trail Making Test is typically thought to
106 measure visual attention, mental flexibility, and executive functioning. The Trail Making Test
107 was contained in the Army General Classification Test, a precursor of the Armed Services
108 Vocational Aptitude Battery used by the United States military. The Trail Making Test involves
109 connecting dots arranged in a numbered sequence in ascending order (Part A) or numbers and
110 letters that alternate (Part B). Traditionally, performance on the Trail Making Test has been
111 timed with a stop-watch and the experimenter has to redirect the participant when they make an
112 error. Unlike the Halstead-Reitan Trail Making Test (Gaudino, Geisler, & Squires, 1995), the path
113 length is equal in Parts A and B of the PEBL Trail Making Test. Behavior on the Trail Making
114 Test is sensitive to wide variety of insults including alcoholism (Chanraud et al. 2009).

115 *Mental Rotation Test*

116 The mental rotation test has been an influential measure in cognitive psychology.
117 Participants must decide whether an image is rotated in space and there is a linear relationship
118 between the angle of rotation and decision time. Males exhibit better performance on spatial
119 ability tests with some evidence indicating that this robust sex difference (e.g. Yassen et al., 2015)
120 is detectable at very young ages (Linn & Petersen, 1985; Moore & Johnson, 2008).

121 *Tower of London*

122 The Tower of London requires planning and judgment to arrive at the most efficient
123 solution and move colored balls from their initial position to a new set of predetermined or goal
124 positions (Shallice, 1982). There are many variations on this “brain teaser” type task including
125 different levels of difficulty and construction (wood versus computerized) (Lezak et al. 2012).
126 An elevation in the number of moves to solve Tower of London type problems has been

127 documented among patients with brain damage and schizophrenia (Morris et al., 1995; Shallice,
128 1982).

129 *Iowa Gambling Task*

130 The Iowa Gambling Task was developed to model real world decision making in a
131 laboratory environment. Participants receive \$2,000 to start and must maximize their profit by
132 choosing cards from among four decks of which two typically result in a net gain (+\$250) and
133 two result in a net loss (-\$250). Although the Iowa Gambling Task has been employed with a
134 wide range of neuropsychiatric disorders, identification of a condition that consistently shows an
135 abnormality on this test has proved difficult with the possible exception of problem gamblers
136 (Buelow & Suhr, 2009; Power, Goodyear, & Crockford, 2012).

137 *Test of Variables of Attention*

138 The Test of Variables of Attention is an index of vigilance and impulsivity in which the
139 participant responds to a target but inhibits responses for non-target stimuli. Although continuous
140 performance tests were intended to discriminate children with, and without, Attention Deficit
141 Hyperactivity Disorder (Greenberg & Waldman, 1993), the Test of Variables of Attention and
142 other similar instruments may have proved even more valuable in measuring attention as a
143 general construct and more specifically in evaluating the efficacy of cognitive enhancing drugs
144 (Huang et al. 2007).

145 Another feature of the PEBL battery is that the key brain structures for these classic tasks
146 are reasonably well characterized based on both lesion studies and more recent neuroimaging
147 investigations (Figure 1). Importantly, as diffuse neural networks are responsible for complex
148 behaviors and the notion of a single neuroanatomical area underlying performance on a test risks
149 oversimplification, more comprehensive information can be found elsewhere (Demakis, 2004;

150 Gerton et al., 2004; Grafton et al., 1992; Hugdahl, Thomsen, & Ersland, 2006; Jacobson et al.,
151 2011; Kaneko et al., 2011; Rogalsky, et al., 2012; Schall et al., 2013; Specht et al., 2009; Tana et
152 al., 2010; Zacks, 2008). Briefly, completing the rotary pursuit with the dominant (right) hand
153 results in a pronounced increase in blood flow in the left primary motor cortex, right cerebellum,
154 the supplementary motor area, and the left putamen (Grafton et al. 1992). Tasks that require
155 sustained attention engage the anterior cingulate and the insula (Tana et al., 2010). Digit Span
156 activates the left prefrontal cortex when examined with near-infrared spectroscopy (Kaneko et al.
157 2011). Whole brain comparison of Digit Span backward, relative for forward, using Position
158 Emission Tomography (PET) revealed blood flow elevations in the dorsal lateral prefrontal
159 cortex, left intraparietal lobule, and in Broca's area (Gerton et al., 2004). The Mental Rotation
160 Test results in a robust activation in the right intraparietal sulcus as well as in the frontal and
161 inferotemporal cortex (Jacobson et al. 2011). Executive function measures like the Trail Making
162 Test, Iowa Gambling Test, Tower of London, and Wisconsin Card Sorting Test have been
163 adapted from their clinical neuropsychological roots to be appropriate in a neuroimaging
164 environment. Part B of the Trail Making Test, relative to Part A, produces Blood Oxygen Level
165 Dependent elevations in the inferior middle frontal gyri (Jacobson et al. 2011). The left middle
166 frontal gyrus and right cerebellar tonsils show Tower of London difficulty dependent activations
167 as determined by both functional magnetic resonance imaging and PET. The left ventral medial
168 prefrontal cortex is engaged during completion of the Iowa Gambling Task (Schall et al. 2013)
169 although lesion studies have produced conflicting evidence regarding the importance of this
170 structure (Shallice, 1982). The Wisconsin Card Sorting Test is a highly cognitively demanding
171 task which involves an extremely diffuse cortical network including the right middle frontal
172 gyrus as well as the left and right parietal lobule (Kaneko et al., 2011).

173 Previously, performance on three of the most prevalent executive function tests including
174 the Wisconsin (Berg) Card Sorting Test, Trail Making Test, and the Tower of London was
175 determined in a lifespan (age 5-87) sample. This investigation identified the anticipated “U-
176 shaped” association between age and performance on these PEBL tests (Piper et al. 2012). One
177 objective of the present report was to extend upon this foundation in a young-adult population by
178 further examining the utility of the three executive function indices as well as six other tests
179 including one (Dexterity) that is completely novel and another (Time-Wall) that is relatively
180 obscure. Each participant in Study I completed all ten measures so that score distributions and
181 the inter-test correlations could be evaluated. This information is necessary because PEBL
182 measures, particularly the indices of executive function, are becoming increasingly utilized. The
183 non-PEBL versions of several tests (Tower of London, Mental Rotation Test, Trail Making Test,
184 rotary pursuit, and digit span) are often conducted using non-computerized methodology (Lezak
185 et al., 2012) so it is currently unclear whether prior data on convergent and discriminant validity
186 will be applicable. Many young adults have extensive experience with computerized measures so
187 it is also crucial to determine whether any measures have ceiling effects.

188 With the exception of a single pilot study (Piper, 2012), there is currently no information
189 about the test-retest reliability of individual PEBL tests or the battery. This dearth of data is
190 unfortunate because the PEBL tests have already been employed in repeated measures designs
191 (Barrett & Gonzalez-Lima, 2013; Premkumar et al., 2013; Wardle et al, 2012) and additional
192 information would aid in the interpretation of those findings. The consistency of measurement is
193 captured by two complementary measures. The correlation between the test and the retest
194 measures the relative consistency, and the effect size quantifies the absolute consistency in
195 performance.

196 There is a vast literature on the reliability of non-PEBL tests (Calamia, Markon, &
197 Tranel, 2013; Lezak et al. 2012) and a few investigations with similar methodology or sample
198 characteristics similar to this report provide some context for the present endeavor. College-
199 students assessed on a computerized target tracking task showed a high correlation ($r = .75$)
200 across sessions separated by two weeks (Fillmore, 2003). Strong correlations ($r > .70$) were also
201 noted on several indices of the Test of Variables of Attention among children completing that
202 vigilance measure with a nine-day inter-test interval (Learck, Wallace, & Fitzgerald, 2004).
203 Veterans in their late-20s exhibited an intermediate ($r = .52$) consistency across three sessions
204 (one/week) of a computerized Digit Span forward (Woods et al., 2010). The percentage
205 selections of the disadvantageous decks showed a moderate correlation ($r \geq .57$) when the Iowa
206 Gambling Task was administered thrice on the same day (Lejuez et al., 2005) but limited
207 information is available at longer intervals (Buelow & Suhr, 2009). The magnitude of practice
208 effects appears to be task dependent with slight changes identified for the Digit Span forward
209 (Woods et al. 2010) and the Test of Variables of Attention (Learck et al. 2004) but pronounced
210 improvements for the Iowa Gambling Task (Bechara, Damasio, & Demasio, 2000; Lejuez et al.
211 2005). Executive function tasks that have a problem solving element may, once solved, have a
212 limited reliability (Lowe & Rabbit, 1998). For example, the correlation of the first with the
213 second 64-trials on the Berg Card Sorting Test was relatively low ($r = .31$) (Fox et al., 2013). In
214 fact, the Wisconsin Card Sorting Test has been referred to as a “one shot test” (Lezak et al.
215 2012).

216 Two secondary objective of this report are also noteworthy. First, these datasets provided
217 an opportunity to identify any sex differences on the PEBL battery. As a general rule, males and
218 females are more similar than dissimilar on most neurocognitive measures. However, as noted

219 previously, the Mental Rotation Test provides a clear exception to this pattern (Linn & Petersen,
220 1985; Moore & Johnson, 2008). A robust male advantage was observed among children
221 completing the PEBL Pursuit Rotor task (Piper, 2010) and similar sex differences have been
222 identified with the non-computerized version (Willingham et al. 1995) of this test. However, sex
223 differences were most pronounced only at older (81+) but not younger (21-80) ages on a
224 computerized task with many similarities to the PEBL Pursuit Rotor (Stirling et al., 2013).

225 A final objective was to evaluate the different card sorting rules on the Berg Card Sorting
226 Task. The PEBL version of the Wisconsin Card Sorting Task has been employed in over a dozen
227 reports (e.g. Danckert et al., 2012; Fox et al., 2013; Piper et al., 2012; Wardle et al., 2011) and
228 may be the most popular of the PEBL tests. Importantly, the Berg Card Sorting Test was
229 programmed based on the definitions of perseverative responses and perseverative errors
230 contained in Esta Berg's 1948 report (Berg, 1948). Alternatively, the Wisconsin Card Sorting
231 Task distributed by Psychological Assessment Resources employs the subsequent definitions of
232 Robert Heaton and colleagues (Heaton et al., 1993).

234

235

MATERIALS AND METHODS

236 *Participants*

237 Participants (Study I: N = 189, 60.3% Female, Age = 18.9 ± 1.0 ; Study II: N = 79, 73.0%

238 Female, Age = 19.1 ± 0.1) were college students receiving course credit. The test sequence in

239 Study I was as follows: written informed consent, Tapping, Pursuit Rotor, Time-Wall, Trail-

240 Making Task, Digit-Span Forward, Berg Card Sorting Test, Mental Rotation, Iowa Gambling

241 Task, Tower-of-London, Dexterity, and the Test of Attentional Vigilance. Due to hardware

242 technical difficulties, data from the tapping motor speed test were unavailable. Half of these

243 measures (Time-Wall, Trail-Making Test, Digit-Span, Mental Rotation and the Test of

244 Attentional Vigilance) contain programming modifications relative to the PEBL battery 0.6

245 defaults and may be found in the Supplemental Materials. All neurobehavioral assessments were

246 completed on one of eight desktop computers running Microsoft Windows. Each of these tests is

247 described further below and screen shots including instructions are in the Supplemental Figure 1.

248 The number of tests was slightly reduced to eight for Study II and the sequence was a written

249 informed consent followed by Pursuit Rotor, Trail-Making Test, Digit-Span, Test of Attentional

250 Vigilance, Tower-of-London, Iowa Gambling Task, and Time-Wall. The interval between the

251 test and retest was two weeks (mean = 14.4 ± 0.2 days, Min = 11, Max = 24). This inter-test

252 interval could be employed to examine the effects of a cognitively enhancing drug. All

253 procedures are consistent with the Declaration of Helinski and were approved by the Institutional

254 Review Board of Willamette University.

255 *PEBL Tests*

256 Pursuit Rotor measures motor-learning and requires the participant to use the computer
257 mouse to follow a moving target on four-fifteen second trials. The target follows a circular path
258 (8 rotations per minute) and the time on target and error, the difference in pixels between the
259 cursor and target, were recorded (González-Giraldo et al., 2015; Piper, 2010).

260 Time-Wall is an attention and decision making task that involves assessing the time at
261 which a target, moving vertically at a constant rate, will have traveled a fixed distance. The
262 primary dependent measure is Inaccuracy, defined as the absolute value of the difference
263 between the participant response time and the correct time divided by the correct time (Minimum
264 = 0.00). The correct time ranged from 2.0 to 9.2 seconds with feedback (“Too short” or “Too
265 long”) provided after each of the ten trials (Piper et al. 2012).

266 The Trail-Making Test is an index of executive function test and assesses set-shifting. In
267 Set A, the participant clicks on an ascending series of numbers (e.g. 1 – 2 – 3 – 4). In Set B, the
268 participant alternates between numbers and letters (e.g. 1 – A – 2 – B). The primary dependent
269 measure from the five trials is the ratio of total time to complete B/A with lower values (closer to
270 1.0) indicative of better performance. Based on the findings of Study I with five trials, only the
271 first two trials were completed in Study II.

272 In the PEBL default Digit Span forward, strings of numbers of increasing length starting
273 with three were presented via headphones and displayed at a rate of one/second. Audio feedback
274 (e.g. “Correct” or “Incorrect”) was provided after each of three trials at each level of difficulty.
275 The primary dependent measure was the number of trials completed correctly.

276 The Berg Card Sort Test measures cognitive flexibility and requires the participant to sort
277 cards into one of four piles based on a rule (color, shape, number) that changes. Feedback
278 (“correct!” or “incorrect”) was displayed for 500 ms after each trial. This test differs somewhat

279 from the version employed previously (Fox et al. 2013; Piper et al. 2011) in that the prior
280 selections were displayed (Supplementary Figure 1E). The primary dependent measure is the
281 percent of the 64 responses that were perseverative errors defined and coded according to the
282 Heaton criteria (Heaton et al., 1993) although the number of categories completed and
283 perseverative responses was also recorded.

284 In the PEBL Mental Rotation Test, the participant must decide whether two 2-
285 dimensional images are identical or if one is a mirror image. There are a total of 64 trials with the
286 angle of rotation varied in 45° increments (-135° to + 180°). The percent correct and response
287 time were the dependent measures.

288 In the PEBL Iowa Gambling Task, the four decks are labeled 1, 2, 3, and 4 rather than A,
289 B, C, and D (Buelow & Suhr, 2009). The primary dependent measure was the \$ at the end and
290 response preference [(Deck 3 + Deck 4) – (Deck 1 + Deck 2)] with Decks 3 and 4 being
291 advantageous and Decks 1 and 2 being disadvantageous. The response to feedback and the
292 frequency different strategies were employed, e.g. payoff and then change piles (Win-Switch),
293 lose money but continue with the same pile (Lose-Stay), was also documented.

294 In the Tower-of-London, the participant must form a plan in order to move colored disks,
295 one at a time, to match a specified arrangement. The number of points to solve twelve problems
296 (3 points/problem) and the average completion time/problem were recorded. Based on some
297 indications of ceiling effects in Study I, Study II employed a more challenging version of this
298 task (Piper et al. 2012) with the primary measure being moves and completion time as a
299 secondary measure.

300 Dexterity is a recently developed test of fine motor function that consists of a circular
301 coordinate plane with the center of the circle (demarcated by a thin black line) at x,y positions

302 0,0. The goal is to move the cursor (depicted as a colored ball) to a target located at various
303 positions. Movement of the cursor is affected by a “noise” component complementing the
304 directional input from the analog mouse to create the effect of interference or “jittering” motion.
305 The effect is such that successful navigation of the coordinate plane using the mouse encounters
306 resistance to purposeful direction, requiring continual adjustment by the participant to maintain
307 the correct path to the target. Visual feedback is given by the use of a color system, wherein the
308 cursor shifts gradually from green to red as proximity to the target becomes lesser. The task
309 consists of 80 trials (10 per “noise” condition), ten seconds maximum in length, with preset noise
310 factors (ranging in intensity) and target locations standardized for consistency between
311 participants. A lack of input from the participant results in a gradual drift towards the center. At
312 the conclusion of each trial, the cursor location is reset to the origin. Completion time and Moves
313 were recorded with Moves defined as the change in the vector direction of the mouse while
314 course correcting toward the target (Supplemental Figure 1H).

315 Finally, in the Test of Attentional Vigilance, participants are presented with “go/no-go”
316 stimuli that they must either respond or inhibit their response. An abbreviated version (6 min)
317 was employed. The primary dependent measures were the reaction time and the variability of
318 reaction times.

319 *Data Analysis*

320 All analyses were conducted using Systat, version 13.0 with figures prepared using
321 Prism, version 6.03. Ceiling and floor effects were determined by examining score distributions
322 for any measure with $\geq 5\%$ of respondents scoring at the maximum or minimum of the obtainable
323 range on that measure. As the PEBL default criteria for perseverative errors on the Berg Card
324 Sorting Test is currently very different than that employed by Heaton et al., 1993 in the

325 Wisconsin Card Sorting Test, secondary analyses were completed with each definition. Sex
326 differences in Study I and the magnitude of practice effects (Study II) were expressed in terms of
327 Cohen's *d* (e.g. [*Absolute value (Mean_{Retest} - Mean_{Test})/SD_{Test}]* with 0.2, 0.5, and 0.8 interpreted
328 as small, medium, and large effect sizes. In Study II, correlation (*r* and *rho*) and paired t-tests
329 were calculated on the test and retest values. Test-retest correlations > 0.7 were interpreted as
330 acceptable (Nunnally, 1994) and < 0.3 as unacceptable. The percent change was determined in
331 order to facilitate comparison across measures.

332

333 RESULTS

334 *Study I: Normative Behavior & Inter-Test Associations*

335 The ten PEBL tests may be organized into the following broad domains: motor function
336 (Pursuit Rotor and Dexterity), Attention (Test of Attentional Vigilance and Time-Wall),
337 Working-Memory (Digit Span), and Executive Functioning/Decision Making (Trail Making
338 Test, Tower of London, Berg Card Sorting Test, Iowa Gambling Test, and the Mental Rotation
339 Test). Table 1 shows that there were substantial individual differences in this sample. With the
340 exception of the Tower of London (Maximum Possible Points = 36), no test showed evidence of
341 ceiling or floor effects. The Berg criteria for coding perseverative responses resulted in a many
342 more than the Heaton criteria ($\text{Mean}_{\text{Berg}} = 30.8 \pm 6.9\%$, $\text{Mean}_{\text{Heaton}} = 11.9 \pm 8.1\%$, $t(172) =$
343 24.10 , $P < .0005$). The difference for perseverative errors was more subtle but still significant
344 ($\text{Mean}_{\text{Berg}} = 12.9 \pm 5.8\%$, $\text{Mean}_{\text{Heaton}} = 11.0 \pm 6.4\%$, $t(172) = 3.79$, $P < .0005$) on the Berg Card
345 Sorting Test.

346 Overall, sex differences were infrequent. On the Pursuit Rotor, the total time on target
347 was greater in males (47.6 ± 6.0) than females (41.8 ± 7.0 sec, $t(182) = 5.79$, $P < .0005$, $d =$
348 0.89). Further analysis determined that target time in males was elevated by over 1,300 msec on
349 each trial relative to females (Figure 2A). On the Mental Rotation Test, there was no sex
350 difference in the number correct (Females = $72.8 \pm 17.9\%$, Males = $74.9\% \pm 19.4\%$, $t(168) =$
351 0.47). Decision time was increased by the angle and the number correct decreased as the rotation
352 angle extended away from zero degrees in either direction (Figure 2B). The sex difference
353 (Males = $2,377.8 \pm 795.8$, Females = $2,638.0 \pm 863.3$) for overall response time was barely
354 significant ($t(168) = 1.99$, $P < .05$, $d = 0.31$) with more pronounced group differences identified

355 at specific angles (e.g. -45° , $d = 0.51$). Further, on the Iowa Gambling Test, total amount earned
356 at the end of the game did not show a sex difference (Males = $\$1,928.95 \pm 707.99$, Females =
357 $1,858.02 \pm 755.47$, $t(180) = 0.64$, $P = .52$) but, following a loss, Males (3.0 ± 3.5) were 73.1%
358 more likely on their following choice to select again from the same deck (Females = 1.7 ± 2.8 ,
359 $t(136.6) = 2.86$, $P < .01$, $d = 0.41$).

360 Table 2 depicts the correlations among the tests. Generally, the association within
361 measures on a single test was moderate to high (e.g. Pursuit Rotor, Test of Attentional Vigilance,
362 Berg Card Sorting Test) whereas between tests Spearman rho values were typically lower. Lower
363 performance on the Pursuit Rotor (i.e. higher Error) was associated with less attentional
364 consistency (i.e. larger Test of Attentional Vigilance variability), longer times to complete
365 Dexterity, more Perseverative Errors on the Berg Card Sorting Test, greater Time-Wall
366 Inaccuracy, and lower Digit Span forward. There were also several correlations on the indices of
367 executive function. Individuals that performed less well on the Trail Making Test (i.e. higher B
368 to A ratios) scored lower on the Tower of London and the Berg Card Sorting Test. The
369 correlation between Berg Card Sorting Test perseverative errors when coded according to the
370 Heaton and default (Berg) criteria was moderately high. More correct Mental Rotation responses
371 also corresponded with higher performance on the Tower of London. Also noteworthy, the B to
372 A Ratio with all five trials showed a strong correspondence with only the first two Trail Making
373 Test trials ($r_s(178) = +0.90$, $P < .0005$, Figure 2C).

374 *Study II: Test-Retest Reliability*

375 Figure 3 shows the test-retest correlations ranked from highest to lowest. Spearman and
376 Pearson correlations ≥ 0.7 were interpreted as acceptable, ≥ 0.3 and $< .7$ as intermediate, and
377 below 0.3 as unacceptable. Acceptable correlations were identified on the Pursuit Rotar and the

378 Test of Attentional Vigilance. Digit Span, Time-Wall, and most measures on the Berg Card
379 Sorting Test were intermediate. Select correlations were below the acceptability cut-off for the
380 Iowa Gambling Task and the Tower of London. The reliability of secondary measures is also
381 listed on Supplemental Table 3. Most notably, reliability coefficients on the Berg Card Sorting
382 Test were equivalent for perseverative errors with the Berg and Heaton definitions.

383 Figure 4 depicts the absolute reliability in terms of effect size from the test to the retest
384 for the primary dependent measures with Supplemental Table 1 also containing secondary
385 indices. Consistent responding (i.e. no significant change) was observed for the number of moves
386 to solve the Tower of London. Slight, but significant ($P \leq .05$) improvements were noted for
387 Digit Span forward and Response Time on the Test of Attentional Vigilance. Significant ($p <$
388 $.01$) practice effects with a small effect size ($d \geq .2$) were identified for the variability of
389 responding on the Test of Attentional Vigilance, the response pattern on the Iowa Gambling
390 Tasks, the B to A ratio on the Trail Making Test as well as time to complete Part A, and
391 perseverative errors on the Berg Card Sorting task defined according to the Berg criteria.
392 Intermediate ($d \geq .6$) practice effects were identified with increased time on target on the Pursuit
393 Rotor, decreased mean time to solve each Tower of London problem, faster completion of Part B
394 of the Trail Making Test, and heightened accuracy on Time-Wall.

395 Further analysis on the Iowa Gambling Tasks determined that the amount earned at the
396 end of each session did not appreciably change from the test ($\$1944.85 \pm 85.04$) to the retest
397 ($2,162.13 \pm 116.03$, $t(67) = 1.59$, $P = .12$; $r(66) = .10$, $P = .40$). However, the number of
398 selections from the disadvantageous decks (1 and 2) decreased 10.3% from the test (45.6 ± 1.4)
399 to the retest (40.9 ± 1.8 , $t(67) = 2.66$, $P < .01$, $d = 3.9$; $r(66) = .41$, $P < .0005$).

400

401 **DISCUSSION**

402 *Study I: Normative Behavior & Inter-Test Associations*

403 The principle objective of the first study was to evaluate the utility of a collection of tests
404 from the PEBL battery including convergent and divergent validity. As also noted in the
405 introduction, there are some methodological differences between the PEBL and non-PEBL tasks.
406 The difference between using a stylus versus a computer mouse to track a moving target in the
407 Pursuit Rotar/Rotary Pursuit may not be trivial. The TOVA, but not the TOAV, includes
408 microswitches to record response time which may result in a higher accuracy than may occur
409 without this hardware. Finally, some of these instruments have a prolonged history (Lezak et al.
410 2012) and the dependent measures for some commercial tests (e.g. the WCST and perseverative
411 errors) have evolved over the past six decades (Berg, 1948; Grant & Berg, 1948; Heaton et al.
412 1993) to be more complex than may be readily apparent based upon reading only the peer-
413 reviewed literature.

414 The ten measures in this dataset were chosen based on a combination of attributes
415 including assessing distinct neurophysiological substrates (Figure 1), theoretically meaningful
416 constructs (Supplemental Table 2), ease and speed of administration, and frequency of use in
417 earlier publications (Mueller & Piper, 2014). Admittedly, a potential challenge that even
418 seasoned investigators have encountered with a young-adult “normal” population is that they can
419 quickly and efficiently solve novel problems which may result in ceiling effects (Yasen et al.
420 2015). However, a substantial degree of individual differences were identified on almost all
421 measures (Table 1). The only test where there might be some concern about score distribution
422 would be the points awarded on the Tower of London. A future study (e.g. testing the efficacy

423 of a cognitive enhancing drug) might consider: 1) using alternative measures like completion
424 time; 2) choosing one of the ten other Tower of London already included, e.g. the test contained
425 in Piper et al., 2012, or, as the PEBL code is moderately well documented for those with at least
426 an intermediate level programming ability, to 3) develop their own more challenging test using
427 one of the existing measures as a foundation.

428 The inter-relationships among tests were characterized to provide additional information
429 regarding validity. For example, indices of attention showed some associations with both motor
430 function and more complex cognitive domains like memory. Overall, the relatively low
431 correlations ($\approx \pm 0.3$) between the Trail Making Test, Tower of London, and Berg Card Sorting
432 Test, are congruent with the sub-component specificity of executive function domains (Miyake et
433 al., 2000). Similarly, the lack of association of the Iowa Gambling Task with other executive
434 function measures is generally concordant with prior findings (Buelow & Suhr, 2009).

435 This dataset also provided an opportunity to examine whether behavior on this battery
436 was sexually dimorphic. Previously, a small ($d = 0.27$) sex difference favoring boys (ages 9 – 13)
437 was identified on the Pursuit Rotor (Piper, 2011). This same pattern was again observed but was
438 appreciably larger ($d = 0.89$) which raises the possibility that completion of puberty in this
439 young-adult sample may be responsible for augmenting this group difference. On the other hand,
440 in a prior study with 3-dimensional Mental Rotation images and a very similar sample (Yasen et
441 al., 2015), sex differences were noted but the effect size was larger ($d = 0.54$) than the present
442 findings ($d = 0.31$). As the PEBL battery currently uses simple 2-dimensional images, image
443 complexity is likely a contributing factor. Sex differences were not obtained on Time-Wall,
444 Berg Card Sorting Tests, Trail-Making, or Tower of London tests which is in-line with earlier
445 findings (Piper et al. 2012).

446 The Berg Card Sorting Test may be the most frequently employed PEBL test in published
447 manuscripts. As both the Berg Card Sorting Test and the Wisconsin Card Sorting Test are based
448 on the same core procedures (Berg, 1948; Grant & Berg, 1948), these tests appear quite similar
449 from the participant's perspective. However, the sorting rules of Heaton et al. (1993) are
450 considerably more complex than those originally developed (Berg, 1948; Grant & Berg, 1948).
451 The finding that five of the correlations with other tests were significant and of the same
452 magnitude with both and Berg and Heaton rules provides some evidence in support of functional
453 equivalence of these tests.

454 *Study II: Test-retest Reliability*

455 The principle objective of the second study was to characterize the test-retest reliability
456 of the PEBL battery with a two-week interval. The correlation between the test and retest is
457 commonly obtained in these types of investigations (Calamia et al., 2013; Fillmore, 2003; Learck
458 et al., 2004; Lejuez et al., 2005; Lezak et al. 2012; Woods et al., 2010). It is also important to be
459 cognizant that the Pearson or the Spearman correlation coefficients may not fully describe the
460 consistency of measurement when the tested participants show an improvement but maintain
461 their relative position in the sample compared to each other. Therefore, a direct comparison
462 between the test and retest scores was also conducted to quantify the extent of any practice
463 effects.

464 The test-retest correlations were high ($\geq .70$) for the Pursuit Rotor and Test of
465 Attentional Vigilance and moderate ($\geq .30$) for Digit Span, the Berg Card Sorting Test. Some
466 measures on the Iowa Gambling Task and the Tower of London have test-retest reliabilities that
467 were low. It is noteworthy that there is no single value that is uniformly employed as the
468 minimum reliability correlation with some advocates of 0.7 or even 0.8 while others reject the

469 notion of an absolute cut-off (Calamia et al., 2013). In general, an extremely thorough meta-
470 analysis concluded that most tests employed by neuropsychologists have correlations above 0.7
471 with lower values observed for measures of memory and executive function (Calamia et al.,
472 2013). Many tests that are widely used clinically and for research have test-retest reliabilities
473 that are in the 0.3 to 0.7 range (Lowe & Rabbitt, 1998). More specifically, the present findings
474 are slightly higher than what has been reported previously for a computerized Rotary Pursuit
475 task⁴². Direct comparison with other psychometric reports is difficult because the test-retest
476 intervals and the participants characteristics were dissimilar but they are generally in line with
477 expectations. Similarly, the degree of improvement from the test to the retest, whether expressed
478 as the percent change or in terms of Cohen's *d*, are in accord with most earlier findings.
479 However, perhaps surprisingly, there is currently very limited reliability data from the non-PEBL
480 computerized versions of the Iowa Gambling Task or the Wisconsin Card Sorting Test for
481 comparative purposes. Overall, it is important to emphasize that reliability is not an inherent
482 characteristic of a test but instead a value that is influenced by the sample characteristics and the
483 amount of time between the test and retest. The two-week interval would be applicable, for
484 example, to assessing the utility of a cognitive enhancing drug but longer intervals should also be
485 examined in the future.

486 Some procedural details of many of the PEBL tasks employed in this study are worthy of
487 consideration. The numbers presented in Digit Span and the cards in the Berg Card Sorting Test
488 are selected from a set of stimuli such that the retest will not be identical to the test. The degree
489 of improvement would likely be even larger without this feature. Although not the goal of this
490 report, we suspect that the magnitude of practice effects would be attenuated if alternative
491 versions of tests were employed for the test and the retest. This possibility is already pre-

492 programmed into the Trail Making Test and Tower of London. Similarly, the direction of
493 rotation could be set at clockwise for the test and counterclockwise for the retest if additional
494 study determined equivalent psychometric properties independent of the direction of target
495 rotation. Another strategy that could attenuate practice effects might be to increase the number of
496 trials, particularly on Time-Wall and the Trail Making Test, until asymptotic performance was
497 observed. Further discussion of the varied parameters and the evolution of the Iowa Gambling
498 Task is available elsewhere (Piper et al., in review).

499 *General Discussion*

500 The information obtained regarding the validity and reliability of the majority of PEBL
501 tests is broadly consistent with expectations (Lezak et al. 2012; Lowe et al., 1998) and indicates
502 that these tests warrant further use for basic and clinical research. The overall profile including
503 the distribution of scores, convergent and divergent validity, practice effects being of the
504 anticipated magnitude, and, where applicable, internal consistency, as well as an expanding
505 evidence base (Mueller & Piper, 2014), demonstrates that the Rotary Pursuit, Test of Attentional
506 Vigilance, Digit Span, and Trail Making Test are particularly appropriate for inclusion in
507 generalized batteries with participants that are similar to those included in this sample.

508 One task where the psychometric properties are concerning is the Iowa Gambling Task.
509 An improvement was noted in the response pattern from the test to the retest which is consistent
510 with what would be expected with this executive function test *a priori*. However, the correlation
511 between the test and retest was not even significant when the more conservative statistic
512 (Spearman rho) was examined. Perhaps, in order to attenuate the practice effect, two alternative
513 forms of the Iowa Gambling Task could be developed (e.g. version A where decks 3 and 4 are
514 advantageous and a version B where decks 3 and 4 are disadvantageous). In fact, even more

515 sophisticated alternative forms of the Iowa Gambling Task which vary based on task difficulty
516 are being developed by others (Xiao et al. 2013). Another modification which might benefit the
517 test-retest correlation would be to increase the salience of feedback that follows each trial. The
518 feedback was very salient in the original (i.e. non computerized) version of this task in that the
519 experimenter would give or take money after each trial (Bechara et al., 1994) . Perhaps, the
520 psychometric properties of the PEBL Iowa Gambling Task would be improved if auditory
521 feedback was presented after each trial or there were a fixed interval between trials which would
522 encourage the participant to reflect on their previous selection. These procedural modifications
523 were made for a subsequent study (Piper et al. in review). Overall, additional study is warranted
524 to better appreciate the present findings as there is no long-term test-retest reliability with the
525 non-PEBL computerized Iowa Gambling Task (Buelow & Suhr, 2009). However, given the
526 limited evidence for convergent validity or test-retest reliability, prior findings with the PEBL
527 Iowa Gambling Task (Lipnicki et al., 2009); may need to be cautiously interpreted.

528 Three limitations of this report should also be acknowledged. First, the PEBL battery also
529 includes many other indices (e.g. Cori's block tapping test of visuospatial working memory, a
530 Continuous Performance Test of vigilance, a Stroop test of executive functioning). Only a subset
531 of the many PEBL tests were utilized due to time constraints (approximately one-hour of
532 availability for each participant). Future investigations may be designed to focus more narrowly
533 on specific domains (e.g. motor function). Second, a future objective would be to provide
534 further information regarding criterion validity, e.g. by determining the similarities, or
535 differences, between the Test of Variables of Attention and PEBL Test of Attentional Vigilance
536 in neurologically intact and various clinical groups as this information is mostly unavailable for
537 the PEBL tests (although see Danckert et al. 2012 which utilized the Berg Card Sorting Tests and

538 brain injured patients). Third, the sample in both studies consisted of young-adult college
539 students, primarily Caucasian and from a middle-class background. There are those that are quite
540 articulate in outlining the limitations of this population (Henrich, Heine, & Norenzayan, 2010;
541 Reynolds, 2010). The data contained in this report should just be viewed as an important first
542 step as further investigations with different ages, socioeconomic, and ethnic groups is needed.

543 There have been several pioneers in the development of new measures which have
544 greatly facilitated our understanding of individual differences in neurobehavioral function (Lezak
545 et al., 2012). We feel that the transparency of the PEBL battery extends upon this earlier work
546 and provides an important alternative to commercial tests. In addition, the ability of anyone with
547 a functional computer, independent of their academic degrees, to use PEBL contributes to the
548 democratization of science.

549 On the other hand, two considerations with PEBL and other similar open-source
550 applications should be acknowledged. First, the flexibility of PEBL also has clear drawbacks in
551 that each investigator can, in theory, modify a test's parameters to meet their own experimental
552 needs. If an investigator reports that they employed a particular test from a specific commercial
553 distributor, there is wide-spread agreement about what this means as many these tests often have
554 only limited modifiability. However, if an investigator changes a PEBL test but fails to make the
555 programming code available, then it is more difficult to critically evaluate research findings. The
556 second potential drawback with PEBL may be ethical. The prohibition against clinical
557 psychologists (American Psychological Association, 2002), but not others, making
558 neurobehavioral tests readily available is discussed elsewhere (Mueller & Piper, 2014). The
559 accessibility of PEBL to anyone, including psychiatrists, neurologists, or cognitive

560 neuroscientists for research or teaching purposes is consistent with the ethos of science (Merton,
561 1979).

562 These findings also begin to aid comparisons with other older neurobehavioral test
563 batteries. Table 3 contrasts PEBL with the Behavioral Assessment and Research System (BARS)
564 and, perhaps the current “gold standard” of batteries, the Cambridge Neuropsychological Test
565 Automated Battery (CANTAB) in terms of intellectual origins, the not insignificant differences
566 in price and transparency, and sample tests. The BARS system is based on the behavioral
567 analysis principles of B.F. Skinner and is designed for testing diverse populations including
568 those with limited education and prior computer experience (Rohlman et al. 2003). The
569 CANTAB battery was designed with an emphasis on translating preclinical findings to humans
570 (Robbins et al. 1994). Each of these platforms have their own advantages and disadvantages with
571 the strength of PEBL being the number of tests, limited cost, and modifiability.

572 **CONCLUSION**

573 In closing, our hope is that thorough, but critical, investigations of the psychometric
574 properties of this novel methodology in normal (present study) and atypical populations will
575 insure that PEBL will continue to be widely used by investigators in basic and applied areas.
576 This will foster further integration between these fields and further advance our understanding of
577 the genetic, biochemical, and neuroanatomical substrates of individual differences in
578 neurocognition.

579

580 **ACKNOWLEDGEMENTS**

581 The technical assistance of Christopher J. Fox, Vera E. Warren, Hannah Gandsey, Sari N.
582 Matisoff, and Donna M. Nolan is gratefully recognized.

584

585

REFERENCES

586

587 American Psychological Association. 2002. Ethical principles of psychologists and code of
588 conduct. *American Psychologist* **47**: 1060–1073.

589

590 Ammons RB, Alprin SI, Ammons CH. 1955. Rotary pursuit performance as related to sex and
591 age of pre-adult subjects. *Journal of Experimental Psychology* **49**:127-133.

592

593 Axelrod BN, Woodard JL, Henry RR. 1992. Analysis of an abbreviated form of the Wisconsin
594 Card Sorting Test. *Clinical Neuropsychology* **6**: 27–31.

595

596 Barrett DW, Gonzalez-Lima F. 2013. Transcranial infrared laser stimulation produces beneficial
597 cognitive and emotional effects in humans. *Neuroscience* **230**: 13-23.

598

599 Bechara A, Damasio H, Damasio A. 2000. Emotion, decision making and the orbitofrontal
600 cortex. *Cerebral Cortex* **10**: 295-307.

601

602 Bechara A., Damasio AR, Damasio H, Anderson SW. 1994. Insensitivity to future consequences
603 following damage to human prefrontal cortex. *Cognition* **50**: 7-15.

604

605 Berg EA. 1948. A simple objective technique for measuring flexibility in thinking. *Journal of*
606 *General Psychology* **39**: 15-22.

607

608 Buelow MT, Suhr JA. 2009. Construct validity of the Iowa Gambling Task. *Neuropsychology*
609 *Review* **19**: 102-114.

610

611 Calamia M, Markon K, & Tranel D. The robust reliability of neuropsychological measures:
612 Meta-analyses of test-retest correlations. *Clinical Neuropsychology* **27**: 1077-1105.

613

614 Carlin D, Bonerba J, Phipps M, Alexander G, Shapiro M, Grafman J. 2000. Planning
615 impairments in frontal lobe dementia and frontal lobe lesion patients. *Neuropsychologia* **38**: 655-
616 665.

617

618 Chanraud S, Reynaud M, Wessa M, Penttila J, Kostogiannia N, Cachia A, et al. 2009. Diffusion
619 Tensor Tractography in Mesencephalic bundles: Relation to mental flexibility in detoxified
620 alcohol-dependent subjects. *Neuropsychopharmacology* **34**: 1223-1232.

621

622 Danckert J, Stöttinger E, Quehl N, Anderson B. 2012. Right hemisphere brain damage impairs
623 strategy updating. *Cerebral Cortex* **22**: 2745-2760.

624

625 Demakis GJ. 2004. Frontal lobe damage and tests of executive processing: A meta-analysis of
626 the category test, stroop test, and trail-making test. *Journal of Clinical & Experimental*
627 *Neuropsychology* **26**: 441-450.

628

629 Fillmore MT. 2003. Reliability of a computerized assessment of psychomotor performance and
630 its sensitivity to alcohol-induced impairment. *Perceptual Motor Skills* **97**: 21-34.
631

632 Fox CJ, Mueller ST, Gray ST, Raber J, Piper BJ. 2013. Evaluation of a short-form of the Berg
633 Card Sorting Test. *PLoS One* **8**: e63885.
634

635 Gaudino EA, Geisler MW, Squires NK. 1995. Construct validity in the Trail Making Test: What
636 makes part B harder? *Journal of Clinical & Experimental Neuropsychology* **17**: 529–535.
637

638 Gerton BK, Brown TT, Meyer-Lindenberg A, Kohn P, Holt JL, Olsen RK, Berman KF. 2004.
639 Shared and distinct neurophysiological components of the digits forward and backward tasks as
640 revealed by functional neuroimaging. *Neuropsychologia* **42**: 1781-1787.
641

642 Grafton ST, Mazziotta JC, Presty S, Friston K.J, Frackowlak SJ, Phelps ME. 1992. Functional
643 anatomy of human procedural learning determined with regional cerebral blood flow and PET.
644 *Journal of Neuroscience* **12**: 2542-2548.
645

646 Grant DA, Berg EA. 1948. A behavioral analysis of degree of reinforcement and ease of shifting
647 to new responses in Weigl-type card-sorting problem. *Journal of Experimental Psychology* **38**:
648 404-411.
649

650 González-Giraldo Y, Rojas J, Novoa P, Mueller ST, Piper BJ, Adan A, Forero DA. 2014.
651 Functional polymorphisms in BDNF and COMT genes are associated with objective differences

652 in arithmetical functioning in a sample of young adults. *Neuropsychobiology* **70**: 152-7. doi:
653 10.1159/000366483.

654

655 González-Giraldo Y, González-Reyes RE, Mueller ST, Piper BJ, Adan A, Forero DA. 2015a.
656 Differences in planning performance, a neurocognitive endophenotype, are associated with a
657 functional variant in PER3 gene. *Chronobiology International* **32**: 591-595. doi:
658 10.3109/07420528.2015.1014096.

659

660 González-Giraldo Y, Rojas J, Mueller ST, Piper BJ, Adan A, Forero DA. 2015b. BDNF
661 Val66Met is associated with performance in a computerized visual-motor tracking test in healthy
662 adults. *Motor Control* in press.

663

664 Greenberg LM, Waldman ID. 1993. Developmental normative data on the Test of Variables of
665 Attention (T.O.V.A.TM). *Journal of Child Psychology & Psychiatry* **34**:1019-1030.

666

667 Heaton RK, Chelune GJ, Talley JL, Kay GG, Curtiss G. 1993. Wisconsin card sorting test
668 manual: Revised and expanded. Odessa: Psychological Assessment Resources.

669

670 Henrich J, Heine SJ, Norenzayan A. 2010. The weirdest people in the world. *Behavioral & Brain*
671 *Sciences* **33**: 61-83. doi: 10.1017/S0140525X0999152X.

672

673 Huang YS, Chao CC, Wu YY, Chen YY, Chen CK. 2007. Acute effects of methylphenidate on
674 performance during the Test of Variables of Attention in children with attention
675 deficit/hyperactivity disorder. *Psychiatry & Clinical Neurosciences* **61**: 219-225.

676

677 Hugdahl K, Thomsen T, Ersland L. 2006. Sex differences in visuo-spatial processing: An fMRI
678 study of mental rotation. *Neuropsychologia* **44**: 1575-1583.

679

680 Jacobson SC, Blanchard M, Connolly CC, Cannon M, Garavan H. 2011. An fMRI investigation
681 of a novel analogue to the Trail-Making Test. *Brain & Cognition* **77**: 60-70.

682

683 Kaneko H, Yoshikawa T, Nomura K, Ito H, Yamauchi H, Ogura M, Honjo S. 2011.

684 Hemodynamic changes in the prefrontal cortex during digit span task: A near-infrared
685 spectroscopy study. *Neuropsychobiology* **63**: 59-65.

686

687 Lezak MD, Howieson DB, Bigler ED, Tranel D. *Neuropsychological Assessment*, 2012, fifth
688 edition, Oxford, New York.

689

690 Learck RA, Wallace DR, Fitzgerald R. 2004. Test-retest reliability and standard error of
691 measurement for the test of variables of attention (T.O.V.A.) with healthy school-age children.

692 *Assessment* **11**: 285-289.

693

694 Lejuez CW, Aklin WM, Richards JB, Strong DR, Karler CW, Read JP. 2005. The Balloon
695 Analogue Risk Task (BART) differentiates smokers and nonsmokers. *Journal of Experimental*
696 *and Clinical Neuropsychology* **11**: 26-33.

697

698 Linn M, Petersen A. 1985. Emergence and characterization of sex differences in spatial ability:
699 A meta-analysis. *Child Development* **56**: 1479-1498.

700

701 Lipnicki DM, Gunga H, Belavy DL, Felsenberg D. Decision making after 50 days of simulated
702 weightlessness. *Brain Research* **1280**: 84-89.

703

704 Lowe C., Rabbitt P. 1998. Test/re-test reliability of the CANTAB and ISPOCD
705 neuropsychological batteries: Theoretical and practical issues. *Neuropsychologia* **36**: 915-923.

706

707 Merton RK. 1979. *The sociology of science: Theoretical and empirical investigations*, University
708 of Chicago: Chicago.

709

710 Miyake A, Friedman NP, Emerson MJ, Witzki AH, Howerter A. 2000. The unity and diversity
711 of executive functions and their contributions to complex “Frontal Lobe” tasks: A latent variable
712 analysis. *Cognitive Psychology* **41**: 49-100.

713

714 Moore DS, Johnson SP. 2008. Mental rotation in human infants: A sex difference. *Psychological*
715 *Science* **19**: 1063-1066.

716

717 Morris RG, Rushe T, Woodruffe PWR, Murray RM. 1995. Problem solving in schizophrenia: A
718 specific deficit in planning ability. *Schizophrenia Research* **14**: 235-246.

719

720 Mueller ST. 2010. A partial implementation of the BICA cognitive decathlon using the
721 Psychology Experiment Building Language (PEBL). *International Journal of Machine*
722 *Consciousness* **2**: 2273-2288.

723

724 Mueller ST. 2014a. The Psychology Experiment Building Language, Version 0.14. Retrieved
725 from <http://pebl.sourceforge.net>.

726

727 Mueller ST. 2014b. The PEBL Manual, Version 0.14, Lulu Press, ISBN 978-0557658176.

728

729 Mueller ST, Piper BJ. 2014. The Psychology Experiment Building Language (PEBL) and PEBL
730 test battery. *Journal of Neuroscience Methods* **222**: 250-259. doi:

731 10.1016/j.jneumeth.2013.10.024.

732

733 Nunnally JC, Bernstein IH. Psychometric theory (3rd ed). New York: McGraw Hill, 1994.

734

735 Piper BJ. 2010. Age, handedness, and sex contribute to fine motor behavior in children. *Journal*
736 *of Neuroscience Methods* **195**: 88-91. doi: 10.1016/j.jneumeth.2010.11.018.

737

738 Piper BJ. 2012. Evaluation of the test-retest reliability of the PEBL continuous performance test
739 in a normative sample. PEBL Technical Report Series [On-line], #2012-05,

740 <http://sites.google.com/site/pebltechnicalreports/home/2012/pebl-technical-report-2012-05>

741

742 Piper BJ, Li V, Eowiz M, Kobel Y, Benice T, Chu A, et al. 2012. Executive function on the

743 Psychology Experiment Building Language test battery. *Behavior Research Methods* **44**: 110-

744 123.

745

746 Power Y, Goodyear B, Crockford D. 2012. Neural correlates of pathological gamblers preference

747 for immediate rewards during the Iowa Gambling Task: An fMRI study. *Journal of Gambling*

748 *Studies* **28**: 623-636.

749

750 Premkumar M, Sable T, Dhanwal D, Dewan, R. 2013. Circadian levels of serum melatonin and

751 cortisol in relation to changes in mood, sleep and neurocognitive performance, spanning a year

752 of residence in Antarctica. *Neuroscience Journal* **2013**; 254090.

753

754 Reynolds CR. 2010. Measurement and assessment. *Psychology Assessment* **22**: 1-4.

755

756 Richardson JTE. 2007. Measures of short-term memory: A historical review. *Cortex* **43**: 635-

757 650.

758

759 Robbins TW, James M, Owen AM, Sahakian BJ, McInnes L, Rabbitt P. 1994. Cambridge
760 Neuropsychological Test Automated Battery (CANTAB): A factor analytic study of a large
761 sample of normal elderly volunteers. *Dementia* **5**: 266-281.

762

763 Rogalsky C, Vidal C, Li X, Damasio H. 2012. Risky decision-making in older adults without
764 cognitive deficits: An fMRI study of VMPFC using the Iowa Gambling Task. *Social*
765 *Neuroscience* **7**: 178-190.

766

767 Rohlman DS, Gimenes LS, Eckerman DA, Kang SK, Farahat F, Anger WK. 2003. Development
768 of the Behavioral Assessment Research System (BARS) to detect and characterize neurotoxicity
769 in humans. *Neurotoxicology* **24**: 523-531.

770

771 Schall U, Johnson P, Lagopoulos J, Juptner M, Jentzen W, Thienel R, et al. 2013. Functional
772 brain maps of Tower of London performance: A positron emission tomography and functional
773 magnetic resonance imaging study. *Neuroimage* **20**: 1154-1161.

774

775 Schmidtke K, Manner H, Kaufmann R, Schmolck H. 2002. Cognitive procedural learning in
776 patients with fronto-striatal lesions. *Learning & Memory* **9**: 419-429.

777

778 Shallice T. 1982. Specific impairments of planning. *Philosophical Transactions of the Royal*
779 *Society of London B: Biological Sciences* **298**: 199-209.

780

781 Specht K, Lie CH, Shah NJ, Fink GR. 2009. Disentangling the prefrontal network for rule
782 selection by means of a non-verbal variant of the Wisconsin Card Sorting Test. *Human Brain*
783 *Mapping* **30**: 1734-1743.
784
785 Stirling LA, Lipsitz LA, Qureshi M, Kelty-Stephan DG, Goldberger AL, Costa MD. 2013. Use
786 of a tracing task to assess visuomotor performance: Effects of age, sex, and handedness. *Journals*
787 *of Gerontology. Series A: Biological Sciences & Medical Sciences* **68**: 938-45. doi:
788 10.1093/gerona/glt003.
789
790 Tana MG, Montin E, Cerutti S, Bianchi AM. 2010. Exploring cortical attentional system by
791 using fMRI during a continuous performance test. *Computational Intelligence & Neuroscience*
792 **329213**. doi: 10.1155/2010/329213.
793
794 Wardle MC, Hart AB, Palmer AA, Wit H. 2012. Does COMT genotype influence the effects of
795 d-amphetamine on executive functioning? *Genes, Brain and Behavior* **12**: 13-20.
796
797 Willingham DB, Hollier J, Joseph J. 1995. A Macintosh analogue of the rotary pursuit task.
798 *Behavior Research Methods* **27**: 491-495.
799
800 Woods DL, Kishiyama MM, Yund EW, Herron TJ, Edwards B, Poliva O, et al. 2010. Improving
801 digit span assessment of short-term verbal memory. *Journal of Clinical & Experimental*
802 *Neuropsychology* **33**: 101-111. doi: 10.1080/13803395.2010.493149.
803

804 Xiao L, Wood SMW, Denburg NL, Moreno GL, Hernandez M, Bechara A. 2013. Is there a
805 recovery of decision-making function after frontal lobe damage? A study using alternative
806 versions of the Iowa Gambling Task. *Journal of Clinical & Experimental Neuropsychology* **35**:
807 518-529. doi: 10.1080/13803395.2013.789484.

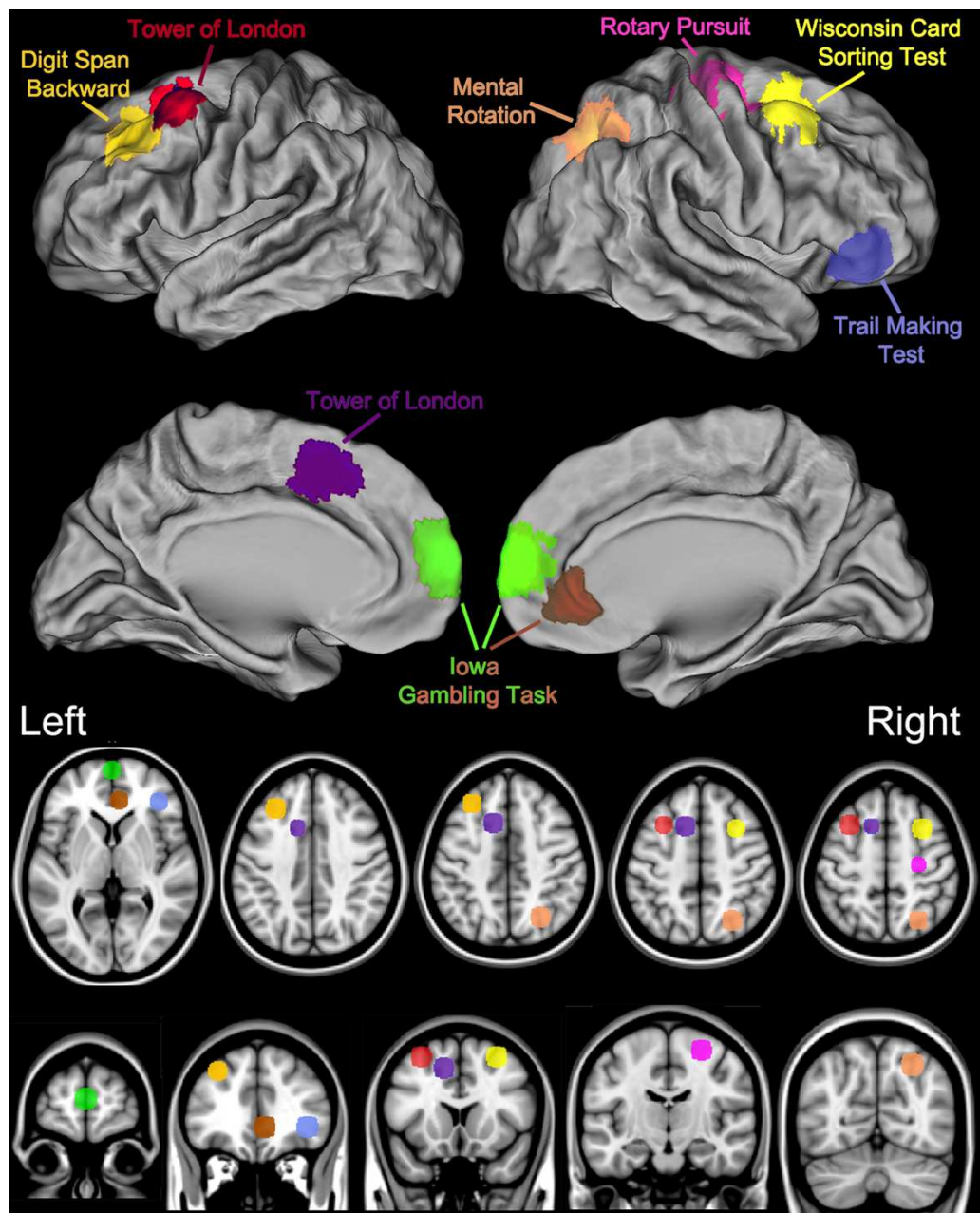
808

809 Yasen AL, Raber J, Miller JK, Piper BJ. 2015. Sex, but not Apolipoprotein E genotype,
810 contributes to spatial performance in young-adults. *Archives of Sexual Behavior* in press.

811

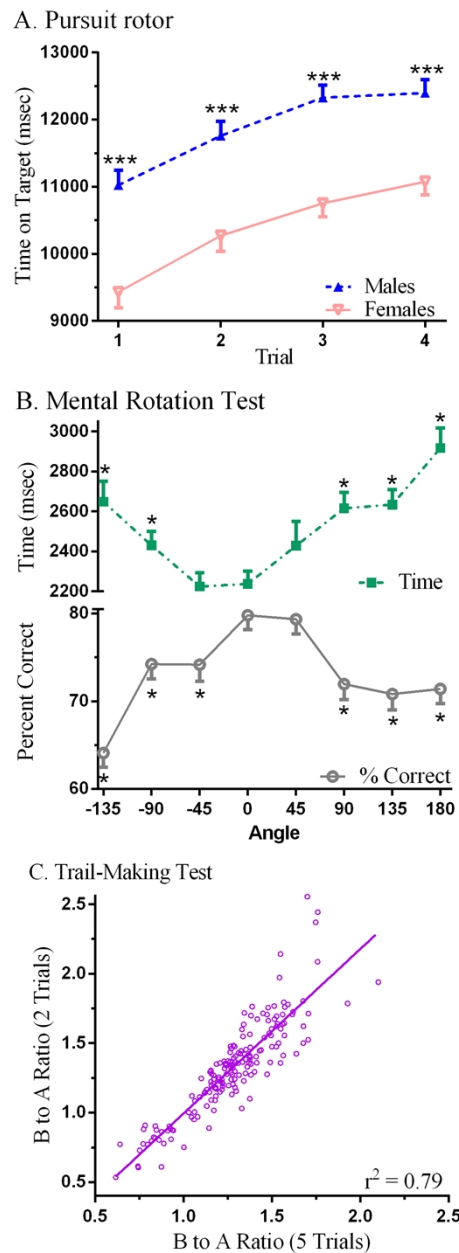
812 Zacks JM. 2008. Neuroimaging studies of mental rotation: A meta-analysis and review. *Journal*
813 *of Cognitive Neuroscience* **20**: 1-19.

815 **Figure 1.** Key brain areas as identified by neuroimaging and lesion studies and the
816 corresponding Psychology Experiment Building Language Tests.



817

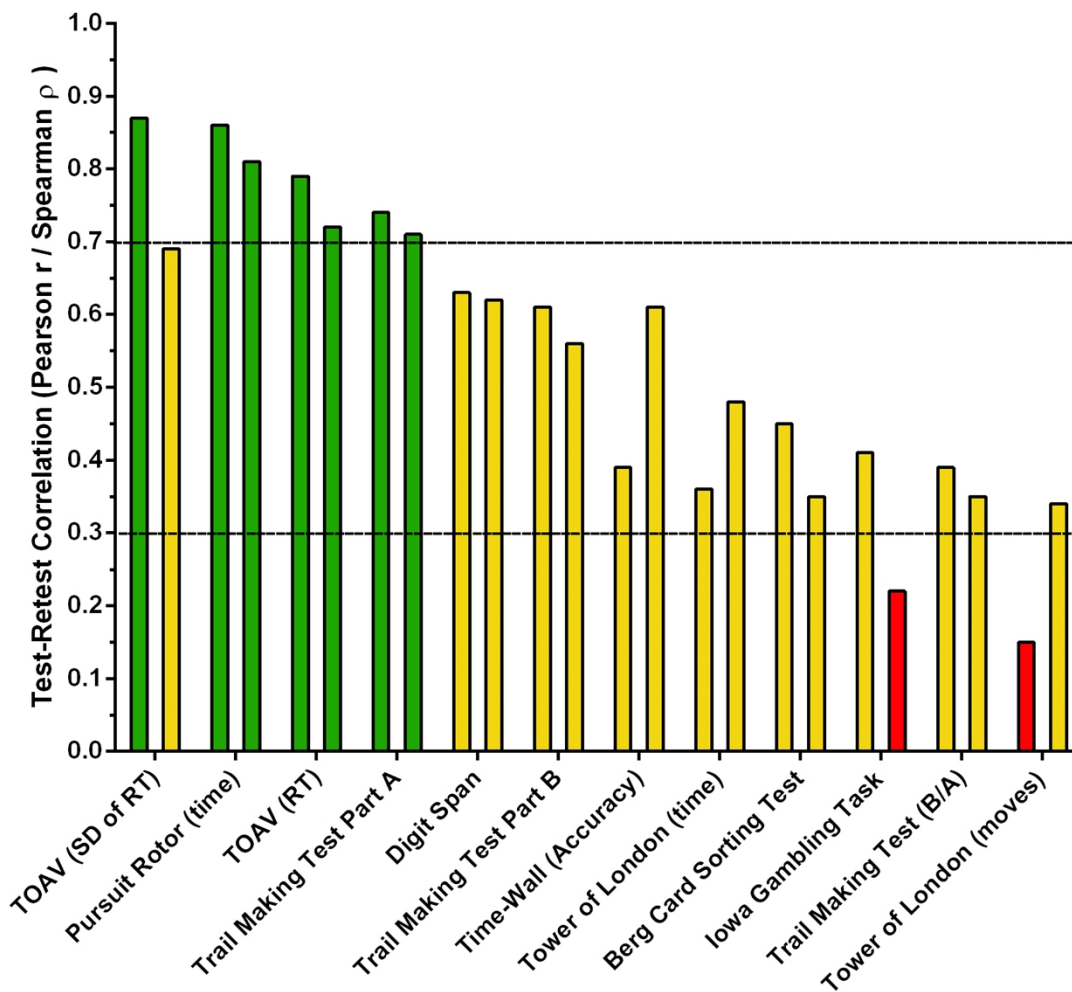
818 **Figure 2.** Neurobehavioral performance on Psychology Experiment Building Language (PEBL)
 819 tests. A) Time on target on the Pursuit Rotor ($***P < .0005$ versus Females); B) Decision time
 820 and percent correct on the Mental Rotation ($*P < .0005$ versus Angle = 0°); C) Scatterplot of
 821 the ratio (Part B/Part A) of times to complete five versus two trials of the PEBL Trail-Making
 822 Test.



823

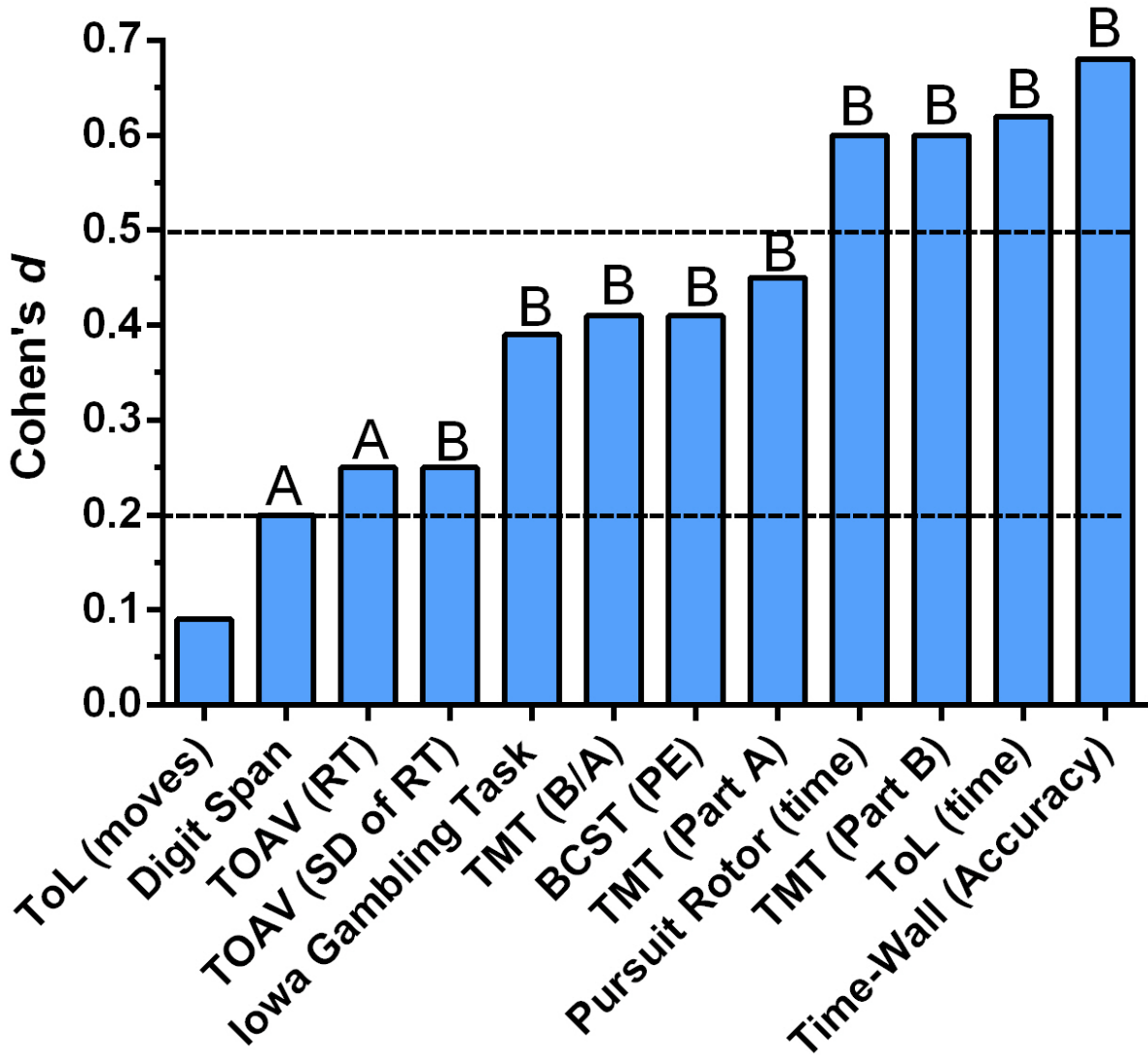
824

825 **Figure 3.** Test-retest correlations ranked from highest to lowest. For each Psychology
826 Experiment Building Language Test, the Pearson r is listed first followed by the Spearman ρ .
827 Correlations $\geq .7$ are acceptable and below 0.3 as unacceptable. RT: Response Time; TOAV:
828 Test of Attentional Vigilance; Trail Making Test Ratio of Completion times for Part B/Part A
829 (B/A).



830

831 **Figure 4.** Change from the test to the retest, expressed as Cohen's *d* measure of effect size,
 832 among young-adults completing the Psychology Experiment Building Language (PEBL)
 833 neurobehavioral test battery. Paired t-test ^A*P* < .05, ^B*P* < .01.

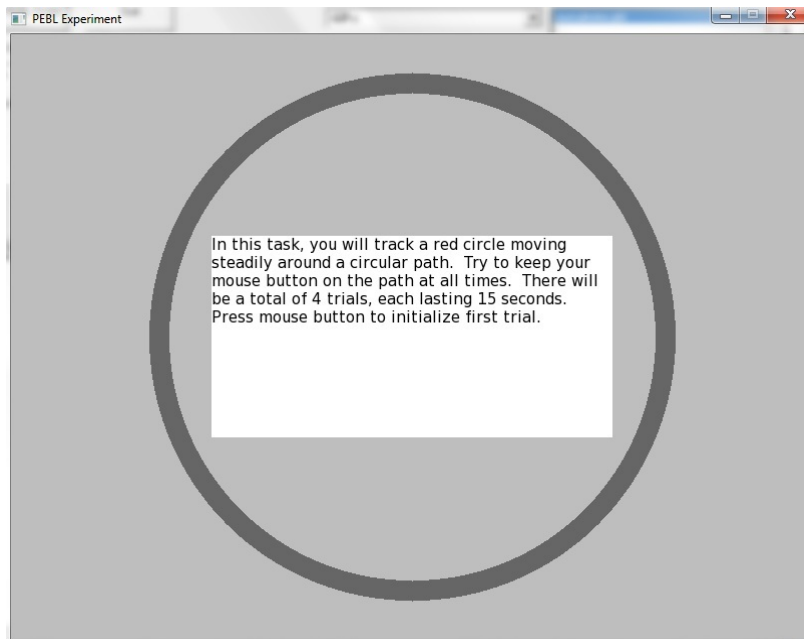


834
 835
 836
 837

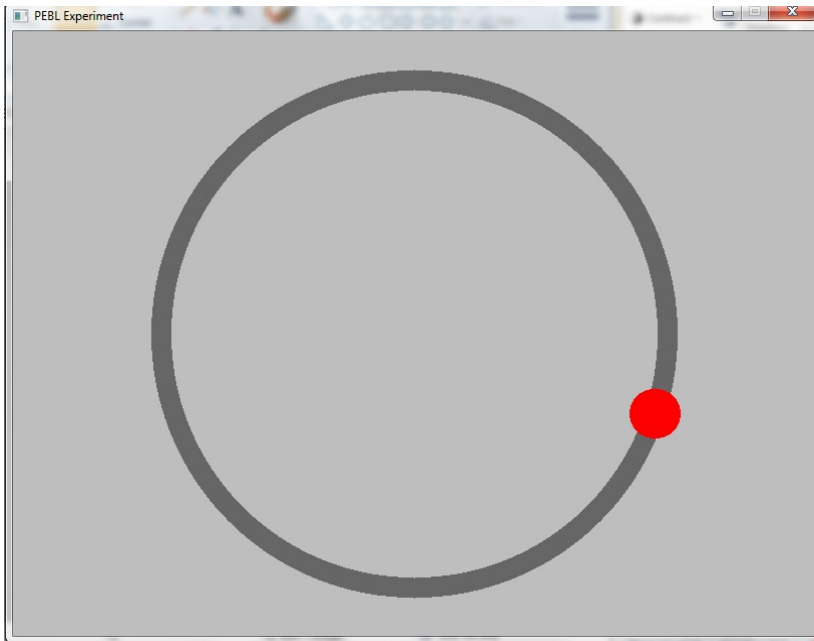
839 **Supplementary Figure 1.** Psychology Experiment Building Language (PEBL) screen-shots
840 including Pursuit Rotor (A), Time-Wall (B), Trail-Making Test (C), Digit Span (D), Wisconsin
841 (Berg) Card Sorting Test (E), Mental Rotation (F), Tower of London (G), Dexterity (H), and the
842 Test of Attentional Vigilance (I).

843

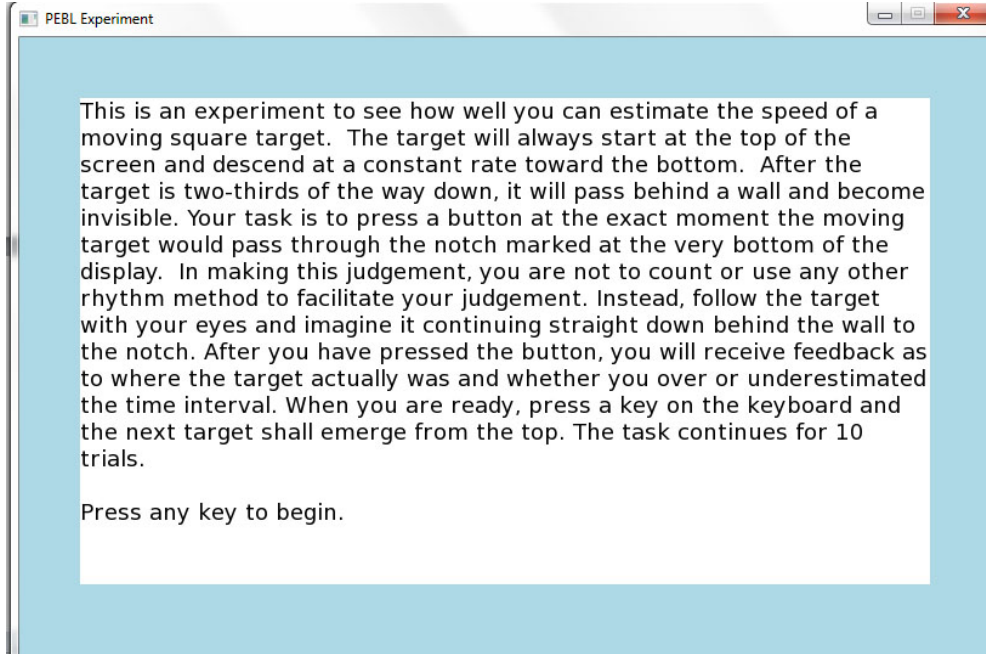
844 A) Pursuit Rotor, a test of fine-motor learning, instructions (top) and example trial (bottom).



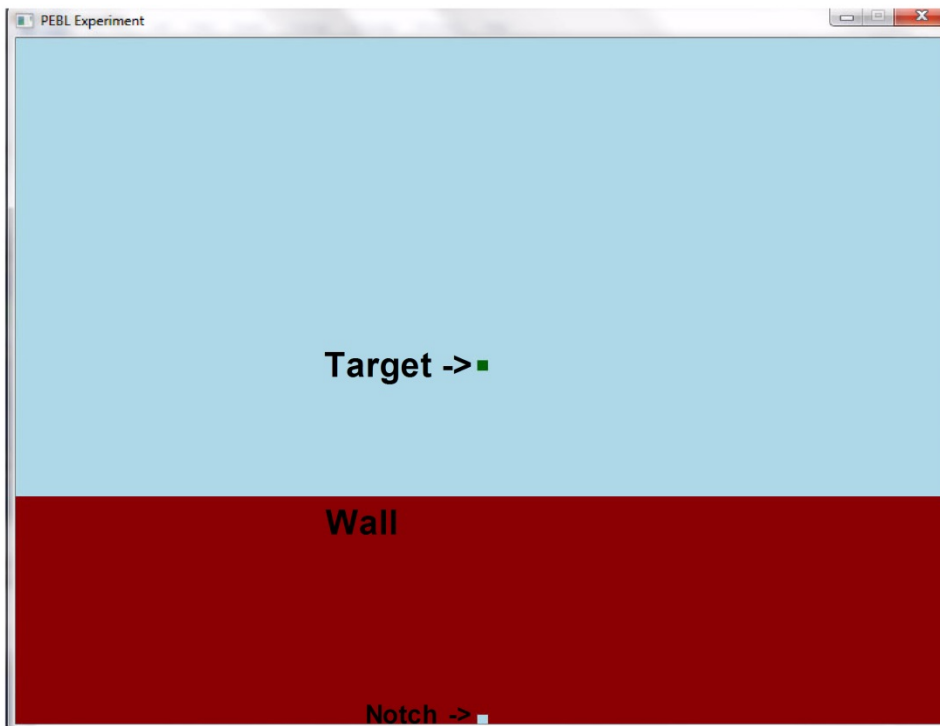
845



848 B. Time-Wall, an index of attention and decision-making, instructions (top) and trial (bottom).



849

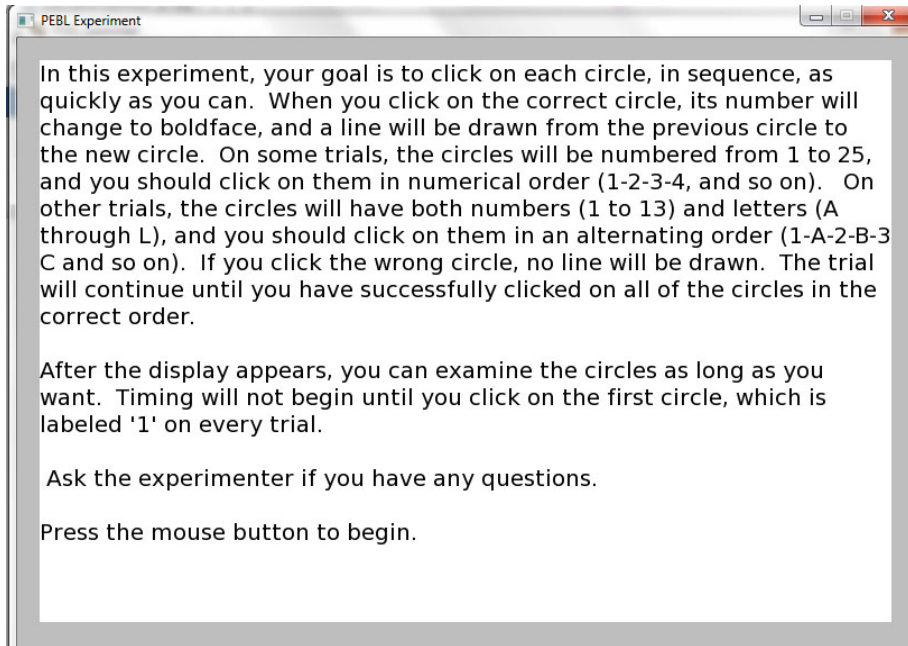


850

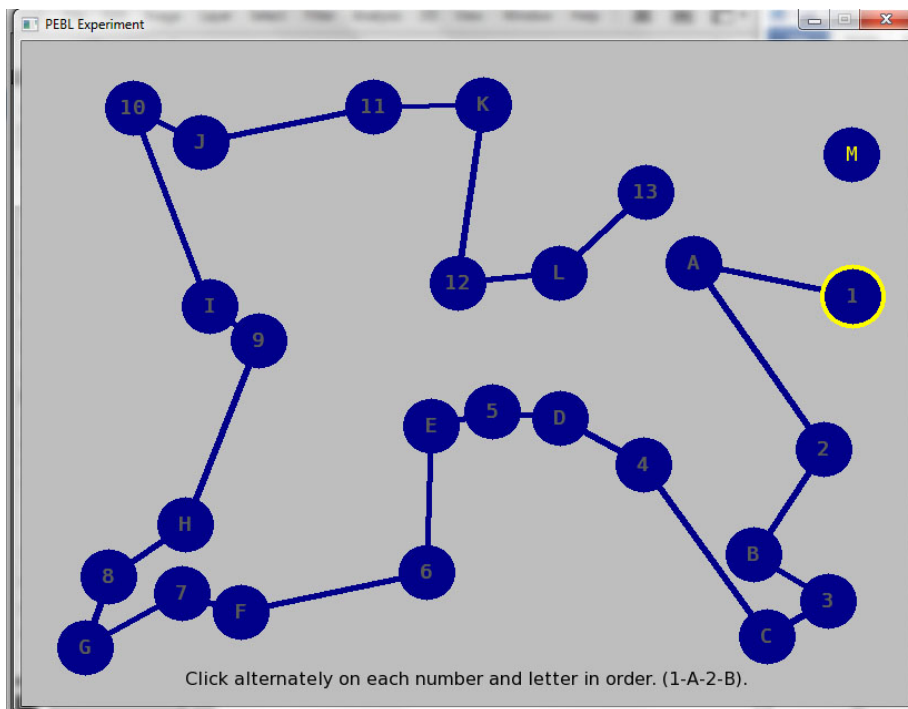
851

852

853 C. PEBL Trail-Making Test, a measure of executive function (set-shifting), instructions (top)
854 and example B trial (bottom).



855



856

857

858

859 **D.** PEBL Digit Span, a measure of working memory, instructions.

You are about to take part in a memory test. You will be presented with a sequence of digits, one at a time on the screen. Each digit will occur only once during a list. You will then be asked to type the list of digits exactly in order. If you do not know what digit comes next, you can skip over it by typing the '-' key. Once entered, you cannot go back to edit your responses. You will start with a list of three items, and will get three different lists at each length. If you are able to recall two out of three lists completely correctly, you will move on to the next longest list length.

860

861

862

863

864

865

866

867

868

869

870

871

872

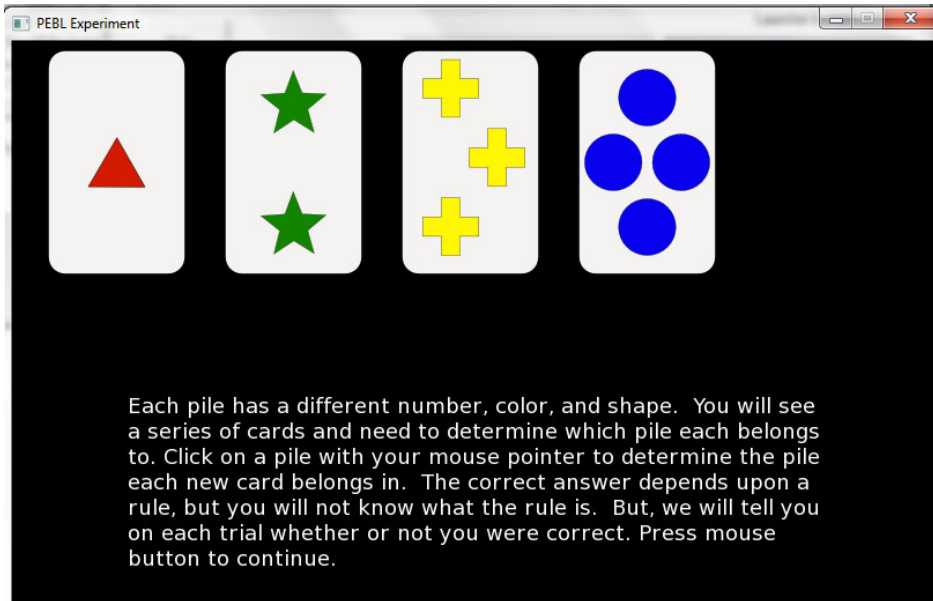
873

874 E. PEBL Berg Card Sort Test, a test of executive function (set shifting), instructions (top and
875 middle) and example trial with color as the current rule (bottom).

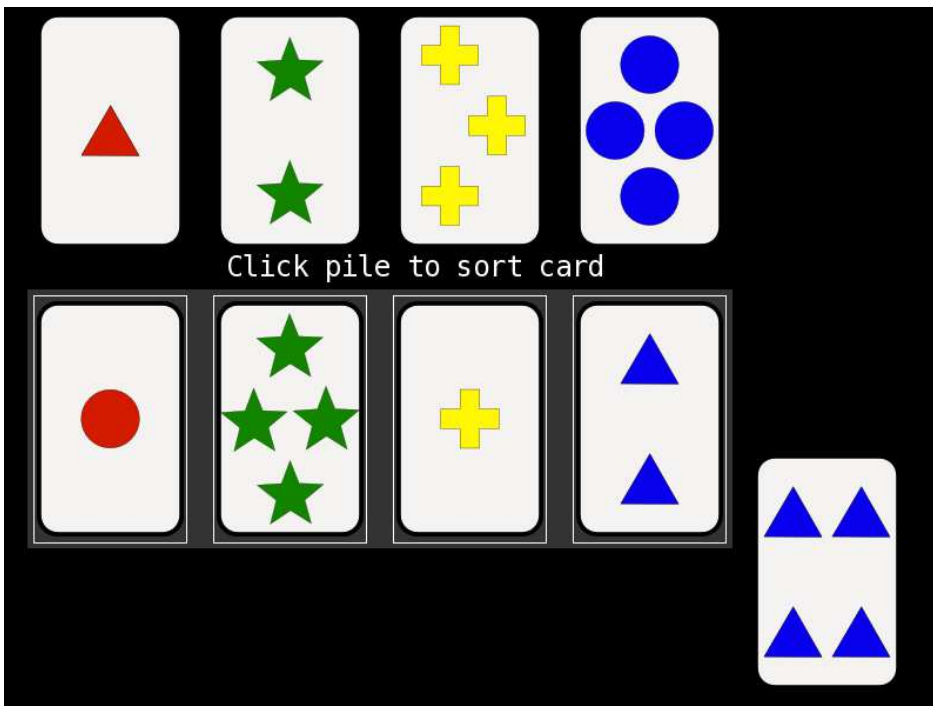
876

You are about to take part in an experiment in which you need to categorize cards based on the pictures appearing on them. To begin, you will see four piles (press the mouse button to see the four piles.)

877



878



879

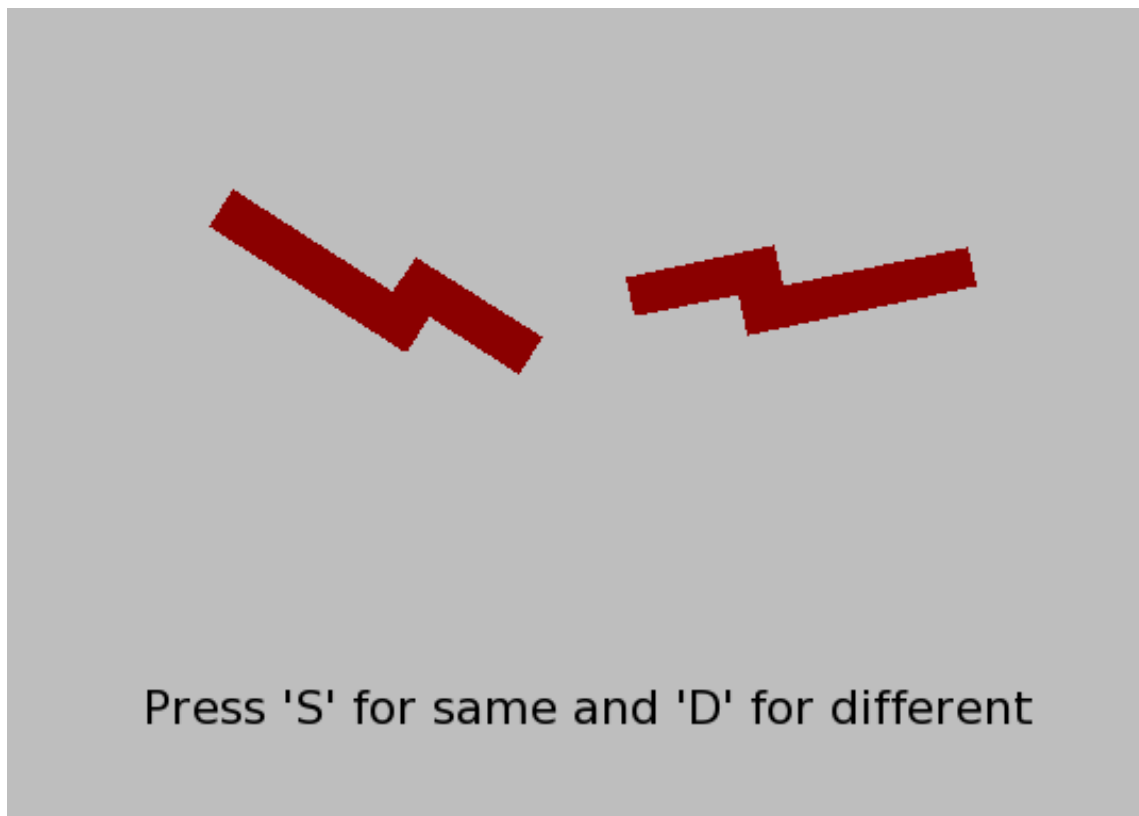
880 **F. Mental Rotation**, an index of executive functioning, instructions (top) and example stimuli

881 (bottom).

This experiment will examine your ability to mentally rotate one figure to compare it with another. You will see a 5 by 5 grid, with five of its cells lighted. You should learn the pattern as quickly and as accurately as possible, and then press a button on the keyboard when you are sure you know the pattern. As soon as you press the key, a new pattern will be presented. If the new pattern is the same as the old pattern, but turned 90 degrees to the left or right, press the left shift key on the keyboard. If the pattern is not a 90 degree left or right rotation of the old pattern, press the right shift key on the keyboard. If you have any questions, please ask the experimenter now.

Press any key to begin

882



883

884

885

886

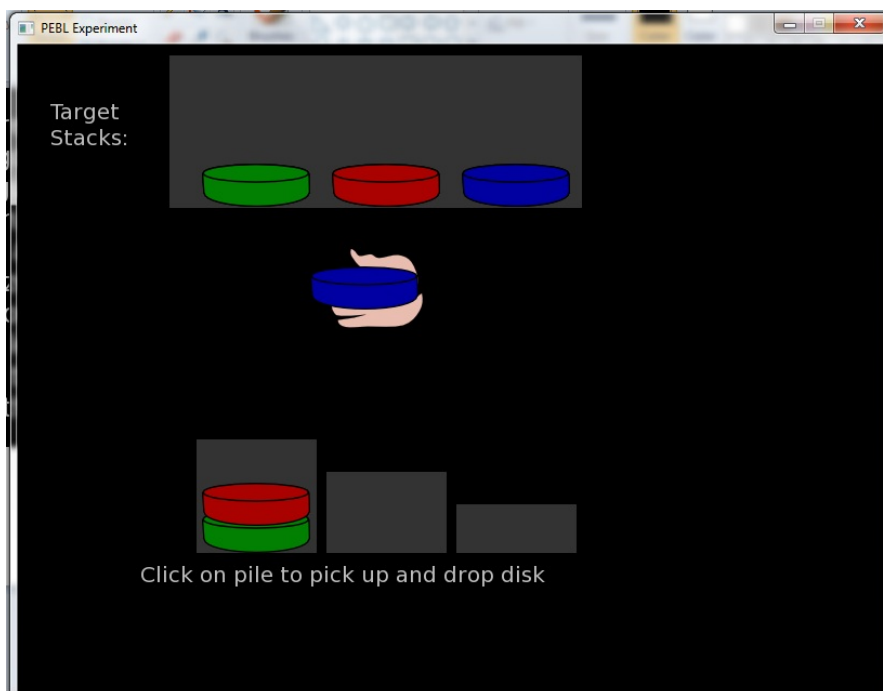
887

888

889 **G.** PEBL Tower of London, a measure of executive-function (planning), instructions (top) and a
890 sample trial (bottom). This version was used in Study II.

You are about to perform a task called the 'Tower of London'.
Your goal is to move a pile of disks from their original
configuration to the configuration shown on the top of the
screen. You can only move one disk at a time, and you cannot
move a disk onto a pile that has no more room (indicated by
the size of the grey rectangle). To move a disk, click on the
pile you want to move a disk off of, and it will move up above
the piles. Then, click on another pile, and the disk will move
down to that pile.
Click the mouse to begin.

891

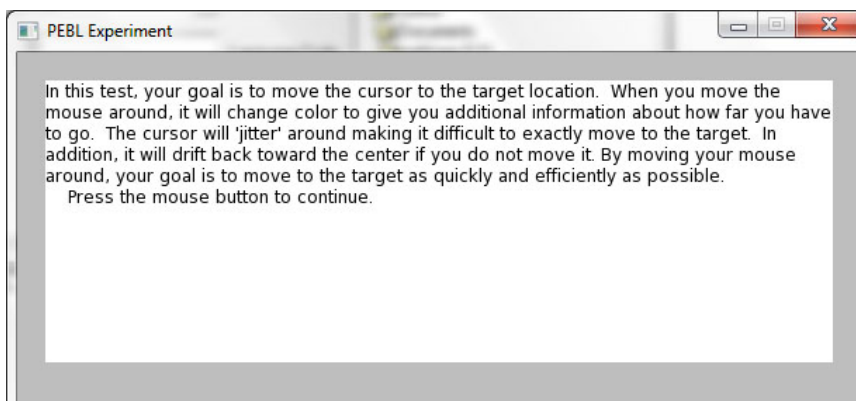


892

893

894

895 **H.** Dexterity, a test of fine motor ability, instructions (top) and example trial (bottom). Dexterity
896 is a recently developed test of fine motor function which consists of a circular coordinate plane
897 with the center of the circle (demarcated by a thin black line) at x,y positions 0,0. The goal is to
898 move the cursor (depicted as a colored ball) to a target located at various positions. Movement of
899 the cursor is affected by a “noise” component complementing the directional input from the
900 analog mouse to create the effect of interference or “jittering” motion. The effect is such that
901 successful navigation of the coordinate plane using the mouse encounters resistance to
902 purposeful direction, requiring continual adjustment by the participant to maintain the correct
903 path to the target. Visual feedback is given by the use of a color system, wherein the cursor shifts
904 gradually from green to red as proximity to the target becomes lesser. The task consists of 80
905 trials (10 per “noise” condition), ten seconds maximum in length, with preset noise factors
906 (ranging in intensity) and target locations standardized for consistency between participants. A
907 lack of input from the participant results in a gradual drift towards the center. At the conclusion
908 of each trial, the cursor location is reset to the origin. Completion time and Moves were recorded
909 with Moves defined as the change in the vector direction of the mouse while course correcting
910 toward the target. The radius of the circular coordinate plane is defined as a function of the
911 screen size and resolution as 300 arbitrary units, the cursor as 2.5 units, and the target as 12.5
912 units. Distance to the target is computed using those arbitrary units. Cursor velocity is defined
913 externally as “intermediate” in the Windows XP Service Pack 2 settings.



914



915

916

917

918

919

921 I. Test of Attentional Vigilance (TOAV), a measure of sustained attention, example trial. The
922 instructions were as follows:

923
924 On each trial, you will see one of two stimuli on the screen. Each will be a white square with a
925 black square inside it. On some trials, this inner square will be near the top of the white square;
926 on other trials it will be near the bottom. Press any key to see the stimuli.

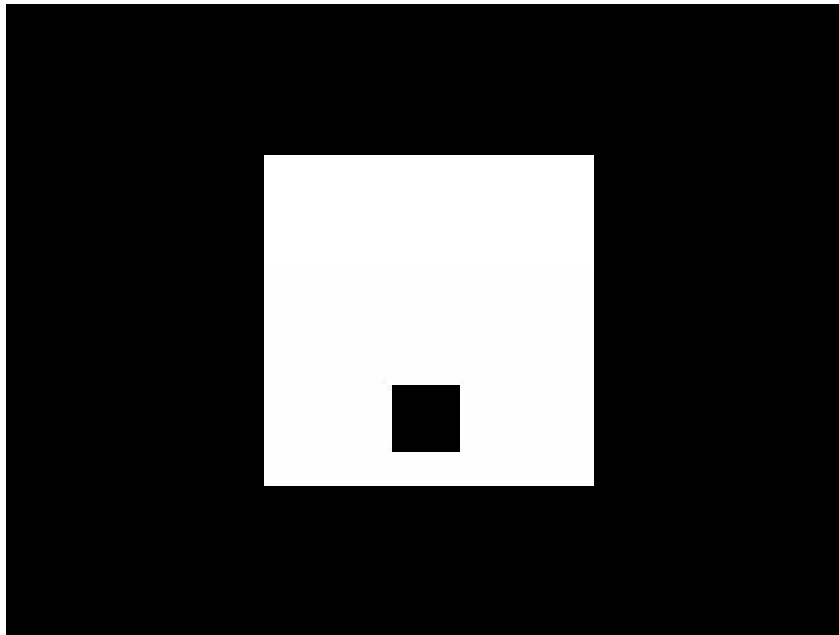
927
928 When the square is on the top, it is a target. During the task, you should press the space bar
929 whenever you see the target stimulus. Press any key to continue.

930
931 When the square is on the bottom, it is NOT a target. During the task, you should not press the
932 space bar when the non-target is displayed. Press any key to continue.

933
934 During the task, you will see a series of targets and non-targets. Press the space bar as quickly as
935 you can whenever you see a target (top square). Do nothing when you see a non-target (bottom
936 square). The task lasts approximately 6 minutes, so you need to concentrate on the task in order
937 to perform well. Press the space bar to begin.

938

939



940

941 **Table 1.** Performance on the Psychology Experimental Building Language (PEBL) battery including total time on target on the
 942 pursuit rotor (PR), Response Time (RT) and RT standard deviation (SD) on the Test of Attentional Vigilance (TOVA), B:A ratio on
 943 the Trail-Making Test (TMT), Tower of London (ToL), Perseverative Errors (PE) on the Wisconsin (Berg) Card Sort Test (BCST),
 944 and Mental Rotation Test (MRT).

945

	<u>Min (N)</u>	<u>Max (N)</u>	<u>Mean</u>	<u>SEM</u>	<u>N</u>
946					
947					
948 A. Pursuit Rotor: Time (sec)	18.8 (1)	56.3 (1)	44.0	0.5	189
949					
950 B. Pursuit Rotor: Error (pixels)	50.5 (1)	322.7 (1)	87.7	2.6	189
951					
952 C. Dexterity (sec)	956.9 (1)	7,276.8 (1)	1,619.4	59.2	175
953					
954 D. Time-Wall (% Inaccuracy)	3.0 (1)	53.0 (1)	10.2	0.5	171
955					
956 E. Test of Attentional Vigilance: RT (msec)	269 (1)	495 (1)	339.6	3.2	150
957					
958 F. Test of Attentional Vigilance: RT SD	42 (1)	288 (1)	100.3	2.6	150
959					
960 G. Digit Span (Points)	7 (7)	21 (3)	13.5	0.3	148
961					
962 H. Trail Making Test (B:A)	0.62 (1)	2.10 (1)	1.28	0.02	180

963						
964	I. Tower of London (Points)	12 (1)	36 (6)	29.0	0.3	182
965						
966	J. Tower of London (sec/trial)	4.8 (1)	32.0 (1)	14.3	0.4	182
967						
968	K. Iowa Gambling Test (\$)	-500 (1)	4,500 (1)	1894	54	184
969						
970	L. Berg Card Sorting Test (% PE Heaton)	3.1 (2)	65.6 (1)	11.0	0.5	173
971						
972	M. Berg Card Sorting Test (% PE Berg)	0.0 (2)	35.9 (1)	12.9	0.5	174
973						
974	N. Mental Rotation Test (% correct)	34.4 (1)	100.0 (1)	73.9	1.4	174
975						
976	O. Mental Rotation Test (msec)	420.4 (1)	5381.6 (1)	2,564.6	66.2	174
977						

979 **Table 2.** Spearman correlations between tests on the Psychology Experimental Building Language (PEBL) battery including Response
 980 Time (RT) and RT standard deviation (SD), Part B to Part A ratio on the Trail-Making Test; Perseverative Errors (PE) on the Berg
 981 Card Sorting Test coded according to the ^BBerg and ^HHeaton criteria. Correlations in **bold** are significant at $P \leq .05$, those in both **bold**
 982 **and italics** are significant at $P < .0005$.

	A.	B.	C.	D.	E.	F.	G.	H.	I.	J.	K.	L.	M.
987 A. Pursuit Rotor: Time (sec)	+1.00												
989 B. Pursuit Rotor: Error	-0.96	+1.00											
991 C. Dexterity (msec)	-0.26	+0.25	+1.00										
993 D. Time-Wall (Inaccuracy)	-0.32	+0.32	+0.13	+1.00									
995 E. TOAV: RT (msec)	-0.13	+0.14	+0.12	+0.13	+1.00								
997 F. TOAV: RT SD	-0.38	+0.38	+0.18	+0.25	+0.46	+1.00							
999 G. Digit Span	+0.14	-0.18	-0.10	-0.08	-0.24	-0.11	+1.00						
1001 H. Trail Making Test (B:A)	-0.16	+0.16	+0.06	+0.09	+0.12	-0.04	-0.01	+1.00					
1003 I. Tower of London (Points)	+0.15	-0.17	-0.13	-0.14	-0.14	-0.13	+0.00	-0.27	+1.00				
1005 J. Iowa Gambling Test	+0.01	-0.02	-0.09	-0.13	+0.16	+0.05	+0.02	-0.10	-0.00	+1.00			

1007	K. Berg Card Sorting Test (% PE ^H)	-0.21	+0.27	-0.03	+0.12	+0.15	+0.07	-0.12	+0.27	-0.34	-0.07	+1.00
1008												
1009	L. Berg Card Sorting Test (% PE ^B)	-0.18	+0.22	+0.02	+0.07	+0.06	+0.06	-0.02	+0.30	-0.19	-0.01	+0.72 1.00
1010												
1011	L. Mental Rotation Test (% correct)	+0.20	-0.20	-0.07	-0.04	-0.08	-0.19	+0.06	-0.08	+0.29	-0.04	-0.20 -0.18 +1.00
1012												
1013	M. Mental Rotation Test (msec)	-0.09	+0.05	-0.07	-0.01	+0.12	-0.05	+0.19	+0.10	+0.13	+0.11	+0.05 +0.03 +0.08
1014												
1015												
1016												

1018 **Table 3.** Comparison of computerized neurobehavioral batteries. Behavioral Assessment and Research System (BARS); Cambridge
 1019 Neuropsychological Test Automated Battery (CANTAB); Continuous Performance Test (CPT), Maximum (Max); Minimum (Min);
 1020 Psychology Experiment Building Language (PEBL); Test of Attentional Vigilance (TOAV).

	<u>BARS</u>	<u>CANTAB</u>	<u>PEBL</u>	
1024	Year developed	1994	1980s	2003
1026	Origins	behavior analysis & cognitive psychology	behavioral neuroscience	experimental & neuropsychology
1029	Philosophy	working populations with different educations & cultures	translational, cultural & language independent	collection of open-source neuropsychological measures
1033	Modifiable	no	no	yes
1035	Cost (Min/Max)	\$950 ^A /\$8,450 ^B per computer	\$1,275 ^C /\$24,480 ^D per computer	free/free for unlimited computers
1037	# Tests	11	25	>100
1039	Example measures	Finger Tapping Reaction time CPT Digit Span	Motor Screening Simple Reaction Time Match to Sample Spatial Span	Tapping Rotary pursuit TOAV, CPT Digit Span, Spatial Span

1043	Selective Attention	Choice Reaction Time	Dexterity
1044	Symbol Digit	Stockings of Cambridge	Tower of London

1045

1046 ^Aone-year preliminary data/student package with 9Button hardware (\$450), ^Bthree-year license with hardware, ^Cone-test with one-year
1047 license, ^Dall tests for 10 year license

1049 **Supplemental Table 1.** Selected peer-reviewed publications using the Psychology Experiment Building Language (PEBL) software.
 1050 Berg/Wisconsin Card Sorting Test (BCST); Delayed Match to Sample (DMS); Implicit Association Task (IAT); Iowa Gambling Task
 1051 (IGT); Psychomotor Vigilance Task (PVT); Situation Awareness Task (SAT); Time-Wall (TW); Tower of London (ToL); and Trail
 1052 Making Test (TMT).
 1053

1054	1055	1056	1057	1058	1059	1060	1061	1062	1063	1064	1065	1066	1067	1068	1069	1070	1071	1072
	<u>Topic</u>	<u>1st Author</u>	<u>Year</u>	<u>Citation</u>														
	Description of PEBL	Mueller	2010	<i>International J of Machine Consciousness</i> 2 : 273-288														
	Pursuit rotor in children	Piper	2010	<i>J Neuroscience Methods</i> 195 :88-91														
	Alcohol & decision making	Lyvers	2010	<i>Addictive Behaviors</i> 35 : 1021-1028														
	Anxiety & decision making	de Visser	2010	<i>Neuropsychologia</i> 48 :1598-1606														
	Caffeine & decision making	Aggarwal	2011	<i>British J Surgery</i> 98 : 1666-1672														
	Behavioral genetics of glutamate	Ness	2011	<i>Neuropharmacology</i> 61 :950-956														
	Heavy drinkers & decision making	Gullo	2011	<i>Drug & Alcohol Dependence</i> 117 :204-210														
	Behavioral genetics & amphetamine	Wardle	2012	<i>Genes, Brain & Behavior</i> 12 :13-20														
	Schizotypy & cognition	Cappe	2012	<i>Psychiatry Research</i> 200 :652-659														
	Brain damage & strategy updating	Danckert	2012	<i>Cerebral Cortex</i> 22 :2745-2760														
	Multiple sclerosis and cognition	Kalinowska	2012	<i>J of Neurological Sciences</i> 321 : 43-48														
	Executive function & lifespan	Piper	2012	<i>Behavior Research Methods</i> 44 : 110-123														
	Aging & executive function	Zebrowitz	2013	<i>Psychology & Aging</i> 28 : 202-212														
	Transcranial Infrared Laser	Barrett	2013	<i>Neuroscience</i> 230 : 13-23														

1073	Obsessive Compulsive Disorder	Tumkaya	2013	<i>Psychiatry Research</i> 209 : 579-588	SAT
1074	Wilson's disease & decisions	Ma	2013	<i>J Clin Exp Neuropsychol</i> 35 : 472-479	BCST
1075	Alcohol consumption & decisions	Bowley	2013	<i>Int J Psychophysiology</i> 89 : 342-348	IAT
1076	Essential tremor & cognition	Bhalsing	2014	<i>Eur J Neurology</i> 21 ; 874-883	BCST
1077	Tourette's & motor skill	Brandt	2014	<i>PLoS One</i> 9 : e98417	Pursuit Rotor
1078	Executive function & Parkinson's	Cohen	2014	<i>J Parkinson's Dis</i> 4 : 111-122	BCST
1079	Hoarders and cognition	Raines	2014	<i>J Affect Dis</i> 166 : 30-35	CPT
1080	Sex differences	Evans	2015	<i>Brain & Cogn</i> 93 : 42-53	IGT
1081	Decision making & alcohol	Lyvers	2015	<i>Addict Behav</i> 41 : 129-135	IGT

1082

1083 * For an updated list, see: http://pebl.sourceforge.net/wiki/index.php/Publications_citing_PEBL

PeerJ PrePrints

1085 **Supplemental Table 2.** Comparison of the antecedents to measures contained in the Psychology Experiment Building Language
 1086 (PEBL) battery including the originator(s), year of key publication, construct measured, and title for PEBL version. Berg Card Sorting
 1087 Test (BCST); Executive Function (EF); Iowa Gambling Task (IGT); Mental Rotation Test (MRT); Tower of London (ToL); Test
 1088 Attentional Vigilance (TOAV); Trail Making Test (TMT)
 1089

1090	<u>Test</u>	<u>Originator</u>	<u>Year</u>	<u>Construct</u>	<u>PEBL Version</u>
1091	Digit Span	Alfred Binet	1905	working memory	Digit Span
1092	Rotary Pursuit	Robert Ammons	1947	procedural learning	Pursuit Rotor
1093	Wisconsin Card Sorting Test	David Grant; Esta Berg	1948	EF: set shifting	BCST
1094	TMT	Ralph Reitan	1955	EF: divided attention	TMT
1095	MRT	Roger Shepard	1978	EF: decision making	MRT
1096	ToL	Tim Shallice	1983	EF: planning	ToL
1097	Test of Variables of Attention	Lawrence Greenberg	1993	sustained attention	TOAV
1098	IGT	Antoine Bechara	1994	decision making	IGT
1099					

1101 **Supplemental Table 3.** Mean performance and correlation between test sessions on Psychology Experiment Building Language
 1102 measures. The number of participants is listed in () after each test; ^rreported previously⁴⁷; Coded according to the ^oOringal Berg
 1103 sorting rules or the ^HHeaton rules; ^At-test $P \leq .01$ versus test; ^Bcorrelation $P < .01$.

	Test		Retest		%	Cohen's	Correlation		
	<u>Mean</u>	<u>SEM</u>	<u>Mean</u>	<u>SEM</u>	<u>Difference</u>	<u>d</u>	<u>Pearson r</u>	<u>rho</u>	
1108	Rotary Pursuit (76)								
1109	Total time (sec)	37.9	1.1	43.7 ^A	0.9	+15.4	.60	.86 ^B	.81 ^B
1110	Error (pixels)	31.7	3.3	24.7	2.7	-22.0	.24	.59 ^B	.77 ^B
1111	Trail Making Test (78)								
1112	Time A (sec) ^F	16.5	0.3	15.4 ^A	0.3	-6.6	.07	.74 ^B	.71 ^B
1113	Time B (sec) ^F	22.1	0.5	19.3 ^A	0.5	-12.8	.13	.61 ^B	.56 ^B
1114	Ratio (B/A)	1.35	0.03	1.25 ^A	0.02	-7.2	.07	.39 ^B	.35 ^B
1115	Digit Span Forward (72)								
1116	Test of Attentional Vigilance (68)								
1117	Response Time	385.4	4.8	395.2	7.7	+2.5	.25	.79 ^B	.72 ^B
1118	SD of Response Time	109.8	5.3	98.9 ^A	6.7	-9.9	.25	.87 ^B	.69 ^B
1119	Omission Errors	8.1	1.9	5.7	1.4	-30.2	.16	.88 ^B	.43 ^B
1120	Commission Errors	16.9	1.1	13.8	1.3	-18.3	.34	.65 ^B	.66 ^B
1121	Tower of London (66)								
1122	Moves	8.9	0.2	9.0	0.3	+1.3	.09	.15	.34 ^B
1123	Time (sec)	16.8	0.5	14.1 ^A	0.5	-16.0	.62	.36 ^B	.48 ^B
1124	Iowa Gambling Task (68)								
1125	Response pattern	8.8	2.9	18.2 ^A	3.5	+106.3	.39	.41 ^B	.22
1126	Money	1944.8	85.0	2162.1	116.0	+11.2	.31	.10	-.01
1127	Berg Card Sorting Test (73 ^O /60 ^H)								
1128	Categories Completed ^O	3.1	0.2	3.6 ^A	0.2	+16.0	.39	.51 ^B	.47 ^B
1129	Errors (%) ^O	30.0	1.7	22.0 ^A	1.6	-26.7	.55	.69 ^B	.68 ^B
1130	Perseverative responses ^O (%)	30.7	0.9	31.1	1.0	-1.3	.05	.03	.00
1131	Perseverative errors ^O (%)	14.9	0.8	12.2 ^A	0.9	-18.2	.41	.45 ^B	.35 ^B
1132									

1133	Perseverative responses ^H (%)	14.5	1.1	12.4	1.4	-24.9	.14	.51 ^B	.35 ^B
1134	Perseverative errors ^H (%)	13.0	0.9	11.2	1.0	-27.7	.14	.53 ^B	.34 ^B
1135	Time-Wall (74)	.175	.02	.070 ^A	.006	-60.0	.68	.39 ^B	.61 ^B
1136									
1137									
1138									
