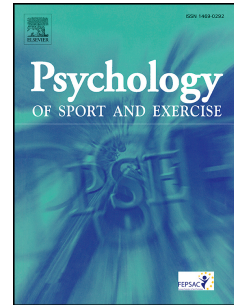


# Journal Pre-proof

Reliable measurement in sport psychology: The case of performance outcome measures

Geoffrey Schweizer, Philip Furley, Nicolas Rost, Kai Barth



PII: S1469-0292(19)30756-3

DOI: <https://doi.org/10.1016/j.psychsport.2020.101663>

Reference: PSYSPO 101663

To appear in: *Psychology of Sport & Exercise*

Received Date: 31 October 2019

Revised Date: 20 January 2020

Accepted Date: 2 February 2020

Please cite this article as: Schweizer, G., Furley, P., Rost, N., Barth, K., Reliable measurement in sport psychology: The case of performance outcome measures, *Psychology of Sport & Exercise* (2020), doi: <https://doi.org/10.1016/j.psychsport.2020.101663>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier Ltd.

Reliable Measurement in Sport Psychology: The Case of Performance Outcome Measures

Geoffrey Schweizer<sup>1</sup>, Philip Furley<sup>2</sup>, Nicolas Rost<sup>3</sup> & Kai Barth<sup>1</sup>

<sup>1</sup>Department of Sport Psychology, Heidelberg University, Heidelberg, Germany

<sup>2</sup>Institute of Cognitive and Team/Racket Sport Research, German Sport University  
Cologne

<sup>3</sup>Max Planck Institute of Psychiatry, Department of Translational Research in Psychiatry,  
Munich, Germany; International Max Planck Research School for Translational  
Psychiatry, Munich, Germany

Author Note

Geoffrey Schweizer, Heidelberg University, Department of Sport Psychology, INF  
720, 69120 Heidelberg, Germany, Phone: +49 (0) 6221 546033, Email:  
[geoffrey.schweizer@issw.uni-heidelberg.de](mailto:geoffrey.schweizer@issw.uni-heidelberg.de)

Philip Furley, German Sport University Cologne, Institute of Cognitive and  
Team/Racket Sport Research, Am Sportpark Müngersdorf 6, 50933 Köln, Germany,  
Email: [p.furley@dshs-koeln.de](mailto:p.furley@dshs-koeln.de)

Nicolas Rost, Max Planck-Institute of Psychiatry, Department of Translational  
Research in Psychiatry, Kraepelinstr. 2-10, 80804 München, Germany; International Max  
Planck Research School for Translational Psychiatry, Munich, Germany, Email:  
[nicolas\\_rost@psych.mpg.de](mailto:nicolas_rost@psych.mpg.de)

Kai Barth, Email: [kai-barth@web.de](mailto:kai-barth@web.de)

Correspondence concerning this article should be addressed to Geoffrey Schweizer,  
Heidelberg University, Department of Sport Psychology, INF 720, 69120 Heidelberg,

Germany, Phone: +49 (0) 6221 546033, Email: [geoffrey.schweizer@issw.uni-heidelberg.de](mailto:geoffrey.schweizer@issw.uni-heidelberg.de)

Declarations of interest: none.

Journal Pre-proof

1

2

3

4

5

6     **Reliable Measurement in Sport Psychology: The Case of Performance Outcome Measures**

7

8

9

Journal Pre-proof

10 Abstract

11 *Objectives:* The present research addresses a neglected aspect within the current Zeitgeist of  
12 improving methodological standards in (sport)psychology: reliable measurement. We discuss  
13 and highlight the importance of reliable measurement from different perspectives and  
14 empirically assess reliability of three commonly used performance outcome measures in order to  
15 give guidelines to researchers on how to increase reliability of measurements of performance  
16 outcomes.

17 *Method:* In three studies we estimate 5 different reliability coefficients for three performance  
18 outcome measures based on 14 golf putts (study 1;  $N = 100$ ), 14 dart throws (study 2;  $N = 200$ ;  
19 100 sports students; 100 non-sports students) and 14 free throws in basketball (study 3;  $N = 192$ ;  
20 100 non-basketball players; 92 basketball players).

21 *Results:* The highest reliability was the odd-even reliability for darts for the whole sample (.888),  
22 followed by golf putts (.714 for distance from the hole, .614 for successful putts) and free throws  
23 (.504 non-basketball players; .62 for basketball players; and .826 for whole sample).

24 *Conclusions:* Based on theoretical considerations and our empirical findings we give practical  
25 guidelines to improve reliability for performance outcome measures in sport psychology.

26 183 words

27 *Keywords:* Classical Test Theory, replicability, research quality, golf putts, darts,  
28 basketball

29

30

31

32       Reliable Measurement in Sport Psychology: The Case of Performance Outcome Measures  
33       “For my money, the #1 neglected topic in statistics is measurement” (Gelman, 2015).  
34       In the past ten years, there has been a controversial discussion regarding the quality of  
35       psychological research (e.g., Nelson, Simmons, & Simonsohn, 2018). Many authors have  
36       criticized what they perceive as problematic research practices, that lead to low rates of  
37       replications and weaker empirical support for psychological theories and interventions than  
38       researchers themselves may believe. Some authors even announced a “crisis of confidence”  
39       (Pashler & Wagenmakers, 2012, p. 528). Currently, the aforementioned discussion has  
40       developed from pointing out existing problems to creating viable practices for strengthening the  
41       quality of empirical research in the future (Lakens & Evers, 2014; Munafò et al., 2017; Nelson et  
42       al., 2018). For example, future research is supposed to benefit from appropriate sample sizes  
43       (Button et al., 2013; Fraley & Vazire, 2014; Schönbrodt & Perugini, 2013; Schweizer & Furley,  
44       2016), methodological advances and new software<sup>1</sup> (e.g., Wagenmakers et al., 2018; Love et al.,  
45       2019), replicability projects (Camerer et al., 2018; OSC, 2015; Soto, 2019), the avoidance of p-  
46       hacking (e.g., Simmons, Nelson, & Simonsohn, 2011), the opportunity to preregister plans for  
47       studies and data analyses (e.g., Nosek, Ebersole, DeHaven, & Mellor, 2018), transparency (e.g.,  
48       Nosek et al., 2015), and, maybe most importantly, a heightened awareness for the importance of  
49       these and other related issues. In the light of these changes, some authors have even gone as far  
50       as to proclaim a “renaissance” of psychological research (Nelson et al., 2018, p. 511).

51       However, as indicated by the quote above, with few exemptions, measurement has long  
52       been absent from the current discussion on research quality<sup>2</sup>. This is particularly remarkable as  
53       measurement is strongly connected to many of the issues that have been debated, as we will  
54       show below. In the present publication, we will focus on one particularly important aspect of

55 measurement, namely reliability. Below, we will outline why we focus on reliability instead of  
56 other aspects of measurement. The main goal of the present paper is twofold: First, we aim to  
57 alert readers to the importance of reliable measurement. We illustrate why reliability is crucial  
58 for high-quality research in general, and particularly so for sport psychology, as we will argue.  
59 Second, we aim to assess reliability estimates for three performance outcome measures (golf  
60 putts, dart throws, and basketball free throws) commonly used in sport psychological research.  
61 Based on these estimates, we make recommendations for future research in sport psychology.

### 62 **What is Reliability?**

63 In Classical Test Theory (CTT), every observed value ( $Y_{\text{observed}}$ ) consists of a true value  
64 ( $T_{\text{true}}$ ) and measurement error ( $Y_{\text{error}}$ ). Measurement error ( $Y_{\text{error}}$ ) is defined as random in this  
65 context (e.g., Bühner, 2011; Lord & Novick, 1968; Steyer & Eid, 1993; Vaughn, Lee, & Kamata,  
66 2012).

67 Equation 1

$$68 \quad Y_{\text{observed}} = T_{\text{true}} + Y_{\text{error}}$$

69 This means that every measurement consists of the respective person's true value and a  
70 random error which is due to the measurement's imperfectness. The smaller a measurement  
71 procedure's error, the closer each single measurement will be on average to the true value. As  
72 measurement error is defined as random, it must have an expectancy value of 0. Furthermore, it  
73 is defined as having a finite variance. Measurement error is usually supposed to be normally  
74 distributed. This means that averaged over a large number of individual measurements the mean  
75 of the measurement error is supposed to be 0. In other words, larger errors are supposed to be  
76 less likely than smaller errors, and errors cancel each other out. With only enough measurements,  
77 therefore, the mean of the observed values will equal the mean of the true values. The more

78 measurement error, the larger this number of measurements has to be (see Appendix A for more  
79 details). However, for every single measurement we cannot say how large its individual error  
80 ( $Y_{\text{error}}$ ) is. In this context measurement error refers to the random error associated to every single  
81 measurement due to the measurement's imperfection. It does not refer to sampling error (i.e.,  
82 differences between samples due to interindividual variation of the true values) or intraindividual  
83 variation (i.e., differences between measurement points due to intraindividual variation of the  
84 true values).

85 Reliability can be understood as the inverse of measurement error. In other words, the  
86 less error a measurement contains, the higher its reliability and vice versa. Formally, reliability is  
87 defined as the proportion of true variation among the entire variation:

88 Equation 2

$$89 \quad r_{yy} = \frac{\sigma_{y_{\text{true}}}^2}{\sigma_{y_{\text{true}}}^2 + \sigma_{y_{\text{error}}}^2}$$

90 Reliability can be estimated via several approaches (e.g., test-retest-reliability; split-half  
91 reliability; parallel-test reliability; Cronbach's alpha; see Appendix A for more information on  
92 different coefficients).

### 93 **Why Does Reliability Matter?**

94 "Measurement error adds noise to predictions, increases uncertainty in parameter  
95 estimates, and makes it more difficult to discover new phenomena or to distinguish among  
96 competing theories" (Loken & Gelman, 2017, p. 584).

97 The question why reliability matters could be answered in one sentence: The higher a  
98 measurement's reliability (i.e., the less measurement error contained on average in a single  
99 observation), the more precise is every single measurement on average. However, it is possible  
100 to look at this feature from several perspectives, and different researchers may be more familiar



101 with some of these perspectives than with others. Therefore, we will explore the benefits of  
102 reliable measurement in more detail below. Importantly, the following arguments are not  
103 independent of each other, but they are all a direct consequence of the above mentioned  
104 observation.

105 **Reliability matters for Type-1 and Type-2 errors.**

106 Both Type-1 errors (false positives) and Type-2 errors (false negatives) can be found in  
107 psychological research. Researchers commit a Type-1 error (or false positive) when they report  
108 finding an effect when truly there is none; or when refuting a null hypothesis when in truth the  
109 null hypothesis is correct. Conversely, researchers commit a Type-2 error when they report not  
110 finding an effect when truly there is one; or when retaining a null hypothesis when in truth the  
111 null hypothesis is false (Fraley & Vazire, 2014). Reliability plays a role for both these errors  
112 (Loken & Gelman, 2017). Type-1 and Type-2 errors are often presented in the framework of  
113 Null-Hypothesis-Significance-Testing (NHST)<sup>3</sup>, however, they are relevant for all statistical  
114 perspectives that contain binary decisions such as accept-or-reject, present-or-absent, retain-or-  
115 dismiss.

116 It has been understood for more than a century that measurement error attenuates effect  
117 sizes (in this case correlations; Spearman, 1904, 1910). This means that when measuring an  
118 effect with less than perfect reliability, the estimated effect size will be smaller than the real-  
119 world effect size. The lower the reliability, the smaller the estimated effect size will be compared  
120 to its real size. For example, when researchers estimate a correlation coefficient between two  
121 measures, then the estimated coefficient will be smaller than the true coefficient to the extent that  
122 one or both of the measurements are less than perfectly reliable. As larger effects are more easily  
123 detected, measurement error thus increases the likelihood of committing a Type-2 error (or false

124 negative): When planning for power, researchers usually plan their sample size with the goal of  
125 achieving a certain power (e.g., .9) based on an expected real-world effect size (e.g.,  $r_{\text{expected}} =$   
126 .7). The less reliable their measurement, the smaller the effect size they estimate in their study  
127 gets (e.g.,  $r_{\text{estimated}} = .5$ ). However, the likelihood of finding this smaller effect size is lower than  
128 the likelihood of finding the real-world (and larger) effect size. Thus, the researchers in this  
129 example risk committing a Type-2 error: They do not find the effect although it truly exists.

130 In order to obtain the same power to detect an effect of the same (true) size, researchers  
131 need smaller samples when using reliable measures than when using unreliable measures.  
132 Therefore, research using more reliable measures is, all else equal, more economical than  
133 research with less reliable measures.

134 Only recently have researchers suggested that in addition to deflating coefficient  
135 estimates, measurement error can also inflate coefficient estimates (Loken & Gelman, 2017).  
136 Whether measurement error deflates or inflates coefficient estimates depends primarily on  
137 sample sizes: In large samples, measurement error nearly always deflates coefficient estimates.  
138 Here, measurement error primarily leads to Type-2 errors. In small samples, measurement error  
139 can deflate and inflate estimates. Therefore, in small samples, measurement error can lead to  
140 Type-1 and Type-2 errors. When researchers are more likely to publish larger effect sizes than  
141 smaller ones (e.g., because they are more likely to be statistically significant or because they  
142 seem to be more impressive), measurement error is likely to contribute to the potentially high  
143 proportion of false-positive findings that have been diagnosed for psychological research,  
144 because it inflates effect sizes in small samples (Loken & Gelman, 2017).

145 Taken together, all else being equal, a research field will benefit from more reliable  
146 measurement, due to less Type-1 and less Type-2 errors. Or stated differently, the field with

147 more reliable measures is both less likely to miss potentially important effects and it is less likely  
148 to report an effect that in reality does not exist.

149 **Reliability matters for replications.**

150 In past replication projects, several definitions and operationalizations of successful or  
151 unsuccessful replication attempts have been employed (Camerer et al., 2018; OSC, 2015; Soto,  
152 2019). What exactly constitutes a successful or unsuccessful replication remains a matter of  
153 debate (for an enlightening discussion and one possible solution see Simonsohn, 2015). Broadly,  
154 one can distinguish two strategies for defining a successful replication. The first strategy  
155 considers a replication as successful when both the original study and the replication study have  
156 significant results in the same direction. The second strategy compares the effect sizes of the  
157 original study and the replication study with each other. An effect is considered to be  
158 successfully replicated when the effect sizes produced by the original and the replication study  
159 do not differ.

160 For the first strategy, reliability matters because it affects the studies' power. The  
161 probability of replicating a true effect depends, among other factors, on the reliability of the  
162 measure with which the to-be-replicated effect is assessed (Stanley & Spence, 2014). The lower  
163 the reliability of a measure, the lower the probability of replicating an effect, even when the  
164 effect is true. Therefore, when measurement is unreliable, unsuccessful replication attempts may  
165 both be a consequence of an effect not being true or of an effect being measured with much  
166 error. For example, Soto finds that the successful replication rate in a large scale replication  
167 project from personality psychology was substantially higher than the respective rates in large  
168 scale replication projects from other behavioral sciences (2019). He goes on to speculate that one  
169 reason for this striking discrepancy (among others, such as sample sizes and type of focal effects)

170 might be that personality psychological research uses more standardized and thus more reliable  
171 measurements than, for example, social psychological research.

172 For the second strategy, reliability matters because, as explained above, measurement  
173 error may both inflate and deflate estimates of effect sizes, making it less informative to compare  
174 two effect sizes to each other. As replication attempts are becoming ever more important in  
175 science, so does the role of reliable measurement, as without reliable measurement replication  
176 attempts are at least hard to interpret or at worst futile.

177 **Reliability matters for the impact of p-hacking.**

178 P-hacking refers to the practice of “selectively reporting data and analyses” or, in other  
179 words, “conducting multiple analyses on the same data set and then reporting only the one(s) that  
180 obtained statistical significance” (Nelson et al., 2018, p. 513). When researchers employ p-  
181 hacking, the likelihood of obtaining a false-positive increases “dramatically” beyond the level  
182 usually assumed by researchers (Nelson et al., 2018, p. 513). Typical examples of p-hacking<sup>4</sup>  
183 include a) having two correlated dependent variables and selectively reporting one of them, b)  
184 adding observations to the sample and stopping once statistical significance has been reached, c)  
185 deciding whether to drop one out of several experimental conditions, d) selectively controlling  
186 for gender or for the interaction of gender with treatment, or e) combinations thereof (Simmons  
187 et al., 2011). In the simulations run by Simmons and colleagues, p-hacking could lead to a  
188 likelihood of obtaining a false-positive (i.e., finding a significant result when truly there is no  
189 effect) of up to 61% (Nelson et al., 2018; Simmons et al., 2011). As a result of p-hacking, there  
190 is supposed to be a high proportion of false and therefore non-replicable findings in the  
191 psychological literature. Nelson and colleagues (2018) suggest that p-hacking (i.e., analyzing the  
192 data until researchers find a significant result) and not publication bias (i.e., simply not

193 publishing non-significant results) is the real answer to the decades-old question how  
194 psychologists manage to publish such a high proportion of significant results when their studies  
195 typically have rather low power. In a recent paper, Friese and Frankenbach (in press) suggest that  
196 p-hacking and publication bias interact: In their simulation study, the extent to which p-hacking  
197 distorts meta-analytic effect size estimates depends on the level of publication bias and on true  
198 effect sizes.

199         Although the problematic influence of p-hacking on research quality has been described  
200 in detail (Simmons et al., 2011), it remains an open question how to reduce its impact in future  
201 research. Again, reliability plays a role as p-hacking exploits random variation. For example, the  
202 strategy of successively increasing the sample size until a certain difference becomes (randomly)  
203 significant works “best” when this very difference is subject to lots of random variation. As  
204 random variation increases with measurement error, so do the opportunities to employ p-hacking.  
205 It follows then that highly reliable measurements should be less vulnerable to p-hacking.  
206 Therefore, one (rather indirect) method among others to reduce the impact of p-hacking would  
207 be to employ highly reliable measurements. Obviously, this would not entirely rule out the  
208 possibility of p-hacking, but it would at least to some extent decrease the potential for employing  
209 them.

#### 210         **Reliability matters for comparisons between measurements.**

211         Whereas the abovementioned benefits of high reliability follow directly from its  
212 definition, namely a measurement with comparably low error, the reliability of measures is  
213 consequential beyond simple estimations of the parameter of interest. Reliability is particularly  
214 important for comparisons between measurements. This can be comparisons between different

215 studies, comparisons between different measurements in a single study or comparisons of some  
216 effect on different measurements in a single study.

217         Suppose some researchers are interested in the question whether some treatment has a  
218 (differential) effect on two different variables. They find several studies reporting an effect on  
219 variable A, and several studies not finding an affect on variable B. They conclude that the  
220 treatment works for variable A, but nor for variable B. However, and unfortunately so, in this  
221 example variable A was assessed with a more reliable measure than variable B. Therefore, the  
222 observed difference might simply be due to measurement error. This hypothetical scenario gets  
223 worse when we assume that preferred measurement and theoretical background of researchers  
224 might be correlated. In this case, different theories might appear to be differentially supported by  
225 evidence, while the only real difference is measurement error. These observations also hold  
226 when researchers compare effects on different variables within one study<sup>5</sup>. For example,  
227 researchers might conclude that their treatment affects variable A (e.g., some symptom of a  
228 disease), but not variable B (e.g., an unwanted side effect). Again, this conclusion is only  
229 legitimate when both variables are measured with the same high reliability.

230         Finally, this also holds for all kinds of multiple regression strategies and related attempts  
231 to control for one variable when estimating associations between two or more additional  
232 variables (Westfall & Yarkoni, 2016). For example, when one predictor significantly predicts the  
233 outcome variable while controlling for another predictor, researchers often interpret this finding  
234 as indicative of incremental validity, which itself may be interpreted as signifying that both  
235 predictors measure “strongly related but conceptually distinct constructs” (Westfall & Yarkoni,  
236 2016, e0152719). However, as Westfall and Yarkoni show, “... a simpler interpretation that is  
237 often equally consistent with the data is that both predictors are simply noisy indicators of the

238 same construct” (2016, e0152719). Westfall and Yarkoni conclude that reliable measurement is  
239 particularly important when trying to assess incremental validity in regression models.

#### 240 **Why reliability matters: A summary.**

241 Taken together, increasing reliability should lead to both less Type-1 and less Type-2  
242 errors, a higher chance of replicating an effect given it is true as well as making replication  
243 attempts more informative in general. Additionally, when measurements are more reliable,  
244 smaller sample sizes are needed in order to safeguard against statistical errors, p-hacking, and  
245 biases. Likewise, comparing between measurements is easier when both measurements are  
246 reliable. Taken together, it seems safe to conclude that research without reliable measurement  
247 does not make much sense in general, and particularly it does not make much sense in the age of  
248 replicability.

#### 249 **Reliability in Sport Psychology**

250 As we have outlined above there are good reasons to make efforts to increase reliability  
251 in science. Nevertheless, it is important to note that not every field of investigation or every  
252 measurement tool faces comparable challenges when it comes to both validity and reliability. In  
253 terms of validity it seems clear that certain psychological measures (e.g., a questionnaire  
254 measuring a person’s tendency to behave aggressively) struggle with more problems regarding  
255 validity as physical measures, for example assessing a person’s weight or body size, which has  
256 led researchers to speak of a validation crisis within the field of psychology (Schimmack, 2010;  
257 2019).

258 While questionnaires can be both problematic in terms of validity and reliability, other  
259 measurement techniques have high face-validity (i.e., there is little doubt as to whether they  
260 actually measure what they claim to measure). For example, most people would probably agree

261 that measuring a person's performance shooting basketball free throws is a valid measure of this  
262 person's ability to shoot free throws. However, it is less clear how reliable this measure is, or  
263 what has to be taken into account when reliably trying to assess perceptual-motor performance in  
264 answering different research questions in sport psychology. Hence, the present research focused  
265 on reliability of commonly used individual sport performance outcome measures.

266         Discussions of reliability have not been absent within sport science (Hopkins, 2000;  
267 2017; Zhu, 2013), and reliability has also been the focus of increasing research endeavors in  
268 some subfields of sport science, for example in determining both the validity and reliability of  
269 new technologies like GPS (Global Positioning System) in assessing sport performance data  
270 (Barbero-Álvarez, Coutts, Granda, Barbero-Álvarez, & Castagna, 2010; Coutts & Duffield,  
271 2010; Jennings, Cormack, Coutts, Boyd, & Aughey, 2010; Johnston, Watsford, Pine, Spurrs,  
272 Murphy, & Pruyn, 2012; Petersen, Pyne, Portus, & Dawson, 2009). In sport and exercise  
273 psychology, Eklund, Tenenbaum and Kamata (2012) provide an extensive overview about nearly  
274 all potentially important aspects of measurement in sport and exercise psychology, from basic  
275 concepts to specific issues, such as cognitive, motivational, emotional and behavioral  
276 measurement. These discussions and analyses have shown that reliable measurement of behavior  
277 in sports, although these measures appear high in face validity, is not a trivial topic. While  
278 measurement is an important topic at all levels of analyses within sports (e.g., biochemical  
279 measures, physiological measures, biomechanical measures, psychological questionnaires,  
280 anthropometric measures, behavioral measures, etc.), some of the most relevant measures (at  
281 least in terms of spectator interest or financial reward) are outcome measures of sports  
282 performance.



283           Somewhat surprisingly, we are not aware of any literature systematically analyzing the  
284 reliability of typically used sport performance outcome measures, although plenty of research  
285 uses sport performance measures as dependent variables. Maybe claims like “sport measures  
286 outcome with a finality of judgment that scientific papers would not pass” (Walsh, 2014, p. 860)  
287 have led researchers to simply assume sport outcome measures are reliable without needing to  
288 pay special attention to this. To address this shortcoming in the literature, we decided to first  
289 identify the most commonly used outcome measures of skilled perceptual-motor performance in  
290 sport psychology and subsequently calculate different reliability indices of these measures in a  
291 series of empirical studies.

292           A literature search identified 40 papers using golf putts as a dependent variable, 37  
293 papers using darts, and 28 using free throws in basketball (see the reference list in the  
294 supplement for an overview). Therefore, it seems safe to argue that these are frequently  
295 employed individual sport performance outcome measures in sport psychology<sup>6</sup>. Reliabilities  
296 were not reported in any of these papers<sup>7</sup>. It is important to note that we are not pointing out or  
297 criticising these papers. We ourselves have not reported reliability coefficients in most of our  
298 papers, when employing other measures than questionnaires.

299           However, this is precisely our point: Whereas everybody cares about reliable  
300 measurement when reporting questionnaire data, hardly anybody does when reporting  
301 performance outcome measures. There are probably several reasons for this. One, whereas it may  
302 seem rather straightforward how to compute reliabilities for questionnaires (e.g., most examples  
303 from the methodological literature refer to questionnaires), it may seem to be more unclear how  
304 to compute reliabilities for performance outcome data. Second, there is a common perception  
305 that a measure must have been reliable when there has been a significant result for this variable

306 (Loken & Gelman, 2017). Therefore, it may seem unnecessary to examine its reliability. As  
307 explained above, this is problematic for several reasons (see Loken & Gelman, 2017, for more  
308 details on this misconception).

### 309 **The Present Research**

310 The main goal of the present research was to estimate reliability coefficients for three  
311 commonly used individual performance outcome measures in sport psychology, namely golf  
312 putts (study 1), darts (study 2) and free throws in basketball (study 3). Furthermore, we aimed to  
313 investigate whether these reliability coefficients are dependent upon different samples in general  
314 and upon participants' experience with the respective task in particular.

### 315 **General Method**

316 Here, we describe the rationale common to all three studies. In the section below we  
317 describe characteristics unique to each study. In all studies, participants provided informed  
318 consent before commencing the study and were thanked and debriefed before receiving some  
319 candy as compensation for participating. Participants were neither paid, nor were they  
320 incentivized dependent on their performance. In all studies, participants first performed 20  
321 training trials, before executing 14 test trials. We only estimated reliability coefficients for the 14  
322 test trials, not for the training trials. The training trials were intended to reduce the influence of  
323 potential short-term learning effects on the reliability estimates (Hopkins, 2000), to get  
324 participants calibrated to the performance context (Ajemian, D'Ausilio, Moorman, & Bizzi,  
325 2010; Wunderlich, Heurer, Furley, & Memmert, 2019), and in turn decrease measurement error.  
326 We assumed potential learning and calibration effects to reach an asymptotical level after the  
327 learning trials and therefore performance to be stable for the test trials. In all studies, the setup  
328 was highly standardized, including videos demonstrating the correct execution of the required

329 movements to all participants (regardless of their familiarity with the task). Participants were told  
330 to try to achieve optimal performance, but experimenters emphasized that we would not evaluate  
331 individual performance in order not to induce pressure (e.g., Baumeister & Showers, 1986).

332 For all variables, we estimated their *split-half reliability* using two different methods, one  
333 method splitting each test in a first half and a second half and the odd-even method (please see  
334 Appendix A for further elaboration on our statistical approach). When using the “first-half vs.  
335 second half” method, for every participant, we computed one mean across the first seven test  
336 trials (i.e., trials 1-7) and one mean across the second seven test trials (i.e., trials 8-14). We then  
337 computed Pearson’s correlation coefficient for the correlation between the first and the second  
338 mean. When using the odd-even method, for every participant, we computed one mean across  
339 the seven odd-numbered test trials (i.e., trials 1, 3, 5, 7, 9, 11, 13) and one mean across the seven  
340 even-numbered test trials (i.e., trials 2, 4, 6, 8, 10, 12, 14). We then computed Pearson’s  
341 correlation coefficient for the correlation between the first and the second mean. According to  
342 Classical Test Theory (CTT), the resulting correlations can be considered one estimate for the  
343 respective measures’ reliability. Reliability depends on the number of items, and split-half  
344 reliabilities thus estimate the reliability for a test of half its original length (i.e., in our case, for  
345 seven instead of 14 trials). Therefore, we used the Spearman-Brown formula to estimate the  
346 reliability of all 14 trials, based on the obtained reliability coefficients (please see Appendix A  
347 for the formula). Additionally, we computed *Cronbach’s Alpha* for all measures.

348 In all studies, we used performance measures (and respective instructions) that allowed  
349 for obtaining continuous measurements, which is a prerequisite for estimating reliability  
350 according to CTT (please see Appendix A for further information). Furthermore, where

351 necessary, we aimed at minimizing the number of missing values, that is attempts that we could  
352 not measure.

353 We planned to total 200 participants per study (please see Appendix A for further  
354 elaboration on sample size planning). In all studies, we planned to collect data in two subsamples  
355 that differ regarding sports experience, each subsample totalling 100 participants. Below, we  
356 describe all studies in detail in order to facilitate interpretation and replication. We encourage  
357 researchers to contact us for more details. When we do not refer to a particular reliability  
358 coefficient, we always refer to the odd-even reliability, as it is usually considered superior to the  
359 “first-half vs. second half” reliability. In all studies and in all subsamples, there is no significant  
360 difference between the mean value for the odd and for the even items, and neither do standard  
361 deviations differ, which is considered a prerequisite for estimating odd-even reliability.

### 362 **Study 1: Golf Putts**

#### 363 **Participants**

364 One hundred students of Heidelberg University participated in the study (58 men and 42  
365 women;  $M_{\text{age}} = 24.9$ ;  $SD_{\text{age}} = 7.9$ ). Sixty-six of them were sports students, 34 were not. None of  
366 them reported having experience playing golf that went beyond participating in one basic course.  
367 Contrary to our plans, we did not collect data from a second subsample (for an explanation why,  
368 please see the Discussion section below).

#### 369 **Apparatus and Procedure**

370 Participants were positioned 200 cm away from the hole and instructed to assume a  
371 typical putting position. They could choose between a putter for left-handers and for right-  
372 handers. In order to maximize standardization, the study was conducted in a laboratory room.

373 Therefore, participants did not perform on a real green but on a putting mat made of plastic, as is  
374 common in sport psychological studies (see Supplement for several examples).

375 Participants were instructed to aim for the hole and informed that performance would be  
376 measured as distance from the hole. This allowed us to measure performance in a continuous  
377 way. Simply counting successful putts is also a commonly used measure and we therefore also  
378 estimated reliabilities for number of successful putts. This allowed us to compare the reliabilities  
379 of two different performance outcome measures constructed from the same task. In order to  
380 obtain continuous measurement, we added up all attempts, as sums of binary variables can be  
381 treated as continuous variables (e.g., Lunney, 1970).

## 382 **Results and Discussion**

383 In the whole sample, the average distance to the hole across all 14 putts was 257 mm (*SD*  
384 = 175; *Md* = 240; *Mode* = 241). On average, 53% of all putts were successful. When looking at  
385 the mean performances for each of the 14 putts separately, no learning trend was apparent (see  
386 Figure 1). Due to failures in data recording, in total the results of 15 putts out of 1400 were not  
387 recorded. Results do not change when excluding the respective participants.

388 For the continuous performance outcome measure (i.e., distance from the hole), the  
389 different reliability coefficients do not differ from each other (see Table 1). Thus, results are not  
390 dependent on a particular coefficient. Reliability coefficients can not be considered acceptable  
391 for seven putts only. When estimating reliabilities for all 14 putts using the Spearman-Brown-  
392 Formula, reliability coefficients expectably get higher, but they are still lower than what is  
393 usually considered acceptable (e.g., Vaughn et al., 2012). The Spearman-Brown corrected odd-  
394 even reliability for the whole sample is .714 ( $CI_{95} [.602; .798]$ ).

395 For the binary performance outcome measure (i.e., number of successful putts), the  
396 different reliability coefficients do not differ from each other (see Table 1). None of them  
397 reaches a level commonly considered as acceptable, with the highest estimate for all trials being  
398 .614 (CI<sub>95</sub> [.475; .723]). Thus, at least descriptively, reliability estimates for the binary  
399 performance outcome measure are lower than for the continuous one. However, confidence  
400 intervals overlap. Based on our data and our sample, golf putts as conducted in the present study  
401 did not possess sufficient reliability to be employed as a performance outcome measure in a sport  
402 psychological study. Reliability estimates for distances from the hole are somewhat better, but  
403 they still do not reach levels usually considered acceptable for other psychological measurement  
404 procedures. Based on the reliability estimates for the number of putts in the current study, the  
405 Spearman-Brown formula allows to calculate what number of putts would be necessary in order  
406 to achieve a certain level of reliability (e.g., .8 or .9). We present these calculations in the section  
407 “*comparison between performance outcome measures*”.

408 Putting distances in sport psychological studies vary from 100 to 400 cm, with 200 cm  
409 being common (see Supplement). Therefore, we decided to use 200 cm in our study. However,  
410 we realized that not only putting performance, but also reliability in golf putts probably strongly  
411 depends on the distance to the hole (as reliability of the performance measurement depends on  
412 the true score of the performance, which probably varies with distance). At the same time, there  
413 is no standard putting distance. Therefore, we decided not to conduct a second study with  
414 another 100 participants, as initially planned, but instead to move on to another performance  
415 measure. Our next two performance measures (darts and free throws) feature standard distances.  
416 Therefore, the problem described above does not apply to them.

417

**Study 2: Darts**

**418 Participants**

419 Study 2 consisted of two subsamples (sample 2a and sample 2b). Both samples consisted  
420 of 100 participants, totalling 200 participants. In sample 2a, there were 50 women and 50 men (4  
421 left-handed). All of them were students at Heidelberg University. In sample 2b, there were 44  
422 women and 56 men (10 left-handed). All of them were sports students at Heidelberg University.  
423 Thus, in the whole sample there were 94 women and 106 men (14 left-handed), 100 non-sports  
424 students and 100 sports students.

**425 Apparatus and Procedure**

426 In line with the World Darts Association's standards, we placed the dart board in such a  
427 way that the centre of the bull (also called bullseye) was at 173 cm (5 ft 8 inches) above ground.  
428 The diameter of the dart board was 400 mm. Participants were positioned behind a line (the so-  
429 called oche) that was 237 cm (7 ft 9.25 inches) away from the board. Behind the dart board we  
430 placed a board made of rigid foam (size: 120 cm x 120 cm; thickness: 40 mm). This setup  
431 allowed us to measure throws that missed the dart board but got stuck in the foam board, in order  
432 to minimize missing values. When a throw did not reach the board, participants were allowed to  
433 repeat the attempt (however, this happened hardly ever, due to the size of the foam board). We  
434 utilized regular tournament darts with a length of circa 155mm and a weight of circa 18g.  
435 Tournament darts come in different variants. Our darts consisted of a steel point, a brass barrel,  
436 an aluminium shaft and a standard shape flight.

437 Participants were instructed to aim for the bull and informed that performance would be  
438 measured as distance from the bull. This allowed us to measure performance in a continuous way  
439 (i.e., to obtain interval scale data), which is a prerequisite for estimating reliability. The usual  
440 scoring system in darts, however, would probably not have produced continuous measurements

441 (Selkirk, 1976; Tibshirani, Price, & Taylor, 2011). Aiming for the bull is a common instruction  
442 in sport psychological studies (see Supplement). In line with our measurement and the  
443 instructions, we did not utilize a dartboard with radial sections and double and triple rings.  
444 Instead, we used a dartboard with concentric rings of equal width.

#### 445 **Results and Discussion**

446 In the whole sample, the average distance to the bull across all 14 darts was 87 mm ( $SD =$   
447  $39$ ;  $Md = 79$ ;  $Mode = 65$ ; see Table 2). The sports students (sample 2b) performed better than the  
448 non-sports students (sample 2a), and the men better than the women (see Table 2). When looking  
449 at the mean performances for each of the 14 darts separately, no clear learning trend was  
450 apparent (see Figure 2).

451 The different reliability coefficients do not differ from each other in each sample (see  
452 Table 1). Thus, results are not dependent on a particular coefficient. Likewise, reliability  
453 coefficients do not differ between both subsamples. Neither does each reliability coefficient for  
454 the whole sample differ from the respective coefficients in both subsamples.

455 Reliability coefficients can already be considered acceptable for seven throws only: The  
456 odd-even reliability for the whole sample based on seven throws is .799 ( $CI_{95} [.743; .844]$ ).  
457 When estimating reliabilities for all 14 throws using the Spearman-Brown-Formula, reliability  
458 coefficients are high: The Spearman-Brown corrected odd-even reliability for the whole sample  
459 is .888 ( $CI_{95} [.855; .914]$ ). At least in this study, both in the subsamples and in the overall  
460 sample, dart throws seemed to capture a substantial proportion of systematic variation as  
461 opposed to random variation and therefore seemed to be able to capture variation in participants'  
462 true score rather well.



463 Furthermore, we investigated into the question whether reliability coefficients varied  
464 between different samples or between different groups of participants. We did so in a more  
465 exploratory manner, based on assumptions that we consider to be common when planning sport  
466 psychological experiments. As different reliability coefficients do not differ from each other,  
467 from now on we refer to the odd-even reliability, as we consider it the most appropriate (see  
468 Appendix A). In our further analyses, we distinguished between a) sports students and non-sports  
469 students, b) women and men, and c) participants who play darts and participants who do not. We  
470 distinguished between darts players and non-darts players based on a median split on  
471 participants' answers to the question "How often did you play darts during the past twelve  
472 months?". All participants who reported never to have played darts in the past twelve months  
473 were assigned to the group of non-darts players ( $n = 106$ ), whereas all other participants were  
474 assigned to the group of darts players ( $n = 94$ ). We conducted a median split in order to obtain  
475 groups of roughly equal sample size, although this approach has some disadvantages. The main  
476 disadvantage here is that the group of darts players did not only contain participants who played  
477 regularly, but also participants who had only played a couple of times. We address this  
478 shortcoming in study 3.

479 First, mean performance differs between the two groups in all three comparisons (see  
480 Table 2). That means, a) sports students were significantly better than non-sports students  
481 ( $t[167.58] = 4.6, p < .001, d = 0.65$ ), b) men were significantly better than women ( $t[137.62] =$   
482  $9.46, p < .001, d = 1.39$ ), and c) darts players were significantly better than non-darts players  
483 ( $t[198] = 4.83, p < .001, d = 0.68$ ).

484 As already reported above, reliability coefficients did not differ between sports students  
485 and non-sports students (i.e., subsample 2b and subsample 2a, see Table 2). Descriptively,

486 reliability was somewhat higher for women than for men, but confidence intervals overlap (see  
487 Table 2). Finally, reliability coefficients do not differ between darts players and non-darts  
488 players (see Table 2). If at all, non-darts players have a slightly higher reliability coefficient, but  
489 again, confidence intervals overlap.

490 Thus, and contrary to what one might have intuitively expected, neither gender, nor  
491 studying sports nor playing darts had an impact on reliability estimates. However, at least the last  
492 finding might be due to the fact that we distinguished darts players from non-darts players based  
493 on a median split, which is not the best method to compare different levels of experience. We  
494 address this issue in study 3.

### 495 **Study 3: Basketball Free Throws**

#### 496 **Participants**

497 As one goal of study three was to further investigate into the role of sports experience for  
498 reliability, we aimed to obtain two different subsamples. One subsample should consist of  
499 experienced basketball players, whereas the other one should consist of comparably  
500 inexperienced players. We assigned potential participants to the sample of experienced players,  
501 when they were active members of a basketball club and reported their free throw success rate to  
502 be at least 30%. We assigned potential participants to the sample of inexperienced players when  
503 they did not fulfil the inclusion criteria for the experienced sample. These rules were defined  
504 prior to data collection. Additionally, participants had to be able to hit the rim or score a basket at  
505 least ten out of 20 times during the practice trials in order to make sure that they were  
506 sufficiently skilled.

507 Therefore, study 3 consisted of two subsamples (sample 3a and sample 3b). Sample 3a  
508 (the inexperienced sample) consisted of 100 participants ( $M_{\text{age}} = 24.8$ ;  $SD_{\text{age}} = 4.2$ ). Sample 3b

509 (the experienced sample) consisted of 92 participants ( $M_{\text{age}} = 25.3$ ;  $SD_{\text{age}} = 7.7$ ), totalling 192  
510 participants. In sample 3a, there were 50 women and 50 men. All of them were sports students at  
511 Heidelberg University. In sample 3b, there were 42 women and 50 men. All of them were  
512 players in regional basketball clubs. Thus, in the whole sample there were 92 women and 100  
513 men, 92 non-sports students and 100 sports students. As a result of the above mentioned criteria  
514 for inclusion of participants into the different subsamples, both subsamples differed considerably  
515 with regard to basketball experience. On average, participants in subsample 3a (the rather  
516 inexperienced) reported to play basketball for 28 minutes per week ( $SD = 52$ ), whereas  
517 participants in subsample 3b (the experienced) reported to play basketball for 314 minutes per  
518 week ( $SD = 206$ ).

### 519 **Apparatus and Procedure**

520 Participants conducted all free throws in line with the regulations of the Fédération  
521 Internationale de Basketball (FIBA, 2018a; b). Participants were positioned behind the free  
522 throw line 422.5 cm from the middle point of the basket. According to FIBA, basketballs for  
523 men are supposed to weigh 567-650 g with a circumference of 74.9-78.0 cm. Basketballs for  
524 women are supposed to weigh 510-567 g with a circumference of 72.4-73.7 cm. We used two  
525 balls: one ball for all men and one ball for all women. From time to time, we made sure that both  
526 balls were still within the limits specified by the regulations. According to FIBA, the basket ring  
527 has to be positioned at a height of 304.8 cm +/- 0.6 cm and have an inside diameter of 45.0-45.9  
528 cm. We made sure that the baskets utilized were within these specifications.

529 We coded each shot as either successful (the ball went through the basket) or not (the ball  
530 did not go through the basket). There were no missing values, as all shots could be coded. We  
531 did not distinguish between different kinds of unsuccessful shots, for example, balls hitting the

532 rim or air balls, as researchers sometimes do. The reason for this is that assessing the difference  
533 between successful shots and different kinds of misses does not produce a continuous (i.e.,  
534 interval scale) measurement. In order to obtain continuous measurement, we added up all  
535 attempts, as sums of binary variables can be treated as continuous variables (e.g., Lunney, 1970).  
536 Examples for this practice can be found, for example, in intelligence tests, where each single  
537 item produces a binary datum, however, items are summed up along scales and then treated in a  
538 continuous manner.

### 539 **Results and Discussion**

540 Experienced participants performed better than inexperienced participants. In the  
541 experienced sample (3b), the average success rate was 74% ( $SD = 17$ ). In the inexperienced  
542 sample (3a), the average success rate was 32% ( $SD = 17\%$ ). When looking at the mean  
543 performances for each of the 14 throws separately, no learning trend was apparent, neither for  
544 sample 3a nor for sample 3b (see Figure 4 and Figure 5; see Appendix A for further elaboration  
545 on this issue).

546 The different reliability coefficients do not differ from each other in each sample (see  
547 Table 1). Thus, results are not dependent on a particular coefficient. At least descriptively,  
548 reliability coefficients are higher for the experienced sample (sample 3b) than for the novice  
549 sample. However, confidence intervals still overlap for both groups. In both subsamples,  
550 reliability coefficients for seven shots only are low, and they are still not acceptable when  
551 correcting for all 14 shots using the Spearman-Brown formula. The Spearman-Brown corrected  
552 odd-even reliability for sample 3a is .504 ( $CI_{95} [.341; .637]$ ) and for sample 3b it is .62 ( $CI_{95}$   
553 [.475; .732]).





599 is a well known property of reliabilities according to CTT. Our results also demonstrate, that  
600 depending on a measurement's reliability, vastly different numbers of items are required in order  
601 to achieve the same acceptable level of reliability. This observation has consequences for the  
602 construction of performance outcome measures as we will discuss in more detail below.

603 Most importantly, our results demonstrate that reliabilities of sport performance outcome  
604 measurements may strongly depend on sample characteristics: Reliabilities estimates for both  
605 subsamples in basketball were very low (.504 and .62), however for the whole sample the  
606 estimated reliability was substantially higher (.826). This observation underlines the necessity of  
607 having samples with true-score variation when researchers want to obtain reliable measurement  
608 (see also Appendix A; Vaughn et al., 2012). At the same time this observation debunks what we  
609 (anecdotally) perceive to be a common misconception in sport psychological research, namely  
610 that reliability will be higher in expert samples than in non-expert samples. This also means that  
611 researchers need to be particularly careful when conducting studies with rather homogenous  
612 expert samples, as this approach might lead to low reliabilities.

613 When interpreting the present results, it is important to keep in mind that all estimates  
614 reported here depend on the respective samples and operationalizations and therefore do not  
615 necessarily generalize to other situations or samples (see Appendix A). Furthermore, it is  
616 important to keep in mind that our results may paint a rosier picture of reliabilities than is  
617 actually warranted when looking at existing studies: One reason why estimates appear to be  
618 rather high in our studies is that we employed 14 trials. To the extent that past studies may have  
619 employed fewer trials, all else equal, they had less reliable measurements.

620 **Limitations and Unintended Consequences**

621 We based our main conclusions on CTT in general and on specific estimators (i.e., odd-  
622 even reliability and the Spearman-Brown formula) in particular. However, alternative approaches  
623 exist. They comprise a) different estimates that have existed within the framework of CTT for a  
624 long time (e.g., the Kristof or the Guttman formulas, see Bühner, 2011); b) novel estimates that  
625 have been proposed only recently (e.g., omega as an alternative to Cronbach's alpha [McNeish,  
626 2018]; weighted kappa [Robinson & O'Donoghue, 2007] to assess agreement amongst observers  
627 in performance analysis; and special coefficients for particular research designs within sport  
628 science [Hopkins, 2017]); and c) estimates computed via structural equation modeling (SEM;  
629 e.g., Raykov, 1997). All of these approaches have advantages and disadvantages and it is  
630 impossible to say that one of them is per se superior. Just as one example, one presumed  
631 advantage of omega above Chronbach's alpha is that omega relies less on modeling assumptions  
632 (McNeish, 2018). However, the advantages of omega have been questioned and currently there is  
633 a controversial discussion regarding its merits (Raykov & Marcoulides, 2019; Savalei & Reise,  
634 2019).

635 Furthermore, CTT itself has some well-known weaknesses, for example its dependency  
636 on sometimes questionable assumptions and on sample characteristics (Bühner, 2011; see also  
637 Appendix A). Item Response Theory (IRT) in turn allows for modeling the probability of a  
638 response to an item as a joint function of both this item's difficulty (the item parameter) and a  
639 person's ability (the person parameter), which is a substantial advantage over CTT (Bühner,  
640 2011). Taken together, we consider it to be important to keep in mind that alternatives to the  
641 particular estimates that we employed and to CTT exist. We hope to stipulate a discussion on  
642 which theoretical approach and which coefficients are best suited for measurement in sport  
643 psychology. In order to foster this discussion, we make all of our raw data public, so that



644 researchers may take these data and calculate other estimates of reliability, SEMs or parameters  
645 from IRT.

646         As we have mentioned above, our estimates depend both on our samples and on our exact  
647 operationalization of the different measures. For example, maybe the reliability estimate for the  
648 free throws would have been different had participants been closer to the basket. Hence,  
649 theoretically, reliability (and its estimates) always refer to a measurement, *not* to an instrument.  
650 Therefore, we cannot say that we estimated *the* reliability of darts, or *the* reliability of free  
651 throws, instead we estimated the reliabilities of our specific measurements. An unintended  
652 consequence of our study would be if from now on researchers in sport psychology would  
653 predominantly use darts as dependent variable, because “it has been proven to be reliable”.  
654 Future studies with different samples might be different in terms of reliability.

655         Moreover, reliability must not be confused with validity. It would be a mistake if  
656 researchers simply used certain measures because they are reliable, and did not care about  
657 validity, a concern that has been raised in psychometrics (e.g., Bühner, 2011). To test theories  
658 that relate theoretical constructs to each other (e.g., construct A influences construct B for  
659 individuals drawn from population P under conditions C), it is necessary to not only have reliable  
660 measures, but also valid measures that actually measure construct A and B and control for P and  
661 C. Validity typically refers to whether a given measure in fact measures what it claims to  
662 measure. Unfortunately, frequently used measures within psychology (e.g., Schimmack, 2019)  
663 and sport science (Fischman, 2015) might not measure what they claim to measure. Although,  
664 the present paper focused on reliability and not validity, high quality measurement in any  
665 scientific field needs to focus on both. However, high reliability is a prerequisite for validity: A  
666 measurement that is not reliable cannot be valid. Finally, we would like to emphasize that our

667 results do not intend to undermine the credibility, quality or replicability of prior studies that  
668 have employed golf putts, darts, or free throws. Instead, they should draw attention to the  
669 importance of reliable measurement in sport psychology with the aim of securing it in the future.

#### 670 **Conclusions**

671 Sport performance outcome measures may substantially differ regarding their reliability  
672 and may have different reliabilities for different samples (and not necessarily in an intuitive  
673 way). Suppose three research teams each used a different one of our measures with their  
674 respective reliabilities to answer a research question (e.g., the effects of pressure or fatigue on  
675 perceptual-motor performance). All else equal, these teams would have substantially different  
676 likelihoods of a) finding an effect, given it exists, of b) replicating an effect found in a prior  
677 study, and c) being able to make meaningful comparisons between studies, variables, and  
678 theories.

679 When conducting studies, we hope that researchers in sport psychology will try to  
680 construct reliable measurements, that they will assess their measurement's reliability, and that  
681 they will interpret their results in the light of these reliabilities. Reliabilities need to be *high*, and  
682 moderate reliabilities may exacerbate methodological problems. For example, Westfall and  
683 Yarkoni (2016) report that the Type-1 error rate when assessing incremental validity via  
684 regression models was highest for moderate reliabilities (at least for certain sample sizes).

685 Regarding conclusions, we hope that researchers will be very careful when comparing  
686 findings to each other that may stem from measurements with different reliabilities. Likewise,  
687 we hope that researchers will consider the role of (different) reliabilities when assessing  
688 replications as being successful or not. If possible, we suggest that researchers pretest their  
689 performance outcome measure and try to determine an optimal number of trials that provides

690 sufficient reliability, but that does not induce threats to validity (such as fatigue or learning  
691 effects) and is still economically feasible (see Appendix A for more information and guidelines  
692 on these issues). Whereas increasing reliability by adding items only works to a certain extent for  
693 common psychological measurement procedures such as questionnaires, for performance  
694 outcome measures, such as discussed in this paper, it seems to be more promising (for more  
695 information see Appendix A). Furthermore, as mentioned above, in order to obtain the same  
696 power to detect an effect of the same (true) size, researchers need smaller samples when using  
697 more reliable measures than when using less reliable measures. Therefore, there is a trade-off  
698 regarding research economy<sup>8</sup>: On the one hand, adding items or trials to a measurement in order  
699 to make it more reliable will make the measurement less economical by increasing its duration.  
700 On the other hand, this approach will make the measurement more reliable and thus more  
701 economical because smaller sample sizes are needed. It seems to be an interesting endeavor for  
702 future research to try and formalize this trade-off depending on its various costs and benefits.

703 In this endeavor, experimenters should attempt to use individual performance outcome  
704 measures that allow for sufficient variation in performance that is indicative of true performance  
705 variation and not random performance fluctuation and measurement error. To this end the  
706 following guiding questions might prove helpful (see also Table 3): a) what is my precise  
707 research question and how well do the variables in my research design measure the constructs in  
708 my research question; b) what is the skill level of my participants or how experienced are  
709 participants with the (or similar) tasks being measured; c) how difficult does the task have to be  
710 (e.g. putting distance in golf); and d) how many trials are sufficient to achieve adequate  
711 reliability, while not threatening validity (e.g. motivation, calibration, learning, fatigue, etc.).

712

713

714

Journal Pre-proof

- 715 References
- 716 Ajemian, R., D'Ausilio, A., Moorman, H., & Bizzi, E. (2010). Why professional athletes need a  
717 prolonged period of warm-up and other peculiarities of human motor learning. *Journal of*  
718 *Motor Behavior*, 42, 381–388. doi:10.1080/00222895.2010.528262
- 719 Amelang, M., & Zielinski, W. (2002). *Psychologische Diagnostik und Intervention*. Heidelberg:  
720 Springer.
- 721 Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk, R., ...  
722 Johnson, V. E. (2018). Redefine statistical significance. *Nature Human Behaviour*, 2, 6–  
723 10. doi:10.1038/s41562-017-0189-z
- 724 Barbero-Álvarez, J. C., Coutts, A., Granda, J., Barbero-Álvarez, V., & Castagna, C. (2010). The  
725 validity and reliability of a global positioning satellite system device to assess speed and  
726 repeated sprint ability (RSA) in athletes. *Journal of Science and Medicine in Sport*, 13,  
727 232-235. doi:10.1016/j.jsams.2009.02.005
- 728 Baumeister, R. F., & Showers, C. J. (1986). A review of paradoxical performance effects:  
729 Choking under pressure in sports and mental tests. *European Journal of Social*  
730 *Psychology*, 16, 361-383. doi:10.1002/ejsp.2420160405
- 731 Bühner, M. (2011). *Einführung in die Test- und Fragebogenkonstruktion*, 3. aktual. Auflage.  
732 München: Pearson.
- 733 Button, K. S., Ioannidis, J. P. A., Mokrysz, C., Nosek, B. A., Flint, J., Robinson, E. S. J., &  
734 Munafò, M. R. (2013). Power failure: Why small sample sizes undermine the reliability  
735 of neuroscience. *Nature Reviews Neuroscience*, 14, 365-376. doi:10.1038/nrn3475
- 736 Camerer, C. F., Dreber, A., Holzmeister, F., Ho, T.-H., Huber, J., Johannesson, M., ... Wu, H.  
737 (2018). Evaluating the replicability of social science experiments in *Nature and Science*

- 738 between 2010 and 2015. *Nature Human Behaviour*, 2, 637-644. doi:10.1038/s41562-018-  
739 0399-z
- 740 Charter, R. A. (1999) Sample size requirements for precise estimates of reliability,  
741 generalizability, and validity coefficients. *Journal of Clinical and Experimental*  
742 *Neuropsychology*, 21, 559-566. doi:10.1076/jcen.21.4.559.889
- 743 Coutts, A. J., & Duffield, R. (2010). Validity and reliability of GPS units for measuring  
744 movement demands of team sports. *Journal of Science and Medicine in Sport*, 13, 33–  
745 135. doi:10.1016/j.jsams.2008.09.015
- 746 Cumming, G. (2012). *Understanding the new statistics: Effect sizes, confidence intervals, and*  
747 *meta-analysis*. New York: Routledge.
- 748 Cumming, G. (2014). The new statistics why and how. *Psychological Science*, 25, 7-29.  
749 doi:10.1177/0956797613504966
- 750 Danner, D. (2015). *Reliabilität: Die Genauigkeit einer Messung. GESIS Survey Guidelines*.  
751 Mannheim: GESIS – Leibniz Institut für Sozialwissenschaften. doi:10.15465/sdm-sg\_011
- 752 Diedenhofen, B. & Musch, J. (2015). cocor: A comprehensive solution for the statistical  
753 comparison of correlations. *PLoS ONE*, 10. e0121945. doi:10.1371/journal.pone.0121945
- 754 Diedenhofen, B., & Musch, J. (2016). cocron: A web interface and R package for the statistical  
755 comparison of Cronbach's alpha coefficients. *International Journal of Internet Science*,  
756 11, 51–60.
- 757 Eklund, R. C., Tenenbaum, G., & Kamata, A. (2012). *Measurement in sport and exercise*  
758 *psychology*. Champaign, IL, Human Kinetics.

- 759 Fischman, M. G. (2015). On the continuing problem of inappropriate learning measures:  
760 Comment on Wulf et al. (2014) and Wulf et al. (2015). *Human Movement Science*, 42,  
761 225-231. doi:10.1016/j.humov.2015.05.011
- 762 Fédération Internationale de Basketball (FIBA) (2018a). *Official Basketball rules 2018*.  
763 Fédération Internationale de Basketball: Mies, Switzerland.  
764 <https://www.fiba.basketball/documents/official-basketball-rules.pdf>
- 765 Fédération Internationale de Basketball (FIBA) (2018b). *Official Basketball rules: Basketball*  
766 *equipment*. Fédération Internationale de Basketball: Mies, Switzerland.  
767 <http://www.fiba.basketball/OBR-2018-Basketball-Equipment-Yellow-Version-2.pdf>
- 768 Fraley, R. C., & Vazire, S. (2014). The N-Pact Factor: Evaluating the quality of empirical  
769 journals with respect to sample size and statistical power. *PLoS One*, 9, e109019. doi:  
770 10.1371/journal.pone.0109019
- 771 Friese, M., & Frankenbach, J. (in press). P-hacking and publication bias interact to distort meta-  
772 analytic effect size estimates. *Psychological Methods*.
- 773 Gaito, J. (1980). Measurement scales and statistics: Resurgence of an old misconception.  
774 *Psychological Bulletin*, 87, 564-567. doi:10.1037/0033-2909.87.3.564
- 775 García-Pérez, M. A. (2017). Thou shalt not bear false witness against null hypothesis  
776 significance testing. *Educational and Psychological Measurement*, 77, 631-662.  
777 doi:10.1177/0013164416668232
- 778 Gelman, A. (2015, April 28). What's the most important thing in statistics that's not in the  
779 textbooks? [Web log post]. Retrieved from  
780 [https://statmodeling.stat.columbia.edu/2015/04/28/whats-important-thing-statistics-thats-](https://statmodeling.stat.columbia.edu/2015/04/28/whats-important-thing-statistics-thats-not-textbooks/)  
781 [not-textbooks/](https://statmodeling.stat.columbia.edu/2015/04/28/whats-important-thing-statistics-thats-not-textbooks/)

- 782 Gelman, A., & Stern, H. (2006). The difference between “significant” and “not significant” is not  
783 itself statistically significant. *American Statistician*, *60*, 328-331.
- 784 Hopkins, W. G. (2017). Spreadsheets for analysis of validity and reliability. *Sportscience*, *21*.
- 785 Hopkins, W. G. (2000). Measures of reliability in sports medicine and science. *Sports Medicine*,  
786 *30*, 1-15. doi:10.2165/00007256-200030010-00001
- 787 Jennings, D., Cormack, S., Coutts, A. J., Boyd, L., & Aughey, R. J. (2010). The validity and  
788 reliability of GPS units for measuring distance in team sport specific running patterns.  
789 *International Journal of Sports Physiology and Performance*, *5*, 328-341.  
790 doi:10.1123/ijsp.5.3.328
- 791 Johnston, R. J., Watsford, M. L., Pine, M. J., Spurrs, R. W., Murphy, A. J., & Pruyn, E. C.  
792 (2012). The validity and reliability of 5-Hz global positioning system units to measure  
793 team sport movement demands. *The Journal of Strength & Conditioning Research*, *26*,  
794 758-765. doi:10.1519/JSC.0b013e318225f161
- 795 Lakens, D., & Evers, E. R. K. (2014). Sailing from the seas of chaos into the corridor of stability:  
796 Practical recommendations to increase the informational value of studies. *Perspectives on*  
797 *Psychological Science*, *9*, 278-292. doi:10.1177/1745691614528520
- 798 Loken, E., & Gelman, A. (2017). Measurement error and the replication crisis. *Science*, *355*,  
799 584-585. doi:10.1126/science.aal3618
- 800 Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Reading:  
801 Addison-Weasley.
- 802 Love, J., Selker, R., Marsman, M., Jamil, T., Dropmann, D., Verhagen, A. J., Ly, A., Gronau, Q.  
803 F., Šmíra, M., Epskamp, S., Matzke, D., Wild, A., Knight, P., Rouder, J. N., Morey, R.



- 804 D., Wagenmakers, E.-J. (2019). JASP: Graphical statistical software for common  
805 statistical designs. *Journal of Statistical Software*, 88. doi:10.18637/jss.v088.i02
- 806 Lunney, G. H. (1970). Using analysis of variance with a dichotomous dependent variable: An  
807 empirical study. *Journal of Educational Measurement*, 7, 263-269.
- 808 McNeish, D. (2018). Thanks coefficient alpha, we'll take it from here. *Psychological Methods*,  
809 23, 412-433. doi:10.1037/met0000144
- 810 McShane, B. B., Gal, D., Gelman, A., Robert, C., & Tackett, J. L. (2019). Abandon statistical  
811 significance. *The American Statistician*, 73, 235-245.  
812 doi:10.1080/00031305.2018.1527253
- 813 Munafò, M. R., Nosek, B. A., Bishop, D. V. M., Button, K. S., Chambers, C. D., Percie du Sert,  
814 N., Simonsohn, U., Wagenmakers, E.-J., Ware, J. J., Ioannidis, J. P. A. (2017). A  
815 manifesto for reproducible science. *Nature Human Behaviour*, 1, 0021.  
816 doi:10.1038/s41562-016-0021
- 817 Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., ...  
818 Yarkoni, T. (2015). Promoting an open research culture. Author guidelines for journals  
819 could help to promote transparency, openness, and reproducibility. *Science*, 348, 1422-  
820 1425. doi:10.1126/science.aab2374
- 821 Nelson, L. D., Simmons, J., & Simonsohn, U. (2018). Psychology's renaissance. *Annual Review*  
822 *of Psychology*, 69, 511-534. doi:10.1146/annurevpsych-122216-011836
- 823 Open Science Collaboration. (2012). An open, large-scale, collaborative effort to estimate the  
824 reproducibility of psychological science. *Perspectives on Psychological Science*, 7, 657-  
825 660.

- 826 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.  
827 *Science*, 349, aac4716. doi:10.1126/science.aac4716
- 828 Pashler, H., & Wagenmakers, E.-J. (2012). Editors' introduction to the special section on  
829 replicability in psychological science: A crisis of confidence? *Perspectives on*  
830 *Psychological Science*, 7, 528-530. doi:10.1177/1745691612465253
- 831 Petersen, C., Pyne, D., Portus, M., & Dawson, B. (2009). Validity and reliability of GPS units to  
832 monitor cricket-specific movement patterns. *International Journal of Sports Physiology*  
833 *and Performance*, 4, 381–393. doi:10.1123/ijsp.4.3.381
- 834 Raykov, T. (1997). Estimation of composite reliability for congeneric measures. *Applied*  
835 *Psychological Measurement*, 21, 173-184.
- 836 Raykov, T., & Marcoulides, G. A. (2019). Thanks coefficient alpha, we still need you!  
837 *Educational and Psychological Measurement*, 79, 200–210.  
838 doi:10.1177/0013164417725127
- 839 Robinson, G., & O'Donoghue, P. (2007). A weighted kappa statistic for reliability testing in  
840 performance analysis of sport. *International Journal of Performance Analysis in Sport*, 7,  
841 12-19. doi:10.1080/24748668.2007.11868383
- 842 Savalei, V., & Reise, S. P. (2019). Don't forget the model in your model-based reliability  
843 coefficients: A reply to McNeish (2018). *Collabra: Psychology*, 5, 36.  
844 doi:10.1525/collabra.247
- 845 Schimmack, U. (2010). What multi-method data tell us about construct validity. *European*  
846 *Journal of Personality*, 24, 241-257. doi:10.1002/per.771
- 847 Schimmack, U. (2019, under review). The validation crisis in psychology. *Meta-Psychology*.  
848 (<https://replicationindex.files.wordpress.com/2019/04/validation.crisis.v3.pdf>)

- 849 Schönbrodt, F. D., & Perugini, M. (2013). At what sample size do correlations stabilize? *Journal*  
850 *of Research in Personality*, 47, 609e612. <http://dx.doi.org/10.1016/j.jrp.2013.05.009>
- 851 Schweizer, G., & Furley, P. (2016). Reproducible research in sport and exercise psychology: The  
852 role of sample sizes. *Psychology of Sport and Exercise*, 23, 114-122.  
853 doi:10.1016/j.psychsport.2015.11.005
- 854 Savalei, V., & Dunn, E. (2015). Is the call to abandon  $p$ -values the red herring of the replicability  
855 crisis? *Frontiers in Psychology*, 6, 1-4. doi:10.3389/fpsyg.2015.00245
- 856 Schönbrodt, F. D., & Perugini, M. (2013). At what sample size do correlations stabilize? *Journal*  
857 *of Research in Personality*, 47, 609e612. doi:10.1016/j.jrp.2013.05.009
- 858 Selkirk, K. (1976). Re-designing the dartboard. *The Mathematical Gazette*, 60, 171-178.
- 859 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-Positive Psychology :  
860 Undisclosed flexibility in data collection and analysis allows presenting anything as  
861 significant. *Psychological Science*, 22, 1359-1366. doi:10.1177/0956797611417632
- 862 Simonsohn, U. (2015). Small telescopes: Detectability and the evaluation of replication results.  
863 *Psychological Science*, 26, 559-569. doi:10.1177/0956797614567341
- 864 Soto, C. J. (2019). How replicable are links between personality traits and consequential life  
865 outcomes? The Life Outcomes Of Personality Replication Project. *Psychological Science*,  
866 30, 711-727. doi:10.1177/0956797619831612
- 867 Spearman, C. (1904). The proof and measurement of association between two things. *The*  
868 *American Journal of Psychology*, 15, 72-101. doi:10.2307/1412159
- 869 Spearman, C. (1910). Correlation calculated from faulty data. *British Journal of Psychology*, 3,  
870 271-295. doi:10.1111/j.2044-8295.1910.tb00206.x

- 871 Stanley, D. J., & Spence, J. R. (2014). Expectations for replications: Are yours realistic?  
872 *Perspectives on Psychological Science*, 9, 305-318. doi:10.1177/1745691614528518
- 873 Steyer, R., & Eid, M. (1993). *Messen und Testen*. Springer: Heidelberg.
- 874 Tibshirani, R. J., Price, A., & Taylor, J. (2011). A statistician plays darts. *Journal of the Royal*  
875 *Statistical Society*, 174, 213-226.
- 876 Vaughn, B. K., Lee, H.-Y., & Kamata, A. (2012). Reliability. In R. C. Eklund, G. Tenenbaum, &  
877 A. Kamata (Eds.), *Measurement in sport and exercise psychology*. Champaign, IL:  
878 Human Kinetics.
- 879 Wagenmakers, E.-J., Marsman, M., Jamil, T., Ly, A., Verhagen, A. J., Love, J., Selker, R.,  
880 Gronau, Q. F., Šmíra, M., Epskamp, S., Matzke, D., Rouder, J. N., Morey, R. D. (2018).  
881 Bayesian inference for psychology. Part I: Theoretical advantages and practical  
882 ramifications. *Psychonomic Bulletin & Review*, 25, 35-57.
- 883 Walsh, V. (2014). Is sport the brain's biggest challenge? *Current Biology*, 24, R859-R860.  
884 doi:10.1016/j.cub.2014.08.003
- 885 Westfall, J., & Yarkoni, T. (2016). Statistically controlling for confounding constructs is harder  
886 than you think. *PLoS ONE*, 11, e0152719. doi:10.1371/journal.pone.0152719
- 887 Wunderlich, F., Heuer, H., Furley, P., & Memmert, D. (2019, iFirst). A serial-position curve in  
888 high-performance darts: The effect of visuomotor calibration on throwing accuracy.  
889 *Psychological Research*, 1-8. doi:10.1007/s00426-019-01205-2
- 890 Zhu, W. (2013). Reliability: what type, please! *Journal of Sport and Health Science*, 2, 62-64.  
891 doi:10.1016/j.jshs.2012.11.001
- 892
- 893

894

**Footnote**

895 <sup>1</sup>For example, the JASP software package includes Bayesian parameter estimation and  
896 Bayes factor hypothesis testing via a graphical user interface (see Love et al., 2019).

897 Furthermore, free and open-source packages for specific procedures in R make it more feasible  
898 for researchers to use these procedures, to name only two examples.

899 <sup>2</sup>Our point is not that measurement is generally absent in the methodological literature in  
900 psychology, quite the contrary. However, in the context of the current debate on methodological  
901 practices (as described by Nelson et al., 2018) only few papers focus on measurement (e.g.,  
902 Loken & Gelman, 2017).

903 <sup>3</sup>There is an ongoing debate in psychology whether researchers should abandon Null  
904 Hypothesis Significance Testing (NHST), and, if they do, which methods they should use  
905 instead. Some authors suggest abandoning not only NHST, but the frequentist perspective  
906 altogether by employing Bayesian methods (e.g., Wagenmakers et al., 2018). Some suggest  
907 abandoning statistical significance as a threshold, but to retain p-values and treat them as one  
908 (albeit continuous) piece of information among others (McShane, Gal, Gelman, Robert, &  
909 Tackett, 2019). Some authors retain a frequentist perspective, but suggest replacing NHST by  
910 focusing on confidence intervals (e.g., Cumming, 2012, 2014). Some authors defend the utility  
911 of NHST (e.g., García-Pérez, 2017; Savalei & Dunn, 2015), and some have even suggested  
912 improving NHST by redefining statistical significance (Benjamin et al., 2018). We would like to  
913 note that in this article, we do not take any position regarding these questions. Instead, we  
914 emphasize that reliable measurement plays a key role for all of the methods discussed above.

915           <sup>4</sup>As these examples refer to choices researchers can make, the underlying construct was  
916 initially called “researcher degrees of freedom” (Simmons et al., 2011, p. 1359). Later, Simmons  
917 and colleagues adopted the term p-hacking (Nelson et al., 2018, p. 513).

918           <sup>5</sup>Generally, one needs to be careful when comparing significant and non-significant  
919 effects to each other: When one effect is significant and the other one is not, this does not mean  
920 that the difference between them is significant. This holds both for differences between groups  
921 (Gelman & Stern, 2006) and for differences between correlations (Diedenhofen & Musch, 2015).

922           <sup>6</sup>First, we looked through the latest issues of sport psychological journals in order to  
923 identify generally used performance outcome variables. This search led us to golf putts, darts and  
924 free throws. Then, we conducted a literature search in google scholar using the key words “golf  
925 putts”; “darts”; “free throws”. We combined these key words with different search terms, such as  
926 “psychology”, “performance”; and “experiment”. Our criterion for inclusion was that the paper  
927 reported a study in which the respective outcome measure had been used (as compared to, for  
928 example, a mathematical model of darts performance).

929           <sup>7</sup>At least in the ones we could access, we could not check the full text of nine articles due  
930 to difficulties acquiring the full text.

931           <sup>8</sup>We thank an anonymous reviewer for this idea.

932

## RELIABILITY

44

933 Table 1

934 *Reliability estimates for different measures*

|                              | N   | split-half<br>(odd-even) | split-half<br>(half-half) | split-half<br>(odd-even) | split-half<br>(half-half) | Cronbach's<br>alpha |
|------------------------------|-----|--------------------------|---------------------------|--------------------------|---------------------------|---------------------|
| Golf (distance)              | 100 | .555 (.402; .678)        | .552 (.399; .675)         | .714 (.602; .798)        | .711 (.598; .796)         | .670 (.567; .758)   |
| Golf (successful<br>putts)   | 100 | .443 (.27; .588)         | .360 (.176; .52)          | .614 (.475; .723)        | .529 (.371; .657)         | .598 (.472; .705)   |
| Darts-I<br>(sample 2a)       | 100 | .797 (.712; .859)        | .700 (.584; .788)         | .887 (.836; .923)        | .824 (.749; .878)         | .855 (.81; .894)    |
| Darts-I<br>(sample 2b)       | 100 | .742 (.639; .819)        | .736 (.631; .815)         | .852 (.787; .898)        | .848 (.782; .895)         | .834 (.782; .878)   |
| Darts total                  | 200 | .799 (.743; .844)        | .732 (.66; .79)           | .888 (.855; .914)        | .845 (.8; .881)           | .863 (.833; .889)   |
| Basketball-I<br>(sample 3a)  | 100 | .337 (.151; .5)          | .300 (.11; .469)          | .504 (.341; .637)        | .462 (.292; .604)         | .502 (.346; .634)   |
| Basketball-II<br>(sample 3b) | 92  | .449 (.269; .599)        | .392 (.204; .552)         | .62 (.475; .732)         | .563 (.405; .688)         | .547 (.399; .672)   |
| Basketball total             | 192 | .703 (.623; .768)        | .677 (.592; .747)         | .826 (.775; .866)        | .807 (.751; .851)         | .812 (.771; .849)   |

935 *Note.* The first two columns for the split-half reliabilities report the simple correlation between the two test halves. The third and the  
 936 fourth columns report the respective coefficients for all 14 items, computed using the Spearman-Brown formula. Cronbach's alpha  
 937 also refers to all 14 items.

938 Numbers in brackets are 95% confidence intervals (rounded to the third decimal place).

939 For estimating the 95% CIs for Cronbach's alpha we used the cocron package in R (via its web interface) (Diedenhofen & Musch,  
 940 2016).

941 Table 2

942 *Darts performance for different subgroups*

|                               | N   | mean<br>performance | split-half<br>(odd-even) |
|-------------------------------|-----|---------------------|--------------------------|
| Darts-I<br>(students)         | 100 | 99 (45)             | .887 (.836; .923)        |
| Darts-II<br>(sports students) | 100 | 75 (28)             | .852 (.787; .898)        |
| Darts total                   | 200 | 87 (39)             | .888 (.855; .914)        |
| Women                         | 94  | 111 (41)            | .871 (.812; .913)        |
| Men                           | 106 | 66 (22)             | .742 (.642; .817)        |
| Darts players                 | 94  | 73 (30)             | .85 (.78; .90)           |
| Non-darts players             | 106 | 99 (43)             | .89 (.84; .92)           |

943

944 *Note.* The third column reports the mean distance from the bulls eye averaged over all 14 throws  
 945 (in mm). Numbers in brackets are standard deviations.

946 The fourth column reports the split-half reliability coefficient for all 14 items, computed using  
 947 the Spearman-Brown formula. Numbers in brackets are confidence intervals (rounded to the  
 948 third decimal place).

949



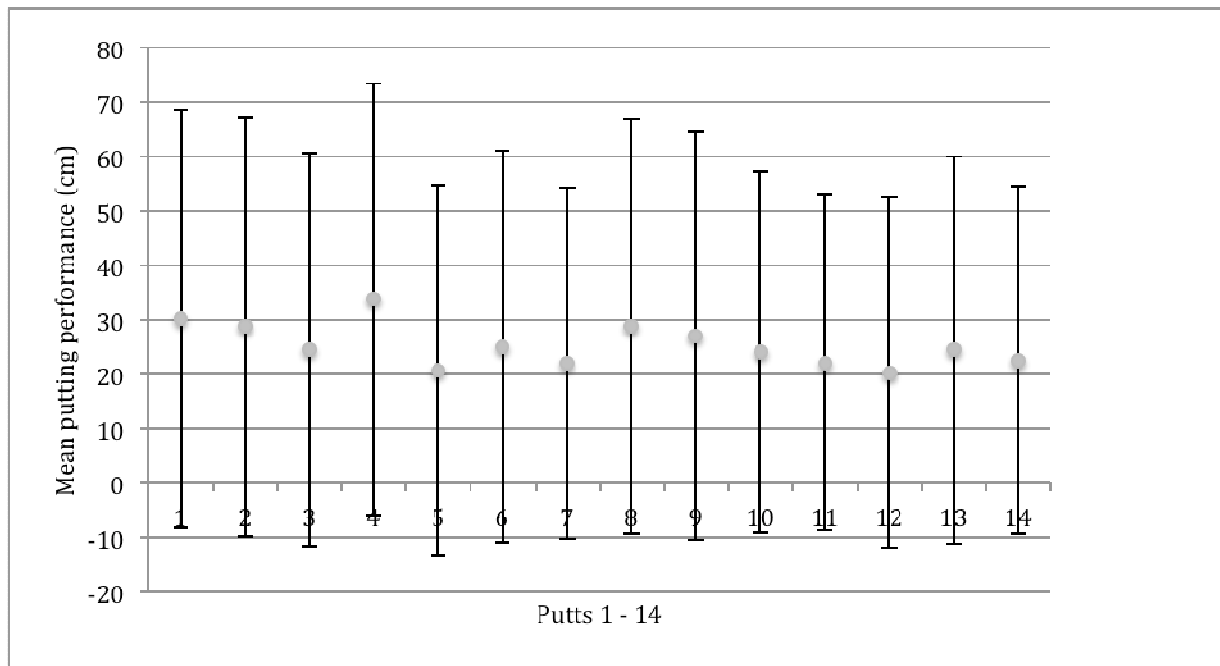
950 Table 3

951 *Guidelines for creating reliable performance outcome measures*

|                        |  |
|------------------------|--|
| Initial considerations | <ul style="list-style-type: none"> <li>• Select a measurement procedure based on theoretical grounds and research goals.</li> </ul>  |
| Pretest                | <ul style="list-style-type: none"> <li>• Pretest the measurement's reliability for a specific number of trials.</li> <li>• Use a sample that is drawn from the same population that you intend to draw your main study's sample from.</li> <li>• When considering sample size, think about precision, not statistical significance (see, for example, Appendix A; Charter, 1999; Schönbrodt and Perugini, 2013).</li> <li>• Consider different estimates of reliability: Which one is best suited for your measurement based on practical and statistical assumptions (see Appendix A for a brief overview)?</li> <li>• Estimate reliability.</li> </ul> |
| Main study             | <ul style="list-style-type: none"> <li>• Based on the above estimate of reliability, calculate the number of trials that you need in order to achieve a certain level of reliability.</li> <li>• Conduct your main study, and estimate reliability again.</li> </ul>   |
| Future studies         | <ul style="list-style-type: none"> <li>• Take into consideration that reliability is dependent (among other factors) on samples.</li> <li>• Different measurements using the same instrument may therefore lead to different estimations of reliability, for example when samples differ regarding true score variation.</li> </ul>  |

952

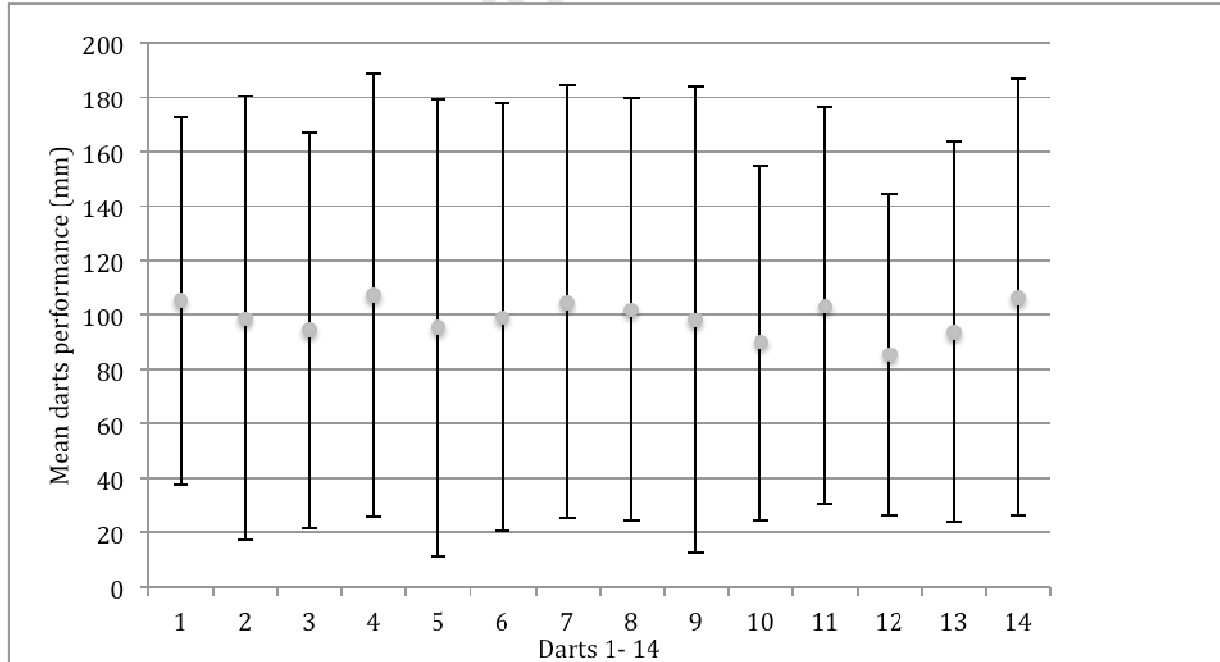
953



954

955 *Figure 1.* Mean putting performance for all 14 putts (distance from the hole in centimeters) in

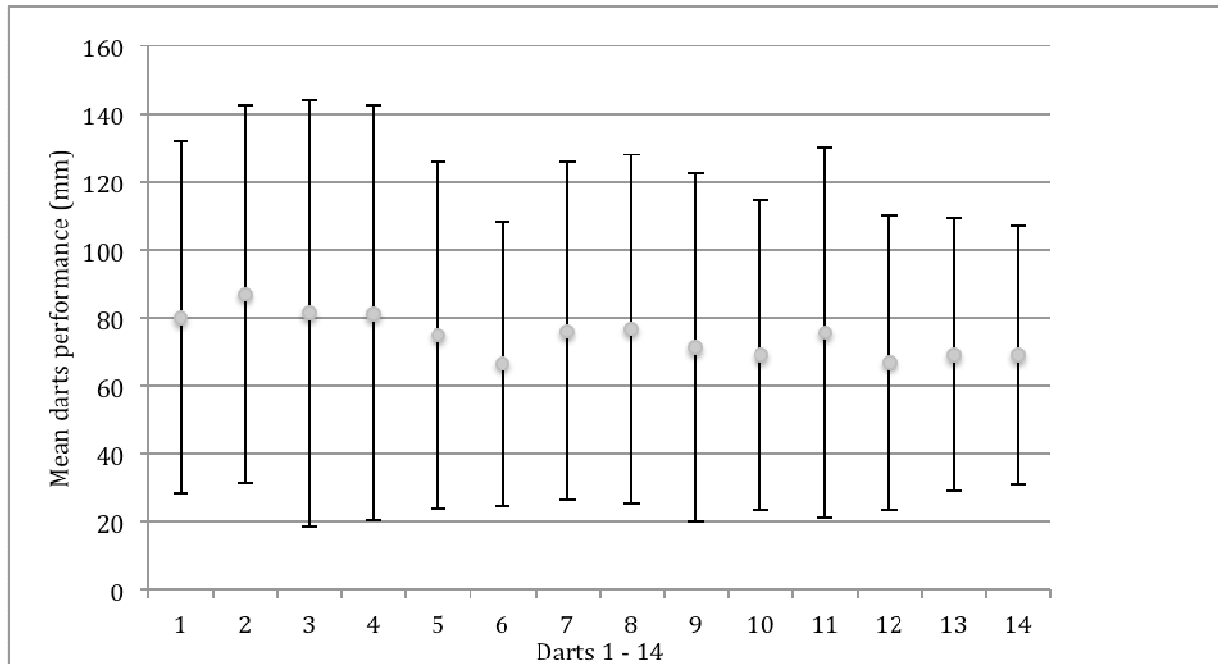
956 Study 1. Error bars are standard deviations.



957

958 *Figure 2.* Mean darts performance (distance from the bull in millimeters) in Study 2, sample 2a

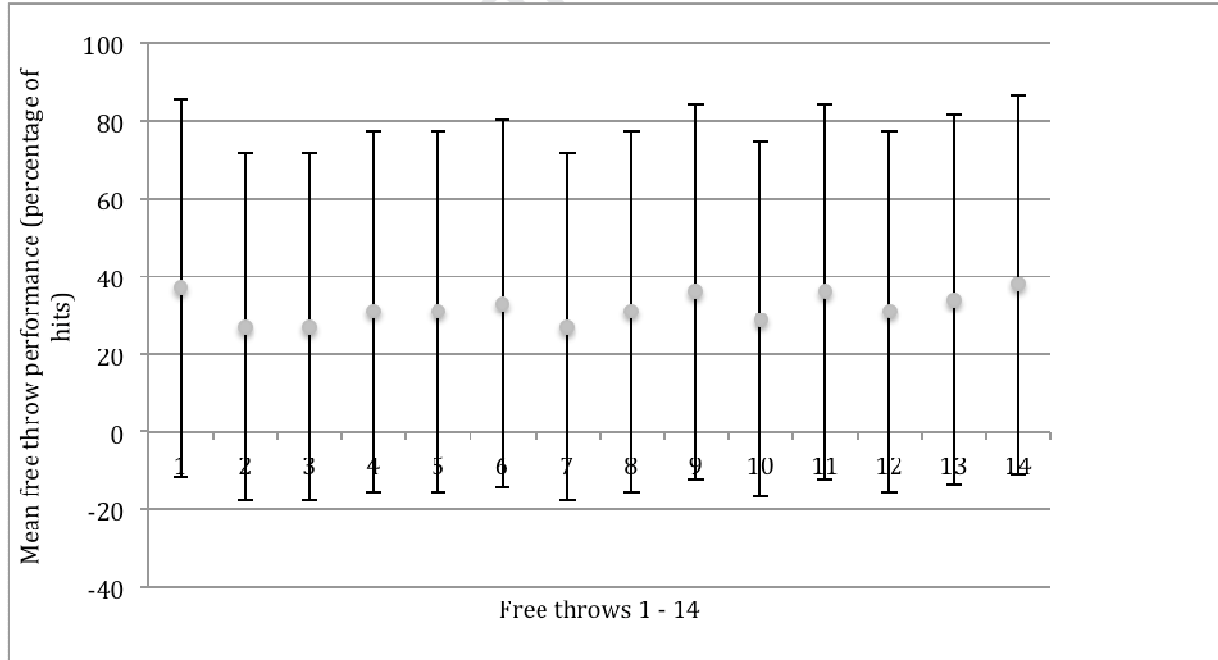
959 (students). Error bars are standard deviations.



960

961 *Figure 3.* Mean darts performance (distance from the bull in millimeters) in Study 2, sample 2b

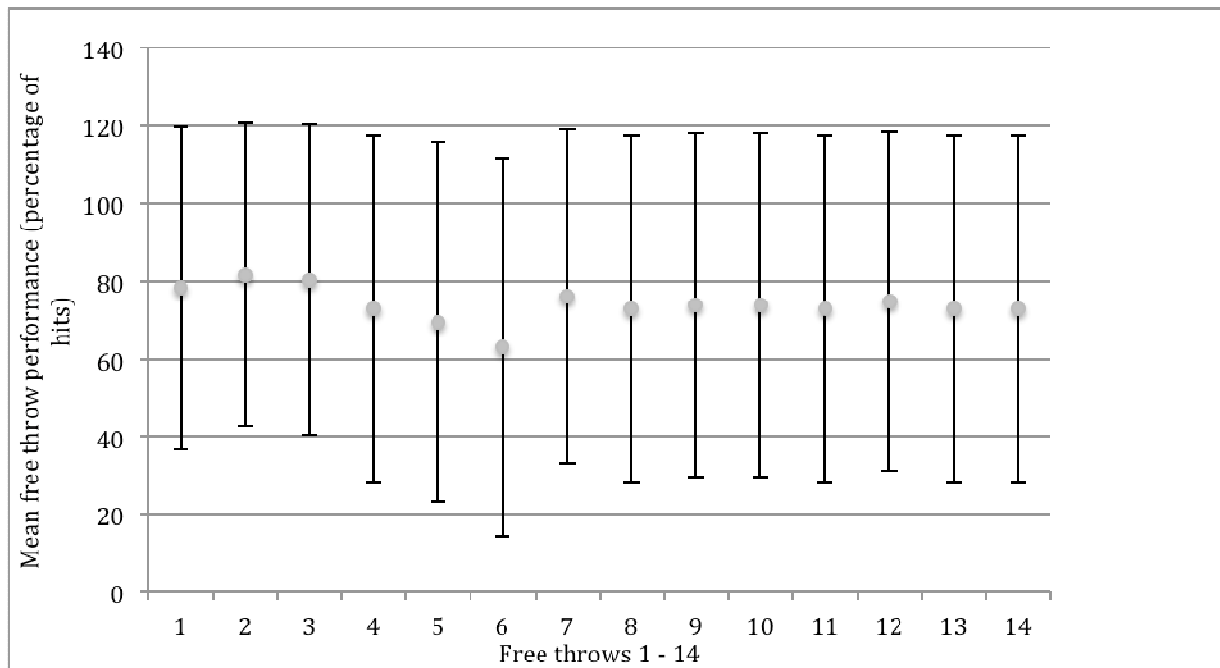
962 (sports students). Error bars are standard deviations.



963

964 *Figure 4.* Mean free throw performance (percentage of hits) in Study 3, sample 3a

965 (inexperienced). Error bars are standard deviations.



966

967 *Figure 5.* Mean free throw performance (percentage of hits) in Study 3, sample 3b (experienced).

968 Error bars are standard deviations.

969

970 Appendix A: Statistical considerations

971 **Justification of sample size planning**

972 In the context of CTT, reliability coefficients are estimated using correlations. Therefore,  
973 their size is independent of the respective sample size (i.e., unlike p-values, correlations do not  
974 change as a mere function of sample size). Thus, the sample size in a reliability analysis does not  
975 affect the estimated reliability per se, but instead the precision of the estimate. The precision of  
976 the estimate can be captured by the confidence interval around the estimate. As correlations  
977 stabilize only at rather high numbers of participants (e.g., Schönbrodt and Perugini [2013]  
978 suggest 250 participants as a reasonable sample size for interpreting single correlation  
979 coefficients), large sample sizes have been suggested for reliability analyses (e.g., Charter,  
980 1999). For example, Charter suggests at least 400 participants in order to conduct reliability  
981 analyses. This is particularly important for studies that aim to estimate a reliability coefficient  
982 which can be interpreted largely independent of the respective sample, for example when  
983 reporting *the* reliability of a certain questionnaire, as is often done in questionnaire construction.  
984 For this kind of analysis, samples are required that are representative of an underlying  
985 population. Still, reliability analyses are possible with fewer participants. They simply lead to  
986 somewhat less precise estimates. For example, with 250 participants, the 95% confidence  
987 interval for a correlation  $r = .85$  ranges from .812 to .881, whereas for 200 participants it ranges  
988 from .806 to .884. Even for 100 participants, the respective CI still ranges from .785 to .897. In  
989 light of these considerations, and given that we had to conduct single-participant sessions, we  
990 considered 200 participants per variable to be acceptable for our approach.

991 **The role of sample homogeneity or heterogeneity for reliability**

992           Whereas sample size does not affect the reliability estimate itself, (i.e., larger sample  
993 sizes do not lead to larger or smaller reliability estimates), true score variation in the respective  
994 sample does (e.g., Bühner, 2011; Steyer & Eid, 1993). That means that more heterogeneous  
995 samples may lead to higher reliability estimates than more homogenous ones. This follows  
996 directly from Equation 2: When variation in the true score increases more strongly than variation  
997 due to measurement error, reliability increases. This observation has some consequences: First,  
998 one and the same measurement instrument may have different reliabilities for different kinds of  
999 samples (or populations, respectively, from which these samples are drawn). For example, an  
1000 instrument assessing political attitudes may be more reliable in a moderate sample (where people  
1001 have different political attitudes) than in an extremist sample (where people have rather similar  
1002 attitudes) (see Danner, 2015, for this example). This transfers to applications in sport  
1003 psychology: When used in a high-performance sample (where there is low variation in athletes'  
1004 performance), a measurement instrument might have lower reliability than when used in a  
1005 sample of more moderate performance (where there is substantial variation in athletes'  
1006 performance). Second, a certain reliability that was estimated based on a representative sample  
1007 of the population may not apply to a more homogenous subsample of that same population.

1008           When one argues that larger samples are more likely to be heterogeneous, then it follows  
1009 that increasing reliability may be an indirect consequence of increasing sample sizes. However,  
1010 this only holds when heterogeneity increases (more technically, as described above, when due to  
1011 increased heterogeneity true score variation increases more strongly than error variation).

#### 1012 **Restrictions and assumptions underlying reliability in CTT**

1013           Importantly, Classical Test Theory can only be applied to measurements that produce  
1014 interval-scale (i.e., continuous) data (Bühner, 2011; Steyer & Eid, 1993; for a disagreement with

1015 this position see Gaito, 1980). The reason for this restriction is that reliability is defined as a  
1016 proportion of variances (see Equation 2), and variances can only be calculated for continuous  
1017 data. Furthermore (following Equation 1), measurement error is defined as the difference  
1018 between the observed value and the true score ( $Y_{\text{error}} = Y_{\text{observed}} - T_{\text{true}}$ ), which again is only  
1019 possible for continuous measurements. This restriction can be misunderstood as meaning that  
1020 reliability coefficients can only be *calculated* for continuous data, however its consequences are  
1021 more far reaching: Indeed, reliability according to CTT is only *defined* for continuous  
1022 measurements. It follows that when one wants to estimate reliability coefficients for a certain  
1023 measurement procedure, this procedure must yield continuous measurement outcomes.

1024 The core of CTT are three definitions, sometimes also called axioms. The first definition  
1025 states that every observed value ( $Y_{\text{observed}}$ ) consists of a true value ( $T_{\text{true}}$ ) and random  
1026 measurement error ( $Y_{\text{error}}$ ). That is,  $Y_{\text{observed}} = T_{\text{true}} + Y_{\text{error}}$ . The second definition states that  
1027 measurement error has an expectancy value of 0 and a finite variance. The third one states that a)  
1028 measurement error of a test  $t_1$  is independent from this test's true values, b) measurement error of  
1029 one test  $t_1$  is independent from measurement error of another test  $t_2$ , and c) that measurement  
1030 error of a test  $t_1$  is independent from the true values of another test  $t_2$ .

1031 Relatedly, CTT contains five models, that describe assumptions that are necessary for  
1032 estimating reliability (Bühner, 2011; Steyer & Eid, 1993). These five models are a) the model of  
1033 parallel measurement, b) the model of essentially parallel measurement, c) the model of tau-  
1034 equivalent measurement, d) the model of essentially tau-equivalent measurement, and e) the  
1035 model of tau-congeneric measurement. These models contain assumptions regarding the true  
1036 scores and (the intercorrelations of) measurement error.

1037           When accepting the axioms of CTT, and when the above described modelling  
1038 assumptions hold, it can be shown that the different reliability coefficients are estimates of the  
1039 measurement instrument's reliability. That is, their respective formulas can be converted (only if  
1040 one assumes that the axioms hold) into the definition of reliability according to Equation 2 (e.g.,  
1041 see Steyer & Eid, 1993). However, when the modelling assumptions do not hold, estimates can  
1042 either under- or overestimate a measurement's reliability (Savalei & Reise, 2019; Steyer & Eid,  
1043 1993). The exact nature of the deviation depends on the exact nature of the violation of the  
1044 assumptions. In cases of extreme violations of the assumptions, reliability estimates can become  
1045 entirely meaningless and unrelated to a measurement's true reliability (Steyer & Eid, 1993).  
1046 Different reliability coefficients require different modelling assumptions.

#### 1047 **Some considerations on different reliability coefficients**

1048           Whereas a measurement only has one reliability (defined by Equation 2), this reliability  
1049 can be assessed or estimated via different reliability coefficients. Reliability itself (and not only  
1050 the estimate) is sample dependant. That means that one and the same measurement may have  
1051 different reliabilities for different samples. When reliability was estimated using a representative  
1052 sample, one may assume that the same reliability holds for samples that are either a) also  
1053 representative or b) drawn randomly from the same population and sufficiently large.

1054           So, why are there several reliability coefficients? First, as mentioned above, different  
1055 reliability coefficients require different modelling assumptions, and only when these are met can  
1056 the respective coefficients be used in order to estimate reliability.

1057           Second, there are some conceptual and practical considerations. Using *test-retest*  
1058 *reliability* only makes sense when the construct to-be-measured is stable between the two  
1059 measurement points and when the measurement is not (differentially) affected by memory effects



1060 or learning. When test-retest reliability is calculated although these conditions are not met, the  
1061 resulting coefficient will underestimate an instrument's reliability. Using *parallel-test reliability*  
1062 only makes sense when two absolutely parallel tests exist for measuring the same construct, in  
1063 other words when two tests exist that measure the same construct with different items. When  
1064 parallel-test reliability is calculated although these conditions are not met, the resulting  
1065 coefficient will underestimate an instrument's reliability. In case that the above mentioned  
1066 conditions are not met, *split-half reliabilities* can be estimated. These tend to be higher the more  
1067 homogenous the measurement instrument is. Calculating split-half reliabilities requires to divide  
1068 all items of the measurement procedure into two equivalent halves. Subsequently, the correlation  
1069 between the two halves is calculated. This only makes sense when the two halves are indeed  
1070 equivalent. In the light of these considerations, we considered split-half reliabilities, and  
1071 particularly odd-even reliabilities to be the most appropriate estimators for our measurements.

### 1072 **The role of the number of items for reliability**

1073 Reliability itself, and not only its estimate, increases when the number of items that  
1074 measure the same underlying construct increases. This is a property of reliability according to  
1075 CTT (Bühner, 2011; Steyer & Eid, 1993). Intuitively, it can be understood when one considers  
1076 that according to CTT, measurement errors cancel each other out, and the more measurements  
1077 one has, the more they cancel each other out. Mathematically, the relationship between the  
1078 length of a measurement procedure and reliability is described by the Spearman-Brown-Formula  
1079 (Bühner, 2011; Steyer & Eid, 1993).

1080 Equation 3: Spearman-Brown-Formula.  $r_{tt}$  is the reliability of the current form of the  
1081 measurement procedure;  $k$  is the factor by which the length of the current measurement  
1082 procedure changes (e.g.,  $k = 0.5$  means half the number of items,  $k = 2$  means twice the number

1083 of items,  $k = 3$  means three times the number of items, and so on);  $r_{ttcorr}$  is the reliability of the  
1084 changed form of the measurement procedure.

$$1085 \quad r_{ttcorr} = \frac{k \cdot r_{tt}}{1 - (k - 1) \cdot r_{tt}}$$

1086 The Spearman-Brown-Formula can be used to predict how reliability will change when  
1087 the number of items of an existing measurement procedure with a known reliability changes.  
1088 Thereby, it can be used to predict how many items researchers need to add (or subtract) in order  
1089 to achieve a certain reliability, once they already know the reliability of a measurement  
1090 procedure with a given number of items. As one can deduce from the Spearman-Brown-Formula,  
1091 increasing the number of items at first leads to relatively high increases in reliability, however,  
1092 further gains in reliability need increasingly more items (Amelang & Zielinski, 2002).

1093 Therefore, practically, researchers who want to obtain a reliable measure can increase the  
1094 number of items. It seems even possible to predefine a certain reliability one wants to achieve  
1095 (say, .90) and then to increase the number of items until this reliability is achieved. However,  
1096 there are three potential problems with this approach (Amelang & Zielinski, 2002). First, the  
1097 single items must all measure the same construct. This seems feasible for performance outcome  
1098 measures as discussed in this paper. However, it may be problematic for other kinds of  
1099 measurement procedures. For example, there may only be a limited number of items that are  
1100 suitable for assessing a certain construct in a questionnaire (e.g., there may only be a limited  
1101 number of items for assessing anxiety). Second, the relationship between reliability and economy  
1102 is an inverse one: When researchers increase reliability by increasing the number of items, they  
1103 also increase the time their measurement procedure takes. Both in applied and in research  
1104 contexts, time is usually limited (and important to consider; e.g. for motivational reasons). This  
1105 is the main reason why constructing short forms of widely used questionnaires has become

1106 common. Furthermore, when a measurement procedure takes more time, adverse effects such as  
1107 fatigue, boredom or concentration problems become more likely to influence the measurement  
1108 outcome, thus limiting the measurement's validity. Still, time constraints do not seem to be a  
1109 major limiting factor for increasing the number of trials in studies employing performance  
1110 outcome measures, as the time needed per trial (e.g., per putt, dart throw or free throw) is very  
1111 short. Third, there may be a seemingly paradoxical relationship between reliability and validity.  
1112 On the one hand, reliability is a prerequisite for validity. That means, a measurement that is not  
1113 reliable cannot be valid. On the other hand, increasing reliability by increasing the number of  
1114 items can decrease validity. The reason for this seemingly paradoxical observation is that adding  
1115 items in order to increase reliability is often likely to make a measure more homogenous. To the  
1116 extent that the construct one intends to assess is rather heterogenous, then, the measurement  
1117 becomes less valid. One way to try and circumvent this problem is by having a measurement  
1118 procedure with several subscales. Each subscale is rather homogenous and constructed in order  
1119 to be highly reliable, whereas the heterogeneity of the construct is captured by the multitude of  
1120 different subscales (e.g., intelligence or personality tests). However, this approach will make the  
1121 measurement procedure less economical again.

1122       Taken together, for the above mentioned reasons increasing reliability by adding items  
1123 only works to a certain extent for common psychological measurement procedures such as  
1124 questionnaires. For performance outcome measures, such as discussed in this paper, it seems to  
1125 be more promising. Therefore, if possible, we suggest that researchers pretest their performance  
1126 outcome measure and try to determine an optimal number of trials: A number, that provides  
1127 sufficient reliability, but that does not induce threats to validity (such as fatigue or learning  
1128 effects) and that is still economically feasible.

1129 As mentioned above, it is possible to use the Spearman-Brown-Formula in order to  
1130 estimate, based on a given reliability ( $r_{tt}$ ) for  $x$  items, how many additional items would be  
1131 necessary in order to achieve a prespecified reliability ( $r_{tt\text{ corr}}$ ). In order to do so, Equation 3 needs  
1132 to be solved for  $k$ , which yields Equation 4.

1133 Equation 4:

$$k = \frac{r_{tt\text{ corr}} \cdot (1 - r_{tt})}{r_{tt} \cdot (1 - r_{tt\text{ corr}})}$$

1135 Importantly,  $k$  is not the number of items but the factor, with which the original number  
1136 of items needs to be multiplied in order to achieve the prespecified reliability. That means, when  
1137  $x$  is the original number of items,  $k \cdot x$  is the new number of items.

1138

|                        |  |
|------------------------|--|
| Initial considerations | <ul style="list-style-type: none"><li>• Select a measurement procedure based on theoretical grounds and research goals.</li></ul>  |
| Pretest                | <ul style="list-style-type: none"><li>• Pretest the measurement's reliability for a specific number of trials.</li><li>• Use a sample that is drawn from the same population that you intend to draw your main study's sample from.</li><li>• When considering sample size, think about precision, not statistical significance (see, for example, Appendix A; Charter, 1999; Schönbrodt and Perugini, 2013).</li><li>• Consider different estimates of reliability: Which one is best suited for your measurement based on practical and statistical assumptions (see Appendix A for a brief overview)?</li><li>• Estimate reliability.</li></ul> |
| Main study             | <ul style="list-style-type: none"><li>• Based on the above estimate of reliability, calculate the number of trials that you need in order to achieve a certain level of reliability.</li><li>• Conduct your main study, and estimate reliability again.</li></ul>  |
| Future studies         | <ul style="list-style-type: none"><li>• Take into consideration that reliability is dependent (among other factors) on samples.</li><li>• Different measurements using the same instrument may therefore lead to different estimations of reliability, for example when samples differ regarding true score variation.</li></ul>   |

- Reliable measurement plays an important yet underrated role for research quality
- Performance outcome measures used in sport psychology differ regarding reliability
- Reliability of performance outcome measures depends on sample characteristics
- Reliability of performance outcome measures depends on item number
- 

Journal Pre-proof

Author statement

Geoffrey Schweizer: Conceptualization; Formal analysis; Supervision; Writing - Original Draft; Writing - Review & Editing

Philip Furley: Conceptualization; Writing - Original Draft; Writing - Review & Editing

Nicolas Rost: Formal analysis ; Investigation; Writing - Review & Editing

Kai-Eric Barth: Formal analysis; Investigation; Writing - Review & Editing

Journal Pre-proof

Declarations of interest: none.

Journal Pre-proof