

1 **Manipulating the Alpha Level Cannot Cure Significance Testing**  
2 **Comments on “Redefine Statistical Significance”**

3  
4 David Trafimow<sup>1a</sup>, Valentin Amrhein<sup>2</sup>, Corson N. Areshenkoff<sup>3</sup>, Carlos Barrera-Causil<sup>4</sup>, Eric J.  
5 Beh<sup>5</sup>, Yusuf Bilgiç<sup>6</sup>, Roser Bono<sup>7</sup>, Michael T. Bradley<sup>8</sup>, William M. Briggs<sup>9</sup>, Héctor A. Cepeda-  
6 Freyre<sup>10</sup>, Sergio E. Chaigneau<sup>11</sup>, Daniel R. Ciocca<sup>12</sup>, Juan Carlos Correa<sup>13</sup>, Denis Cousineau<sup>14</sup>,  
7 Michiel R. de Boer<sup>15</sup>, Subhra Sankar Dhar<sup>16</sup>, Igor Dolgov<sup>1</sup>, Juana Gómez-Benito<sup>7</sup>, Marian  
8 Grendar<sup>17</sup>, James Grice<sup>18</sup>, Martin E. Guerrero-Gimenez<sup>12</sup>, Andrés Gutiérrez<sup>19</sup>, Tania B. Huedo-  
9 Medina<sup>20</sup>, Klaus Jaffe<sup>21</sup>, Armina Janyan<sup>22,23</sup>, Ali Karimnezhad<sup>24</sup>, Fränzi Korner-Nievergelt<sup>25</sup>,  
10 Koji Kosugi<sup>26</sup>, Martin Lachmair<sup>27</sup>, Rubén Ledesma<sup>28</sup>, Roberto Limongi<sup>29</sup>, Marco Tullio Liuzza<sup>30</sup>,  
11 Rosaria Lombardo<sup>31</sup>, Michael Marks<sup>1</sup>, Gunther Meinlschmidt<sup>32,33,34</sup>, Ladislav Nalborczyk<sup>35,36</sup>,  
12 Hung T. Nguyen<sup>37</sup>, Raydonal Ospina<sup>38</sup>, Jose D. Perezgonzalez<sup>39</sup>, Roland Pfister<sup>40</sup>, Juan José  
13 Rahona<sup>27</sup>, David A. Rodríguez-Medina<sup>41</sup>, Xavier Romão<sup>42</sup>, Susana Ruiz-Fernández<sup>27</sup>, Isabel  
14 Suarez<sup>43</sup>, Marion Tegethoff<sup>44</sup>, Mauricio Tejo<sup>45</sup>, Rens van de Schoot<sup>46,47</sup>, Ivan Vankov<sup>22</sup>,  
15 Santiago Velasco-Forero<sup>48</sup>, Tonghui Wang<sup>49</sup>, Yuki Yamada<sup>50</sup>, Felipe C. M. Zoppino<sup>12</sup> &  
16 Fernando Marmolejo-Ramos<sup>51</sup>

17  
18 Authorship order is alphabetical, except for the first, second, and last author.

- 19  
20 1. Department of Psychology, New Mexico State University, U.S.A.  
21 2. Zoological Institute, University of Basel, Basel, Switzerland  
22 3. Centre for Neuroscience Studies, Queens University, Ontario, Canada  
23 4. Faculty of Applied and Exact Sciences, Metropolitan Technological Institute, Medellín,  
24 Colombia  
25 5. School of Mathematical & Physical Sciences, University of Newcastle, Australia  
26 6. Department of Mathematics, State University of New York at Geneseo, U.S.A.  
27 7. Quantitative Psychology Unit, Faculty of Psychology, University of Barcelona, Barcelona,  
28 Spain  
29 8. Department of Psychology, Faculty of Arts, University of New Brunswick, Canada  
30 9. Independent Researcher, New York, U.S.A.  
31 10. School of Psychology, Benemérita Universidad Autónoma de Puebla, México

- 32 11. Center for Cognition Research, CINCO, School of Psychology, Universidad Adolfo Ibáñez,  
33 Santiago, Chile
- 34 12. Oncology Laboratory, IMBECU, CCT CONICET Mendoza, Argentina
- 35 13. School of Statistics, Faculty of Sciences, National University of Colombia, Medellín,  
36 Colombia
- 37 14. School of Psychology, University of Ottawa, Ottawa, Canada
- 38 15. Department of Health Sciences, Vrije Universiteit Amsterdam and Amsterdam Public Health  
39 research institute, Amsterdam, The Netherlands
- 40 16. Department of Mathematics and Statistics, IIT Kanpur, India
- 41 17. Biomedical Center Martin, Jessenius Faculty of Medicine, Comenius University, Slovakia,  
42 and Institute of Measurement Science, Slovak Academy of Sciences, Slovakia
- 43 18. Department of Psychology, Oklahoma State University, U.S.A.
- 44 19. Faculty of Statistics, Saint Thomas University, Colombia
- 45 20. Department of Allied Health Sciences, College of Health, Agriculture, and Natural  
46 Resources, University of Connecticut, U.S.A.
- 47 21. Simón Bolívar University, Caracas, Venezuela
- 48 22. Department of Cognitive Science and Psychology, New Bulgarian University, Sofia,  
49 Bulgaria
- 50 23. National Research Tomsk State University, Tomsk, Russia
- 51 24. Department of Biochemistry, Microbiology, and Immunology, University of Ottawa, Ottawa,  
52 Canada
- 53 25. Oikostat GmbH, Ettiswil, Switzerland
- 54 26. Department of Education, Yamaguchi University, Japan
- 55 27. Leibniz Knowledge Media Research Center, Multimodal Interaction Lab, Tübingen,  
56 Germany
- 57 28. Consejo Nacional de Investigaciones Científicas y Técnicas, and Universidad Nacional de  
58 Mar del Plata, Argentina
- 59 29. School of Pedagogy, Pontifical Catholic University of Valparaíso, Chile
- 60 30. Department of Medical and Surgical Sciences, “Magna Graecia” University of Catanzaro,  
61 Catanzaro, Italy
- 62 31. Economics Department, University of Campania “Luigi Vanvitelli”, Capua, Italy

- 63 32. Department of Psychosomatic Medicine, University Hospital Basel and University of Basel,  
64 Basel, Switzerland
- 65 33. Division of Clinical Psychology and Cognitive Behavioral Therapy, International  
66 Psychoanalytic University, Berlin, Germany
- 67 34. Division of Clinical Psychology and Epidemiology, Department of Psychology, University of  
68 Basel, Basel, Switzerland
- 69 35. Univ. Grenoble Alpes, CNRS, LPNC UMR 5105, F-38000, Grenoble, France
- 70 36. Department of Experimental Clinical and Health Psychology, Ghent University, Belgium
- 71 37. Faculty of Economics, Chiang Mai University, Thailand
- 72 38. Department of Statistics, Computational Statistics Laboratory (CAST), Universidade Federal  
73 de Pernambuco, Brazil
- 74 39. Business School, Massey University, New Zealand
- 75 40. Department of Psychology III, University of Würzburg, Germany
- 76 41. School of Psychology, National Autonomous University of Mexico, Mexico
- 77 42. CONSTRUCT-LESE, Faculty of Engineering, University of Porto, Portugal
- 78 43. Department of Psychology, Universidad del Norte, Barranquilla, Colombia
- 79 44. Division of Clinical Psychology and Psychiatry, Department of Psychology, University of  
80 Basel, Basel, Switzerland
- 81 45. Facultad de Ciencias Naturales y Exactas, Universidad de Playa Ancha, Valparaíso, Chile
- 82 46. Utrecht University, Faculty of Social and Behavioural Sciences, Department of Methods and  
83 Statistics, Utrecht, The Netherlands
- 84 47. North-West University, Optentia Research Focus Area, Vanderbijlpark, South Africa
- 85 48. MINES Paristech, PSL Research University, Centre for Mathematical Morphology, France
- 86 49. Department of Mathematical Sciences, New Mexico State University, U.S.A.
- 87 50. Faculty of Arts and Science, Kyushu University, Japan
- 88 51. School of Psychology, The University of Adelaide, Australia
- 89
- 90 ☐ Corresponding author: David Trafimow, Department of Psychology, New Mexico State  
91 University, U.S.A. Email: dtrafimo@nmsu.edu

92 **Acknowledgments:** GM has been acting as consultant for Janssen Research & Development,  
93 LLC. MG acknowledges support from VEGA 2/0047/15 grant. RvdS was supported by a grant  
94 from the Netherlands organization for scientific research: NWO-VIDI-45-14-006.

95

96

97 **One sentence summary:** We argue that depending on  $p$ -values to reject null hypotheses,  
98 including a recent call for changing the canonical alpha level for statistical significance from .05  
99 to .005, is deleterious for the finding of new discoveries and the progress of cumulative science.

100

101

102 Many researchers have criticized null hypothesis significance testing, though many have  
103 defended it too (see Balluerka, Gómez, & Hidalgo, 2005 for a review). Sometimes, there is a  
104 recommendation that the alpha level be reduced to a more conservative value, to reduce the Type  
105 I error rate. For example, Melton (1962), the editor of *Journal of Experimental Social*  
106 *Psychology* from 1950–1962, favored an alpha level of .01 over the typical .05 alpha level. More  
107 recently, Benjamin and 71 scientists (2017) recommended shifting to .005—consistent with  
108 Melton’s comment that even the .01 level might not be “sufficiently impressive” to warrant  
109 publication (p. 554). In addition, Benjamin et al. (2017) stipulated that the .005 criterion should  
110 be for new findings but were vague about what to do with findings that are not new. Though not  
111 necessarily endorsing significance testing as the preferred inferential statistical procedure,<sup>1</sup>  
112 Benjamin et al. (2017) did argue that using a .005 criterion would fix much of what is wrong  
113 with significance testing. Unfortunately, as we will demonstrate, the problems with significance  
114 tests cannot be importantly mitigated merely by having a more conservative rejection criterion,  
115 and some problems are exacerbated by adopting a more conservative criterion.

116 We commence with some claims on the part of Benjamin et al. (2017). For example, they  
117 wrote “...changing the  $p$ -value threshold is simple, aligns with the training undertaken by many  
118 researchers, and might quickly achieve broad acceptance.” If significance testing—at any  $p$ -  
119 value threshold—is as badly flawed as we will maintain it is (see also Amrhein, Korner-  
120 Nievergelt, & Roth, 2017; Greenland, 2017), these reasons are clearly insufficient to justify it.

---

<sup>1</sup> Many of the authors favor Bayesian procedures.

121 Consider another claim: “The new significance threshold will help researchers and readers to  
122 understand and communicate evidence more accurately.” But if researchers have understanding  
123 and communication problems with a .05 threshold, it is unclear how using a .005 threshold will  
124 eliminate these problems. And consider yet another claim: “Authors and readers can themselves  
125 take the initiative by describing and interpreting results more appropriately in light of the new  
126 proposed definition of statistical significance.” Again, it is not clear how adopting a .005  
127 threshold will allow authors and readers to take the initiative with respect to better data  
128 interpretation. Thus, even prior to a discussion of our main arguments, there is reason for the  
129 reader to be suspicious of hasty claims with insufficient support.

130 With the foregoing out of the way, consider that a basic problem with tests of  
131 significance is that the goal is to reject the null hypothesis. This goal seems to demand—if one is  
132 a Bayesian—that the posterior probability of the null hypothesis should be low given the  
133 obtained finding. But the  $p$ -value one obtains is the probability of the finding (or a more extreme  
134 finding) given the null hypothesis (and the assumptions underlying the test), and one would need  
135 to make an invalid inverse inference to draw a conclusion about the probability of the null  
136 hypothesis given the finding. And if one is a frequentist, there is no way to traverse the logical  
137 gap from the probability of the finding given the null hypothesis to a decision about whether one  
138 should accept or reject the null hypothesis (Briggs, 2016; Trafimow, 2017). We accept that, by  
139 frequentist logic, the probability of a Type I error really is lower if  $p = .005$  than if  $p = .05$ , all  
140 else being equal. We also accept the Bayesian argument by Benjamin et al. (2017) that the null  
141 hypothesis is less likely if  $p = .005$  than if  $p = .05$ , all else being equal (although determining  $p$ -  
142 values via Bayes Factors is problematic; see Appendix).<sup>2</sup> Finally, we acknowledge that Benjamin  
143 et al. (2017) performed a service for science by further stimulating debate about significance  
144 testing. But there are important issues Benjamin et al. (2017) seem not to have considered,  
145 discussed in the following sections.

---

<sup>2</sup> Depaoli and van de Schoot (2017) provided a critique showing how it is possible to abuse Bayesian statistics, and provided potential solutions to such abuse. Konijn, van de Schoot, Winter, and Ferguson (2015) suggested a way to use Bayesian statistics to reduce publication bias.

146 *Regression and Reliability*

147       Trafimow and Earp (2017) argued against the general notion of setting an alpha level to  
148 make decisions to reject or not reject null hypotheses, and the arguments retain their force even if  
149 the alpha level is reduced to .005. In some ways, the reduction worsens matters. One problem is  
150 that  $p$ -values have a sampling distribution,<sup>3</sup> as do other statistics (Cumming, 2012). Whether the  
151  $p$ -value obtained in any experiment passes the alpha level is partly a matter of luck (which  $p$ -  
152 value one happens to sample), with the caveat that large effect and sample sizes, and small  
153 variation, should decrease  $p$ -values. Absent the caveat, the researcher is unlikely to re-sample a  
154  $p$ -value below a significance threshold upon replication, as there may be many more  $p$ -values  
155 above than below the threshold in the  $p$ -value distribution. Thus, the phenomenon of regression  
156 to the mean suggests that the  $p$ -value obtained in a replication experiment is likely to regress to  
157 whatever the mean  $p$ -value would be if many replications were performed to obtain a distribution  
158 of  $p$ -values for the experiment. How much regression should occur? That depends on the  
159 reliability of  $p$ -values.

160       Based on data placed online by the Open Science Collaboration (2015;  
161 <https://osf.io/fgjvw>), Trafimow and de Boer (2017) calculated the correlation between  $p$ -values  
162 obtained in the original cohort of studies with  $p$ -values obtained in the replication cohort, and  
163 obtained the dismal value of .004.<sup>4</sup> Clearly, then, the obtained  $p$ -value in the original study has  
164 little to do with the  $p$ -value obtained in a replication experiment. The best prediction would be a  
165  $p$ -value for the replication experiment being vastly closer to the mean of the  $p$ -value distribution  
166 than to the  $p$ -value obtained in the original experiment. Under the null hypothesis, the lower the  
167  $p$ -value published in the original experiment (e.g., .005 rather than .05), the greater the amount of  
168 distance of the  $p$ -value from the  $p$ -value mean, implying increased regression to the mean.<sup>5</sup> Thus,  
169 even using the .05 value is problematic, with exacerbation using the .005 value (Amrhein &  
170 Greenland, 2017). When studies have low power, it is not rare to obtain large sample effects that

---

<sup>3</sup> For a test of the difference between two normal means, the  $p$ -value is uniformly distributed on [0,1] under the point null hypothesis. Under a range alternative hypothesis, the distribution may be unknowable.

<sup>4</sup> There are several possible reasons for the low value. These could include the nonlinear relation between  $p$ -values and effect sizes, mixing cases where the null hypothesis is true (or close to true) with cases where it is not, publication bias, and imperfect replication methodology, as well as random sampling error.

<sup>5</sup> Recall, the  $p$ -value distribution under the alternative hypothesis often is not knowable.

171 are overestimates, and using the .005 threshold instead of .05 would guarantee that statistically  
172 significant results are even larger overestimates of population effect sizes (Button et al. 2013).

173 In addition, from a measurement point of view, where reliability is a prerequisite for  
174 validity, the  $p$ -value correlation (reliability) of .004 obtained by Trafimow and de Boer (2017)  
175 indicates that as a basis for binary decisions,  $p$ -values are incapable of measuring anything  
176 validly, including the strength of the evidence (Fisher, 1925; 1973) or the severity of the test  
177 (Mayo, 1996).<sup>6</sup> This could be argued to be a good reason not to use  $p$ -values at all. Alternatively,  
178 the dismal  $p$ -value reliability as evidenced by the Open Science Collaboration could be  
179 attributed, in part, to the publication bias caused by having a publishing criterion (Locascio,  
180 2017a). But if one wishes to make such an attribution, although it provides a justification for  
181 using  $p$ -values in a hypothetical scientific universe where  $p$ -values are more reliable because of a  
182 lack of publication bias, the attribution provides yet another important reason to avoid publishing  
183 criteria based on  $p$ -values.

184

#### 185 *Type I and Type II Errors*

186 Another disadvantage of using any set criterion level for publication is that the relative  
187 importance of Type I and Type II errors might differ across studies within or between areas and  
188 researchers (Trafimow & Earp, 2017). Setting a blanket level of either .05 or .005, or anything  
189 else, forces researchers to pretend that the relative importance of Type I and Type II errors is  
190 constant.<sup>7</sup> Benjamin et al. (2017) pointed out that a few areas of science use very low criterion  
191 levels to justify their recommendation to reduce to the .005 level, but this justification seems to  
192 tacitly admit that a blanket level across many areas is undesirable. It seems obvious that a wide  
193 variety of factors can influence the relative importance of Type I and Type II errors, thereby  
194 rendering any blanket recommendation undesirable (indeed Miller & Ulrich, 2016, show how  
195 these and other factors have a direct bearing on the final research payoff). These factors may  
196 include the clarity of the theory or auxiliary assumptions, practical or applied concerns, or  
197 experimental rigor. There is an impressive literature attesting to the difficulties in setting a

---

<sup>6</sup> “Correcting” the correlation for attenuation due to restriction of range, in the original cohort of studies, increases the correlation to .01, which is still low.

<sup>7</sup> Another problem is that for different sample sizes the same  $p$ -value may imply a different extent of the evidence against the null hypothesis (Royall, 1986).



198 blanket recommendation (e.g., Buhl-Mortensen, 1996; Lemons, Shrader-Frechette, & Cranor,  
199 1997; Lemons & Victor, 2008; Lieberman & Cunningham, 2009; Mudge, Baker, Edge, &  
200 Houlahan, 2012; Myhr, 2010; Rice & Trafimow 2010). This argument is not a recommendation  
201 that every researcher should get to set her own criterion, as that has obvious problems too (as  
202 Trafimow & Earp, 2017, showed).<sup>8</sup> Rather, given that blanket and variable criterion levels both  
203 are problematic, it is sensible to dispense with significance testing altogether.

204

### 205 *Defining Replicability*

206 Yet another disadvantage pertains to what Benjamin et al. (2017) touted as the main  
207 advantage of their proposal, that published findings will be more replicable using the .005 than  
208 .05 alpha level. This depends on what is meant by “replicate” (see Lykken, 1968, for some  
209 definitions). If one insists on the same alpha level for the original study and the replication study,  
210 then we see no reason to believe that there will be more successful replications using the .005  
211 level than using the .05 level. In fact, the statistical regression argument made earlier suggests  
212 that the regression issue is made even worse using .005 than using .05. Alternatively, as  
213 Benjamin et al. (2017) seem to suggest, one could use .005 for the original study and .05 for the  
214 replication study. In this case, we agree that the combination of .005 and .05 will create fewer  
215 unsuccessful replications than the combination of .05 and .05, for the initial and replication  
216 studies, respectively. However, this comes at a high price in arbitrariness. Suppose that two  
217 studies come in at  $p < .005$  and  $p < .05$ , respectively. This would count as a successful  
218 replication. In contrast, suppose that the two studies come in at  $p < .05$  and  $p < .005$ ,  
219 respectively. Only the second study would count, and the combination would not qualify as  
220 indicating a successful replication. The arbitrariness of declaring the combination of .005 and .05  
221 as being a successful replication, whereas the combination of .05 and .005 is not, adds to the  
222 myriad difficulties researchers have interpreting their data. More generally, insisting that setting

---

<sup>8</sup> In addition to creating new issues of how researchers should decide on the criteria for each experiment, how editors and reviewers should evaluate different criteria proposed by different authors, and losing what many consider to be the point of NHST—which is to have a consistent threshold level across a scientific domain: with variable thresholds, many old problems with NHST remain unsolved, such as the problems of regression to the mean, unreliability of  $p$ -values, inflation of effect sizes, publication bias, and the general disadvantage of forcing decisions too quickly rather than considering cumulative evidence across experiments.



223 a criterion of .005 renders research more replicable demands much more specificity with respect  
224 to how to conceptualize replicability. In addition, we do not see a single replication success or  
225 failure as definitive. If one wishes to make a strong case for replication success or failure,  
226 multiple replication attempts are desirable.<sup>9</sup>

227

### 228 *Questioning the Assumptions*

229 The discussion thus far is under the pretense that the assumptions underlying the  
230 computation of  $p$ -values are true. But how likely is this? Berk and Freedman (2003) have made a  
231 strong case that the assumptions of random sampling from a population and independence are  
232 rarely true. The problems are particularly salient in the clinical sciences, where the falsity of the  
233 assumptions, as well as the divergences between statistical and clinical significance, are  
234 particularly obvious and dramatic (Bhardwaj, Camacho, Derrow, Fleischer, & Feldman 2004;  
235 Ferrill, Brown, & Kyle, 2010; Fethney, 2010; Page, 2014). The problem of likely false  
236 assumptions underlying the computation of  $p$ -values, in combination with the other problems  
237 already discussed, render the illusory garnering of truth from  $p$ -values yet more dramatic.

238

### 239 *The Population Effect Size*

240 Let us continue with the significance and replication issues, reverting to the pretense that  
241 significance testing assumptions are correct, while keeping in mind that this is unlikely. Consider  
242 that as matters now stand using tests of significance with the .05 criterion, the population effect  
243 size plays an important role both in obtaining statistical significance (all else being equal, the  
244 sample effect size will be larger if the population effect size is larger) and in obtaining statistical  
245 significance twice for a successful replication. Switching to the .005 criterion would not lessen  
246 the importance of the population effect size, and would increase its importance unless sample  
247 sizes increased substantially from those commonly used.<sup>10</sup> And there is good reason to reject that

---

<sup>9</sup> The present NHST focus should not detract from the importance of the quality of the theory and auxiliary assumptions for replication, as is attested to by recent successful replication studies in cognitive psychology (Zwaan et al., 2017) and social sciences (Mullinix et al., 2015).

<sup>10</sup> In addition, with an alpha level of .005, large effect sizes would be more important for publication, and researchers might lean much more towards “obvious” research than in testing creative ideas where there is more of a risk of weak effects and  $p$ -values that fail to meet the .005 bar.

248 replicability should depend on the population effect size. To see this quickly, consider one of the  
249 most important science experiments of all time, by Michelson and Morley (1887). They used  
250 their interferometer to test whether the universe is filled with a luminiferous ether that allows  
251 light to travel to Earth from the stars. Their sample effect size was very small, and physicists  
252 accept that the population effect size is zero because there is no luminiferous ether. Using  
253 traditional tests of significance with either a .05 or .005 criterion, replicating Michelson and  
254 Morley would be problematic (see Sawilowsky, 2003, for a discussion of this experiment in the  
255 context of hypothesis testing). And yet physicists consider the experiment to be highly replicable  
256 (see also Meehl, 1967).<sup>11</sup> More generally, an experiment's replicability should not depend on the  
257 population effect size. Any proposal that features  $p$ -value rejection criteria forces the replication  
258 probability to be impacted by the population effect size, and should be rejected.

259

#### 260 *Accuracy of Published Effect Sizes*

261 It is desirable that published facts in scientific literatures accurately reflect reality.  
262 Consider again the regression issue. The more stringent the criterion level for publishing, the  
263 more distance there is from a finding that passes the criterion to the mean, and so there is an  
264 increasing regression effect. Even at the .05 level, researchers have long recognized that  
265 published effect sizes likely do not reflect reality, or at least not the reality that would be seen if  
266 there were many replications of each experiment and all were published (see Briggs, 2016;  
267 Grice, 2017; Hyman, 2017; Kline, 2017; Locascio, 2017a; 2017b; and Marks, 2017 for a recent  
268 discussion of this problem). Under reasonable sample sizes and reasonable population effect  
269 sizes, it is the abnormally large sample effect sizes that result in  $p$ -values that meet the .05 (or  
270 .005) criterion, as is obvious from the standpoint of statistical regression. Moreover, with  
271 typically low sample sizes, statistically significant effects often require overestimates of  
272 population effect sizes. Effect size overestimation was empirically verified by the Open Science  
273 Collaboration project (2015), where the average effect size in the replication cohort of studies  
274 was dramatically reduced from the average effect size in the original cohort (from .403 to .197).  
275 Changing to a more stringent .005 criterion merely would result in yet worse effect size

---

<sup>11</sup> Very likely, a reason null results are so difficult to publish in sciences such as psychology is because the tradition of using  $p$ -value cutoffs is so ingrained. It would be well to terminate this tradition.

276 overestimation (Button et al. 2013). The importance of having published effect sizes accurately  
277 reflect population effect sizes contradicts the use of significance tests, at any criterion.

278

### 279 *Sample size and Alternatives to Significance Testing*

280 We stress that replication depends largely on sample size, but there are factors that  
281 interfere with researchers using the large sample sizes necessary for good sampling precision and  
282 replicability. In addition to the obvious costs of obtaining large sample sizes, there may be an  
283 underappreciation of how much sample size matters (Vankov, Bowers, & Munafo, 2014), of the  
284 importance of incentives to favor novelty over replicability (Nosek, Spies, & Motyl, 2012) and  
285 of a prevalent misperception that the complement of  $p$ -values measures replicability (Cohen,  
286 1994; Thompson, 1996; Greenland et al. 2016). A focus on sample size suggests an alternative to  
287 significance testing. Trafimow (2017; Trafimow & MacDonald, 2017) suggested a procedure as  
288 follows. The researcher specifies how close she wishes the sample statistics to be to their  
289 corresponding population parameters, and the desired probability of being that close. Trafimow's  
290 equations can be used to obtain the necessary sample size to meet specifications. The researcher  
291 then obtains the necessary sample size, computes the descriptive statistics, and takes them as  
292 accurate estimates of population parameters (provisionally on new data, of course).<sup>12</sup> This *a*  
293 *priori* procedure stresses (a) deciding what it takes to believe that the sample statistics are good  
294 estimates of the population parameters before data collection rather than afterwards, and (b)  
295 obtaining a large enough sample size to be confident that the obtained sample statistics really are  
296 within specified distances of corresponding population parameters. The procedure also does not  
297 promote publication bias because there is no cutoff for publication decisions.<sup>13</sup>

298 The larger point is that there are creative alternatives to significance testing that confront  
299 the sample size issue much more directly than significance testing does. The “statistical toolbox”  
300 (Gigerenzer & Marewski, 2015) further includes, for example, confidence intervals, equivalence

---

<sup>12</sup> An optimal way to obtain reliable estimation is via robust methods (Erceg-Hurn, Wilcox, & Keselman, 2013; Field & Wilcox, 2017; Huber, 1972; Portnoy & He, 2000; Rousseeuw, 1991; Tukey, 1979).

<sup>13</sup> The foregoing description may make the *a priori* procedure seem to be the same as traditional power analysis, but this is not so. First, the goal of traditional power analysis is to find the sample size needed to have a good chance of obtaining a statistically significant  $p$ -value. Second, traditional power analysis is strongly influenced by the expected effect size whereas this *a priori* procedure is completely uninfluenced by the expected effect size.

301 tests, alternative ways of dealing with  $p$ -values as continuous indices, Bayesian methods, or  
302 information criteria; but none of those tools should replace conventional significance testing as  
303 the new magic method giving clear-cut mechanical answers (Cohen, 1994). In fact, inference  
304 should not be based on single studies at all (Neyman & Pearson, 1933; Fisher, 1937; Greenland,  
305 2017), nor on replications from the same lab, but on cumulative evidence from multiple  
306 independent studies. It is desirable to obtain precise estimates in those studies, but the more  
307 important goal may be to publish also our wide confidence intervals and small effects, without  
308 which the cumulative evidence will be distorted (Amrhein, Korner-Nievergelt, & Roth, 2017;  
309 Amrhein & Greenland, 2017). Along these lines, Briggs (2016) argues for abandoning  
310 parameter-based inference and adopting purely predictive, and therefore verifiable, probability  
311 models, and Greenland (2017) sees “a dire need to get away from inferential statistics and hew  
312 more closely to descriptions of study procedures, data collection [...], and the resulting data.”  
313

#### 314 *Conclusion*

315         It seems appropriate to conclude with the basic issue that has been with us from the  
316 beginning. Should  $p$ -values and  $p$ -value thresholds be used as the main criterion for making  
317 publication decisions? The mere fact that researchers are concerned with replication, however it  
318 is conceptualized, indicates an appreciation that single studies are rarely definitive and rarely  
319 justify a final decision. Thus,  $p$ -value criteria may not be very sensible. A counterargument  
320 might be that researchers often make decisions about what to believe, and using  $p$ -value criteria  
321 formalize what otherwise would be an informal process. But this counterargument is too  
322 simplistic. When evaluating the strength of the evidence, sophisticated researchers consider, in  
323 an admittedly subjective way, theoretical considerations such as scope, explanatory breadth, and  
324 predictive power; the worth of the auxiliary assumptions connecting nonobservational terms in  
325 theories to observational terms in empirical hypotheses; the strength of the experimental design;  
326 or implications for applications. To boil all this down to a binary decision based on a  $p$ -value  
327 threshold of .05, .01, .005, or anything else, is not acceptable.

## References

- 328  
329  
330 Amrhein, V. & Greenland, S. (2017). Remove, rather than redefine, statistical significance.  
331 *Nature Hum. Behav.* **1**, 0224.
- 332 Amrhein, V., Korner-Nievergelt, F., & Roth, T. (2017). The earth is flat ( $p>0.05$ ): significance  
333 thresholds and the crisis of unreplicable research. *PeerJ* **5**, e3544.
- 334 Balluerka, N., Gómez, J., Hidalgo, D. (2005). The controversy over null hypothesis significance  
335 testing revisited. *Methodology* **1**, 55–77.
- 336 Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk, R.,  
337 Bollen, K. A., Brembs, B., Brown, L., Camerer, C., Cesarini, D., Chambers, C. D., Clyde,  
338 M., Cook, T. D., De Boeck, P., Dienes, Z., Dreber, A., Easwaran, K., Efferson, C., Fehr,  
339 E., Fidler, F., Field, A. P., Forster, M., George, E. I., Gonzalez, R., Goodman, S., Green,  
340 E., Green, D. P., Greenwald, A., Hadfield, J. D., Hedges, L. V., Held, L., Ho, T.-H.,  
341 Hoijtink, H., Jones, J. H., Hruschka, D. J., Imai, K., Imbens, G., Ioannidis, J. P. A., Jeon,  
342 M., Kirchler, M., Laibson, D., List, J., Little, R., Lupia, A., Machery, E., Maxwell, S. E.,  
343 McCarthy, M., Moore, D., Morgan, S. L., Munafó, M., Nakagawa, S., Nyhan, B., Parker,  
344 T. H., Pericchi, L., Perugini, M., Rouder, J., Rousseau, J., Savalei, V., Schönbrodt, F. D.,  
345 Sellke, T., Sinclair, B., Tingley, D., Van Zandt, T., Vazire, S., Watts, D. J., Winship, C.,  
346 Wolpert, R. L., Xie, Y., Young, C., Zinman, J., & Johnson, V. E. (2017). Redefine  
347 statistical significance. *Nature Hum. Behav.* **1**, 0189.
- 348 Berk, R. A., & Freedman, D. A. (2003). Statistical assumptions as empirical commitments. In T.  
349 G. Blomberg & S. Cohen (Eds). *Law, punishment, and social control: Essays in honor of*  
350 *Sheldon Messinger* (2<sup>nd</sup> Ed, pp. 235–254). Aldine de Gruyter.
- 351 Bhardwaj, S., Camacho, F., Derrow, A., Fleischer, A., & Feldman, S. (2004). Statistical  
352 significance and clinical relevance. *Archives of Dermatology* **140**, 1520–1523.
- 353 Briggs, W. M. (2016). *Uncertainty: The Soul of Modeling, Probability & Statistics*. New York:  
354 Springer.
- 355 Buhl-Mortensen, L. (1996). Type-II statistical errors in environmental science and the  
356 precautionary principle. *Marine Pollution Bulletin* **32**, 528–531.

- 357 Button, K. S., Ioannidis, J. P. A., Mokrysz, C., Nosek, B. A., Flint, J., Robinson, E. S. J., &  
358 Munafo, M. R. (2013). Power failure: why small sample size undermines the reliability of  
359 neuroscience. *Nature Reviews Neuroscience* **14**, 365376.
- 360 Cohen, J. (1994). The earth is round ( $p < .05$ ). *American Psychologist* **49**, 997–1003.
- 361 Cumming, G. (2012). *Understanding the new statistics: Effect sizes, confidence intervals, and*  
362 *meta-analysis*. New York: Routledge.
- 363 Depaoli, S., & van de Schoot, R. (2017). Improving transparency and replication in Bayesian  
364 statistics: The WAMBS-Checklist. *Psychological Methods* **22**, 240–261.
- 365 Erceg-Hurn, D. M., Wilcox, R. R., & Keselman, H. J. (2013). Robust statistical estimation. In T.  
366 Little (Ed.), *The Oxford Handbook of Quantitative Methods*, Vol. 1, 388–406. New York:  
367 Oxford University Press.
- 368 Ferrill, M., Brown, D., & Kyle, J. (2010). Clinical versus statistical significance: Interpreting p  
369 values and confidence intervals related to measures of association to decision making.  
370 *Journal of Pharmacy Practice* **23**, 344–351.
- 371 Fethney, J. (2010). Statistical and clinical significance, and how to use confidence intervals to  
372 help interpret both. *Australian Critical Care* **23**, 93–97.
- 373 Field, A. & Wilcox, R. (2017). Robust statistical methods: a primer for clinical psychology and  
374 experimental psychopathology researchers. *Behavior Research and Therapy* **98**, 19–38.
- 375 Fisher, R. A. (1925). *Statistical methods for research workers* (1<sup>th</sup> ed.). Edinburgh: Oliver and  
376 Boyd.
- 377 Fisher, R. A. (1937). *The design of experiments* (2<sup>nd</sup> ed.). Edinburgh: Oliver and Boyd.
- 378 Fisher, R. A. (1973). *Statistical methods and scientific inference* (3<sup>rd</sup> ed.). London: Macmillan.
- 379 Gigerenzer, G. & Marewski, J. N. (2015). Surrogate science: the idol of a universal method for  
380 scientific inference. *Journal of Management* **41**, 421–440.
- 381 Greenland, S. (2017). The need for cognitive science in methodology. *American Journal of*  
382 *Epidemiology* **186**, 639–645.
- 383 Greenland, S., Senn, S. J., Rothman, K. J., Carlin, J. B., Poole, C., Goodman, S. N., & Altman,  
384 D. G. (2016). Statistical tests, P values, confidence intervals, and power: a guide to  
385 misinterpretations. *European Journal of Epidemiology* **31**, 337–350.

- 386 Grice, J. W. (2017). Comment on Locascio's results blind manuscript evaluation proposal. *Basic*  
387 *and Applied Social Psychology*. In press.
- 388 Huber, P. J. (1972). Robust statistics: a review. *The Annals of Mathematical Statistics* **43**, 1041–  
389 1067.
- 390 Hyman, M. (2017). Can 'results blind manuscript evaluation' assuage 'publication bias'? *Basic*  
391 *and Applied Social Psychology*, In press.
- 392 Kline, R. (2017). Comment on Locascio, results blind science publishing. *Basic and Applied*  
393 *Social Psychology*, In press.
- 394 Konijn, E. A., van de Schoot, R., Winter, S. D., & Ferguson, C. J. (2015) Possible solution to  
395 publication bias through Bayesian statistics, including proper null hypothesis testing,  
396 *Communication Methods and Measures* **9**, 280–302.
- 397 Lemons, J., & Victor, R. (2008) Uncertainty in river restoration. In Darby, S., Sear, D. (Eds.),  
398 *River restoration: managing the uncertainty in restoring physical habitat*. John Wiley &  
399 Sons.
- 400 Lemons, J., Shrader-Frechette, K., Cranor, C. (1997) The precautionary principle: Scientific  
401 uncertainty and type I and type II errors. *Foundations of Science* **2**, 207–236.
- 402 Lieberman, M. D., Cunningham, W. A. (2009). Type I and Type II error concerns in fMRI  
403 research: re-balancing the scale. *Social Cognitive and Affective Neuroscience* **4**, 423–428.
- 404 Locascio, J. (2017a). Results blind publishing. *Basic and Applied Social Psychology*. In press.
- 405 Locascio, J. (2017b). Rejoinder to responses to "results blind publishing." *Basic and Applied*  
406 *Social Psychology*. In press.
- 407 Lykken, D. T. (1968). Statistical significance in psychological research. *Psychological Bulletin*  
408 **70**, 151–159.
- 409 Marks, M. J. (2017). Commentary on Locascio 2017. *Basic and Applied Social Psychology*. In  
410 press.
- 411 Mayo, D. (1996). *Error and the growth of experimental knowledge*. Chicago: The University of  
412 Chicago Press.
- 413 Meehl, P. E. (1967). Theory-testing in psychology and physics. A methodological paradox.  
414 *Philosophy of Science* **34**, 103–115.
- 415 Melton, A. (1962). Editorial. *Journal of Experimental Psychology* **64**, 553–557.



- 416 Michelson, A. A., & Morley, E. W. (1887). On the relative motion of earth and luminiferous  
417 ether. *American Journal of Science, Third Series*, *34*, 203, 233–245.  
418 <http://history.aip.org/exhibits/gap/PDF/michelson.pdf>
- 419 Miller, J., & Ulrich, R. (2016). Optimizing research payoff. *Perspectives on Psychological*  
420 *Science* **11**, 664–691.
- 421 Mudge, J.F., Baker, L.F., Edge, C.B., Houlahan, J.E. (2012). Setting an Optimal  $\alpha$  That  
422 Minimizes Errors in Null Hypothesis Significance Tests. *PLoS ONE* **7**, e32734.
- 423 Mullinix, K. J., Leeper, T. J., Druckman, J. N., & Freese, J. (2015). The generalizability of  
424 survey experiments. *Journal of Experimental Political Science* **2**, 109–138.
- 425 Myhr, A. I. (2010). A precautionary approach to genetically modified organisms: challenges and  
426 implications for policy and science. *Journal of Agricultural and Environmental Ethics*  
427 **23**, 501–525.
- 428 Neyman, J. & Pearson, E. S. (1933). On the problem of the most efficient tests of statistical  
429 hypotheses. *Philosophical Transactions of the Royal Society of London, Series A* **231**,  
430 289–337.
- 431 Nosek, B. A., Spies, J. R., & Motyl, M. (2012). Scientific utopia: II. Restructuring incentives and  
432 practices to promote truth over publishability. *Perspectives in Psychological Science* **7**,  
433 615–631.
- 434 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.  
435 *Science* **349** (6251). aac4716. doi: 10.1126/science.aac4716.
- 436 Page, P. (2014). Beyond statistical significance: Clinical interpretation of rehabilitation research  
437 literature. *The International Journal of Sports Physical Therapy* **9**, 72.
- 438 Portnoy, S., & He, X. (2000). A robust journey in the new millennium. *Journal of the American*  
439 *Statistical Association* **95**, 1331–1335.
- 440 Rice, S., & Trafimow, D. (2010). How many people have to die for a type II error? *Theoretical*  
441 *Issues in Ergonomics Science* **11**, 387–401.
- 442 Rousseeuw, P. J. (1991). Tutorial to robust statistics. *Journal of Chemometrics* **5**, 1–20.
- 443 Royall, R. M. (1986). The Effect of Sample Size on the Meaning of Significance Tests. *The*  
444 *American Statistician* **40**, 313–315.

- 445 Sawilowski, S. (2003). Deconstructing arguments from the case against hypothesis testing.  
446 *Journal of Modern Applied Statistical Methods* **2**, 467–474.
- 447 Thompson, B. (1996). AERA editorial policies regarding statistical significance testing: Three  
448 suggested reforms. *Educational Researcher* **25**, 26–30.
- 449 Trafimow, D. (2017). Using the coefficient of confidence to make the philosophical switch from  
450 *a posteriori* to *a priori* inferential statistics. *Educational and Psychological Measurement*  
451 **77**, 831–854.
- 452 Trafimow, D., & de Boer, M. (2017). Measuring the strength of the evidence. Under submission.
- 453 Trafimow, D., & Earp, B. D. (2017). Null hypothesis significance testing and the use of P values  
454 to control the Type I error rate: The domain problem. *New Ideas in Psychology* **45**, 19–  
455 27. <http://dx.doi.org/10.1016/j.newideapsych.2017.01.002>.
- 456 Trafimow, D., & MacDonald, J. A. (2017). Performing inferential statistics prior to data  
457 collection. *Educational and Psychological Measurement* **77**, 204–219.
- 458 Tukey, J. W. (1979). Robust techniques for the user. In R. L. Launer and G. N. Wilkinson (Eds.),  
459 *Robustness in Statistics* (pp. 103–106). Academic Press, New York.
- 460 Vankov, I., Bowers, J., & Munafò, M. R. (2014). On the persistence of low power in  
461 psychological science. *Quarterly Journal of Experimental Psychology* **67**, 1037–1040.
- 462 Zwaan, R., Pecher, D., Paolacci, G., Bouwmeester, S., Verkoeijen, P., Dijkstra, K., &  
463 Zeelenberg, R. (2017). Participant Nonnaiveté and the reproducibility of cognitive  
464 psychology. *Psychonomic Bulletin and Review*. [doi.org/10.3758/s13423-017-1348-y](https://doi.org/10.3758/s13423-017-1348-y)

## Appendix

465

466

467 The Bayes Factor is an approach to model selection that attempts to quantify the posterior  
468 probability of one model relative to another, given set of observed data (Kass & Raftery, 1995).  
469 Formally, given a model  $H$  and observed data  $D$ , the posterior probability of  $H$  is given by Bayes  
470 theorem:

$$471 \quad P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

472 where  $P(D|H)$  is the likelihood (determined by the statistical model), and  $P(H)$  is the prior on the  
473 model  $H$ . Given two competing models  $H_1$  and  $H_2$ , the ratio of the posterior probabilities is given  
474 by

$$475 \quad \frac{P(H_1|D)}{P(H_2|D)} = \frac{P(D|H_1)P(H_1)}{P(D|H_2)P(H_2)}$$

476 Here, the quantity

$$477 \quad B = \frac{P(D|H_1)}{P(D|H_2)}$$

478 which multiplies the ratio of prior probabilities to obtain the posterior, is traditionally called the  
479 Bayes factor, and is generally interpreted as reflecting the relative weight of evidence provided  
480 by the data for the models  $H_1$  and  $H_2$ . Quantities greater than 1 suggest that the data favor  $H_1$ ,  
481 while quantities less than 1 suggest support for  $H_2$ .

482 Ostensibly a form of model comparison, the Bayes factor has been widely employed in  
483 the sciences to perform null-hypothesis significance testing by specifying the null model  $H_1$  to be  
484 the alternative  $H_2$  with some parameter(s) set to zero (e.g. Wagenmakers, 2007; Wetzels &  
485 Wagenmakers, 2012; Ly, Verhagen, & Wagenmakers, 2016). In this way, Bayes factors find use  
486 as a sort of “Bayesian hypothesis test”. Although increasingly popular, Bayes factors have  
487 several problems which complicate their use as a form of hypothesis testing. Most notably,  
488 Bayes factors are highly sensitive to the choice of prior in a way that a true Bayesian analysis  
489 (one which returns a posterior distribution over the model parameters) is not, as we demonstrate  
490 below. Further, this prior sensitivity can often behave unintuitively; for instance, Gelman, et al  
491 (2014; pp. 183–184) provide an example in which the goal is to estimate a mean treatment effect  
492 within several groups. A Bayes factor is employed to compare a null model in which all groups  
493 have the same mean value, which has a normal prior; and a model with no shrinkage, in which

494 the means are independent draws from the same normal prior. In this case, the resulting Bayes  
495 factor is highly sensitive to the variance of the prior, and will always select the null model as the  
496 prior variance goes to infinity. A researcher, mistakenly believing that they are constructing a  
497 non-informative prior, might choose a very large prior variance, unknowingly forcing the Bayes  
498 factor to select the null model, regardless of what the data say.

499 Moreover, in a fully Bayesian model the likelihood begins to dominate the prior as the  
500 sample size goes to infinity, as we might expect (since a larger sample provides more  
501 information about the model parameters). The prior thus has less effect with increasing sample  
502 size. This is not, in general, true of Bayes factors, which retain their prior sensitivity even with  
503 large samples.

504 As an example, consider a set of data which are assumed to be Poisson distributed with  
505 rate  $\lambda$ , with competing models  $H_1: \lambda \leq 1$  vs  $H_2: \lambda > 1$ . A simulation study considering different  
506 sample sizes (20, 30 and 50) out with 10,000 replicates. The prior distributions for  $H_1$  and  $H_2$ ,  
507 and mean and standard deviations are presented in Table 1.

n	Prior $H_1$	Prior $H_2$	$mean(BF)$	$sd(BF)$
20	Gamma(1,2)	Gamma(3,3)	7.0	2.7
	Gamma(2,2)		202.0	74.7
	Gamma(3,2)		3726.9	1331.9
	Gamma(2,2)	Gamma(1,3)	11570.1	3139.9
		Gamma(2,3)	1496.0	472.5
		Gamma(4,3)	29.1	12.6
		Gamma(10,3)	<b>0.0019</b>	<b>0.0020</b>
	Laplace	Laplace	2.5	13.6
	Jeffreys	Jeffreys	<b>0.7</b>	<b>4.0</b>
	30	Gamma(1,2)	Gamma(3,3)	357.1
Gamma(2,2)			15626.6	5500.1
Gamma(3,2)			428190.9	151300.7
Gamma(2,2)		Gamma(1,3)	1936147.3	597773.2
		Gamma(2,3)	167656.2	54433.1
		Gamma(4,3)	1539.0	603.0
		Gamma(10,3)	<b>0.0113</b>	<b>0.0088</b>
Laplace		Laplace	1.7	11.5
Jeffreys		Jeffreys	<b>0.4</b>	<b>2.6</b>
50		Gamma(1,2)	Gamma(3,3)	2161996.0
	Gamma(2,2)		156358014.0	58652977.0
	Gamma(3,2)		7218780790.0	2682940095.0
	Gamma(2,2)	Gamma(1,3)	53119836572.0	18718219507.0
		Gamma(2,3)	2776659384.0	993973186.0
		Gamma(4,3)	9361043.0	3659097.0
		Gamma(10,3)	3.8	2.4
	Laplace	Laplace	<b>0.7</b>	<b>4.7</b>
	Jeffreys	Jeffreys	<b>0.1</b>	<b>0.8</b>

508

509 *Table 1: Mean and standard deviations for Bayes Factor (BF) when different prior distributions*510 *for the hypothesis  $H_1$  and  $H_2$  are considered. The BF's that provide support to  $H_2$  are shown in*511 *bold.*

512 It is evident that the BF changes considerably when prior distributions and sample sizes  
513 change. In many cases,  $H_1$  is more strongly supported by the data than  $H_2$ , which is not correct  
514 since  $\lambda = 1.5$ . And yet, the posterior distribution over  $\lambda$  is largely insensitive to these choices.  
515 For  $\lambda = 1.5$  and a sample size of 50, the expected sum over all observations for is 75. Since a  
516 Gamma prior is conjugate for a Poisson likelihood, we can compute the expected posterior  
517 directly: for prior shape  $\alpha$  and rate  $\beta$ , the expected posterior shape and rate are  $\alpha + 75$  and  $\beta +$   
518 50, giving a posterior mean of

$$519 \frac{(\alpha + 75)}{(\beta + 50)}.$$

520 For Gamma(1,3) and Gamma(2,3) priors, we have posterior means of 1.43 and 1.45,  
521 respectively—a negligible difference. And yet, the Bayes factors resulting from these priors  
522 differ by an order of magnitude.

523 Additionally, Bayes factors may exhibit strange behavior when used to test point-nulls  
524 for continuous models (e.g. testing that a mean difference is exactly zero). Aitkin, Boys, and  
525 Chadwick (2005) study an example in which a hypothesis test for a binomial probability,  
526 conducted via Bayes factor, returns strong support for a point null  $H_0: p = p_0$  when, in fact, the  
527 posterior distribution (and a data itself) overwhelmingly support a value  $p \neq p_0$ . The authors  
528 note that this is an example where two competing approaches—estimation and hypothesis  
529 testing—are in clear conflict. In general, researchers should be aware that the question they are  
530 asking—that is, do the data support the null value, or is *some other value* better supported by the  
531 data—is not necessarily the question being answered by the Bayes factor, nor is a Bayes factor  
532 used to test a point null on a parameter consistent in general with the posterior distribution over  
533 that same parameter. For these reasons, we do not feel that the Bayes factor is a satisfactory  
534 substitute for traditional hypothesis testing, nor does it address the fundamental problems  
535 associated with such approaches: namely, that they ignore uncertainty in favor of binary decision  
536 making. Using a threshold for the BF will result in a similar dilemma as with a threshold for the  
537  $p$ -value. As Konijn et al. (2015) suggested, “God would love a Bayes Factor of 3.01 nearly as  
538 much as a BF of 2.99”. Moreover, the Bayes factor does not provide a good measure of statistical  
539 evidence, as it fails the coherence desideratum (see Lavine & Schervish, 1999).

## References

- 540  
541
- 542 Aitkin, M., Boys, R. J., & Chadwick, T. (2005). Bayesian point null hypothesis testing via the  
543 posterior likelihood ratio. *Statistics and Computing* **15**, 217–230.
- 544 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B.  
545 (2014). *Bayesian data analysis, Third Edition*. Boca Raton, FL: CRC press.
- 546 Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical*  
547 *Association* **90**, 773–795.
- 548 Konijn, E. A., van de Schoot, R., Winter, S. D., & Ferguson, C. J. (2015) Possible solution to  
549 publication bias through Bayesian statistics, including proper null hypothesis testing.  
550 *Communication Methods and Measures* **9**, 280–302.
- 551 Lavine, M. and Schervish, M. J. (1999). Bayes factors: what they are and what they are not. *The*  
552 *American Statistician* **53**, 119–122.
- 553 Ly, A., Verhagen, J., & Wagenmakers, E. J. (2016). Harold Jeffreys’s default Bayes factor  
554 hypothesis tests: Explanation, extension, and application in psychology. *Journal of*  
555 *Mathematical Psychology* **72**, 19–32.
- 556 Wagenmakers, E. J. (2007). A practical solution to the pervasive problems of p  
557 values. *Psychonomic Bulletin & Review* **14**, 779–804.
- 558 Wetzels, R., & Wagenmakers, E. J. (2012). A default Bayesian hypothesis test for correlations  
559 and partial correlations. *Psychonomic Bulletin & Review* **19**, 1057–1064.