

1 Transparency in ecology and evolution: real problems, real solutions

2

3 Timothy H. Parker, Whitman College, Walla Walla, USA

4 Wolfgang Forstmeier, Max Planck Institute for Ornithology, Seewiesen, Germany

5 Julia Koricheva, Royal Holloway University of London, Egham, UK

6 Fiona Fidler, University of Melbourne, Melbourne, Australia

7 Jarrod Hadfield, University of Edinburgh, Edinburgh, UK

8 Yung En Chee, University of Melbourne, Melbourne, Australia

9 Clint Kelly, University of Quebec, Montreal, Canada

10 Jessica Gurevitch, Stony Brook University, New York, USA

11 Shinichi Nakagawa, University of New South Wales, Sydney, Australia

12

13 Keywords:

14 confirmation bias; inflated effect size; *p*-hacking; pre-registration; replication; selective reporting

15

16

17

18

19

20

21 Abstract

22

23 To make progress scientists need to know what other researchers have found and how they found it.
24 However, transparency is often insufficient across much of ecology and evolution. Researchers often fail
25 to report results and methods with detail sufficient to permit interpretation and meta-analysis, and
26 many results go entirely unreported. Further, these unreported results are often a biased subset. Thus
27 the conclusions we can draw from the published literature are themselves often biased and sometimes
28 might be entirely incorrect. Fortunately there is a movement across empirical disciplines, and now
29 within ecology and evolution, to shape editorial policies to better promote transparency. This can be
30 done by either requiring more disclosure by scientists or by developing incentives to encourage
31 disclosure.

32

33

34 Science needs transparency

35

36 Science is a uniquely effective tool for understanding the world, and ecologists and evolutionary
37 biologists have built a robust body of scientific knowledge over the past century. However, several
38 common practices are limiting progress in these fields. For science to progress, results and clear
39 explanations of methods must be shared with other scientists. Although this fundamental principle is
40 widely understood, practices that cloud transparency of methods and results, such as selective reporting
41 (see glossary), appear far more common than they should be. This is unlikely to be an issue of deliberate
42 dishonesty, which we assume is rare in ecology and evolution. Instead, we believe that the unintended
43 negative consequences of insufficient transparency are often unrecognized by many members of the
44 scientific community. In addition, the institutions that shape our choices often inadvertently encourage
45 or reward choices that obstruct transparency [1]. Without sufficient transparency, we are hindered in
46 our ability to interpret published findings, conclusions based on published literature can be biased or
47 wrong, and meta-analytical syntheses are weakened [2]. Although these challenges to transparency vary
48 across disciplines and sub-disciplines, evidence suggests they are often common and present very real
49 problems for the advancement of ecology and evolutionary biology. In this paper, we first review
50 evidence of insufficient transparency in ecology and evolutionary biology, and then discuss new efforts
51 in these fields and in empirical science in general to improve transparency and thus improve scientific
52 progress.

53

54

55 Evidence of insufficient transparency

56

57 Selective reporting - Once researchers have collected and analyzed data, they commonly publish only a
58 portion of the results derived from these data (Fig. 1). Such selective reporting can lead to publication
59 bias (see glossary) if researchers preferentially publish certain types of results, such as those with the
60 strongest or the most surprising patterns. However, selective reporting is not limited to the classic 'file-
61 drawer' problem in which a study that does not produce the hoped-for result goes unpublished (e.g.,
62 [3]). For instance, researchers might conduct multiple alternative forms of an analysis and report only
63 the one with the strongest relationships or lowest p -values. This practice has become known as ' p -
64 hacking' (see glossary) [4, 5]. P -hacking and other forms of selective reporting can be masked by
65 'HARKing', or Hypothesizing After Results are Known (see glossary)[6]. We might convince ourselves of
66 the validity of selective reporting in various ways. For instance, human cognitive tendencies, such as
67 confirmation bias (see glossary) (Box 1)[7], can lead researchers to select evidence that lends the
68 clearest support for a pre-existing hypothesis. Alternatively, selective reporting might not seem
69 problematic as researchers often tend to be more interested in the existence of patterns than in their
70 absence. However, ignoring weak, negative, or absent patterns is a major hindrance to our
71 understanding of the biological world. First, the absence of an effect or the presence of only a weak
72 effect is itself important as we sort through explanations of how biological systems work. Second, any
73 observed statistical relationship is an estimate of a true biological relationship, and as an estimate, it is
74 inherently uncertain. Sampling variance results in some estimates being higher than the true value, and
75 some lower (Type M errors; see glossary), and some being even opposite in sign (Type S error; see
76 glossary) [8]. If we systematically eliminate the smaller or contradictory effect sizes (see glossary) from
77 publication, we get a biased picture of the size of the true underlying effect, and under some
78 circumstances this bias can be extreme [2]. Methods exist for estimating the effect of publication bias in
79 meta-analysis, but these methods are imperfect because most are indirect and thus must make major
80 assumptions about missing unpublished results whose true values we can never know [9]. Therefore,
81 the clearest path towards a reliable average is minimizing bias in the original sample of statistical effects
82 [2]. The selective reporting behind much publication bias clearly varies among sub-disciplines and with
83 the type of data reported, but evidence suggests it is common in many areas of ecology and evolution,
84 as in many other scientific disciplines. Most authors of this manuscript have engaged in selective
85 reporting at one or more points in their pasts, sometimes at the request of reviewers or editors, and

86 anecdotal evidence from conversations with others suggest it could be widespread and frequent.
87 However, it is not just our personal experience that suggests selective reporting is common. There is
88 considerable published empirical evidence for publication bias in ecology and evolutionary biology.

89
90 Under-reporting (see glossary) is the easiest form of selective reporting to document because we know
91 the analysis was completed; the paper just fails to provide all the details of results or statistical methods.
92 For instance, studies sometimes include means with no indication of uncertainty around those means, p -
93 values with no indication of the direction of the trend, or statistical results without the sample size for
94 the particular subset of data examined. These practices all limit readers' abilities to build an unbiased
95 understanding of a system and severely limit the usefulness of data for meta-analysis. A long and
96 growing list of surveys and meta-analyses has documented widespread under-reporting across many of
97 our sub-disciplines. Studies in fields including conservation [10], plant ecology [11], behavioral ecology
98 [12], ecosystem ecology [13, 14], community ecology [15], and others [16, 17] often find that around
99 half of published articles lack at least one key piece of information regarding statistical relationships
100 (Table 1). Further, where it has been examined these under-reported results were more likely to come
101 from non-significant comparisons or patterns contradictory to the primary hypothesis [18]. Finally, even
102 if authors report statistical results, they often do not report how the analyses were conducted in
103 sufficient detail, which makes it impossible for readers to critique the statistical methodology and to
104 replicate the analyses.

105
106 Estimating the rate at which results go completely unreported is more challenging. Results could remain
107 hidden from comparisons that authors decided were uninteresting. Unreported results might also come
108 from alternative versions of analyses conducted with, for instance, different covariates, interactions, or
109 subsets of data, as we might expect from p -hacking. One proposed method for identifying p -hacking is
110 ' p -curve' analysis, which predicts a clumping of p -values just below 0.05 if p -hacking is common [5].
111 Recently p -curve analysis was used to argue that p -hacking was having only modest impacts on biology
112 [4]. Regrettably, this reassuring conclusion is unwarranted. First, when researchers can include or
113 exclude covariates depending on their effects on p -values, p -values much smaller than 0.05 can often be
114 generated in the absence of a real effect [19, 20]. Thus, p -curve analysis focused on a 0.05 threshold can
115 dramatically underestimate p -hacking in fields where multiple covariates are common [19], such as
116 much of ecology and evolutionary biology. In fact, p -values have been shown to clump under lower
117 thresholds (0.01, 0.001, etc.) as well [21], as would be expected if p -hacking often ended with
118 calculation of a "highly significant" p -value. However, the second problem with these analyses is that
119 assumptions about the expected distribution of a collection of published p -values are almost certainly
120 incorrect, and thus inferring bias from the ' p -curve' is untenable under most conditions [22].

121
122 There are, however, other ways to estimate the magnitude of selective reporting. We can compare rates
123 of publication of statistically significant results with the observed distribution of statistical power (see
124 glossary) and estimates of average strength of effect. Rates of publication of statistically significant
125 effects are very high. In "Environment/Ecology" and "Plant and Animal Sciences", 74% of 150 and 78% of
126 200 statistical tests, each from a different randomly selected paper, were statistically significant and
127 supported the researchers' putative *a priori* hypotheses [23]. Similarly, in a cross-section of biological
128 journals, many from the disciplines of ecology and evolution, only 8.6% presented non-significant tests
129 of the main hypothesis [24]. Part of the explanation for these numbers is likely to be HARKing, in which
130 authors choose their strongest patterns and build the paper around those results, either de-emphasizing
131 or leaving out other results. While in some sub-fields of ecology and evolution researchers might often
132 test hypotheses that are likely to be true, this is probably not the case across all of ecology and
133 evolution. Further, even if most of our hypotheses were true, the proportion of statistically significant
134 results should be much lower since many of our studies have low statistical power. This low power
135 results from sample sizes that are often small, and average effect sizes that are also relatively small ($|r|$
136 = 0.19 [25], which should actually be an overestimate [26]) and thus difficult to detect (Box 2). The
137 resulting statistical power to detect effects of this observed average magnitude in the behavior, ecology,

138 and evolution literature is in the neighborhood of 20% [27, 28] (Box 2). If we thus conclude that typical
139 power is about 20% and we assume that 74% of tested hypotheses are true, we would still expect only
140 16% of findings to be statistically significant (Box 3) rather than 74%. This is a strong indication of
141 HARKing and selective reporting. Further, we discuss evidence below which suggests that published
142 statistically significant results might often be false or inflated relative to the true effect.

143
144 Sources of bias - The proportion of significant results that are false positives is, somewhat counter-
145 intuitively, increased in studies with small samples and low power [29]. This increase happens because
146 the probability of detecting a true positive declines as power is reduced but the probability of detecting
147 a false positive remains fixed (typically at 0.05). As a consequence a greater proportion of positives will
148 be false as power decreases (Box 3). This means that reports of significant findings with low sample size
149 should be disproportionately likely to be incorrect [30], and of course such underpowered studies are
150 common in much of ecology and evolutionary biology [27].

151
152 Insufficient statistical power also hinders detection of real effects, and Type II errors (see glossary)
153 should thus also be common in ecology and evolution [31]. In fact, we predict that Type II error, when
154 they occur, will often go hand and hand with Type I error, as *p*-hacking extracts false positives from data
155 while true relationships go undetected. As described above, the rarity of negative results in the
156 literature suggests that Type II error is often concealed by HARKing, selective reporting, or both.

157
158 Much of our focus in this paper is on null hypothesis tests because these tests remain the most common
159 type of statistical analyses in ecology and evolution. However, it is important to note that most of the
160 choices related to sample size and selective reporting that can bias null hypothesis tests can bias other
161 threshold tests (e.g., Akaike information criterion: $\Delta AIC > 2$ [32]) and can also generate misleading and
162 inflated effect sizes. For instance, large effects reported from studies with small samples are likely to
163 often be inflated, or even of the wrong sign [30]. Examination of 3867 ecological studies from 52
164 previously published meta-analyses showed that studies with the largest effect sizes tended to have the
165 lowest samples sizes [33]. Further, '*p*-hacking' could also be considered 'effect-size hacking' since the
166 same practices produce inflated effect sizes, and if combined with selective reporting, produce a
167 distribution of published effects that is biased upwards.

168
169 Given that studies with larger effects could be more likely to end up in journals with higher impact
170 scores [34], perhaps high impact journals are often publishing studies with large effects despite their
171 small samples and unreliability. Although there is evidence that in some subsets of the published
172 literature sample size and journal impact factor are negatively correlated, this trend appears to vary
173 across study types, and when averaged across a large number of studies ($n = 3867$), impact factor was
174 uncorrelated with sample size [33]. While this lack of correlation is certainly better than a consistent
175 negative correlation, given that studies with larger samples produce more reliable results, it would
176 actually be preferable to see a positive relationship between sample size and journal impact factor.
177 Further, it is effect size, not sample size, that predicts the number of citations a study receives [33]. So,
178 not only are published studies with small sample sizes more likely to report inflated effects (i.e. more
179 prone to Type M errors), the unreliability of these studies does not dependably deter their publication in
180 high impact journals or their accumulation of citations.

181
182 It has long been established that as the number of statistical comparisons increases, the probability of
183 observing patterns that result only from chance (i.e., false positives) also increases [35]. This happens
184 both with multiple separate tests or if, instead of alternative tests, we combine multiple possible
185 predictors in the same model [36]. Within a single model we might include a set of different equally
186 plausible predictors of the variable of interest, or we might include multiple alternative interaction
187 terms between our predictor of interest and different covariates. In a survey of 50 randomly selected
188 studies from ecology and evolution, 28 studies (56%) used GLMs with two or more predictors [36], and
189 none of these 28 considered any type of correction for multiple comparisons to counter the risk of

190 inflated significance. We could not locate other attempts to quantify failures to correct for multiple
191 comparisons, but uncorrected multiple comparisons appear common in at least some portions of the
192 literature [12]. Although false positives from multiple comparisons in exploratory analyses need not be a
193 major problem if we recognize the provisional nature of the results [35], two current practices in our
194 disciplines make uncorrected multiple comparisons a severe issue. First, multiple comparisons are often
195 hidden, with researchers conducting multiple tests but only reporting a subset of them. Thus the
196 likelihood that a result is a false positive is concealed and the scientific community is misled about the
197 probability that the result is true. Second, calls for tolerating a high false positive rate (to reduce Type II
198 errors) emphasize the importance of validating findings with replication studies [35], but replications or
199 other types of independent evaluation are currently far too rare to sort out the false from the true
200 positives [37, 38].

201
202 The role of institutions - The problems outlined above are heavily influenced by the institutions that
203 shape the decisions of researchers, including journals, funding bodies, and employers. Calls for
204 individual scientists to improve transparency are not uncommon [e.g., 39, 40, 41], and scientists
205 sometimes respond to these calls. However, individual scientists, like all people, make decisions in
206 response to the institutions in which they operate [1]. Funders reward novelty, typically to the complete
207 exclusion of replication, and journals preferentially publish statistically significant findings, especially if
208 those findings are surprising. These factors alone would influence researchers' decisions, but these
209 incentives are even more influential because universities and research institutes often hire and promote
210 scientists based on their record of acquiring grant money and the number and impact factors of their
211 publications [1]. Thus to increase transparency, we should identify components of this incentive
212 structure amenable to improvement.

213
214 Some solutions to improve transparency
215

216 There is growing recognition of the problems hindering empirical progress and of the role that
217 institutions must play in shaping science in ecology, evolutionary biology, and beyond [42-44]. In
218 November 2015, representatives (mostly editors-in-chief) from nearly 30 journals in ecology and
219 evolution joined funding agency panelists and other researchers to identify ways to improve
220 transparency in these disciplines. At this workshop, strong support emerged for the recently introduced
221 Transparency and Openness Promotion (TOP) framework (<https://cos.io/top/>)[45]. TOP currently
222 consists of eight guidelines that can be implemented by journals and funding agencies. Institutions can
223 adopt whichever of the eight guidelines they choose, and they can implement these guidelines along a
224 gradient of stringency. The rapid and extensive spread of support for TOP (>500 journals in < 1 year)
225 across scientific disciplines appears to herald a revolution in transparency standards.

226
227 Several TOP guidelines simply request or require more thorough reporting of methods, results, data, or
228 analysis code. Ecologists and evolutionary biologists made important progress in this regard several
229 years ago when a growing number of journals began requiring the archiving of data [46]. Calls for more
230 expanded archiving are growing in ecology and evolution [47], and the TOP guidelines can facilitate the
231 expansion of these types of disclosures. Interestingly, an incentive to archive in the form of a badge
232 appears similarly effective [48] as requiring archiving [49] and could therefore eliminate much of the
233 controversy regarding archiving [e.g., 50]. However, challenges remain, such as ensuring inclusion of
234 sufficient metadata [49]. The TOP guideline titled 'analysis and design transparency' calls for discipline-
235 specific guidance regarding what information should be disclosed in publications, and to that end, the
236 workshop produced a document 'Tools for Transparency in Ecology and Evolution' (TTEE;
237 <https://osf.io/g65cb/>) that provides checklist questions that journals can provide to authors, reviewers,
238 and editors to facilitate transparent reporting. Promoting more thorough and consistent reporting of
239 results and methods through TOP and TTEE should dramatically improve transparency, but here we also
240 highlight two other TOP components that could have transformative impacts on our field.

241

242 Pre-registration (see glossary), in which researchers register their study and data analysis plan prior to
243 collecting data, can greatly improve transparency. Although requiring pre-registration (as in clinical trial
244 research)[51] might thwart publication of valuable exploratory and serendipitous findings in ecology and
245 evolution, encouraging pre-registration where appropriate has large potential benefits. Most obviously,
246 it makes unpublished results more discoverable [45], thus helping to reduce publication bias. Potentially
247 more important, however, pre-registration of analysis plans ensures that we can identify genuine *a*
248 *priori* planned tests, helping to improve confidence in results because they are unlikely to derive from
249 hidden multiple hypothesis testing and selective reporting. As pre-registration becomes more common,
250 results that do not come from pre-registered analysis plans become viewed as exploratory, and thus
251 provisional and less convincing than pre-registered results [52], providing a strong incentive to pre-
252 register studies. We acknowledge that exploratory work is hugely important in ecology and evolutionary
253 biology and we do not wish to impede it, but it should be more consistently identifiable and it should be
254 follow-up with planned, ideally pre-registered, tests [35]. A common concern is that pre-registration
255 ignores the inevitable tweaking of methods that occurs as field projects evolve. However, alterations to
256 methods or analysis plans can be justified in the published study [e.g., 48]. Reviewers and editors can
257 decide if the reported methods and analyses adhered closely enough to the pre-registration to earn a
258 pre-registration badge (<https://osf.io/tvyxz/wiki/home/>). Further, pre-registered analyses and
259 exploratory results can be published in the same paper when the distinction between them is made
260 clear. In an effort to further jump start the pre-registration process, the Center for Open Science
261 recently announced the Pre-registration Challenge, in which the first thousand researchers to publish
262 pre-registered research will be awarded US\$1000 each (<https://cos.io/prereg/>). Independently,
263 institutions promoting systematic reviews in ecology and conservation have also been encouraging pre-
264 registration (<http://www.environmentalevidence.org/>; <http://cebc.bangor.ac.uk/>).

265
266 The final TOP guideline promotes replications (see glossary) of previously published studies. Replication
267 to assess validity and generality of prior results is a core practice of science. Exact replication is not
268 possible, especially in field studies, but various forms of replication, especially when combined with
269 meta-analysis, are powerful tools for establishing the applicability of hypotheses [37]. Unfortunately,
270 institutional incentive structures often work strongly against replication in ecology and evolution,
271 especially replications that seek to closely match methods as part of the process of assessing validity
272 [37]. Journals and funding bodies explicitly favor novelty. Of course progress requires novelty, but
273 progress also requires rigorous evaluation of prior findings. Not all studies are of high priority for
274 replication. The more interesting or important a finding, however, the more important it is to replicate
275 that study. Allocating funding to replication would certainly increase its frequency, as would journals
276 adopting policies that explicitly encourage submission of replications (e.g.,
277 <http://biotropica.org/reproducibility-repeatability/>). As with any other articles, journals can reject less
278 valuable replication studies. For instance, journals might require sample sizes larger than in the original
279 study, review of methods prior to conducting the research (i.e., ‘registered reports’; see glossary) [53],
280 or replications only of original studies that cross some threshold of impact or interest. Replication is an
281 essential part of doing science in other fields, as, for example, anyone who remembers the ‘cold fusion
282 in a jar’ debacle of 1989 can attest [54].

283
284 As institutions in ecology and evolutionary biology more vigorously promote transparency, we will
285 become better able to evaluate the results we read, the average result will be more reliable, and there
286 will be clearer paths for empirical progress (Fig. 1). We need to deliberately shape the institutions in
287 which we operate to best facilitate scientific progress. Not all institutions will be equally responsive to
288 attempts at reform. However, we already know that journals can take deliberate steps to increase
289 transparency [46], and in response to the TTEE workshop mentioned above, nearly 30 ecology and
290 evolution journals are engaged in ongoing discussions about adopting TOP guidelines or have already
291 adopted these guidelines. Funding agencies have also implemented data archiving policies [46] and
292 could promote transparency in multiple other ways as guided by TOP. The proposals we review here are
293 only a subset of possible solutions to insufficient transparency. We hope to stimulate a continuing

294 exploration of these issues. This is an historic crossroads for the practice of science in ecology and
295 evolutionary biology, and for empirical disciplines in general [45].

296 Acknowledgements

297

298 We thank Mark Elgar for requesting an aggregation of evidence regarding the current state of
299 transparency in ecology and evolution. This request was made at the November 2015 workshop titled
300 “Improving Inference in Evolutionary Biology and Ecology.” Other participants at this workshop
301 (complete list: <https://osf.io/dhp3t/>) were also vital contributors to discussions that inspired this paper.
302 Financial support for the workshop was provided by the US National Science Foundation (DEB: 1548207)
303 and The Laura and John Arnold Foundation, and logistical support was provided by the Center for Open
304 Science. ARC Future Fellowships supported S. N. (FT130100268) and F. F. (FT150100297). We also thank
305 Losia Lagisz for helping to make Figure 1. Comments from an anonymous reviewer significantly
306 improved the manuscript.

307 Glossary

308

309 **Blind observation:** The observer (person making measurements) is unaware of the group membership
310 (e.g., treatment condition) of the subject being measured

311

312 **Confirmation bias:** The widespread human tendency to interpret observations as consistent with one's
313 belief about how the world works or to preferentially search for and recall such observations

314

315 **Effect size:** A measure of study outcome that indicates the magnitude and direction of the outcome of
316 each study. Effect sizes can be based on the magnitude of difference between groups or the strength of
317 the correlation between variables. Effect sizes can be unstandardized (e.g., mean difference or
318 covariance) or standardized (e.g., Cohen's *d* or correlation coefficient).

319

320 **Exploratory analysis:** conducting many graphical and/or statistical comparisons in an effort to identify
321 previously unidentified relationships among variables in a data set

322

323 **False positive:** In null hypothesis testing, a rejection of the null hypothesis when the null hypothesis is
324 actually true (Type I error)

325

326 **HARKing:** Hypothesizing After Results are Known – presenting a *post hoc* explanation for an exploratory
327 result as though it were an *a priori* hypothesis. Many of us were taught to HARK and to write papers as
328 though we were testing *a priori* hypotheses even if we were conducting exploratory analyses. Although
329 philosophers debate the importance of distinguishing between *a priori* and *post hoc* hypotheses,
330 HARKing is problematic even if one discounts this distinction. This is because HARKing often serves to
331 conceal selective reporting of exploratory analyses (often without a deliberate attempt to deceive), and
332 thus skews the distribution of reported results.

333

334 **Inflated effect size:** An estimated effect size that is larger than the actual effect size, for instance
335 because the researcher selected the covariate that led to the largest effect in the target relationship
336 after testing multiple covariates

337

338 **Meta-analysis:** The quantitative synthesis of the outcomes of different studies, based on combining
339 effect sizes, to determine overall results across studies and sources of heterogeneity in outcomes among
340 studies. Generally study outcomes are weighted by the precision with which the effects are estimated.
341 Meta-regression is a variant of meta-analysis in which the effects of covariates are modeled statistically.

342

343 ***p*-hacking:** A variety of practices that increase the odds of finding a statistically significant result by, for
344 instance, conducting multiple versions of an analysis with different covariates, interactions, or subsets of
345 data. Some processes that contribute to *p*-hacking, such as conducting multiple versions of an analysis
346 with different interaction terms, might be pursued out of a sincere desire to discover the story the data
347 have to tell. However, each additional version of the analysis increases the risk of a false positive or of
348 an inflated effect, and unless we disclose all results from all versions of analyses and all decisions
349 regarding data gathering and analyses, we will contribute to the biased distribution of effects in the
350 literature.

351

352 **Pre-registration:** A process by which planned studies, including methods and an analysis plan, are
353 registered in a secure and accessible platform (e.g. website such as Open Science Framework;
354 <https://osf.io/>) before commencement of the research. Once a pre-registration has been submitted, it
355 cannot be altered. Pre-registrations can be embargoed to protect ideas prior to publication.

356

357 **Publication bias:** A bias in the distribution of published effect sizes resulting from any number of factors,
358 including selective reporting by authors and rejection of non-significant results by editors

359
360 **Registered report:** A study in which the rationale, methods, and analysis plan are submitted to a journal
361 for review, and possible revision, with the objective of achieving in-principle acceptance based on the
362 importance of the question and the quality of the study design, not the outcome, prior to initiation of
363 the study.
364
365 **Replication:** a study designed to replicate a previously published result, either by closely following the
366 original methods in an effort to assess validity ('direct' or 'close' replication) or by designing a study
367 inspired by the original concept in an effort to assess generality ('conceptual replication')
368
369 **Selective reporting:** Reporting only a subset of analyses conducted. In medicine, a similar concept is
370 often referred to as reporting bias.
371
372 **Statistical power:** The probability of detecting a statistically significant effect if that effect actually exists.
373 This probability is a function of the significance threshold, sample size, and strength of statistical effect.
374
375 **Type I error:** Rejection of a null hypothesis when the null hypothesis is true (a 'false positive').
376
377 **Type II error:** a failure to reject a null hypothesis when the null hypothesis is false (a 'false negative')
378
379 **Type M error:** an error in estimating the magnitude of an effect
380
381 **Type S error:** an error in estimating the sign of an effect
382
383 **Under-reporting:** Reporting an analysis without sufficient details of analytical methods or results to
384 allow for interpretation
385
386
387
388

389 Text boxes

390

391 Text Box 1

392

393 Confirmation bias

394

395 People have a strong tendency to interpret observations as supporting their existing worldview and to
396 seek out evidence in support of this worldview [7]. This can play out in various forms of selective
397 reporting as we convince ourselves that we are simply focusing our reporting on the real phenomena.
398 Confirmation bias can thus help rationalize p -hacking and selective reporting, often by preventing us
399 from recognizing our own subtle HARKing. Confirmation bias can also influence data gathering. Studies
400 in ecology and evolution in which individuals gathering data were not blind to the treatment condition
401 or the predicted outcomes showed stronger effects and higher rates of significance than studies with
402 blinded observers [55, 56]. Blind observation (see Glossary) is quite rare in ecology and evolutionary
403 biology [57] in part because in some studies blinding is nearly impossible. However, in a large sample of
404 recent studies, 56% that could have benefited from blinding could also have implemented it with little
405 difficulty (e.g., no additional personnel), and an additional 22% could have adopted blinding by
406 employing an observer naïve to certain details of the study [57].

407

408

409 Text Box 2

410

411 Evidence of low power

412

413 In a sample of 1362 statistical tests from 697 papers published in 2000 in 10 behavior, evolution, and
414 ecology journals, the average power to detect a small effect ($|r| = 0.1$) was only 13-16% [27]. In other
415 words, studies would only be expected to reject a false null hypothesis 13-16% of the time in the case of
416 weak effects. Power to detect medium ($|r| = 0.3$) and large ($|r| = 0.5$) effects, though of course higher
417 (40-47% and 65-72%, respectively), was still typically well below the commonly recommended threshold
418 of 80%. Examined another way, the proportion of studies reaching this 80% power threshold to detect
419 weak effects was 2-3%, 13-21% for medium effects, and 37-50% for strong effects [27]. Other analyses
420 of power find similar results. For example, an analysis of studies published in *Animal Behaviour* in 1996,
421 2003, and 2009 found, across all three years, an average power of just 23-26% for detection of medium
422 effects and 1-2% for weak effects [28]. It thus appears that studies in ecology and evolution often lack
423 power to detect small and medium effects, and this is particularly problematic because effects in
424 ecology and evolution tend to be weak. Average effects across 43 meta-analyses in ecology and
425 evolutionary biology were found to be weak to moderate ($|r| = 0.18-0.19$) [25]. Further, these rather
426 low values are actually overestimates because averages of estimated absolute values of effect size are
427 upwardly biased [26]. To detect these relatively small effects requires large samples (e.g., $n = 207$ to
428 obtain an 80% probability of detecting a true effect of $r = 0.193$) [25], but obtaining sufficient power
429 through large samples is rare [27].

430

431 Text Box 3

432 False-positive report probability (FPRP)

433

434 In many sub-fields of evolution and ecology it remains common to use a significance threshold of 5%.
 435 This means that if our null hypothesis were true we would incorrectly reject it 5% of the time. However,
 436 we often incorrectly attribute a frequency of 5% to a different phenomenon: the chance that a
 437 significant finding is a false positive. This is incorrect because the probability that a positive result is a
 438 false positive depends on three factors (1) the proportion of our hypotheses that are in fact true (π , the
 439 probability that a hypothesis is true), (2) the significance threshold (α), and (3) statistical power ($1 - \beta$,
 440 where β is the probability of making a type II error; Table 1): $FPRP = (\alpha(1 - \pi)/[\alpha(1 - \pi) + (1 - \beta)\pi]$. With
 441 50% of our hypotheses true and statistical power of 20% (a power typical in ecology and evolution [25]),
 442 the chance that a significant finding is a false positive is 20%. This value is known as the false positive
 443 report probability [58]. This number is notably larger than 5%, but it becomes dramatically larger when,
 444 in pursuit of novelty, we turn our interest towards testing relatively unlikely hypotheses, those that in
 445 the Bayesian sense could be said to have a low prior probability. For instance, when only 10% of tested
 446 hypotheses are in fact true, the expected false positive report probability rises to 69% $((0.05(1 -$
 447 $0.1)/((0.05(1 - 0.1) + (0.2)0.1))$ [58]! In fact, false positives could be even more prevalent. The above
 448 calculations assume complete and transparent reporting of the full set of analyses conducted, as
 449 promoted by pre-registration and other recently proposed transparency tools. If, in contrast,
 450 researchers make their choices of analysis strategy conditional on the outcome as with *p*-hacking (i.e.
 451 preferring test variants that yield significance or stronger effects) then the false-report probability
 452 increases further.

453

454 I. Four possible outcomes from a null hypothesis statistical test together with the probabilities of
 455 each outcome depending on whether the null-hypothesis is true

	Null Hypothesis True	Alternate Hypothesis True
Significant Finding	False Positive: α	True Positive: $1 - \beta$
Non-Significant Finding	True Negative: $1 - \alpha$	False Negative: β

456

457

458

459

460 Tables

461

462 Table 1. A sample of studies in ecology and evolution that quantify rates of under-reporting of important
463 details of methods or results in the published literature.

464

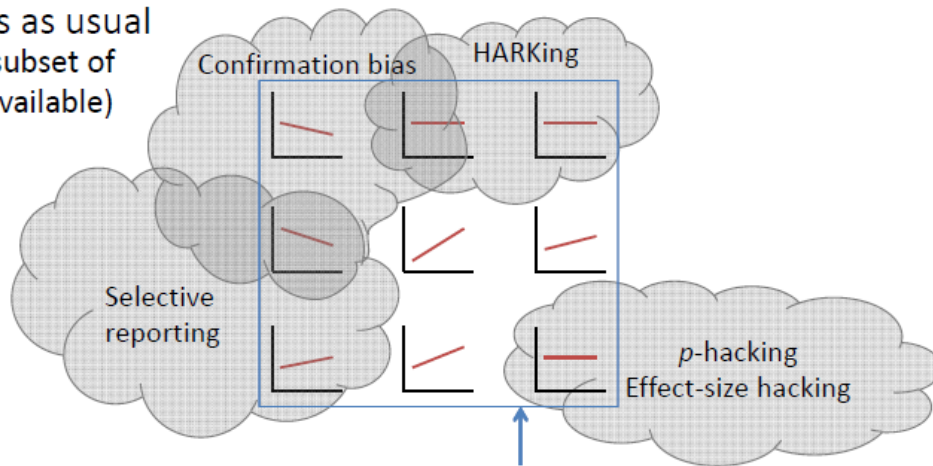
Citation	Studies reviewed	finding
Ferreira et al. (2015)	99 studies of litter decomposition in streams as an effect of nutrient enrichment	Estimates of decomposition rate presented without estimate of uncertainty in 54% of studies (even after requesting details directly from authors)
Fidler et al. (2006)	78 articles published in 2005 in Conservation Biology and Biological Conservation	58% missing at least one effect size 51% missing at least one sample size 85% missing at least one SE or SD
Parker (2013)	48 studies of plumage color in a well-studied European songbird species	409 of 997 main-effect relationships lacked information to estimate the strength and/or direction of the effect
Zhang et al. (2012)	54 studies of forest productivity as a function of tree diversity	29 studies failed to provide either estimates of variance associated with means or corresponding sample sizes

465

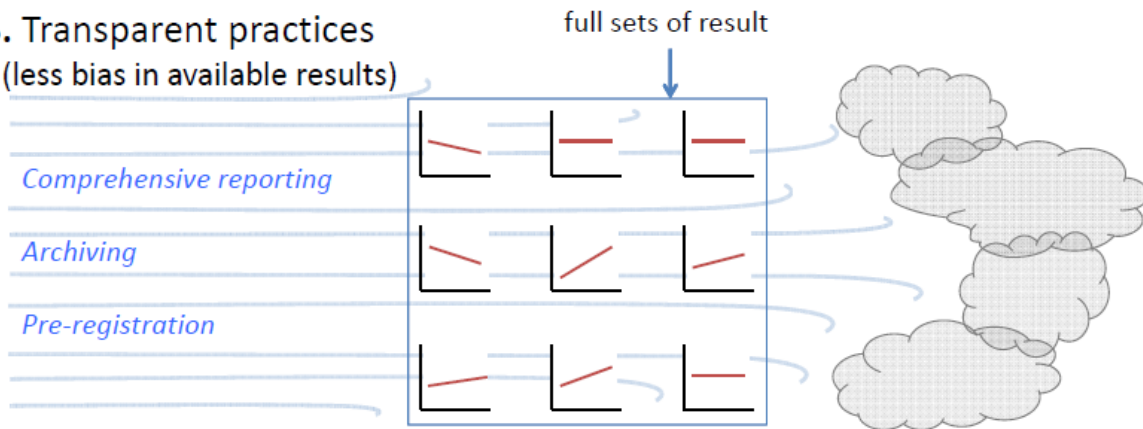
466

467

A. Business as usual
(biased subset of results available)



B. Transparent practices
(less bias in available results)



470
471 Figure 1. ‘Business as usual’ in ecology and evolution allows and often promotes practices that keep
472 many analyses hidden and this leads to biases in the published literature. For example, current practices
473 (A) could result in only the three ‘unclouded’ graphs making it to publication, leaving the impression that
474 all results were consistently positive. However, full transparency (B) will sometimes leave a very
475 different impression of results. In this illustration, we see results that are more complicated and less
476 consistent, and suggest a much smaller average effect, if any.
477

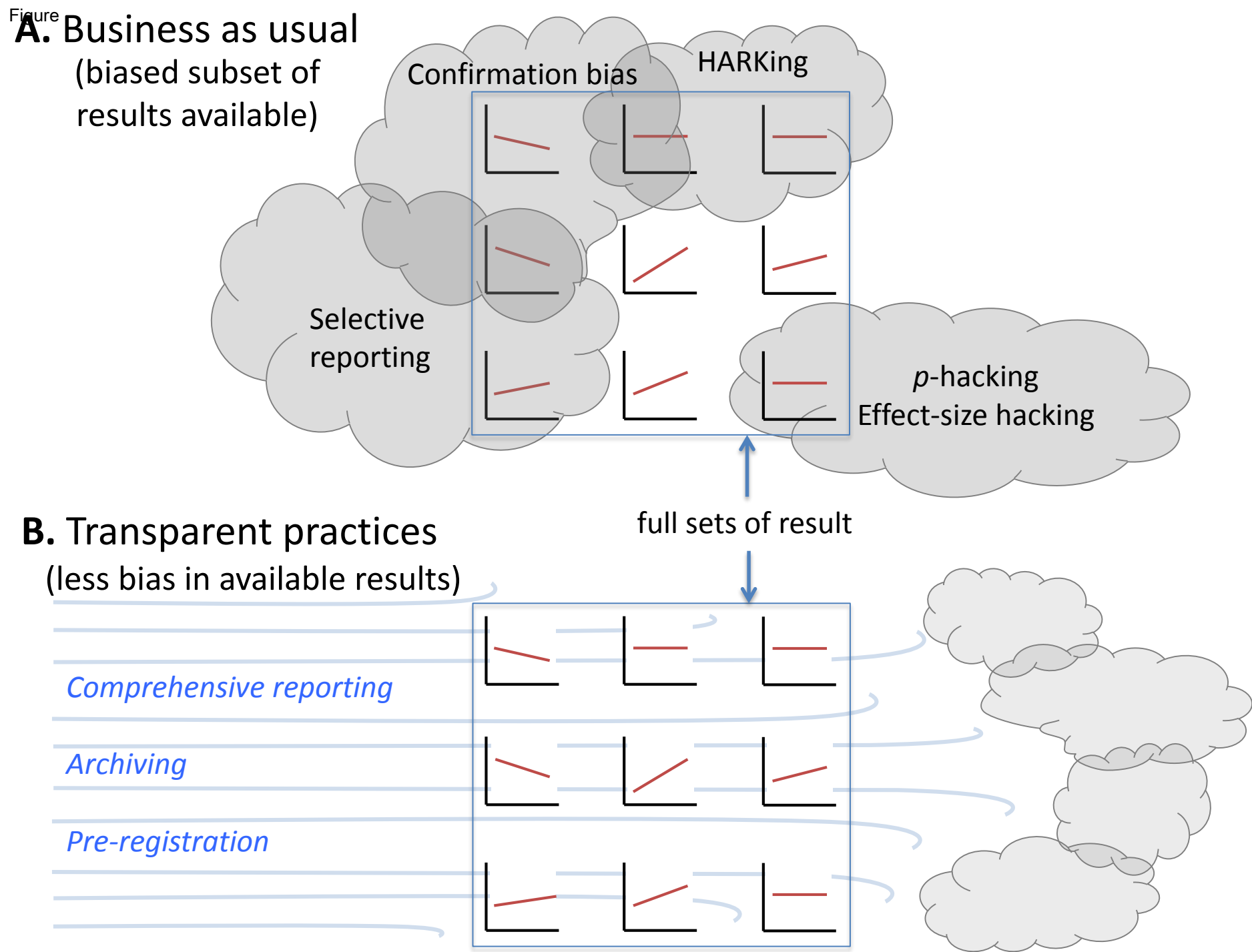
478 Literature Cited

- 479
- 480 1. Smaldino, P.E. and McElreath, R. (2016) The natural selection of bad science. *arXiv*,
481 1605.19511v19511.
 - 482 2. Møller, A.P. and Jennions, M.D. (2001) Testing and adjusting for publication bias. *Trends in*
483 *Ecology & Evolution* 16, 580-586.
 - 484 3. Godefroid, S., *et al.* (2011) How successful are plant species reintroductions? *Biological*
485 *Conservation* 144, 672-682.
 - 486 4. Head, M.L., *et al.* (2015) The extent and consequences of *P*-Hacking in science. *PLoS Biol* 13,
487 e1002106.
 - 488 5. Simonsohn, U., *et al.* (2014) *P*-curve: a key to the file drawer. *Journal of Experimental*
489 *Psychology: General* 143, 534-547.
 - 490 6. Kerr, N.L. (1998) HARKing: hypothesizing after the results are known. *Personality and Social*
491 *Psychology Review* 2, 196-217.
 - 492 7. Nickerson, R.S. (1998) Confirmation bias: A ubiquitous phenomenon in many guises. *Review of*
493 *General Psychology* 2, 175-220.
 - 494 8. Gelman, A. (2015) Working through some issues. *Significance* 12, 33-35.
 - 495 9. Rothstein, H.R., *et al.*, eds (2005) *Publication bias in meta-analysis: prevention, assessment and*
496 *adjustments*. John Wiley & Sons, Lt.
 - 497 10. Fidler, F., *et al.* (2006) Impact of criticism of null-hypothesis significance testing on statistical
498 reporting practices in conservation biology. *Conservation Biology* 20, 1539-1544.
 - 499 11. Koricheva, J. and Gurevitch, J. (2014) Uses and misuses of meta-analysis in plant ecology. *Journal*
500 *of Ecology* 102, 828-844.
 - 501 12. Parker, T.H. (2013) What do we really know about the signalling role of plumage colour in blue
502 tits? A case study of impediments to progress in evolutionary biology. *Biological Reviews* 88,
503 511-536.
 - 504 13. Ferreira, V., *et al.* (2015) A meta-analysis of the effects of nutrient enrichment on litter
505 decomposition in streams. *Biological Reviews* 90, 669-688.
 - 506 14. Menge, D.N.L. and Field, C.B. (2007) Simulated global changes alter phosphorus demand in
507 annual grassland. *Global Change Biology* 13, 2582-2591.
 - 508 15. Zhang, Y., *et al.* (2012) Forest productivity increases with evenness, species richness and trait
509 variation: a global meta-analysis. *Journal of Ecology* 100, 742-749.
 - 510 16. Leisner, C.P. and Ainsworth, E.A. (2012) Quantifying the effects of ozone on plant reproductive
511 growth and development. *Global Change Biology* 18, 606-616.
 - 512 17. Moles, A.T., *et al.* (2011) Assessing the evidence for latitudinal gradients in plant defence and
513 herbivory. *Functional Ecology* 25, 380-388.
 - 514 18. Cassey, P., *et al.* (2004) A survey of publication bias within evolutionary ecology. *Proceedings of*
515 *the Royal Society of London B: Biological Sciences* 271, S451-S454.
 - 516 19. Bruns, S.B. and Ioannidis, J.P.A. (2016) *p*-Curve and *p*-Hacking in observational research. *PLoS*
517 *ONE* 11, e0149144.
 - 518 20. Bishop, D.V.M. and Thompson, P.A. (2016) Problems in using *p*-curve analysis and text-mining to
519 detect rate of *p*-hacking and evidential value. *PeerJ* 4, e1715.
 - 520 21. Ridley, J., *et al.* (2007) An unexpected influence of widely used significance thresholds on the
521 distribution of reported *P*-values. *J. Evol. Biol.* 20, 1082-1089.
 - 522 22. Gelman, A. and O'Rourke, K. (2014) Discussion: Difficulties in making inferences about scientific
523 truth from distributions of published *p*-values. *Biostatistics* 15, 18-23.
 - 524 23. Fanelli, D. (2010) "Positive" results increase down the hierarchy of the sciences. *PLoS ONE* 5,
525 e10068.
 - 526 24. Csada, R.D., *et al.* (1996) The "file drawer problem" of non-significant results: does it apply to
527 biological research? *Oikos* 76, 591-593.

- 528 25. Møller, A.P. and Jennions, M.D. (2002) How much variance can be explained by ecologists and
529 evolutionary biologists? *Oecologia* 132, 492-500.
- 530 26. Hereford, J., *et al.* (2004) Comparing strengths of directional selection: how strong is strong?
531 *Evolution* 58, 2133-2143.
- 532 27. Jennions, M.D. and Møller, A.P. (2003) A survey of the statistical power of research in behavioral
533 ecology and animal behavior. *Behav. Ecol.* 14, 438-445.
- 534 28. Smith, D.R., *et al.* (2011) Power rangers: no improvement in the statistical power of analyses
535 published in *Animal Behaviour*. *Animal Behaviour* 81, 347-352.
- 536 29. Button, K.S., *et al.* (2013) Power failure: why small sample size undermines the reliability of
537 neuroscience. *Nat Rev Neurosci* 14, 365-376.
- 538 30. Gelman, A. and Weakliem, D. (2009) Of beauty, sex, and power. *American Scientist* 97, 310-316.
- 539 31. Eberhardt, L.L. and Thomas, J.M. (1991) Designing environmental field studies. *Ecological*
540 *Monographs* 61, 53-73.
- 541 32. Murtaugh, P.A. (2014) In defense of *P* values. *Ecology* 95, 611-617.
- 542 33. Barto, E.K. and Rillig, M.C. (2012) Dissemination biases in ecology: effect sizes matter more than
543 quality. *Oikos* 121, 228-235.
- 544 34. Murtaugh, P.A. (2002) Journal quality, effect size, and publication bias in meta-analysis. *Ecology*
545 83, 1162-1166.
- 546 35. Pike, N. (2011) Using false discovery rates for multiple comparisons in ecology and evolution.
547 *Methods in Ecology and Evolution* 2, 278-282.
- 548 36. Forstmeier, W. and Schielzeth, H. (2011) Cryptic multiple hypotheses testing in linear models:
549 overestimated effect sizes and the winner's curse. *Behavioral Ecology and Sociobiology* 65, 47-
550 55.
- 551 37. Nakagawa, S. and Parker, T.H. (2015) Replicating research in ecology and evolution: feasibility,
552 incentives, and the cost-benefit conundrum. *BMC Biology* 13, 88.
- 553 38. Kelly, C.D. (2006) Replicating empirical research in behavioral ecology: how and why it should be
554 done but rarely ever is. *Q. Rev. Biol.* 81, 221-236.
- 555 39. Birkhead, T.R. (2002) Of Moths and Men (book review). *International Society for Behavioral*
556 *Ecology Newsletter* 14, 15-16.
- 557 40. Nakagawa, S. and Cuthill, I.C. (2007) Effect size, confidence interval and statistical significance: a
558 practical guide for biologists. *Biological Reviews* 82, 591-605.
- 559 41. Belovsky, G.E., *et al.* (2004) Ten suggestions to strengthen the science of ecology. *BioScience* 54,
560 345-351.
- 561 42. Baker, M. (2016) 1,500 scientists lift the lid on reproducibility. *Nature* 533, 452-454.
- 562 43. Parker, T.H. and Nakagawa, S. (2014) Mitigating the epidemic of type I error: ecology and
563 evolution can learn from other disciplines. *Frontiers in Ecology and Evolution* 2.
- 564 44. Open_Science_Collaboration (2015) Estimating the reproducibility of psychological science.
565 *Science* 349.
- 566 45. Nosek, B.A., *et al.* (2015) Promoting an open research culture. *Science* 348, 1422-1425.
- 567 46. Whitlock, M.C. (2011) Data archiving in ecology and evolution: best practices. *Trends in Ecology*
568 *& Evolution* 26, 61-65.
- 569 47. Mislán, K.A.S., *et al.* (2016) Elevating the status of code in ecology. *Trends in Ecology & Evolution*
570 31, 4-7.
- 571 48. Kidwell, M.C., *et al.* (2016) Badges to acknowledge open practices: a simple, low cost, effective
572 method for increasing transparency. *PLOS Biology* 14, e1002456.
- 573 49. Roche, D.G., *et al.* (2015) Public data archiving in ecology and evolution: how well are we doing?
574 *PLoS Biology* 13, e1002295.
- 575 50. Mills, J.A., *et al.* (2015) Archiving primary data: solutions for long-term studies. *Trends in Ecology*
576 *& Evolution* 30, 581-589.
- 577 51. Ross, J.S., *et al.* (2009) Trial publication after registration in ClinicalTrials.gov: a cross-sectional
578 analysis. *PLoS Med* 6, e1000144.

- 579 52. Wagenmakers, E.-J., *et al.* (2012) An agenda for purely confirmatory research. *Perspectives on*
580 *Psychological Science* 7, 632-638.
- 581 53. Chambers, C., D. (2013) *Registered Reports*: a new publishing initiative at *Cortex*. *Cortex* 49, 609-
582 610.
- 583 54. Huizenga, J.R. (1994) *Cold Fusion: The Scientific Fiasco of the Century*. Oxford University Press.
- 584 55. van Wilgenburg, E. and Elgar, M.A. (2013) Confirmation bias in studies of nestmate recognition:
585 a cautionary note for research into the behaviour of animals. *PLoS ONE* 8, e53548.
- 586 56. Holman, L., *et al.* (2015) Evidence of experimental bias in the life sciences: why we need blind
587 data recording. *PLoS Biol* 13, e1002190.
- 588 57. Kardish, M.R., *et al.* (2015) Blind trust in unblinded observation in ecology, evolution and
589 behavior. *Frontiers in Ecology and Evolution* 3, 51.
- 590 58. Wacholder, S., *et al.* (2004) Assessing the probability that a positive report is false: an approach
591 for molecular epidemiology studies. *Journal of the National Cancer Institute* 96, 434-442.

592



The instructions for authors state that an “outstanding questions box” is not required, but the submission website required a document in this category.

The main purpose of this article is not to stimulate further research into transparency, but to make the case to ecologists and evolutionary biologists that editorial policies that promote transparency are desirable.

That said, there is certainly room to build our empirical understanding of publication bias and other obstacles to transparency, so we will be happy to provide “outstanding questions” if requested.