

# Empowering peer reviewers with a checklist to improve transparency

Timothy H. Parker<sup>1,2\*</sup>, Simon C. Griffith<sup>3</sup>, Judith L. Bronstein<sup>3</sup>, Fiona Fidler<sup>4,5</sup>, Susan Foster<sup>6</sup>, Hannah Fraser<sup>4</sup>, Wolfgang Forstmeier<sup>7</sup>, Jessica Gurevitch<sup>8</sup>, Julia Koricheva<sup>9</sup>, Ralf Seppelt<sup>10,11,12</sup>, Morgan W. Tingley<sup>13</sup> and Shinichi Nakagawa<sup>14</sup>

**Peer review is widely considered fundamental to maintaining the rigour of science, but it often fails to ensure transparency and reduce bias in published papers, and this systematically weakens the quality of published inferences. In part, this is because many reviewers are unaware of important questions to ask with respect to the soundness of the design and analyses, and the presentation of the methods and results; also some reviewers may expect others to be responsible for these tasks. We therefore present a reviewers' checklist of ten questions that address these critical components. Checklists are commonly used by practitioners of other complex tasks, and we see great potential for the wider adoption of checklists for peer review, especially to reduce bias and facilitate transparency in published papers. We expect that such checklists will be well received by many reviewers.**

Two important tasks facing peer reviewers are assessing the soundness of study design and evaluating the reporting of methods and results. Study soundness and reporting both bear directly on the reliability of the inferences that can be drawn from the papers that are ultimately published<sup>1</sup>. Other reviewing tasks include considering the placement of the study in a broader context, the writing and the importance of the research, but these vary by journal and the expertise of the reviewer, and are often more subjective. We therefore focus on only the first two reviewing tasks. Our goal here is to explain particular components of this assessment process that we believe are too frequently ignored by peer reviewers, ultimately to the detriment of the scientific literature. We combine these components in a checklist that reviewers can use to improve transparency and reduce bias, and thus improve the reliability of scientific inferences.

We present this checklist as a series of ten questions (summarized in Box 1 and Supplementary Information), each accompanied by suggestions for how the reviewer should proceed depending on the answer to that question. The checklist is not meant to be comprehensive. A longer checklist to help reviewers in ecology and evolutionary biology promote transparency was created as part of Tools for Transparency in Ecology and Evolution (TTEE; <https://osf.io/y8aqx/>) in an effort to help journals in ecology and evolutionary biology adopt transparency and openness promotion guidelines<sup>2</sup>. TTEE checklists, for both reviewers and authors, were designed to cover a broad swath of transparency issues. In contrast, the short checklist we present in this paper focuses on the subset of practices that we think are critically in need of improvement, and on which

we think a concise checklist can achieve greatest impact. Our checklist provides reviewers with an efficient tool for promoting transparency in empirical research papers.

## Why a checklist?

The use of checklists is well established among skilled practitioners working in complex systems. Checklists make flying complicated aircraft safer, they free architects to devote their mental energy to creativity and they help surgeons focus on applying their skill without forgetting vital tasks<sup>3,4</sup>. Good checklists do not replace complex thought; they facilitate it. Of course, effective peer review requires expertise and critical thinking skills that no practical checklist can provide. However, this does not mean that checklists cannot be used to improve peer review, even dramatically, by calling attention to essential elements that are often overlooked.

Checklists can be of use to peer reviewers in two primary ways related to creating a more transparent and less biased literature: to help reviewers check (1) mundane but important details, and (2) both their own and the authors' potential biases. With regard to the first point, incomplete reporting of information hinders interpretation of studies and effective synthesis, and thus scientific progress<sup>1,5</sup>. We know from surveys of subsets of the ecology literature that approximately half of published papers omit important information such as sample size or variability associated with estimates<sup>6–8</sup>. Nearly all papers omitting this information were peer reviewed, suggesting that reviewers either overlooked these details, or felt that it was someone else's job to monitor them. Whether we notice omissions as reviewers depends on scrutiny that may vary unconsciously with

<sup>1</sup>Department of Biology, Whitman College, Walla Walla, WA, USA. <sup>2</sup>Department of Biological Sciences, Macquarie University, North Ryde, New South Wales, Australia. <sup>3</sup>Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ, USA. <sup>4</sup>School of BioSciences, University of Melbourne, Melbourne, Victoria, Australia. <sup>5</sup>History & Philosophy of Science, School of Historical & Philosophical Studies, University of Melbourne, Melbourne, Victoria, Australia. <sup>6</sup>Department of Biology, Clark University, Worcester, MA, USA. <sup>7</sup>Department of Behavioural Ecology and Evolutionary Genetics, Max Planck Institute for Ornithology, Seewiesen, Germany. <sup>8</sup>Department of Ecology and Evolution, Stony Brook University, Stony Brook, NY, USA. <sup>9</sup>School of Biological Sciences, Royal Holloway University of London, Egham, Surrey, UK. <sup>10</sup>Department of Computational Landscape Ecology, UFZ – Helmholtz Centre for Environmental Research, Leipzig, Germany. <sup>11</sup>Institute of Geoscience and Geography, Martin-Luther-University Halle-Wittenberg, Halle (Saale), Germany. <sup>12</sup>Div – German Centre for Integrative Biodiversity Research Halle-Jena-Leipzig, Leipzig, Germany. <sup>13</sup>Department of Ecology and Evolutionary Biology, University of Connecticut, Storrs, CT, USA. <sup>14</sup>Evolution & Ecology Research Centre, School of Biological, Earth and Environmental Sciences, University of New South Wales, Randwick, New South Wales, Australia. \*e-mail: [parkerth@whitman.edu](mailto:parkerth@whitman.edu)

**Box 1 | Concise version of ten questions reviewers can use to improve transparency and reduce bias in the empirical literature****Questions to promote transparent reporting of methods and results**

- 1 Were all sample sizes fully reported, including exact values for all subsets of data (for example, each treatment group), and for all statistical analyses?
- 2 Are the methods reported in sufficient detail to allow another researcher to gather the same data and run the identical analyses?
- 3 Are statistical results reported completely (considered in two parts below)?
- 3a Are statistical results for each test reported in sufficient detail? What qualifies as 'sufficient detail' will differ among analyses.
- 3b Are results from all variables and from all models reported? Complete reporting should include results related to all variables examined in preliminary models and all results from exploratory analyses.

**Questions to check biases of reviewers and authors**

- 4 Were observers kept unaware of the experimental treatment imposed on the samples (for example, organisms, plots) when recording observations or measurements so as to minimize unconscious bias?
- 5 Did the authors explain how sample size was decided (for example, based on a priori power analysis or logistical constraints), or when an experiment with pre-set sample sizes was terminated? If sample size or the end of the experiment was not decided prior to the initiation of the study, was there a decision rule for when to cease data collection?
- 6 Did the authors develop their analysis plan, including choices of variables, without looking at the data, for instance prior to gathering data or with a dummy data set? This is most easily determined by the existence of a pre-registered analysis plan. In the absence of pre-registration, a statement from the authors about the development of their analysis plan is still important.
- 7 How suitable do you find the research methods without considering the outcome? Evaluate the design and methods regardless of whether or not there was a finding of 'statistical significance', or whether or not the results conform to a predicted pattern.
- 8 Are the sample sizes large enough to justify the authors' conclusions? If presenting significance tests, how much power would this study have to detect statistically significant weak, moderate and strong effects? Expectation of effect size can best be derived from average effect sizes presented in meta-analyses of similar topics. The effect size reported in the manuscript under review can be a poor estimate of the underlying effect size, especially if the sample size is small, which elevates sampling uncertainty. Statistical significance is a poor indicator of the reliability of an estimate across a wide range of sample sizes and common effect sizes.
- 9 What does the size of the estimated effect (for example, slope, correlation coefficient, difference in means) suggest about its biological or practical importance, and what does uncertainty around that effect estimate suggest about the estimate's precision?
- 10 How unexpected would you judge these results to be in light of prior empirically derived understanding? Effects that are more surprising in light of robust prior information are those that had a lower prior probability of being correct.

See the main text for details.

factors such as whether we agree with the study's conclusions, our perception of the expertise of the authors, or whether we have used similar research designs. Regardless, the frequency of these omissions

in the literature is evidence of a systematic problem, but one that could be resolved with the help of an appropriate checklist.

With regard to the second point, we need to explicitly address potential bias from authors and reviewers because all people, scientists included, are subject to biases that influence the information we notice and how we interpret that information<sup>9,10</sup>. Such biases have been shown to have major impacts on the content of scientific papers<sup>11–13</sup>, and so we expect them also to influence the opinions we form when reviewing such papers. In fact, evidence suggests that peer review often suffers from a multitude of complex, systematic biases<sup>14</sup>.

We hope that reviewers will find the questions in this checklist useful for most reviews. To facilitate the checklist's use, we provide some suggestions for reviewer responses to individual checklist questions, although we cannot provide a set of all possible answers to each one. Occasionally reviewers will be uncertain about answers to one or more of these questions. Sometimes this uncertainty can be resolved by asking for additional information from the authors, and sometimes the reviewer should simply notify the editor so that she or he can seek additional reviewer expertise if needed. Of course some questions may not apply to some papers; it will be up to the reviewer to determine the relevance of each question. Determining its relevance may be aided by the explanation and justification that we provide following each particular checklist item.

**Questions to promote transparent reporting of methods and results****1. Were all sample sizes fully reported, including exact values for all subsets of data (for example, each treatment group), and for all statistical analyses?**

- If 'no', request that authors provide this information.

Knowledge of sample size is essential for understanding the power of analyses (see below) and the reliability of estimates, and thus for interpreting results. It is also essential for later meta-analytic synthesis<sup>5</sup>. Yet, researchers fail to report sample sizes with troubling frequency<sup>7,8</sup>. Reporting a range (for example, '9–12 replicates per treatment') is inadequate.

**2. Are the methods for carrying out the study and analysing the results reported in sufficient detail to allow another researcher to gather the same data and run the identical analyses?**

When not in the paper itself, methodological details should be included in a supplement, or in many cases, archived in a publicly accessible and curated repository.

- If 'no', request that authors provide the relevant information.
- If you are uncertain about some aspect of the methods, state your uncertainty to the editor so that she or he can seek appropriate expertise as needed.

By keeping replicability in mind while reading the methods, the reviewer can determine if methods have been reported in sufficient detail. Necessary details vary among studies with different methods. For instance, in the case of Bayesian analyses, authors should explicitly define their priors and report how their posterior distributions were derived, if applicable including Markov chain Monte Carlo specifications, and method of convergence (mixing) assessment. Archiving of details such as analysis code is essential if others are to understand how results were derived<sup>15</sup> (see also <http://www.britishecologicalsociety.org/wp-content/uploads/2017/12/guide-to-reproducible-code.pdf>) and, at least theoretically, be able to replicate the study, including the analyses. This information should be stored in curated archives. Temporary and uncured repositories, including personal websites and the version-control site GitHub, are

not viable for long-term storage. There will occasionally be valid justifications for not reporting certain information (for example, population locations for species threatened by illegal collection), but in most cases these exceptions should be explicitly addressed in the manuscript.

### 3. Are statistical results reported completely (considered in two parts below)?

**3a. Are statistical results for each test reported in sufficient detail?** What qualifies as ‘sufficient detail’ will differ among analyses. For most analyses, however, this will include (but not be limited to) basic parameter estimates of central tendency (for example, means) or other basic estimates (for example, regression or correlation coefficients) and variation (for example, standard deviation) or associated estimates of uncertainty (for example, confidence/credible intervals). For null hypothesis tests, reporting *P* values and test statistics by themselves is almost always insufficient.

- If ‘no’, request that authors provide this information.
- If you are uncertain, state your uncertainty to the editor so that he or she can seek appropriate statistical expertise as needed. Remember that you may be the only reviewer looking carefully at this aspect of the manuscript.

**3b. Are results from all variables and from all models reported?** Complete reporting should include results related to all variables examined in preliminary models and all results from exploratory analyses. It will sometimes be appropriate to include these as supplementary materials. For analysis types that generate vast sets of results, it may be appropriate to place results in data archives.

- If ‘no’, request that authors provide this information
- If you are uncertain, ask the authors to declare in the paper that all exploratory analyses are reported in full. We recommend using the ‘Standard Reviewer Statement for Disclosure of Sample, Conditions, Measures, and Exclusions’: “I request that the authors add a statement to the paper confirming whether, for all experiments, they have reported all measures, conditions, data exclusions, and how they determined their sample sizes. The authors should, of course, add any additional text to ensure the statement is accurate. This is the standard reviewer disclosure request endorsed by the Center for Open Science [see <http://osf.io/hadz3>]”.

Insufficient reporting of results is one of the largest obstacles to an unbiased understanding of empirical progress<sup>1,16</sup>. Sometimes authors state that an analysis was conducted, but fail to provide all the relevant statistical outcomes such as slope estimates or estimates of variability<sup>6–8,17</sup>. At other times, authors conduct multiple analyses but do not explicitly acknowledge that they are reporting results from only a subset. Both practices may sometimes result from a direct request by the journal to shorten the text because of space limits or a desire for a concise story. Regardless, they weaken our ability to draw unbiased conclusions from the published literature. The failure to provide all relevant details from a reported analysis is often easily recognized by reviewers. In contrast, analyses that have been conducted but are completely unreported are more difficult, and sometimes even impossible, to recognize. However, there can be signs of unreported analyses: for instance, different variables may be included in different models without obvious *a priori* justification, a subset of potential interactions may be provided without clear justification for the choice, or the authors may have failed to examine obvious predictions that are testable with available data. Each of these signs was found in a sample of literature in behavioural ecology, providing circumstantial evidence of unreported analyses<sup>17</sup>. Reviewers can prompt authors to include missing information in

Supplementary materials or in searchable, curated data archives. Asking authors to state whether all results from all analyses have been reported should lead authors to be more transparent about their exploratory work<sup>18</sup>. If necessary, authors should be directed to consult published recommendations regarding thorough reporting of results (and methodological choices) from the type of analysis they have conducted<sup>5</sup>. Finally, it may help to remind authors that ‘not statistically significant’ does not mean ‘not interesting or not important’.

### Questions to check biases of reviewers and authors

#### 4. Were observers kept unaware of the experimental treatment imposed on the samples (for example, organisms, plots) when recording observations or measurements so as to minimize unconscious bias?

- If not stated, then request clarification in the manuscript of whether methods were adopted that reduced the possibility of unconscious bias influencing observations.
- If no steps were taken to prevent observer bias, request an explanation to appear in the manuscript of how unconscious bias could have influenced observations.

It is now well demonstrated that researchers’ observations are often influenced by what they expect to see<sup>12,13</sup>. For instance, when researchers were unaware of the colony of origin of the ants they were observing, they were more than three times more likely to report aggression between colony mates than were researchers who knew the ants’ colony of origin<sup>12</sup>. Keeping observers unaware of treatment categories or expected outcomes is not always possible or reasonable, but researchers should at least discuss the possibility of unconscious bias<sup>19</sup>.

#### 5. Did the authors explain how sample size was decided (for example, based on *a priori* power analysis or logistical constraints), or when an experiment with pre-set sample sizes was terminated?

If sample size or the end of the experiment was not decided prior to the initiation of the study, was there a decision rule for when to cease data collection?

- If not reported, request that authors provide this information.
- If the stopping rule included iterative statistical tests or examination of patterns as data accumulated, request that authors acknowledge the bias resulting from this process.

Cessation of data collection should never be made in response to reaching some threshold of statistical significance or effect. Such a practice leads to strong bias in favour of effects inflated by sampling error<sup>20,21</sup>. An explanation for the choice of stopping point should be provided (for example, ‘we planned to harvest samples at the end of the second growing season’).

#### 6. Did the authors develop their analysis plan, including choices of variables, without looking at the data, for instance prior to gathering data or with a dummy data set?

This is most easily determined by the existence of a pre-registered analysis plan. In the absence of pre-registration, a statement from the authors about the development of their analysis plan is still important.

- If no, request that authors acknowledge the exploratory nature of their analyses and declare that they are reporting the complete set of results from all exploratory analyses.
- If authors deviated from their analysis plan, request an explanation of why and how they deviated from the plan.

**Table 1 | Power to detect a true biological effect as statistically significant ( $P < 0.05$ ) as a function of sample size and actual effect size for two types of simple analysis**

Effect size			Sample size						
			10	20	50	100	200	500	
Correlation	<i>r</i>	Power (to detect a true effect)							
		0.1	Small	0.06	0.07	0.11	0.17	0.29	0.61
		0.3	Medium	0.14	0.26	0.57	0.86 <sup>a</sup>	>0.99 <sup>a</sup>	>0.99 <sup>a</sup>
		0.5	Large	0.33	0.64	0.97 <sup>a</sup>	>0.99 <sup>a</sup>	>0.99 <sup>a</sup>	>0.99 <sup>a</sup>
			Sample size (summed across both treatments in balanced design)						
			10	20	50	100	200	500	
Comparison of means (for example, <i>t</i> -test)	Hedge's <i>d</i>	Power (to detect a true effect)							
		0.2	Small	0.06	0.07	0.11	0.17	0.29	0.61
		0.5	Medium	0.11	0.18	0.41	0.7	0.94 <sup>a</sup>	>0.99 <sup>a</sup>
		0.8	Large	0.2	0.4	0.79 <sup>a</sup>	0.98 <sup>a</sup>	>0.99 <sup>a</sup>	>0.99 <sup>a</sup>

<sup>a</sup>High power, which is typically considered 0.8, or an 80% chance of detecting an effect, if the effect exists. Note that obtaining high power to detect small to medium effects (those most common in ecology and evolution) requires sample sizes much larger than are typical.

Choosing the analyses to present based on the strength of the effects derived from those analyses or models biases the distribution of presented results and can lead to presentation of entirely spurious relationships<sup>20,22</sup>. An ideal solution is to develop an analysis plan before examining the data and file it in a pre-registration archive such as offered by the Open Science Framework (<https://cos.io/pre-reg/>). One plausible alternative is an unusually detailed and publicly available grant proposal. Either way, the pre-registration or proposal should be cited in the manuscript. Researchers will sometimes need to deviate from pre-registered analysis plans. The pre-registration simply makes this transparent and gives the reviewer, and later the reader, the opportunity to assess whether deviations were sufficiently justified. Regardless of the availability of an analysis plan, reporting all versions of all analyses is essential for avoiding bias.

### 7. How suitable do you find the research methods without considering the outcome?

Evaluate the design and methods regardless of whether or not there was a finding of 'statistical significance', or whether or not the results conform to a predicted pattern.

- If the methods seem to have been flawed, call attention to the problems and, if possible, recommend a better design. Deciding whether the problems with the methods are sufficient to justify a recommendation of rejection will require your expert judgement.
- If uncertain about the suitability of some aspect of the methods, state your uncertainty to the editor so that she or he can seek appropriate methodological expertise as needed.

One driver of bias in the published literature is that we often evaluate the suitability of a study's methods based on the direction and strength of results<sup>23</sup>. This is especially true in cases of smaller samples or weaker study designs. In such cases, studies producing statistically significant or strong effects are sometimes incorrectly viewed as more plausible than those reporting weak or statistically non-significant results. There is a tendency among people we have talked with to assume that if a study found statistically significant results, sample sizes were sufficient or methodological weakness was not much of a problem. However, as 'strong' or 'significant' effects can often arise by chance<sup>24</sup> or be selected for reporting from other unreported results<sup>20,22</sup>, such results cannot be taken as proof that a study's methodological limitations were not a problem. Instead, the

quality of the methods must be judged independent of the results. (Of course, some studies include tests designed to assess a method's effectiveness rather than to assess the biological effect of primary interest, and those tests should be used to determine the quality of methods.) Doubts about the reliability of the methods should be given equal strength regardless of the primary outcome.

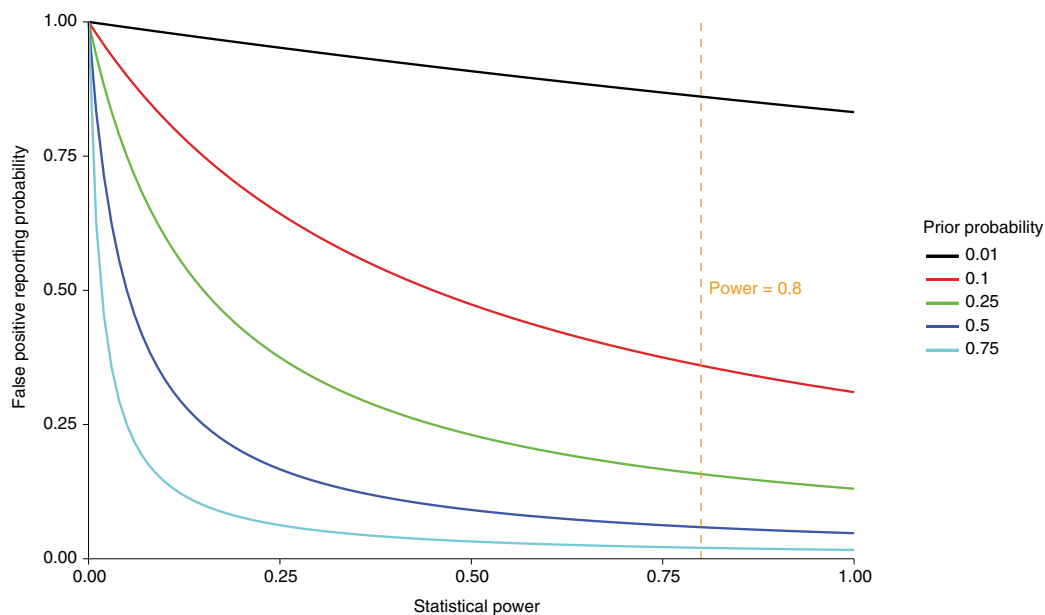
### 8. Are the sample sizes large enough to justify the authors' conclusions?

If presenting significance tests, how much power would this study have to detect statistically significant weak, moderate and strong effects? (See Table 1 for examples of how sample size and effect size combine to determine power in two types of simple analysis.) Expectation of effect size can best be derived from average effect sizes presented in meta-analyses of similar topics. The effect size reported in the manuscript under review can be a poor estimate of the underlying effect size, especially if the sample size is small thus elevating sampling uncertainty. Statistical significance is a poor indicator of the reliability of an estimate across a wide range of sample sizes and common effect sizes (Table 1 provides insight into statistical power).

- If sample sizes are small in a system where effects are expected to be weak to moderate, request that authors avoid inferences based on threshold  $P$  values, acknowledge uncertainty in effect size estimates and acknowledge the need for further study.
- Do not use sample size as a criterion for recommending publication unless you do so regardless of study outcome (that is, regardless of reported effect size and regardless of the outcomes of tests for significance).
- Do not use the failure to surpass a significance threshold as a reason to recommend rejection.

Presumably, nearly all ecologists and evolutionary biologists understand that there are problems with low power. However, it is clear that most of us would benefit from a reminder that type II error (false negatives) is only one of these problems. Because effect sizes are more variable with small samples, inflated effect sizes are more likely, and thus large effects derived from small samples can be unreliable<sup>25–27</sup>. In fact, with low power caused by some combination of a small sample and relatively weak biological effect, studies are likely to reach statistical significance only if sampling error drives the observed effect size much higher than the true effect<sup>25,27</sup>.





**Fig. 1 | Relationship between prior probability, statistical power and the false positive report probability.** The false positive report probability is the probability of a statistically significant result being a false positive (in other words, the probability that, in the case of a statistically significant rejection of the null, the null hypothesis is actually true). Note that for unlikely hypotheses, large portions of statistically significant findings will be false positives even with high power. This figure is based on a significance threshold of  $P < 0.05$ .

Unfortunately the weak- to moderate-strength biological effects that contribute to low power are common in ecology and evolutionary biology, at least in some sub-disciplines<sup>27,28</sup>. However, biological effects can be larger in some types of study and in some systems<sup>27,29</sup>, and so what qualifies as a small sample in one study may be sufficiently large in another. Thus, evaluating sample size and power will benefit from knowledge of the typical effect sizes for the type of study in question, and this can most reliably be learned by consulting meta-analyses. If the study under review seems to have low power, we should not consider meeting a threshold  $P$  value to be a reliable index of the validity of a pattern or a given effect size. In general, the reviewer should treat conclusions derived from low-powered studies as tentative, whether or not some significance threshold was met. However, studies with low power may often be worthy of publication, as some studies face major logistical obstacles regarding sample size, and it is only through publication and subsequent meta-analysis of a series of studies with small samples that we build a robust understanding of the true effect size<sup>27</sup>.

**9. What does the size of the estimated effect (for example, slope, correlation coefficient, difference in means) suggest about its biological or practical importance, and what does uncertainty around that effect estimate suggest about the estimate's precision?** Depending on the biological question, weak effects may be either biologically important or of limited interest; authors should justify their interpretation accordingly. Uncertainty around effects can be estimated with standard error of the mean (s.e.m.), 95% confidence intervals (approximately  $2 \times \text{s.e.m.}$ ), or with other statistics. As sample size increases (see checklist question 8 above) and variance decreases, s.e.m. decreases and we gain confidence in the precision of the effect estimate.

- If the authors do not interpret their results in terms of the biological relevance of the effect and the uncertainty surrounding their effect estimate, request that they consider doing so.

Evaluating results based on the size of the effect estimate and the associated uncertainty rather than based on a  $P$  value provides more direct insight into the biological phenomenon of interest<sup>30</sup>.

Too often, interpretation of results focuses on statistical significance rather than on biological significance, and thus we can be led astray regarding our understanding of their relevance.

#### 10. How unexpected would you judge these results to be in light of prior empirically derived understanding?

Effects that are more surprising in light of robust prior information are those that had a lower prior probability of being correct. When testing unlikely hypotheses, the chance that a statistically significant result is a false positive rises dramatically (Table 2, Fig. 1).  $P < 0.05$  is a poor threshold for evaluating the significance of an unexpected discovery and should be presented as no more than suggestive evidence for such discoveries.

- If a result is unexpected in light of prior evidence and is not supported by very strong new evidence (for example, multiple lines of convincing evidence), do not recommend against publication on these grounds, but request that the authors acknowledge the tentative nature of their results.

**Table 2 | False positive report probability (the probability that a statistically significant result is a false positive — in other words, the probability that, in the case of a statistically significant rejection of the null, the null hypothesis is actually true) as a function of prior probability and statistical power**

	Power			
	0.1	0.2	0.5	0.8
Prior	False positive report probability			
0.01	0.98	0.96	0.91	0.86
0.1	0.82	0.69	0.47	0.36
0.25	0.60	0.43	0.23	0.16
0.5	0.33	0.20	0.09	0.06
0.75	0.14	0.08	0.03	0.02

Note that for unlikely hypotheses, larger portions of statistically significant findings will be false positives. This table assumes a significance threshold of  $P < 0.05$ .

Findings should be interpreted in light of previously published information, and the more robust the body of pre-existing information, the more caution authors should exercise when interpreting the implications of their contradictory results. To quote Carl Sagan, “Extraordinary claims require extraordinary evidence”. For instance, many researchers in biology are unaware that the strength of evidence presented by a *P* value depends on the prior probability of the outcome. When testing moderately unlikely hypotheses (those with a 10% chance of being true) in a test with high statistical power, more than one-third of statistically ‘significant’ effects below the *P* < 0.05 threshold will be false positives (Table 2, Fig. 1)<sup>21</sup>. Thus, if robust pre-existing information makes a result unlikely, that result should be held to a higher standard of evidence than would be appropriate for a hypothesis that has already been empirically supported and thus has a higher prior probability<sup>31</sup>. For instance, a finding that parental diet influenced offspring phenotype is consistent with previously published findings and theory, but a finding that parental diet influenced grand-offspring phenotype more strongly than it influenced offspring phenotype would be extraordinary. Extraordinary results may be correct, but relative to results with a high prior probability, the extraordinary results are more likely to be false positives. We are not suggesting that reviewers estimate prior probabilities. However, a qualitative consideration of this issue is important for thoroughly evaluating the link between evidence and inference presented in a manuscript.

## Conclusions

We have designed this checklist for the use of reviewers, but we also hope that editors and authors will find it useful. Of course, reviewers are also authors, and many editors are also researchers, and understanding of the issues raised here can contribute to excellence in scientific publication in ecology and evolution in many ways. Currently, a small number of journals where ecologists and evolutionary biologists publish have adopted checklists for authors that rigorously address some of the issues we raise here (for example, Nature journals (<https://www.nature.com/authors/policies/ReportingSummary.pdf>), Conservation Biology ([https://mc.manuscriptcentral.com/societyimages/conbio/checklist\\_26.08.2016.docx](https://mc.manuscriptcentral.com/societyimages/conbio/checklist_26.08.2016.docx))). These are important steps forward. As such checklists become more widespread, they should reduce the need for separate reviewer checklists. However, until rigorous author checklists designed to promote transparency and reduce bias are standard across journals, checklists such as this one will continue to play an important role. And even when author checklists become widespread, reviewers will still have an important function because, in their role as reviewers, they are not subject to the incentives that might lead authors or editors to publish biased subsets of results or be insufficiently transparent.

How will peer review checklists be received by peer reviewers? Journal editors often struggle to recruit the necessary two or three reviewers per manuscript, and so editors are legitimately reluctant to do anything that makes reviewing seem more burdensome. However, even if journal editors decide not to make review checklists mandatory, they can still make them readily available to reviewers. Our discussions with new reviewers (for example, senior PhD students and post-docs) suggest that there is strong demand for this sort of guidance in peer reviewing papers.

Our checklist questions are practical tools. We hope that these questions, along with peer review checklists that address a broader set of topics (for example, TTEE), will improve transparency in the published literature and thus reduce bias therein. More transparency and less bias should mean more reliable inferences in published papers and later in the meta-analyses based on those

published papers<sup>1</sup>. As we improve peer review, we improve the quality of science.

Received: 25 October 2017; Accepted: 26 March 2018;

Published online: 22 May 2018

## References

- Parker, T. H. et al. Transparency in ecology and evolution: real problems, real solutions. *Trends Ecol. Evol.* **31**, 711–719 (2016).
- TTEE Working Group *Tools for Transparency in Ecology and Evolution (TTEE)* (Open Science Framework, 2016); <https://osf.io/g65cb/>
- Arriaga, A. F. et al. Simulation-based trial of surgical-crisis checklists. *New Engl. J. Med.* **368**, 246–253 (2013).
- Gawande, A. A. *The Checklist Manifesto: How to Get Things Right* (Metropolitan Books, New York, 2009).
- Gerstner, K. et al. Will your paper be used in a meta-analysis? Make the reach of your research broader and longer lasting. *Methods Ecol. Evol.* **8**, 777–784 (2017).
- Ferreira, V. et al. A meta-analysis of the effects of nutrient enrichment on litter decomposition in streams. *Biol. Rev.* **90**, 669–688 (2015).
- Fidler, F., Burgman, M. A., Cumming, G., Buttrose, R. & Thomason, N. Impact of criticism of null-hypothesis significance testing on statistical reporting practices in conservation biology. *Conserv. Biol.* **20**, 1539–1544 (2006).
- Zhang, Y., Chen, H. Y. H. & Reich, P. B. Forest productivity increases with evenness, species richness and trait variation: a global meta-analysis. *J. Ecol.* **100**, 742–749 (2012).
- Nickerson, R. S. Confirmation bias: a ubiquitous phenomenon in many guises. *Rev. Gen. Psychol.* **2**, 175–220 (1998).
- Fischhoff, B. Hindsight not equal to foresight – effect of outcome knowledge on judgment under uncertainty. *J. Exp. Psychol. Human.* **1**, 288–299 (1975).
- Kozlov, M. V., Zverev, V. & Zvereva, E. L. Confirmation bias leads to overestimation of losses of woody plant foliage to insect herbivores in tropical regions. *PeerJ* **2**, e709 (2014).
- van Wilgenburg, E. & Elgar, M. A. Confirmation bias in studies of nestmate recognition: a cautionary note for research into the behaviour of animals. *PLoS ONE* **8**, e53548 (2013).
- Holman, L., Head, M. L., Lanfear, R. & Jennions, M. D. Evidence of experimental bias in the life sciences: why we need blind data recording. *PLoS Biol.* **13**, e1002190 (2015).
- Lee, C. J., Sugimoto, C. R., Zhang, G. & Cronin, B. Bias in peer review. *Adv. Inf. Sci.* **64**, 2–17 (2013).
- Mislan, K. A. S., Heer, J. M. & White, E. P. Elevating the status of code in ecology. *Trends Ecol. Evol.* **31**, 4–7 (2016).
- Fidler, F. et al. Meta-research for evaluating reproducibility in ecology and evolution. *BioScience* **67**, 282–289 (2017).
- Parker, T. H. What do we really know about the signalling role of plumage colour in blue tits? A case study of impediments to progress in evolutionary biology. *Biol. Rev.* **88**, 511–536 (2013).
- Simmons, J. P., Nelson, L. D. & Simonsohn, U. A 21 word solution. *Dialogue* **26**, 4–7 (2012).
- Kardish, M. R. et al. Blind trust in unblinded observation in ecology, evolution and behavior. *Front. Ecol. Evol.* **3**, 51 (2015).
- Simmons, J. P., Nelson, L. D. & Simonsohn, U. False positive psychology: undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychol. Sci.* **22**, 1359–1366 (2011).
- Forstmeier, W., Wagenmakers, E.-J. & Parker, T. H. Detecting and avoiding likely false-positive findings – a practical guide. *Biol. Rev.* **92**, 1941–1968 (2017).
- Forstmeier et al. present insights that can help reviewers recognize and guide authors away from potentially biased and unreliable reporting.
- Forstmeier, W. & Schielzeth, H. Cryptic multiple hypotheses testing in linear models: overestimated effect sizes and the winner’s curse. *Behav. Ecol. Sociobiol.* **65**, 47–55 (2011).
- Palmer, A. R. Quasireplication and the contract of error: lessons from sex ratios, heritabilities and fluctuating asymmetry. *Annu. Rev. Ecol. Syst.* **31**, 441–480 (2000).
- Halsey, L. G., Curran-Everett, D., Vowler, S. L. & Drummond, G. B. The fickle *P* value generates irreproducible results. *Nat. Methods* **12**, 179–185 (2015).
- Gelman, A. & Weakliem, D. Of beauty, sex, and power. *Am. Sci.* **97**, 310–316 (2009).
- Barto, E. K. & Rillig, M. C. Dissemination biases in ecology: effect sizes matter more than quality. *Oikos* **121**, 228–235 (2012).
- Barto and Rillig provide evidence that various forms of bias, rather than concerns about data quality, have often influenced publication patterns in ecology.

27. Lemoine, N. P. et al. Underappreciated problems of low replication in ecological field studies. *Ecology* **97**, 2554–2561 (2016).  
**Lemoine et al. discuss how bias can emerge from low-powered studies, and also how bias can be avoided, even in systems where low power is inevitable due to logistical constraints.**
28. Møller, A. P. & Jennions, M. D. How much variance can be explained by ecologists and evolutionary biologists? *Oecologia* **132**, 492–500 (2002).
29. Duffy, J. E., Godwin, C. M. & Cardinale, B. J. Biodiversity effects in the wild are common and as strong as key drivers of productivity. *Nature* **549**, 261–264 (2017).
30. Nakagawa, S. & Cuthill, I. C. Effect size, confidence interval and statistical significance: a practical guide for biologists. *Biol. Rev.* **82**, 591–605 (2007).
31. Benjamin, D. J. et al. Redefine statistical significance. *Nat. Hum. Behav.* **2**, 6–10 (2018).

## Acknowledgements

We thank A. Moore for suggestions that improved the manuscript.

## Author contributions

T.H.P. composed the original draft of this manuscript in consultation with S.C.G. and S.N. S.N. made the figure. The manuscript was edited substantially over multiple rounds with input from all co-authors.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41559-018-0545-z>.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Correspondence** should be addressed to T.H.P.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.