

# Jury Theorems for Peer Review

Marcus Arvan, Liam Kofi Bright,  
and Remco Heesen

## Abstract

Peer review is often taken to be the main form of quality control on academic research. Usually journals carry this out. However, parts of maths and physics appear to have a parallel, crowd-sourced model of peer review, where papers are posted on the arXiv to be publicly discussed. In this paper we argue that crowd-sourced peer review is likely to do better than journal-solicited peer review at sorting papers by quality. Our argument rests on two key claims. First, crowd-sourced peer review will lead on average to more reviewers per paper than journal-solicited peer review. Second, due to the wisdom of the crowds, more reviewers will tend to make better judgements than fewer. We make the second claim precise by looking at the Condorcet jury theorem as well as two related jury theorems developed specifically to apply to peer review.

- 1 *Introduction*
- 2 *Assumptions of Peer Review*
- 3 *Crowd-Sourcing More Reviewers*
- 4 *The Basic Condorcet Jury Theorem*
- 5 *A Jury Theorem for Reviewer Scores*
- 6 *A Jury Theorem for Reviewer Reasons*
- 7 *Replies to Potential Objections*
  - 7.1 *Manipulation of reviewer scores?*
  - 7.2 *Greater average competence in journal-solicited peer review?*
  - 7.3 *Failures of independence in crowd-sourced peer review?*
- 8 *Conclusion*

## 1 Introduction

Peer review is supposed to secure an epistemic benefit. By ensuring that only work that has been validated by multiple experts is allowed into the academic literature, peer review is thought to

function as a quality control that prevents us from wasting time on poor work. Rather than have to wade through every half-baked flight of fancy, a discerning researcher may simply peruse peer-reviewed journals and read only that which passes peer review. However, in this essay we argue that an open, ‘crowd-sourced’ approach to peer review is more likely to reliably identify high-quality work compared to traditional, journal-solicited peer review.

The widely practised journal-solicited approach to peer-review filters the quality of academic work through a small number of experts—typically, a few editors and one to three outside referees. The normative assumption that appears to underlie this practice is the belief that a small number of experts (reading anonymized submissions) are the best mechanism for distinguishing between high-quality and low-quality work. In this system, quality assessment occurs in two stages: first, in pre-publication peer review, which sorts papers into journals; and second, by a journal’s readership post-publication. We take it that the first stage is intended to provide a proxy for the second, which is the more long-term assessment by the field. In this paper, we argue that a ‘crowd-sourced’ approach to peer review that bypasses the first stage—immediately opening up papers for evaluations by large numbers of readers—is likely to more reliably evaluate paper quality than the traditional model of peer review. In brief, we use the Condorcet jury theorem (Condorcet [1785]) and some closely related mathematical results to argue that a large number of evaluators is more likely to produce an accurate quality assessment of a paper than a smaller number of evaluators.

While we offer up some specifics of an open, crowd-sourced peer-review system, this paper is not intended to provide a full outline of such a system. In Sections 3–6 we offer enough of the details to make our comparative case. In Section 7, while considering some objections, we gesture to some further features one may wish to incorporate into such a system. For instance, one possibility we note is that in addition to crowd-sourcing from the academic community at large, it might be desirable to have a core of expert reviewers whose assessments are recorded separately. So, while we do not rest our case on such specifics of a crowd-sourced model, we are often (as in this case) supportive of particular proposals for how such a system might work. In our view, an experimental attitude to crowd-sourcing peer review will be a much better way to work out the details than any argument we could provide here.

We close this introduction by relating our argument to the previous literature. There is renewed interest in the epistemic benefits secured by large numbers of diverse agents (List and Goodin [2001]; Hartmann and Sprenger [2012]), including in the social epistemology of science (Heesen *et al.* [2019]; O'Connor and Bruner [2019]; Singer [2019]). Our intent is to bring this literature to bear on a concrete problem in the social epistemology of science, namely, peer review. Further, given the replication crisis, there has recently been interest in systematic failures of peer review (Romero [2016]; Heesen [2018]). Our paper offers a thoroughgoing solution to these problems. Like those who argue we should eliminate peer review in the context of project funding (Avin [2019]), we think we should abandon the idea that a small number of experts can reliably predict which ideas will be worth reading. We are not the first to suggest opening up peer review (see Gibson [2007]; Nosek and Bar-Anan [2012]; Heesen and Bright [2021]). However, we offer a novel argument for its epistemic benefits, claiming in particular that crowd-sourced peer review is better equipped than journal-solicited peer review to play the role that Heesen and Bright ([2021]) called 'epistemic sorting' (from the perspective of that paper, the present paper may be viewed as arguing that epistemic sorting should be viewed as a factor that favours abolishing prepublication peer review, rather than a neutral factor as Heesen and Bright argued). If our argument is sound, nothing should get a potentially deceptive stamp of authority through journal-solicited peer review. Instead, to gauge the quality of academic work, we should rely on the long-run and aggregated views of many diverse researchers.

## 2 Assumptions of Peer Review

Our purpose is to compare the present system of journal-solicited prepublication peer review against a crowd-sourced model. However, we expressly do not engage in an all-things-considered comparison. Rather, we focus on one goal of peer review (central to its defenders), namely, the selection of high-quality papers. Thus we set aside other goals, such as improving the quality of papers.

In this section we argue that if the present system of peer review really helps us pick out high-quality papers (however imperfectly), then research quality and the peer reviewers who assess it must satisfy certain assumptions. Our argument's structure is loosely analogous to a

transcendental deduction: we argue that without satisfying these assumptions, the idea that any form of peer review could successfully select for quality makes no sense. The two assumptions are 'competency' and 'intersubjectivity'.

We start with competency. We assume that researchers are at least better than chance at picking out papers of high quality or ranking papers according to quality. If this were false, then the current system could do no better than a system of random publication. So quality is the sort of thing that researchers can and do discern, and respond appropriately to given their reviewing task.

Our second assumption is that there exists broad (if rough) intersubjective agreement about what constitutes quality. The idea is that for any given academic paper, there is a specific notion of quality—which tracks what it means to be a good contribution to that particular topic—which is to some extent shared among readers and reviewers. Note that we are not saying that there is only one type of quality for all papers. Rather, we are assuming that once you fix facts about a paper's topic and the type of impact it is intended to have, then there are better and worse ways of fulfilling the paper's purpose. Nor are we saying that being the best possible version of a given idea is either necessary or sufficient for being a high-quality paper. A perfectly executed paper on a trivial topic may rightly be judged low-quality and a timely and important article may be judged high-quality despite a mediocre execution (one could think of the choice of topic as setting the range of the quality scale, and the execution of the paper as determining where in that range the paper lands). We are merely assuming that what constitutes quality may depend on the field, methodology, intended audience, importance of the specific topic to broader issues in its field, and other contextual features of the paper.

Even when relativized in this way, a unidimensional notion of quality invites scepticism. Kuhn ([1977], Chapter 13) emphasized that there are many respects in which a scientific theory (and by extension a scientific paper) might be judged good or bad, and different reviewers will weigh these differently (for a more formal approach to this issue, see also Okasha [2011]; Heesen [forthcoming]).

However, consider what would be the case if our assumption were false. If there is no intersubjectively agreed unidimensional notion of quality, then it is unclear what peer review is

doing. Why should the fact that reviewers like a paper give me any reason to think it worth my time? The fact that we assume that peer reviewers can assess quality belies a presupposition that relevant experts know what it would take to be a more or less worthy version of a given paper, and that we can reasonably expect some agreement on this point. We allow, and in fact it will be essential to our argument, that this agreement may be partial and accompanied by substantive and persistent disagreement on particular points.

Moreover, academic practices including hiring, tenure, and promotion decisions are at least partially based on the journals researchers have published in. This assumes that better journals are more likely to publish better papers, which in turn assumes that something worth caring about can be learned from peer reviewers' assessment. Thus, anyone whose opinion on anything has been influenced by where something was published is implicitly committed to intersubjectivity and competence. Hence, if peer review is to make sense as a system of quality control, then the following must hold: Paper quality is such that for any given paper there are better and worse ways it could fulfil its own potential, and the research community contains people who are competent to assess this.

While our focus is exclusively on the role of peer review in quality control, we feel it is useful at this stage to briefly discuss fairness considerations. Proponents of journal-solicited peer review often suggest that our current practices are the fairest method available, as a particular feature of the existing peer review process—*anonymization*—is vital for protecting against reviewer bias. However, *anonymization* also has downsides, as it can facilitate bullying and mercilessly harsh critique. Moreover, the credit incentive to put work out there depends heavily on one's name being associated with the work (Merton [1957]). We take no stance on *anonymization*. Our arguments instead support the following conditional claims: if *anonymization* is important for fairness in peer review, then an *anonymized crowd-sourced peer review process* would be superior to current processes; and if *anonymization* is unimportant, then a *non-anonymized crowd-sourced peer review process* would be superior. We can imagine experimenting with hybrid systems, for instance by mandating that new manuscripts (or new reviews) would be anonymous for an initial period of evaluation and response, before revealing the name and letting authors take credit for their contribution. There are a variety of ways this

could work, and experience will be the best guide in working out the details.

### 3 Crowd-Sourcing More Reviewers

One notable feature of journal-solicited peer review is that paper quality is judged by a small number of evaluators. First, papers are often read and 'desk-rejected' by a single editor. Second, when papers are sent out for review, they are typically reviewed by anywhere from one to three referees. Contrast this to the publication model already utilized in maths and physics. In these disciplines, it is standard for unpublished papers to be posted on individuals' professional websites and on central repositories, such as the arXiv. For this arXiv publication model to fully play its part in the crowd-sourced peer review system we imagine, there would have to be a disciplinary norm for peers to read and publicly evaluate new submissions.

Evidence for such a norm in maths and physics is anecdotal at best, as there are no formal studies of how widely arXiv preprints are discussed prior to journal publication. However, such discussions have not been altogether limited to preprints from well-known figures. For example, in 2007, a PhD researcher in physics who had left academia posted an arXiv preprint entitled 'An Exceptionally Simple Theory of Everything' (Lisi [unpublished]). Despite having no academic post, Lisi's paper was discussed at major physics blogs,<sup>1</sup> and in further arXiv preprints (Distler and Garibaldi [unpublished]). This process of online discussion led to a disciplinary consensus that Lisi's paper is flawed.<sup>2</sup> Likewise, this format does not inevitably give prominent figures an easy ride or render people unwilling to take on lengthy review tasks. In 2012, a well-known mathematician posted drafts of four preprints on his website, totaling around 500 pages, claiming to prove the *abc* conjecture (Mochizuki [unpublished]). A flurry of discussion on blogs and in follow-up preprints followed, identifying and debating potential flaws with the proof.<sup>3</sup> So at least in some instances the open disciplinary conversation seems to work well.

However, given the lack of systematic evidence, all we can confidently say is more evid-

<sup>1</sup> <[backreaction.blogspot.com/2007/11/theoretically-simple-exception-of.html](http://backreaction.blogspot.com/2007/11/theoretically-simple-exception-of.html)>.

<sup>2</sup> See <[www.scientificamerican.com/article/wipeout-theory](http://www.scientificamerican.com/article/wipeout-theory)>.

<sup>3</sup> For discussion of the mathematical content, see (Scholze and Stix [2018]) and <[www.galoisrepresentations.com/2017/12/17/the-abc-conjecture-has-still-not-been-proved](http://www.galoisrepresentations.com/2017/12/17/the-abc-conjecture-has-still-not-been-proved)>; and for a timeline up until April 2020 when the four preprints (Mochizuki [2021a], [2021b], [2021c], [2021d]) were accepted for publication, see <[twitter.com/andrewaberdein/status/1246553878553939980](https://twitter.com/andrewaberdein/status/1246553878553939980)>.

ence is needed as to what would happen if our proposal were implemented. We think it is at least plausible that informal social norms would arise and be generally adhered to requiring that people take part in crowd-sourced peer review. After all, participation in journal-solicited review is also largely a result of soft pressure from such a social norm, and there is no particular reason to think this would be different. Should the system allow reviews to be signed, then this would make it more observable whether one was actually complying with the social norm to contribute review labour, which may increase uptake among less conscientious academics (especially if coupled with targeted incentives, such as recognizing ‘top reviewers’, like the film review website Rotten Tomatoes does). In any case, we put some faith in the proposition that academics are opinionated about research in their field, and would not shirk the opportunity to voice their thoughts.

It is worth noting that our argument does not require the overall time spent reviewing papers to increase. Suppose that the current disciplinary norm to volunteer one’s time to review papers for journals is shifted over to the new system of crowd-sourced peer review, such that each member of the academic community volunteers exactly the same amount of time and reviews the same number of papers. Under journal-solicited peer review, journals base their decision about the quality of a paper only on the reviews solicited by that journal. For those papers that have already been rejected from other journals, all previous reviews are normally ignored. In contrast, under crowd-sourced peer review all reviews are public. Thus, if the total number of reviews remains constant, the average number of actually available reviews per paper under crowd-sourced peer review will be higher than under journal-solicited peer review.

A subtlety with respect to the point of the previous paragraph arises when we consider the possibility of authors revising their papers in response to reviews. We very much hope and expect that authors will continue to do this. To keep an accurate scholarly record, it would be best if each revision is clearly marked (with its own version number and DOI, as is already standard on the arXiv and other preprint servers) and reviews are dated and keyed to specific versions of the paper. But if reviews are specific to versions, this dilutes the number of reviews, potentially undermining our claim that crowd-sourcing increases the number of available reviews per paper. We offer two points to address this worry. First, if the average number of revisions (of a

paper in our crowd-sourced system) is lower than the average number of journal submissions (of a paper in the current model), then our claim from the previous paragraph still goes through even if 'reviews per paper' is weakened to 'reviews per revision of a paper'. Second, it seems too extreme to assume that reviews written for previous versions of a paper are necessarily completely irrelevant to the evaluation of later versions. If the revisions are fairly minor (or the scope of the review sufficiently sweeping), old reviews may still be useful. We might, for example, send automated emails to reviewers when a revision is uploaded, inviting them to resubmit their review (with or without amendments), if it is still relevant. Here again we call for experimentation regarding the details of implementation.

This subtlety notwithstanding, we will assume throughout the rest of this paper that moving to crowd-sourced peer review increases the average number of reviews per paper. For a reader who thinks that crowd-sourced peer review will lower the number of reviews per paper, the arguments presented below will favour journal-solicited peer review over crowd-sourced peer review. Alternatively, they could be read as a normative argument in favour of increasing the average number of evaluators of any given paper.

#### **4 The Basic Condorcet Jury Theorem**

The Condorcet jury theorem shows that subject to three assumptions, the judgements of a jury—a group charged with voting on the truth of a proposition (where a majority vote wins)—have a greater probability of accuracy the greater the number of people in the jury. The first assumption is that there is a correct answer: the proposition the jury is judging is either true or false. The second assumption is that every member of the jury has some probability of voting for the correct truth-value of the proposition that is probabilistically independent of the other jury members' vote. Finally, the third assumption is that the average probability that any individual in the jury votes correctly is greater than 0.5.<sup>4</sup> The theorem then says that adding more voters to the pool (keeping the average probability of voting correctly constant) makes it progressively less likely that the majority vote for the wrong conclusion.

Here is an intuitive illustration of the theorem. Suppose the average probability that a juror

---

<sup>4</sup> This generalizes the original theorem by allowing individual probabilities to vary.

votes for the right answer is 0.51. If only 100 people serve on the jury, then the most likely result is that 51 jury members will vote for the correct answer and 49 for the wrong answer. If, however, just one additional jury member votes wrongly, then the result will be a tie. And if two additional jury members vote wrongly, then the jury will vote 51–49 for the wrong verdict. So it is not unlikely for this jury to go wrong (this happens with a probability of approximately 0.38, with an additional 0.08 probability of a tie). Now consider a jury of 100,000.<sup>5</sup> If the average probability of a correct vote remains 0.51, the most likely result is that 51,000 jury members will vote for the right verdict and 49,000 for the wrong verdict. Consequently, a thousand additional jury members would have to make a mistake to shift the jury from the single most likely outcome to the wrong result. This occurs only with a probability of around one in ten billion.

The theorem also shows that in the limit (an infinite-sized jury), the majority will vote for the correct answer with probability one (that is, 100% of the time). The relevant point for our purposes, however, is the comparative claim: the more jury members there are, the more likely it is that a majority of them will vote for the correct answer. This is important in light of our assumption that the typical article in a crowd-sourced peer review model will be read and evaluated by more people than under journal-solicited peer review.

To see how the Condorcet jury theorem plausibly supports crowd-sourced peer review, compare the assumptions of the theorem to those discussed in Section 2. The theorem's first assumption is that the proposition being judged is true or false. In the case of peer review this proposition would be something like 'This paper is of high quality'. This aligns closely with the second assumption we argued peer review must satisfy: that there is an intersubjective quality standard for a paper on a particular topic. For the moment we are assuming that peer reviewers give (only) a binary judgement of quality: thumbs up or thumbs down. One of the motivations of the models in Sections 5–6 is to consider more informative, graded reviewer judgements.

Now consider the Condorcet theorem's second assumption: that every jury member has an

---

<sup>5</sup> These numbers are wildly unrealistic for peer review as currently practised. Exaggerated numbers work better for illustrating how the theorem works. However, the theorem applies equally whether we are comparing a jury of three to a jury of five or comparing juries of 100 to 100,000: the larger the jury the higher the probability that the majority judgement is accurate. The only nuance to this is that differences in competence among jury members are more of a concern with smaller juries. We address this issue in Section 7.2.

independent probability of voting for the correct result. This assumption does not correspond directly to one of the assumptions from Section 2. Under crowd-sourced peer review, we can imagine reviewers' judgements becoming correlated due to reviewers being able to read other reviews, whereas under journal-solicited peer review, the active hand of an editor may likewise induce correlation of reviewers' judgements. So whether the independence assumption is satisfied may well depend on the mechanism by which the different peer review systems are implemented. We will say more about steps a crowd-sourced peer review model could take to ensure reviewer independence in Section 7.1 and Section 7.3. For now we emphasize that independence is assumed in the Condorcet jury theorem, and hence the real-world applicability of our argument in this section hinges on providing a mechanism to guarantee it.

Finally, consider the Condorcet theorem's third assumption: that on average voters' probability of voting for the correct answer is better than chance. This corresponds to the first assumption we argued that peer review must satisfy: that reviewers are competent at picking out high-quality papers.

Our claim is that for peer review to reliably select papers for publication on the basis of quality, two of the three assumptions of the Condorcet jury theorem must be satisfied. Moreover, the third assumption (independence) will be satisfied by crowd-sourced peer review if the latter is carefully implemented (we defer our discussion of this to Section 7). But then a crowd-sourced peer review model is more reliable than a journal-solicited peer review model. For we have argued in Section 3 that crowd-sourced peer review will tend to base evaluations on the judgement of a larger jury than journal-solicited peer review. And by the Condorcet jury theorem a larger jury is more likely to arrive at an accurate evaluation.

## 5 A Jury Theorem for Reviewer Scores

While we find the argument of the previous section convincing, we recognize that the model is highly idealized and thereby open to objections. In this section we provide a model intended to be more tailored to the specifics of peer review, and show that an analogous theorem holds. This shows that the jury theorem is robust against certain changes, thus strengthening our argument.

The basic Condorcet model assumes that agents make a binary judgement on a single pro-

position. In contrast, real peer reviewers provide more nuanced judgements. These may come in the form of numerical scores or qualitative reasons for the reviewer's verdict. This section considers a model of peer review where reviewers only provide a numerical score; we add qualitative reasons in the next section.

Whereas in the previous section the goal was to evaluate the truth value of 'This paper is of high quality', now the goal is to rank papers, with those ranked highest most recommended to the attention of other researchers. By the intersubjectivity assumption, for any two papers, one can accurately be said to be of higher quality than another.

Each review consists of a numerical score, which is the reviewer's estimate of the paper's quality. We write  $q_i$  for the quality estimate provided by reviewer  $i$ . By the competence assumption, reviewers tend to give higher scores to better papers. But as in the previous section, we assume that there is some random variation in reviewer scores, reflecting individual reviewer biases and idiosyncrasies. Also as in the previous section, we assume that this variation is independent across reviewers, so reviewer scores can be modelled as independent random draws from a large pool of potential reviewers or reviewer scores (we refer again to Section 7.3 for more discussion of the independence assumption).

In this setup, we can represent the competence assumption by assuming that, on average, reviewers agree on the quality of a paper (that is, for each paper there is a number  $q$  such that  $\mathbb{E}[q_i] = q$  for all  $i$ ). And we can represent reviewer biases and idiosyncrasies by assuming that there is some random variation around this average ( $\text{Var}[q_i] = \sigma^2 > 0$  for all  $i$ ). The intersubjectivity assumption is reflected in the different average for each paper (if  $q_i$  and  $r_i$  are reviewer  $i$ 's scores for distinct papers, the former is intersubjectively better if  $\mathbb{E}[q_i] > \mathbb{E}[r_i]$ ).

Given differing quality estimates from reviewers that are each taken to be competent, a journal editor or arXiv reader may take the average of these estimates to be her best estimate of the relative quality of a paper. Averaging in this way has been defended in the literature on combining forecasts (Clemen [1989]; Armstrong [2001], especially p. 422) and peer disagreement (Elga [2007]; Christensen [2007]; Cohen [2013]), while (weighted) linear averaging more generally has been widely defended by formal epistemologists (Lehrer and Wagner [1981]; Klein and Sprenger [2015]; Martini and Sprenger [2017]; Pettigrew [2019]). So the quantity of in-

terest that will be used to make decisions under either journal-solicited or crowd-sourced peer review is the average of  $n$  reviewer scores  $q_1, q_2, \dots, q_n$ , which may be written  $\frac{1}{n} \sum_{i=1}^n q_i$ .<sup>6</sup>

Because individual reviewer scores are equal in expectation, so is the average reviewer score (that is,  $\mathbb{E}[\frac{1}{n} \sum_{i=1}^n q_i] = \mathbb{E}[q_1]$ ). Perhaps more importantly, the random variation in the average reviewer score will decrease as the number of reviewers increases, according to the formula  $\text{Var}[\frac{1}{n} \sum_{i=1}^n q_i] = \sigma^2/n$ . This means that the more reviewers there are, the smaller the probability that the average reviewer score will be much different from its expectation. Since better papers have a higher expectation, the probability that two papers are ranked incorrectly will similarly decrease.

This gives a clear analogy to the Condorcet jury theorem. Previously, the probability of a correct verdict increased, whereas here the probability that papers are ranked correctly increases with the number of reviewers. To complete the analogy, note that the random variation goes to zero in the limit with an infinite number of reviewers, meaning that papers are ranked correctly with probability one.

Once again, granted the assumption that crowd-sourced peer review will have on average more reviewers per paper than journal-solicited peer review, this yields an argument in favour of crowd-sourced peer review as more likely to yield accurate quality judgements.

## 6 A Jury Theorem for Reviewer Reasons

In this section we expand on the previous section's model by including reviewers' reasons for giving a particular (numerical) quality judgement. We represent these reasons by thinking of papers as having a number of features and peer reviewers as having opinions on which combinations of features make for a high-quality paper. More specifically, we assume there

---

<sup>6</sup> A potential objection here is that the average is only a meaningful quantity if reviewer scores are assumed to be measured on a cardinal scale, whereas arguably such judgements only have ordinal significance (for background on this classification, see Tal [2020], Section 3.2). This need not always be a problem in practice. For example, if the underlying distribution is symmetric, the mean coincides with the median, which is ordinally meaningful, so the average will 'accidentally' track something meaningful. But to address this worry more fully we might take the median of reviewer scores as the quantity of interest instead. If we change our assumptions appropriately (individual reviewer scores share a median, with a higher median for better papers), essentially the same argument goes through, as the variance of the median of reviewer scores will also decrease with more reviewers. Alternatively, if we assume reviewer scores are discrete rather than continuous, the jury theorem for the median proven by Morreau ([2021]) applies, so our argument goes through in that setting as well. Working with the median has another advantage: it is robust against outliers. We recommend using both where feasible (see Section 7.1). We thank Michael Morreau for suggesting both the objection and the response.

are  $m$  features that peer reviewers evaluate for a paper on a certain topic (recall that quality standards are paper-specific). A paper is represented by its feature coordinates  $x_1, x_2, \dots, x_m$ , which provide a numerical 'score' for how that paper does on each feature. We imagine that for each feature there is a 'golden mean' (possibly relative to the value of the other features) such that both more and less of that feature would make the paper worse in the eyes of the reviewer.

For example, say that for a given paper feature 1 concerns the paper's discussion of the external validity of its results, so  $x_1$  indicates how the paper scores on this feature. A low value of  $x_1$  might indicate that the discussion unnecessarily constrains the external validity (compared to what is justified per the scores on other features); a high value then says that the discussion generalizes the study's results too widely (in ways not sufficiently supported by the evidence). A medium value indicates a sensible discussion that avoids unsupported claims.

What might the set of features look like? Since our argument goes through regardless, we can remain agnostic between the following suggestions. First, the features might be Kuhn's criteria for theory choice: empirical adequacy, simplicity, and so on. Second, the features might be some variation on those that peer reviewers are explicitly asked to score papers on by journals, such as novelty and methodological soundness. Third, our preferred option, the features might be anything and everything peer reviewers use to evaluate papers, at as fine-grained a level as possible.

According to our intersubjectivity assumption, any given paper has an (intersubjectively agreed) quality. We conceive of both a paper's quality and any given reviewer's opinion of its quality as a function of that paper's feature coordinates  $x_1, \dots, x_m$ . Quality can then be characterized as something that looks like an epistemic landscape (in the sense of Weisberg and Muldoon [2009]; Alexander *et al.* [2015]; Thoma [2015]): each  $m$ -dimensional point  $x = (x_1, \dots, x_m)$  represents a possible paper, and the height of the landscape at that point is the quality of that paper. We define the function  $f : \mathbb{R}^m \rightarrow [0, \infty)$  to describe this epistemic landscape. That is,  $f(x)$  is the intersubjective quality of a paper with characteristics  $x$ .

In accordance with our competence assumption, peer reviewers (whether crowd-sourced or journal-solicited) assess the quality of a paper. However, they are not perfect. First, they may be biased in that the combinations of features they value highly are different from the combinations

that really constitute intersubjective quality. Second, there may be measurement error: peer reviewers may make mistakes in evaluating a paper on some or all features. We roll these two types of errors into a single bias  $b_i$  for reviewer  $i$ . The bias  $b_i$  is an  $m$ -dimensional point representing the total distortion in reviewer  $i$ 's assessment of quality, such that the reviewer's evaluation of a paper with characteristics  $x$  will be  $f(x + b_i)$ .

We denote by  $\mu$  the centre of mass of the epistemic landscape of intersubjective quality, and assume that it exists.<sup>7</sup> Consequently, for any reviewer  $i$ , the epistemic landscape characterizing how that reviewer assesses quality also has a centre of mass located at  $\mu - b_i$ .

As before, we assume that an editor or reader averages reviewer assessments to obtain her best estimate of the quality of a paper. So for a paper  $x$  reviewed by  $n$  reviewers the quality estimate is  $f_n(x) = \frac{1}{n} \sum_{i=1}^n f(x + b_i)$ .

The question then is whether the quality estimate improves with the number of reviewers. This is the question whether, for a paper  $x$ ,  $f_n(x)$  gets closer to  $f(x)$  as  $n$  increases. Depending on the shape of the landscape and reviewers' biases, this may be true for some values of  $x$  and false for others. What we would like to know, then, is whether for an arbitrary paper the quality estimate gets closer to the intersubjective quality with more reviewers, that is, whether the function  $f_n$  as a whole becomes more similar to the function  $f$  as  $n$  increases.

How do we characterize the similarity of two functions? Here we take the following approach: compare the centres of mass. The centre of mass measures the central tendency of a function, giving some indication of where in the landscape the highest peaks of quality occur. This is a fairly crude measure of similarity: two functions may have the same centre of mass but be dissimilar in other respects. However, it has the advantage of giving us a single number (or  $m$ -dimensional point, rather) for each function. This measure works well when the landscapes are single-peaked and mostly smooth, as two such landscapes with similar centres of mass will usually agree on relative judgements (which of two papers is better). The centre of mass of  $f_n$  is  $\mu - \frac{1}{n} \sum_{i=1}^n b_i$ , which we compare to  $\mu$ , the centre of mass of  $f$ .

<sup>7</sup> More formally, we assume  $\int_{\mathbb{R}^m} x_j f(x) dx$  is finite for each feature  $j$  and define  $\mu$  coordinate-wise by  $\mu_j = \int_{\mathbb{R}^m} x_j f(x) dx / \int_{\mathbb{R}^m} f(x) dx$ . Given our 'golden mean' approach to paper quality, this assumption is fairly innocent. It holds in particular if  $f$  has a finite maximum (as it does under any reasonable formalization of the 'golden mean' approach) and the features are measured on finite scales. If the features are measured on infinite scales, whether the assumption holds depends on how quickly quality drops off away from the maximum.

It remains to worry about how the reviewer biases are distributed. We assume that we can treat these as random variables. This need only be true in a subjective sense: you do not know in advance which reviewer and hence which bias will be selected. We assume that expected bias is zero (loosely speaking, this says that bias is equally likely to be in any direction) and expected variation in bias is finite.<sup>8</sup>

It follows that in expectation the centre of mass of estimated quality is equal to the centre of mass of intersubjective quality ( $\mathbb{E}[\mu - \frac{1}{n} \sum_{i=1}^n b_i] = \mu$ ), so on average there will be no bias at all. But assuming the biases of different reviewers are independent (see Section 7.3 for discussion), we also get that the probabilistic variation in the centre of mass of estimated quality decreases with the number of reviewers ( $\text{Cov}[\mu - \frac{1}{n} \sum_{i=1}^n b_i] = \Sigma/n$ ). This means that the centre of mass of estimated quality is more likely to be far away from the centre of mass of intersubjective quality if there are fewer reviewers, and more likely to be close if there are more. Moreover, since the variation reduces to zero in the limit,  $\mu - \frac{1}{n} \sum_{i=1}^n b_i$  probabilistically converges to  $\mu$ .

These results provide another close parallel to the Condorcet jury theorem. Estimated quality is likely to be closer to intersubjective quality as the number of reviewers increases, and they coincide (by our crude measure) in the infinite limit. We conclude that crowd-sourced peer review, insofar as it tends to involve a greater number of reviewers, outperforms journal-solicited peer review under the basic assumptions (competency and intersubjectivity) that are required for journal-solicited peer review to make sense.

## 7 Replies to Potential Objections

### 7.1 Manipulation of reviewer scores?

Our assumptions of reviewer competence and independence of reviewer judgements would not be plausible for crowd-sourced peer review if it is overwhelmed by internet trolls with a political agenda or other forms of organized manipulation. If people base their judgement of a paper on reasons orthogonal to its quality, then our reviewer competence assumption will not be satisfied.

---

<sup>8</sup> More formally, we assume (for all  $i$ ) that  $\mathbb{E}[b_i] = 0$  and  $\text{Cov}[b_i] = \Sigma$ . Here,  $0$  is the  $m$ -dimensional origin, and  $\Sigma$  gives the covariances between a particular reviewer's bias in each of the  $m$  features (not the covariances between different reviewers' biases). In virtue of being the covariance matrix,  $\Sigma$  is a symmetric and positive semi-definite  $m \times m$  matrix. We make no assumptions on  $\Sigma$  except that it is not the zero matrix.

If groups of people are mobilized to leave reviews of a paper without much thought, then our independence assumption fails. Note that the latter will be a problem for our independence assumption regardless of whether such a 'mass reviewing campaign' is ultimately motivated by scientific (as when a large research programme gangs up on a smaller one), political, or other reasons. This is an important worry for crowd-sourced peer review given the extent to which various social media have recently been overwhelmed by such phenomena.

It is tempting to address this issue by putting tight restrictions on who is allowed to review. There are a number of ways of doing this. We might use formal requirements such as possession of a doctorate or academic employment, or social requirements such as endorsement from existing reviewers or reviews rated sufficiently helpful by other reviewers. We might apply such requirements at a system-wide level (that is, to decide whether a given person is allowed to review anything at all) or at a subfield-specific level, by requiring, say, a doctorate in a specific subfield or endorsement from subfield specialists to be allowed to review papers in that subfield.

As explained below, there may be ways of testing whether such restrictions are necessary, and if so, what types of restrictions would function best. However, our current opinion is that such restrictions go against the spirit of the proposal of crowd-sourced peer review and will limit its advantages. Restricting who is allowed to review will lower the average number of reviews per paper, thus reducing the benefits from large numbers we have discussed. Moreover, such restrictions may reinforce existing disciplinary boundaries and subfield-level group-think (where it exists), whereas one of the key envisioned strengths of crowd-sourced peer review is that it will be easier for disparate fields to cross-pollinate, benefit from each other's insights, and correct each other's biases.

For these reasons we favour a system in which anyone is allowed to review anything, regardless of whether they are a recognized expert or even an academic. But this does not mean giving free rein to trolls, mobs, and other manipulation. We think there are various possible measures to guard against these.

Here is a relatively simple one. For each subfield, curate a set of expert reviewers in one of the ways suggested above, such as by having reviewers endorse each other's expertise in the

given subfield. Here we imagine subfields to be relatively small, say, more than twenty but less than a hundred endorsed reviewers per subfield. Then, for each paper, report both the overall average reviewer score and the average expert reviewer score when taking into account only reviewers endorsed for that particular subfield.

This system, familiar from the website Rotten Tomatoes—which reports qualified reviewer scores and general audience scores—has a number of advantages. Researchers who prefer something close to journal-solicited peer review can use the expert reviewer average, whereas those favouring the wisdom of the crowds can focus on the overall average. More importantly, one can look at both. When they are similar, either there were little or no non-expert reviewers, or the non-expert reviewers tended to agree with the expert reviewers. It gets interesting when there is significant divergence between the expert reviewer average and the overall average. This could be evidence of a mob coming in to manipulate the score. However, it could also be evidence of group-think within the subfield, exposed by the independent insights of outsiders. In any case, the divergent scores will be a signal that something is up (Nosek and Bar-Anan [2012], p. 238). Individual readers will be alerted that at least one of the scores is misleading and that blind reliance on averages is not advisable for this particular paper. Such readers would be encouraged to read and judge the paper for themselves, potentially leaving new, genuine reviews clarifying the epistemic contribution of the paper.

More generally, while we argue that using the overall average reviewer score from crowd-sourced peer review will give better quality judgements than journal-solicited peer review, it is emphatically not part of our proposal that overall average scores should be the only thing available. A lot of additional information should be made available to potential readers, so they can freely choose which metrics they think are more informative. This includes the median and mode of reviewer scores, the content and score of each individual review, the total number of reviews and the number of reviews by endorsed reviewers, the ranking of the paper relative to other papers in its subfield, and ratings of the helpfulness of individual reviewers.

Combining metrics will provide additional insight relevant to the problem of manipulation. For example, the median reviewer score is robust to manipulation as long as fewer than half of the reviewer scores are manipulated. Thus, papers with big gaps between the average score and

the median score should be treated with care. For another example, one would typically expect better papers to receive more reviews, as readers are attracted by the high score. Thus, papers with an unusually high number of reviews but a low average score should raise suspicions. The same thing goes for a paper where most of the reviews come from reviewers who have never reviewed anything else. Using this information, we think academics will be able to make use of crowd-sourced peer review to identify and read high-quality papers with minimal interference from manipulation.

There may also be technological or procedural ways to monitor and prevent internet mobbing and the like. First, statistical software might be used to detect highly correlated votes (as might be the case when a particular article initially receives a good proportion of positive reviews only to be followed by a quick succession of overwhelmingly negative reviews). Second, reviewers might be afforded the ability to 'flag' particular reviews as suspicious.

Third, the proprietors or 'section editors' of the arXiv-like site we propose could be alerted by their site's software to a large number of reviews for a given paper being posted by visitors from particular outside websites (such as a Reddit or Twitter thread advocating for 'review bombing' the paper). Although we would not necessarily advocate deleting suspicious reviews (which Rotten Tomatoes did in early 2019)—as suspicious reviews could well be genuine—proprietors could flag individual papers as potential victims of illicit reviewer behaviour, alerting other readers and reviewers to the possibility that the paper's scores may be corrupted.

Fourth, moderation to remove abusive (for instance, sexist or otherwise bigoted) reviews is within the spirit of our proposal. Such reviews would not add to the evaluation of the paper. However, in addition to their intrinsic cultural or moral harm, they could harm the epistemic performance of science by contributing to some groups being less able to have their claims fairly assessed by the scientific community. Whatever potential for abuse of power exists in this content moderation system is surely no worse than what present editors and reviewers have.

To be clear, although we think these interventions might be effective, we do not commit to any particular scheme for addressing illicit reviewing practices. These matters are probably best addressed through practical experimentation of the sort that review aggregators like Rotten

Tomatoes have and continue to do.

On that note, we are optimistic that our model's viability and any implementation issues could be examined empirically. One possibility would be for a central repository to roll out a 'beta' version of the system, implementing the model we outline above in a specific subfield, recruiting a batch of expert reviewers in that subfield, randomly selecting new unpublished manuscripts to receive reviews and ratings, and permitting those manuscripts to be reviewed and rated by expert reviewers and others.

If this beta is rolled out prominently, it might encourage participation from significant numbers of people in the profession. Such a trial could run for six months or longer. In addition to gaining feedback from the academic community on what works well and what does not, the implementers could collect data for statistical analysis over several years (say, one to three years). They might use citation counts and other measures of engagement to look at whether papers rated well by expert reviewers and general readership tend to have a greater impact on the field than low-rated papers. Or they might use textual analysis to see whether follow-up work by other authors is largely positive (constructive) or negative (critical), how this correlates with reviewer ratings, and how this compares to papers published using journal-solicited peer review. Although scientific quality is difficult to gauge, this could provide systematic statistical evidence (albeit defeasible) of whether the trial system tends to select higher-quality work than traditional peer review. Further betas (implementing tweaks in response to feedback) could be carried out and examined if the data appears promising.

## **7.2 Greater average competence in journal-solicited peer review?**

Our arguments assumed that reviewer competence is randomly distributed throughout the population of possible reviewers. More specifically, in Section 4 we assumed that the average probability that a reviewer in journal-solicited peer review will arrive at an accurate judgement of a paper's merit is the same as the average probability of an accurate judgement from a crowd-sourced reviewer. Similarly, in Sections 5–6 we assumed that the judgement of a randomly selected peer reviewer follows the same probability distribution whether the reviewer is journal-solicited or crowd-sourced.

However, some may doubt this. First, some might suggest that journal editors are likely to select substantially more competent reviewers than the average reviewer in the population, perhaps because they commission reviews from the most accomplished figures in the field. These reviewers, due to their exceptional achievements, may perhaps be expected to judge a given paper's merits more accurately or with less bias.<sup>9</sup> Second, some might argue that insofar as journal-solicited peer review commissions reviews by specialists in the paper's field, those specialists are likely to have higher accuracy or be less biased than a pool of reviewers that includes non-specialists. If as a result of either mechanism journal-solicited peer review reliably selects reviewers who are more competent than crowd-sourced peer review, then our arguments do not go through. In order to successfully defend a Condorcet-style argument in this case, we would need to show that the accuracy increase generated by increasing the size of the jury pool (through crowd-sourced peer review) is greater than the accuracy increase generated by how journal editors select reviewers.

Our reply to this concern is three-fold. First, the balance of present evidence does not support the empirical claim that journals select better-than-average reviewers. Second, there are a number of *prima facie* reasons to believe that journal-solicited reviewers are likely to be at least as biased as the population from which crowd-sourced peer review might draw. These sources of evidence together provide some grounds for thinking journal-solicited reviewers may be no more reliable, on average, than reviewers in a crowd-sourced model. Third, if it turns out that journal-solicited reviewers are systematically better than crowd-sourced ones after all, crowd-sourced peer review can still outperform journal-solicited review by using weighted averaging.

Arguments that journal-solicited reviewers are likely to be more competent to evaluate papers than readers at large tend to come from the armchair. However, two sources of empirical evidence collectively cast doubt on this intuition. First, empirical studies on the quality of journal-solicited reviews suggest very low interrater reliability (Lee *et al.* [2013], pp. 5–6; Bornmann [2011], p. 207). Interrater reliability measures the level of agreement between different reviewers judging the same paper, which is relevant here because disagreement imposes an

---

<sup>9</sup> See [philosopherscooon.typepad.com/blog/2018/12/incentivizing-better-reviewer-behavior.html](https://philosopherscooon.typepad.com/blog/2018/12/incentivizing-better-reviewer-behavior.html), where David Bourget notes that journal editors may aim to select more accomplished, senior scholars as reviewers for this reason.

upper bound on reviewer accuracy. In one study, interrater reliability barely exceeded chance (Kravitz *et al.* [2010]). In terms of the basic Condorcet model, this corresponds to probabilities of voting correctly barely exceeding 0.5. In short, we know of no clear empirical evidence that journal-solicited reviewers evaluate paper quality more reliably than crowd-sourced reviewers and the evidence that does exist judges the reliability of journal-solicited reviewers so poorly that it leaves little room for crowd-sourced reviewers to be less reliable.

Second, anecdotal reports suggest that reviewers at highly selective journals routinely misjudge papers that larger audiences have judged more accurately. Gans and Shepherd ([1994]) report how a variety of classic (including Nobel Prize-winning) economics articles were systematically rejected by top-ranked journals in the field. Anecdotally, this also happens in academic philosophy—for instance, Jason Stanley reported that four of his articles rejected from multiple highly ranked journals are now among the twenty most-cited articles in those very journals since 2000.<sup>10</sup> This phenomenon is further illustrated by a study in which twelve articles were submitted to the same highly ranked psychology journals that already published them (Peters and Ceci [1982]). Of the nine that made it past desk-rejection, eight of the papers were rejected without reviewers or editors realizing that the journal had already published them—in many cases on the basis of ‘serious methodological flaws’.

How can we square the intuitive thought that prestigious journals will select the most competent reviewers with the empirical research and anecdotes indicating that journal-solicited reviewers are highly unreliable? Although we can only speculate, there are reasons to think that journal-solicited reviewers are likely to have a variety of biases that can be expected to interfere with them reliably evaluating paper quality that crowd-sourced reviewers may fail to have, or may not have to the same extent.

First, insofar as journals tend to commission comparatively accomplished reviewers, these reviewers may have biases favouring their own views and work (as they have a vested interest in their views remaining influential) and biases for particular arguments (for example, by authors they admire or are personally acquainted with). Relatively accomplished reviewers thus plausibly have reasons to be biased in favour of the *status quo*, judging papers conservat-

---

<sup>10</sup> <[philosopherscocoon.typepad.com/blog/2015/09/stanley-on-peer-review.html](https://philosopherscocoon.typepad.com/blog/2015/09/stanley-on-peer-review.html)>.

ively to (perhaps subconsciously) preserve the prominence of their own favoured arguments and theories. The anecdotal cases above plausibly lend at least some support to this, as Nobel Prize-winning and highly cited works tend to make dramatic contributions, while a significant number of them have been rejected by journal reviewers. Further, empirical studies on formal peer review support the hypothesis of conservative bias (Luukkonen [2012]; Lee et al. [2013], pp. 9–10; Hug and Ochsner [2022]). In contrast, this particular kind of bias seems likely to be less prominent in the general reviewer population, which can be expected to have a more evenly distributed mix of accomplished and less-accomplished authors.

Second, incentives in the journal-solicited peer review system arguably buttress the above biases while plausibly introducing additional ones. First, when reviewers know that the journal they are reviewing for has a high rejection rate, their presumption may be that they should recommend rejection unless they are convinced the paper is excellent. This potentially lowers the risk of accepting bad papers but increases the risk of rejecting good papers. Anecdotally, this is borne out in journal-solicited peer review—as illustrated again by significant numbers of seminal economics papers being rejected by journals. Second, journal editors plausibly have grounds to err on the side of rejection as well, as publishing a bad paper may harm the journal's reputation. Finally, by explicitly selecting 'specialists' (people who have already published in the area) to review papers, journal-solicited peer review runs a serious risk of 'group-think'. Indeed, consider two causal antecedents of group-think: group cohesiveness and insularity (Janis [1972]). Both are arguably embedded in journal-solicited peer review insofar as editors seek out specialist reviewers—individuals who have chosen to work in a similar area, attend conferences together, and share unpublished work among each other. Conversely, journal-solicited peer review does not appear to regularly involve a practice empirical evidence suggests serves to prevent group-think: the stimulation of intellectual conflict (Turner and Pratkanis [1994]) through the inclusion of outside perspectives (Janis [1972], pp. 209–15). Because crowd-sourced peer review would not only invite more people to review papers, but also give reviewers the opportunity to contest each other's reviews (see Section 7.3), it would be more likely to generate the kind of intellectual conflict necessary for combating group-think. Similar arguments have been made by philosophers under the label of 'epistemic diversity': researchers actively pursu-

ing opposing theories or methodologies is often fruitful (Feyerabend [1975]; Lakatos [1978]; Longino [1990]; Kitcher [1993]; Zollman [2010]).

Thus, there are reasons to believe that even if journal editors recruit ‘the best reviewers’ (which we have no clear empirical evidence for), journal-solicited peer review introduces biases that are likely to be less pronounced or more evenly distributed in a general population of readers in the discipline. The ‘beta’ experiment suggested at the end of Section 7.1 could help support (or undermine) these claims.

To be clear, readers in the general academic population will tend to have biases of their own. Like journal-solicited reviewers, they plausibly have vested interests in advancing their own views, idiosyncrasies, and so on. Ideally, we would have some empirical information on how reliable crowd-sourced reviewers are, which we could compare directly to the reliability of journal-solicited reviewers. But no such evidence exists, to our knowledge. However, we think the considerations we have adduced in this section add up to a strong case against the claim that journal-solicited reviewer are more reliable. We have argued that (1) there is no evidence that journal-solicited reviewers are particularly reliable (in an absolute sense), which severely limits the extent to which they could possibly be more reliable than crowd-sourced reviewers, and (2) there are some positive reasons (supported by anecdotal evidence) to think that journal peer review processes introduce some systematic biases favouring conservatism and group-think. In contrast, crowd-sourced reviewers may be individually biased, but because these biases are not systematic, they will cancel each other out. Our three jury theorems suggest that this cancelling of biases will be more effective in a larger jury. With a small jury (say, two reviewers), distorting biases are more likely to produce the wrong verdict.

Consequently, we submit that there is no clear support for the proposition that journal-solicited reviewers are more accurate in their judgements of paper quality than academics at large. Given that the evidence is unclear, we believe it is more appropriate to assume that accuracy and bias are randomly distributed in an academic population—unless and until clear evidence is provided to the contrary.

We now briefly revisit the concern (Footnote 5) that in our initial presentation of the jury theorem we used wildly unrealistic numbers (comparing juries of 100 to 100,000). If, as we

just argued, it is reasonable to assume no systematic differences in reviewer accuracy between a journal-solicited and a crowd-sourced approach to peer review, then our argument supports the crowd-sourced approach even if the number of crowd-sourced reviews is only modestly higher on average (say, five reviews) than in traditional peer review (say, three reviews). This is because each additional reviewer increases the expected accuracy of the group.

If, on the other hand, in spite of the argument of this section, we did acquire clear evidence that journal-solicited peer reviewers are superior to crowd-sourced ones, then the overall argument of this paper would fail. However, even in this case there is room to benefit from a crowd-sourced approach, if we allow some additional assumptions. Here we take our cue from Klein and Sprenger ([2015]) and Martini and Sprenger ([2017]), who argue that with sufficient information about the relative expertise of reviewers, a competence-weighted average of the reviewer scores outperforms the straight (unweighted) average we used in Section 5. In this scenario, a version of our jury theorem still goes through. Assume that (a) we make sure the reviewers that would be solicited by the journal(s) still submit their review to the crowd-sourced system and (b) we have accurate information about the individual competence of reviewers. Now if we add additional crowd-sourced reviewers, no matter how incompetent relative to the journal-solicited ones (though they should still meet the minimal competency requirement of Section 5), the variance of the weighted average reviewer score will decrease.<sup>11</sup> That is, additional reviewers can still be used to improve the accuracy with which paper quality is assessed. However, this move only works with sufficient information about reviewer competence. In particular, if journal-solicited reviewers are significantly more competent than crowd-sourced ones, then the unweighted crowd-sourced average performs worse than the journal-solicited average. In such scenarios, competence-weighting is necessary rather than just helpful.

---

<sup>11</sup> Following Klein and Sprenger's notation, let  $\text{Var}[q_i] = \sigma_i^2$  denote the variance in reviewer  $i$ 's reviewer score (with lower variance indicating higher competence) and let  $c_i$  be the weight given to reviewer  $i$  so that the competence-weighted average is  $\sum_{i=1}^n c_i q_i$ . If the competences  $\sigma_i^2$  are known, the optimal weights are  $c_i^* = (\sum_{j=1}^n \sigma_j^2 / \sigma_i^2)^{-1}$  and the variance of the weighted average is  $\text{Var}[\sum_{i=1}^n c_i^* q_i] = \sum_{i=1}^n c_i^{*2} \sigma_i^2 = (\sum_{i=1}^n 1 / \sigma_i^2)^{-1}$ . Keeping the competences of the first  $n$  reviewers constant, adding an  $n + 1$ -st reviewer will decrease this variance regardless of the value of  $\sigma_{n+1}^2$ .

### 7.3 Failures of independence in crowd-sourced peer review?

Another worry is that the jury theorems hold only when votes are probabilistically independent. This assumption seems plausible for journal-solicited peer review: each reviewer judges a given paper without knowledge of what other reviewers think. Conversely, with crowd-sourced peer review, one (influential) reviewer's evaluation of a paper may affect others—potentially generating a snowball effect. When votes in a jury are correlated, the collective competence of a jury may, in the worst case, be lower than the competence of individual jurors (Kaniowski and Zaigraev [2011]). More generally, our theorems offer no guarantee that the collective competence of a larger group of correlated (crowd-sourced) reviewers will be higher than that of a smaller group of uncorrelated (journal-solicited) reviewers.

Our reply begins with a conceptual note. As Estlund ([1994]) points out, the mere fact that early reviewers might influence the opinion of later reviewers is not necessarily inconsistent with probabilistic independence of reviewer judgements. Estlund thus shows that the 'correlation' (in the intuitive sense) induced by early influential reviewers need not entail the specific type of correlation that would undermine our argument. It is important not to confuse the two notions of independence and correlation that are at stake here.

Questions about reviewer competence and correlation between their opinions would ideally be settled empirically. As discussed earlier, the empirical finding of low interrater reliability suggests that reviewer competence may indeed be low. We are not aware of any empirical work on correlations among reviewer opinions in a crowd-sourced model of peer review (perhaps the 'beta' experiment discussed in Section 7.1 could provide some). However, potentially relevant evidence comes from recent studies using surveys and prediction markets to see whether academics can predict the results of replications (Dreber et al. [2015]; Camerer et al. [2016], [2018]; Forsell et al. [2019]). There is a close analogy here, as the surveys measure individual academics' opinions without information about what others think (as in journal-solicited peer review) while the prediction markets measure opinions in the presence of information about others (like crowd-sourced peer review). The results provide reason for optimism. Prediction markets tend to do at least as well as surveys, suggesting that information about other academics' opinions does not introduce the sort of correlation that undermines the benefits of large

numbers.

We add some speculative reasons to doubt that the correlation of votes under crowd-sourced peer review would be high. Academic training and incentives encourage academics to evaluate and counter arguments they find unpersuasive. Consider the kinds of discussions that occur in journals after an article or book is published—say, John Rawls's book *A Theory of Justice*. Some commentators were highly critical of Rawls's arguments (Hare [1973a], [1973b]; Nowell-Smith [1973]; Barry [1973]); others more sympathetic (Mandelbaum [1973]). Scholars then began debating each other on particular 'flaws' in Rawls's book (for example, Bedau [1975]). Eventually, a general consensus emerged that Rawls's work is important yet flawed in particular ways. Crowd-sourced reviewers could be expected to do something similar: present their own evaluations of a given piece of work and contest others. Assuming such practices are central to crowd-sourced peer review, 'votes' will tend to be poorly correlated.

Moreover, we expect expert reviewers in particular to form their opinions independently. Recall that in Section 7.1 we introduced expert reviewers (based on colleagues' endorsements) whose scores potential readers may want to consider separately as a way to guard against internet trolls. Now a genuine expert cannot just be somebody who happens to have more true beliefs than average about a domain: such a person would lack what has been called 'contributory expertise' in the literature (Collins *et al.* [2016]). The contributory expert must also know the methods and heuristics one ought to adopt to reliably arrive at true beliefs about the topic (Licon [2012], p. 451). We take it that knowledge of these methods is meant to make the contributory expert someone whose beliefs are relatively independent of their peers, conditional on the truth. This point is strengthened by Estlund ([1994]), who observes that in the presence of even fairly high degrees of deference to influential reviewers, the kind of reviewer independence needed for our theorems to apply can be maintained as long as reviewers add at least a modicum of their own (truth-tracking) insight.

If this is granted, then at least the set of expert reviewers will satisfy the independence condition, and hence all the conditions for our jury theorems to apply. Assuming that the average number of expert reviewers under crowd-sourced peer review is at least as high as the average number of reviewers under journal-solicited peer review, it follows (at minimum) that

basing one's opinion on the average expert reviewer score is better than basing one's opinion on journal-solicited peer review.

We also think that, barring cases of internet mobs, correlation among non-expert reviewers will tend to be low, and so basing one's opinion on the overall average score is even better than basing one's opinion on the average expert reviewer score (as this will make for an even larger jury). Whether we are right about this will become clear over time as we experiment with crowd-sourced peer review.

## 8 Conclusion

We have argued that if the presuppositions which peer review is based upon are correct, then three jury theorems suggest that an open, crowd-sourced model of post-publication peer review would do better at directing researchers towards better work. This leaves two major questions open. First, are the presuppositions of peer review correct? The brief arguments we gave for them here should be supplemented with more sustained, and often empirical, socio-epistemic inquiry. However, it should be noted that if these presuppositions turn out to be false, journal-solicited peer review is arguably even less defensible than our argument suggests. Second, is it actually desirable to direct researchers towards the best work? This might seem good for individual researchers but that is not yet to argue for its socio-epistemic optimality (Mayo-Wilson *et al.* [2011]). We leave both these projects for future work.

We conclude by noting that the practical difficulties with implementing our proposal are substantial but not insurmountable. Online forums for public peer review (such as the arXiv in physics) already exist, even serving as a primary point of publication in some fields. It is not beyond our capacities to add the necessary features on some expanded version of these venues. What is presently lacking is the will. We thus hope that our paper goes towards building this will, such that researchers will become able to take further and more systematic advantage of the combined wisdom of the academic community.

## Acknowledgements

All authors contributed equally. We thank Justin Bruner, Allard Tamminga, Boudewijn de Bruin, Kevin Zollman, Michael Morreau, two anonymous referees, and audiences at the University of Groningen, Eindhoven University of Technology, and the University of Bayreuth for valuable comments and discussion. This research was supported by the Leverhulme Trust (Philip Leverhulme Prize 2020 to Liam Kofi Bright) and the Dutch Research Council (016.Veni.195.141 to Remco Heesen).

Marcus Arvan

*Department of Philosophy and Religion*

*University of Tampa*

*Tampa, FL, United States*

*marvan@ut.edu*

Liam Kofi Bright

*Department of Philosophy, Logic and Scientific Method*

*London School of Economics*

*London, United Kingdom*

*l.k.bright@lse.ac.uk*

Remco Heesen

*Department of Philosophy*

*University of Western Australia*

*Crawley, WA, Australia*

and

*Faculty of Philosophy*

*University of Groningen*

*Groningen, The Netherlands*

*remco.heesen@uwa.edu.au*

## References

- Alexander, J. M., Himmelreich, J. and Thompson, C. [2015]: 'Epistemic Landscapes, Optimal Search, and the Division of Cognitive Labor', *Philosophy of Science*, **82**, pp. 424–53.
- Armstrong, J. S. [2001]: 'Combining Forecasts', in J. S. Armstrong (ed.), *Principles of Forecasting: A Handbook for Researchers and Practitioners*, New York: Kluwer, pp. 417–40.
- Avin, S. [2019]: 'Centralized Funding and Epistemic Exploration', *British Journal for the Philosophy of Science*, **70**, pp. 629–56.
- Barry, B. M. [1973]: *The Liberal Theory of Justice: A Critical Examination of the Principal Doctrines in a Theory of Justice by John Rawls*, Oxford: Clarendon Press.
- Bedau, H. A. [1975]: 'Review of *The Liberal Theory of Justice: A Critical Examination of the Principal Doctrines in a Theory of Justice* by John Rawls by Brian Barry', *Philosophical Review*, **84**, pp. 598–603.
- Bornmann, L. [2011]: 'Scientific Peer Review', *Annual Review of Information Science and Technology*, **45**, pp. 197–245.
- Camerer, C. F. et al. [2016]: 'Evaluating Replicability of Laboratory Experiments in Economics', *Science*, **351**, pp. 1433–36.
- Camerer, C. F. et al. [2018]: 'Evaluating the Replicability of Social Science Experiments in *Nature* and *Science* between 2010 and 2015', *Nature Human Behaviour*, **2**, pp. 637–44.
- Christensen, D. [2007]: 'Epistemology of Disagreement: The Good News', *The Philosophical Review*, **116**, pp. 187–217.
- Clemen, R. T. [1989]: 'Combining Forecasts: A Review and Annotated Bibliography', *International Journal of Forecasting*, **5**, pp. 559–83.
- Cohen, S. [2013]: 'A Defense of the (Almost) Equal Weight View', in D. Christensen and J. Lackey (eds), *The Epistemology of Disagreement: New Essays*, Oxford: Oxford University Press, pp. 98–119.

- Collins, H., Evans, R. and Weinel, M. [2016]: 'Expertise Revisited, Part II: Contributory Expertise', *Studies in History and Philosophy of Science Part A*, **56**, pp. 103–10.
- de Condorcet, N. [1785]: *Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix*, Paris: Imprimerie Royale.
- Distler, J. and Garibaldi, S. [unpublished]: 'There Is No "Theory of Everything" inside  $E_8$ ', available at <[arxiv.org/abs/0905.2658](https://arxiv.org/abs/0905.2658)>.
- Dreber, A. *et al.* [2015]: 'Using Prediction Markets to Estimate the Reproducibility of Scientific Research', *Proceedings of the National Academy of Sciences*, **112**, pp. 15343–47.
- Elga, A. [2007]: 'Reflection and Disagreement', *Noûs*, **41**, pp. 478–502.
- Estlund, D. M. [1994]: 'Opinion Leaders, Independence, and Condorcet's Jury Theorem', *Theory and Decision*, **36**, pp. 131–62.
- Feyerabend, P. [1975]: *Against Method*, London: New Left Books.
- Forsell, E. *et al.* [2019]: 'Predicting Replication Outcomes in the Many Labs 2 Study', *Journal of Economic Psychology*, **75**, p. 102117.
- Gans, J. S. and Shepherd, G. B. [1994]: 'How Are the Mighty Fallen: Rejected Classic Articles by Leading Economists', *Journal of Economic Perspectives*, **8**, pp. 165–79.
- Gibson, T. A. [2007]: 'Post-publication Review Could Aid Skills and Quality', *Nature*, **448**, p. 408.
- Hare, R. M. [1973a]: 'Rawls' Theory of Justice—I', *The Philosophical Quarterly*, **23**, pp. 144–55.
- Hare, R. M. [1973b]: 'Rawls' Theory of Justice—II', *The Philosophical Quarterly*, **23**, pp. 241–52.
- Hartmann, S. and Sprenger, J. [2012]: 'Judgment Aggregation and the Problem of Tracking the Truth', *Synthese*, **187**, pp. 209–21.

- Heesen, R. [2018]: 'Why the Reward Structure of Science Makes Reproducibility Problems Inevitable', *The Journal of Philosophy*, **115**, pp. 661–74.
- Heesen, R. [forthcoming]: 'The Necessity of Commensuration Bias in Grant Peer Review', *Ergo*, available at <[philsci-archive.pitt.edu/15930/](http://philsci-archive.pitt.edu/15930/)>.
- Heesen, R. and Bright, L. K. [2021]: 'Is Peer Review a Good Idea?', *British Journal for the Philosophy of Science*, **72**, pp. 635–63.
- Heesen, R., Bright, L. K. and Zucker, A. [2019]: 'Vindicating Methodological Triangulation', *Synthese*, **196**, pp. 3067–81.
- Hug, S. E. and Ochsner, M. [2022]: 'Do Peers Share the Same Criteria for Assessing Grant Applications?', *Research Evaluation*, **31**, pp. 104–17.
- Janis, I. L. [1972]: *Victims of Groupthink: A Psychological Study of Foreign-Policy Decisions and Fiascoes*, Boston: Houghton Mifflin.
- Kaniovski, S. and Zaigraev, A. [2011]: 'Optimal Jury Design for Homogeneous Juries with Correlated Votes', *Theory and Decision*, **71**, pp. 439–59.
- Kitcher, P. [1993]: *The Advancement of Science: Science without Legend, Objectivity without Illusions*, Oxford: Oxford University Press.
- Klein, D. and Sprenger, J. [2015]: 'Modelling Individual Expertise in Group Judgements', *Economics and Philosophy*, **31**, pp. 3–25.
- Kravitz, R. L. *et al.* [2010]: 'Editorial Peer Reviewers' Recommendations at a General Medical Journal: Are They Reliable and Do Editors Care?', *PLOS ONE*, **5**, p. e10072.
- Kuhn, T. S. [1977]: *The Essential Tension: Selected Studies in Scientific Tradition and Change*, Chicago, IL: University of Chicago Press.
- Lakatos, I. [1978]: *The Methodology of Scientific Research Programmes*, Cambridge: Cambridge University Press.

- Lee, C. J., Sugimoto, C. R., Zhang, G. and Cronin, B. [2013]: 'Bias in Peer Review', *Journal of the American Society for Information Science and Technology*, **64**, pp. 2–17.
- Lehrer, K. and Wagner, C. [1981]: *Rational Consensus in Science and Society: A Philosophical and Mathematical Study*, Dordrecht: D. Reidel.
- Licon, J. A. [2012]: 'Sceptical Thoughts on Philosophical Expertise', *Logos & Episteme*, **3**, pp. 449–58.
- Lisi, A. G. [unpublished]: 'An Exceptionally Simple Theory of Everything', available at <[arxiv.org/abs/0711.0770](https://arxiv.org/abs/0711.0770)>.
- List, C. and Goodin, R. E. [2001]: 'Epistemic Democracy: Generalizing the Condorcet Jury Theorem', *Journal of Political Philosophy*, **9**, pp. 277–306.
- Longino, H. E. [1990]: *Science as Social Knowledge: Values and Objectivity in Scientific Inquiry*, Princeton, NJ: Princeton University Press.
- Luukkonen, T. [2012]: 'Conservatism and Risk-Taking in Peer Review: Emerging ERC Practices', *Research Evaluation*, **21**, pp. 48–60.
- Mandelbaum, M. [1973]: 'Review of *A Theory of Justice* by John Rawls', *History and Theory*, **12**, pp. 240–50.
- Martini, C. and Sprenger, J. [2017]: 'Opinion Aggregation and Individual Expertise', in T. Boyer-Kassem, C. Mayo-Wilson and M. Weisberg (eds), *Scientific Collaboration and Collective Knowledge*, Oxford: Oxford University Press, pp. 180–201.
- Mayo-Wilson, C., Zollman, K. J. S. and Danks, D. [2011]: 'The Independence Thesis: When Individual and Social Epistemology Diverge', *Philosophy of Science*, **78**, pp. 653–77.
- Merton, R. K. [1957]: 'Priorities in Scientific Discovery: A Chapter in the Sociology of Science', *American Sociological Review*, **22**, pp. 635–59.
- Mochizuki, S. [2021a]: 'Inter-universal Teichmüller Theory I: Construction of Hodge Theaters', *Publications of the Research Institute for Mathematical Sciences*, **57**, pp. 3–207.

- Mochizuki, S. [2021b]: 'Inter-universal Teichmüller Theory II: Hodge-Arakelov-Theoretic Evaluation', *Publications of the Research Institute for Mathematical Sciences*, **57**, pp. 209–401.
- Mochizuki, S. [2021c]: 'Inter-universal Teichmüller Theory III: Canonical Splittings of the Log-Theta-Lattice', *Publications of the Research Institute for Mathematical Sciences*, **57**, pp. 403–626.
- Mochizuki, S. [2021d]: 'Inter-universal Teichmüller Theory IV: Log-Volume Computations and Set-Theoretic Foundations', *Publications of the Research Institute for Mathematical Sciences*, **57**, pp. 627–723.
- Mochizuki, S. [unpublished]: 'Inter-universal Teichmüller Theory I–IV', available at [www.kurims.kyoto-u.ac.jp/motizuki/papers-english.html](http://www.kurims.kyoto-u.ac.jp/motizuki/papers-english.html).
- Morreau, M. [2021]: 'Democracy without Enlightenment: A Jury Theorem for Evaluative Voting', *Journal of Political Philosophy*, **29**, pp. 188–210.
- Nosek, B. A. and Bar-Anan, Y. [2012]: 'Scientific Utopia: I. Opening Scientific Communication', *Psychological Inquiry*, **23**, pp. 217–43.
- Nowell-Smith, P. H. [1973]: 'Review Symposium: I—A Theory of Justice?', *Philosophy of the Social Sciences*, **3**, pp. 315–29.
- O'Connor, C. and Bruner, J. [2019]: 'Dynamics and Diversity in Epistemic Communities', *Erkenntnis*, **84**, pp. 101–19.
- Okasha, S. [2011]: 'Theory Choice and Social Choice: Kuhn versus Arrow', *Mind*, **120**, pp. 83–115.
- Peters, D. P. and Ceci, S. J. [1982]: 'Peer-review Practices of Psychological Journals: The Fate of Published Articles, Submitted Again', *Behavioral and Brain Sciences*, **5**, pp. 187–95.
- Pettigrew, R. [2019]: 'On the Accuracy of Group Credences', in T. S. Gendler and J. Hawthorne (eds), *Oxford Studies in Epistemology*, Vol. 6, Oxford: Oxford University Press, pp. 137–60.

- Romero, F. [2016]: 'Can the Behavioral Sciences Self-Correct? A Social Epistemic Study', *Studies in History and Philosophy of Science Part A*, **60**, pp. 55–69.
- Scholze, P. and Stix, J. [2018]: 'Why *abc* Is Still a Conjecture', available at [www.kurims.kyoto-u.ac.jp/motizuki/SS2018-08.pdf](http://www.kurims.kyoto-u.ac.jp/motizuki/SS2018-08.pdf).
- Singer, D. J. [2019]: 'Diversity, Not Randomness, Trumps Ability', *Philosophy of Science*, **86**, pp. 178–91.
- Tal, E. [2020]: 'Measurement in Science', in E. N. Zalta (ed.), *Stanford Encyclopedia of Philosophy*, available at [plato.stanford.edu/archives/fall2020/entries/measurement-science/](http://plato.stanford.edu/archives/fall2020/entries/measurement-science/).
- Thoma, J. [2015]: 'The Epistemic Division of Labor Revisited', *Philosophy of Science*, **82**, pp. 454–72.
- Turner, M. E. and Pratkanis, A. R. [1994]: 'Social Identity Maintenance Prescriptions for Preventing Groupthink: Reducing Identity Protection and Enhancing Intellectual Conflict', *International Journal of Conflict Management*, **5**, pp. 254–70.
- Weisberg, M. and Muldoon, R. [2009]: 'Epistemic Landscapes and the Division of Cognitive Labor', *Philosophy of Science*, **76**, pp. 225–52.
- Zollman, K. J. S. [2010]: 'The Epistemic Benefit of Transient Diversity', *Erkenntnis*, **72**, pp. 17–35.