

Match made by BERT?

Towards Interpretable Paper-Reviewer Assignments in NLP

Anonymous ACL submission

Abstract

Both scientific progress and individual researcher careers depend on the quality of peer review, which in turn depends on paper-reviewer matching. Surprisingly, this problem has been mostly approached simply as an automated recommendation problem rather than as a matter where different stakeholders (authors, reviewers, area chairs) have accumulated experience worth taking into account. We present the results of the first survey of the NLP community, identifying common issues and perspectives on what factors should be considered in paper-reviewer matching. This study contributes actionable recommendations for improving future NLP conferences, and desiderata for interpretable peer review assignments.

1 Introduction

Peer review is increasingly coming under criticism for its arbitrariness. Two NeurIPS experiments (Price, 2014; Cortes and Lawrence, 2021; Beygelzimer et al., 2021) have shown that the reviewers are good at identifying papers that are clearly bad, but the agreement on the “good” papers appears to be close to random. Among the likely reasons for that are cognitive and social biases of NLP reviewers (see overview by Rogers and Augenstein, 2020), fundamental disagreements in such an interdisciplinary field as NLP, and acceptance rates that are kept low¹ irrespective of the ratio of high-quality submissions.

Such arbitrariness leads to understandable frustration on the part of the authors whose jobs and graduation depend on publications, and it also means lost time and opportunities (Aczel et al., 2021; Gordon and Poulin, 2009) for science overall. Reviews written by someone who does not have the requisite expertise, or does not even consider the given type of research as a contribution, it is a

¹<https://twitter.com/tomgoldsteincs/status/1388156022112624644>

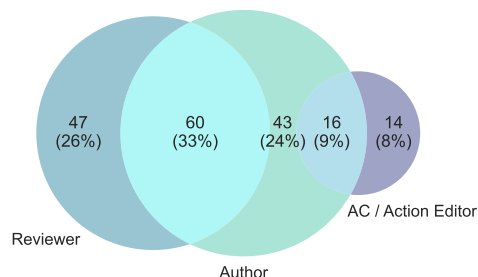


Figure 1: Overview of all respondents and overlap of their roles for their last experience at NLP venues.

loss for all parties: the authors do not get the intellectual exchange that could improve their projects and ideas, and reviewers simply lose valuable time without learning something they could use. It is also a loss for the field overall: less popular topics could be systematically disadvantaged, leading to ossification of the field (Chu and Evans, 2021).

This paper contributes a snapshot of this problem in NLP venues, based on a survey of authors, reviewers and area chairs (ACs). We collected 180 responses, which is comparable to the volume of feedback collected for implementing the ACL Rolling Review (ARR). The overall distribution of respondents’ roles is shown in fig. 1. We present the commonly reported issues and community preferences for different paper assignment workflows (section 4). We derive actionable recommendations to how peer review in NLP could be improved (section 5), discuss the limitations of survey methodology (section 6.2), and conclude with desiderata for interpretable peer review assignments (section 6.3).

2 Background: Peer Review in NLP

Paper-reviewer assignments are matches between submissions to conferences or journals, and their available pool of reviewers, taking into account the potential conflicts of interest (COI) and reviewer assignment quotas.

066 Among the systems used in recent NLP con- 112
067 ferences, the Softconf matching algorithm takes 113
068 into account bidding, quotas, and manual assign- 114
069 ments, and randomly assigns the remaining papers 115
070 as evenly as possible². ARR switched to Open- 116
071 Review and currently uses SPECTRE-MFR sys- 117
072 tem (OpenReview, 2021) for computing affinity 118
073 scores³, which the ACs may fully or partially rely 119
074 on. NAACL and ACL 2021 used SoftConf, but also 120
075 provided their ACs with affinity scores produced by 121
076 a “paraphrastic similarity” system (Wieting et al., 122
077 2019; Neubig et al., 2021). In our AC experience, 123
078 when program chairs provide affinity scores, they 124
079 caution against fully relying on them. 125

080 The problem of paper-reviewer assignment is 126
081 by itself an active area of research in several ar- 127
082 eas of computer science (Hartvigsen et al., 1999; 128
083 Wang et al., 2010; Li and Watanabe, 2013, inter 129
084 alia), including the problem of “fair” assignments 130
085 (Long et al., 2013; Stelmakh et al., 2019). Such 131
086 studies tend to be hypothesis-driven: they make an 132
087 assumption about what criteria should be taken into 133
088 account, design a system and evaluate it. To the 134
089 best of our knowledge, ours is the first study in the 135
090 field to address the opposite question: what criteria 136
091 should be taken into account, given the diversity of 137
092 perspectives in an interdisciplinary field? We take 138
093 that question to the community. 139

094 3 Methodology: survey structure and 140 095 distribution 141

096 We developed three separate surveys for the main 142
097 groups of stakeholders in the peer review pro- 143
098 cess: authors, reviewers and ACs. They follow 144
099 the same basic structure: consent to participation 145
100 (see Impact Statement), background information, 146
101 questions on most recent experiences in the role 147
102 which the survey pertains to, and how the respon- 148
103 dents believe paper-reviewer matching should be 149
104 performed. Most questions are asked to respon- 150
105 dents in all three roles, reformulated to match their 151
106 different perspectives. The responses were col- 152
107 lected late 2021 and all respondents are required 153
108 to confirm that their most recent experience as an 154
109 AC/reviewer/author is in 2019-2021. 155

110 The full surveys and response data are included 156
111 in the materials accompanying the submission, and 157

²<https://www.softconf.com/about/index.php/start/administration-view>

³Scores indicating how well a given submission matches a given reviewer, typically computed on the basis of their publication history. 158
159
160

will be made public upon acceptance. 112

Participant background. All surveys include 113
questions on career status and the number of times 114
the respondents have been authors/reviewers/ACs 115
at NLP venues. We ask what venues they have 116
experience with (as broad categories) and what 117
types of contributions they make in their work. 118

Participant experience with peer review. We 119
further ask the respondents a range of questions 120
about their experience as AC/reviewer/author: how 121
satisfied they are with the process, what issues they 122
have experienced, what was the assignment load 123
(ACs and reviewers), how paper-reviewer match- 124
ing was done, how they would prefer it to be done, 125
and which factors they believe to be important for 126
paper-review matching. Most of the questions are 127
multiple-choice, with addition of some open-ended 128
questions where appropriate, so that respondents 129
can elaborate their answers or add to the available 130
options. Whenever possible, the question formu- 131
lations were taken from the question bank of UK 132
Data Service (Hyman et al., 2006). Attitude ques- 133
tions use a 5-point Likert scale. 134

Limited memory is an important concern in sur- 135
veys (Sudman and Bradburn, 1973; Öztas Ayhan 136
and Isiksal, 2005), and we cannot expect the re- 137
spondents to accurately recall all their experience 138
with peer review. To reduce memory recall errors, 139
the survey focuses on the respondent’s most recent 140
experience, but they also have a chance to reflect 141
on prior experience in open-ended questions, and 142
to report whether they experienced certain issues 143
at any time in their career. 144

Survey distribution. We distributed the surveys 145
via three channels: by handing out flyers at 146
EMNLP 2021, through mailing lists (ML-news, 147
corpora list, Linguist list), and through Twitter with 148
the hashtag #NLProc. Participation was voluntary, 149
with no incentives beyond potential utility of this 150
study for improving NLP peer review. 151

Data validation. Given that links to surveys were 152
distributed openly and that we did not ask for any 153
identifiable information, the surveys needed to in- 154
clude other means of validation to ensure that the 155
responses included in the analysis were from atten- 156
tive, relevant individuals. Our approach for validat- 157
ing the data quality follows *satisficing theory* (Liu 158
and Wronski, 2018), with the main safeguards be- 159
ing 1) the checking of response consistency, includ- 160

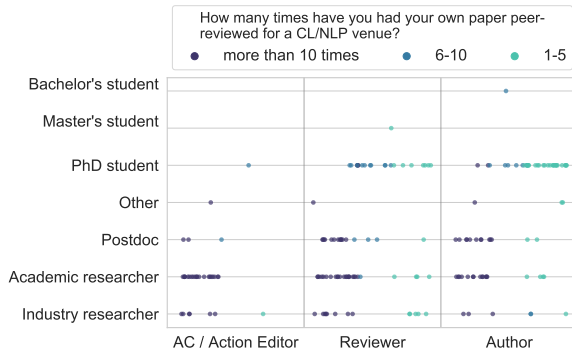


Figure 2: Career status of the respondents vs their experience receiving peer review. Numerical data is available in table 1 in the appendix.

ing a few “traps” where inconsistency or illogical responses can be exposed, and 2) the inclusion of open-ended questions.

73% ACs, 40% reviewers, 33% authors have provided at least one response to our open-ended questions, all found relevant in manual analysis. For consistency checks, all respondents were asked:

- How many times they have been an AC/reviewer/author. One of the options was “0”, contradicting the earlier confirmation of experience in a given role.
- When was the last time they were an AC/reviewer/author. One of the options was “earlier than 2019”, contradicting the earlier confirmation of peer review experience in 2019-2021.
- Whether they have performed the other roles. New authors may have not reviewed or AC-ed, but reviewers should also have been authors, and ACs should have experience with all roles.

4 Results

Overall we received 38 responses from ACs, 87 from reviewers and 81 from authors (206 in total). After removing 20 incomplete responses and 8 responses inconsistent with the “trap” questions, we report the results for 30 responses from ACs, 77 from reviewers and 73 from authors (180 in total).

4.1 Who are the respondents?

According to the past conference statistics, we could expect that many submissions would be primarily authored by the students, and reviewers are generally expected to be relatively senior, which should correspond to their going through peer review more often. We can use this expected pattern as an extra validation step for the survey responses.

Figure 2 shows that the responses are in line with this expected pattern. We received the most responses from academic researchers (62), PhD students (54), and postdocs (32). Most academic researchers and postdocs, but not PhD students, have had their work reviewed more than 10 times. At the same time 65% of the PhD students who served as reviewers went through peer review more than 5 times, as opposed to 24.2% of PhD students in the author role. Fewer industry than academic researchers responded to the survey.

4.2 Paper types

The next question is to see what kinds of research papers the respondents to our surveys have authored: engineering experiment, survey, position paper etc., according to the COLING taxonomy by Bender and Derczynski (2018). We expect that more senior researchers will have more experience with different types of work. Indeed, on average the authors have worked with 2.5 types of papers, vs. 3.0 for reviewers and 3.6 for ACs. The distribution is shown in fig. 3. The most respondents have authored engineering experiment papers (with the authors reporting the most such work).

Note that this only indicates whether the respondents to our surveys have or have not authored certain types of papers, rather than how many. In terms of volume, the engineering papers are a lot more prevalent: e.g. at ACL 2021 the “Machine learning” track had 332 submissions, vs 168 in the “Resources and evaluation” track (Xia et al., 2021).

4.3 What kinds of problems do people report?

As with any voluntary feedback, our surveys were likely to receive more responses from people who

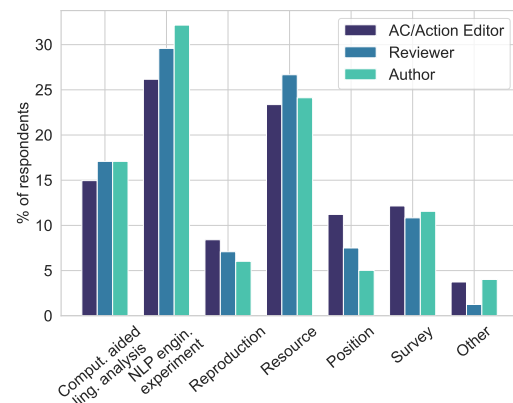


Figure 3: Types of research performed by respondents (multiple options could be selected).

had a grievance with the current process. Indeed, we find that only 6.7% of ACs, 20.5% of authors, and 22.1% of reviewers say that they have not had any issues in their last encounter with NLP venues.

The overall distribution for the types of problems reported by the authors, reviewers and ACs in their last and overall experience is shown in fig. 5, and the key takeaways are as follows:

- Two of the most frequent complaints of ACs (about 50% of the respondents) are insufficient information about reviewers and clunky interfaces;
- Many paper-reviewer mismatches (about 30%, if the report of the last experience is representative) are *avoidable*: they should have been clear from the reviewers’ publication history;
- Over a third of the author respondents in their last submission (about 50% over all history) received reviews from reviewers lacking either expertise or interest, and that is supported by the reviewers’ reports of being assigned papers that were mismatched on one of these dimensions;
- The authors report that many reviews (over a third in last submission, close to 50% overtime) are biased or shallow, which might be related to the above mismatches in expertise or interest.
- Two patterns are exclusive to ARR: insufficient time for ACs, and zero authors with no issues.

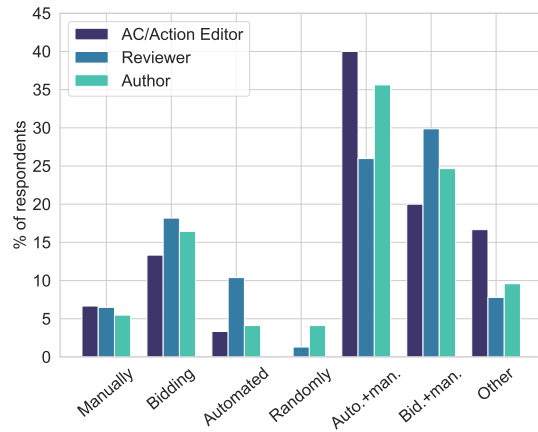
4.4 Knowledge of the workflow

Our next question is what methods NLP venues use to match submissions to reviewers, and to what extent the stakeholders (authors and reviewers) are aware of how it is done. We find that relatively few authors (23.3%) and reviewers (23.4%) know for sure what process was used, which begs for *more transparency in the conference process*. The ACs report that the most frequent case (37%) is a combination of automated and manual assignments. Interestingly, most reviewers believe that their assignments were automated (36%), and only (28%) believe they were automated+manual. See App. Figure 8 for full distribution.

5 The Ideal Process

5.1 Ideal workflow

When asked about what paper-assignment process they would prefer (given that fully manual matching is impractical for large conferences), most ACs and authors opted for automated+manual process, but for the reviewers this is the second preferred process (26%), with 30% opting for bidding + man-



Topics mentioned in the open-ended comments
(See supplementary materials for full categorized comments)

ACs: bidding (2), similarity+manual (1), similarity+bidding+manual (5), keyword-based filtering + bidding (2), similarity (1), tracks (1), other info (2), ARR (1), interface (1)

Reviewers: manual (2), similarity + bidding (3), similarity+bidding+manual (3), keywords (1), keywords+similarity (1), tracks (2), tracks+bidding (1), other (4)

Authors: against similarity (2), similarity + bidding (2), similarity+bidding+manual (2), ARR (2), random (2)

Figure 4: Which of the following options would you consider best for assigning reviewers to submissions?

ual checks (see fig. 4). There was also relatively large support for pure bidding (13-18% of respondents in all roles), and cumulatively pure bidding and bidding with manual adjustments have as much or more support from all respondent categories than the automated matching + manual assignments.

The analysis of open-ended comments suggested that the respondents were aware that bidding is quite labor-intensive on the part of the reviewers. 5 ACs, 3 reviewers and 2 authors suggested using affinity scores to filter the papers on which bids would be requested, followed with manual checking. Another suggestion was keywords or more fine-grained areas/tracks, potentially as alternative to affinity scores for filtering down the list of papers to bid on. One AC suggested “an extensive, but still finite, set of tags (e.g. an ACL-version of ACM CCS concepts, or FAccT’s submission tags”. One reviewer stressed that the keywords should be provided by the authors, to match what they perceive to be the salient aspects of the paper.

1 reviewer and 1 author suggested looking at whether the paper *cites* the potential reviewer⁴, as

⁴We believe this is an interesting idea, but it could lead to authors strategically placing citations to maximize the chances of acceptance, or being punished for citing work that they may criticize or claim to improve upon.



Topics mentioned in the open-ended comments
(The full comments categorized by these topics can be found in the ??)

Area chairs: interface issues (7), bad reviewers/reviews (5), workload issues (6), issues with ARR (4), lacking information on reviewers (4), communication issues between both systems and other human agents (4), lack of qualified reviewers in the pool (3), issues with meta-reviews (2), affinity score complaints (2), affinity score for finding reviewers the AC does not know personally (1), preference for manually recruited reviewers (1), papers assigned to ACs outside their area of expertise (1), too many declines (1), mismatch in goals of reviewers and authors (1), emergency reviews (1), bidding enabling bias (1).

Reviewers: choices forced by ACs (5), preference for bidding (4), areas of past expertise not currently of interest (4), lack of interest in the paper (3), methodological mismatch between generations of NLP researchers (3), mismatch in research methods (2), publication records as an unreliable indicator for assignments (1), mismatch in languages (1), time issues (1), reviewer bias (1)

Authors: reviewer expectation for a certain kind of research (6), inattentive reviews (5), short reviews (3), mismatch between the score and the text of the review (3), requests for irrelevant citations (2), confirmation bias (1), non-constructive criticism (1), shallow reviews (1), lack of reviewer competence (2), missing reviews (2), requests for irrelevant comparisons (1), "wild" estimates of impact (1), unannounced policy changes (1)

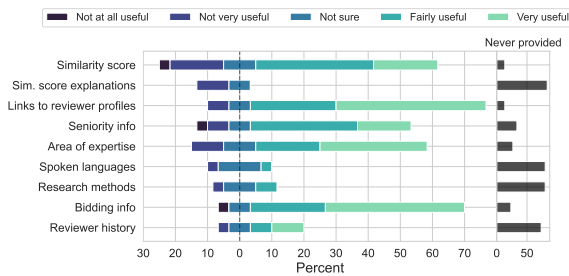
Figure 5: The issues with peer review process, reported by ACs, reviewers and authors, in their last (on the left) versus historical (on the right) experience with CL/NLP venues.

this could be a good indicator for the reviewer’s interest. 1 reviewer and 2 authors voiced support for some randomness in the assignments (given a track-level match): *“Bidding + some random assignment to ensure diversity in the matching. We don’t want reviewers to review only papers they *want* to review. However these random assignments should be clearly indicated to all, and treated accordingly.”*

5.2 Ideal assignment criteria

AC past experience. Figure 5 shows that one of the most common problems for the ACs is that they were not provided with enough information to facilitate the paper-reviewer matching. The follow-up question is what information they are provided with, and how useful they find it.

Figure 6 shows that the types of information with the highest utility information are links to reviewer profiles, bidding information, and affinity scores.



Topics mentioned in the open-ended comments: reviewer history (2), number of assigned papers (1), being able to ask SACs for advice (1), reviewer affiliation (e.g. academic or industry) (1), correct area match for both ACs and reviewers (1).

Figure 6: The utility of different kinds of information about reviewers that ACs may have been presented with to assist in manual checks of paper-reviewer matches.

But affinity scores are also the most controversial: it is the type of information that the most ACs find “not very useful” or “not useful at all” (20%).

Overall the results suggest that ACs are presented with little structured information about reviewers, and have to identify the information they need from a glance at the reviewers’ publication record. Seniority, expertise, and reviewer history notes from other ACs are all reported to be useful, but they were never provided directly to many ACs.

An avenue for future research is offered by three types of information that the most ACs are not sure about, presumably because they are rarely provided: structured information about the methods that the reviewers were familiar with, the languages they spoke, and affinity score explanations. We will show below that there is much support for taking such methods into account. For the languages, this might be due to the “default” status of English (Bender, 2019). We hypothesize that providing this information would make it easier to provide better matches for papers on other languages, which would in turn encourage the authors to submit more such work. Affinity will be discussed in section 6.3.

Stakeholder preferences. We then asked the respondents what factors they believe are important for paper-reviewer assignments. The overall mean importance rankings (on scale 0-5) are as follows:

- 3.95 Reviewer has worked on the same task
- 3.85 Reviewer bid on the paper
- 3.72 Reviewer has worked with the same method
- 3.32 Reviewer has authored the same type of paper
- 3.11 AC knows & trusts the reviewer
- 2.81 Reviewer has worked with the same kind of data
- 1.99 The affinity score is high

The fact that affinity scores rank the least important is an interesting reflection on how much NLP researchers would rely on their own techniques in a process with high stakes for themselves.

Our data suggests that within groups of stakeholders the individual variation in importance of different factors is higher for some factors and stakeholders than others: e.g. ACs vary within 1 point on the importance of knowing the data, but only within 0.74 points on importance of knowing the tasks. This has implications for approaches who would rely on AC assignments as ground truth for automated assignment systems: they could end up modeling the annotator instead of the task (Geva et al., 2019) See App. table 2 for full data.

We then explored the question of whether the experience of having authored research of a certain type correlates with any changes in the attitude towards some of these paper-reviewer matching factors. For each pair of type of research and matching factor, we ran two-sided Fisher’s Exact tests for all respondents who have authored (or not) the types of research and the importance they attached to different factors in paper-reviewer assignment (binning on less than moderately important and more than moderately important). For some pairs there were statistically significant differences: e.g. the respondents who have authored reproduction papers were significantly more likely to believe it important that the reviewer has worked with the same kind of data ($p = 0.008$), and less likely to believe it important that the reviewer has worked with the same method ($p = 0.036$). See table 3 in the appendix for all p -values and more details on the tests. We note that the relationship is not necessarily causal.

We conclude that our sample does provide evidence (the first, to our knowledge) that researchers in interdisciplinary fields who perform different kinds of research may have differing preferences for what information should be taken into account for paper-reviewer assignments. If that effect is robust, it should be considered in assignment systems for interdisciplinary fields. We hope that this finding would be explored in a larger study, taking into account both the experience of authoring a given type of paper and how central that type of research is for a given researcher (a factor that we did not consider). Another direction for future work is exploring this question from the perspective of demographic characteristics and the type of institution the respondents work in. Should there be

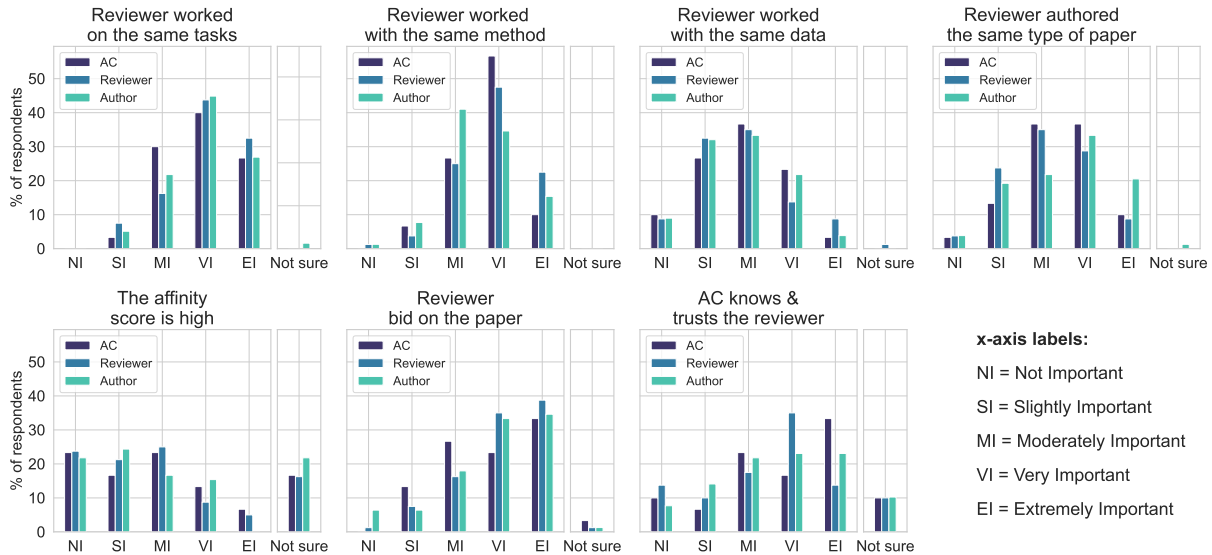


Figure 7: Question: *How important do you think the following factors are for a good paper-reviewer match?*

380 significant differences, more targeted assignments
 381 could be a powerful tool for diversifying the field.

382 5.3 Ideal workload

383 We asked our reviewer and AC respondents how
 384 many assignments they received at their most re-
 385 cent NLP venue, and what would be the optimal
 386 number (given a month to review, and a week for
 387 AC assignments). For ACs, the mean optimal num-
 388 ber of assignments is 8.5 ± 4.2 vs. 9.1 ± 5.1 they
 389 received at the most recent venue, and for review-
 390 ers it is 2.8 ± 1.0 vs. 3.3 ± 1.8 . Whether this is an
 391 issue depends on how much time a given venue
 392 allows. The ARR reviewers have even less than a
 393 month, and they indicated preference for fewer as-
 394 signments than they received (2.4 ± 1.0 vs 3.3 ± 1.9).
 395 See App. fig. 10 for data on other venues.

396 The lack of reviewers is a well-known problem.
 397 One of the possible causes is that many authors are
 398 students not yet ready to be reviewers. To investi-
 399 gate that, we asked the authors if they also reviewed
 400 for the venues where they last submitted a paper,
 401 and the reviewers and ACs - if they also submitted.
 402 If the core problem is that many authors are not
 403 qualified, we would expect more non-student au-
 404 thors to also be reviewers. Among all respondents
 405 there are 24% authors who submit to a venue but do
 406 not review there or help in some other role (fig. 1),
 407 but if we consider only non-student respondents
 408 that ratio is still 18% (see non-student role distri-
 409 bution in App. fig. 9). This suggests that *many*
 410 *qualified people do not review.*

411 6 Discussion

412 6.1 Reviewer interests

413 Our results suggest the lack of interest is one of the
 414 most common problems in paper-reviewer match-
 415 ing, for both authors and reviewers. The authors
 416 are aware of this problem and sometimes try to op-
 417 timize for it by pursuing the “safe”, popular topics.
 418 Unenthusiastic reviewers will likely produce shal-
 419 low, heuristic-based reviews, essentially penalizing
 420 non-mainstream research. Both tendencies con-
 421 tribute to ossification of the field (Chu and Evans,
 422 2021), and generally need to be minimized.

423 It is in the AC’s interest to find interested re-
 424 viewers, since that minimizes late reviews, but they
 425 need to know who finds what interesting. That is
 426 not as simple as a match by topic/methodology,
 427 clear from the publication record. Interests change
 428 not only gradually over time but also according
 429 to what is popular or *salient* at the given moment
 430 (Tversky and Kahneman, 1974; Dai et al., 2020),
 431 or even in seemingly irrational ways (e.g. by be-
 432 ing sensitive to the framing) (Tversky and Kahne-
 433 man, 1981). But although experience and knowl-
 434 edge may provide more stable descriptions of a
 435 reviewer, looking into dated publication records
 436 may be counter-productive. According to one of
 437 our respondents: *“I prefer the conferences who*
 438 *offer bidding processes to select the papers to re-*
 439 *view... I am more enthusiastic to review the papers*
 440 *compared to conferences that assign papers based*
 441 *on what my interests were x years ago.”*

442 Bidding however has its own set of problems,

including the practical impossibility to elicit all preferences over a big set of papers, the possibility of collusion rings (Littman, 2021), and, as one of our respondents put it, “biases towards/against certain paper types when bidding is enabled”. But these problems potentially have solutions: there is work on detecting collusion rings (Boehmer et al., 2022), and several respondents suggested that bidding could be facilitated by subsampling with either keyword- or affinity-score-based approaches.

We support some of our respondents’ recommendation for a combination of interest-based and non-interest-based (within a matching area) assignments, with the latter clearly marked as such for ACs and reviewers, and separate playbooks for the two cases. The reviewer training programs should aim to develop the expectation that peer review is something that combines utility and exploration.

6.2 Limitations

We readily acknowledge that, like with any surveys with voluntary participation, our sample of respondents may not be representative of the field overall, since the people who have had issues with peer review system are more incentivized to respond. However, precisely for that reason this methodology can be used to learn about the commonly reported types of problems, which was our goal. Our response rate turned out to be comparable to the response rate of the official ACL survey soliciting feedback on its peer review reform proposal (Neubig, 2020), which received 199 responses.

It is an open problem how future conferences could systematically improve, if they cannot rely on surveys to at least reliably estimate at what scale an issue occurs. Asking about satisfaction with reviews does not seem to produce reliable results (Daumé III, 2015; Cardie et al.). Our survey included a question about satisfaction with the paper-reviewer matching, and whether the most recent experience was better or worse than on average. Both reviewers and authors were more satisfied than dissatisfied, and considered the recent experience better than on average, despite reporting so many issues (see App. fig. 12 for the distribution).

6.3 Interpretable Paper-Reviewer Matching: Problem Formulation

There already are many proposed solutions for paper-reviewer matching (see section 2), but their evaluation is the more difficult problem. The obvious approach would be to use bidding information

or real assignments made by ACs as ground truth, but this data is typically not shared to protect reviewer anonymity. It would also provide a very noisy signal not just due to different assignment strategies between ACs, but also different quality of assignments depending on how much time they have on a given day. Both ACs and bidding reviewers are also likely⁵ to favor top-listed candidates. And, as our findings suggest, the optimal assignment strategies in an interdisciplinary field might genuinely vary between different types of papers and tracks. A system unaware of that might systematically disadvantage whole research agendas.

Given that even the human experts cannot tell what the best possible assignments are, we propose to reformulate the problem as *interpretable paper-reviewer matching*. That problem is *not* the same as the problem of faithfully explaining why a given paper-reviewer matching system produced a certain score, for which we have numerous interpretability techniques (Søgaard, 2021). The AC goal is fundamentally different: not to understand the system, but to quickly find the information that the AC⁶ considers relevant for making the best possible match. Therefore *the task of interpretable paper-reviewer matching is rather to help to identify the information that the stakeholders wish the decisions to be based on, and to provide that information as justification for the decisions*.

7 Conclusion

We present the results of the first survey on paper-reviewer assignment from the perspective of three groups of stakeholders in the NLP community: authors, reviewers, and ACs. The results point at a host of issues, some immediately actionable (e.g. providing the ACs with better information), some normative (e.g. different kinds of research may need different assignment strategies), and some open (e.g. how do we evaluate the effect of any changes to peer review process?)

A big issue for both authors and reviewers is mismatches due to lack of interest, which is in tension with explorative aspects of peer review. We recommend to address this issue with a combination of assignments based on bidding and random matches within area, backed up by reviewer training.

⁵Position bias is well documented in search & recommendation systems (Craswell et al., 2008; Collins et al., 2018).

⁶Or the program chairs, should the conference aim to have consistent policies for all ACs.

Impact Statement

Broader impact. The study identifies types of information that could be used to provide better paper-reviewer matches. Used strategically by a conference, it could be a powerful tool for diversifying the field, by helping the non-mainstream papers find the reviewers more open to them. By the same token, if the entity organizing the review process aimed for suppressing such research, deprioritising this information could harm such papers. Our proposal of interpretable paper-reviewer assignments would mitigate this potential risk by requiring the organizers to disclose their rationale for any given match.

Personal data. The surveys are designed to not solicit any personally identifiable information (including comments about individual peer review cases in the past conferences), or demographic information about participants.

Potential risks. The respondents are participants in anonymous peer review process, and as such being tracked back to individual peer review cases could expose them to retaliation. The survey therefore did not solicit information about specific venues (only broader categories such as “*ACL conferences”), and we manually verified that the open-ended comments also do not contain references to specific cases. We thus foresee no potential risks from deanonymization of the respondents.

Informed consent. The respondents are informed about the organizers and the objective of the study: to identify current practises of paper-reviewer assignments in CL/NLP conferences and ways in which this process can be improved. Responses are anonymous and respondents consent to the use and sharing of their responses for research purposes. Respondents must give consent to continue the survey.

Intended use. The survey data and forms will be made publicly available for research purposes.

Institutional approval. The study was approved by the Research Ethics Committee at the authors’ institution.

References

Balazs Aczel, Barnabas Szasz, and Alex O. Holcombe. 2021. [A billion-dollar donation: Estimating the cost](#)

[of researchers’ time spent on peer review](#). *Research Integrity and Peer Review*, 6(1):14.

Emily M. Bender. 2019. [The #BenderRule: On Naming the Languages We Study and Why It Matters](#).

Emily M. Bender and Leon Derczynski. 2018. [Paper Types](#).

Alina Beygelzimer, Yann Dauphin, Percy Liang, and Jennifer Wortman Vaughan. 2021. [The NeurIPS 2021 Consistency Experiment](#).

Niclas Boehmer, Robert Bredereck, and André Nichterlein. 2022. [Combating Collusion Rings is Hard but Possible](#).

Claire Cardie, Iryna Gurevych, and Yusuke Miyao. [ACL 2018: A Report on the Review Process of ACL 2018](#).

Johan S. G. Chu and James A. Evans. 2021. [Slowed canonical progress in large fields of science](#). *Proceedings of the National Academy of Sciences*, 118(41).

Andrew Collins, Dominika Tkaczyk, Akiko Aizawa, and Joeran Beel. 2018. [A Study of Position Bias in Digital Library Recommender Systems](#). *arXiv:1802.06565 [cs]*.

Corinna Cortes and Neil D. Lawrence. 2021. [Inconsistency in Conference Peer Review: Revisiting the 2014 NeurIPS Experiment](#). *arXiv:2109.09774 [cs]*.

Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. 2008. [An experimental comparison of click position-bias models](#). In *Proceedings of the 2008 International Conference on Web Search and Data Mining, WSDM ’08*, pages 87–94, New York, NY, USA. Association for Computing Machinery.

Jane Dai, Jeremy Cone, and Jeff Moher. 2020. [Perceptual salience influences food choices independently of health and taste preferences](#). *Cognitive research: principles and implications*, 5(1):2–13.

Hal Daumé III. 2015. [Some NAACL 2013 statistics on author response, review quality, etc](#).

Mor Geva, Yoav Goldberg, and Jonathan Berant. 2019. [Are We Modeling the Task or the Annotator? An Investigation of Annotator Bias in Natural Language Understanding Datasets](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1161–1166, Hong Kong, China. Association for Computational Linguistics.

Richard Gordon and Bryan J. Poulin. 2009. [Cost of the NSERC Science Grant Peer Review System Exceeds the Cost of Giving Every Qualified Researcher a Baseline Grant](#). *Accountability in Research*, 16(1):13–40.

David Hartvigsen, Jerry C. Wei, and Richard Czuchlewski. 1999. [The Conference Paper-Reviewer Assignment Problem*](#). *Decision Sciences*, 30(3):865–876.

639	Laura Hyman, Julie Lamb, and Martin Bulmer. 2006.	Fan Wang, Ning Shi, and Ben Chen. 2010. A comprehensive survey of the reviewer assignment problem.	691
640	The use of pre-existing survey questions: Implications for data quality. In <i>Proceedings of the European Conference on Quality in Survey Statistics</i> , pages 1–8. Cardiff Wales, UK.	<i>International Journal of Information Technology & Decision Making</i> , 09(04):645–668.	692
641			693
642			694
643			
644	Xinlian Li and Toyohide Watanabe. 2013. Automatic Paper-to-reviewer Assignment, based on the Matching Degree of the Reviewers. <i>Procedia Computer Science</i> , 22:633–642.	John Wieting, Kevin Gimpel, Graham Neubig, and Taylor Berg-Kirkpatrick. 2019. Simple and effective paraphrastic similarity from parallel translations. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4602–4608, Florence, Italy. Association for Computational Linguistics.	695
645			696
646			697
647			698
648	Michael L. Littman. 2021. Collusion rings threaten the integrity of computer science research. <i>Communications of the ACM</i> , 64(6):43–44.		699
649			700
650		Fei Xia, Wenjie Li, and Roberto Navigli. 2021. ACL-IJCNLP 2021 Program Chair Report.	701
651	Mingnan Liu and Laura Wronski. 2018. Trap questions in online surveys: Results from three web survey experiments. <i>International Journal of Market Research</i> , 60(1):32–49.		702
652			703
653		H. Öztas Ayhan and Semih Isiksal. 2005. Memory recall errors in retrospective surveys: A reverse record check study. <i>Quality and Quantity</i> , 38(5):475–493. Copyright - Kluwer Academic Publishers 2004; Last updated - 2021-09-09.	704
654			705
655	Cheng Long, Raymond Chi-Wing Wong, Yu Peng, and Liangliang Ye. 2013. On Good and Fair Paper-Reviewer Assignment. In <i>2013 IEEE 13th International Conference on Data Mining</i> , pages 1145–1150.		706
656			707
657			708
658			
659			
660	Graham Neubig. 2020. ACL Rolling Review Proposal.		
661	Graham Neubig, John Wieting, Arya McCarthy, Amanda Stent, Natalie Schluter, and Trevor Cohn. 2021. ACL Reviewer Matching Code. Association for Computational Linguistics.		
662			
663			
664			
665	OpenReview. 2021. Paper-reviewer affinity modeling for OpenReview.		
666			
667	Eric Price. 2014. The NIPS experiment.		
668	Anna Rogers and Isabelle Augenstein. 2020. What Can We Do to Improve Peer Review in NLP? In <i>Findings of EMNLP</i> , pages 1256–1262, Online. Association for Computational Linguistics.		
669			
670			
671			
672	Anders Søgaard. 2021. Explainable Natural Language Processing. <i>Synthesis Lectures on Human Language Technologies</i> , 14(3):1–123.		
673			
674			
675	Ivan Stelmakh, Nihar B. Shah, and Aarti Singh. 2019. PeerReview4All: Fair and Accurate Reviewer Assignment in Peer Review. In <i>Proceedings of the 30th International Conference on Algorithmic Learning Theory</i> , pages 828–856.		
676			
677			
678			
679			
680	Seymour Sudman and Norman M. Bradburn. 1973. Effects of time and memory factors on response in surveys. <i>Journal of the American Statistical Association</i> , 68(344):805–815.		
681			
682			
683			
684	A Tversky and D Kahneman. 1981. The framing of decisions and the psychology of choice. <i>Science (American Association for the Advancement of Science)</i> , 211(4481):453–458.		
685			
686			
687			
688	Amos Tversky and Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases. <i>Science</i> , 185(4157):1124–1131.		
689			
690			

709
710
711

A Appendix

In this appendix we introduce supplementary figures and tables.

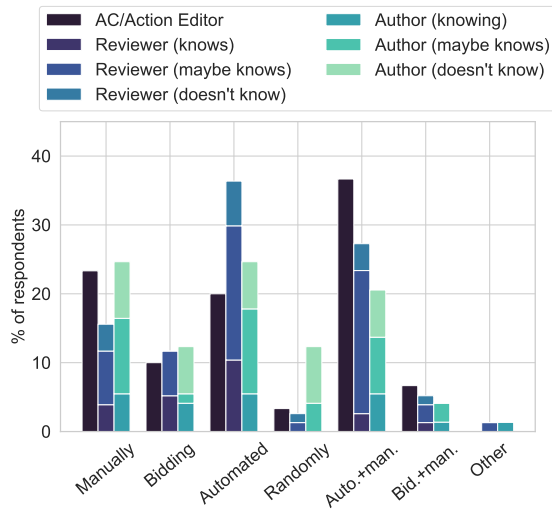


Figure 8: We ask reviewers and authors whether they know for certain, or maybe knows, or do not know how the paper-reviewer matching was done for their last CL/NLP venue. We then ask both reviewers, authors and ACs what they believe (or knows in the case of some) was the process for this venue. The guesses, and knowledge herof, are much different from *best* options in fig. 4, discussed in section 5.

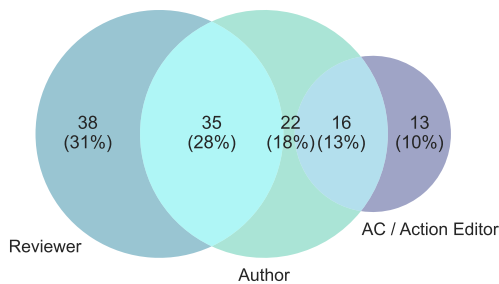
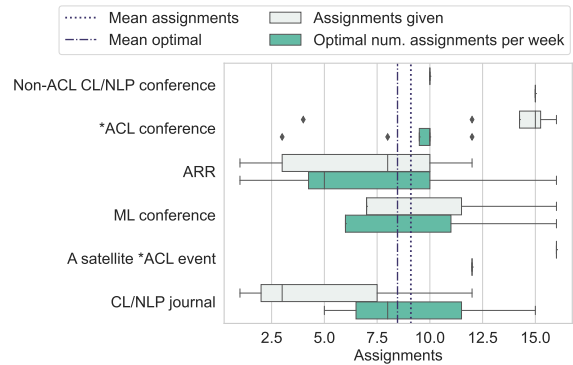
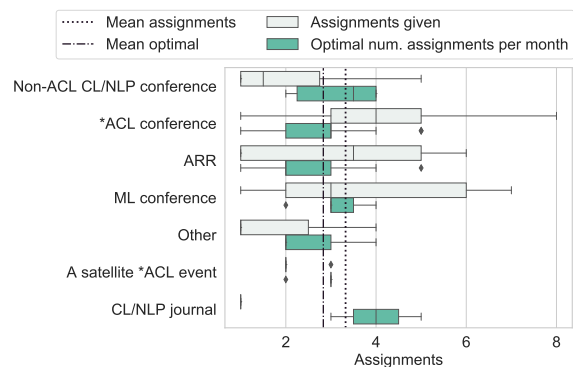


Figure 9: Distribution of **non-students** in the three roles, with overlap derived from asking the question *Did you also serve as a reviewer/author?* for their last CL/NLP venue.



(a) AC/Action Editors



(b) Reviewers

Figure 10: The boxplots shows the number of assignments given and optimal for a) ACs and b) Reviewers, discussed in section 5.3. Number of given assignments are reported for the venue in which the respondent last served as AC/reviewer, and optimal number of assignments are reported for time periods one week for ACs and one month for reviewers. Mean given and optimal number of assignments, across all respondents/venues, are shown with vertical striped lines.

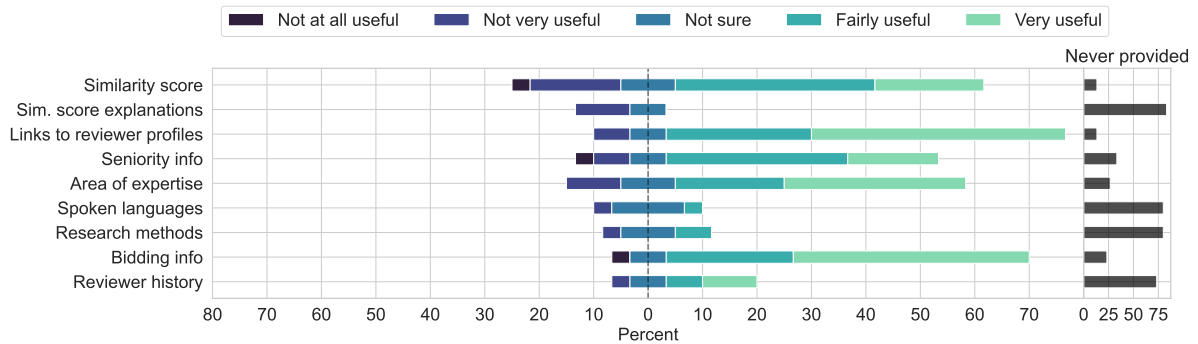
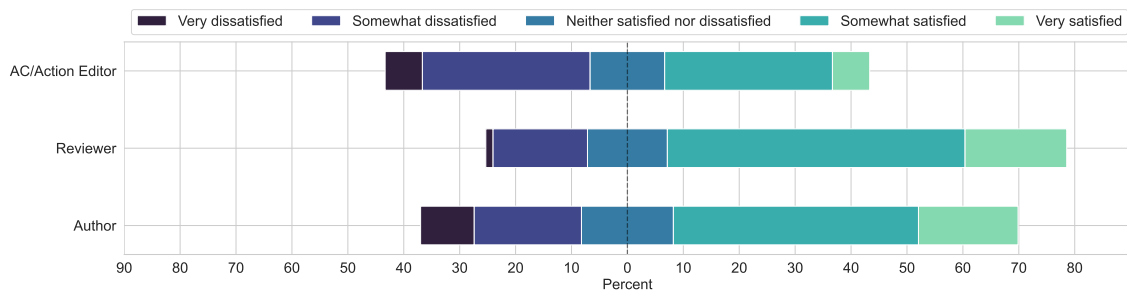
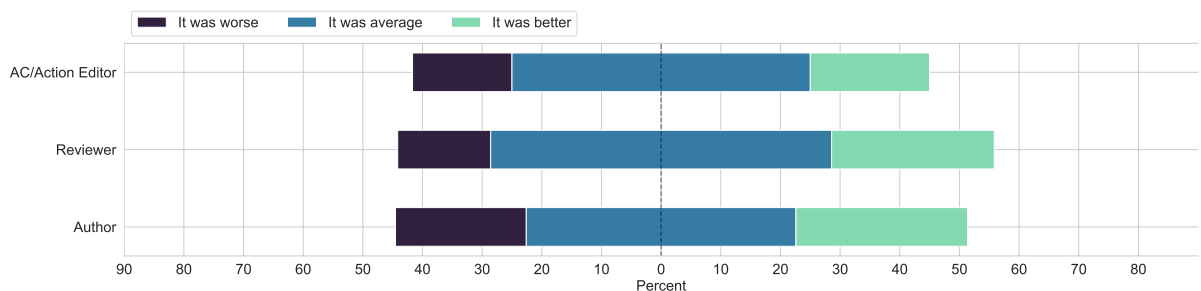


Figure 11: This is a larger (perhaps more easily readable) version of fig. 6. The diverging bars shows the experienced utility of different kinds of information about reviewers that ACs may have been presented with to assist in manual checks of paper-reviewer matches. If the respondent had never been presented with the specific kind of information they chose “Never provided”.



(a) Reported satisfaction with most recent experience. AC / Action Editor question: (...) how satisfied were you with the support provided to you to improve the paper-reviewer matching? Reviewer question: (...) How satisfied were you with the paper-reviewer matching? Author question: (...) How satisfied were you with the amount of constructive criticism in the reviews you received?



(b) Question: would you say your most recent experience with paper-reviewer matching/paper assignment(s)/set of reviews described above, was better or worse than on average?

Figure 12: We ask all respondents general questions about their satisfaction with their last experience as an AC/reviewer/author and their overall satisfaction in this role. We discuss these results in the limitations section 6.2.

	AREA CHAIR			REVIEWER			AUTHOR		
	>10	6-10	1-5	>10	6-10	1-5	>10	6-10	1-5
Bachelor's student	0	0	0	0	0	0	0	1	0
Master's student	0	0	0	0	0	1	0	0	0
PhD student	0	1	0	3	10	7	2	6	25
Other	0	0	0	1	0	0	1	0	2
Postdoc	2	1	0	12	3	1	11	0	2
Academic researcher	19	0	0	21	1	4	13	0	4
Industry researcher	6	0	1	8	0	5	3	2	1
Total	27	3	1	45	14	18	30	9	34

Table 1: This table shows the count of respondents from each role (AC/reviewer/author) reporting one of 7 career statuses and an amount of times having had their own papers reviewed. The numbers reflects those plotted in [fig. 2](#), [section 4.1](#).

	Tasks	Method	Data	Type of paper	Affinity score	Bidding	Trust
AC / Action Editor	3.90±0.83	3.70±0.74	2.83±1.00	3.37±0.95	2.13±1.50	3.67±1.25	3.27±1.67
Reviewer	4.00±0.90	3.86±0.85	2.75±1.09	3.16±0.97	1.97±1.38	4.03±1.07	2.96±1.57
Author	3.95±0.86	3.59±0.87	2.84±0.98	3.45±1.18	1.85±1.32	3.86±1.22	3.11±1.60
Grand mean	3.95	3.72	2.81	3.32	1.99	3.85	3.11

Table 2: Mean importance with 0=Not sure, 1=Not important, 2=Slightly important, 3=Moderately important, 4=Very important and 5=Extremely important, for the seven paper-reviewer matching factors shown in [fig. 7](#). Removing "Not sure" does not change the overall ranking. The grand mean is the unweighted mean of ACs', reviewers' and authors' mean scores. The mean absolute difference is greatest between ACs and reviewers (0.20) and smallest between ACs and authors (0.12), while between reviewers and authors it is 0.16. These results are discussed in [section 5.2](#) under "Stakeholder preferences".

	Tasks	Methods	Data	Paper type	Sim. score	Bidding	Trust
Computationally-aided linguistic analysis	0.763	0.242	0.024 ^{<}	0.097	1.000	0.487	0.707
NLP engineering experiment paper	1.000	0.072	1.000	0.227	1.000	0.732	0.427
Reproduction paper	0.457	0.036 ^{<}	0.008 ^{>}	1.000	0.587	1.000	1.000
Resource paper	0.728	0.757	1.000	0.164	0.808	0.610	0.658
Position paper	0.227	0.264	0.657	0.227	0.001 ^{<}	0.767	0.251
Survey paper	0.751	1.000	0.842	0.426	0.656	0.135	0.847
Other	1.000	1.000	0.205	0.020 ^{<*}	0.063	0.067	0.740

Table 3: P-values of two-sided Fisher Exact tests, discussed in [section 5.2](#). For each contribution type, we test the null hypothesis that there is no difference in whether respondents find a paper-match factor (from [fig. 7](#)) more than or less than moderately important, depending on whether or not individuals have worked on the specific types of papers (contribution types). For each combination of contribution type i and paper-match factor j , a 2×2 contingency table is made with the counts of a) respondents having worked with type i and finding factor j *less* than moderately important, b) having worked with type i and finding factor j *more* than moderately important, c) having *not* worked with i and finding j *less* than moderately important, d) having *not* worked with i and finding j *more* than moderately important. The p-values reflect the probability of observing the given counts or something more imbalanced between types i . Significant p-values, $p < 0.05$, are in bold, and for these, superscript $>$ denotes that respondents having worked with i believe factor j is *more* than moderately important, and the superscript $<$ denotes the opposite.

*Unfortunately, we did not ask the respondents what "other" paper types they worked on, so this result is difficult to interpret.