

## DIVERSITY

# Can rubrics combat gender bias in faculty hiring?

Some bias persisted, but rubric use should be encouraged

By **Mary Blair-Loy<sup>1</sup>, Olga V. Mayorova<sup>1</sup>,  
Pamela C. Cosman<sup>2</sup>, Stephanie I. Fraley<sup>3</sup>**

**R**esearch has documented the presence of bias against women in hiring, including in academic science, technology, engineering, and mathematics (STEM). Hiring rubrics (also called criterion checklists, decision support tools, and evaluation tools) are widely recommended as a precise, cost-effective remedy to counteract hiring bias, despite a paucity of evidence that they actually work (see table S8). Our in-depth case study of rubric usage in faculty hiring in an academic engineering department in a very research-active university found that the rate of hiring women increased after the department deployed rubrics and used them to guide holistic discussions. Yet we also found evidence of substantial gender bias persisting in some rubric scoring categories and evaluators' written comments. We do not recommend abandoning rubrics. Instead, we recommend a strategic and sociologically astute use of rubrics as a department self-study tool within the context of a holistic evaluation of semifinalist candidates.

Although academic STEM aspires to be a meritocracy, its taken-for-granted cultural schemas of merit smuggle in biases (1), which contribute to a dearth of diversity that undermines scientific innovation and impact (see table S8). In academic engineering, one of the most male-dominated STEM fields, on average 17.6% of engineering faculty positions are held by women (2). Although the percent of women engineering doctorates increased from 15.8% in 2000 to 24% in 2019 (3), these increases will not be matched by gains in the professoriate if women face unfair barriers at hiring.

Academic policy-makers and EDI (equity, diversity, and inclusion) specialists

strongly encourage rubric usage, in which faculty evaluators systematically rate each candidate on a set of previously agreed-on criteria. This process is believed to counteract the bias of individual evaluators by promoting slower, more deliberative, and analytical thinking and by focusing them on skill sets that directly affect job performance rather than on impressions and intuitions (4, 5). However, we are aware of only one study, conducted in a laboratory setting, in which participants rated candidate summaries, which shows that agreeing on rubric criteria in advance reduces evaluation bias (6). We are aware of no studies that analyze the effect of rubric use on bias in real-world hiring, in which actual evaluators assess voluminous candidate files to make actual high-stakes decisions. Real-world case studies are important because the effectiveness of interventions depends on the social context and the identities of all involved (7).

Despite this paucity of evidence, many fields have developed rubrics to standardize candidate assessment and have promoted rubrics as a best practice for EDI in hiring. In policy guides for academic hiring, several applied treatises and websites provide sample rubrics (see table S8) (5). A recent influential review of faculty hiring lists "mandatory use" of rubrics as one of the interventions that the authors "view as having the most promise for [university] institutions seeking to improve inclusivity in hiring across disciplines" (4).

## CASE STUDY OF RUBRIC USAGE

To help address this knowledge and policy gap, we developed an in-depth case study of an engineering department in a research-intensive (Carnegie classification R1), highly ranked university. Like other R1s, this department strongly values research productivity when evaluating faculty candidates (8). Like most academic engineering departments, our case department was male dominated; women composed 18% of the faculty, which is close to the 17.6% national average (2).

## Rubric usage and hiring patterns

We started with a faculty candidate evaluation template from the University of Michigan STRIDE program (Strategies and Tactics for Recruiting to Improve Diversity and Excellence), which is funded by the National Science Foundation (NSF). This template includes widely accepted criteria for faculty at research-intensive universities (5).

We worked with the department under study to adapt the template to fit its searches. Its rubric evaluated faculty candidates across six dimensions: research productivity, research impact, teaching ability, contributions to diversity, potential for collaboration, and overall impression (see fig. S1). Rubric scores ranged from excellent to poor for each evaluation category (we translated these ratings into a numerical variable: excellent = 4, good = 3, neutral = 2, fair = 1, and poor = 0). Written commentary was also encouraged.

Department faculty agreed to fill out the rubric as a tool in their evaluation of the semifinalist list compiled by the recruitment committee. In four faculty search cycles over four recent academic years, faculty used the rubric to evaluate written materials supporting applications of 62 semifinalists (32 women and 30 men; gender was self-reported by candidates in their application file). At the beginning of each faculty meeting that was focused on selecting finalists, a faculty member summarized and presented rubric scoring results and commentary, with evaluators anonymized, filtering out any inaccurate or off-topic content.

Our analysis next compared the proportion of women hired during the 8-year period immediately before rubric use (which we refer to as "Phase One") with the proportion hired during the four academic years of rubric use ("Phase Two"). Near the outset of Phase One, the campus implemented three EDI interventions: Faculty serving as equity advisers took on administrative oversight of shortlisted candidates, applicants' "Contributions to Diversity Statements" (C2D) were added to application files, and equity advisers began giving diversity training to search committee members. Such training focused on evaluating C2D statements and an overview of research on implicit biases. Nonetheless, during Phase One, the department conducted eight searches and hired eight men and one woman.

At the outset of Phase Two, the department introduced an additional intervention: rubrics used as described above. During Phase Two, the department conducted four faculty searches and hired three women and six men. The number of women hired increased from only one per nine hires in Phase One to three per nine hires in Phase

<sup>1</sup>Department of Sociology, University of California, San Diego, La Jolla, CA, USA. <sup>2</sup>Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA, USA. <sup>3</sup>Department of Bioengineering, University of California, San Diego, La Jolla, CA, USA. Email: mblairloy@ucsd.edu; sifraley@ucsd.edu

Two. The phase with the increased hiring of women coincides with the period when rubrics were used. (We cannot fully attribute this increase to rubric utilization because unmeasured factors may also have changed.)

### Analysis of rubric scoring

All 62 semifinalist candidates received rubric scores from 6 to 21 faculty, with a mean of 13.5 and a median of 12. There was no statistically significant gender difference in the number of scores received.

Analysis of the rubric scoring patterns for men and women candidates revealed statistically significant differences in three of the six evaluation categories (see table S1.A): Women were scored lower than men in research productivity and research impact but higher than men in contributions to diversity. In the other categories, including the overall impression category, scores for men and women were not significantly different.

To determine whether gender bias was incorporated into rubric scores, we analyzed the research productivity category because it can be most directly compared with external metrics. We chose two metrics calculated for the candidate's application year. First, we tallied from candidate curricula vitae the number of articles published (and confirmed this tally in the Web of Science database). Second, we pulled from the Web of Science the H-index, a dominant measure of researcher output that incorporates productivity and impact in a single number that can be compared across faculty of all seniority levels (9). We call research productivity, a rubric category that can be measured independently, a "calibration category" (see table S2, footnote).

To test for gender bias, we constructed ordinary least squares (OLS) regression models to predict rubric scores of research productivity, controlling for the independently measured categories of seniority and number of articles or H-index (see table S2). We found that women candidates, on average, received statistically significantly lower productivity rubric scores than those of men, even after controlling for seniority (measured as number of years since PhD) and number of articles published [unstandardized  $\beta$  coefficient ( $B$ ) =  $-0.36$ ,  $P \leq 0.01$ ] (see table S2, model 1). Similarly, women receive significantly lower scores on average than men while controlling for seniority and H-index ( $B = -0.29$ ,  $P < 0.05$ ) (table S2, model 3). Thus, rubric scoring alone did not appear to fully mitigate gender bias. These findings mirror the gender bias detected in other academic peer-review processes (10, 11). Because the H-index itself has been found to incorporate bias against women (see table S8), our findings should be interpreted as additional bias.

Social psychology literature on double standards finds that among more junior candidates, women are often held to higher standards of competence than men in ways that sometimes change for candidates with more experience (12). We thus tested whether the effect of gender on rubric scores is contingent on the value of the external metric (see table S2, models 2 and 4).

Women's productivity rubric scores are consistently below those of men who have the same number of articles and seniority (see the first figure, left). Women face an average 0.36-point penalty, which remains the

same across the range of number of articles published. (The rubric scores are mostly clustered at the middle to high portion of the scale, between 2 and 4; within this range, there is an approximately 18% penalty for being a woman.)

Yet when controlling for H-index and seniority, the gender penalty is harshest among candidates with the lowest H-indices (see the first figure, right), who are disproportionately junior (see table S4, footnote). At the lowest tail of the H-index distribution, men receive research productivity rubric scores that are on average 0.7 points higher than the scores of women with the same seniority. At this end of the H-index distribution, the rubric scores are mostly clustered between 1 and 3; for these candidates, there is an approximately 35% penalty for being a woman.

The gender difference in rubric scoring gradually decreases by about 0.01 point for each 1 point of H-index gained. Yet women do not catch up to men in how productive they are rated in the rubrics until reaching an H-index of 17.5, a productivity index well above the 12.8 average and achieved by only a handful of candidates. Because rubric scores for men and women in the overall impression category were not statistically significantly different (see table S1), evaluators may have combined category scores so that women's higher average scores on contributions to diversity offset their lower average scores on productivity and impact.

### Content analysis of qualitative comments

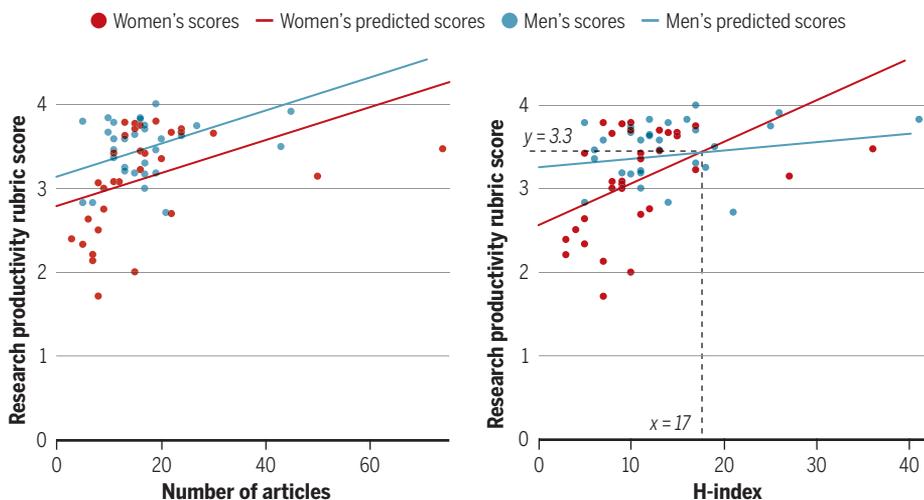
Next, we conducted what to our knowledge is the first content analysis of qualitative rubric comments in a real faculty search context. Candidates received written commentary alongside their rubric scores. An average of three and a maximum of nine evaluators wrote comments on each candidate. The number of comments received did not differ by candidate gender.

In our content analysis, we prepared a dataset of comments, in which candidate gender was concealed, by removing gender indicators such as pronouns. We then combined an inductive exploration for emergent themes with deductive searches for specific patterns found in previous literature on letters of recommendation. We conducted hand coding, which some research suggests is superior to computer-assisted coding for studies such as ours with new analyses (13).

Many comments contained evaluative notes on the quality, number, authorship order, or impact of the publications. Inductively, we coded these as either negative (for example, "some gaps in the pubs." and "only one paper from a postdoc of three-plus years") or positive (for example, "strong publication record, letters attest to research

## Indication of gender bias in rubric scores on research productivity

Graphs show predicted values by gender from ordinary least squares regression models regressing rubric scores of research productivity on independent productivity metrics, controlling for seniority.  $n = 62$  semifinalist candidates (see table S2, models 1 and 4).



output and collaboration”). We also searched deductively for the presence of two themes mentioned in previous literature on other types of evaluation. One code is “standout” language (13, 14). Examples in the rubric comments include “outstanding productivity and quality” and “probably the best [search specialty] candidate out there this cycle.” The last code is “doubt raisers” (15), when a seemingly positive or neutral comment is accompanied by language that minimizes the accomplishment or raises concerns (for example, “several publications, but... some impact factors are very modest”). We turned these four codes into four dichotomous measures of the presence or absence of comments in each category. (See table S3 for details on coding methods, inter-coder reliability, dichotomous measure rationale, and robustness check.)

Candidate gender was subsequently unmasked, and the percentages of women and men who had received at least one positive comment, negative comment, standout term, and doubt raiser were calculated (see table S3). We found that 86% of men but only 63% of women candidates received at least one positive comment. Men were half as likely to receive a negative comment (25%) compared with women (50%). Men were 3.5 times more likely to receive standout language (32%) compared with women (9%). A  $\chi^2$  test indicates statistically significant gender differences for these three variables. Thirteen percent of women and 25% of men received a doubt-raising comment. This pattern was not expected, yet this gender difference is not statistically significant (see table S3). Overall, the gendered patterns in rubric quantitative scores align with the qualitative comments and with previous research on other evaluative language (13–15).

### Survey of faculty

At the conclusion of Phase Two, we conducted an anonymous survey of department faculty, to which 56% of the professors responded. Our check for selection bias in responses revealed no statistically significant difference in the two indicators we have for respondents and nonrespondents: gender and tenure status (see table S6). Most reported that rubric usage helped them evaluate candidates in a more organized fashion (78%) and may have helped them be more objective (78%), potentially reducing individual bias (see fig. S2 and table S5).

## Using rubrics for EDI hiring

Recommended process for rubric use to improve EDI in faculty hiring. Steps 1, 3, and 6 are conducted by faculty evaluators while steps 2, 4, and 5 are conducted by faculty on the search committee with administrative support.

### Evaluators

Work as a group to determine rubric categories and weights, in accordance with disciplinary schemes of merit, to mitigate evaluation bias.

Individually complete rubrics for all candidates to mitigate individual bias.

Use findings to guide holistic discussion of candidates, further reducing individual bias and mitigating interactional bias.

### Process managers

Ensure that rubric includes a calibration category, i.e., one that is independently quantifiable.

Compile and analyze rubric results, comparing calibration category scores to independent metrics to detect any remaining individual bias.

Start the meeting of evaluators by presenting rubric results to mitigate first speaker bias, noting high inter-evaluator variance and any large inconsistencies of scores with independent metrics.

At the beginning of the meeting to discuss the semifinalist list and choose finalists to invite for interviews, a faculty member presented rubric numeric scores and comments with the identity of the evaluators anonymized. It was explained to faculty that this step aimed to enable all viewpoints to be heard while avoiding first-speaker bias, in which initial speakers set the tone for discussion. However, faculty were not discouraged from making additional comments or claiming or echoing their support for particular comments after the presentation. Most survey respondents said that this practice prompted meeting attendees to focus on more objective criteria (78%), improved the climate of the meeting (80%), and reduced the first speaker effect (67%) (see fig. S2).

Additionally, some commented in open-ended survey responses that meeting time was used efficiently to quickly identify candidates with strong consensus and then spend more time on those with high variance in ratings [whom we determined were disproportionately women (see table S1.B)]. Taken together, these results suggest that beginning the faculty meeting with rubric results reduced interactional bias emergent in the faculty meeting. By opening with a neutral reading of the full set of both positive and negative rubric comments, the impact was blunted of any first speakers, often senior men, attempting to vociferously promote or shoot down a candidate. The faculty meeting format may have mitigated gender bias in the research productivity scores and the selection of finalists; 47% of women semifinalists and 37% of men semifinalists advanced to the finalist stage.

## POLICY TEMPLATE

In light of our findings that gender bias remains endemic even in this seemingly objective evaluation process, it is vital that rubric usage be accompanied by strategic application in departmental meetings to counteract individual bias and check interactional bias during the discussion of candidates. Our results suggest that using rubrics according to this process framework can improve diversity in hiring. Thus, we recommend a strategic and sociologically astute use of rubrics as a department self-study tool within the context of a holistic evaluation of the short-listed candidates (see the second figure).

We have studied this process with regard to gender. Given the otherwise limited diversity among candidates in our study, we were unable to address whether rubrics could also be a tool to promote and check on the fairness of evaluations with regard to race/ethnicity or other minoritized identities. This suggests priorities for future research. ■

## REFERENCES AND NOTES

1. M. Blair-Loy, E. Cech, *Misconceiving Merit: Paradoxes of Excellence and Devotion in Academic Science and Engineering* (Univ. Chicago Press, 2022).
2. J. Roy, C. Wilson, A. Erdiaw-Kwasie, C. Stuppard, “Engineering and engineering technology by the numbers 2019” (American Society for Engineering Education, 2020).
3. J. Falkenheim, “Doctoral Recipients from U.S. Universities: 2019” (National Science Foundation, 2019).
4. K. O’Meara, D. Culpepper, L. L. Templeton, *Rev. Educ. Res.* **90**, 311 (2020).
5. University of Michigan, in *Advance Program*, University of Michigan, Ed. (University of Michigan Office of the Provost, 2018).
6. E. Uhlmann, G. L. Cohen, *Psychol. Sci.* **16**, 474 (2005).
7. R. H. Thaler, K. R. Sunstein, *Nudge: Improving Decisions About Health, Wealth, and Happiness* (Yale Univ. Press, 2008).
8. National Research Council, *Gender Differences at Critical Transitions in the Careers of Science, Engineering, and Mathematics Faculty* (National Academies Press, 2010).
9. V. Koltun, D. Hafner, *PLOS ONE* **16**, e0253397 (2021).
10. C. Wenneras, A. Wold, *Nature* **387**, 341 (1997).
11. E. R. Andersson, C. E. Hagberg, S. Hägg, *Front. Res. Metr. Anal.* **6**, 594424 (2021).
12. M. Foschi, *Annu. Rev. Sociol.* **26**, 21 (2000).
13. S. J. Correll, K. R. Weisshaar, A. T. Wynn, *Am. Sociol. Rev.* **85**, 1022 (2020).
14. T. Schmader, J. Whitehead, V. H. Wysocki, *Sex Roles* **57**, 509 (2007).
15. J. M. Madera, M. R. Hebl, H. Dial, R. Martin, V. Valian, *J. Bus. Psychol.* **34**, 287 (2019).

## ACKNOWLEDGMENTS

This work was supported by the National Science Foundation, grant 1661306 (M.B.-L.). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation. S.I.F. is a cofounder and scientific adviser for MelloLabs and holds equity in the company.

## SUPPLEMENTARY MATERIALS

science.org/doi/10.1126/science.abm2329

10.1126/science.abm2329

## Can rubrics combat gender bias in faculty hiring?

Mary Blair-Loy Olga V. Mayorova Pamela C. Cosman Stephanie I. Fraley

*Science*, 377 (6601), • DOI: 10.1126/science.abm2329

### View the article online

<https://www.science.org/doi/10.1126/science.abm2329>

### Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)