

1 **Original Manuscript**

2

3

4 **How to best evaluate applications for junior fellowships?**

5 **Remote evaluation and face-to-face panel meetings**

6 **compared**

7

8 Marco Bieri^{1¶*}, Katharina Roser^{1,2¶}, Rachel Heyard¹, Matthias Egger^{1,3,4}

9

10 ¹Swiss National Science Foundation, Berne, Switzerland

11 ²Department of Health Sciences and Medicine, University of Lucerne, Lucerne, Switzerland

12 ³Institute of Social and Preventive Medicine, University of Berne, Berne, Switzerland

13 ⁴Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK

14 ¶These authors contributed equally to this work

15

16

17

18

19 *Corresponding author:

20 Marco Bieri

21 Swiss National Science Foundation

22 Wildhainweg 3

23 CH-3001 Berne

24 marco.bieri@snf.ch

25

26 Word count: abstract 291 words, box on strengths and limitations 119 words, main text 2835

27 words, 29 references, 4 tables, 1 figure.

28

29

30

31 **Abstract**

32 **Objectives**

33 To test a simplified evaluation of fellowship proposals by analyzing the agreement of funding
34 decisions with the official evaluation, and to examine the use of a lottery-based decision for
35 proposals of similar quality.

36

37 **Design**

38 The study involved 134 junior fellowship proposals (Postdoc.Mobility). The official method
39 used two panel reviewers who independently scored the application, followed by triage and
40 discussion of selected applications in a panel. Very competitive/uncompetitive proposals were
41 directly funded/rejected without discussion. The simplified procedure used the scores of the
42 two panel members, with or without the score of an additional, third expert. Both methods
43 could further use a lottery to decide on applications of similar quality close to the funding
44 threshold. The same funding rate was applied, and the agreement between the two methods
45 analyzed.

46

47 **Setting**

48 Swiss National Science Foundation (SNSF).

49

50 **Participants**

51 Postdoc.Mobility panel reviewers and additional expert reviewers.

52

53 **Primary outcome measure**

54 Per cent agreement between the simplified and official evaluation method with 95%
55 confidence intervals (95% CI).

56

57 **Results**

58 The simplified procedure based on three reviews agreed in 80.6% (95% CI 73.9-87.3) with
59 the official funding outcome. The agreement was 86.6% (95% CI 80.8-92.4) when using the
60 two reviews of the panel members. The agreement between the two methods was lower for
61 the group of applications discussed in the panel (64.2% and 73.1%, respectively), and higher

62 for directly funded/rejected applications (range 96.7% to 100%). The lottery was used in eight
63 (6.0%) of 134 applications (official method), 19 (14.2%) applications (simplified, three
64 reviewers) and 23 (17.2%) applications (simplified, two reviewers). With the simplified
65 procedure, evaluation costs could have been halved and 31 hours of meeting time saved for
66 the two 2019 calls.

67

68 **Conclusion**

69 Agreement between the two methods was high. The simplified procedure could represent a
70 viable evaluation method for the Postdoc.Mobility early career instrument at the SNSF.

71

72

73 **Strengths and limitations of this study**

- 74 ■ The study examined the outcome between a simplified and the official evaluation
75 procedure for junior fellowship applications for different research disciplines.
- 76 ■ The study discussed the agreement between the two evaluation methods in the context of
77 the general uncertainty around peer review and estimated the costs that could have been
78 saved with the simplified evaluation procedure.
- 79 ■ It is the first study to provide insight into lottery-based decisions in the context of the
80 evaluation of junior fellowship applications.
- 81 ■ The study lacks statistical power because the sample size of applications was relatively
82 small.
- 83 ■ The study addressed the specific context and evaluation of the SNSF Postdoc.Mobility
84 funding scheme, results may thus not be generalizable to other funding programs.

85

86

87 **Introduction**

88 Peer review of grant proposals is costly and time-consuming. The burden on the scientific
89 system is increasing, affecting funders, reviewers, and applicants [1,2]. In response,
90 researchers have studied the review process and examined simplifications. For example, Snell
91 [3] studied the number of reviewers and consistency of decisions and found that five
92 evaluators represented an optimal tradeoff. Graves et al. [4] assessed the reliability of
93 decisions made by evaluation panels of different sizes. They concluded that reliability was
94 greatest with about ten panel members. Herbert et al. [5] compared smaller panels and shorter
95 research proposals with the standard review procedure. The agreement was about 75%
96 between simplified and standard procedures. As an alternative to face-to-face (FTF) panels,
97 the use of virtual, online meetings has also been examined. Bohannon [6] reported that at the
98 National Science Foundation (NSF) and National Institutes of Health (NIH), virtual meetings
99 could reduce costs by one-third. Gallo et al. [7] compared teleconferencing with FTF
100 meetings and found only few differences in the scoring of the applications. Later studies also
101 found that virtual and FTF panels produce comparable outcomes [8–10].

102 With virtual formats, panel members still need to attend time-consuming meetings.
103 Using the reviewers' written assessments without FTF or virtual panel discussions would
104 simplify the process further. Fogelholm et al. [11] reported that results were similar when
105 using panel consensus or the mean of reviewer scores. Obrecht et al. [12] noted that panel
106 review changed the funding outcome of only 11% of applications. Similarly, Carpenter et al.
107 [8] found that the impact of discussions was small, affecting the funding outcome of about
108 10% of applications. Pina et al. [13] studied Marie Curie Actions applications and concluded
109 that ranking applications based on reviewer scores might work for some but not all
110 disciplines. In the humanities, social and economic sciences, an exchange between reviewers
111 may be particularly relevant. The triaging of applications has also been examined: after an
112 initial screening, noncompetitive and very competitive proposals are either directly rejected or
113 funded. Vener et al. [14] validated the triage model of the NIH and found that the likelihood
114 of erroneously discarding a competitive proposal was very small. Bornmann et al.'s [15]
115 findings on a multi-stage fellowship selection process also supported the use of a triage.

116 Mandated by the government, the Swiss National Science Foundation (SNSF) is
117 Switzerland's foremost funding agency, supporting scientific research in all disciplines.
118 Following innovations in career funding, the SNSF will experience a significant increase of
119 applications for the junior "Postdoc.Mobility" fellowship scheme, which offers postdoctoral

120 researchers a stay at a research institution abroad for up to 24 months. The aim of this work
121 was to compare the evaluation of applications by expert review, triage, and discussion in an
122 evaluation panel with expert reviews only.

123

124 **Methods**

125 **Sample**

126 We included applications submitted for the August 2019 Postdoc.Mobility fellowship call.
127 We also included applications by Postdoc.Mobility fellows for a return grant to facilitate their
128 return to Switzerland. Both, fellowship and return grants were evaluated according to the
129 same criteria by the Humanities panel, the Social Sciences panel, Science, Technology,
130 Engineering, Mathematics (STEM) panel, the Biology or Medicine panels.

131

132 **Study design**

133 We compared funding outcomes based on the official, legally binding evaluation with a
134 simulated, hypothetical evaluation. The official evaluation was based on the triage of
135 applications based on expert reviews, followed by a discussion of the meritorious applications
136 in an FTF panel: the Triage-Panel Meeting (TPM) format ([Figure 1](#)). In a first step, each
137 proposal was independently reviewed and scored by two panel members. For the assessment,
138 the evaluation criteria defined in the Postdoc.Mobility regulations [16] were applied. The
139 criteria address different aspects of the applicant, the proposed research project, and the
140 designated research location. Panel members used a 6-point scale: outstanding=6 points,
141 excellent=5 points, very good=4 points, good=3 points, mediocre=2 points, poor=1 point.
142 Applications were then allocated to three groups based on the ranking of the mean scores
143 given to each proposal: Fund without further discussion (F in [Figure 1](#)), Discuss in panel
144 meeting (D), and Reject (R). Panel members could request that applications in the F or R
145 group are reallocated to D and discussed. In a second step, the D proposals were discussed in
146 the FTF panel meeting, ranked and funded or rejected. Random Selection (RS in [Figure 1](#))
147 could be used to fund or reject proposals of similar quality close to the funding threshold if
148 the panel could not reach a decision. Funding decisions were based on the standard two-stage
149 method, which included FTF panel meetings (TPM).

150 The simulated alternative procedure consisted only of the first step, i.e., was entirely
151 based on the ranking of proposals based on the expert reviews: the Expert Review-Based
152 (ERB) evaluation. In addition to the two panel members, a third expert reviewer who was not

153 a member of the panel assessed the proposal. The same 6-point scale was used. The proposals
154 were then allocated to one of three groups based on the mean scores (F, RS, and R in Figure
155 1). Random selection was used whenever the funding line went through a group of two or
156 more applications with identical scores. The funding rate of the TPM was applied to the
157 simulated ERB method.

158

159 **Data analysis**

160 To determine the agreement between the two evaluation methods, we used 2x2 contingency
161 tables. We calculated the simple agreement with 95% Wald confidence intervals (CI) for
162 proportions. We also examined the agreement between the TPM and the ERB approach using
163 only the assessments from the two panel members, thus excluding the assessment from the
164 third reviewer. We calculated discipline- and gender-specific levels of agreement and tested
165 for differences in agreement between disciplines and gender using chi-squared tests for
166 categorical data.

167

168 **Costs**

169 We determined the costs related to the evaluation. The costs comprised expenses related to
170 the scientific assessment of the individual applications and the panel meetings. The SNSF
171 compensates panel reviewers with USD 275 per scientific assessment. Panel reviewers further
172 receive a meeting allowance of up to USD 550 depending on the duration of the meeting.
173 Further, the SNSF reimburses travel expenses and accommodation costs. The five panels
174 included 96 members and met twice in 2019.

175

176 **Ethics approval**

177 The Ethics Committee of the Canton of Berne confirmed that the study does not fall under the
178 Federal Act on Research involving Human Beings. No reviewer, applicant or application can
179 be identified from this study.

180

181

182 **Results**

183 **Study sample and success rates**

184 The sample consisted of 134 applications, including 124 fellowship applications and ten
185 requests for a return grant. The mean age of applicants was 32.7 years (SD 3.2 years) in men
186 and 33.5 years (SD 2.8 years) in women. Each reviewer received a mean of 2.5 (SD 1.4)
187 applications to evaluate.

188 Table 1 shows the distribution of applications and success rates across disciplines,
189 genders and the three evaluation methods: the legally binding TPM format and the simulated
190 ERB evaluations with three or two reviewers. Most applications came from biology, followed
191 by the STEM disciplines and the social sciences. Almost two-thirds of applications came from
192 men. With TPM, success rates were slightly higher in women (60.4%) than in men (50.0%).
193 This was driven by the middle group of applications that were discussed in the panels, where
194 the success rates of women overall was 66.7% (24 of 67 applicants were women in this
195 group). Success rates were similar across disciplines, ranging from 56.2% in the humanities to
196 52.2% in the social sciences. By design, overall success rates were the same with the ERB
197 evaluations; however, the difference between genders was smaller with ERB than with TPM
198 (Table 1).

199

200 **Agreement between evaluation by ERB or TPM**

201 Comparing the ERB evaluation based on three reviewers with the standard TPM format, the
202 agreement overall was 80.6% (95% CI 73.9-87.3). The agreement was highest in the
203 Medicine panel (90.0%; CI 76.9-100), and lowest in the Social Sciences panel (73.9%; CI
204 56.0-91.8). However, the statistical evidence for differences in agreement between panels was
205 weak ($P=0.58$, Table 2). As expected, the agreement was higher when comparing the ERB
206 evaluation based on the two panel members with TPM. Overall, for two reviews, the
207 agreement was 86.6% (95% CI 80.8-92.4). It ranged from 75.0% (CI 53.8-96.2) in the
208 Humanities panel to 91.3% (CI 79.8-100) in the Social Sciences panel ($P=0.51$). Both for
209 ERB evaluation with three and two reviewers, the agreement was slightly higher for women
210 than for men ($P>0.70$, Table 3).

211 In Table 4, we calculated agreement separately for the triage categories: Fund (F),
212 Discuss (D), Reject (R). With the ERB evaluation based on three reviewers, agreements for F
213 and R were close to 100% (97.3% and 96.7%, respectively) but considerably lower for D:

214 64.2% (95% CI 52.7-75.7), with $P < 0.001$ for differences in agreement across categories from
215 chi-squared test. For ERB evaluation with two reviewers (the two panel members), the
216 agreement was 100% for F and R, but 73.1% (95% CI 62.5- 83.7) for D, with $P < 0.001$ for
217 differences in agreement.

218

219 **Random selection in TPM and ERB evaluation**

220 With the standard TPM evaluation, only eight (11.9%) of the 67 applicants in the D group, or
221 eight (6.0%) of 134 applicants were entered into a lottery of whom four were funded. With
222 the simulated ERB evaluation based on three reviewers, 19 (14.2%) of the 134 applicants
223 would have entered the lottery, and with the ERB with two reviewers 23 (17.2%) applications
224 would have been subjected to random selection.

225

226 **Cost savings**

227 We determined the resources that could be saved with the use of an ERB evaluation compared
228 to the TPM. By comparison with the current valid TPM evaluation procedure for the
229 Postdoc.Mobility, we calculated that about USD 91,000 related to the holding of meetings
230 could have been saved if an ERB evaluation had been used for the two Postdoc.Mobility calls
231 in 2019. This saving corresponds to 55% of total costs. Moreover, the holding of all panel
232 sessions in 2019 amounted to a total duration of 31 hours. This represents a significant
233 workload that could have been eliminated with the use of the ERB approach.

234

235

236

237 **Discussion**

238 In this comparative study of the evaluation of early-career funding applications, we found that
239 the simulated funding outcomes of a simplified, expert review-based (ERB) approach agreed
240 well with the official funding outcomes based on the standard, time-tested triage and panel
241 meeting (TPM) format. Applications for fellowships covered a wide range of disciplines,
242 from the humanities and social sciences to STEM, biology and medicine. The agreement was
243 very high for proposals which, in the TPM evaluation, were either allocated to the Fund or
244 Reject categories, but lower in the middle category of proposals that were discussed by the
245 panels. More applicants entered the lottery with the simplified ERB approach than with TPM
246 evaluation. Finally, the simplified ERB evaluation approach was associated with a substantial
247 reduction in costs. Overall, our results support the notion that a sound evaluation of early-
248 career funding applications is possible with an ERB approach.

249 Although panel review is considered as a “de facto” standard, the consistency of
250 decisions from panels has been shown to be limited. For example, previous work by Cole
251 [17], Hodgson [18], Fogelholm [11] and Clarke [19] found an agreement of 65% to 83%
252 between two independent panels evaluating the same set of applications. Thus, in these
253 studies, the funding outcome also depended on the panel that evaluated the application, and
254 not only on the scientific content. Against this background, the agreement of over 80%
255 between ERB and TPM in this study is remarkable. Among the different discipline-specific
256 review panels, our results showed a slightly lower agreement in the humanities and social
257 sciences compared to life sciences and medicine. These differences did not reach
258 conventional levels of statistical significance but were in line with previous findings reported
259 by Pina et al. [13].

260 In the middle group of applications based on the triage step of TPM, the agreement
261 was lower; 64% with three reviewers and 73% with the two reviewers. This is not surprising
262 considering the results from previous studies that suggest that peer review has difficulties in
263 discriminating between applications that are neither clearly competitive nor noncompetitive
264 [20–22]. Agreement between ERB and TPM was also generally lower with ERB using three
265 reviewers than with ERB with two reviewers. An additional reviewer may introduce a
266 different viewpoint. Also, the third reviewer was not a member of the corresponding panel,
267 and not involved in previous panel discussions, which have led to some degree of calibration
268 between assessments of panel members. Such calibration is more difficult to achieve with a
269 remote, ERB approach. However, information and briefing sessions could be held to

270 compensate for the lack of FTF panel meetings. Of note, previous studies reported that
271 reviewers appreciated the social aspects and the camaraderie in FTF settings and that physical
272 meetings are important for building trust among the evaluators [8,9].

273 We found that the panel discussions in the TPM format resulted in higher success rates
274 for women compared to the ERB format. Gender equality is a key concern at the SNSF,
275 which is committed to promoting women in research. The panels will have been aware of the
276 under-representation of female researchers in certain areas, for example, the STEM
277 disciplines, and the SNSF's agenda to promote women. It is, therefore, possible that during
278 the panel deliberations and for funding decisions, the gender of applicants was taken into
279 account in addition to the quality of the proposal.

280 We estimated that about USD 91,000 could have been saved for the two
281 Postdoc.Mobility calls in 2019 if they had been evaluated by ERB rather than by TPM. The
282 meeting costs represented about 55% of the total evaluation costs. In other words, the ERB
283 evaluation based on the two panel reviewers would have cut the expenses by more than half.
284 The experience described here with the junior Postdoc.Mobility fellowship scheme indicates
285 that substantial cost savings could also result from simplifications in the evaluation of other
286 funding instruments at the SNSF. However, any such changes need to be considered carefully.
287 The quality of the evaluation should not be allowed to be compromised because costs may be
288 reduced.

289 To the best of our knowledge, the Health Research Council of New Zealand (HRC-
290 NZ) [23], the Volkswagen Foundation [24], and recently the Austrian Research Fund FWF
291 [25] are the only funders that have used or examined the use of a random selection element in
292 the evaluation process of funding instruments, with a focus on transformative research or
293 unconventional research ideas. The random selection for decisions on applications close to the
294 funding threshold could avoid bias if evaluation criteria do not allow any further
295 differentiation for a small set of similarly qualified applications [22,26]. The applicants were
296 informed about the possible random selection and the evaluation process thus complied with
297 the San Francisco Declaration on Research Assessment (DORA) [27], which states that
298 funders must be explicit about assessment criteria. There was some reservation on the random
299 selection approach among some panel members, but acceptance grew over time. Of note,
300 panels applied the random selection only in a few cases, in eight (6.0%) of 134 applications.
301 In the context of the Explorer Grant scheme of the HRC-NZ, Liu et. al [28] recently reported
302 that most applicants agreed with the use of a random selection. In this study, no negative or
303 positive reactions to the use of random selection were received from applicants.

304 Our study has several limitations. It addressed the specific context of the SNSF
305 Postdoc.Mobility funding scheme and results may not be generalizable to other funding
306 instruments. The sample size was relatively small, and the study lacked statistical power, for
307 example, to examine differences in agreement between TPM and ERB evaluation across
308 disciplines. The two evaluation methods were not independent, since the two assessments of
309 the panel reviewers were used for both methods. We were relying on reviewer evaluation
310 scores which might not always perfectly reflect the quality of the proposed project, might be
311 biased, and depend on the reviewers' previous experience with grant evaluation. However,
312 our study design allowed us to investigate the impact of panel meetings on funding outcomes
313 compared to an ERB approach. This study provides further insights into peer review and a
314 modified lottery approach selection in the context of the evaluation of fellowship applications.
315 More research on the limitations inherent in peer review and grant evaluation is urgently
316 needed. Funders should be creative when investigating the merit of different evaluation
317 strategies [29].

318

319 **Conclusions**

320 In conclusion, we simulated an ERB approach in the evaluation of the junior
321 Postdoc.Mobility funding scheme at the SNSF and compared the funding outcome to the
322 standard TPM format, which has been in use for many years. We found an overall high
323 agreement between the two methods. Discrepancies were mainly observed in the middle
324 group of applications that were discussed in the panel meetings. Based on the evidence that
325 peer review has difficulties in making fine-grained differentiations between meritorious
326 applications [20–22], we are unsure which method performs better. Our findings indicate that
327 the ERB approach represents a viable evaluation method for the Postdoc.Mobility selection
328 process that could save cost and time which could be invested in science and research.

329

330

331 **Acknowledgements**

332 We thank the Management of the SNSF Administrative Offices for approving additional
333 resources for the conduction of this study. We also thank the SNSF Postdoc.Mobility staff of
334 the Administrative Offices for their excellent support in implementing the additional
335 reviewers used for the study.

336

337 **Contributors**

338 Conceived and designed the experiments: MB KR ME. Performed the experiments: MB KR.
339 Analyzed the data: KR RH. Contributed reagents/materials/analysis tools: MB KR ME RH.
340 Wrote the initial draft: MB. Contributed to writing: KR ME RH.

341

342 **Funding**

343 This research received no specific grant from any funding agency in the public, commercial
344 or not-for-profit sectors.

345

346 **References**

- 347 1 Guthrie S, Ghiga I, Wooding S. What do we know about grant peer review in the health
348 sciences? *F1000Res* 2018;**6**:1335. doi:10.12688/f1000research.11917.2
- 349 2 Guthrie S, Rodriguez Rincon D, McInroy G, *et al.* Measuring bias, burden and
350 conservatism in research funding processes. *F1000Res* 2019;**8**:851.
351 doi:10.12688/f1000research.19156.1
- 352 3 Snell RR. Menage a Quoi? Optimal Number of Peer Reviewers. *PLoS ONE*
353 2015;**10**:e0120838. doi:10.1371/journal.pone.0120838
- 354 4 Graves N, Barnett AG, Clarke P. Funding grant proposals for scientific research:
355 retrospective analysis of scores by members of grant review panel. *BMJ*
356 2011;**343**:d4797–d4797. doi:10.1136/bmj.d4797
- 357 5 Herbert DL, Graves N, Clarke P, *et al.* Using simplified peer review processes to fund
358 research: a prospective study. *BMJ Open* 2015;**5**:e008380. doi:10.1136/bmjopen-2015-
359 008380
- 360 6 Bohannon J. Meeting for Peer Review at a Resort That’s Virtually Free. *Science*
361 2011;**331**:27–27. doi:10.1126/science.331.6013.27
- 362 7 Gallo SA, Carpenter AS, Glisson SR. Teleconference versus Face-to-Face Scientific Peer
363 Review of Grant Application: Effects on Review Outcomes. *PLoS ONE* 2013;**8**:e71693.
364 doi:10.1371/journal.pone.0071693
- 365 8 Carpenter AS, Sullivan JH, Deshmukh A, *et al.* A retrospective analysis of the effect of
366 discussion in teleconference and face-to-face scientific peer-review panels. *BMJ Open*
367 2015;**5**:e009138. doi:10.1136/bmjopen-2015-009138
- 368 9 Pier EL, Raclaw J, Nathan MJ, *et al.* Studying the Study Section: How Group Decision
369 Making in Person and via Videoconferencing Affects the Grant Peer Review Process.
370 2015.https://wcer.wisc.edu/docs/working-papers/Working_Paper_No_2015_06.pdf
- 371 10 Vo NM, Quiggle GM, Wadhvani K. Comparative outcomes of face-to-face and virtual
372 review meetings. *International Journal of Surgery Open* 2016;**4**:38–41.
373 doi:10.1016/j.ijso.2016.07.002
- 374 11 Fogelholm M, Leppinen S, Auvinen A, *et al.* Panel discussion does not improve reliability
375 of peer review for medical research grant proposals. *Journal of Clinical Epidemiology*
376 2012;**65**:47–52. doi:10.1016/j.jclinepi.2011.05.001
- 377 12 Obrecht M, Tibelius K, D’Aloisio G. Examining the value added by committee discussion
378 in the review of applications for research awards. *res eval* 2007;**16**:70–91.
379 doi:10.3152/095820207X223785
- 380 13 Pina DG, Hren D, Marušić A. Peer Review Evaluation Process of Marie Curie Actions
381 under EU’s Seventh Framework Programme for Research. *PLoS ONE* 2015;**10**:e0130753.
382 doi:10.1371/journal.pone.0130753

- 383 14 Vener KJ, Feuer EJ, Gorelic L. A statistical model validating triage for the peer review
384 process: keeping the competitive applications in the review pipeline ¹. *FASEB j*
385 1993;**7**:1312–9. doi:10.1096/fasebj.7.14.8224604
- 386 15 Bornmann L, Mutz R, Daniel H-D. Latent Markov modeling applied to grant peer review.
387 *Journal of Informetrics* 2008;**2**:217–28. doi:10.1016/j.joi.2008.05.003
- 388 16 SNSF Postdoc.Mobility Regulations.
389 http://www.snf.ch/SiteCollectionDocuments/Reglement_PM_ab2021_en.pdf
- 390 17 Cole S, Cole, Simon G. Chance and consensus in peer review. *Science* 1981;**214**:881–6.
391 doi:10.1126/science.7302566
- 392 18 Hodgson C. How reliable is peer review? An examination of operating grant proposals
393 simultaneously submitted to two similar peer review systems. *Journal of Clinical*
394 *Epidemiology* 1997;**50**:1189–95. doi:10.1016/S0895-4356(97)00167-4
- 395 19 Clarke P, Herbert D, Graves N, *et al.* A randomized trial of fellowships for early career
396 researchers finds a high reliability in funding decisions. *Journal of Clinical Epidemiology*
397 2016;**69**:147–51. doi:10.1016/j.jclinepi.2015.04.010
- 398 20 Scheiner SM, Bouchie LM. The predictive power of NSF reviewers and panels. *Frontiers in*
399 *Ecology and the Environment* 2013;**11**:406–7. doi:10.1890/13.WB.017
- 400 21 Fang FC, Bowen A, Casadevall A. NIH peer review percentile scores are poorly predictive
401 of grant productivity. *eLife* 2016;**5**:e13323. doi:10.7554/eLife.13323
- 402 22 Klaus B, del Alamo D. Talent Identification at the limits of Peer Review: an analysis of the
403 EMBO Postdoctoral Fellowships Selection Process. *Scientific Communication and*
404 *Education* 2018. doi:10.1101/481655
- 405 23 Health Research Council of New Zealand Explorer Grants.
406 [https://gateway.hrc.govt.nz/funding/researcher-initiated-proposals/2021-explorer-](https://gateway.hrc.govt.nz/funding/researcher-initiated-proposals/2021-explorer-grants)
407 [grants](https://gateway.hrc.govt.nz/funding/researcher-initiated-proposals/2021-explorer-grants)
- 408 24 Volkswagen Foundation Experiment!
409 [https://www.volkswagenstiftung.de/en/funding/our-funding-portfolio-at-a-](https://www.volkswagenstiftung.de/en/funding/our-funding-portfolio-at-a-glance/experiment)
410 [glance/experiment](https://www.volkswagenstiftung.de/en/funding/our-funding-portfolio-at-a-glance/experiment)
- 411 25 FWF 1000 Ideas Programme. [https://www.fwf.ac.at/en/research-funding/fwf-](https://www.fwf.ac.at/en/research-funding/fwf-programmes/1000-ideas-programme/)
412 [programmes/1000-ideas-programme/](https://www.fwf.ac.at/en/research-funding/fwf-programmes/1000-ideas-programme/)
- 413 26 Fang FC, Casadevall A. Research Funding: the Case for a Modified Lottery. *mBio*
414 2016;**7**:e00422-16, /mbio/7/2/e00422-16.atom. doi:10.1128/mBio.00422-16
- 415 27 San Francisco Declaration on Research Assessment (DORA). <https://sfedora.org/>
- 416 28 Liu M, Choy V, Clarke P, *et al.* The acceptability of using a lottery to allocate research
417 funding: a survey of applicants. *Res Integr Peer Rev* 2020;**5**:3. doi:10.1186/s41073-019-
418 0089-z

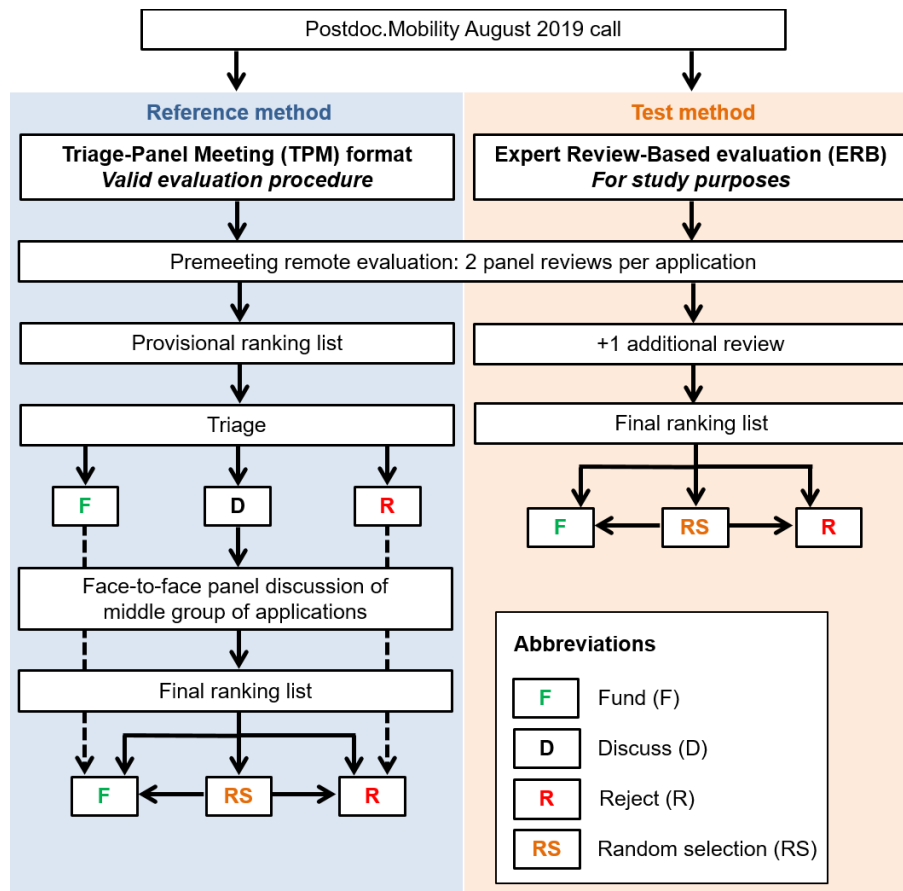
419 29 Severin A, Egger M. Research on research funding: an imperative for science and society.
420 *Br J Sports Med* 2020;:bjssports-2020-103340. doi:10.1136/bjssports-2020-103340

421

422

423 Figures and Tables

424



425

426 **Figure 1. Design of the study comparing the Expert Review-Based evaluation (ERB) with the**

427 **Triage-Panel Meeting (TPM) format.** The ERB and the TPM were dependent in terms of the two

428 assigned panel reviewers per application. The third reviewers were only added for the ERB, their

429 assessments were not considered for the TPM and therefore the official funding outcome.

430

431

432 **Table 1. Success rates by gender of applicants, by discipline and type of evaluation.**
 433

Discipline	All applicants		Women		Men	
	N	N funded (%)	N	N funded (%)	N	N funded (%)
TPM						
All disciplines	134	72 (53.7)	48	29 (60.4)	86	43 (50.0)
Humanities	16	9 (56.2)	9	4 (44.4)	7	5 (71.4)
Social Sciences	23	12 (52.2)	10	7 (70.0)	13	5 (38.5)
STEM	35	19 (54.3)	10	6 (60.0)	25	13 (52.0)
Biology	40	21 (52.5)	14	8 (57.1)	26	13 (50.0)
Medicine	20	11 (55.0)	5	4 (80.0)	15	7 (46.7)
ERB (3 reviewers*)						
All disciplines	134	72 (53.7)	48	27 (56.3)	86	45 (52.3)
Humanities	16	9 (56.3)	9	5 (55.6)	7	4 (57.1)
Social Sciences	23	12 (52.2)	10	6 (60.0)	13	6 (46.2)
STEM	35	19 (54.3)	10	4 (40.0)	25	15 (60.0)
Biology	40	21 (52.5)	14	8 (57.1)	26	13 (50.0)
Medicine	20	11 (55.0)	5	4 (80.0)	15	7 (46.7)
ERB (2 reviewers[‡])						
All disciplines	134	72 (53.7)	48	25 (52.1)	86	47 (54.7)
Humanities	16	9 (56.3)	9	5 (55.6)	7	4 (57.1)
Social Sciences	23	12 (52.2)	10	6 (60.0)	13	6 (46.2)
STEM	35	19 (54.3)	10	4 (40.0)	25	15 (60.0)
Biology	40	21 (52.5)	14	7 (50.0)	26	14 (53.8)
Medicine	20	11 (55.0)	5	3 (60.0)	15	8 (53.3)

434 Abbreviations: N: Number of applications; STEM: Science, Technology, Engineering, Mathematics; TPM: Triage-panel meeting
 435 format; ERB: Expert review-based evaluation.

436 *Two of the three expert reviewers were also members of the evaluation panel.

437 [‡]Both expert reviewers were also members of the evaluation panel.

438

439

440 **Table 2. Agreement between the simulated expert review-based (ERB) evaluation and the**
 441 **triage-panel meeting (TPM) format, by discipline.**
 442

Discipline	N	Funded by TPM		Agreement (%) (95% Wald CI)
		Yes	No	
Funded by ERB (3 reviewers*)				
All disciplines	Yes	59	13	80.6
	No	13	49	(73.9-87.3)
Humanities	Yes	7	2	75.0
	No	2	5	(53.8-96.2)
Social Sciences	Yes	9	3	73.9
	No	3	8	(56.0-91.8)
STEM	Yes	15	4	77.1
	No	4	12	(63.2-91.0)
Biology	Yes	18	3	85.0
	No	3	16	(73.9-96.1)
Medicine	Yes	10	1	90.0
	No	1	8	(76.9-100)
<i>P</i> -value				0.58
Funded by ERB (2 reviewers[‡])				
All disciplines	Yes	63	9	86.6
	No	9	53	(80.8-92.4)
Humanities	Yes	7	2	75.0
	No	2	5	(53.8-96.2)
Social Sciences	Yes	11	1	91.3
	No	1	10	(79.8-100)
STEM	Yes	16	3	82.9
	No	3	13	(70.4-95.4)
Biology	Yes	19	2	90.0
	No	2	17	(80.7-99.3)
Medicine	Yes	10	1	90.0
	No	1	8	(76.9-100)
<i>P</i> -value				0.51

443 Abbreviations: N: Number of applications; CI: Confidence interval; STEM: Science, Technology, Engineering, Mathematics;

444 ERB: Expert review-based evaluation; TPM: Triage-panel meeting format.

445 *P*-values for differences in agreement across disciplines from chi-squared test.

446 *Two of the three expert reviewers were also members of the evaluation panel.

447 [‡]Both expert reviewers were also members of the evaluation panel.

448

449 **Table 3. Agreement between the simulated expert review-based (ERB) evaluation and the**
 450 **triage-panel meeting (TPM) format, by gender.**
 451

Gender		Funded by TPM		Agreement (%) (95% Wald CI)
		Yes	No	
Funded by ERB (3 reviewers*)				
Women	Yes	24	3	83.3
	No	5	16	(72.7-93.9)
Men	Yes	35	10	79.1
	No	8	33	(70.5-87.7)
<i>P</i> -value				0.71
Funded by ERB (2 reviewers[‡])				
Women	Yes	24	1	87.5
	No	5	18	(78.1-96.9)
Men	Yes	39	8	86.0
	No	4	35	(78.7-93.3)
<i>P</i> -value				0.99

452 Abbreviations: N: Number of applications; CI: Confidence interval; STEM: Science, Technology, Engineering, Mathematics;

453 ERB: Expert review-based evaluation; TPM: Triage-panel meeting format.

454 *P*-values for differences in agreement across genders from chi-squared test.

455 *Two of the three expert reviewers were also members of the evaluation panel.

456 [‡]Both expert reviewers were also members of the evaluation panel.

457

458

459 **Table 4. Agreement between the simulated expert review-based (ERB) evaluation and the**
 460 **triage-panel meeting (TPM) format, by triage results.**
 461

Triage result		Funded by TPM		Agreement (%) (95% Wald CI)
		Yes	No	
Funded by ERB (3 reviewers*)				
Fund (F)	Yes	36	0	97.3
	No	1	0	(92.1-100)
Discuss (D)	Yes	23	12	64.2
	No	12	20	(52.7-75.7)
Reject (R)	Yes	0	1	96.7
	No	0	29	(90.3-100)
<i>P</i> -value				<0.001
Funded by ERB (2 reviewers[‡])				
Fund (F)	Yes	37	0	100
	No	0	0	
Discuss (D)	Yes	26	9	73.1
	No	9	23	(62.5-83.7)
Reject (R)	Yes	0	0	100
	No	0	30	
<i>P</i> -value				<0.001

462 Abbreviations: N: Number of applications; CI: Confidence interval; STEM: Science, Technology, Engineering, Mathematics;

463 ERB: Expert review-based evaluation; TPM: Triage-panel meeting format.

464 *P*-values for differences in agreement across triage groups from chi-squared test.

465 *Two of the three expert reviewers were also members of the evaluation panel.

466 [‡]Both expert reviewers were also members of the evaluation panel.

467