

# UC Santa Cruz

## UC Santa Cruz Previously Published Works

### Title

The reproducibility crisis meets stock assessment science: Sources of inadvertent bias in the stock assessment prioritization and review process

### Permalink

<https://escholarship.org/uc/item/2q73z9ss>

### Author

Satterthwaite, William H

### Publication Date

2023-10-01

### DOI

10.1016/j.fishres.2023.106763

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NoDerivatives License, available at <https://creativecommons.org/licenses/by-nd/4.0/>

Peer reviewed

1 *Accepted manuscript version of Satterthwaite, W. H. 2023. The reproducibility crisis meets stock*  
2 *assessment science: Sources of inadvertent bias in the stock assessment prioritization and*  
3 *review process. Fisheries Research 266:106763. <https://doi.org/10.1016/j.fishres.2023.106763>*  
4

5 **Title: The reproducibility crisis meets stock assessment science: Sources of inadvertent**  
6 **bias in the stock assessment prioritization and review process**

7 Author: William H. Satterthwaite

8 Fisheries Ecology Division, Southwest Fisheries Science Center, National Marine Fisheries  
9 Service, National Oceanic and Atmospheric Administration, Santa Cruz, CA, USA

10  
11 Corresponding author. 110 McAllister Way, Santa Cruz, California, USA, 95060.  
12 will.satterthwaite@noaa.gov

13  
14 **Highlights**

- 15 • Prioritization, review, and adoption of assessments is subject to inadvertent bias
- 16 • Many of the factors introducing bias are analogous to "p-hacking" in science broadly
- 17 • This compromises the interpretation of risk metrics based on probability statements
- 18 • Bias may be comparable to differences between commonly-applied uncertainty buffers
- 19 • Solutions/mitigations proposed for p-hacking broadly have analogs for use here

20  
21 **Abstract**

22 The broader scientific community is struggling with a reproducibility crisis brought on by  
23 numerous factors, including "p-hacking" or selective reporting that may increase the rate of false  
24 positives or generate misleading effect size estimates from meta-analyses. This results when  
25 multiple modeling approaches or statistical tests may be brought to bear on the same problem,  
26 and there are pressures or rewards for finding "significant" results. Fisheries science is unlikely  
27 to be immune to this problem, with numerous opportunities for bias to inadvertently enter into  
28 the process through the prioritization of stocks for assessment, decisions about competing  
29 model approaches or data treatments within complex assessment models, and decisions about  
30 whether to adopt assessments for management after they are reviewed. I present a simple  
31 simulation model of a system where many assessments are performed each management cycle  
32 for a multi-stock fishery, and show how asymmetric selection of assessments for extra scrutiny  
33 or re-assessment within a cycle can turn a process generating unbiased advice on fishing limits  
34 into one that is biased high. I show similar results when sequential assessments receive extra  
35 scrutiny if they show large proportional decreases in catch limits compared to a prior  
36 assessment for the same stock, especially if there are only small changes in true stock size or  
37 status over the interval between assessments. The level of bias introduced by a plausible level  
38 of asymmetric scrutiny is unlikely to fundamentally undermine scientific advice, but may be  
39 sufficient to compromise the nominal "overfishing probabilities" used in a common framework for  
40 accommodating uncertainty, and introduce a level of bias comparable to the difference between  
41 buffers corresponding to commonly-applied levels of risk tolerance.

42  
43 **Keywords:** Reproducibility, bias, p-hacking, unintended consequences  
44

## 45 1. Introduction

46

47 While reproducibility has long been a cornerstone of science, in the last decade there  
48 has been an explosion of references to a reproducibility “crisis” (Fanelli 2018). This has led to  
49 alarming suggestions that most published scientific studies are false (Ioannidis 2005), or that  
50 meta-analyses synthesizing average effect sizes across studies may yield misleading results  
51 (Schooler 2011). Other authors have taken more nuanced views suggesting that the problem,  
52 while real, may cause fewer false positives (Jager and Leek 2014) or confound meta-analyses  
53 less (Head et al. 2015) than the worst-case predictions. Nevertheless, reliability and  
54 reproducibility remain major concerns across all fields of science (Baker 2016), with the  
55 potential to affect both the advancement of knowledge and specific policy and management  
56 decisions guided by scientific advice. Recent challenges to reproducibility of hot topics in marine  
57 science (e.g., Clark et al. 2020, Provencher et al. 2020) suggest that fisheries science, which is  
58 often closely linked to economically and culturally impactful policy and management decisions,  
59 is unlikely to be immune to these problems.

60 Although some instances of non-reproducible science may reflect data fabrication or  
61 other scientific misconduct (Viglione 2020), this is widely believed to be uncommon relative to  
62 much more prevalent issues resulting from selective reporting of statistically “significant” results,  
63 pressure to publish, and improper application of statistics or experimental design (Baker 2016).  
64 The data and models used in fisheries stock assessments vary by region and stock, but are  
65 often subject to multiple layers of oversight and review, which should make outright data  
66 fabrication or other forms of deliberate misconduct unlikely in systems with robust review  
67 processes. However, I will argue that fisheries science and stock assessments can be, and  
68 likely are, affected by forces analogous to other sources of the reproducibility crisis.

69 Much of the lack of reproducibility in the scientific literature may be attributable to “p-  
70 hacking” on the part of researchers (Benjamin et al. 2018), and/or selective reporting on the part  
71 of authors and journals that may only, or at least preferentially, publish statistically “significant”  
72 results (Schooler 2011). Statistical “significance” is often established in a frequentist setting via  
73 a null hypothesis testing framework, wherein results are deemed “significant” if a null model  
74 suggests less than a 5% probability (p-value) of generating a pattern at least as strong as the  
75 one observed in the data by chance alone (Wasserstein and Lazar 2016). In “p-hacking”,  
76 multiple statistical tests are performed, but only the “significant” results are retained and  
77 reported, meaning that false positives are likely to be reported at a higher rate than the nominal  
78 value assigned for a single test. This is a particular problem when analysts have large datasets

79 with multiple potential predictors and/or response variables (leading to decisions on which to  
80 include and how to weight them) and a variety of plausible statistical model forms to test.  
81 However, even in simpler situations where only a single testing approach is considered,  
82 interpretation of p-values is clouded when additional data are accumulated and tests are  
83 repeated for larger datasets without clear stopping rules (Wicherts et al. 2016).

84 P-hacking is not necessarily done with ill intent, rather it may reflect an innocent lack of  
85 understanding of statistics (Peng 2015) or the strong ability of post-hoc reasoning to convince  
86 scientists that whatever course of decisions led to the “significant” (and publishable, thus  
87 rewarding) analysis were the correct ones (Simmons et al. 2011). Even when the scientists  
88 performing individual studies do everything right, if journals tend to publish “significant” results  
89 while rejecting inconclusive studies, this may lead to larger mean effects in the published  
90 literature compared to the means that would be concluded from all valid studies performed,  
91 including those that did not produce “significant” results.

92

## 93 **2. P-hacking and selective reporting: parallels in fisheries and stock assessments**

94

95 Stock assessments are fundamental to scientifically informed fisheries management  
96 (Hilborn and Walters 1992, Hilborn et al. 2020) for purposes such as setting allowable catch  
97 limits. These catch limits are often set using control rules intended to maintain population  
98 abundance and spawning output levels near those expected to produce maximum sustainable  
99 yield (Melnychuk et al. 2013, Methot et al. 2014). Stock assessments are also the product of  
100 often-complex models based on imperfectly-measured data and that require numerous choices  
101 on the part of the analysts about data to include versus exclude, how to weight different data  
102 sets, parameters to fix versus estimate, and functional forms to assume (Maunder and Piner  
103 2015, 2017). As a result, stock assessment outputs are unavoidably uncertain (Hilborn and  
104 Walters 1992, Mildenerger et al. 2022), and different but more or less equally defensible  
105 decisions on the part of the analysts (and/or reviewers, who often drive final model form in  
106 conjunction with the original analysts) could lead to different results (Ralston et al. 2011).

107 One approach to dealing with this uncertainty that has been adopted in multiple regions  
108 of the United States, and in somewhat similar forms in other countries, is the  $P^*/\sigma$  approach  
109 (Shertzer et al. 2008). This approach assumes that overfishing limits (OFLs) are estimated  
110 without bias, but with uncertainty expressed by assuming a lognormally distributed ratio  
111 between the true OFL and the assessed OFL, where the median is equal to one and the log-  
112 scale standard deviation is  $\sigma$  (see Ralston et al. 2011 and Privitera-Johnson and Punt

113 2020a for approaches to estimating sigma). Then, an acceptable biological catch (ABC) is  
114 determined by multiplying the OFL by the  $P^*$  quantile of the distribution. If all statistical  
115 assumptions were met, this would result in a  $P^*$  probability that the ABC exceeds the “true” OFL  
116 that would have been estimated given perfect knowledge. When applied to multiple stocks  
117 simultaneously, this implies that a fraction  $P^*$  of the ABCs established for a multi-stock fishery  
118 would be larger than the OFLs that would have been estimated given perfect knowledge,  
119 loosely analogous to the expectation that 5% of scientific results reported as significant at the  
120  $p < 0.05$  level would be false positives. Just as p-hacking and selective reporting may lead to a  
121 higher than nominal false positive rate in the scientific literature, and over-estimate mean effect  
122 size, it would seem unavoidable that if analogous pressures operate in the analysis and  
123 adoption of assessments of stock status and fishing limits, the interpretation of  $P^*$  may be  
124 similarly clouded.

125

## 126 *2.1 Scope of the problem*

127

128 I suggest that just like academic scientists may face pressure to produce studies with  
129  $p < 0.05$  that can be published (with resultant career benefits), and journals may be more likely to  
130 publish “significant” results, stock assessors (and the review bodies that can influence the final  
131 structure of stock assessments) may face pressures to produce “favorable” results, and/or  
132 management bodies may be more likely to adopt and use assessments perceived as “less  
133 pessimistic” (Seagraves and Collins 2012).

134

### 135 *2.1.1 Assessment prioritization*

136

137 In most regions, there are far more stocks in need of assessment than there are  
138 resources to assess them. As a result, management agencies in the United States (primarily the  
139 National Marine Fisheries Service that implements many assessments and Fishery  
140 Management Councils that lead development of management responses) have adopted a  
141 comprehensive prioritization process (Methot 2015), which typically considers many factors  
142 (NMFS 2022) including economic and ecological importance, trends in survey data, time since  
143 last assessment, and status at most recent assessment. Of note, although this is only one of the  
144 many factors considered, stocks which were last assessed to be in high status are given the  
145 lowest score for the status component of their overall prioritization score, while stocks recently  
146 assessed to be in poor status are given the highest score – higher even than stocks which lack

147 recent (or even any) assessments and with attributes associated with high vulnerability (i.e.,  
148 stocks with low productivity and high susceptibility to the fishery [Cope et al. 2011]). This  
149 strategy runs the risk that a new assessment of a low-status stock means a new error in the  
150 inevitably uncertain assessment; thus, a more favorable assessment is not conclusive proof of  
151 better status. When considering repeated assessments of multiple stocks overall, this strategy  
152 could lead to an asymmetry where there may be more chances to incorrectly reverse an  
153 assignment of poor status than there are to incorrectly re-assign a stock from good status to  
154 poor. At the same time, it may divert resources from other stocks in need of assessment where  
155 an assessment might do more to address important management uncertainties (Cadrin et al.  
156 2015).

157         Given limited resources for assessments and the large number of species/stocks to  
158 assess, stock assessment analysts have also developed “data-moderate” approaches that  
159 consider fewer types of input data, have less flexibility in choices among alternative  
160 assumptions, and have fewer model structure alternatives (e.g., Rudd et al. 2021), with the  
161 hope of increasing throughput. When this approach was first proposed to the Pacific Fishery  
162 Management Council (PFMC), it was suggested that the output of a “data-moderate”  
163 assessment might be acceptable for use in management if it returned a favorable estimate of  
164 stock status but not be used as the basis for determining that a stock was overfished (PFMC  
165 2013), or that there be an option for an “out-of-cycle” assessment to provide a second estimate  
166 of status before adopting an overfished status from a data-moderate assessment (NMFS 2013).  
167 Such an asymmetric standard of proof, especially when confounded with the lack of priority  
168 given to re-assessing stocks with favorable assessment outputs, seems likely to introduce bias  
169 at the level of the suite of stocks subject to the same management process, such as all the  
170 stocks in a Fishery Management Plan (FMP).

171

### 172 *2.1.2 Conduct of assessments and reviews*

173

174         Once a stock has been chosen for assessment, stock assessment analysts still face  
175 numerous decisions about the specific datasets to include as well as the treatment of putative  
176 “outliers” within accepted datasets, the weightings applied to different data sources, potential  
177 use of priors, parameters to fix versus estimate, and various functional forms. Assessments are  
178 typically subjected to review panels where data treatments and other modeling choices are  
179 scrutinized and alternatives are explored. In the vast majority of cases, the model endorsed at  
180 the end of the review process has some differences from the initially proposed base model,

181 reflecting the combined deliberations of the assessors and reviewers. Ideally, the reviewers  
182 would be guided solely by scientific considerations, but just as with academic scientists  
183 (Simmons 2011), there is the potential for conscious or subconscious considerations leading to  
184 post-hoc reasoning to support the outcome perceived as likely to be most palatable to  
185 managers at the next step in the review and adoption process. There may also be incentives to  
186 be more critical of proposed model changes that reduce status, or to be less likely to  
187 recommend ultimate acceptance of an assessment yielding low status. These pressures may be  
188 most acute when the reviewers are drawn from bodies whose members are dependent on  
189 managers or politicians for their appointment or renewal (Crosson 2013).

190

### 191 *2.1.3 Adoption of assessments for management*

192

193           Following the initial peer review, adoption of stock assessments by Fishery Management  
194 Councils in the United States comes after review by a Scientific and Statistical Committee  
195 (SSC). Ideally, the SSC would apply equal standards of proof for acceptance of any  
196 assessment, but consciously or subconsciously they may apply extra scrutiny to assessments  
197 outputting poor status, although assessments yielding unexpectedly positive outcomes have  
198 also faced extra scrutiny. Once the SSC has endorsed an assessment, although Councils “may  
199 not exceed” the fishing level recommendations of the SSC, the Council must act to formally  
200 adopt the assessments recommended by the SSC, and this does not always happen (Crosson  
201 2013, Nies 2022). Councils are intended to represent the public interest, but with an emphasis  
202 on the fishing industry. For example, as of 2022 the 72 appointed seats on U.S. Fishery  
203 Management Councils consisted of 29 representatives of commercial fishing interests, 25  
204 representatives of recreational fishing interests, and 18 representatives of "other" interests  
205 (NMFS 2021), which can include tribal fishery representatives and individuals formerly more  
206 closely associated with a fishing interest. This in itself is of course not a definitive basis for  
207 concluding that Councils are more likely to reject a “pessimistic” assessment and accept an  
208 “optimistic one”, but suggests that it could be a reasonable expectation. Indeed, in response to  
209 concerns expressed about increased and repeated scrutiny of poor-status assessments by the  
210 PFMC in 2021 (SSC 2021), a voting Council member stated that from the perspective of a  
211 manager and/or Council member, it is logical to expect that more attention will be given to more  
212 pessimistic assessments, consistent with the need to instill confidence for managers and  
213 stakeholders that the results are robust (SSC 2022, p. 13), suggesting that managers do indeed  
214 apply different standards of proof depending on management implications and may not

215 appreciate the parallels with p-hacking or the potential biases introduced by shifting standards  
216 of proof.

217

## 218 *2.2 Quantifying the problem*

219

220 To quantify the likely magnitude of bias that might be introduced by various approaches  
221 to selecting assessments remanded for further scrutiny, revised, and/or re-assessed on a very  
222 short timeline (i.e., before another year is added to the assessment), I developed a simple  
223 simulation model. As with any model, it requires numerous simplifying assumptions, and does  
224 not incorporate all of the potential qualitative sources of bias described above. Nevertheless, I  
225 hope it is useful in demonstrating the potential scope of the problems introduced by asymmetric  
226 scrutiny during the review and adoption steps of the stock assessment process.

227 I simulated a system in which a large number of stocks are assessed in each  
228 assessment cycle, with the assumption that (prior to any additional selective scrutiny at the  
229 review and adoption stage) assessments are uncertain but median-unbiased. I then examined  
230 the distributional properties of the outputs of the full suite of assessments after a selected  
231 subset of assessments had been re-done within the same cycle, with the assumption that the  
232 redone assessments were also median unbiased, exploring various scenarios for the degree of  
233 independence between the initial and redone assessment. I also explored a scenario where  
234 sequential assessments are performed for a stock, and the degree of scrutiny applied to the  
235 assessment done at the later timestep depends on the proportional difference in OFL compared  
236 to the outcome of the assessment from the first timestep (Section 2.2.3).

237

### 238 *2.2.1 Model structure – assessments redone in same year*

239

240 Following the assumptions at the heart of the P\*/sigma approach, I assumed that for the  
241 initial version of each assessment:

242

$$243 \quad 1) \quad \log\left(\frac{OFL_{true}}{OFL_{assessed}}\right) = \epsilon_s + \epsilon_a$$

244

245 where

246

$$247 \quad 2) \quad \epsilon_s \sim Normal(0, \sigma_s)$$



248

249 represents persistent errors for a given stock (arising from underlying issues in the primary data,  
250 assumptions shared across all candidate models [e.g., assumptions about steepness (Thorson  
251 et al. 2019) or natural mortality (Hamel 2014)], and if applicable persistent assessor and/or  
252 reviewer effects – and how all of these interact with the biology of the stock in question), and  
253

254 3)  $\epsilon_a \sim Normal(0, \sigma_a)$

255

256 represents assessment-specific errors arising from choices about datasets used, data  
257 weightings, and selection of model assumptions from within the general scope of acceptable  
258 assumptions.

259

260 To preserve the assumption that if all assessments are performed only once,

261

262 4)  $\log\left(\frac{OFL_{true}}{OFL_{assessed}}\right) \sim Normal(0, \sigma)$

263

264 I used the equation for variance of a sum of random variables (assuming independence, which  
265 may be reasonable given how the respective errors were defined, but potentially problematic if  
266 certain features of the data lead to a tendency to select certain modeling options) to obtain

267

268 5)  $\sigma_a = \sqrt{\sigma^2 - \sigma_s^2}$ .

269

270 To simulate an assessment cycle that assesses  $N$  total stocks, I simulated the  
271 distribution of true versus assessed OFLs by drawing vectors of length  $N$  for  $\epsilon_s$  and  $\epsilon_a$  across  
272 various values of  $\sigma_s$  (chosen such that anywhere from 0% to 100% of the variance of  
273 assessment outputs was driven by assessment-specific factors) while holding  $\sigma$  constant at 0.5  
274 (and thus determining  $\sigma_a$  via equation 5), then exponentiated the sum of these vectors. 0.5 is  
275 the default value of sigma applied by the PFMC for “category 1” assessments, the most data-  
276 rich and complex models used (PFMC 2022). To simulate a distribution of true versus assessed  
277 OFLs after some stocks were re-assessed, I retained all draws of  $\epsilon_s$ , and all draws of  $\epsilon_a$  for  
278 stocks that were not re-assessed, while drawing new values of  $\epsilon_a$  for re-assessed stocks. I then  
279 exponentiated the summed vectors as before. For realistic values of  $N$ , stochastic variation from

280 run-to-run is expected to predominate, thus I chose  $N=2,000,000$  to approximate the asymptotic  
281 expectation. I focused on  $P^*$  values of 0.45 and 0.40, as they are the most commonly used by  
282 the PFMC.

283 I explored various scenarios for the selection of assessments to be redone within a  
284 single assessment cycle: 1) selecting  $x\%$  of assessments at random (exploring varying levels of  
285  $x$ ), 2) assuming skillful selection of problematic assessments by including the  $x/2\%$  highest and  
286  $x/2\%$  lowest ratios of true versus assessed OFLs, or 3) assuming skillful selection of  
287 problematic assessments aimed only at the pessimistic ones by including the  $x\%$  highest ratios  
288 of true:assessed OFLs.

289

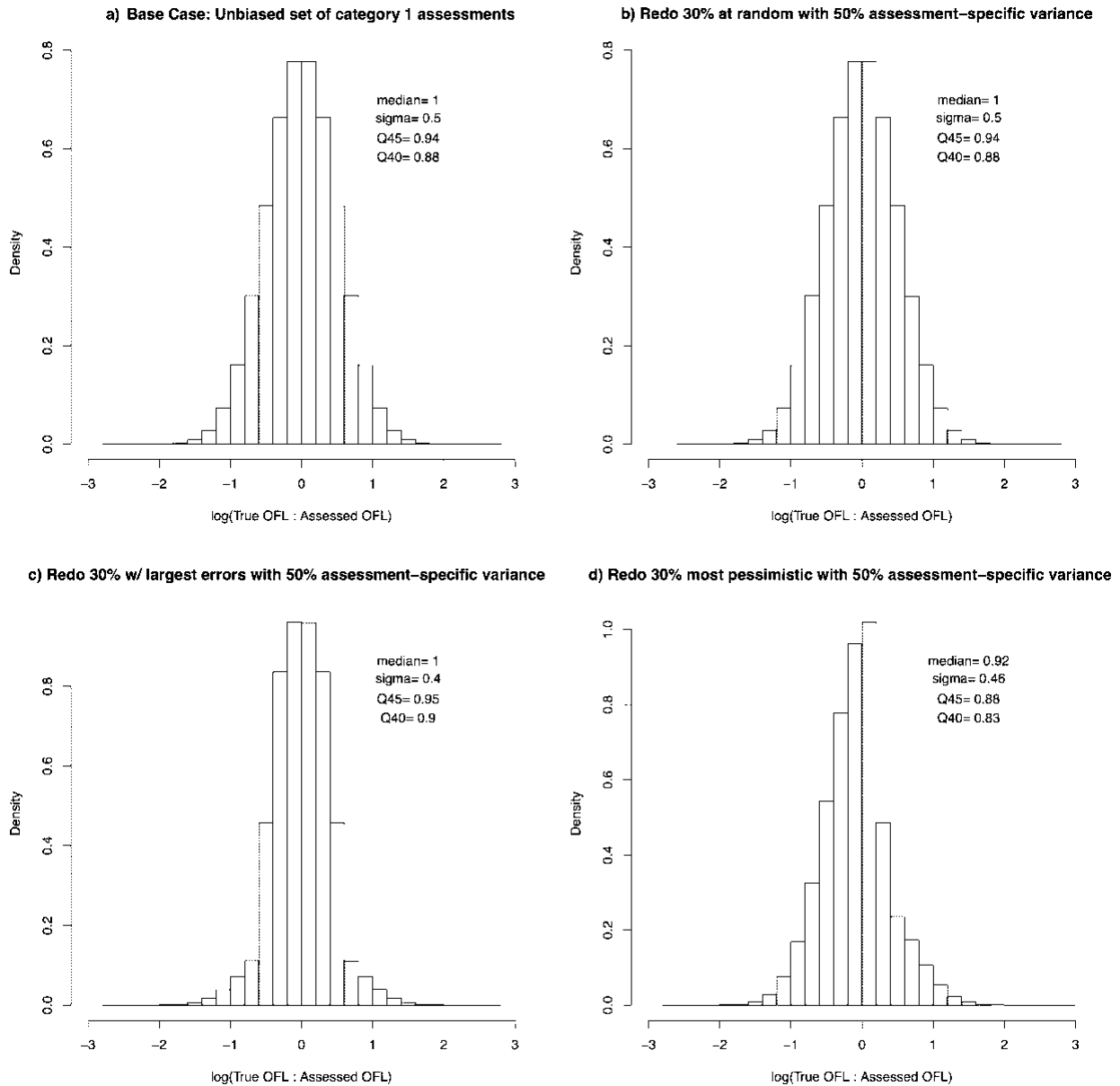
### 290 *2.2.2 Model outputs for within-cycle scrutiny*

291

292 As expected, redoing assessments at random does not change the FMP-wide  
293 distribution of ratios between true and assessed OFLs (compare Figure 1a versus Figure 1b). If  
294 inaccurate assessments are skillfully selected regardless of the direction of error, the median  
295 ratio remains fixed at 1.0 (i.e., there is still no bias at the FMP-wide level) while the distribution  
296 narrows (but remains symmetric while no longer lognormal) and sigma becomes smaller (Figure  
297 1c). If the least accurate assessments with errors in the direction of poor status are redone, the  
298 median ratio drops below 1.0 (indicating FMP-wide bias) and the distribution becomes non-  
299 symmetric (Figure 1d). The degree of bias introduced increases with the fraction of  
300 assessments redone and with the fraction of variance in assessment outputs attributable to  
301 assessment-specific factors (Figure 2). Note that this FMP-wide bias occurs even though the  
302 individual assessments and re-assessments are assumed to provide unbiased estimates.

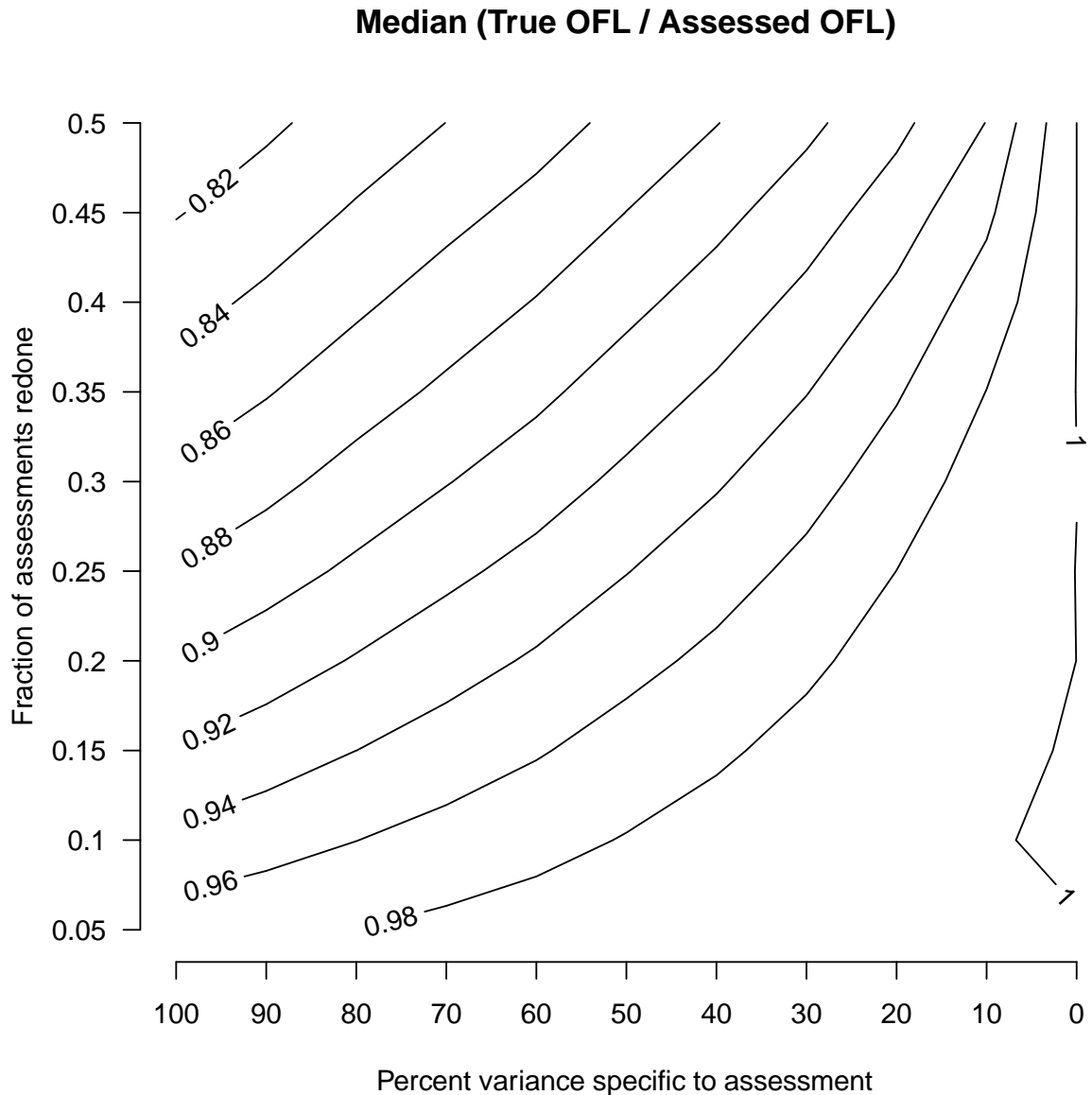
303

304 **Figure 1.** Example outputs showing the ratio between true and assessed overfishing limits  
 305 (OFLs) for the initial set of unbiased assessments (a), after redoing some assessments within-  
 306 cycle at random (b), redoing the least accurate assessments within-cycle regardless of the  
 307 direction of error (c), or redoing (within-cycle) the assessments with the largest errors in the  
 308 direction of low status (d). Q45 and Q40 denote the 45<sup>th</sup> and 40<sup>th</sup> quantiles, respectively.



309  
 310

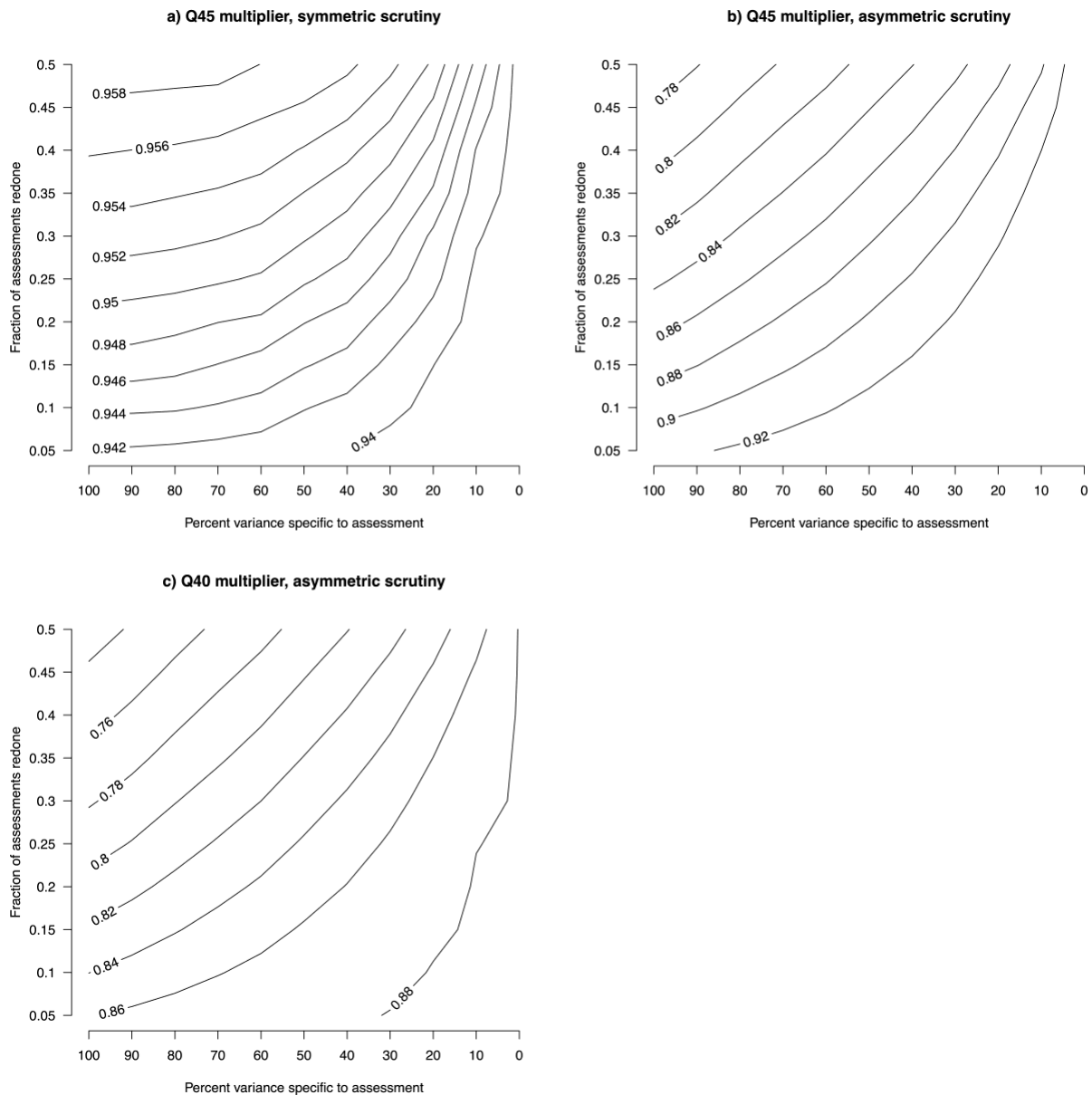
311 **Figure 2.** Median ratio between the true and assessed OFL following the skillful selection of  
 312 assessments with the largest proportional errors in the direction of low status to be redone  
 313 within a single assessment cycle.



314  
 315 These changes in the distribution of true versus assessed OFLs in the skillful selection  
 316 scenarios for redoing assessments lead to changes in the multiplier needed to result in a  
 317 specified probability that the ABC calculated by applying the multiplier to the assessed OFL will  
 318 be higher than the true OFL. In the base case of lognormal distribution with  $\sigma=0.5$ , a buffer  
 319  $Q45=0.945$  results in the expectation that the ABC will be greater than the true OFL 45% of the  
 320 time. If assessments are skillfully selected to redo, but without attention to the direction of error,

321 a slightly larger multiplier (i.e., reduced buffer between ABC and OFL) can be used to achieve  
322 the same expectation (Figure 3a), although the change in multiplier is small. In contrast, skillfully  
323 selecting the most pessimistic assessments to be redone requires larger changes in the  
324 multiplier to preserve the nominal probability of ABCs exceeding the true OFL, and the multiplier  
325 needs to be decreased (Figure 3b). Even larger changes in the multiplier are required to  
326 achieve a nominal 40% probability of ABCs exceeding the true OFL (Figure 3c), which requires  
327 a multiplier of 0.881 in the base case where no assessments are redone.  
328

329 **Figure 3.** Multiplier required to achieve a 45% (Q45, panels a and b) or 40% (Q40, panel c)  
 330 probability that the ABC obtained by applying the multiplier to the assessed OFL is larger than  
 331 the true OFL, under skillful and symmetric selection of the least accurate assessments to be  
 332 redone within an assessment cycle (a) or under skillful selection of the most pessimistic  
 333 assumptions to be redone within an assessment cycle (b and c). Note that the contour spacings  
 334 are different in panel a versus b and c.



335  
 336

337 2.2.3 Responses to changes between consecutive assessments

338

339 So far, my quantitative analysis focused on actions taken within a single assessment  
340 cycle, such that both the original and redone assessment are estimating status for the same  
341 terminal year and setting catch limits for the same management year(s). It may also be the case  
342 that managers and reviewers would give extra scrutiny to an assessment of a stock that differed  
343 substantially in its status estimate or OFL determination compared to the previous assessment  
344 of that stock, although such changes could be unsurprising given long intervals between  
345 assessments or significant changes in the environment or management affecting the stock.

346 To simulate a system where consecutive assessments of the same stock are compared  
347 and unexpected changes may trigger further scrutiny of the more recent assessment, I assumed  
348 a set of OFLs was determined for a suite of stocks in both timestep 1 and timestep 2. In  
349 timestep 1 I assumed the ratios between true and assessed OFLs were determined as in  
350 equation 1, but with  $\epsilon_a$  subscripted by timestep to reflect its potential to vary between timestep 1  
351 and timestep 2:

352 6) 
$$\log\left(\frac{OFL_{true,1}}{OFL_{assessed,1}}\right) = \epsilon_s + \epsilon_{a,1}.$$

353 I assumed that the dynamics of the stock between timestep 1 and timestep 2 changed  
354 the true OFL by a proportion given by a lognormal distribution:

355 7) 
$$\log\left(\frac{OFL_{true,2}}{OFL_{true,1}}\right) = \partial,$$

356 where

357 8) 
$$\partial \sim Normal(0, \sigma_d)$$

358 and  $\sigma_d$  represents the variability in stock dynamics. For each stock, the proportional difference  
359 between assessed OFLs in timestep 2 and timestep 1 is

360 9) 
$$\log\left(\frac{OFL_{assessed,2}}{OFL_{assessed,1}}\right) = \partial + \epsilon_{a,2,initial} - \epsilon_{a,1}$$

361 where  $\epsilon_{a,2,initial}$  is the assessment-specific error associated with the initial iteration of the  
362 timestep 2 assessment (with the “initial” subscript reflecting the potential that some timestep 2  
363 assessments will be redone), noting that  $\epsilon_s$  affects the assessed OFLs in both timesteps equally  
364 and so it drops out of the comparison.

365 For the initial set of timestep 2 assessments,

366 10) 
$$\log\left(\frac{OFL_{true,2}}{OFL_{assessed,2,initial}}\right) = \epsilon_s + \epsilon_{a,2,initial}$$

367 and I simulated scenarios where the timestep 2 assessments with the largest proportional  
368 changes in OFL compared to timestep 1 (as determined in equation 9) were redone, resulting in  
369 new errors:

$$370 \quad 10) \quad \log\left(\frac{OFL_{true,2}}{OFL_{assessed,2,redone}}\right) = \epsilon_s + \epsilon_{a,2,redone}$$

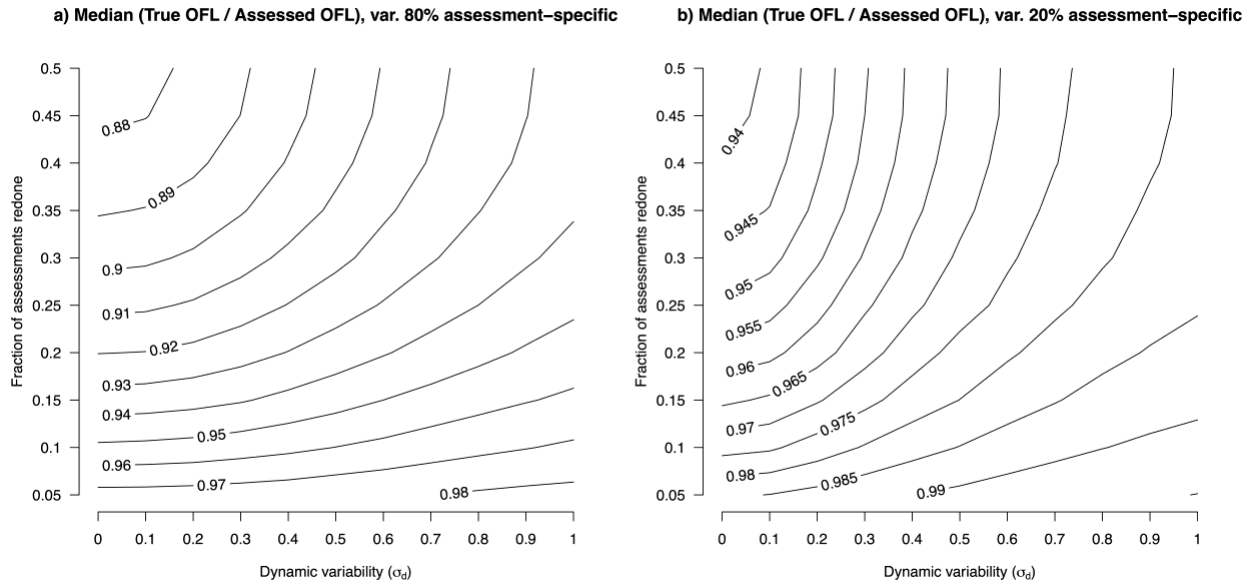
371 and explored the distribution of the ratio between true and assessed OFLs given different levels  
372 of  $\sigma_d$  and different proportions of assessments being redone within timestep 2. I explored  
373 scenarios where the  $x\%$  of largest proportional decreases in OFL or the  $x/2\%$  of largest  
374 proportional changes in either direction led to re-assessment. For these simulations, I held  $\sigma_a$   
375 constant at either 0.447 (80% of total variance given  $\sigma=0.5$ ) or 0.225 (20% of total variance) to  
376 reduce the dimensionality of the simulations. Given the meta-analytic approach to estimating  $\sigma$   
377 based on repeated assessments (Ralston et al. 2011, Privitera-Johnson and Punt 2020b), it is  
378 likely that estimates of  $\sigma$  are dominated by assessment-specific factors, with intrinsic factors  
379 largely constant across assessments and so not revealed by comparison of sequential  
380 assessments. Therefore, the larger value of  $\sigma_a$  may be more appropriate, since  $\sigma_a$  likely drives  
381 current estimates of  $\sigma$ .

382         Applying extra scrutiny to large proportional decreases in OFL from one assessment to  
383 the next can yield similar biases to redoing assessments within-cycle if the changes in true  
384 stock size/status between assessment periods is small relative to the error in assessments  
385 (compare Figure 4 to Figure 2, comparing Figure 4a to the part of Figure 2 where a large  
386 proportion of variance is assessment-specific and 4b to the part of Figure 2 where a small  
387 proportion of variance is assessment-specific). As the variation in true stock size/status between  
388 assessment periods becomes larger, the bias introduced is less (Figure 4) and results become  
389 more similar to picking assessments to be redone at random, because when changes in true  
390 stock size/status between assessments are large, assessment error has relatively little effect on  
391 the probability of observing a large overall change in the estimated OFL. As expected, redoing  
392 assessments that show large proportional changes in OFLs regardless of the direction of  
393 change does not introduce a bias (results not shown, but are produced by the code available  
394 online). Implications for appropriate uncertainty buffers or multipliers are similar to the effects on  
395 bias, requiring larger buffers in the case of directional scrutiny and having minimal effects on  
396 appropriate buffers in the case of symmetric scrutiny (results not shown, but are produced by  
397 the code available online).

398



399 **Figure 4.** Median ratio between the true and assessed OFL after redoing the second in a pair of  
 400 sequential assessments if there was a large proportional reduction in the OFL for the second  
 401 assessment compared to the first. In panel a, 80% of the variance associated with assessment  
 402 error is due to assessment-specific factors, whereas in panel b 80% of the variance is  
 403 associated with intrinsic factors that do not vary between assessment iterations.



404  
 405 **2.2.4 Caveats and potential extensions**

406  
 407 The models presented here are admittedly oversimplifications of complicated processes  
 408 where there may be no bright line between “redoing” an assessment within a cycle and the  
 409 revisions that normally occur during the process of model development and review. Numerous  
 410 factors could lead to the expectation that true OFLs would change between assessments  
 411 performed at different times, with the direction of change depending on the scenario (e.g., OFL  
 412 likely to increase through time for a rebuilding stock, or decrease for a newly-targeted stock  
 413 currently assessed to be well above its biomass target). The assumption that first-pass  
 414 assessments are median unbiased, or that a single distribution can describe the uncertainty  
 415 associated with each assessment, is also a gross oversimplification. Numerous factors could  
 416 affect the covariance between estimates from initial and revised assessments, and the  
 417 expectation that revised assessments would be median-unbiased is questionable. Requests for  
 418 model changes in revised assessments may have anticipated directional effects, or new  
 419 research projects may be funded with the anticipation of directional changes in assessment  
 420 outcomes (Terceiro 2018, Lynn et al. 2022). Nevertheless, I chose to model a scenario in which

421 redone assessments were median unbiased to illustrate the potential for inadvertent bias to be  
422 introduced even when requests for assessment revisions were not made with the intent or  
423 anticipation of driving the results in a particular direction. It may often be the case that proposed  
424 revisions to historical datasets, addition or removal of datasets, changes in data treatments,  
425 and/or revised prior specifications would have predictable effects. These sorts of predictable  
426 biases are not included in my simulations, and might be guarded against by restricting the  
427 opportunities for requesting such changes to early in the assessment process, before results  
428 are known.

429 In addition, stock assessments output numerous other quantities of scientific and  
430 management importance beyond the OFL, including estimates of status or depletion and  
431 estimates of biological parameters that affect productivity and so influence the projections  
432 needed for multi-year catch advice and for identifying sustainable fishing rates. The uncertainty  
433 in many of these estimates may also be reasonably described by a lognormal distribution (Bi et  
434 al. 2023) amenable to exploration with a similar approach, although there may be less of an  
435 established basis for the value of sigma to assume. I focused on OFLs and ABCs given the  
436 clear frequentist interpretation of  $P^*$  and a well-established existing framework for using this  
437 approach to characterize uncertainty. However, while incentives associated with p-hacking and  
438 publication of novel results is expected to lead to increased Type I error rates (i.e., false  
439 positives), it might be argued that a tendency to be suspicious of low-status assessments and  
440 favor model variants giving more moderate status could be more akin to increasing the rate of  
441 Type II errors (false negatives or incorrectly rejecting an assessment of poor status).

442 It is important to realize that full ABCs may not be attained, and thus an ABC higher than  
443 the true OFL does not necessarily mean that biological overfishing will or is likely to occur (i.e.,  
444 the fishing mortality rate actually estimated for recent years is often well below the proxy for the  
445 rate expected to produce maximum sustainable yield [e.g., Figure 3.4.1 of Harvey et al. (2022)]).  
446 Thus, even if the calculated  $P^*$  is lower than the true probability of an ABC exceeding the true  
447 OFL, biological overfishing may still be acceptably unlikely given expected attainment levels.  
448 This might reduce concern about the accuracy of  $P^*$  calculations, but this complication would be  
449 better addressed through a framework where  $P^*$  represents the probability that the expected  
450 harvest, as opposed to the ABC, would result in overfishing. Such an approach has not yet been  
451 developed.

452 Empirically quantifying the degree to which assessments are “redone” prematurely, and  
453 the extent to which low status predicts assessments receiving extra scrutiny, would be a  
454 formidable task requiring extensive review of the grey literature and likely a fair amount of

455 informed speculation about the motivation underlying incompletely documented decisions.  
456 Silvar-Viladomiu et al. (2021) and Bi et al. (2023) did not detect evidence of bias when  
457 comparing year-specific status estimates between repeated assessments of the same stocks,  
458 but each assessment was the product of a similar process that might exert similar effects on  
459 each iteration of the assessment. In addition, these analyses may not have sufficient power to  
460 detect small biases. Note also that my simulations assumed the re-assessments were  
461 themselves unbiased, and the potential bias arises when considering advice on the managed  
462 suite of stocks as whole rather than an expectation that the redone assessments are  
463 themselves biased. Some reviews of the scientific literature have used techniques like p-curves  
464 to test whether distributions of critical test statistics of published papers have discontinuities at  
465 critical “significance” levels that could be indicative of p-hacking or publication bias (Simonsohn  
466 et al. 2015). A similar approach might be used to examine the frequency of assessments  
467 indicating status just above versus below target or limit reference points, although one might  
468 expect that successful management would naturally lead to discontinuities around such  
469 reference points. Additionally, one could test whether assessments of status just above limits or  
470 targets in the terminal year of an earlier assessment tended to be consistent with updated  
471 perceptions of status for that year based on future assessments (Bi et al. 2023). Future  
472 simulation work could relax the assumption that initial or redone assessments are median  
473 unbiased, model the effects of directed requests for assessment revisions, and/or partition the  
474 assessment-specific error into additional components such as individual assessor effects,  
475 institutional effects, modeling platform effects, and the like. This might be addressed through a  
476 hierarchical modeling framework.

477

#### 478 *2.2.5 Likely magnitude of the problem*

479

480 In practice, reviewers and decision makers likely have some skill in identifying less  
481 accurate assessments, but do not have perfect knowledge, suggesting that the outcome in  
482 practice is likely to be somewhere in between the random selection and skillful selection  
483 scenarios. Note that the effects of moderate changes in the fraction of assessments redone or  
484 the proportion of variance attributable to assessment-specific factors can change the multiplier  
485 by amounts comparable to or larger than the changes needed to achieve 45% versus 40%  
486 probability of establishing an ABC higher than the true OFL. The bias in the OFL estimate  
487 introduced by asymmetric scrutiny is larger than can be countered by the default Q45 when as  
488 few as about 10% of assessments are redone if a high proportion of variance in OFL estimates

489 is attributable to assessment-specific factors, or about 20% of assessments if 50% of OFL  
490 variance is assessment-specific (Figure 2). If assessments are not redone within-cycle, but  
491 large relative decreases in OFLs between sequential assessments prompt extra scrutiny, this  
492 can introduce comparable levels of bias if the interval between assessments is short relative to  
493 the rate of change in true stock size/status. Thus, these problems may be more acute for long-  
494 lived, frequently assessed species and less acute for short-lived, infrequently assessed species.

495 Overall, these simulations suggest that randomly selecting assessments to be redone  
496 within a cycle is a waste of time and resources. Skillfully and symmetrically selecting  
497 assessments to be redone will not introduce bias, but is likely an inefficient use of resources,  
498 given the small changes in suitable multipliers. Redoing only the most pessimistic assessments  
499 within a cycle would introduce a bias that may be comparable in magnitude to the differences  
500 between choice of  $P^*=0.45$  versus 0.40. Similar biases could result when applying extra scrutiny  
501 to reductions in OFL or status between consecutive assessments, although the bias is reduced  
502 when changes in true stock size/status over the relevant interval are likely to be large, such as  
503 for infrequent assessment of a short-lived and dynamic stock. Differences on the order of a few  
504 percent may be acceptable relative to other uncertainties in the process, but the selection of  
505 assessments to be redone should be judicious to avoid the introduction of larger biases.

506

### 507 **3. Potential solutions**

508

509 Similar to several broad reviews of the scientific literature, the examples and models  
510 explored here suggest that various processes operating at the analysis and “publication” or  
511 adoption stage for stock assessment science are introducing a bias into fishery-wide OFL  
512 specification, and the rate at which ABCs exceed the OFLs that would have been established  
513 given perfect knowledge is likely higher than the nominal value of  $P^*$ , similar to how publications  
514 of “significant” results likely have a higher false positive rate than the nominal p-value. While the  
515 magnitude of the bias may not be sufficient to call the scientific or assessment enterprise as a  
516 whole into question, it does seem sufficient to warrant caution and efforts to limit known sources  
517 of bias as much as possible.

518 For the broader scientific enterprise, several courses of action have been proposed  
519 (Wicherts et al. 2016), most of which have clear analogs in the assessment process. 1)  
520 Analyses should be driven by clear *a priori* hypotheses that lead directly to a parsimonious set  
521 of candidate explanatory covariates, with objective methods for model and variable selection. 2)  
522 There should be transparency in statistical model selection and significance criteria. 3) Pre-

523 registration should be considered and employed to the extent possible. At the funding or pre-  
524 publication stage, clear statements should be made of the motivation for a study, the  
525 hypotheses it will test, the data to be collected and the analyses to be performed (ideally with a  
526 power analysis indicating sufficient power to detect meaningful effects if present), and the  
527 statistical tests to be performed along with the significance criterion and its justification. If all of  
528 these are satisfactory, publication should be assured regardless of the p-value obtained. 4)  
529 When multiple statistical tests or model formulations are applied to the same dataset, some  
530 adjustments like the Holm-Bonferroni procedure or Šidák correction should be applied. 5) A  
531 more stringent “significance” standard than  $p < 0.05$  should be considered for novel findings  
532 (Benjamin et al. 2018).

533         The stock assessment prioritization and review process along with the application of the  
534  $P^*/\sigma$  system for developing ABCs from OFLs offers analogies to all of these  
535 recommendations. 1) At the beginning of each assessment cycle, the species to be assessed,  
536 the spatial boundaries in assessment and management units, the data sources to be considered  
537 for inclusion, and the standards for review should all be specified in advance. 2) Review criteria  
538 should be clearly tied to the strength of scientific evidence, not management implications. 3)  
539 Assessments should not be aborted or rejected for use in management based on a politically  
540 unfavorable outcome. 4) The buffer between the ABC and OFL should be increased beyond  
541 that implied by the nominal choice of  $P^*$ , based on an approach similar to the models explored  
542 in Section 2.3. Similar adjustments to estimates of depletion and status may also be warranted.  
543 5) Strict standards should be adopted for revising, further reviewing, or rapidly revisiting an  
544 assessment that has been endorsed by scientific reviewers. For example, the SSC of the Mid-  
545 Atlantic Council will only reconsider a recommendation if new data are found or an error is  
546 discovered in an assessment (Crosson 2013), and the New England Fishery Management  
547 Council has similar limits on when an SSC recommendation can be remanded (Nies 2022).  
548 There may also be benefits in determining *a priori* criteria for how large a change would be  
549 required to deem a revised assessment sufficiently different from the initial assessment to revisit  
550 the adoption of the original assessment (SSC 2022), which might be based on evaluating the  
551 magnitude of the difference between two model alternatives relative to the overall level of  
552 uncertainty (Cope and Gertseva 2020).

553  
554  
555  
556

557 **4. Conclusions**

558

559 Overall, there seem to be numerous pathways by which inadvertent bias may be  
560 introduced into the stock assessment prioritization, review, and adoption process; and these  
561 pathways have commonalities in concerns raised about “p-hacking” in the scientific enterprise  
562 more broadly. Fortunately, the broader scientific literature also poses potential solutions or at  
563 least steps to reduce the influence of p-hacking, and many of these steps have direct analogs  
564 that can be applied to reduce the chances of introducing inadvertent bias into the fisheries stock  
565 assessment process. Simply raising awareness of the issue may go a long way toward fostering  
566 more careful work that is less likely to create bias (Peng 2015).

567

568 **Data Availability**

569

570 No original data were used in this paper. R code to run the simulations is available in a  
571 Mendeley archive at <https://data.mendeley.com/datasets/d49ct4fypr/2>.

572

573 **Funding**

574

575 This research did not receive any specific grant from funding agencies in the public,  
576 commercial, or not-for-profit sectors.

577

578 **Declaration of Competing Interest**

579

580 The scientific results and conclusions, as well as any views or opinions expressed herein, are  
581 those of the author and do not necessarily reflect the views of NOAA or the Department of  
582 Commerce. The author declares that he has no known competing financial interests or personal  
583 relationships that could have appeared to influence the work reported in this paper. The author  
584 is a member of the Pacific Fishery Management Council's Scientific and Statistical Committee at  
585 the time of writing, and his experiences there may have affected his perception of issues  
586 discussed in this paper, but he does not speak on their behalf.

587

588 **Acknowledgments**

589

590 This paper was informed and improved by feedback from members of the PFMC's SSC and  
591 stock assessment teams, attendees at a University of Washington Think Tank seminar, Eric  
592 Ward, Corey Ridings, John Field, Michael O'Farrell, Olaf Jensen, and an anonymous reviewer.

593

594 **References**

595

596 Baker, M., 2016. Is there a reproducibility crisis? Nature 533, 542-454.

597 <https://doi.org/10.1038/533452a>.

598 Benjamin, D.J., Berger, J.O., Johannesson, M.J., Nosek, B.A., Wagenmakers, E.-J., et al., 2018.  
599 Redefine statistical significance. *Nature Human Behaviour* 2, 6-10.  
600 <https://doi.org/10.1038/s41562-017-0189-z>.

601 Bi, R., Collier, C., Mann, R., Mills, K.E., Saba, V., Wiedenmann, J., Jensen, O.P., 2023. How  
602 consistent is the advice from stock assessments? Empirical estimates of inter-  
603 assessment bias and uncertainty for marine fish and invertebrate stocks. *Fish and*  
604 *Fisheries* 24, 126-141. <https://dx.doi.org/10.1111/faf.12714>.

605 Cadrin, S., Henderschedt, J., Mace, P., Mursalski, S., Powers, J., Punt, A.E., Restrepo, V., 2015.  
606 Addressing Uncertainty in Fisheries Science and Management. National Aquarium.  
607 <http://www.fao.org/3/a-bf336e.pdf>.

608 Clark, T.D., Raby, G.D., Roche, D.G., Binning, S.A., Speers-Roesch, B., Jutfelt, F., Sundin, J. 2020.  
609 Ocean acidification does not impair the behaviour of coral reef fishes. *Nature* 577, 370-  
610 375. <https://doi.org/10.1038/s41586-019-1903-y>.

611 Cope, J.M., DeVore, J., Dick, E.J., Ames, K., Budrick, J., Erickson, D.L., et al., 2011. An approach to  
612 defining stock complexes for U.S. West Coast groundfishes using vulnerabilities and  
613 ecological distributions. *N. Am. J. Fish. Manag.* 31, 589–604.  
614 <https://dx.doi.org/10.1080/02755947.2011.591264>.

615 Cope, J.M., Gertseva, V., 2020. A new way to visualize and report structural and data  
616 uncertainty in stock assessments. *Can. J. Fish. Aquat. Sci.* 77, 1275-280.  
617 <https://doi.org/10.1139/cjfas-2020-0082>.

618 Crosson, S., 2013. The impact of empowering scientific advisory committees to constrain catch  
619 limits in US fisheries. *Science and Public Policy* 40, 261-273.  
620 <https://doi.org/10.1093/scipol/scs104>.

621 Fanelli, D., 2018. Is science really facing a reproducibility crisis, and do we need it to? *Proc. Nat.*  
622 *Acad. Sci.* 115, 2628-2631. <https://doi.org/10.1073/pnas.1708272114>.

623 Harvey, C., et al. 2022. 2021-2022 California Current Integrated Ecosystem Assessment (CCIEA)  
624 California Current ecosystem status report. Report to the Pacific Fishery Management  
625 Council. [https://www.pcouncil.org/documents/2022/02/h-2-a-cciea-team-report-1-  
626 2021-2022-california-current-ecosystem-status-report-and-appendices.pdf/](https://www.pcouncil.org/documents/2022/02/h-2-a-cciea-team-report-1-2021-2022-california-current-ecosystem-status-report-and-appendices.pdf/).

627 Hamel, O.S., 2014. A method for calculating a meta-analytical prior for the natural mortality  
628 rate using multiple life history correlates. *ICES J. Mar. Sci.* 72, 62-69.  
629 <https://doi.org/10.1093/icesjms/fsu131>.

630 Head, M.L., Holman, L., Lanfear, R., Kahn, A.T., Jennions, M.D., 2015. The Extent and  
631 Consequences of P-Hacking in Science,. *PLoS Biol.* 13, e1002106.  
632 <https://dx.doi.org/doi:10.1371/journal.pbio.1002106>.

633 Hilborn, R., Amoroso, R.O., Anderson, C.M., Baum, J.K., Branch, T.A., Costello, C., De Moor, C.L.,  
634 Faraj, A., Hively, D., Jensen, O.P., Kurota, H., 2020. Effective fisheries management  
635 instrumental in improving fish stock status. *Proc. Nat. Acad. Sci.* 117, 2218-2224.

636 Hilborn, R., Walters, C., 1992. *Quantitative Fisheries Stock Assessment: Choice, Dynamics and*  
637 *Uncertainty*. Chapman and Hall.

638 Ioannidis, J.P.A., 2005. Why most published research findings are false. *PLoS Medicine* 2, e124.  
639 <https://dx.doi.org/10.1371/journal.pmed.0020124>.

640 Jager, L.R., Leek, J.T., 2014. An estimate of the science-wise false discovery rate and application  
641 to the top medical literature. *Biostatistics* 15, 1-12.  
642 <https://dx.doi.org/10.1093/biostatistics/kxt007>.

643 Lynn, K., Dorval E., Porzio, D., Nguyen, T., 2022. A collaborative survey of coastal pelagic species  
644 in nearshore California waters. *Fisheries* 47, 500-508.  
645 <https://doi.org/10.1002/fsh.10840>.

646 Maunder, M.N., Piner, K.R., 2015. Contemporary fisheries stock assessment: many issues still  
647 remain. *ICES J. M. Sci.* 72, 7-18. <https://doi.org/10.1093/icesjms/fsu015>.

648 Maunder, M.N., Piner, K.R., 2017. Dealing with data conflicts in statistical inference of  
649 population assessment models that integrate information from multiple diverse data  
650 sets. *Fish. Res.* 192, 16-27. <https://doi.org/10.1016/j.fishres.2016.04.022>.

651 Melnychuk, M.C., Banobi, J.A., Hilborn, R., 2013. Effects of management tactics on meeting  
652 conservation objectives for western North American groundfish fisheries. *PLoS One* 8,  
653 e56684. <https://doi.org/10.1371/journal.pone.0056684>.

654 Methot Jr., R.D. (editor). 2015. Prioritizing fish stock assessments. U.S. Dep. Commer., NOAA  
655 Tech. Memo. NMFS-F/SPO-152. <https://repository.library.noaa.gov/view/noaa/12874>.

656 Methot Jr, R.D., Tromble, G.R., Lambert, D.M. Greene, K.E., 2014. Implementing a science-  
657 based system for preventing overfishing and guiding sustainable fisheries in the United  
658 States. *ICES J. Mar. Sci.* 71, 183-194. <https://doi.org/10.1093/icesjms/fst119>.

659 Mildenberger, T., Berg, C.W., Kokkalis, A., Hordyk, A.R., Wetzel, C., Jacobsen, N.S., Punt, A.E.,  
660 Nielsen, J.R., 2022. Implementing the precautionary approach into fisheries  
661 management: Biomass reference points and uncertainty buffers. *Fish Fish.* 23, 73-92.  
662 <https://doi.org/10.1111/faf.12599>.

663 NMFS (National Marine Fisheries Service), 2013. Use of stock status results from data-moderate  
664 assessments. Presentation to Pacific Fishery Management Council.  
665 [https://www.pcouncil.org/documents/2013/03/h-groundfish-management-march-  
666 2013.pdf/#page=40](https://www.pcouncil.org/documents/2013/03/h-groundfish-management-march-2013.pdf/#page=40).

667 NMFS (National Marine Fisheries Service), 2021. 2021 Report to Congress on the Regional  
668 Fishery Management Councils and Scientific and Statistical Committee members'  
669 financial interest disclosure and recusal requirements and on the Regional Fishery  
670 Management Councils membership apportionment developed pursuant to section  
671 302(b)(2)(b) and section 302(j)(9) of the Magnuson-Stevens Fishery Conservation and  
672 Management Act. [https://media.fisheries.noaa.gov/2022-03/Final-  
673 2021%20Report%20to%20Congress%20on%20Councils.pdf](https://media.fisheries.noaa.gov/2022-03/Final-2021%20Report%20to%20Congress%20on%20Councils.pdf).

674 NMFS (National Marine Fisheries Service), 2022. Prioritizing groundfish stock assessments:  
675 Preliminary analysis and ranking of Pacific coast groundfish species for assessment in  
676 2023. Presentation to Pacific Fishery Management Council.  
677 [https://www.pcouncil.org/documents/2022/03/e-8-a-supplemental-noaa-fsc-  
678 presentation-1-preliminary-analysis-and-ranking-of-pacific-coast-groundfish-species-for-  
679 assessment-in-2023-dr-hastie-and-dr-wetzel.pdf/](https://www.pcouncil.org/documents/2022/03/e-8-a-supplemental-noaa-fsc-presentation-1-preliminary-analysis-and-ranking-of-pacific-coast-groundfish-species-for-assessment-in-2023-dr-hastie-and-dr-wetzel.pdf/)

680 Nies, T. 2022. SSC Roles and Responsibilities. SSC Planning Meeting, presentation of the NEFMC  
681 Executive Director. [https://s3.us-east-1.amazonaws.com/nefmc.org/SSC-roles-and-  
682 responsibilities\\_for-web.pdf](https://s3.us-east-1.amazonaws.com/nefmc.org/SSC-roles-and-responsibilities_for-web.pdf).



683 Peng, R., 2015. The reproducibility crisis in science: A statistical counterattack. Significance 6,  
684 30-32. <https://doi.org/10.1111/j.1740-9713.2015.00827.x>.

685 PFMC (Pacific Fishery Management Council), 2013. Status determination criteria for data-  
686 moderate stocks. [https://www.pcouncil.org/documents/2013/03/h-groundfish-  
687 management-march-2013.pdf/#page=27](https://www.pcouncil.org/documents/2013/03/h-groundfish-management-march-2013.pdf/#page=27).

688 PFMC (Pacific Fishery Management Council), 2022. Terms of reference for the groundfish stock  
689 assessment review process for 2023-2024.  
690 [https://www.pcouncil.org/documents/2022/06/terms-of-reference-for-the-groundfish-  
691 stock-assessment-review-process-for-2023-2024-june-2022.pdf/](https://www.pcouncil.org/documents/2022/06/terms-of-reference-for-the-groundfish-stock-assessment-review-process-for-2023-2024-june-2022.pdf/).

692 Privitera-Johnson, K.M., Punt, A.E., 2020a. A review of approaches to quantifying uncertainty in  
693 fisheries stock assessments. Fish. Res. 226, 105503.  
694 <https://doi.org/10.1016/j.fishres.2020.105503>.

695 Privitera-Johnson, K.M., Punt, A.E., 2020b. Leveraging scientific uncertainty in fisheries  
696 management for estimating among-assessment variation in overfishing limits. ICES J.  
697 Mar. Sci. 77, 515-526. <https://doi.org/10.1093/icesjms/fsz237>.

698 Provencher, J.F., Covernton, G.A., Moore, R.C., Horn, D.A., Conkle, J.L., Lusher, A.L. 2020.  
699 Proceed with caution: The need to raise the publication bar for microplastics research.  
700 Sci. Tot. Env. 748, 141426. <https://doi.org/10.1016/j.scitotenv.2020.141426>.

701 Ralston, S., Punt, A.E., Hamel, O.S., Devore, J.D., Conser, R.J., 2011. A meta-analytic approach to  
702 quantifying scientific uncertainty in stock assessments. Fish. Bull. 109, 217–231.  
703 [https://spo.nmfs.noaa.gov/content/meta-analytic-approach-quantifying-scientific-  
704 uncertainty-stock-assessments](https://spo.nmfs.noaa.gov/content/meta-analytic-approach-quantifying-scientific-uncertainty-stock-assessments).

705 Rudd, M.B., Cope, J.M., Wetzel, C.R., Hastie, J., 2021. Catch and length models in the stock  
706 synthesis framework: expanded application to data-moderate stocks. Front. Mar. Sci. 8,  
707 663554. <https://doi.org/10.3389/fmars.2021.663554>.

708 Schooler, J., 2011. Unpublished results hide the decline effect. Nature 470, 437.  
709 <https://doi.org/10.1038/470437a>.

710 Seagraves, R., Collins, K. (eds), 2012. Fourth National Meeting of the Regional Fishery  
711 Management Council’s Scientific and Statistical Committees. Report of a National SSC  
712 Workshop on Scientific Advice on Ecosystem and Social Science Considerations in U.S.  
713 Federal Fishery Management. Mid-Atlantic Fishery Management Council, Williamsburg,  
714 VA. <http://www.fisherycouncils.org/ssc-workshops/fourth-national-ssc-workshop-2011>.

715 Shertzer, K.W., Prager, M.H., Williams, E.H., 2008. A probability-based approach to setting  
716 annual catch levels. Fish. Bull. 106, 225–232. [https://media.fisheries.noaa.gov/dam-  
717 migration/ns1-shertzer-et-al-2008.pdf](https://media.fisheries.noaa.gov/dam-migration/ns1-shertzer-et-al-2008.pdf).

718 Silvar-Viladomiu, P., Minto, C., Halouani, G., Betts, L., Brophy, D., Lordan, C., Reid, D.G., 2021.  
719 Moving reference point goalposts and implications for fisheries sustainability. Fish Fish.  
720 22, 1345-1358. <https://doi.org/10.1111/faf.12591>.

721 Simmons, J.P., Nelson L.D., Simonsohn, U., 2011. False-positive psychology: undisclosed  
722 flexibility in data collection and analysis allows presenting anything as significant. Psych.  
723 Sci. 22, 1359-1366. <https://dx.doi.org/10.1177/0956797611417632>.

724 Simonsohn, U., Simmons, J. P., Nelson, L. D., (2015). Better P-curves: Making P-curve analysis  
725 more robust to errors, fraud, and ambitious P-hacking, a Reply to Ulrich and Miller

726 (2015). *J. Experimental Psychology: General* 144, 1146–1152.  
727 <https://doi.org/10.1037/xge0000104>.

728 SSC (Scientific and Statistical Committee of the Pacific Fishery Management Council), 2021.  
729 Scientific and Statistical Committee Report on Harvest Specifications for 2023-2024  
730 Including Final Overfishing Limits and Acceptable Biological Catches.  
731 <https://www.pcouncil.org/documents/2021/11/e-3-a-supplemental-ssc-report-1-2.pdf/>.

732 SSC (Scientific and Statistical Committee of the Pacific Fishery Management Council), 2022. SSC  
733 Groundfish Subcommittee Report on Groundfish Stock Assessment Process Review  
734 Webinar held on January 25, 2022. <https://www.pcouncil.org/documents/2022/02/e-8-attachment-6-ssc-groundfish-subcommittee-report-on-groundfish-stock-assessment-process-review-webinar-held-on-january-25-2022.pdf/>.

737 Terceiro, M., 2018. The summer flounder chronicles III: struggling with success, 2011-2016. *Rev. Fish Biol. Fisheries* 28, 381-404. <https://doi.org/10.1007/s11160-017-9506-x>.

739 Thorson, J.T., Dorn, M.W., Hamel, O.S., 2019. Steepness for West Coast rockfishes: results from  
740 a twelve-year experiment in iterative regional meta-analysis. *Fish. Res.* 217, 11-20.  
741 <https://doi.org/10.1016/j.fishres.2018.03.014>.

742 Viglione, G., 2020. 'Avalanche' of spider-paper retractions shakes behavioural-ecology  
743 community. *Nature* 578, 199-200. [https://www.nature.com/articles/d41586-020-00287-](https://www.nature.com/articles/d41586-020-00287-y)  
744 [y](https://www.nature.com/articles/d41586-020-00287-y).

745 Wasserstein, R.I., Lazar, N. A., 2016. The ASA statement on p-values: Context, process, and  
746 purpose. *Am. Stat.* 70, 129-133. <https://dx.doi.org/10.1080/00031305.2016.1154108>.

747 Wicherts, J.M., Veldkamp, C.L.S., Augusteijn, H.E.M., Bakker, M., van Aert, R.C.M. and van  
748 Assen, M.A.L.M., 2016. Degrees of freedom in planning, running, analyzing, and  
749 reporting psychological studies: A checklist to avoid p-Hacking. *Front. Psychol.* 7,1832.  
750 <https://dx.doi.org/10.3389/fpsyg.2016.01832>.