

1 **The quality of open datasets shared by researchers in ecology and evolution is moderately**
2 **repeatable and slow to change**

3
4 Dominique G Roche^{1,2,3}, Ilias Berberi¹, Fares Dhane², Félix Lauzon^{2,4}, Sandrine Soeharjono²,
5 Roslyn Dakin¹, Sandra A Binning²

6
7 ¹ *Department of Biology, Carleton University, Ottawa, Canada*

8 ² *Département de science biologiques, Université de Montréal, Montréal, Canada*

9 ³ *Institut de Biologie, Université de Neuchâtel, Neuchâtel, Switzerland*

10 ⁴ *Department of Biology, McGill University, Montréal, Canada*

11
12 Correspondence to: dominique.roche@mail.mcgill.ca

13
14 **Key words:** data sharing, FAIR data, metascience, open science, public data archiving,
15 reproducibility

16
17 **We assessed the quality of 362 open datasets shared by 100 principal investigators (PIs) in**
18 **ecology and evolution to identify predictors of data quality. Datasets generally scored low**
19 **on completeness and reusability, but these metrics were slightly higher for more recently**
20 **archived datasets and PIs with less seniority. Journal data sharing policies had no effect on**
21 **data quality, whereas PI identity explained the largest proportion of the variance in both**
22 **data completeness (27.8%) and reusability (22.0%), suggesting that a PI's training and lab**
23 **culture are key determinants of data quality. Thus, greater incentives and training for**
24 **individual researchers could help improve data sharing practices.**

25
26 The debate is over regarding the value of open and FAIR data (Findable, Accessible,
27 Interoperable, and Reusable data¹). Many journals, funding agencies, and policymakers agree
28 that the societal benefits of publicly sharing (non-sensitive) research data outweigh any
29 perceived or reported costs to individual researchers²⁻⁶. Indeed, while some researchers remain
30 reluctant to share the data underlying their published results^{3,7}, most view open data
31 positively^{2,8,9}. Since 2010, when a handful of journals began requiring open data in ecology and
32 evolution (E&E)¹⁰, policies encouraging this practice have grown rapidly. Now, over 60 journals
33 publishing research in E&E mandate open data². Strong journal policies are effective at ensuring
34 that more datasets are shared^{11,12}, which is often touted as a win for open science¹³. Yet,
35 problems persist^{12,14}. For instance, more than half of open datasets associated with 100 E&E
36 studies published in 2012 and 2013 were incomplete and/or archived in ways that prevented
37 reuse¹⁵. Similar issues have been documented in psychology¹⁶ and cognition research¹⁷, pointing
38 to the inherent problem with journals mandating open data without appropriate oversight or
39 quality control¹⁸⁻²⁰: datasets get archived but the majority are incomplete and challenging to
40 reuse. Developing effective strategies to promote good data sharing practices requires that we
41 first identify which factors predict high-quality, FAIR data.

42
43 We assessed the quality (completeness and reusability) of open datasets associated with
44 publications by tenured or tenure-track E&E faculty members (PIs) in biology departments at the
45 21 highest-ranked universities in Canada. PIs were necessarily first or last author on the
46 publications assessed (see pre-registered methods). Data completeness (availability of data

47 allowing computational reproducibility) and reusability (ease with which data can be reused by
48 third parties) were assessed following Roche et al 2015¹⁵. Scores above 3 on two 5-point scales
49 indicate complete or reusable data (Table S1). In total, we examined 362 datasets shared by 44
50 women and 56 men (Table S2). We ran a Bayesian bivariate mixed-effects linear model with PI
51 identity and their institution of employment as random factors. We tested whether article
52 publication date, journal open data policy, and the seniority, gender and corresponding author
53 status of PIs predicted the quality of their open datasets.

54
55 The completeness and reusability scores of datasets varied considerably within and among PIs
56 (Fig. 1). Overall, 56.4% of datasets were complete (mean completeness score of 3.4 ± 1.3 SD),
57 and 45.9% were reusable (mean reusability score of 3.1 ± 1.4 SD) (Fig. S1) This represents a
58 moderate improvement of approximately 10% above the completeness and reusability of datasets
59 associated with E&E studies published in 2012 and 2013¹⁵. Data completeness and reusability
60 were strongly correlated within ($R^2=0.79$, 95% CI: 0.72–0.85) and among ($R^2=0.77$, 95% CI:
61 0.56–0.92) PIs.

62
63 Open data is a relatively recent concept in E&E, having been introduced in earnest a decade
64 ago⁶. As such, PIs who developed their research skills prior to this period may be less likely to
65 have incorporated these principles into their research workflow. We assessed datasets as far back
66 as 2013 in our analysis but included faculty members hired as recently as 2019, using year of
67 first scientific publication as a proxy for PI seniority. Therefore, our study likely includes
68 datasets published by new PIs during their PhD and post-doc years, when they might have had
69 access to various training opportunities. For instance, a growing number of biology departments
70 recognize the value of data science and initiatives such as Data Carpentry (datacarpentry.org)
71 and FOSTER (fosteropenscience.eu) now routinely offer workshops in data management across
72 North America and Europe. We found that PIs with less seniority tended to share slightly more
73 reusable data than PIs with more seniority, suggesting that early training initiatives may be
74 bearing fruit (Figs 2, S2 I). This result is good news because younger researchers tend to be more
75 fearful and reluctant to share their data than senior researchers⁸ despite reporting a more
76 favorable attitude towards open data^{8,21-23} (but see²⁴). Our study included datasets spanning six
77 years (2013-2019). We found that datasets associated with more recent studies were slightly
78 more reusable than those of older studies (Figs 2, S2 J). In contrast, dataset completeness was
79 independent of PI seniority and publication date (Fig. 2, S2 D,F). These results suggest that,
80 while there have been slight improvements to data sharing practices through time, these are slow
81 to change. Training and increased exposure to open science practices are no doubt contributing
82 to this slow improvement, but additional work is needed at all career stages to enhance the
83 quality of open data.

84
85 We examined whether the quality of open datasets was influenced by PIs being corresponding
86 author on the published article (in addition to being first or senior author). We took PI
87 corresponding author status as an indicator that the PIs themselves archived the data. Assuming
88 that PIs are highly competent at managing their research data, we expected that open datasets
89 with PIs as corresponding author to be of higher in quality, on average, than those archived by
90 presumably less experienced researchers (likely students or post-docs). We found no support for
91 this hypothesis: data quality was unrelated to corresponding author status (Figs 2, S2 C,D).
92 Conventions regarding who is corresponding author on a published study may vary among sub-

93 disciplines and research labs. However, the corresponding author is ultimately responsible for
94 compliance with journal policies, including open data²⁵. The fact that corresponding author status
95 had no bearing on data quality is worrying and suggests that PIs do not understand the
96 responsibilities associated with this role, do not have the tools or training to ensure that open data
97 are compliant with journal policies, or simply do not care. Education, capacity-building, and
98 incentives targeted at individuals are needed to address these issues^{2,26}.

99
100 We had no a priori hypothesis for why PI gender might influence data quality. However, we
101 included this predictor in the model because we believe that gender differences are important to
102 consider in scientific research. We found no evidence to suggest that PI gender influences the
103 quality of open data (Figs 2, S2 A,B). This result is encouraging given that men in E&E self-
104 identify as experiencing more costs than women as a result of sharing open data². Women
105 accounted for almost half of the PIs assessed in our study, yet far fewer than 50% of PIs in
106 biology departments at Canadian institutions identify as women²⁷. While gender does not appear
107 to be a barrier to open data quality, we nevertheless encourage diversity incentives aimed at
108 promoting training and support in open science practices for minority gender identities as a way
109 of enhancing the quality, visibility, and impact of research published by underrepresented groups
110 in E&E.

111
112 When journals have a mandatory open data policy, the number of archived datasets underlying
113 published research articles increases¹¹. We tested whether such policies also translate into more
114 complete and reusable data. We hypothesized that data quality would be higher for studies
115 published in journals requiring open data. Alternatively, it is also possible that researchers
116 voluntarily archiving datasets in journals without a policy share higher quality data than
117 researchers who are forced to do so. Contrary to this logic, a journal's open data policy had no
118 bearing on data quality (Fig. 2, S2 B,G), indicating that policies alone do little to ensure that
119 shared data are complete and reusable. Some of the world's largest funding agencies (e.g., ERC,
120 NSF, NERC, Canadian Tri-Council) now require that PIs specify data management and/or
121 sharing plans to obtain funding. However, compliance with these policies is low¹⁴. Unless
122 resources are invested in training, technical support and policy oversight², data risk not being
123 archived or not contributing to advancing knowledge in instances where they are made available.

124
125 We assessed multiple datasets by the same researchers, which allowed us to calculate
126 repeatability scores for both data completeness and reusability. Repeatability (R) ranges between
127 0 and 1 and is the proportion of the total variance in scores attributed to inter-individual
128 differences: high R values indicate large score differences among individuals and consistent
129 scores within individuals. Data completeness was moderately repeatable with $R_{adj} = 0.28$ (95%
130 CI: 0.16–0.39), and data reuse with $R_{adj} = 0.22$ (95% CI: 0.14–0.34), revealing differences within
131 and among PIs (Fig. 1). This variability reflects several realities that merit discussion. On the one
132 hand, PIs in academia often conduct research through students and post-docs, who may
133 ultimately be responsible for data collection, curation and archiving. Thus, in some cases, the PI
134 may not have performed quality control over the archived files, potentially explaining the
135 considerable within-individual variation in data completeness and reusability scores we observed
136 (variable scores in Fig. 1). In these cases, data quality may be a better reflection of data
137 management not by the PI, but by the person within the PI's research group who was responsible
138 for archiving the dataset. On the other hand, some PIs consistently scored high or low in both

139 data completeness and reuse (low and high scores in Fig. 1). This consistency within research
140 groups suggests both a robust lab culture promoting good research data management or,
141 alternatively, a PI's lack of competence or reluctance to engage in data sharing and student
142 training in this regard. PIs who oppose open data initiatives^{2,7} are unlikely to respond positively
143 to incentives or training opportunities to improve data quality and FAIRness. However, our
144 results suggest that only a minority of PIs potentially fall within this category (approximately
145 10%; Fig. 1). Rather, most PIs were inconsistent in how they shared data associated with their
146 publications, or consistently shared highly complete and reusable data. This finding points to the
147 importance of facilitating sound research data management practices within research groups to
148 achieve complete and FAIR data sharing.

149
150 Overall, our data suggest that journal policies are ineffective at ensuring that open data in E&E
151 are complete and reusable. We also found that data quality is slow to improve over time.
152 However, most PIs did share high-quality data, either consistently or occasionally. Striking
153 variation in the quality of open data within PIs suggests that education, training, and technical
154 support could help raise the bar by enabling good data sharing practices to become the rule rather
155 than the exception²⁶.

156 157 **Methods**

158
159 Our methods were pre-registered at [https://doi.org/ 10.17605/OSF.IO/A492M](https://doi.org/10.17605/OSF.IO/A492M). We assessed open
160 datasets from research faculty members in biology departments at the 21 highest-ranked
161 Canadian universities based on the 2019 Times Higher Education World University Rankings.
162 Although we initially planned to select the top 20 Canadian universities, we retained 21
163 universities due to a three-way tie for rank 19. Our study focused on Canadian academic faculty.
164 However, our findings are likely to be representative of the broader population of PIs in ecology
165 and evolution given Canada's diverse academic institutions as well as the high degree of PI
166 mobility in today's globalized academic landscape. Furthermore, many granting agencies in
167 Europe and the USA require that data from funded research be publicly archived within a
168 specified timeframe of publishing. This is not yet the case in Canada: the Tri-council Granting
169 agencies now require a data management plan for grants submitted in 2021 and beyond, but this
170 does not yet include a requirement for open data. This allowed us to assess the effect of journal
171 policies on archiving practices by Canadian PI's independent of requirements from funding
172 agencies on the same practices.

173
174 We reviewed the biology department website at each of the 21 selected universities in a random
175 order and identified all researchers primarily conducting research in the fields of ecology and/or
176 evolution (E&E) with a rank of assistant, associate or full professor. Adjunct professors and
177 researchers who primarily focus on molecular biology, genetics, genomics, bioinformatics,
178 theoretical biology, comparative physiology and paleontology were excluded given our focus on
179 researchers in E&E. Each researcher's primary fields of study were determined from public
180 information on the university websites and cross-checked by a minimum of two people (IB, FD,
181 SAB, RD, DGR). This criterion served to limit the scope of the study to E&E and facilitate
182 consistent assessment of datasets, given our shared expertise. To standardize the selection of
183 researchers across universities and avoid bias, we omitted E&E researchers who are primarily
184 affiliated with departments other than biology (e.g., environmental sciences, natural resources,

185 fisheries and ocean sciences, veterinary sciences). In total, we identified 351 researchers that met
186 these criteria Table S2).

187
188 To be included in our study, candidate researchers had to have at least two articles containing a
189 data availability statement that were published in a peer-reviewed scientific journal between
190 January 2013 and June 2019. The researcher had to be first or senior (last) author on these
191 articles, ensuring that they were one of the primary intellectual contributors (in E&E, the
192 convention is that the first and last author are primarily responsible for the work). We used
193 Google Scholar and the researchers' personal and/or institutional websites to identify articles.
194 When researchers did not have a Google Scholar profile, we verified their publication list using
195 Web of Science. Researchers at each university were screened in a random order. Articles for
196 each researcher were manually searched in a reverse chronological order (i.e., starting in 2019,
197 ending in 2013) to determine whether a data availability statement was present, either stated
198 explicitly at the end of the article, or embedded in the main text. If there was an absence of an
199 open data statement but presence of electronic supplementary material (ESM), we looked for
200 evidence of open data in the ESM (i.e., raw or processed data as opposed to summary statistics).
201 Reviews, commentaries, and theoretical or simulation studies were excluded. The article search
202 for every researcher was completed when seven articles containing a data availability statement
203 or open data were identified, or when the reverse-chronological scan reached January 2013. In
204 total, 4,322 articles were examined, 928 of which contained a data availability statement and/or
205 associated open data.

206
207 The strength of a journal's open data policy and date of implementation was determined by
208 reviewing each journal's author guidelines and relevant editorials. When necessary, we contacted
209 journal editors for clear information on whether open data were required (i.e., mandatory open
210 data) or encouraged (i.e., optional open data) as a condition of publication at the time a paper
211 was published. Journals without an open data policy were categorized as optional open data.

212
213 We identified 19 researchers with at least two articles containing a data availability statement
214 and/or open data across the 21 universities (Table S2). Of these, we randomly selected up to four
215 women and four men at each university to evaluate the quality of their shared datasets. We aimed
216 to randomly select three women and three men at each university but some universities had
217 fewer than three researchers per gender (Table S2). The departments of biology at two
218 institutions had no researchers that met our selection criteria. One researcher identified as gender
219 non-binary but was not part of our random sample. We made assumptions about gender based on
220 names and pronouns used in public profiles on university websites or social media. We recognize
221 that gender presentation, names, and pronouns are not necessarily indications of a person's
222 gender and that, in the absence of additional information from the individuals, we may have
223 unintentionally made incorrect assumptions about individuals' genders. In total we assessed 362
224 datasets published in 97 journals by 100 PIs. We scored the completeness and reusability of
225 shared datasets on a scale from 0 (min score) to 5 (max score) following Roche et al. (2015)¹⁵
226 (Table S1). The number of datasets assessed per researcher ranged from two to five; if a
227 researcher had more than five shared datasets in the period from 2013-2019, we selected the
228 most recent five.

229
230 *Statistical analysis*

231
232 We used a Bayesian bivariate mixed-effects regression model (R package MCMCglmm v2.32²⁸)
233 to identify factors influencing the quality of shared data and estimate its repeatability (i.e., the
234 proportion of the total variance attributable to differences among individuals)²⁹. Data
235 completeness and reusability scores were included as two dependent variables in the model;
236 researcher ID and university were specified as random effects, with researcher nested within
237 university; PI gender, PI seniority (measured as the year of their first peer-reviewed publication,
238 assessed on Google Scholar or Web of Science), PI author status (corresponding author or not),
239 journal open data policy at the time of publication (mandatory, optional), and year of study
240 publication were included as fixed effects.

241
242 The two dependent variables were mean-centered and standardized to one standard deviation unit
243 prior to inclusion in the model (i.e., mean=0, standard deviation=1). The numerical (PI seniority,
244 year of study publication) and categorical predictors (gender, corresponding author, journal open
245 data policy) were mean-centered and standardized to two standard deviation units (i.e., mean=0,
246 standard deviation=0.5) following Araya-Ajoy *et al*³⁰. Categorical predictors were treated as
247 binary variables (values of 0 and 1) to allow centering and standardization. The advantage of
248 mean-centering the predictors is that it ensures model intercepts are estimated for the average
249 value of the predictors, facilitating interpretation of the results. Mean centering allows the
250 estimate of the intercepts to be calculated for the average ‘environmental’ conditions³¹; the use of
251 two standard deviations for predictor standardization allows for direct comparison of the
252 variance explained by categorical and continuous predictors³⁰.

253
254 We specified a mildly informative inverse-Wishart prior and tested the sensitivity of the model to
255 prior specification by examining how the posterior means and 95% credible intervals changed
256 when specifying a parameter-expanded prior (see ²⁸). We checked the model by plotting the
257 traces of the parameters, examining autocorrelation among samples drawn by MCMCglmm, and
258 computing the Gelman-Rubin statistic to evaluate convergence (see archived script). Model
259 diagnostics were satisfactory and conclusions were not sensitive to the choice of prior (Fig. S2).

260
261 We calculated the adjusted repeatability (R_{adj}) for a researcher’s data completeness and data
262 reusability as the proportion of the total variance due to differences among individuals when
263 accounting for fixed and random effects in the statistical model ³². Within- and among-individual
264 correlations between data completeness and reusability were calculated as outlined in Roche *et*
265 *al*²⁹. All analyses were done in R version 4.0.3.

266
267
268

269 **Acknowledgements**

270

271 We thank Redouan Bshary for the idea that led to this study and Ellen Bledsoe for helpful
272 comments on the manuscript. We acknowledge funding from the Natural Sciences and
273 Engineering Research Council of Canada (grant no. UIF-537860–2018). DGR was supported by
274 the European Union’s Horizon 2020 research and innovation program under Marie Skłodowska-
275 Curie grant agreement no. 838237-OPTIMISE.

276

277 **Author Contributions**

278

279 DGR, RD, and SAB designed the study. IB, FD, FL, SS, and DGR collected and managed the
280 data. DGR analyzed the data. DGR and FD made the figures. DGR and SAB wrote the paper. All
281 authors revised and approved the paper.

282

283 **Competing Interests statement**

284

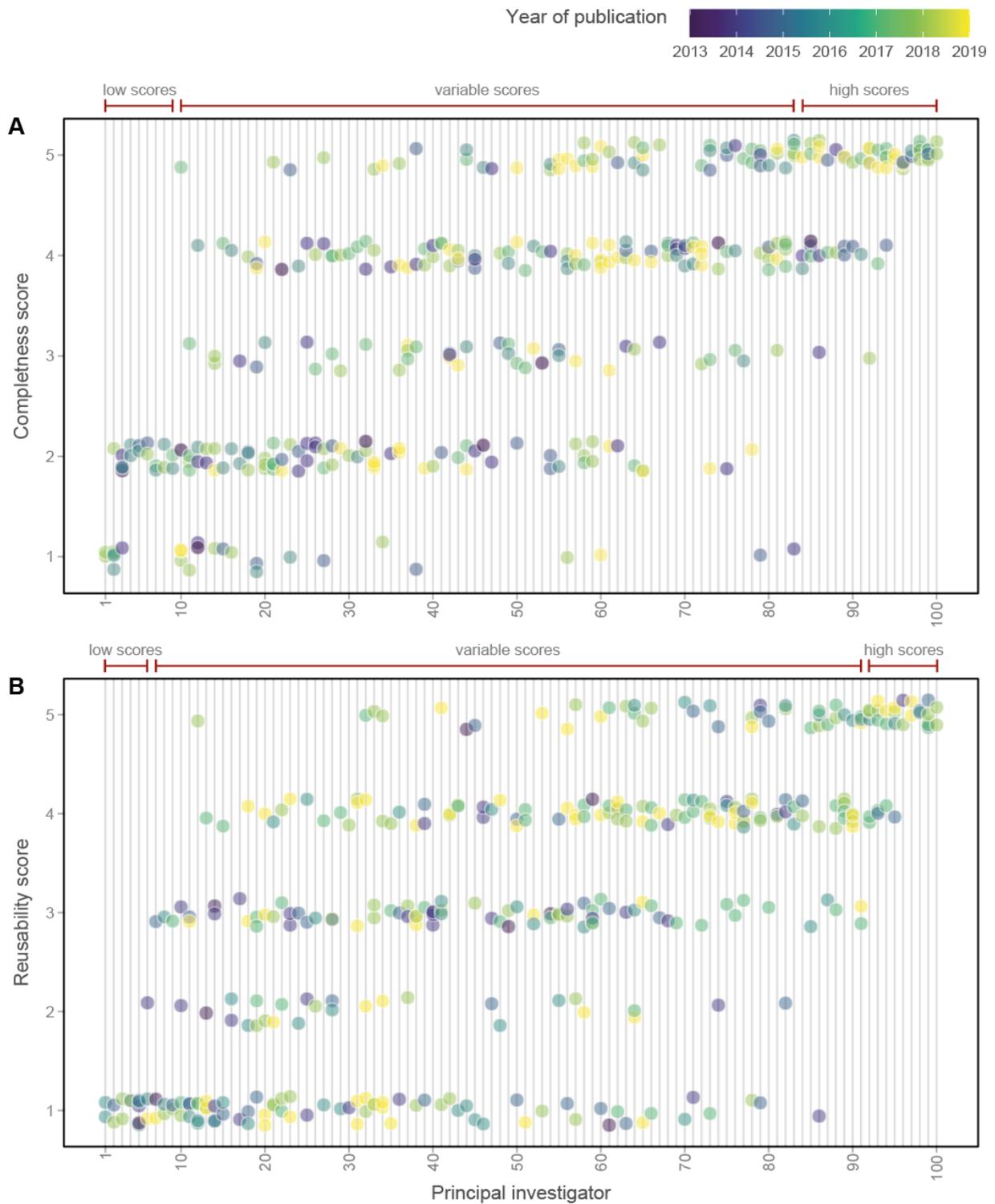
285 The authors declare no competing interests.

286

287 **Data and code availability**

288

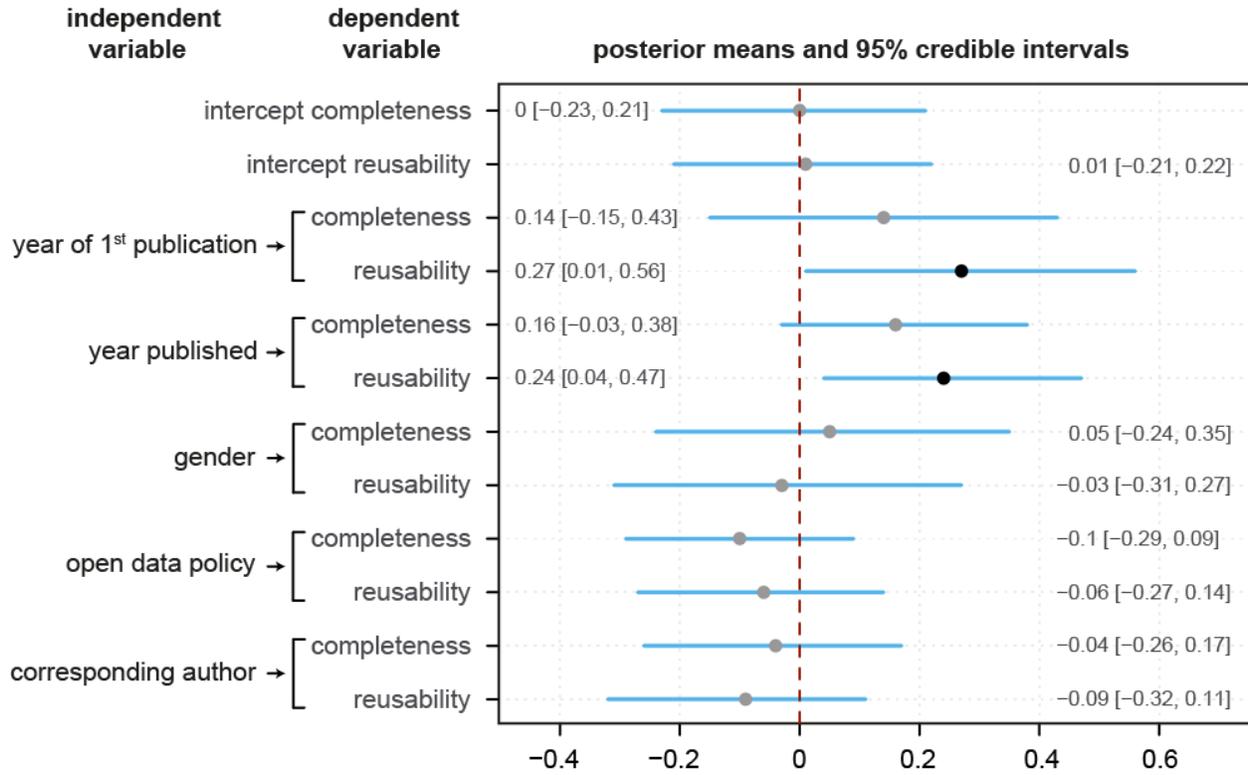
289 The anonymized data and the analysis script for this study are available on the Open Science
290 Framework (<https://doi.org/10.17605/OSF.IO/P2YWM>) and were shared with the
291 editors/reviewers upon submission. This study was pre-registered
292 (<https://doi.org/10.17605/OSF.IO/A492M>)



293
 294
 295
 296
 297
 298
 299

Fig. 1 Individual differences in the quality of open data published by principal investigators (PIs) in ecology and evolution. Each caterpillar plot shows (A) the completeness and (B) the reusability of data from 362 datasets published by 100 PIs. PIs are identified by a vertical grey line and ordered from lowest to highest individual mean score. The colour of the data points indicates the year in which a study was published.

300
301
302



303
304
305
306
307
308
309
310
311
312
313
314

Fig. 2 Only publication year and PI seniority predict data reusability. Posterior means and 95% credible intervals from a Bayesian bivariate mixed-effects model to examine predictors of data completeness and reusability (n = 362 open datasets shared by 100 principal investigators [PIs]). The predictor variables included in the model include: the year of the PI's first publication as a measure of seniority, the year in which the study was published, the PI's gender, the journals' open data policy, and whether the PI was the corresponding author on the published study. Black dots indicate weak relationships and grey dots indicate posteriors that overlap zero.

315 **References**

316

- 317 1 Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
318 stewardship. *Sci. Data* **3**, 160018 (2016).
- 319 2 Soeharjono, S. & Roche, D. G. Reported individual costs and benefits of sharing open
320 data among Canadian faculty members in ecology and evolution. *BioScience*,
321 doi:10.1093/biosci/biab1024 (2021).
- 322 3 Roche, D. G. *et al.* Troubleshooting public data archiving: suggestions to increase
323 participation. *PLoS Biol.* **12**, e1001779, doi:10.1371/journal.pbio.1001779 (2014).
- 324 4 Culina, A., van den Berg, I., Evans, S. & Sánchez-Tójar, A. Low availability of code in
325 ecology: A call for urgent action. *PLoS Biol.* **18**, e3000763,
326 doi:10.1371/journal.pbio.3000763 (2020).
- 327 5 Costello, M. J. Motivating online publication of data. *BioScience* **59**, 418-427 (2009).
- 328 6 Parr, C. S. & Cummings, M. P. Data sharing in ecology and evolution. *Trends Ecol. Evol.*
329 **20**, 362-363, doi:10.1016/j.tree.2005.04.023 (2005).
- 330 7 Mills, J. A. *et al.* Archiving primary data: solutions for long-term studies. *Trends Ecol.*
331 *Evol.* **30**, 581-589 (2015).
- 332 8 Tenopir, C. *et al.* Changes in data sharing and data reuse practices and perceptions among
333 scientists worldwide. *PloS One* **10**, e0134826 (2015).
- 334 9 Digital Science & Figshare. The State of Open Data Report 2019 [online]: Available
335 from [https://www.digital-science.com/resources/portfolio-reports/the-state-of-open-data-](https://www.digital-science.com/resources/portfolio-reports/the-state-of-open-data-2019/)
336 [2019/](https://www.digital-science.com/resources/portfolio-reports/the-state-of-open-data-2019/). (2019).
- 337 10 Moore, A. J., McPeck, M. A., Rausher, M. D., Rieseberg, L. & Whitlock, M. C. The need
338 for archiving data in evolutionary biology. *J. Evol. Biol.* **23**, 659-660, doi:10.1111/j.1420-
339 9101.2010.01937.x (2010).
- 340 11 Vines, T. H. *et al.* Mandated data archiving greatly improves access to research data.
341 *FASEB J.* **27**, 1304-1308 (2013).
- 342 12 Federer, L. M. *et al.* Data sharing in PLOS ONE: An analysis of data availability
343 statements. *PloS one* **13**, e0194768 (2018).
- 344 13 O’Dea, R. E. *et al.* Towards open, reliable, and transparent ecology and evolutionary
345 biology. *BMC Biology* **19**, 68, doi:10.1186/s12915-021-01006-3 (2021).
- 346 14 Couture, J. L., Blake, R. E., McDonald, G. & Ward, C. L. A funder-imposed data
347 publication requirement seldom inspired data sharing. *PLOS ONE* **13**, e0199789,
348 doi:10.1371/journal.pone.0199789 (2018).
- 349 15 Roche, D. G., Kruuk, L. E., Lanfear, R. & Binning, S. A. Public data archiving in
350 ecology and evolution: how well are we doing? *PLoS Biol.* **13**, e1002295 (2015).
- 351 16 Towse, J., Ellis, D. A. & Towse, A. S. Opening Pandora’s Box: Peeking inside
352 psychology’s data sharing practices, and recommendations for change. doi:
353 10.31234/osf.io/k31236rux (2020).
- 354 17 Hardwicke, T. E. *et al.* Data availability, reusability, and analytic reproducibility:
355 Evaluating the impact of a mandatory open data policy at the journal *Cognition*. *Royal*
356 *Society Open Science* **5**, 180448 (2018).
- 357 18 Sholler, D., Ram, K., Boettiger, C. & Katz, D. S. Enforcing public data archiving policies
358 in academic publishing: A study of ecology journals. *Big Data Soc.* **6**,
359 2053951719836258, doi:10.1177/2053951719836258 (2019).

360 19 Christian, T.-M., Gooch, A., Vision, T. & Hull, E. Journal data policies: Exploring how
361 the understanding of editors and authors corresponds to the policies themselves. *PLOS*
362 *ONE* **15**, e0230281, doi:10.1371/journal.pone.0230281 (2020).

363 20 Roche, D. G. Open data: policies need policing. *Nature* **538**, 41 (2016).

364 21 Tenopir, C. *et al.* Data sharing by scientists: practices and perceptions. *PLoS One* **6**,
365 e21101, doi:10.1371/journal.pone.0021101 (2011).

366 22 Chawinga, W. D. & Zinn, S. Global perspectives of research data sharing: A systematic
367 literature review. *Library & Information Science Research* **41**, 109-122 (2019).

368 23 Abele-Brehm, A. E., Gollwitzer, M., Steinberg, U. & Schönbrodt, F. D. J. S. P. Attitudes
369 toward Open Science and public data sharing: A survey among members of the German
370 Psychological Society. **50**, 252 (2019).

371 24 Campbell, H. A., Micheli-Campbell, M. A. & Udyawer, V. Early career researchers
372 embrace data sharing. *Trends Ecol. Evol.* **34**, 95-98 (2019).

373 25 McNutt, M. K. *et al.* Transparency in authors' contributions and responsibilities to
374 promote integrity in scientific publication. *Proceedings of the National Academy of*
375 *Sciences* **115**, 2557-2560, doi:10.1073/pnas.1715374115 (2018).

376 26 Mons, B. Invest 5% of research funds in ensuring data are reusable. *Nature* **578**, 491
377 (2020).

378 27 Natural Sciences and Engineering Research Council of Canada. Women in Science and
379 Engineering in Canada. Available at: [https://www.nserc-crsng.gc.ca/doc/Reports-](https://www.nserc-crsng.gc.ca/doc/Reports-Rapports/WISE2017_e.pdf)
380 [Rapports/WISE2017_e.pdf](https://www.nserc-crsng.gc.ca/doc/Reports-Rapports/WISE2017_e.pdf). (2017).

381 28 Hadfield, J. D. MCMC methods for multi-response generalized linear mixed models: the
382 MCMCglmm R package. *Journal of Statistical Software* **33**, 1-22 (2010).

383 29 Roche, D. G., Careau, V. & Binning, S. A. Demystifying animal 'personality' (or not):
384 why individual variation matters to experimental biologists. *J. Exp. Biol.* **219**, 3832-3843,
385 doi:10.1242/jeb.146712 (2016).

386 30 Araya-Ajoy, Y. G., Mathot, K. J. & Dingemanse, N. J. An approach to estimate short-
387 term, long-term and reaction norm repeatability. *Methods in Ecology and Evolution* **6**,
388 1462-1473 (2015).

389 31 Dingemanse, N. J. & Dochtermann, N. A. Quantifying individual variation in behaviour:
390 mixed-effect modelling approaches. *J. Anim. Ecol.* **82**, 39-54 (2013).

391 32 Nakagawa, S. & Schielzeth, H. Repeatability for Gaussian and non-Gaussian data: a
392 practical guide for biologists. *Biol. Rev.* **85**, 935-956 (2010).

393
394

395
396
397
398
399

SUPPLEMENTARY INFORMATION

Table S1. Data completeness and reusability assessment. Reproduced from Roche et al. (2015) <https://doi.org/10.1371/journal.pbio.1002295>

Data Completeness		
Score	Description	Criteria
5	Exemplary	All the data necessary to reproduce the analyses and results (in practice) are archived. There is informative metadata with a legend detailing column headers, abbreviations, and units.
4	Good	All the data necessary to reproduce the analyses and results (in practice) are archived. The metadata are limited or absent, but column headings, abbreviations, and units can be understood from reading the paper.
3	Small omission	Most of the data necessary to repeat the analyses are archived except for a small amount (e.g., for a supporting or exploratory analysis). The metadata are informative OR the archived data can be interpreted from reading the paper.
2	Large omission	The main analyses in the paper cannot be redone because essential data are missing AND/OR insufficient metadata or information in the paper precludes interpreting the data AND/OR the authors archived summary statistics (e.g., means), but not the raw data used in the analyses.
1	Poor	The data are not archived OR the wrong data are archived OR insufficient information is provided in the metadata or paper for the data to be intelligible.
Data Reusability		
Score	Description	Criteria
5	Exemplary	The data are archived in a nonproprietary, human- and machine-readable file format that facilitates data aggregation and can be processed with both free and proprietary software (e.g., csv, text; see Table 3). The metadata are highly informative (such that column headings, abbreviations, and units can be understood in isolation from the original paper). Raw data are presented (perhaps in combination with processed data such as means). ^a
4	Good	The data are archived in a format that is designed to be machine readable with proprietary software (e.g., Excel), and the metadata are highly informative (such that column headings, abbreviations, and units can be understood in isolation from the original paper). [OR] The data are archived in a nonproprietary, human- and machine-readable file format, and the metadata are sufficiently informative to be understood when combined with information from the associated paper. Raw data are presented (perhaps in combination with processed data such as means). ^a
3	Average	The data are archived in a format that is designed to be machine readable with proprietary software (e.g., Excel). The metadata are sufficiently informative to be understood when combined with information from the associated paper. Raw data are presented (perhaps in combination with processed data such as means). ^a
2	Poor	The data are archived in a human- but not machine-readable format. The metadata are highly informative OR sufficiently informative to be understood with information from the associated paper. Raw data are presented (perhaps in combination with processed data such as means). ^a
1	Very poor	The metadata are insufficient for the data to be intelligible even when combined with information from the associated paper AND/OR processed but not raw data are presented. ^a

N.B. Reusability was assessed for archived data independently of completeness. One point was subtracted when data were included as supplementary material on the journal website, except when the reusability score was 1 to avoid zero values (see S1 Text).

^a Raw data were considered unprocessed data (e.g., trait values used in a principal component analysis rather than principle component scores, values underlying means presented in figures). Studies that did not archive duplicate or triplicate measurements to account for measurement error were not considered as missing raw data.

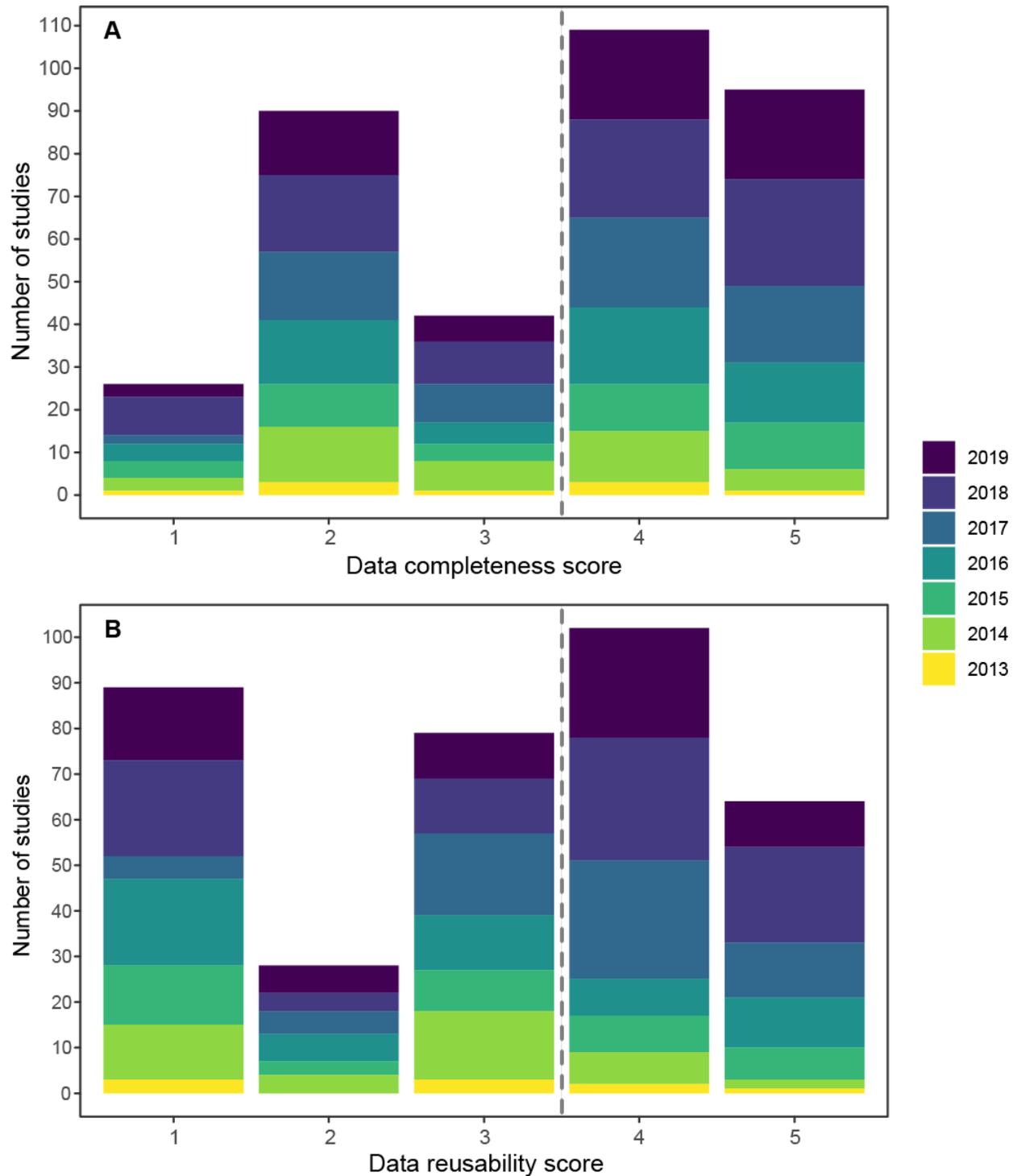
400

doi:10.1371/journal.pbio.1002295.t002

401 **Table S2.** The number of ecologists and evolutionary biologists in Canada's top-21 ranked
 402 universities according to the 2019 THE World University Rankings. Universities are ordered
 403 alphabetically. Indicated are the number of PIs in ecology and evolution at each university (E&E
 404 PIs), PIs with at least two journal articles published between Jan 2013-June 2019 containing a
 405 data availability statement and/or associated open data (Open data PIs), and PIs randomly
 406 selected for analysis in this study (Selected PIs).
 407

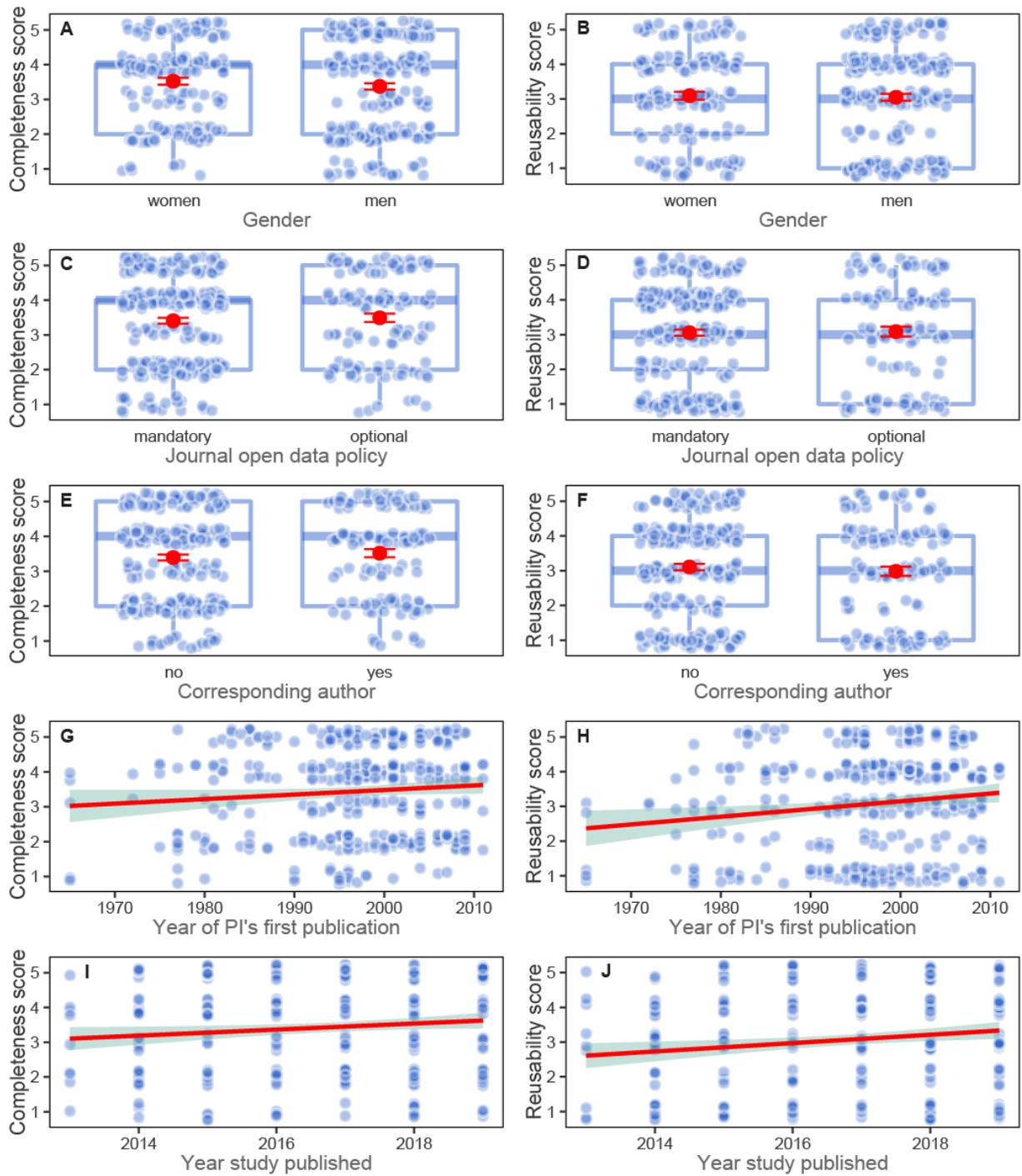
University	E&E PIs	Open data PIs	Selected PIs	
			Women	Men
Carleton University	11	9	3	3
Dalhousie University	12	8	1	4
Laval University	18	8	1	4
McGill University	17	16	3	3
McMaster University	7	0	0	0
Memorial University	13	5	2	2
Queen's University	12	5	1	4
Simon Fraser University	18	13	3	3
University of Alberta	23	11	3	3
University of British Columbia	41	21	3	3
University of Calgary	17	8	3	3
University of Guelph	19	8	2	3
University of Manitoba	14	4	0	2
University of Montreal	16	9	2	4
University of Ottawa	13	9	3	3
University of Saskatchewan	14	4	2	2
University of Toronto	35	31	3	3
University of Victoria	11	7	3	3
University of Waterloo	8	2	0	0
Western University	18	11	2	3
York University	14	5	4	1
Total	351	194	44	56

408
 409



411
412
413
414
415
416

Fig. S1 Frequency distribution of the (A) completeness and (B) reusability scores for open datasets associated with 362 studies shared by 100 researchers between 2013-2019. A score of 5 indicates exemplary archiving, and a score of 1 indicates poor archiving (Table S1). Studies with scores of 3 or lower (left of the grey dashed lines) are incomplete or difficult to reuse.

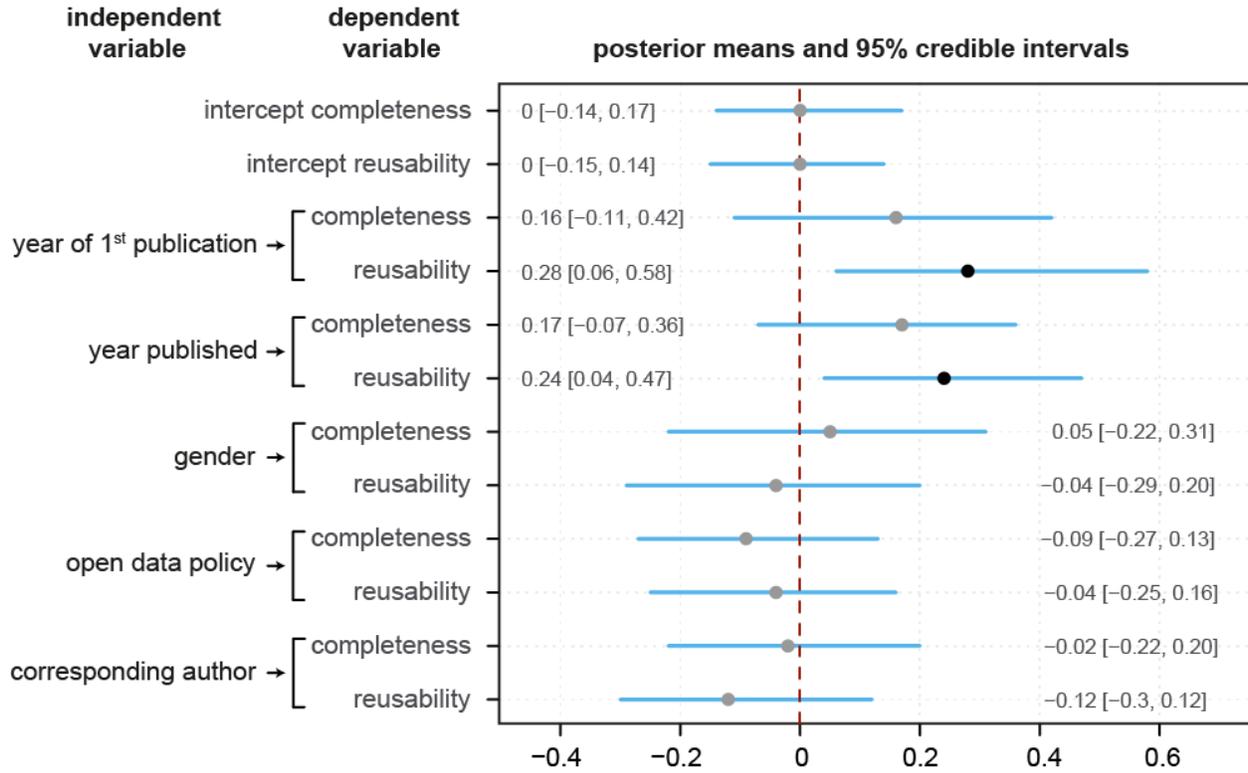


417
 418
 419
 420
 421
 422
 423
 424
 425

Fig. S2 The relationship between five independent variables and the completeness (A, C, E, G, I) and reusability score (B, D, F, H, J) of 362 open datasets shared by 100 principal investigators (PIs). Independent variables include: the gender of the PI, the journals' open data policy, whether the PI was the corresponding author on the associated paper, the year of the PI's first publication (i.e., seniority), and the year in which the study was published. Red dots are means and error bars represent 95% confidence intervals (CIs). Red lines are least square regressions and shaded areas represent 95% CIs. Note that these graphs are included for visualization purposes only to show

426 the raw data. The relationships depicted do not control for other predictors in the analysis, nor do
427 they control for repeated measurements (i.e., the random effects PI and university), which are
428 included in the main analysis reported in the main text. For posterior means and 95% credible
429 intervals from the Bayesian bivariate mixed-effects model, see Fig. 2.
430

431
432



433
434

435 **Fig. S3** Posterior means and 95% credible intervals from a Bayesian bivariate mixed-effects
 436 model investigating the effect of five independent variables on the completeness and reusability
 437 scores of 362 open datasets shared by 100 principal investigators (PIs). The prior was specified
 438 as a parameter-expanded prior (in contrast to an inverse-Wishart prior; see Fig. 2). The predictor
 439 variables included in the model include: the year of the PI's first publication as a measure of
 440 seniority, the year in which the study was published, the PI's gender, the journals' open data
 441 policy, and whether the PI was the corresponding author on the published study. Black dots
 442 indicate weak relationships and grey dots indicate posteriors that overlap zero.