

1 **Slow improvement to the archiving quality of open datasets shared by researchers in**
2 **ecology and evolution**

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, K1S 5B6*

8 ² *Département de science biologiques, Université de Montréal, Montréal, Canada, H3C 3J7*

9 ³ *Institut de Biologie, Université de Neuchâtel, Neuchâtel, Switzerland; 2000*

10 ⁴ *Department of Biology, McGill University, Montréal, Canada, H3A 1B1*

11

12 Correspondence to: Dominique G. Roche

13 209 Nesbit Biology Building, Carleton University

14 1125 Colonel By Drive

15 Ottawa, Ontario, Canada, K1S 5B6

16 E: dominique.roche@mail.mcgill.ca

17 T: 1-438-889-5822

18

19 **Key words:** data sharing, FAIR data, metascience, open science, public data archiving,
20 reproducibility

21 **Abstract**

22 Many leading journals in evolution and ecology now mandate open data upon publication. Yet,
23 there is very little oversight to ensure the completeness and reusability of archived datasets, and
24 we currently have a poor understanding of the factors associated with high quality (FAIR) data-
25 sharing. We assessed 362 open datasets linked to first- or senior-authored papers published by
26 100 principal investigators (PIs) in the fields of evolution and ecology over a period of seven
27 years to identify predictors of data completeness and reusability ('data archiving quality').
28 Datasets scored low on these metrics: 56.4% were complete and 45.9% were reusable. Data
29 reusability, but not completeness, was slightly higher for more recently archived datasets and PIs
30 with less seniority. Journal open data policy, PI gender, and PI corresponding author status were
31 unrelated to data archiving quality. However, PI identity explained a large proportion of the
32 variance in data completeness (27.8%) and reusability (22.0%), indicating consistent inter-
33 individual differences in data sharing practices by PIs across time and contexts. Several PIs
34 consistently shared data of either high or low archiving quality, but most PIs were inconsistent in
35 how well they shared. One explanation for the high intra-individual variation we observed is that
36 PIs often conduct research through students and post-docs, who may be responsible for the data
37 collection, curation and archiving. Levels of data literacy vary among trainees and PIs may not
38 regularly perform quality control over archived files. Our findings suggests that research data
39 management training and culture within a PI's group are likely to be more important
40 determinants of data archiving quality than other factors such as a journal's open data policy.
41 Greater incentives and training for individual researchers at all career stages could improve data
42 sharing practices and enhance data transparency and reusability.

43

44 **Main text**

45 The debate is over regarding the value of open and FAIR data (Findable, Accessible, Interoperable,
46 and Reusable data [1]). Many journals, funding agencies, and policymakers agree that the societal
47 benefits of publicly sharing (non-sensitive) research data outweigh any perceived or reported costs
48 to individual researchers [2-6]. Making data underlying scientific studies publicly available
49 facilitates exploring, validating, and building on published results [7], with few researchers in
50 evolution and ecology (E&E) reporting negative outcomes from sharing [2]. Not only do open data

51 accelerate scientific discovery, as illustrated during the Covid-19 pandemic [8], but sharing data
52 also promotes an academic value system that is more equitable, diverse, and inclusive [9-11].

53
54 Some researchers are reluctant to share the data underlying their published results [6, 12, 13], yet
55 most researchers view open data positively [2, 14, 15]. Since 2010, when a handful of journals
56 began requiring open data in E&E [16], policies encouraging this practice have grown rapidly.
57 Now, 20 percent of journals publishing research in E&E mandate open data [17]. Strong journal
58 policies are effective at ensuring that more datasets are shared [18-20], which is often touted as a
59 win for open science [21]. Yet, problems persist [19, 22, 23]. For instance, more than half of
60 open datasets associated with 100 E&E studies published in 2012 and 2013 were incomplete
61 and/or archived in ways that prevented reuse [24]. Similar issues have been documented in
62 psychology [25] and cognition research [26], pointing to the inherent problem with journals
63 mandating open data without appropriate oversight or quality control [27-29]: datasets get
64 archived but the majority are incomplete and challenging to reuse. Developing effective
65 strategies to promote good data sharing practices requires that we first identify which factors are
66 associated with complete and reusable open data [30].

67
68 We assessed the archiving quality (completeness and reusability) of open datasets associated
69 with publications by tenured or tenure-track E&E faculty members (PIs) in biology departments
70 at the 21 highest-ranked universities in Canada. PIs were necessarily first or last author on the
71 publications assessed (see pre-registered methods). Data completeness (availability of data
72 allowing computational reproducibility) and reusability (ease with which data can be reused by
73 third parties) were assessed following Roche et al. [24]. Scores above 3 on two 5-point scales
74 indicate complete or reusable data (Table 1). In total, we examined 362 datasets shared by 44
75 women and 56 men (Table S1). We ran a Bayesian bivariate mixed-effects linear model with PI
76 identity and their institution of employment as random factors. We tested whether article
77 publication date, journal open data policy, and the seniority, gender and corresponding author
78 status of PIs predicted the archiving quality of their open datasets. A post-hoc exploratory
79 analysis was also carried out to examine the relationship between the complexity of datasets
80 (estimated as the number and size of archived data files) and their completeness and reusability.

81

82 The completeness and reusability scores of datasets varied considerably within and among PIs
83 (Fig. 1). Overall, 56.4% of datasets were complete (mean completeness score of 3.4 ± 1.3 SD),
84 and 45.9% were reusable (mean reusability score of 3.1 ± 1.4 SD) (Fig. S1) This represents a
85 moderate improvement of approximately 10% above the completeness and reusability of datasets
86 associated with E&E studies published in 2012 and 2013 [24]. Data completeness and reusability
87 were strongly correlated within ($R^2=0.79$, 95% CI: 0.72–0.85) and among ($R^2=0.77$, 95% CI:
88 0.56–0.92) PIs (see [31] for an explanation of among and within individual correlations).

89
90 Open data is a relatively recent concept in E&E, having been introduced in earnest a decade ago
91 [5, 32]. As such, PIs who developed their research skills prior to this period may be less likely to
92 have incorporated these principles into their research workflow. We assessed datasets as far back
93 as 2013 in our analysis but included faculty members hired as recently as 2019, using year of
94 first scientific publication as a proxy for PI seniority. Therefore, our study likely includes
95 datasets published by new PIs during their PhD and post-doc years, when they might have had
96 access to various training opportunities. For instance, a growing number of biology departments
97 recognize the value of data science and initiatives such as Data Carpentry (datacarpentry.org)
98 and FOSTER (fosteropenscience.eu) now routinely offer workshops in data management across
99 North America and Europe. We found that PIs with less seniority tended to share slightly more
100 reusable data than PIs with more seniority, suggesting that early training initiatives may be
101 bearing fruit (Figs 2, S2 H). This result is good news because younger researchers tend to be
102 more fearful and reluctant to share their data than senior researchers [33] despite reporting a
103 more favorable attitude towards open data [33-36] [but see 37]. Our study included datasets
104 spanning seven years (2013-2019). We found that datasets associated with more recent studies
105 were slightly more reusable than those of older studies (Figs 2, S2 J). In contrast, dataset
106 completeness was independent of PI seniority and publication date (Fig. 2, S2 G,I). These results
107 suggest that, while there have been slight improvements to data sharing practices through time,
108 these are slow to change. Training and increased exposure to open science practices are no doubt
109 contributing to this slow improvement, but additional work is needed at all career stages to
110 enhance data archiving quality.

111

112 We examined whether data archiving quality was influenced by PIs being corresponding author
113 on the published article (in addition to being first or senior author). We took PI corresponding
114 author status as an indicator that the PIs themselves archived the data. Assuming that PIs are
115 highly competent at managing their research data, we expected open datasets with PIs as
116 corresponding author to be of higher archiving quality, on average, than those archived by
117 presumably less experienced researchers (likely students or post-docs). We found no support for
118 this hypothesis: data archiving quality was unrelated to corresponding author status (Figs 2, S2
119 E,F). Conventions regarding who is corresponding author on a published study may vary among
120 sub-disciplines and research labs. However, the corresponding author is ultimately responsible
121 for compliance with journal policies, including open data [38]. The fact that corresponding
122 author status had no bearing on data archiving quality is worrying and suggests that PIs do not
123 understand the responsibilities associated with this role, do not have the tools or training to
124 ensure that open data are compliant with journal policies, or simply do not care. Education,
125 capacity-building, and incentives targeted at individuals are needed to address these issues [2, 7,
126 39, 40].

127
128 We had no a priori hypothesis for why PI gender might influence data archiving quality.
129 However, we included this predictor in the model because we believe that gender differences are
130 important to consider in scientific research. We found no evidence to suggest that PI gender
131 influences the quality of open data (Figs 2, S2 A,B). This result is encouraging given that men in
132 E&E self-identify as experiencing more costs than women as a result of sharing open data [2].
133 Women accounted for almost half of the PIs assessed in our study, yet far fewer than 50% of PIs
134 in biology departments at Canadian institutions identify as women [41].

135
136 When journals have a mandatory open data policy, the number of archived datasets underlying
137 published research articles increases [18]. We tested whether such policies also translate into
138 more complete and reusable data. We hypothesized that data archiving quality would be higher
139 for studies published in journals requiring open data. Alternatively, it is also possible that
140 researchers voluntarily archiving datasets in journals without a policy share higher quality data
141 than researchers who are forced to do so. Contrary to this logic, a journal's open data policy had
142 no bearing on data archiving quality (Fig. 2, S2 C,D), indicating that policies alone do little to

143 ensure that shared data are complete and reusable. Some of the world's largest funding agencies
144 (e.g., ERC, NSF, NERC, Canadian Tri-Council) now require that PIs specify data management
145 and/or sharing plans to obtain funding. However, compliance with these policies is low [22].
146 Unless resources are invested in training, technical support and policy oversight [2, 15], data risk
147 not being archived or not contributing to advancing knowledge in instances where they are made
148 available.

149
150 We assessed multiple datasets by the same researchers, which allowed us to calculate
151 repeatability scores for both data completeness and reusability. Repeatability (R) ranges between
152 0 and 1 and is the proportion of the total variance in scores attributed to among (or inter-)
153 individual differences: high R values indicate large score differences among individuals and
154 consistent scores within individuals [42]. Data completeness was moderately repeatable with R_{adj}
155 = 0.28 (95% CI: 0.16–0.39), and data reuse with R_{adj} = 0.22 (95% CI: 0.14–0.34), revealing
156 differences within and among PIs (Fig. 1). This variability reflects several realities that merit
157 discussion. On the one hand, PIs in academia often conduct research through students and post-
158 docs, who may ultimately be responsible for data collection, curation and archiving. Thus, in
159 some cases, the PI may not have performed quality control over the archived files, potentially
160 explaining the considerable within-individual variation in data completeness and reusability
161 scores we observed (variable scores in Fig. 1). In these cases, data archiving quality may be a
162 better reflection of data management not by the PI, but by the person within the PI's research
163 group who was responsible for archiving the dataset. On the other hand, some PIs consistently
164 scored high or low in both data completeness and reuse (low and high scores in Fig. 1). This
165 consistency within research groups suggests both a robust lab culture promoting good research
166 data management or, alternatively, a PI's lack of competence or reluctance to engage in data
167 sharing and student training in this regard. PIs who oppose open data initiatives (e.g., [2, 12]) are
168 unlikely to respond positively to incentives or training opportunities to improve data archiving
169 quality and FAIRness. However, our results suggest that only a minority of PIs potentially fall
170 within this category (approximately 10%; Fig. 1). Rather, most PIs were inconsistent in how they
171 shared data associated with their publications, or consistently shared highly complete and
172 reusable data. We found no indication that dataset complexity influenced archiving quality,
173 suggesting that PI often struggle to share even simple datasets (Fig. S4). These finding points to

174 the importance of facilitating sound research data management practices within research groups
175 to achieve high-quality, FAIR data sharing.

176

177 Overall, our data suggest that journal policies are ineffective at ensuring that open data in E&E
178 are complete and reusable. We also found that data archiving quality is slow to improve over
179 time. However, most PIs did share high-quality open data, either consistently or occasionally.
180 Striking variation in data archiving quality within PIs suggests that education, training, and
181 technical support could help raise the bar by enabling good data sharing practices to become the
182 rule rather than the exception [39].

183

184 *Materials and Methods*

185 Our methods were pre-registered at <https://doi.org/10.17605/OSF.IO/A492M>. We assessed open
186 datasets from research faculty members in biology departments at the 21 highest-ranked
187 Canadian universities based on the 2019 Times Higher Education World University Rankings.
188 Although we initially planned to select the top 20 Canadian universities, we retained 21
189 universities due to a three-way tie for rank 19. Our study focused on Canadian academic faculty.
190 However, our findings are likely to be representative of the broader population of PIs in E&E
191 given Canada's diverse academic institutions as well as the high degree of PI mobility in today's
192 globalized academic landscape. Furthermore, many granting agencies in Europe and the USA
193 require that data from funded research be publicly archived within a specified timeframe of
194 publishing. This is not yet the case in Canada: the Tri-council Granting agencies now require a
195 data management plan for grants submitted in 2021 and beyond, but this does not yet include a
196 requirement for open data. This allowed us to assess the effect of journal policies on archiving
197 practices by Canadian PI's independent of requirements from funding agencies on the same
198 practices.

199

200 We reviewed the biology department website at each of the 21 selected universities in a random
201 order and identified all researchers primarily conducting research in the fields of ecology and/or
202 evolution (E&E) with a rank of assistant, associate or full professor. Adjunct professors and
203 researchers who primarily focus on molecular biology, genetics, genomics, bioinformatics,
204 theoretical biology, comparative physiology and paleontology were excluded given our focus on

205 researchers in E&E. Each researcher's primary fields of study were determined from public
206 information on the university websites and cross-checked by a minimum of two people (IB, FD,
207 SAB, RD, DGR). This criterion served to limit the scope of the study to E&E and facilitate
208 consistent assessment of datasets, given our shared expertise. To standardize the selection of
209 researchers across universities and avoid bias, we omitted E&E researchers who are primarily
210 affiliated with departments other than biology (e.g., environmental sciences, natural resources,
211 fisheries and ocean sciences, veterinary sciences). In total, we identified 351 researchers that met
212 these criteria Table S1).

213

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

232

233 The strength of a journal's open data policy and date of implementation was determined by
234 reviewing each journal's author guidelines and relevant editorials. When necessary, we contacted
235 journal editors for clear information on whether open data were required (i.e., mandatory open

236 data) or encouraged (i.e., optional open data) as a condition of publication at the time a paper
237 was published. Journals without an open data policy were categorized as optional open data.

238

239 We identified 194 researchers with at least two articles containing a data availability statement
240 and/or open data across the 21 universities (Table S1). Of these, we randomly selected up to four
241 women and four men at each university to evaluate the data archiving quality of their open
242 datasets. We aimed to randomly select three women and three men at each university but some
243 universities had fewer than three researchers per gender (Table S1). The departments of biology
244 at two institutions had no researchers that met our selection criteria. One researcher identified as
245 gender non-binary but was not part of our random sample. We made assumptions about gender
246 based on names and pronouns used in public profiles on university websites or social media. We
247 recognize that gender presentation, names, and pronouns are not necessarily indications of a
248 person's gender and that, in the absence of additional information from the individuals, we may
249 have unintentionally made incorrect assumptions about individuals' genders. In total we assessed
250 362 datasets published in 97 journals by 100 PIs. We scored the completeness and reusability of
251 shared datasets on a scale from 0 (min score) to 5 (max score) following Roche *et al.* 2015 [24]
252 (Table 1). The number of datasets assessed per researcher ranged from two to five; if a
253 researcher had more than five shared datasets in the period from 2013-2019, we selected the
254 most recent five.

255

256 *Statistical analysis*

257 We used a Bayesian bivariate mixed-effects regression model (R package MCMCglmm v2.32
258 [43]) to identify factors influencing data archiving quality and estimate its repeatability (i.e., the
259 proportion of the total variance attributable to differences among individuals) [42, 44, 45]. Data
260 completeness and reusability scores were included as two dependent variables in the model;
261 researcher ID and university were specified as random effects, with researcher nested within
262 university; PI gender, PI seniority (measured as the year of their first peer-reviewed publication,
263 assessed on Google Scholar or Web of Science), PI author status (corresponding author or not),
264 journal open data policy at the time of publication (mandatory, optional), and year of study
265 publication were included as fixed effects. Journal impact factor (JIF) was not included in the

266 model because MCMCglmm does not tolerate missing values in the fixed predictors (this is a
267 deviation from the preregistered methodology doi:10.17605/OSF.IO/A492M).

268
269 The two dependent variables were mean-centered and standardized to one standard deviation unit
270 prior to inclusion in the model (i.e., mean=0, standard deviation=1). The numerical (PI seniority,
271 year of study publication) and categorical predictors (gender, corresponding author, journal open
272 data policy) were mean-centered and standardized to two standard deviation units (i.e., mean=0,
273 standard deviation=0.5) following Araya-Ajoy et al. [46]. Categorical predictors were treated as
274 binary variables (values of 0 and 1) to allow centering and standardization. The advantage of
275 mean-centering the predictors is that it ensures model intercepts are estimated for the average
276 value of the predictors, facilitating interpretation of the results. Mean centering allows the
277 estimate of the intercepts to be calculated for the average ‘environmental’ conditions [44]; the
278 use of two standard deviations for predictor standardization allows for direct comparison of the
279 variance explained by categorical and continuous predictors [46].

280
281 We specified a mildly informative inverse-Wishart prior and tested the sensitivity of the model to
282 prior specification by examining how the posterior means and 95% credible intervals changed
283 when specifying a parameter-expanded prior [see 43]. We checked the model by plotting the
284 traces of the parameters, examining autocorrelation among samples drawn by MCMCglmm, and
285 computing the Gelman-Rubin statistic to evaluate convergence (see archived script). Model
286 diagnostics were satisfactory and conclusions were not sensitive to the choice of prior (Fig. S3).

287
288 We calculated the adjusted repeatability (R_{adj}) for a researcher’s data completeness and data
289 reusability as the proportion of the total variance due to differences among individuals when
290 accounting for fixed and random effects in the statistical model [45]. Within- and among-
291 individual correlations between data completeness and reusability were calculated as outlined in
292 Roche et al. [42].

293
294 Following a reviewer suggestion, we conducted an exploratory (i.e., non-registered) analysis to
295 examine whether dataset complexity could explain variation in data archiving quality. The
296 rationale for this analysis is that simple datasets (e.g., simple experimental design, few variables,

297 low sample size) might be easier to share in a complete and reusable fashion than complex
298 datasets containing many different experiments or observational studies, a large number of
299 variables, and many measurements. We estimated dataset complexity as the number and size (in
300 KB) of data files and plotted these variables against data completeness and reusability scores for
301 each archived dataset (Fig. S4).

302

303 All analyses were done in R version 4.0.3.

304

305 **Acknowledgements and funding sources**

306 We thank Redouan Bshary for the idea that led to this study and Ellen Bledsoe and the associate
307 editor and two anonymous reviewers for helpful comments on the manuscript. We acknowledge
308 funding from the Natural Sciences and Engineering Research Council of Canada (grant no. UIF-
309 537860–2018) and the Canada Research Chair program. DGR was supported by the European
310 Union’s Horizon 2020 research and innovation program under Marie Skłodowska-Curie grant
311 agreement no. 838237-OPTIMISE.

312

313 **Author Contributions**

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

317

318 **Competing Interests statement**

319 DGR is a member of Research Data Canada’s Policy Committee, a member of the Canadian
320 National Committee for CODATA, and the president of the Society for Open, Reproducible and
321 Transparent Ecology and Evolutionary biology (sortee.org).

322

323 **Data and code availability**

324 This study was pre-registered (<https://doi.org/10.17605/OSF.IO/A492M>). The anonymized data
325 and the analysis script to reproduce our results are available on the Open Science Framework
326 (<https://doi.org/10.17605/OSF.IO/P2YWM>) and were shared with the editors and reviewers upon
327 submission.

328 **Table and Figure captions**

329

330 **Table 1. Scoring system and criteria used to assess data completeness and reusability.**

331 Reproduced from Roche et al. (2015) <https://doi.org/10.1371/journal.pbio.1002295>

332

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

338

339 **Fig. 2 Only publication year and PI seniority predict data reusability (but not completeness).**

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

347 **References**

- 348 1. Wilkinson M.D., Dumontier M., Aalbersberg I.J., Appleton G., Axton M., Baak A., Blomberg
349 N., Boiten J.-W., da Silva Santos L.B., Bourne P.E. 2016 The FAIR Guiding Principles
350 for scientific data management and stewardship. *Sci. Data* **3**, 160018.
- 351 2. Soeharjono S., Roche D.G. 2021 Reported individual costs and benefits of sharing open data
352 among Canadian faculty members in ecology and evolution. *BioScience* **71**, 750–756.
353 (doi:10.1093/biosci/biab024).
- 354 3. Culina A., van den Berg I., Evans S., Sánchez-Tójar A. 2020 Low availability of code in
355 ecology: A call for urgent action. *PLoS Biol.* **18**, e3000763.
356 (doi:10.1371/journal.pbio.3000763).
- 357 4. Costello M.J. 2009 Motivating online publication of data. *BioScience* **59**, 418-427.
- 358 5. Parr C.S., Cummings M.P. 2005 Data sharing in ecology and evolution. *Trends Ecol. Evol.* **20**,
359 362-363. (doi:10.1016/j.tree.2005.04.023).
- 360 6. Roche D.G., Lanfear R., Binning S.A., Haff T.M., Schwanz L.E., Cain K.E., Kokko H.,
361 Jennions M.D., Kruuk L.E. 2014 Troubleshooting public data archiving: suggestions to
362 increase participation. *PLoS Biol.* **12**, e1001779. (doi:10.1371/journal.pbio.1001779).
- 363 7. Westoby M., Falster D.S., Schrader J. 2021 Motivating data contributions via a distinct career
364 currency. *Proc. R. Soc. Biol. Sci. B* **288**, 20202830. (doi:10.1098/rspb.2020.2830).
- 365 8. Xu B., Gutierrez B., Mekaru S., Sewalk K., Goodwin L., Loskill A., Cohn E.L., Hswen Y.,
366 Hill S.C., Cobo M.M., et al. 2020 Epidemiological data from the COVID-19 outbreak,
367 real-time case information. *Sci. Data* **7**, 106. (doi:10.1038/s41597-020-0448-0).
- 368 9. Hood A.S.C., Sutherland W.J. The data-index: An author-level metric that values impactful
369 data and incentivizes data sharing. *Ecology and Evolution* **11**, 14344–14350.
370 (doi:10.1002/ece3.8126).
- 371 10. Buxton R., Nyboer E.A., Pigeon K., Raby G.D., Rytwinsky T., Gallagher A.J., Schuster R.,
372 Lin H.S., Fahrig L., Bennett J.R., et al. 2021 Avoiding wasted research resources in
373 conservation science. *Conserv. Sci. Prac.* **3**, e329. (doi:10.1111/csp2.329).
- 374 11. Roche D.G., O'Dea R.E., Kerr K.A., Rytwinski T., Schuster R., Nguyen V.M., Young N.,
375 Bennett J.R., Cooke S.J. Closing the knowledge-action gap in conservation with open
376 science. *Conserv. Biol.* (doi:10.1111/cobi.13835).

- 377 12. Mills J.A., Teplitsky C., Arroyo B., Charmantier A., Becker P.H., Birkhead T.R., Bize P.,
378 Blumstein D.T., Bonenfant C., Boutin S. 2015 Archiving primary data: solutions for
379 long-term studies. *Trends Ecol. Evol.* **30**, 581-589.
- 380 13. Tedersoo L., Küngas R., Oras E., Köster K., Eenmaa H., Leijen Ä., Pedaste M., Raju M.,
381 Astapova A., Lukner H., et al. 2021 Data sharing practices and data availability upon
382 request differ across scientific disciplines. *Sci. Data* **8**, 192. (doi:10.1038/s41597-021-
383 00981-0).
- 384 14. Tenopir C., Rice N.M., Allard S., Baird L., Borycz J., Christian L., Grant B., Olendorf R.,
385 Sandusky R.J. 2020 Data sharing, management, use, and reuse: Practices and perceptions
386 of scientists worldwide. *PLOS ONE* **15**, e0229003. (doi:10.1371/journal.pone.0229003).
- 387 15. Science Digital, Simons N., Goodey G., Hardeman M., Clare C., Gonzales S., Strange D.,
388 Smith G., Kipnis D., Iida K., et al. 2021 The State of Open Data 2021. Digital Science.
389 Report. <https://doi.org/10.6084/m9.figshare.17061347.v1>
- 390 16. Moore A.J., McPeck M.A., Rausher M.D., Rieseberg L., Whitlock M.C. 2010 The need for
391 archiving data in evolutionary biology. *J. Evol. Biol.* **23**, 659-660. (doi:10.1111/j.1420-
392 9101.2010.01937.x).
- 393 17. Berberi I., Roche D.G. 2021 Living database of journal data policies in E&E. Available on
394 the Open Science Framework: <https://doi.org/10.17605/OSF.IO/D6SP3>.
- 395 18. Vines T.H., Andrew R.L., Bock D.G., Franklin M.T., Gilbert K.J., Kane N.C., Moore J.-S.,
396 Moyers B.T., Renaut S., Rennison D.J. 2013 Mandated data archiving greatly improves
397 access to research data. *FASEB J.* **27**, 1304-1308.
- 398 19. Federer L.M., Belter C.W., Joubert D.J., Livinski A., Lu Y.-L., Snyders L.N., Thompson H.
399 2018 Data sharing in PLOS ONE: An analysis of data availability statements. *PloS One*
400 **13**, e0194768.
- 401 20. Roche D.G., Raby G.D., Norin T., Ern R., Scheuffele H., Skeeles M., Morgan R.,
402 Andreassen A.H., Clements J.C., Louissaint S., et al. 2022 Paths towards greater
403 consensus building in experimental biology. *J. Exp. Biol.* **225**, jeb243559.
404 (doi:10.1242/jeb.243559).
- 405 21. O’Dea R.E., Parker T.H., Chee Y.E., Culina A., Drobniak S.M., Duncan D.H., Fidler F.,
406 Gould E., Ihle M., Kelly C.D., et al. 2021 Towards open, reliable, and transparent

- 407 ecology and evolutionary biology. *BMC Biology* **19**, 68. (doi:10.1186/s12915-021-01006-
408 3).
- 409 22. Couture J.L., Blake R.E., McDonald G., Ward C.L. 2018 A funder-imposed data publication
410 requirement seldom inspired data sharing. *Plos One* **13**, e0199789.
411 (doi:10.1371/journal.pone.0199789).
- 412 23. Vines T.H. 2014 Scientific community: Journals must boost data sharing. *Nature* **508**, 44-44.
- 413 24. Roche D.G., Kruuk L.E., Lanfear R., Binning S.A. 2015 Public data archiving in ecology and
414 evolution: how well are we doing? *PLoS Biol.* **13**, e1002295.
- 415 25. Towse J., Ellis D.A., Towse A.S. 2020 Opening Pandora's Box: Peeking inside psychology's
416 data sharing practices, and recommendations for change. *Beh. Res. Meth.*
417 (doi:10.31234/osf.io/k6rux).
- 418 26. Hardwicke T.E., Mathur M.B., MacDonald K., Nilsonne G., Banks G.C., Kidwell M.C.,
419 Hofelich Mohr A., Clayton E., Yoon E.J., Henry Tessler M. 2018 Data availability,
420 reusability, and analytic reproducibility: Evaluating the impact of a mandatory open data
421 policy at the journal *Cognition*. *R. Soc. Open Sci.* **5**, 180448.
- 422 27. Sholler D., Ram K., Boettiger C., Katz D.S. 2019 Enforcing public data archiving policies in
423 academic publishing: A study of ecology journals. *Big Data Soc.*, 1-18.
424 (doi:10.1177/2053951719836258).
- 425 28. Christian T.-M., Gooch A., Vision T., Hull E. 2020 Journal data policies: Exploring how the
426 understanding of editors and authors corresponds to the policies themselves. *PLOS One*
427 **15**, e0230281. (doi:10.1371/journal.pone.0230281).
- 428 29. Roche D.G. 2016 Open data: policies need policing. *Nature* **538**, 41.
- 429 30. Wass M.N., Ray L., Michaelis M. 2019 Understanding of researcher behavior is required to
430 improve data reliability. *GigaScience* **8**, 1-8. (doi:10.1093/gigascience/giz017).
- 431 31. Careau V., Wilson R.S. 2017 Of Uberfleas and Krakens: Detecting Trade-offs Using Mixed
432 Models. *Integrative and Comparative Biology* **57**, 362-371. (doi:10.1093/icb/ix015).
- 433 32. Fairbairn D.J. 2011 The advent of mandatory data archiving. *Evolution* **65**, 1-2.
- 434 33. Tenopir C., Dalton E.D., Allard S., Frame M., Pjesivac I., Birch B., Pollock D., Dorsett K.
435 2015 Changes in data sharing and data reuse practices and perceptions among scientists
436 worldwide. *PloS One* **10**, e0134826.

- 437 34. Tenopir C., Allard S., Douglass K., Aydinoglu A.U., Wu L., Read E., Manoff M., Frame M.
438 2011 Data sharing by scientists: practices and perceptions. *PLoS One* **6**, e21101.
439 (doi:10.1371/journal.pone.0021101).
- 440 35. Chawinga W.D., Zinn S. 2019 Global perspectives of research data sharing: A systematic
441 literature review. *Libr. Inf. Sci. Res.* **41**, 109-122.
- 442 36. Abele-Brehm A.E., Gollwitzer M., Steinberg U., Schönbrodt F.D.J.S.P. 2019 Attitudes
443 toward Open Science and public data sharing: A survey among members of the German
444 Psychological Society. *Soc. Psychol.* **50**, 252-260. (doi:10.1027/1864-9335/a000384).
- 445 37. Campbell H.A., Micheli-Campbell M.A., Udyawer V. 2019 Early career researchers embrace
446 data sharing. *Trends Ecol. Evol.* **34**, 95-98.
- 447 38. McNutt M.K., Bradford M., Drazen J.M., Hanson B., Howard B., Jamieson K.H., Kiermer
448 V., Marcus E., Pope B.K., Schekman R., et al. 2018 Transparency in authors'
449 contributions and responsibilities to promote integrity in scientific publication. *Proc.*
450 *Natl. Acad. Sci.* **115**, 2557-2560. (doi:10.1073/pnas.1715374115).
- 451 39. Mons B. 2020 Invest 5% of research funds in ensuring data are reusable. *Nature* **578**, 491.
- 452 40. Renaut S., Budden A.E., Gravel D., Poisot T., Peres-Neto P. 2018 Management, archiving,
453 and sharing for biologists and the role of research institutions in the technology-oriented
454 age. *BioScience* **68**, 400-411. (doi:10.1093/biosci/biy038).
- 455 41. Natural Sciences and Engineering Research Council of Canada. 2017 Women in Science and
456 Engineering in Canada. Available at: [https://www.nserc-crsng.gc.ca/doc/Reports-
457 Rapports/WISE2017_e.pdf](https://www.nserc-crsng.gc.ca/doc/Reports-Rapports/WISE2017_e.pdf).
- 458 42. Roche D.G., Careau V., Binning S.A. 2016 Demystifying animal 'personality' (or not): why
459 individual variation matters to experimental biologists. *J. Exp. Biol.* **219**, 3832-3843.
460 (doi:10.1242/jeb.146712).
- 461 43. Hadfield J.D. 2010 MCMC methods for multi-response generalized linear mixed models: the
462 MCMCglmm R package. *J. Stat. Softw.* **33**, 1-22.
- 463 44. Dingemanse N.J., Dochtermann N.A. 2013 Quantifying individual variation in behaviour:
464 mixed-effect modelling approaches. *J. Anim. Ecol.* **82**, 39-54.
- 465 45. Nakagawa S., Schielzeth H. 2010 Repeatability for Gaussian and non-Gaussian data: a
466 practical guide for biologists. *Biol. Rev.* **85**, 935-956.

- 467 46. Araya-Ajoy Y.G., Mathot K.J., Dingemanse N.J. 2015 An approach to estimate short-term,
468 long-term and reaction norm repeatability. *Methods Ecol. Evol.* **6**, 1462-1473.

SUPPLEMENTARY INFORMATION

for: Roche DG, Berberi I, Dhane F, Lauzon F, Soeharjono S, Dakin R, Binning SA (in revision)
Slow improvement to the archiving quality of open datasets shared by researchers in ecology and evolution (in revision)

Table S1. The number of ecologists and evolutionary biologists in Canada's top-21 ranked universities according to the 2019 THE World University Rankings. Universities are ordered alphabetically. Indicated are the number of PIs in ecology and evolution at each university (E&E PIs), PIs with at least two journal articles published between Jan 2013-June 2019 containing a data availability statement and/or associated open data (Open data PIs), and PIs randomly selected for analysis in this study (Selected PIs).

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

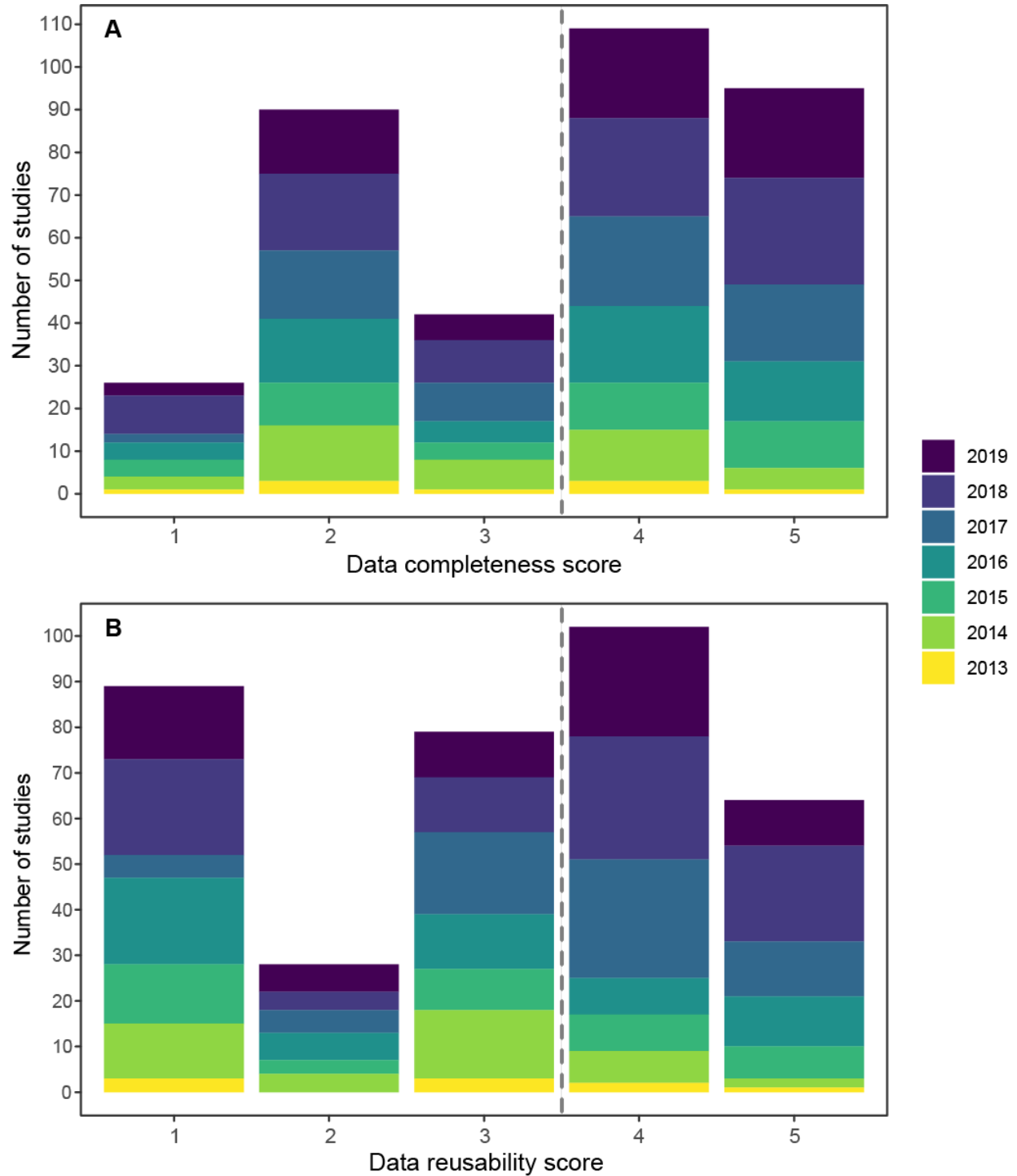


Fig. S1 Frequency distribution of the (A) completeness and (B) reusability scores for open datasets associated with 362 studies shared by 100 researchers between 2013-2019. A score of 5 indicates exemplary archiving, and a score of 1 indicates poor archiving (Table 1). Studies with scores of 3 or lower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

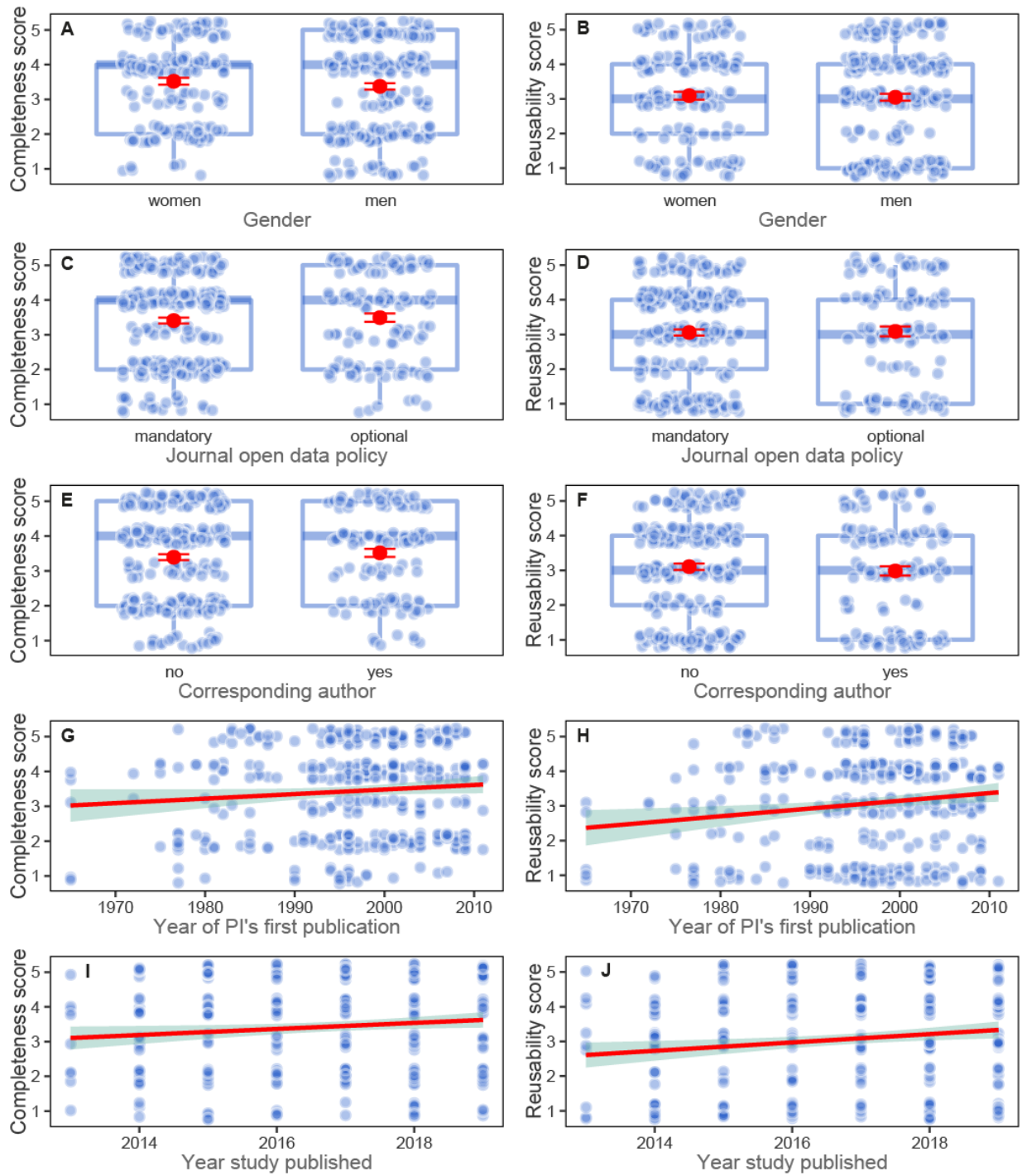


Fig. S2 The relationship between five independent variables and the completeness (A, C, E, G, I) and reusability score (B, D, F, H, J) of 362 open datasets shared by 100 principal investigators (PIs). Independent variables include: the gender of the PI, the journals' open data policy, whether the PI 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 the

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

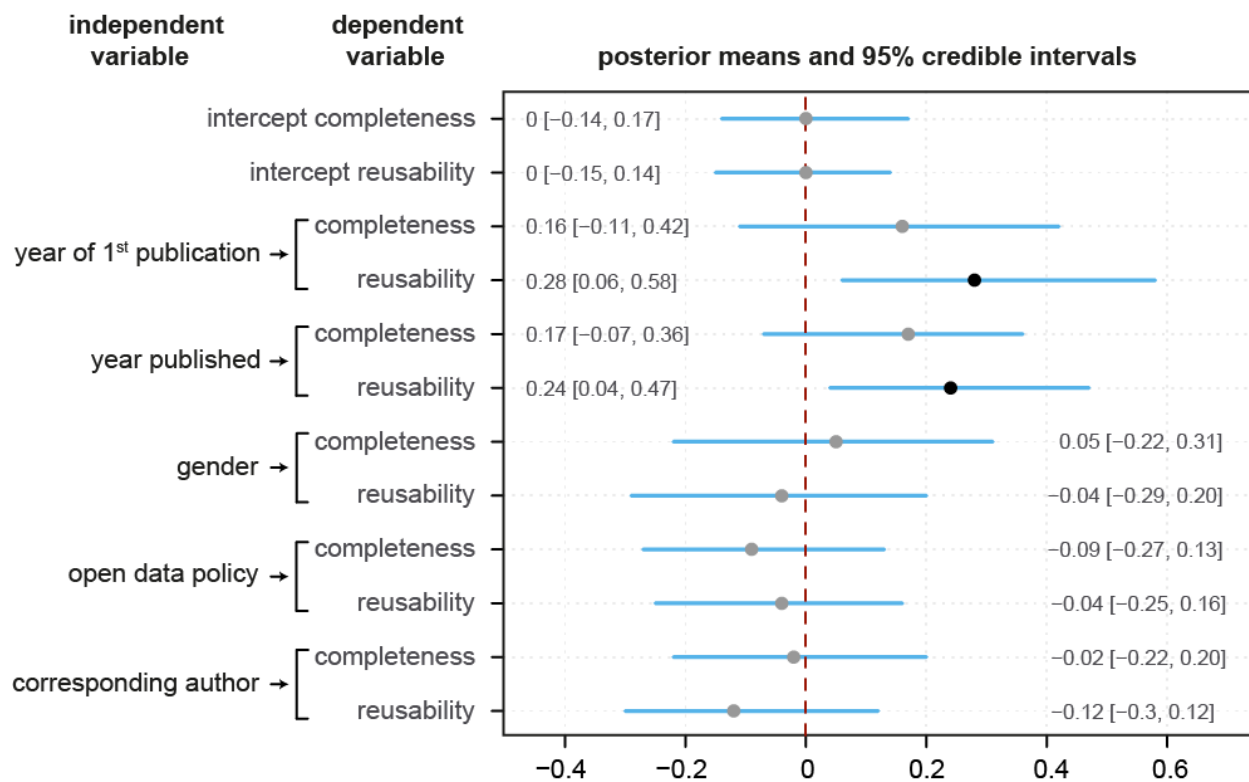


Fig. S3 Posterior means and 95% credible intervals from a Bayesian bivariate mixed-effects model investigating the effect of five independent variables on the completeness and reusability scores of 362 open datasets shared by 100 principal investigators (PIs). The prior was specified as a parameter-expanded prior (in contrast to an inverse-Wishart prior; see Fig. 2). 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.

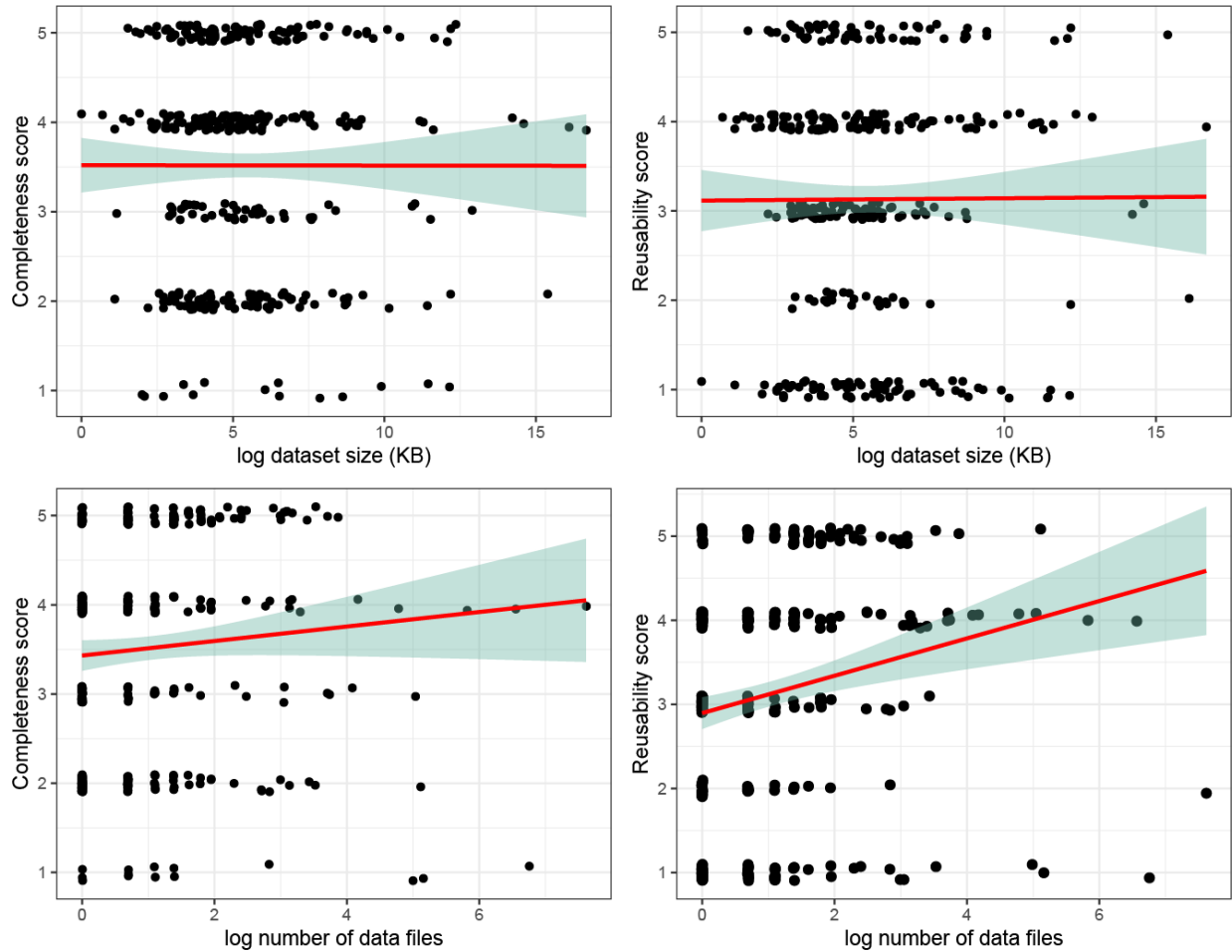


Fig. S4 The relationship between the completeness and reusability scores of open datasets and their complexity as estimated by the number and size (in KB) of archived files. Note that both variables used to estimate dataset complexity should be interpreted with caution as datasets were not archived using a standard file type (some file types [e.g. .xlsx] are inherently larger than others [e.g., .txt] and some files may contain multiple spreadsheets (e.g., .xlsx files can contain multiple tabs or spreadsheets as opposed to .txt or .csv files, which only contain one).