# A game theoretic analysis of research data sharing

Tessa E Pronk, Paulien H Wiersma, Anne van Weerden, Feike Schieving

While reusing research data has evident benefits for the scientific community as a whole, decisions to archive and share these data are primarily made by individual researchers. In this paper we analyse, within a game theoretical framework, how sharing and reuse of research data affect individuals who share or do not share their datasets. We construct a model in which there is a cost associated with sharing datasets whereas reusing such sets implies a benefit. In our calculations conflicting interests appear for researchers. Individual researchers are *always* better off not sharing and omitting the sharing cost, at the same time both sharing and not sharing researchers are better off if (almost) all researchers share. Namely, the more researchers share, the more benefit can be gained by the reuse of those datasets. We simulated several policy measures to increase benefits for researchers sharing or reusing datasets. Results point out that, although policies should be able to increase the rate of sharing researchers, and increased discoverability and dataset quality could partly compensate for costs, a better measure would be to directly lower the cost for sharing, or even turn it into a (citation-) benefit. Making data available would in that case become the most profitable, and therefore stable, strategy. This means researchers would willingly make their datasets available, and arguably in the best possible way to enable reuse.

# A GAME THEORETIC ANALYSIS OF RESEARCH DATA SHARING

Tessa E. Pronk, [1]*

Paulien H. Wiersma*

Anne van Weerden*

Feike Schieving[#]

*Utrecht University Library , Utrecht University, Heidelberglaan 3, Utrecht, the Netherlands

[#] Ecology and Biodiversity, Utrecht University, Padualaan 8, Utrecht, the Netherlands

[1] Corresponding author: T.E.Pronk@uu.nl

## Introduction

20

21 While sharing datasets has group benefits for the scientific community and society as a whole,

22 decisions to archive datasets are made by individual researchers. It is less obvious that the

23 benefits of sharing outweigh the costs for all individuals [Tenopir et al., 2011; Roche et al.,

24 2014]. Many researchers are reluctant to share their dataset publicly because of real or

25 perceived individual costs [Pitt and Tang, 2013]. This probably explains why sharing datasets  is

26 no daily practice [Roche et al., 2014], especially when compared to sharing knowledge and

27 information in the form of a scientific paper. Costs to individual researchers include time

28 investment, money, the chance of being scooped by others on any future publications on the

29 dataset, a chance that results from published papers will be over-scrutinized, misinterpretation

30 of data resulting in faulty conclusions [Atici et al., 2013], misuse [Bezuidenhout, 2013], and

31 possible infringement of the privacy of test subjects [Antman, 2014]. Also, datasets are

32 perceived as intellectual property and researchers simply do not want others to benefit from it

33 [Vickers, 2011].

34 In contrast, the act of sharing research data could have advantageous consequences.

35 Scientific outreach might be extended into other than the original research areas [Chao, 2011],

36 and researchers' reputations could grow by the publicity of good sharing practices, possibly

37 initiating new collaborations. In genetics [Botstein, 2010; Piwowar and Vision, 2013] it was

38 calculated that papers with open data were cited more than studies without the data available.

39 This citation advantage was also found in other disciplines like astronomy [Henneken E.A.,

40 2011; Dorch, 2012] and oceanography [Sears, 2011]. As citations to papers for many disciplines

41 are a the key metric by which impact of researchers is measured, this could mean a very

42 important incentive to researchers for sharing their data. Moreover, there is a tendency to

43 regard datasets as research output that can be used as a citeable reference or source in their

44 own right [Costello et al., 2013; Neumann and Brase, 2014]. For the field of oceanography it

45 was found that  datasets can be cited even more than most papers [Belter, 2014]. This would

46 mean that sharing datasets in the near future could have a direct positive influence on a

47 researcher's scientific impact.

48 On the other side of the coin, a researcher who reuses a dataset that was shared can

49 gain several advantages. Time is saved in not having to collect or produce the data, which can

50 be put to use to produce more papers. Papers can be enhanced with a comparison or meta-

51 analysis based on an extra dataset. If the added dataset merits publication in a higher impact

52 journal, the paper could be cited more often. In more general terms, the scientific community

53 can benefit from reuse of datasets. Sharing data enables open scientific inquiry, encourages

54 diversity of analysis and opinion, promotes new research, facilitates the education of new

55 researchers, enables novel applications to data not envisioned by the initial investigators,

56 permits the creation of new datasets when data from multiple sources are combined, and

57 provides a basis for new experiments [Ascoli, 2007; Kim, 2013; Pitt and Tang, 2013]. It also is a

58   way to prevent scientific fraud; with the dataset provided one should be able to reproduce
59   scientific results.
60        To summarize, data sharing implies costs and/or benefits for the individual researcher,
61   but are of clear benefit to researchers that reuse the dataset, and to the scientific community
62   as a whole. In this context, the problem of data sharing can be studied as a game theoretic
63   problem. The strength of game theory lies in the methodology it provides for structuring and
64   analysing problems of strategic choice. The players, their strategic options, the external factors
65   of influence on those decisions, all have to be made explicit. With the model we show how
66   research data sharing fits the definition of a typical 'tragedy of the commons', in which
67   cooperating is the best strategy but cheating is the evolutionary stable strategy. In addition, we
68   assess measures for altering costs and benefits with sharing and reuse and analyse how each
69   measure would turn the balance towards *more* sharing and *more* benefits from sharing,
70   benefitting the community, society and the individual researcher.
71

## Methods
72

73

**A Model for Impact**
74
75   We assume a community of researchers who publish papers. We consider two types of
76   researchers: those sharing and not sharing research data associated with those papers. We
77   make the simplifying assumption that the goal for both types of researchers is to perform well
78   by making a significant contribution to science, i.e. to have a large impact on science. We
79   assume that produced papers, $P_s$ for sharers and $P_{ns}$ for non sharers, create impact by getting
80   cited a number of times $c$. We assume c is constant, which means we do not distinguish
81   between low and highly cited papers. To increase their performance, researchers need to be
82   efficient, i.e. they should try to minimize the time spent on producing a paper, so more papers
83   can be produced within the same timeframe. Papers from which the dataset is shared gain an
84   extra citation advantage, increasing the impact of that paper by a factor $b$. In our model we
85   consider only papers with a dataset as a basis, i.e. no review or opinion papers. So, the
86   performance of researchers is expressed as an impact rate, in terms of citations per year, i.e.
87   the impact for sharing and non-sharing researchers is defined as

88   $$E_s = P_s \cdot c \cdot (1 + b) \qquad E_{ns} = P_{ns} \cdot c \qquad\qquad\qquad (1)$$

89   From the above expressions it is clear that the difference in impact between sharing and not
90   sharing researchers is to a large extent dependent on the number of publications $P$ per year.
91   These publications can be expressed in terms of an average time to write a paper $T_s$ for sharers
92   and $T_{ns}$ for not sharers.

93 $\quad P_s = \dfrac{1}{T_s} \qquad P_{ns} = \dfrac{1}{T_{ns}}$ $\hfill$ (2)

94 The time $T$ consists of several elements that we make explicit here. Each paper costs time $t_a$ to

95 produce. Producing the associated dataset costs a certain time $t_d$. Sharing a dataset implies a

96 time cost $t_c$. We do not distinguish between large and small efforts to prepare a dataset for

97 sharing; all datasets take the same amount of time. We assume there is a certain probability $f$

98 to find an appropriate dataset for a paper from the pool of shared datasets $X$, in which case the

99 time needed to produce a dataset $t_d$ is avoided. We do acknowledge that some time is needed

100 for a good 'getting to know' the external dataset and to process it, resembled in the time cost

101 $t_r$. We calculate the time to produce a paper by

102 $\quad T_s = t_a + \dfrac{t_d}{1+f\cdot X} + \left( t_r - \dfrac{t_r}{1+f\cdot X} \right) + t_c \qquad T_{ns} = t_a + \dfrac{t_d}{1+f\cdot X} + \left( t_r - \dfrac{t_r}{1+f\cdot X} \right)$ $\hfill$ (3)

103 In these formulae, the pool of available datasets X determines the value of the terms with $t_d$

104 and $t_r$. When $X$ is close to zero, the term with $t_d$ approaches $t_d$. This implies that everybody has

105 to produce their own dataset with time cost $t_d$. In contrast, when $X$ is very large the term

106 approaches zero, implying almost everyone can reuse a dataset and almost no time is spent in

107 the community to produce datasets. Between these two extremes, the term first rapidly

108 declines with increasing $X$ and then ever more slowly approaches zero (see the plots in the last

109 column in the figure in Appendix 2). This is under the assumption that at a small number of

110 available datasets, adding datasets will have a profound influence on the reuse possibilities. If

111 datasets are already superfluous, adding extra datasets will have less influence on the reuse

112 rate. The term representing the effort to reuse a paper $t_r$ works opposite to the term

113 representing $t_d$. When $X$ is close to zero, the term approaches zero, implying nobody spends

114 time to prepare a set for reuse. When $X$ is very large the term approaches $t_r$; everyone spends

115 this time because everyone has found a set for reuse.

116 $\qquad$ While the pool of datasets $X$ determines the values of the terms with $t_d$ and $t_r$ and with

117 that the number of shared datasets, at the same time the shared datasets accumulate in the

118 pool of shared datasets $X$. To come to a specification of this pool size $X$ we formulate a

119 differential equation for the pool size. A change in the pool of available, shared datasets $X$

120 depends on adding datasets belonging to papers $P_s$ from sharing researchers $Y_s$, minus the

121 decay $q_x \cdot X$ of the datasets. Such a decay rate could be a result from a fixed storage time after

122 which datasets would be disposed of or by a loss of data value, for instance by outdated

123 techniques.

124 $\quad \dfrac{dX}{dt} = Y_s \cdot P_s - q_x \cdot X$ $\hfill$ (4)

125 Using Formula (2) and (3) with the system at steady state i.e. $dX/dt = 0$, the pool size $X$ as

126 function of the publication parameters and the size of the group of sharing researchers is given

127 by

128 $$X = \frac{-\left(q_x\left(t_a + t_c + t_d\right) - Y_s f\right) + \sqrt{\left(q_x\left(t_a + t_c + t_d\right) - Y_s f\right)^2 - 4\left(q_x \cdot f\left(t_a + t_c + t_r\right)\right) \cdot \left(-Y_s\right)}}{2\left(q_x f\left(t_a + t_c + t_r\right)\right)} \quad (5)$$

129 (Formula (5) is derived in Appendix 1). So, for each parameter setting, we calculate $X$, and

130 consequently, we calculate the impact in terms of citation rates $E_s$ and $E_{ns}$ with Formulae (1-3).

131 Table 1 gives the default parameter settings that we use for our simulations.

132

133 **An Individual Based Model**

134 In addition to the model for impact we set up an individual based model to assess the impact

135 for individual researchers depending on their personal publication rate, sharing and reuse

136 habits, rather than to work with averages. We use the 'model for impact' as a basis for the

137 calculations and then assign characteristics to individuals. First, a publication rate $P_r$ per

138 researcher is assigned at random to individual researchers. $P_r$ is based on the distribution as

139 seen in Figure 1, fitted with the function

140 $$P_r = Y \cdot e^{-(t_a + t_d)} \qquad (6)$$

141 As a next step we introduce parameters that have to do with sharing. The percentage of sharing

142 researchers is a fixed parameter in this model. The researchers sharing type is assigned at

143 random to individuals. The actual reuse of a dataset, based on the probability to find an

144 appropriate dataset for a paper, is assigned at random to publications. The portion of papers $R$

145 for which an appropriate dataset for reuse is found is calculated as

146 $$R = 1 - \frac{1}{1 + f \cdot X} \qquad (7)$$

147 We now have a mix of individual researchers that share or do not share, find a dataset for reuse

148 or not for any of their papers, and publish different number of papers in a year. Based on the

149 parameters in Table 1 we assign costs and benefits with these traits. These factors determine

150 the performance of researchers in terms of impact by citations.

151     To determine the publication rate distribution in Figure 1, we sampled the bibliographic

152 database Scopus. We selected the first four papers for each of the 26 subject areas in Scopus-

153 indexed papers, published in 2013. If a paper appeared within the first four in more than one

154 subject area, it was replaced by the next paper in that subject area. For each of the selected

155 papers we noted down all authors and checked how many papers each author (co-) authored in

156 total in 2013. We came to 366 unique authors in our selected papers. Authors that were

157 ambiguous, because they seemingly published many papers, were checked individually and

158 excluded if it was a group of authors publishing under the same name with different affiliations

159    between the papers. For the data see [Pronk et al., 2015]. This distribution, based on our

160    sampling, implies that most researchers publish one- and a few researchers publish many

161    papers in a given year. We fitted an exponential distribution through the sampled population

162    (Formula 6). The average for the distribution is close to three papers per researcher in a given

163    year.

164

165    **Simulations**

166    For the R-scripts to generate the plots for all simulations, see [Pronk et al., 2015].

167         We start with a set of simulations regarding performances per sharing type, with the

168    model for impact. We calculate the impact for the two types of researchers  over a range of

169    sharing from zero to a hundred percent of all researchers. In addition to the default values (see

170    Table 1), we change parameters to assess their influence on the publication rate and associated

171    impact by citations for sharing and not sharing researchers. In Table 2 we list the parameters

172    changed in the simulations and a score of the measures that would have these effects in a 'real

173    world' scientific community [Chan et al., 2014].

174         To have a closer look on individual performance, we perform the same set of

175    simulations with the individual based model. For each setting we calculate the difference

176    between the publication rate assigned in Formula (6) at no costs or benefits with sharing or

177    reuse, and a new, calculated publication rate based on sharing and reuse traits per researcher

178    under the assumption that half of the researchers share. So, again we change the parameters in

179    Table 2 and assess their influence, as in the first simulation.

180         We end by zooming out to community performance with the model for impact. We

181    calculate the average impact over all researchers in the community, now at more extreme

182    settings of the citation benefit $b$ and in a second simulation at even higher cost $t_c$ for preparing

183    a dataset for sharing. This is to provide a broader range of results. Citation benefit $b$ and the

184    sharing rate are changed within their range in one hundred equal steps.

185

186    # Results

187    Shown in Figure 2 are the simulations with the model for impact (Formulae 1-5). The simulation

188    in (a) is at default parameter values (Table 1). In (b-f) we simulated measures to improve upon

189    impact. There are two important observations. First, in all (but the last) subfigure of Figure 2 (a-

190    e) the average impact of not sharing researchers exceeds that of sharing researchers

191    irrespective of how many sharing researchers there are. This means that *not sharing* is the best

192    option, at all percentages sharing researchers. In this scenario it would be logical if all individual

193    researchers would choose not to share and eventually end up getting the average impact by

194    citations depicted at zero percent sharing. So we see here a classical example of the tragedy of

195    the commons or prisoners dilemma phenomenon. What is important to note though is that the

196    measures in (b) (c) (d) and (e) ascertain a key effect when compared to the default in (a). The

197 average impact of sharing researchers at the highest percentage sharing researchers  (straight
198 horizontal light-grey line; stripes) is increasingly higher with the measures than the average
199 impact for not sharing researchers at zero percentage sharing researchers (straight horizontal
200 dark-grey line; dots-stripes). Should a policy enforce the sharing, or all would agree to
201 cooperate and share, a higher gain is achieved than in the case that researchers would all
202 choose not to share. This illustrates the conflicting interest for individual researchers, who are
203 better off not sharing, while they would do better if all of them did share. Subfigure (f) of Figure
204 2 shows the potential of the citation benefit with sharing. In the picture it is profitable to share
205 at low sharing rates, and profitable not to share at high sharing rates, leading to a stable
206 coexistence of sharing and not sharing researchers. This means that the community would exist
207 of researchers from both strategies. Hypothetically, should the citation benefit be even higher,
208 the sharing strategy would outperform the not sharing strategy at all sharing percentages.
209 Researchers would in this case choose to share even without measures to promote sharing,
210 simply because it directly increases their impact.
211 　　　　Second, it can be noted that in some subfigures of Figure 2 (a, b, c, e) the average
212 citations are the highest at intermediate sharing. This means that if sharing increases further, it
213 has a detrimental effect on average community impact. This is because the model is formulated
214 in Formula (3) in a way that total costs for sharing increase for the community as more
215 researchers share, whereas total benefits cease to increase at high sharing rate. The extra
216 datasets do not contribute much to the benefits, or in other words, the research community
217 has become saturated with datasets. Compared to the average community citations, which are
218 highest at intermediate sharing, for both sharing and not sharing researchers the highest
219 impact by citations is at the point at which everyone is sharing.
220 　　　　Results from the individual based in Figure 3 model show that the individual researchers
221 have various gains depending on their publication rate, reuse, and dataset sharing habits. In (a)
222 are the gains and losses in impact, at default parameter values (Table 1). In (b-f) we simulated
223 measures to improve gains or limit losses. A possible desired effect of sharing of datasets would
224 be that every individual researcher can benefit, sharing or not sharing. It can be observed that
225 in Figure 3 (a-e) most of the sharing researchers have lower benefits or even costs compared to
226 not sharing researchers. This logically is in line with the lower averages for sharing researchers
227 in Figure 2. Also, it can be noted in all subfigures of Figure 3 that there are always sharing
228 researchers that do not benefit from the availability of datasets by the reuse of datasets. These
229 researchers were not (fully) able to compensate for the cost to share their data. It is notable
230 that in (b) individual researchers are left with lower costs than in (c). This is because in (b) the
231 probability of finding an appropriate dataset for reuse $f$ is set higher, compensating the sharing
232 costs for many of the researchers. In (c) the time cost $t_r$ with reuse per paper is lower,
233 benefitting only those few researchers that do find a reusable set. In (d) the lowering of the
234 time cost $t_c$ for preparing a dataset for sharing improves the situation for *all* researchers

235     compared to the default in (a), but still some researchers are not fully compensated. In (e) the
236     introduction of the citation benefit *b* does not help much to improve the benefits for sharing
237     researchers. Only when in (f) a substantial citation benefit *b* is introduced for sharing
238     researchers, the costs associated with sharing are (more than) compensated for, for all sharing
239     researchers.
240         When simulating community impact in Figure 4 (a) and (b) it can be seen that, as the
241     benefits *b* for sharing increase towards the right of the plot, the average community impact
242     increasingly starts to rise with more sharing in both plots. Even the drop after the initial
243     increase at increased sharing caused by the datasets saturation is eventually compensated for
244     with the increase of the citation benefit with sharing. In subfigure (b) at the left side of the plot,
245     without a citation benefit and with the very high cost for sharing $t_c$, there appears an alarming
246     effect. At these parameter values the average impact becomes lower at high sharing than at no
247     sharing at all. Policies increasing sharing would, if successful, in this case backfire and reduce
248     scientific community impact.
249

## Discussion

251     We analysed the effect of sharing and not sharing research data on scientific community
252     impact. We found that there is a conflicting interest for individual researchers, who are *always*
253     better off not sharing and omitting the sharing cost while they would have higher impact when
254     sharing as a community. With our model we assessed some measures to improve the costs and
255     benefits with sharing and reuse of data, to make most researchers profit from the sharing of
256     datasets. We simulated policies to increase sharing, measures to stimulate reuse by reducing
257     reuse costs or increasing discoverability of datasets, and measures to stimulate sharing by
258     lowering costs associated with sharing or introducing a citation benefit with each shared
259     dataset. These simulations concretize the notion in literature that improving spontaneous
260     participation in sharing datasets will require lowering costs and/or increasing benefits for
261     sharing [Smith, 2009; Roche et al., 2014] and values different measures to do so.
262         A policy is a straightforward measure to increase community impact simply by enforcing
263     higher percentages of sharing researchers. Moreover, policies are pivotal for establishing
264     acceptable data sharing practices and community-level standards. Such policies can be
265     enforced on the level of institutions, funders, or journals. In the model these do increase
266     community impact, as long as the community is not already saturated with datasets. In real life,
267     at least for journals,  policies have not been enough to convince researchers to actually make
268     their dataset publicly available [Wicherts et al., 2006; Savage and Vickers, 2009; Alsheikh-Ali et
269     al., 2011; Wicherts and Bakker, 2012; Vines et al., 2013]. This could be exemplary for the
270     reluctance of individual researchers to share datasets because of real, or perceived costs. The
271     inequality in costs between sharing researchers and not sharing researchers remains with
272     mandated sharing, and the researcher that does not share a dataset but does reuse a dataset

273  will have the highest impact compared to all others. Of course there are many factors for
274  researchers to decide to share data or not, but simply said this could predispose a researcher
275  towards not sharing. The 'reuse-don't share' strategy is a true current sentiment towards using:
276  according to a survey in 2011 of about 1,300 scientists, more than 80 percent said they would
277  use other researchers' datasets but only few wanted to make their dataset available to others,
278  for a variety of reasons [Tenopir et al., 2011; Fecher et al., 2015].

279  Stimulating reuse by reducing reuse costs or increasing discoverability of datasets in the
280  model increases average community impact, though not equally for all individuals within the
281  community. Only the researchers that actually reuse a dataset profit from these measures, and
282  the costs for those who share, although partly compensated, still exist.  Again, although helpful,
283  the inequality in costs between sharing and not sharing researchers is not addressed with such
284  measures.

285  A direct reduction of the time costs with sharing a dataset in our model improved the
286  situation for all sharing researchers. Only a small inequality between sharing and not sharing
287  researchers remained. The best solution is however to introduce a 'citation benefit' for papers
288  with the dataset shared, to directly balance the costs of sharing individuals. The citation benefit
289  in real life can not only come from increased citations to the paper [Botstein, 2010; Sears, 2011;
290  Dorch, 2012; Piwowar and Vision, 2013] but also from citations to the shared dataset itself
291  [Costello et al., 2013; Belter, 2014; Neumann and Brase, 2014]. With a relatively high citation
292  benefit, sharing datasets even becomes more profitable than not sharing, at any percentage of
293  sharing researchers. Sharing then is not only optimal for maximizing community impact, but
294  also for the individual researcher.

295  All in all, enhancement of the citation benefit would bring about better incentives to
296  share datasets than imposing an obligation to share by funders, institutes or journals, or partly
297  compensating for costs by enabling reuse. Better incentives arguably also lead to better sharing
298  practices as researchers would strive to present their dataset as such that its reuse potential is
299  optimal.

300  All models come with simplifications and assumptions. A central assumption of the
301  model is the gain of scientific impact by citations to papers, and implicitly datasets. For some
302  communities the concept of impact by citations is less applicable overall [Krell, 2002]. These fall
303  outside the scope of this model. It also should be noted that there are other ways to count
304  scientific impact such as Altmetrics [Roemer and Borchardt, 2012]. Additionally, we derived
305  general phenomena for the scientific community, whereas (perceived) costs and benefits with
306  sharing will differ between scientific communities [Vickers, 2011; Tenopir et al., 2011; Kim,
307  2013] and attitudes towards sharing can differ largely between disciplines [Kirwan, 1997; Huang
308  et al., 2012; Pitt and Tang, 2013; Anagnostou et al., 2013]. This means that the measures taken
309  to make sharing worthwhile will have to differ in their focus in each scientific community
310  [Borgman et al., 2007; Acord and Harley, 2013]. To apply the current model to any specific

311 situation or community, parameter values for that community should be carefully determined
312 and, where necessary, the model should be adjusted or expanded. Additional factors that may
313 influence the outcome of this model and that could possibly be incorporated in community
314 specific versions or future refinements of this model include: differences in quality of papers
315 leading to differences in citation rates, heterogeneity in the costs of sharing (small and easy
316 versus big and complicated datasets to document), heterogeneity in the contribution of a
317 papers' dataset to the available pool of datasets, introducing and allowing for heterogeneity in
318 search time for datasets, feedback between the number of times a dataset is reused and the
319 citation benefit for that dataset. A focal point to assess in the current model would also be the
320 pool of available datasets. What is the relation between available datasets and reuse rate for
321 researchers, do these datasets overlap in content, will all new datasets contribute to science,
322 does the pool become saturated, are all datasets reused, what is the decay rate of datasets in
323 the pool for that specific community?

324 Lastly, it is clear that not all data can or should be made fully or immediately publicly
325 available for a variety of practical reasons (e.g., lack of interest, sheer volume and lack of
326 storage, cheap-to-recreate data, high time costs to prepare the data for reuse, the wish to
327 publish later perhaps, patents pending, privacy sensitive data) [Kim, 2013; Cronin, 2013]. With
328 our simulations we show that if costs for sharing are too high relative to the benefits of reuse,
329 in theory sharing policies to increase sharing could even backfire and reduce scientific
330 community impact. It should be carefully considered whether the alleged benefits of storage
331 for the scientific community will outweigh the costs for each data type and set. For easily
332 obtainable data such as the data underlying this paper, recreating it is probably cheaper than
333 storing and interpreting the datasheet.

334 In conclusion, we performed a game-theoretic analysis to provide structure and to
335 analyse problems of strategic data sharing. In the simulations there appeared a conflicting
336 interest for individual researchers, who are *always* better off not sharing and omitting the
337 sharing cost, while they are ultimately better off all sharing as a community. Although policies
338 are indispensable and should be able to increase the rate of sharing researchers, and increased
339 discoverability and dataset quality could partly compensate for costs, a better measure to
340 promote sharing would be to lower the cost for sharing, or even turn it into a (citation-) benefit.

341

347

348 REFERENCES

349    Acord, S. K. and D. Harley (2013), Credit, time, and personality: The human challenges to

350    sharing scholarly work using Web 2.0, New Media and Society, 15(3), 379-397,

351    doi:10.1177/1461444812465140.

352    Alsheikh-Ali, A. A., W. Qureshi, M. H. Al-Mallah, and J. P. Ioannidis (2011), Public availability of

353    published research data in high-impact journals, PLoS One, 6(9), e24357,

354    doi:10.1371/journal.pone.0024357.

355    Anagnostou, P., M. Capocasa, N. Milia, and G. D. Bisol (2013), Research data sharing: Lessons

356    from forensic genetics, Forensic. Sci. Int. Genet., 7(6), e117-9, doi:10.1016/j.fsigen.2013.07.012.

357    Antman, E. (2014), Data sharing in research: benefits and risks for clinicians, BMJ, 348, g237,

358    doi:10.1136/bmj.g237.

359    Ascoli, G. A. (2007), Successes and rewards in sharing digital reconstructions of neuronal

360    morphology, Neuroinformatics, 5(3), 154-160, doi:NI:5:3:154.

361    Atici, L., S. W. Kansa, J. Lev-Tov, and E. C. Kansa (2013), Other People's Data: A Demonstration

362    of the Imperative of Publishing Primary Data, J. Archaeol. Method and Theory, 20(4), 663-681,

363    doi:10.1007/s10816-012-9132-9.

364    Belter, C. W. (2014), Measuring the value of research data: a citation analysis of oceanographic

365    data sets, PLoS One, 9(3), e92590, doi:10.1371/journal.pone.0092590.

366    Bezuidenhout, L. (2013), Data sharing and dual-use issues, Sci. Eng. Ethics, 19(1), 83-92,

367    doi:10.1007/s11948-011-9298-7.

368   Borgman, C. L., J. C. Wallis, and N. Enyedy (2007), Little science confronts the data deluge:

369   Habitat ecology, embedded sensor networks, and digital libraries, Int. J. Digital Libr., 7(1-2), 17-

370   30, doi:10.1007/s00799-007-0022-9.


371   Botstein, D. (2010), It's the data! Mol. Biol. Cell, 21(1), 4-6.


372   Chan, A. W., F. Song, A. Vickers, T. Jefferson, K. Dickersin, P. C. Gotzsche, H. M. Krumholz, D.

373   Ghersi, and H. B. van der Worp (2014), Increasing value and reducing waste: addressing

374   inaccessible research, Lancet, 383(9913), 257-266, doi:10.1016/S0140-6736(13)62296-5.


375   Chao, T. C. (2011), Disciplinary reach: Investigating the impact of dataset reuse in the earth

376   sciences, Proc. ASIST Ann. Meet., 48, doi:10.1002/meet.2011.14504801125.


377   Costello, M. J., W. K. Michener, M. Gahegan, Z. -. Zhang, and P. E. Bourne (2013), Biodiversity

378   data should be published, cited, and peer reviewed, Trends Ecol. Evol., 28(8), 454-461,

379   doi:10.1016/j.tree.2013.05.002.


380   Cronin, B. (2013), Thinking about data, J. Am. Soc. Inf. Sci. Technol., 64(3), 435-436,

381   doi:10.1002/asi.22928.


382   Dorch, B. (2012), On the Citation Advantage of linking to data. hprints-00714715, version 2,

383   Hprints, http://hprints.org/hprints-00714715.


384   Fecher, B., S. Friesike, M. Hebing, S. Linek, and A. Sauermann (2015), A Reputation Economy:

385   Results from an Empirical Survey on Academic Data Sharing. DIW Berlin Discussion Paper No.

386   1454., doi:http://dx.doi.org/10.2139/ssrn.2568693.

387  Henneken E.A., A. A. (2011), Linking to data - effect on citation rates in astronomy.

388  arXiv:1111.3618v1, http://arxiv.org/abs/1111.3618.

389  Huang, X., B. A. Hawkins, F. Lei, G. L. Miller, C. Favret, R. Zhang, and G. Qiao (2012), Willing or

390  unwilling to share primary biodiversity data: Results and implications of an international survey,

391  Conserv. Lett., 5(5), 399-406, doi:10.1111/j.1755-263X.2012.00259.x.

392  Kim, J. (2013), Data sharing and its implications for academic libraries, New Libr. World,

393  114(11), 494-506, doi:10.1108/NLW-06-2013-0051.

394  Kirwan, J. R. (1997), Making original data from clinical studies available for alternative analysis,

395  J. Rheumatol., 24(5), 822-825.

396  Krell, F. T. (2002), Why impact factors don't work for taxonomy. Nature, 415(6875), 957.

397  Neumann, J. and J. Brase (2014), DataCite and DOI names for research data, J. Comput. Aided

398  Mol. Des., doi:10.1007/s10822-014-9776-5.

399  Pitt, M. A. and Y. Tang (2013), What should be the data sharing policy of cognitive science? Top.

400  Cogn. Sci., 5(1), 214-221, doi:10.1111/tops.12006.

401  Piwowar, H. A. and T. J. Vision (2013), Data reuse and the open data citation advantage, PeerJ,

402  1, e175, doi:10.7717/peerj.175.

403     Pronk, T. E., P. H. Wiersma, and v. A. Weerden (2015), Replication data for: GAMES WITH

404     RESEARCH DATA SHARING", http://hdl.handle.net/10411/20328 V4 [Version], Dutch Dataverse

405     Network.

406     Roche, D. G., R. Lanfear, S. A. Binning, T. M. Haff, L. E. Schwanz, K. E. Cain, H. Kokko, M. D.

407     Jennions, and L. E. Kruuk (2014), Troubleshooting public data archiving: suggestions to increase

408     participation, PLoS Biol., 12(1), e1001779, doi:10.1371/journal.pbio.1001779.

409     Roemer, R. C. and R. Borchardt (2012), From bibliometrics to altmetrics. A changing scholarly

410     landscape. College & Research Libraries News, 73, 596.

411     Savage, C. J. and A. J. Vickers (2009), Empirical study of data sharing by authors publishing in

412     PLoS journals, PLoS ONE, 4(9), doi:10.1371/journal.pone.0007078.

413     Sears, J. R. L. (2011), Data Sharing Effect on Article Citation Rate in Paleoceanography, Fall

414     Meeting, AGU, San Francisco, Calif., 5-9 Dec., Abstract IN53B-

415     1628,http://adsabs.harvard.edu/abs/2011AGUFMIN53B1628S.

416     Smith, V. S. (2009), Data publication: towards a database of everything, BMC Res. Notes, 2, 113-

417     0500-2-113, doi:10.1186/1756-0500-2-113.

418     Tenopir, C., S. Allard, K. Douglass, A. U. Aydinoglu, L. Wu, E. Read, M. Manoff, and M. Frame

419     (2011), Data sharing by scientists: practices and perceptions, PLoS One, 6(6), e21101,

420     doi:10.1371/journal.pone.0021101.

421    Vickers, A. J. (2011), Making raw data more widely available, BMJ, 342, d2323,

422    doi:10.1136/bmj.d2323.

423    Vines, T. H., R. L. Andrew, D. G. Bock, M. T. Franklin, K. J. Gilbert, N. C. Kane, J. S. Moore, B. T.

424    Moyers, S. Renaut, D. J. Rennison, T. Veen, and S. Yeaman (2013), Mandated data archiving

425    greatly improves access to research data, FASEB J., 27(4), 1304-1308, doi:10.1096/fj.12-218164.

426    Wicherts, J. M. and M. Bakker (2012), Publish (your data) or (let the data) perish! Why not

427    publish your data too? Intelligence, 40(2), 73-76, doi:10.1016/j.intell.2012.01.004.

428    Wicherts, J. M., D. Borsboom, J. Kats, and D. Molenaar (2006), The poor availability of

429    psychological research data for reanalysis, Am. Psychol., 61(7), 726-728, doi:10.1037/0003-

430    066X.61.7.726.

431

**Table 1**(on next page)

Parameters, variables, and their values.

Table 1. Overview of parameters, variables, and their standard values used in the model. Grey rows indicate the parameters that are varied in the model to assess their influence (examples for real-world measures to change these are explained in Table 2).

1  # Table 1.

2

3  Table 1. Overview of parameters, variables, and their standard values used in the model. Grey rows indicate the
4  parameters that are varied in the model to assess their influence (examples for real-world measures to change
5  these are explained in Table 2).

6

| Parameter | Meaning | Value | Source | Unit |
|---|---|---|---|---|
| $t_a$ | Time-cost to produce a paper | 0.13 | Derived: $t_a+t_d$ amount to 121 days; leading to ~3 papers a year (similar to the average in Figure 1) | Year/ Paper |
| $t_d$ | Time-cost to produce a dataset | 0.2 | Derived: $t_a+t_d$ amount to 121 days; leading to ~3 papers a year (similar to the average in Figure 1) | Year/ Paper |
| $t_c$ | Time-cost to prepare a dataset for sharing | 0.1 | Estimated: 36.5 days | Year/ Paper |
| $t_r$ | Time-cost to prepare a dataset to reuse | 0.05 | Estimated: 18.25 days | Year/ Paper |
| $q_x$ | Decay rate of shared datasets | 0.1 | Derived: based on a storage time of 10 years | 1 / Year |
| $b$ | Citation benefit (sharing researcher) | 0 | Estimated: percent extra citations | Percent |
| $f$ | Probability to find an appropriate dataset | 0.00001 | Fitted | 1 / Dataset |
| $c$ | Citations per paper produced | 3.4 | Derived: approximate from 'baselines'; average citation rate by year three, Thompson Reuters | Citation / Paper |
| **State Variables** | **Meaning** | **Value** | | **Unit** |
| $E$ | Efficiency of researchers | See formula (1) | Calculated | Citation / Year |
| $P$ | Number of papers | See formula (2) | Calculated | Paper / Year |
| $T$ | Time for a publication | See formula (3) | Calculated | Year / Paper |
| $X$ | Pool of shared datasets | See formula (5) | Calculated | Dataset |
| $Y$ | Number of researchers | 10000 | Defined | n.a. |

7

8

**Table 2**(on next page)

Changed parameters and associated measures

Table 2. Overview of considered parameters determining reuse and sharing habits of researchers, and possible measures to improve these in a realistic setting.

1   Table 2.

2

3   Table 2. Overview of considered parameters determining reuse and sharing habits of researchers, and possible

4   measures to improve these in a realistic setting.

5

| Parameters investigated in the model | Possible associated measures to improve this |
|---|---|
| Time '$t_r$' spent to assess and include an external dataset | • Improve data quality, for instance by the use of data journals [Costello et al., 2013; Atici et al., 2013; Gorgolewski et al., 2013], or peer review of datasets (i.e. a 'comment' field in data repositories).<br>• Offer techniques or tools for easy assessment of dataset quality [Eijssen et al., 2013], faster pre-processing or data cleaning (i.e. 'OpenRefine' or 'R statistical language'). |
| Chance '$f$' to find an external dataset | • Harvest databases through data portals to reduce 'scattering' of datasets.<br>• Standardization of metadata and documentation<br>• Advanced community and project-specific databases<br>• Library assistance in finding and using appropriate datasets. |
| Time '$t_c$' associated with sharing of research data | • Offer a good storing & sharing IT infrastructure.<br>• Assistance with good data management planning at the early stages of a research project. |
| Benefit in citation per paper '$b$' associated with sharing of research data | • Provide a permanent link between paper and dataset.<br>• Increase attribution to datasets by citation rules .<br>• Establish impact metrics for datasets. |
| Percentage of scientists sharing their research data | • Promote sharing by a top down policy from an institute, funder, or journal.<br>• Promote sharing bottom up by offering education on the benefits of sharing, to change researchers' mind set. |

6

**Figure 1**(on next page)

Publication distribution

Figure 1. The sampled (bars) and fitted (line) distribution of published papers per researcher in a given year, in this case 2013. For reasons of visualisation the distribution is shown up to thirty publications, whereas the sampling sporadically included more publications per researcher. The fitted line is used as the published papers' distribution for the simulated community.
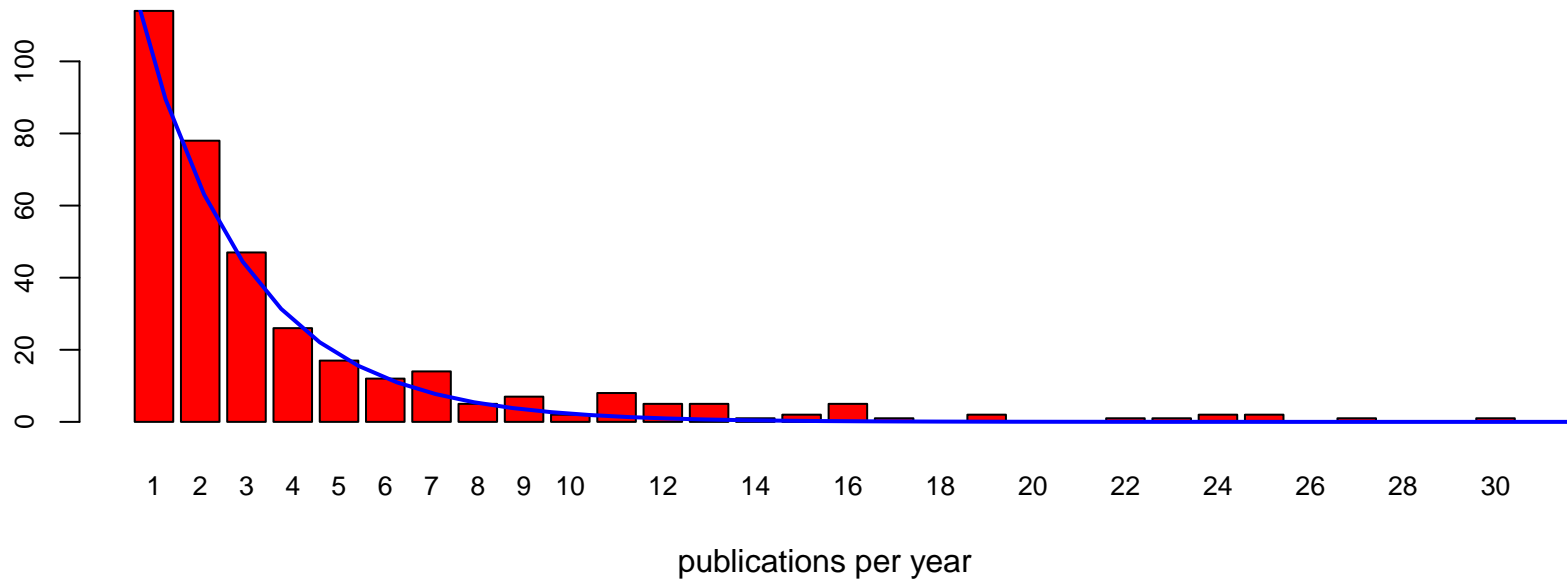
publications per year

# Figure 2 (on next page)

Impact per sharing type

Figure 2. Citations ('impact') per year for researchers sharing and not sharing, at different percentages of sharing researchers. The simulations are done at parameter settings a) default (see Table 1), b) default but with $f$ increased threefold c) default but with $t_r$ decreased threefold d) default but with $t_c$ decreased threefold e) default but with $b$ set to 0.1 f) default but with $b$ set to 0.4. The curved light-grey line depicts the impact of the sharing researchers. The curved dark-grey line depicts the impact of the not sharing researchers. The thin dotted curved black line is the averaged community impact. The straight black vertical dotted line depicts the percentage of sharing researchers at which community impact is maximized. The straight horizontal lines respectively depict the impact at zero percent researchers sharing (dark-grey line; dots-stripes) and hundred percent sharing researchers (light-grey line; stripes).

## Figure 3(on next page)

Individual gains with sharing

Figure 3. Gains from sharing in number of citations per individual researcher. These gains are calculated for the situation with fifty percent sharing researchers compared to the same situation without sharing researchers. For visualization purposes the researchers are sorted according to sharing habitat: not sharing researchers (dark grey circles) to the left, sharing researchers (light grey circles) to the right. See the legend of Figure 2 for parameter settings in all subfigures.
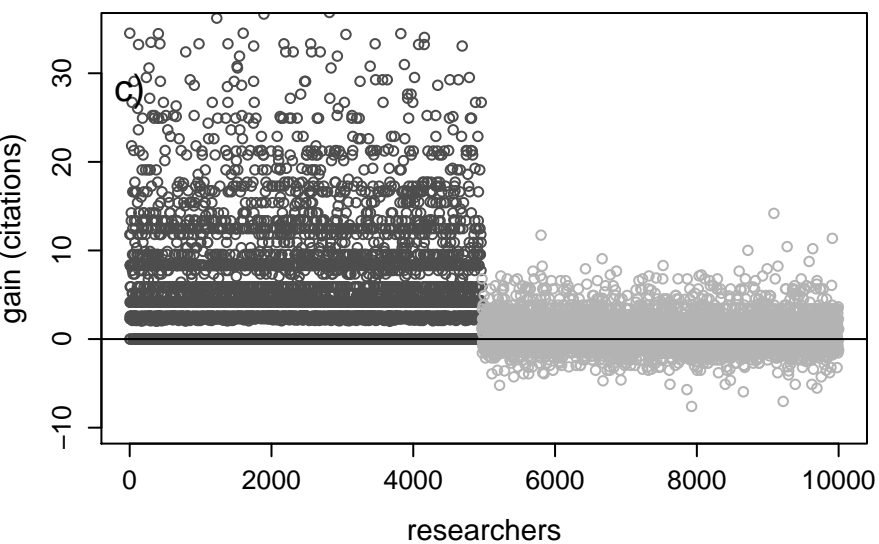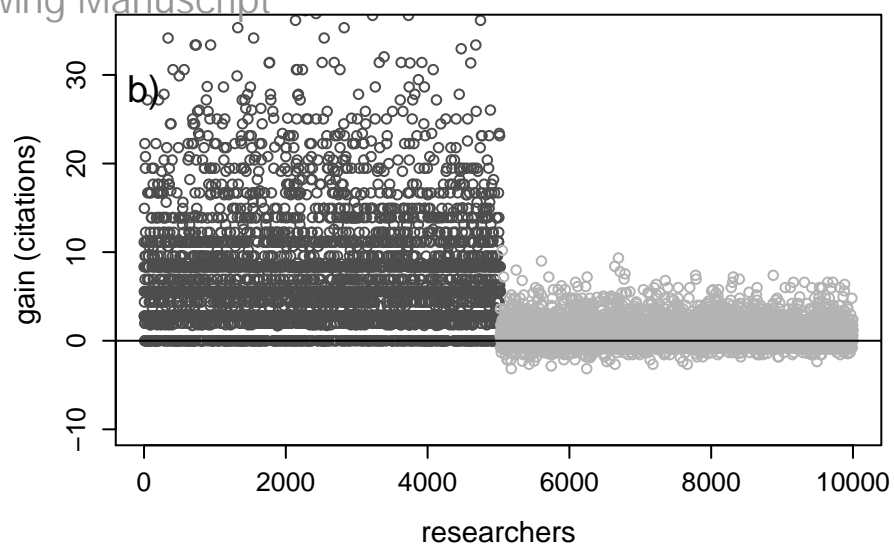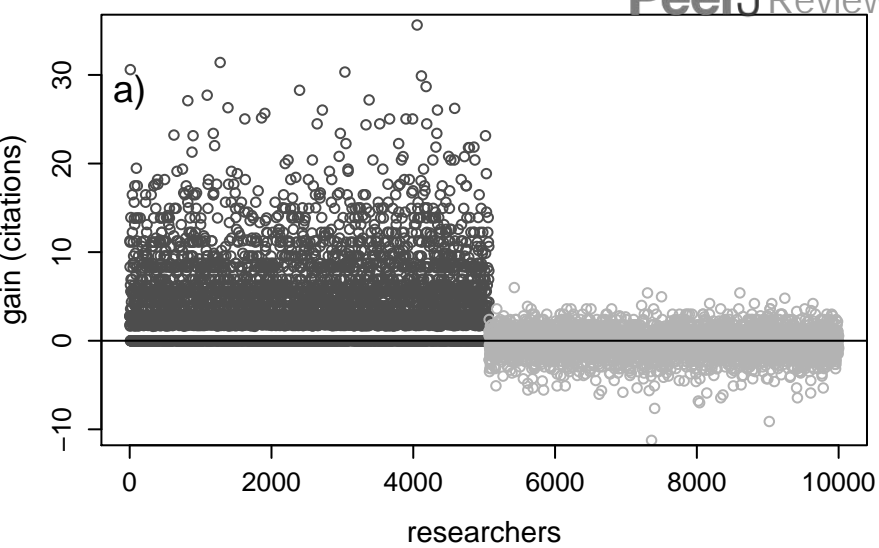
# Figure 4(on next page)

Community impact

Figure 4. Average community impact with varying percentage of sharing researchers and varying sharing benefit $b$. Figures are calculated at default parameter values (see Table 1) with the exception of $b$ which is varied, and for subplot (b) $t_c$, of which the value was set from 0.1 to 0.2. On the z-axis is t he average community impact. On the x and y axes respectively increasing benefits $b$ for sharing from 0 to 0.8 (0 to 80% citation benefit with sharing) and increasing percentage of sharing researchers from 0 to100%.