

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/338562114>

Validity of machine learning in biology and medicine increased through collaborations across fields of expertise

Article in *Nature Machine Intelligence* · January 2020

DOI: 10.1038/s42256-019-0139-8

CITATIONS

11

READS

436

20 authors, including:



Maria Littmann

Technische Universität München

16 PUBLICATIONS 49 CITATIONS

[SEE PROFILE](#)



Liel Cohen-Lavi

Ben-Gurion University of the Negev

10 PUBLICATIONS 51 CITATIONS

[SEE PROFILE](#)



Evans Kataka

Technische Universität München

5 PUBLICATIONS 13 CITATIONS

[SEE PROFILE](#)



Anja Mösch

Technische Universität München

10 PUBLICATIONS 34 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Identifying bacteria associated with tumor development in a UPR-model of colonic tumorigenesis [View project](#)



Optoacoustic Monitoring of Thermal Treatments [View project](#)

Validity of machine learning in biology and medicine increased through collaborations across fields of expertise

Maria Littmann^{1,*}, Katharina Selig^{2,*}, Liel Cohen^{3,4}, Yotam Frank⁵, Peter Hönigschmid⁶, Evans Kataka⁶, Anja Mösch⁶, Kun Qian^{7,8}, Avihai Ron⁹, Sebastian Schmid¹⁰, Adam Sorbie¹¹, Liran Szlak¹², Ayana Dagan-Wiener¹³, Nir Ben-Tal¹⁴, Masha Y. Niv^{13,15}, Daniel Razansky^{9,16,17}, Björn W. Schuller¹⁸, Donna Ankerst², Tomer Hertz^{3,4,19} & Burkhard Rost^{1,20}

- 1 TUM (Technical University of Munich) Department of Informatics, Bioinformatics & Computational Biology - i12, Boltzmannstr. 3, 85748 Garching/Munich, Germany
- 2 TUM Department of Mathematics, Boltzmannstr. 3, 85748 Garching/Munich, Germany
- 3 Department of Microbiology, Immunology and Genetics, Faculty of Health Sciences, Ben-Gurion University of the Negev, Beer-Sheva, Israel National Institute for Biotechnology in the Negev
- 4 National Institute of Biotechnology in the Negev, Beer-Sheva, Israel.
- 5 The Blavatnik School of Computer Science, Tel-Aviv University, Ramat Aviv 69978, Israel
- 6 TUM Department of Bioinformatics, Wissenschaftszentrum Weihenstephan, Maximus-von-Imhof-Forum 3, 85354 Freising, Germany
- 7 TUM Chair of Human-Machine Communication, Theresienstraße 90, 80333 Munich, Germany
- 8 Educational Physiology Laboratory, Graduate School of Education, The University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo, 113-0033, Japan
- 9 Institute for Biological and Medical Imaging, TUM and Helmholtz Center Munich, Germany
- 10 TUM Chair of Food Chemistry and Molecular Sensory Science, Lise-Meitner Str. 34, 85354 Freising, Germany
- 11 TUM Chair of Nutrition and Immunology, Gregor-Mendel-Str. 2, 85354 Freising, Germany
- 12 Weizmann Institute of Science, Rehovot 7610001, Israel
- 13 The Institute of Biochemistry, Food and Nutrition, The Robert H Smith Faculty of Agriculture, Food and Environment, The Hebrew University, 76100 Rehovot, Israel
- 14 Department of Biochemistry and Molecular Biology, George S. Wise Faculty of Life Sciences, Tel Aviv University, Ramat Aviv, 69978 Tel Aviv, Israel
- 15 The Fritz Haber Center for Molecular Dynamics, The Hebrew University, Jerusalem, Israel
- 16 Faculty of Medicine and Institute of Pharmacology and Toxicology, University of Zurich, Switzerland
- 17 Institute for Biomedical Engineering and Department of Information Technology and Electrical Engineering, ETH Zurich, Switzerland
- 18 Group on Language, Audio & Music (GLAM), Imperial College London, UK
- 19 Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA.
- 20 Institute for Advanced Study (TUM-IAS), Lichtenbergstr. 2a, 85748 Garching/Munich, Germany & TUM School of Life Sciences Weihenstephan (WZW), Alte Akademie 8, Freising, Germany & Columbia University, Department of Biochemistry and Molecular Biophysics, 701 West, 168th Street, New York, NY 10032, USA New York Consortium on Membrane Protein Structure (NYCOMPS) & Department of Biochemistry and Molecular Biophysics, Columbia University, 701 West, 168th Street, New York, NY 10032, USA

* Corresponding authors: Maria Littmann: littmann@rostlab.org, <http://www.rostlab.org/>, Tel: +49-289-17-814 (email rost: assistent@rostlab.org); Katharina Selig: katharina.selig@tum.de

◇ These authors share first authorship

Running Title: ML in Life Sciences

Doc stats 1st: Abstract: 215 words, Text: 2878 words, References: 23, figures: 5; tables: 0; pages: 12

Doc stats R1: Abstract: 193 words, Text: 2893 words, References: 23, figures: 5; tables: 0; pages: 22

Journal: Nature Machine Intelligence

1 **Submitted:** 2019/08/09, R1: 2019/11/08

Abstract

Machine Learning (ML) has become an essential asset for the life sciences and medicine. We selected 300 articles describing ML applications from 17 journals sampling 26 different fields between 2011 and 2018. Independent evaluation by two readers highlighted three results. First, only half of the articles shared software, 64% shared data, and 81% applied any kind of evaluation. Although these aspects are crucial to ensure validity and reliability of ML applications, they were met more by publications in lower-ranked journals. Second, the authors' scientific background highly influenced how technical aspects were addressed: reproducibility and computational evaluation methods were more prominent with computational co-authors; experimental proofs more with experimentalists. Third, 73% of the ML applications resulted from interdisciplinary collaborations comprising authors from at least two of the three disciplines: computational sciences, experimental biology, medicine. ~~The data suggested collaborations between computational and experimental scientists to generate more computationally sound and impactful work integrating knowledge. Furthermore, such collaborations provide opportunities to both sides: computational scientists are given access to novel and challenging real-world biological data increasing the scientific impact of their research, and experimentalists benefit from more in-depth computational analyses improving the technical correctness of work.~~

Key words: machine learning, life sciences, medicine, open access, open data, interdisciplinary research, sustainable research, standardization

Abbreviations used: **AI**, Artificial Intelligence; **ML**, Machine Learning; **NAR**, Nucleic Acids Research; **NC**, number of citations; **NC/year**: number of citations normalized by number of years since publication; **NEJM**, New England Journal of Medicine; **PNAS**, Proceedings of the National Academy of Sciences.

NOTE editor & reviewers:

Major changes in red font, major deletions marked as ~~DELETED~~

Introduction

Growing importance of Machine Learning (ML). Large amounts of experimental data triggered by technological advances are increasing the interaction between biology, medicine and quantitative sciences¹⁻³. For instance, the amount of genome sequencing data is growing exponentially while data storage capacity only grows linearly⁴. Numerous large databases in molecular biology and large clinical datasets increasing through electronic health records call for novel ways to interrogate, analyze, and process biological and biomedical data for gaining biological and medical insights⁵.

Machine Learning (ML) automatically identifies patterns and regularities in existing data to accurately predict for unseen data⁶. Despite the complexity of the underlying mathematical concepts, ML has attracted broad attention even outside the research community: querying Google Trends⁷ with “machine learning” demonstrated an exponential increase over the last decade (01/2010-02/2019, data not shown). This general rise has been mirrored in many fields of biology and medicine, i.e. the life sciences⁸⁻¹¹ although keeping track with the rapid evolution of artificial intelligence (AI) challenges even those applying ML¹². Typically, large biological or medical datasets enable the development of ML models that can be used to predict biological or clinical phenotypes through measurements from novel samples.

Quality and validity of ML models hinge upon two primary factors: (1) size, quality and universal validity of data, and (2) the correct development and assessment of the resulting methods^{5,13}. Successful ML applications extract generic principles from today’s data, allowing the generalization, i.e. accurate prediction, for tomorrow’s data. This needs proper extraction and processing of data and features often requiring expert knowledge¹⁴⁻¹⁶. The development and application of ML models to the life sciences needs expertise from both computational and biological/medical fields. In contrast, ML applications to areas such as object and speech recognition or complex games (including chess and Go) for which task and success are more clearly defined and thus require mainly expertise in ML.

Interdisciplinary research might have more impact. In many fields of science, interdisciplinarity has become crucial to produce groundbreaking results through the integration of approaches from different disciplines^{17,18}. Several recent studies^{17,19-24} have been investigating the role of interdisciplinarity by automatically extracting tens and hundreds of thousands of publications (e.g. from WoS²⁵ or PNAS). Toward this end, one definition of *interdisciplinarity* is through the field of the journal in which they are published compared to the journal in which they are cited (the US NSF²⁶ classifies journals into 14 different disciplines and 143 subdisciplines^{17,21}; if published and cited in different fields or subfields, the article is deemed “interdisciplinary”^{17,21,24}). Other definitions^{19,20} focus on the author’s field and define interdisciplinary articles as work published by authors from different fields. So far this definition has been limited to Italian scientists for who there is a public directory mapping researcher and field^{19,20}.

The scientific impact of an article is usually measured by the number of citations for this article^{17,24}. To correct for field- and journal-specific effects that number is often normalized by taking average citation rates and a journal’s impact factor into account^{23,24}. Since the impact factor is calculated from the number of citations of articles published in this journal²⁷, articles from higher-ranked journals are expected to have higher citation counts.

All those automated studies allowed the assessment of many articles while being limited to the extraction of only particular type of information. The studies disagree in their findings regarding the importance of interdisciplinary collaborations: one²⁴ finds no consistent

1 correlation between impact and interdisciplinarity sampling over 750k publications: for some
2 disciplines interdisciplinarity were proportional to citations, for others the relation was reversed
3 xxbr: where is physics in this? Important for next sentence. Another work, focusing on
4 xxbr_number publications from physics²³ found interdisciplinarity was proportional to citation
5 rates but only when published in journals with citation rates below average. Yet other studies
6 xxbr_number¹⁷ and xxbr_number publications²⁰ agreed that interdisciplinarity creates higher
7 impact than non-interdisciplinary work. Also, specific collaborations between scientists from
8 related fields leads to higher-impact publications than generic collaborations between
9 scientists from very different fields²⁰. Clearly, there is no simple red line leading through all
10 those findings. However, what made us re-open the can and begin our analysis were three
11 other reasons: (1) the focus on the life sciences, not explicitly covered by others, (2) the aim
12 of separating the analysis of scientific quality (soundness) and of impact, and (3) the
13 introduction of a more rigorous definition of interdisciplinarity: instead of proxying by the
14 number of disciplines citing a work, we require experts from different disciplines to co-author
15 a work (incidentally, the same sort of definition was used for the analysis of Italian authors
16 ^{19,20}).

17
18 **Focus of this work.** Here, we assessed several aspects of ML applications in the life
19 sciences. We started with the selection of 17 journals representing computational/experimental
20 biology and medicine (Materials & Methods: Supporting Online Material, SOM). Amongst all
21 papers published in those 17 journals in the years 2011-2016, keyword searches (Table S1)
22 matched in 4,306 articles, about 2,100 of those were deemed correct hits after a quick expert
23 analysis. From those, initially 250 were randomly selected (*Materials & Methods* SOM;
24 complete list in additional file paper_table.csv, list of identified falsely extracted articles is
25 provided in additional file false_articles.csv). Subsequently, we applied the same selection
26 process and chose another 50 papers from 2018 to verify that the major findings have not
27 changed through the most recent advent of deep learning (xxbr: cite review here). In contrast
28 to previous studies^{17,19-24}, our assessment focused on ML applications in the life sciences and
29 all information we analyzed was manually extracted from the articles. This allowed, for
30 instance, to correct the 50% false positives from the keyword searches, and also allowed to
31 define interdisciplinarity through the author's background for non-Italians (simply by reading
32 partial CVs for all 1,918 authors of the 250 papers). Each article was classified independently
33 by two of us. These investments limited the number of papers analyzed but allowed a more
34 fine-grained assessment not accessible to automatic extraction. Our focus had several
35 implications, including that all papers reported applications of machine learning to the life
36 sciences, as opposed to more theoretical treatments. In some sense the application of ML
37 (computational sciences) to the life sciences is by definition interdisciplinary. However, we
38 sharpened the perspective by distinguishing expertise from three different disciplines:
39 computational sciences, experimental biology and medicine (expertise of author verified
40 through CV, not through affiliation). Thus, papers could maximally be co-authored by authors
41 from three disciplines, and minimally by one. To simplify, we loosely referred to the case of
42 $N=1$ as to "non-interdisciplinary" and the case of $N>1$ as to "interdisciplinary". For some results,
43 we also showed differences between $N=2$ and $N=3$.

44
45 The correct application of ML requires expertise from those familiar with ML and those
46 familiar with the life sciences, i.e. different disciplines. Thus, we hypothesized articles written

1 by research teams from different disciplines to be more likely to report the necessary
2 evaluation methods ensuring proper implementation of ML methods, to make their data
3 publicly available so others could validate their results, and subsequently, to be accepted in
4 higher-ranked journals and have more citations.

5

Results and Discussion

Coverage of machine learning varies between journals and fields. 58% of the chosen 250 papers (*Material & Methods* in Supporting Online Material – SOM – for more details on how these articles were selected) appeared in only four of the 17 journals (by occurrence: *Bioinformatics*, *PNAS*, *PLOS Computational Biology* and *BMC Bioinformatics*), i.e. were 2.5–fold over-represented (xxbr: can we compute the surprise, simplest model: $0.58/(4/17)$ – but I am sure you ladies might come up with something more intelligent;). Most articles were cited fewer than 100 times, and the number of citations was proportional to time passed since publication (Spearman correlation coefficient $\rho=-0.22$, p-value = 0.03; **Fig. 1**, Fig. S1). The average number of citations for articles from Nature and Science (2011-2016) showed the same trend as that for all 250 articles (Fig. S1). Since the time-dependency obfuscated inter-year comparisons, we normalized by the number of years (*SOM Material & Methods*). As the number of citations correlated with the journal impact factor^{28,29} ($\rho=0.52$, p-value<0.001, **Fig. 1**), all aspects correlating with the impact factor trivially correlated with the number of citations. Normalizing by year and impact factor, removed this correlation. We continued also using the impact factor to assess the visibility of an article as publications in higher-ranked journals tend to be downloaded more often from bioRxiv³⁰. Xxbr: do we really need this addition in red? I for one cannot immediately see what exactly you are aiming at, and adding more than one sentence seems a lot. Do you mean that number of citations (btw. should we introduce “names” NC/NCperanno/NCcorrected?) reflect deflection by bioRxiv more than else)?

>>>

Fig. 1

<<<

The number of articles differed highly between fields: the top five (*molecular biology* 26%, *genetics* 24%, *medicine* 14%, *oncology* 10% and *neuroscience* 9%) accounted for 76% of the 250 articles (Table S2, Fig. S2). Numbers varied even more by disciplines (author expertise): (xxbr: make sure we are NOT confusing discipline=1,2,3 interdisciplinary and field=previous sentence; whatever we settle on: we HAVE to stick to it to avoid further complications): Computational scientists co-authored 88% of all articles, and 95% of those from *genetics* (Fig. S3). Experimental biologists co-authored 70% of all and 59% in *medicine*. Physicians were primarily involved in articles from *medicine* and *oncology* (Fig. S3). Numbers of citations were largely similar for all fields (Fig. S4) but articles focusing on *medicine*, *neuroscience*, and *oncology* tended to be published in higher impact journals (Fig. S4). While the disciplines *experimental biologist* and *physician* correlated positively with impact factor ($\rho=0.30$ /p-value<0.001, $\rho=0.26$ /p-value<0.001, respectively), *computational science* correlated negatively ($\rho=-0.30$ /p-value<0.001; Fig. 1). Computational scientists might focus more on methods, experimental biologists and physicians more on new data that tend to be highly cited in the life sciences.

Three levels of interdisciplinarity. By definition, all the papers analyzed applied methods from computational fields to the life sciences, i.e. were intrinsically interdisciplinary. Most likely all 300 papers analyzed would have been considered “interdisciplinary” by automated analyses checking from which field/discipline the article was quoted. To generate a more detailed lens, we distinguished three disciplines (computational scientists, experimental biologists, and physicians) and introduced interdisciplinarity as a number ranging from 1-3 depending on how many disciplines were represented by the authors of the work. Most of the 250 papers were co-authored by two disciplines (one: 27%, two: 53%, three: 20%). Given

1 these levels, we could classify all papers according to their level of interdisciplinarity and
2 differentially analyze the key indicators: validity (evaluation and sharing) and impact (NC:
3 number of citations, NC/year, NC/year*journal impact factor).

4
5 **Scientific validity higher with experts participating in collaboration.** We proxied the
6 **validity** of papers describing the application of machine learning (ML) methods to biology &
7 medicine (the life sciences) through six different indicators. The first four relate to whether or
8 not the method was assessed in ways needed to ascertain that it works as promised (or at all).
9 We asked: did the authors use cross-validation (V1: binary value), more than one single
10 measure for performance (V2: integer), additional test sets (V3: binary value), and additional
11 experimental verification (V4: binary value). While method evaluation might correctly estimate
12 performance for unseen data without V4, it appears impossible to accomplish this simple
13 objective without V1-V3. The second two indicators related to sharing methods and results.
14 These were sharing data (V5), programs and codes (V6) through publicly available sites.
15 Typically, reviewing ML applications by journal reviewers and the public at large requires
16 availability of data and programs in a form beyond what is available through what can be
17 squeezed into writing. In ML it is almost impossible to imagine the development of the best
18 possible method without any assessment (V1-V3). What if someone might have decided not
19 to publish that assessment? On top, should the aspects of sharing (V5, V6) better be termed
20 “reproducibility” than “validity”. Given the rules of proper scientific conduct, we answer both
21 questions in the negative arguing that without making the evaluation available or making the
22 content of a publication reproducible, the work should either not be published as an application
23 of ML or should be considered as invalid.

24 Evaluation methods (e.g. cross-validation), usage of independent test sets, and/or
25 independent experimental proofs reduce the chance of overfitting and enhance the
26 applicability of the model to future data. Indeed, 80% of the articles with only computational
27 authors, applied some evaluation methods or independent tests; compared to 41% of those
28 written by “experimentalists” (experimental biologists & physicians; Fig. 4). However, most
29 articles written solely by experimentalists provided independent experimental proof (55%), so
30 did 16% of those from only computational co-authors (Fig. 4). The corresponding numbers for
31 interdisciplinary collaborations between computational and experimental scientists (level of
32 interdisciplinarity \geq 2) were between these two extremes: 67% evaluated their methods, 43%
33 provided independent experimental proof, suggesting that such collaborations facilitate
34 experimental and computational validation. On the flip side: 19% of all articles did not provide
35 any evaluation; this number rose as high as 34% without computational co-authors (Fig. 4). To
36 put this most clearly: 19-34% of the papers should have never been accepted, because
37 applications of ML without evaluation resemble “experiments” with no output and/or no
38 measurement.

39 >>>

Fig. 4

<<<

40
41 Several evaluation metrics are required to assess the performance of ML applications
42 (e.g. precision, recall, accuracy or confusion matrices). 6% of all articles used no evaluation
43 metric, 53% used one or two, and 6% used over five (Fig. S7). Although, more metrics do not
44 necessarily imply better assessment, even for binary predictions (separation of two
45 classes/classifications), we have to consider the predictive power of the model for both classes

1 separately, i.e. minimally need two evaluation metrics. More complex problems require more
2 evaluation metrics. Typically, clearly more than two metrics are needed to show different
3 strengths and weaknesses of a prediction method. To put the number five metrics (6% used
4 ≥ 5), none of the ML applications with >400 citations published by the most senior co-author of
5 this manuscript used fewer than eight different metrics.

6 Slightly more than half (52%) of the articles compared their method to others; this again
7 dropped to 21% without computational co-authors (p-value = 0.001; Fig. 2). Although method
8 comparisons are crucial for validation, they might add complications leading to acceptance in
9 lower-ranked journals (Fig. 3C) and possibly to lower impact (Fig. 3A; although adjusting also
10 by impact factor suggested a slight pay-off from method comparisons: Fig. 3B).

>>>

Fig. 3

<<<

13 Reproducibility is a major pillar of science³¹⁻³³ partially relying on making data and
14 methods publicly available. It is particularly critical for ML applications because many minor
15 technical details may invalidate results²⁷. Overall, 64% of the articles shared their data (V5;
16 with large variation between journals: from NAR=89% to NEJM=8%, Fig. S5) reflecting the
17 general trend that articles from *medicine* shared data the least (Fig. S6). We could not establish
18 whether or not this related to sensitive patient data. While all journals encourage data sharing,
19 many do not enforce it.

20 Overall, 68 % of the articles with computational scientists shared data, opposed to 31%
21 without (p-value < 0.001; Fig. 2 _deleted_). 57% of the articles relied on data extracted from
22 public resources or previous articles. However, 22% of those who did, did not publish their
23 data. Data sharing was highest for collaborations with computer scientists (xxbr: as discussed:
24 have not quite seen this plot, we should look separately at data sharing -ds- for computer
25 scientists CS, exp biol EB, clinicians/physicians CP, and interdisciplinary=2 (EB+CP vs EB+CS
26 + CP+CS), and interdisciplinary=3 (all), here I argue that we see interdisciplinary+CS always
27 best, if not: we'll have to refine the statement).

>>>

Fig. 2

<<<

29
30 ***Collaborations of scientists with different expertise somehow cited more often.***
31 Interdisciplinary collaborations of researchers from different fields seem increasingly important
32 to generate new ideas and results^{34,35}. The higher the level of interdisciplinarity, the higher the
33 NC/year (number of citations divided by number of years since publication; $\rho=0.22$, p-
34 value=0.02; Fig. 1, Fig. S8A) and the higher the impact factor ($\rho=0.24$, p-value=0.002; Fig. 1,
35 Fig. S8C). When adjusting the number of citations also by impact factor, the correlation was
36 no longer significant (Fig. 1, Fig. S8B) suggesting that interdisciplinary articles were cited more
37 mainly because they were published in higher-ranked journals (Fig. S8C). The correlation
38 between impact factor and level of interdisciplinarity (Fig. S8C) suggested that authors profit
39 from collaborations.

>>>

Fig. 5

<<<

41
42 Closer analysis of the correlation between interdisciplinarity and impact refined the
43 message: distinguishing just two groups (*computational* and *experimental*), revealed NC to be
44 higher for research teams of only experimental scientists (Fig. 5A), and this outcome was
45 largely caused by physicians (xxbr: point to a figure that makes that point or keep it out?). The

1 results for impact factor and NC adjusted by impact factor suggested that the higher NC
2 originated essentially from physicians publishing in higher-ranked journals (Fig. 5B and C).

3 Did scientific validity (evaluation and sharing) correlate with impact? Computational
4 evaluations correlated negatively with the impact factor ($\rho=-0.31$, $p\text{-value}<0.001$); using no
5 evaluation method correlated positively with the impact factor ($\rho=0.23$, $p\text{-value}=0.004$), but we
6 could not detect a significant relationship between impact factor and experimental proof (Fig.
7 1). Since all articles analyzed here focus on applications, the absence of proper evaluation -
8 independent of the focus of a paper - clearly contradicts good scientific conduct.

9 Data sharing was not rewarded by increases in NC or NC/year (Fig. 3A), although
10 adjusting also by impact factor hinted at a tendency that sharing leads to more citations (Fig.
11 3B). Thus, although data sharing is crucial to ascertain validity and reproducibility, it is not
12 incentivized by increased visibility. In fact, there was no significant difference in the impact
13 factor (Fig. 3C).

14 Software sharing also did not correlate with NC/year (Fig. 3A; the trend changed toward
15 more cited when adjusting NC by impact factor: Fig. 3B). On the contrary, **not** sharing software
16 seemed to lead to acceptance of articles in higher-ranked journals, but again the difference
17 was not significant (Fig. 3C). Certainly, method sharing is crucial for reproducibility and for the
18 impact of a method on science. Therefore, we were surprised that program sharing appeared
19 neither crucial for visibility nor acceptance in the research community as proxied by citations
20 and journal rank. Ultimately, this might shed light on the limitations of such measures to
21 evaluate scientific impact.

22
23 **More computational scientists involved in 2018.** Artificial intelligence (AI) and ML are so
24 rapidly evolving that papers published from 2011-2016 might simply not be up-to-date enough
25 to capture the newest trends. We attempted to address this worry by analyzing another 50
26 articles describing ML applications to the life sciences published in 2018 (selected and
27 analyzed largely by the same criteria as the other 250, Material & Methods in SOM for details).
28 The major differences were: fewer publications without computational scientists (6% 2018 vs.
29 12% 2011-2016), and program sharing rose (70% vs. 50%). Although data sharing did not
30 change significantly (68% vs. 64%), those papers that shared data were cited more often and
31 accepted to higher-ranked journals, we could not detect a significant difference (Fig. S9). T
32 Other aspects did not change significantly, neither program sharing, nor the fact that papers
33 sharing programs tended to be published in lower-ranked journals (Fig. S9), nor the proxies
34 for impact (e.g. NC, NC/year, impact factor, NC normalized by impact factor). Overall, the most
35 substantial change was that computational scientists contributed more often in 2018. This
36 might reflect the increasing complexity of realizing ever more popular deep learning-type
37 solutions of ML.

38
39 **Limitations.** Although our analysis revealed interesting insights, some issues remain to be
40 addressed in the future. First of all, thoroughly analyzing more than 300 articles will render the
41 conclusions more valid. The problem might be one between the Scylla of two few papers and
42 Charybdis of too unreliable analysis. Our solution fell victim of the first, while other solutions
43 xxbr quote-may be NOT incl. the Italians here fell victim of the later. Secondly, we proxied
44 impact and visibility through number of citations and the impact factor. However, other factors
45 are also influencing the number of citations that can seem superficial and can be controlled by
46 the authors³⁶ and it is hard to compensate for these factors. Using the impact factor for
47 measuring scientific impact has been criticized in the literature and the increasing use of social

1 media might increase the visibility of research independent of the journals impact factor^{37,38}.
2 Thirdly, the scope of a journal might influence the description of ML applications. Journals
3 focusing on methodologies are more likely to require certain standards in ML, those focusing
4 on biological and medical relevant novelties are less likely to specifically ask for methodological
5 details. Xxbr: remove (argument: this is peer-reviewed; put up or shut up!) ~~Also, the~~
6 ~~assumption that articles not reporting evaluations did not evaluate is over-simplified.~~
7 ~~_DELETED_~~ Fourthly, we considered any publicly available information to assign author fields,
8 but could not account for paid statisticians not listed as authors. A variety of medical scientists
9 from pathologists to clinicians were all simplified as *physician* ignoring large differences in
10 scientific training. These simplifications might lead to under-estimate computational expertise
11 in publications. Furthermore, we considered data and program availability as stated in the
12 articles, but did not attempt to contact authors to obtain those if not available. Finally, since
13 several aspects in our analysis correlated with the impact factor and they also correlated with
14 each other, confounding factors might influence the results and these interrelationships are
15 difficult to separate.

16
17 Xxbr: the following I suggest to remove, NOT because I don't like it, but because (1) it
18 doesn't really fit here and I do not know how to replace it, and (2) too long anyway. Reason to
19 keep: is simple and some of it is bla & others like more bla than !! ~~For research teams with~~
20 ~~only computational expertise, contributions from physicians or colleagues with expertise in wet~~
21 ~~lab experiments can help to add new data, find biologically relevant applications and~~
22 ~~interpretations of the results, and increase the relevance of ML applications leading to more~~
23 ~~visibility of conducted research because it might be accepted in higher-ranked journals.~~
24 ~~Involving computational scientists in their work does not increase the visibility of research for~~
25 ~~physicians or experimental biologists because this work is rather accepted in lower-ranked~~
26 ~~journals. However, they might benefit from colleagues with knowledge in computer science to~~
27 ~~add evaluation methods, bring a greater variety of tools, and help with the interpretation of the~~
28 ~~scientific and statistical significance of results. Therefore, the results focus more on technical~~
29 ~~aspects making it possibly less intuitive for a broader research community but increasing its~~
30 ~~scientific value by achieving more technical correctness.~~

31 Most likely, with the introduction of new high-dimensional datasets and high-throughput
32 technologies, the need for collaborations will increasingly grow. As the enforcement of data
33 and program transparency will increase, ML methods in biology and medicine will have to be
34 implemented more carefully. While using the impact factor to measure the success of a
35 scientific article currently does not show an advantage of collaborations for experimental
36 scientists (Fig. 5C), we suggest that these collaborations will become more frequent and
37 impactful in the near future.

Conclusions

We analyzed 250 articles describing ML applications to the life sciences published 2011-2016 and another 50 published in 2018 in 17 journals from 26 different biological/medical fields (SOM). This diversity of fields was mirrored by the diversity of how machine learning was applied. Reproducibility and correct evaluation of results are crucial to ascertain validity and reliability of ML applications. Surprisingly, many articles did not focus on these aspects: 50% shared no software, 36% shared no data, and 19% applied no evaluation. In fact, an entire third (34%) of the articles only written by experimentalists described no evaluation. While we hypothesized that ensuring validity of ML applications would be necessary to achieve high visibility of the research, we found the opposite: more valid work was often published in lower-ranked journals attracting fewer citations (Fig. 1, Fig. 3).

In general, how these technical aspects were addressed was highly influenced by the authors' scientific background: Reproducibility and evaluation were more prominent with computational scientists as co-authors (Fig. 2, Fig. 4 _deleted_), while articles co-authored by experimentalists more frequently provided independent experimental proof (Fig. 4). Thus, collaborations of authors from different disciplines provided more opportunity for higher quality results integrating knowledge from various fields of expertise.

We hypothesized that collaborative research should also be cited more often and be accepted in higher-ranked journals. However, this was only true for computational scientists who profited from collaborating with experimentalists, in particular physicians, by getting accepted in higher impact factor journals (Fig. 5C).

One of the most substantial challenges for AI and ML is a comprehensive, adequate evaluation; incorrect application of such tools can lead to drawing false conclusions or to overestimating the predictive power of a method. Collaborations between computational and experimental scientists substantially increased the correctness of evaluations and the likelihood of reproducibility. Thus, increased the scientific validity of published research, a good incentive to focus on such collaborations to improve ML applications that will advance the life sciences in the future.

Acknowledgments

1
2 Thanks to T. Karl and I. Weise (both TUM) for invaluable help with technical and administrative
3 aspects of this work. Thanks to the TUM Graduate School (in particular Z. Zhang) for
4 organizing the Summer School, to the TUM (in particular Dr. H. Keidel and President W.
5 Herrmann) for substantial support on several levels, to the Weizmann Institute, Tel Aviv
6 University, Technion, and Hebrew University for financial and general support; thanks also to
7 the enlightening talks by D. Cremers (TUM), M. Linial (IAS Israel, Hebrew University), Y. Ofran
8 (Bar-Ilan University); thanks to PubMed for providing easy access to published articles and
9 supporting automatic access; thanks to the maintainers of Biopython for providing excellent
10 code to access various databases and process biological data. Particular thanks also to both
11 reviewers for their extremely helpful and plentiful criticisms substantially improving this work.
12 Last, not least, thanks to all maintainers of public databases and to all experimentalists who
13 enabled this analysis by making their data publicly available. This work was supported by grant
14 2003403.
15

Author's contribution

1
2 M.L. and K.S. performed the major part of data analysis and of writing the manuscript. M.L.
3 created and adapted the pre-defined list of articles. K.S. generated figures and performed
4 statistical tests. L.C. assisted in finding interesting correlations in the data by performing
5 complex analyses and statistical test and in generating figures. M.L, K.S., L.C., Y.F., P.H, E.K.,
6 A.M., K.Q., A.R., S.S., A.S., L.S., and A. D.-W. participated in the summer school where the
7 idea for this work was developed, were involved in agreeing on the goals and analysis methods
8 of this work, were involved in data analysis by collecting data from the pre-defined list of
9 articles, and assisted in writing the manuscript. M.L., K.S., and A.M. collected the data for
10 2018. N.B.-T., M.Y.N, D.R., and B.W.S. supervised the work over the entire time and proofread
11 the manuscript. D.A. provided valuable comments especially regarding statistical analysis and
12 was involved in manuscript writing. T.H. and B.R. initiated and supervised the summer school
13 where the idea for this project was developed. T.H. provided important comments to refine the
14 analysis and contributed to manuscript writing. B.R. supervised and guided the work over the
15 entire time and proofread the manuscript. All authors read and approved the final manuscript.

References

- 1 Bleicher, K. H., Bohm, H. J., Muller, K. & Alanine, A. I. Hit and lead generation: beyond high-throughput screening. *Nat Rev Drug Discov* **2**, 369-378, doi:10.1038/nrd1086 (2003).
- 2 Sulakhe, D. *et al.* High-throughput translational medicine: challenges and solutions. *Adv Exp Med Biol* **799**, 39-67, doi:10.1007/978-1-4614-8778-4_3 (2014).
- 3 Howard, J. Quantitative cell biology: the essential role of theory. *Mol Biol Cell* **25**, 3438-3440, doi:10.1091/mbc.E14-02-0715 (2014).
- 4 Cook, C. E. *et al.* The European Bioinformatics Institute in 2016: Data growth and integration. *Nucleic Acids Res* **44**, D20-26, doi:10.1093/nar/gkv1352 (2016).
- 5 Chicco, D. Ten quick tips for machine learning in computational biology. *BioData Min* **10**, 35, doi:10.1186/s13040-017-0155-3 (2017).
- 6 Cios, K. J., Kurgan, L. A. & Reformat, M. Machine learning in the life sciences. *IEEE Eng Med Biol Mag* **26**, 14-16 (2007).
- 7 *Google Trends*, <<https://trends.google.de/trends>> (
- 8 Rost, B., Radivojac, P. & Bromberg, Y. Protein function in precision medicine: deep understanding with machine learning. *FEBS Lett* **590**, 2327-2341, doi:10.1002/1873-3468.12307 (2016).
- 9 Webb, S. Deep learning for biology. *Nature* **554**, 555-557, doi:10.1038/d41586-018-02174-z (2018).
- 10 Min, S., Lee, B. & Yoon, S. Deep learning in bioinformatics. *Brief Bioinform* **18**, 851-869, doi:10.1093/bib/bbw068 (2017).
- 11 Larranaga, P. *et al.* Machine learning in bioinformatics. *Brief Bioinform* **7**, 86-112 (2006).
- 12 Frank, M. R., Wang, D., Cebrian, M. & Rahwan, I. The evolution of citation graphs in artificial intelligence research. *Nature Machine Intelligence* **1**, 79-85 (2019).
- 13 Domingos, P. A few useful things to know about machine learning. *Commun. ACM* **55**, 78-87, doi:10.1145/2347736.2347755 (2012).
- 14 Chou, K.-C. Some remarks on protein attribute prediction and pseudo amino acid composition. *Journal of Theoretical Biology* **273**, 236-247 (2011).
- 15 Ioannidis, J. P. *et al.* Increasing value and reducing waste in research design, conduct, and analysis. *Lancet* **383**, 166-175, doi:10.1016/S0140-6736(13)62227-8 (2014).
- 16 Gron, A. *Hands-On Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. (O'Reilly Media, Inc., 2017).
- 17 Chen, S., Arsenault, C. & Larivière, V. Are top-cited papers more interdisciplinary? *Journal of Informetrics* **9**, 1034-1046 (2015).
- 18 Cummings, J. & Kiesler, S. Organization Theory and the Changing Nature of Science. *Journal of Organization Design* **3** (2014).
- 19 Abramo, G., D'Angelo, C. A. & Di Costa, F. Authorship analysis of specialized vs diversified research output. *Journal of Informetrics* **13**, 564-573 (2019).
- 20 Abramo, G., D'Angelo, C. A. & Di Costa, F. Do interdisciplinary research teams deliver higher gains to science? *Scientometrics* **111**, 317-336, doi:10.1007/s11192-017-2253-x (2017).
- 21 Chen, S., Arsenault, C., Gingras, Y. & Larivière, V. Exploring the interdisciplinary evolution of a discipline: the case of Biochemistry and Molecular Biology. *Scientometrics* **102**, 1307-1323, doi:10.1007/s11192-014-1457-6 (2015).

1 22 Xie, Z., Li, M., Li, J., Duan, X. & Ouyang, Z. Feature analysis of multidisciplinary
2 scientific collaboration patterns based on PNAS. *EPJ Data Science* **7**, 5,
3 doi:10.1140/epjds/s13688-018-0134-z (2018).

4 23 Rinia, E. J., van Leeuwen, T. N. & van Raan, A. F. J. Impact measures of
5 interdisciplinary research in physics. *Scientometrics* **53**, 241-248,
6 doi:10.1023/a:1014856625623 (2002).

7 24 Larivière, V. & Gingras, Y. On the relationship between interdisciplinarity and scientific
8 impact. *Journal of the American Society for Information Science and Technology* **61**,
9 126-131, doi:10.1002/asi.21226 (2010).

10 25 Thomson Reuter's Web of Science, <<http://wokinfo.com>> (
11 26 U. S. National Science Foundation, <<https://www.nsf.gov/>> (
12 27 Wallach, J. D., Boyack, K. W. & Ioannidis, J. P. A. Reproducible research practices,
13 transparency, and open access data in the biomedical literature, 2015-2017. *PLoS Biol*
14 **16**, e2006930, doi:10.1371/journal.pbio.2006930 (2018).

15 28 Impact Factors, <<https://www.scijournal.org/>> (
16 29 Journal Citation Reports, <<https://clarivate.com/products/journal-citation-reports/>> (
17 30 Abdill, R. J. & Blekhnman, R. Tracking the popularity and outcomes of all bioRxiv
18 preprints. *Elife* **8**, doi:10.7554/eLife.45133 (2019).

19 31 Berger, B. et al. ISCB's Initial Reaction to The New England Journal of Medicine
20 Editorial on Data Sharing. *PLoS Comput Biol* **12**, e1004816,
21 doi:10.1371/journal.pcbi.1004816 (2016).

22 32 Drazen, J. M. Data Sharing and the Journal. *N Engl J Med* **374**, e24,
23 doi:10.1056/NEJMe1601087 (2016).

24 33 Longo, D. L. & Drazen, J. M. Data Sharing. *N Engl J Med* **374**, 276-277,
25 doi:10.1056/NEJMe1516564 (2016).

26 34 Mind meld. *Nature* **525**, 289-290, doi:10.1038/525289b (2015).

27 35 Nissani, M. Ten cheers for interdisciplinarity: The case for interdisciplinary knowledge
28 and research. *The Social Science Journal* **34**, 201-216 (1997).

29 36 van Wesel, M., Wyatt, S. & ten Haaf, J. What a difference a colon makes: how
30 superficial factors. *Scientometrics* **98**, 1601-1615 (2014).

31 37 Fitzgerald, R. T. & Radmanesh, A. Social media and research visibility. *AJNR Am J*
32 *Neuroradiol* **36**, 637, doi:10.3174/ajnr.A4054 (2015).

33 38 Patton, R. M., Stahl, C. G. & Wells, J. C. Measuring Scientific Impact Beyond Citation
34 Counts. *D-Lib Magazine* **22**, 5-5 (2016).

35

Figure legends

Xxbr: I might want to go over the following captions in the next round....

Fig. 1: Spearman correlation coefficients for numeric and binary variables. We assessed the correlation between the different criteria using the Spearman correlation and tested the significance at a level of 0.05. Significant p-values are displayed using * for p-value < 0.05, ** for p-value < 0.01 and *** for p-value < 0.001 after adjusting for multiple testing using the Benjamin-Hochberg procedure. Blank squares denote that the correlation is non-significant. Citations adj. (Year + Imp. Fct.) denote the citations adjusted by year and impact factor.

Fig. 2: More sharing and method comparison with computational scientists. The involvement of a computational scientist was highly correlated with ~~sharing the data~~, ~~making the program available~~, or ~~performing a comparison with other methods~~. Percentage of articles with data or program available or performance of a comparison with other methods with 95% percentile bootstrap confidence intervals split by whether a computational scientist was involved.

Fig. 3: Sharing and method comparison hardly impact citations. A. Number of citations adjusted by year were not influenced by data or program availability. Comparing the developed method to others led to a small, decrease in the number of citations. B. Adjusting also by impact factor showed a small trend towards higher citations when data or program were available, or a comparison to other methods was performed. C. The impact factor was higher for articles that did not make data or program available, or compared their method to others.

Fig. 4: Method testing depends on author expertise. ~~Articles involving a computational scientist applied a computational evaluation method more often than articles with only an experimentalist (physician or biologist).~~ ~~Articles co-authored by an experimentalist provided experimental proof more often than without such a co-author.~~ ~~Providing no evaluation method was more common among articles written solely by experimentalists.~~ Percentage of articles with computational evaluation methods, experimental proof or no evaluation methods are shown with 95% percentile bootstrap confidence intervals split by author background.

Fig. 5: Adjusted number of citations and impact factor for different collaborations. A. The number of citations adjusted by year is slightly higher for articles solely written by experimentalists compared to articles involving computational scientists. B. Adjusting also by impact factor removes this difference. This suggests that the higher number of citations for experimentalists was mainly caused by the fact that their work got accepted in higher-ranked journals. C. ~~Impact factor was higher for articles only published by experimentalists (biologists and/or physicians) than for articles involving also computational scientists.~~

Fig. 1: Spearman correlation coefficients for numeric and binary variables

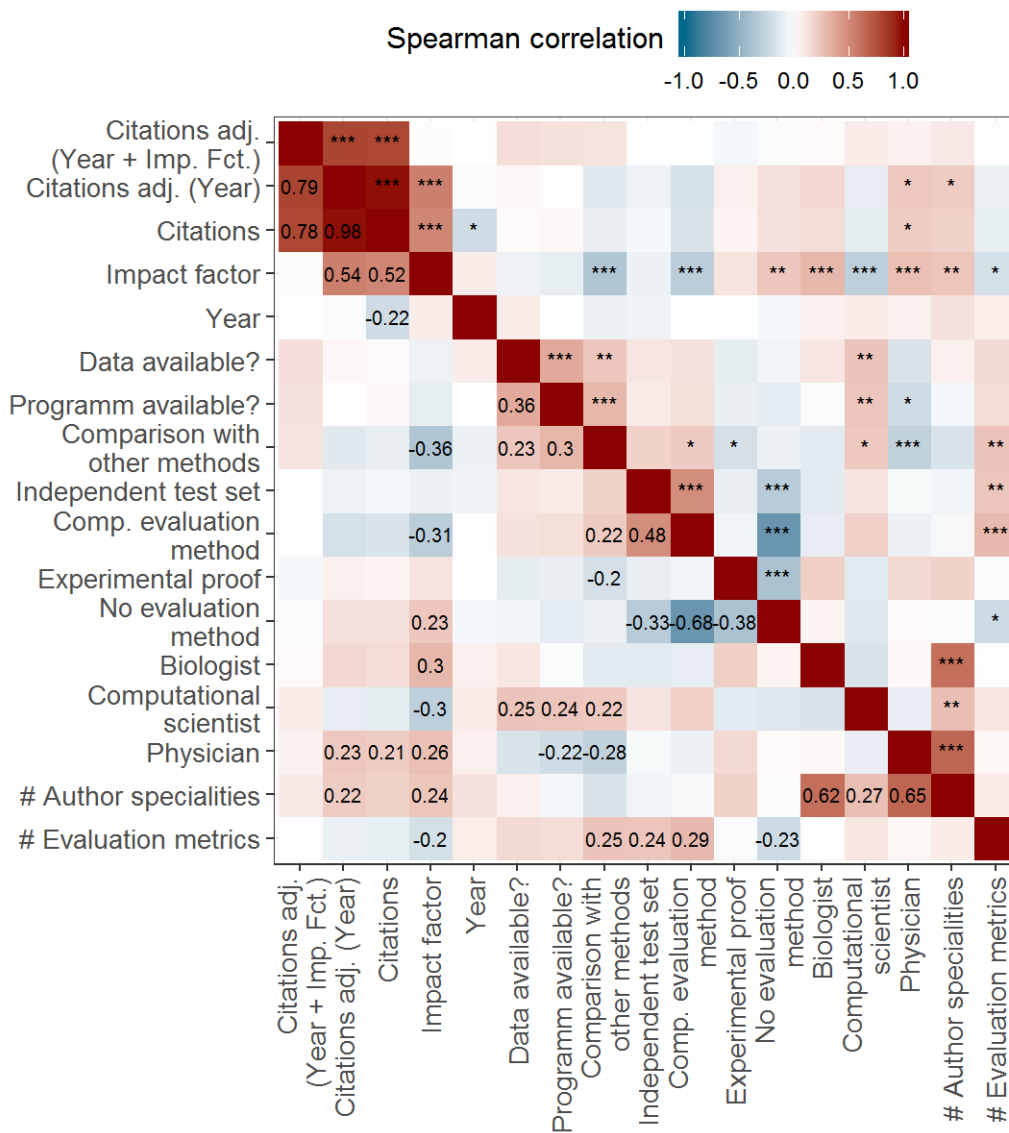


Fig. 1: Spearman correlation coefficients for numeric and binary variables. We assessed the correlation between the different criteria using the Spearman correlation and tested the significance at a level of 0.05. Significant p-values are displayed using * for p-value < 0.05, ** for p-value < 0.01 and *** for p-value < 0.001 after adjusting for multiple testing using the Benjamin-Hochberg procedure. Blank squares denote that the correlation is non-significant. Citations adj. (Year + Imp. Fct.) denote the citations adjusted by year and impact factor.

Fig. 2: Method validation depends on author expertise

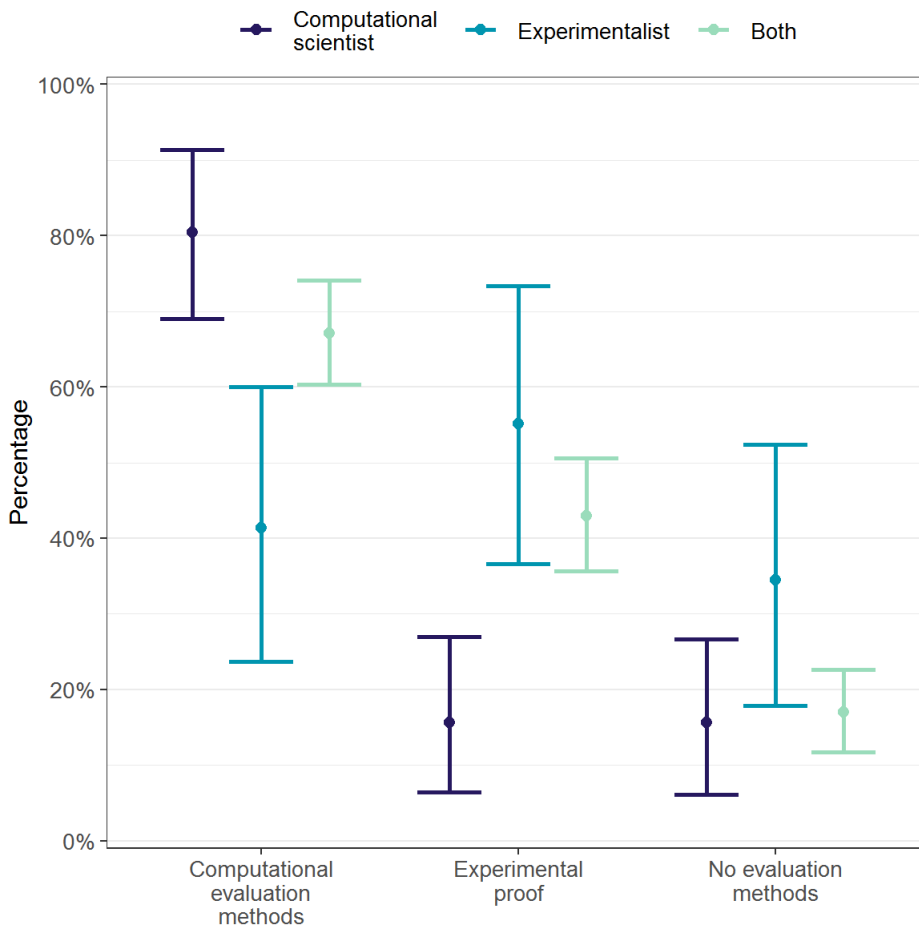
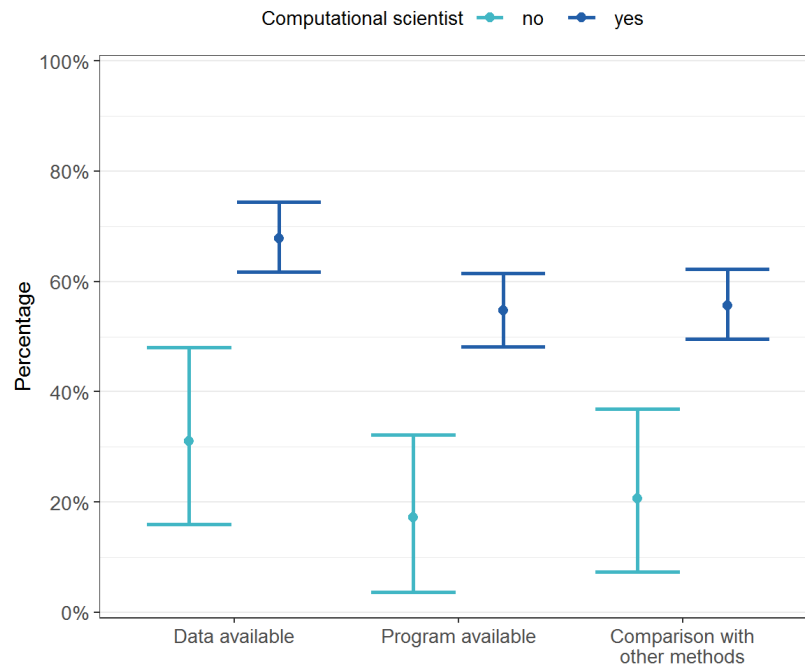


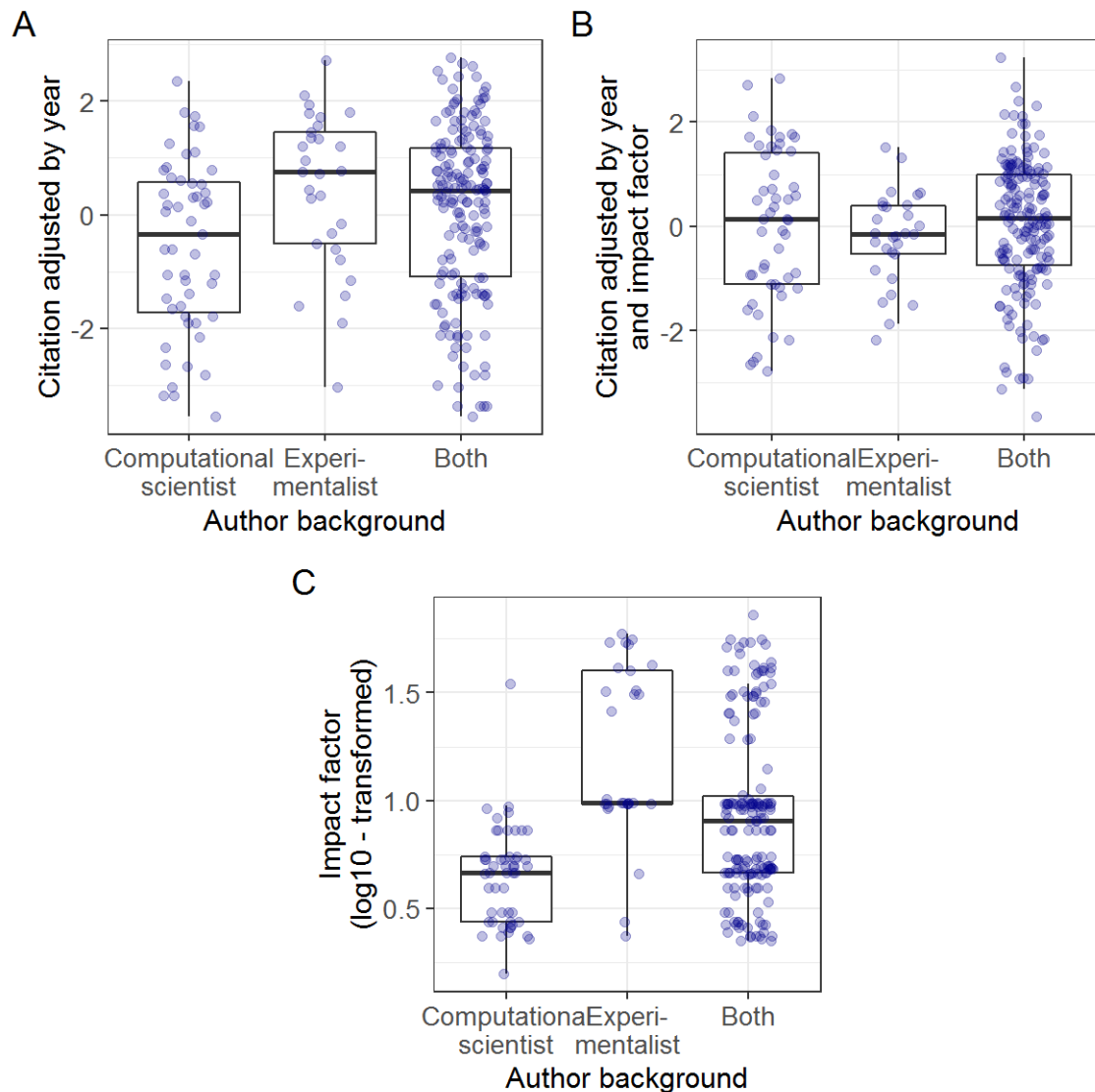
Fig. 2: Method validation depends on author expertise. ~~Articles involving a computational scientist applied a computational evaluation method more often than articles with only an experimentalist (physician or experimental biologist).~~ ~~Articles co-authored by experimentalists provided experimental proof more often than those without.~~ ~~Providing no evaluation method was more common among articles written solely by experimentalists.~~ Percentage of articles with computational evaluation methods, experimental proof or no evaluation methods are shown with 95% percentile bootstrap confidence intervals split by author background.

Xxbr: may be join Fig. 2+3 into one with two panels (assuming you can compress the two to fit together onto a page, i.e. each being a column wide (7.5 cm))

Fig. 3: More sharing and method comparison with computational scientists

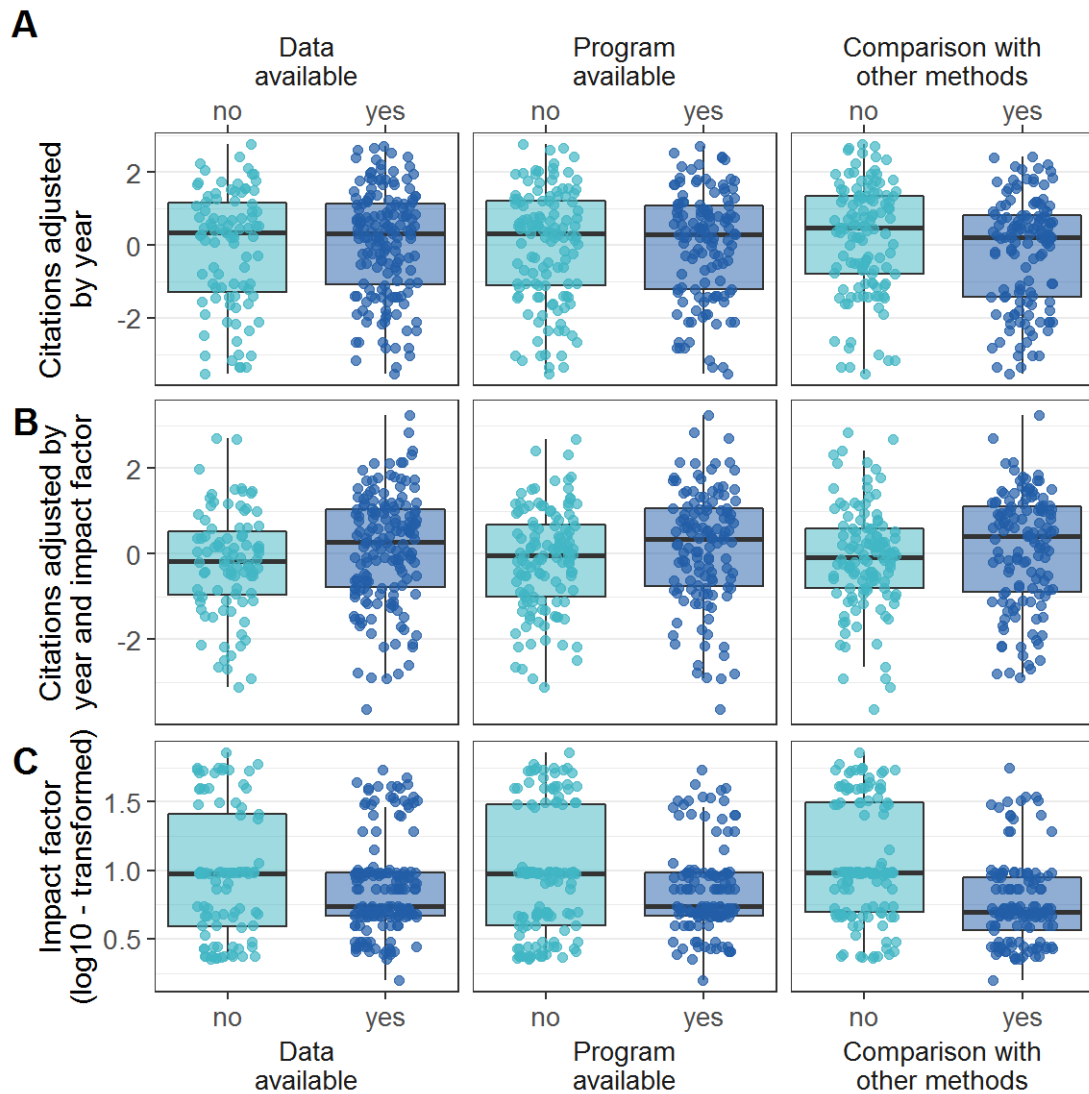
Xxbr: we really ought to optimize this figure: it will be totally fine as a one-column (column=7.5 cm including legends and all); if I shrank it down to that size, we'd lose the font: too small, and there is way too much white space: make thicker lines, bigger dots, less wide bars for errors, points closer to each other in x-axis. Coloring: 1 be careful to use color consistently between figures, already inconsistent between 2 and 3; maybe full black line (yes) and dashed gray: no; you could use the same scheme in Fig. 2 but would have to come up with a new one (blue?) for "both". Btw. Fig. 5 ought to follow that concept and Fig. 4 should NOT interfere – which might speak against blue for both in Fig. 2...

Fig. 3: More sharing and method comparison with computational scientists. The involvement of a computational scientist was highly correlated with DELETED sharing the data, DELETED making the program available, or DELETED performing a comparison with other methods. Percentage of articles with data or program available or performance of a comparison with other methods with 95% percentile bootstrap confidence intervals split by whether a computational scientist was involved.

Fig. 4: NC and impact factor not consistently higher for collaborations

Xxbr: I assume that this will be either one figure spanning two columns (then have them all 3 on one, guess that might work), or 3 under each other as one column.

Fig. 4: NC and impact factor not consistently higher for collaborations. **A.** The number of citations adjusted by year was slightly higher for articles solely written by experimentalists compared to articles involving computational scientists. **B.** Adjusting also by impact factor removed this difference. This suggests that the higher number of citations for experimentalists was mainly caused by the fact that their work got accepted in higher-ranked journals. **C.** ~~DELETED~~ Impact factor was higher for articles only published by experimentalists (biologists and/or physicians) than for articles with computational scientists.

Fig. 5: Sharing and method comparison hardly impact citations

A. Number of citations adjusted by year were not influenced by data or program availability. Comparing the developed method to others led to a small, decrease in the number of citations. B. Adjusting also by impact factor showed a small trend towards higher citations when data or program were available, or a comparison to other methods was performed. C. The impact factor was higher for articles that did not make data or program available, or compared their method to others.