
Managing the Whole Research Process on GitHub

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This is a position paper proposing the idea of managing the entire research process
2 on GitHub. The current machine learning research community faces a variety of
3 problems, such as poor quality and low reproducibility of peer review at interna-
4 tional conferences. These problems are caused by a lack of transparency in the
5 research process and a lack of accessibility, where not everyone can participate in
6 any given process of research. Thus, we propose that any information that arises
7 in the research process be posted on GitHub and that contributions to the research
8 be managed like those in an open-source software project. This could provide a
9 springboard for solving the challenges of machine learning through clarifying con-
10 tributors, allowing fine-grained contributions, improving reproducibility, enabling
11 post-publication peer review, enhancing diversity, and protecting ideas.

12 1 Introduction

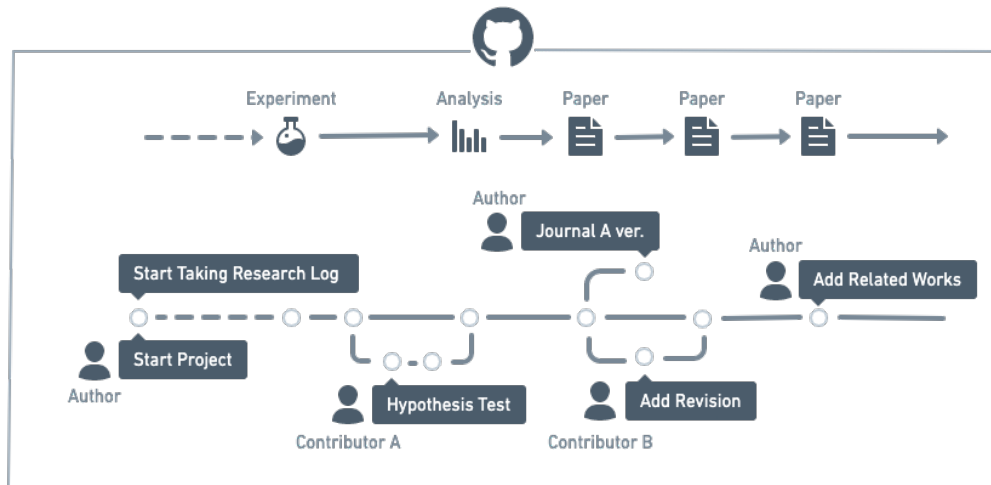


Figure 1: Conceptual diagram of managing research on GitHub. All intermediate outputs of the research process are placed on GitHub. Each contribution in the study is represented as a commit on GitHub, and authors can incorporate contributions from any contributor by accepting pull requests.

13 The machine learning community faces several challenges, such as replication difficulty [1] and
14 the low quality of peer reviews [2]. Solving these challenges will ensure more reliable knowledge

15 production and quality of knowledge. Also, making it easier for researchers to share knowledge and
16 contribute to the others' research can accelerate knowledge production. Enabling these is essential
17 for optimized human knowledge production.

18 An idea to make these possible is to manage research in a public repository on GitHub, a Git repository
19 hosting service ¹. We know that experimental code is published on GitHub in machine learning
20 research community. On the other hand, we propose to manage all intermediate outputs of the
21 research process on GitHub, from the determination of the research topic to the post-publication
22 paper. Literally, all intermediate products, including information about how researchers found the
23 literature, what they thought, what initial experiments they did, etc. In addition, we propose that the
24 progress of the research itself and the maintenance/management of the research results be done in a
25 manner that mimics the way software open source projects proceed.

26 Managing research on GitHub is a very simple proposal, but it would provide one prototype for solving
27 many of the problems facing the machine learning research community. Specifically, managing
28 research on GitHub would bring about the following benefits: clarifying contributors, allowing
29 for contributions in a finer scale, improving reproducibility, enabling post-publication peer review
30 (PPPR), embracing diversity, and protecting research ideas.

31 **2 Proposed Idea**

32 As mentioned in the previous section, we suggest that literally the entire research process be done on
33 GitHub whenever possible. The first step is to create a research log. This can be in any format, e.g.
34 markdown file. Then, since the stage of deciding on a research topic, you will publish it to the public
35 repository on GitHub. Then, via *commits* in GitHub, you record each process of the research in the
36 research log, e.g., wherein the material you consulted and what ideas you got.

37 You also put on GitHub all intermediate outputs generated during the research. For example, the
38 experimental conditions, the log files output for debugging, the results of numerical calculations, the
39 results of pilot studies, and all other intermediate products.

40 When you start writing a paper (some people write it from the beginning of the research, others at the
41 end), you also upload its latex file and all the files and data needed to compile it. Once the paper has
42 passed peer review, you create a branch with the name of the journal/international conference. Then,
43 in the main branch, you will continue to revise and update the paper. Issues that were not addressed in
44 the peer review or that should be improved upon will be created as *issues*. You will then continue to
45 revise the paper by addressing these issues even after the paper is published. You accept contributions
46 from anyone as *pull requests*, both during the paper is written and after the paper is published.

47 **3 The Benefits of Managing Research on GitHub**

48 **3.1 GitHub clarifies contributors**

49 The first advantage is that managing research on GitHub clarifies who contributed and how. Currently,
50 the most common way to express research contributions is to indicate the name of the author on the
51 paper and to give implicit meaning to the order of the authors. However, this makes it difficult for a
52 reader to know who contributed and how much; the contributions of the second and third authors
53 may be far apart, or they may be almost the same. It may also specify what each author did but still
54 does not show enough specific information to reproduce what was done in the research process. At
55 best, they may say that they "did the calculations" or "wrote the text".

56 Suppose that all discussions and revisions of the manuscript are on GitHub. Then, the commit history
57 can be traced to show who contributed to what part of the research process at a glance. Even if
58 the authors' names are not listed in the paper, the history on GitHub makes it clear to everyone the
59 appropriate allocation of credit for the research.

60 This might, for example, eliminate gift-authorship problems [3]. Or it may solve the problem of how
61 to order author names in a collaborative research project. It may also address the plagiarism issue [4].
62 This is because these are the problems that one's contribution is not recognized in the publication.

¹<https://github.com/>

63 **3.2 GitHub allows for finer contributions**

64 The second advantage is that on a more granular basis you can contribute to research and publish
65 research results. Currently, those with smaller contributions are not included as authors in research
66 papers. Even if they are included, it is common to mention them just in acknowledgments. However,
67 if it were possible to contribute to research by commits on GitHub, these small contributions that
68 have not been visualized so far would be visualized on the log of the research.

69 For example, someone good at statistical analysis may be able to commit only to the statistical
70 analysis part. Or someone good at writing may contribute to the paper writing. Furthermore, minor
71 typo corrections, short advice, and help for calculation are also visualized as contributions.

72 Furthermore, enabling contribution on a per-commit basis means that citations will be available on a
73 per-commit hash basis, not just on a per-article basis. This may lead to more direct recognition of
74 individual contributions rather than papers. Citing a commit hash might make it possible to refer to
75 specific ideas, results of experiments, problem formulation, etc. These would allow researchers to
76 utilize multiple talents more effectively in their research.

77 **3.3 GitHub enables post-publication peer review**

78 The third advantage is that it may address the review crisis, which is a problem in machine learning
79 research [5]. To begin with, the number of scientific papers published is increasing rapidly [6]. In
80 particular, because of the recent machine learning boom, there are not enough reviewers for the
81 increasing number of researchers and papers in machine learning research. Because of an insufficient
82 number of reviewers, non-experts often review [5]. Especially in the case of international conferences
83 on machine learning, it is more difficult to ensure the quality of peer review because of the limited
84 review period. Also, because of the fast research cycle in machine learning, several important
85 international conferences are held throughout the year [7]. Therefore, the peer review period for other
86 conferences often overlaps with the preparation period for submitting your new paper to another
87 conference, meaning that you have even less time to spend on peer review. These factors result in the
88 decline and variation in the quality of reviews.

89 One possible solution to this problem of review quality is PPPR [8]. PPPR is an attempt to ensure the
90 validity of papers by evaluating the results after publication. The problem with the current machine
91 learning peer review system is that papers are only evaluated by a specific person in a period of time.
92 On the contrary, in a PPPR system, the quality of a paper is continuously evaluated by several peers
93 for a long period of time after publication. This system seems to make sense, given that science is the
94 activity of making certain that knowledge is irrefutable.

95 Managing the research on GitHub would help to achieve this PPPR. As we mentioned above, you
96 place the LaTeX file or markdown file of published paper on GitHub. Then, you accept corrections
97 and comments on the file as pull requests. This will allow more people to participate in the evaluation
98 of the published paper at any time and in any form they wish. This allows us to continually check the
99 quality of research results on an ongoing basis.

100 Furthermore, this will lead to a return of the peer review process itself from "review to accept or
101 not" to what it should be, "evaluation of the quality of the research results. Managing research on
102 GitHub might be the beginning of a shift away from the competition for the top journal to a system
103 that evaluates what is truly useful for the production of human knowledge.

104 **3.4 GitHub improves reproducibility**

105 The fourth advantage is that it may lead to better reproducibility of research. As mentioned above,
106 machine learning researchers publish their code on GitHub. However, information during the research
107 process, such as "when it didn't work" and "what hyperparameters were important," is lost in some
108 final repositories, decreasing reproducibility of the support for the main argument [1]. Even slight
109 change in the code sometimes makes your experiment fails. This means you should repeat the same
110 process of searching for conditions the original proposer probably went through. In addition, many
111 heuristics are not described or emphasized in the paper.

112 We suspect this is due to the current pressure in the machine learning community to publish positive
113 results to get into good international conferences. Publishing all intermediate thoughts and outputs of

114 the research process would allow researchers to avoid making the same mistake twice. This helps
115 another researcher to assess the work's soundness since all intermediate outputs are public.

116 In addition, the logging research process may contribute to reducing questionable research practices.
117 For example, it might reduce the number of HarKing by making it possible to compare the time
118 stamp when the hypothesis was determined with the time stamp when the experimental results were
119 obtained. These attempts could lead to more reproducible and transparent research practices.

120 Furthermore, since the experimental design is publicly available, something corresponding to pre-
121 registration could also be done on GitHub. Pre-registration is also considered an effective way to
122 reduce QRP, so this is another way to help advance science.

123 **3.5 GitHub embraces diversity**

124 The fifth advantage is that more diverse people can participate in a research project. Open collabora-
125 tion does not care who a person is, what institution they belong to, or what their background is. It
126 does not matter what race, social class, or educational background they have. Research collaboration
127 through GitHub would mitigate the diversity problem many research communities face [9].

128 One of the reasons diversity is an issue is researchers must belong to an organization, and selection
129 must be made to determine the organization's membership. If research is conducted through a
130 web-based platform such as GitHub, researchers can collaborate without any kind of affiliation. For
131 these reasons, the generalization of research on GitHub would be effective in increasing the diversity
132 of the research ecosystem.

133 **3.6 GitHub protects research ideas**

134 The sixth advantage is that more people may be able to publish unpublished content without fear of
135 plagiarism [4]. Currently, researchers generally do not publish their ideas until they are in a paper.
136 This is because they fear that their ideas will be stolen and published first. However, withholding
137 information increases the risk of multiple people duplicating and doing the same thing at the same
138 time, and issues that could be solved immediately by others may not be solved and may not lead to
139 results. Therefore, withholding information slows down the progress of science.

140 As mentioned above, a research log on GitHub shows who was thinking of what ideas and when.
141 Therefore, even if a similar research idea appears later, you can track who states it first by checking
142 the commit history. Hence, more people would be willing to share knowledge with others without
143 fear of having their idea stolen by someone else.

144 Recently, there has been a growing discussion about tying research contributions to the blockchain to
145 guarantee uniqueness (e.g. DeSci [10]). However, it will take time for this to be practical. Starting by
146 visualizing research contributions on GitHub would be a good idea since it is an easy way to start.

147 **4 Discussion**

148 In this paper, we introduced the idea of managing research on GitHub. Then, we explained that doing
149 research on GitHub can solve a variety of machine learning problems.

150 To put this into practice, it would be good to start by putting published papers on GitHub and
151 accepting PPPRs. This is because it may provide a prescription for the serious problem of lack of
152 reviewers despite the psychologically and practically low barrier to entry compared to others.

153 In addition, the use of OSF (Open Science Framework) [11] to disclose the research process is
154 gradually becoming more prevalent in other research fields, such as psychology. And you can connect
155 OSF with GitHub [12]. So, it may be a good idea to start by making OSF pervasive in machine
156 learning research, and then gradually shift to managing it on GitHub.

157 This could be effectively extended to other scientific domains, not just machine learning. The machine
158 learning research community is a desirable environment to experiment with new research styles
159 because experimental code is already available, most papers are open access, and research cycles are
160 fast. By testing various research methods in machine learning research and exporting the good ones
161 to other scientific fields, the development of science as a whole would be accelerated.

162 References

- 163 [1] Edward Raff. A step toward quantifying independently reproducible machine learning research.
164 *Advances in Neural Information Processing Systems*, 32, 2019.
- 165 [2] Corinna Cortes and Neil D Lawrence. Inconsistency in conference peer review: revisiting the
166 2014 neurips experiment. *arXiv preprint arXiv:2109.09774*, 2021.
- 167 [3] Maria Christina Anna Grieger. Authorship: an ethical dilemma of science. *Sao Paulo Medical*
168 *Journal*, 123:242–246, 2005.
- 169 [4] Melissa S Anderson and Nicholas H Steneck. The problem of plagiarism. In *Urologic Oncology:*
170 *Seminars and Original Investigations*, volume 29, pages 90–94. Elsevier, 2011.
- 171 [5] Alessio Russo. Some ethical issues in the review process of machine learning conferences.
172 *arXiv preprint arXiv:2106.00810*, 2021.
- 173 [6] Lutz Bornmann and Rüdiger Mutz. Growth rates of modern science: A bibliometric analysis
174 based on the number of publications and cited references. *Journal of the Association for*
175 *Information Science and Technology*, 66(11):2215–2222, 2015.
- 176 [7] Paper With Code. Ai conference deadlines. [https://aideadlin.es/?sub=ML,CV,CG,NLP,](https://aideadlin.es/?sub=ML,CV,CG,NLP,RO,SP,DM)
177 [RO,SP,DM](https://aideadlin.es/?sub=ML,CV,CG,NLP,RO,SP,DM), 2022. Accessed: 2022-09-29.
- 178 [8] Paul Knoepfler. Reviewing post-publication peer review. *Trends in Genetics*, 31(5):221–223,
179 2015.
- 180 [9] Sarah Myers West, Meredith Whittaker, and Kate Crawford. Discriminating systems. *AI Now*,
181 2019.
- 182 [10] Sarah Hamburg et al. Call to join the decentralized science movement. *Nature*, 600(7888):221–
183 221, 2021.
- 184 [11] Erin D Foster and Ariel Deardorff. Open science framework (osf). *Journal of the Medical*
185 *Library Association: JMLA*, 105(2):203, 2017.
- 186 [12] Center for Open Science. Connect github to a project. [https://help.osf.io/article/](https://help.osf.io/article/211-connect-github-to-a-project)
187 [211-connect-github-to-a-project](https://help.osf.io/article/211-connect-github-to-a-project), 2022. Accessed: 2022-09-29.