# A large-scale analysis of bioinformatics code on

# GitHub

Pamela H Russell[1]*, Rachel L Johnson[1], Shreyas Ananthan[2], Benjamin Harnke[3], Nichole E

Carlson[1]


[1] Department of Biostatistics and Informatics, Colorado School of Public Health, Aurora, CO, USA

[2] High-Performance Algorithms and Complex Fluids, National Renewable Energy Laboratory, Golden,

CO, USA

[3] Health Sciences Library, University of Colorado Anschutz Medical Campus, Aurora, CO, USA


* Corresponding author

E-mail: pamela.russell@ucdenver.edu

18 **Abstract**

19

20 In recent years, the explosion of genomic data and bioinformatic tools has been accompanied

21 by a growing conversation around reproducibility of results and usability of software. However,

22 the actual state of the body of bioinformatics software remains largely unknown. The purpose of

23 this paper is to investigate the state of source code in the bioinformatics community, specifically

24 looking at relationships between code properties, development activity, developer communities,

25 and software impact. To investigate these issues, we curated a list of 1,720 bioinformatics

26 repositories on GitHub through their mention in peer-reviewed bioinformatics articles.

27 Additionally, we included 23 high-profile repositories identified by their popularity in an online

28 bioinformatics forum. We analyzed repository metadata, source code, development activity, and

29 team dynamics using data made available publicly through the GitHub API, as well as article

30 metadata. We found key relationships within our dataset, including: certain scientific topics are

31 associated with more active code development and higher community interest in the repository;

32 most of the code in the main dataset is written in dynamically typed languages, while most of

33 the code in the high-profile set is statically typed; developer team size is associated with

34 community engagement and high-profile repositories have larger teams; the proportion of

35 female contributors decreases for high-profile repositories and with seniority level in author lists;

36 and, multiple measures of project impact are associated with the simple variable of whether the

37 code was modified at all after paper publication. In addition to providing the first large-scale

38 analysis of bioinformatics code to our knowledge, our work will enable future analysis through

39 publicly available data, code, and methods. Code to generate the dataset and reproduce the

40 analysis is provided under the MIT license at https://github.com/pamelarussell/github-

41 bioinformatics. Data are available at https://doi.org/10.17605/OSF.IO/UWHX8.

42

43

## Author summary

45

46  We present, to our knowledge, the first large-scale analysis of bioinformatics source code. The

47  purpose of our work is to contribute data to the growing conversation in the bioinformatics

48  community around reproducibility, code quality, and software usability. We analyze a large

49  collection of bioinformatics software projects, identifying relationships between code properties,

50  development activity, developer communities, and software impact. Throughout the work, we

51  compare the large set of projects to a small set of highly popular bioinformatics tools,

52  highlighting features associated with high-profile projects. We make our data and code publicly

53  available to enable others to build upon our analysis or generate new datasets. The significance

54  of our work is to (1) contribute a large base of knowledge to the bioinformatics community about

55  the state of their software, (2) contribute tools and resources enabling the community to conduct

56  their own analyses, and (3) demonstrate that it is possible to systematically analyze large

57  volumes of bioinformatics code. This work and the provided resources will enable a more

58  effective, data-driven conversation around software practices in the bioinformatics community.

59

60

61

## Introduction

Bioinformatics is broadly defined as the application of computational techniques to analyze biological data. Modern bioinformatics can trace its origins to the 1960s, when improved access to digital computers coincided with an expanding collection of amino acid sequences and the recognition that macromolecules encode information [1]. The field underwent a transformation with the advent of large-scale DNA sequencing technology and the availability of whole genome sequences such as the draft human genome in 2001 [2]. Since 2001, not only the volume but also the types of available data have expanded dramatically. Today, bioinformaticians routinely incorporate whole genomes or multiple whole genomes, high-throughput DNA and RNA sequencing data, large-scale genetic studies, data addressing macromolecular structure and subcellular organization, and proteomic information [3].

Some debate has centered around the difference between "bioinformatics" and "computational biology". One common opinion draws a distinction between bioinformatics as tool development and computational biology as science [4]. However, no consensus has been reached, nor is it clear whether one is needed. The terms are often used interchangeably, as in the "Computational biology and bioinformatics" subject area of *Nature* journals, described as "an interdisciplinary field that develops and applies computational methods to analyse large collections of biological data" [5]. In this article we use the umbrella term "bioinformatics" to refer to the development of computational methods and tools to analyze biological data.
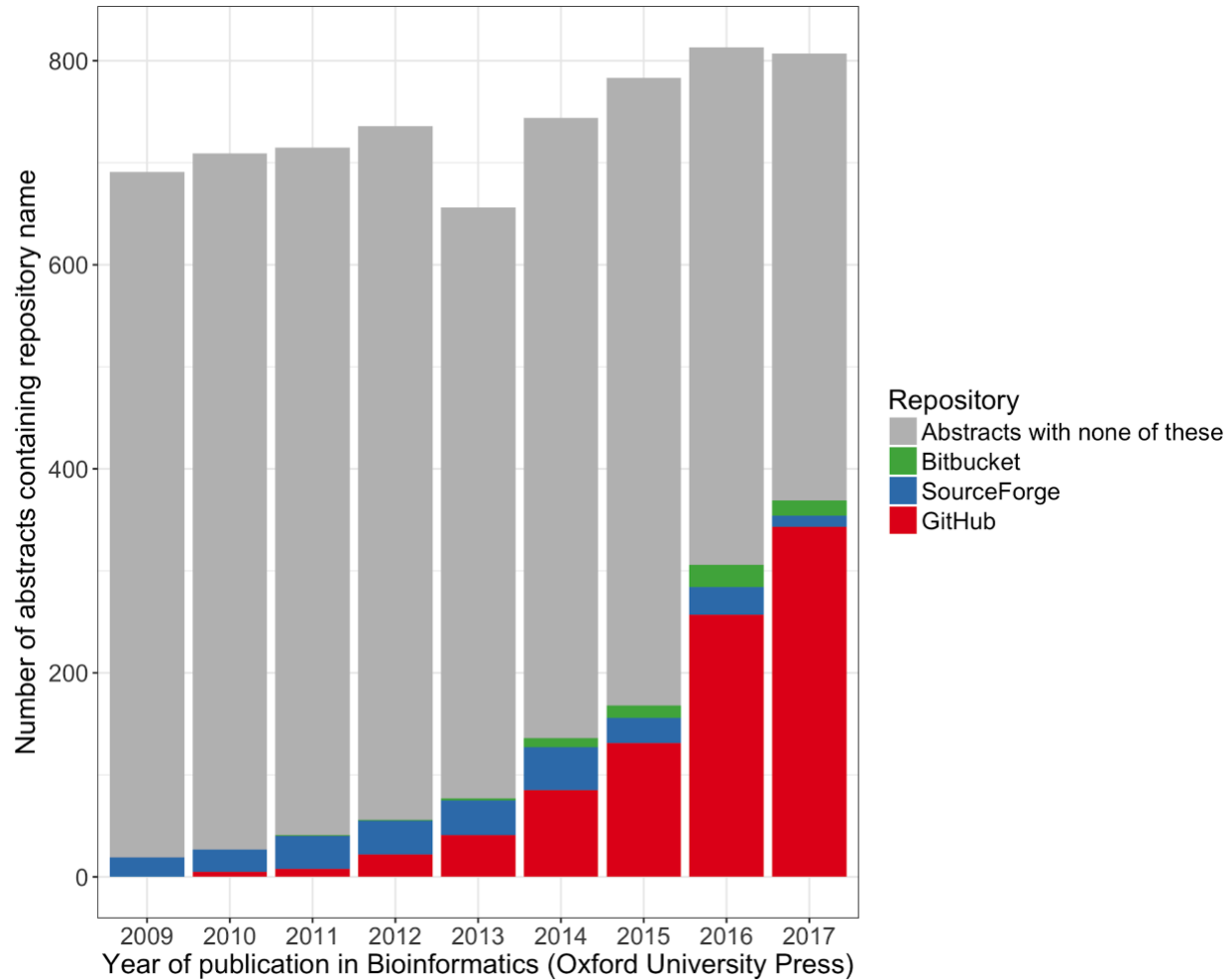
In recent years, the explosion of genomic data and bioinformatic tools has been accompanied by a growing conversation around reproducibility of results and usability of software [6–9]. Reproducibility requires that authors publish original data and a clear protocol to allow repetition

87    of the analysis in a paper [7]. Usability refers to ease and transparency of installation and

88    usage. Version control systems such as Git and Subversion, which allow developers to track

89    changes to code and maintain an archive of all old versions, are widely accepted as essential to

90    the effective development of all non-trivial modern software. In particular, transparent version

91    control is important for long-term reproducibility and usability in bioinformatics [6–9].

92

93    The dominant version control system today is the open source distributed system Git [10], used

94    by 87.2% of respondents to the 2018 Stack Overflow Developer Survey [11]. A Git "repository"

95    is a directory that has been placed under version control, containing files along with all tracked

96    changes. A "commit" is a snapshot of tracked changes that is preserved in the repository;

97    developers create commits each time they wish to preserve a snapshot. Many online sharing

98    sites host Git repositories, allowing developers to share code publicly and collaborate effectively

99    with team members. GitHub [12] is a tremendously popular hosting service for Git repositories,

100   with 24 million users across 200 countries and 67 million repositories in 2017 [13]. Since its

101   initial launch in 2008, GitHub has grown in popularity within the bioinformatics field, as

102   demonstrated by the proportion of articles in the journal *Bioinformatics* mentioning GitHub in the

103   abstract (Fig 1). For an excellent explanation of Git and GitHub including additional definitions,

104   see [14].

105

**Fig 1. Source code repositories in the journal *Bioinformatics*.** Here the term "repository"

refers to online code hosting services. The journal *Bioinformatics* publishes new developments

in bioinformatics and computational biology. If a paper focuses on software development,

authors are required to state software availability in the abstract, including the complete URL

[15]. URLs for software hosted on the popular services GitHub, Bitbucket, and SourceForge

contain the respective repository name except in rare cases of developers referring to the

repository from a different URL or page. The figure shows the results of PubMed searches for

the repository names in the title or abstract of papers published in *Bioinformatics* between 2009

and 2017. The category "Abstracts with none of these" captures all remaining articles published

in *Bioinformatics* for the year, and likely includes many software projects hosted on organization

117    websites or featuring their own domain name, as well as any articles that did not publish

118    software.

119

120    The bioinformatics field embraces a culture of sharing — for both data and source code — that

121    supports rapid scientific and technical progress. In this paper, we present, to our knowledge, the

122    first large-scale study of bioinformatics source code, taking advantage of the popularity of code

123    sharing on GitHub. Our analysis data include 1,720 GitHub repositories published along with

124    bioinformatics articles in peer-reviewed journals. Additionally, we have identified 23 "high-

125    profile" GitHub repositories containing source code for popular and highly respected

126    bioinformatic tools. We analyzed repository metadata, source code, development activity, and

127    team dynamics using data made available publicly through the GitHub API [16]. We provide all

128    scripts used to generate the dataset and perform the analysis, along with detailed instructions.

129    We work within the GitHub Terms of Service [17] to make all data except personal identifying

130    information publicly available, and provide instructions to reconstruct the removed columns if

131    needed. Our main analysis results are provided as a table with over 400 calculated features for

132    each repository.

133

134    Although the software engineering literature describes many analyses of GitHub data [18–24],

135    bioinformatics software has not been looked at specifically. These software engineering studies

136    often look only at highly active projects in wide community use, with many contributors utilizing

137    the collaborative features of GitHub. Public bioinformatics software serves a variety of purposes,

138    from analysis code supporting scientific results to polished tools intended for adoption by a wide

139    audience. With exceptions, code bases published along with bioinformatics articles tend to be

140    small, with one or a few contributors, and use GitHub mostly for its version control and public

141    sharing features. Additionally, the interdisciplinary nature of bioinformatics creates a unique

142    culture around programming, with developers bringing experience from diverse backgrounds

143    [25]. The projects in our dataset treat a variety of scientific topics, use many different

144    programming languages, and show a diverse range of team dynamics.

145

146    We describe our dataset from the perspective of the articles announcing the repositories, the

147    source code itself, and the teams of developers. We observe several features that are

148    associated with overall project impact. Our analysis points to simple recommendations for

149    selecting bioinformatic tools from among the thousands available. Our dataset also contributes

150    to and highlights the importance of the ongoing conversation around reproducibility and

151    software quality.

152

153

154    **Results**

155

156    **A dataset of 1,740 bioinformatics repositories on GitHub**

157

158    We curated a set of 1,720 GitHub repositories mentioned in bioinformatics articles in peer-

159    reviewed journals (referred to throughout the paper as the "main" dataset), as well as 23 high-

160    profile repositories that were not necessarily on GitHub at the time of publication or are not

161    published in journals. Three repositories overlapped between the two sets. As a resource for the

162    community, we provide the full pipeline to extract all repository data from the GitHub API, all

163    extracted data except personal identifying information, scripts to perform all analysis, and

164    citations for the articles announcing each repository.
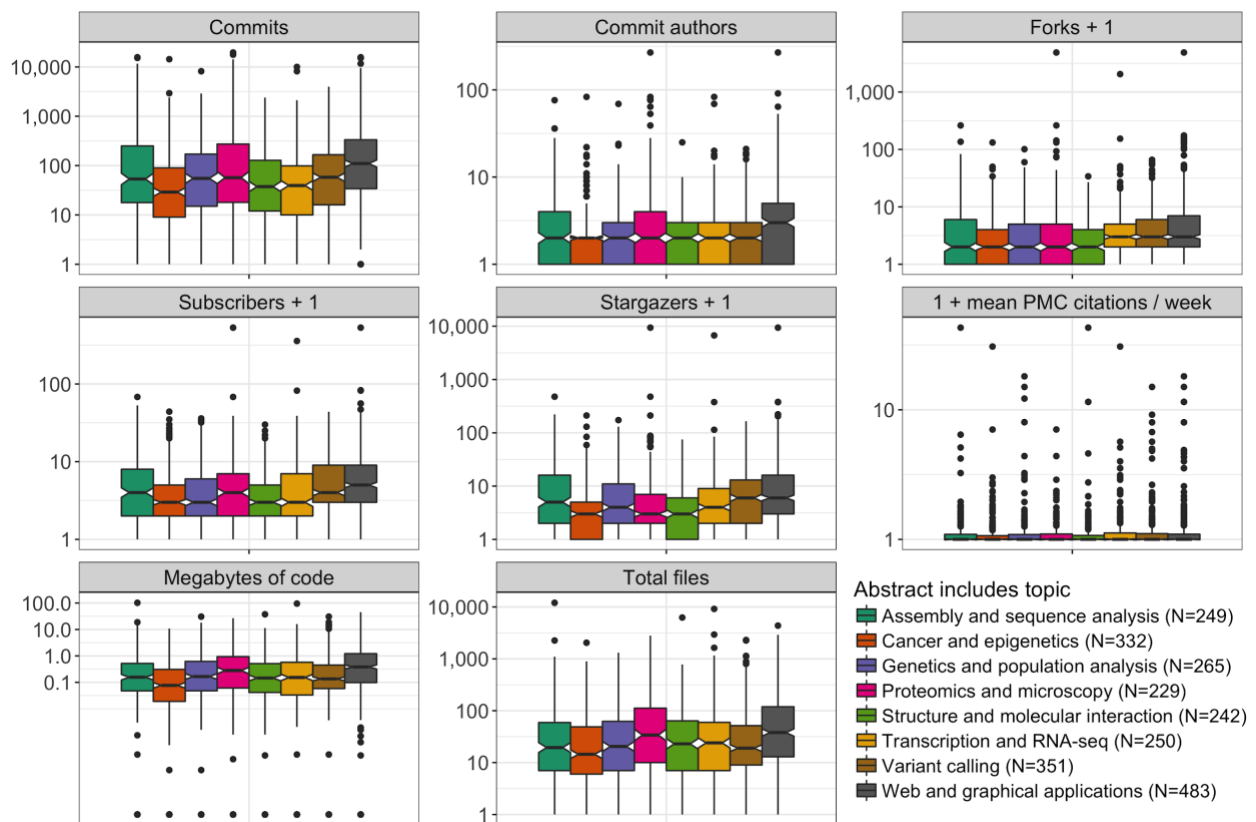
165

166    **Article topics**

167

168    We performed topic modeling [26] on the abstracts of the articles announcing each repository in

169    the main dataset, associating each article with one or more topics. We manually assigned labels

170    to each topic based on top associated terms (Fig S1); for example, the topic "Transcription and

171    RNA-seq" is associated with the terms "rna", "seq", and "transcript". We found that the topic

172    "Web and graphical applications" was positively associated with several measures of project

173    size and activity, as were, to a lesser extent, some other topics (Fig 2). We found that code for

174    articles about certain topics was disproportionately written in certain languages; for example, the

175    greatest amount of code for "Assembly and sequence analysis" was in C and C++, while the

176    greatest amount of code for "Web and graphical applications" was in JavaScript (Fig S2).

177    *Bioinformatics* was the most common journal for all topics, probably due in part to the relative

178    ease of finding relevant projects in this journal (Fig S3). Fig S4 shows topic distribution by year

179    of initial commit and article publication.

180



181

182 **Fig 2. Project features by article topic.** Projects are broken into groups according to whether

183 the accompanying paper abstract is associated with each topic category. Projects that are

184 associated with multiple topics are counted separately for each topic. Topic labels were

185 assigned manually after examining top terms associated with each category. We added one to

186 several variables to facilitate plotting on a log scale; these are noted in the variable name. All

187 variables refer to the GitHub repository except "1 + mean PMC citations / week", which refers to

188 the paper and looks at citations in PubMed Central per week starting two years after the initial

189 publication of the paper. Commits is the total number of commits to the default branch. Commit

190 authors have created commits but do not necessarily have push access to the main branch; we

191 attempted to collapse individuals with multiple aliases. Forks are individual copies of the

192 repository made by community members. Subscribers are users who have chosen to receive

193 notifications about repository activity. Stargazers are users who have bookmarked the

194 repository as interesting. Megabytes of code and total files include source code only, excluding

195 data file types such as JSON and HTML. The horizontal line at the center of the notch

196 corresponds to the median. The lower and upper limits of the colored box correspond to the first

197 and third quartiles. The whiskers extend beyond the hinges by at most an additional 1.5 times

198 the inter-quartile range. Outliers are plotted individually. The notches correspond to roughly a

199 95% confidence interval for comparing medians [27]. The table of repository features is provided
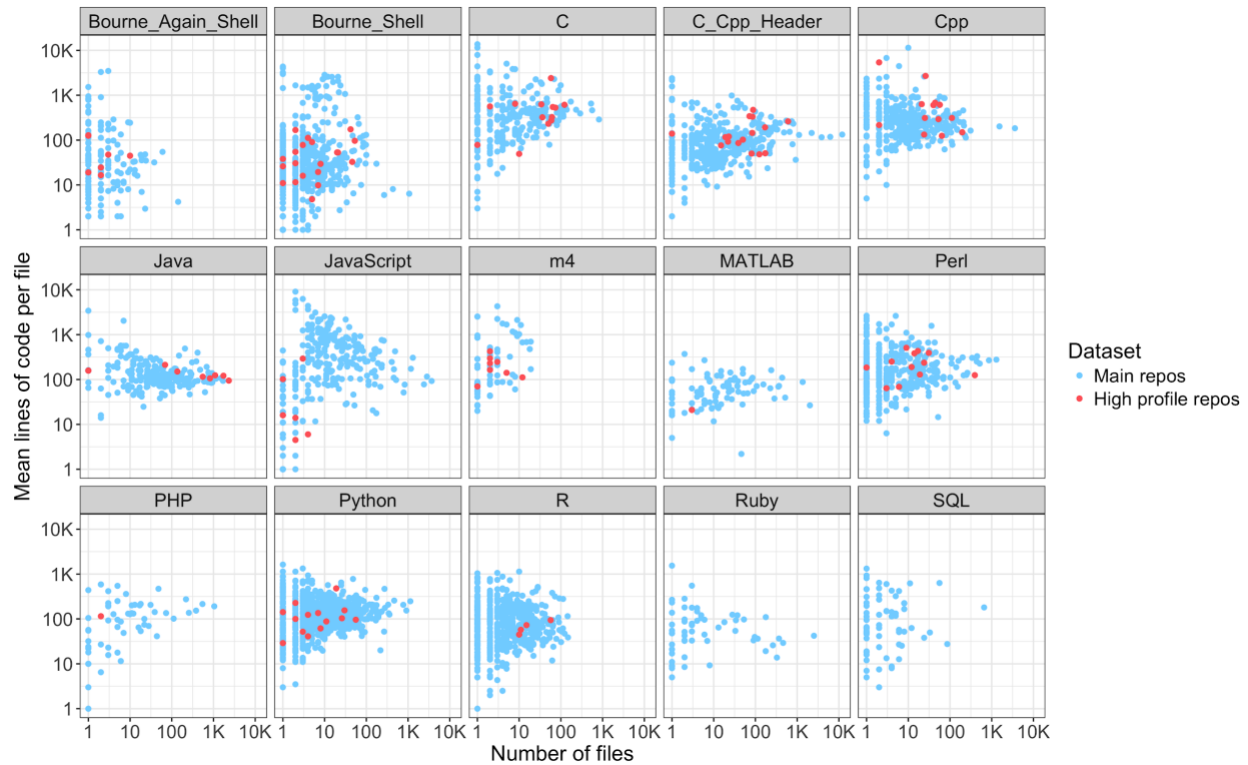
200 as Table S8.

201

202 **Programming languages**

203

204 We identified a programming language for each source file and analyzed the prevalence of

205 languages along several dimensions including total number of source files, lines of code, and

206 size of source files in bytes. In high-profile repositories, the greatest amount of code in bytes

207 was in Java, followed by C and C++. In the main dataset, two repositories contained entire

208    copies of the large C++ Boost libraries [28]. Ignoring those copies of Boost, the greatest amount

209    of code in the main dataset was in Javascript, followed by Java, Python, C++, and C (Fig S5).

210

211    We analyzed language features including primary execution mode (interpreted or compiled),

212    type system (static or dynamic, strong or weak), and type safety. High-profile repositories

213    tended to emphasize compiled, statically typed languages, with the largest contribution being

214    from Java. The main dataset contained a greater proportion of code written in interpreted or

215    hybrid interpreted/compiled (such as Python) and dynamically typed languages (Fig 3, Fig S6,

216    Table S6, Table S7). This difference could reflect the fact that interpreted and dynamically typed

217    languages provide a powerful platform to quickly design prototypes for small projects, while

218    static typing provides important safety checks for larger projects. Indeed, there was a

219    relationship between project size (total lines of code) and amount of statically typed code

220    (percentage of bytes in statically typed languages): the Spearman correlation between these

221    variables over the entire dataset was 0.41 (P=2.2e-16) (Table S8). Our data support the intuition

222    that Java, Python and R are more succinct than lower-level languages such as C and C++, as

223    the former group tended to have fewer lines of code per source file in the presumably

224    sophisticated high-profile repositories (Fig 3).

225

**Fig 3. Number and length of source files by programming language.** Languages included in at least 50 main repositories are shown. Each dot corresponds to one repository and indicates the number of files in the language and the mean number of lines of code per file not including comments. The data are provided as Table S8.
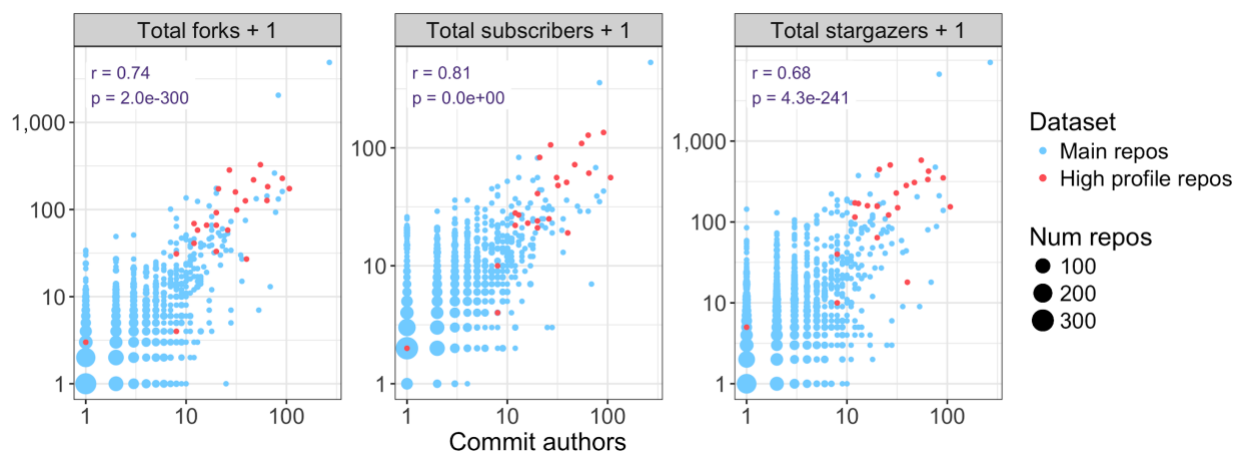
**Developer communities**

For version control systems such as Git, "commits" refer to batches of changes contributed by individual users; each commit causes a snapshot of the repository to be saved along with records of all changes. Each GitHub repository has a core team of developers with commit access; these developers can push changes directly to the repository. In addition, GitHub facilitates community collaboration through a system of forks and pull requests. Anyone can create a personal copy of a public repository, called a "fork", and make arbitrary changes to their fork. If an outside developer feels their changes could benefit the main project, they can create

241     a "pull request": a request for members of the core team to review and possibly merge their

242     changes into the main project. In that case, the commit records for the main project would show

243     the outside contributor as the commit author and the core team member who merged the

244     changes as the committer.

245

246     We looked at the size of each developer team (including users with commit access and outside

247     contributors) as well as other measures of community engagement, including number of forks,

248     subscribers, and stargazers. Subscribers are users who have chosen to receive notifications

249     about repository activity. Stargazers are users who have bookmarked the repository as

250     interesting. Neither subscribers nor stargazers necessarily touch any code, though in practice

251     they are likely to include the developer team. Not surprisingly, the size of the developer team (all

252     commit authors) was strongly associated with the number of forks, subscribers, and stargazers.

253     High-profile repositories tended to have larger teams and more community engagement by

254     these measures (Fig 4). The number of outside contributors was also associated with these

255     measures, though less strongly, perhaps because only 14% of main repositories had any

256     outside contributors and these already tended to be within the highly active subset; 70% of high-

257     profile repositories had outside contributors (Fig S7).

258



259

260   **Fig 4. Size of developer community.** Various measures of community engagement are plotted

261   against the number of commit authors. Each dot represents one repository or a set of

262   repositories with identical values for the variables. We added one to the vertical axis variables to

263   facilitate plotting on a log scale due to many zero values. The pearson correlation and

264   associated p-value are displayed for each variable versus number of commit authors. Commit

265   authors refers to the number of unique commit authors to the default branch. The high-profile

266   repository with a single contributor is s-andrews/FastQC [29]. This repository appears to have

267   been created by a single developer importing a previously existing code base to GitHub. The

268   table of repository features is provided as Table S8.

269

270   **Gender distribution of developers and article authorships**

271

272   We analyzed the gender distribution of developers and article authorships in the dataset as a

273   whole and within teams. Developer and author first names were submitted to the Genderize.io

274   API [30] and high-confidence gender calls were counted. We found that the proportion of female

275   authors decreased with seniority in author lists and the proportion of female developers was

276   lower in high-profile repositories compared to the main dataset. In the main dataset, 12% of

277   developers were women while only 6% of commits were contributed by women; these numbers

278   were lower in the high-profile dataset (7% and 2%, respectively). In biology articles, it is

279   customary to list the lead author first and the senior author last, with additional authors in the

280   middle. We found that in the articles announcing each repository, middle authors included the

281   greatest proportion of women. Women comprised 22% of all authorships in the main dataset

282   and 21% in the high-profile dataset, compared to 18% and 0% for first authors and 14% and 8%

283   (representing only one person) for the most senior last authors (Fig 5). A separate study of

284   author gender in computational biology articles found a similar trend of decreased

285    representation of women with increased seniority in author lists; the authors additionally

286    identified a pattern of more female authors on papers with a female last author [31].
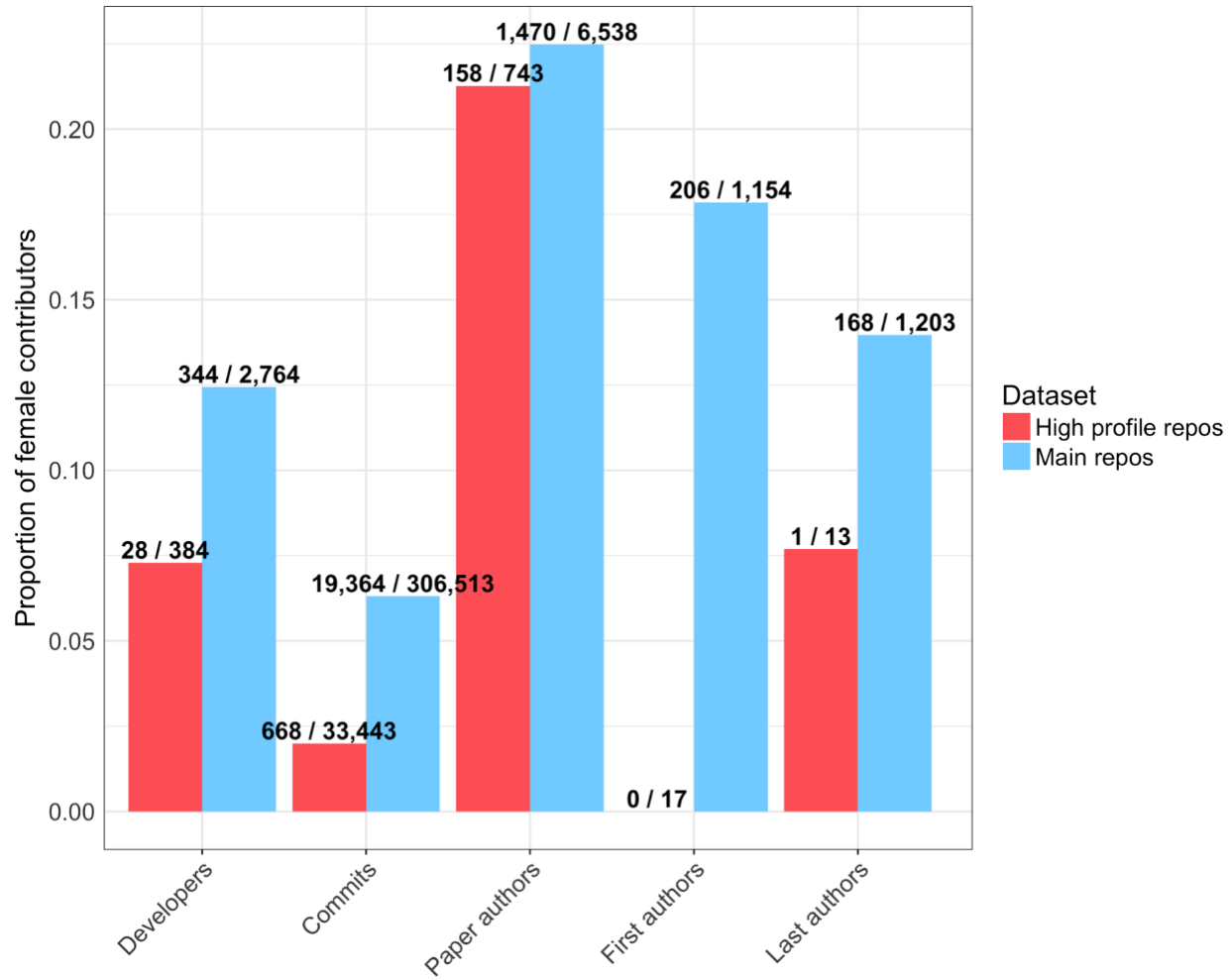
287

288    We analyzed the gender composition of each team of developers and paper authors. The most

289    common type of team in the main dataset was a single male developer and an all-male author

290    list. The most common type of team in the high-profile dataset was a majority-male developer

291    team and an all-male author list. Only ten main repositories and no high-profile repositories had

292    all or majority female developer and author teams; all ten of these developer teams consisted of

293    a single female developer (Fig S8).

294

295    We quantified gender diversity within teams using the Shannon index of diversity [32]. A

296    Shannon index of 0 means all members have the same gender, while the maximum value of the

297    Shannon index with two categories is $\ln(2) = 0.69$, achieved with equal representation of both

298    categories. We found that 13% of main repositories and 62% of high-profile repositories had a

299    nonzero Shannon index for the developer team. There were no high-profile repositories with a

300    Shannon index greater than 0.4; the percentage of main repositories with Shannon index

301    greater than 0.4 was 12% (Fig S9).

302

**Fig 5. Distribution of developers, commits and paper authorships by gender.** "Developers" are unique commit authors or committers over the entire dataset; we attempted to collapse individuals with multiple aliases. "Commits" are individual commits to default branches of repositories. "Paper authors" are individual authorships on papers, not necessarily unique people. For each repository, the one paper announcing the repository is included; papers were then deduplicated because some papers announced multiple repositories. First and last authors are only counted for papers with at least two authors. Names for which a gender could not be inferred are excluded. Bar height corresponds to the number of female contributors divided by the number of contributors with a gender call; these numbers are labeled above each bar. The features for each repo are provided in Table S8.
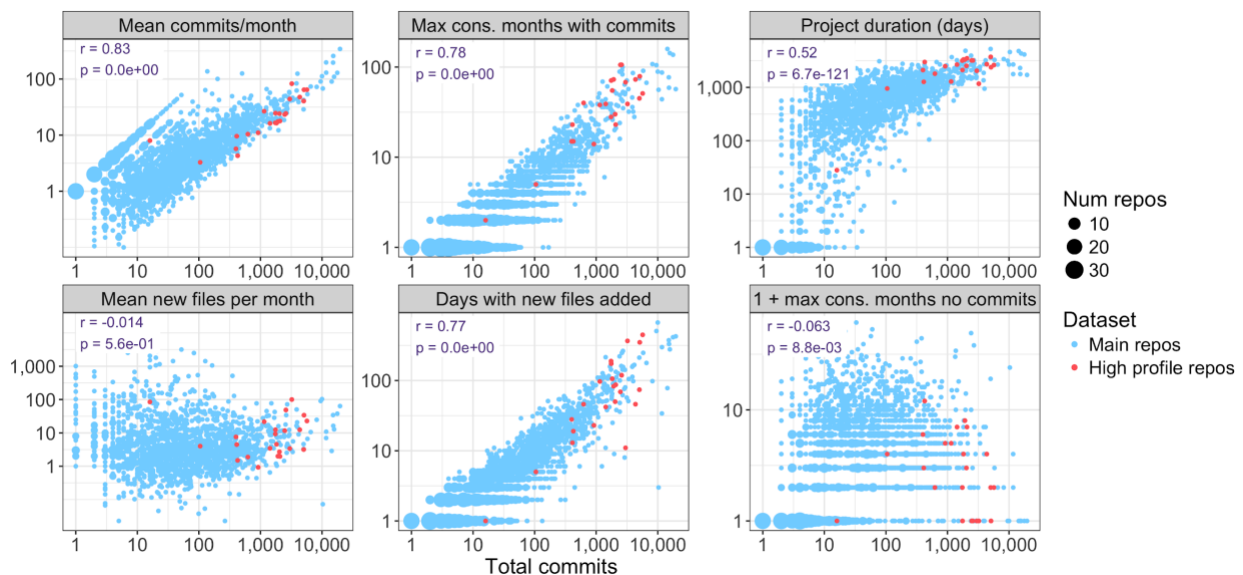
314

315     **Commit dynamics**

316

317     We looked at several measures of commit timing along with total number of commits to each

318     repository. Not surprisingly, the total number of commits was strongly associated with density of

319     activity (commits per month and maximum consecutive months with commits) and overall

320     project duration. High-profile repositories tended to have longer project duration and greater

321     density of commit activity (Fig 6).

322



323

324     **Fig 6. Commit timing versus total commits.** Various timing dynamics are plotted versus total

325     commits to the default branch. Each dot represents one repository or a set of repositories with

326     identical values for the variables. For each variable, the total time interval covered by the project

327     is the interval starting with the first commit and ending with the last commit at the time we

328     accessed the data. For example, "Mean new files per month" counts only months from the first

329     to last commit. The high-profile repository with only 16 commits and all files added on a single

330     day is s-andrews/FastQC [29]. This repository appears to have been created by importing a

331     previously existing code base to GitHub. The data are provided as Table S8.
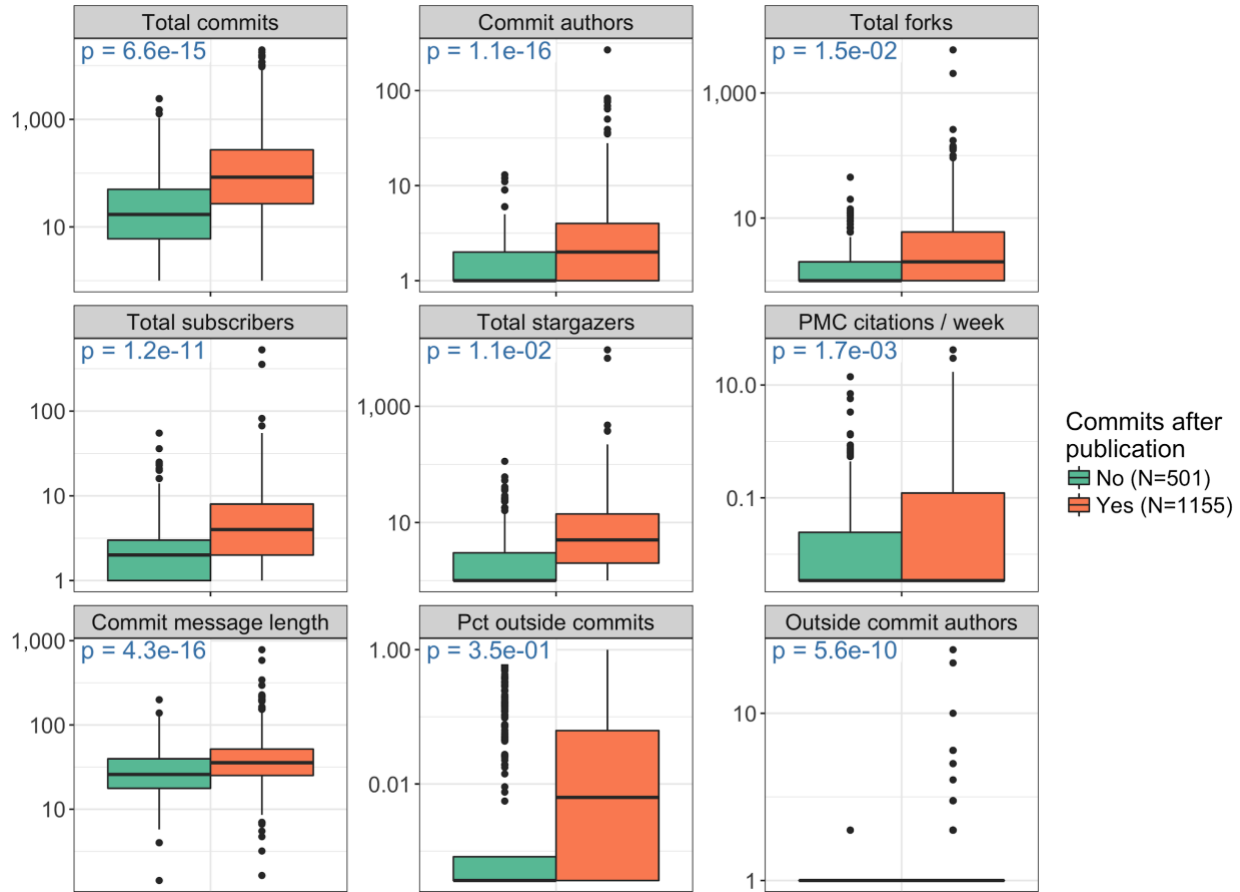
332

### A simple proxy for project impact

334

335    We looked at the simple binary feature of whether any commits were contributed to each

336    repository after the associated article appeared in PubMed. We found that this simple feature

337    was associated with several measures of project activity and impact (Fig 7). Not surprisingly, it

338    was strongly associated with the total number of commits and size of the developer team.

339    Presumably, larger projects tend to be those that are useful to many people and for which

340    development continues after the paper is published. The metric was also associated with

341    measures of community engagement such as forks, stargazers, and outside contributors. This

342    could be explained in part by the previous point and in part by outside community members

343    voluntarily becoming involved in the project after reading the paper. However, interestingly, the

344    association with the proportion of commits contributed by outside authors was not statistically

345    significant, suggesting that overall team size may be the principal feature driving the relationship

346    with the number of outside commit authors. Additionally, the metric was associated with

347    frequency of citations in PubMed Central, which could indicate that people are discovering the

348    code through the paper and using it, and the code is therefore being maintained. Interestingly,

349    repositories with commits after the paper was published had longer commit messages

350    (explanations included by commit authors along with their changes to the repository). This could

351    be due to a relationship between both variables and the size of the developer team; perhaps

352    members of larger teams tend to write longer commit messages to meet the increased burden

353    of communication with more team members. Indeed, there was a moderate linear relationship (r

354    = 0.14, p = 1.9e-09) between total number of commit authors and mean commit message length

355    in the main dataset.

356

357

**Fig 7. Commits after paper publication.** Each data point contributing to each box plot is one repository in the main dataset. Repositories are separated by whether the last commit timestamp at the time we accessed the data was after the date the corresponding publication appeared in PubMed. Repositories for which we do not have a publication date in PubMed are excluded. See Fig 2 legend for the explanation of "Total commits", "Commit authors", "Total forks", "Total subscribers", "Total stargazers", and "PMC citations / week". "Commit message length" is the mean number of characters in a commit message. "Pct outside commits" is the proportion of commits with an author who is never a committer. "Outside commit authors" is the number of commit authors who are never committers. The p-value refers to the two-sided *t*-test for different means between the two groups. The data used to compute the p-value include zero values, but for the plot, we replaced zeros by the minimum positive value of each variable to facilitate plotting on a log scale. The horizontal line across the box corresponds to the median.

370     The lower and upper limits of the box correspond to the first and third quartiles. The whiskers

371     extend beyond the box by at most an additional 1.5 times the inter-quartile range. Outliers are

372     plotted individually. The table of repository features is provided as Table S8.

373

374

375     **Discussion**

376

377     We have presented the first large-scale analysis of bioinformatics code to our knowledge. Our

378     analysis gives a high-level picture of the current state of software in bioinformatics, summarizing

379     scientific topics, source code features, development practices, community engagement, and

380     team dynamics. The culture of sharing in bioinformatics will continue to enable deeper study of

381     software practices in the field. Our hope is that readers will uncover additional insights in our

382     tables of hundreds of calculated features for each repository (Table S8), many of which were

383     not analyzed in this paper, and that some readers will use or adapt our code to generate data

384     and analyze repositories in unanticipated ways.

385

386     Interestingly, despite being made public on GitHub, nearly half of all repositories in our dataset

387     do not feature explicit licenses (Fig S10), in most cases likely unintentionally restricting the

388     rights of others to reuse and modify the code. Nonetheless, the type of research described here

389     may proceed under the GitHub Terms of Service [17] and Privacy Statement [33].

390

391     With the overwhelming variety of public bioinformatics software available, users are constantly

392     faced with the question of which tool to use. Several features of our analysis point to simple

393     heuristics based on information available on GitHub. We observed relationships between

394     community engagement and various measures of project size and activity level (Fig 4, Fig 6,

395     Fig S7). Our final analysis looked at the simple question of whether the developers had revisited

396     their code at all after the paper was published; we found that this feature is associated with

397     several measures of impact (Fig 7). Intuitively, these points suggest that users should prioritize

398     software that is being consistently maintained by an active team of developers. The GitHub web

399     interface prominently displays the total number of commits, number of contributors, and time of

400     latest commit on the front page for each repository. Additionally, GitHub provides a full-featured

401     mechanism, called Issues, that allows the developer team or any user to create tracked

402     requests within the project. We did not analyze issues because these are a relatively advanced

403     feature that is rarely used in our dataset; nonetheless, a consistent flow of issues can help

404     identify sophisticated projects under active development.

405

406     Bioinformatics is a hybrid discipline combining biology and computer science. There are three

407     major paths into the field: (1) computer scientists and programmers can become familiar with

408     the relevant biology, (2) biologists can learn programming and data analysis, or (3) students can

409     train specifically in increasingly popular bioinformatics programs [25]. Our dataset likely includes

410     developers from all three major paths. However, our analysis of developer gender demonstrates

411     that the gender distribution in bioinformatics more closely resembles that of computer science

412     than biology. Indeed, the underrepresentation of women in our dataset was more extreme than

413     among students awarded PhDs in computer science in the United States in 2016 [34]. A

414     possible reason for this could be that, despite relatively high numbers of women in biology,

415     biologists who make the transition to bioinformatics tend to be male. Another possible

416     explanation could be that the subset of bioinformaticians who publish code on GitHub are

417     disproportionately those from the computer science side. Importantly, our analysis does not

418     address other intersections of identity and demographics that affect individuals' experience

419     throughout the academic life cycle. Beyond simply pushing for fair treatment of all scientists,

420    researchers have argued that team diversity leads to increased productivity of software

421    development and higher quality science [35–37].

422

423    **Limitations**

424

425    Our dataset represents a large cross section of bioinformatics code bases, but many projects

426    are excluded for various reasons. First of all, due to the challenges of full-text literature search,

427    we did not identify all articles in the biomedical literature that mention GitHub. In particular, we

428    did not use the open access set of articles in PubMed Central because these included too many

429    mentions of GitHub to manually curate for both bioinformatics topics and code being announced

430    with the respective articles, and efforts to train automated classifiers left too many false

431    positives that tended to skew the picture of repository properties compared to true

432    announcements of bioinformatics code. We therefore selected a search strategy that was

433    limited enough to generate a high-quality hand-curated set and could include papers that were

434    not open access. Second, we are missing repositories that were not on GitHub at the time of

435    publication or are primarily described on a main project website other than GitHub, with the

436    exception of the high-profile repositories we added manually. Third, we have not included large

437    open source collaborations such as Bioconductor [38], BioJava [39], and Biopython [40], due to

438    project-specific substructure making it unfair to compare them to the rest of the dataset. Finally,

439    our dataset could be biased due to our use of GitHub itself: it is possible that developers with

440    certain backgrounds are disproportionately likely to host code on GitHub, while we have not

441    analyzed any code not hosted on GitHub.

442

443    The spirit of sharing has led to an increase in popularity of preprints: advance versions of

444    articles that have not yet been published in peer-reviewed journals. Preprints can allow scientific

445    progress to continue during the sometimes extensive review process. However, we chose not to

446     include preprints in our literature search for three main reasons. First, we believed that

447     successful peer review was a fair criterion on which to identify serious code bases. Second, we

448     wanted to analyze article metadata that would only be available from databases such as

449     PubMed. Third, the most popular preprint server for biology, bioRxiv [41], does not currently

450     provide an API, putting programmatic access out of reach.

451

452     **Future research**

453

454     Several interesting future analyses are possible with our dataset or extensions to it. First, we did

455     not examine the important topic of software documentation, either within source code or for

456     users. The myriad forms of user documentation (README files, help menus, wikis, web pages,

457     forums, and so on) make this a difficult but important topic to study. Second, static code

458     analysis would provide deep insight into software quality and style. While impractical for a large

459     heterogeneous set of code bases written in many different languages, future studies could

460     uncover valuable insights through focused static analysis of repositories sharing common

461     features. Third, we did not study the behavior of individual developers in depth. Future studies

462     could analyze the social and coding behavior of individuals across all their projects and interests

463     on GitHub. Finally, our analysis does not address the important question of software validity:

464     whether a program correctly implements its stated specification and produces the expected

465     results. The complexity of bioinformatic analysis makes validity testing a very challenging

466     problem. Nevertheless, progress has been made in this area [42–44]. Our hope is that others

467     will leverage our work to answer further important questions about bioinformatics code.

468

469     **Toward better bioinformatics software**

470

471    Our work provides data to enhance the ongoing community-wide conversation around

472    reproducibility and software quality in bioinformatics. Several features of our data suggest a

473    need for community-wide software standards, including the widespread absence of open source

474    licenses (46% of main repositories have no detectable license), the number of repositories not

475    appearing to use version control effectively (12% of main repositories added all new files on a

476    single day, while 40% have a median commit message length less than 20 characters), and the

477    apparent lack of reuse of the software (28% of papers in the main dataset have never been

478    cited by articles in PubMed Central, while 68% have fewer than five citations) (Table S8).

479    Similarly, a study based on text mining found that over 70% of bioinformatics software

480    resources described in PubMed Central were never reused [45]. These orthogonal lines of

481    evidence support the need for the already growing efforts toward supporting better software in

482    bioinformatics and scientific research in general.

483

484    Existing efforts to improve research software include the Software Sustainability Institute

485    [46,47], which works toward a mission of improving software to enable more effective research;

486    Better Scientific Software [48], a project that provides resources to improve scientific and

487    engineering software; and Software Carpentry [49–51], which provides highly practical training

488    for research computing. In addition, several reviews recommend specific practices for the

489    software development lifecycle in academic science. In [8], the author provides specific

490    recommendations to improve usability of command line bioinformatics software. The authors of

491    [52] recommend specific software engineering practices for scientific computing. In [9], the

492    authors outline several practices for the entire software development lifecycle. In [53], members

493    of a small biology lab describe their efforts to bring better software development practices to

494    their lab. In [54], the author advocates for changes at the institutional and societal levels that

495    would lead to better software and better science.

496

497  Our contribution to this conversation, in addition to the specific conclusions from our analysis, is

498  to demonstrate that it is possible to study bioinformatics software at the atomic level using hard

499  data. With continued updates, this paradigm will enable a more effective, data-driven

500  conversation around software practices in the bioinformatics community.

501

502

503  **Methods**

504

505  **Identification of bioinformatics repositories on GitHub**

506

507  GitHub repositories containing bioinformatics code were found through their mention in

508  published journal articles pertaining to bioinformatics topics. Briefly, a literature search identified

509  articles that were likely to pertain to bioinformatics topics and contained mentions of GitHub.

510  Manual curation identified the subset of these articles treating bioinformatics topics, using a

511  detailed definition of bioinformatics. GitHub repository names were automatically extracted from

512  the bioinformatics articles. Mentions of each repository in each article were manually examined

513  to identify repositories containing code for the paper, as opposed to mentions of outside

514  repositories. Repository names were manually deduplicated and fixed for other noticeable

515  issues such as inclusion of extra text due to the automatic parsing of context around the

516  repository name. Repository names were automatically checked for validity using the GitHub

517  API, and repositories with issues in this check were manually fixed or removed if the repository

518  no longer existed. The final set included 1,720 repositories. In addition to the 1,720 repositories

519  identified through the literature search, we also curated a separate set of 23 high-profile

520  repositories — highly popular and respected tools in the bioinformatics community — based on

521  the high volume of posts about these projects on the online forum Biostars [55]. The two

522    datasets are referred to throughout the paper as the "main" and "high-profile" datasets. See

523    Supplemental Section 2 for details. The repositories are listed in Table S4 and Table S5.

524

525    **Extraction of repository data from GitHub API**

526

527    Repository data were extracted from the GitHub REST API v3 [16] and saved to tables on

528    Google BigQuery [56] for efficient downstream analysis. Data extracted for each repository

529    include repository-level metrics, file information, file creation dates, file contents, commits, and

530    licenses. GitHub API responses were obtained using the PycURL library [57]. The JSON

531    responses were converted to database records and pushed to tables on BigQuery using the

532    BigQuery-Python library [58]. See Supplemental Section 3 for details.

533

534    **Topic modeling of article abstracts**

535

536    We used latent Dirichlet allocation (LDA) [59] to infer topics for abstracts of the articles

537    announcing each repository in the main dataset. From the LDA model, we identified terms that

538    were primarily associated with a single topic. We chose a model with eight topics due to its

539    maximal coherence of concepts within the top topic-specialized terms. We manually assigned a

540    label to each of the eight topics that captures a summary of the top terms. We then classified

541    each article abstract into one or more topics. Details are in Supplemental Section 4.

542

543    **Programming languages**

544

545    We identified 515,017 total files files among the repositories in the main dataset and 22,396

546    total files in the high-profile dataset. Contents of 425,967 and 18,501 files respectively (349,834

547    and 16,917 with unique contents) with size under 999KB were saved to tables in BigQuery for

548   further analysis. (See Supplemental Section 3.) We used cloc (Count Lines of Code) version

549   1.72 [60] to identify the programming language, count lines of code and comments, and extract

550   comment-stripped source code for each file. A total of 221,343 unique files in the main dataset

551   and 11,425 in the high-profile dataset had an identifiable programming language. Language

552   execution modes were obtained from [61]. Type systems were obtained from [62]. Further

553   details are presented in Supplemental Section 5.

554

555   **Developer communities**

556

557   We identified the number of commit authors and outside contributors for each repository. For

558   commit authors, we attempted to count unique people by collapsing users with the same name

559   or login. For outside contributors, we counted commit authors whose author ID is never a

560   committer ID for the repository. The counts of forks, subscribers and stargazers were returned

561   directly from the GitHub API. Further details are presented in Supplemental Section 6.

562

563   **Gender analysis**

564

565   We attempted to infer a gender for each commit author, committer, and article author using the

566   Genderize.io API [30], which returns a gender call and probability of correctness for a given first

567   name. Names were first cleaned to remove noise such as single-word handles or organization

568   names, and then the first word of each cleaned full name was submitted to Genderize. We

569   accepted gender calls whose reported probability was 0.8 or greater. We proceeded with

570   analysis of "female" and "male" categories only. We assume that transgender and non-binary

571   contributors have names that reflect their gender identity. There may be erroneous calls for

572   individuals who do not identify with a binary gender. The gender calls are also expected to

573     include a few errors for cisgender individuals as we accept calls with global probability of 0.8 or

574     higher.

575

576     To analyze the gender breakdown of developers, we counted unique full names of authors and

577     committers. For commits, we joined commit records to genders by the full name of the commit

578     author and counted individual commits. For paper authors, we counted individual authorships on

579     papers instead of unique individuals, reasoning that multiple different authorships for the same

580     individual should be counted separately. We analyzed team composition for the 504 projects in

581     the main dataset for which we could infer a gender for at least 75% of developers and 75% of

582     paper authors (Fig S8). We calculated the Shannon index of diversity [32] for the 602

583     repositories in the main dataset for which we could infer a gender for at least 75% of developers

584     (Fig S9). Details are described in Supplemental Section 7.

585

586     **Commit dynamics**

587

588     We defined project duration as the time span between the first and last commit timestamps for

589     the repository. Metrics describing monthly activity are with respect to the number of months in

590     the project duration. We identified the initial commit time for each file by taking the earliest

591     timestamp of all commits touching the file. Details are described in Supplemental Section 8.

592

593     **Proxy for project impact**

594

595     We defined "commits after publication" to be true if the latest commit timestamp at the time we

596     accessed the data was after the day the associated article appeared in PubMed. Articles were

597     identified and article metadata were extracted as described in Supplemental Section 2.

598    Repository data were extracted from the GitHub API as described in Supplemental Section 3.

599    Details are described in Supplemental Section 9.

600

601    **Availability of data and software**

602

603    All repository data extracted from the GitHub API, except file contents, are available at

604    https://doi.org/10.17605/OSF.IO/UWHX8. For file contents, in the absence of explicit open

605    source licenses for the majority of repositories studied, we recorded the Git URL for the specific

606    version of each file so that the exact dataset can be reconstructed using our downstream

607    scripts. Additionally, we have removed personal identifying information from commit records, but

608    have included API references for each commit record so that the full records can be

609    reconstructed. Software to generate the dataset and replicate the results in the paper is

610    available at https://github.com/pamelarussell/github-bioinformatics. See Supplemental Section 1

611    for details on the data and software.

612

613

614    # Acknowledgements

615

616    We thank Debashis Ghosh, Wladimir Labeikovsky, and Matthew Mulvahill for helpful

617    conversations and comments on the manuscript. We thank the GitHub support staff for their

618    effort in determining how we could work within the GitHub Terms of Service to publish a

619    reproducible study.

620

621

622    # Author contributions

623

624  PR: Conceptualization, Data Curation, Formal Analysis, Methodology, Project Administration,

625  Resources, Software, Supervision, Visualization, Writing - Original Draft Preparation, Writing -

626  Review & Editing.

627

628  RJ: Data Curation, Writing - Review & Editing.

629

630  SA: Conceptualization, Writing - Review & Editing.

631

632  BH: Data Curation, Investigation, Writing - Original Draft Preparation, Writing - Review &

633  Editing.

634

635  NC: Funding Acquisition, Project Administration, Supervision, Writing - Review & Editing.

636

637

638  **References**

639

640  1. Hagen JB. The origins of bioinformatics. Nat Rev Genet. 2000;1: 231–236.
641  doi:10.1038/35042090

642  2. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial
643  sequencing and analysis of the human genome. Nature. 2001;409: 860–921.
644  doi:10.1038/35057062

645  3. Scope Guidelines | Bioinformatics | Oxford Academic [Internet]. [cited 19 Mar 2018].
646  Available: https://academic.oup.com/bioinformatics/pages/scope_guidelines

647  4. Searls DB. The roots of bioinformatics. PLoS Comput Biol. 2010;6: e1000809.
648  doi:10.1371/journal.pcbi.1000809

649  5. Computational biology and bioinformatics - Latest research and news | Nature [Internet].
650  7 Mar 2018 [cited 24 Mar 2018]. Available: https://www.nature.com/subjects/computational-
651  biology-and-bioinformatics

652   6.  Hothorn T, Leisch F. Case studies in reproducibility. Brief Bioinform. 2011;12: 288–300.
653       doi:10.1093/bib/bbq084

654   7.  Sandve GK, Nekrutenko A, Taylor J, Hovig E. Ten simple rules for reproducible
655       computational research. PLoS Comput Biol. 2013;9: e1003285.
656       doi:10.1371/journal.pcbi.1003285

657   8.  Seemann T. Ten recommendations for creating usable bioinformatics command line
658       software. Gigascience. 2013;2: 15. doi:10.1186/2047-217X-2-15

659   9.  Leprevost F da V, Barbosa VC, Francisco EL, Perez-Riverol Y, Carvalho PC. On best
660       practices in the development of bioinformatics software. Front Genet. 2014;5: 199.
661       doi:10.3389/fgene.2014.00199

662   10.     Git [Internet]. [cited 24 Mar 2018]. Available: https://git-scm.com/

663   11.     Stack Overflow Developer Survey 2018. In: Stack Overflow [Internet]. [cited 18
664       Mar 2018]. Available: https://insights.stackoverflow.com/survey/2018/

665   12.     Build software better, together [Internet]. Github; Available: https://github.com

666   13.     GitHub Octoverse 2017 [Internet]. Github; Available:
667       https://octoverse.github.com/

668   14.     Blischak JD, Davenport ER, Wilson G. A Quick Introduction to Version Control
669       with Git and GitHub. PLoS Comput Biol. 2016;12: e1004668.
670       doi:10.1371/journal.pcbi.1004668

671   15.     Instructions for Authors | Bioinformatics | Oxford Academic [Internet]. [cited 27
672       Apr 2018]. Available:
673       https://academic.oup.com/bioinformatics/pages/instructions_for_authors

674   16.     GitHub API v3 | GitHub Developer Guide [Internet]. Github; Available:
675       https://developer.github.com/v3/

676   17.     GitHub Terms of Service - User Documentation [Internet]. Github; Available:
677       https://help.github.com/articles/github-terms-of-service/

678   18.     Ray B, Posnett D, Filkov V, Devanbu P. A large scale study of programming
679       languages and code quality in github. Proceedings of the 22nd ACM SIGSOFT
680       International Symposium on Foundations of Software Engineering. ACM; 2014. pp. 155–
681       165. doi:10.1145/2635868.2635922

682   19.     Kochhar PS, Bissyandé TF, Lo D, Jiang L. An Empirical Study of Adoption of
683       Software Testing in Open Source Projects. 2013 13th International Conference on Quality
684       Software. 2013. pp. 103–112. doi:10.1109/QSIC.2013.57

685   20.     Hu Y, Zhang J, Bai X, Yu S, Yang Z. Influence analysis of Github repositories.
686       Springerplus. 2016;5: 1268. doi:10.1186/s40064-016-2897-7

687   21.     Borges H, Hora A, Valente MT. Understanding the Factors that Impact the
688       Popularity of GitHub Repositories [Internet]. arXiv [cs.SE]. 2016. Available:
689       http://arxiv.org/abs/1606.04984

690    22.      Blincoe K, Sheoran J, Goggins S, Petakovic E, Damian D. Understanding the
691    popular users: Following, affiliation influence and leadership on GitHub. Information and
692    Software Technology. 2016/2;70: 30–39. Available:
693    http://www.sciencedirect.com/science/article/pii/S0950584915001688

694    23.      Ma W, Chen L, Zhou Y, Xu B. What Are the Dominant Projects in the GitHub
695    Python Ecosystem? 2016 Third International Conference on Trustworthy Systems and their
696    Applications (TSA). 2016. pp. 87–95. doi:10.1109/TSA.2016.23

697    24.      Sheoran J, Blincoe K, Kalliamvakou E, Damian D, Ell J. Understanding
698    "Watchers" on GitHub. Proceedings of the 11th Working Conference on Mining Software
699    Repositories. New York, NY, USA: ACM; 2014. pp. 336–339.
700    doi:10.1145/2597073.2597114

701    25.      Spotlight on Bioinformatics. NatureJobs. Nature Publishing Group; 2016;
702    doi:10.1038/nj0478

703    26.      Blei DM. Probabilistic Topic Models. Commun ACM. New York, NY, USA: ACM;
704    2012;55: 77–84. doi:10.1145/2133806.2133826

705    27.      McGill R, Tukey JW, Larsen WA. Variations of Box Plots. Am Stat. [American
706    Statistical Association, Taylor & Francis, Ltd.]; 1978;32: 12–16. doi:10.2307/2683468

707    28.      Boost C++ Libraries [Internet]. [cited 18 Mar 2018]. Available:
708    http://www.boost.org/

709    29.      Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput
710    Sequence Data [Internet]. [cited 18 Mar 2018]. Available:
711    http://www.bioinformatics.babraham.ac.uk/projects/fastqc/

712    30.      Strømgren C. Genderize.io | Determine the gender of a first name [Internet].
713    [cited 25 Jan 2018]. Available: https://genderize.io/

714    31.      Bonham KS, Stefan MI. Women are underrepresented in computational biology:
715    An analysis of the scholarly literature in biology, computer science and computational
716    biology. PLoS Comput Biol. 2017;13: e1005134. doi:10.1371/journal.pcbi.1005134

717    32.      Shannon CE. A mathematical theory of communication. The Bell System
718    Technical Journal. 1948;27: 379–423. doi:10.1002/j.1538-7305.1948.tb01338.x

719    33.      GitHub Privacy Statement - User Documentation [Internet]. Github; Available:
720    https://help.github.com/articles/github-privacy-statement/

721    34.      National Science Foundation, National Center for Science and Engineering
722    Statistics. Doctorate Recipients from U.S. Universities: 2016 [Internet]. Alexandria, VA.:
723    National Science Foundation; 2017. Report No.: Special Report NSF 18-304. Available:
724    https://www.nsf.gov/statistics/2018/nsf18304/

725    35.      Ortu M, Destefanis G, Counsell S, Swift S, Tonelli R, Marchesi M. How diverse is
726    your team? Investigating gender and nationality diversity in GitHub teams. J Softw Eng Res
727    Dev. Springer Berlin Heidelberg; 2017;5: 9. doi:10.1186/s40411-017-0044-y

728    36.      Nielsen MW, Alegria S, Börjeson L, Etzkowitz H, Falk-Krzesinski HJ, Joshi A, et

729     al. Opinion: Gender diversity leads to better science. Proc Natl Acad Sci U S A. 2017;114:
730     1740–1742. doi:10.1073/pnas.1700616114

731     37.      Vasilescu B, Posnett D, Ray B, van den Brand MGJ, Serebrenik A, Devanbu P,
732     et al. Gender and Tenure Diversity in GitHub Teams. Proceedings of the 33rd Annual ACM
733     Conference on Human Factors in Computing Systems. New York, NY, USA: ACM; 2015.
734     pp. 3789–3798. doi:10.1145/2702123.2702549

735     38.      Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, et al.
736     Bioconductor: open software development for computational biology and bioinformatics.
737     Genome Biol. 2004;5: R80. doi:10.1186/gb-2004-5-10-r80

738     39.      Prlić A, Yates A, Bliven SE, Rose PW, Jacobsen J, Troshin PV, et al. BioJava: an
739     open-source framework for bioinformatics in 2012. Bioinformatics. 2012;28: 2693–2695.
740     doi:10.1093/bioinformatics/bts494

741     40.      Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython:
742     freely available Python tools for computational molecular biology and bioinformatics.
743     Bioinformatics. 2009;25: 1422–1423. doi:10.1093/bioinformatics/btp163

744     41.      bioRxiv.org - the preprint server for Biology [Internet]. [cited 24 Mar 2018].
745     Available: https://www.biorxiv.org/

746     42.      Kamali AH, Giannoulatou E, Chen TY, Charleston MA, McEwan AL, Ho JWK.
747     How to test bioinformatics software? Biophys Rev. 2015;7: 343–352. doi:10.1007/s12551-
748     015-0177-3

749     43.      Yang A, Troup M, Ho JWK. Scalability and Validation of Big Data Bioinformatics
750     Software. Comput Struct Biotechnol J. 2017;15: 379–386. doi:10.1016/j.csbj.2017.07.002

751     44.      Chen TY, Ho JWK, Liu H, Xie X. An innovative approach for testing
752     bioinformatics programs using metamorphic testing. BMC Bioinformatics. 2009;10: 24.
753     doi:10.1186/1471-2105-10-24

754     45.      Duck G, Nenadic G, Filannino M, Brass A, Robertson DL, Stevens R. A Survey
755     of Bioinformatics Database and Software Usage through Mining the Literature. PLoS One.
756     2016;11: e0157989. doi:10.1371/journal.pone.0157989

757     46.      The Software Sustainability Institute | Software Sustainability Institute [Internet].
758     [cited 2 May 2018]. Available: https://www.software.ac.uk/

759     47.      The Software Sustainability Institute: changing research software attitudes and
760     practices | Software Sustainability Institute [Internet]. [cited 2 May 2018]. Available:
761     https://www.software.ac.uk/software-sustainability-institute-changing-research-software-
762     attitudes-and-practices

763     48.      Better Scientific Software [Internet]. [cited 2 May 2018]. Available:
764     https://bssw.io/pages/about

765     49.      Software Carpentry. In: Software Carpentry [Internet]. [cited 2 May 2018].
766     Available: https://software-carpentry.org/

767     50.      Wilson G. Software Carpentry: Getting Scientists to Write Better Code by Making

768  Them More Productive. Computing in Science and Engg. Piscataway, NJ, USA: IEEE
769  Educational Activities Department; 2006;8: 66–69. doi:10.1109/MCSE.2006.122

770  51.      Wilson G. Software Carpentry: lessons learned. F1000Res. 2014;3.
771  doi:10.12688/f1000research.3-62.v1

772  52.      Wilson G, Aruliah DA, Brown CT, Chue Hong NP, Davis M, Guy RT, et al. Best
773  practices for scientific computing. PLoS Biol. 2014;12: e1001745.
774  doi:10.1371/journal.pbio.1001745

775  53.      Crusoe MR, Brown CT. Walking the Talk: Adopting and Adapting Sustainable
776  Scientific Software Development processes in a Small Biology Lab. J Open Res Softw.
777  2016;4. doi:10.5334/jors.35

778  54.      Goble C. Better Software, Better Research. IEEE Internet Comput. 2014;18: 4–8.
779  doi:10.1109/MIC.2014.88

780  55.      Parnell LD, Lindenbaum P, Shameer K, Dall'Olio GM, Swan DC, Jensen LJ, et
781  al. BioStar: an online question & answer resource for the bioinformatics community. PLoS
782  Comput Biol. 2011;7: e1002216. doi:10.1371/journal.pcbi.1002216

783  56.      BigQuery - Analytics Data Warehouse | Google Cloud Platform. In: Google Cloud
784  Platform [Internet]. [cited 19 Mar 2018]. Available: https://cloud.google.com/bigquery/

785  57.      Kjetil Jacobsen MFXJO. PycURL Home Page [Internet]. [cited 19 Mar 2018].
786  Available: http://pycurl.io/

787  58.      Treat T. BigQuery-Python [Internet]. Github; Available:
788  https://github.com/tylertreat/BigQuery-Python

789  59.      Blei DM, Ng AY, Jordan MI. Latent Dirichlet Allocation. J Mach Learn Res.
790  2003;3: 993–1022. Available: http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf

791  60.      cloc [Internet]. Github; Available: https://github.com/AlDanial/cloc

792  61.      Wikipedia contributors. List of programming languages by type. In: Wikipedia,
793  The Free Encyclopedia [Internet]. 12 Dec 2017 [cited 15 Mar 2018]. Available:
794  https://en.wikipedia.org/w/index.php?title=List_of_programming_languages_by_type&oldid=
795  814994307

796  62.      Wikipedia contributors. Comparison of type systems. In: Wikipedia, The Free
797  Encyclopedia [Internet]. 5 Sep 2017 [cited 15 Mar 2018]. Available:
798  https://en.wikipedia.org/w/index.php?title=Comparison_of_type_systems&oldid=799049191

799

## Supporting information captions

801

**Supplemental Information. Supplemental information, methods, and figures.**

803

804    **Table S1. Definition of bioinformatics topics.**

805

806    **Table S2. Manual classification of articles as bioinformatics or not.**

807

808    **Table S3. Automatic identification of GitHub repository names in articles.**

809

810    **Table S4. Manual curation of GitHub repository names.**

811

812    **Table S5. High-profile repositories.**

813

814    **Table S6. Programming language type systems.**

815

816    **Table S7. Programming language execution modes.**

817

818    **Table S8. Calculated repository features.**

819