

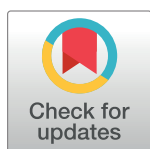
PERSPECTIVE

# Challenges and recommendations to improve the installability and archival stability of omics computational tools

Serghei Mangul<sup>1,2</sup>\*, Thiago Mosqueiro<sup>1,2</sup>, Richard J. Abdill<sup>3</sup>, Dat Duong<sup>1</sup>, Keith Mitchell<sup>1</sup>, Varuni Sarwal<sup>4</sup>, Brian Hill<sup>1</sup>, Jaqueline Brito<sup>5</sup>, Russell Jared Littman<sup>1</sup>, Benjamin Statz<sup>1</sup>¶, Angela Ka-Mei Lam<sup>1</sup>, Gargi Dayama<sup>3</sup>, Laura Grieneisen<sup>3</sup>, Lana S. Martin<sup>2</sup>, Jonathan Flint<sup>6</sup>, Eleazar Eskin<sup>1,7</sup>, Ran Blekhman<sup>3,8</sup>

**1** Department of Computer Science, University of California Los Angeles, Los Angeles, California, United States of America, **2** Institute for Quantitative and Computational Biosciences, University of California Los Angeles, Los Angeles, California, United States of America, **3** Department of Genetics, Cell Biology, and Development, University of Minnesota, Minneapolis, Minnesota, United States of America, **4** Indian Institute of Technology Delhi, Hauz Khas, New Delhi, India, **5** Institute of Mathematics and Computer Science, University of São Paulo, São Paulo, Brazil, **6** Center for Neurobehavioral Genetics, Semel Institute for Neuroscience and Human Behavior, University of California Los Angeles, Los Angeles, California, United States of America, **7** Department of Human Genetics, University of California Los Angeles, Los Angeles, California, United States of America, **8** Department of Ecology, Evolution, and Behavior, University of Minnesota, Minnesota, United States of America

\* These authors contributed equally to this work.  
 ¶ Authorship confirmed by corresponding author.  
 \* [smangul@ucla.edu](mailto:smangul@ucla.edu)



**OPEN ACCESS**

**Citation:** Mangul S, Mosqueiro T, Abdill RJ, Duong D, Mitchell K, Sarwal V, et al. (2019) Challenges and recommendations to improve the installability and archival stability of omics computational tools. *PLoS Biol* 17(6): e3000333. <https://doi.org/10.1371/journal.pbio.3000333>

**Published:** June 20, 2019

**Copyright:** © 2019 Mangul et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** Raw data and parsed HTTP information for each link are available at <https://doi.org/10.6084/m9.figshare.7641083>. Our protocol to check the archival stability of published software tools is available at <https://github.com/smangul1/good.software>.

**Funding:** SM acknowledges support from a QCB Collaboratory Postdoctoral Fellowship and the QCB Collaboratory community directed by Matteo Pellegrini. SM and EE are supported by National Science Foundation grants 0513612, 0731455, 0729049, 0916676, 1065276, 1302448, 1320589, 1331176, and 1815624 and National Institutes of Health grants K25-HL080079, U01-DA024417, P01-HL30568, P01-HL28481, R01-GM083198, R01-ES021801, R01-MH101782, and R01-ES022282. RB is grateful for support from the

## Abstract

Developing new software tools for analysis of large-scale biological data is a key component of advancing modern biomedical research. Scientific reproduction of published findings requires running computational tools on data generated by such studies, yet little attention is presently allocated to the installability and archival stability of computational software tools. Scientific journals require data and code sharing, but none currently require authors to guarantee the continuing functionality of newly published tools. We have estimated the archival stability of computational biology software tools by performing an empirical analysis of the internet presence for 36,702 omics software resources published from 2005 to 2017. We found that almost 28% of all resources are currently not accessible through uniform resource locators (URLs) published in the paper they first appeared in. Among the 98 software tools selected for our installability test, 51% were deemed “easy to install,” and 28% of the tools failed to be installed at all because of problems in the implementation. Moreover, for papers introducing new software, we found that the number of citations significantly increased when authors provided an easy installation process. We propose for incorporation into journal policy several practical solutions for increasing the widespread installability and archival stability of published bioinformatics software.

National Institutes of General Medicine (R35-GM128716) and a McKnight Land-Grant Professorship from the University of Minnesota.

**Competing interests:** The authors have declared that no competing interests exist.

**Abbreviations:** CRAN, Comprehensive R Archive Network; FTP, file transfer protocol; HTTP, hypertext transfer protocol; HTTPS, hypertext transfer protocol secure; NCBI, National Center for Biotechnology Information; OTU, operational taxonomic unit; SV, structural variant; URL, uniform resource locator; UX, user experience; WGS, whole-genome sequencing; XML, extensible markup language; YUM, Yellowdog Updater, Modified.

**Provenance:** Commissioned; Externally peer reviewed.

## Introduction

During the past decade, the rapid advancement of genomics and sequencing technologies has inspired a large and diverse collection of new algorithms in computational biology [1,2]. In the last 15 years, the amount of available genomic sequencing data has doubled every few months [3,4]. Life-science and biomedical researchers are leveraging computational tools to analyze this unprecedented volume of genomic data [3,4], which has been critical in solving complex biological problems and subsequently laying the essential groundwork for the development of novel clinical translations [5]. The exponential growth of genomic data has reshaped the landscape of contemporary biology, making computational tools a key driver of scientific research [6,7].

As computational and data-enabled research become increasingly popular in biology, novel challenges arise, accompanied by standards that attempt to remedy them. One such challenge is computational reproducibility—the ability to reproduce published findings by running the same computational tool on the data generated by the study [8–10]. Although several journals have introduced requirements for the sharing of data and code, there are currently no effective requirements to promote installability and long-term archival stability of software tools, creating situations in which researchers share source code that either doesn't run or disappears altogether. These issues can limit the applicability of the developed software tools and impair the community's ability to reproduce results generated by software tools in the original publication.

The synergy between computational and wet-lab researchers is especially productive when software developers distribute their tools as packages that are easy to use and install [11]. Though many new tools are released each year, comparatively few incorporate adequate documentation, presentation, and distribution, resulting in a frustrating situation in which existing tools address every problem except how to run them [12].

Widespread support for software installability promises to have a major impact on the scientific community [13], and practical solutions have been proposed to guide the development of scientific software [14–17]. Although the scale of this issue in computational biology has yet to be estimated, the bioinformatics community warns that poorly maintained or improperly implemented tools will ultimately hinder progress in data-driven fields like genomics and systems biology [3,7,18].

## Challenges to effective software development and distribution in academia

Successfully implementing and distributing software for scientific analysis involves numerous unique challenges that have been previously outlined by other scholars [11,15,16,19]. In particular, fundamental differences between software development workflows in academia and in industry challenge the installability and archival stability of novel tools developed by academics. These differences can be broken down into three broad categories:

1. Software written by researchers tends to be written with the idea that users will be knowledgeable about the code and appropriate environment and dependencies. This sometimes results in tools that are difficult to install, with instructions and command-line options that are unclear and confusing but are also critical for the tool's function.
2. Academic journals are a primary source for information and documentation of noncommercial scientific software, even though the static nature of publications means this vital information quickly falls out of date.

3. Incentives in academia heavily favor the publication of new software, not the maintenance of existing tools.

First, software developers in industrial settings receive considerably more resources for developing user-friendly tools than their counterparts in academic settings [20]. Commercial software is developed by large teams of software engineers that include specialized user experience (UX) developers. In academic settings, software is developed by smaller groups of researchers who may lack formal training in software engineering, particularly UX and cross-platform design. Many computational tools lack a user-friendly interface to facilitate the installation or execution process [12]. Developing an easy-to-use installation interface is further complicated when the software relies on third-party tools that need to be installed in advance, called “dependencies.” Installing dependencies is an especially complicated process for researchers with limited computational knowledge. Well-defined UX standards for software development could help software developers in computational biology promote widespread implementation and use of their newly developed computational tools.

Second, companies efficiently distribute industry-produced software using dedicated company units or contractors—services that universities and scientific funding agencies do not typically provide for academic-developed software. The computational biology community has adopted by default a pragmatic, short-term framework for disseminating software development [21], which generally consists of publishing a paper describing the software tool in a peer-reviewed journal. So-called methods papers are dedicated to explaining the rationale behind the novel computational tool and demonstrating its efficacy with sample datasets. Supplemental materials such as detailed instructions, tutorials, dependencies, and source code are made available on the internet and included in the published paper as a uniform resource locator (URL), but they generally exist in a location out of the journal’s direct control. The quality, format, and long-term availability of supplemental materials varies among software developers and is subject to less scrutiny in the peer-review process compared with the published paper itself. This approach limits the installability of software tools for use in research and hinders the community’s ability to evaluate the tools themselves [22].

Third, the academic structures of funding, hiring, and promotion offer little reward for continuous, long-term development and maintenance of tools and databases [23], and software developers can lose funding for even the most widely used tools. Loss of external funding can slow and even discontinue software development, potentially impacting the research productivity of studies that depend on these tools [24]. Interrupted development also hinders the ability to reproduce results from published studies that use discontinued tools. In general, industry-developed software is supported by teams of software engineers dedicated to developing and implementing updates for as long as the software is considered valuable. Many software developers in academia do not have access to mechanisms that could ensure a similar level of maintenance and stability.

We combined two approaches to determine the effects of these challenges on the proportion of bioinformatics tools that could be considered user friendly. First, we investigated tens of thousands of URLs corresponding to bioinformatics tools and resources to determine whether they are archivally stable—whether users can even reach the websites described in the papers evaluated. Next, we investigated the number of tools that provided an easy-to-use installation interface to download and install the software and any required dependencies.

### **Archival stability of published computational tools and resources**

The World Wide Web provides a platform of unprecedented scope for data and software archival stability, yet long-term preservation of online resources remains a largely unsolved

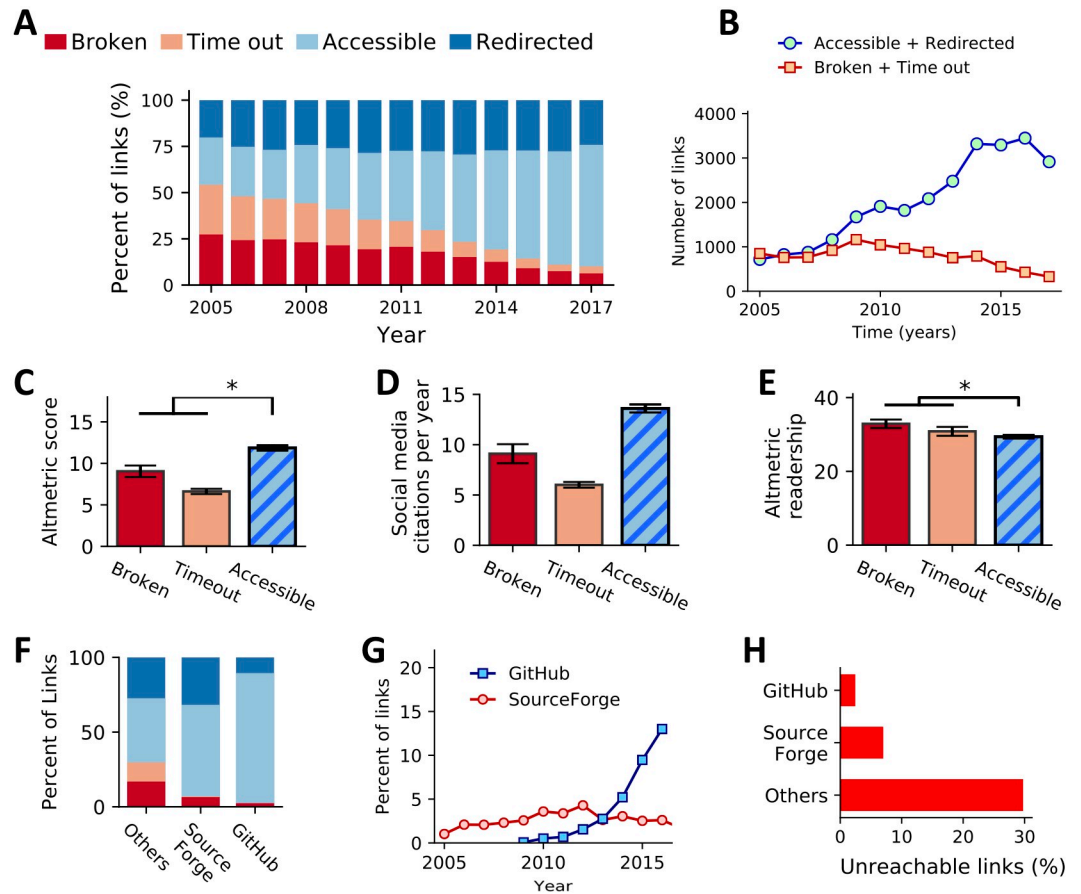
problem [25]. Published software tools are made accessible through the URL, which is typically provided in the abstract or main text of the paper and is often assumed to be a practically permanent locator. However, a URL may become inactive because of the removal or reconfiguration of web content. The “death of URLs” [26] has been described for decades in various terms, including “link rot” [27] and “lost internet references” [28]. At the onset, the World Wide Web promised the virtually infinite availability of digital resources; in practice, many are lost. For example, many tools in computational biology are hosted on academic web pages that become inactive with time, sometimes only months after their initial publication. These software packages are typically developed by small groups of graduate students or postdoctoral scholars who, considering the temporary nature of such positions, cannot maintain such websites and software for longer periods of time.

Multiple studies have identified the deterioration of long-term archival stability of published software tools [18,26,28–31]. In order to begin assessing the magnitude of these issues in computational biology, we comprehensively evaluated the archival stability of computational biology tools used in 51,236 biomedical papers published across 10 relevant peer-reviewed journals over a span of 18 years, from 2000 to 2017 (S1 Table). Out of the 51,236 examined papers, 13.6% contained at least one URL in their abstracts, and another 38.3% contained URLs in the body of the paper. To increase the likelihood that the identified URL corresponds to a software tool or database, we inspected 10 neighboring words for specific keywords commonly used, including “pipeline,” “code,” “software,” “available,” “publicly,” and others (see Methods section). Complete details on our methodology for extracting the URLs, including all options and thresholds, are provided in the Methods section.

We used a web-mining approach to test 36,702 published URLs that our survey identified. We categorized unreachable URLs into two groups: unreachable due to connection time-out and unreachable due to error (“broken” links, i.e., 404 hypertext transfer protocol [HTTP] status). We separately categorized accessible URLs that returned immediately and those that used redirection—that is, URLs to which servers responded by pointing the user to a new URL that then connects successfully. We found that 26.7% of evaluated URLs are successfully redirected to new URLs. (Some URLs were redirected to pages that subsequently returned an error; these were not considered successful.) Of all identified URLs, 11.9% were unreachable because of connection time-outs, and 15.9% were “broken.” To prevent erroneous classification caused by configuration of our automated tests, we manually verified more than 900 URLs reported as “time-outs,” or requests that did not receive a response within an acceptable amount of time (S1 Fig).

Next, we grouped the URLs by the year in which the computational biology tool was first referenced in a publication. As expected, the time since publication is a key predictor of URL archival stability (Fisher exact test,  $p$ -value  $< 10^{-15}$ ). In total, 41.9% of the software referenced before 2012 ( $n = 15,439$ ) is unavailable, whereas only 17.5% of the recent software referenced in 2012 and later ( $n = 21,263$ ) is unavailable (Fig 1A). After 2013, we observe a drop in the absolute number of archivally unstable resources (Fig 1B). Despite the strong decline in the percentage of missing resources over time, there are still 200 resources published every year with links that were broken by the time we tested them. The data and scripts for reproducing the plots in Fig 1 are available at <https://github.com/smangul1/good.software>.

Prior research demonstrates that the availability of published bioinformatics resources has a significant impact on citation counts [30]. In addition to those generally accepted measures of scientific impact, we assessed the effect of software availability on complementary metrics of impact, such as measures of social media mentions, media coverage, and public attention (Fig 1C–1E). We found that papers with accessible links exhibit increased engagement by readers in social media, reflected in a significantly higher number of citations in social media platforms



**Fig 1. Archival stability of 36,702 published URLs across 10 systems and computational biology journals over the span of 13 years.** An asterisk (\*) denotes categories that have a difference that is statistically significant. Error bars, where present, indicate SEM. (A) Archival stability status of all links evaluated from papers published between 2005 and 2017. Percentages of each category (y-axis) are reported over a 13-year span (x-axis). (B) A line graph comparing the overall numbers (y-axis) of functional (green circles) and nonfunctional (orange squares) links observed in papers published over time (x-axis). (C) A bar chart showing the mean Altmetric “attention score” (y-axis) for papers, separated by the status of the URL (x-axis) observed in that paper. (D) A bar chart showing the mean number of mentions of papers in social media (blog posts, Twitter feeds, etc.) according to Altmetric, divided by the age of the paper in years (y-axis). Papers are separated by the status of the URL (x-axis) found in the paper. (E) A bar chart illustrating the mean Altmetric readership count per year of papers (y-axis) containing URLs in each of the categories (x-axis). (F) The proportion of unreachable links (due to connection time-out or due to error) stored on web services designed to host source code (e.g., GitHub and SourceForge) and “Other” web services. (G) A line plot illustrating the proportion (y-axis) of the total links observed in each year (x-axis) that point to GitHub or SourceForge. (H) A bar chart illustrating the proportion of links hosted on GitHub or SourceForge (vertical axis) that are no longer functional (horizontal axis) compared with links hosted elsewhere. SEM, standard error of the mean; URL, uniform resource locators.

<https://doi.org/10.1371/journal.pbio.3000333.g001>

per year (Fig 1D; Kruskal–Wallis  $p$ -value =  $1.75 \times 10^{-161}$ , Dunn’s test  $p$ -value =  $9.66 \times 10^{-103}$  for accessible versus broken,  $p$ -value =  $3.37 \times 10^{-76}$  for accessible versus time-out, adjusted for multiple tests using the Benjamini–Hochberg procedure) and an increased Altmetrics score [32] when compared with papers with “broken” and “time-out” links (Fig 1C; Kruskal–Wallis,  $p$ -value =  $1.66 \times 10^{-25}$ , Dunn’s test  $p$ -value =  $2.47 \times 10^{-17}$  for accessible versus broken,  $p$ -value =  $4.16 \times 10^{-14}$  for accessible versus Time-out). Although the difference is small, we found the readership of papers with accessible links differed significantly from papers with links that are classified as broken or time-outs—surprisingly, the median reader count per year (according to Altmetric) was lower for papers with accessible links (Fig 1E; Dunn’s test

$p$ -value =  $1.17 \times 10^{-6}$  for accessible versus broken,  $p$ -value =  $8.42 \times 10^{-15}$  for accessible versus time-out).

In addition, we tested the impact of using websites designed to host source code, such as GitHub and SourceForge, on the archival stability of bioinformatics software. These websites have been used by the bioinformatics community since 2001, and the proportion of software tools hosted on these sites has grown substantially, from 1.6% in 2012 to 13% in 2016 (Fig 1G). We find that URLs pointing to these websites have a high rate of archival stability: 97.6% of the links to GitHub and 93.0% of the links to SourceForge are accessible, whereas only 70.3% of links hosted elsewhere are accessible (Fig 1H)—a significant difference (Fisher exact test GitHub versus Others,  $p$ -value =  $2.70 \times 10^{-106}$ ; SourceForge versus Others,  $p$ -value =  $1.31 \times 10^{-5}$ ).

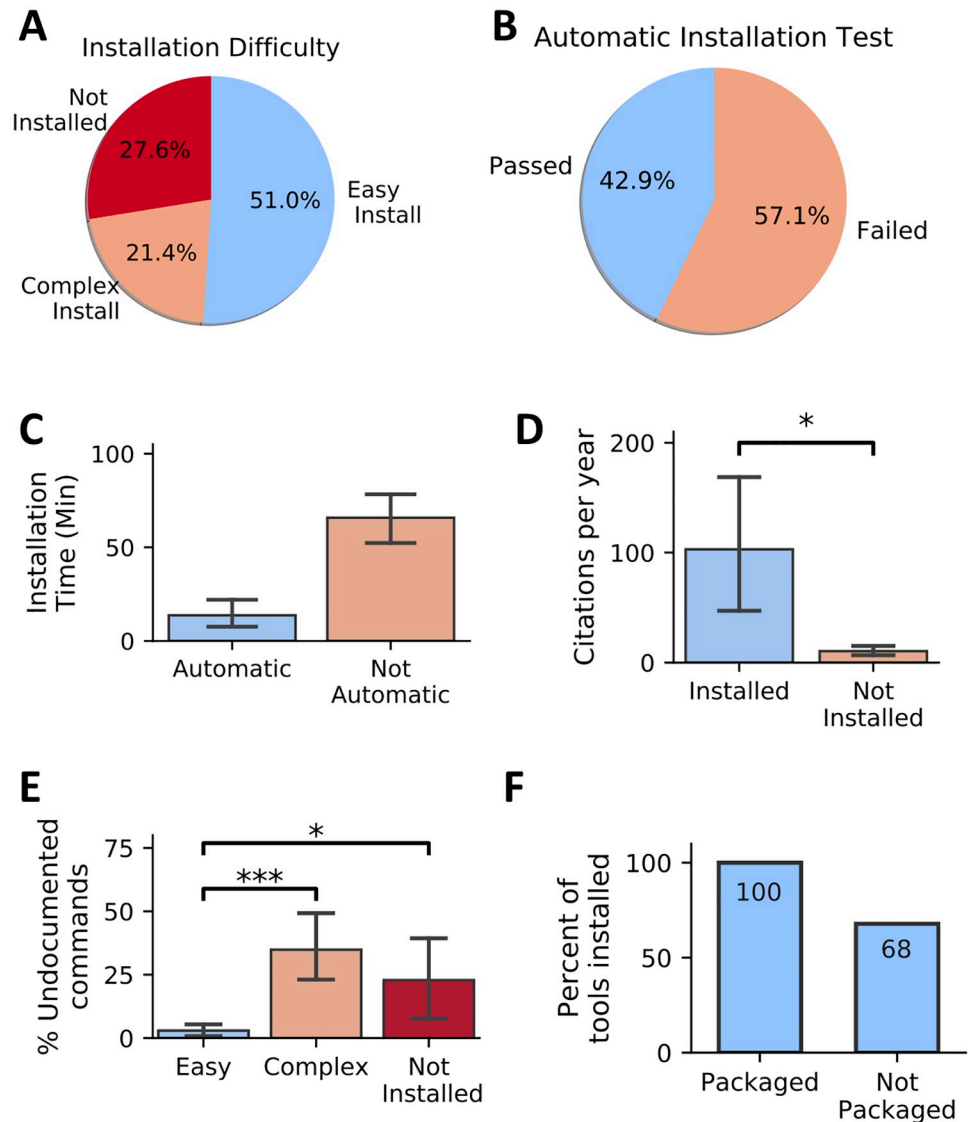
Our results suggest that the computational biology community would benefit from such approaches, which effectively guarantee permanent access to published scientific URLs. Specifically, several key principles emerge that promise to positively impact the availability of published bioinformatics resources, including the number of citations and social media references. In addition, bioinformatics tools and resources stored on web services designed to host source code have a significantly higher chance of remaining accessible.

## Installability of published software tools

We developed a computational framework capable of systematically verifying the archival stability and installability of published software tools. We applied this framework to 98 randomly selected tools across various domains of computational biology (Methods section). Those tools were selected independently from the 36,702 URLs used previously (Archival stability of published computational tools and resources). We engaged undergraduate and graduate students to run the installation test using a standardized protocol (S2 Fig); we recorded the time required to install the tools and other important features, allowing up to 2 hours per software package. In total, 71 hours of installation time was required in attempts to install 98 tools. We categorized a tool as “easy to install” if it could be installed in 15 minutes or less, “complex” if it required more than 15 minutes but was successfully installed before the 2-hour limit, and “not installed” if the tool could not be successfully installed within 2 hours (S2 Table and Fig 2).

The most stringent evaluation was the “automatic installation test,” in which the tester is required to strictly follow the instructions provided in the manual of the software tool (Methods; Fig 2A)—we determined that 57.1% of the selected tools failed this test. The vast majority (39 out of 42) of the tools that passed this test finished in fewer than 15 minutes and were classified as “easy to install” (S2 Table). For the tools failing the test, we performed manual intervention during which the tester was allowed to install missing dependencies and modify code to resolve installation errors. On average, it took an additional 70 minutes to install tools failing the “automatic installation test” (Mann–Whitney  $U$  test,  $p$ -value =  $4.7 \times 10^{-9}$ ; Fig 2C). Manual intervention was unsuccessful for 66% of the tools that initially failed the automatic installation test. Failed manual installation was due to numerous issues, including hard-coded options, invalid folder paths or header files, and usage of unavailable software dependencies.

Next, we assessed the effect of the ease of installation on the popularity of tools in the computational biology community by investigating the number of citations for the paper describing the software tools. We find that tools that we were able to install had significantly more citations compared with tools that we were not able to successfully install within 2 hours (Fig 2D; Mann–Whitney  $U$  test,  $p$ -value = 0.032). These results suggest, perhaps not surprisingly, that tools that are easier to install are more likely to be adopted by the community.



**Fig 2. Installability of 98 randomly selected published software tools across 22 life-science journals over a span of 15 years.** Error bars, where present, indicate SEM. (A) Pie chart showing the percentage of tools with various levels of installability. (B) A pie chart showing the proportion of evaluated tools that required no deviation from the documented installation procedure. (C) Tools that require no manual intervention (pass automatic installation test) exhibit decreased installation time. (D) Tools installed exhibit increased citation per year compared with tools that were not installed (Kruskal–Wallis,  $p$ -value = 0.035). (E) Tools that are easy to install include a decreased portion of undocumented commands (Not Installed versus Easy Install: Mann–Whitney  $U$  test,  $p$ -value = 0.01, Easy Install versus Complex Install: Mann–Whitney  $U$  test,  $p$ -value =  $8.3 \times 10^{-8}$ ). (F) Tools available in well-maintained package managers such as Bioconda were always installable, whereas tools not shipped via package managers were prone to problems in 32% of the studied cases. SEM, standard error of the mean.

<https://doi.org/10.1371/journal.pbio.3000333.g002>

In addition, we aimed to see whether the accuracy of a tool’s installation instructions affects its installation time. Considering the proportion of commands that are undocumented (estimated as a ratio between the executed commands and commands in the manual), we find that tools with easier installation have a significantly lower percentage of undocumented commands (Fig 2E; Mann–Whitney  $U$  test,  $p$ -value = 0.04). Considering a significant increase of installation time and a low rate of success for tools failing automatic installation test, we argue

that reliance on manual intervention to successfully install and run computational biology tools is an unsustainable practice. Software developers would benefit from ensuring a simple installation process and providing adequate installation instructions.

The vast majority of surveyed tools fail to provide one-line solutions for installation, instead providing step-by-step instructions. On average, eight commands were required to install surveyed tools, whereas only 3.9 commands were provided in the manual. Among the surveyed software tools, 23 tools provide one-line installation solutions that worked successfully, of which nine were available via the Bioconda package manager [33] (S2 Table). A package manager is a system that automates the installation, upgrade, and configuration of a collection of software tools in a consistent manner. Tools with single-command installations require on average 6 minutes of installation time, which is significantly faster when compared with tools that require multicommand installation (Kruskal–Wallis,  $p$ -value =  $4.7 \times 10^{-6}$ ) (S3 Fig). Tools available in well-maintained package managers (e.g., Bioconda) were always installable, whereas tools not shipped via package managers failed to install in 32% of the studied cases (Fig 2F). The results from our study point to several specific opportunities for establishing an effective software development and distribution practice (Box 1).

## Discussion

Our study assesses a critical issue in computational biology that is characterized by lack of standards regarding installability and long-term archival stability of omics computational tools and resources. Despite recent requirements on the behalf of journals to impose data and code sharing on published authors' work, 27.8% of 36,702 omics software resources examined in this study are not currently accessible via the original published URLs. Among the 98 software packages selected for our installability test, 49.0% of omics tools failed our “easy-to-install” test. In addition, 27.6% of surveyed tools could not be installed because of severe problems in the implementation process. One-quarter of examined tools are easy to install and use; in these cases, we identify a set of good practices for software development and dissemination.

Reviewers assessing the papers that present new software tools could begin addressing this problem with the adoption of a rigorous, standardized approach during the peer-review process. Feasible solutions for improving the installability and archival stability of peer-reviewed software tools include requirements for providing installation scripts, test data, and functions that allow automatic checks for the plausibility of installing and running the tool. For example, “forking” is a simple procedure that ensures the version of cited code within an article may persist beyond initial publication [40]. Academic journals recently took a major step toward improving archival stability by permanently forking published software on GitHub (e.g., [41]).

The current workflow of computational biology software development in academia encourages researchers to develop and publish new tools, but this process does not incentivize long-term maintenance of existing tools. Results from this study provide a strong argument for the development of standardized approaches capable of verifying and archiving software. Furthermore, our results suggest that funding agencies should emphasize support for maintenance of existing tools and databases.

Manual interventions and long installation times are unappealing to many users, especially to those with limited computational skills. Many life-science and medical researchers lack formal computational training and may be unable to perform manual interventions (e.g., installing dependencies or editing computer code during installation). Users could leverage advanced knowledge of the time and computational skills required to properly



## Box 1. Principles to increase installability and archival stability of omics computational tools and resources

Here, we present eight principles to increase the installability and archival stability of omics computational tools and resources. The tool was considered installable if the tool and its corresponding dependencies can be installed on Linux/UNIX-based operating systems and the tool can produce expected results from the input data with no errors. The majority of surveyed software tools and resources address only a portion of these principles.

### 1. Host software and resources on archivally stable services

Selecting the appropriate service to host your software and resources is critical. A simple solution is to use web services designed to host source code (e.g., GitHub [34,35] or SourceForge). In our study, we have determined that more than 96% of software tools and resources stored at GitHub or SourceForge are accessible, and tools hosted on these services remain stable for longer periods of time (S3 Table). Ideally, the repositories storing code should also be more permanently archived using a service such as Zenodo (<https://zenodo.org>), which is designed to provide long-term stability for scientific data.

### 2. Provide easy-to-use installation interface

Use sustainable and comprehensive software distribution. One example of a sustainable package manager is Bioconda [33], which is language agnostic and available on Linux and Mac operating systems. The package manager is a collection of software tools that automate the installation of a tool's core package and updates in a consistent manner. Package managers also help solve the "dependencies problem" by automatically installing required third-party software packages. Bioconda, technically a "channel" within the broader Conda project, is one the most popular package managers in bioinformatics, currently covering 2,900 software tools that are continuously maintained, updated, and extended by a growing global community [33]. Bioconda provides a one-line solution for downloading and installing a tool.

### 3. Take care of all the dependencies the tool needs

Even the most widely used tools rely on dependencies. To facilitate simple installation, provide an easy-to-use interface to download and install all required dependencies. Ideally, all necessary installation instructions should be included in a single script, especially when the number of installation commands is large. Package managers can potentially make this problem easier to solve. Bioconda also automatically generates containers for each Bioconda "recipe" [36], which provides all files and information needed to install a package. Other implementations of containerized software (e.g., Docker [<https://www.docker.com/>] and Singularity [<https://www.sylabs.io/docs/>]) also usually have all dependencies preinstalled. Often, language-specific solutions are also available (e.g., Bioconductor [37] and the Comprehensive R Archive Network [CRAN]). One drawback of Bioconda is that the existing tools in portable package managers are manually updated by the team or community, often delaying such updates. For example, as of August 10, 2018, R 3.5 was unavailable under Bioconda despite being released almost 4 months prior. If possible, one should design an installation script combining the commands for installing dependencies and developed software tools into a single script. Ideally, these dependencies should be installed in a user-configurable directory, as with Python "virtual environments," which can help avoid conflicts with existing software on the system.

#### 4. Provide an example dataset

Provide an example dataset inside the software package with a description of the expected results. Similar to unit- and integration-testing practices in software engineering, example datasets allow the user to verify that the tool was successfully installed and works properly before running the tool on experimental data. A tool may be installed with no errors, yet it may still fail to successfully run on the input data. Only 68% of examined tools provide an example dataset ([S2 Table](#)).

#### 5. Provide a “quick start” guide

Allow the user to verify the installation and performance of the tool. Providing a “quick start” guide is the best way for the user to validate that the tools are installed and working properly. The guide should provide the commands needed to download, install, and run the software tool on the example dataset. An example of a “quick start” guide is provided in [S1 Text](#). In addition to the “quick start,” a detailed manual must be provided with information on options, advanced features, and configuration. Best practices for creating bioinformatics software documentation are discussed elsewhere [38].

#### 6. Choose an adequate name

Choose a software name that best reflects the developed tool or resource. Today’s “age of Google” places new demands on the function of tool names, which should be memorable and unique, yet easily searchable. In addition, there are no regulations on tool names. For example, there are at least six tools named “Prism,” making it challenging to find the right tool ([S3 Text](#)). Scout the web to check the uniqueness of a name before publishing a new tool.

#### 7. Assume no root privileges

Tools are often installed on high-performance computing clusters in which users do not have administrative (root/superuser) privileges to install software into system directories. When developing instructions for installation of the proposed software tool, avoid commands that require root access. Examples of such commands include those that use package managers that require root/superuser privileges, such as “apt-get install” or “yum install.”

#### 8. Make platform-agnostic decisions when possible

Create tools that will work on as many systems as possible—specification of various versions of UNIX-like systems may limit the installability of software. Design your software to minimize reliance on operating system-specific functionality to make it easier for users to use your tools in diverse environments. Platform-specific installation commands (e.g., Homebrew [39]) should also be avoided.

install a software package. We propose a prototype of a badge server that runs an automated installation test, thus introducing to the peer-review process explicit assessment of a tool’s installability ([S2 Text](#)). This badge server would be particularly useful in computational biology, an interdisciplinary field composed of reviewers who often lack the skills and time to

verify the installability of software tools. Many benchmarking studies already routinely report relative ease of installation and use of new tools as components of their performance metrics [42].

## Methods

### Protocol to check the archival stability of published software tools

We downloaded open-access papers via PubMed from 10 systems and computational biology journals from the National Center for Biotechnology Information (NCBI) file transfer protocol (FTP) server (<ftp://ftp.ncbi.nlm.nih.gov/pub/pmc/>). We included the following journals: *Nature Biotechnology*, *Genome Medicine*, *Nature Methods*, *Genome Biology*, *BMC Systems Biology*, *Bioinformatics*, *PLOS Computational Biology*, *BMC Bioinformatics*, *BMC Genomics*, and *Nucleic Acids Research*.

Papers were downloaded in extensible markup language (XML) format, which contains name tags for field extraction. (Raw data from PubMed are available at <https://doi.org/10.6084/m9.figshare.7641083>.) Specifically, we focused on three tags: <abstract>, <body>, and <text-link>. Each paper's abstract is enclosed inside the <abstract> tag (S1 Fig). The <body> tag contains the key contents, like introduction, methods, results, and discussion. The <ext-link> tags contain internet addresses for external sources (e.g., supplementary data and directions for downloading data sources and software packages). We have prepared a folder containing a small set of papers in XML format for testing purposes, available at [https://github.com/smangul1/good.software/blob/master/download.parse.data/Nat\\_Methods.tar.gz?raw=true](https://github.com/smangul1/good.software/blob/master/download.parse.data/Nat_Methods.tar.gz?raw=true).

We deployed a heuristic approach to limit links to software produced by each paper's authors. We assumed that these links are in <ext-link> tags whose neighbor words contain one of the following keywords: "here," "pipeline," "code," "software," "available," "publicly," "tool," "method," "algorithm," "download," "application," "apply," "package," and "library." We searched for these words in a neighborhood that extended 75 characters from both the start and end of each <ext-link> tag.

For each extracted link, we initially used the `HTTPError` class of the Python library `urllib2` to get the HTTP status. Status numbers 400 and above indicate broken links; for example, the well-known 404 code indicates "Page Not Found." Some URLs point at servers that did not respond at all. Because the threshold for the allotted time to wait for a response may bias the results, we manually verified 931 URLs reported with the time-out error code (S1 Fig).

Multiple attempts were made to validate each extracted URL: First, an HTTP request was sent to each URL; if that was not successful, an FTP request was sent to avoid marking URLs as "broken" if they used this older method of transferring files instead. HTTP requests that received "redirect" responses (status codes 300–399) were followed to the end point specified by the redirection (or redirections) to determine the final destination of the request. If the request ultimately completed successfully, the initial redirect code was recorded, and that link appears in our data as a redirection. However, some requests eventually resulted in errors—for example, if a server rewrites a received URL according to a formula, but the rewritten URL points to a file that doesn't exist. Redirections that eventually resulted in an error were recorded with that error code instead. There is only one exception to this classification: if a server responded with a redirection status, but the redirection pointed at a URL that only changes the URL's protocol from "HTTP" to "hypertext transfer protocol secure (HTTPS)," we classified this as a "success" rather than a "redirection." Our protocol to check the archival

stability of published software tools is available at <https://github.com/smangul1/good.software>. Parsed HTTP information for each link is available at <https://doi.org/10.6084/m9.figshare.7738901>.

### Protocol to check the installability of published software tools

To standardize the operating system environment for each tool installation, we used a CentOS 7 (v1710.01) Vagrant virtual machine. CentOS is an open-source operating system that is widely used in research computing. To prevent dependency mismatches caused by previously installed packages, we installed each tool in a new Vagrant virtual machine. Our virtual machine was provisioned with several commonly used software tools using the Yellowdog Updater, Modified (YUM) package manager to accommodate low-level dependencies that many developers would assume were already installed: epel-release, java (version 1.8.0), wget, vim, unzip, gcc and gcc-devel, python, and R. Users seeking to replicate this environment can use the Vagrant provisioning script found here: <https://github.com/smangul1/good.software/blob/master/toolInstall/Vagrantfile>.

We present a summary of our protocol in [S3 Fig](#). Tools were classified into three categories: (1) easy to install, when installation took fewer than 15 minutes; (2) hard to install, when installation took between 15 minutes and 2 hours; and (3) not installed, meaning installation took longer than 2 hours or could not be completed. We tested a total of 98 tools across various categories and fields as described in the following sections: Tools for microbiome profiling, Tools for read alignment, Tools for variant calling tools, Tools for structural variants tools, and Additional omics tools. Information on the tools tested and the results of the test are available in [S2 Table](#) and are shared at <https://doi.org/10.6084/m9.figshare.7738949>.

### Tools for microbiome profiling

The installability of 10 common tools for microbiome analysis was tested. To develop a list of popular tools, two coauthors independently made lists of 30 tools currently used for microbiome data processing based on a literature survey and identified those present on both lists. Microbiome tools can vary in their specificity of use; we limited the final tool list to five tools that process raw sequences into a final operational taxonomic unit (OTU) table and five tools capable of broad downstream analysis functions.

### Tools for read alignment

We tested the installability of 10 tools for read alignment. We randomly selected a total of 20 tools—10 tools from a recent survey [43] and 10 tools from PubMed (<https://www.ncbi.nlm.nih.gov/pubmed/>). The full list of extracted URLs is available at <https://github.com/smangul1/good.software>. To confirm that the installation process indeed worked, we used reads generated from the complete genome of *Enterobacteria phage lambda* (NC\_001416.1).

### Tools for variant calling tools

We tested the installability of seven randomly sampled tools designed for variant calling [44]. We confirmed successful software installation when the core functionality of each package could be executed with an example dataset. Only one of the tools was not packaged with an example dataset, in which case we randomly chose an open example dataset. We discarded from our study the tools for which papers could not be located.

## Tools for structural variants tools

We examined the installability of 52 common tools used for the structural variant (SV) calling from whole-genome sequencing (WGS) data. First, we compiled a list of tools that use read alignment, in which reads aligned to the locations are inconsistent with the expected insert size of the library or expected read depth at a specific locus. We randomly selected 50 tools out of 70 programs designed to detect SVs from WGS data and published after 2011. We confirmed the successful installation of each software package by executing its core functionality with an example dataset.

## Additional omics tools

Lastly, we randomly selected 20 published tools based on the URL present in the abstract or the body of the publications available in PubMed (<https://www.ncbi.nlm.nih.gov/pubmed/>). The full list of extracted URLs is available at <https://github.com/smangul1/good.software>.

## Statistical analysis

Once the archival information was recorded, variance analysis was performed to assess the differences among the links categorized as “accessible,” “redirected,” “broken,” and “time out” as they related to four paper-level metrics: the number of citations in the original paper in which the tool was published; number of citations per year in social media platforms such as blogs and Twitter feeds; total readership per year, as measured by Altmetrics; and the Altmetric “attention score.” Because the distributions of all five measures presented heavy tails and deviated from a bell-shaped distribution, we performed a Kruskal–Wallis test on ranks followed by pairwise Dunn’s tests to confirm which groups presented significant differences with a significance level of 0.01. We provide all  $p$ -values and test statistics from these experiments in our electronic supplemental material on GitHub (<https://github.com/smangul1/good.software>).

## Supporting information

**S1 Text. An example of the “quick start”.**  
(PDF)

**S2 Text. Automatic verification of software installability.**  
(PDF)

**S3 Text. List of bioinformatics tools with the name Prism.**  
(PDF)

**S1 Table. The names of the 10 journals that were used to retrieve the URLs.** We reported the total number of papers with URLs in the abstract or body of the paper (“Number of URLs”) and the number of accessible URLs, which were not broken or time-out (“Number of accessible URLs”). URL, uniform resource locator.  
(XLSX)

**S2 Table. Installability of 98 published software tools between 2004 and 2018.**  
(XLSX)

**S3 Table. List of earliest published software tools and resources stored on <https://sourceforge.net> and <https://github.com/>.**  
(XLSX)

**S1 Fig. Protocol to check the archival stability of a published software tool or resource.**  
Numbers are provided for illustrative purposes and correspond to the link presented in the abstracts of the published papers considered in this study.

(TIFF)

**S2 Fig. Protocol to verify the installability of a published software tool.**

(TIFF)

**S3 Fig. A box plot showing the time required to install tools that required a single command compared with tools that required multiple commands (Mann–Whitney *U* test, *p*-value =  $4.7 \times 10^{-6}$ ).**

(TIFF)

## Acknowledgments

We thank John Didion (<https://twitter.com/jdidion>) for an interesting discussion over Twitter about the issue of software installability.

## References

1. Van Noorden R, Maher B, Nuzzo R. The top 100 papers. *Nature*. 2014; 514: 550–553. <https://doi.org/10.1038/514550a> PMID: 25355343
2. Wren JD. Bioinformatics programs are 31-fold over-represented among the highest impact scientific papers of the past two decades. *Bioinformatics*. 2016; 32: 2686–2691. <https://doi.org/10.1093/bioinformatics/btw284> PMID: 27153671
3. Greene AC, Giffin KA, Greene CS, Moore JH. Adapting bioinformatics curricula for big data. *Brief Bioinform*. 2016; 17: 43–50. <https://doi.org/10.1093/bib/bbv018> PMID: 25829469
4. Stephens ZD, Lee SY, Faghri F, Campbell RH, Zhai C, Efron MJ, et al. Big Data: Astronomical or Genomic? *PLoS Biol*. 2015; 13: e1002195. <https://doi.org/10.1371/journal.pbio.1002195> PMID: 26151137
5. Ahn W-Y, Busemeyer JR. Challenges and promises for translating computational tools into clinical practice. *Current Opinion in Behavioral Sciences*. 2016; 11: 1–7. <https://doi.org/10.1016/j.cobeha.2016.02.001> PMID: 27104211
6. Markowitz F. All biology is computational biology. *PLoS Biol*. 2017; 15: e2002050. <https://doi.org/10.1371/journal.pbio.2002050> PMID: 28278152
7. Marx V. The big challenges of big data. *Nature*. 2013; 498: 255–260. <https://doi.org/10.1038/498255a> PMID: 23765498
8. Stodden V, Seiler J, Ma Z. An empirical analysis of journal policy effectiveness for computational reproducibility. *Proc Natl Acad Sci U S A*. 2018; 115: 2584–2589. <https://doi.org/10.1073/pnas.1708290115> PMID: 29531050
9. Gertler P, Galiani S, Romero M. How to make replication the norm. *Nature*. 2018; 554: 417–419. <https://doi.org/10.1038/d41586-018-02108-9> PMID: 29469135
10. Beaulieu-Jones BK, Greene CS. Reproducibility of computational workflows is automated using continuous analysis. *Nat Biotechnol*. 2017; 35: 342–346. <https://doi.org/10.1038/nbt.3780> PMID: 28288103
11. List M, Ebert P, Albrecht F. Ten Simple Rules for Developing Usable Software in Computational Biology. *PLoS Comput Biol*. 2017; 13: e1005265. <https://doi.org/10.1371/journal.pcbi.1005265> PMID: 28056032
12. Baxter SM, Day SW, Fetrow JS, Reisinger SJ. Scientific Software Development Is Not an Oxymoron. *PLoS Comput Biol*. 2006; 2: e87. <https://doi.org/10.1371/journal.pcbi.0020087> PMID: 16965174
13. Carpenter AE, Kamensky L, Eliceiri KW. A call for bioimaging software usability. *Nat Methods*. 2012; 9: 666–670. <https://doi.org/10.1038/nmeth.2073> PMID: 22743771
14. Leprevost F da V, Barbosa VC, Francisco EL, Perez-Riverol Y, Carvalho PC. On best practices in the development of bioinformatics software. *Front Genet*. 2014; 5. <https://doi.org/10.3389/fgene.2014.00199> PMID: 25071829
15. Plić A, Procter JB. Ten simple rules for the open development of scientific software. *PLoS Comput Biol*. 2012; 8: e1002802. <https://doi.org/10.1371/journal.pcbi.1002802> PMID: 23236269

16. Altschul S, Demchak B, Durbin R, Gentleman R, Krzywinski M, Li H, et al. The anatomy of successful computational biology software. *Nat Biotechnol.* 2013; 31: 894–897. <https://doi.org/10.1038/nbt.2721> PMID: 24104757
17. Jiménez RC, Kuzak M, Alhamdoosh M, Barker M, Batut B, Borg M, et al. Four simple recommendations to encourage best practices in research software. *F1000Res.* 2017; 6. <https://doi.org/10.12688/f1000research.11407.1> PMID: 28751965
18. Ósz Á, Pongor LS, Szirmai D, Gyórfy B. A snapshot of 3649 Web-based services published between 1994 and 2017 shows a decrease in availability after 2 years. *Brief Bioinform.* 2017. <https://doi.org/10.1093/bib/bbx159> PMID: 29228189
19. Gewaltig M-O, Cannon R. Current practice in software development for computational neuroscience and how to improve it. *PLoS Comput Biol.* 2014; 10: e1003376. <https://doi.org/10.1371/journal.pcbi.1003376> PMID: 24465191
20. Guellec D, Van Pottelsberghe De La Potterie B. The impact of public R&D expenditure on business R&D\*. *Economics of Innovation and New Technology.* 2003; 12: 225–243.
21. Ahmed Z, Zeeshan S, Dandekar T. Developing sustainable software solutions for bioinformatics by the “Butterfly” paradigm. *F1000Res.* 2014; 3: 71. <https://doi.org/10.12688/f1000research.3681.2> PMID: 25383181
22. Kanitz A, Gypas F, Gruber AJ, Gruber AR, Martin G, Zavolan M. Comparative assessment of methods for the computational inference of transcript isoform abundance from RNA-seq data. *Genome Biol.* 2015; 16: 150. <https://doi.org/10.1186/s13059-015-0702-5> PMID: 26201343
23. Support Model Organism Databases [Internet]. [cited 11 Aug 2018]. <http://www.genetics-gsa.org/MODsupport>.
24. Database under maintenance. *Nat Methods.* 2016; 13: 699–699.
25. Chen S-S. Digital Preservation: Organizational Commitment, Archival Stability, and Technological Continuity. *Journal of Organizational Computing and Electronic Commerce.* 2007; 17: 205–215.
26. Carnevale RJ, Aronsky D. The life and death of URLs in five biomedical informatics journals. *Int J Med Inform.* 2007; 76: 269–273. <https://doi.org/10.1016/j.ijmedinf.2005.12.001> PMID: 16458066
27. Markwell J, Brooks DW. “Link rot” limits the usefulness of web-based educational materials in biochemistry and molecular biology. *Biochemistry and Molecular Biology Education.* 2003; 31(1): 69–72. <https://doi.org/10.1002/bmb.2003.494031010165>
28. Dellavalle RP, Hester EJ, Heilig LF, Drake AL, Kuntzman JW, Graber M, et al. Information science. Going, going, gone: lost Internet references. *Science.* 2003; 302: 787–788. <https://doi.org/10.1126/science.1088234> PMID: 14593153
29. Ducut E, Liu F, Fontelo P. An update on Uniform Resource Locator (URL) decay in MEDLINE abstracts and measures for its mitigation. *BMC Med Inform Decis Mak.* 2008; 8. <https://doi.org/10.1186/1472-6947-8-23> PMID: 18547428
30. Wren JD, Georgescu C, Giles CB, Hennessey J. Use it or lose it: citations predict the continued online availability of published bioinformatics resources. *Nucleic Acids Res.* 2017; 45: 3627–3633. <https://doi.org/10.1093/nar/gkx182> PMID: 28334982
31. Wren JD. URL decay in MEDLINE—a 4-year follow-up study. *Bioinformatics.* 2008; 24: 1381–1385. <https://doi.org/10.1093/bioinformatics/btn127> PMID: 18413326
32. Piwowar H. Altmetrics: Value all research products. *Nature.* 2013; 493: 159. <https://doi.org/10.1038/493159a> PMID: 23302843
33. Grüning B, The Bioconda Team, Dale R, Sjödin A, Chapman BA, Rowe J, et al. Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nat Methods.* 2018; 15: 475–476. <https://doi.org/10.1038/s41592-018-0046-7> PMID: 29967506
34. Perez-Riverol Y, Gatto L, Wang R, Sachsenberg T, Uszkoreit J, Leprevost F da V, et al. Ten Simple Rules for Taking Advantage of Git and GitHub. *PLoS Comput Biol.* 2016; 12: e1004947. <https://doi.org/10.1371/journal.pcbi.1004947> PMID: 27415786
35. Perkel, J. When it comes to reproducible science, Git is code for success. 2018 Jun 11 [cited 11 Aug 2018]. In: *Nature Index* [Internet]. <https://www.natureindex.com/news-blog/when-it-comes-to-reproducible-science-git-is-code-for-success>.
36. da Veiga Leprevost F, Grüning BA, Alves Aflitos S, Röst HL, Uszkoreit J, Barsnes H, et al. BioContainers: an open-source and community-driven framework for software standardization. *Bioinformatics.* 2017; 33: 2580–2582. <https://doi.org/10.1093/bioinformatics/btx192> PMID: 28379341
37. Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, et al. Bioconductor: open software development for computational biology and bioinformatics. *Genome Biol.* 2004; 5: R80. <https://doi.org/10.1186/gb-2004-5-10-r80> PMID: 15461798

38. Karimzadeh M, Hoffman MM. Top considerations for creating bioinformatics software documentation. *Brief Bioinform.* 2018; 19: 693–699. <https://doi.org/10.1093/bib/bbw134> PMID: 28088754
39. Howell M. Homebrew. [software]. [cited 17 Aug 2018]. <https://brew.sh/>.
40. Guerreiro M. Forking software used in eLife papers to GitHub. 2017 Apr 14. In: eLife [Internet]. eLife Sciences Publications Limited; 2017. <https://elifesciences.org/inside-elife/dbcb6949/forking-software-used-in-elife-papers-to-github>.
41. Mosqueiro T, Cook C, Huerta R, Gadau J, Smith B, Pinter-Wollman N. Task allocation and site fidelity jointly influence foraging regulation in honeybee colonies. *R Soc Open Sci.* 2017; 4: 170344. <https://doi.org/10.1098/rsos.170344> PMID: 28878985
42. Hunt M, Newbold C, Berriman M, Otto TD. A comprehensive evaluation of assembly scaffolding tools. *Genome Biol.* 2014; 15: R42. <https://doi.org/10.1186/gb-2014-15-3-r42> PMID: 24581555
43. Fonseca NA, Rung J, Brazma A, Marioni JC. Tools for mapping high-throughput sequencing data. *Bioinformatics.* 2012; 28: 3169–3177. <https://doi.org/10.1093/bioinformatics/bts605> PMID: 23060614
44. Pabinger S, Dander A, Fischer M, Snajder R, Sperk M, Efremova M, et al. A survey of tools for variant analysis of next-generation genome sequencing data. *Brief Bioinform.* 2014; 15: 256–278. <https://doi.org/10.1093/bib/bbs086> PMID: 23341494