

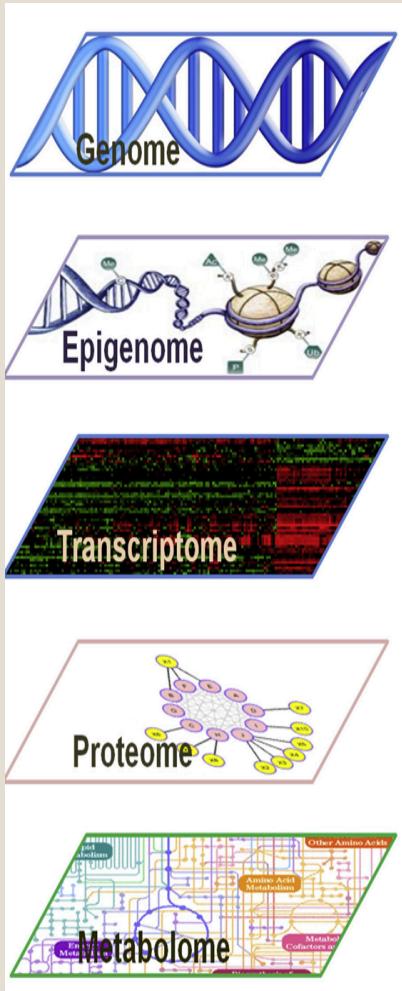
# ANALYSIS OF THE USABILITY AND ARCHIVAL STABILITY OF OMICS COMPUTATIONAL TOOLS AND RESOURCES

Serghei Mangul  
University of California, Los Angeles



@serghei\_mangul

# Flood of omics software tools



Bioinformatics tools

# Many bioinformatics tools are hard to install

⤤ You Retweeted



**Emily Madden** @\_emadden · Oct 27

The tweet I come across as I'm taking a break from trying to properly install and run a bioinformatics tool. [#perfecttiming](#)

**Ran Blekhman** @blekhman

Say you pick 100 random bioinformatics software tools -- how many will you actually be able to access, install, and run?

Our new paper: [biorxiv.org/content/early/](https://biorxiv.org/content/early/)... [twitter.com/serghei\\_mangul...](https://twitter.com/serghei_mangul)

Show this thread

# Is user-friendly scientific software an oxymoron?

OPEN ACCESS Freely available online

PERSPECTIVES

PLOS COMPUTATIONAL BIOLOGY

## Scientific Software Development Is Not an Oxymoron

Susan M. Baxter\*, Steven W. Day, Jacquelyn S. Fetrow, Stephanie J. Reisinger

*"Many scientists and engineers spend much of their lives writing, debugging, and maintaining software, but only a handful have ever been taught how to do this effectively: after a couple of introductory courses, they are left to rediscover (or reinvent) the rest of programming on their own. The result? Most spend far too much time wrestling with software, instead of doing research, but have no idea how reliable or efficient their programs are." —Greg Wilson [1]*

software development life cycle onto computational biology projects to build a solid foundation for success.

Two of us are card-carrying software engineers; two of us are formally trained scientists. We are all battle-scared veterans of large scientific software development projects, while working in business, nonprofit, government, and academic settings. Many of those projects were successful; some were not. We think that the best practices learned and

PLOS COMPUTATIONAL BIOLOGY

BROWSE PUBLISH

OPEN ACCESS

EDITORIAL

## Ten Simple Rules for Developing Usable Software in Computational Biology

Markus List, Peter Ebert, Felipe Albrecht

OPEN ACCESS Freely available online

EDITORIAL

PERSPECTIVE

PLOS COMPUTATIONAL BIOLOGY

## Ten Simple Rules for the Open Development of Scientific Software

Andreas Prlic\*, James B. Procter<sup>2</sup>

1 San Diego Supercomputer Center, University of California San Diego, La Jolla, California, United States of America, 2 School of Life Sciences Research, College of Life Sciences, University of Dundee, Dundee, Scotland, United Kingdom

OPEN ACCESS Freely available online

PERSPECTIVE

PLOS COMPUTATIONAL BIOLOGY

## Current Practice in Software Development for Computational Neuroscience and How to Improve It

Marc-Oliver Gewaltig<sup>1\*</sup>, Robert Cannon<sup>2</sup>

1 Blue Brain Project, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland, 2 Textensor Ltd., Edinburgh, United Kingdom

nature biotechnology

Feature | Published: 08 October 2013

## The anatomy of successful computational biology software

Stephen Altschul, Barry Demchak, Richard Durbin, Robert Gentleman, Martin Krzywinski, Heng Li, Anton Nekrutenko, James Robinson, Wayne Rasband, James Taylor & Cole Trapnell

Nature Biotechnology 31, 894–897 (2013) | Download Citation ↴

frontiers in Genetics

Front Genet. 2014; 5: 199.

Published online 2014 Jul 2. doi: [10.3389/fgene.2014.00199](https://doi.org/10.3389/fgene.2014.00199)

PMCID: PMC4078907

PMID: 25071829

## On best practices in the development of bioinformatics software

Felipe da Veiga Leprevost,<sup>1,2,\*</sup> Valmir C. Barbosa,<sup>3</sup> Eduardo L. Francisco,<sup>2</sup> Yasset Perez-Rivero,<sup>4</sup> and Paulo C. Carvalho<sup>1</sup>

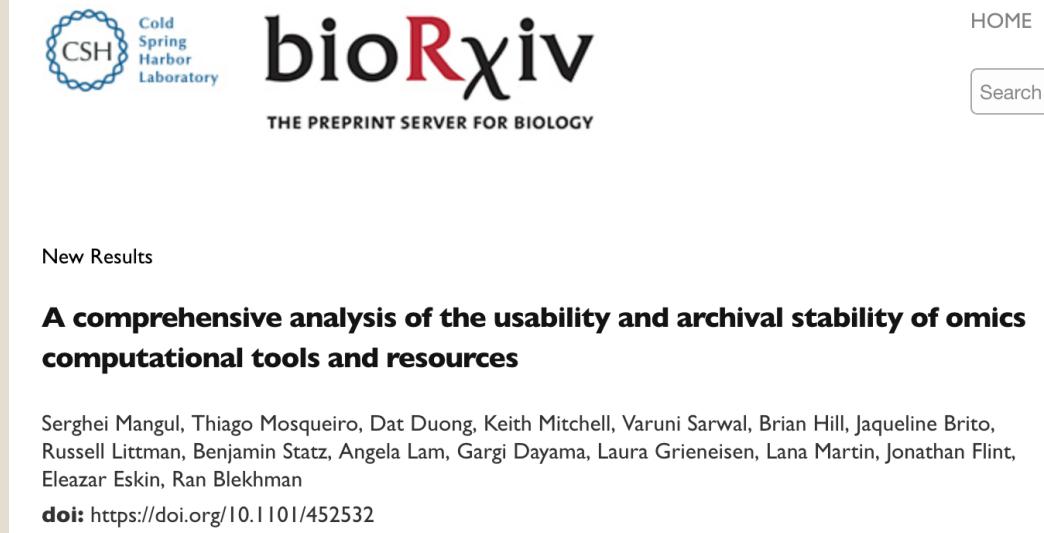
## Challenges to effective software development and distribution in academia

- Software written by researchers tends to be written with the idea that all users will know as much about the code as its original authors
- Incentives in academia heavily favor the publication of new software, not the maintenance of existing tools
- There is lack of protocols to check insatiability of published software tools in academia

# Software development in academia versus industry

- Industry receive considerably more resources for developing user-friendly tools
- Companies efficiently distribute industry-produced software using dedicated company units or contractors
- There is little reward for continuous, long-term development and maintenance of tools in academia

# Poorly implemented tools will ultimately hinder progress in big data-driven fields



The image shows the bioRxiv preprint server interface. At the top left is the CSHL logo with the text "Cold Spring Harbor Laboratory". Next to it is the bioRxiv logo with the tagline "THE PREPRINT SERVER FOR BIOLOGY". To the right are navigation links for "HOME" and "Search". Below the header, there is a section titled "New Results" featuring a specific preprint. The title of the preprint is "A comprehensive analysis of the usability and archival stability of omics computational tools and resources". The authors listed are Serghei Mangul, Thiago Mosqueiro, Dat Duong, Keith Mitchell, Varuni Sarwal, Brian Hill, Jaqueline Brito, Russell Littman, Benjamin Statz, Angela Lam, Gargi Dayama, Laura Grieneisen, Lana Martin, Jonathan Flint, Eleazar Eskin, and Ran Blekhman. The DOI provided is <https://doi.org/10.1101/452532>.



The image shows a page from the journal Genome Biology. At the top right is the journal title "Genome Biology". Below it is an "EDITORIAL" article by Mangul et al. The article is titled "Improving the usability and archival stability of bioinformatics software". It is marked as "Open Access". The authors listed are Serghei Mangul<sup>1,2\*</sup>, Lana S. Martin<sup>2</sup>, Eleazar Eskin<sup>1,3</sup> and Ran Blekhman<sup>4,5</sup>. There is also a "Check for updates" button.

# Software crisis is a reproducibility crisis

- Limited software usability and archival stability of computational tools leads to (**computational reproducibility crisis**)
- **Computational reproducibility** is the ability to replicate published findings by running the same computational tool on the data generated by the published study

# Data: archival stability of omics computational tools and resources

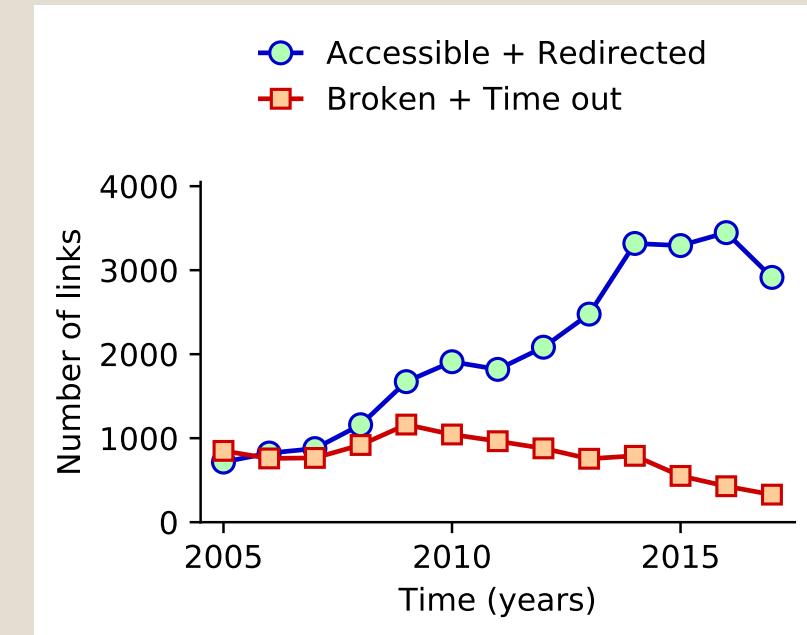
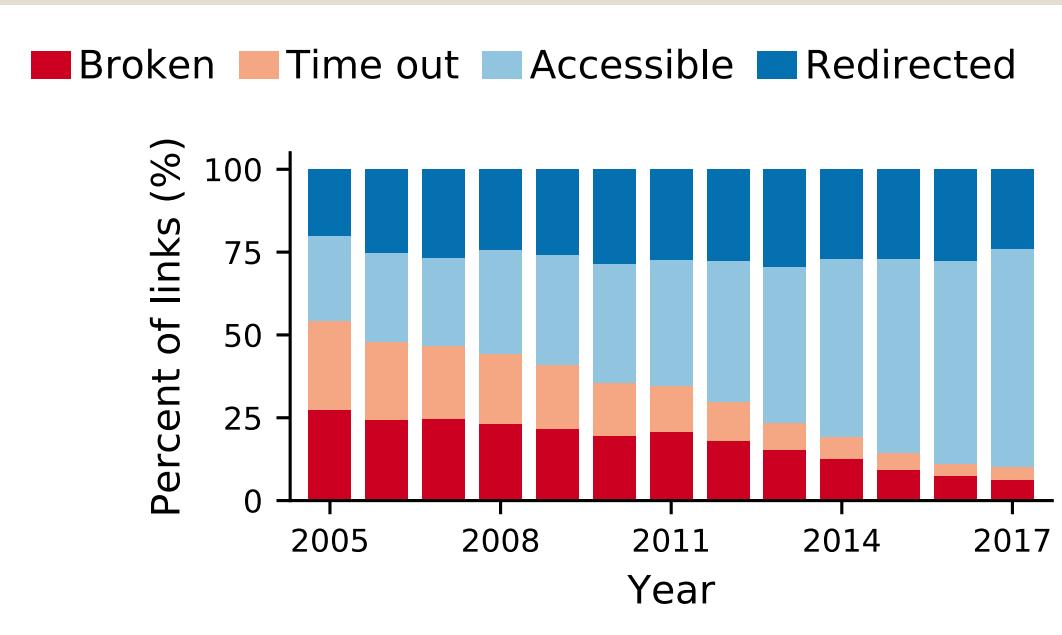
- 51,236 biomedical papers from PubMed
- 99 randomly selected bioinformatics tools

Journal name	Number of URLs
Nature Biotechnol	180
Genome Medicine	352
Nature Methods	403
Genome Biology	904
BMC Systems Biology	912
Bioinformatics	3131
PLoS Comp. Biology	3226
BMC Bioinformatics	6840
BMC Genomics	7651
Nucleic Acids Research	13103

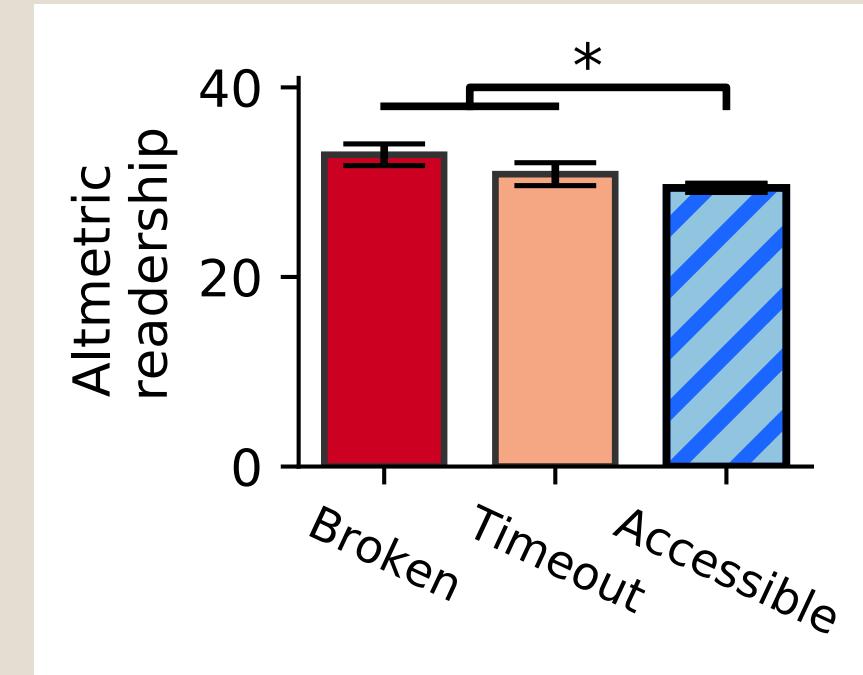
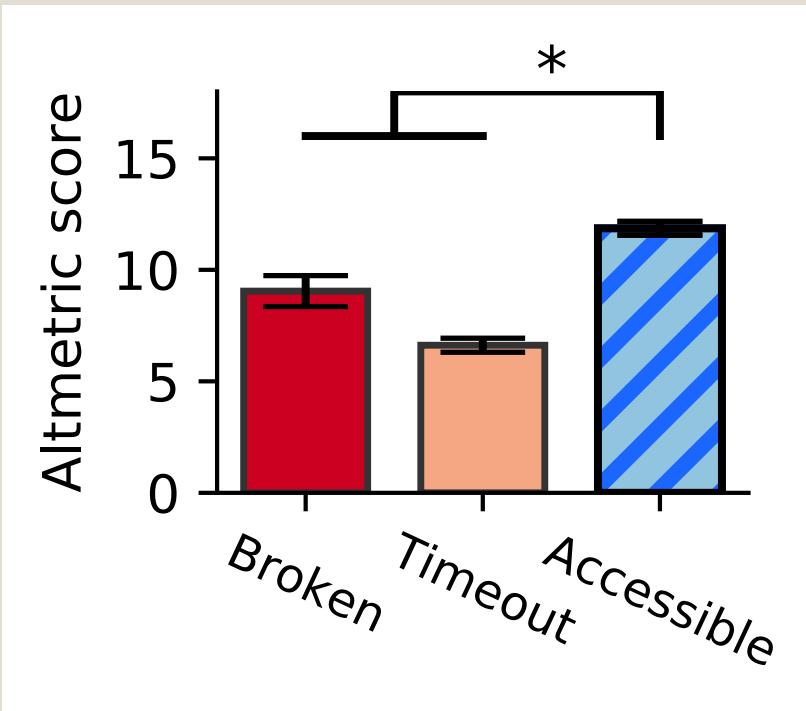
## Analysis of archival stability

- We downloaded **51,236** open access papers via PubMed from 10 systems and computational biology journals (Raw data in XML format)
- We developed an approach to extract **36,702** software links from the downloaded papers and verify the archival stability of links
- Timeout links were manually verified

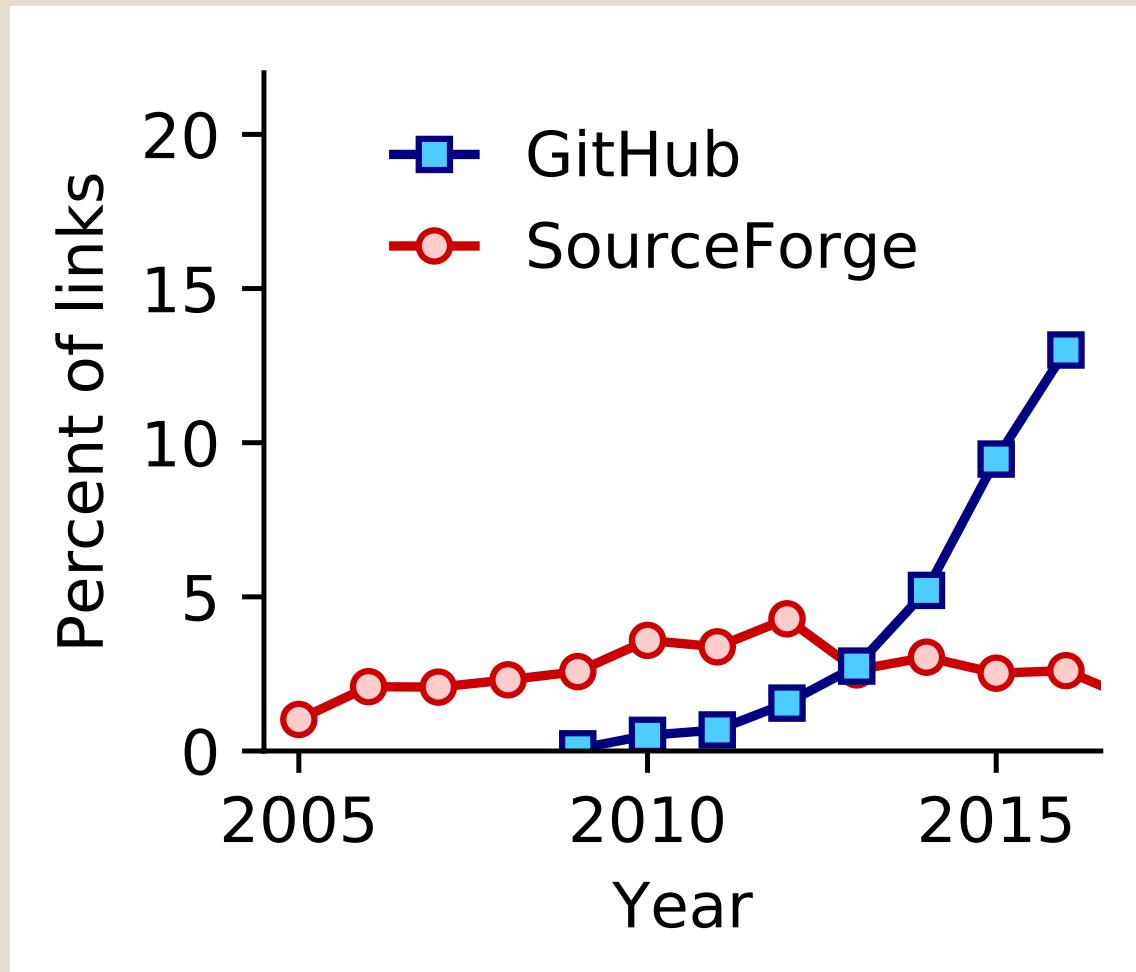
# Archival stability of abstract URLs



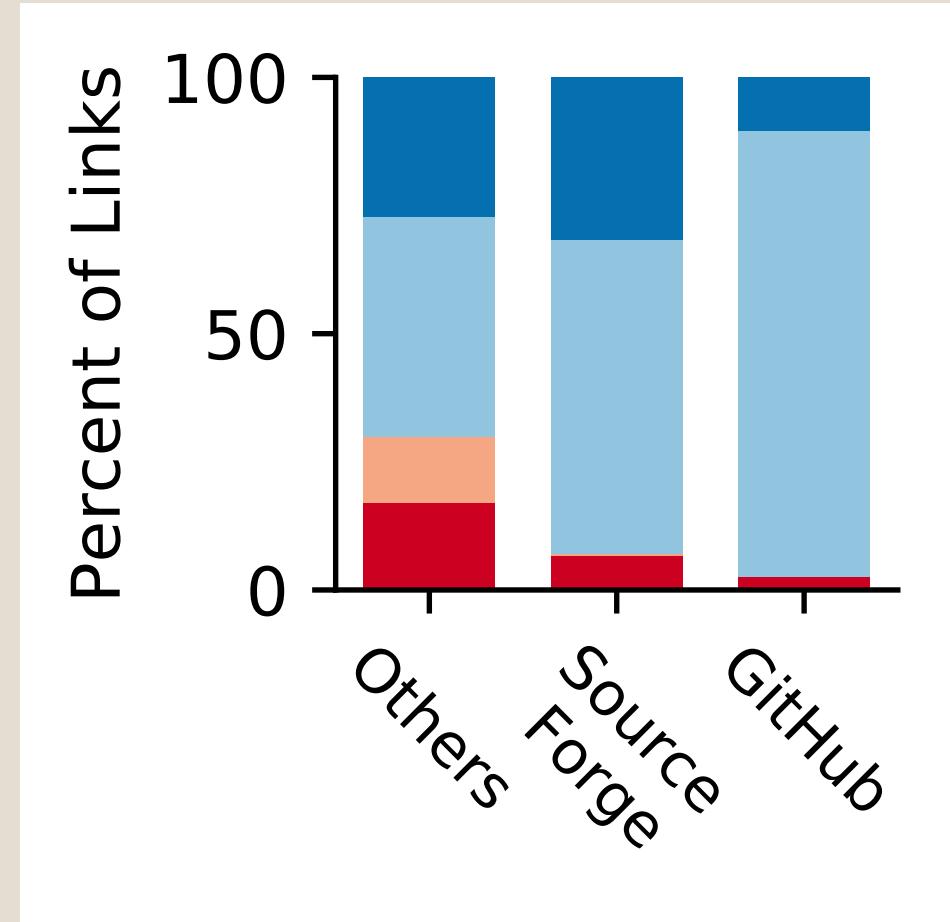
# Effect of social media



# GitHub is the most popular platform to host scientific code



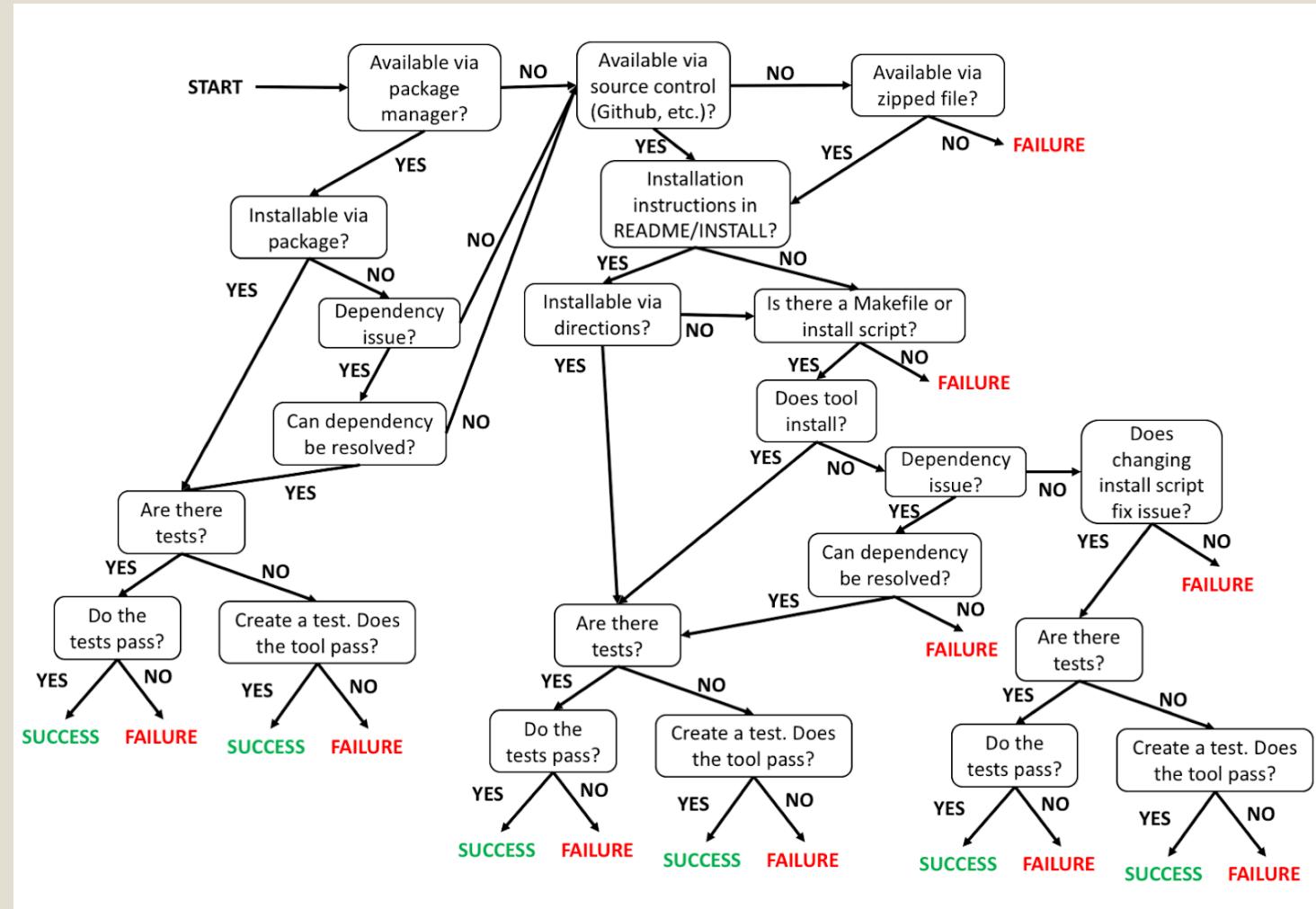
# Tools hosted on GitHub have <3% of unreachable links



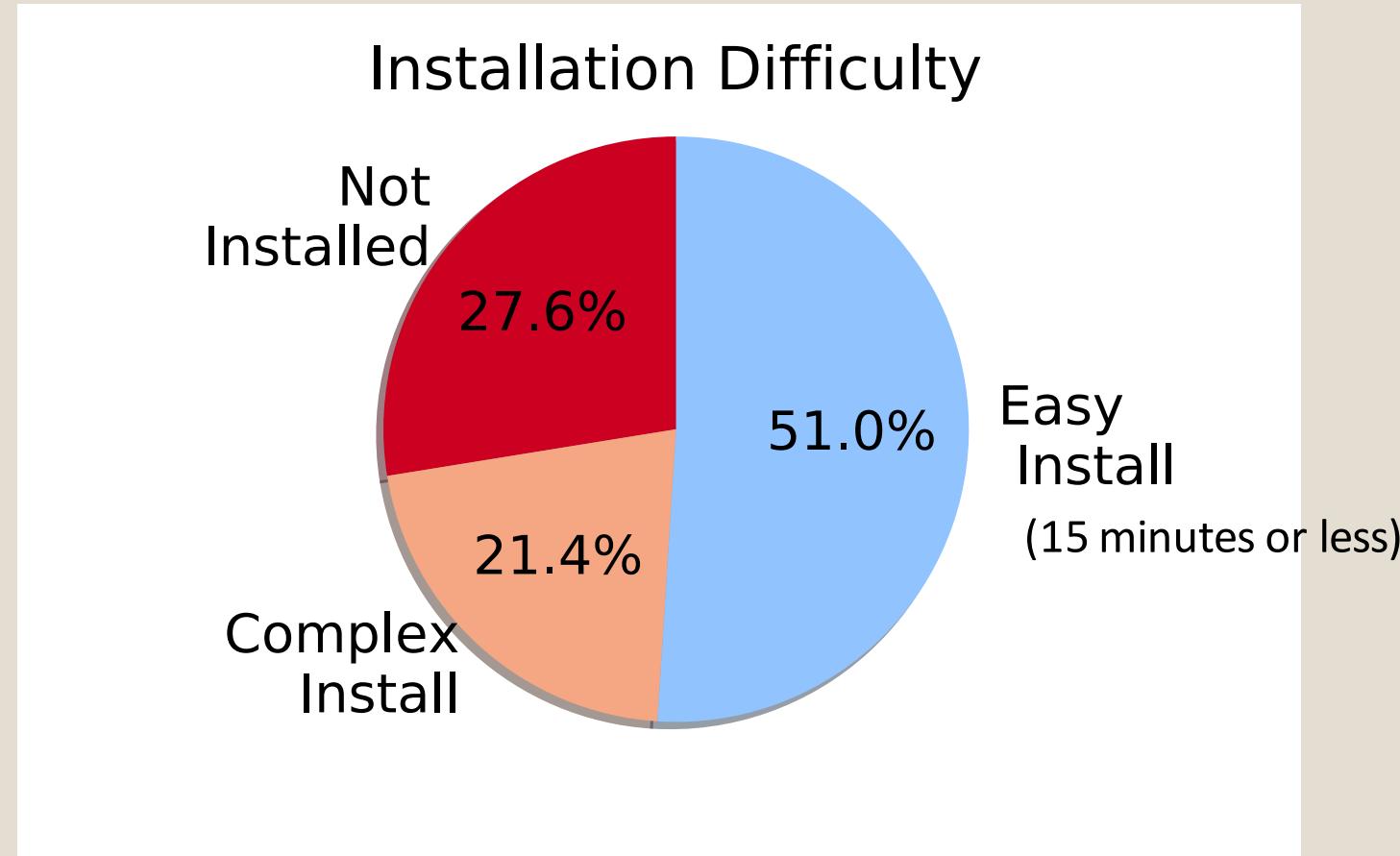
# Analysis of the software usability

- We have randomly selected **99** tools across various domains of computational biology.
- Assign one undergraduate student to install 10 tools or more
- Record total time, number of commands to install the tools
- Total human time to install 99 tools -- **72 hours**

# Protocol to check the insatiability of the tool

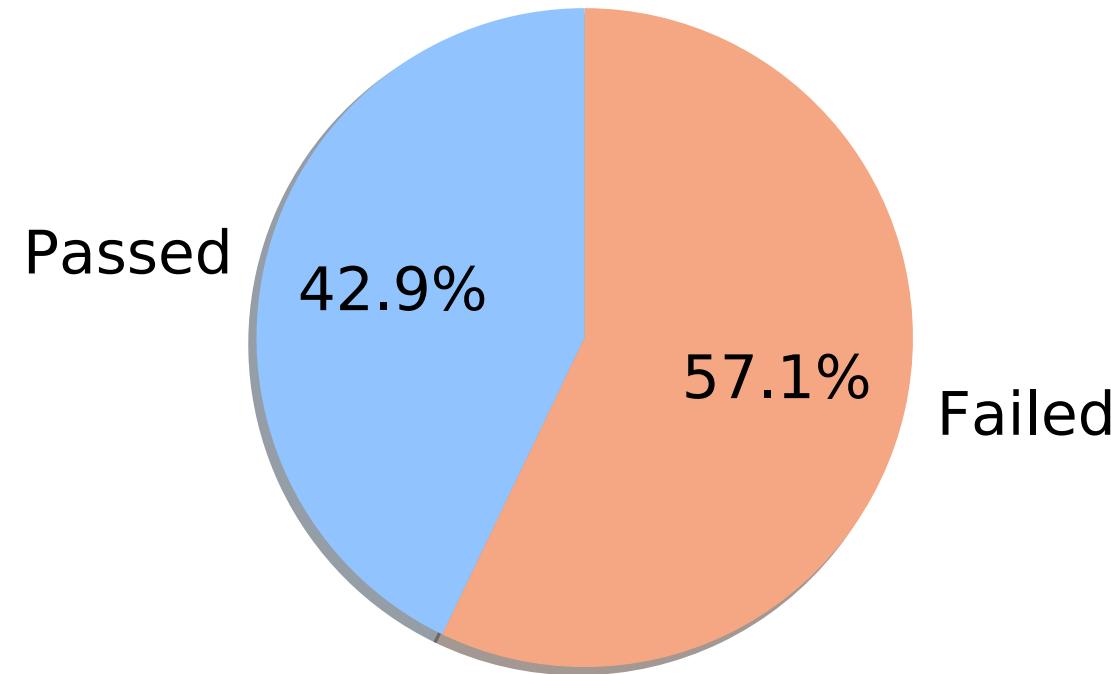


# Many tools are hard or impossible to install



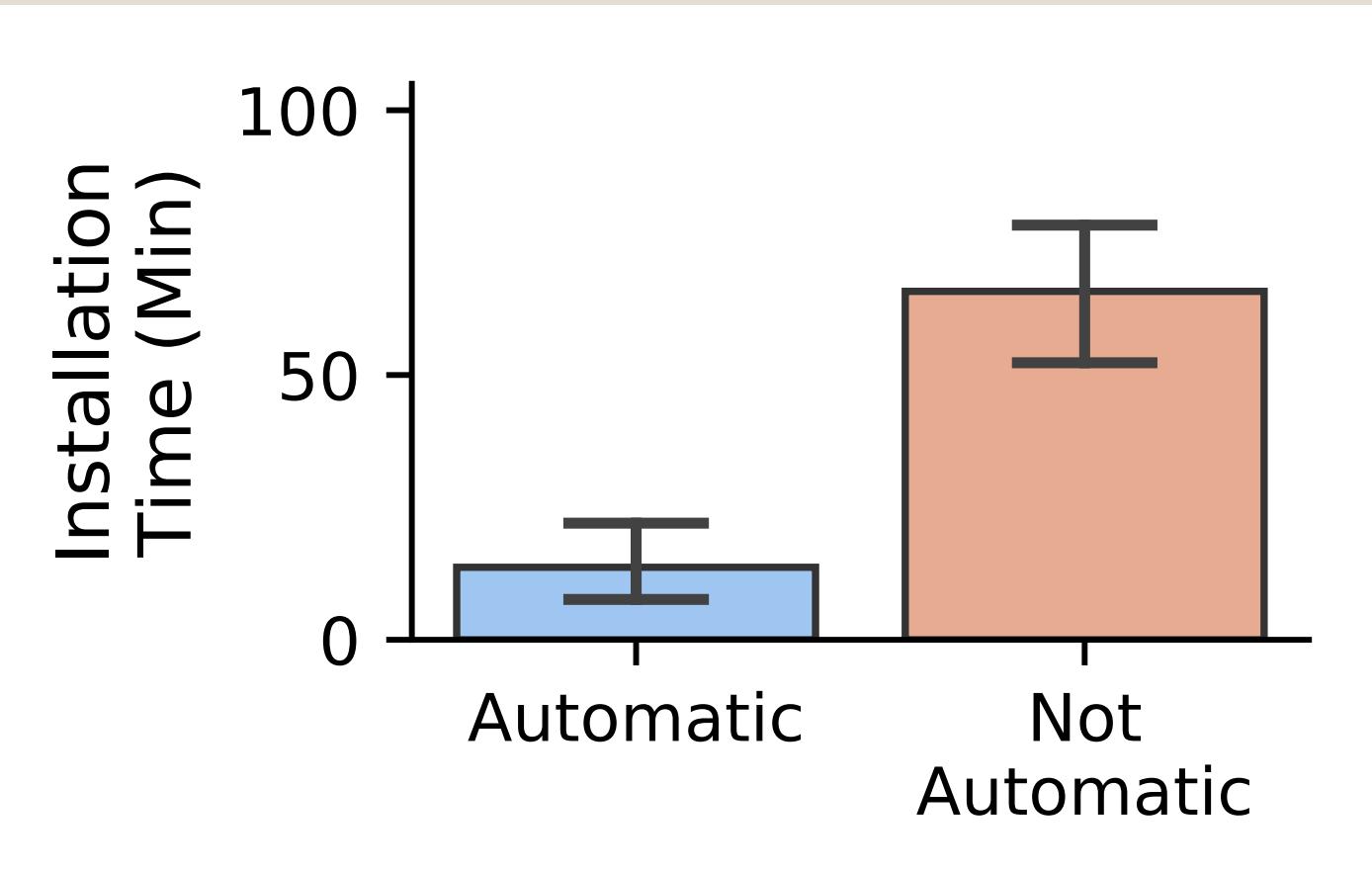
# Automatic installation test

Automatic Installation Test

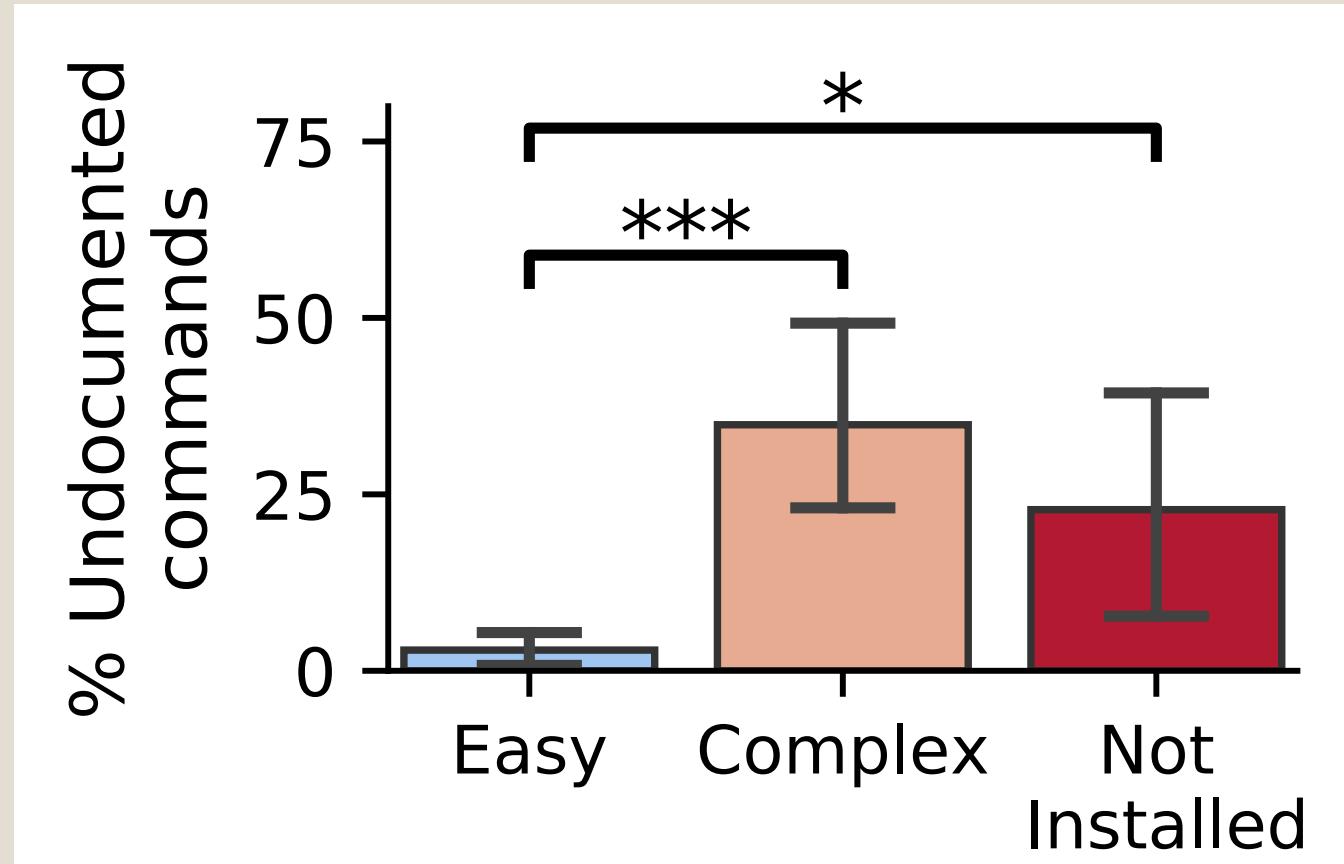


Passed: no manual intervention

# Manual interventions are time consuming

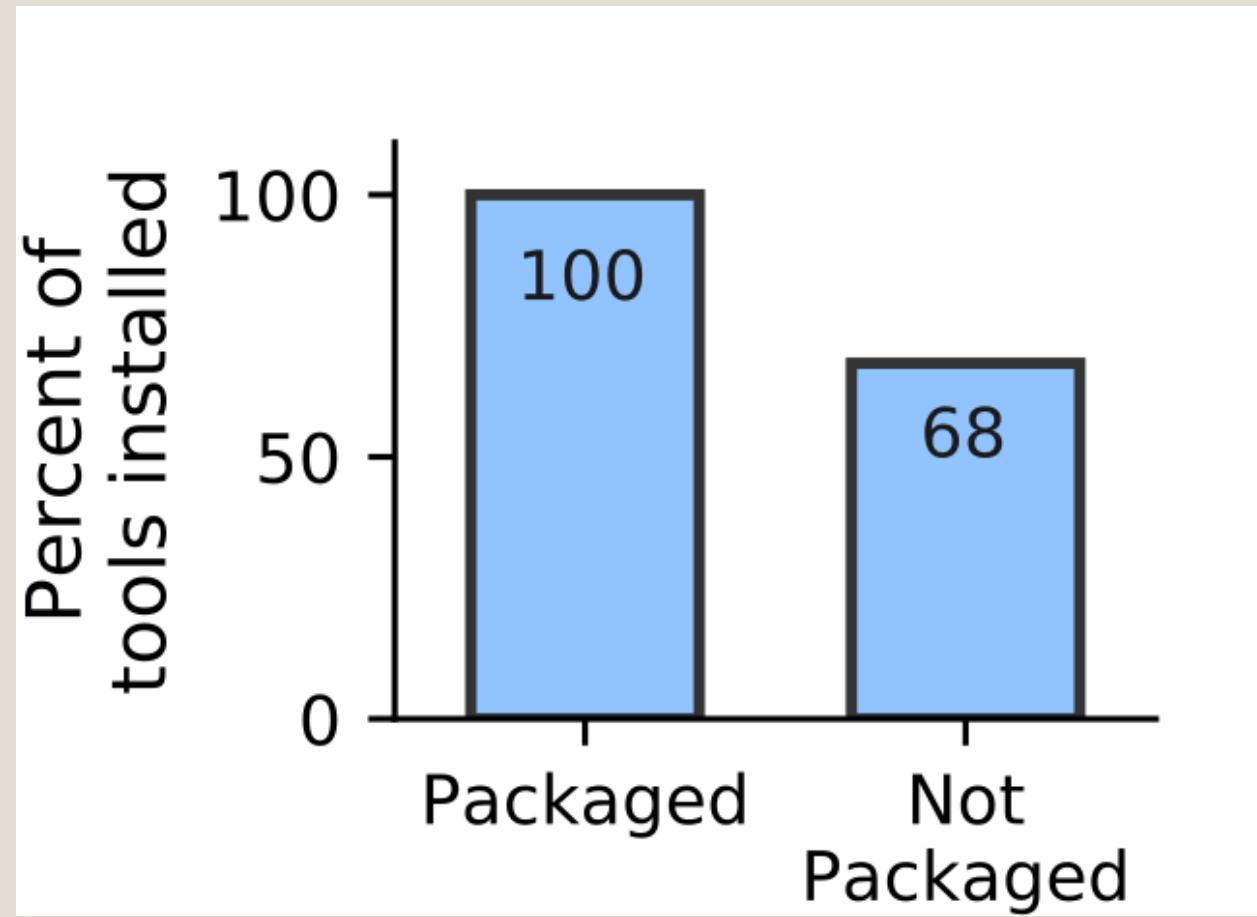


# Easy-to-install\* tools have fewer undocumented commands

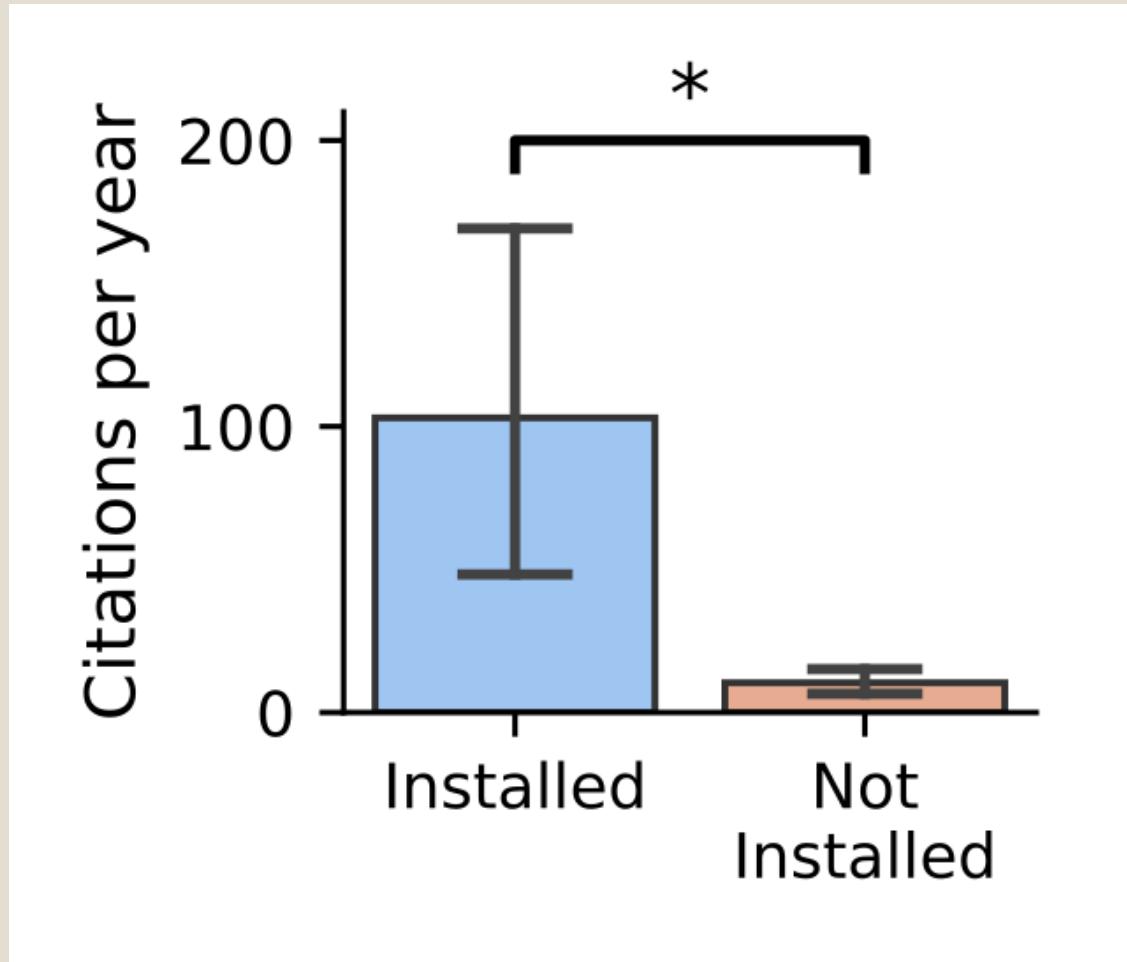


\*<15min

# Bioconda tools were always installable

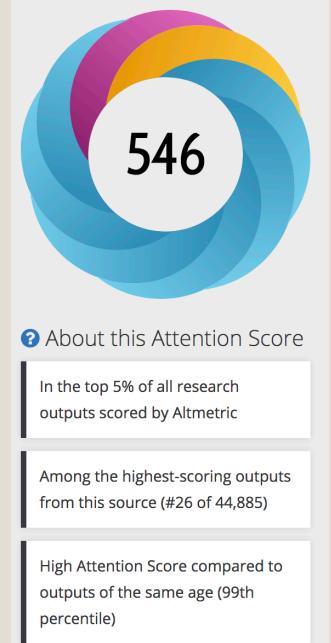


# Installable tools have more citations



# Conclusions

- Host software on archaically stable services e.g. **GitHub**
- Provide easy-to-use installation interface and get more citations
- Tools wrapped in **package managers (e.g. Bioconda)** are 100% installable!



Abstract	Pdf
13,248	3,434

Below the table, there is a section with social media metrics:

- Blogged by 1
- Tweeted by 1087
- Mentioned in 1 Google+ posts
- 40 readers on Mendeley

# Future work

 **Serghei Mangul** @serghei\_mangul · Oct 30

1/2 CALL FOR PARTICIPANTS: Given the interest in our study about bioinformatics tools usability, we want to work on the follow up a paper about the package managers for bioinformatics tools. Please reply in this tweet if you are interested to participate! [@thmosqueiro](#) [@blekhman](#)



5 5 10 ||

## Acknowledgment

- Thiago Mosqueiro, Dat Duong, Keith Mitchell, Varuni Sarwal, Brian Hill, Jaqueline Brito, Russell Littman, Benjamin Statz, Angela Lam, Gargi Dayama, Laura Grieneisen, Lana Martin, Jonathan Flint, Eleazar Eskin, Ran Blekhman

We have prepared Jupyter Notebooks that utilize the raw data described above to reproduce the results and figures presented in our paper

→ <https://github.com/smangul1/good.software>

[smangul@ucla.edu](mailto:smangul@ucla.edu)

<http://www.sergheimangul.com/>

THANK YOU.