# ~~Science~~ - Broken Science

Paywalled, proprietary tools & software, non shared data, non-reproducible, anonymous peer-review, publish or perish (alone)!

**PAYWALLED (32$)**

Nature 171 (1953)
Molecular Structure of Nucleic Acids: A Structure for Deoxyribose Nucleic Acid
J.D. WATSON & F. H. C. CRICK.

I have heard from graduate students opting out of academia, assistant professors afraid to come up for tenure, mid-career people wondering how to protect their labs, and senior faculty retiring early, all because of methodological terrorism.

APS Observer (2016)

**METHODOLOGICAL TERRORISM**

A second concern held by some is that a new class of research person will emerge — people who had nothing to do with the design and execution of the study but use another group's data for their own ends, possibly stealing from the research productivity planned by the data gatherers, or even use the data to try to disprove what the original investigators had posited.There is concern among some front-line researchers that the system will be taken over by what some researchers have characterized as research parasites

The New England Journal of Medicine (2016)

**RESEARCH PARASITES**

# ~~Open Science~~ - Science

Open Access, Open Source, Open Data, Open Methodology,  Open Education, Open Peer-review, Much more fun & efficient!

Original Article

Companion Notebook

# Once upon a time, there was a post-doc…

_Interaction between cognitive and motor cortico-basal ganglia loops during decision making: a computational study._ M. Guthrie, A. Leblois, A. Garenne, and T. Boraud, Journal of Neurophysiology, 109, 2013

Nice paper, good results, but…

- No public repository, no version control
- Sources were mixing actual computation and GUI code
- Model was split into a hundred files, main file 6,000 lines long
- Several configuration files, no data saved
- Model description included ambiguous information

Model was hardly reproducible.

_You can download our code from the URL supplied. Good luck downloading the only postdoc who can get it to run, though…_

Ian Holmes



HOW TO WRITE GOOD CODE:

# Ma'am, the dog ate my program

_We describe a study into the extent to which Computer Systems researchers share their code and data and the extent to which such code builds. Starting with 601 papers from ACM conferences and journals, we examine 402 papers whose results were backed by code. For 32.3% of these papers we were able to obtain the code and build it within 30 minutes; for 48.3% of the papers we managed to build the code, but it may have required extra effort; for 54.0% of the papers either we managed to build the code or the authors stated the code would build with reasonable effort._

```
From: Christian Collberg <ccollberg@gmail.com>
To:   first-or-corresponding-author
Cc:   remaining-authors
Subject: Your  conference-name  paper

Dear Dr.  first-or-corresponding-author ,

I've been looking at your   conference-name  paper
   paper-title
and would like to try out the implementation. However,
I haven't been able to find it online. Would you please
let me know how I can obtain the source code so that I
can try to build and run it?

Thank you very much for your help!

Christian Collberg
ccollberg@gmail.com
```

# Ma'am, the dog ate my program

---

Reasons why code cannot be shared:

→ Versioning Problems
→ Code Will be Available Soon
→ No Intention to Release
→ Programmer Left
→ Bad Backup Practices
→ Commercial Code
→ Proprietary Academic Code
→ Industrial Lab Issues
→ Unavailable Subsystems
→ Multiple Reasons
→ Intellectual Property
→ Research vs. Sharing
→ Security and Privacy
→ Design Issues
→ Too Busy to Help

```
⟨STUDENT⟩ was a graduate student in our program but he left a
while back so I am responding instead.  For the paper we used a
prototype that included many moving pieces that only ⟨STUDENT⟩
knew how to operate and we did not have the time to integrate them in
a ready-to-share implementation before he left.  Still, I hope you can
build on the ideas/technique of the paper. Regards,
```

```
Since this work has been done at ⟨COMPANY⟩ we don't open-source
code unless there is a compelling business reason to do so.  So
unfortunately I don't think we'll be able to share it with you.
```

```
Thank you for your interest in our work.  Unfortunately the current
system is not mature enough at the moment, so it's not yet publicly
available. We are actively working on a number of extensions and
things are somewhat volatile. However, once things stabilize we plan
to release it to outside users. At that point, we would be happy to
send you a copy.
```

```
Thanks for your interest in the implementation of our paper. The good
news is that I was able to find some code. I am just hoping that it is
a stable working version of the code, and matches the implementation
we finally used for the paper. Unfortunately, I have lost some data
when my laptop was stolen last year. The bad news is that the code is
not commented and/or clean. So, I cannot really guarantee that you
will enjoy playing with it.
```

```
The code used to implement the ⟨CONFERENCE⟩ paper is
complete, but hardly usable by anyone other than the authors.  This is
due in large part due to our decision to use Template Haskell for the
input language. The error messages which are produced by the compiler
are useless to anyone not fluent in both Haskell, BSV, and the
compiler architecture.
```

# A brand new implementation

Remember? *Interaction between cognitive and motor cortico-basal ganglia loops during decision making: a computational study.* M. Guthrie, A. Leblois, A. Garenne, and T. Boraud, Journal of Neurophysiology, 109, 2013.→ 100 files, 6,000 lines of Delphi
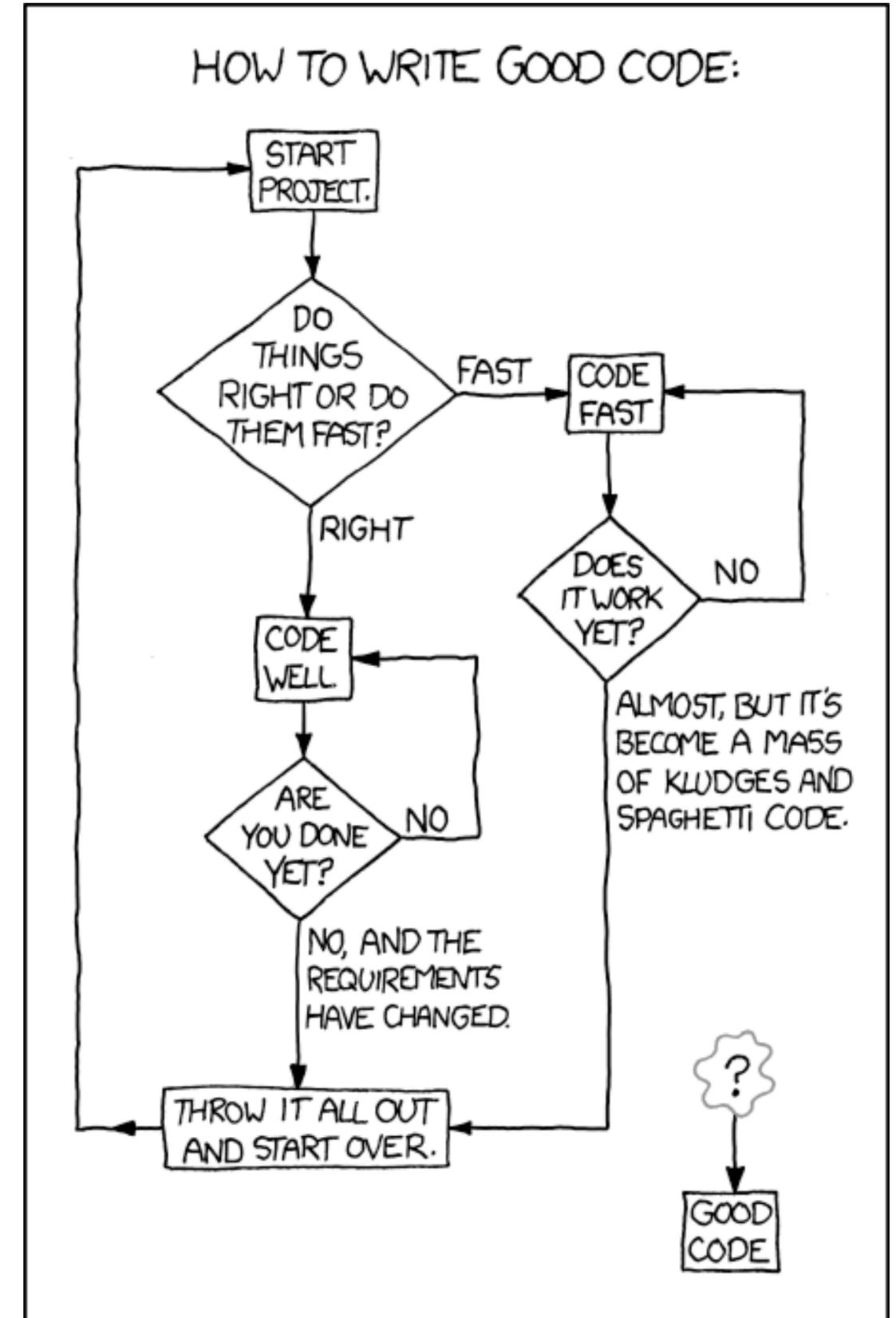
I asked my PhD student (M. Topalidou) to write a brand new implementation. Together, it took us three months of hard work to replicate the model using

- Python language and numerical libraries
- DANA library for intuitive description
- IPython notebook for interactive sessions

Source is now a single file of 200 readable notebook available on GitHub. Without this replication effort, original model would have been useless for our research.

Because of strong incentives for innovation and weak incentives for confirmation, direct replication is rarely practiced or published… Innovative findings produce rewards of publication, employment, and tenure; replicated findings produce a shrug.

Brian Nosek, The Reproducibility Project, 2012

# Reproducible computational neuroscience

Any model in Science is doomed to be proved wrong or incomplete and replaced by a more accurate one. In the meantime, for such replacement to happen, we have first to make sure that models are actually reproducible such that they can be tested, evaluated, criticized and ultimately modified, replaced or even rejected.

This is where the shoe pinches.

If we cannot reproduce a model in the first place, we're doomed to re-invent the wheel again and again, preventing us from building an incremental computational knowledge.

My field of research is quite different from computational neuroscience, but I recognize the problem described in this paper very well. The core issue has in my opinion been identified in the comment by Jan Moren: there is no obvious way to publish complex scientific models other than as part of simulation software.
Konrad Hinsen, 2015

404 code not found

# WHAT DO WE DO

NEXT ?

# 'cause we are $cientific publi$her$

— Elsevier, can I publish my replication in your journal?
— Nope!

— Hi Springer, interested in replication?
— Failure or success?
— Success!
— Nope!

— Hello Mr Wiley, did you hear about reproducible Science?
— tut…. tut…. tut…

— Dear beloved Frontiers, can you review this?
— Ha ha ha…. No.

— Well, well, well…

# The ReScience journal

———

ReScience is an open peer-reviewed journal that target any computational research and encourage the explicit replication of already published research promoting new and open-source implementations.

ReScience lives on github where each new implementation is made available together with explanations (article).

Each published article is archived on Zenodo and code is saved by Software Heritage

**ReScience in numbers:**

4 editors-in-chief
12 associate editors
110 registered reviewers
72 published articles
100% replication rate (strong bias)



We redo Science !

# ReScience

Reproducible science is good. Replicated science is better.

# The R quintuplet (R⁵)

Rerunnable
  Can you re-run your program ?
  One day, one week, one month, one year (just kidding) apart ?

Repeatable
  Can you re-run your program and get same results ?
  Did you save everything, including random seed ?

Reproducible
  Can someone re-run your program and get same results ?
  Did you save the software stack ?

Replicable
  Can someone reimplement your model and get same results ?
  Did you describe everything ?

Reusable
  Can someone reuse your program using different data ?
  Is your software data-dependent ?

Re-run, Repeat, Reproduce, Reuse, Replicate: Transforming
Code into Scientific Contributions - Benureau & Rougier, 2018

rerun
Robust

variations on
experiment and set up

repeat
Defend

same experiment,
same set up, same lab

replicate
Certify

same experiment,
same set up, independent lab

reproduce
Compare

variations on experiment,
on set up, independent labs

reuse
Transfer

# Replications in the wild

**What is a replication?**

Bob reads Alice's paper, takes note of all model properties and then implements the model himself using a method of his choice.

Bob confirms Alice's result by obtaining qualitatively the same results.

Alice's model has been replicated.

**Who wants to write replication?**

During the course of a PhD, it is often the case that a student will try to replicate results from the literature, possibly interacting with the original authors.

Such replication generally lives inside the hard-drive of the computer's student while it would be actually useful for the whole scientific community.
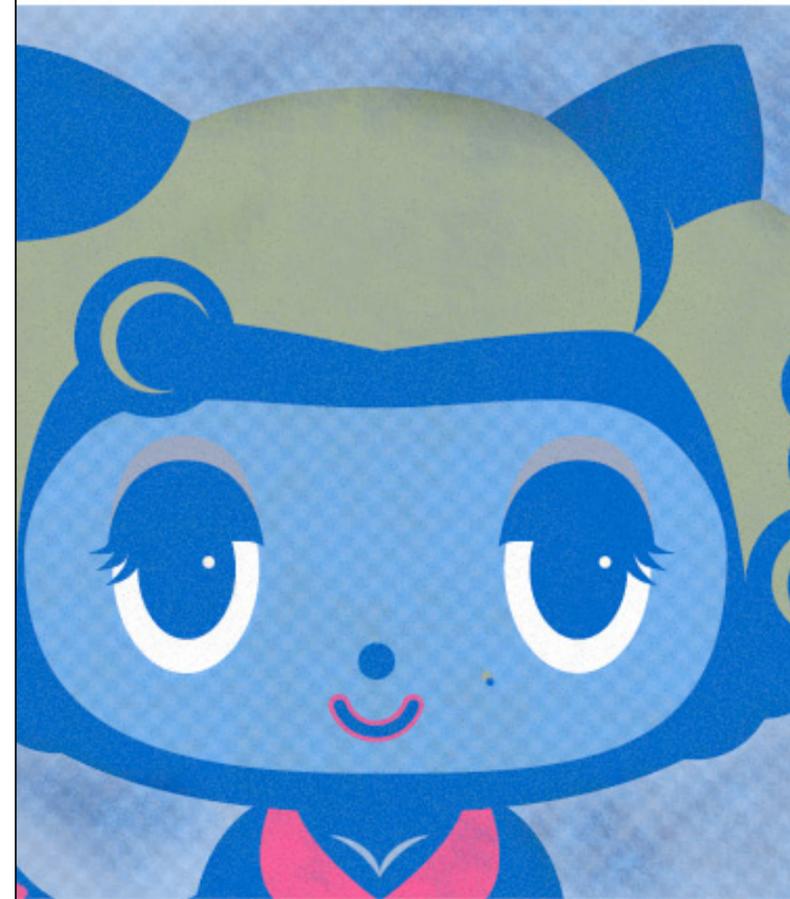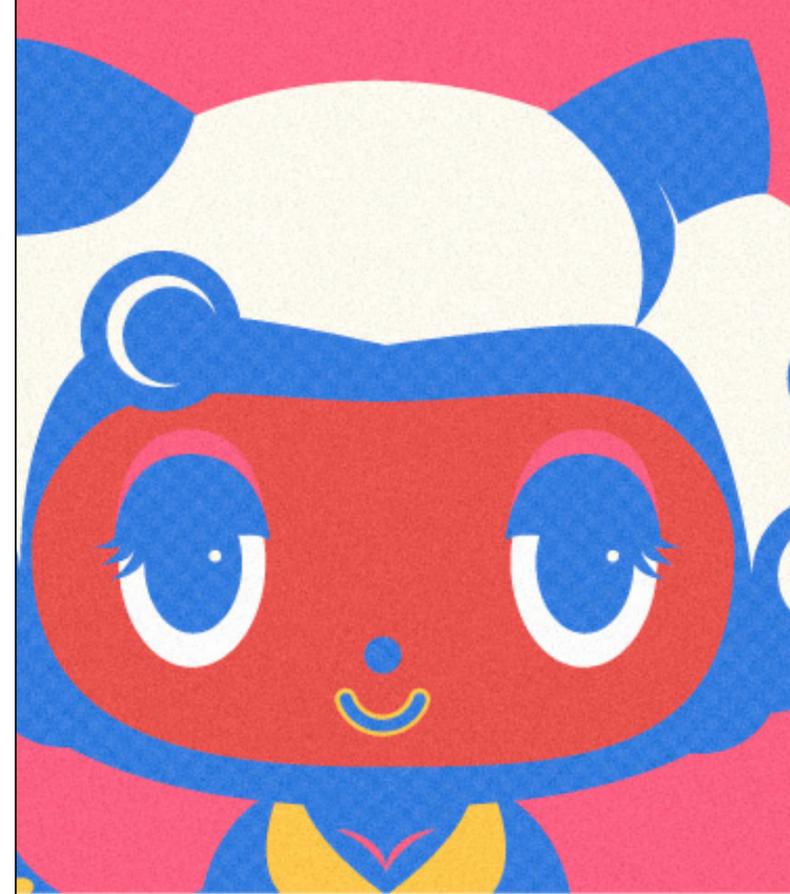
**Who wants to review & publish such replication?**

We do!

# Why GitHub ?

GitHub offers a web-based git repository hosting service with great specific features (issue, pull request, etc).

→ Version control
→ Public repositories
→ Transparency and verifiability
→ Easy exploration of new ideas

A kind of modern lab for the computer scientist.

→ Popular among developers (Google, Microsoft, etc.)
→ Ergonomic & efficient
→ Free (as in beer)

But

→ Closed sources
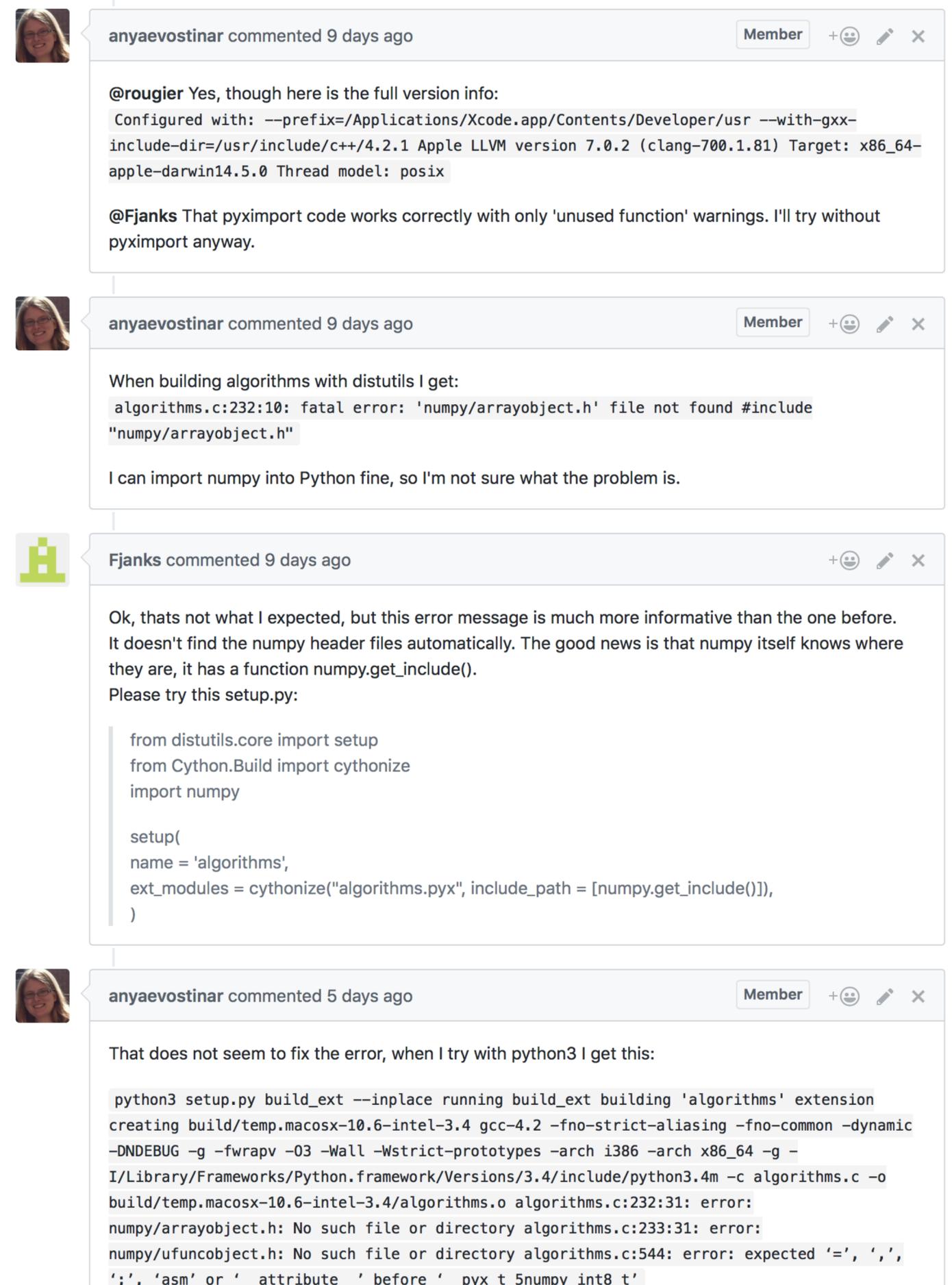→ Ran by a private company
→ Can close tomorrow

# Open peer-review

_____

Editor is publicly assigned by editor-in-chief.

Reviewers are publicly invited to review (they can decline the invitation of course)

The actual review takes place in the discussion area of the issue. Anybody can enter the discussion unless this discussion is locked.

This means anybody can give advice and/or comment because this discussion is public.

---

**anyaevostinar** commented 9 days ago    `Member`  +☺ ✎ ✕

@**rougier** Yes, though here is the full version info:

```
Configured with: --prefix=/Applications/Xcode.app/Contents/Developer/usr --with-gxx-
include-dir=/usr/include/c++/4.2.1 Apple LLVM version 7.0.2 (clang-700.1.81) Target: x86_64-
apple-darwin14.5.0 Thread model: posix
```

@**Fjanks** That pyximport code works correctly with only 'unused function' warnings. I'll try without pyximport anyway.

---

**anyaevostinar** commented 9 days ago    `Member`  +☺ ✎ ✕

When building algorithms with distutils I get:

```
algorithms.c:232:10: fatal error: 'numpy/arrayobject.h' file not found #include
"numpy/arrayobject.h"
```

I can import numpy into Python fine, so I'm not sure what the problem is.

---

**Fjanks** commented 9 days ago    +☺ ✎ ✕

Ok, thats not what I expected, but this error message is much more informative than the one before. It doesn't find the numpy header files automatically. The good news is that numpy itself knows where they are, it has a function numpy.get_include().
Please try this setup.py:

```
from distutils.core import setup
from Cython.Build import cythonize
import numpy

setup(
name = 'algorithms',
ext_modules = cythonize("algorithms.pyx", include_path = [numpy.get_include()]),
)
```
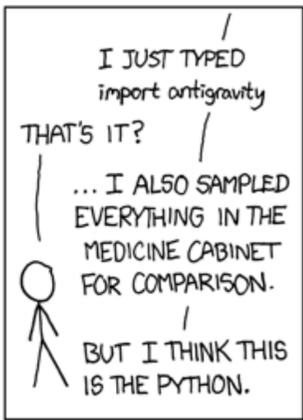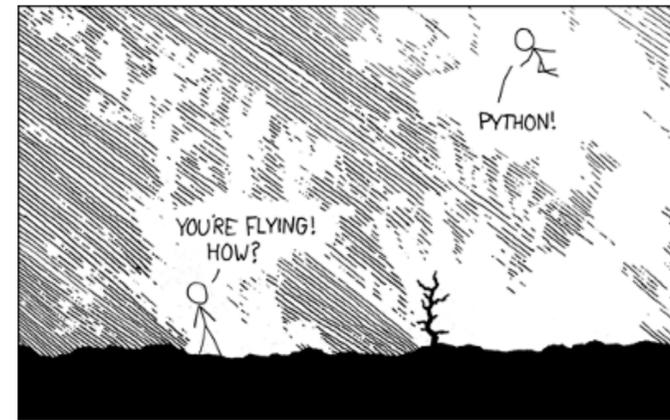
---

**anyaevostinar** commented 5 days ago    `Member`  +☺ ✎ ✕

That does not seem to fix the error, when I try with python3 I get this:

```
python3 setup.py build_ext --inplace running build_ext building 'algorithms' extension
creating build/temp.macosx-10.6-intel-3.4 gcc-4.2 -fno-strict-aliasing -fno-common -dynamic
-DNDEBUG -g -fwrapv -O3 -Wall -Wstrict-prototypes -arch i386 -arch x86_64 -g -
I/Library/Frameworks/Python.framework/Versions/3.4/include/python3.4m -c algorithms.c -o
build/temp.macosx-10.6-intel-3.4/algorithms.o algorithms.c:232:31: error:
numpy/arrayobject.h: No such file or directory algorithms.c:233:31: error:
numpy/ufuncobject.h: No such file or directory algorithms.c:544: error: expected '=', ',',
':', 'asm' or '__attribute__' before '__pyx_t_5numpy_int8_t'
```
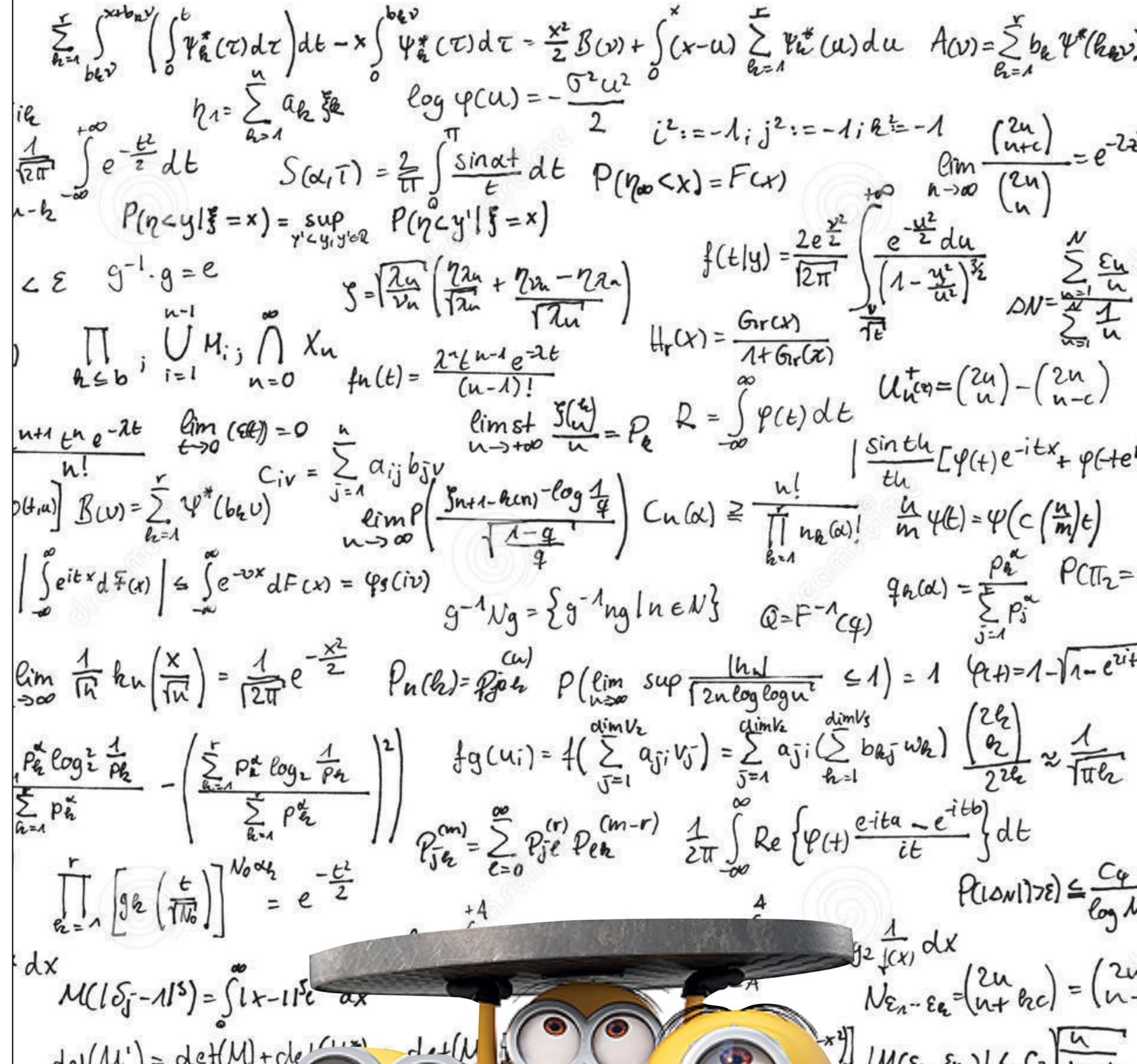
# The Lazarus effect

# Frequently Asked Questions

___

## What kind of research can I replicate?

Any computational research in any domain of science as long as there is an editor from the Board who has the expertise to edit your submission. The editorial board is growing to increase the scientific domains being covered. If no editor is able to edit your submission, you can also propose a guest editor (who must be willing to work with our GitHub-based editorial processes). about replication of my own work?

## I'm a student, can I submit?

Yes ! Students are strongly encouraged to submit their work. Although the ReScience publishing model is a bit different from other academic journals, it can give students a first experience at peer-reviewed scholarly publishing, including meeting standards of scientific rigor and addressing reviewers' comments. Publishing in ReScience is also a way to actively contribute to open science while adding to one's publication record.

# Frequently Asked Questions

_____

## What if I cannot replicate a result?

Some research may not be replicable. Before declaring a research result non-replicable, we require extra caution to be taken. In addition to scrutiny of your submission by reviewers and editors, we will contact the authors of the original research, and issue a challenge to the ReScience community to spot and report (using the issue tracker) errors in your implementation.
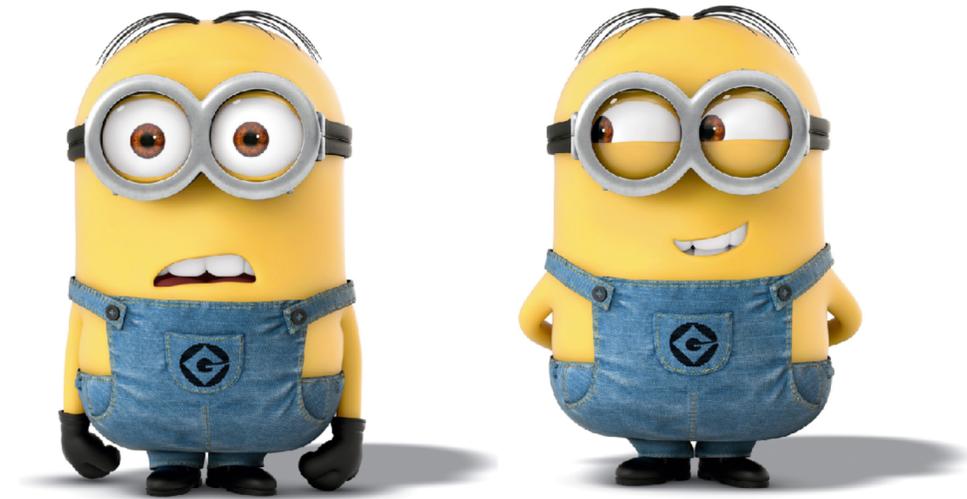
If no errors are found, your work will be accepted and the original research will be declared non-replicable.

## What about replication of my own work?

No. Mistakes in the implementation of research questions and methods are often due to biases authors invariably have, consciously or not. One's biases will inevitably carry over to how one approaches a replication.

Perhaps even more importantly, we aim at the cross-fertilization of research and trying to replicate the work of one's peers might pave the way for a future collaboration, or may give rise to new ideas as a result of the replication effort.
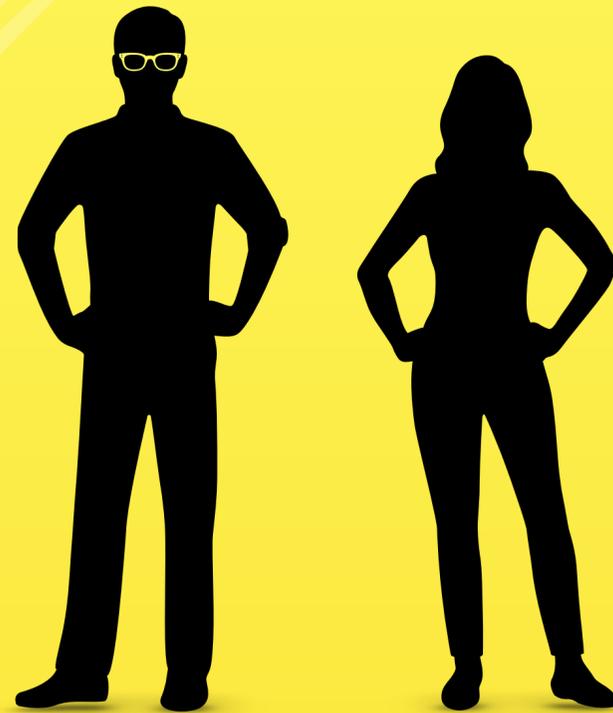
**This changed in 2020…**

WHAT'S NEXT?

# Re|Science X

FREE TO READ - FREE TO PUBLISH

SINCE JANUARY 2021