# Research software preservation: a publisher's perspective

Naomi Penfold, Innovation Officer

Email: n.penfold@elifesciences.org

Twitter: @eLifeInnovation

## Maria Guerreiro

Journal Development Editor

Summary and slides available at elifesci.org/software-preservation
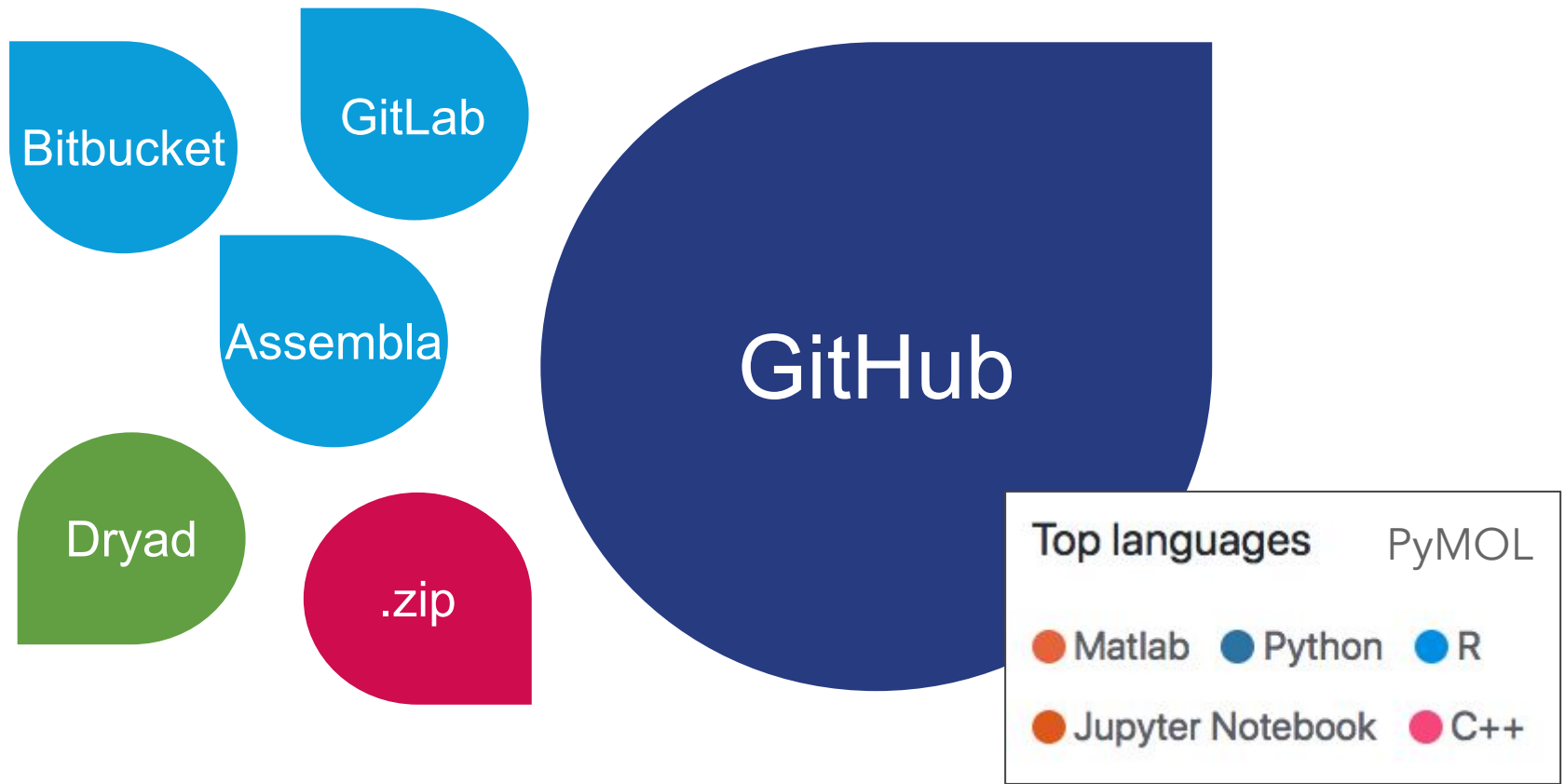
eLIFE

Helping scientists accelerate discovery by operating a platform for research communication that encourages and recognises the most responsible behaviours in science
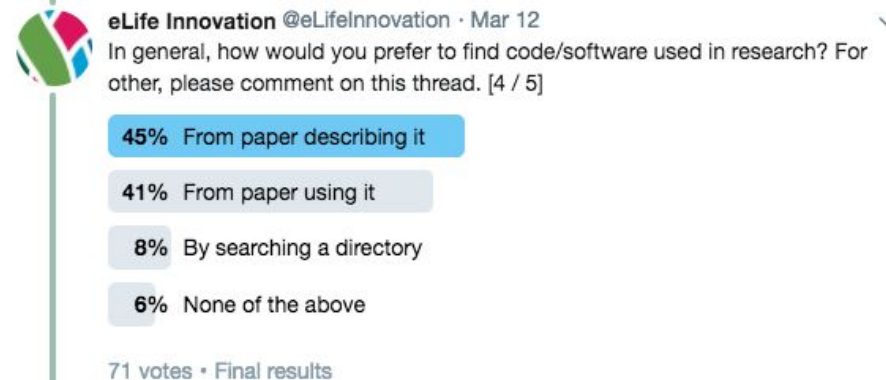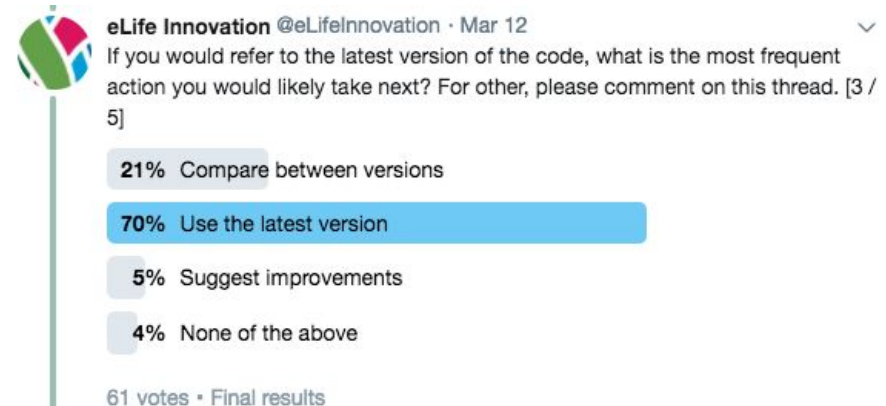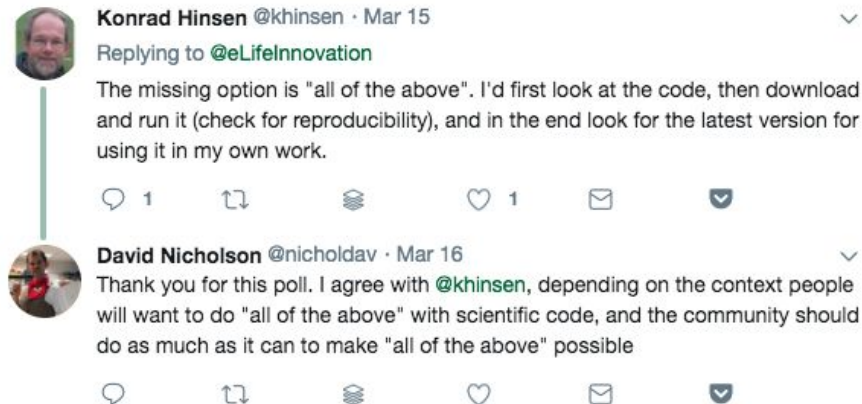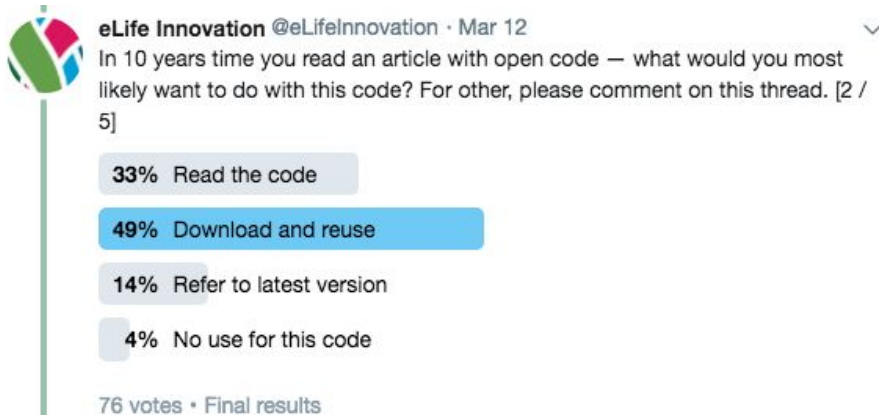
# Today…

- Software shared with eLife
- How we cite and preserve software
- Cost-benefit of developing?
- Key requirements
- Opportunity to encourage best practice?

# Research software shared with eLife



Bitbucket

GitLab

Assembla

Dryad

.zip

GitHub

Top languages    PyMOL

● Matlab   ● Python   ● R

● Jupyter Notebook   ● C++

# Researchers want to keep options open for reuse



**eLife Innovation** @eLifeInnovation · Mar 12
In 10 years time you read an article with open code — what would you most likely want to do with this code? For other, please comment on this thread. [2 / 5]

33% Read the code
49% Download and reuse
14% Refer to latest version
4% No use for this code

76 votes · Final results

**Konrad Hinsen** @khinsen · Mar 15
Replying to @eLifeInnovation
The missing option is "all of the above". I'd first look at the code, then download and run it (check for reproducibility), and in the end look for the latest version for using it in my own work.

💬 1   ⟲   ⩘   ♡ 1   ✉   ⌄

**David Nicholson** @nicholdav · Mar 16
Thank you for this poll. I agree with @khinsen, depending on the context people will want to do "all of the above" with scientific code, and the community should do as much as it can to make "all of the above" possible

💬   ⟲   ⩘   ♡   ✉   ⌄

**eLife Innovation** @eLifeInnovation · Mar 12
If you would refer to the latest version of the code, what is the most frequent action you would likely take next? For other, please comment on this thread. [3 / 5]

21% Compare between versions
70% Use the latest version
5% Suggest improvements
4% None of the above

61 votes · Final results

**eLife Innovation** @eLifeInnovation · Mar 12
In general, how would you prefer to find code/software used in research? For other, please comment on this thread. [4 / 5]

45% From paper describing it
41% From paper using it
8% By searching a directory
6% None of the above

71 votes · Final results

Dr. Rachael Tatman
@rctatman

Follow

Reproducibility tip of the day: If you're sharing research code for a paper, make sure to double-check that you've 1) included a link in the paper 2) the link works and 3) your code is actually at the link. (Too many papers have links to "coming soon" empty repos 😭 😭 😭)
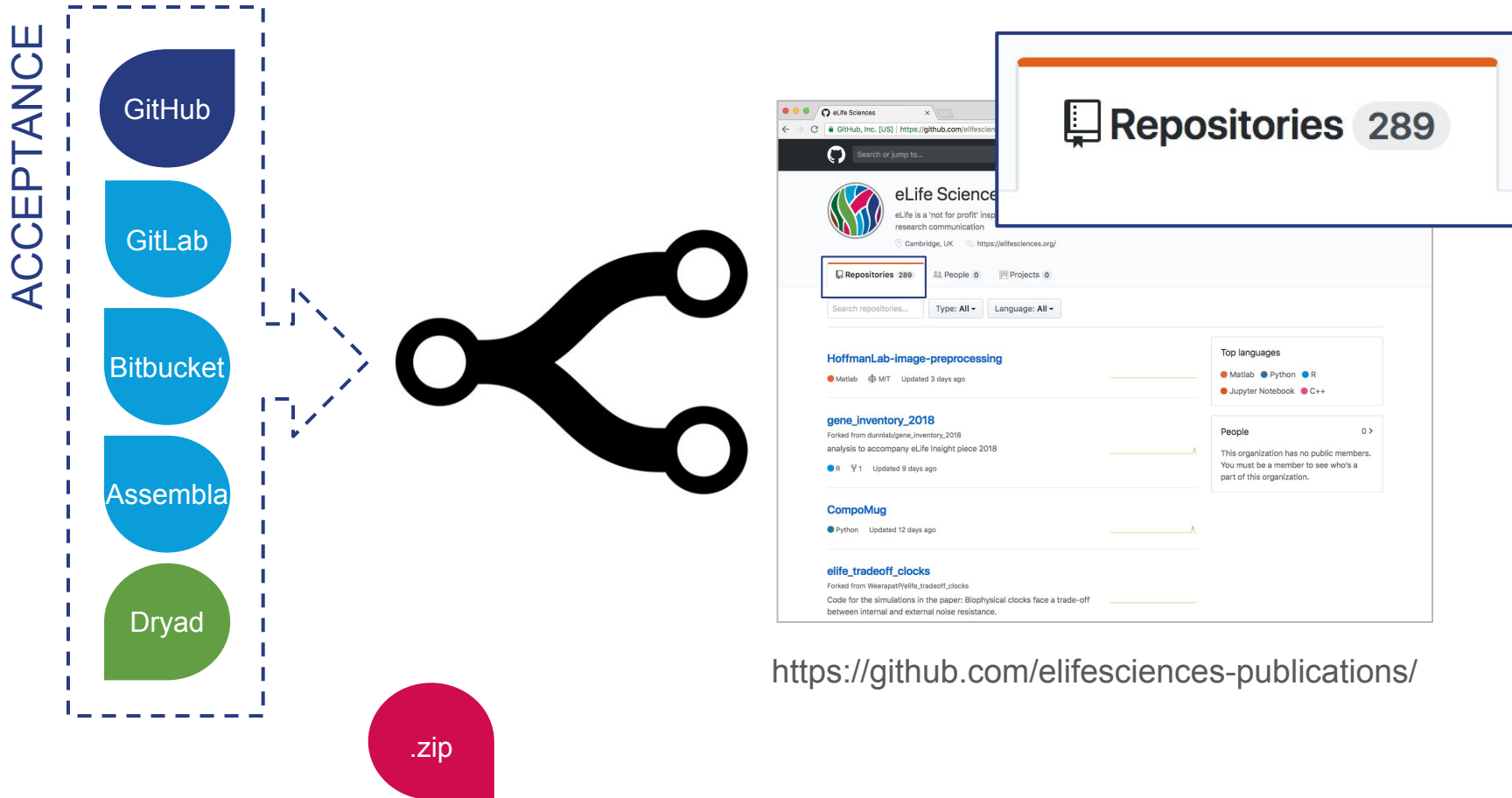
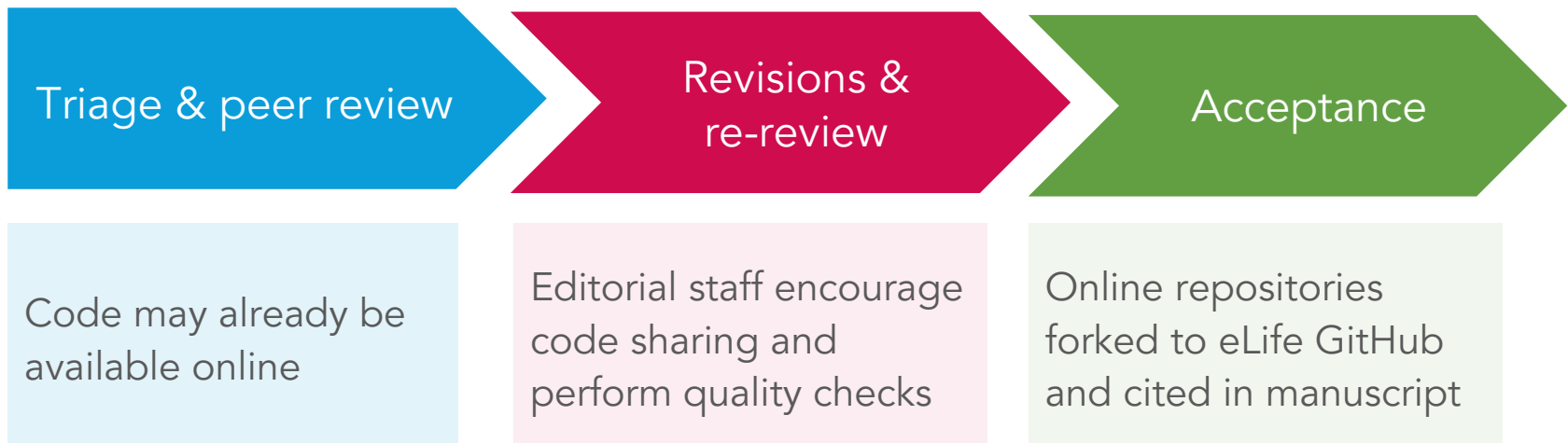10:18 AM - 6 Jul 2018

13 Retweets  86 Likes

https://twitter.com/rctatman/status/1015283853131304960
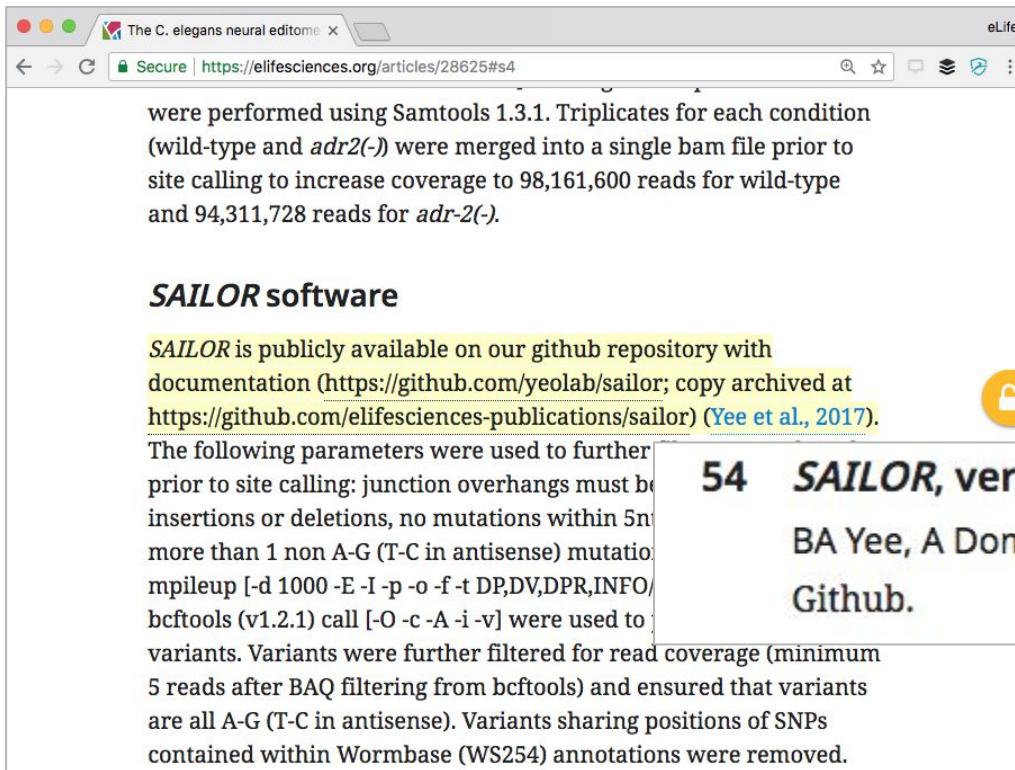
@eLifeInnovation
elifesci.org/software-preservation

# We fork to our own GitHub repository

ACCEPTANCE

GitHub

GitLab

Bitbucket

Assembla

Dryad

.zip

https://github.com/elifesciences-publications/

# The workflow in detail

| Triage & peer review | Revisions & re-review | Acceptance |
|---|---|---|
| Code may already be available online | Editorial staff encourage code sharing and perform quality checks | Online repositories forked to eLife GitHub and cited in manuscript |

# Citation in the text and references



were performed using Samtools 1.3.1. Triplicates for each condition (wild-type and *adr2(-)*) were merged into a single bam file prior to site calling to increase coverage to 98,161,600 reads for wild-type and 94,311,728 reads for *adr-2(-)*.

## *SAILOR* software

*SAILOR* is publicly available on our github repository with documentation (https://github.com/yeolab/sailor; copy archived at https://github.com/elifesciences-publications/sailor) (Yee et al., 2017). The following parameters were used to further [...] prior to site calling: junction overhangs must be [...] insertions or deletions, no mutations within 5n[...] more than 1 non A-G (T-C in antisense) mutatio[...] mpileup [-d 1000 -E -I -p -o -f -t DP,DV,DPR,INFO[...] bcftools (v1.2.1) call [-O -c -A -i -v] were used to [...] variants. Variants were further filtered for read coverage (minimum 5 reads after BAQ filtering from bcftools) and ensured that variants are all A-G (T-C in antisense). Variants sharing positions of SNPs contained within Wormbase (WS254) annotations were removed.

54    *SAILOR*, version 9a57b4b
BA Yee, A Domissy, EC Wheeler, GW Yeo (2017)
Github.

eLife 2017;6:e28625 DOI: 10.7554/eLife.28625

# Research software continues to develop

# Preserving research software at eLife

## Benefits

- Scientist-driven
- Less work for authors
- Reuse is facilitated
- **Encourages best practice: powerful when combined with data and other resources**

## Limitations

- No DOI
- Reliant on Github
- Requires staff time
- Not for source code files shared directly with the journal; what are the advantages/disadvantages of hosting code on the journal website?

# Going further: is it worth it?

- How much are we prepared to invest as a community in process innovation or development? Or to support an archive?
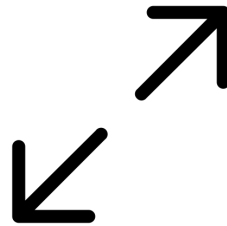- For how much added value? To whom?

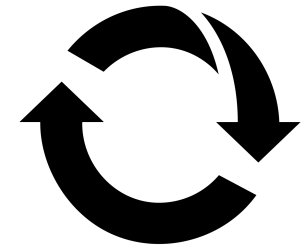How do we evaluate cost-benefit for software preservation?
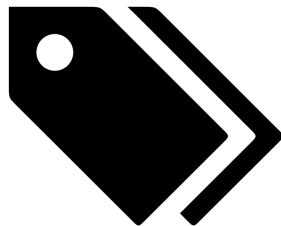
# Requirements as a publisher

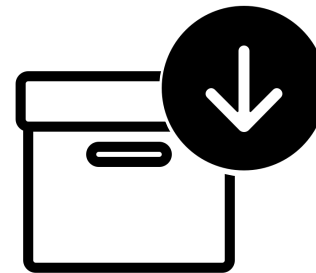Minimise work for authors, make it sustainable for publishers

Scalable and agnostic to platform, format or tool

For reusability wherever possible

Metadata collected at source, compliant with citation guidelines

Persistent and retrievable for as long as is reasonable

# Can we encourage best practice?

☑ LICENSE.md

☑ CITATION.cff or codemeta.json

⏩ Include event-driven process in open source publishing platform

Together with:
- Text
- Data
- Key Resources Table

→ open and reproducible research package

What if the software is **<u>not</u>** open?

How do we do the best for all cases?

# Further investigations

- Software citations, interactions, and activity at eLife
- Other publisher workflows

<span style="color:#d0006f">Collaboration welcome</span>

# In summary…

- Software shared with eLife is mainly on Github but we need to cater for any source
- We fork online repositories to our Github upon acceptance but have no process to preserve source code files
- We would like to minimise burden and cost, encourage best practice, and support researchers to reuse where reasonable

Today we ask:

- What do other publishers do?
- Can we encourage best practice?
- Can we help you test new process(es)?
- **Is it worth it?**

# Questions?

Slides: elifesci.org/software-preservation
Email: n.penfold@elifesciences.org
Twitter: @eLifeInnovation