# `commit` -ment issues with Git: investigating & archiving y'alls work

Sarah Nguyen, Vicky Steeves (Presenting)
& Genevieve Milliken

csv,conf,v5 | 2020.05.14

Slides: osf.io/uftkn

# About IASGE
`(ice-age)`

# Who are we?

New York University, Division of Libraries!

In addition to the core team (below), this project wouldn't be possible without the time and efforts of colleagues in Digital Library Technology Services



Vicky Steeves
Project Lead

Genevieve Milliken
Research Scientist

Sarah Nguyen
Research Scientist

# Project overview

An Alfred P. Sloan Foundation funded project, IASGE (pronounced `ice-age`) has two main streams of work:

1) Study how academics/folks in academia are using Git and Git hosting platforms and how these tools could be better aligned with their needs
2) Evaluate the extent to which the scholarship on Git hosting platforms is being preserved by professionals, and write an archival spec

The results of this project aim to inform the way code and annotations on Git hosting platforms move from a phase where they are highly active and collaborative, to a state where they are stable, permanently citable, and under **active, professional preservation**.

# What is Git?

- Git is a revision control system, a program to manage your source code history
- It is strictly a command-line tool
- Revision control systems let us compare, restore, and merge changes to our [plain-text] stuff

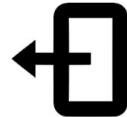This is hugely important for collaboration and transparency!

In case of fire

1. git commit
2. git push
3. leave building

# What are "Git Hosting Platforms"

Literally, places that host git repositories on the web

They are NOT the same as Git, but rather are places where you can upload Git repositories with some additional features

The most popular include:

1. GitHub
2. GitLab
3. Bitbucket
4. Sourceforge

| Name | Manager | Established | Server side: all free software | Client side: all-free JS code |
|------|---------|-------------|-------------------------------|-------------------------------|
| Assembla | Assembla, Inc | 2005 | No | Unknown |
| Azure DevOps Services | Microsoft | 2012[1] | No | No |
| Bitbucket | Atlassian | 2008 | No | No |
| Buddy | Buddy, LLC. | 2015 | No | No |
| CloudForge | CollabNet | 2012 | No | Unknown |
| Gitea | Gitea organization (open source community)[4] | 2016 | Yes | Yes |
| GForge | The GForge Group,Inc. [5] | 2006 | Partial | Yes |
| GitHub | GitHub, Inc | 2008-04 | No | No |
| GitLab | GitLab Inc. | 2011-09[6] | Partial[7] | Yes[8] |
| GNU Savannah | Savannah Administration | 2001-01 | Yes | Yes |
| Helix TeamHub | Perforce Software | 1995 | No | No |
| Launchpad | Canonical | 2004 | Yes | No |
| OSDN | OSDN K.K. (Q11237954) | 2002-04 | Unknown | Yes |
| Ourproject.org | Comunes Collective | 2002 | Yes | Yes |
| OW2 Consortium | OW2 Consortium | Unknown | Unknown | No |
| Phabricator | Phacility, Inc | 2010 | Yes | Yes |
| Rosetta Code | Unknown | 2007 | Unknown | Unknown |
| SEUL | Unknown | 1997-05 | Unknown | No |
| SourceForge | BizX LLC | 1999-11 | Yes[14][15] | Yes |

# Examples of "scholarly Git" usage

1.  Publishing code and data as supplementary materials
2.  Quality assurance workflows for data analysis
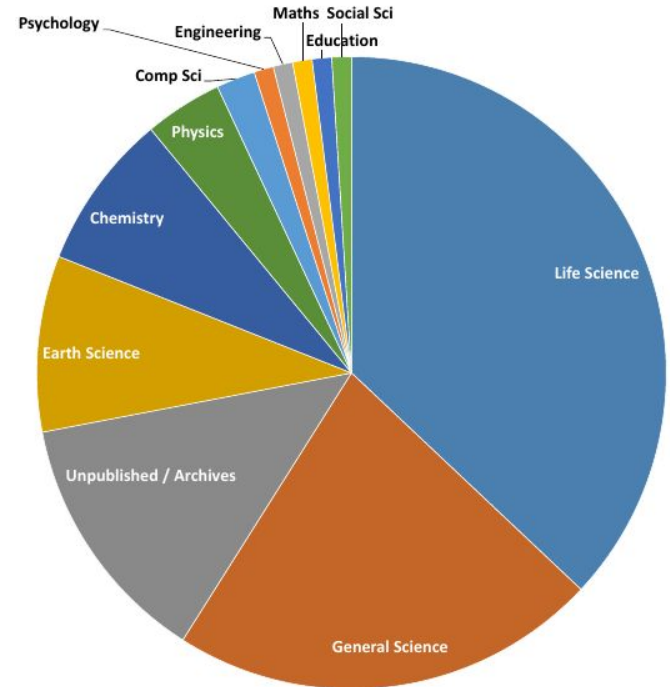3.  Journal infrastructure with peer review

# Estimated scope of scholarship in GHPs

"Over 5,000 Github software repositories have been identified as research software according to the criteria explained previously: either a research publication referenced the software repository, or the software repository referenced a research publication."
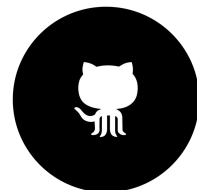
-- Hasselbring, Wilhelm, et al. "FAIR and Open Computer Science Research Software." ArXiv:1908.05986 [Cs], Aug. 2019. http://arxiv.org/abs/1908.05986.



Research areas of publications cited from Github repositories from arXiv: 1909.05986

# What's the problem we want to solve

- Researchers use a variety of scholarly tools on the web during the research process, which includes designing, developing, and refining (through versioning) source code
- This source code is contextualized by the "scholarly ephemera" associated with it (e.g. issue disc.)
- No project currently captures both source code and scholarly ephemera

# 🐘 in the room: GitHub's Archiving Program

There are a lot of open questions about GitHub's Archive Program which are probably shared by people in this room

The fact remains that none of the partners or solutions here capture the ephemera + source code reliably together, which we posit as important for usability in the long-term

Also, read DSHR's takedown:
blog.dshr.org/2019/11/seeds-or-code.html

# Gap Analysis to Understand Scholars Using Git

# Researching the Scholarly Git Experience



Literature Review

Systematic Review

Focus Group

Broad Survey

User Experience Interviews
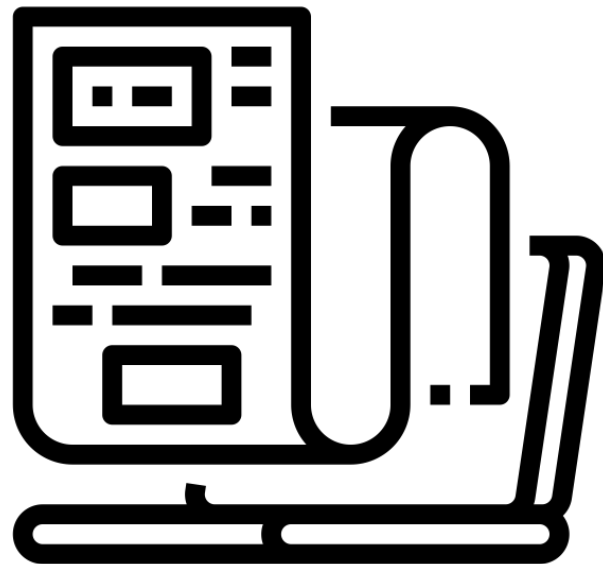
# Research Questions

How are scholars currently using this toolkit of
Git + GHPs?

How can features in Git + GHPs serve scholarship?

How can we improve teaching Git for minimal users?

Follow our IASGE blog  for updates

investigating-archiving-git.gitlab.io/updates/

# Examples of scholarly Git* usage

| Git Experiences | Related GHP Feature | Related Git Commands |
|---|---|---|
| **Version control** | - Commit logs<br>- Branches | `git log`<br>`git diff` |
| **Community & collaboration** | - Issue Tracker<br>- Pull requests | `git add <files>`<br>`git commit -m "[message]"`<br>`git push` |
| **Method tracking** | - README<br>- Wiki<br>- Posts<br>- Commit logs | `git commit logs` |
| **Education** | - README<br>- Wiki<br>- Issue Tracker<br>- Pull requests | `open-issues`<br>`close-issues`<br>`list-issues`<br>`check-review` |
| **Data processing** | - Continuous integration | (various) |
| **Reproducibility** | - README<br>- Continuous integration | `git clone`<br>`git pull` |
| **Publishing** | - Pages services<br>- README | (various) |

# A graduate student perspective on overcoming barriers to interacting with open-source software

Oihane Cereceda ✉, Danielle E.A. Quinn ✉

⬇ PDF      ⬇ Citation (RIS )      ⬇ Citation (BibTeX )

# Abstract

Computational methods, coding, and software are important tools for conducting research. In both academic and industry data analytics, open-source software (OSS) has gained massive popularity. Collaborative source

Many academic institutions are using GitHub to share COVID-19 Data

Johns Hopkins University's Center for Systems Science & Engineering

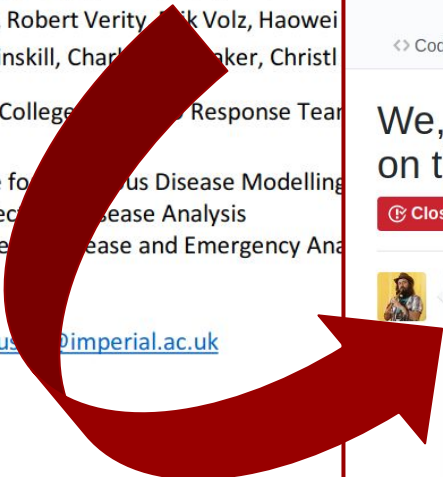Academic institutions are using GitHub to share COVID-19 models

**Report 9: Impact of non-pharmaceutical interventions (NPIs) to reduce COVID-19 mortality and healthcare demand**

Neil M Ferguson, Daniel Laydon, Gemma Nedjati-Gilani, Natsuko Imai, Kylie Ainslie, Marc Baguelin, Sangeeta Bhatia, Adhiratha Boonyasiri, Zulma Cucunubá, G Dorigatti, Han Fu, Katy Gaythorpe, Will Green, Arran Ham Elsland, Hayley Thompson, Robert Verity, ik Volz, Haowei Caroline Walters, Peter Winskill, Char aker, Christl

On behalf of the Imperial Colle Response Tear

WHO Collaborating Centre fo us Disease Modelling MRC Centre for Global Infec sease Analysis Abdul Latif Jameel Institute ease and Emergency Ana Imperial College London

Correspondence: neil.fergu @imperial.ac.uk

Imperial College of London's MRC Centre for Global Infectious Disease Analysis (MRC GIDA)

mrc-ide / **covid-sim**

⊙ Watch ▾  75      ★ Star  823      Fork  163

<> Code      ⊙ Issues  29      Pull requests  8      ⊙ Actions      🛡 Security  0      Insights

We, the undersigned software engineers, call for any papers based on this codebase to be immediately retracted. #165

🔴 Closed    **jMyles** opened this issue 5 days ago · 43 comments

**jMyles** commented 5 days ago · edited ▾

The tests in this project, being limited to broad, "smoke test"-style assertions, do not support an assurance that the equations are being executed faithfully in discrete units of logic, nor that they are integrated into the application in such a way that the accepted practices of epidemiology are being modeled in accordance with the standards of that profession.

Billions of lives have been disrupted worldwide on the basis that the study produced by the logic contained in this codebase is accurate, and since there are no tests to show that, the findings of this study (and any others based on this codebase) are not a sound basis for public policy at this time.

I want to be clear that this Issue is not meant to denigrate the authors of this code - we've all written code that isn't our best work and code that is untested. But when a codebase is used to craft scholarly publications that are in turn used to influence public policy, the authors of those publications (and ultimately policy) need to ensure that the science is verifiable in a public sense. The lack of tests makes that an impossibility. So closure of this Issue, by retraction of studies based on it, is meant as a critique of the publication and policy authors, not the contributors to this repo.

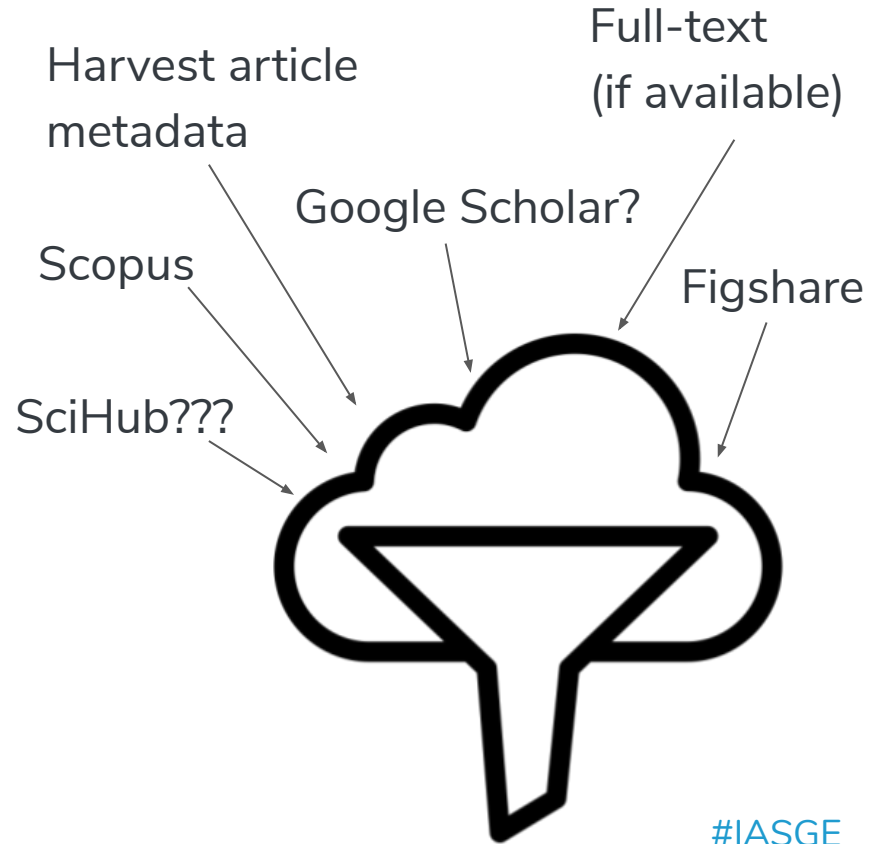Assignees
No one assigned

Labels
None yet

Projects
None yet

Milestone
No milestone

Linked pull requests

# Part I: Systematic Review (Quant)

An approach to understand the landscape of published scholarly articles that reference Git repositories.

"All our source code is available on [GitLab], to allow community to reproduce our results, from the training of the networks, until the statistical analyses." (Perez, 2019)
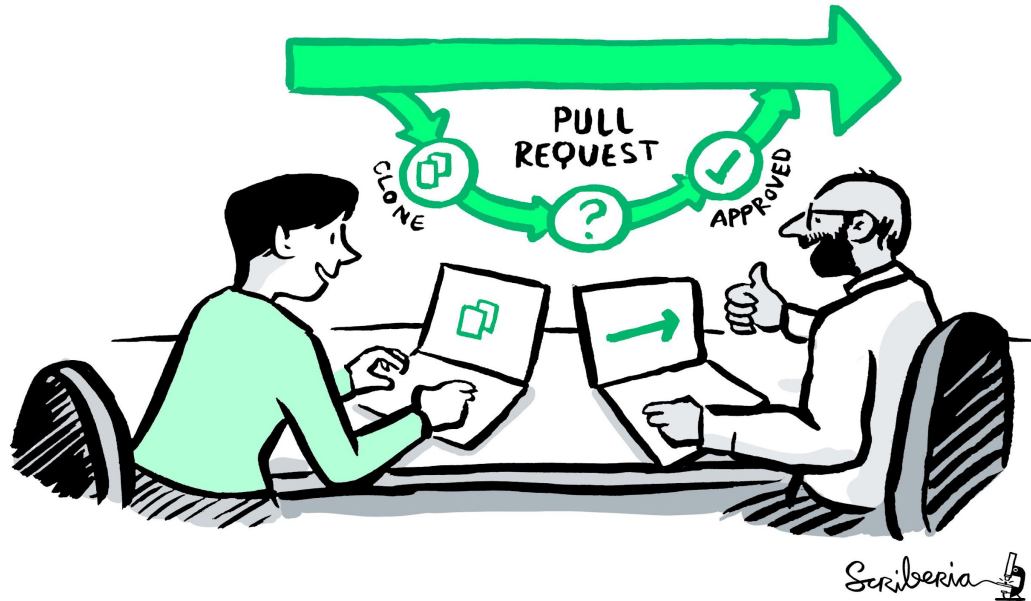
Harvest article metadata

Full-text (if available)

Google Scholar?

Scopus

Figshare

SciHub???

# Part II: Focus Group (Qual)

*. . . teachers were surprised by how overwhelming the student enthusiasm is for adopting VCS, [but] they also discovered that they lacked understanding about the system or having confidence in their ability to use it effectively beyond the course.*

──Glassey, R. (2019). Adopting Git/Github Within Teaching: A Survey of Tool Support. *Proceedings of the ACM Conference on Global Computing Education*, 143–149. https://doi.org/10.1145/3300115.3309518
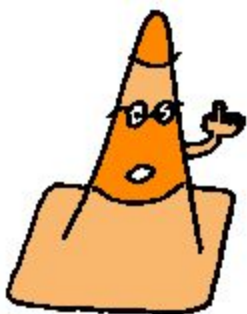
#IASGE

**Kirstie Whitaker** @kirstie_j · 2h

I've done this a few times in the last month and it really makes me happy. A GitHub pull request is truly one of the biggest barriers to more folks working collaboratively online. We have these great tools, but we don't teach them often or compassionately enough.

2        1            ❤ 10

# Part III: Survey (Quant)

Target population:
Scholars who use Git across all disciplines & statuses

Goal:
To gather a wide-ranging & comparable census

Please participate & share widely!
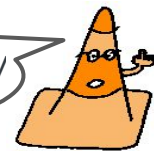bit.ly/3aO0ykQ

**Closes June 22nd**

Themes

❏ Learning Git

❏ Teaching Others

❏ Daily Use

❏ Features on hosting platforms

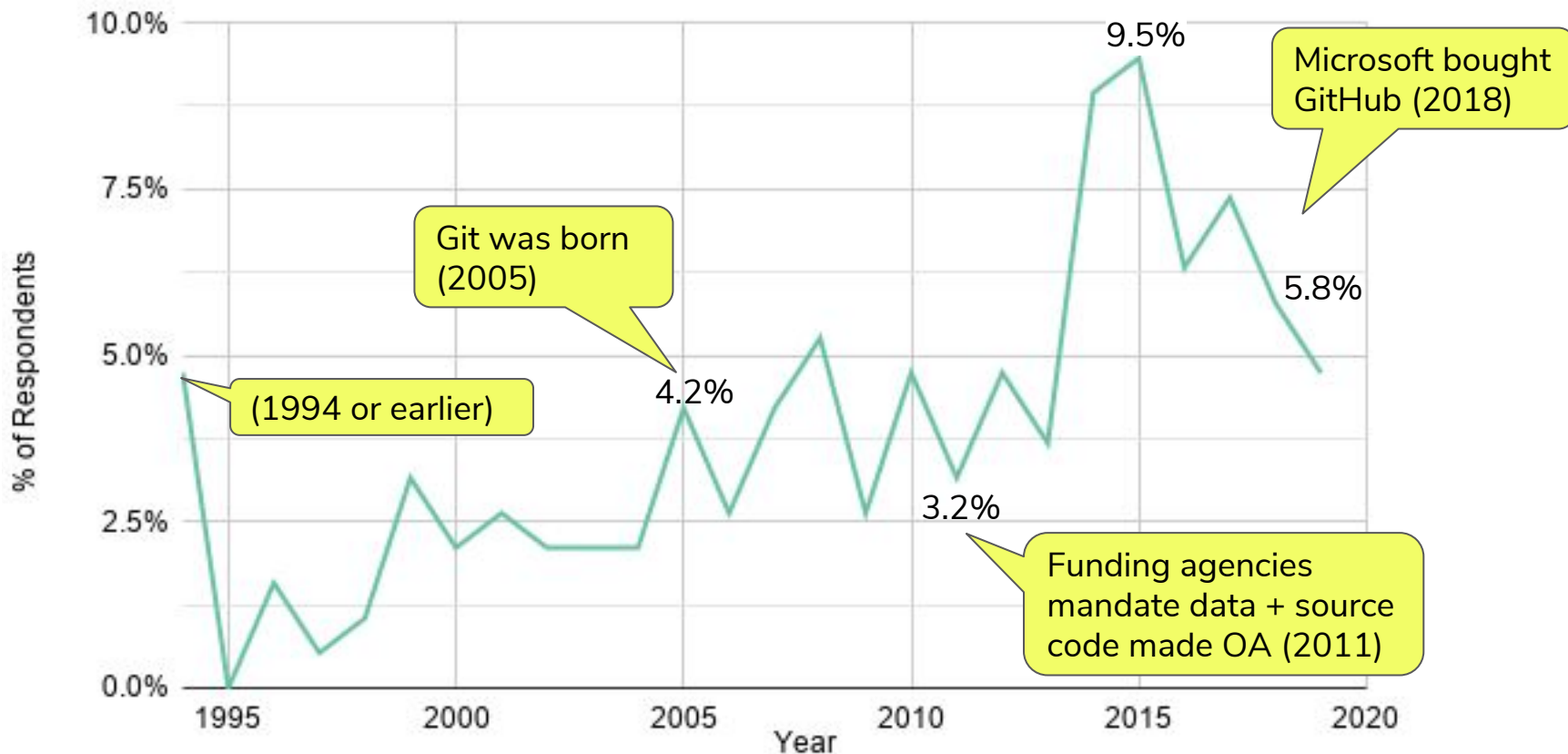❏ Scholarship

❏ Follow-up

# Survey - Preliminary Findings
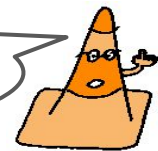
Please participate & share widely
bit.ly/3aO0ykQ

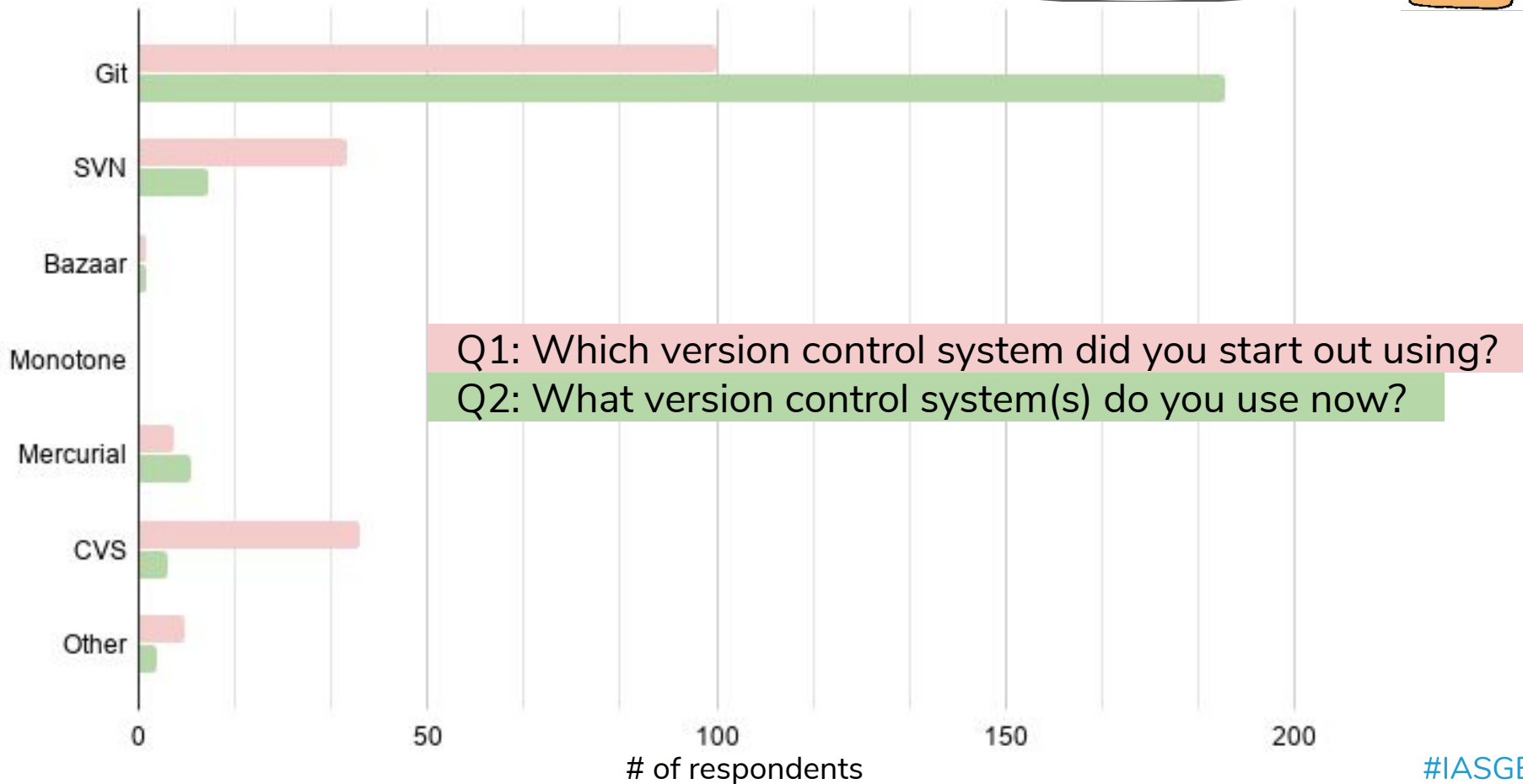Q: When did you first start using a version control system?

# Survey - Preliminary Findings

Please participate & share widely!
bit.ly/3aO0ykQ



Q1: Which version control system did you start out using?
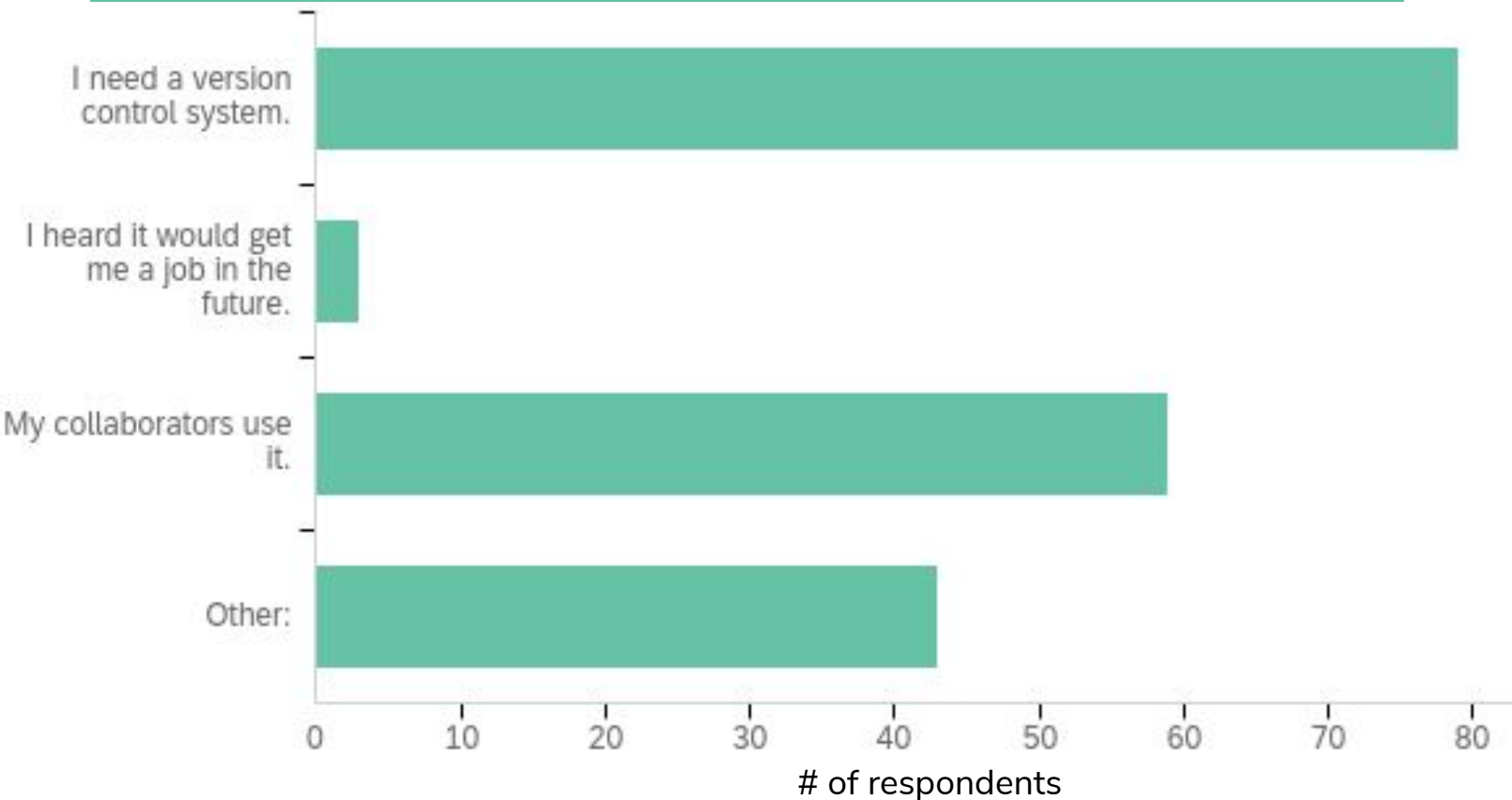Q2: What version control system(s) do you use now?

#IASGE

# Survey - Preliminary Findings

Please participate & share widely!
bit.ly/3aO0ykQ

Q: Why did you first enter the world of git and version control?
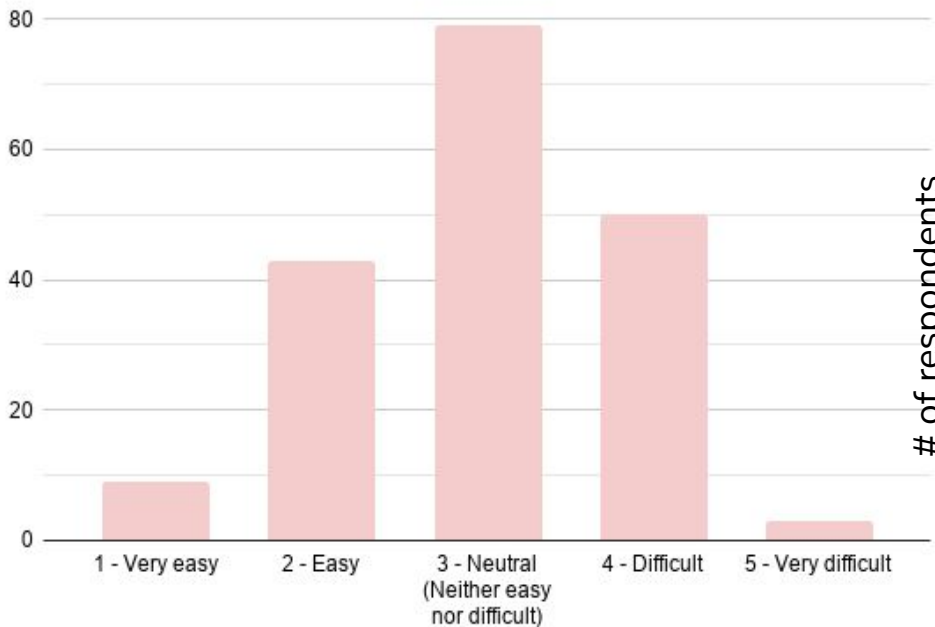


# of respondents

#IASGE

# Survey - Preliminary Findings

Please participate & share widely!
bit.ly/3aO0ykQ
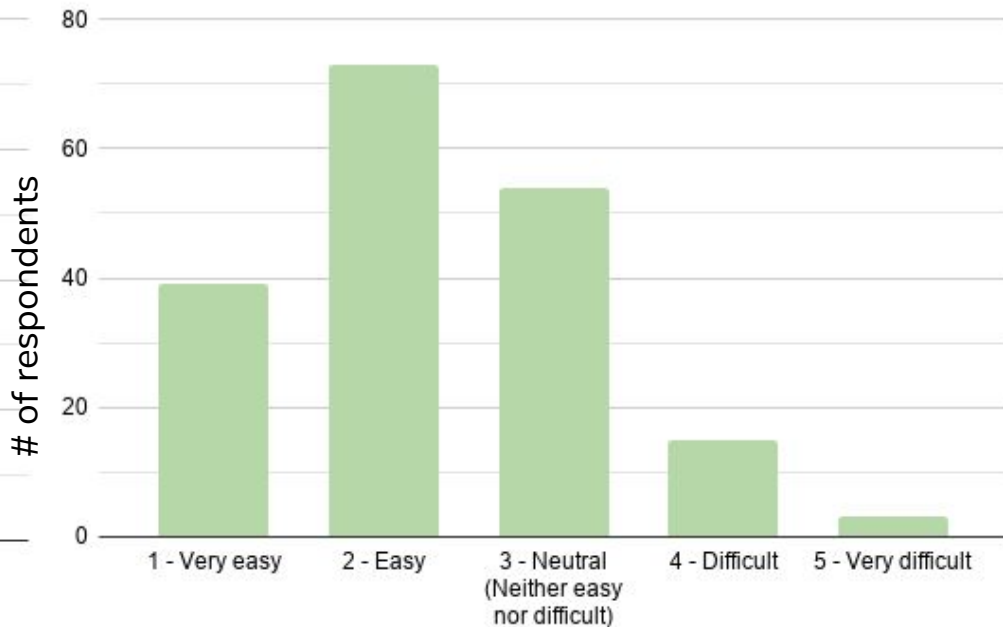
Q1: How difficult was it for you to learn how to use git on your local computer?

Q2: How difficult was it for you to learn how to use the git hosting platform?
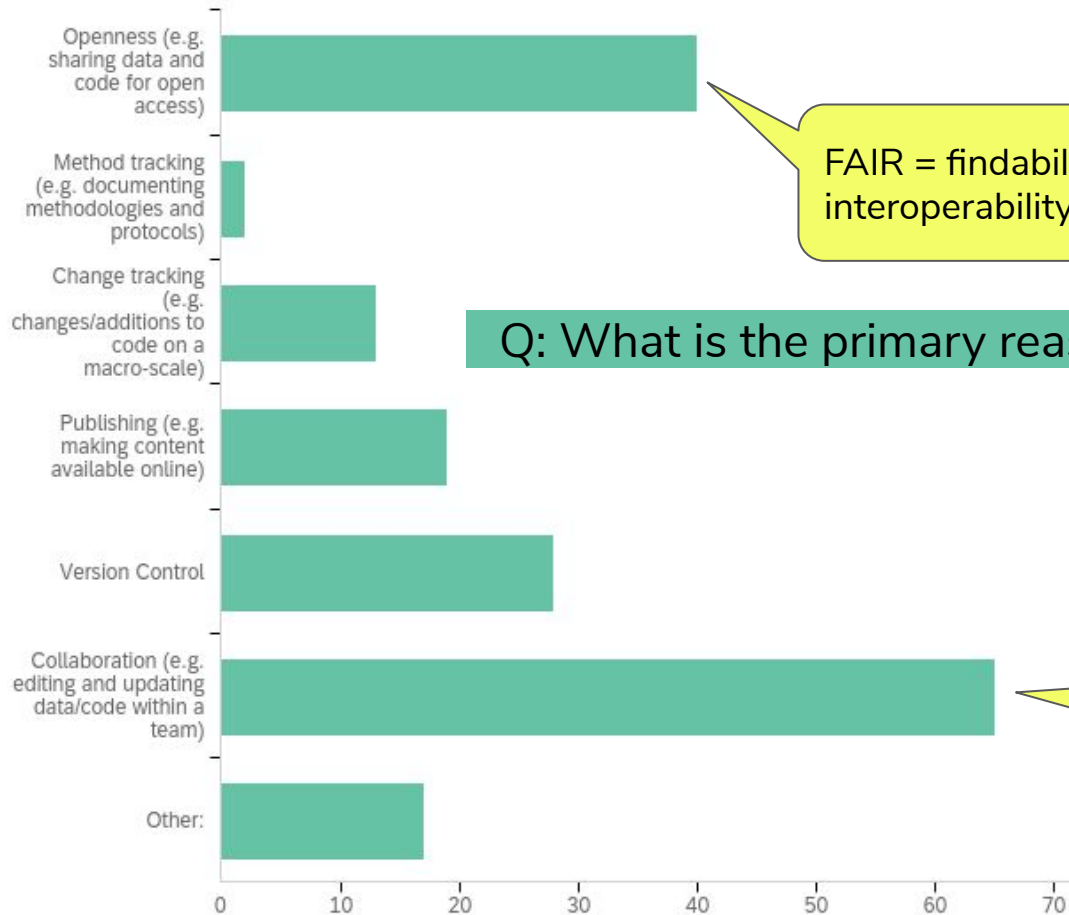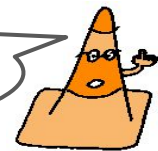


Difficulty Scale

Difficulty Scale

#IASGE

# Survey - Preliminary Findings

Please participate & share widely!
bit.ly/3aO0ykQ



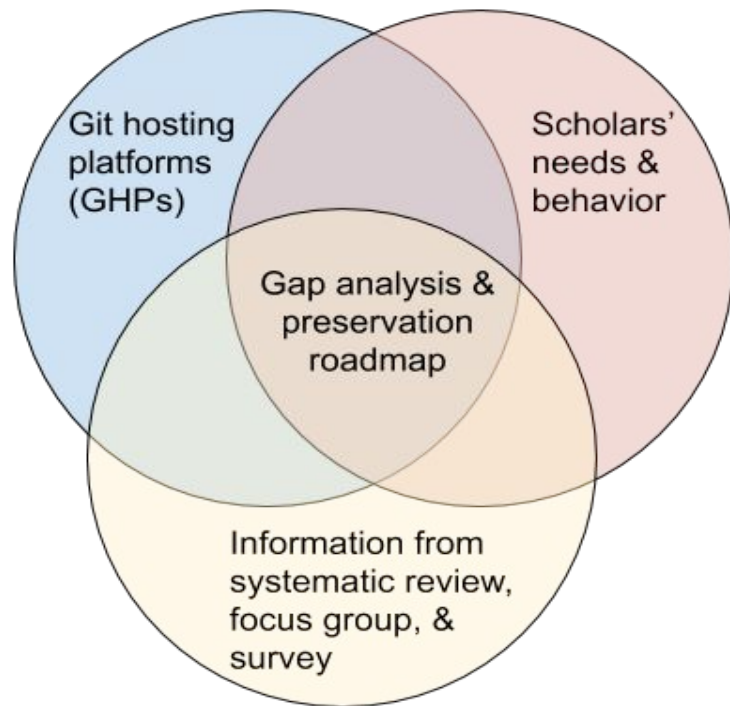FAIR = findability, accessibility, interoperability and reusability

Q: What is the primary reason you use git hosting platforms?

Community, collaboration, & peer production

# Part IV: User Experience Interviews (Qual)

Semi-structured interviews with 50 scholars to understand their behaviours

- Why did folks stop using Git & GHPs?
- If/how are they versioning their work without these toolkits?
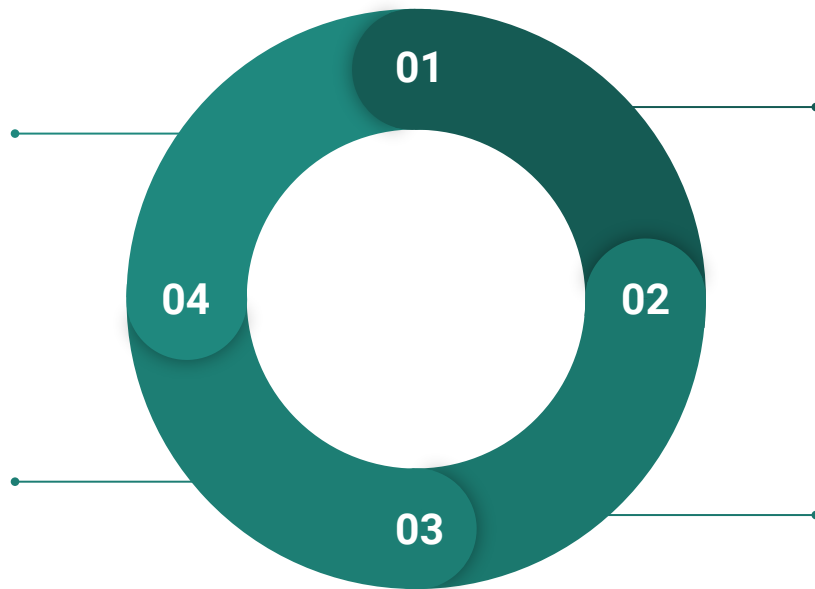- What features & workflows are in practice that we have not heard of before?



Git hosting platforms (GHPs)

Scholars' needs & behavior

Gap analysis & preservation roadmap

Information from systematic review, focus group, & survey

# Current
# Archival
# Approaches

# IASGE Environmental Scan



**Software Preservation**

Best practices, software curation and description

**01**

**Web Archiving**

State-of-the-art web archiving tools and technologies; Could they be used for software capture?

**02**

**Self-Archiving**

Motivations to self-archive; appeal of general repositories vs. institutional repositories

**03**

**04**

**Programmatic Captures**

large-scale archiving of GitHub API data; large-scale archiving of source code from Git hosting platforms; select archiving of repos
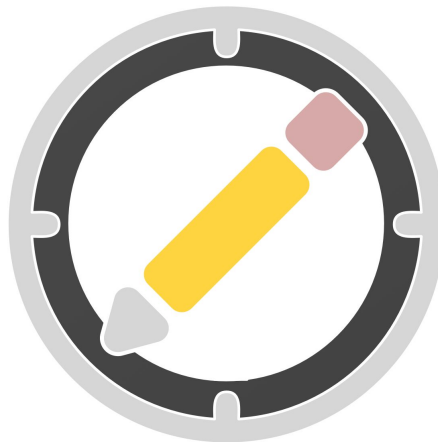
#IASGE

# Web Archiving

- The use of web crawlers/software to capture web-based content
- Understand who, if anyone, is currently using this technology for code repositories

- Projects of Interest:
  - Archive-It
  - Webrecorder
  - Memento Tracer

# Self-Archiving

1. Identify motivations for self-archiving in general repositories
2. Identify gaps regarding IRs and hosting/describing software
3. Understand platform integration (GitHub <> Zenodo)
4. Should we keep rehabbing IRs or move to more flexible models?



Software in Zenodo with Integrations with GitHub and indexed by OpenAIRE

# Programmatic Captures

1. Software to capture software
2. The use of indexing, cloning, and APIs to capture software and contextual information
3. Scholarly ephemera in one place, software in another
4. Projects of interest
    a. Software Heritage
    b. GHTorrent and GH Archive
    c. SARA

The GHTorrent project

# Software Preservation



1. Understanding current communities of practice
2. Software metadata & citation
3. Software curation
4. Projects of interest
   a. Software Preservation Network
   b. Software Emulation (e.g. EaaSI)
   c. Software Sustainability Institute
   d. US Research Software Sustainability Institute

# Towards an Archival Spec

Culmination of this research will be an archival spec that can be used by institutions.

Details will include:

1. The capture of both source code and its ephemera (commit messages, merge requests, issues, wikis)
2. Description and curation
3. Sample preservation workflows

# WA Tool Testing Phase - Preliminary Findings

## Archive-It

a. Test crawl four git repositories using the standard crawler and Brozzler

b. Found issues with capturing GitLab (rendering issues) and Bitbucket (only got a white page!)

c. Promising results with standard crawler on GitHub

d. Limitations: patch crawling issues, time needed, and "uncharted territory" using this method
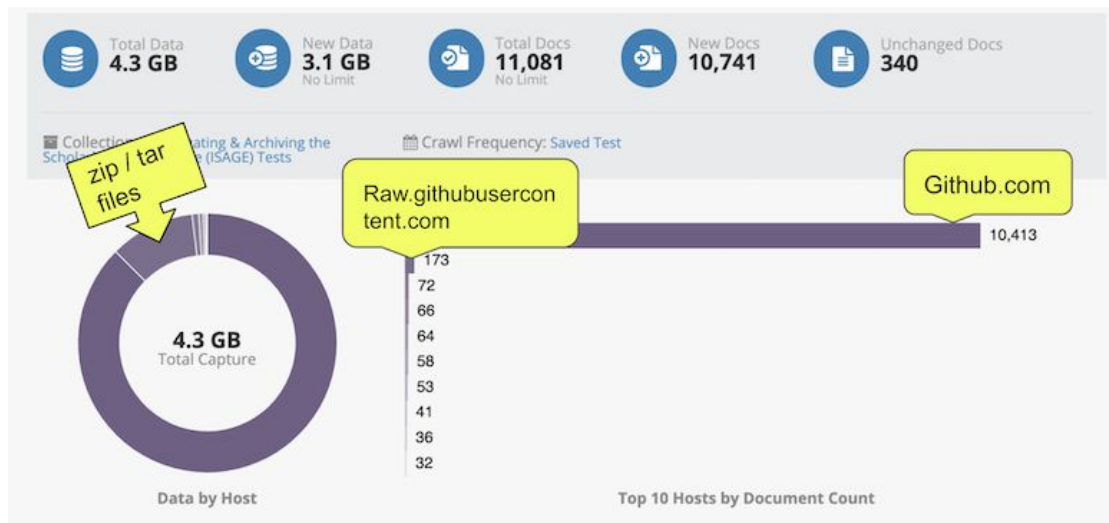


Read more: Lab Notes detailing Archive-it testing: investigating-archiving-git.gitlab.io/updates/lab-notes-archive-it

#IASGE

# WA Tool Testing Phase - Preliminary Findings

What was captured:

- Zip of source code with past versions
- 23 open issues (at time of capture) and 111 closed issues, and their labels.
- PRs and messages were also captured, except for indiv. Commit messages
- 4 pages of commits were saved, but older commit messages were missing



Screenshot of Archive-It crawl report using Standard crawler for single GitHub repository link from: investigating-archiving-git.gitlab.io/updates/lab-notes-archive-it/
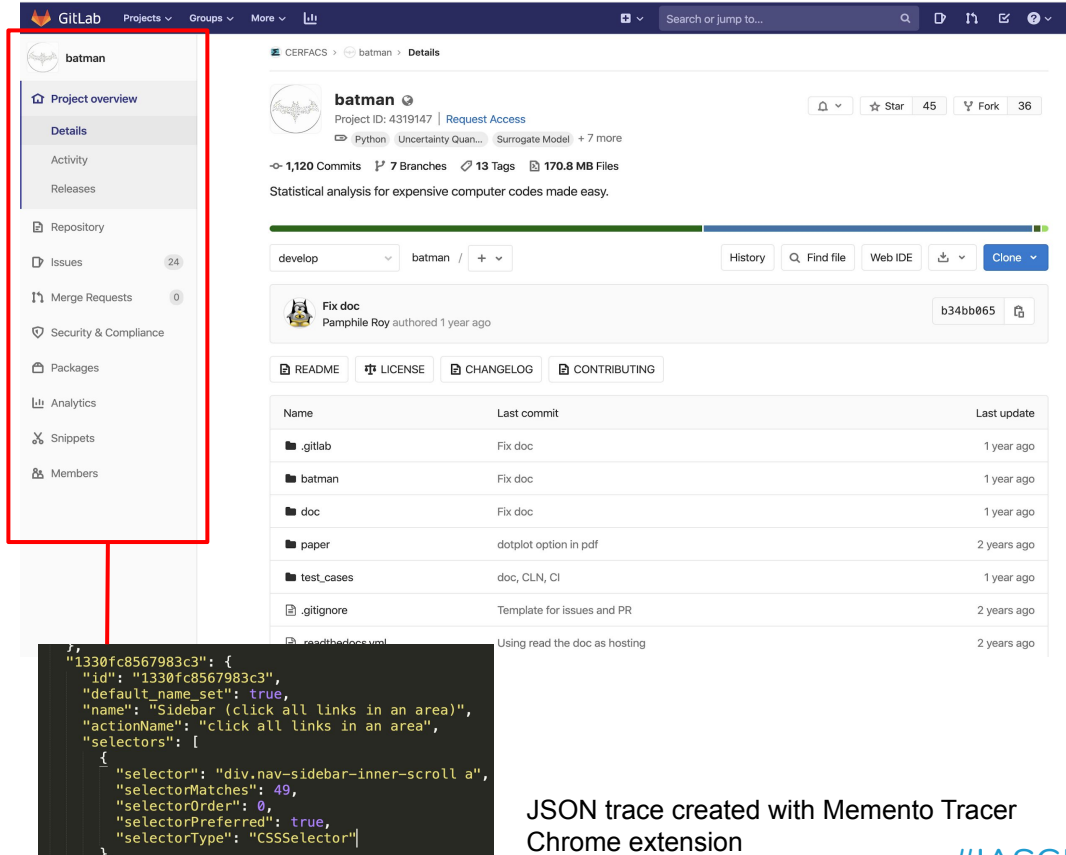
#IASGE

# What's next? Further testing!

- ## Memento Tracer
  - JSON traces serve as "instructions" for the crawler, which can be shared and reused
  - JSON file can be applied to any repositories on that domain (i.e. any repository on GitHub).
  - http://tracer.mementoweb.org/

- ## Webrecorder
  - Uses user interactions (clicking, scrolling) to create high quality captures
  - Limited scalability
  - Potential to create an Autopilot behaviors for GitHub? feature request ;-)
  - https://webrecorder.io/



JSON trace created with Memento Tracer Chrome extension

#IASGE

# Summary, calls to action, & all our contact info!

- Code/software and the contextual ephemera are worth saving
- Scholarship in Git format & in Git hosting platforms is at risk because there is no preservation plan
- Understanding behaviour patterns of academics using GHPs will help in T&L and archiving work
- Be on the lookout for research-centered blog posts
- Send us feedback on posts and/or resources you think we should know about!

Project website:
https://investigating-archiving-git.gitlab.io

GitLab repo:
gitlab.com/investigating-archiving-git

Emails:
vicky.steeves@nyu.edu
genevieve.milliken@nyu.edu
sarahtnguyen@nyu.edu

Twitters:
@VickySteeves
@gen_milliken
@snewyuen

#IASGE

# Survey - Preliminary Findings

Summary

Q33 - Do you use git hosting platforms as a storage place to backup your code?
- % said Yes

Q37 - How is your research or scholarship funded? Check all that apply.
- Mostly public

Q42 - How often do you use git to collaborate on authoring code?
- Mostly daily

Q43 - When a new collaborator joins your team, is there an onboarding process or protocol specifically for introducing them to your coding practices and use of version control?
- Mostly No

Q47 - Do you copy your repositories to external long-term storage services or platforms (e.g. Zenodo, OSF, institutional repository)?
- Mostly No

#IASGE