



Identifying Billions of Source Code Artifacts: the SWHID in Publication Workflows

Miguel Colom
Centre Borelli
ENS Paris-Saclay



Morane Gruenpeter
Software Heritage
Inria



Co-funded by
the European Union



The importance of referencing source code

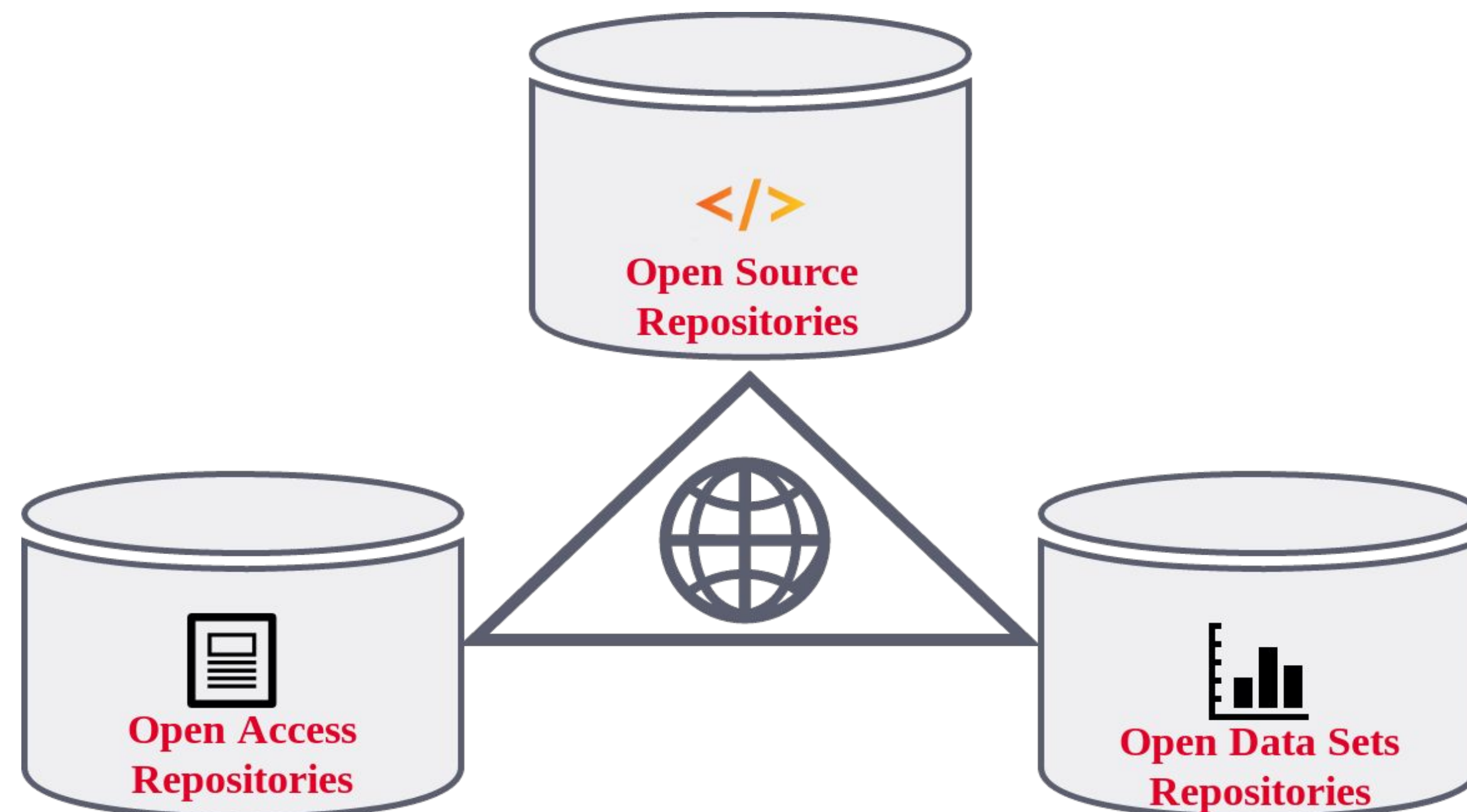
Open science: open access, open data, open source

The three **pillars** of **open science**:

- → **Open source repositories**
- Open datasets repositories
- Open access repositories

Software has **multiple facets**:

- a **tool**
- a research **outcome** or result
- the **object** of research



*Three pillars of Open Science
Software Heritage CC-BY 4.0 2019*

The importance of archiving source code permanently

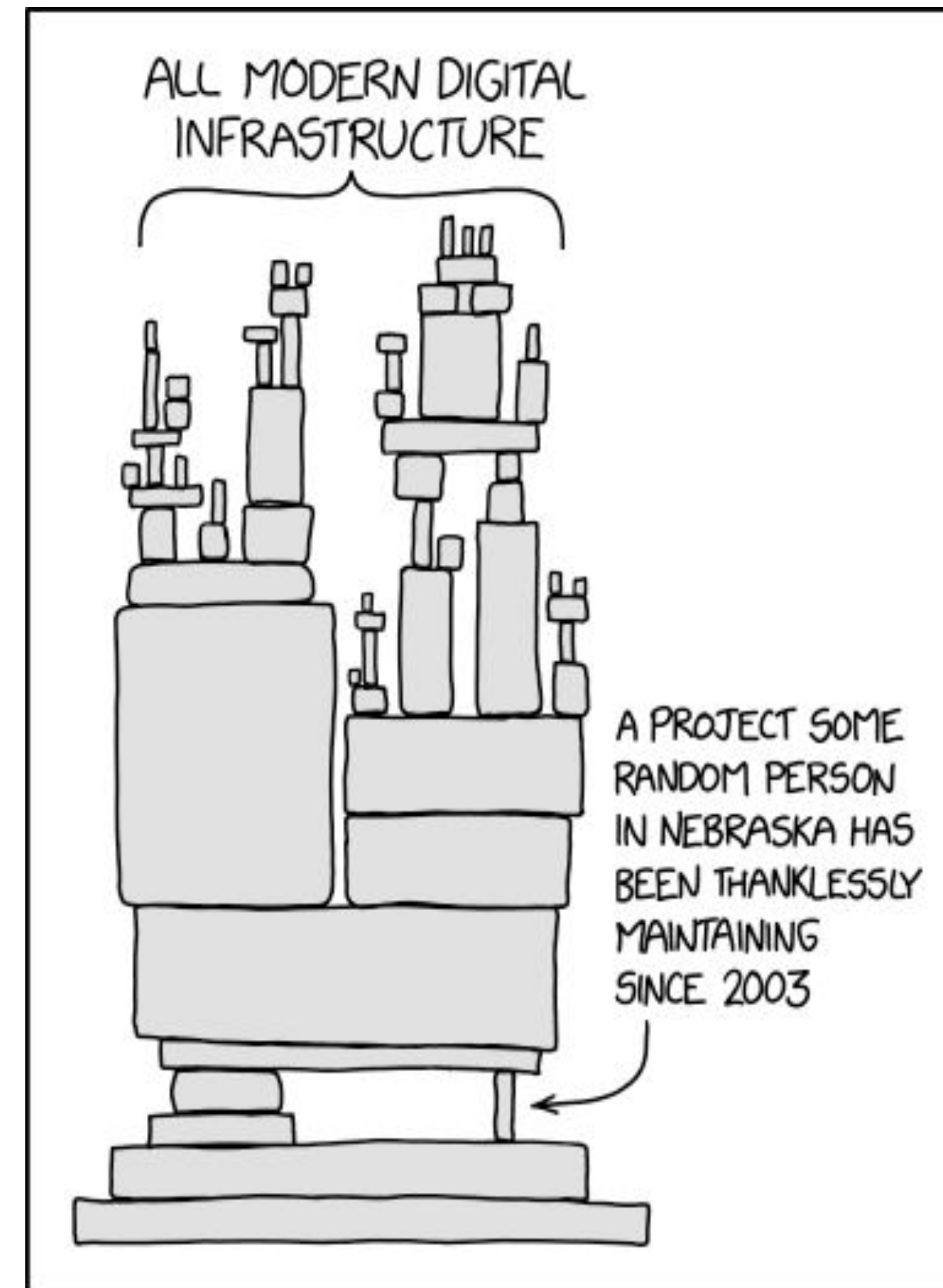
Software evolves over time

- projects may last decades
- the **development history** is key to its understanding

Contains part of the **knowledge of humanity!**

Complexity

- **millions of lines** of code
- large web of **dependencies**
 - easy to **break**, difficult to **maintain**
- sophisticated developer communities



<https://xkcd.com/2347/>

The Software Heritage initiative

One common infrastructure
to **collect, preserve and share** *all source code*



Bitbucket 2,036,916 origins	git 7,300 origins	R 21,620 origins
debian 129,507 origins	Guix 6,463 origins	GitHub 156,095,789 origins
GitLab 3,990,594 origins	launchpad 368,181 origins	GNU 354 origins
heptapod 1,098 origins	npm 1,799,296 origins	Maven™ 93,710 origins
NixOS 12,466 origins	python Package Index 432,092 origins	SOURCEFORGE 308,965 origins
Phabricator 184 origins		





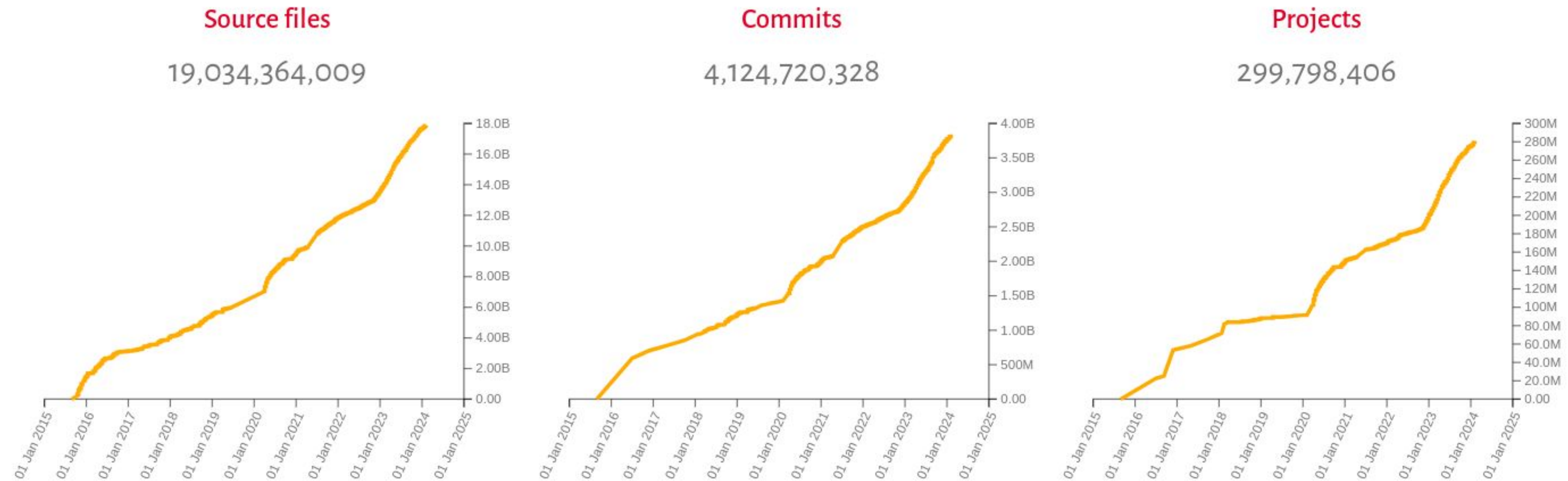
Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Are you familiar with Software Heritage?

Size

As of today the archive already contains and keeps safe for you the following amount of objects:



An international and non-profit initiative launched in 2016



built for the long term



Sharing the vision



Donors, members, sponsors



<http://www.softwareheritage.org/support/testimonials>

Save code feature

<https://archive.softwareheritage.org/save/>

Saving ~~your~~ any code now!

Software Heritage Archive

Save code now

Enter a SWHID to resolve or keyword(s) to search for in origin URLs

Features

- Search
- Downloads
- Save code now
- Help

You can contribute to extend the content of the Software Heritage archive by submitting an origin save request. To do so, fill the required info in the form below:

Origin type	Origin url
git 1	2

Submit 3

Help

Browse save requests

A "Save code now" request takes the form

Advantages

- All dev history is also saved
- URLs from **different platforms** are accepted
- PID to **reference** specific pieces of code

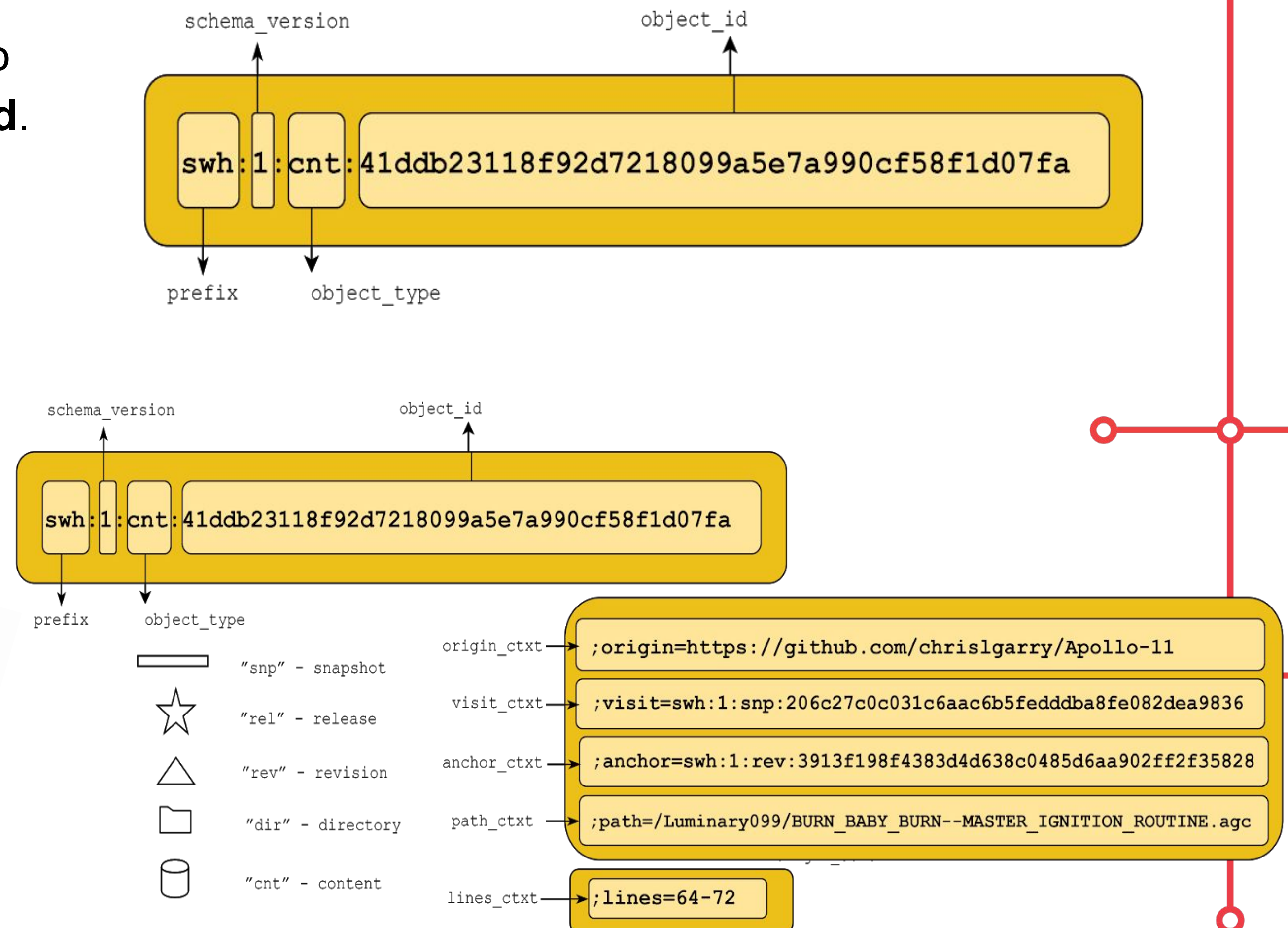
Meet the SWHID - The Software Hash Identifier

SWHIDs are **intrinsic identifiers** which are intimately **bound** to the **designated object**, they do **not need a register**, only agreement on a **standard**.

- **Intrinsic**: compute a unique **digital fingerprint**
- **decentralised**: do not need a registry, only agreement on a standard
- **cryptographically strong** identifiers

[Intrinsic vs. extrinsic blog post](#)

Go to [API endpoint](#)



Choose a SWHID on Software Heritage

Full width Home Development Documentation Donate Operational login

Software Heritage Archive

Browse the archive

Enter a SWHID to resolve or keyword(s) to search for in origin URLs

https://github.com/CGAL/cgal

23 October 2020, 13:31 UTC

Code Branch

Revision: b86a5

Permalinks

Choose a - `directory`

Select below a type

directory revision snapshot

archived repository archived swh:1:dir:abc0e2cbbfdfee8de52f0842263fbadf65f5b211

```
swh:1:dir:abc0e2cbbfdfee8de52f0842263fbadf65f5b211;
origin=https://github.com/CGAL/cgal;
visit=swh:1:snp:78e145aa8174e576786284475a76cf6f187b3475;
anchor=swh:1:rev:b86a5018c7f5f733c80fe40eee65803c112f2685;
path=/Hyperbolic_triangulation_2/
```

Add contextual information

Copy identifier Copy permalink

1.6 KB

Add context to SWHID

Copy identifier

The need of a standard PID to reference and fully describe SW

Why do we *need* a **standard**?

- In **Open Science** there are plenty of different journals based on **different formats**
- To **define an agreed upon format**, we need the **standardisation** of mechanisms to **reference** software ⇒ **SWHID** is **recommended** in the [EOSC SRIA](#) and in the [EOSC SIRS](#) reports
- A **link** to Github or another **forge** is **not enough** to ensure **access** and **interoperability**

→ An open process has been put in place in order to produce a publicly available specification

The SWHID Working Group

Working group and governance

- In **2023** the WG has been **established**
<https://swhid.org/>
- **Governance:**
<https://www.swhid.org/governance/>
- **Specification completed:** v1.0
(21/06/2023), v1.1 (6/11/2023)

Current status of SWHID

- 30+B SWHIDs in the archive
- Mention in Linux Foundation's [SPDX 2.2](#)
- IANA registered
- [WikiData P6138](#)



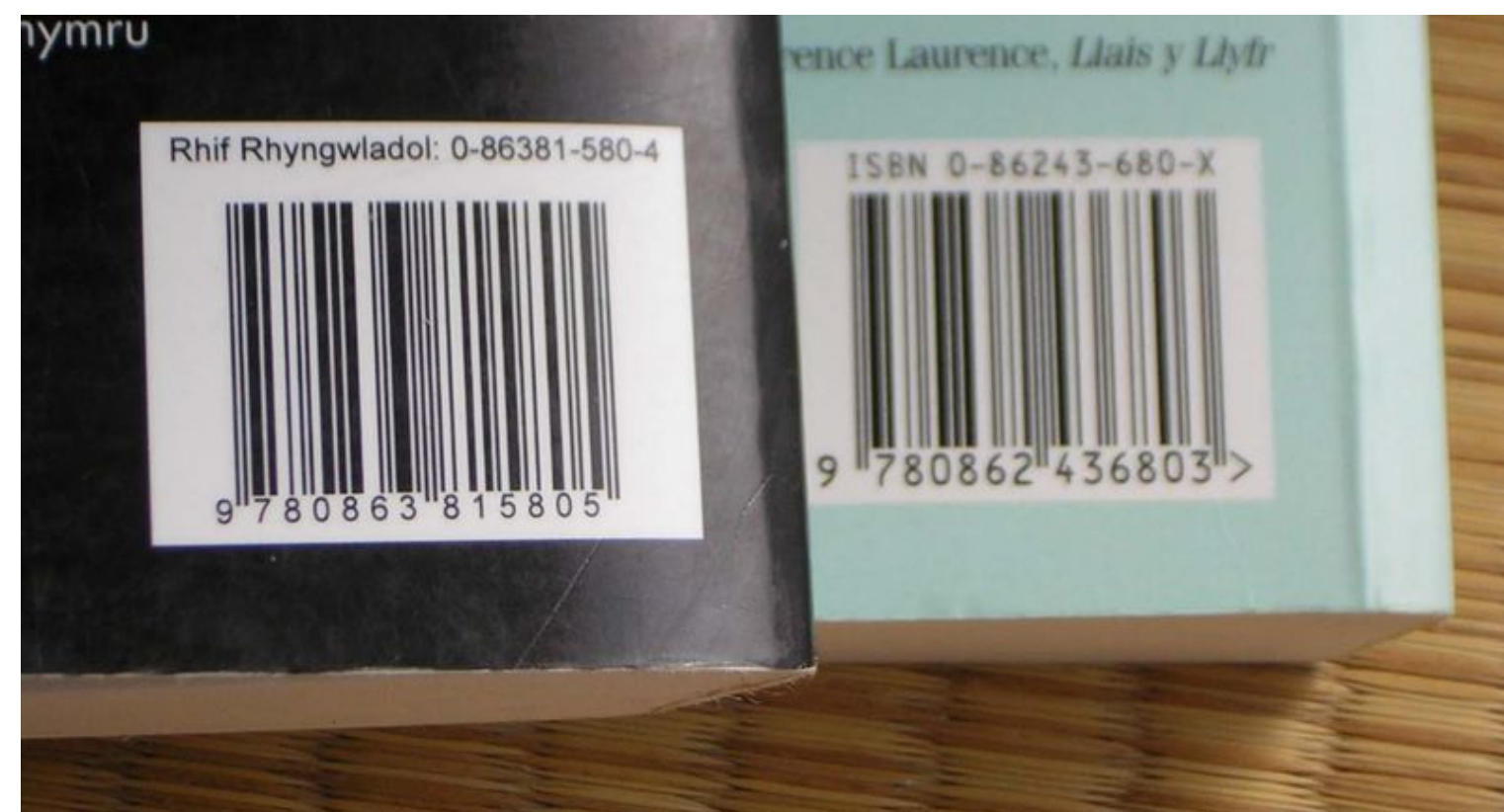
Go to [API endpoint](#)

Identification vs. location

Identification of a book

Goal: attribute the same code to the same type of objects. *What is it?*

- One ISBN number per published book
- ISO 2108 Standard specification



ISBN.JPG,
<https://commons.wikimedia.org/w/index.php?title=File:ISBN.JPG&oldid=456239269> (last visited January 26, 2024).

Location of a copy of a book

Goal: find (a copy of) a book. *Where is it?*

- Many locations (locations can change!)
- Many approaches



https://commons.wikimedia.org/wiki/File:Art_Books_on_Library_Shelf.JPG

→ We are mainly interested in **identification** and not that much with *location*

Extrinsic vs. Intrinsic identifiers

Main difference:

- how the **relation** between **identifier** and designated **object** is **created** and **maintained**.

Persistence is a key **desired property**.

	Extrinsic	Intrinsic
Relation	Register	Convention - agreed protocol/naming
Pre-internet	Passport number, ISBN, SSN, etc.	Music / Chemistry notation
Internet era	DOI, Handle, ARK, etc.	Cryptographic hashes, e.g git, bitcoin, SWHID

[Intrinsic vs. extrinsic blog post](#)

Which part of the software do we want to *identify*? (granularity)

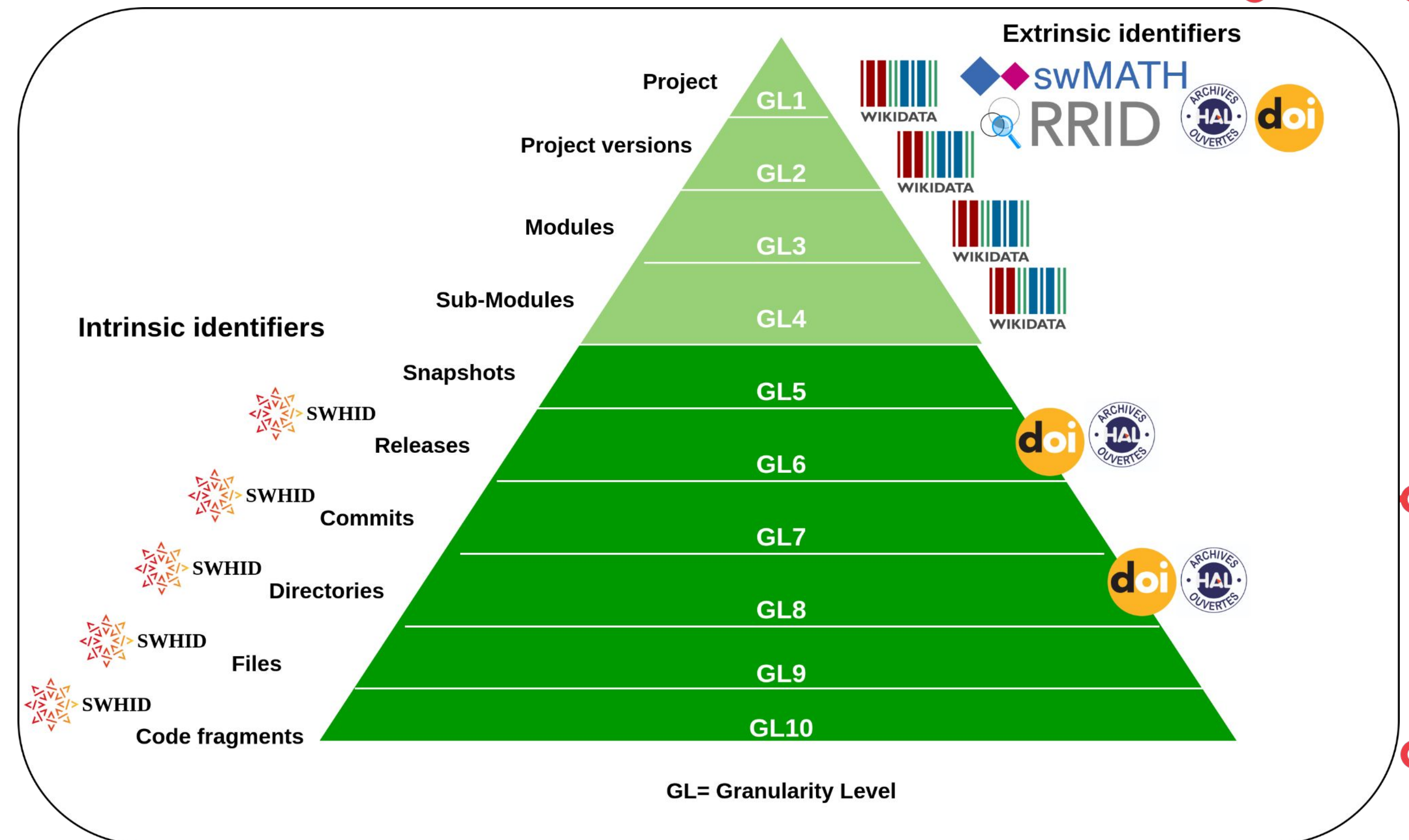
Software concept / project / collection

Description in registry, a homepage or any other form of metadata record

- Project versions (for example Python2 and Python3)
- Modules
- Sub-modules

Software artifact

- Executable (download link)
- Software source code
 - ◆ Dynamic artifact - current development code
 - ◆ Archived copy
 - Snapshot (all branches, all dev history)
 - Release / Package
 - Commit- a specific point in development history
 - Directory
 - File
 - Algorithm

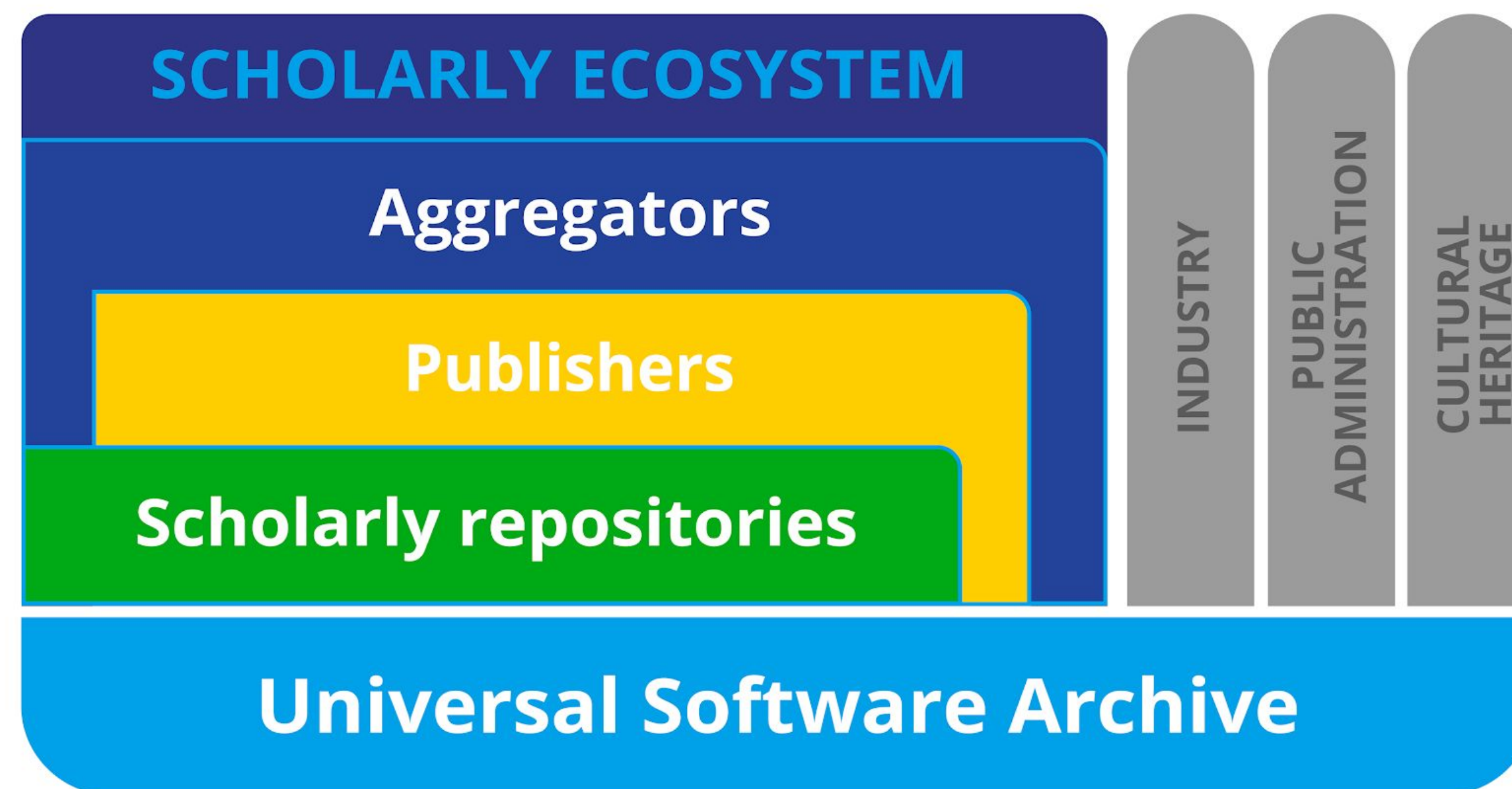


Research Data Alliance/FORCE11 Software Source Code Identification WG et al. (2020). Use cases and identifier schemes for persistent software source code identification (V1.1). *Research Data Alliance*. <https://doi.org/10.15497/RDA00053>

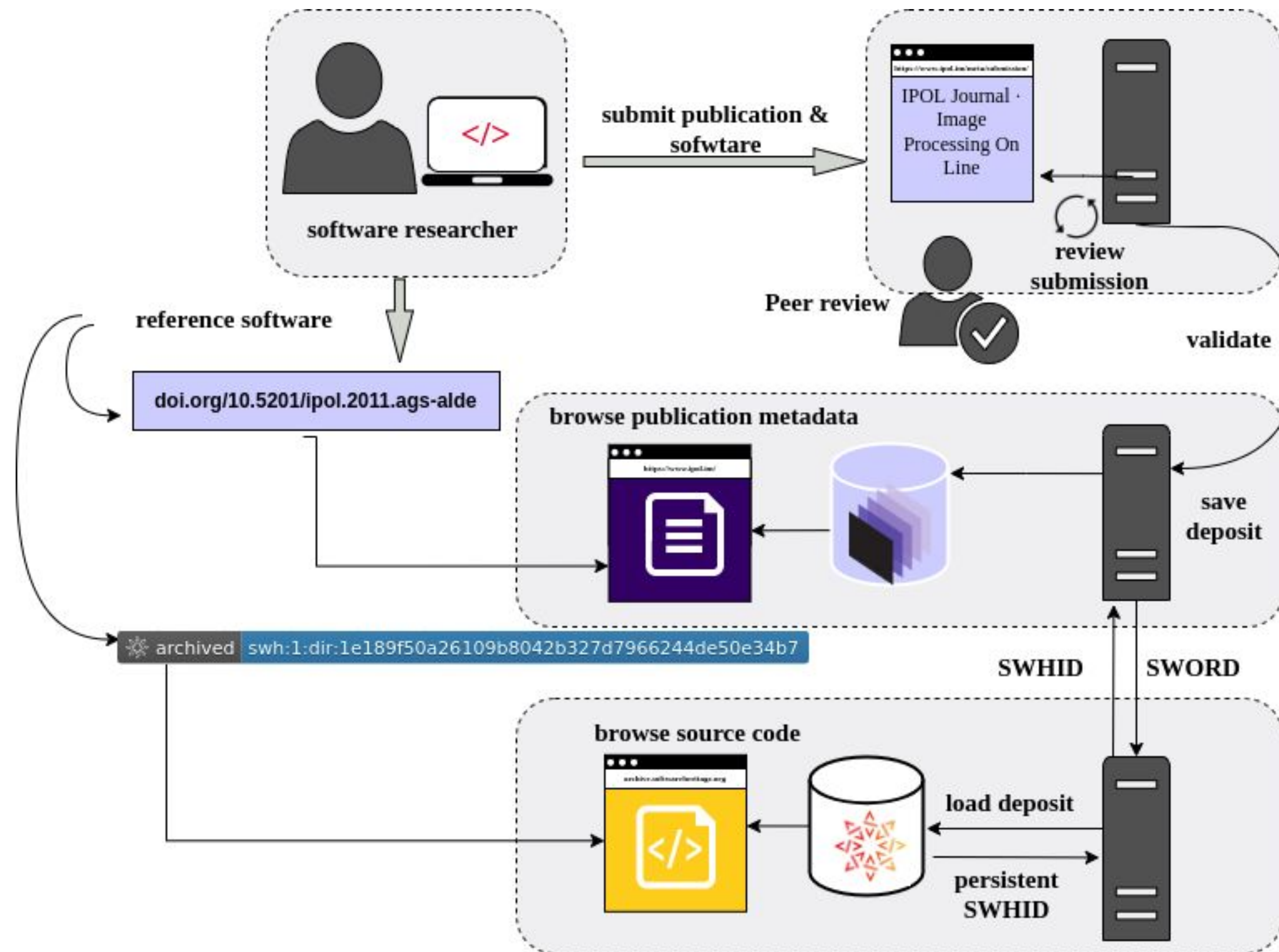
Use cases - connecting an ecosystem with SWHIDs

SIRS report: European Commission,
Directorate-General for Research and Innovation,
*Scholarly infrastructures for research software :
report from the EOSC Executive Board Working
Group (WG) Architecture Task Force (TF) SIRS,*
Publications Office, 2020,
<https://data.europa.eu/doi/10.2777/28598>

Video: [EOSC Software Infrastructures for
Research Software: J. B. Gonzalez Lopez
\(CERN\)](#)



IPOL publishing use case



IPOL journal publishing use case

- **Journal** Image Processing On Line (IPOL, <https://www.ipol.im/>)
- **Research software** packages are **identified** with:
 - The article **DOI**: (<https://doi.org/10.5201/ipol.2021.286>)
 - The **software SWHID**: the publisher deposits the software in Software Heritage with the DOI as an origin (<https://archive.softwareheritage.org/swh:1:dir:2cb75d8c95eb61d047d89428d0ec40a2286c0311;origin=https://doi.org/10.5201/ipol.2021.286;visit=swh:1:snp:23a5f7ee209b593e9b3e60ebe2bc42f1e6b76ff3;anchor=swh:1:rel:2de235c8fc3dd527cfaaba5cbf1d8144fee14f40>)
- **Links from the paper and metadata DOI to**:
 - the **software deposit** and its **SWHID**,
 - the **live demo** of the **software** (in the *demo* tab)

IPOL Journal · Image Processing On Line
HOME · ABOUT · ARTICLES · PREPRINTS · WORKSHOPS · NEWS · SEARCH

Image Inpainting using Patch Consensus and DCT Priors

Ignacio Ramírez Paulino, Ignacio Hounie

article demo archive

published • 2021-01-09
reference • IGNACIO RAMÍREZ PAULINO, AND IGNACIO HOUNIE, *Image Inpainting using Patch Consensus and DCT Priors*, Image Processing On Line, 11 (2021), pp. 1–17. <https://doi.org/10.5201/ipol.2021.286>

BibTeX info

Communicated by Pablo Arias
Demo edited by Pablo Arias

Abstract

We present an implementation of the PACO-DCT inpainting algorithm. This method is based on maximizing the likelihood of image patches in terms of their DCT coefficients, while requiring consensus on the overlapping patches. The resulting problem is solved as an instance of the PACO framework.

Download

- full text manuscript: PDF low-res. (577.7kB) PDF (6.6MB) [?]
- source code: ZIP SWHID info </> </> Software Heritage Archive

```
@softwareversion{sw-ipol.2021.286,  
  title = {{Image Inpainting using Patch Consensus and DCT Priors}},  
  author = {Ignacio Ramírez Paulino, Ignacio Hounie},  
  date = {2021-01-01},  
  license = {GPL-3.0-or-later},  
  version = {1.0},  
  swhid =  
{swh:1:dir:2cb75d8c95eb61d047d89428d0ec40a2286c0311;origin=https://doi.org/10.5201/ipol.2021.286;vis
```

Copy to clipboard

Preview

Loading takes a few seconds. Images and graphics are degraded here for faster rendering. See the downloadable PDF documents for original high-quality versions.

HAL (CCSD) x SWH = combined solution for researchers

The **institutional** Research Software **deposit** in the French National Archive - **HAL**

- ★ Two options for researchers:
 - 🚀 Copy-pasting a PID
 - 📦 Adding files
- ★ Moderation of the metadata

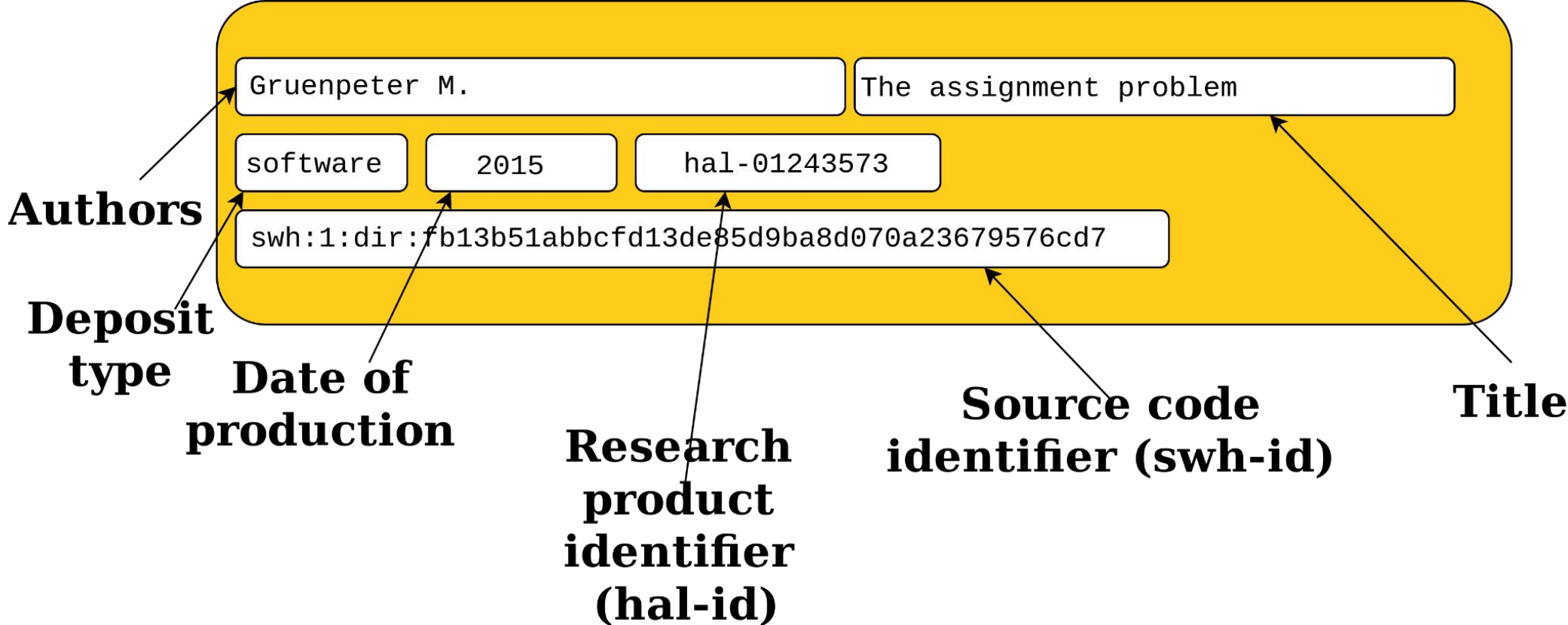
The screenshot displays the HAL interface for the LinBox software. On the left, the HAL page shows metadata for 'hal-02130801, version 1' (11-06-2019), including identifiers, metadata (version 1.6.3, GNU Lesser General Public License v2.1), and a citation. An orange arrow points from the HAL page to the Software Heritage interface on the right. The Software Heritage page shows the source code for 'config-blas.h' with a commit hash 'e8e18328952266b7875c692963b11963b1496107'. A red box highlights the commit hash and the file path '393b611 / linbox-1.6.3 / linbox / config-blas.h'. Below this, a blue box contains the SWHID: `swh:1:dir:393b611a1424f032e83569bf6762502371cfcf65`. At the bottom right, a blue box contains the URL: `https://hal.science/hal-02130801`.

Cite

The Linbox Group. LinBox. 2019, (swh:1:dir:393b611a1424f032e83569bf6762502371cfcf65;origin=https://hal.archives-ouvertes.fr/hal-02130801;visit=swh:1:snp:19c29b988fe02623c70c7dc8bc97c42481eb691b;anchor=swh:1:rev:e8e18328952266b7875c692963b11963b1496107;path=/). (hal-02130801)

Open Science series: the software source code deposit

For which use case: *reference vs. citation*



Archive & Index

- metadata record (extrinsic)
- artifact itself (intrinsic)

Credit & Attribution

- a metadata record
- all authors & contributors

Reuse & Reproducibility

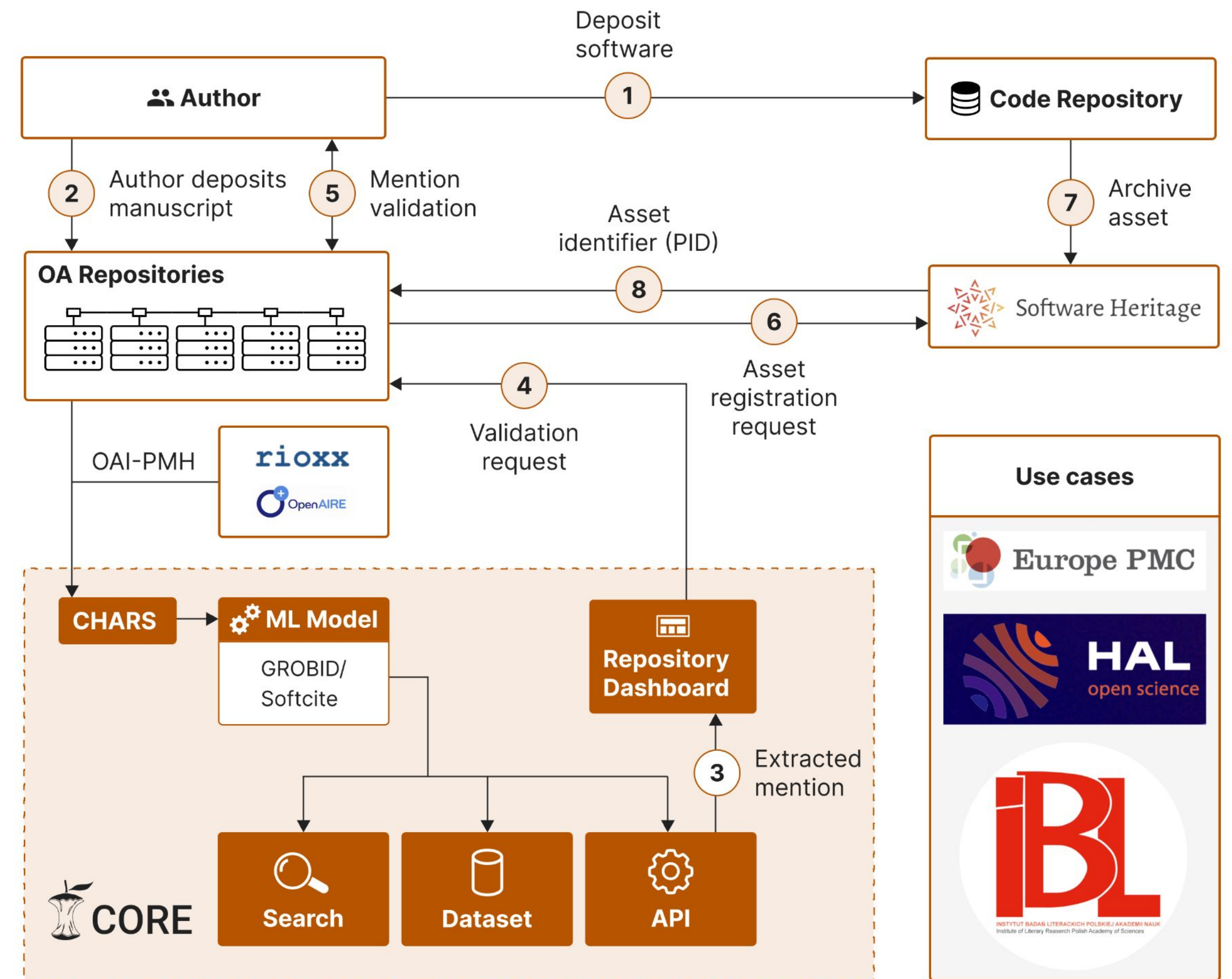
- a specific artifact
- appropriate metadata/docs

Building bridges between infrastructures

SoFAIR project, launched in January 2024

- **Goal:** establish a machine-assisted workflow embedded into widely used open scholarly infrastructures to **assist researchers** in
 - **identifying,**
 - **describing,**
 - **registering,**
 - **linking,**
 - **and archiving research software.**
- Funded by **CHIST-ERA**
- <https://sofair.org/>

Stay tuned...



Conclusions (1 / 2)

- **SW is valuable** by itself
 - Not just a *supplementary material* of a scientific article, for example
- **SW** needs to be **preserved** for the long term
 - It's part of the knowledge heritage of humanity
- Preservation of SW and proper description ⇒ allow **reproducible research**
- **Good news**: **SWH** is **preserving all** publicly available **source codes** :-)

Conclusions (2 / 2)

- **SW** needs to be **properly referenced**
 - A DOI is not enough. We need a **complete PID**
 - We propose the **SWHID** to point to **permanently-stored source code** and properly **describe** it
- Our goal is that the **SWHID** is a **general standard to all kind of software** for:
 - **Preservation**
 - **Proper citation**
 - **Interoperability**

Thank you for your attention

<https://www.softwareheritage.org/>



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE



FAIRCORE4EOSC
Core Components Supporting a FAIR EOSC



Google begins shutdown of its code repository

Latest News Published: March 12th, 2015 - Michael Pehel

After nine years, Google's open-source code repository, Google Code, started closing shop today by disabling new projects and announcing the permanent shut down of the service by January 26, 2016.

Google Code started as Google's answer to SourceForge, the predominant code repository back in 2006. The reliability of SourceForge was brought into question that year when SourceForge.net's database was hacked and user data was compromised. Problems continued the following summer in 2007 with a temporary service outage in August. The appeal of a repository backed by Google only increased as a result.



Latest Deep tech Sustainability Ecosystems Data and security Fintech and ecommerce Future of work More

This article was published on March 3, 2015

INSIDER

Code collaboration platform GitLab acquires rival Gitorious, will shut it down on June 1

March 3, 2015 - 4:11 pm



Gabriel Altay
@gabrielaltay

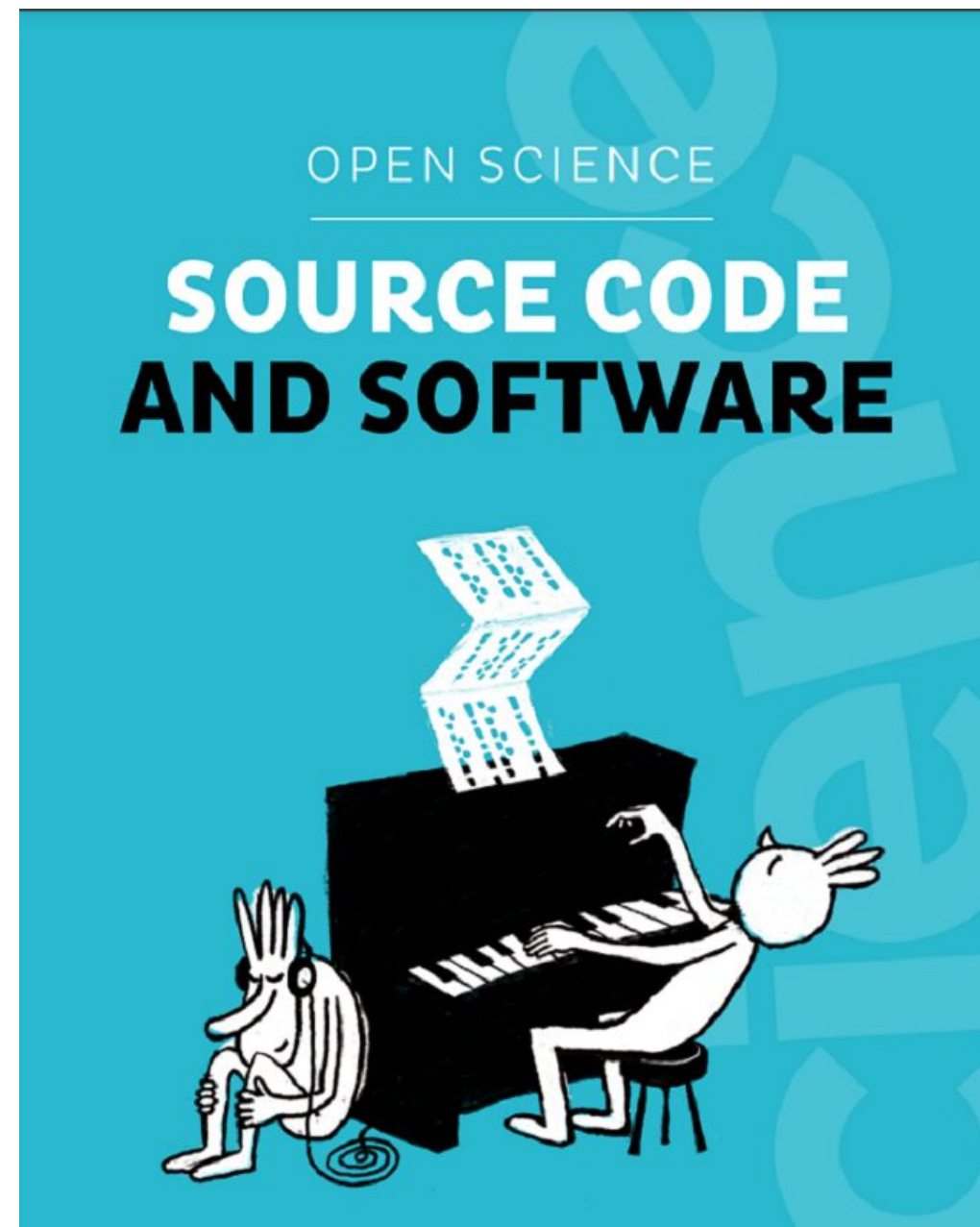
Just realized [@Bitbucket](#) disabled all mercurial repositories when the [@asclnet](#) informed me that a link associated with an old paper of mine was down. Thought all was lost, but someone archived all the repos! very classy move by [@octobus_net](#) and [@SWHeritage](#).

[Traduire le Tweet](#)

1:48 AM · 31 août 2020 · Twitter Web App

Adoption in Academia: a few indicators

Policy



- [Funding agencies recommendations ANR 2023 guidelines \(p. 17\)](#)
- HAL+SWH in [the Open Science software booklet](#) from the French Ministry of Higher Education and Research

Citation & ACM

- November 2017: first meeting with [ACM Digital Library](#)
- May 2020 [BibLaTeX](#) released
- April 2022: inclusion in [ACMART](#) Class for typesetting publications of ACM



Replicability stamp

the [Computer Graphic Replicability Stamp Initiative \(GRSI\)](#)



The first free software distribution backed by a stable archive - making [reproducible research papers possible](#)