

# The Software Pillar of Open Science

policy, needs, and how to address them

Roberto Di Cosmo

Director, Software Heritage  
Inria and Université de Paris Cité

May 2022



# Software Heritage

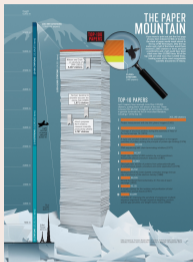
THE GREAT LIBRARY OF SOURCE CODE

- 1 Software and Open Science
- 2 Policy framework and growing needs
- 3 Can you address these needs?
- 4 Yes you can!
- 5 Call to action



# Software is a pillar of Open Science

## Software powers modern research



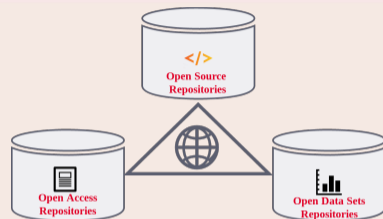
*[...] software [...] essential in their fields.*

*Top 100 papers (Nature, 2014)*

*Sometimes, if you don't have the software, you don't have the data*

*Christine Borgman, Paris, 2018*

## A key pillar: software (source code)



The links in the picture are **important**

## Nota Bene

software may be a *tool*, a *research outcome* and a *research object*

access to the *source code* is essential!

Preserving (the history of) source code is necessary for *reproducibility*

# Software Source Code is Precious Knowledge

Harold Abelson, Structure and Interpretation of Computer Programs (1st ed.)

1985

*“Programs must be written for people to read, and only incidentally for machines to execute.”*

## Apollo 11 source code (excerpt)

```
P63SP0T3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND   CHAN33
              EXTEND
              BZF    P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF    CODE500      # ASTRONAUT:  PLEASE CRANK THE
              TC     BANKCALL      #              SILLY THING AROUND
              CADR   GOPERF1
              TCF    GOTOP00H      # TERMINATE
              TCF    P63SP0T3      # PROCEED    SEE IF HE'S LYING

P63SP0T4      TC     BANKCALL      # ENTER      INITIALIZE LANDING RADAR
              CADR   SETPOS1

              TC     POSTJUMP      # OFF TO SEE THE WIZARD ...
              CADR   BURNBABY
```

## Quake III source code (excerpt)

```
float Q_rsqrt( float number )
{
    long i;
    float x2, y;
    const float threehalfs = 1.5F;

    x2 = number * 0.5F;
    y = number;
    i = * ( long * ) &y; // evil floating point bit level hacking
    i = 0x5f3759df - ( i >> 1 ); // what the fuck?
    y = * ( float * ) &i;
    y = y * ( threehalfs - ( x2 * y * y ) ); // 1st iteration
    // y = y * ( threehalfs - ( x2 * y * y ) ); // 2nd iteration, this
    // can be removed

    return y;
}
```

Len Shustek, Computer History Museum

2006

*“Source code provides a view into the mind of the designer.”*

- 1 Software and Open Science
- 2 Policy framework and growing needs
- 3 Can you address these needs?
- 4 Yes you can!
- 5 Call to action



# The Paris Call on Software Source code (2019, UNESCO)

Experts call for greater recognition of software source code as heritage for sustainable development

6 November 2018



UNESCO, Inria, Software Heritage invite  
40 international experts to meet in Paris



The call is published on Feb 2019

*"[We call to] promote software development as a valuable research activity, and research software as a key enabler for Open Science/Open Research, sharing good practices and recognising in the careers of academics their contributions to high quality software development, in all their forms"*

<https://en.unesco.org/foss/paris-call-software-source-code>



## 2nd National Plan for Open Science (6/7/2021)

### Open and promote research software source code

- actions (selection)
  - charter for research software policy
  - recognize software development (see [announcement of the 2021 prize](#))
  - coordinate communities of practice
  - connected ecosystem of research outputs
- recommendations (selection)
  - archive in Software Heritage
  - standardise and use SWHID
  - build a national catalog of research software
  - leverage ADAC network

See [official announcement](#)



[Accueil](#) > [Recherche](#) > [Science ouverte](#)

Publié le 05.02.2022

## Sommaire

- [The Coq proof assistant](#) : lauréat de la catégorie Scientifique et technique
- [Scikit-learn](#) : lauréat de la catégorie Communauté
- [Faust](#) : lauréat de la catégorie Documentation
- [Gammapy](#) : prix du jury
- [Jury](#)

# Remise des prix science ouverte du logiciel libre de la recherche

**Le ministère de l'Enseignement supérieur, de la Recherche et de l'Innovation remet pour la première année les Prix science ouverte du logiciel libre de la recherche. Dix logiciels mis au point par des équipes françaises sont récompensés pour leur contribution à l'avancée de la connaissance scientifique.**



# A plurality of needs that we must address

## Researchers

- **archive** and **reference** software used in articles
- **find** useful software
- get **credit** for developed software
- verify, **reproduce**, improve results

## Laboratories/teams

- **track** software contributions
- produce reports
- maintain web page

## Research Organization

know its **software assets**

- technology **transfer**
- impact **metrics**
- funding **strategy**
- career **evaluation**

## Archive

Research software artifacts must be properly **archived**  
make sure we can *retrieve* them (*reproducibility*)

## Reference

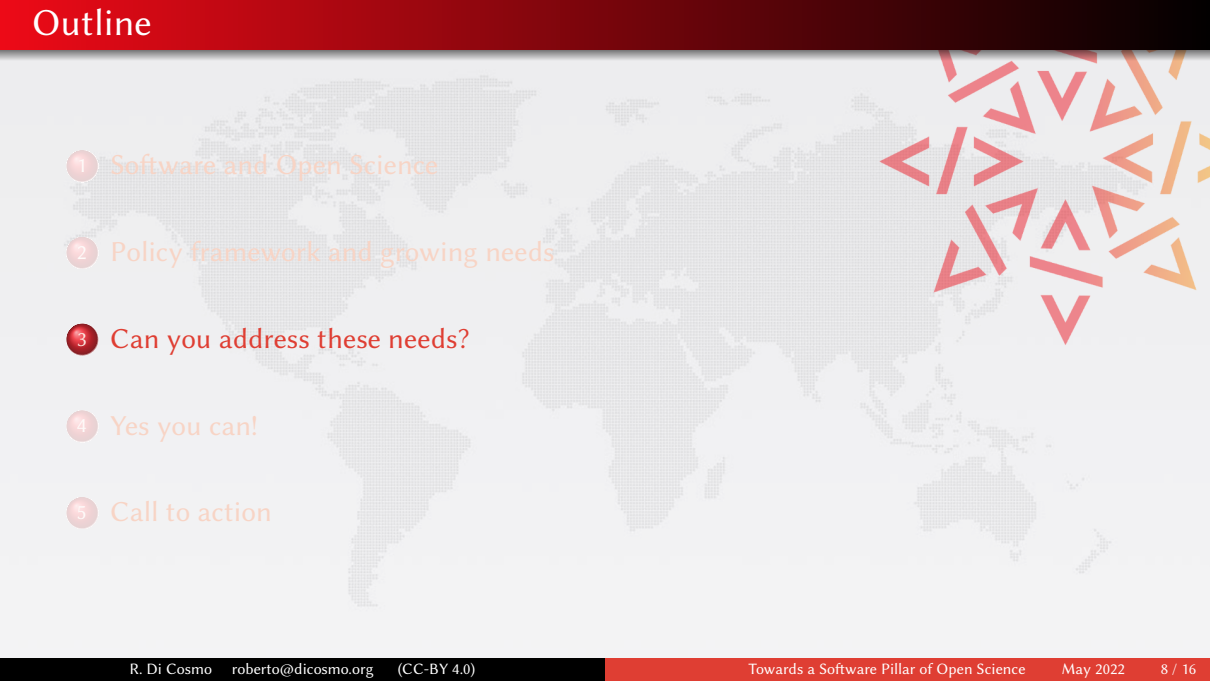
Research software artifacts must be properly **referenced**  
make sure we can *identify* them (*reproducibility*)

## Describe

Research software artifacts must be properly **described**  
make it easy to *discover* and *reuse* them (*visibility*)

## Cite/Credit

Research software artifacts must be properly **cited** (*not the same as referenced!*)  
to give *credit* to authors (*evaluation!*)

- 
- 1 Software and Open Science
  - 2 Policy framework and growing needs
  - 3 Can you address these needs?
  - 4 Yes you can!
  - 5 Call to action

# A word of warning: forges are *not* archives!

## 2015: the first big bad news

Google Code and Gitorious.org shutdown: ~1M endangered repositories

- broken links in the web of knowledge (my papers too)

## 2019: big bad news keep coming in

- summer 2019: BitBucket announces Mercurial VCS sunset
- july 2020: BitBucket erases 250.000+ repositories (including research software)

## 2021: ... in Academia too

- october 2021: Inria's old gforge is unplugged
  - **breaks the build chain** of the OCaml package manager (Opam)

## Bottomline

we need a universal archive of software source code: now we have one!



# Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Collect, preserve and share *all* software source code

Preserving our heritage, enabling better software and better science for all

## Reference catalog



**find** and **reference** all software source code

## Universal archive



**preserve** all software source code

## Research infrastructure



**enable analysis** of all software source code

# The largest software archive, a shared infrastructure

Cultural Heritage



Industry



Research



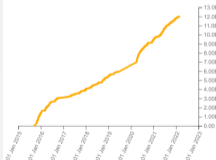
Public Administration



## Software Heritage

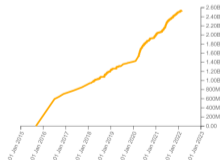
Source files

12,032,627,304



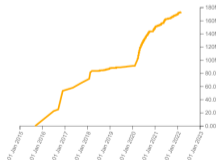
Commits

2,536,918,821



Projects

173,242,749



Directories

9,946,192,395

Authors

47,334,620

Releases

31,763,605

## Sharing the vision



United Nations  
Educational, Scientific and  
Cultural Organization



And many more ...

[www.softwareheritage.org/support/testimonials](http://www.softwareheritage.org/support/testimonials)

## Donors, members, sponsors

*Inria*

Diamond sponsor



Platinum sponsors



Gold sponsors

openinventionnetwork



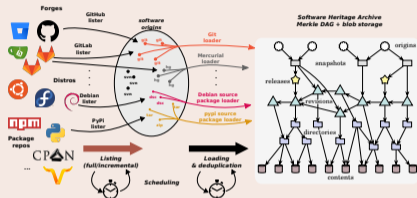
Silver sponsors



Bronze sponsors



## Archive (12B+ files, 170M+ projects)



- [save.softwareheritage.org](https://save.softwareheritage.org)
- [deposit.softwareheritage.org](https://deposit.softwareheritage.org)

## Describe

- *Intrinsic metadata* from source code
- Contributed the [Codemeta generator](#)

## Reference (20 billion SWHIDs)

Intrinsic, decentralised, cryptographically strong identifiers, SWHIDs



Now supported in [SPDX 2.2](#), [Wikidata](#) etc.

## Cite/Credit

- Contributed *software citation* style  
[biblatex-software](#), v 1.2-2 now on [CTAN](#)



# HAL and Software Heritage: building a curated software catalog

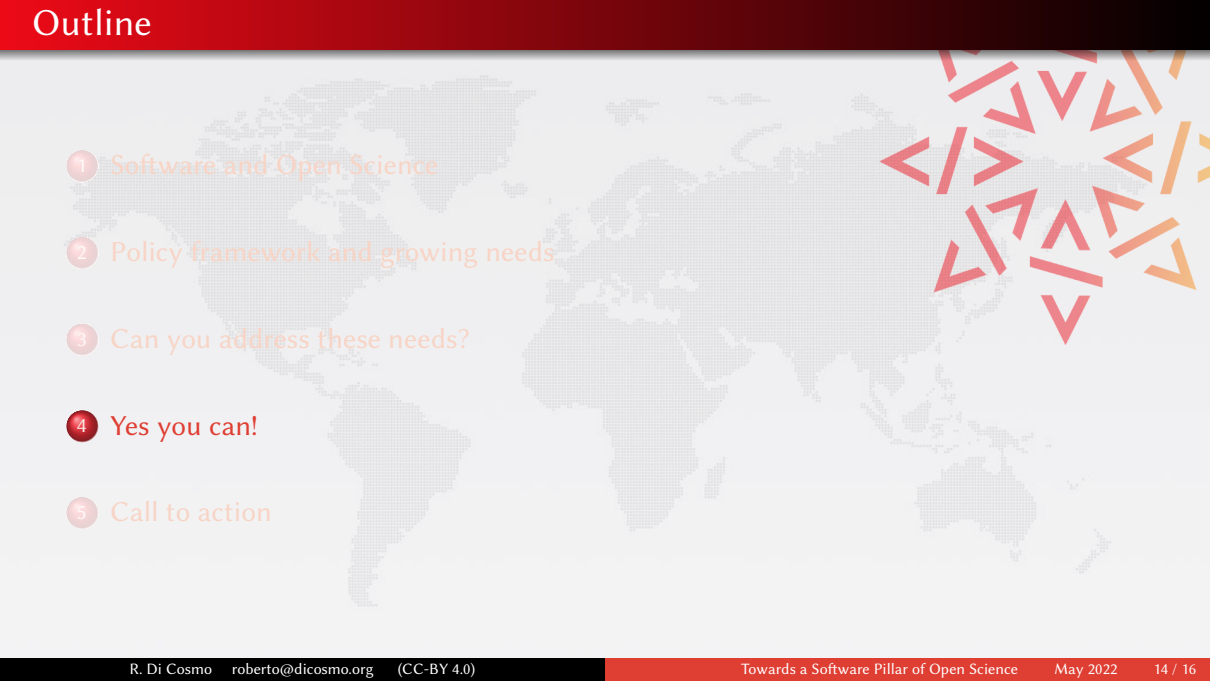
The diagram on the left illustrates the workflow for depositing software into HAL and Software Heritage:

- work on a YCS**: A researcher works on a system using tools like Git and Docker.
- submit repository url**: The researcher provides a repository URL to the SWHIB system.
- Save code now**: A button to save code to [softwareheritage.org](https://softwareheritage.org).
- digital archive**: The code is archived in a digital archive.
- validate**: The deposit is validated.
- publish**: The code is published to the public.
- cite & browse deposit url**: The deposit URL is cited and browsed.

The screenshots show the HAL website interface:

- HAL Archives-ouvertes.fr**: URL <https://hal.archives-ouvertes.fr/hal-02130801>. Slogan: "Free and accessible knowledge".
- LinBox**: Metadata for the LinBox library, including authors (Eric Grönroos, Ridaa Alges, etc.) and abstract.
- Metadata**: Details for version 1.6.3, including the GNU Lesser General Public License v2.1 or later, C++ programming language, and code repository URL.
- Export**: Options to export the metadata in various formats (Checklist, BibTeX, etc.).
- Permalink**: A unique URL for the deposit: <https://hal.archives-ouvertes.fr/hal-02130801>.
- Revision**: e8e18328952266b7875c692963b11965b1496107.
- Tip revision**: e8e18328952266b7875c692963b11965b1496107 authored by Software Heritage on 11 June 2019, 08:12 UTC.
- config-blas.h**: Content of the config file, including copyright information and license details.

Software Heritage URL: [swh:1:dir:393b611a1424f032e83569bf6762502371cfc65](https://swh:1:dir:393b611a1424f032e83569bf6762502371cfc65)

- 
- 1 Software and Open Science
  - 2 Policy framework and growing needs
  - 3 Can you address these needs?
  - 4 **Yes you can!**
  - 5 Call to action

# An example is worth a thousand words

- Browse [the archive](#) (your work [may be already there](#) !)
- [Trigger archival](#) of your preferred software in a breeze
- Get and use SWHIDs ([full specification available online](#))
- Cite software using the [biblatex-software](#) package from CTAN
- Example in a journal: [an article from IPOL](#)
- Example with Parmap: [devel on Github](#), [archive in SWH](#), [curated deposit in HAL](#)
- Extracting all the software products [for Inria](#), [for CNRS](#), [for LIRMM](#) or [for Rémi Gribonval](#) using HalTools
- [Curated deposit in SWH via HAL](#), see for example: [LinBox](#), [SLALOM](#), [Givaro](#), [NS2DDV](#), [SumGra](#), [Coq proof](#), ...
- Example use in a research article: compare Fig. 1 and conclusions
  - in [the 2012 version](#)
  - in [the updated version](#) using SWHIDs and Software Heritage
- Example use in a research article: extensive use of SWHIDs in [a replication experiment](#)

- 1 Software and Open Science
- 2 Policy framework and growing needs
- 3 Can you address these needs?
- 4 Yes you can!
- 5 Call to action



## Archiving and referencing

For **all source code** used in research (*yes, even small scripts!*)

- ensure it is archived in Software Heritage (see [save code now](#))
- get the proper **SWHID** for your software (see [detailed HOWTO](#))
- add it to research articles for reproducibility (see [detailed HOWTO](#))

## Describing and Citing/Crediting

For **software you want to put forward** (*mention in your CV, reports, etc., get citations and credit for it*), do the following **extra steps**:

- add **codemeta.json** with description (see the [codemeta generator](#))
- reference in the HAL portal (french partners, see [online HAL documentation](#))
- cite software using the [biblatex-software](#) package (in CTAN and TeXLive)

HAL+SWH let you address all the needs at once...

- *researcher, engineer*: archival, reference, credit, CV etc. *with a little effort from them*
- *labs, organizations*: track and report software production in a simple way
- *technology transfer offices*: view the software production
- *national level*: a *curated* catalog of the software production

... with a little effort from your side

- Update the Open Science policy to include software
- Train on the use of SWH and HAL for software
- Join the network of HAL moderators for software

it's a long road, but together we can make it

## Questions?