

# Building Infrastructures for Open Source Software

a short introduction

Roberto Di Cosmo  
UN Open Source Retreat 2024

Director, Software Heritage  
Inria and Université de Paris Cité

May 14th 2024



# Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Introduction
- 2 Selected actions of interest (France)
- 3 Meet Software Heritage
- 4 An example is worth a thousand words
- 5 A lot more lies ahead
- 6 Call to action



# Short Bio: Roberto Di Cosmo

Computer Science professor in Paris, now working at INRIA

- 35+ years of research (Theor. CS, Programming, Software Engineering, Erdos #: 3)
- 25+ years of Free and Open Source Software
- 15+ years building and directing structures for the common good



1999 *DemoLinux* – first live GNU/Linux distro

2007 *Free Software Thematic Group*

150 members 40 projects 200Me

2008 *Mancoosi project* [www.mancoosi.org](http://www.mancoosi.org)

2010 *IRILL* [www.irill.org](http://www.irill.org)

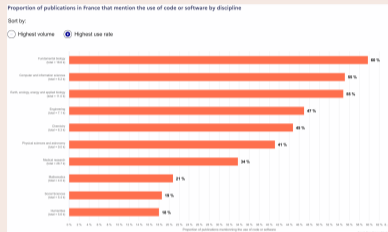
2015 *Software Heritage* at INRIA

2018 *National Committee for Open Science*, France

2021 *EOSC Task Force on Infrastructures for Software*,  
European Union

# Software is a pillar of Open Science

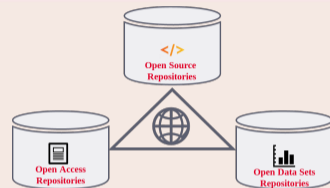
## Software powers modern research



Over 20% of articles using software across all disciplines share it

2024 French Open Science Monitor

## Key pillar: software



Links are **important**

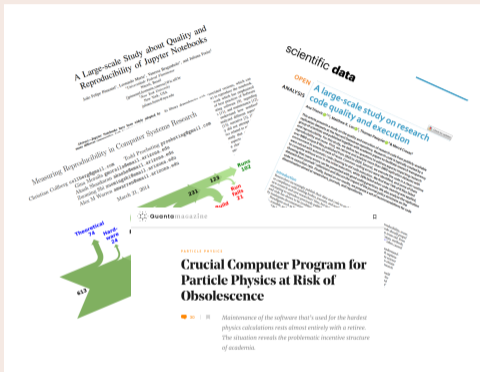
## Nota Bene

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

access to the *source code* is essential!

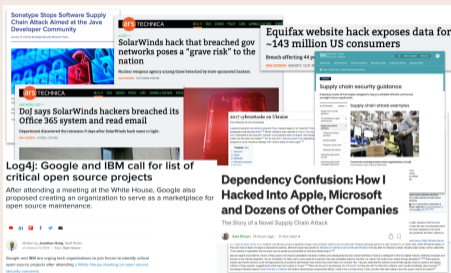
# How are we managing our (open source) software ?

## Reproducibility, maintenance in Academia



(articles: [here](#), [here](#), [here](#) and [here](#))

## Security, integrity, traceability in Industry



Can they track the software that they

- ship, use, acquire
- has that bug or vulnerability

awareness is raising at the level of public policy

- 1 Introduction
- 2 Selected actions of interest (France)
- 3 Meet Software Heritage
- 4 An example is worth a thousand words
- 5 A lot more lies ahead
- 6 Call to action



# French National plan for Open Science, 2021-2024



## SECOND FRENCH PLAN FOR OPEN SCIENCE

Generalising open science in France 2021-2024



1

## Second French Plan for Open Science



Launch on 6 July 2021 by Frédérique Vidal, Minister for Higher Education, Research and Innovation

- Multiplying the **levers for change** in order to **generalise open science practices**
- Structuring the **policy for opening up or sharing research data**
- New commitments to the **opening of source code** produced by research
- **European and international inclusion** in the context of the French Presidency of the European Union
- **Disciplinary and thematic variations**: open science policies must be adapted to disciplinary specificities

2

## Path Three : Opening up and promoting source code produced by research

7

Recognize and support the dissemination under an open source licence of software produced by publicly funded research programmes

« The opening of software source code is a major challenge for the **reproducibility** of scientific results. »

8

Highlight the production of source code from higher education, research and innovation

« Distribution of software products under **open source licence** will be preferred. »

9

Define and promote an open source software policy

3

### Define and promote an open source software policy

- Produce a **National Charter for Open Source Software** coming from higher education, research and innovation
- Develop the **link between data and software** through a network of **Chief Data Officers** in the various universities and research performing organisations.
- Develop the **economic models of open source software** and make them known within commercialization services
- **Support Software Heritage** and recommend it for the archiving and referencing of source code

### Recognise source code as a contribution to research

- Create an **open source research software prize**
- **Provide greater recognition** for software production in the career of researchers, research support staff

### Build an ecosystem that connects code, data and publications

- Develop **proper coordination** between software forges, open publication archives, data repositories and the scientific publishing sector.

4

## Software College in the National Committee for Open Science

- Five action lines (licence is just the tip of the iceberg, see [details online](#))
- Co-Chairs: Roberto Di Cosmo, François Pellegrini
- Members: 20+ from all research fields
- Key productions:

- Source Code primer



- Report on software forges (needs, limitations, options):



- Address GitLab instance limitations for global contribution:

<https://code.gouv.fr/fr/bluehats/outils-de-forge-2024/>

- *National research software award*. Open Research Europe 2023





# French National plan for Open Science awards 2022 and 2023

## First edition, 2022 prize



MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE

Accueil > Recherche > Science ouverte

Publié le 05.02.2022

### 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.

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

- 129 projects
- 4 awards
- 6 accessit
- first edition
- Coq proof assistant
- Scikit-Learn ML/AI
- Faust music
- Gammapy astronomy

## Second edition, 2023 prize



MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE

Accueil > Recherche > Science ouverte

Publié le 29/11/2023

### Remise des prix science ouverte du logiciel libre de la recherche 2023

Le ministère de l'Enseignement supérieur et de la Recherche remet pour la deuxième édition les Prix science ouverte du logiciel libre de la recherche. Huit logiciels développés par des équipes françaises sont récompensés pour leur contribution à l'avancée de la connaissance scientifique ou pour le caractère prometteur de leurs travaux.

Sommaire

- PfanGOOLIN : espoir de la catégorie « Scientifique et technique »
- Sireli : lauréat de la catégorie « Scientifique et technique »
- NoiseCapture : espoir de la catégorie « Communauté »
- OCaml : lauréat de la catégorie « Communauté »
- KoOps : espoir de la catégorie « Documentation »
- Brian : lauréat de la catégorie « Documentation »
- Fink : espoir de la catégorie « Coup de cœur » du jury
- Hylio : lauréat de la catégorie « Coup de cœur » du jury
- Jury

- 66 projects
- 4 awards
- 4 "espoirs"
- will be run annually

# Blueprint and data from the first edition (2021-2022)

## Blueprints and analysis available



Open Research Europe

99 Views | 12 Downloads | 0 Citations

OPEN LETTER

### Establishing a national research software award

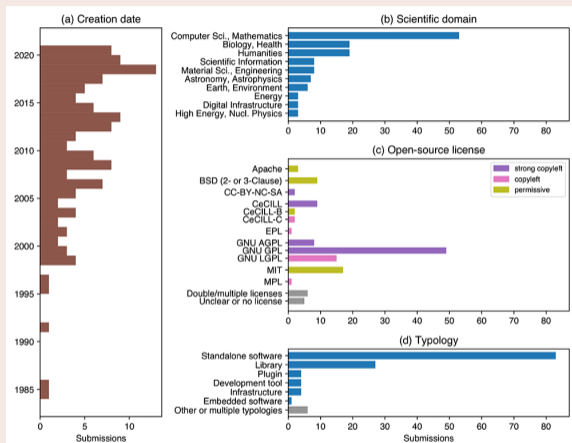
[version 1; peer review: 2 approved]

Isabelle Blanc Catala, Roberto Di Cosmo, Mathieu Graud, Daniel Le Berre, Violaine Louvet, Sophie Renaudin

College of experts for source code and software Committee for Open Science

- goals, design decisions
- challenges and solutions
- lessons learned
- detailed data

## A look at the data



## Emulation is working

Australia: [here](#) and [here](#); Germany: [Helmholtz](#); European Commission: [a CSA](#); ...

# A few key shared needs across all areas



We need a universal shared infrastructure to address them!

- 1 Introduction
- 2 Selected actions of interest (France)
- 3 Meet Software Heritage**
- 4 An example is worth a thousand words
- 5 A lot more lies ahead
- 6 Call to action





## 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** and **share** all software source code

### Research infrastructure



**enable analysis** of all software source code

## 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



Diamond sponsor



Platinum sponsors



Gold sponsors



Silver sponsors



Bronze sponsors



# The largest software archive, a shared infrastructure

One infrastructure  
open and shared

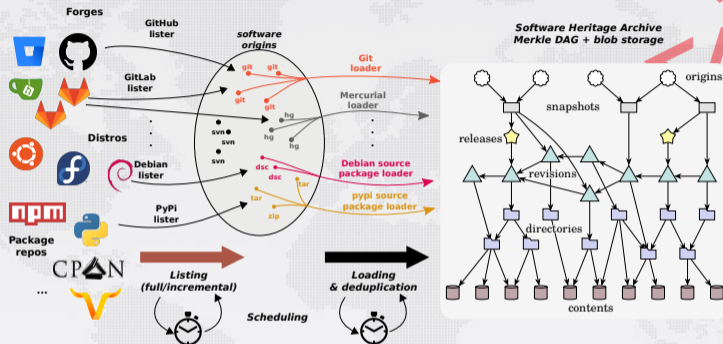


The largest archive ever built



Bitbucket 2,509,402 origins	debian 56,983 origins	git 24,600 origins
GitHub 26,599 origins	gitleaks 136,338 origins	GitLab 53,297 origins
git 197,883,004 origins	Gogs 10,171 origins	GO 4,216,298 origins
git 2,926 origins	GNU 172 origins	heptapod 971,549 origins
Guix 14,482 origins	GNU 354 origins	NixOS 1,207 origins
launchpad 503,631 origins	Maven 312,461 origins	NixOS 14,482 origins

# Address common Open Science and Open Source needs: archival



Global development history permanently archived in a uniform data model

- over 18 billion unique source files from over 290 million software projects
- ~1.5PB (compressed) blobs, ~35 B nodes, ~500 B edges



# Address common Open Science and Open Source needs: reference

## Software Heritage Identifiers (SWHID)

[link to full docs](#)



40+B  
intrinsic,  
decentralised,  
cryptographic

## Full fledged *source code references* for traceability, integrity and reproducibility

- Linux Foundation [SPDX 2.2](#)
- IANA-registered "swh:"
- WikiData property [P6138](#)

Examples: [Apollo 11 AGC excerpt](#), [Quake III rsqrt](#)  
Guidelines available, see [the HOWTO](#)

**Breaking news: standardisation**, see [swhid.org](#)

- 1 Introduction
- 2 Selected actions of interest (France)
- 3 Meet Software Heritage
- 4 An example is worth a thousand words**
- 5 A lot more lies ahead
- 6 Call to action



# A walkthrough

- Browse and Reference (e.g. [Apollo 11 \[excerpt\]](#), your work [may be already there](#) !)
- Trigger archival, use the [updateswh](#) browser extension, configure the [webhooks](#)
  - example: French public open source at: [code.gouv.fr](#)
- Cite with [biblatex-software](#) (CTAN, [Overleaf ACMART template](#))
- Describe with Codemeta (use [codemeta generator](#))
- Curated deposit in SWH via HAL, see for example: [LinBox](#), [SLALOM](#), [Givaro](#), [NS2DDV](#), [SumGra](#), [Coq proof](#), ...
- Extracting all the software products [for Inria](#), [for CNRS](#), [for CNES](#), [for LIRMM](#) or [for Rémi Gribonval](#) using [HalTools](#)

- 1 Introduction
- 2 Selected actions of interest (France)
- 3 Meet Software Heritage
- 4 An example is worth a thousand words
- 5 A lot more lies ahead
- 6 Call to action





DALL-E's view of an institutional portal

**Problem:** institutions need to identify **projects of interest...**  
... projects spread across 100s of platforms

**Solution:** personalised institution portals on SWH

- deposit and curation of (meta)data
- presentation
- extraction

**Why Software Heritage:**

- platform agnostic, metadata from multiple institutions
- shared knowledge base



DALL-E's view of SWH as CD

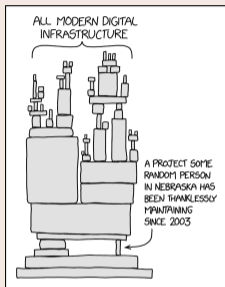
**Problem:** organizations need trusted timestamps on software deposits...

**Solution:** personalised deposit portals with a secure timestamp conforming to ISO/IEC 18014

**Why Software Heritage:**

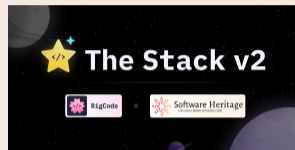
- SWH's Merkle tree is *half* a blockchain
- SWHID guarantees integrity
- universal archive, standardised
- minimal disruption to current deposit workflow

## Cybersecurity: tracing vulnerabilities



See the [SWHSec project](#)

## (Generative) AI



[The Stack V2](#) (subset of Software Heritage) used to build [StarCoder2](#), *best open AI model for coding today* (3, 7 and 15B parameters)

## Big Data challenges



- big data infrastructure
- efficient queries
- integration with other knowledge graphs

- 1 Introduction
- 2 Selected actions of interest (France)
- 3 Meet Software Heritage
- 4 An example is worth a thousand words
- 5 A lot more lies ahead
- 6 Call to action





# A rally flag for a grand vision

Bring together academia, industry, governments, communities

*"to build a reference, global infrastructure for open and better software"*

Software Heritage is the first brick ...

- vendor neutral
- open source
- a worldwide initiative
- a long term initiative

... that will enable

- archival, reference, integrity
- qualification, sharing and reuse
- a global software knowledge base
- test and deploy world class tooling

A lot more is needed

Software Heritage can be the *catalyser* of a way bigger undertaking

Let's work together!

## Annual report



## 5 years in 5 minutes

[Link](#)



## Evolution of our codebase

[Link](#)

