

# Software Heritage

Infrastructure for Open Science

Roberto Di Cosmo  
Director  
roberto@dicosmo.org

March 24, 2026



Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 Adoption and ecosystem
- 4 Building reliable indicators
- 5 Strategic perspectives
- 6 A plan for Spain

## Invisible fabric of digital society



## Knowledge is in the source

```
/**
 * @brief The basic unit of the simulation and is associated to a geographical location.
 *
 * Interventions (e.g., school closures) are tracked at this level. It contains a list of its
 * members (people), places (schools, universities, workplaces etc.), road networks, links to
 * airports etc.
 */
struct Microcell
{
    /* Note use of short int here limits max run time to USHRT_MAX*ModelTimeStep - e.g. 65536*0.25=16384 days=
     * Global search and replace of 'unsigned short int' with 'int' would remove this limit, but use more mem
     */

    int n; // Number of people in microcell
    int adunit; // admin unit microcell belongs to
    int* members; // array of members/hosts of microcell

    int* places[MAX_NUM_PLACE_TYPES]; // list of places (of various place types) within microcell
    unsigned short int NumPlacesByType[MAX_NUM_PLACE_TYPES]; // number of places (of various place types) with
    unsigned short int keyworkerproph, move_trig, place_trig, socdist_trig, keyworkerproph_trig;
    unsigned short int move_start_time, move_end_time;
    unsigned short int place_end_time, socdist_end_time, keyworkerproph_end_time;
    TreatStat moverest, treat, vacc, socdist, placeclose;
    unsigned short int treat_trig, vacc_trig;
    unsigned short int treat_start_time, treat_end_time;
    unsigned short int vacc_start_time;
    IndexList* AirportList;
};
```

## Covid Sim ( [excerpt](#) )

Len Shustek, Computer History Museum

2006

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

# A Global Undertaking

## From all continents

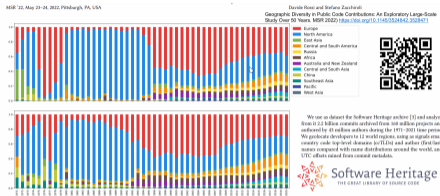
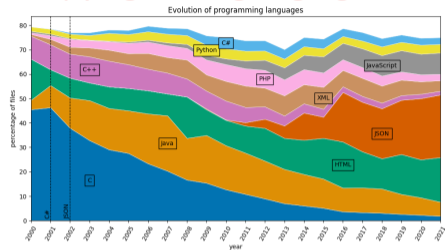
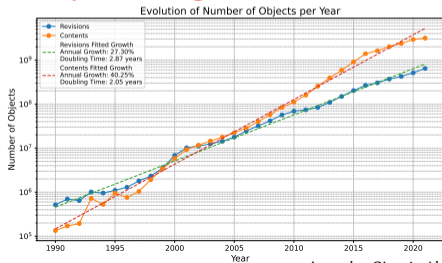


Figure 3: Ratio of commits (above) and active authors (below) by world zone over the 1971-2020 period.

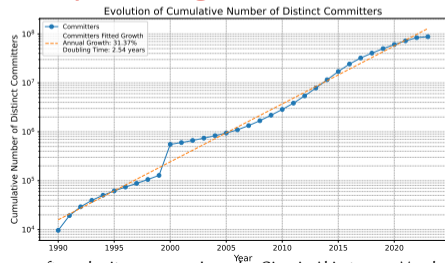
## Many programming languages



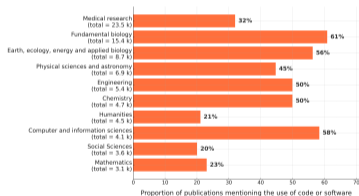
## An exponential growth (code)



## An exponential growth (contributors)



## Pillar of Science across all research areas



## A plurality of needs

### Researcher

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

### Laboratory/team

- track software contributions
- produce reports / web page

### Research Organization

- know its software assets
- technology transfer
- impact metrics

we need a dedicated infrastructure: now we have it!

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder**
- 3 Adoption and ecosystem
- 4 Building reliable indicators
- 5 Strategic perspectives
- 6 A plan for Spain



Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

Inria with  **unesco**





Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

Inria with  unesco



# The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good *built in France since 2015*

Cultural Heritage



Industry



Research



Public Administration



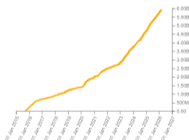
Source files

28,152,651,759



Commits

5,920,787,012



Projects [↗](#)

429,963,770





Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

Inria with  unesco



The largest open source code archive: one infrastructure, open, shared, non profit  
Unique digital common good *built in France since 2015*

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

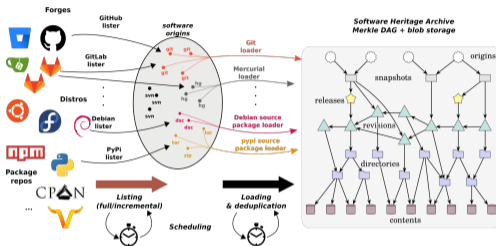


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history development in a single graph



Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

Inria with  unesco



The largest open source code archive: one infrastructure, open, shared, non profit  
Unique digital common good *built in France since 2015*

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

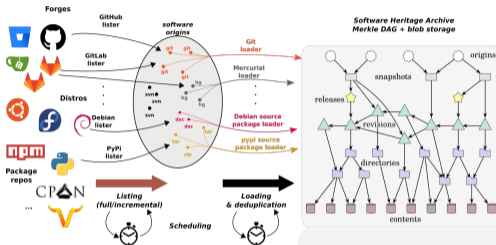


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history  
development in a single graph

- 50 × 10<sup>9</sup> nodes
- 1000 × 10<sup>9</sup> edges
- ~ 3 PB of storage



Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

Inria with  unesco



# The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good *built in France since 2015*

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

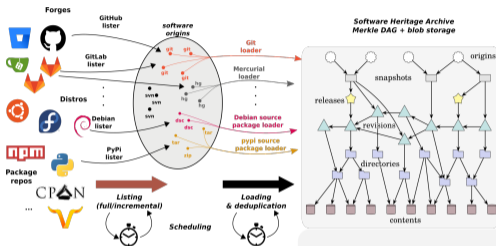


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history  
development in a single graph

- 50 × 10<sup>9</sup> nodes  
- 1000 × 10<sup>9</sup> edges  
~ 3 PB of storage

A revolutionary **infrastructure** ensures **availability** guarantees **integrity** enables **traceability**





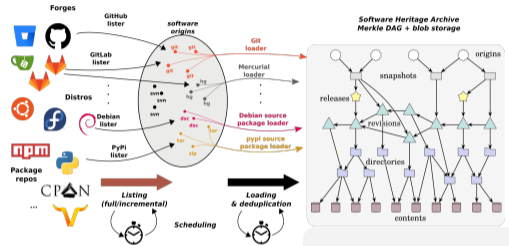
**Software Heritage**  
THE GREAT LIBRARY OF SOURCE CODE

*Inria* with unesco



**The largest open source code archive: one infrastructure, open, shared, non profit**  
Unique digital common good *built in France since 2015*

- Cultural Heritage
- Industry
- Research
- Public Administration



**5000+ platforms**

**All versions, all history development in a single graph**

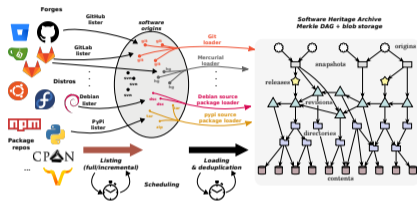
- 50 × 10<sup>9</sup> nodes
- 1000 × 10<sup>9</sup> edges
- ~ 3 PB of storage

A revolutionary **infrastructure** ensures **availability** guarantees **integrity** enables **traceability**



# Addressing key needs in (Open) Science

## Archive (28B+ files, 430M+ projects)



- save now, updateswh, webhooks
- [deposit.softwareheritage.org](https://deposit.softwareheritage.org)

## Describe

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

## Reference (50 billion SWHIDs)

## Intrinsic, cryptographically strong IDs



## Now in SPDX 2.2, Wikidata

<https://swhid.org> - ISO/IEC 18670

## Cite/Credit

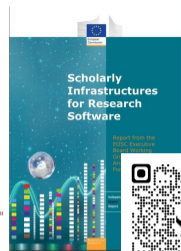
- Contributed [biblatex-software style](#)
- Software Citation from the archive!

# An example is worth a thousand words

- Browse + Reference [ISO 18670] (Apollo 11 [excerpt], your work may be already there !)
- Trigger archival, use the [updateswh](#) browser extension, configure the webhooks
- Cite [from the archive](#) with [biblatex-software](#) (CTAN, [ACMART](#))
- 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](#)
- Example with Parmap: [devel on Github](#), [archive in SWH](#), [curated deposit in HAL](#)
- Example research articles:
  - compare Fig. 1 and conclusions in [the 2012 version](#) and [the updated version](#)
  - SWHID in [a replication experiment](#)

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 Adoption and ecosystem**
- 4 Building reliable indicators
- 5 Strategic perspectives
- 6 A plan for Spain

# Connections and mutualization for the scholarly ecosystem



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

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

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

9 Define and promote an open source software policy

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

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



- [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

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 Adoption and ecosystem
- 4 Building reliable indicators**
- 5 Strategic perspectives
- 6 A plan for Spain

## Software Heritage

Report on the public  
software collected,  
preserved and  
referenced for  
Acad-Hal-Fr

Version 1.0 – 2025-12-29

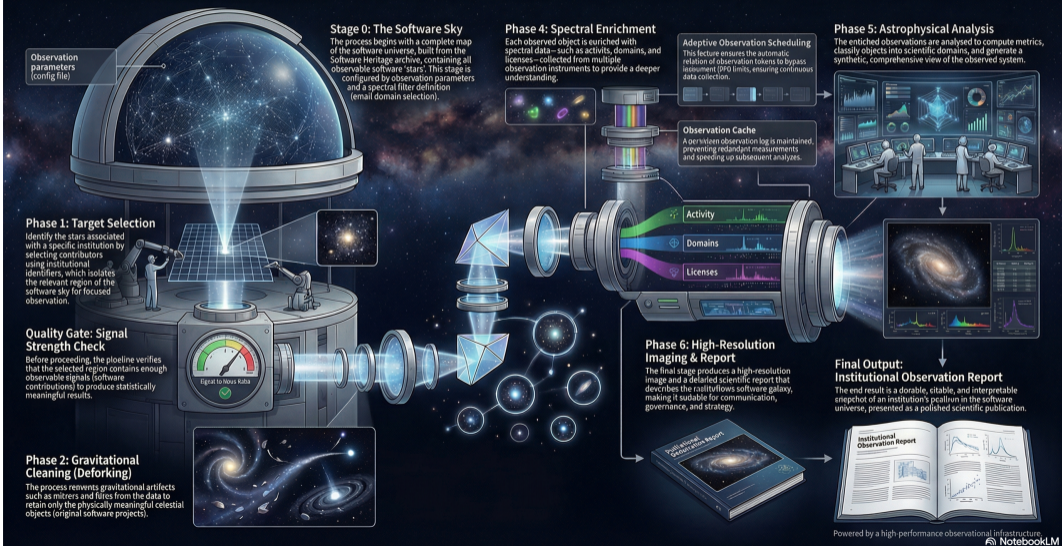
## Contents

<b>1 Introduction</b>	<b>3</b>
1.1 Bird's-Eye View: the key numbers	3
1.2 Technical Information	4
<b>2 Analysis of Contributions</b>	<b>5</b>
2.1 Density of contributions per project	5
2.2 Distribution of first contributions over time	7
2.3 Timespan of institutional contributions	9
2.4 Distribution of project lifetime (based on Software Heritage information)	11
2.5 Distribution of project lifetime (GitHub-based)	13
2.6 Distribution of Citation Files	15
2.6.1 Projects with Citation Files	15
2.6.2 Projects with Citation Files Only in Software Heritage	17
2.7 Top Projects by Repository Host	19
<b>3 Analysis of the Projects to which Acad-Hal-Fr contributes</b>	<b>34</b>
3.1 Distribution by number of stars	34
3.2 Main organizations	36
3.3 Top 40 projects by number of contributions	37
3.4 Top 40 projects by GitHub stars	39
3.5 Top 40 projects by repository age	41
<b>4 Programming Languages</b>	<b>53</b>
<b>5 License Analysis</b>	<b>55</b>
5.1 License Distribution	55
5.2 License Status Overview	55
5.3 License Statistics	55
5.4 Top License Types	55
<b>6 Platform Distribution Analysis</b>	<b>57</b>
6.1 Platform Distribution Overview	57
6.2 Platform Distribution (Log Scale)	57
6.3 Platform Statistics	57
6.4 Top Platforms	57
6.5 Summary	58
<b>7 Tentative Classification in Research Areas/Domains</b>	<b>59</b>
7.1 Distribution of domains of contributors (all projects)	60
7.2 Distribution of domains of contributors (selected projects with more than 10 contributions)	61
<b>8 Methodology: From Source Graph to Analysis-Ready Dataset</b>	<b>61</b>
<b>9 Conclusion</b>	<b>63</b>



# Software Heritage is the Very Large Telescope for source code

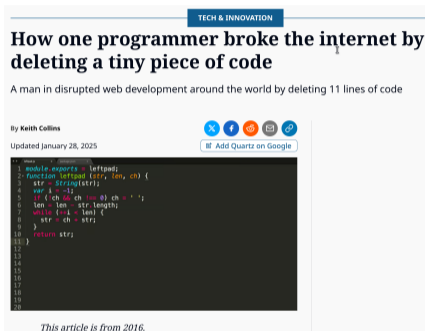
## From the Software Sky to Institutional Insight: Anatomy of an Observational Pipeline



- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 Adoption and ecosystem
- 4 Building reliable indicators
- 5 Strategic perspectives**
- 6 A plan for Spain

# The next frontier: digital resilience

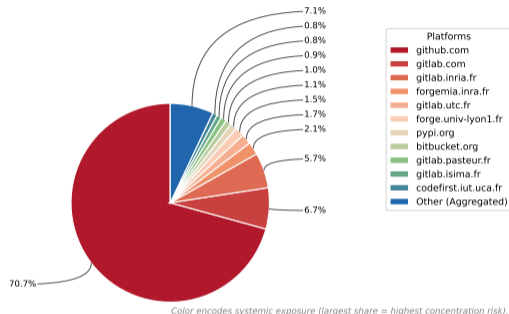
## An early warning: left-pad (March 2016)



See the [Wikipedia](#) article:

- Facebook, PayPal, Netflix, Spotify affected
- web broken worldwide ~2.5 hours

## A clear and present danger



French Academia devs (Source: Software Heritage)

*"If GitHub or PyPI disappeared tomorrow,  
all research would not just slow down  
it would stop."*

## Connecting reproducible deployment to a long-term source code archive



Ludovic Courtès — March 29, 2019

GNU Guix can be used as a “package manager” to install and upgrade software packages as is familiar to GNU/Linux users, or as an environment manager, but it can also provision containers or virtual machines, and manage the operating system running on your machine.

One foundation that sets it apart from other tools in these areas is *reproducibility*. From a high-level view, Guix allows users to *declare* complete software environments and instantiate them. They can share those environments with others, who can replicate them or adapt them to their needs. This aspect is key to reproducible computational experiments: scientists need to reproduce software environments before they can reproduce experimental results, and this is one of the things we are focusing on in the context of the [Guix-HPC](#) effort. At a lower level, the project, along with others in the [Reproducible Builds](#) community, is working to ensure that software build outputs are [reproducible, bit for bit](#).

Work on reproducibility at all levels has been making great progress. Guix, for instance, allows you to [travel back in time](#). That Guix can travel back in time *and* build software reproducibly is a great step forward. But there’s still an important piece that’s missing to make this viable: a stable source code archive. This is where [Software Heritage](#) (SWH for short) comes in.

**When source code vanishes**

Next step: scale up a sovereign fallback for Europe’s build chains

When a platform or registry fails,

Guix fallbacks

to Software Heritage



and continues to build!

*(since 2019)*

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 Adoption and ecosystem
- 4 Building reliable indicators
- 5 Strategic perspectives
- 6 A plan for Spain

## ENCA 2023-2027



- Objective 1 (page 11): "[...] infraestructuras digitales interoperables suficientemente robustas y bien articuladas"
- Objective 2 (page 11): "Establecer nuevos mecanismos de evaluación de la investigación..."

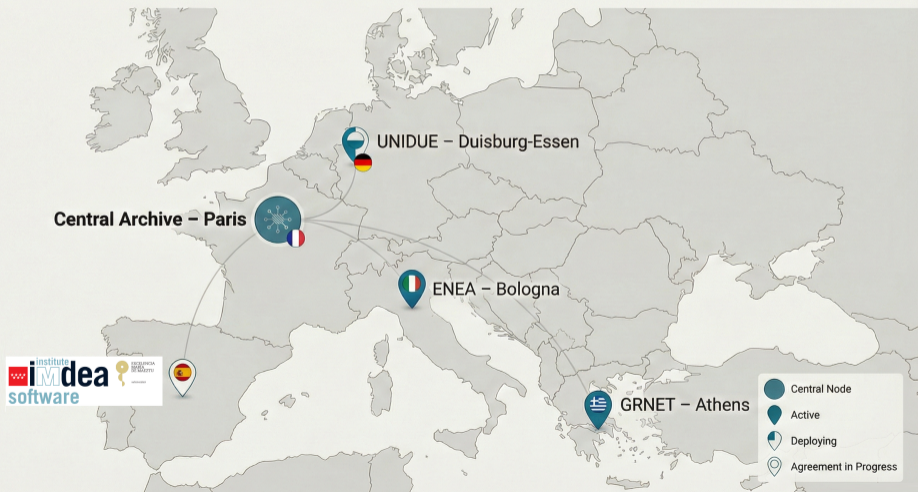
## French National Plan for Open Science 2021

- Define and promote an open source software policy
- Recognise source code as a contribution to research
- Create an open source research software prize
- Support Software Heritage for archiving and referencing source code

## Key Actions (selection)



## The **Software Heritage European Network:** A Distributed Infrastructure for Global Source Code Preservation



## Foster adoption of best practices

archive, reference, describe, cite software: see [the Software Heritage HOWTO](#)

## Establish incentives: funding, evaluation and recognition

- count *quality* software contributions (all aspects!)
- establish National/Regional/Institutional *awards* ([get the blueprints from France](#))

## Build objective indicators: *Strategic Insights Report, not surveys*

## Avoid *balkanisation*, build upon shared, open, non profit infrastructures (POSI, UNESCO)

- [FAIRCORE4EOSC](#) interconnected infrastructures with Software Heritage
- [OSPO-Radar](#) institutional portals/reporting via Software Heritage outside France

## Avoid proprietarization of public research result

*publicly funded research software should be open source*, exceptions **must be justified**