

Software Heritage

Sponsor and Members General Assembly

Roberto Di Cosmo
Director
roberto@dicosmo.org

February 18th, 2026



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Software as a pillar of modern society
- 2 Software Heritage's mission
- 3 Adoption and ecosystem
- 4 A new phase: products and services for all stakeholders
- 5 Strategic Software Insights Report

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 memco
     */

    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](#))

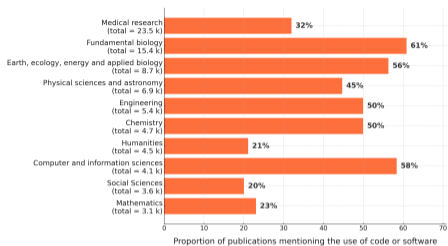
Len Shustek, Computer History Museum

2006

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

A Global Undertaking

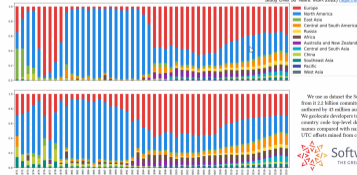
Pillar of Science across all research areas



From all continents

NSR '22, May 21-24, 2022, Pittsburgh, PA, USA

David R. Rossi and Stefano Zacchiroli
Geographic Diversity in Public Code Contributions: An Exploratory Large-Scale Study Over 50 Years. MSR 2022 (<https://doi.org/10.1145/3539462.3539471>)



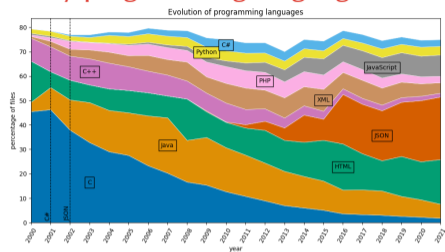
We use as dataset the Software Heritage archive [1] and analyze from it 2.2 billion commits archived from 500 active projects and authored by 41 million authors during the 1971-2021 time period. We analyze developers in 12 world regions, using as regional proxy country code top-level domains (ccTLDs) and authors (first last) names compared with name distributions around the world, and UTC offsets raised from commit metadata.



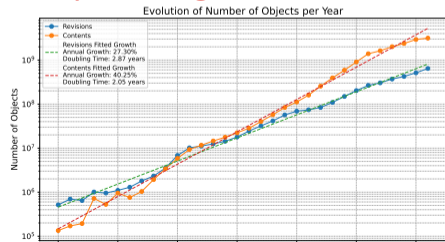
Figure 3: Ratio of commits (above) and active authors (below) by world zone over the 1971-2020 period.



Many programming languages



An exponential growth



Nature, October 2025

nature View all journals Q Search Login
Explore content About the journal Publish with us Subscribe Sign up for alerts RSS feed

nature > comment > article

COMMENT | 17 October 2025

Stop treating code like an afterthought: record, share and value it

Scientists, research institutions, funders, libraries and publishers must all improve software practices.

By Roberto Di Cosmo, Sabrina Grange, Konrad Hinsen, Nicolas Julien, Daniel Le Berre, Violaine Louvet, Camille Mauzet, Chloémette Maurice, Raehel Morat & Nicolas P. Rouquier



Illustration: Phil Wheeler

Software Heritage sponsor meeting

www.softwareheritage.org

ICSE FoSE 2026

Reclaiming Software Engineering as the Enabling Technology for the Digital Age

Tanja E. J. Vos (1), Tijs van der Storm (2), Alexander Serebrenik (3), Lionel Briand (4), Roberto Di Cosmo (5), J.-M. Bruel (6), Benoit Combemale (7, 6)

Show details

- 1 Open University of the Netherlands [Heerlen]
- 2 CWI - Centrum Wiskunde & Informatica
- 3 TU/e - Eindhoven University of Technology [Eindhoven]
- 4 LERO - The Irish Software Engineering Research Centre
- 5 Software Heritage
- 6 IIRIT-SM@RT - Smart Modeling for softw@re Research and Technology
- 7 DiverSe - Diversity-centric Software Engineering

Abstract en

Software engineering is the invisible infrastructure of the digital age. Every breakthrough in artificial intelligence, quantum computing, photonics, and cybersecurity relies on advances in software engineering, yet the field is too often treated as a supportive digital component rather than as a strategic, enabling discipline. In policy frameworks, including major European programmes, software appears primarily as a building block within other technologies, while the scientific discipline of software engineering remains largely absent. This position paper argues that the long-term sustainability, dependability, and sovereignty of digital technologies depend on investment in software engineering research. It is a call

Keywords en

open source sustainability
digital sovereignty research policy
enabling technology
software engineering

Domains

Software Engineering [cs.SE]

- 1 Software as a pillar of modern society
- 2 Software Heritage's mission
- 3 Adoption and ecosystem
- 4 A new phase: products and services for all stakeholders
- 5 Strategic Software Insights Report

Unveiled in 2016



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

Software Heritage sponsor meeting

Universal archive



preserve and share all software source code

www.softwareheritage.org @swheritage

Research infrastructure



enable analysis of all software source code

Contact: roberto@dicosmo.org February 18th, 2026

A universal software archive, as a shared infrastructure

One infrastructure
open and shared

Cultural Heritage



Industry



Research

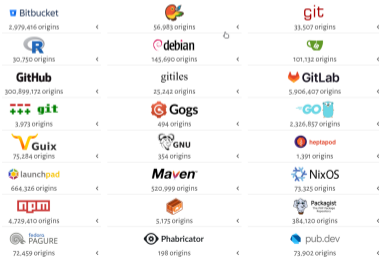


Public Administration



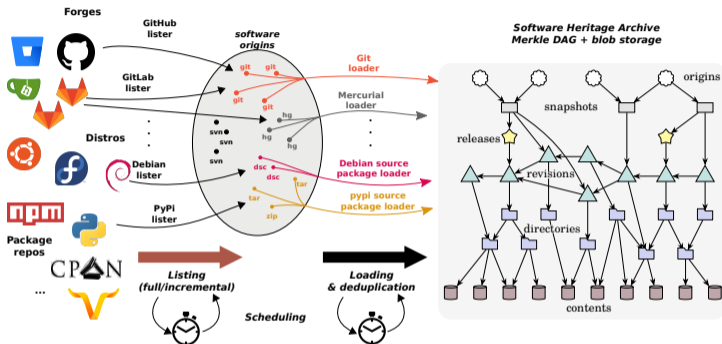
Software Heritage

The largest archive ever built



figures as of January 8 2026

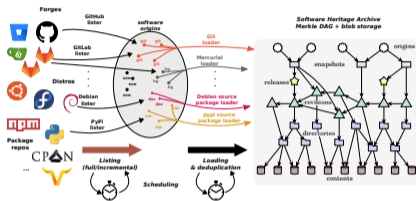
From a Babel tower to a giant graph



Global development history permanently archived in a uniform data model

- over 27 billion unique source files from over 420 million software projects
- ~2PB (compressed) blobs, ~50 B nodes, ~1 trillion edges

Archive (25B+ files, 400M+ projects)



- save now, updateswh, webhooks
- deposit.softwareheritage.org

Describe

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

Reference (50 billion SWHIDs)

Intrinsic, cryptographically strong IDs



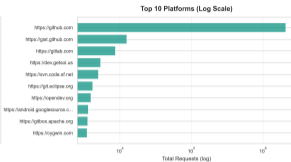
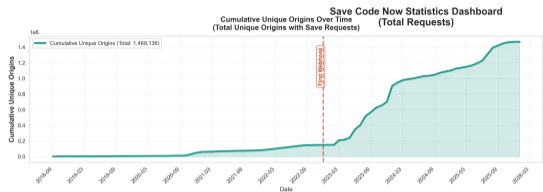
Now in SPDX 2.2, Wikidata

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

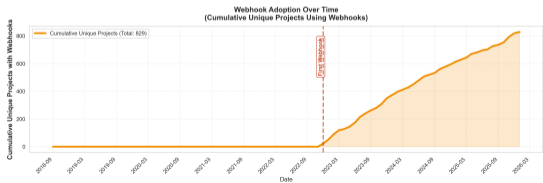
Cite/Credit

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

A peek at Save Code Now adoption



Request Type Distribution



Summary Statistics (Total Requests Mode)

Total: 2,219,897
Manual: 2,183,218 (98.4%)
Webhook: 35,887 (1.6%)

Date Range:
2018-09-07
to
2026-01-06

Unique Origins:
Total: 1,469,136
Avg/Origin: 1.51

Webhook:
Projects: 829
Platforms: 2,991

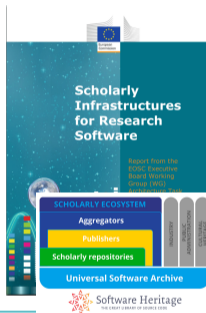
Top Platform:
https://github.com
(2,886,212)

- 1 Software as a pillar of modern society
- 2 Software Heritage's mission
- 3 Adoption and ecosystem**
- 4 A new phase: products and services for all stakeholders
- 5 Strategic Software Insights Report

A few adoption indicators



Policy



- [Recommendations in ANR 2023 guidelines \(p. 17\)](#)
- HAL+SWH in [the Open Science software booklet](#)

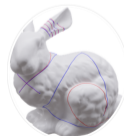
Users and collaborations



What are they “referencing”?

source	n	percentage
Not available	2868	46.22
GitHub	1151	18.55
software heritage	387	6.24
zenodo	142	2.29
r package	70	1.13
cran	56	0.90
r package version	54	0.87
gitlab	35	0.56

Graphics Replicability Stamp Initiative



b/Surf: Interactive Bézier Splines on Surface Meshes

Claudio Mancinelli, Giacomo Nazzaro, Fabio Pellacini, Enrico Puppo
IEEE Transactions on Visualization and Computer Graphics (TVCG)



Repository



Projects



FAIRCORE4EOSC
Core Components Supporting a FAIR EOSC

The CodeMeta Project



FAIR-IMPACT
Expanding FAIR solutions across EOSC

- 1 Software as a pillar of modern society
- 2 Software Heritage's mission
- 3 Adoption and ecosystem
- 4 A new phase: products and services for all stakeholders
- 5 Strategic Software Insights Report

Institutional portal



Vulnerability graph



Strategic Insights



Dataset Factory



See annual report for more 



- 1 Software as a pillar of modern society
- 2 Software Heritage's mission
- 3 Adoption and ecosystem
- 4 A new phase: products and services for all stakeholders
- 5 **Strategic Software Insights Report**

Software is a strategic asset, but

- ... tracking and analysing your software assets is **difficult**
- ... assessing your open source presence is **hard**
- ... knowing how you compare to other is **complex**

Software Heritage helps you gain deep insight into your software assets

A preview of the Strategic Software Insights report

Software Heritage

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

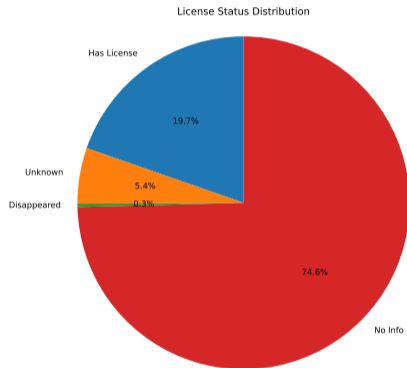
Version 1.0 – 2025-12-29

Contents

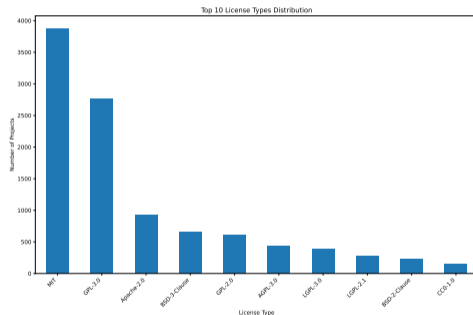
1 Introduction	3
1.1 Bird's-Eye View: the key numbers	3
1.2 Technical Information	4
2 Analysis of Contributions	5
2.1 Density of contributions per project	5
2.2 Distribution of first contributions over time	7
2.3 Timespan of institutional contributions	9
2.4 Distribution of project lifetime (based on 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
3 Analysis of the Projects to which Acad-Hal-Fr contributes	34
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
4 Programming Languages	53
5 License Analysis	55
5.1 License Distribution	55
5.2 License Status Overview	55
5.3 License Statistics	55
5.4 Top License Types	55
6 Platform Distribution Analysis	57
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
7 Tentative Classification in Research Areas/Domains	59
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
8 Methodology: From Source Graph to Analysis-Ready Dataset	61
9 Conclusion	63

A peek inside the French Academic Software report, cont'd

How many projects have a license?



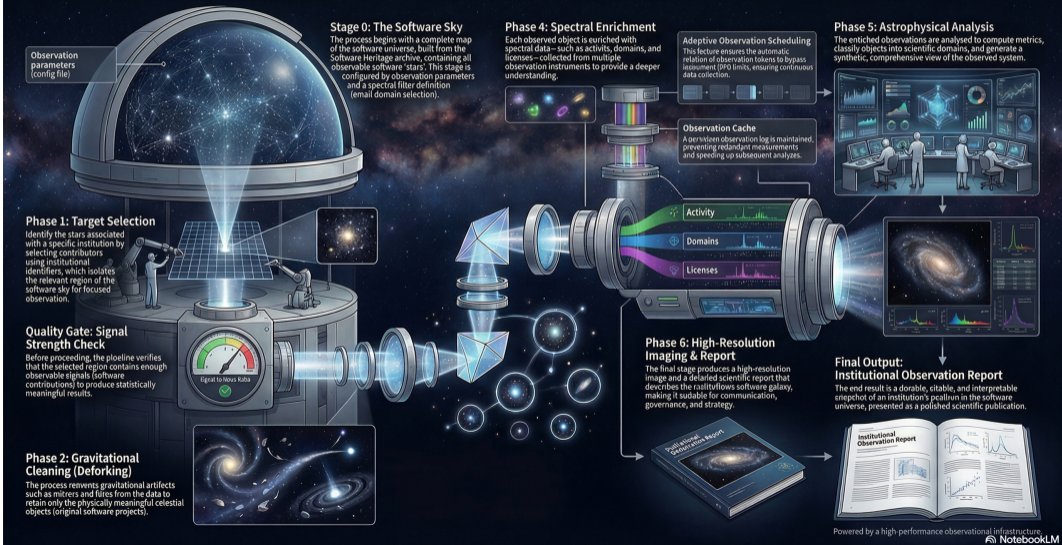
What are the top licences used?



how do we build these indicators?

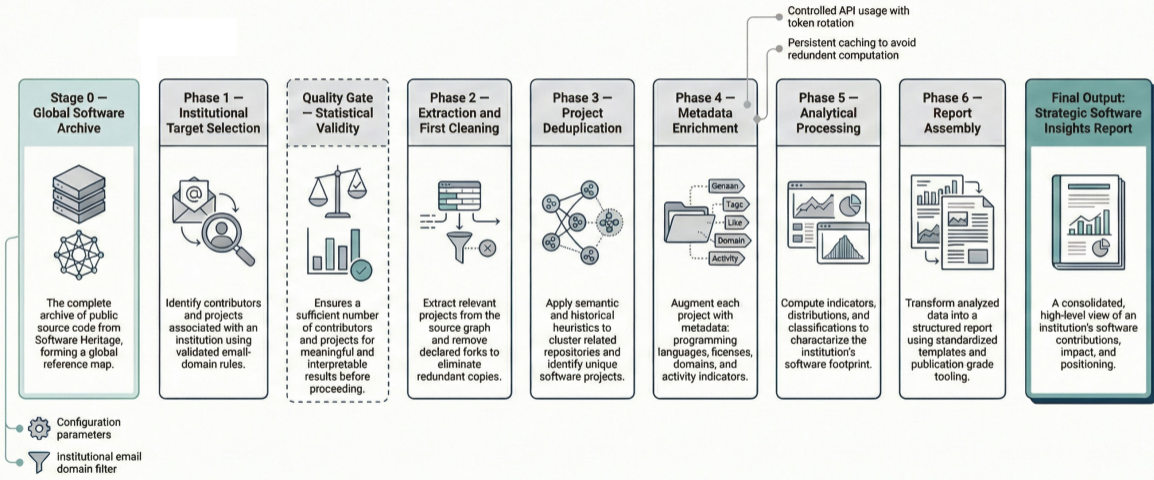
Software Heritage is the Very Large Telescope for source code

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



Strategic Software Insights Report — From Global Archive to Institutional Insight

A multi-stage analytical pipeline built on the Software Heritage archive



A complex, costly process scans the whole Software Heritage archive to deliver **key**

The report in context, and your input welcome

CodeCommons (BPI project, Datasets for AI and beyond)

- broad spectrum of metadata: licence, language, vulnerabilities, issues, PR...
- connections with context:
 - forge metadata: GitHub, GitLab, etc.
 - curated metadata: HAL, OSPO-Radar, catalogue, etc.
 - qualified software mentions in publications: Hal, swMath, SoFAIR, OpenAlex, etc.
 - vulnerabilities

Enriched Data Model

Question: what extra (meta)data is needed/useful for you?

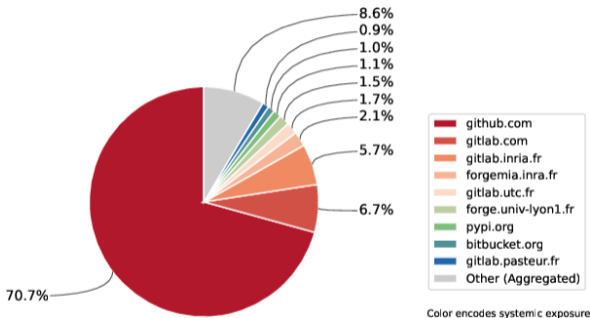
Report generation constraints

- first phase = filtering 100+M personal emails (cannot open this)
- Software Heritage must do the heavy lifting for you

Important: it will be a professional service model, with a cost

The next frontier: from preservation and analysis to resilience

Infrastructure concentration (French academia)



Collaborative development depends **massively**
on **very few** *single points of failures*

Software Heritage is the systemic response

Today

- All the code, also niche and legacy code
- ISO 18670 for robust traceability
- Enables long term reproducibility

With significant more funding

- automatic recovery (registry failure, network partition)
- resilience (mirrors)
- strategic autonomy