

Software as Critical Infrastructure

addressing Europe's Strategic Blind Spot

Roberto Di Cosmo

Director, Software Heritage
Inria and Université Paris Cité

March 18 2026
DG RTD / DG CONNECT, Brussels



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Why Software Source Code matters
- 2 How are we managing [our] software source code?
- 3 Regulation and Infrastructure
- 4 Software Heritage: open, non profit, multistakeholder
- 5 From archive to strategic backbone
- 6 The way forward

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](#))

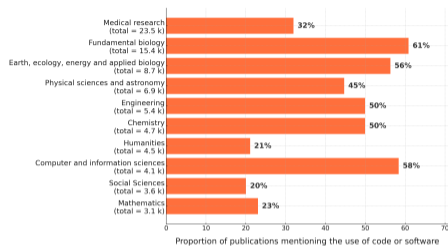
Len Shustek, Computer History Museum

2006

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

A Global Undertaking...

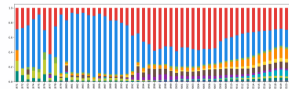
Pillar of Science across all research areas



From all continents

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

David R. Rosen and Stefano Zachariu
Geographic Diversity in Public Code Contributions: An Exploratory Large-Scale Study Over 50 Years. MSR 2022 (<https://doi.org/10.1145/3534462.3534471>)



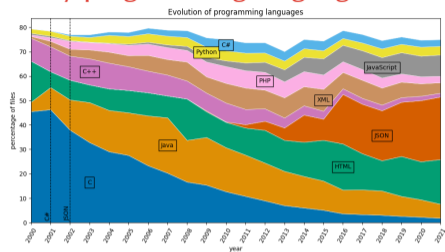
We use as dataset the Software Heritage archive [1] and extract from it 2.2 billion commits archived from 500 active projects and authored by 41 million authors during the 1971-2021 time period. We geocode developers to 12 world regions, using as signals email country code top-level domains (ccTLDs)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



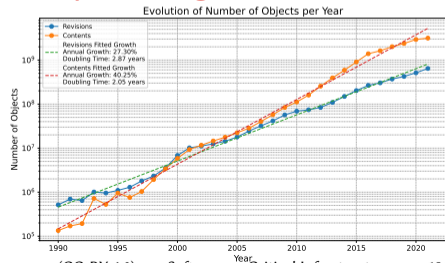
Figure 5: Ratio of commits (above) and active authors (below) by world zone over the 1971-2020 period.



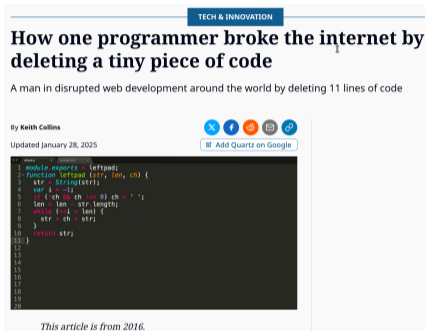
Many programming languages



An exponential growth



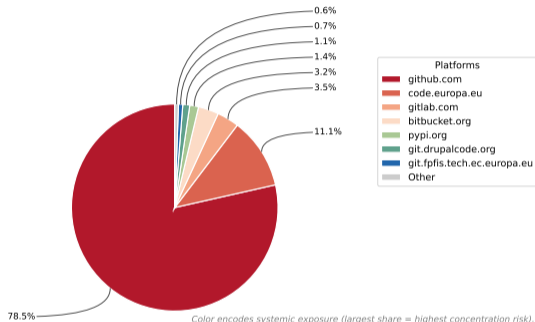
An early warning: left-pad (March 2016)



See the [Wikipedia](#) article:

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

A clear and present danger



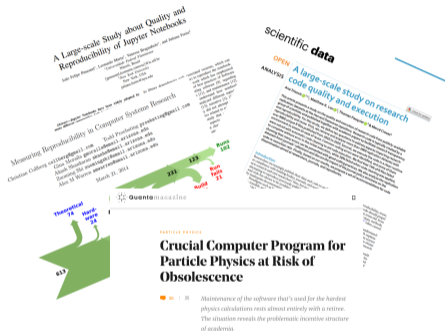
Used by europa.eu devs (Source: Software Heritage)

"If GitHub or PyPI disappeared tomorrow, Europe would not just slow down — it would stop."

- 1 Why Software Source Code matters
- 2 How are we managing [our] software source code?
- 3 Regulation and Infrastructure
- 4 Software Heritage: open, non profit, multistakeholder
- 5 From archive to strategic backbone
- 6 The way forward

Different angles, same bad news

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

Policy is evolving in response to this challenge

- 1 Why Software Source Code matters
- 2 How are we managing [our] software source code?
- 3 Regulation and Infrastructure**
- 4 Software Heritage: open, non profit, multistakeholder
- 5 From archive to strategic backbone
- 6 The way forward

Regulation

Europe's Software Sovereignty

REGULATION



- **Cyber Resilience Act (CRA) & NIS2**
Mandate strict software security, long-term maintenance, and end-to-end supply-chain resilience.
- **AI Act Transparency**
Legal requirements for the transparency and reproducibility of training data and the underlying algorithms.
- **Open Science Policies**
Policies requiring long-term availability and reproducibility of research software for scientific integrity.

“Legal obligations assume permanent access to source code.”

Regulation Without Infrastructure: Europe's Software Sovereignty Gap

REGULATION



- **Cyber Resilience Act (CRA) & NIS2**
Mandate strict **software security**, long-term maintenance, and end-to-end supply-chain resilience.
- **AI Act Transparency**
Legal requirements for the **transparency** and **reproducibility** of training data and the underlying algorithms.
- **Open Science Policies**
Policies requiring long-term availability and reproducibility of research software for scientific integrity.

INFRASTRUCTURE



- **Dependence on Non-European Platforms**
Relies on **proprietary platforms** like GitHub, GitLab.com, and package managers (PyPI, npm, Maven Central).
- **No EU-Wide Continuity Guarantee**
Existing platforms provide **no Service Level Agreements (SLA)** for preservation and operate under **no legal obligation** to protect European sovereignty.
- **Fragmented National Initiatives**
Current efforts are often national, not interoperable, and lack a shared global governance structure.

"Legal obligations assume permanent access to source code."

"Critical dependencies outside European control."

Regulation Without Infrastructure: Europe's Software Sovereignty Gap

REGULATION



- **Cyber Resilience Act (CRA) & NIS2**
Mandate strict **software security**, long-term maintenance, and end-to-end supply-chain resilience.
- **AI Act Transparency**
Legal requirements for the **transparency** and **reproducibility** of training data and the underlying algorithms.
- **Open Science Policies**
Policies requiring long-term availability and reproducibility of research software for scientific integrity.

INFRASTRUCTURE



- **Dependence on Non-European Platforms**
Relies on **proprietary platforms** like GitHub, GitLab.com, and package managers (PyPI, npm, Maven Central).
- **No EU-Wide Continuity Guarantee**
Existing platforms provide **no Service Level Agreements (SLA)** for preservation and operate under **no legal obligation** to protect European sovereignty.
- **Fragmented National Initiatives**
Current efforts are often national, not interoperable, and lack a shared global governance structure.

⚠️ **THE STRUCTURAL GAP** ⚠️

“Legal obligations assume permanent access to source code.”

Europe regulates software as critical infrastructure – but does not operate a software continuity infrastructure

“Critical dependencies outside European control.”

- 1 Why Software Source Code matters
- 2 How are we managing [our] software source code?
- 3 Regulation and Infrastructure
- 4 Software Heritage: open, non profit, multistakeholder**
- 5 From archive to strategic backbone
- 6 The way forward



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



Source files

28,152,651,759

Industry



Commits

5,920,787,012

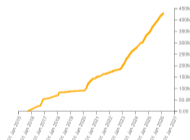
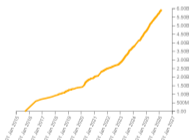
Research



Projects

429,963,770

Public Administration





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



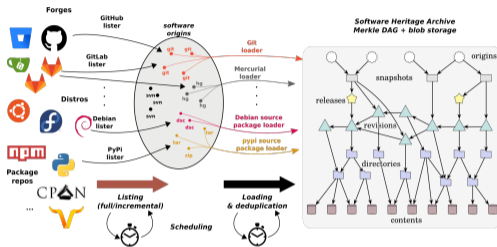
Commits

5,920,787,012



Projects

429,963,770



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

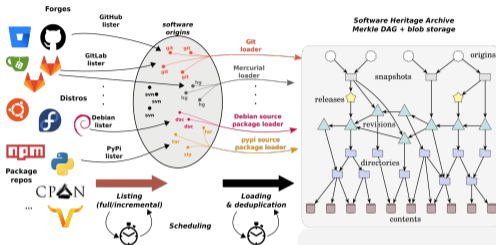


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph

- 50 × 10⁹ nodes
- 1000 × 10⁹ 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

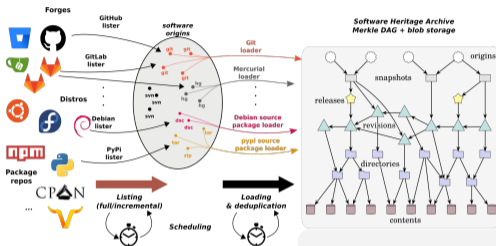


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph

- 50 × 10⁹ nodes
- 1000 × 10⁹ edges
~ 3 PB of storage

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





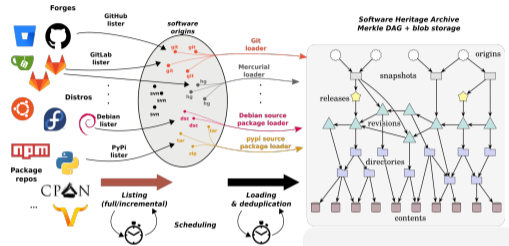
Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

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 × 10⁹ nodes
- 1000 × 10⁹ edges
~ 3 PB of storage

A revolutionary infrastructure ensures availability guarantees integrity enables traceability



- 1 Why Software Source Code matters
- 2 How are we managing [our] software source code?
- 3 Regulation and Infrastructure
- 4 Software Heritage: open, non profit, multistakeholder
- 5 From archive to strategic backbone**
- 6 The way forward

Standards (Software Hash ID, SWHID)



Intrinsic identifier
integrity, traceability
50B+ SWHID in archive

Software Pillar of (Open) Science EOSC SIRS report (2020)



National Open Science Plan
(2021)
French National RI Roadmap
(2022)

Pillar of transparent AI



StarCoder2/The Stack v2
(Hugging Face)
CodeCommons (BPI, France)
AI Factory (AI2F, EU)

Next gen Cybersecurity



Global Vulnerabilities
(Sec4AI4Sec, SWHSec)
Software supply chain
(CRA, SBOM,...)

A (version of a) dependency of a critical software vanishes from GitHub, Gitlab, ...:

- how to **audit** it five years later?
- how to **reproduce** the results?
- how to **ensure** continuity of service?
- how to **retrieve** them?

Proof of capability: Guix integration

Software Heritage and GNU Guix join forces
to enable long term reproducibility



Guix demonstrates automated recovery using the Software Heritage archive

When a platform or registry fails...

Guix fallbacks to Software Heritage

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

and continues to build!



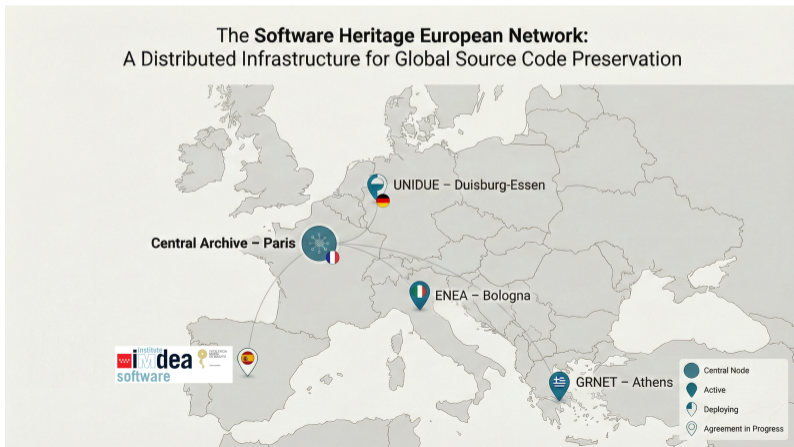
Since 2019

Next step: Sovereign fallback for Europe's build chains

Extend to key package managers

npm, PyPI, Maven Central, Cargo, CPAN, RubyGems, Packagist, ...

Realize ecumenical resilience



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

Software Heritage: one infrastructure, six strategic levers

Research infrastructure

two facets:

- support Open Science
- enable Software Science

Traceability and compliance

- Software supply chain (SBOM)
- integrity (SWHID)
- availability
- support CRA/NIS2 mandates

Cybersecurity

Cyber forensics, vulnerability tracking at world scale

Resilience

when npm/PyPI/Maven **fail**,
Europe **continues to build**

Transparent AI

- **CodeCommons** : provenance and attribution for LLMs
- AI Act, Art. 11, 12, 15:
 - technical documentation
 - record-keeping
 - reproducibility

Metrics we can count on

- Objective Strategic Insights
- Support policy implementation

Regulation Without Infrastructure: Bridging Europe's Software Sovereignty Gap

REGULATION



- **Cyber Resilience Act (CRA) & NIS2**
Mandate strict software security, long-term maintenance, and end-to-end supply-chain resilience.
- **AI Act Transparency**
Legal requirements for the transparency and reproducibility of training data and the underlying algorithms.
- **Open Science Policies**
Policies requiring long-term availability and reproducibility of research software for scientific integrity.

"Legal obligations assume permanent access to source code."

INFRASTRUCTURE



- **Dependence on Non-European Platforms**
Relies on proprietary platforms like GitHub, GitLab.com, and package managers (PyPI, npm, Maven Central).
- **No EU-Wide Continuity Guarantee**
Existing platforms provide no Service Level Agreements (SLA) for preservation and operate under no legal obligation to protect European sovereignty.
- **Fragmented National Initiatives**
Current efforts are often national, not interoperable, and lack a shared global governance structure.

"Critical dependencies outside European control."

! THE STRUCTURAL GAP !

Europe regulates software as critical infrastructure – but does not operate a software continuity infrastructure

SOFTWARE HERITAGE

Over
24
billion

Universal Archive of Source Code

An operational global archive containing over 24 billion files and 375 million projects (and growing).



Intrinsic Identifiers (SWHID – ISO 18670)

Provides a standardized "fingerprint" for software to ensure absolute traceability, auditability, and citation.



Sovereign Mirror Network

Resilient, neutral network of nodes (e.g., Italy, Greece, Spain) ensuring archive cannot be destroyed or controlled by a single entity.

Ready for Scaling, Not Invention

Already operational; requires political and financial scaling to fully support European regulatory ambitions.

"Infrastructure that makes European regulation enforceable."

- 1 Why Software Source Code matters
- 2 How are we managing [our] software source code?
- 3 Regulation and Infrastructure
- 4 Software Heritage: open, non profit, multistakeholder
- 5 From archive to strategic backbone
- 6 The way forward

A plan for Europe

Europe **does not need to invent** this infrastructure...

it can **adopt and scale** what took 10+years to build

Next steps

All recognise **software continuity** as critical enabler for

- research and open science
- digital sovereignty
- regulation enforcement (CRA, NIS2, AI Act, etc.)

All **avoid balkanisation**, support **mutualisation**

- build on common, shared, open, non profit infrastructures

DG RTD / DG CONNECT coordinate a joint effort

- **grow** Software Heritage into the EU strategic software backbone
- **unlock** funding *now* to close the gap (~20 to 40Me table stakes)
- **build** the long term resilience and continuity strategy