

Software Heritage

Towards A New Era for Software Engineering, Cybersecurity, and AI

Roberto Di Cosmo

Director, Software Heritage
Inria and Université Paris Cité

May 21st, 2026

Colloquium Sorbonne Université



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

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.”

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)

```
P63SPOT3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND    CHAN33
              EXTEND
              BZF     P63SPOT4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF     CODE500      # ASTRONAUT:  PLEASE CRANK THE
              TC      BANKCALL     #              SILLY THING AROUND
              CADR    GOPERF1
              TCF     GOTOP00H      # TERMINATE
              TCF     P63SPOT3      # PROCEED     SEE IF HE'S LYING

P63SPOT4      TC      BANKCALL     # ENTER      INITIALIZE LANDING RADAR
              CADR    SETPOS1

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

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)

```
P63SPOT3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND     CHAN33
              EXTEND
              BZF      P63SPOT4 # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF      CODE500 # ASTRONAUT: PLEASE CRANK THE
              TC       BANKCALL #
              CADR     GOPERF1 #
              TCF      GOTOP00H # TERMINATE
              TCF      P63SPOT3 # PROCEED SEE IF HE'S LYING

P63SPOT4      TC       BANKCALL # ENTER INITIALIZE LANDING RADAR
              CADR     SETPOS1

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

Covid Sim (excerpt)

```
/**
 * @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=44 yrs.
    Global search and replace of 'unsigned short int' with 'int' would remove this limit, but use more memory.
    */

    int n; // Number of people in microcell
    int admin; // 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) within microcell
    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;
};
```

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)

```
P63SPOT3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND    CHAN33
              EXTEND
              BZF     P63SPOT4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF     CODE500      # ASTRONAUT: PLEASE CRANK THE
              TC      BANKCALL     # SILLY THING AROUND
              CADR    GOPERF1
              TCF     GOTOP00H     # TERMINATE
              TCF     P63SPOT3     # PROCEED SEE IF HE'S LYING

P63SPOT4      TC      BANKCALL     # ENTER INITIALIZE LANDING RADAR
              CADR    SETPOS1

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

Covid Sim (excerpt)

```
/**
 * @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=44 yrs.
     * Global search and replace of 'unsigned short int' with 'int' would remove this limit, but use more memory.
     */

    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) within microcell
    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;
};
```

Len Shustek, *Computer History Museum*

2006

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

Yuval Noah Harari (on COVID 19)

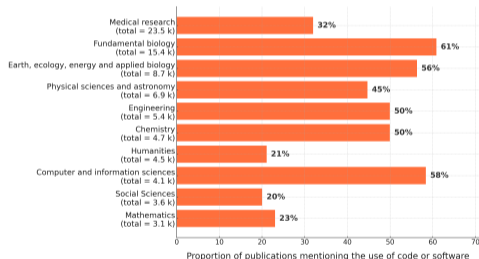
“The real antidote [to epidemic] is scientific knowledge and global cooperation.”

(Open Source) Software is *precious technical and scientific knowledge*

Yuval Noah Harari (on COVID 19)

“The real antidote [to epidemic] is scientific knowledge and global cooperation.”

Software powers modern research



20%+ articles use software, all disciplines

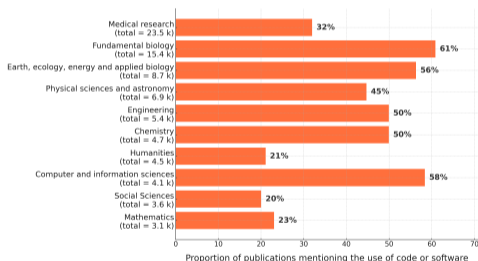
2025 French Open Science Monitor

(Open Source) Software is *precious technical and scientific knowledge*

Yuval Noah Harari (on COVID 19)

"The real antidote [to epidemic] is scientific knowledge and global cooperation."

Software powers modern research



20%+ articles use software, all disciplines
2025 French Open Science Monitor

We can still talk to the early inventors



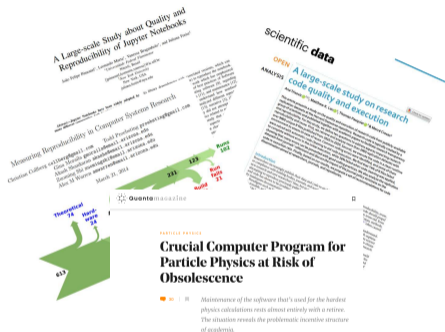
"Telling historical stories is the best way to teach. It's much easier to understand something if you know the threads it is connected to."

Donald E. Knuth
Len Shustek

CACM, January 2021

How are we managing our (open source) software?

Reproducibility, maintenance in Academia

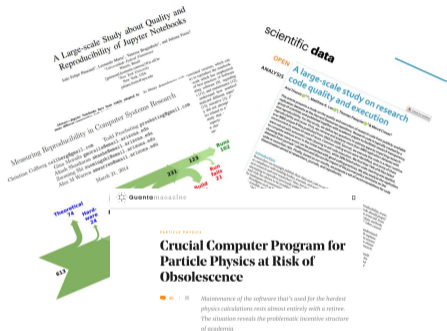


(articles: [here](#), [here](#), [here](#) and [here](#))

How are we managing our (open source) software?

Reproducibility, maintenance in Academia

Security, integrity, traceability in Industry



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



Can they track the software that they

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

We need a *dedicated infrastructure* to adress *all* this!

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion



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

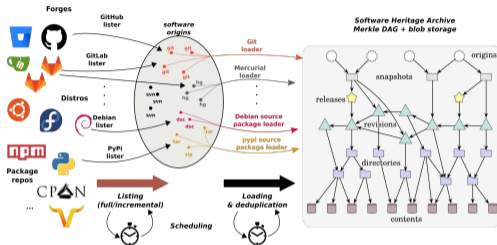
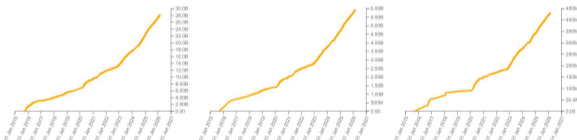
Research



Projects

429,963,770

Public Administration



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

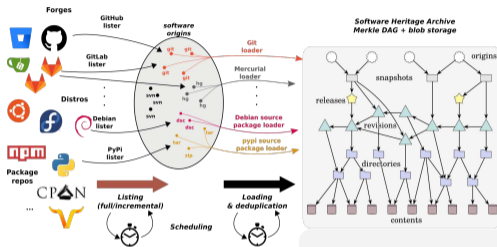
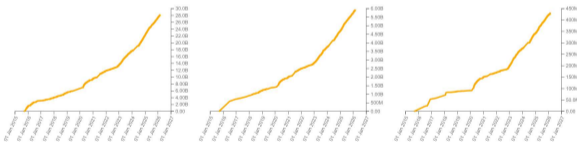
Research



Projects

429,963,770

Public Administration



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

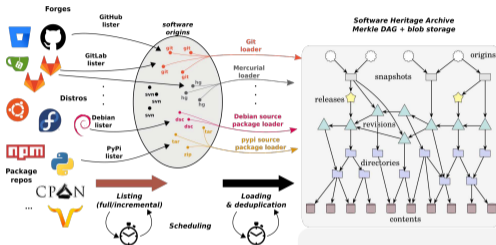
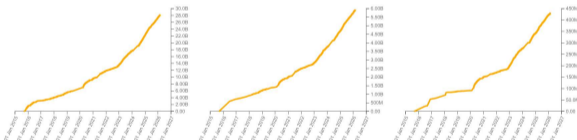


Public Administration



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



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

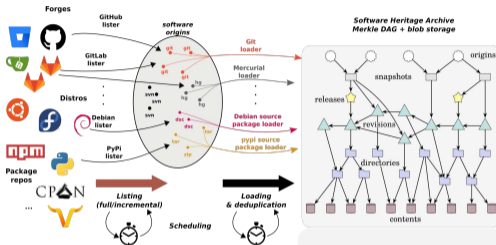
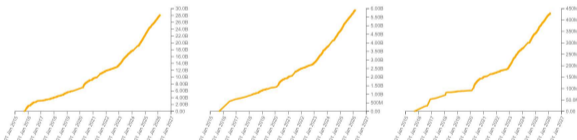


Public Administration



Projects

429,963,770



5000+ platforms

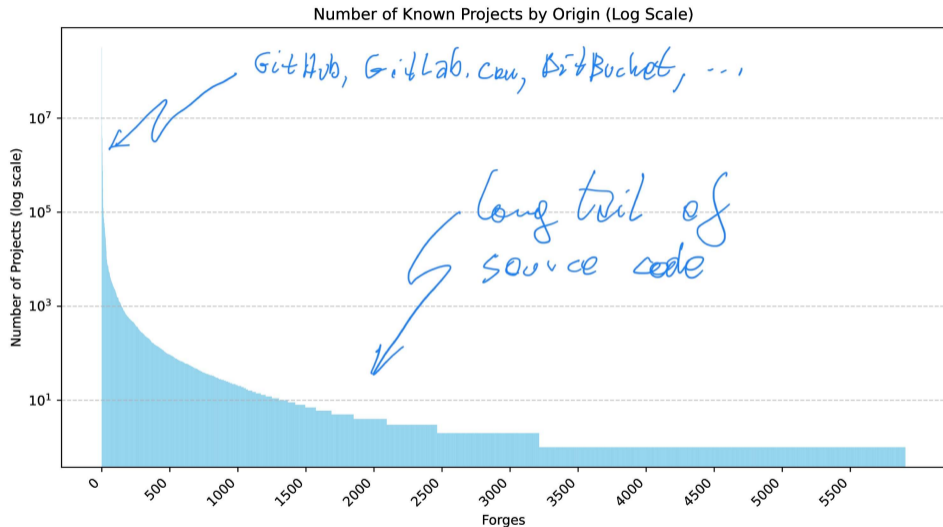
All versions, all history development in a single graph

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

A revolutionary infrastructure ensures availability guarantees integrity enables traceability



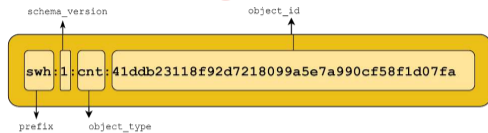
A peek at the code hosting landscape



Over 2 million projects in the long tail!

The Software Hash identifier (SWHID)

Software Heritage Identifiers (SWHID)



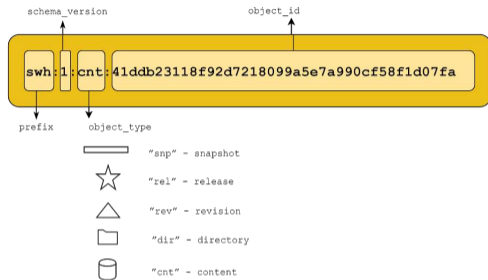
see swhid.org

50+B
intrinsic,
decentralised,
cryptographic

The Software Hash identifier (SWHID)

Software Heritage Identifiers (SWHID)

see swhid.org



50+B
intrinsic,
decentralised,
cryptographic

The Software Hash identifier (SWHID)

Software Heritage Identifiers (SWHID)

see swhid.org



50+B
intrinsic,
decentralised,
cryptographic

The Software Hash identifier (SWHID)

Software Heritage Identifiers (SWHID)

see swhid.org



50+B
intrinsic,
decentralised,
cryptographic

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

- Linux Foundation [SPDX 2.2](https://spdx.org/licenses/)
- IANA-registered "swh:"
- WikiData property [P6138](https://www.wikidata.org/wiki/P6138)

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

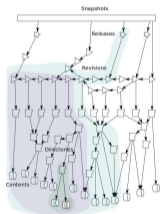
[ISO/IEC 18670](#), see swhid.org

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure**
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

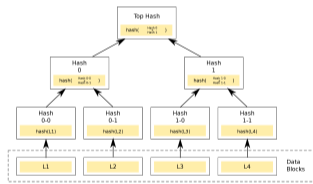
From archive to revolutionary infrastructure

Modern "Library of Alexandria", *international, non profit, long term* initiative addressing the needs of *industry, research, culture and society as a whole*

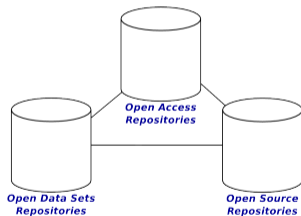
Software Graph



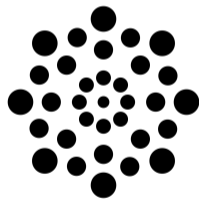
Software Blockchain



Open Science pillar



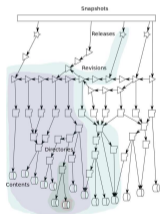
Big Code



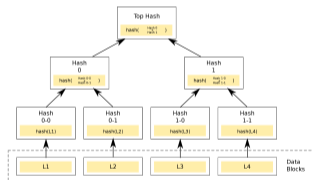
From archive to revolutionary infrastructure

Modern "Library of Alexandria", *international, non profit, long term* initiative addressing the needs of *industry, research, culture and society as a whole*

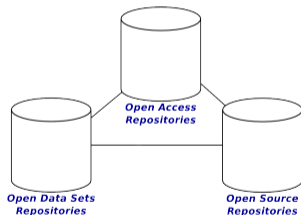
Software Graph



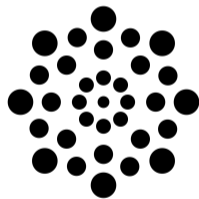
Software Blockchain



Open Science pillar



Big Code



Today we focus on

Open Science • **Large-scale analysis** • **Cybersecurity** • **AI**

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive**
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

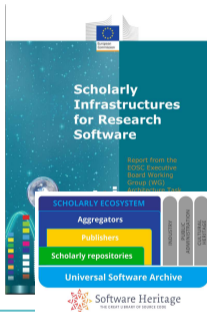
An example is worth a thousand words

- Browse (e.g. [Apollo 11 \[excerpt\]](#), your work [may be already there](#) !)
- Trigger archival, use the [updateswh](#) browser extension, configure the [webhooks](#)
- Get and use SWHIDs ([full specification available online](#))
- Cite software with [biblatex-software](#) package from CTAN
 - [Overleaf ACMART template](#) available
- Example in journals: [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 CNES](#), [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 research articles:
 - compare Fig. 1 and conclusions in [the 2012 version](#) and [the updated version](#)
 - SWHID in [a replication experiment](#)

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)



Projects



FAIRCORE4EOSC
Core Components Supporting a FAIR EOSC

The CodeMeta Project



FAIR-IMPACT
Expanding FAIR solutions across EOSC

Artifact Evaluation Committees

- current status
 - many guidelines developed in the 2010's ignore Software Heritage
 - ACM DL or Zenodo deposit of a zip + DOI still required
- state of the art
 - archive on Software Heritage
 - use the qualified SWHID

Artifact Evaluation Committees

- current status
 - many guidelines developed in the 2010's ignore Software Heritage
 - ACM DL or Zenodo deposit of a zip + DOI still required
- state of the art
 - archive on Software Heritage
 - use the qualified SWHID

Software citation

- current status
 - cursory mention in the text, sometimes citation of a research article, or manual ad hoc entries
- state of the art
 - include codemeta.json (or citation.cff, better if generated from codemeta.json)
 - generate biblatex entry, and use biblatex-software (CTAN, acmart.cls, etc.)

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive**
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

<https://registry.opendata.aws/software-heritage/>

Registry of Open Data on AWS



Software Heritage Graph Dataset

[digital preservation](#) [free software](#) [open source software](#) [source code](#)

Description

[Software Heritage](#) is the largest existing public archive of software source code and accompanying development history. The Software Heritage Graph Dataset is a fully deduplicated Merkle DAG representation of the Software Heritage archive. The dataset links together file content identifiers, source code directories, Version Control System (VCS) commits tracking evolution over time, up to the full states of VCS repositories as observed by Software Heritage during periodic crawls. The dataset's contents come from major development forges (including GitHub and GitLab), FOSS distributions (e.g., Debian), and language-specific package managers (e.g., PyPI). Crawling information is also included, providing timestamps about when and where all archived source code artifacts have been observed in the wild.

Update Frequency

Data is updated yearly

License

Creative Commons Attribution 4.0 International. By accessing the dataset, you agree with the Software Heritage [Ethical Charter for using the archive data](#) and the [terms of use for bulk access](#).

Documentation

<https://docs.softwareheritage.org/devel/swh-dataset/graph/athena.html>

Managed By

Software Heritage

See all datasets managed by [Software Heritage](#).

Contact

R. Di Cosmo roberto@dicosmo.org (CC-BY 4.0)

Resources on AWS

Description

Software Heritage Graph Dataset

Resource type

S3 Bucket

Amazon Resource Name (ARN)

```
arn:aws:s3:::softwareheritage
```

AWS Region

```
us-east-1
```

AWS CLI Access (No AWS account required)

```
aws s3 ls --no-sign-request s3://softwareheritage/
```

Description

S3 Inventory files

Resource type

S3 Bucket

Amazon Resource Name (ARN)

```
arn:aws:s3:::softwareheritage-inventory
```

AWS Region

```
us-east-1
```

AWS CLI Access (No AWS account required)

```
aws s3 ls --no-sign-request s3://softwareheritage-inventory/
```

Software Heritage for SE, CyberSecurity and AI

May 21st, 2026

11 / 41

Example: most popular commit verbs (stemmed)

Query using Amazon Athena

```
SELECT COUNT(*) AS C, word FROM (  
    SELECT word_stem(lower(split_part(  
        trim(from_utf8(message)), ' ', 1)))  
    AS word FROM revision  
    WHERE length(message) < 1000000)  
WHERE word != ''  
GROUP BY word  
ORDER BY C  
DESC LIMIT 20;
```

Example: most popular commit verbs (stemmed)

Query using Amazon Athena

```
SELECT COUNT(*) AS C, word FROM (  
  SELECT word_stem(lower(split_part(  
    trim(from_utf8(message)), ' ', 1)))  
  AS word FROM revision  
  WHERE length(message) < 1000000)  
WHERE word != ''  
GROUP BY word  
ORDER BY C  
DESC LIMIT 20;
```

Results

Completed Time in queue: 272 ms Run time: 33.545 sec Data scanned: 94.51 GB

Results (20) [Copy](#) [Download results](#)

| # | c | word |
|----|-----------|--------|
| 1 | 271573294 | updat |
| 2 | 163328012 | merg |
| 3 | 140044381 | add |
| 4 | 105800317 | fix |
| 5 | 103646653 | ad |
| 6 | 52891401 | bump |
| 7 | 50067041 | initi |
| 8 | 45609622 | creat |
| 9 | 42633225 | remov |
| 10 | 32230842 | chang |
| 11 | 23110410 | delet |
| 12 | 20734745 | new |
| 13 | 16644508 | commit |
| 14 | 15651821 | test |

State-of-the-art graph compression from social networks



Paolo Boldi, Antoine Pietri, Sebastiano Vigna, Stefano Zacchioli

Ultra-Large-Scale Repository Analysis via Graph Compression

SANER 2020, 27th Intl. Conf. on Software Analysis, Evolution and Reengineering. IEEE

Results

Full graph structure (50 B nodes, 800 B edges) in 400 GiB RAM

- traversal time is tens of ns per edge
- bidirectional traversals implemented
- **beware:** metadata access is still *off RAM*

Java Rust and gRPC APIs available ...

docs.softwareheritage.org/devel/swh-graph/grpc-api.html

What it does

- User-space POSIX filesystem (FUSE) exposing the **entire SWH archive** as a directory tree.
- Navigate by SWHID:
archive/swh:1:rev:...,
archive/swh:1:dir:...,
archive/swh:1:cnt:...
- New backend on top of **sw-h-graph**: traverse history, parents, snapshots at full graph speed.
- Lazy content fetch + local cache — you only pay for what you read.

Use case: **Whisper team / Coccinelle**

Inria Whisper (LIP6, Sorbonne Université)

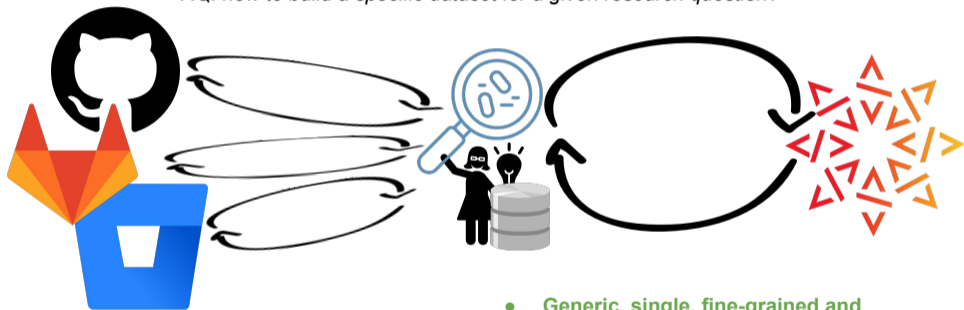
Coccinelle — semantic-patch engine for the Linux kernel. With sw-h-fuse, Coccinelle can run *across the full history of public code*, not just the latest checkout: study patch propagation, vulnerability fix coverage, code-evolution patterns at SWH scale.

docs.softwareheritage.org/devel/sw-h-fuse • gitlab.softwareheritage.org/sw-h/devel/sw-h-fuse

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies**
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

Mining Android Applications on Software Heritage

RQ: how to build a specific dataset for a given research question?



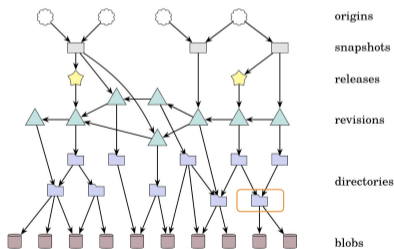
- **Specific and limited API**
- **Hardly reproducible**

- **Generic, single, fine-grained and unlimited API**
- **Growing number of source codes**
- **Easy to update the dataset**

(from the Inria/IRISA DiverSE team)

Using the SWH merkle dag to identify android repositories

Identify android application repositories = Find the AndroidManifest.xml among the sources

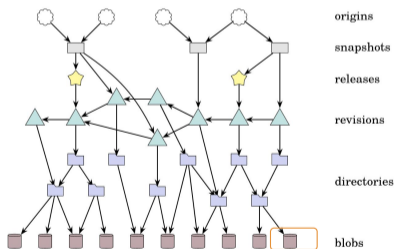


SWH Merkle DAG, Antoine Pietri

1) Iterate over the graph nodes until you find a directory node containing a file named "AndroidManifest.xml".

Using the SWH merkle dag to identify android repositories

Identify android application repositories = Find the AndroidManifest.xml among the sources

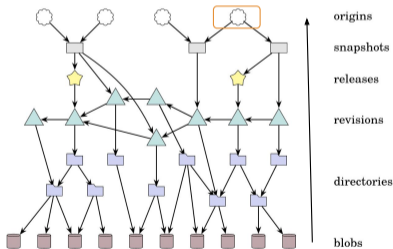


SWH Merkle DAG, Antoine Pietri

2) Extract the SWH identifier of the blob corresponding to the AndroidManifest.xml and download the corresponding file through the SWH Web API

Using the SWH merkle dag to identify android repositories

Identify android application repositories = Find the AndroidManifest.xml among the sources



SWH Merkle DAG, Antoine Pietri

3) Traverse the graph in backward direction to the origin node and get the repository url

Broad variety of sources in *one open dataset*

reduces usual GH bias

Reference simple *standard data format*

VCS and forge details are abstracted away

Simplifies reproducibility packages

no need to create a full copy, *just list the SWHIDs!*

Software Heritage does the heavy lifting for you

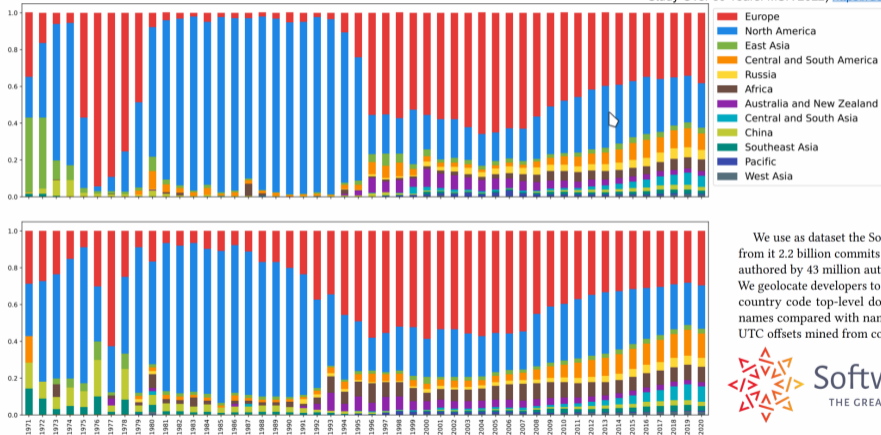
no need to scrape/download repositories all over again

Example: (Open) Source Code comes from all over the world

MSR '22, May 23–24, 2022, Pittsburgh, PA, USA

Davide Rossi and Stefano Zacchiroli

Geographic Diversity in Public Code Contributions: An Exploratory Large-Scale Study Over 50 Years. MSR 2022) <https://doi.org/10.1145/3524842.3528471>



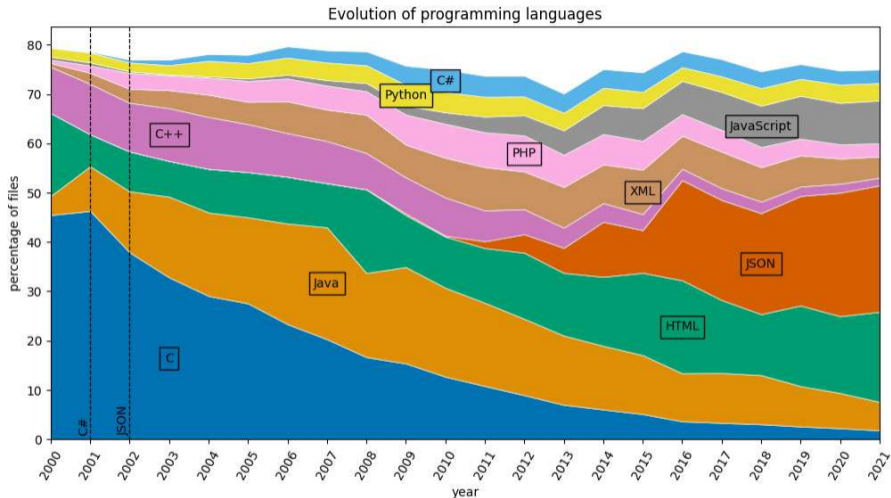
We use as dataset the Software Heritage archive [3] and analyze from it 2.2 billion commits archived from 160 million projects and authored by 43 million authors during the 1971–2021 time period. We geolocate developers to 12 world regions, using as signals email country code top-level domains (ccTLDs) and author (first/last) names compared with name distributions around the world, and UTC offsets mined from commit metadata.



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

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

Example: Programming language evolution in (Open) Source



A. Desmazières, R. Di Cosmo, V. Lorentz

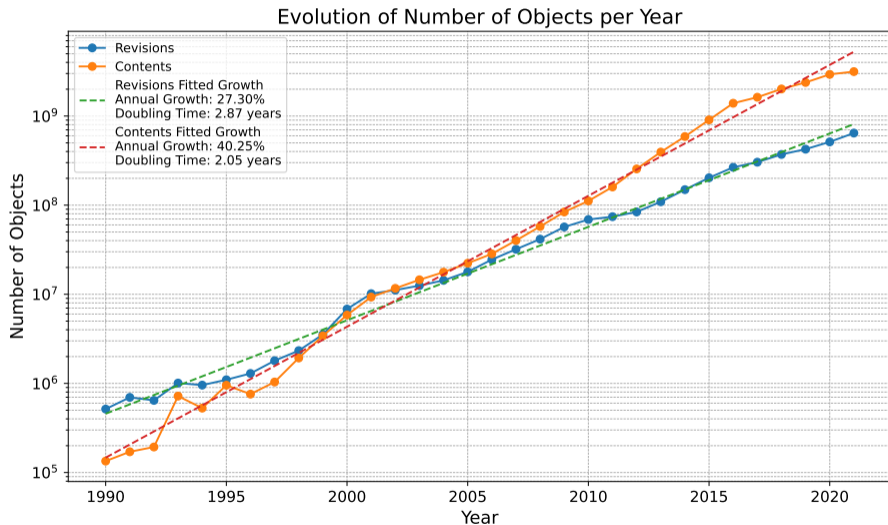
50 years of programming language evolution through the Software Heritage Looking Glass

International Conference on Mining Software Repositories (MSR) 2025.

<https://hal.science/hal-04924849/>



Example: (Open) Source Code grows at an exponential rate



Software
Heritage

Get the data
and the code



Question:

can we leverage the Software Heritage graph to map contributions from a research institution in the galaxy of software development?

Question:

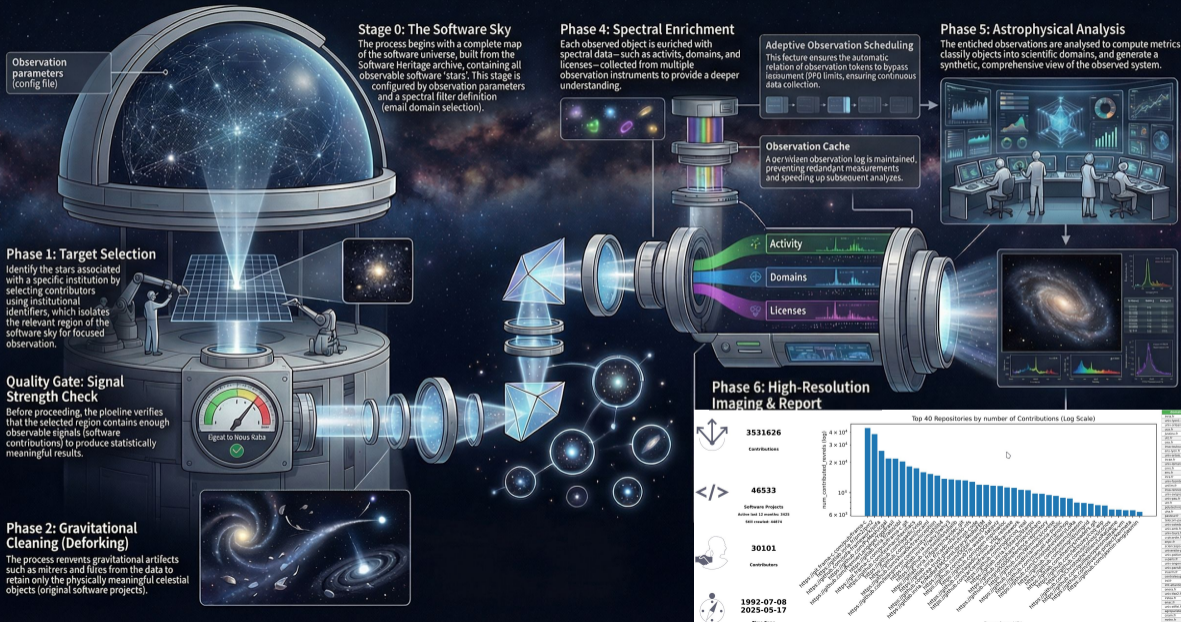
can we leverage the Software Heritage graph to map contributions from a research institution in the galaxy of software development?

Answer: let's try

- **Extract:** projects and contributions using email domain matching in the Software Heritage graph.
- **Deduplicate:** Cluster forks/copies using commit metadata, select a canonical repository.
- **Enrich:** Add GitHub metadata (topics, README, etc.) for context and visibility.
- **Analyze:** Compute key metrics on impact, contributor domains, and research relevance.

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy**
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

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



Observation parameters (config file)

Stage 0: The Software Sky

The process begins with a complete map of the software universe, built from the Software Heritage archive, containing all observable software 'stars'. This stage is configured by observation parameters and a spectral filter definition (email domain selection).

Phase 1: Target Selection

Identify the stars associated with a specific institution by selecting contributors using institutional identifiers, which isolates the relevant region of the software sky for focused observation.

Quality Gate: Signal Strength Check

Before proceeding, the pipeline verifies that the selected region contains enough observable signals (software contributions) to produce statistically meaningful results.



Phase 2: Gravitational Cleaning (Deforking)

The process renvents gravitational artifacts such as mirrors and forks from the data to retain only the physically meaningful celestial objects (original software projects).



Phase 4: Spectral Enrichment

Each observed object is enriched with spectral data—such as activities, domains, and licenses—collected from multiple observation instruments to provide a deeper understanding.



Adaptive Observation Scheduling

This feature ensures the automatic relation of observation tokens to bypass instrument (PPO) limits, ensuring continuous data collection.

Observation Cache

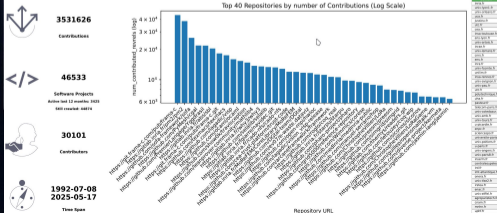
A per-woken observation log is maintained, preventing redundant measurements and speeding up subsequent analyzes.

Phase 6: High-Resolution Imaging & Report



Phase 5: Astrophysical Analysis

The enticed observations are analysed to compute metrics, classify objects into scientific domains, and generate a synthetic, comprehensive view of the observed system.



First SWH VLT pictures : SU contribution to Open Source



155686

Contributions



2739

Software Projects

Active last 12 months: 91

Still crawled: 2469



1096

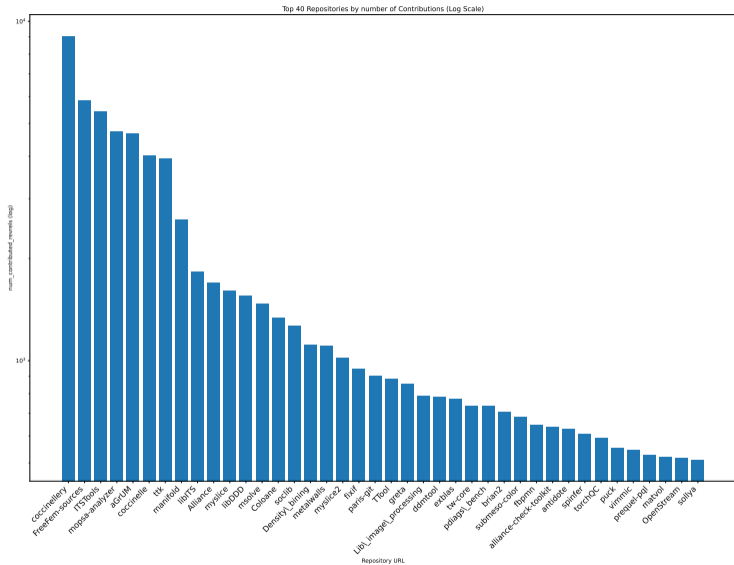
Contributors



1998-10-02
2026-03-04

Time Span

First SWH VLT pictures: SU contribution to Open Source



- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)**
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion

Programme

France 2030 PTCC — 5 years (2024–2028) — 8 research teams

Co-leaders: Olivier Barais (Inria, DiverSE) · Stefano Zacchiroli (Télécom Paris, ACES)

Partners — academic

- **APR** — LIP6, **Sorbonne Université** (*prog. languages, static analysis MOPSA, verification*)
- **ACES** — Télécom Paris (*cybersecurity, software supply chain*)
- **CEA List** (*static analysis: Frama-C, Binsec*)

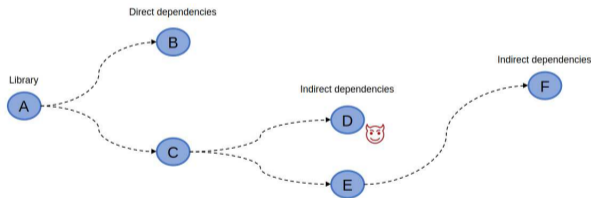
Partners — Inria

- **DiverSE** — supply chain security
- **RMOD** — reverse engineering (Moose)
- **Software Heritage** — archive, infrastructure
- **Spirals** — self-protective systems, web security
- **Whisper** — Linux kernel, Coccinelle

Software supply chain attacks

Reusing OSS via dependencies

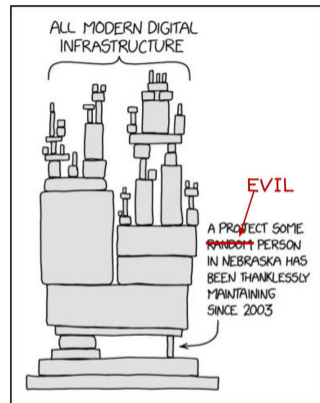
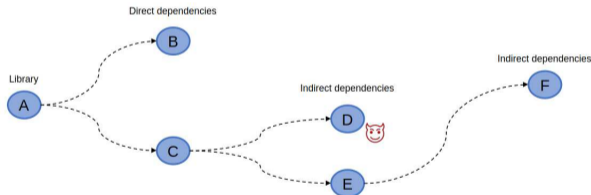
- **Software dependencies:** a popular way of reusing open source software.
- Software product *A* uses functionalities implemented in OSS product *B* ... and so on.



Software supply chain attacks

Reusing OSS via dependencies

- **Software dependencies**: a popular way of reusing open source software.
- Software product *A* uses functionalities implemented in OSS product *B* ... and so on.



based on xkcd.com/2347

Attacking the software supply chain

- Attacking **undermaintained "leaf" packages** (e.g., D) → efficient attack strategy
- Many documented attacks: event-stream (2018), node-ipc (2022), XZ utils (2024), ...

An universal knowledge base about public code vulnerabilities

Vision

- **Software Heritage** is the perfect (and only) place where to build an universal knowledge base that maps known vulnerabilities to public code artifacts.
- SWH can provide an **open data API mapping SWHIDs to CVEs**, that knows about *all public commits* and can be leveraged to increase open source security.

EU Cyber Resilience Act (CRA)

- Key helper to abide to CRA obligations, coming into effect ~Q3 2026.
- SWH funding member of the **Open Regulatory Compliance Working Group**.



Roadmap

- Context: SWHSec project (PTCC-funded, 2023-2027).
- Current status: working prototype that processes OSV.dev data and use it to "color" the entire SWH commit graph (~5 billion commits) with vulnerability information.



Syful Islam, Stefano Zacchioli.

On the Informativeness of Security Commit Messages: A Large-scale Replication Study.

To appear in EASE 2026, ACM, Glasgow, UK.



Syful Islam, Stefano Zacchioli.

On the Use of Commit Messages for Corrective Software Maintenance: A Systematic Mapping Study.

To appear in EASE 2026, ACM, Glasgow, UK.



Romain Robbes, Théo Matricon, Thomas Degueule, André Hora, Stefano Zacchioli.

Promises, Perils, and (Timely) Heuristics for Mining Coding Agent Activity.

To appear in MSR 2026, ACM, Rio de Janeiro, Brazil.



Luís Soeiro, Thomas Robert, Stefano Zacchioli.

Finding Software Supply Chain Attack Paths with Logical Attack Graphs.

FPS 2025, Brest, France. Springer, 2026.



Solal Rapaport, Laurent Pautet, Samuel Tardieu, Stefano Zacchioli.

Altered Histories in Version Control System Repositories: Evidence from the Trenches.

ASE 2025, IEEE, Seoul, South Korea. Pages 2183–2194.

- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs**
- 10 One last thing
- 11 Conclusion

Software Heritage and Generative AI, first contacts

October 19, 2023

Software Heritage Statement on Large Language Models for Code



Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting machine learning models must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The initial training data extracted from the Software Heritage archive must be fully and precisely identified by, for example, publishing the corresponding SWHID identifiers (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Software Heritage and Generative AI, first contacts

October 19, 2023

Software Heritage Statement on Large Language Models for Code

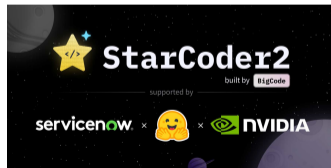
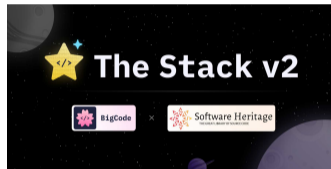


Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting machine learning models must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The initial training data extracted from the Software Heritage archive must be fully and precisely identified by, for example, publishing the corresponding SWHID identifiers (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

February 2024

Yes, it's possible!



Software Heritage and Generative AI, first contacts

October 19, 2023

Software Heritage Statement on Large Language Models for Code

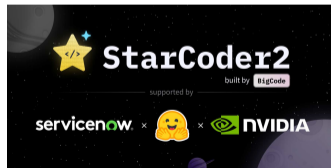
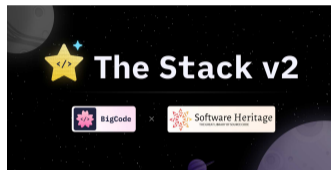


Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting machine learning models must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The initial training data extracted from the Software Heritage archive must be fully and precisely identified by, for example, publishing the corresponding SWHID identifiers (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

February 2024

Yes, it's possible!



But it's hard...

Lessons learned

Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data extracted from the Software Heritage archive* must be fully and precisely identified by, for example, publishing the corresponding **SWHID identifiers** (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Lessons learned

Transparency is easy: [SWHID](#) (undergoing ISO standardisation) and Software Heritage
N.B. : may be mandated by regulations!

Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data extracted from the Software Heritage archive* must be fully and precisely identified by, for example, publishing the corresponding **SWHID identifiers** (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Lessons learned

Transparency is easy: [SWHID](#) (undergoing ISO standardisation) and Software Heritage
N.B. : may be mandated by regulations!

Opt out is complex: who is *the real right owner*?
(similar issues to license compliance)

Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data extracted from the Software Heritage archive* must be fully and precisely identified by, for example, publishing the corresponding **SWHID identifiers** (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Lessons learned

Transparency is easy: [SWHID](#) (undergoing ISO standardisation) and Software Heritage
N.B. : may be mandated by regulations!

Opt out is complex: who is *the real right owner?*
(similar issues to license compliance)



Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data extracted from the Software Heritage archive* must be fully and precisely identified by, for example, publishing the corresponding **SWHID identifiers** (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Lessons learned

Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data* extracted from the Software Heritage archive must be fully and precisely identified by, for example, publishing the corresponding SWHID identifiers (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Transparency is easy: [SWHID](#) (undergoing ISO standardisation) and Software Heritage
N.B. : may be mandated by regulations!

Opt out is complex: who is *the real right owner*?
(similar issues to license compliance)

● **Building the training set is complex:** e.g. includes **license compliance** alike work **at massive scale**



Lessons learned

Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data* extracted from the Software Heritage archive must be fully and precisely identified by, for example, publishing the corresponding SWHID identifiers (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Transparency is easy: [SWHID](#) (undergoing ISO standardisation) and Software Heritage
N.B. : may be mandated by regulations!

Opt out is complex: who is *the real right owner*?
(similar issues to license compliance)



- **Building the training set is complex:** e.g. includes **license compliance** alike work **at massive scale**
- **Generating attribution information on model output is more complex** than license compliance

Lessons learned

Principles

1. Knowledge derived from the Software Heritage archive must be given back to humanity, rather than monopolized for private gain. The resulting *machine learning models* must be made available under a suitable open license, together with the documentation and toolings needed to use them.
2. The *initial training data* extracted from the Software Heritage archive must be fully and precisely identified by, for example, publishing the corresponding SWHID identifiers (note that, in the context of Software Heritage, public availability of the *initial training data* is a given: anyone can obtain it from the archive). This will enable use cases such as: studying biases (fairness), verifying if a code of interest was present in the training data (transparency), and providing appropriate attribution when generated code bears resemblance to training data (credit), among others.
3. Mechanisms should be established, where possible, for authors to exclude their archived code from the training inputs before model training begins.

Transparency is easy: [SWHID](#) (undergoing ISO standardisation) and Software Heritage
N.B. : may be mandated by regulations!

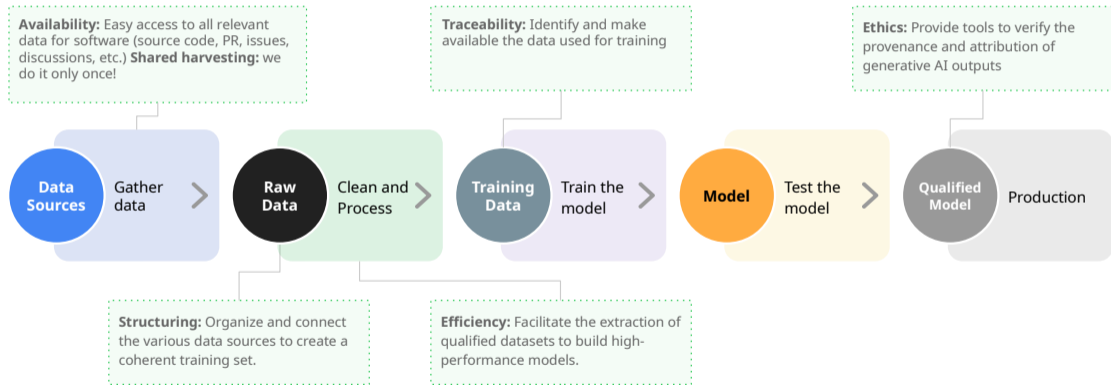
Opt out is complex: who is *the real right owner?*
(similar issues to license compliance)



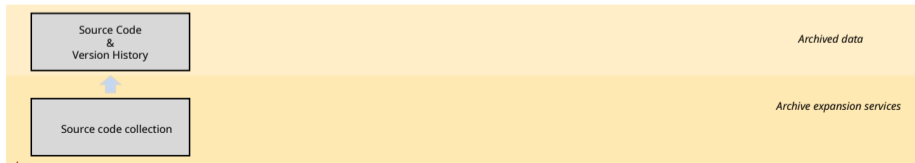
- **Building the training set is complex:** e.g. includes **license compliance** alike work **at massive scale**
- **Generating attribution information on model output is more complex** than license compliance

We need a coordinated effort to ensure fully open models will succeed!

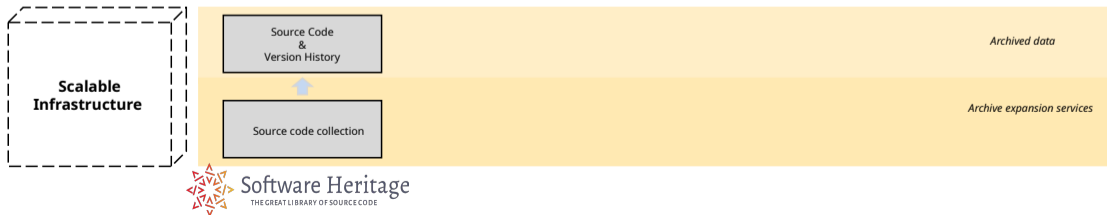
A STEP FORWARD: CodeCommons



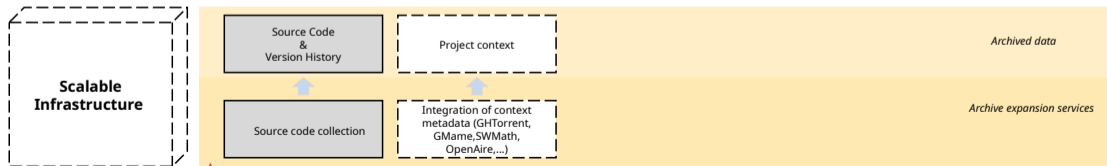
CodeCommons: bird's eye view (technical focus)



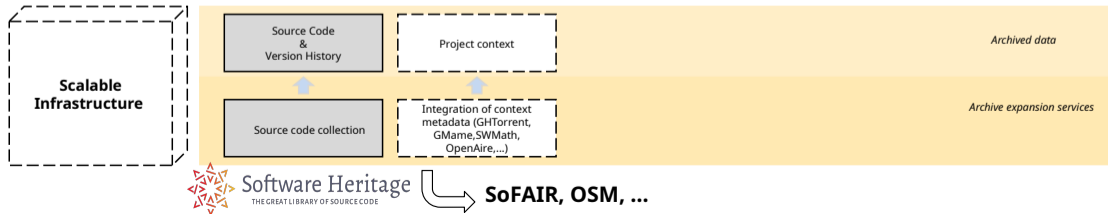
CodeCommons: bird's eye view (technical focus)



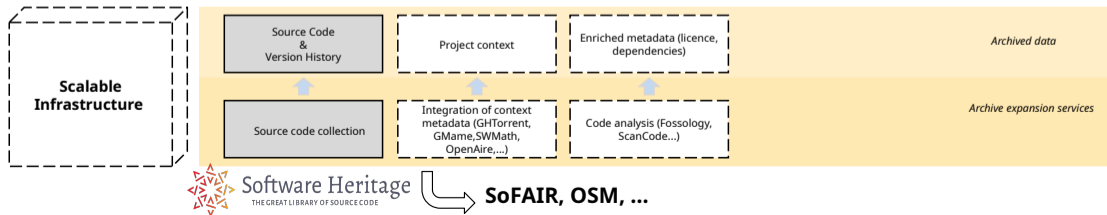
CodeCommons: bird's eye view (technical focus)



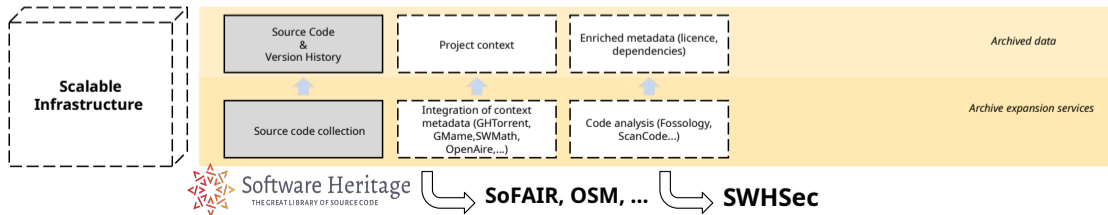
CodeCommons: bird's eye view (technical focus)



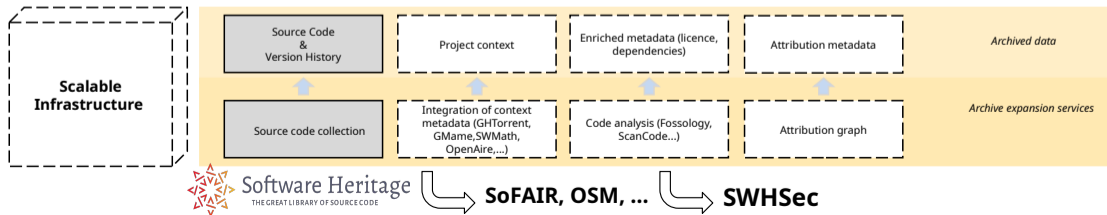
CodeCommons: bird's eye view (technical focus)



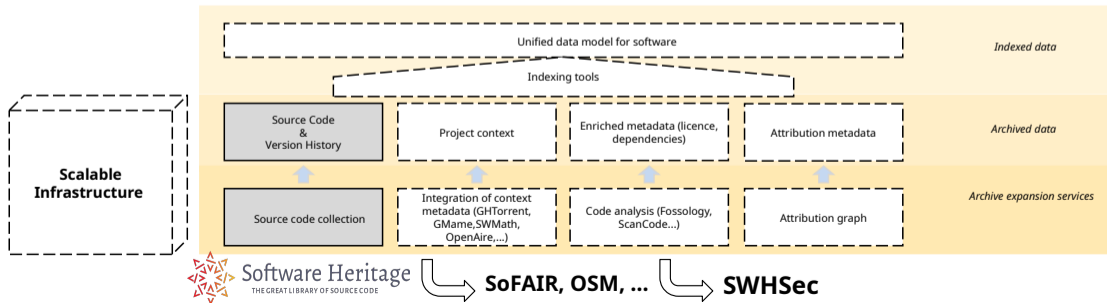
CodeCommons: bird's eye view (technical focus)



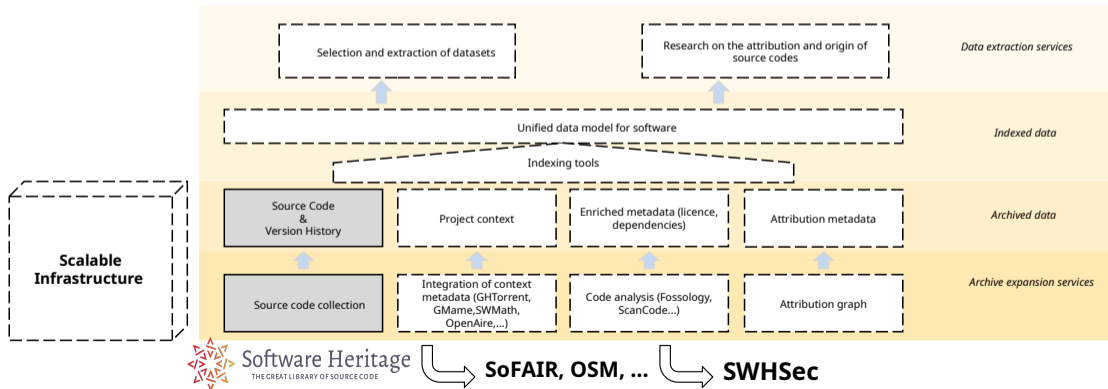
CodeCommons: bird's eye view (technical focus)



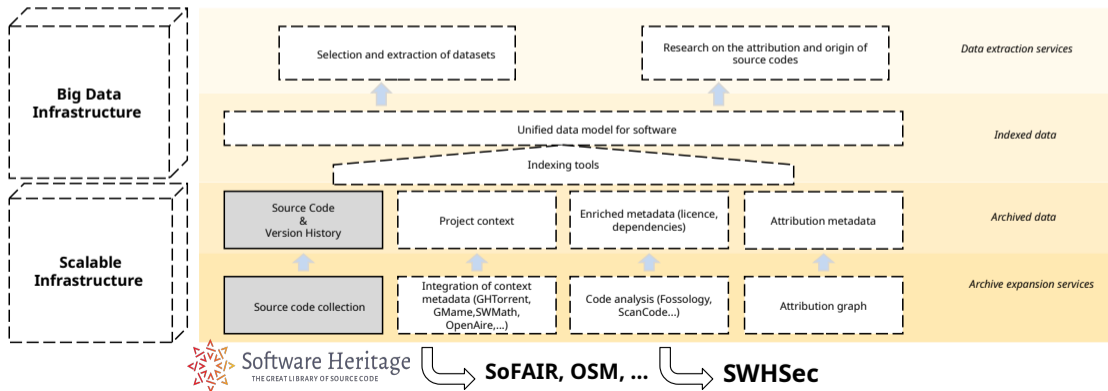
CodeCommons: bird's eye view (technical focus)



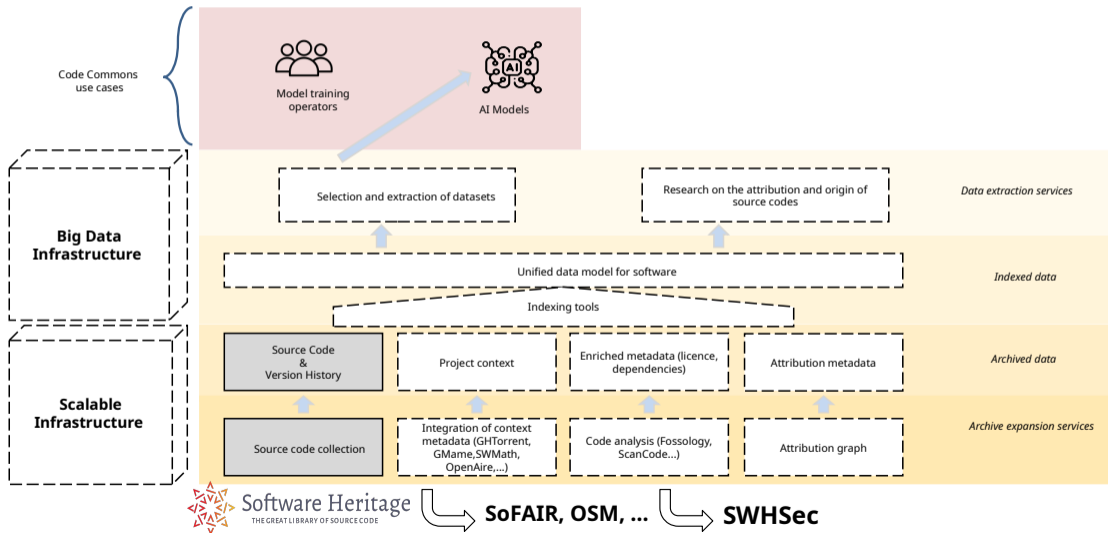
CodeCommons: bird's eye view (technical focus)



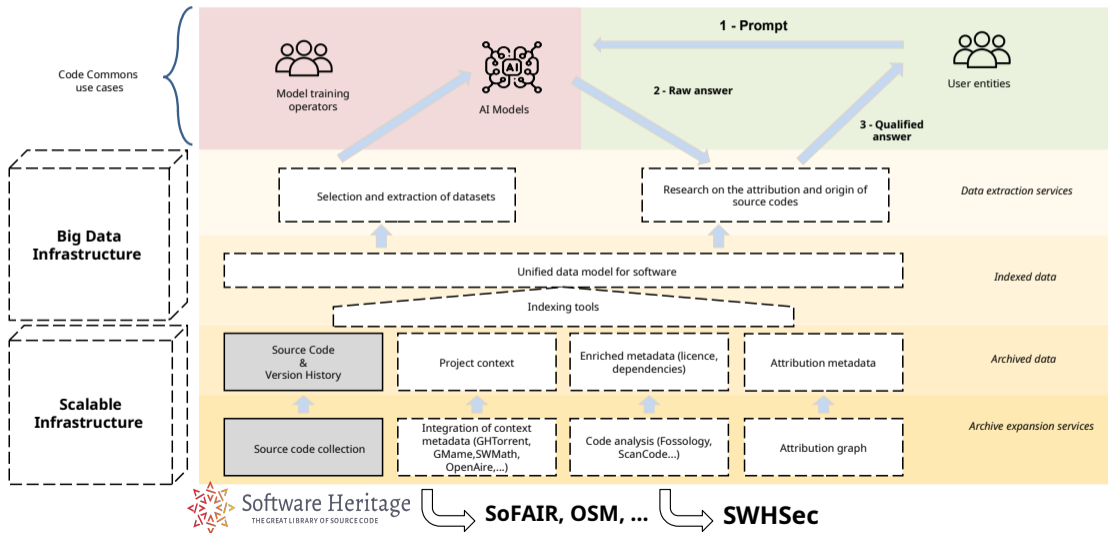
CodeCommons: bird's eye view (technical focus)



CodeCommons: bird's eye view (technical focus)



CodeCommons: bird's eye view (technical focus)



CodeCommons

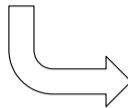
Open, responsible, and transparent AI: Our shared goal

CodeCommons is an ambitious project to create the world's most comprehensive digital commons for code

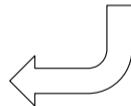
Building on the existing foundation of Software Heritage, the largest publicly available source code archive, CodeCommons aims to bring into one place all the **critical** and **qualified** information needed to create **smaller, better** datasets for the next generation of AI tools.

At its core, the project prioritizes transparency and traceability, enabling model builders and users to **respect creators' rights** while promoting **sovereign** and **sustainable** AI.

Learn more



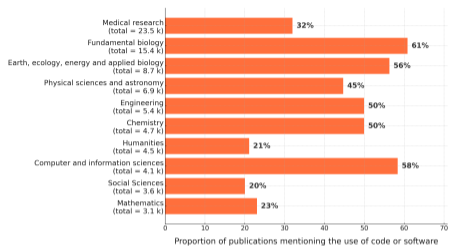
Meet the teams



- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing**
- 11 Conclusion

Software is a global undertaking...

Pillar of Science across all research areas



From all continents

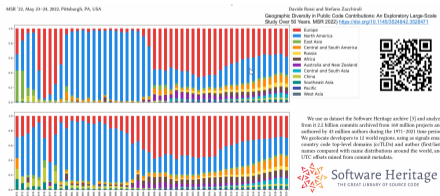
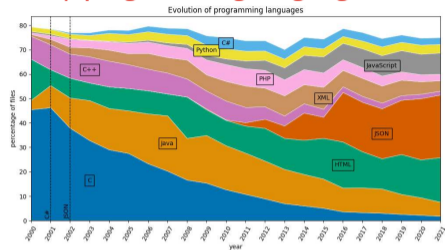
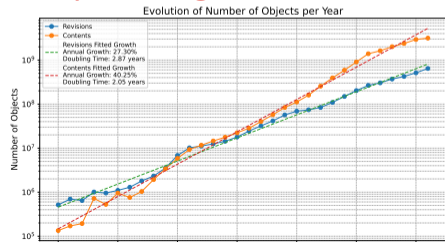


Figure 3: Ratio of commits (above) and active authors (below) by world zone over the 1971-2020 period.

Many programming languages

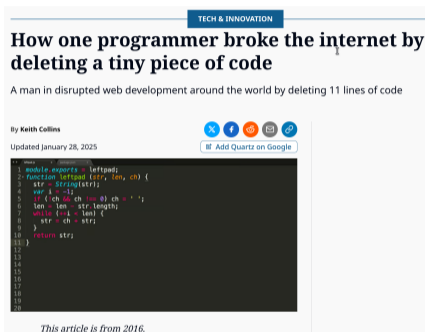


An exponential growth



... with deep (infra)structural dependencies!

An early warning: left-pad (March 2016)



See the [Wikipedia](#) article:

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

Resilience, demonstrated: Guix (and Nix) 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!



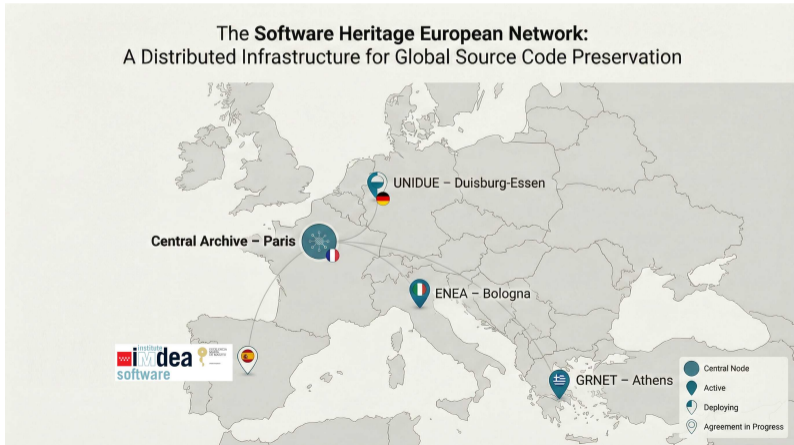
Since 2019

An ecumenical approach to sovereign fallback of build chains is possible

Extend to key package managers

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

Realize ecumenical resilience



- 1 Introduction
- 2 Meet Software Heritage
- 3 From archive to revolutionary infrastructure
- 4 Software for research: reference archive
- 5 Research on software: datasets from the archive
- 6 Selected highlight: Impact on software studies
- 7 A Very Large Telescope to explore the software development galaxy
- 8 Selected highlight: Improving Open Source Security with SWH(Sec)
- 9 Selected highlight: AI and transparent LLMs
- 10 One last thing
- 11 Conclusion**

A growing and active community

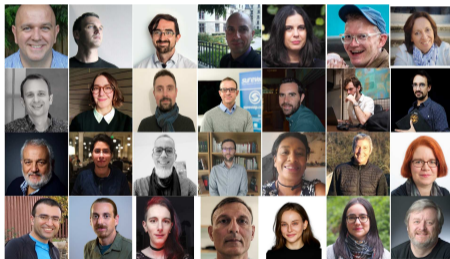
Core Team



All together, 2026 Symposium



Ambassadors



A rally flag for a grand vision

Bring together academia, industry, governments, communities

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

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

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

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

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

Software Heritage enables

- archival, reference, integrity
- qualification, sharing and reuse

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

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

Call to action

- leverage for AI, Cyber, SwEng
- get involved research, code

Software Heritage enables

- archival, reference, integrity
- qualification, sharing and reuse

- promote for AEC, Journals, etc.
- fund the long term mission

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

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

Call to action

- leverage for AI, Cyber, SwEng
- get involved research, code

Software Heritage enables

- archival, reference, integrity
- qualification, sharing and reuse

- promote for AEC, Journals, etc.
- fund the long term mission

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

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

Call to action

- leverage for AI, Cyber, SwEng
- get involved research, code

Software Heritage enables

- archival, reference, integrity
- qualification, sharing and reuse

- promote for AEC, Journals, etc.
- fund the long term mission

Annual report 2025



5 years in 5 minutes



Get these slides



Appendix

12 Resilience: closing the regulation / infrastructure gap

13 SWHID

14 Scanner

15 SWHSec

16 Demo time!

Security

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

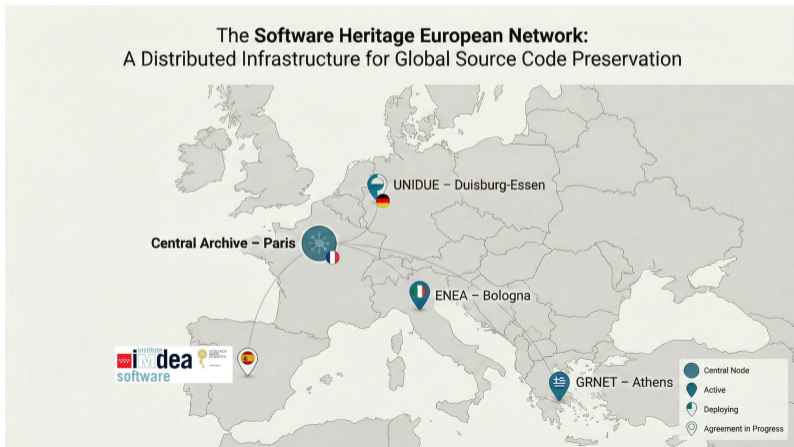
"Critical dependencies outside European control."

A sovereign fallback for Europe's build chains is possible

Extend to key package managers

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

Realize ecumenical resilience



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

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**

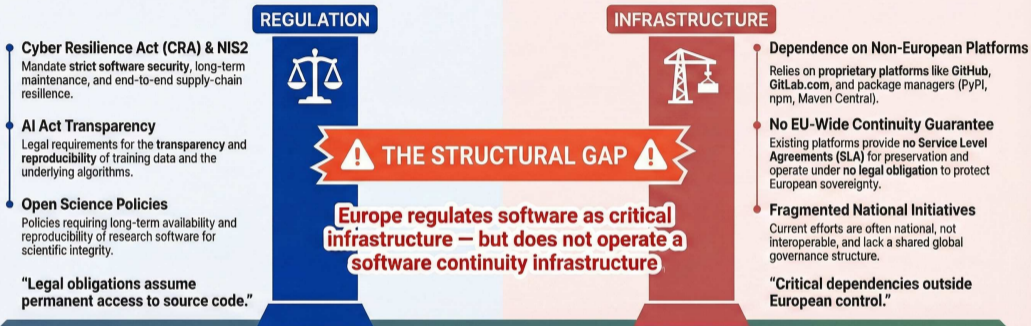
Transparent AI

- **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



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.”



Julien Malka, Théo Zimmermann, Stefano Zacchioli.

Does Functional Package Management Enable Reproducible Builds at Scale? Yes.

MSR 2025, IEEE, Ottawa, Canada. Pages 775–787.

ACM SIGSOFT Distinguished Paper Award.

12 Resilience: closing the regulation / infrastructure gap

13 SWHID

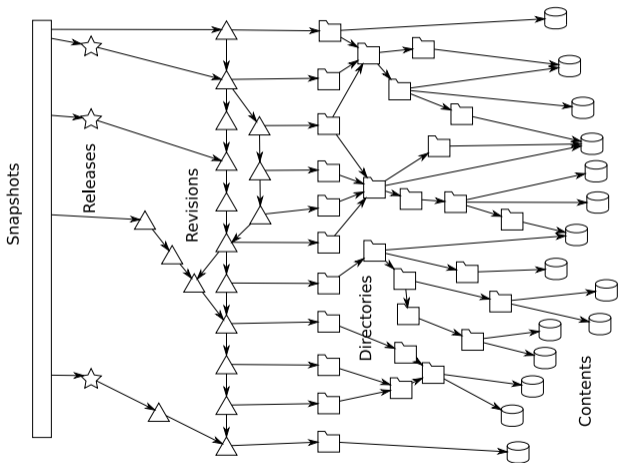
14 Scanner

15 SWHSec

16 Demo time!

Security

A giant (extended) Merkle DAG



A giant (extended) Merkle DAG

Contents

GNU GENERAL PUBLIC LICENSE
Version 3, 29 June 2007

Copyright (C) 2007 Free Software Foundation, Inc. <<http://fsf.org/>>
Everyone is permitted to copy and distribute verbatim copies
of this license document, but changing it is not allowed.

Preamble

The GNU General Public License is a free, copyleft license for
software and other kinds of works.

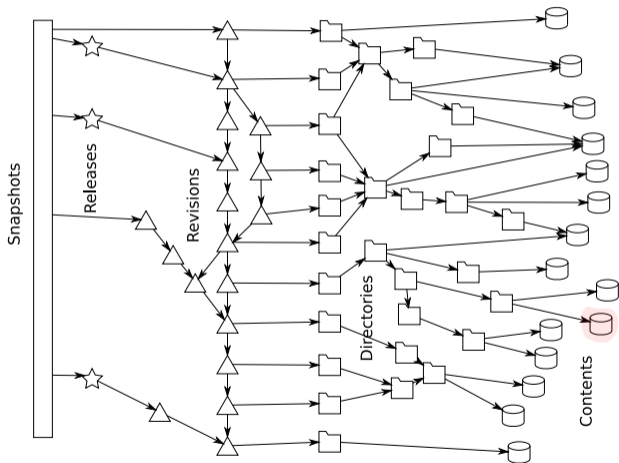
The licenses for most software and other practical works are designed
to take away your freedom to share and change the works. By contrast,
the GNU General Public License is intended to guarantee your freedom to
share and change all versions of a program--to make sure it remains free
software for all its users. We, the Free Software Foundation, use the
GNU General Public License for most of our software; it applies also to
any other work released this way by its authors. You can apply it to
your programs, too.

When we speak of free software, we are referring to freedom, not
price. Our General Public Licenses are designed to make sure that you
have the freedom to distribute copies of free software (and charge for
them if you wish), that you receive source code or can get it if you
want it, that you can change the software or use pieces of it in new
free programs, and that you know you can do these things.

To protect your rights, we need to prevent others from denying you
these rights if you have paid for software. Therefore, you have the right
to copy, modify and redistribute copies of free software for your
own private use, on the same conditions surrounding copy-
ing the original work.

sha1: 8624bcdae55baeef...
sha256: 8ceb4b9ee5aded...
sha1_git: 94a9ed024d385...
length: 35147

A giant (extended) Merkle DAG



A giant (extended) Merkle DAG

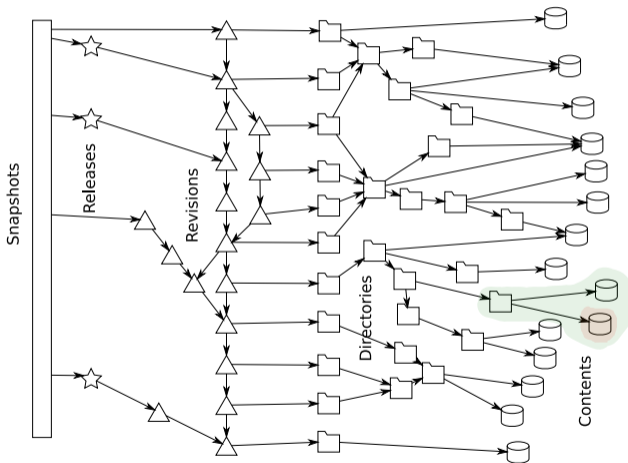


Directories

```
100644 blob c5baade4c44766042186ef858c0fd63d587ebf09 .gitignore
100644 blob 2d0a34af6f52cf3cf6b0c2f7bd0648fbd255e77f AUTHORS
100644 blob 94a9ed024d3859793618152ea559a168bbcbb5e2 LICENSE
100644 blob d9b2665a435a43f8a79a84e0867751dfb095c7bb MANIFEST.in
100644 blob 524175c2bad0b35b975f79284c2f5a6d5eaf2eb4 Makefile
100644 blob 5c7e3a5bbddb038682ba7793f440492ed9678bb3 Makefile.local
100644 blob 8617980629cd24e6080404f09aa749b085b3e07b README.db_testing
100644 blob 76b29f94cf815e0869c414d38d78d7ce08ec514e README.dev
040000 tree e1e10ecef948af0b93adb0372afc89f12e92618a bin
040000 tree 83e56d0beaf7793c77a45a34c80fcb8af503013 debian
040000 tree a34c9c4ba213f0cedc67f9816348d27955577af5 docs
100644 blob f2a6d32c6135aa7287bbd76167b01df2ae4f1539 requirements.txt
100755 blob eee147c36caf1bbc2d820da8dc026cb5b68180bc setup.py
040000 tree 224bb4c1f4c67fca1d160bff2d06094e7e1abf3 sql
040000 tree 8631c9cd77bbe993168107ab5baf51f40c6300be swh
040000 tree 8fb905b56ba8ed692f1209b2773b474c6c1d66c1 utils
```

id: 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d


A giant (extended) Merkle DAG



A giant (extended) Merkle DAG

Revisions

Details Changes Files

SHA: 963634dca6ba5dc37e3ee426ba091092c267f9f6 


Author: [Nicolas Dandrimont <nicolas@dandrimont.eu>](mailto:nicolas@dandrimont.eu) (Thu Sep 1 14:26:13 2016)

Committer: [Nicolas Dandrimont <nicolas@dandrimont.eu>](mailto:nicolas@dandrimont.eu) (Thu Sep 1 14:26:13 2016)

Subject: provenance.tasks: add the revision -> origin cache task

Parent: [fc3a8b59ca1df424d860f2c29ab07fee4dc35d10](https://sw.hq.mozilla.org/repo/v2/963634dca6ba5dc37e3ee426ba091092c267f9f6) : test_storage: property pipeline origin and cont...

provenance.tasks: add the revision -> origin cache task

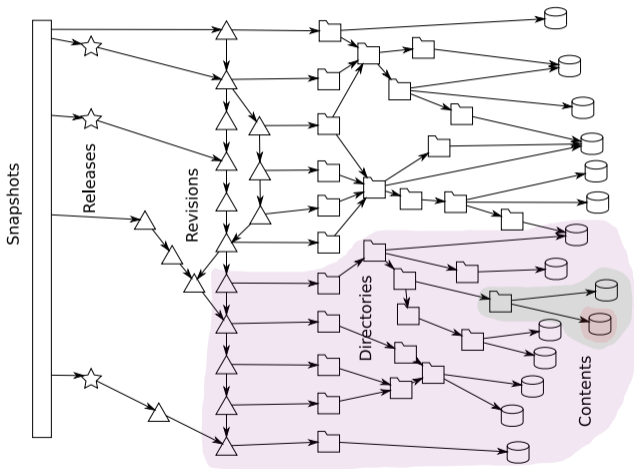
[sw/h/storage/provenance/tasks.py](#)  77

tree 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d
parent fc3a8b59ca1df424d860f2c29ab07fee4dc35d10
author Nicolas Dandrimont <nicolas@dandrimont.eu> 1472732773 +0200
committer Nicolas Dandrimont <nicolas@dandrimont.eu> 1472732773 +0200

provenance.tasks: add the revision -> origin cache task

id: 963634dca6ba5dc37e3ee426ba091092c267f9f6

A giant (extended) Merkle DAG



A giant (extended) Merkle DAG

Releases

tag v0.0.51
Tagger: Nicolas Dandrimont <nicolas@dandrimont.eu>
Date: Wed Aug 24 14:36:03 2016 +0200

Release swh.storage v0.0.51

- Add new metadata column to origin_visit
- Update swh-add-directory script for updated API
[...]

commit c0c9f16b1e134f593e7567570a1761b156e6eb1d

```
object c0c9f16b1e134f593e7567570a1761b156e6eb1d
type commit
tag v0.0.51
tagger Nicolas Dandrimont <nicolas@dandrimont.eu> 1472042163 +0200
```

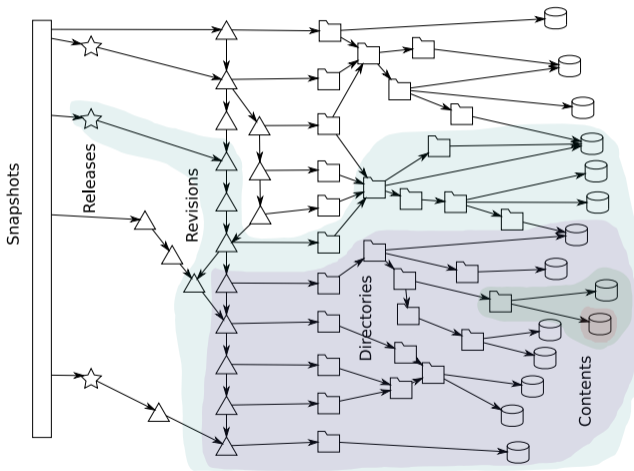
Release swh.storage v0.0.51

- Add new metadata column to origin_visit
- Update swh-add-directory script for updated API
-----BEGIN PGP SIGNATURE-----

```
iQIzBAABCAAdBQJXvZTNFhxuaWNvbGFzOGRhbmRyaW1vbnQuZlUACgkQ7AWLMO2+
neqorw//aq6S0b5DijzEa+kWN3rXgV5+1K1vEVh1wNKAwxBeKJ7aX2kEILDHt7uf
ahpZ6pz3q8nqs6aC1+YrxBFcih3L2YtrdZeWxWqr8xWNMaEoYDb8qaphwhBAD5t2
ICBIh2ujXuCrD193eKpWvzZxg+hB9sMWy35Dr6jWZ7K4MuJGgJyHPf55yo
lCEndWno7VfH1Vm6e1+5aB7l5mXRaqk+becqduhTZz.vjj+jpglRqC8cyqN3nmfL
qsjZmu8kyz3t8tG/H1jpv+15OwBlnPoS5TH0tujjEYqPKlghSP790QuHDH2FkCao
klj6kAWyU80Mxb+nKVjjeLbr3+yWBFj3Qp5a1/V8oOTh6E1dALCNMpEakCoKtMt
d/gMRax1l1/g0EDfnsW67G6sDwKPKPHhgTVLQ3nV3GaQQTnu1RpMz006H9/tawzC
Gg/K1PdHT4hzOI46wYPZyjeDU2VXGFu6vVU9vFQ4ZR/Wjn+DzMzdcRdrJSUOMn
RpTTFusbXUeXHGOpgkXhSYTnvp1gdPc76U5TsK0aGe84AZm1lk0mGrwXCvPqjfo
nhhibBSHBNMoqyF6yTSOpUbyK70tpYRRUGKWDeRK0wK5xkWKUZGtKzy6jYqjjo29
gulwzQif5qWQCB0OontAL2+HvPfaVyckMejUhg62cP/+EHlvUk=
=kOxP
-----END PGP SIGNATURE-----
```

id: [85083a5cc14a441c89dea73f5bdf67c3f9c6afdb](https://sw.hq.mozilla.org/swid/85083a5cc14a441c89dea73f5bdf67c3f9c6afdb)

A giant (extended) Merkle DAG



12 Resilience: closing the regulation / infrastructure gap

13 SWHID

14 Scanner

15 SWHSec

16 Demo time!

Security

Vision

swh-scanner is an **open source** and **open data** source code scanner for **open compliance** workflows, backed by the **largest public archive** of FOSS source code.

Design

- Query Software Heritage as source of truth about public code
- Leverages the Merkle DAG model and SWHIDs for maximum scanning efficiency
 - E.g., no need to query the back-end for files contained in a known directory
- File-level granularity
- Output: source tree partition into known (= published before) v. unknown

Source: gitlab.softwareheritage.org/swh/devel/swh-scanner

License: GPL-3+

Package: pypi.org/project/swh.scanner

swh-scanner demo — Efficiency

```
$ du -sh --exclude=.git /srv/src/linux/git
4,1G /srv/src/linux/git
```

```
$ time swh scanner scan /srv/src/linux
```

```
Files:                78277
      known:          78267 ( 99%)
directories:          5085
      fully-known:    5081 ( 99%)
      partially-known: 4 ( 0%)
```

```
38,65s user 4,71s system 81% cpu 53,127 total
```

```
$ swh scanner scan --output-format ndjson /srv/src/linux/git | grep false
```

```
...
{"scripts/kconfig/symbol.o": {"swhid": "swh:1:cnt:874f19...", "known": false}}
```

12 Resilience: closing the regulation / infrastructure gap

13 SWHID

14 Scanner

15 SWHSec

16 Demo time!

Security



PTCC
Axe 1 : Programme R&D



SWHSec

Leveraging Software Heritage
to Enhance Cybersecurity

Co-porteur

Nom : Barais

Prénom : Olivier

Email : olivier.barais@irisa.fr

Co-porteur

Nom : Di Cosmo

Prénom : Roberto

Email : roberto@dicosmo.org

Co-porteur

Nom : Zacchioli

Prénom : Stefano

Email : stefano.zacchioli@telecom-paris.fr



Enjeux stratégiques du projet



Aujourd'hui, **seulement les hyperscalers** fournissent les plateformes qui hébergent les codes sources et distribuent les binaires des logiciels Open Source

- Github (Microsoft)
- Gitlab (gitlab.com et un large ensemble de repositories privés)

On assiste à la **concentration des outils d'analyse de la supply chain open source** entre des acteurs **non européens**, e.g.:

- Sur les dépendances (Dependatbot <https://github.com/dependabot>)
- Sur le code (rough-auditing-tool-for-security RATS, CodeQL by github, ...)

Opportunité: un outil unique au monde (SWH) qui fournit une source de données vaste et peut être outillé pour la cybersécurité

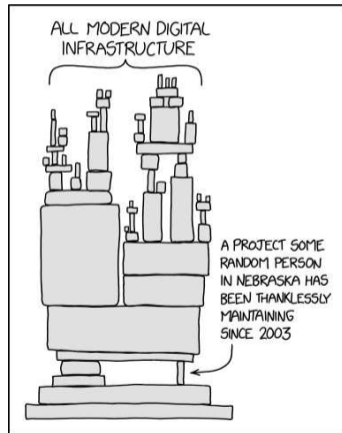


Les objectifs du projet



Construire par-dessus l'infrastructure de Software Heritage (SWH), une chaîne d'analyse et de remédiation unique scalable dédiée à la cybersécurité

- **Analyseur du code source**, capable de bénéficier de l'architecture de SWH
- Infrastructure de gestion des dépendances permettant une **analyse temporelle de l'évolution des dépendances** dans le domaine de l'open-source
- Analyse de **l'impact d'une vulnérabilité** à l'aide de SWH
- Remédiation sur un ensemble de projets d'une vulnérabilité découverte
- **Extension de Software Heritage** pour la prise en compte des outils d'analyse de sécurité

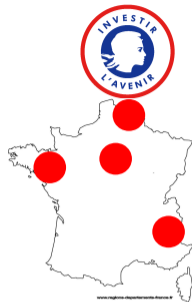


10



SWHSec

Leveraging Software Heritage
to Enhance Cybersecurity



12 Resilience: closing the regulation / infrastructure gap

13 SWHID

14 Scanner

15 SWHSec

16 Demo time!

Security

Get SWHID of compromised openssl files

```
$ for f in openssl-1.0.1*/ssl/d1_both.c; do swh-identify $f; done
swh:1:cnt:0a84f957118afa9804451add380eca4719a9765e  openssl-1.0.1-beta1/ssl/d1_both
swh:1:cnt:7a5596a6b373aeabbd6d8d674f0e20b1618c5012  openssl-1.0.1f/ssl/d1_both.c
swh:1:cnt:2e8cf681ed0976e2b16460170fda27c77cfec6cc  openssl-1.0.1g/ssl/d1_both.c
swh:1:cnt:04aa23107ec53c184505e98091306c7391091bb5  openssl-1.0.1h/ssl/d1_both.c
swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7  openssl-1.0.1/ssl/d1_both.c
```

Get SWHID of compromised openssl files

```
$ for f in openssl-1.0.1*/ssl/d1_both.c; do swh-identify $f; done
swh:1:cnt:0a84f957118afa9804451add380eca4719a9765e  openssl-1.0.1-beta1/ssl/d1_both
swh:1:cnt:7a5596a6b373aeabbd6d8d674f0e20b1618c5012  openssl-1.0.1f/ssl/d1_both.c
swh:1:cnt:2e8cf681ed0976e2b16460170fda27c77cfec6cc  openssl-1.0.1g/ssl/d1_both.c
swh:1:cnt:04aa23107ec53c184505e98091306c7391091bb5  openssl-1.0.1h/ssl/d1_both.c
swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7  openssl-1.0.1/ssl/d1_both.c
```

Look up one origin that contains it using the graph

```
$ swh-graph-lookup.py -c swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7
swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7;
path=ssl/d1_both.c;
anchor=swh:1:rev:86628df45f9eec5b2d46aeb77644ae8f544d1291;
visit=swh:1:snp:6163a539c30011303b5162931fdafd84af8d1c09;
origin=https://github.com/taptipalit/openssl
```

let's check [this occurrence](#)

Find all origins...

```
$ swh-graph-lookup.py --all-origins \  
  -c swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7 \  
  | cut -d ; -f 3 | sort -u | grep swh | sed 's/anchor=/' > allrevs  
$ head -3 allrevs  
swh:1:rev:005b61176f3f72c6c31a2c9431dbc8b5730023ed  
swh:1:rev:01b47383d76f8b9653c6418b0fe1c36043b83ea1  
swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd
```

Find all origins...

```
$ swh-graph-lookup.py --all-origins \  
  -c swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7 \  
  | cut -d ; -f 3 | sort -u | grep swh | sed 's/anchor=/' > allrevs  
$ head -3 allrevs  
swh:1:rev:005b61176f3f72c6c31a2c9431dbc8b5730023ed  
swh:1:rev:01b47383d76f8b9653c6418b0fe1c36043b83ea1  
swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd
```

One of them is pretty late!

```
$ getrevdate.py --swhid swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd  
{'revision': {'swhid': 'swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd',  
              'date': {'date': '2016-02-23T16:22:12+08:00'}}}
```

Find all origins...

```
$ swh-graph-lookup.py --all-origins \  
  -c swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7 \  
  | cut -d ; -f 3 | sort -u | grep swh | sed 's/anchor=/' > allrevs  
$ head -3 allrevs  
swh:1:rev:005b61176f3f72c6c31a2c9431dbc8b5730023ed  
swh:1:rev:01b47383d76f8b9653c6418b0fe1c36043b83ea1  
swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd
```

One of them is pretty late!

```
$ getrevdate.py --swhid swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd  
{'revision': {'swhid': 'swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd',  
              'date': {'date': '2016-02-23T16:22:12+08:00'}}}
```

Let's see where it comes from

```
$ grep "swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd" allrevs  
swh:1:cnt:de8bab873f2cf114d0d1b3e49acfa09bb9d0e4f7;path=ssl/d1_both.c;  
anchor=swh:1:rev:03487116266297d1611556910515f7a3cd7f5fcd;
```

let's check [that occurrence](#)

The screenshot shows a GitHub repository page for `https://github.com/Harrisono984/openssl-android`. The page title is "Browse the archive". A search bar at the top right contains the text "Enter a SWHID to resolve or keyword(s) to". Below the repository name, there is a commit history section showing a commit on 05 April 2019, 18:40:36 UTC. The commit message is "Tip revision: 03487116266297d1611556910515f7a3cd7f5fcd authored by heyunpeng on 23 February 2016, 08:22:12 UTC fix build script". The file `d1_both.c` is selected, and its content is displayed in a code editor. The code is a C file with a header comment that reads: "/* ssl/d1_both.c */", "/*", "* DTLS implementation written by Nagendra Modadugu", "* (nagendra@cs.stanford.edu) for the OpenSSL project 2005.", "*/", "/* =====", "* Copyright (c) 1998-2005 The OpenSSL Project. All rights reserved.", "*".

Navigation options include: `<> Code`, `Branches (2)`, `Releases (0)`, and `Visits`. The current revision is `03487116266297d1611556910515f7a3cd7f5fcd`, and the file path is `5675cb2 / ssl / d1_both.c`. A `Raw File` button is also visible.

```
1 /* ssl/d1_both.c */
2 /*
3  * DTLS implementation written by Nagendra Modadugu
4  * (nagendra@cs.stanford.edu) for the OpenSSL project 2005.
5  */
6 /* =====
7  * Copyright (c) 1998-2005 The OpenSSL Project. All rights reserved.
8  *
```