

Software Heritage

key infrastructure for Open Science and Software Science

Roberto Di Cosmo

Director, Software Heritage
Inria and Université de Paris Cité

12 April 2024



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 Conclusion



Short Bio: Roberto Di Cosmo

Computer Science professor in Paris, now working at INRIA

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



1999 *DemoLinux* – first live GNU/Linux distro

2007 *Free Software Thematic Group*
150 members 40 projects 200Me

2008 *Mancoosi project* www.mancoosi.org

2010 *IRILL* www.irill.org

2015 *Software Heritage* at INRIA

2018 *National Committee for Open Science*, France

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

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 Conclusion



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*)

```
P63SP0T3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND   CHAN33
              EXTEND
              BZF    P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF    CODE500      # ASTRONAUT:  PLEASE CRANK THE
              TC     BANKCALL     #              SILLY THING AROUND
              CADR   GOPERF1
              TCF    GOTOP00H     # TERMINATE
              TCF    P63SP0T3     # PROCEED   SEE IF HE'S LYING

P63SP0T4      TC     BANKCALL     # ENTER      INITIALIZE LANDING RADAR
              CADR   SETPOS1

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

Quake III source code (*excerpt*)

```
float Q_rsqrt( float number )
{
    long i;
    float x2, y;
    const float threehalfs = 1.5F;

    x2 = number * 0.5F;
    y = number;
    i = * ( long * ) &y; // evil floating point bit level hacking
    i = 0x5f3759df - ( i >> 1 ); // what the fuck?
    y = * ( float * ) &i;
    y = y * ( threehalfs - ( x2 * y * y ) ); // 1st iteration
    // y = y * ( threehalfs - ( x2 * y * y ) ); // 2nd iteration, this
    // can be removed

    return y;
}
```

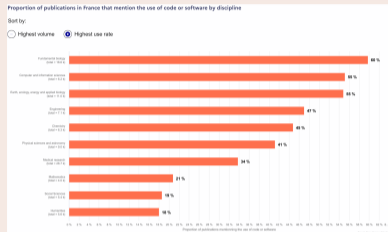
Len Shustek, Computer History Museum

2006

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

Software is a pillar of Open Science

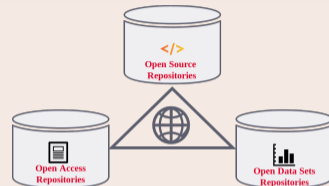
Software powers modern research



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

2024 French Open Science Monitor

Key pillar: software



Links are **important**

Nota Bene

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

access to the *source code* is essential!

Preserving (the history of) source code is necessary for *reproducibility*

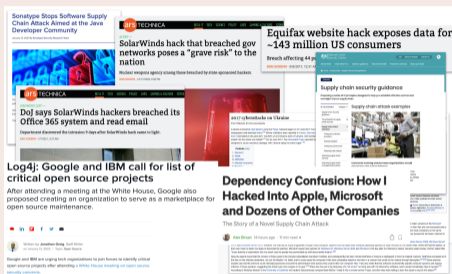
How are we managing our software ?

Reproducibility, maintenance in Academia



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

Security, integrity, traceability in Industry



Can they track the software that they

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

awareness is raising at the level of public policy

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 Conclusion



International highlights

Paris Call on Software Source code (2019, UNESCO)

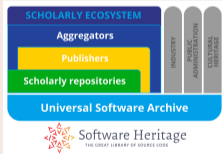


40 international experts call to “promote software development as a valuable research activity, and research software as a key enabler for Open Science/Open Research, [...] recognising in the careers of academics their contributions to high quality software development, in all their forms”



Open Source in UNESCO [recommendations](#) for Open Science, 2021

Software in the EOSC



2020 [EOSC SIRS](#) connect scholarly ecosystem via Software Heritage
2021 [EOSC Task Force](#) on Infrastructures for Research Software
2022 [FAIRCORE4EOSC project](#) WP6 implements SIRS report

And much more

Software track in [OSEC 2022](#), Software working group launched in Science Europe, DFG adds software [to model CV \(9/22\)](#), NASA unveils [Open Science policy \(12/22\)](#), ...

Fundamental needs for software in Open Science (selection)

Archive

Research software artifacts must be properly **archived**
make sure we can *retrieve* them (*reproducibility*)

Reference

Research software artifacts must be properly **referenced**
make sure we can *identify* them (*reproducibility*)

Describe

Research software artifacts must be properly **described**
make it easy to *discover* and *reuse* them (*visibility*)

Cite/Credit

Research software artifacts must be properly **cited** (*not the same as referenced!*)
to give *credit* to authors (*evaluation!*)

Archive and reference: some popular approaches that do not fit the bill

A - Since the 1970's 1990's

.zip or .tar file on:

- ftp server (e.g. [gnu](#))
- web page ([example](#))
- document archive (+ DOI [sample](#))

B - Since the 2000's

Rely on *software forges*

- institutional/project (e.g. [example](#))
- free commercial ones: BitBucket, GitHub, GitLab, ... (e.g. [parmap](#))

C: a mix of the two

Artifacts Available Artifacts Evaluated & Functional

Authors/Contributors: [Authors Info & Affiliations](#)

DOI: <https://doi.org/10.1145/> [redacted] Version: 1.0

Description

A source archive of [redacted], and the version of [redacted] used in the paper eval. A more up-to-date version of [redacted] can be found at [github.com/\[redacted\]/\[redacted\].git](https://github.com/[redacted]/[redacted].git)

Assets

Read Me [redacted]

[Download \(3.5 KB\)](#)

Can get no satisfaction...

- A *Poor user experience*
- B *No preservation guarantee*
- C *Can do so much better*

Forges are *not* archives!

2015: the first big bad news

Google Code and Gitorious.org shutdown: ~1M endangered repositories

- broken links in the web of knowledge (my papers too)

Big bad news keep coming in

- summer 2019: BitBucket announces Mercurial VCS sunset
- july 2020: BitBucket erases 250.000+ repositories (including research software)
- summer 2022: GitLab.com considers erasing **all** projects that are **inactive for a year**

In Academia too!

- 2021: Inria's old gforge is unplugged... **breaks the Opam build chain** for OCaml

We need a universal archive of software source code: now we have one!

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility**
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 Conclusion





Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Collect, preserve and share *all* software source code

Preserving our heritage, enabling better software and better science for all

Reference catalog



find and **reference** all software source code

Universal archive



preserve and **share** all software source code

Research infrastructure



enable analysis of all software source code

Sharing the vision



United Nations
Educational, Scientific and
Cultural Organization



And many more ...

www.softwareheritage.org/support/testimonials

Donors, members, sponsors

Inria

Diamond sponsor



Platinum sponsors



Gold sponsors



Silver sponsors



Bronze sponsors



The largest software archive, a shared infrastructure

One infrastructure
open and shared

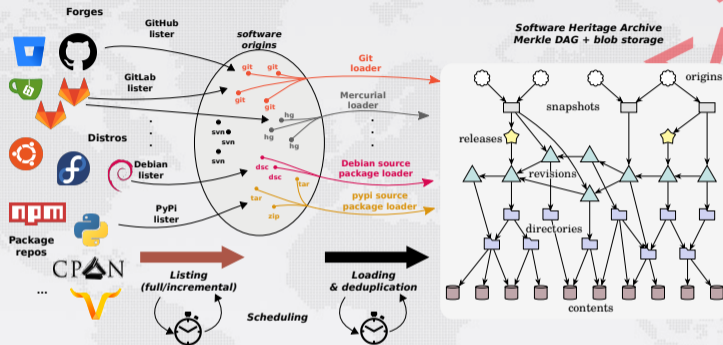


The largest archive ever built



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

Software Heritage: a *radically different* approach to archiving



Global development history **permanently archived** in a **uniform data model**

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

Meet the SWHID identifier

Software Hash Identifiers (SWHID)

see swhid.org

35+B **intrinsic, decentralised, cryptographically strong identifiers, SWHIDs**



In **SPDX 2.2**; IANA registered "swh : "; WikiData [P6138](#); ISO standard ([ongoing](#))

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

Examples: [Apollo 11 AGC](#), [Quake III rsqrt](#); Guidelines available: [HOWTO](#) and [ICMS 2020](#)

A quick tour as a user

- **designed for source code:** [Browse](#) (e.g. [Apollo 11 excerpt](#), see also [Apollo 11 blog post](#)) like on a developer platform, not a document archive!
- **reference source code:** all granularities, using SWHIDs ([full specification available online](#))
 - compare Fig. 1 and conclusions in [the 2012 version](#) and [the updated version](#)
 - SWHID in [a replication experiment](#)
 - SWHIDs *guarantee integrity* like in *blockchains*
demo if time left:
 - 1 download a version of a project for a given SWHID
 - 2 compute locally the SWHID with `swh-identify`
 - 3 check that the computed id match the given one

Getting software archived

- **automated harvesting**: over **290 million software origins**, your researchers' work may already be there (actually, [here](#))!
- **universal archive**: *all* source code **from all platforms** (BitBucket, GitHub, GitLab, your own forge, etc.)
 - **trigger archival** of *any code* in one click with **the updateswh browser extension**
 - **use webhooks** to automatically archive *your code* (a **GitHub action** is available too)
 - **journals, libraries, open access portals** may *deposit sourcecode and metadata*
 - Example [article from IPOL](#)
 - Example [article from eLife](#)

A look at some adoption indicators

From [Melissa Harrison's OSEC 2022 talk](#)



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

- 6205 "software" references identified
- Top 8 listed, then long tail of 1055 other sites – 932 are unique "source"

Use on replicabilitystamp.org

Lightweight Curvature Estimation on Point Clouds with Randomized Corrected Curvature Measures

Jacques-Olivier Lachaud, David Coeurjolly, Céline Labart, Pascal Romon, Boris Thibert
Wiley Computer Graphics Forum (CGF)



HAL+SWH in the Open Science software booklet

Funding agencies recommendations [ANR 2023 guidelines](#) (p. 17)

Enfin, conformément au 2^{ème} Plan national pour la science ouverte, L'ANR recommande que les logiciels développés durant le projet soient mis à disposition sous une licence libre³⁰ et que les codes sources soient stockés dans l'archive Software Heritage³¹ en indiquant la référence au financement ANR.

Call to action: best practices for ARDC are available... today!

Archiving and referencing

For **all source code** used in research (*yes, even small scripts!*)

- ensure it is archived in Software Heritage (see [save code now](#))
- get the proper **SWHID** for your software (see [detailed HOWTO](#))
- add it to research articles for reproducibility (see [detailed HOWTO](#))

Describing and Citing/Crediting

For **software you want to put forward** (*mention in your CV, reports, etc., get citations and credit for it*), do the following **extra steps**:

- add **codemeta.json** with description (see the [codemeta generator](#))
- reference in the HAL portal (french partners, see [online HAL documentation](#))
- cite software using the [biblatex-software](#) package (in CTAN and TeXLive)

- train students, colleagues
- engage journals, conferences, learned societies

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software**
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 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](#).

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-
```


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;
```

Total cost: approximately .5 euros

Results

Completed

Time in queue: 272 ms

Run time: 33.545 sec

Data scanned: 94.51 GB

Results (20)

Copy

Download results

Search rows

< 1 > ⚙

#	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

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph**
- 7 A lot more lies ahead
- 8 Conclusion



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



Paolo Boldi, Antoine Pietri, Sebastiano Vigna, Stefano Zacchiroli

Ultra-Large-Scale Repository Analysis via Graph Compression

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

Results

Full graph structure (35 B nodes, 500 B edges) in 300 GiB RAM

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

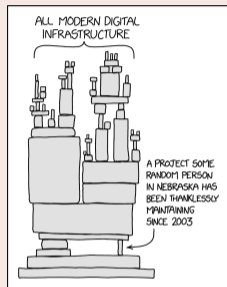
Java and gRPC APIs available, Rust is coming!

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

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 Conclusion

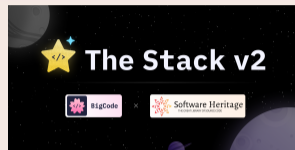


Cybersecurity: tracing vulnerabilities



See the [SWHSec project](#)

(Generative) AI



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

Big Data challenges



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

- 1 Introduction
- 2 Software and Source Code
- 3 An emerging policy framework for Open Science
- 4 Software Heritage for Open Science and Reproducibility
- 5 Software Heritage for Research on Software
- 6 Efficient traversal of the full graph
- 7 A lot more lies ahead
- 8 Conclusion**

A rally flag for a grand vision

Bring together academia, industry, governments, communities

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

Software Heritage is the first brick ...

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

... that will enable

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

A lot more is needed

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

You can help!

use, disseminate, contribute, build&adapt research tools, ...

Join a growing and active community

Team

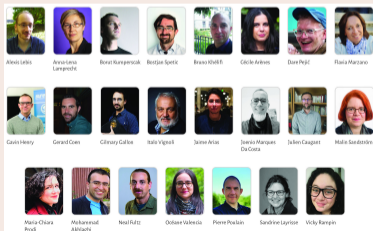


Contributors to the platform

```
13:21:32 <sew> from last time I ran it? it very likely is
13:21:47 <sew> we had a X2 on the edges in a single year
13:22:14 <vlorenzo> ah
13:53:44 <acko> sure, i think i was remembering the LLP time on granet rather than the one (on the previous
13:54:01 <acko> wasn't it something like 10-14 days (on granet)?
13:55:11 <sew> i think it depends on the number of weights you use
13:55:23 <sew> i had something like that to do the parameter sweep
13:56:31 <sew> but then i settled on a few good gamma values
13:55:44 <sew> and afterwards it was only over ~3-4 days
14:02:57 <acko> ok
15:19:35 <jelmer> vlorenzo: when is jenkins meant to kick in ? I didn't think the CI would mean you passing test
15:19:59 <jelmer> alternatively, I could try to get it working locally - for some reason tox doesn't run here,
15:20:48 <jelmer> completing it can't find swb.scheduler
15:20:48 <vlorenzo> jenkins is down until tomorrow evening (pairs time)
15:20:59 <vlorenzo> bad day for submitting your code :D
15:21:18 <vlorenzo> er yeah, i just fixed that issue
15:21:31 <vlorenzo> but the fixed swb.scheduler is not pushed to pypi because jenkins
15:23:25 <jelmer> ah
15:23:40 <vlorenzo> in the meantime, you can change apply this patch: https://github.com/softwareheritage/org/
15:23:44 <vlorenzo> s/pip/pypi/546
15:24:13 <vlorenzo> as an ugly workaround
15:24:13 <vlorenzo> actually, just adding 'typed-progress' = 4.0.0* should do it
15:25:00 <vlorenzo> when jenkins is back online I'll push a new release of swb.scheduler without the missing
15:25:00 <vlorenzo> dependency on typed-progress
```

- Nicks
- Alphare
- arnes
- arij
- arombo
- ar-jan
- bchauret[m]
- cmarric[m]
- dani
- deuts
- ericson2514
- franc3dre1
- guythly
- GuyMS2
- hultsted
- hpior
- jpyshov
- jelmer
- KSHivendu
- landry[m]
- marcouste

Ambassadors



Work with us!

Big Data Development and Architecture Engineer

The Software Heritage project Software Heritage is a universal software source code archive project, whose aim is to recover, preserve for the very long term and share all publicly available source co..

November 24, 2023

March 1, 2024

[Read More](#)

DevOps Engineer

The Software Heritage project Software Heritage is a universal software source code archive project, whose aim is to recover, preserve for the very long term and share all publicly available source co..

November 24, 2023

[Read More](#)

Fullstack Python Developer

The Software Heritage project Software Heritage is a universal software source code archive project, whose aim is to recover, preserve for the very long term and share all publicly available source co..

November 13, 2023

[Read More](#)

<https://softwareheritage.org/jobs/>

Annual report



5 years in 5 minutes

[Link](#)



Evolution of our codebase

[Link](#)

