# Software Heritage

## key infrastructure for Open Science and Software Science

Jaime Arias

Research Engineer
CNRS, LIPN, Université Sorbonne Paris Nord

November 27, 2024

## Software Heritage

### THE GREAT LIBRARY OF SOURCE CODE

# Hello!

I am **Jaime Arias**

- CNRS Research Engineer @ LIPN
- Member @ Collège Codes Sources et Logiciels
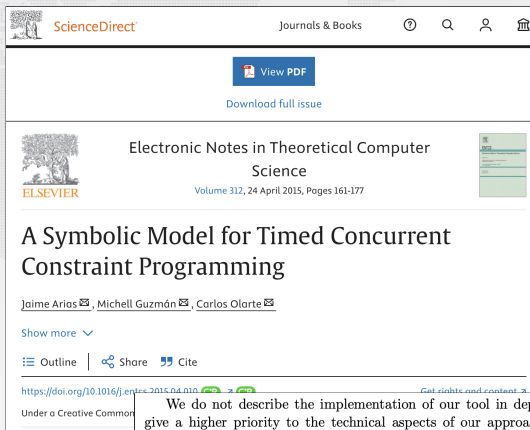- Ambassador @ Software Heritage

You can find me at:

✉ `arias@lipn.fr`

🔗 `https://www.jaime-arias.fr`

We do not describe the implementation of our tool in depth here in order to give a higher priority to the technical aspects of our approach. The reader can find the details of the implementation as well as the execution of the examples described in this paper at http://www.labri.fr/perso/jarias/symbolicMC.

# Software *source code* is precious knowledge

## 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     GOTOPOOH      # 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
```

# Software *source code* is precious knowledge

## 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    GOTOPOOH    # 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
```

## 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;
}
```

# 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."*

# 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."*

> ### Len Shustek, Computer History Museum    2006
> *"Source code provides a view into the mind of the designer."*

# 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."*

**Len Shustek, Computer History Museum** 2006

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

**Sonatype Survey** 2017

80% to 90% of a new application is ... just to **reuse**!

# 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."*

**Len Shustek, Computer History Museum** — 2006

*"Source code provides a view into the mind of the designer."*
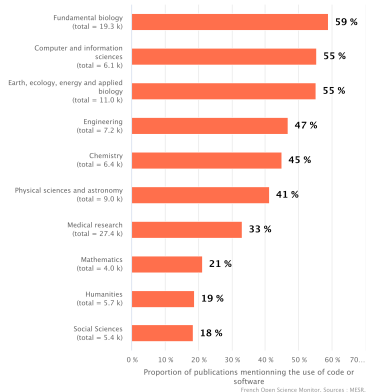
**Sonatype Survey** — 2017

80% to 90% of a new application is . . . just to **reuse**!

**Art. L. 112-2 du Code de la Propriété Intellectuelle** — 1994

*"Sont considérés notamment comme œuvres de l'esprit au sens du présent code: . . .
13o «Les logiciels, y compris le matériel de conception préparatoire»; . . . "*

## Software powers modern research



Proportion of publications in France published in 2022 that mention the use of code or software by discipline

| Discipline | Percentage |
|---|---|
| Fundamental biology (total = 19.3 k) | 59 % |
| Computer and information sciences (total = 6.1 k) | 55 % |
| Earth, ecology, energy and applied biology (total = 11.0 k) | 55 % |
| Engineering (total = 7.2 k) | 47 % |
| Chemistry (total = 6.4 k) | 45 % |
| Physical sciences and astronomy (total = 9.0 k) | 41 % |
| Medical research (total = 27.4 k) | 33 % |
| Mathematics (total = 4.0 k) | 21 % |
| Humanities (total = 5.7 k) | 19 % |
| Social Sciences (total = 5.4 k) | 18 % |

Proportion of publications mentionning the use of code or software

French Open Science Monitor, Sources : MESR.

*Over 20% of articles using software across all disciplines share it*
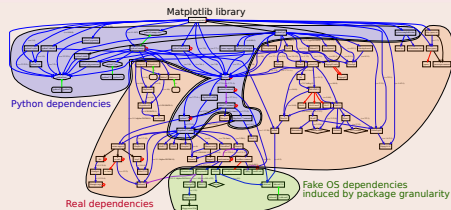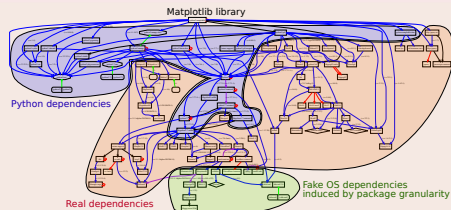*2024 French Open Science Monitor*

# Source code is *special* (software is *not* data)

## Software *evolves* over time

- projects may last decades
- the *development history* is key to its *understanding*

# Source code is *special* (software is *not* data)

## Software *evolves* over time

- projects may last decades
- the *development history* is key to its *understanding*

## Complexity

- *millions* of lines of code
- large *web of dependencies*
  - easy to break, difficult to maintain
  - *research software* a thin top layer
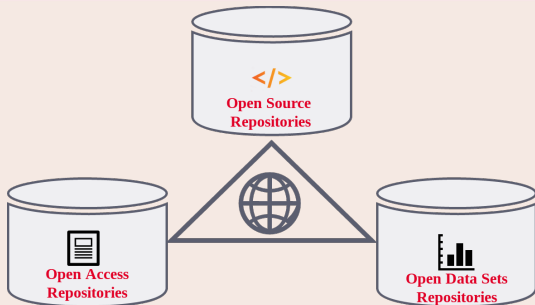- sophisticated *developer communities*

# Source code is *special* (software is *not* data)

## Software *evolves* over time
- projects may last decades
- the *development history* is key to its *understanding*

## Complexity
- *millions* of lines of code
- large *web of dependencies*
  - easy to break, difficult to maintain
  - *research software* a thin top layer
- sophisticated *developer communities*



Matplotlib library

Python dependencies

Real dependencies

Fake OS dependencies
induced by package granularity

## The human side
design, algorithm, code, test, documentation, community, funding

and so many more facets . . .

**Key pillar: software**

Open Source Repositories

Open Access Repositories

Open Data Sets Repositories

Links are important

**Nota Bene**

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

# Software is a pillar of Open Science

**Key pillar: software**

</>
**Open Source Repositories**

**Open Access Repositories**

**Open Data Sets Repositories**

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*

## Archive

Research software artifacts must be properly archived

make sure we can *retrieve* them (*reproducibility*)

## 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*)

# 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*)

# 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!*)

# Where is the source code?

## Collaborative development platforms (aka "forges")

- BitBucket, GitLab(.com), GitHub, etc.
- support for version control, issues, etc.
- example:
  - https://depot.lipn.univ-paris13.fr/cosyverif/cosydraw
  - https://gitlab.inria.fr/gt-sw-citation/bibtex-sw-entry/

# Where is the source code?

## Collaborative development platforms (aka "forges")

- BitBucket, GitLab(.com), GitHub, etc.
- support for version control, issues, etc.
- example:
  - https://depot.lipn.univ-paris13.fr/cosyverif/cosydraw
  - https://gitlab.inria.fr/gt-sw-citation/bibtex-sw-entry/

## Distribution platforms

- CTAN, CRAN, PyPi, Debian, etc.
- example: https://ctan.org/pkg/biblatex-software

# Where is the source code?

## Collaborative development platforms (aka "forges")

- BitBucket, GitLab(.com), GitHub, etc.
- support for version control, issues, etc.
- example:
    - `https://depot.lipn.univ-paris13.fr/cosyverif/cosydraw`
    - `https://gitlab.inria.fr/gt-sw-citation/bibtex-sw-entry/`

## Distribution platforms

- CTAN, CRAN, PyPi, Debian, etc.
- example: `https://ctan.org/pkg/biblatex-software`

## Archives

- Software Heritage
- example: archived version of biblatex-software

## A - Since the ~~1970's~~ 1990's

.zip or .tar file on:

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

# 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. imitator)

# 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. imitator)

## C: a mix of the two

# 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. imitator)

## C: a mix of the two



## 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

# 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

## 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

# 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

## 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

# 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

## 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:**

# 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

## 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!**

# 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

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

## Collect, preserve and share *all* software source code

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

## Reference catalog



find and reference all
software source code

# Software Heritage
## 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

# 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



UNESCO
United Nations
Educational, Scientific and
Cultural Organization

And many more ...
www.softwareheritage.org/support/testimonials

## Sharing the vision



UNESCO
United Nations
Educational, Scientific and
Cultural Organization



And many more ...
www.softwareheritage.org/support/testimonials

## Donors, members, sponsors

Inría

**Diamond sponsor**

cea

**Platinum sponsors**

cnrs · intel · MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE · Microsoft · HUAWEI

**Gold sponsors**

Hugging Face · open invention network · servicenow · SOCIETE GENERALE · SORBONNE UNIVERSITÉ · Université Paris Cité

**Silver sponsors**

AdaCore · MINISTÈRE DES ARMÉES DGA · RÉPUBLIQUE FRANÇAISE · GitHub · Google · UNIVERSITÀ DI PISA

**Bronze sponsors**

SCANOSS · SCUOLA NORMALE SUPERIORE · Université de Lorraine

One infrastructure
open and shared

Cultural Heritage   Industry   Research   Public Administration

Software Heritage

**One** infrastructure **open** and **shared**

Cultural Heritage · Industry · Research · Public Administration

**Software Heritage**

The largest archive ever built

| Source files | Commits | Projects |
|---|---|---|
| 17,798,218,376 | 3,802,143,973 | 278,187,495 |

| Directories | Authors | Releases |
|---|---|---|
| 14,364,868,206 | 69,923,710 | 82,196,102 |

One infrastructure open and shared

**Cultural Heritage**  **Industry**  **Research**  **Public Administration**

Software Heritage

The largest archive ever built

| | | |
|---|---|---|
| **Source files** | **Commits** | **Projects** |
| 17,798,218,376 | 3,802,143,973 | 278,187,495 |
| **Directories** | **Authors** | **Releases** |
| 14,364,868,206 | 69,923,710 | 82,196,102 |

| | | |
|---|---|---|
| **Bitbucket** 2,509,402 origins | 56,983 origins | **git** 24,600 origins |
| **R** 26,599 origins | **debian** 136,338 origins | 53,297 origins |
| **GitHub** 197,883,004 origins | **gitiles** 10,171 origins | **GitLab** 4,216,298 origins |
| **git** 2,926 origins | **Gogs** 172 origins | **GO** 971,549 origins |
| **Guix** 14,482 origins | **GNU** 354 origins | **heptapod** 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 20 billion unique source files from over 300 million software projects
- ~2PB (compressed) blobs, ~50 B nodes, ~700 B edges

# Software Heritage is *radically different*, cont'd

## Software Hash Identifiers (SWHID)

50+B intrinsic, decentralised, cryptographically strong identifiers, SWHIDs

## Software Hash Identifiers (SWHID)

see swhid.org

50+B intrinsic, decentralised, cryptographically strong identifiers, SWHIDs



schema_version

object_id

`swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa`

prefix   object_type

"snp" - snapshot

"rel" - release

"rev" - revision

"dir" - directory

"cnt" - content

## Software Hash Identifiers (SWHID)

50+B intrinsic, decentralised, cryptographically strong identifiers, SWHIDs

# Software Heritage is *radically different*, cont'd

## Software Hash Identifiers (SWHID)

50+B intrinsic, decentralised, cryptographically strong identifiers, SWHIDs



In SPDX 2.2; IANA registered `"swh:"`; WikiData P6138; ISO standard

# Software Heritage is *radically different*, cont'd

## Software Hash Identifiers (SWHID)                                    see swhid.org

50+B intrinsic, decentralised, cryptographically strong identifiers, SWHIDs



In SPDX 2.2; IANA registered `"swh:"`; WikiData P6138; ISO standard

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

Examples: Apollo 11 AGC, Quake III rsqrt; Guidelines available: HOWTO and ICMS 2020

# Software Heritage is *radically different*, cont'd

A quick tour as a user

- **designed for source code**: Browse (e.g. Apollo 11 excerpt) like on a developer platform, not a document archive!

- reference source code: all granularities, using SWHIDs (full specification available online)
  - SWHIDs *guarantee integrity* like in *blockchains*



Figure: Compare Fig. 1 and conclusions in the 2012 version and the updated version

## Getting software archived

- **automated harvesting**: over 290 million software origins, your researchers' work may already be there (actually, here)!

## 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 walkthrough

- Browse (e.g. Imitator [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 adt2amas: code source, 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, …

# An example of long term reproducibility for HPC

(re)create fully reproducible binaries from source...   https://guix.gnu.org/

- functional package manager
- bit by bit reproductibility
- *from the source code*

# An example of long term reproducibility for HPC

## (re)create fully reproducible binaries from source... `https://guix.gnu.org/`



- functional package manager
- bit by bit reproductibility
- *from the source code*

## ... with a focus on HPC `https://hpc.guix.info/`


Reproducible software deployment for high-performance computing.

- environment control
- support cluster deployment
- *from the source code*

# An example of long term reproducibility for HPC

## (re)create fully reproducible binaries from source... `https://guix.gnu.org/`



- functional package manager
- bit by bit reproductibility
- *from the source code*

## ... with a focus on HPC `https://hpc.guix.info/`



Reproducible software deployment for high-performance computing.

- environment control
- support cluster deployment
- *from the source code*

## connection with Software Heritage

- source code *archival and identification* for `guix` and `nix`
- automatic fallback for missing sources (see experience report)

https://hal.archives-ouvertes.fr/hal-02130801

swh:1:dir:393b611a1424f032e83569bf6762502371cfcf65

with minimal user overhead!

# 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)

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

## Archiving and referencing

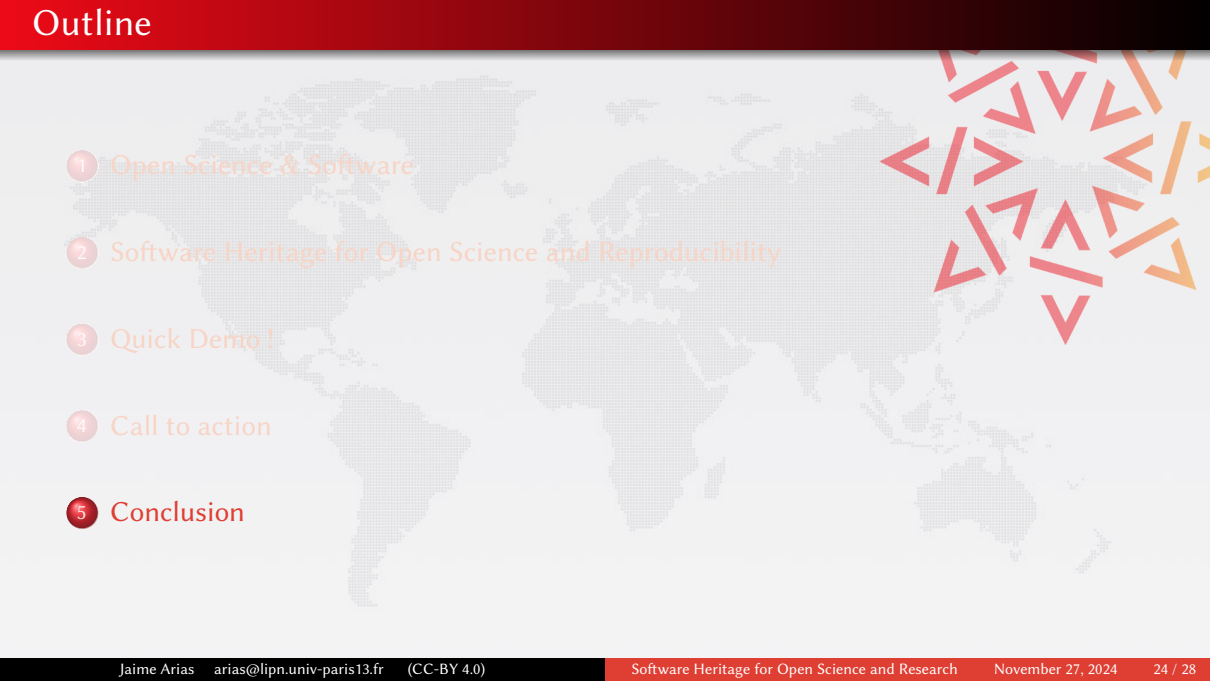For all source code used in research (*yes, even small scripts!*)

- 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)

# 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

# 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 the first brick …

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

# 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, …

## Team

# Join a growing and active community

## Team



## Ambassadors

## Team



## Contributors to the platform



## Ambassadors

# Join a growing and active community

## Team



## Contributors to the platform



## Ambassadors



Alexis Lebis    Anna-Lena Lamprecht    Borut Kumperscak    Bostjan Jpetic    Bruno Khelifi    Cécile Antzen    Dare Pejić    Flavia Marzano

Gavin Henry    Gerard Coen    Gihozary Calline    Italo Vignoli    Jaime Arias    Joenio Marques Da Costa    Julien Caupart    Malo Sandström

Maria-Chiara Prodi    Mohannoud Alhraghi    Neal Fultz    Octave Valecela    Pierre Poulain    Sandrine Layrisse    Vicky Rampin

## 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...

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

## Annual report



## 5 years in 5 minutes

## Annual report



## 5 years in 5 minutes — Link



## Evolution of our codebase — Link

it's a long road, but together we can make it

# Thank you

This presentation reuses material from Roberto di Cosmo's presentations.