

Vers un pilier logiciel de la Science Ouverte

défis et opportunités pour la reproductibilité et pour la science ouverte

Roberto Di Cosmo

Director, Software Heritage
Inria and Université de Paris Cité

18 Janvier 2023



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

- 1 Introduction
- 2 Software and Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions



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 Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions



Why Open Science?

Open Science ([Second National Plan for Open Science](#), France, 2021)

Unhindered dissemination of results, methods and products from scientific research. It draws on *the opportunity provided by recent digital progress* to develop *open access to publications* and – as much as possible – *data, source code and research methods*.

Why Open Science?

Open Science ([Second National Plan for Open Science](#), France, 2021)

Unhindered dissemination of results, methods and products from scientific research. It draws on *the opportunity provided by recent digital progress* to develop *open access* to *publications* and – as much as possible – *data, source code* and *research methods*.

Jean-Eric Paquet (EU DGRI, [on the objective of Open Science](#))

“Increase scientific quality, the pace of discovery and technological development, as well as societal trust in science.”

Why Open Science?

Open Science ([Second National Plan for Open Science](#), France, 2021)

Unhindered dissemination of results, methods and products from scientific research. It draws on *the opportunity provided by recent digital progress* to develop *open access to publications* and – as much as possible – *data, source code and research methods*.

Jean-Eric Paquet (EU DGRI, [on the objective of Open Science](#))

“Increase *scientific quality*, the *pace of discovery and technological development*, as well as *societal trust in science*.”

Mariya Gabriel ([EU Commissioner](#) for Research)

The COVID-19 crisis has also shown that cooperation at international level in research and innovation is more important than ever, including through *open access to data and results*. *No nation, no country can tackle any of these global challenges alone*.

Why Open Science?

Open Science ([Second National Plan for Open Science](#), France, 2021)

Unhindered dissemination of results, methods and products from scientific research. It draws on the opportunity provided by recent digital progress to develop open access to publications and – as much as possible – data, source code and research methods.

Jean-Eric Paquet (EU DGRI, [on the objective of Open Science](#))

“Increase scientific quality, the pace of discovery and technological development, as well as societal trust in science.”

Mariya Gabriel ([EU Commissioner](#) for Research)

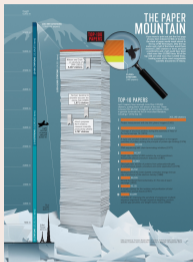
The COVID-19 crisis has also shown that cooperation at international level in research and innovation is more important than ever, including through *open access to data and results. No nation, no country can tackle any of these global challenges alone.*

Yuval Noah Harari (on COVID 19)

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

Software is a pillar of Open Science

Software powers modern research



[...] software [...] essential in their fields.

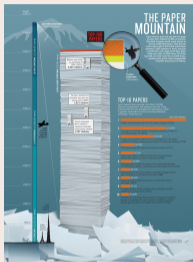
Top 100 papers (Nature, 2014)

Sometimes, if you don't have the software, you don't have the data

Christine Borgman, Paris, 2018

Software is a pillar of Open Science

Software powers modern research



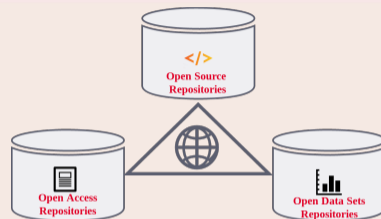
[...] software [...] essential in their fields.

Top 100 papers (Nature, 2014)

Sometimes, if you don't have the software, you don't have the data

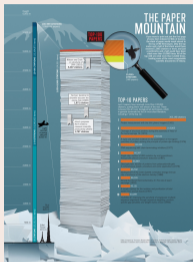
Christine Borgman, Paris, 2018

A key pillar: software (source code)



Software is a pillar of Open Science

Software powers modern research



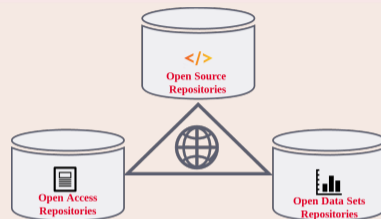
[...] software [...] essential in their fields.

Top 100 papers (Nature, 2014)

Sometimes, if you don't have the software, you don't have the data

Christine Borgman, Paris, 2018

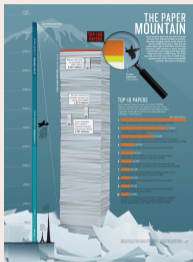
A key pillar: software (source code)



The links in the picture are **important**

Software is a pillar of Open Science

Software powers modern research



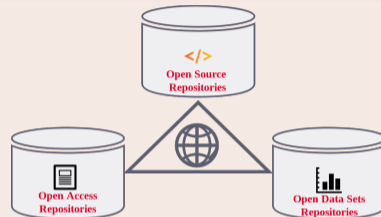
[...] software [...] essential in their fields.

Top 100 papers (Nature, 2014)

Sometimes, if you don't have the software, you don't have the data

Christine Borgman, Paris, 2018

A key pillar: software (source code)



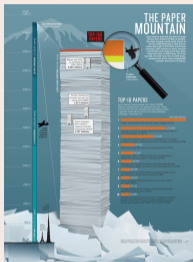
The links in the picture are **important**

Nota Bene

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

Software is a pillar of Open Science

Software powers modern research



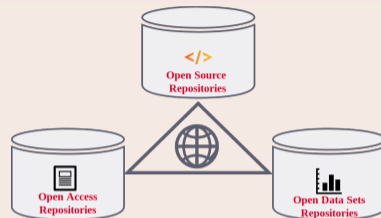
[...] software [...] essential in their fields.

Top 100 papers (Nature, 2014)

Sometimes, if you don't have the software, you don't have the data

Christine Borgman, Paris, 2018

A key pillar: software (source code)



The links in the picture are **important**

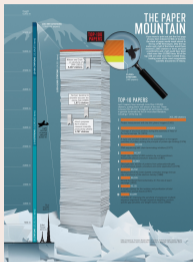
Nota Bene

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

access to the *source code* is essential!

Software is a pillar of Open Science

Software powers modern research



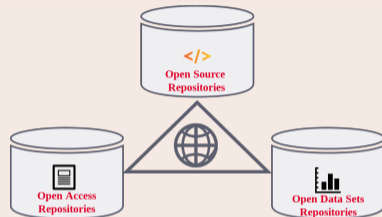
[...] software [...] essential in their fields.

Top 100 papers (Nature, 2014)

Sometimes, if you don't have the software, you don't have the data

Christine Borgman, Paris, 2018

A key pillar: software (source code)



The links in the picture 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*

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.”

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
```


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;
}
```

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 source code is *not* data

Software *evolves* over time

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

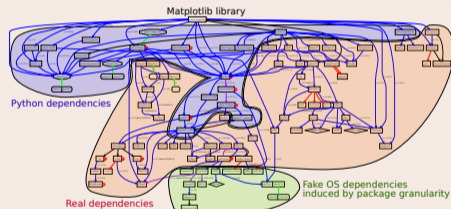
Software source code is *not* data

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*



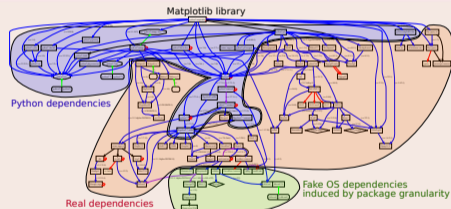
Software source code is *not* data

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*



The human side

design, algorithm, code, test, documentation, community, funding

and so many more facets ...

How are we managing our software ?

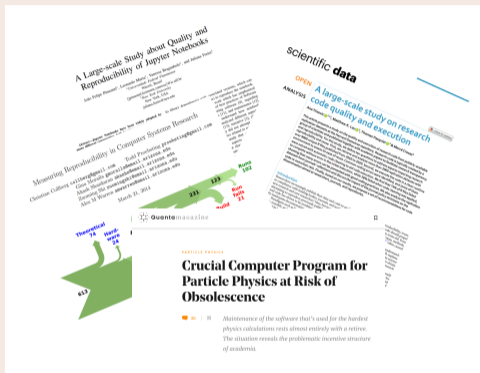
Reproducibility, maintenance in Academia



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

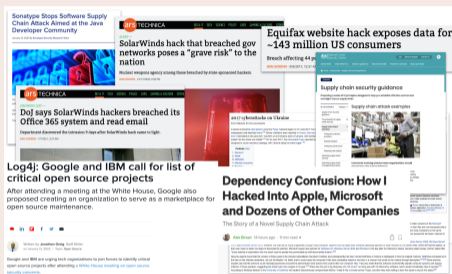
How are we managing our software ?

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

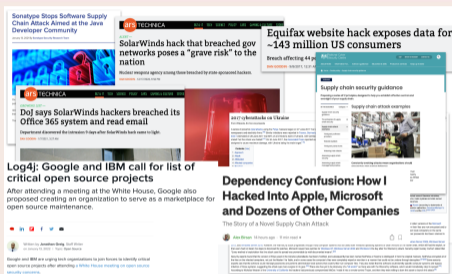
How are we managing our software ?

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 Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions



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

International highlights

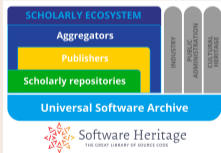
Paris Call on Software Source code (2019, UNESCO)



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
- 2023 [INFRAEOSC call](#) on quality of scientific software

International highlights

Paris Call on Software Source code (2019, UNESCO)

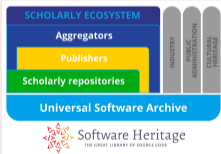


40 international experts call to “promote software development as a valuable research activity, and research software as a key enabler for Open Science/Open Research, [...] recognising in the careers of academics their contributions to high quality software development, in all their forms”



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

2023 [INFRAEOSC call](#) on quality of scientific software

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 \(12/22\)](#), ...

French National plan for Open Science, 2021-2024


MINISTÈRE
DE L'ENSEIGNEMENT
SUPÉRIEUR,
DE LA RECHERCHE
ET DE L'INNOVATION
*Liberté
Égalité
Fraternité*



SECOND FRENCH PLAN FOR OPEN SCIENCE

Generalising open science in France 2021-2024



1


MINISTÈRE
DE L'ENSEIGNEMENT
SUPÉRIEUR,
DE LA RECHERCHE
ET DE L'INNOVATION
*Liberté
Égalité
Fraternité*

Second French Plan for Open Science



GENERALISING
OPEN SCIENCE
IN FRANCE 2021-2024

Launch on 6 July 2021 by Frédérique Vidal, Minister for Higher Education, Research and Innovation

- Multiplying the **levers for change** in order to **generalise open science practices**
- Structuring the **policy for opening up or sharing research data**
- New commitments to the **opening of source code** produced by research
- **European and international inclusion** in the context of the French Presidency of the European Union
- **Disciplinary and thematic variations**: open science policies must be adapted to disciplinary specificities

2

French National plan for Open Science, 2021-2024



SECOND FRENCH PLAN FOR OPEN SCIENCE

Generalising open science in France 2021-2024



1

Second French Plan for Open Science



Launch on 6 July 2021 by Frédérique Vidal, Minister for Higher Education, Research and Innovation

- Multiplying the **levers for change** in order to **generalise open science practices**
- Structuring the **policy for opening up or sharing research data**
- New commitments to the **opening of source code** produced by research
- **European and international inclusion** in the context of the French Presidency of the European Union
- **Disciplinary and thematic variations**: open science policies must be adapted to disciplinary specificities

2

Path Three : Opening up and promoting source code produced by research

7

Recognize and support the dissemination under an open source licence of software produced by publicly funded research programmes

« The opening of software source code is a major challenge for the **reproducibility** of scientific results. »

8

Highlight the production of source code from higher education, research and innovation

« Distribution of software products under **open source licence** will be preferred. »

9

Define and promote an **open source software policy**

3

Define and promote an open source software policy

- Produce a **National Charter for Open Source Software** coming from higher education, research and innovation
- Develop the **link between data and software** through a network of **Chief Data Officers** in the various universities and research performing organisations.
- Develop the **economic models of open source software** and make them known within commercialization services
- **Support Software Heritage** and recommend it for the archiving and referencing of source code

Recognise source code as a contribution to research

- Create an **open source research software prize**
- **Provide greater recognition** for software production in the career of researchers, research support staff

Build an ecosystem that connects code, data and publications

- Develop **proper coordination** between software forges, open publication archives, data repositories and the scientific publishing sector.

4

Five action lines (see [details online](#))

- Identifying and highlighting research software production
- Technical and social tools and best practices
- Valorization and sustainability
- Liaison and animation at national, European, and international levels
- Recognition and careers

Five action lines (see [details online](#))

- Identifying and highlighting research software production
- Technical and social tools and best practices
- Valorization and sustainability
- Liaison and animation at national, European, and international levels
- Recognition and careers

Leveraging experience and connections

- Open Source thematic group in Systematic (since 2007, more on demand)
- Collaboration with DINUM, Eclipse Foundation, OW2, ...

Five action lines (see [details online](#))

- Identifying and highlighting research software production
- Technical and social tools and best practices
- Valorization and sustainability
- Liaison and animation at national, European, and international levels
- Recognition and careers

Leveraging experience and connections

- Open Source thematic group in Systematic (since 2007, more on demand)
- Collaboration with DINUM, Eclipse Foundation, OW2, ...

Composition

Chairs: Roberto Di Cosmo and François Pellegrini

20+ active members from a broad panel of institutions and fields

- 1 Introduction
- 2 Software and Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them**
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions



ARDC

- **Archive** for retrieval
(*reproducibility*)
- **Reference** for
identification
(*reproducibility*)
- **Describe** for discovery
and reuse
- **Cite/Credit** for credit
and evaluation

ARDC

- **Archive** for retrieval (*reproducibility*)
- **Reference** for identification (*reproducibility*)
- **Describe** for discovery and reuse
- **Cite/Credit** for credit and evaluation

Before ARDC

- **Development** practices and tools (VCS, build system, test suites, CI, code quality, ...)
- **Opening up** towards a community (documentation, organization, communication)

Need training, tooling, infrastructures, best practices

ARDC

- **Archive** for retrieval (*reproducibility*)
- **Reference** for identification (*reproducibility*)
- **Describe** for discovery and reuse
- **Cite/Credit** for credit and evaluation

Before ARDC

- **Development** practices and tools (VCS, build system, test suites, CI, code quality, ...)
- **Opening up** towards a community (documentation, organization, communication)

Need training, tooling, infrastructures, best practices

Beyond ARDC

- **Policies** (dissemination, reuse, careers, ...)
- **Sustainability** (legal, financial, etc.)
- Technology transfer
- Advanced technologies and tools (quality, traceability, etc.)

ARDC

- **Archive** for retrieval (*reproducibility*)
- **Reference** for identification (*reproducibility*)
- **Describe** for discovery and reuse
- **Cite/Credit** for credit and evaluation

Before ARDC

- **Development** practices and tools (VCS, build system, test suites, CI, code quality, ...)
- **Opening up** towards a community (documentation, organization, communication)

Need training, tooling, infrastructures, best practices

Beyond ARDC

- **Policies** (dissemination, reuse, careers, ...)
- **Sustainability** (legal, financial, etc.)
- Technology transfer
- Advanced technologies and tools (quality, traceability, etc.)

Here we will focus on ARDC

Some popular approaches, and why they do not fit the bill

A - Since the 1970's 1990's

.zip or .tar file on:

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

Some popular approaches, and why they do not fit the bill

A - Since the 1970's 1990's

.zip or .tar file on:

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

B - Since the 2000's

Rely on *software forges*

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

Some popular approaches, and why they do not fit the bill

A - Since the 1970's 1990's

.zip or .tar file on:

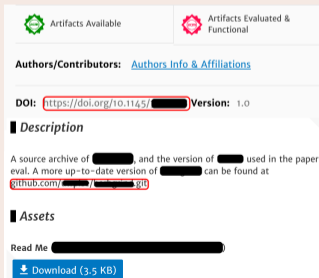
- ~~ftp server~~ (e.g. [gnu](#))
- **web page** (e.g. [myself](#))
- **document archive** (+ DOI [sample](#))

B - Since the 2000's

Rely on *software forges*

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

C: a mix of the two



The screenshot shows a software artifact page with the following details:

- Artifacts Available (green icon)
- Artifacts Evaluated & Functional (red icon)
- Authors/Contributors: [Authors Info & Affiliations](#)
- DOI: <https://doi.org/10.1145/1234567890> (highlighted with a red box)
- 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/1234567890 (highlighted with a red box).
- Assets: Read Me [redacted]
- Download (3.5 KB) button

Some popular approaches, and why they do not fit the bill

A - Since the 1970's 1990's

.zip or .tar file on:

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

B - Since the 2000's

Rely on *software forges*

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

C: a mix of the two

The screenshot shows a software artifact page with the following details:

- Two status indicators: "Artifacts Available" (green) and "Artifacts Evaluated & Functional" (red).
- Section: "Authors/Contributors: [Authors Info & Affiliations](#)"
- DOI: <https://doi.org/10.1145/...> Version: 1.0
- Section: "Description" containing text: "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)"
- Section: "Assets" with a "Read Me" file and a "Download (3.5 KB)" button.

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)

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**

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

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:

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 Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection**
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions





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



United Nations
Educational, Scientific and
Cultural Organization



And many more ...

www.softwareheritage.org/support/testimonials

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

openinventionnetwork



Silver sponsors



Bronze sponsors



The largest software archive, a shared infrastructure

Cultural Heritage



Industry



Research



Public Administration



Software Heritage

The largest software archive, a shared infrastructure

Cultural Heritage



Industry



Research



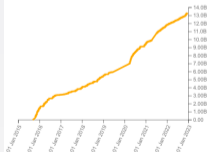
Public Administration



Software Heritage

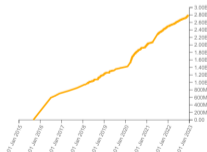
Source files

13,338,879,609



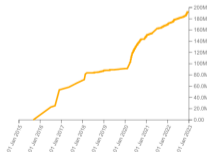
Commits

2,801,474,518



Projects

193,279,669



Directories

10,905,976,959

Authors

51,691,686

Releases

35,938,300

The largest software archive, a shared infrastructure

Cultural Heritage



Industry



Research



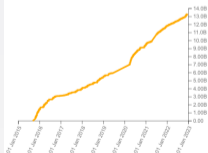
Public Administration



Software Heritage

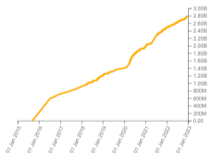
Source files

13,338,879,609



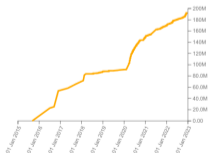
Commits

2,801,474,518



Projects

193,279,669



Directories

10,905,976,959

Authors

51,691,686

Releases

35,938,300

Bitbucket

1,925,997 origins

git

21,603 origins

R

21,113 origins

debian

128,719 origins

gnome

5,947 origins

GitHub

137,564,899 origins

GitLab

3,982,586 origins

Guix

12,032 origins

GNU

354 origins

heptapod

1,068 origins

launchpad

329,908 origins

Maven

93,738 origins

NixOS

12,032 origins

npm

1,799,296 origins

Python

4,083 origins

Phabricator

192 origins

python

410,582 origins

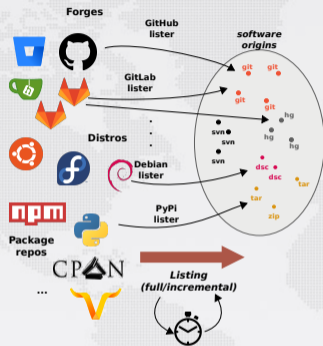
SOURCEFORGE

308,990 origins

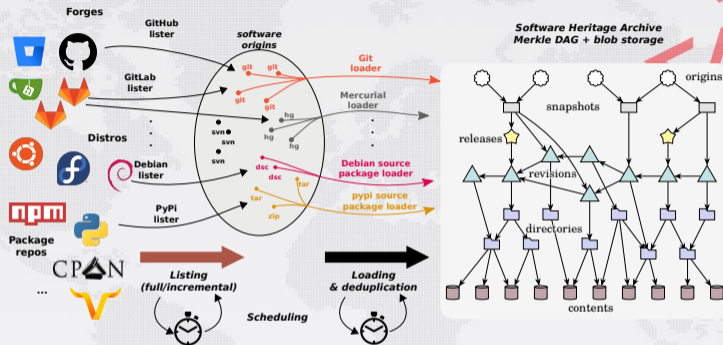
- 1 Introduction
- 2 Software and Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference**
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions



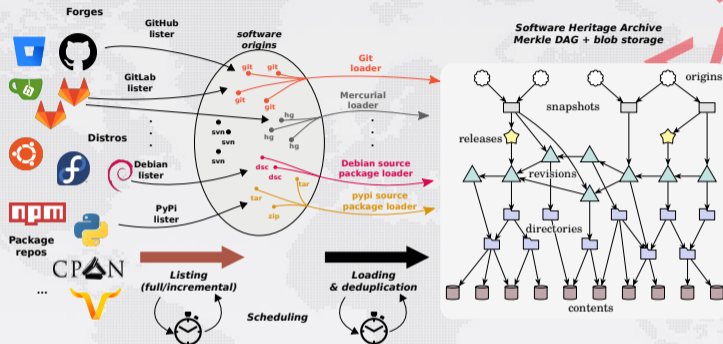
Address common Open Science and Open Source needs: archival



Address common Open Science and Open Source needs: archival



Address common Open Science and Open Source needs: archival



Global development history permanently archived in a uniform data model

- over 13 billion unique source files from over 200 million software projects
- ~1PB (uncompressed) blobs, ~25 B nodes, ~350 B edges

Address common Open Science and Open Source needs: reference

Software Heritage Identifiers (SWHID)

[link to full docs](#)

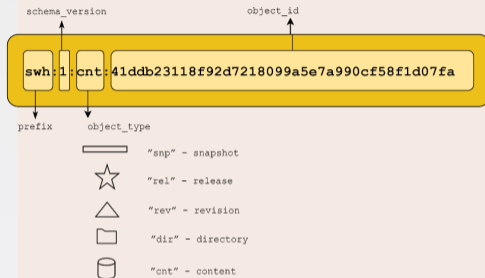


25+B
intrinsic,
decentralised,
cryptographic

Address common Open Science and Open Source needs: reference

Software Heritage Identifiers (SWHID)

[link to full docs](#)



25+B
intrinsic,
decentralised,
cryptographic

Address common Open Science and Open Source needs: reference

Software Heritage Identifiers (SWHID)

[link to full docs](#)



25+B
intrinsic,
decentralised,
cryptographic

Address common Open Science and Open Source needs: reference

Software Heritage Identifiers (SWHID)

[link to full docs](#)



25+B
intrinsic,
decentralised,
cryptographic

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

- Linux Foundation [SPDX 2.2](#)
- IANA-registered "swh:"
- WikiData property [P6138](#)

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

Breaking news: standardisation, see [swhid.org](#)

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!

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](#)
 - [guidelines](#) and [a full article](#)
 - 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 **200 million software origins**, your researchers' work may already be there (actually, [here](#))!

Getting software archived

- **automated harvesting:** over 200 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"

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"

Upcoming on replicabilitystamp.org (please do not spoil :-))

Preview

Large Growth Deformations of Thin Tissue using Solid Shells

Daniel Hwang and Ilan Pass
IEEE Transactions on Visualization and Computer Graphics (TVCG)

doi [Repository](#)

10.21203/rs.3.rs-3184856/v1

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"

Upcoming on replicabilitystamp.org (please do not spoil :-))

Large Growth Deformations of Thin Tissue using Solid Shells

Daniel Heiligman Iliescu
IEEE Transactions on Visualization and Computer Graphics (TVCG)

doi [Repository](#)

archived swin1.snp.5f1848561119a066354416eb54ad94011865fc81

HAL+SWH in the Open Science software booklet

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"

Upcoming on replicabilitystamp.org (please do not spoil :-))

Preview

Large Growth Deformations of Thin Tissue using SolidShells

Daniel Harrigan, Ilja Blass
IEEE Transactions on Visualization and Computer Graphics (TVCG)

doi

Repository

archived swf1.snp.5f1848561119a066354416eb54ad94011865fc81

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.

- 1 Introduction
- 2 Software and Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 Actions



Software metadata: `codemeta.json`

- example from [Parmap](#), created using the [Codemeta generator](#)

Software metadata: codemeta.json

- [example from Parmap](#), created using the [Codemeta generator](#)

Integration with the HAL national french open access archive

- **Curated deposit:** metadata quality due to moderation
 - all pieces of the puzzle together: one researcher does all the steps (Parmap)

Software metadata: codemeta.json

- [example from Parmap](#), created using the [Codemeta generator](#)

Integration with the HAL national french open access archive

- **Curated deposit:** metadata quality due to moderation
 - all pieces of the puzzle together: one researcher does all the steps (Parmap)
- export of citation information for [biblatex-software](#)
- examples: [LinBox](#), [SLALOM](#), [Givaro](#), [NS2DDV](#), [SumGra](#), [Coq proof](#), ...
- generation of reports, cv, web pages: [for Inria](#), [for CNRS](#), [for CNES](#), [for LIRMM](#) or [for Rémi Gribonval](#) using [HalTools](#)

Software Heritage + a *curated* metadata repository allows to address all needs ...

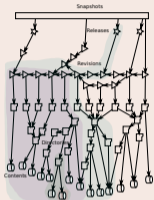
- *researcher, engineer*: archival, reference, credit, CV etc. *with a little effort from them*
- *labs, organizations*: track and report software production in a simple way
- *technology transfer offices*: view the software production
- *national level*: a *curated* catalog of the software production

- 1 Introduction
- 2 Software and Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more**
- 9 Actions



A few words about what we did not see

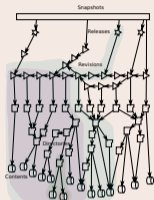
The *graph* of Software Development



All software development
in a **single graph** ...

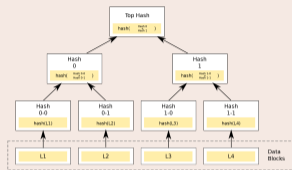
A few words about what we did not see

The *graph* of Software Development



All software development
in a **single graph** ...

The *blockchain* of Software Development



... a single
Merkle graph!

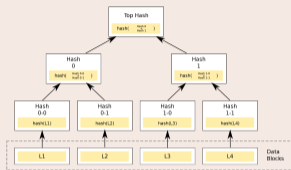
A few words about what we did not see

The *graph* of Software Development



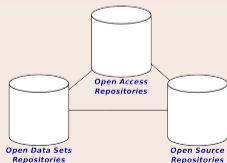
All software development
in a **single graph** ...

The *blockchain* of Software Development



... a single
Merkle graph!

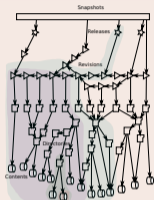
A *pillar* of Open Science



Reference **archive** of
Research Software

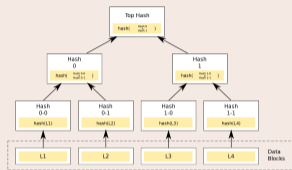
A few words about what we did not see

The *graph* of Software Development



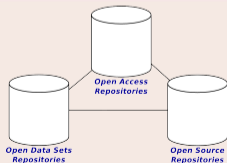
All software development
in a **single graph** ...

The *blockchain* of Software Development



... a single
Merkle graph!

A *pillar* of Open Science



Reference **archive** of
Research Software

Reference platform for *Big Code*



A **single, uniform** data structure

- 1 Introduction
- 2 Software and Open Science
- 3 An emerging policy framework
- 4 Assessing the needs and a strategy to address them
- 5 Meet Software Heritage and the HAL connection
- 6 Archive and reference
- 7 Describe, cite, credit
- 8 There is much more
- 9 **Actions**



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)

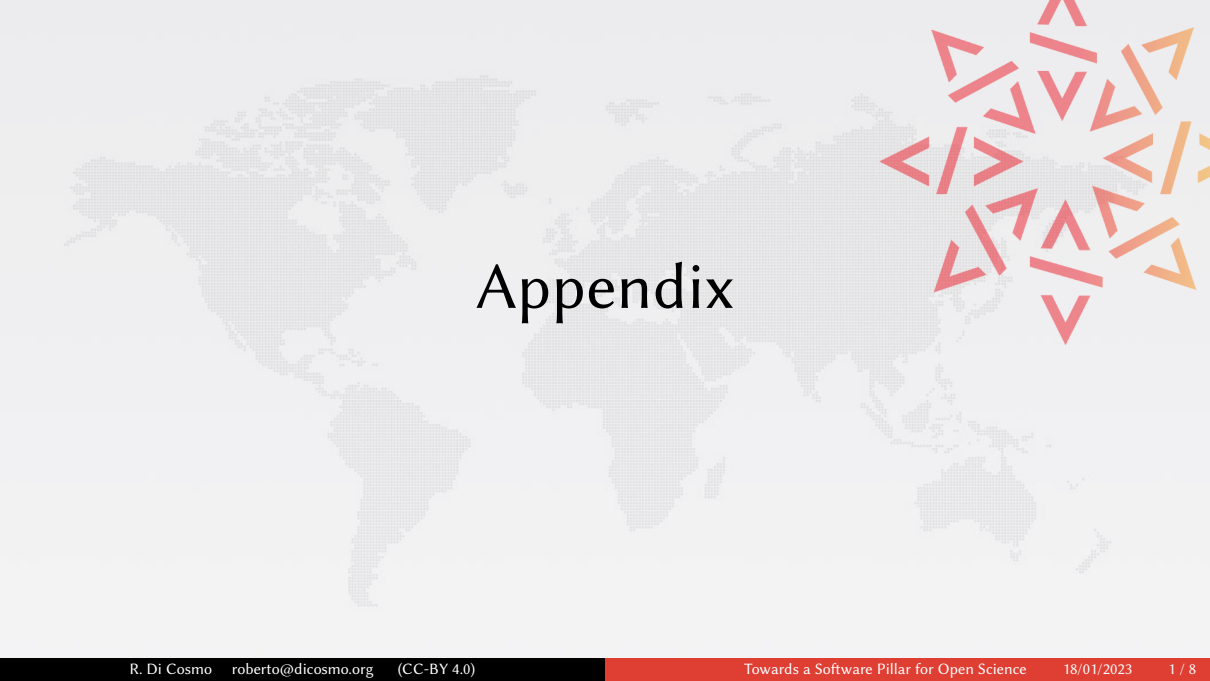
- train students and colleagues
- engage journals, conferences, learned societies

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

Questions?

References

-  UNESCO, *Draft recommendations on Open Science* 2021, ([online](#))
-  French Ministry of Research, *Second National Plan for Open Science* 2021, ([online](#))
-  EOSC SIRS Task Force, *Scholarly Infrastructures for Research Software* 2020, Publications office of the European Commission, ([10.2777/28598](#))
-  R. Di Cosmo, *Archiving and Referencing Source Code with Software Heritage* International Conference on Mathematical Software 2020 ([10.1007/978-3-030-52200-1_36](#))
-  J.F. Abramatic, R. Di Cosmo, S. Zacchiroli, *Building the Universal Archive of Source Code* CACM, October 2018 ([10.1145/3183558](#))



Appendix



10 SWHIDs by the example

11 Public code, mirrors

12 FAIR

13 END

A word on the trust model for systems of identifiers

Two general classes of systems of identifiers

intrinsic *computed from the object (no registry required, fully decentralised)*
(e.g.: chemical notation, music notation, hashes, SWHIDs)

A word on the trust model for systems of identifiers

Two general classes of systems of identifiers

intrinsic *computed from the object (no registry required, fully decentralised)*
(e.g.: chemical notation, music notation, hashes, SWHIDs)

extrinsic *assigned by an authority (need a registry)*
(e.g.: passport number, DOI, ARK, RRID, etc.)

A word on the trust model for systems of identifiers

Two general classes of systems of identifiers

intrinsic *computed from the object (no registry required, fully decentralised)*
(e.g.: chemical notation, music notation, hashes, SWHIDs)

extrinsic *assigned by an authority (need a registry)*
(e.g.: passport number, DOI, ARK, RRID, etc.)

See [the dedicated blog post](#) for more details

A word on the trust model for systems of identifiers

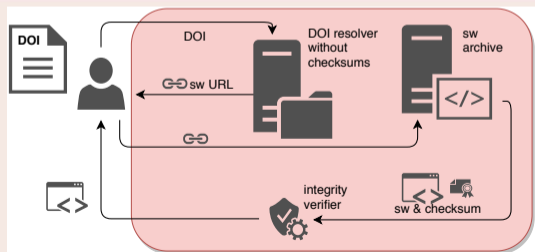
Two general classes of systems of identifiers

intrinsic *computed from the object (no registry required, fully decentralised)*
(e.g.: chemical notation, music notation, hashes, SWHIDs)

extrinsic *assigned by an authority (need a registry)*
(e.g.: passport number, DOI, ARK, RRID, etc.)

See [the dedicated blog post](#) for more details

Trust model, extrinsic (e.g. DOIs)



A word on the trust model for systems of identifiers

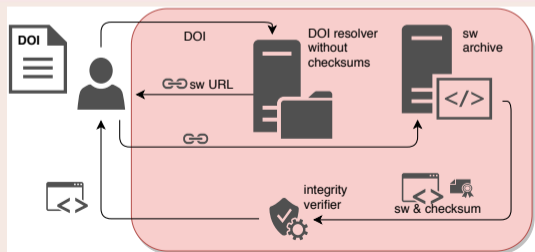
Two general classes of systems of identifiers

intrinsic *computed from the object (no registry required, fully decentralised)*
(e.g.: chemical notation, music notation, hashes, SWHIDs)

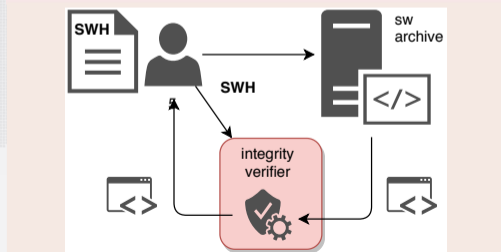
extrinsic *assigned by an authority (need a registry)*
(e.g.: passport number, DOI, ARK, RRID, etc.)

See [the dedicated blog post](#) for more details

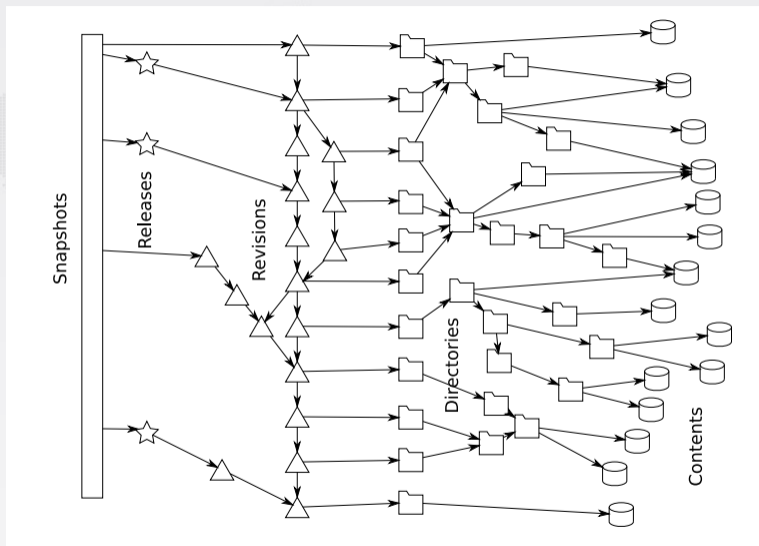
Trust model, extrinsic (e.g. DOIs)



Trust model, intrinsic (e.g. SWHIDs)



A worked example



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.

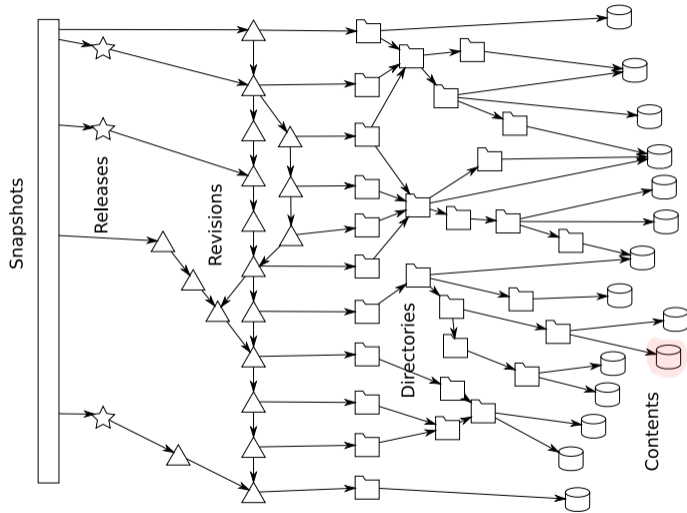
The licenses for most software and other practical works are designed
to take away your freedom to share and change the works. By contrast,
the GNU General Public License is intended to guarantee your freedom to
share and change all versions of a program--to make sure it remains free
software for all 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 and to make sure you have received the full text of the
```

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

A worked example



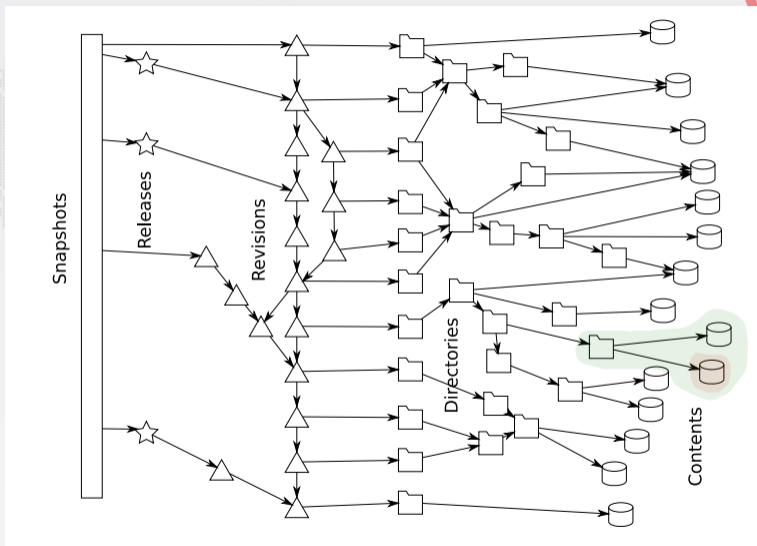


Directories


```
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 e1e10ecf948af0b93adb0372afc89f12e92618a bin
040000 tree 83e56d0beaf7793c77a45a345c80fcb8af503013 debian
040000 tree a34c9c4ba213f0cedc67f9816348d27955577af5 docs
100644 blob f2a6d32c6135aa7287bbd76167b01df2ae4f1539 requirements.txt
100755 blob eee147c36caf1bbc2d820da8dc026cb5b68180bc setup.py
040000 tree 224bb4c1f4c67fca1d160bfd2d06094e7e1abf3 sql
040000 tree 8631c9cd77bbe993168107ab5baf51f40c6300be swl
040000 tree 8fb905b56ba8ed692f1209b2773b474c6c1d66c1 utils
```

id: 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d

A worked example



Revisions

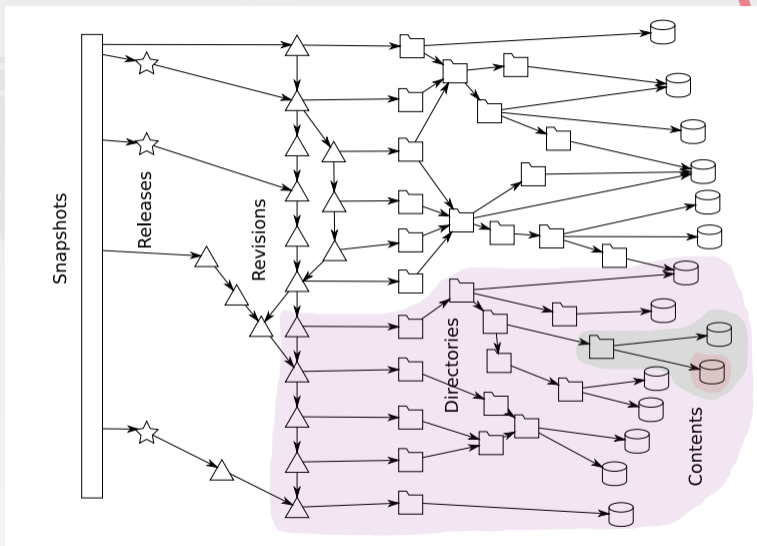
Details	Changes	Files
SHA: 963634dca6ba5dc37e3ee426ba091092c267f9f6		
Author: Nicolas Dandrimont <nicolas@dandrimont.eu> (Thu Sep 1 14:26:13 2016)		
Committer: Nicolas Dandrimont <nicolas@dandrimont.eu> (Thu Sep 1 14:26:13 2016)		
Subject: provenance.tasks: add the revision -> origin cache task		
Parent: fc3a8b59ca1df424d860f2c29ab07fee4dc35d10 : test...storage: properly pipeline origin and cont...		
provenance.tasks: add the revision -> origin cache task		
swh/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 worked example



Releases

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

```
Release sw-h.storage v0.0.51
```

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

```
commit c0c9f16b1e134f593e7567570a1761b156e6eb1d
```

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

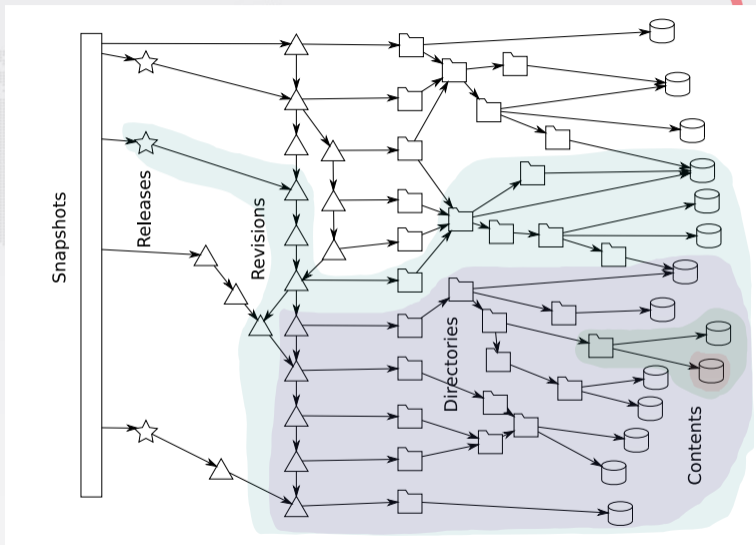
```
Release sw-h.storage v0.0.51
```

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

```
iQIzBAABCAAdBQJXvZTNFhXuaWNvbGFzQGRhbmRyaW1vbnQuZXUACgkQ7AWLMo2+
neqorw/aaq65Ob5DijzEa+kWN3rXgV5+1K1vEVh1wNKAwX8eKJ7aX2kEiLdt7uf
ahpZ6pz3q8nqs6aC1+YrxBfcih3L2YtrdZeWXXWqr8xWNMaEoYDb8qaphwh8AD5t2
ICBlit2ujtXuCrDt93eKKPwvzZxg+h80sMWy35Dr6jW7Z7K4Mu/PgGlyLHPY55yo
IGEndWno7VfH1Vm6t1n5qB7l5mXRaqA+becqddbTZ2xij+jpIUqC8cyqN3hm/fL
qsj2mu8kyz3t8tG/H1/pV+I5OwBlNpO5STH0tujojEVgPK/dHSP79QuHDHZFkCao
kij6kAWyU80Mxb+nKVjeLbrR3+yWBFj3Qp5a1/V8o0Th6E1dALcNMpEaKCoKtMt
d/gMRax1l1/g0EDfnsW67G6sDwKPKPhngfVLQ3nV3GaQQTnu1RpMz006H9/tAwzC
Gg/K1PdHT4hz0iI46wYPZyje0U2VXGFu6vVU9vFQ4ZR/Wjn+0zMzdcRdrJlSUOMn
RpTTfUbsXUeXHGOpkgXhSYTnvp1gdPc76U5TsK0aGe84AZm1Ik0mGrwXCvFPqYo
nhhibB5HBNMoqyF6yTSOpUbYK70tpYRRUGKwDeRk0wKSxkWKUZGtKzy6jYqJjo29
gulwgZQif5qWQC80ontL2+HvFfaVyckMejUhg62cP/+EHlvUk=
=kOxP
---END PGP SIGNATURE---
```

id: 85083a5cc14a441c89dea73f5bdf67c3f9c6afdb

A worked example



Snapshots

git show-refs

```
commit 08ffeb25770109525eb3ce21691466c53a1d9158 refs/heads/atime
commit ba5443a24e3f9fe323a46c292cec4fcbe61c67eb refs/heads/directory-listing-arrays
commit d69e0dbf892383ff6589b27fbelc05d27238d9c5 refs/heads/foo
commit c7ff9eea0eb22f8946908f5a8019f67de468e08 refs/heads/master
commit 7eca197fc66d2024047e54b1ed9e8b44361a0fc2 refs/heads/tmp-directory-add
commit 642a205f37de85005a85d427b53ee4fb2252e82e refs/heads/tmp/generic-releases
tag 20f043b1379cf768d966597799fd4907c757f755 refs/tags/v0.0.1
tag 72a21991a384e539996dbb867bfb0bee72aee2cd refs/tags/v0.0.10
tag 3590e0ca0ebbb070e5b376705fa230bbfa4ffa5cc refs/tags/v0.0.11
tag 33378427a403ba569a67777b8d58f6674fbc6556 refs/tags/v0.0.12
tag 06f74652755b327cf590311c2bfa036cf3b4b35d refs/tags/v0.0.13
tag 5a6325fe86ab854b581d7442667d92a11e32f3bd refs/tags/v0.0.14
tag 586fba4e580b4f5fab05f599367643c3cb1a9c7f refs/tags/v0.0.15
tag 8cd8b885f4098bf36317742bd289f660e5be51c refs/tags/v0.0.16
tag a542444ee3f0fbcd35efb202fee035c809abc7d6 refs/tags/v0.0.17
tag 228a2f1650dd12222e556559462e1e06fc4993d9 refs/tags/v0.0.18
tag 606979a4ca05d497fc0d24aad0dce82636ef47c refs/tags/v0.0.19
tag 32bf5a59fc2a323baa6d5f15a6ad5382ec275a67 refs/tags/v0.0.2
tag 3147c3d31ec46cf6492f881e908b1237ebdff2c7 refs/tags/v0.0.20
tag 215ea50daball1e082e0b72e76eb4b6073a87908 refs/tags/v0.0.21
tag 3fb168c2072a5d6252124257a1e5dfc0f5ffa1df refs/tags/v0.0.22
tag 8cdbee8da4d73fc5d262789e460a16ac3c72aba4 refs/tags/v0.0.23
...
```

id: b464cad1b66fff266a37b46ea6e7a04b545e904b

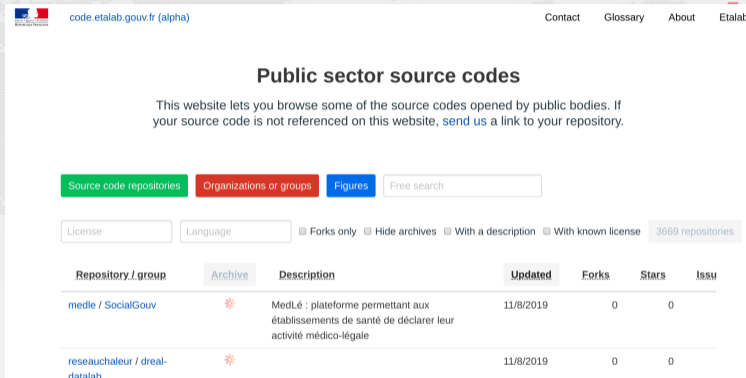


10 SWHIDs by the example



11 Public code, mirrors

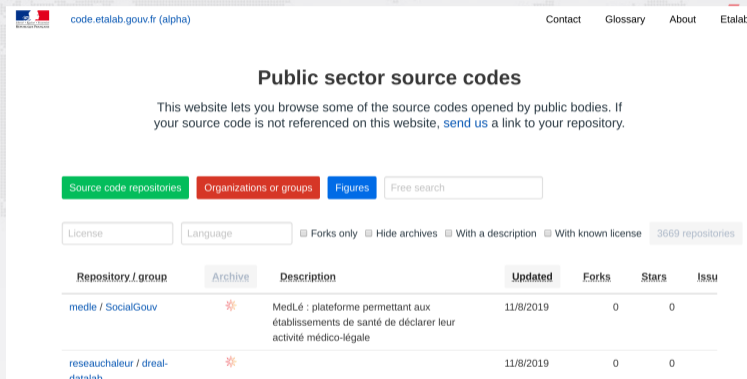
12 FAIR

13 END



The screenshot shows the website code.etalab.gouv.fr (alpha). The main heading is "Public sector source codes". Below it, a paragraph states: "This website lets you browse some of the source codes opened by public bodies. If your source code is not referenced on this website, [send us](#) a link to your repository." There are three colored buttons: "Source code repositories" (green), "Organizations or groups" (red), and "Figures" (blue). A search bar labeled "Free search" is present. Below these are filters for "License" and "Language", and checkboxes for "Forks only", "Hide archives", "With a description", and "With known license". A badge indicates "3669 repositories". A table lists two repositories:

Repository / group	Archive	Description	Updated	Forks	Stars	Issu
medle / SocialGouv		MedLé : plateforme permettant aux établissements de santé de déclarer leur activité médico-légale	11/8/2019	0	0	
reseauchaleur / dreal-datalab			11/8/2019	0	0	



The screenshot shows the homepage of code.etalab.gouv.fr (alpha). The page features a navigation bar with links for Contact, Glossary, About, and Etalab. The main heading is "Public sector source codes", followed by a descriptive paragraph. Below this are several filters: "Source code repositories" (green), "Organizations or groups" (red), "Figures" (blue), and a "Free search" input field. Further down, there are filters for "License" and "Language", along with checkboxes for "Forks only", "Hide archives", "With a description", and "With known license". A badge indicates "3669 repositories". The main content is a table with columns for Repository / group, Archive, Description, Updated, Forks, Stars, and Issu.

Repository / group	Archive	Description	Updated	Forks	Stars	Issu
medle / SocialGouv	*	MedLé : plateforme permettant aux établissements de santé de déclarer leur activité médico-légale	11/8/2019	0	0	
reseauchaleur / dreal-datalab	*		11/8/2019	0	0	

<https://code.etalab.gouv.fr>

Thomas Jefferson, February 18, 1791

... let us save what remains: not by vaults and locks which fence them from the public eye and use in consigning them to the waste of time, but by such a multiplication of copies, as shall place them beyond the reach of accident.

Thomas Jefferson, February 18, 1791

... let us save what remains: not by vaults and locks which fence them from the public eye and use in consigning them to the waste of time, but by such a multiplication of copies, as shall place them beyond the reach of accident.

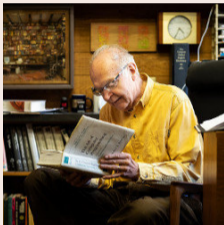
Welcoming ENEA



Italian National Agency for New Technologies,
Energy and Sustainable Economic Development

- first **institutional** mirror
- increased resilience
- **AI infrastructure** for researchers
- stepping stone to
an European joint effort

Communications of the ACM, February 2021



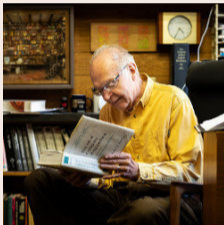
"Telling historical stories is the best way to teach. It's much easier to understand something if you know the threads it is connected to."

Let's Not Dumb Down the History of Computer Science

Donald E. Knuth, Len Shustek

<https://doi.org/10.1145/3442377>

Communications of the ACM, February 2021



"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."

Let's Not Dumb Down the History of Computer Science

Donald E. Knuth, Len Shustek

<https://doi.org/10.1145/3442377>

A unique opportunity

most of the creators are still here: we can talk to them!

but the clock is ticking...

Source code history for Security and Transparency

Where does reused software come from?



A word cloud featuring various source code hosting and distribution platforms. The most prominent words are 'Git Hub' and 'Sourceforge'. Other visible words include 'Debian', 'CPAN', 'Maven', 'Bitbucket', 'GoogleCode', 'Gitlab', 'CTAN', 'BerliOs', 'Adullact', 'Inria', and 'Gitorious'. The words are rendered in different colors and fonts, creating a dynamic and colorful composition.

Source code history for Security and Transparency

Where does reused software come from?



Do you know where it comes from?

- the software you ship
- the software you use
- the software you acquire
- the software that
 - has that bug
 - has that vulnerability

Source code history for Security and Transparency

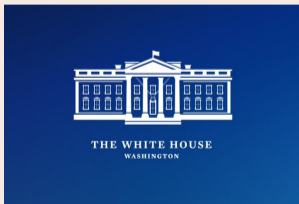
Where does reused software come from?



Do *you* know where it comes from?

- the software you ship
- the software you use
- the software you acquire
- the software that
 - has that bug
 - has that vulnerability

KYSW: Know Your SoftWare



Like KYC in banking, KYSW is now essential all over IT...

Sec. 4. Enhancing Software Supply Chain Security
ensuring and attesting, to the extent practicable, to the integrity and provenance of open source software

May 2021 POTUS Executive Order

10 SWHIDs by the example

11 Public code, mirrors

12 FAIR

13 END



What about FAIR? (Findable, Accessible, Interoperable, Reusable)

FAIR data principles *for data*

in a nutshell: metadata, metadata, metadata all over the place (makes sense for data)

What about FAIR? (Findable, Accessible, Interoperable, Reusable)

FAIR data principles *for data*

in a nutshell: metadata, metadata, metadata all over the place (makes sense for data)

But software is *not data* ...

the terms *interoperability* and *reusability* have precise technical meaning for software, and *differ significantly* from what is intended by the I and R of FAIR;

- see the entries for [software interoperability](#) and [software reusability](#)
- it is *very difficult* to achieve these properties even for commercial software developed by multi billion dollars corporations

What about FAIR? (Findable, Accessible, Interoperable, Reusable)

FAIR data principles *for data*

in a nutshell: metadata, metadata, metadata all over the place (makes sense for data)

But software is *not data* ...

the terms *interoperability* and *reusability* have precise technical meaning for software, and *differ significantly* from what is intended by the I and R of FAIR;

- see the entries for [software interoperability](#) and [software reusability](#)
- it is *very difficult* to achieve these properties even for commercial software developed by multi billion dollars corporations

FAIR for software is a distraction

let's focus on the real issues at stake: ARDC a good starting point



10 SWHIDs by the example

11 Public code, mirrors

12 FAIR

13 END