

Software Heritage

The universal source code archive

Roberto Di Cosmo, Morane Guenpeter

Director, Software Heritage

Computer Science full professor, Inria and IRIF

`roberto@dicosmo.org`

July 11th, 2018



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 
- 1 Software Heritage
 - 2 Relevance for research software publishing



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Collect, preserve and share the source code of all the software

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



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Collect, preserve and share the source code of all the software

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

Reference catalog



find and reference **all** the
source code



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Collect, preserve and share the source code of all the software

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

Reference catalog



find and reference **all** the
source code

Universal archive



preserve **all** the source
code



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Collect, preserve and share the source code of all the software

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

Reference catalog



find and reference **all** the
source code

Universal archive



preserve **all** the source
code

Research infrastructure



enable analysis of **all** the
source code

Cultural Heritage



Industry



Research



Education



Software Heritage



Cultural Heritage



Industry



Research



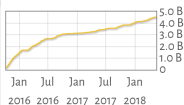
Education



Software Heritage

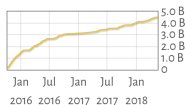
Source files

4,536,067,027



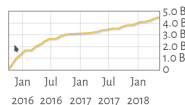
Commits

1,024,675,748



Projects

83,801,775



Cultural Heritage



Industry



Research



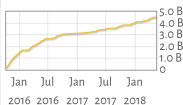
Education



Software Heritage

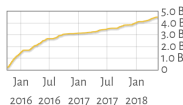
Source files

4,536,067,027



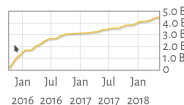
Commits

1,024,675,748



Projects

83,801,775



Open approach

- open source
- transparency

In for the long haul

- non profit, **mirrors**
- **intrinsic** identifiers

Exhaustive

- **all** software
- open to **all** communities

Growing Support

Raising awareness: landmark agreement, 3/4/2017; grand opening, 7/6/2018



Sharing the vision



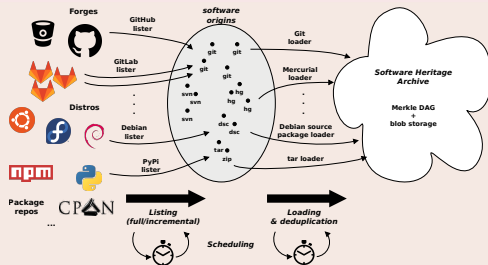
Sponsoring our work



- 
- 1 Software Heritage
 - 2 Relevance for research software publishing

Zoom on the collection phase

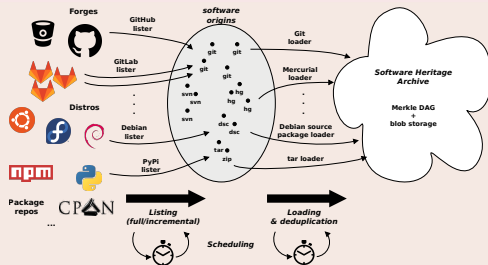
Much more than an archive!



- GitHub
- Debian, GNU
- Gitorious, Google Code
- WIP: Bitbucket, FusionForge, GitLab.com
- *add your own plugins!*

Zoom on the collection phase

Much more than an archive!



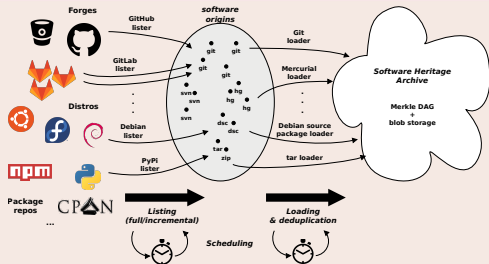
- GitHub
- Debian, GNU
- Gitorious, Google Code
- WIP: Bitbucket, FusionForge, GitLab.com
- *add your own plugins!*

Important properties

- mission: **exhaustive** and **up to date** collection of **source code**, *specifically*
- strategy: **automatic** harvesting + *deposit* from *selected* sources

Zoom on the collection phase

Much more than an archive!



- GitHub
- Debian, GNU
- Gitorious, Google Code
- WIP: Bitbucket, FusionForge, GitLab.com
- *add your own plugins!*

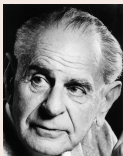
Important properties

- mission: **exhaustive** and **up to date** collection of **source code**, *specifically*
- strategy: **automatic** harvesting + *deposit* from *selected* sources

The *richest* source code archive already, ... and growing daily!

How we built our scientific knowledge

Reproducibility is the key

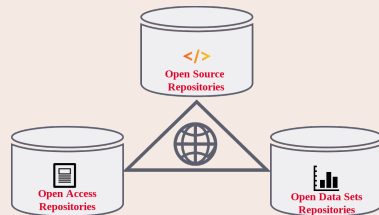


non-reproducible single occurrences are of no significance to science

Karl Popper, The Logic of Scientific Discovery, 1934

the software source code is special

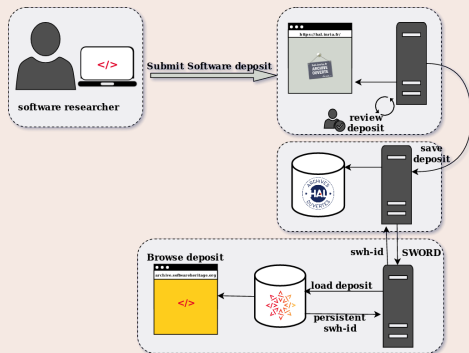
- It **embodies the logic** of the data transformation.
- It **must** be included in strategies for *scientific knowledge preservation*.
- knowing the **exact version** of the software used is essential for reproducibility.



The research software (deposit) use case

Deposit software in HAL

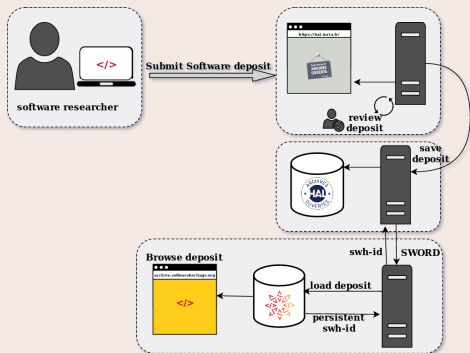
<http://hal.inria.fr/hal-01738741>



The research software (deposit) use case

Deposit software in HAL

<http://hal.inria.fr/hal-01738741>



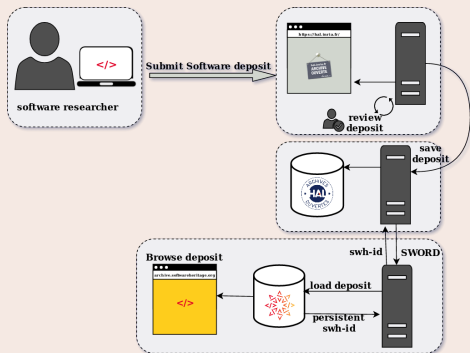
Generic mechanism:

- SWORD based
- review process
- versioning

The research software (deposit) use case

Deposit software in HAL

<http://hal.inria.fr/hal-01738741>



Generic mechanism:

- SWORD based
- review process
- versioning

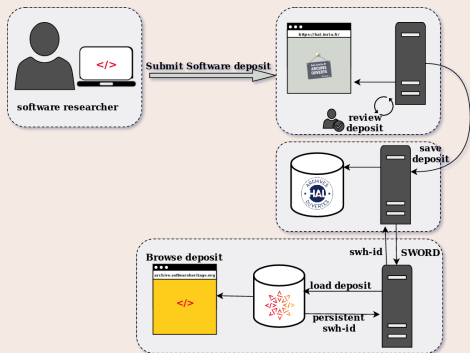
How to do it:

- **today:** deposit .zip or .tar.gz file (*guide*)
- **tomorrow:**
 - provide *SWH id* and metadata
 - include *metadata file* for automatic metadata extraction
 - ...

The research software (deposit) use case

Deposit software in HAL

<http://hal.inria.fr/hal-01738741>



Generic mechanism:

- SWORD based
- review process
- versioning

How to do it:

- **today:** deposit .zip or .tar.gz file (*guide*)
- **tomorrow:**
 - provide *SWH id* and metadata
 - include *metadata file* for automatic metadata extraction
 - ...

September 2018: **open to all** on <https://hal.archives-ouvertes.fr/>

Identifying and retrieving source code

Intrinsic identifiers

(spec: <http://bit.ly/swhpids>)

- provide **integrity** guarantees
- **all software and VCS** (not just git or GitHub)
- use for identifying **a precise version** of source code
- learn more in the forthcoming iPres 2018 paper

e.g: **swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa**

Identifying and retrieving source code

Intrinsic identifiers

(spec: <http://bit.ly/swhpids>)


- provide **integrity** guarantees
- **all software and VCS** (not just git or GitHub)
- use for identifying **a precise version** of source code
- learn more in the forthcoming iPres 2018 paper

e.g: **swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa**

"Wayback-machine-style" identifiers

- point to software **origins**
- expose the SWH crawling history

use **when no precise version is known**



Demo Time!

- example deposits in HAL
- example use of `https://archive.softwareheritage.org`

Access using intrinsic IDs

Click on the links in the paper and view the source code!

"Our **Parmap.parmap** and **Parmap.parfold** functions may be used to seamlessly ..."

The screenshot shows the Software Heritage website interface. On the left is a sidebar with icons for 'Software Heritage Archive', a book icon, a gear icon, a magnifying glass icon, a download icon, and a question mark icon. The main content area displays the source code for the 'parmap.parmap' function in Scala. The code is highlighted in green. A red vertical bar labeled 'Permalinks' is positioned over the code. To the right of the code, there is a text box explaining that to reference or cite objects in the Software Heritage archive, permalinks based on persistent identifiers must be used instead of copying and pasting the URL from the address bar. Below this text box, there is a section titled 'Select below a type of object currently browsed in order to display its associated persistent identifier and permalink.' with three buttons: 'content', 'revision', and 'snapshot'. The 'content' button is selected. Below these buttons, there is a text input field containing the permalink: 'swh:1:cnt:52dba04fcffb3b7c0206b45a3f0640c841a2c459;lines=92-101;origin=https://github.com/rdicos'. Below the input field, there are two buttons: 'Copy identifier' and 'Copy permalink'. At the bottom, there is a section titled 'Options' with two checkboxes: 'Add origin info' and 'Add selected lines info', both of which are checked.

```
90 {6 Parallel map}
91
92 val parmap : ?init:(int -> unit) -> ?finalize:(unit -> unit) -> ?ncores:int -> ?chunksize:int -> ('a -> 'b) ->
93 (** [parmap ~ncores:n f (l l)] computes [List.map f l]
94 by forking [n] processes on a multicore machine.
95 [parmap ~ncores:n f (A a)] computes [Array.map f a]
96 by forking [n] processes on a multicore machine.
97 If the optional [chunksize] parameter is specified
98 the process will be executed on blocks of size [chunksize]
99 on blocks of size [chunksize]
100 load balancing is performed on the processes
101 of the result
102
103 (** {6 Parallel iterator}
104
105 val pariter : ?init:(int -> unit) -> ?finalize:(unit -> unit) -> ?ncores:int -> ?chunksize:int -> ('a -> 'b) ->
106 (** [pariter ~ncores:n f (l l)] computes [List.map f l]
107 by forking [n] processes on a multicore machine.
108 [pariter ~ncores:n f (A a)] computes [Array.map f a]
109 by forking [n] processes on a multicore machine.
110 If the optional [chunksize] parameter is specified
111 the process will be executed on blocks of size [chunksize]
112 on blocks of size [chunksize]
113 load balancing is performed on the processes
114 of the result
115
116 (** {6 Parallel map}
117
```

To reference or cite the objects present in the Software Heritage archive, permalinks based on persistent identifiers must be used instead of copying and pasting the url from the address bar of the browser (as there is no guarantee the current URI scheme will remain the same over time).

Select below a type of object currently browsed in order to display its associated persistent identifier and permalink.

content revision snapshot

swh:1:cnt:52dba04fcffb3b7c0206b45a3f0640c841a2c459;lines=92-101;origin=https://github.com/rdicos

Options

☒ Add origin info ☒ Add selected lines info

Copy identifier Copy permalink

All features of Software Heritage *for free*

- **intrinsic IDs** (integrity, not dependent on resolvers!)
 - specification: <http://bit.ly/swhpids>
 - **iPres2018** paper: <http://bit.ly/swhpidpaper>
- browse, download (now)
- metadata, licenses, provenance (plagiarism detection), classification (wip), ...

All features of Software Heritage *for free*

- **intrinsic IDs** (integrity, not dependent on resolvers!)
 - specification: <http://bit.ly/swhpids>
 - **iPres2018** paper: <http://bit.ly/swhpidpaper>
- browse, download (now)
- metadata, licenses, provenance (plagiarism detection), classification (wip), ...

Coverage and uniformity

- **one** archive for **all** domains (industry included)
- reference *any* software, not just the deposited ones
- **git-compatible** identifiers greatly simplify workflows

All features of Software Heritage *for free*

- **intrinsic IDs** (integrity, not dependent on resolvers!)
 - specification: <http://bit.ly/swhpids>
 - **iPres2018** paper: <http://bit.ly/swhpidpaper>
- browse, download (now)
- metadata, licenses, provenance (plagiarism detection), classification (wip), ...

Coverage and uniformity

- **one** archive for **all** domains (industry included)
- reference *any* software, not just the deposited ones
- **git-compatible** identifiers greatly simplify workflows

Sustainability

... doors are open!

one infrastructure

independent non profit foundation

worldwide mirrors

Operational adoption

June 7th swMath.org points into SWH for the source code see "Code" link in, e.g.
<http://swmath.org/software/7116>

ongoing OpenAire detects source code links in articles, resolves them to SWH

September HAL opens the software deposit doors on all portals

Institutional adoption

July 4th Software Heritage is part of the french National Plan for Open Science

Organiser

- Soutenir la **Research data alliance** (RDA) et créer le chapitre français de l'alliance (RDA France).
- Soutenir **Software heritage**, la bibliothèque des codes sources

Adopt and reuse

deposit/reference research software in SWH

morane.gg@gmail.com

become a SWH mirror

roberto@dicosmo.org

Standardise

metadata

RDA source code IG

identifiers

RDA source code identification WG

Support

sponsoring / partnership

sponsorship.softwareheritage.org

donations

softwareheritage.org/donate

our own code

forge.softwareheritage.org

- 
- 3 Intrinsic PID
 - 4 Our role in the publication workflow
 - 5 The Metadata challenge
 - 6 Collection strategies

Our challenge in the PID arena

Long term

Identifiers must be there for the long term

No middle man

Identifiers must be meaningful even if resolvers go away

Integrity, not just naming

Identifier must ensure that the retrieved object is the intended one

Uniqueness by design

one name identifies a single object, and each object has only one name

Exploring the PID landscape

A lot of options out there...

URL, URI, PURL, URN, ARK, DOI, ...

... some are widely used

- articles
- data
- even software artefacts!

Exploring the PID landscape

A lot of options out there...

URL, URI, PURL, URN, ARK, DOI, ...

... some are widely used

- articles
- data
- even software artefacts!

We can get no satisfaction

of all the key criteria

Exploring the PID landscape

A lot of options out there...

URL, URI, PURL, URN, ARK, DOI, ...

... some are widely used

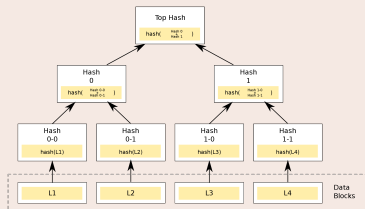
- articles
- data
- even software artefacts!

We can get no satisfaction

of all the key criteria

we adopted something radically different

Merkle tree (R. C. Merkle, Crypto 1979)



Combination of

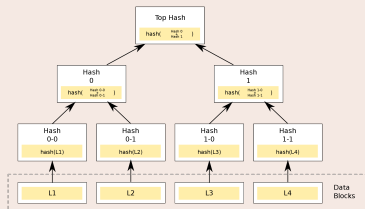
- tree
- hash function

Classical cryptographic construction

fast, parallel signature of large data structures, built-in deduplication

Intrinsic identifiers in Software Heritage

Merkle tree (R. C. Merkle, Crypto 1979)



Combination of

- tree
- hash function

Classical cryptographic construction

fast, parallel signature of large data structures, built-in deduplication

- satisfies all three criteria
- widely used in industry (e.g., Git, nix, blockchains, IPFS, ...)

DIO (digital identifier of an object)

- digital identifiers for traditional (non digital) objects
 - epistemic complications (manifestations, versions, locations, etc.)
 - significant governance issues, ...

DIO (digital identifier of an object)

- digital identifiers for traditional (non digital) objects
 - epistemic complications (manifestations, versions, locations, etc.)
 - significant governance issues, ...

IDO (identifier of a digital object)

- (digital) identifier for digital objects
 - much simpler to build/handle
 - can (and must) be intrinsic

Back to basics: DIOs vs. IDOs

DIO (digital identifier of an object)

- digital identifiers for traditional (non digital) objects
 - epistemic complications (manifestations, versions, locations, etc.)
 - significant governance issues, ...

IDO (identifier of a digital object)

- (digital) identifier for digital objects
 - much simpler to build/handle
 - can (and must) be intrinsic

Separation of concerns

- yes, we **need both** DIOs and IDOs
- no, we **must not mistake** DIOs for IDOs (and viceversa)

DIO (digital identifier of an object)

- digital identifiers for traditional (non digital) objects
 - epistemic complications (manifestations, versions, locations, etc.)
 - significant governance issues, ...

IDO (identifier of a digital object)

- (digital) identifier for digital objects
 - much simpler to build/handle
 - can (and must) be intrinsic

Separation of concerns

- yes, we **need both** DIOs and IDOs
- no, we **must not mistake** DIOs for IDOs (and viceversa)

Example: links to *software source code* in an article

Leveraging the Software Heritage universal archive:

set of files `swh:1:dir:06741c8c37c5a384083082b99f4c5ad94cd0cd1f`
id of directory object listing all the files in a project (at a given time)

revision `swh:1:rev:7598fb94d59178d65bd8d2892c19356290f5d4e3`
id of commit object which a tree and (a pointer to) the history

Full specification available online:

<https://docs.softwareheritage.org/devel/swh-model/persistent-identifiers.html>

Example: links to *software source code* in an article

Leveraging the Software Heritage universal archive:

set of files `swh:1:dir:06741c8c37c5a384083082b99f4c5ad94cd0cd1f`
id of directory object listing all the files in a project (at a given time)

revision `swh:1:rev:7598fb94d59178d65bd8d2892c19356290f5d4e3`
id of commit object which a tree and (a pointer to) the history

Full specification available online:

<https://docs.softwareheritage.org/devel/swh-model/persistent-identifiers.html>

metadata this *will* involve some form of DIO

- and we get all the complications back

- 
- 3 Intrinsic PID
 - 4 Our role in the publication workflow
 - 5 The Metadata challenge
 - 6 Collection strategies

Our role : handle *all* the *software source code*

At the end of the process

Explicit deposit, coordinated with the publisher

- store the *final* source code (no garbage)
- store only public source code
- **N.B.:** no embargo or access control (yet)

During the review

Access to the largest available source code base

- provenance, plagiarism detection (for new code)
- metrics (for long standing projects)

Our role : handle *all* the *software source code*

At the end of the process

Explicit deposit, coordinated with the publisher

- store the *final* source code (no garbage)
- store only public source code
- **N.B.:** no embargo or access control (yet)

During the review

Access to the largest available source code base

- provenance, plagiarism detection (for new code)
- metrics (for long standing projects)

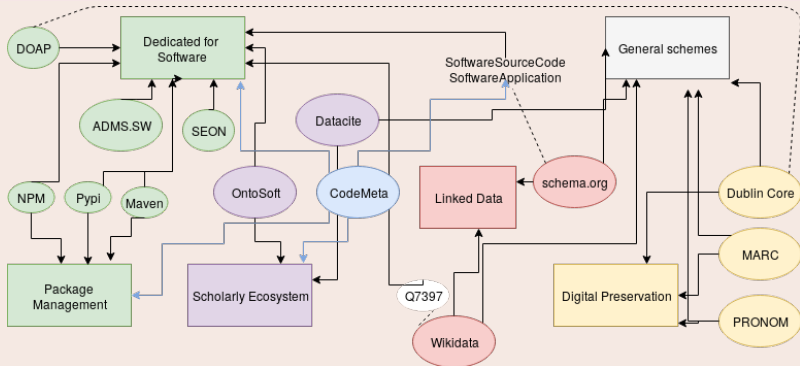
Later on

- Support embargo/access control

- 
- 3 Intrinsic PID
 - 4 Our role in the publication workflow
 - 5 The Metadata challenge
 - 6 Collection strategies

Collecting metadata for 60+ million projects

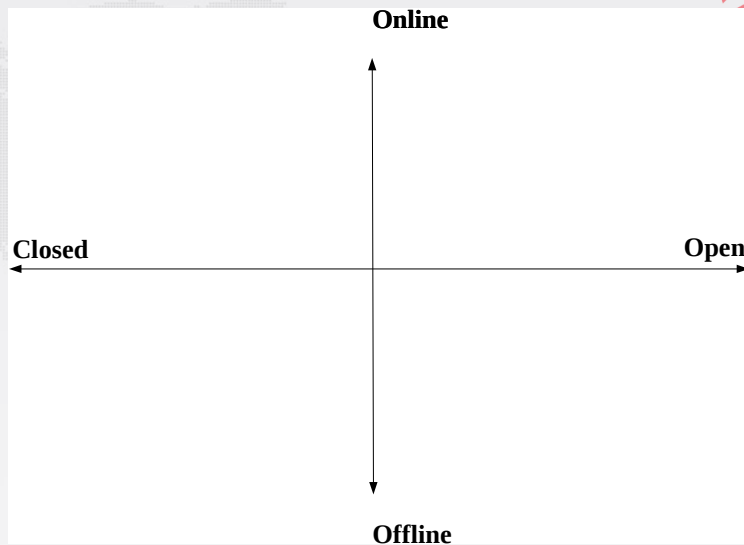
Landscape of Software Ontologies

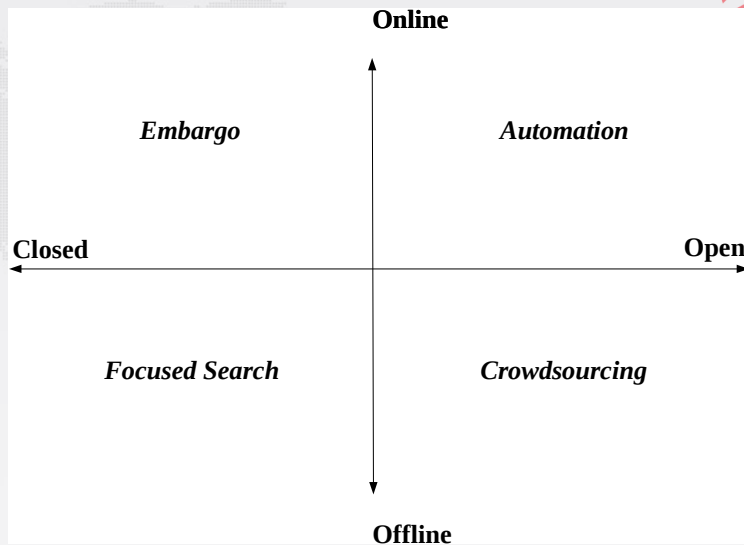


It's the real world!

reconcile metadata from different origins, handle conflicts, synthesise missing information, classify (automatically) the projects, etc.

- 
- 3 Intrinsic PID
 - 4 Our role in the publication workflow
 - 5 The Metadata challenge
 - 6 Collection strategies





Online, open source code: automation overview

