

Software Heritage: why and how

Building the Universal Archive of Source Code

Roberto Di Cosmo

`roberto@dicosmo.org`

May 30th, 2018



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we loosing trace of our knowledge?
- 4 The Software Heritage initiative
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion



Software is everywhere



Source code is *executable* and *human readable* knowledge

a growing part of our *Cultural Heritage*

Harold Abelson, Structure and Interpretation of Computer Programs

“Programs must be written for people to read, and only incidentally for machines to execute.”

Quake 2 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;
}
```

Net. queue in Linux (excerpt)

```
/*
 * SFB uses two B[l][n] : L x N arrays of bins (L levels, N bins per level)
 * This implementation uses L = 8 and N = 16
 * This permits us to split one 32bit hash (provided per packet by rxhash or
 * external classifier) into 8 subhashes of 4 bits.
 */
#define SFB_BUCKET_SHIFT 4
#define SFB_NUMBUCKETS (1 << SFB_BUCKET_SHIFT) /* N bins per Level */
#define SFB_BUCKET_MASK (SFB_NUMBUCKETS - 1)
#define SFB_LEVELS (32 / SFB_BUCKET_SHIFT) /* L */

/* SFB also uses a virtual queue, named "bin" */
struct sfb_bucket {
    u16      qlen; /* length of virtual queue */
    u16      p_mark; /* marking probability */
};
```

Len Shustek, Computer History Museum

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

~ 50 years, a lightning fast growth

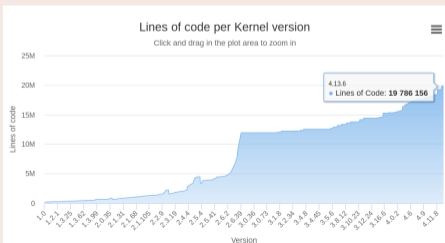
Apollo 11 Guidance Computer (~60.000 lines), 1969



"When I first got into it, nobody knew what it was that we were doing. It was like the Wild West."

Margaret Hamilton

Linux Kernel



... now in your pockets!

are we taking care of all this?

Software is spread all around

Debian CPAN
Sourceforge Gitorious
Maven Inria
Bitbucket
Git GitHub
BerliOs CTAN
GoogleCode GitLab Adullact CRAN



A word cloud of terms related to software fragility and security. The words are arranged in various orientations and colors (purple, blue, green, brown) against a background of a world map and a decorative geometric pattern of triangles in the top right corner.

damage
disaster
malicious
deletion
reference
storage
obsolete
dependencies
attack
aging
media
tear
dangling
wear
corruption
encryption
format

Software lacks its own research infrastructure



Photo: ALMA(ESO/NAOJ/NRAO), R. Hills

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we loosing trace of our knowledge?
- 4 The Software Heritage initiative
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion



How we built our scientific knowledge

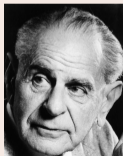
The experimental method



- make an *observation*
- formulate an *hypothesis*
- set up an **experiment**
- formulate a *theory*

And then we **reproduce** and **verify**.

Reproducibility is the key



non-reproducible single occurrences are of no significance to science

Karl Popper, The Logic of Scientific Discovery, 1934

For an experiment involving software, we need

- open access** to the scientific article describing it
- open data sets** used in the experiment
- source code** of all the components
- environment** of execution
- stable references** between all this

Remark

The first two items are already widely discussed!

... what about *software*?

Collberg's report from the trenches

Analysis of 613 papers

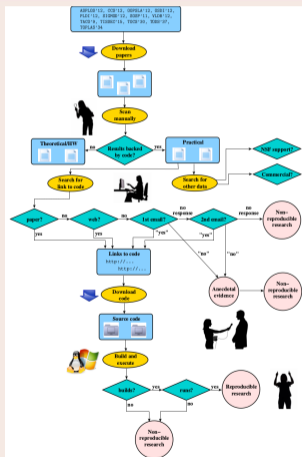
- 8 ACM conferences: ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12
- 5 journals: TACO'9, TISSEC'15, TOCS'30, TODS'37, TOPLAS'34

all very practical oriented

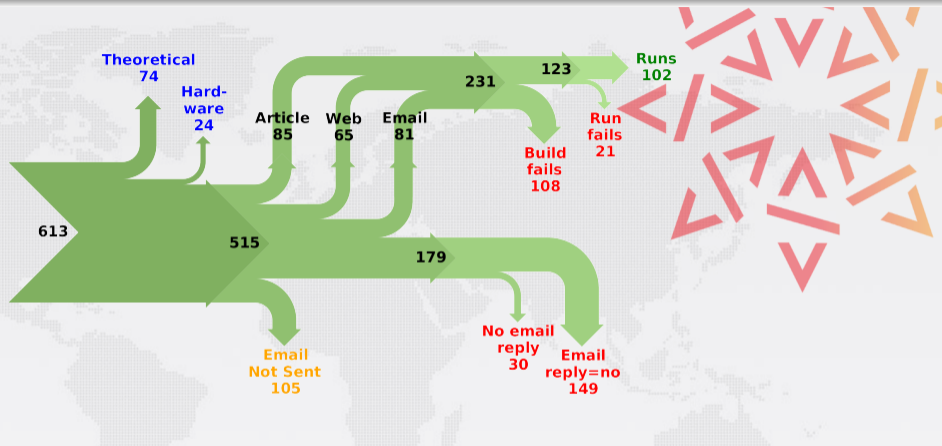
The basic question

can we get the code to build and run?

The workflow



The result



This can be debated (see <http://cs.brown.edu/~sk/Memos/Examining-Reproducibility/>), but...

... that's a whopping 81% of **non reproducible** works!

What about our Software Engineering community?

Even higher expectations, and yet similarly disappointing results <http://fr.slideshare.net/carloghezzi18/icse-2009-keynote-15919951>

Reference journal

ACM Transactions on Software Engineering and Methodology (TOSEM)

- analysis by Carlo Ghezzi, in 2009, of TOSEM from 2001 to 2006
- 60% of papers refer to a tool
- 20% only are *installable*

Reference conference

International Conference on Software Engineering (ICSE)

- analysis by Zannier, Melrik, Maurer 2006
- complete absence of replication studies

Evaluation of software artefacts (optional)



- tools are usable, in line with expectations
- started as a contest in 2011 (ESEC/FSE) (winner *Vouillon and Di Cosmo*)
- now going mainstream: POPL'17, POPL'16, ECOOP'16, OOPSLA'16, CGO'16, VISSOFT'16, PLDI'16, CGO'15, PPOPP'15, VISSOFT'15, ISSTA'15, OOPSLA'15, PLDI'15, POPL'15, CAV'15, ECOOP'15, FSE'15, ISSTA'14, OOPSLA'14, PLDI'14, ECOOP'14, FSE'14, SAS'13, OOPSLA'13, ECOOP'13, FSE'13, FSE'11

Some people claim that having (all) the source of the code used in an experiment is *not worth the effort* (see “Replicability is not Reproducibility: Nor is it Good Science”, Chris Drummond, ICML 2009)

Sure, diversity *is* important, but:

- Source code is like the proof used in a theorem: can we really accept *Fermat statements* like “the details are omitted due to lack of space”?
- modern complex systems makes even the simplest experiment depend on a wealth of components and configuration options
- access to *all* the source code is not just necessary to *reproduce*, it is also useful to *evolve and modify*, to *build new experiments* from the old ones

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we losing trace of our knowledge?
- 4 The Software Heritage initiative
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion



URL decay disrupts the *web of reference*

Web links *are not* permanent (even *permalinks*)

there is no general guarantee that a URL... which at one time points to a given object continues to do so

T. Berners-Lee et al. Uniform Resource Locators. RFC 1738.

404

URLs used in articles *decay!*

Analysis of *IEEE Computer* (Computer), and the *Communications of the ACM* (CACM): 1995-1999

- the *half-life* of a referenced URL *is approximately 4 years* from its publication date
D. Spinellis. The Decay and Failures of URL References.

Communications of the ACM, 46(1):71-77, January 2003.

Similar findings in Lawrence, S. et al. *Persistence of Web References in Scientific Research*, *IEEE Computer*, 34(2), pp. 26-31, 2001.

An example from Astronomy

Domain	links (broken)	.html	.txt	.dat	.gz	.tar	.fits	tilde
oac.harvard.edu	802 (110)	336 (70)	0	0	4 (2)	5 (4)	1	0
heasarc.gsfc.nasa.gov	640 (33)	423 (27)	1	0	0	0	0	0
www.stsci.edu	498 (61)	205 (29)	3	0	0	0	0	15 (10)
esc.harvard.edu	471 (152)	212 (99)	0	0	0	0	0	1 (1)
ssc.spitzer.caltech.edu	427 (194)	125 (76)	3 (3)	0	0	0	0	0
cfa-www.harvard.edu	352 (68)	277 (52)	1	0	0	0	0	54 (17)
archive.stsci.edu	308 (58)	57 (9)	2	1 (0)	0	0	0	0
www.ipac.caltech.edu	285 (14)	209 (12)	0	0	0	0	0	0
www.atnf.csiro.au	211 (21)	12 (6)	0	0	0	0	0	7 (5)
space.mit.edu	193 (10)	58 (5)	1	0	0	0	0	2 (1)
www.astro.psu.edu	186 (4)	103 (1)	1	10 (1)	1	1	0	2
www.eso.org	186 (58)	54 (22)	1 (1)	0	0	0	0	4 (1)
isa.ipac.caltech.edu	163 (5)	38	0	0	1	0	0	0
www.sdss.org	156 (2)	106 (1)	0	0	0	0	0	0
hea-www.harvard.edu	125 (37)	42 (17)	1	0	0	1	0	26 (16)
physics.nist.gov	125 (3)	63 (2)	0	0	0	0	0	0
www.noao.edu	120 (3)	50 (2)	0	0	0	0	0	0
rmm.vilspa.esa.es	118 (35)	23 (19)	0	0	8 (1)	0	0	1 (1)
www.astro.princeton.edu	115 (31)	43 (14)	0	0	0	0	0	53 (12)
ad.usno.navy.mil	110 (27)	98 (22)	3 (3)	0	0	0	0	1 (1)

This table lists total number of links and broken links (HTTP status codes 3xx, 4xx, and 5xx) to top domains (domains with over 100 links) found within articles published in the four main astronomy journals between 1997 and 2008. The table also shows, for each domain, the portion of links to common filename extensions, as well as links that contain the tilde character.

doi:10.1371/journal.pone.0104798.t001

How Do Astronomers Share Data?

Pepe, Goodman, Muench, Crosas, Erdmann

[dx.doi.org/10.1371/journal.pone.0104798](https://doi.org/10.1371/journal.pone.0104798)

PLOS August 28, 2014

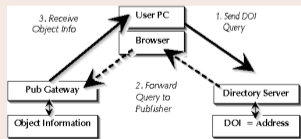
DOI limitations

Example: `doi:10.1109/MSR.2015.10`

- to find what `10.1109/MSR.2015.10` is, go to a *resolver* (e.g. `doi.org`)
- this returns `http://ieeexplore.ieee.org/document/7180064/`
- at this URL we find ...

The screenshot shows a web page with the title "Mining Component Dependencies for Installability Issues". It features a navigation bar with "View Document" (1 page), "145 Citations", and "145 Text Issues". Below the navigation bar, there is a section for "Abstract" with tabs for "Abstract", "Authors", "Figures", "References", "Citations", "Keywords", "Metrics", and "Media". The abstract text discusses component repositories and their role in software life cycle management. It also includes a "Published in:" section mentioning "Mining Software Repositories (MSR), 2015 IEEE/ACM LBNV Workshop Conference on". At the bottom, there are links for "Download PDF" and "Read the Full Document".

Architecture of the DOI infrastructure



- DOI resolution *can change*
- content at URL *can change*
- no *intrinsic* way of noticing
- persistence based on *good will* of *multiple parties*

No catalog, no archive, no references, ... and we are at a turning point

Looking at the past

- a lot of old software misplaced, lost, or behind barriers, but...
- most founding fathers are still here, and willing to share
- **urgent** to collect their knowledge

Only a few years left.

Looking at the future

- software development and use skyrockets: more programmers, and more code!
- **essential** to provide a **universal** platform for all the future software source code

Every year that goes by makes the problem worse.

it is **urgent** to take action!

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we loosing trace of our knowledge?
- 4 The Software Heritage initiative**
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion





Software Heritage



Our mission

Collect, preserve and share the *source code* of *all the software* that is available

Past, present and future

Preserving the past, enhancing the present, preparing the future

Cultural Heritage



Industry



Research



Education



Software Heritage

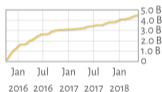
Source files

4,536,067,027



Commits

1,024,675,748



Projects

83,801,775



Technology

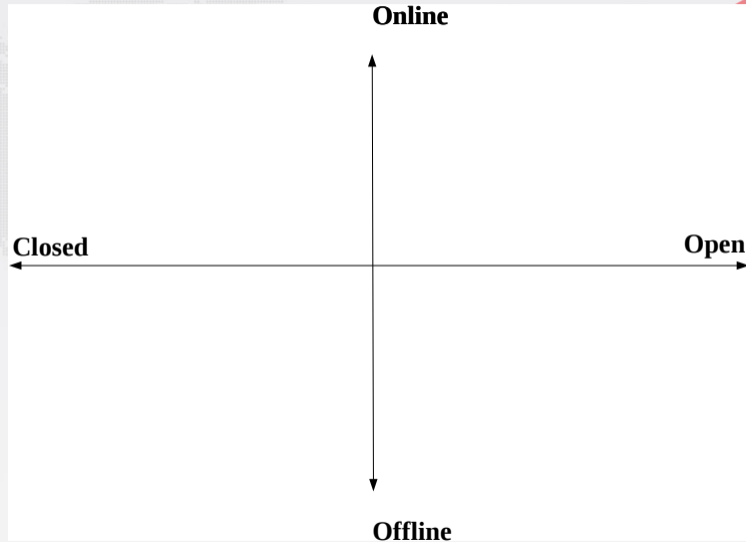
- transparency and FOSS
- replicas all the way down

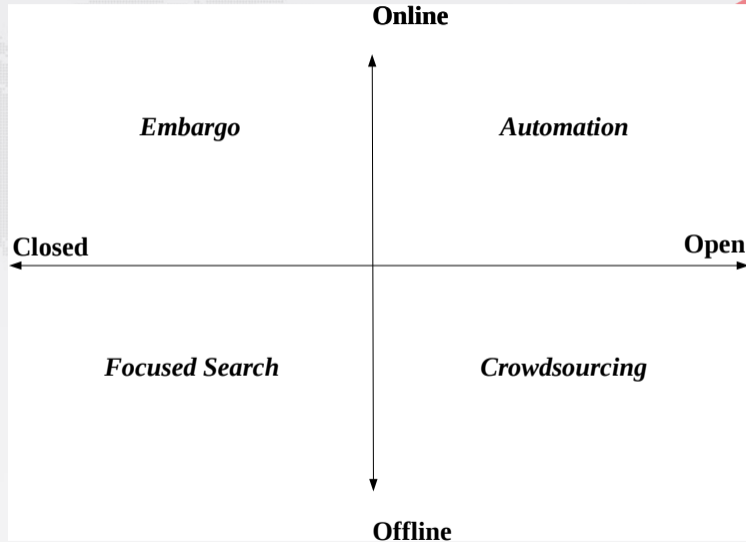
Content

- intrinsic identifiers
- facts and provenance

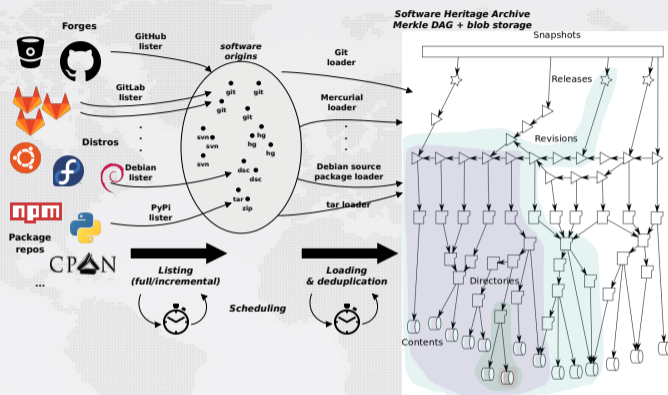
Organization

- non-profit
- multi-stakeholder





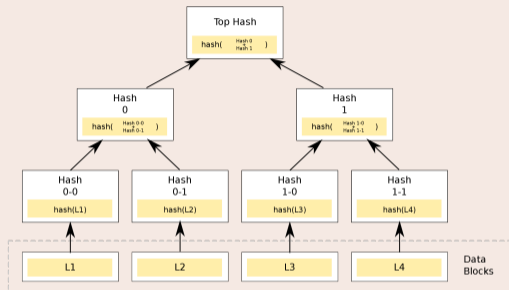
Architecture (simplified, first quadrant)



- full development history permanently archived
- origins: GitHub (automated), Debian (automated), Gitorious, Google Code, GNU
- ~200Tb raw contents, ~10Tb graph (7+Bn nodes, 60+Bn edges)

Much more than an archive!

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



Combination of

- tree
- hash function

Classical cryptographic construction

- fast, parallel signature of large data structures
- widely used (e.g., Git, blockchains, IPFS, ...)
- built-in deduplication

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we loosing trace of our knowledge?
- 4 The Software Heritage initiative
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion



Reference archive for all software

A "wayback machine" for software source code ...

and **intrinsic identifiers!**

- <http://archive.softwareheritage.org/browse> (icse / 2018)
- <http://bit.ly/swhpids> for persistent identifiers

Demo time: let's highlight some features...

Origin search

origin name	description
01 BrewerySoftware.org http://archive.softwareheritage.org/browse	
02 BrewerySoftware.org http://archive.softwareheritage.org/browse	
03 BrewerySoftware.org http://archive.softwareheritage.org/browse	
04 BrewerySoftware.org http://archive.softwareheritage.org/browse	
05 BrewerySoftware.org http://archive.softwareheritage.org/browse	
06 BrewerySoftware.org http://archive.softwareheritage.org/browse	
07 BrewerySoftware.org http://archive.softwareheritage.org/browse	
08 BrewerySoftware.org http://archive.softwareheritage.org/browse	
09 BrewerySoftware.org http://archive.softwareheritage.org/browse	
10 BrewerySoftware.org http://archive.softwareheritage.org/browse	
11 BrewerySoftware.org http://archive.softwareheritage.org/browse	
12 BrewerySoftware.org http://archive.softwareheritage.org/browse	
13 BrewerySoftware.org http://archive.softwareheritage.org/browse	
14 BrewerySoftware.org http://archive.softwareheritage.org/browse	
15 BrewerySoftware.org http://archive.softwareheritage.org/browse	
16 BrewerySoftware.org http://archive.softwareheritage.org/browse	
17 BrewerySoftware.org http://archive.softwareheritage.org/browse	
18 BrewerySoftware.org http://archive.softwareheritage.org/browse	
19 BrewerySoftware.org http://archive.softwareheritage.org/browse	
20 BrewerySoftware.org http://archive.softwareheritage.org/browse	
21 BrewerySoftware.org http://archive.softwareheritage.org/browse	
22 BrewerySoftware.org http://archive.softwareheritage.org/browse	
23 BrewerySoftware.org http://archive.softwareheritage.org/browse	
24 BrewerySoftware.org http://archive.softwareheritage.org/browse	
25 BrewerySoftware.org http://archive.softwareheritage.org/browse	
26 BrewerySoftware.org http://archive.softwareheritage.org/browse	
27 BrewerySoftware.org http://archive.softwareheritage.org/browse	
28 BrewerySoftware.org http://archive.softwareheritage.org/browse	
29 BrewerySoftware.org http://archive.softwareheritage.org/browse	
30 BrewerySoftware.org http://archive.softwareheritage.org/browse	
31 BrewerySoftware.org http://archive.softwareheritage.org/browse	
32 BrewerySoftware.org http://archive.softwareheritage.org/browse	
33 BrewerySoftware.org http://archive.softwareheritage.org/browse	
34 BrewerySoftware.org http://archive.softwareheritage.org/browse	
35 BrewerySoftware.org http://archive.softwareheritage.org/browse	
36 BrewerySoftware.org http://archive.softwareheritage.org/browse	
37 BrewerySoftware.org http://archive.softwareheritage.org/browse	
38 BrewerySoftware.org http://archive.softwareheritage.org/browse	
39 BrewerySoftware.org http://archive.softwareheritage.org/browse	
40 BrewerySoftware.org http://archive.softwareheritage.org/browse	

Directory browsing

File	Modification	Size	Version
0 glibc	2005-01-11 10:26:45	1415713	2.2.5
1 bin	2005-01-11 10:26:45	1415713	2.2.5
2 libc	2005-01-11 10:26:45	1415713	2.2.5
3 ld	2005-01-11 10:26:45	1415713	2.2.5
4 ld-linux.so.2	2005-01-11 10:26:45	1415713	2.2.5
5 ld.so	2005-01-11 10:26:45	1415713	2.2.5
6 ld.so.1	2005-01-11 10:26:45	1415713	2.2.5
7 ld.so.2	2005-01-11 10:26:45	1415713	2.2.5
8 ld.so.3	2005-01-11 10:26:45	1415713	2.2.5
9 ld.so.4	2005-01-11 10:26:45	1415713	2.2.5
10 ld.so.5	2005-01-11 10:26:45	1415713	2.2.5
11 ld.so.6	2005-01-11 10:26:45	1415713	2.2.5
12 ld.so.7	2005-01-11 10:26:45	1415713	2.2.5
13 ld.so.8	2005-01-11 10:26:45	1415713	2.2.5
14 ld.so.9	2005-01-11 10:26:45	1415713	2.2.5
15 ld.so.10	2005-01-11 10:26:45	1415713	2.2.5
16 ld.so.11	2005-01-11 10:26:45	1415713	2.2.5
17 ld.so.12	2005-01-11 10:26:45	1415713	2.2.5
18 ld.so.13	2005-01-11 10:26:45	1415713	2.2.5
19 ld.so.14	2005-01-11 10:26:45	1415713	2.2.5
20 ld.so.15	2005-01-11 10:26:45	1415713	2.2.5
21 ld.so.16	2005-01-11 10:26:45	1415713	2.2.5
22 ld.so.17	2005-01-11 10:26:45	1415713	2.2.5
23 ld.so.18	2005-01-11 10:26:45	1415713	2.2.5
24 ld.so.19	2005-01-11 10:26:45	1415713	2.2.5
25 ld.so.20	2005-01-11 10:26:45	1415713	2.2.5
26 ld.so.21	2005-01-11 10:26:45	1415713	2.2.5
27 ld.so.22	2005-01-11 10:26:45	1415713	2.2.5
28 ld.so.23	2005-01-11 10:26:45	1415713	2.2.5
29 ld.so.24	2005-01-11 10:26:45	1415713	2.2.5
30 ld.so.25	2005-01-11 10:26:45	1415713	2.2.5
31 ld.so.26	2005-01-11 10:26:45	1415713	2.2.5
32 ld.so.27	2005-01-11 10:26:45	1415713	2.2.5
33 ld.so.28	2005-01-11 10:26:45	1415713	2.2.5
34 ld.so.29	2005-01-11 10:26:45	1415713	2.2.5
35 ld.so.30	2005-01-11 10:26:45	1415713	2.2.5
36 ld.so.31	2005-01-11 10:26:45	1415713	2.2.5
37 ld.so.32	2005-01-11 10:26:45	1415713	2.2.5
38 ld.so.33	2005-01-11 10:26:45	1415713	2.2.5
39 ld.so.34	2005-01-11 10:26:45	1415713	2.2.5
40 ld.so.35	2005-01-11 10:26:45	1415713	2.2.5

Revisions as diffs

```
diff --git a/src/main/java/org/apache/struts2/struts2-plugin.jar b/src/main/java/org/apache/struts2/struts2-plugin.jar
index 123456789..987654321
--- a/src/main/java/org/apache/struts2/struts2-plugin.jar
+++ b/src/main/java/org/apache/struts2/struts2-plugin.jar
@@ -123,13 +123,13 @@
- public class MyServlet extends HttpServlet {
+ public class MyServlet extends HttpServlet {
     public void doGet() {
@@ -136,13 +136,13 @@
     public void doPost() {
@@ -149,13 +149,13 @@
     }
 }
```

A glimpse at the technical roadmap

Features...

- (done) **lookup** by content hash
- **browsing**: "wayback machine" for archived code
 - (done) <http://archive.softwareheritage.org/api>
 - (done) <http://archive.softwareheritage.org/browse/search>
- (done) **download**: `wget / git clone` from the archive
- (done) **deposit** of source code bundles directly to the archive
- (todo) **provenance** lookup for all archived content
- (todo) **full-text search** on all archived source code files

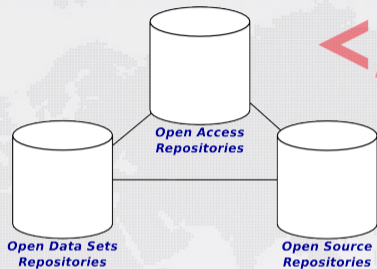
... and much more ...

you have the world's software development graph at your hands!

your tools could be here!

Paper points to lost source code on gitorious

- https://www.openaire.eu/search/publication?articleId=dedup_wf_001::cd996f0b6236b90659f84f99feb62bcc
- <https://gitorious.org/parmap>
- <https://archive.softwareheritage.org/browse/search/?url=%22gitorious.org/parmap%22>

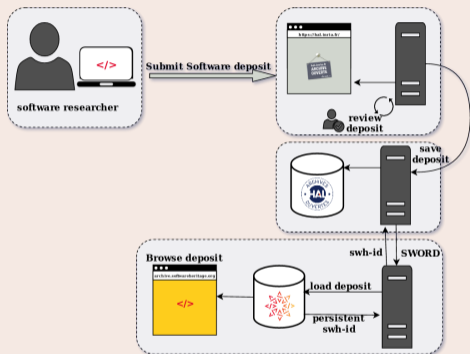


A global library referencing all software used in all research fields

- completes the infrastructure for **Open Access** in science
- provides intrinsic persistent identifiers needed for scientific **reproducibility**
- enables large scale, verifiable **software studies**

Deposit software in HAL

<http://bit.ly/swhdeposithalen>



Generic mechanism:

- SWORD based
- review process
- versioning

How to do it:

- today: deposit .zip file
- tomorrow:
 - *provide SWH id and metadata*
 - *provide SWH id, metadata is extracted*
 - ...

Intrinsic PIDS for referencing content now available

see <http://bit.ly/swhpids> and the forthcoming iPres 2018 article

The way to go to archive and reference scientific software

All features of Software Heritage *for free*

- **intrinsic IDs** (integrity, not just DIOs!), browse, download (now)
- metadata, licenses, provenance analysis (plagiarism detection), classification (wip)
- and many more (powerful connections with SE and Industry)

Coverage and uniformity

- **one** archive for **all** domains (industry included)
- you can reference *any* software, not just the deposited one
(thanks D. Katz for pointing this out)
- **git-compatible** identifiers greatly simplify workflows

Sustainability

... doors are open!

one infrastructure

independent non profit foundation

worldwide mirrors



Large scale *repeatable* software studies...

- vulnerability detection
- dependency analysis
- pattern elicitation
- automatic classification ...

... need a uniform representation

Software Heritage has **one data model** for all forges/VCS...

... yes, we do **data normalization** of software evolutiona!

Coming soon to a platform near you!

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we loosing trace of our knowledge?
- 4 The Software Heritage initiative
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion



Landmark Inria Unesco agreement, April 3rd, 2017



Sharing the vision



Contributing to the mission



The Software Heritage Foundation

- independent
- long term mission
- multistakeholder

The community

- academia: Open Access, research
- industry: better software
- cultural heritage: **all** the software history

The mirror network

- resilience
- biodiversity

“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

You can help!

Take the research challenges

- efficient tracking of development streams
- machine learning/classification
- ...

Contribute to the development see <http://forge.softwareheritage.org>

★★ listers/loaders for other unsupported forges, VCS

★★ Web UI improvements

Funding

- pester *companies* to become sponsors :
sponsorship.softwareheritage.org
- give *your own contribution* :
www.softwareheritage.org/donate

Spread the word!

- help research teams *use* the archive
- tell everybody about Software Heritage

- 1 Software is everywhere and nowhere
- 2 Software source code in Science
- 3 Are we loosing trace of our knowledge?
- 4 The Software Heritage initiative
- 5 Using the Software Heritage archive
- 6 Building for the long term
- 7 Conclusion



Come in, we're open!



Software Heritage

www.softwareheritage.org

@swheritage

Grand opening, June 7th, UNESCO headquarters!

Library of Alexandria of code



- recover the past
- structure the future

A CERN for Software



- build better software
 - for industry
 - for society as a whole