

Software Heritage

Collecting, preserving and sharing the software source code of Mankind

Roberto Di Cosmo
Inria

roberto@dicosmo.org

November 7, 2017
Evry



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood
- 7 Building for the long term
- 8 Conclusion



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

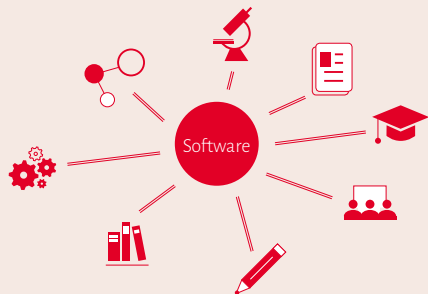
2010 *IRILL* www.irill.org

2015 *Software Heritage* at INRIA

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood
- 7 Building for the long term
- 8 Conclusion

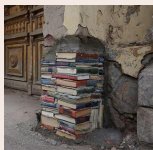


At the heart of our society



- communication, entertainment
- administration, finance
- health, energy, transportation
- education, research, politics
- ...

Key mediator for accessing all information (c) Banski



Information is a main pillar of our modern societies.

Absent an ability to correctly interpret digital information, we are left with [...] "rotting bits" [...] of no value.

Vinton G. Cerf IEEE 2011

Software is an essential component of modern scientific research

[...] the vast majority describe experimental methods or software that have become essential in their fields.

Top 100 papers (Nature, October 2014)



Software embodies our collective Knowledge and Cultural Heritage

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood
- 7 Building for the long term
- 8 Conclusion



Source code matters!



"The source code for a work means the preferred form of the work for making modifications to it."
— GPL Licence

Hello World

Program (excerpt of binary)

```
4004e6: 55
4004e7: 48 89 e5
4004ea: bf 84 05 40 00
4004ef: b8 00 00 00 00
4004f4: e8 c7 fe ff ff
4004f9: 90
4004fa: 5d
4004fb: c3
```

Program (source code)

```
/* Hello World program */
#include<stdio.h>

void main()
{
    printf("Hello World");
}
```


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.”

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.”

Distinguishing features

- *executable and human readable knowledge (an all time new)*
 - even hardware is... software! (VHDL, FPGA, ...)
 - *text files are forever*
- naturally *evolves* over time
 - the *development history* is key to its *understanding*
- complex: large *web of dependencies*, millions of SLOCs

In a word

- software *is not just another* sequence of bits
- a software archive *is not just another* digital archive

~ 50 years, a lightning fast growth

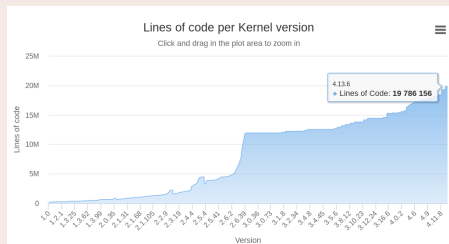
Apollo 11 Guidance Computer (~60.000 lines), 1969



"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?

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood
- 7 Building for the long term
- 8 Conclusion



Software is spread all around



Fashion victims

- many disparate development platforms
- a myriad places where distribution may happen
- projects tend to migrate from one place to another over time

Where is the place ...

where we can find, track and search *all* source code?



A word cloud of terms related to software fragility and digital information loss. The most prominent words are 'damage', 'disaster', 'malicious', 'obsolete', 'deletion', and 'format'. Other visible words include 'attack', 'dependencies', 'aging', 'media', 'tear', 'dangling', 'wear', 'corruption', 'encryption', 'reference', and 'storage'. The words are arranged in a cluster, with some overlapping.



Like all digital information, FOSS is fragile

- inconsiderate and/or malicious code loss (e.g., Code Spaces)
- business-driven code loss (e.g., Gitorious, Google Code)
- for obsolete code: physical media decay (data rot)

Where is the archive...

where we go if (a repository on) GitHub or GitLab.com goes away?



A wealth of software research on crucial issues...

- safety, security, test, verification, proof
- software engineering, software evolution
- big data, machine learning, empirical studies

If you study the stars, you go to Atacama...

... where is the *very large telescope* of source code?

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 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative**
- 6 Under the hood
- 7 Building for the long term
- 8 Conclusion





Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Our mission

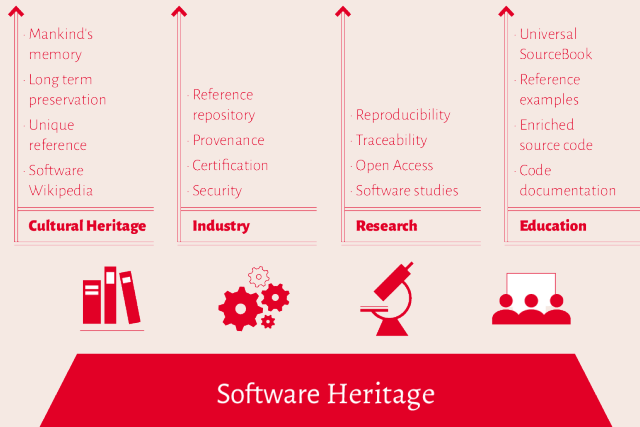
Collect, preserve and **share** the *source code* of *all the software* that is publicly available.

Past, present and future

Preserving the past, enhancing the present, preparing the future.

We are working on the foundations

One infrastructure to build them all





A structured archive of all of the world's software

- preserve humanity's technological and scientific **knowledge**
- enable continued **access** to all digital documents and information
- building block for **thematic portals** and collections



A unique reference catalog of all industrial software components

- a single entry point to discover, explore and reuse source code
- eases vulnerability tracking for more secure software
- simplifies **traceability** for better software integration
- ensures long term preservation of critical software



A global source referencing all software

- the ultimate **source book** for computer science and programming classes
- intrinsic persistent identifiers for stable **course materials**
- enables real-world, semi-automated **documentation**

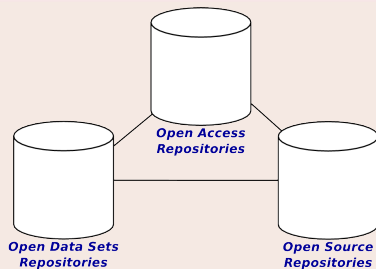


A global library referencing all software used in all research fields

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

The Knowledge Conservancy Magic Triangle

The Knowledge Conservancy Magic Triangle

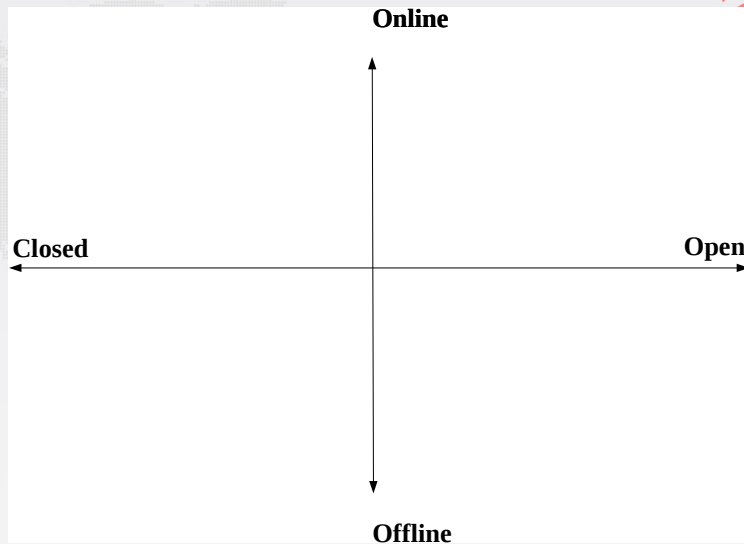


Legenda (links are important!)

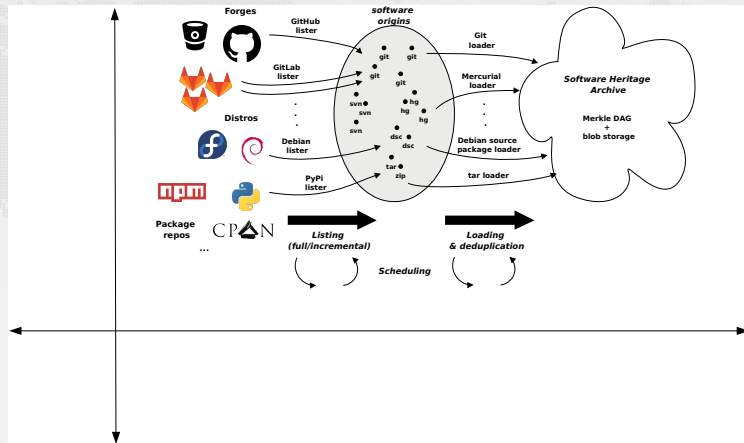
- articles: ArXiv, HAL, ...
- data: Zenodo, ...
- software: *Software Heritage* to the rescue

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood**
- 7 Building for the long term
- 8 Conclusion





Online, open source code: automation overview



Archive coverage



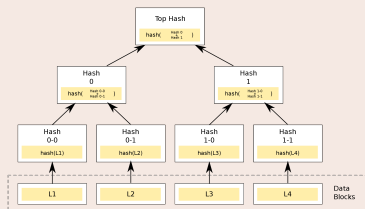
~150 TB blobs, ~5 TB database (as a graph: ~7 B nodes + ~60 B edges)

Our sources

- GitHub — full, up-to-date mirror
- Debian — automation in progress; GNU
- Gitorious, Google Code — processing (Archive Team & Google)
- Bitbucket — WIP

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

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
- widely used in industry (e.g., Git, nix, blockchains, IPFS, ...)

Features...

- (done) **lookup** by content hash
- **browsing**: "wayback machine" for archived code
 - (done) <http://archive.softwareheritage.org/api>
 - (in progress) via Web UI
- (in progress) **download**: `wget / git clone` from the archive
- (in progress) **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 than one could possibly imagine

all the world's software development history in a single graph!

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood
- 7 Building for the long term**
- 8 Conclusion



Cultural Heritage



Industry



Research



Education



Software Heritage

Open approach

- Transparency
- Free Software
- User and contributor community building

Objectiveness

- Facts and provenance
- *Intrinsic* identifiers
- Full development history

Long term

- Multi-stakeholder
- Nonprofit
- Replication *at all layers*

Three pillars

Science and technology

- build on sound basis
- fantastic playground for research

Resources

- fund the effort
- transfer to industry and society

Awareness

- promote public and private policies
- community building

Inria
INVENTEURS DU MONDE NUMÉRIQUE

>= 100Ke/year	 Microsoft	 intel	 SOCIETE GENERALE	
>= 50Ke/year				
>= 25Ke/year		 HUAWEI		
>= 10Ke/year	<small>Data Archiving and Networked Services</small> DANS	NOKIA Bell Labs	 <small>ALMA MATER STUDIORUM UNIVERSITA' DI BOLOGNA DIPARTIMENTO DI INFORMATICA - SCIENZA E INGEGNERIA</small>	GitHub

Sharing the Software Heritage vision



See more

<http://www.softwareheritage.org/support/testimonials>

April 3rd, 2017: landmark Inria Unesco agreement...

Inria
INVENTEURS DU MONDE NUMÉRIQUE



<https://www.softwareheritage.org/blog>

September 28th, 2017

Mauritius Call on information access

- 1 Introductions
- 2 Software is everywhere
- 3 Source code is essential...
- 4 ... but we are not taking care of it!
- 5 The Software Heritage initiative
- 6 Under the hood
- 7 Building for the long term
- 8 Conclusion



You can help!

Coding

- www.softwareheritage.org/community/developers/
- forge.softwareheritage.org – our own code

Current development priorities

- ★★★ listers for unsupported forges, distros, pkg. managers
- ★★★ loaders for unsupported VCS, source package formats
- ★ content indexing and search
- ★ efficient data representation

... *all* contributions equally welcome!

Join us

- www.softwareheritage.org/jobs – job openings
- wiki.softwareheritage.org – internships

Come in, we're open!

Learn more

A white paper is available
<http://bit.ly/swpaper>

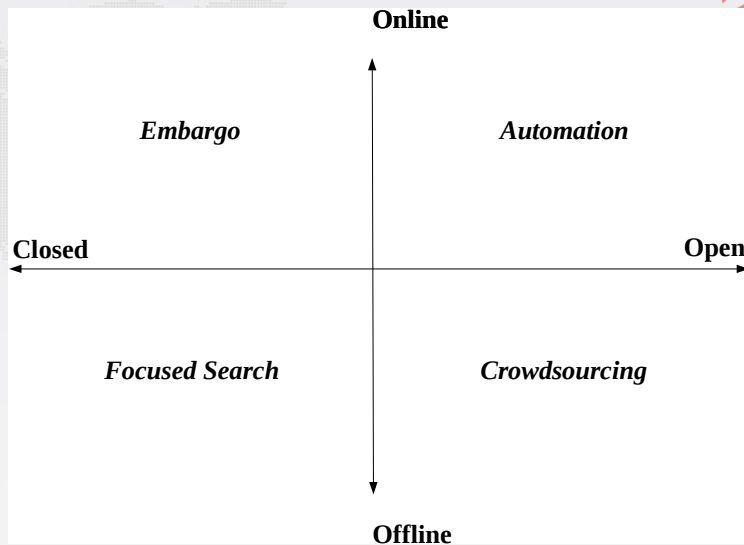
Presented at iPres 2017



Get involved

sponsoring / partnership
working groups, leads
our own code

sponsorship.softwareheritage.org
wiki.softwareheritage.org
forge.softwareheritage.org

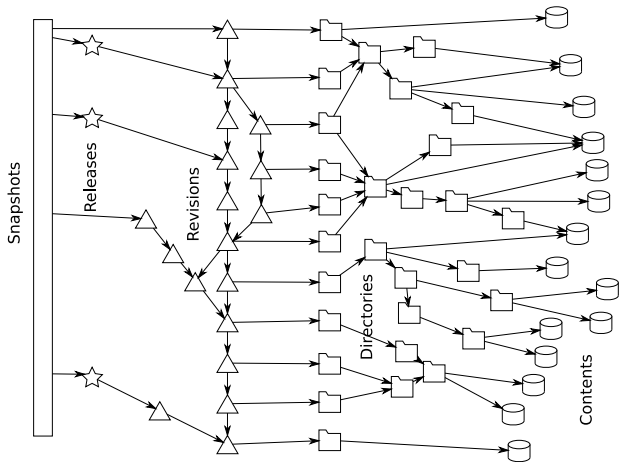




9 A glimpse of the archive

The archive in a few pictures

A giant (extended) Merkle DAG



First public version of our Web API (Feb 2017)

<https://archive.softwareheritage.org/api/>

Features

- pointwise **browsing** of the Software Heritage archive
 - ... releases → revisions → directories → contents ...
- full access to the **metadata** of archived objects
- **crawling** information
 - *when have you last visited this Git repository I care about?*
 - *where were its branches/tags pointing to at the time?*

Complete endpoint index

<https://archive.softwareheritage.org/api/1/>

A tour of the Web API — origins & visits

```
GET https://archive.softwareheritage.org/api/1/origin/ \
    git/url/https://github.com/hylang/hy
{ "id": 1,
  "origin_visits_url": "/api/1/origin/1/visits/",
  "type": "git",
  "url": "https://github.com/hylang/hy"
}
```

```
GET https://archive.softwareheritage.org/api/1/origin/ \
    1/visits/
[ ...,
  { "date": "2016-09-14T11:04:26.769266+00:00",
    "origin": 1,
    "origin_visit_url": "/api/1/origin/1/visit/13/",
    "status": "full",
    "visit": 13
  }, ...
]
```



A tour of the Web API — snapshots

```
GET https://archive.softwareheritage.org/api/1/origin/ \
  1/visit/13/
{ ...,
  "occurrences": { ...,
    "refs/heads/master": {
      "target": "b94211251...",
      "target_type": "revision",
      "target_url": "/api/1/revision/b94211251.../"
    },
    "refs/tags/0.10.0": {
      "target": "7045404f3...",
      "target_type": "release",
      "target_url": "/api/1/release/7045404f3.../"
    }, ...
  }, ...
},
"origin": 1,
"origin_url": "/api/1/origin/1/",
"status": "full",
"visit": 13
}
```



A tour of the Web API — revisions

```
GET https://archive.softwareheritage.org/api/1/revision/ \
    6072557b6c10cd9a21145781e26ad1f978ed14b9/
{
  "author": {
    "email": "tag@pault.ag",
    "fullname": "Paul Tagliamonte <tag@pault.ag>",
    "id": 96,
    "name": "Paul Tagliamonte"
  },
  "committer": { ... },
  "date": "2014-04-10T23:01:11-04:00",
  "committer_date": "2014-04-10T23:01:11-04:00",
  "directory": "2df4cd84e...",
  "directory_url": "/api/1/directory/2df4cd84e.../",
  "history_url": "/api/1/revision/6072557b6.../log/",
  "merge": false,
  "message": "0.10: The Oh f*ck it's PyCon release",
  "parents": [ {
    "id": "10149f66e...",
    "url": "/api/1/revision/10149f66e.../"
  }
]
```



```
GET https://archive.softwareheritage.org/api/1/content/ \
    adc83b19e793491b1c6ea0fd8b46cd9f32e592fc/
{
  "data_url": "/api/1/content/sha1:adc83b19e.../raw/",
  "filetype_url": "/api/1/content/sha1:.../filetype/",
  "language_url": "/api/1/content/sha1:.../language/",
  "length": 1,
  "license_url": "/api/1/content/sha1:.../license/",
  "sha1": "adc83b19e...",
  "sha1_git": "8b1378917...",
  "sha256": "01ba4719c...",
  "status": "visible"
}
```

Caveats

- rate limits apply throughout the API
- blob download available for selected contents