

Browsing the Free Software Commons

Stefano Zacchioli

University Paris Diderot & Inria – zack@upsilon.cc

22 August 2018

OpenSym 2018

Paris, France



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE



1 The Software Commons

2 Software Heritage

3 Accessing the archive

4 Getting involved

(Free) Software is everywhere



Software source code is *special*

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.”

Definition (Commons)

The **commons** is the cultural and natural resources accessible to all members of a society, including natural materials such as air, water, and a habitable earth. These resources are held in common, not owned privately. <https://en.wikipedia.org/wiki/Commons>

Definition (Software Commons)

The **software commons** consists of all computer software which is available at little or no cost and which can be altered and reused with few restrictions. Thus *all open source software and all free software are part of the [software] commons. [...]*

https://en.wikipedia.org/wiki/Software_Commons

Definition (Commons)

The **commons** is the cultural and natural resources accessible to all members of a society, including natural materials such as air, water, and a habitable earth. These resources are held in common, not owned privately. <https://en.wikipedia.org/wiki/Commons>

Definition (Software Commons)

The **software commons** consists of all computer software which is available at little or no cost and which can be altered and reused with few restrictions. Thus *all open source software and all free software are part of the [software] commons.* [...]

https://en.wikipedia.org/wiki/Software_Commons

Source code is a precious part of our commons

are we taking care of it?



Fashion victims

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

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, including: damage, disaster, malicious, deletion, reference, storage, attack, obsolete, dependencies, dangling, wear, corruption, encryption, format, aging, media, and tear.

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)



A word cloud of terms related to software fragility, including: damage, disaster, malicious, deletion, reference, storage, attack, obsolete, dependencies, dangling, wear, corruption, encryption, format, aging, media, and tear.

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

Software lacks its own research infrastructure



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?



1 The Software Commons

2 Software Heritage

3 Accessing the archive

4 Getting involved



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.

Cultural Heritage



Industry



Research



Education



Software Heritage

Cultural Heritage



Industry



Research



Education



Software Heritage

Open approach

- 100% Free Software
- transparency

In for the long haul

- replication
- non profit

Archiving goals

Targets: VCS repositories & source code releases (e.g., tarballs)

We DO archive

- file **content** (= blobs)
- **revisions** (= commits), with full metadata
- **releases** (= tags), ditto
- where (**origin**) & when (**visit**) we found any of the above

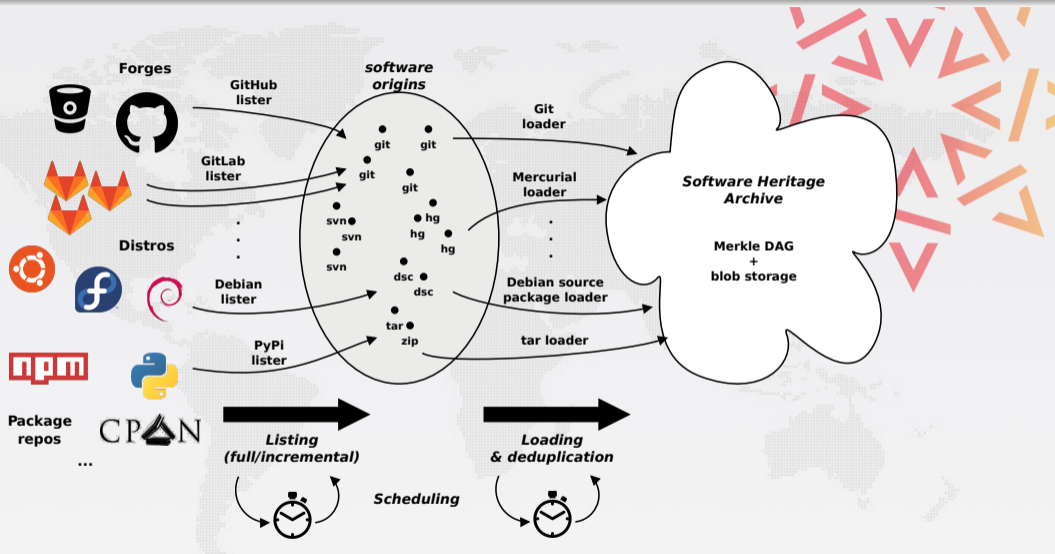
... in a VCS-/archive-agnostic **canonical data model**

We DON'T archive

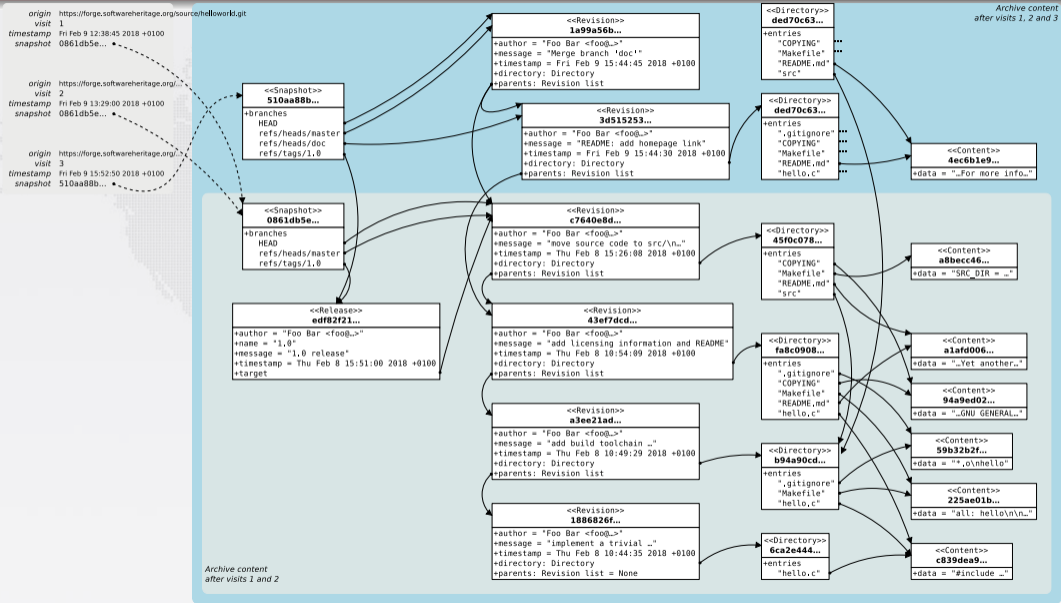
- homepages, wikis
- BTS/issues/code reviews/etc.
- mailing lists

Long term vision: play our part in a *"semantic wikipedia of software"*

Data flow



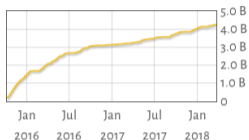
The archive: a (giant) Merkle DAG



Archive coverage

Source files

4,290,063,587



Commits

980,310,191



Projects

83,797,945



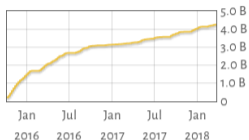
Current sources

- live: GitHub, Debian
- one-off: Gitorious, Google Code, GNU
- WIP: GitLab, PyPI, Bitbucket

Archive coverage

Source files

4,290,063,587



Commits

980,310,191



Projects

83,797,945



Current sources

- live: GitHub, Debian
- one-off: Gitorious, Google Code, GNU
- WIP: GitLab, PyPI, Bitbucket

175 TB (compressed) blobs, 6 TB database (as a graph: 10 B nodes + 100 B edges)



Current sources

- live: GitHub, Debian
- one-off: Gitorious, Google Code, GNU
- WIP: GitLab, PyPI, Bitbucket

175 TB (compressed) blobs, 6 TB database (as a graph: 10 B nodes + 100 B edges)

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

- 
- 1 The Software Commons
 - 2 Software Heritage
 - 3 Accessing the archive
 - 4 Getting involved

RESTful API to programmatically access the Software Heritage archive

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

Features

- pointwise **browsing** of the archive
 - ... snapshots → 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?*

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.../"
  }
]
```



A tour of the Web API — contents

```
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"
}
```



```
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

Vault service

- source code is thoroughly deduplicated within the Software Heritage archive
- bulk download of large artefacts (e.g., a Linux kernel release) requires collecting millions of objects
- the **Software Heritage Vault** cooks and caches source code bundles for bulk download needs

Tech bits


- **RESTful API** to request downloads, notifications, and monitoring
- `docs.softwareheritage.org/devel/swh-vault`

Browser-based interface to browse the Software Heritage archive

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

Features

- all **REST API features**, but good looking :-)
 - browsing: snapshots → revisions → directories → contents ...
 - access to metadata and crawling information
- **origin search**, as full text indexing of origin URLs
- bulk **download**, via integration with the Vault

- 
- 1 The Software Commons
 - 2 Software Heritage
 - 3 Accessing the archive
 - 4 Getting involved

Features...

- (done) **lookup** by content hash
- (done) **browsing**: "wayback machine" for source code (API + UI)
- (early access) **deposit** of source code bundles directly to the archive
- (done) **download**: `wget / git clone` from the archive
- (todo) **provenance** lookup for all archived content
- (todo) **full-text search** on all archived source code files

Features...

- (done) **lookup** by content hash
- (done) **browsing**: "wayback machine" for source code (API + UI)
- (early access) **deposit** of source code bundles directly to the archive
- (done) **download**: `wget / git clone` from 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 at hand's reach!

You can help!

Coding

- ★★ Web UI improvements
- ★ loaders/listers for unsupported VCS/forges
- ★★★ developer documentation

<https://docs.softwareheritage.org/devel/>

You can help!

Coding

- ★★ Web UI improvements
- ★ loaders/listers for unsupported VCS/forges
- ★★★★ developer documentation

<https://docs.softwareheritage.org/devel/>

Community

- ★★★★ spread the world, help us with sustainability
- ★★ document endangered source code

wiki.softwareheritage.org/index.php?title=Suggestion_box

You can help!

Coding

- ★★ Web UI improvements
- ★ loaders/listers for unsupported VCS/forges
- ★★★★ developer documentation

<https://docs.softwareheritage.org/devel/>

Community

- ★★★★ spread the world, help us with sustainability
- ★★ document endangered source code

wiki.softwareheritage.org/index.php?title=Suggestion_box

Join us

- www.softwareheritage.org/jobs – **job openings**
- wiki.softwareheritage.org/index.php?title=Internship – **internships**

Software Heritage is

- a reference archive of **all Free Software** ever written
- an international, open, nonprofit, **mutualized infrastructure**
- **now accessible** to developers, users, researchers
- at the service of our community, **at the service of society**

Come in, we're open!

`www.softwareheritage.org` – general information

`wiki.softwareheritage.org` – internships, leads

`forge.softwareheritage.org` – our own code