# Software Heritage, the universal archive at the crossroads of Software Forges

Benoit Chauvet

Software Heritage
Software Engineering Manager

29 June 2023
GT données GitLab 2023
Paris, France

# Software Heritage

## THE GREAT LIBRARY OF SOURCE CODE

# Outline

- Software Heritage
- Started in 2015 (Roberto Di Cosmo & Stefano Zacchiroli)
- Universal archive of source code
- Initiated by Inria, in collaboration with Unesco

# Outline

# A few basic needs for research software

## Archive

Research software artifacts must be properly *archived*

make sure we can *retrieve* them (*reproducibility*)

# A few basic needs for research software

## Archive

Research software artifacts must be properly archived

make sure we can *retrieve* them (*reproducibility*)

## Reference

Research software artifacts must be properly referenced

make sure we can *identify* them (*reproducibility*)

## Archive

Research software artifacts must be properly *archived*

make sure we can *retrieve* them (*reproducibility*)

## Reference

Research software artifacts must be properly *referenced*

make sure we can *identify* them (*reproducibility*)

## Describe

Research software artifacts must be properly *described*

make it easy to *discover* and *reuse* them (*visibility*)

# A few basic needs for research software

## Archive

Research software artifacts must be properly archived

make sure we can *retrieve* them (*reproducibility*)

## Reference

Research software artifacts must be properly referenced

make sure we can *identify* them (*reproducibility*)

## Describe

Research software artifacts must be properly described

make it easy to *discover* and *reuse* them (*visibility*)

## Cite/Credit

Research software artifacts must be properly cited *(not the same as referenced!)*

to give *credit* to authors (*evaluation!*)

# Outline

- many disparate development platforms, with a few dominant players (e.g., GitHub)
- a myriad places where distribution may happen
- most of them operated by for-profit companies

# Software source code is fragile

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

## Like all digital information, FOSS is fragile

- link rot: projects are created, moved around, removed
- business-driven code loss (e.g., Gitorious, Google Code, Bitbucket)
- data rot: physical media with legacy software decay

# Software source code is fragile

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

## Like all digital information, FOSS is fragile

- link rot: projects are created, moved around, removed
- business-driven code loss (e.g., Gitorious, Google Code, Bitbucket)
- data rot: physical media with legacy software decay

## If a website disappears you go to the Internet Archive…

where do you go if (a repository on) GitHub or GitLab goes away?

# Outline

## 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

Cultural Heritage | Industry | Research | Public Administration

Software Heritage

| | Source files | Commits | Projects |
|---|---|---|---|
| | 13,974,813,954 | 2,912,845,019 | 207,160,527 |

archive.softwareheritage.org

Cultural Heritage | Industry | Research | Public Administration

**Software Heritage**

| Source files | Commits | Projects |
| --- | --- | --- |
| 13,974,813,954 | 2,912,845,019 | 207,160,527 |

archive.softwareheritage.org
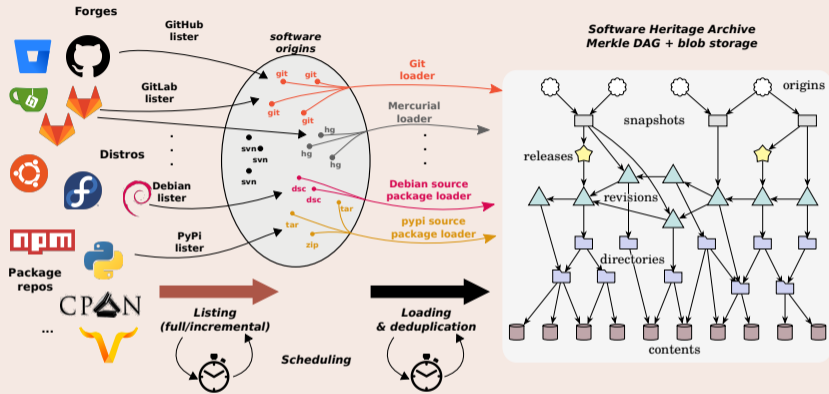
## Technology
- transparency and FOSS
- replicas all the way down

## Content (billions!)
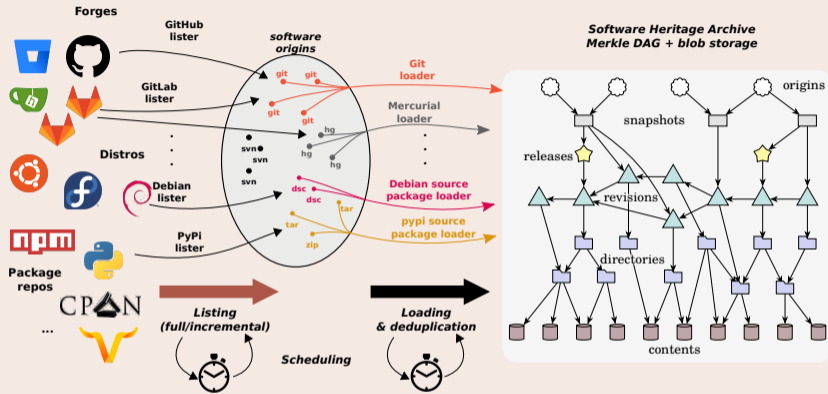- intrinsic identifiers
- facts and provenance

## Organization
- non-profit
- multi-stakeholder

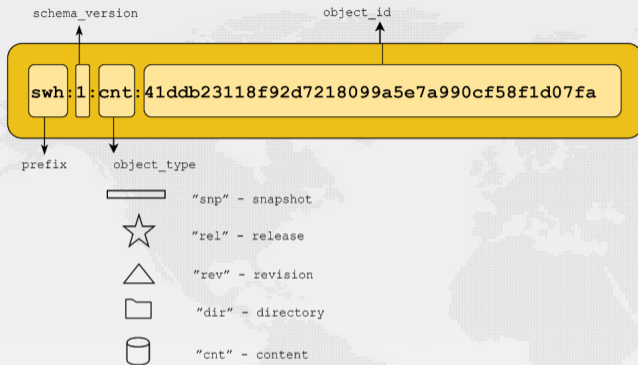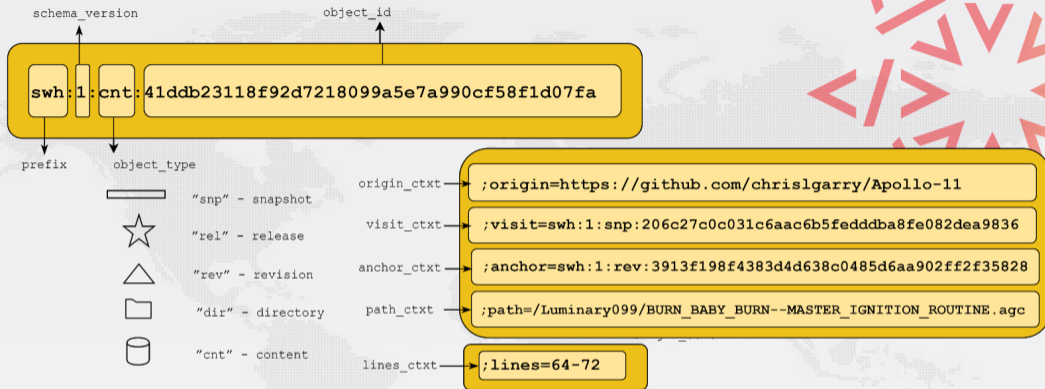A global graph linking together fully deduplicated source code artifact (files, commits, directories, releases, etc.) to the places that distribute them (e.g., Git repositories), providing a unified view on the entire *Software Commons*. (Size: ~30 B nodes, ~300 B edges, ~1 PiB blobs)

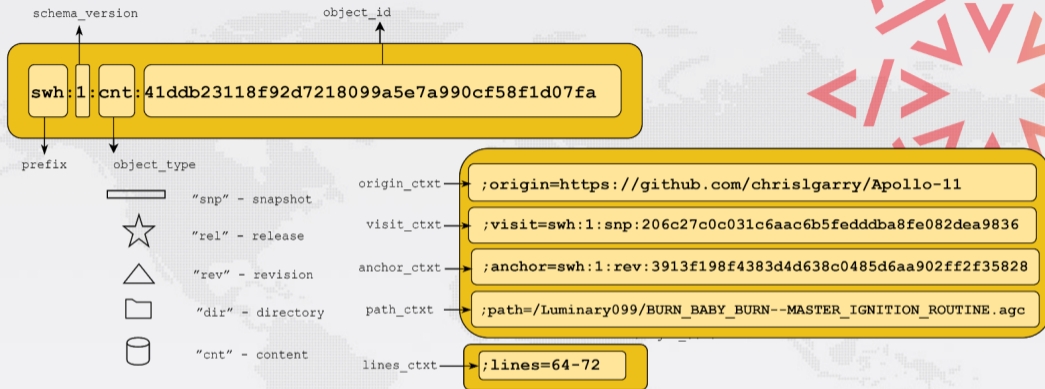# Outline

schema_version                    object_id

```
swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa
```

prefix          object_type

schema_version · object_id

```
swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa
```

prefix · object_type

▭    "snp" – snapshot

☆    "rel" – release

△    "rev" – revision

▢    "dir" – directory

⬭    "cnt" – content

schema_version

object_id

```
swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa
```

prefix

object_type

⬜ "snp" – snapshot

☆ "rel" – release

△ "rev" – revision

▭ "dir" – directory

⬭ "cnt" – content

origin_ctxt → `;origin=https://github.com/chrislgarry/Apollo-11`

visit_ctxt → `;visit=swh:1:snp:206c27c0c031c6aac6b5fedddba8fe082dea9836`

anchor_ctxt → `;anchor=swh:1:rev:3913f198f4383d4d638c0485d6aa902ff2f35828`

path_ctxt → `;path=/Luminary099/BURN_BABY_BURN--MASTER_IGNITION_ROUTINE.agc`

lines_ctxt → `;lines=64-72`

schema_version

object_id

```
swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa
```

prefix

object_type

⎯ "snp" – snapshot

☆ "rel" – release

△ "rev" – revision

▢ "dir" – directory

⬗ "cnt" – content

origin_ctxt → `;origin=https://github.com/chrislgarry/Apollo-11`

visit_ctxt → `;visit=swh:1:snp:206c27c0c031c6aac6b5fedddba8fe082dea9836`

anchor_ctxt → `;anchor=swh:1:rev:3913f198f4383d4d638c0485d6aa902ff2f35828`

path_ctxt → `;path=/Luminary099/BURN_BABY_BURN--MASTER_IGNITION_ROUTINE.agc`

lines_ctxt → `;lines=64-72`

## An emerging standard

- in Linux Foundation's SPDX 2.2
- IANA registered, WikiData property P6138

schema_version

object_id

```
swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa
```

prefix

object_type

```
▭    "snp" - snapshot
☆    "rel" - release
△    "rev" - revision
▢    "dir" - directory
⬡    "cnt" - content
```

origin_ctxt → `;origin=https://github.com/chrislgarry/Apollo-11`

visit_ctxt → `;visit=swh:1:snp:206c27c0c031c6aac6b5fedddba8fe082dea9836`

anchor_ctxt → `;anchor=swh:1:rev:3913f198f4383d4d638c0485d6aa902ff2f35828`

path_ctxt → `;path=/Luminary099/BURN_BABY_BURN--MASTER_IGNITION_ROUTINE.agc`

lines_ctxt → `;lines=64-72`

## An emerging standard

- in Linux Foundation's SPDX 2.2
- IANA registered, WikiData property P6138

## Examples:

- Apollo 11 AGC excerpt
- Quake III rsqrt

# Outline

## Archive a single repository

## Automate archival from your forge

- Automatically trigger archival in Software Heritage
- Triggering events:
  - Tags or branch creation
  - Release creation
  - …
- Available for:
  - Bitbucket
  - Gitea
  - GitHub
  - GitLab
  - Sourceforge
- Webhooks howto available here…

## Archive a whole forge



**Software Heritage**

≡ **Request the addition of a forge into the archive**

Enter a SWHID to resolve or keyword(s) to search for in origin

Features
🔍 Search
⬇ Downloads
📷 Save code now
🏹 Add forge now
❓ Help

Administration
📷 Save code now
🏹 Add forge now
📁 Deposit

"Add forge now" provides a service for Software Heritage users to save a complete forge in the Software Heritage archive by requesting the addition of the forge URL into the list of regularly visited forges.

| Submit a Request | Browse Requests | Help |

**Forge type** *

bitbucket

Supported forge types in software archive.

**Forge URL** *

Remote URL of the forge.

**Forge contact name** *

Name of the forge administrator.

**Forge contact email** *

Email of the forge administrator. The given email address will not be used for any purpose outside the "add forge now" process.

☐ I consent to add my username in the communication with the forge.

**Comment**

Optionally, leave a comment to the moderator regarding your request.

Submit Add Request

Once an add-forge-request is submitted, its status can be viewed in the **submitted requests list**. This process involves a moderator approval and might take a few days to handle (it primarily depends on the response time from the forge).

## How to deposit

- Local method: deposit .zip/.tar.gz file
- SWHID method: deposit SWHID + metadata

## Embedded metadata

- README file (Mardown or plain text)
- AUTHORS files (plain text)
- LICENCE file (plain text) or LICENCES directory
- codemeta.json file CodeMeta generator

# CodeMeta generator

Most fields are optional. Mandatory fields will be highlighted when generating Codemeta.

## The software itself

**Name**
My Software

the software title

**Description**
My Software computes ephemerides and orbit propagation. It has been developed from early '80.

**Creation date**
YYYY-MM-DD

**First release date**
YYYY-MM-DD

**License(s)**

from SPDX licence list

## Discoverability and citation

**Unique identifier**
10.151.xxxxx

such as ISBNs, GTIN codes, UUIDs etc.. http://schema.org/identifier

**Application category**
Astronomy

**Keywords**
ephemerides, orbit, astronomy

**Funding**
PRA_2018_73

grant funding software development

**Funder**
Università di Pisa

organization funding software development

Authors and contributors can be added below

## Development community / tools

**Code repository**
git+https://github.com/You/RepoName.git

**Continuous integration**
https://travis-ci.org/You/RepoName

**Issue tracker**
https://github.com/You/RepoName/issues

**Related links**

## Run-time environment

**Programming Language**
C#, Java, Python 3

**Runtime Platform**
.NET, JVM

**Operating System**
Android 1.6, Linux, Windows, macOS

**Other software requirements**
Python 3.4
https://github.com/psf/requests

## Current version of the software

**Version number**
1.0.0

**Release date**
YYYY-MM-DD

**Download URL**
https://example.org/MySoftware.tar.gz

**Release notes**
Change log: this and that;
Bugfixes: that and this.

## Additional Info

**Reference Publication**
https://doi.org/10.1000/xyz123

**Development Status**

see www.repostatus.org for details

**Is part of**
http://The.Bigger.Framework.org

## Authors
Add one   Remove last

## Contributors

## SWHID Deposit in HAL

- https://www.softwareheritage.org/2023/04/04/
  swhid-deposit-hal/

## Search and expose software publications

- https://haltools.archives-ouvertes.fr

# Outline

# Questions ?

# Outline

# Appendix

## Software Heritage

- Our mission
- Browse the archive
- Save and reference research software
- HOWTO archive and reference your code
- Save Code Now
- Add Forge Now
- Webhooks for auto-archival

## Misc

- CodeMeta generator
- SWHID deposit in HAL
- HAL Tools