# Software Heritage: why and how we preserve our Software Commons
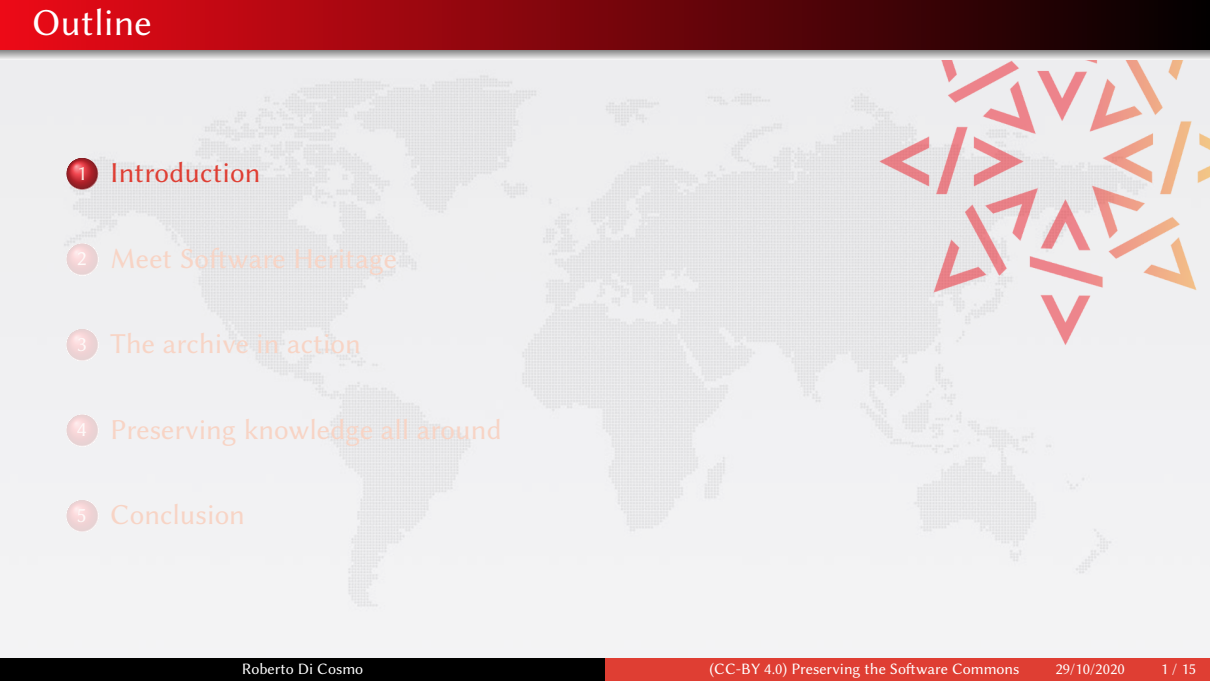
Roberto Di Cosmo
Director, Software Heritage
Inria and Université de Paris

29/10/2020, *49 JAIIO*

Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

# Outline

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

2015   *Software Heritage* at INRIA

2018   *National Committee for Open Science*, France

# Software source code: a precious part of our heritage

## 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."*

## Apollo 11 source code (excerpt)

```
P63SPOT3     CA      BIT6        # IS THE LR ANTENNA IN POSITION 1 YET
             EXTEND
             RAND    CHAN33
             EXTEND
             BZF     P63SPOT4    # BRANCH IF ANTENNA ALREADY IN POSITION 1

             CAF     CODE500     # ASTRONAUT:   PLEASE CRANK THE
             TC      BANKCALL    #              SILLY THING AROUND
             CADR    GOPERF1
             TCF     GOTOPOOH    # TERMINATE
             TCF     P63SPOT3    # PROCEED     SEE IF HE'S LYING

P63SPOT4     TC      BANKCALL    # ENTER       INITIALIZE LANDING RADAR
             CADR    SETPOS1

             TC      POSTJUMP    # OFF TO SEE THE WIZARD ...
             CADR    BURNBABY
```

## Quake III source code ( excerpt )

```c
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;
}
```

## Len Shustek, Computer History Museum

*"Source code provides a view into the mind of the designer."*

## Yuval Noah Harari (on COVID 19)

*"The real antidote [to epidemic] is* scientific knowledge *and* global cooperation."
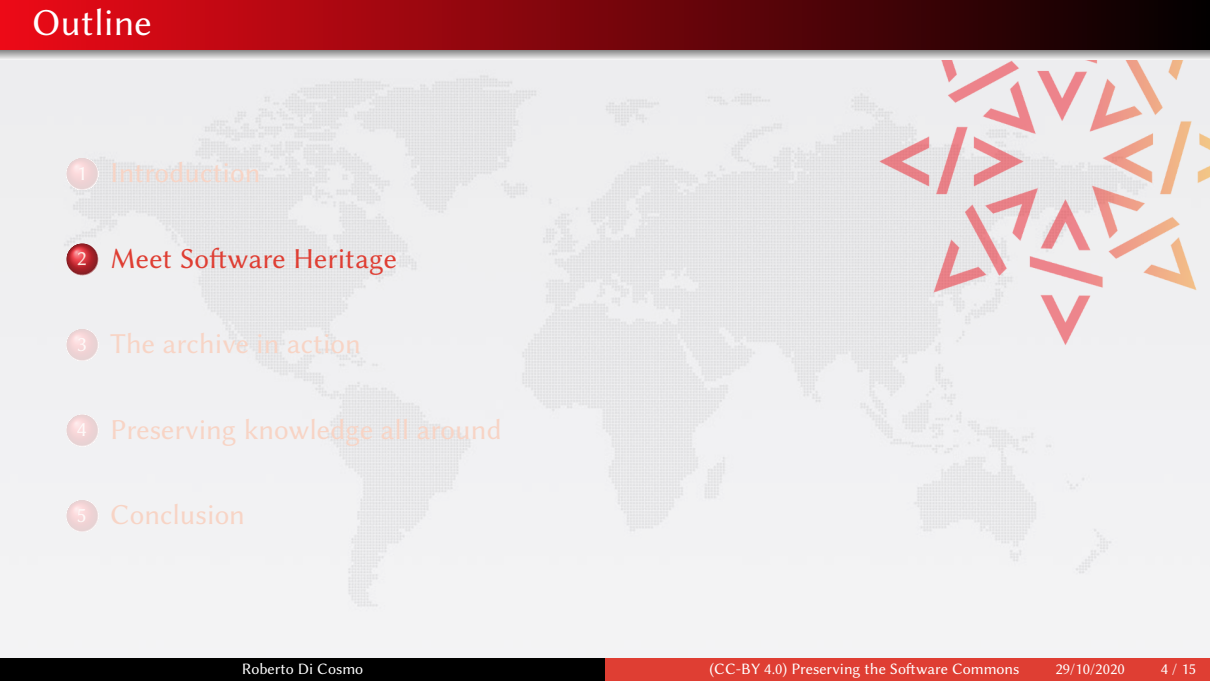
## Definition (Software Commons)

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: part of our commons ... pillar of Open Science! *(would require another talk)*

## Precious, endangered *executable* and *human readable* knowledge

key people passing away, platforms (GoogleCode, Gitorious, etc.) closing down ...

# 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

### Reference catalog



**find** and **reference** all software source code

### Universal archive



**preserve** all software source code

### Research infrastructure



**enable analysis** of all software source code

Cultural Heritage | Industry | Research | Education

## Software Heritage

As of today the archive already contains and keeps safe for you the following amount of objects:

| Source files | Commits | Projects |
|---|---|---|
| 8,846,381,610 | 1,880,663,008 | 140,348,311 |

| Directories | Authors | Releases |
|---|---|---|
| 7,506,954,410 | 38,603,337 | 15,051,940 |

## Technology
- transparency and FOSS
- replicas all the way down

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

## Organization
- non-profit
- multi-stakeholder

## Sharing the vision



And many more …
www.softwareheritage.org/support/testimonials

## Donors, members, sponsors



**Platinum sponsors**

**Gold sponsors**

**Silver sponsors**

**Bronze sponsors**

# A peek under the hood



*Global development history* permanently archived in a *unique* git-like Merkle DAG

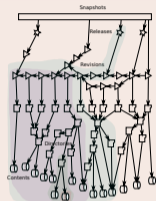- ~400 TB (uncompressed) blobs, ~20 B nodes, ~280 B edges

```
                    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 "swh:" URI prefix
- WikiData property P6138

## Examples

- Apollo 11 AGC excerpt
- Quake III rsqrt

# A revolutionary infrastructure for software source code

## The *graph* of Software Development



All software development with its history, in a single graph …

## The *blockchain* of Software Development



… a single Merkle graph, with *intrinsic ids* for traceability

## A *pillar* of Open Science



Reference archive of Research Software

## Reference platform for *Big Code*



One uniform data structure enables *massive* machine learning for quality, cybersecurity, etc.

# Outline

# An example is worth a thousand words

- Browse the archive
- Trigger archival of your preferred software in a breeze
- Get and use SWHIDs (full specification available online)
- curated deposit in SWH via HAL, see for example: LinBox, SLALOM, Givaro, NS2DDV, SumGra, Coq proof, …
- cite software with the biblatex-software style from CTAN
- Example use in a research article: compare Fig. 1 and conclusions
  - in the 2012 version
  - in the updated version using SWHIDs and Software Heritage
- Example use in a research article: extensive use of SWHIDs in a replication experiment

# Outline

Experts call for greater recognition of software source code as heritage for sustainable development

.6 November 2018





At UNESCO, Inria, and Software Heritage invitation,
40 international experts meet in Paris …

Their call is published on Feb 2019
it's open for signatures!

## Paris Call on Software Source Code

"[We call to] support efforts to gather and preserve the artifacts and narratives of the history of computing, while the earlier creators are still alive"

Rescue Legacy Software, Curate the code, Archive it



The Software Heritage Acquisition Process

UNESCO, UniPi and Software Heritage collaboration                    worldwide scope

- The control code for the music synthesizer TAU2
- FORTRAN II and TAUmus command language
- Istituto di Elaborazione dell'Informazione CNR
- Group led by the late M° P. Grossi
  - Le Sacre du Printemps (ABSTRACT)

# ... and massive amounts of present source code!

## Paris Call on Software Source Code

"[We call to] support a universal archive, as part of a broad effort at digital preservation, that will ensure persistence of and universal access to software source code"

## Summer 2019: BitBucket announce Mercurial VCS phase out!

- fall 2019: Software Heritage teams up with Octobus (funded by NLNet, thanks!)
- july 2020: BitBucket erases *250.000* repositories
- august 2020: bitbucket-archive.softwareheritage.org is live

## Preserving the web of knowledge                    (Tweet is here )

**Gabriel Altay**
@gabrielaltay

Just realized @Bitbucket disabled all mercurial repositories when the @asclnet informed me that a link associated with an old paper of mine was down. Thought all was lost, but someone archived all the repos! very classy move by @octobus_net and @SWHeritage.

Traduire le Tweet

Bottomline
*explicit deposit* is important, ...
                  ... and we must promote it...
                           ... but will never be enough.

# Outline

# Software Heritage

www.softwareheritage.org          @swheritage

## Everybody is concerned, everybody can help build

### The 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

# Appendix

**Online**

**Escrow**

**Automation**

**Closed**

**Open**

**Focused Search**

**Crowdsourcing**

**Offline**

# A word on the trust model for systems of identifiers

## Two general classes of systems of identifiers

**intrinsic** *computed* from the object *(no registry required, fully decentralised)*
*(e.g.: chemical notation, music notation, hashes, SWHIDs)*

**extrinsic** *assigned* by an authority *(need a registry)*
*(e.g.: passport number, DOI, ARK, RRID, etc.)*

See the dedicated blog post for more details

### Trust model, extrinsic (e.g. DOIs)



### Trust model, intrinsic (e.g. SWHIDs)

# A worked example

# Contents



```
              GNU GENERAL PUBLIC LICENSE
               Version 3, 29 June 2007

Copyright (C) 2007 Free Software Foundation, Inc. <http://fsf.org/>
Everyone is permitted to copy and distribute verbatim copies
of this license document, but changing it is not allowed.

                     Preamble

  The GNU General Public License is a free, copyleft license for
software and other kinds of works.

  The licenses for most software and other practical works are designed
to take away your freedom to share and change the works.  By contrast,
the GNU General Public License is intended to guarantee your freedom to
share and change all versions of a program--to make sure it remains free
software for all its users.  We, the Free Software Foundation, use the
GNU General Public License for most of our software; it applies also to
any other work released this way by its authors.  You can apply it to
your programs, too.

  When we speak of free software, we are referring to freedom, not
price.  Our General Public Licenses are designed to make sure that you
have the freedom to distribute copies of free software (and charge for
them if you wish), that you receive source code or can get it if you
want it, that you can change the software or use pieces of it in new
free programs, and that you know you can do these things.

  To protect your rights, we need to pre
```

sha1: 8624bcdae55baeef...
sha256: 8ceb4b9ee5aded...
sha1_git: 94a9ed024d385...
length: 35147

# A worked example

# Revisions

# Releases

tag v0.0.51
Tagger: Nicolas Dandrimont <nicolas@dandrimont.eu>
Date:   Wed Aug 24 14:36:03 2016 +0200

Release swh.storage v0.0.51

 - Add new metadata column to origin_visit
 - Update swh-add-directory script for updated API
[...]

commit c0c9f16b1e134f593e7567570a1761b156e6eb1d

object c0c9f16b1e134f593e7567570a1761b156e6eb1d
type commit
tag v0.0.51
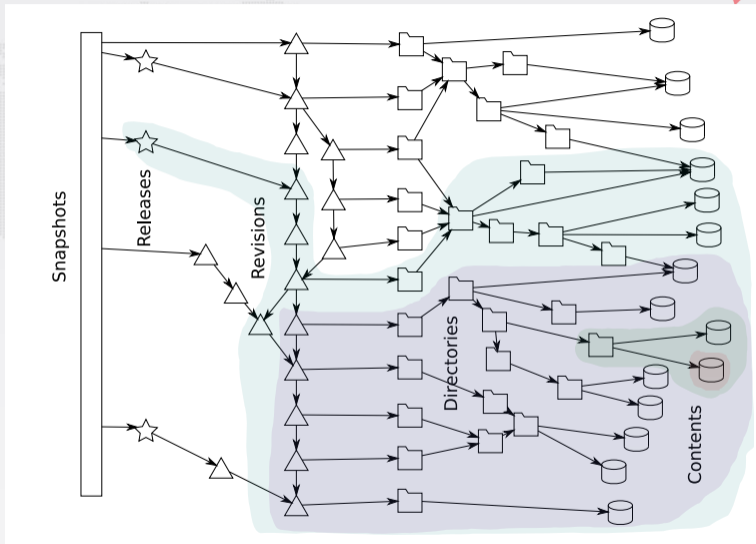tagger Nicolas Dandrimont <nicolas@dandrimont.eu> 1472042163 +0200

Release swh.storage v0.0.51

 - Add new metadata column to origin_visit
 - Update swh-add-directory script for updated API
-----BEGIN PGP SIGNATURE-----

iQIzBAABCAAdBQJXvZTNFhxuaWNvbGFzQGRhbmRyaWw1vbnQuZXUACgkQ7AWLMo2+
neqorw//aq6SOb5DijzEa+kWN3rXgVS+1K1vEVh1wNKAwx8eKJ7aX2kEiLDtt7uf
ahpZ6pz3q8nqs6aC1+YrxBfcih3L2YtrdZeWXWqr8xWNMaEoYDb8qaphwh8AD5t2
ICBIit2uJtXuCrDt93eKKPwvzZXg+hB0sMWy35Dr6jW7Z7K4Mu/PGglyIHPYS5yo
IGEndWno7VfH1Vm6t1n5qB7I5mXRaqA+becqddubTZ2xjj+jplUqC8cyqN3hm/tL
qsj2mu8kyz3t8tG/H1/pV+I5OwBlnPoS5TH0tujojEVgPK/dHSP79QuHDHZFkCao
kIj6kAWyU80Mxb+nKV/jeLbrR3+yWBFj3Qp5a1/V8oOTh6E1dALcNMpEaKCoKtMt
d/gMRax1l1g0EDfnsW67G6sDwKPKPHhgfVLQ3nV3GaQQTnu1RpMz006H9/tAwzC
Gg/K1PdHT4hzOi46wYPZyje0U2VXGFu6vVU9vFQ4ZR/Wjn+0zMzdcRdrlJSUOMn
RpTTfUsbXUeXHGOpkgXhSYTnvp1gdPc76US1sK0aGe84AZm1lk0mGrwXCVfPqlYo
nhhibBSHBNMoqyF6yTSOpUbYK70tpYRRUGKWDeRK0wKSxkWKUZGtKzy6JYqJjo29
gulwqZQif5qWQCB0OontAL2+HvPFaVyckMejUhg62cP/+EHIvUk=
=kOxP
-----END PGP SIGNATURE-----

id: 85083a5cc14a441c89dea73f5bdf67c3f9c6afdb

## Snapshots

```
commit 08ffeb25770109525eb3ce21691466c53a1d9158 refs/heads/atime
commit ba5443a24e3f9fe323a46c292cec4fcbe61c67eb refs/heads/directory-listing-arrays
commit d69e0dbf892383ff6589b27fbe1c05d27238d9c5 refs/heads/foo
commit cf7ff9eea0eb22f0946908f5a8019f67de468e08 refs/heads/master
commit 7eca197fc66d2024047e54b1ed9e8b44361a0fc2 refs/heads/tmp-directory-add
commit 642a205f37de85005a85d427b53ee4fb2252e82e refs/heads/tmp/generic-releases
tag 20f043b1379cf768d966597799fd4907c757f755 refs/tags/v0.0.1
tag 72a21991a384e539996dbb867bfb0bee72aee2cd refs/tags/v0.0.10
tag 3590e0ca0ebb070e5b376705fa230bbfa4ffa5cc refs/tags/v0.0.11
tag 33378427a403ba569a67777b8d58f6674fbc6556 refs/tags/v0.0.12
tag 06f74652755b327cf590311c2bfa036cf3b4b35d refs/tags/v0.0.13
tag 5a6325fe86ab854b581d7442667d92a11e32f3bd refs/tags/v0.0.14
tag 586fba4e580b4f5fab05f599367643cbcb1a9c7f refs/tags/v0.0.15
tag 8cd8b885f4098bf363177742bd289f660e5be51c refs/tags/v0.0.16
tag a542444ee3f0fbed35efb202fee035c809abc7d6 refs/tags/v0.0.17
tag 228a2f1650dd12222e556559462e1e06fc4993d9 refs/tags/v0.0.18
tag 606979a4ca05d497fc0d24aad00dce82636ef47c refs/tags/v0.0.19
tag 32bf5a59fc2a323baa6d5f15a6ad5382ec275a67 refs/tags/v0.0.2
tag 3147c3d31ec46cf6492f881e908b1237ebdff2c7 refs/tags/v0.0.20
tag 215ea50daba111e082e0b72e76eb4b6073a87908 refs/tags/v0.0.21
tag 3fb168c2072a5d62521242257a1e5dfc0f5ffa1df refs/tags/v0.0.22
tag 8cdbee8da4d73fc5d262789e460a16ac3c72aba4 refs/tags/v0.0.23
...
```
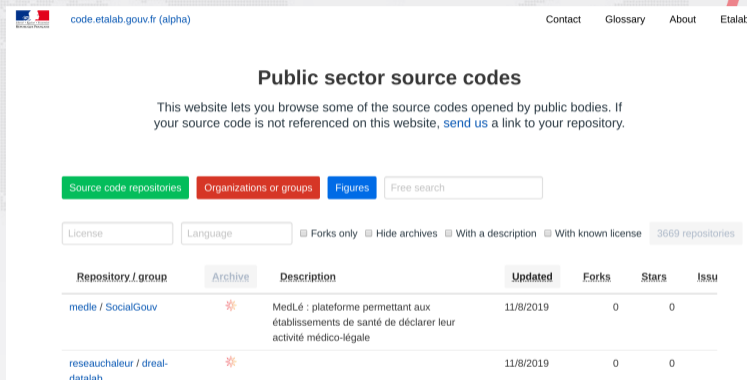
git show-refs

id: b464cad1b66fff266a37b46ea6e7a04b545e904b

# Outline

https://code.etalab.gouv.fr