

Reflections on the Future of Software Development

Roberto Di Cosmo
Director, Software Heritage

Apr 7th, 2020



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

Software source code: *human readable and executable knowledge*

Harold Abelson, Structure and Interpretation of Computer Programs

(1985)

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

Apollo 11 source code (excerpt)

```
P63SP0T3      CA      BIT6          # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND     CHAN33
              EXTEND
              BZF      P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF      CODE500      # ASTRONAUT: PLEASE CRANK THE
              TC       BANKCALL     # SILLY THING AROUND
              CADR     GOPERF1
              TCF      GOTOP00H     # TERMINATE
              TCF      P63SP0T3     # PROCEED SEE IF HE'S LYING

P63SP0T4      TC       BANKCALL     # ENTER INITIALIZE LANDING RADAR
              CADR     SETPOS1

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

Quake III 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;
}
```

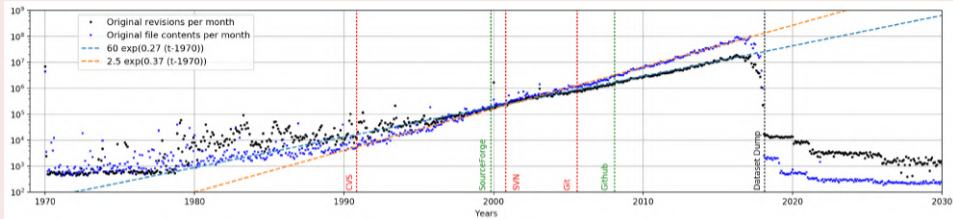
Len Shustek, Computer History Museum

(2006)

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

Exponential growth and increasing complexity

Growth of *globally original known content* (source: Rousseau, Di Cosmo, Zacchioli 2019)

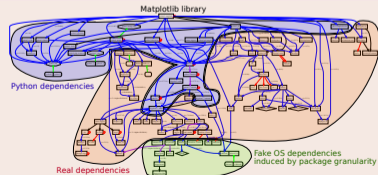


trend over 20 years: *original content doubles every 22 months* (i.e. **~45% per year**)

Complexity of the software ecosystem

8000+ langs

- *millions* of lines of code
- large *web of dependencies*
 - easy to break, difficult to maintain
- decentralised *developer communities*



Precious asset, *endangered heritage*

We are at a turning point

software **now a critical asset for society**, but key scientists/developers **are passing away**, and source code is getting **lost or misplaced** while software development **skyrockets!**

To enable next generation research and software development we need
a common, non profit, long term, shared infrastructure that provides

A catalog



A word cloud featuring various software development projects and platforms. The most prominent words are 'Git Hub', 'Sourceforge', 'Debian', 'CRAN', 'Maven', 'Bitbucket', 'GoogleCode', 'GitLab', 'CRAN', 'eTAN', 'CRAN', 'Adalact', 'BoltOz', 'Ivoria', 'Cibonius', and 'Cibonius'.

An archive



A word cloud of terms related to data loss and preservation. The most prominent words are 'damage', 'disaster', 'malicious', 'deletion', 'attack', 'obsolete', 'dependencies', 'format', 'encryption', 'corruption', 'wear', 'dangling', 'reference', 'storage', 'tear', 'aging', 'media', and 'attack'.

A research infrastructure





Software Heritage

Our mission

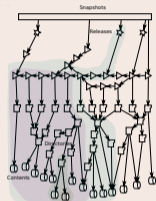
Collect, preserve and share the *source code* of *all the software* that is available



Preserving the past, enhancing the present, preparing the future

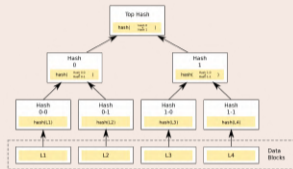
A revolutionary infrastructure for software source code

The *graph* of Software Development



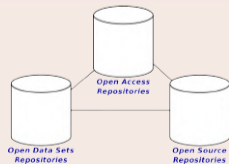
All software development with its history, in a **single graph** ...

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.

Sharing the vision



United Nations
Educational, Scientific and
Cultural Organization



And many more ...

www.softwareheritage.org/support/testimonials

Donors, members, sponsors



Platinum sponsors



Gold sponsor



Silver sponsors



Bronze sponsors



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

6 November 2018



UNESCO, Inria, Software Heritage invite
40 international experts meet in Paris ...

It's an important *policy tool*, already referenced and used ...

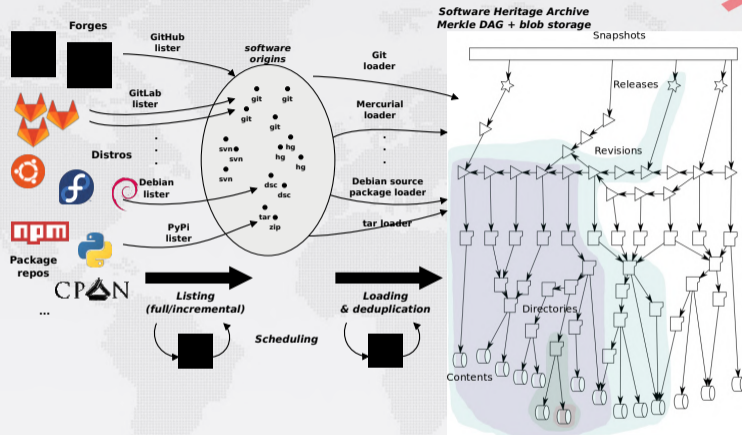
yes, you can sign it!

<https://en.unesco.org/foss/paris-call-software-source-code>



Their call is published on Feb 2019

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

Questions?

Learn more

www.softwareheritage.org/publications

-  **Jean-François Abramatic, Roberto Di Cosmo, Stefano Zacchioli**
Building the Universal Archive of Source Code, Communications of the ACM, October 2018
-  **P. Alliez, R. Di Cosmo, B. Guedj, A. Girault, M. Hacid, A. Legrand, N. Rougier**
Attributing and Referencing (Research) Software: Best Practices and Outlook From Inria, Computing in Science & Engineering, 22 (1), pp. 39-52, 2020, ISSN: 1558-366X
-  **Roberto Di Cosmo, Morane Gruenpeter, Stefano Zacchioli**
Referencing Source Code Artifacts: a Separate Concern in Software Citation, Computing in Science & Engineering, 2020, ISSN: 1521-9615
-  **Roberto Di Cosmo, Stefano Zacchioli**
Software Heritage: Why and How to Preserve Software Source Code, iPRES 2017



1 More about Software Heritage

2 The SWH-ID: the source code fingerprint

Highlights from the launching phase

Summer 2015



The collection starts: first server, (very) early prototype

June 30th 2016



Public unveiling, with the first sponsors: Microsoft and DANS

April 3rd 2017



Unesco - Inria agreement on software access and preservation.

June 7th 2018



Opening the archive to the world

December 7th 2018



Starting the mirror network

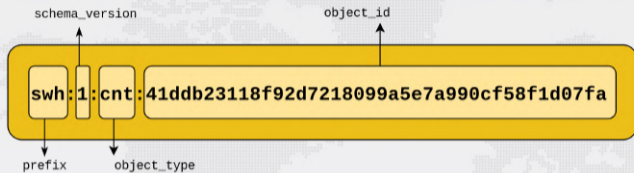
February 26th 2019



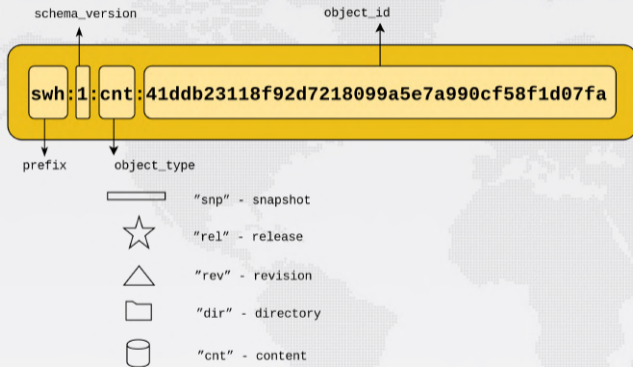
Publication of the expert meeting
Paris Call on
Software Source Code

- 
- 1 More about Software Heritage
 - 2 The SWH-ID: the source code fingerprint

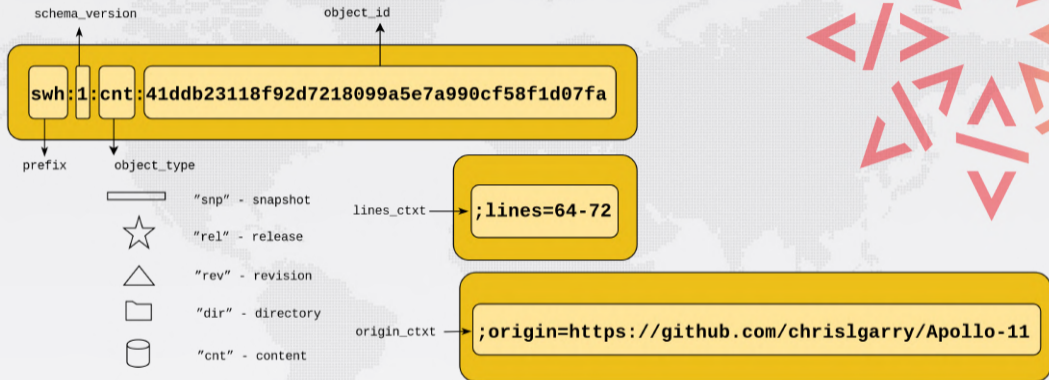
The SWH-ID schema



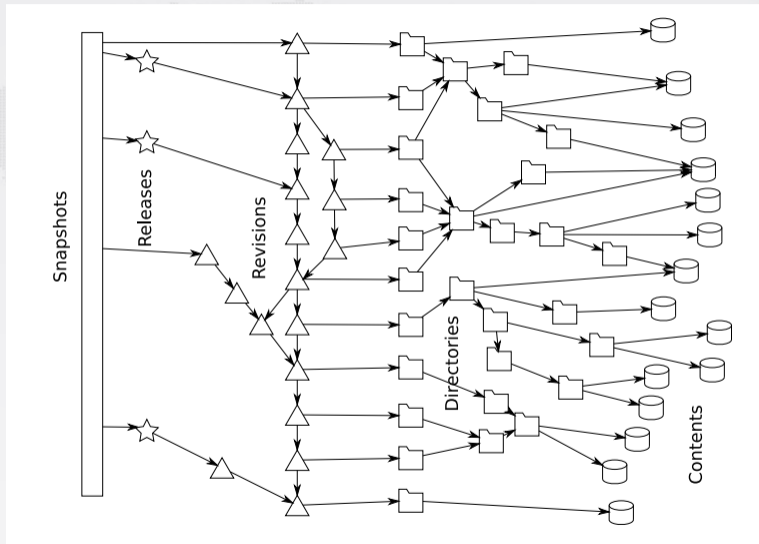
The SWH-ID schema



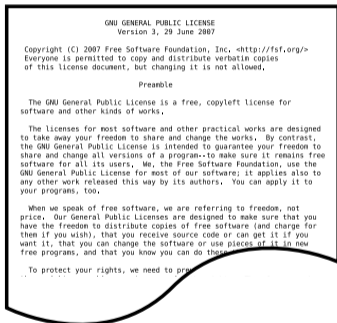
The SWH-ID schema



A worked example

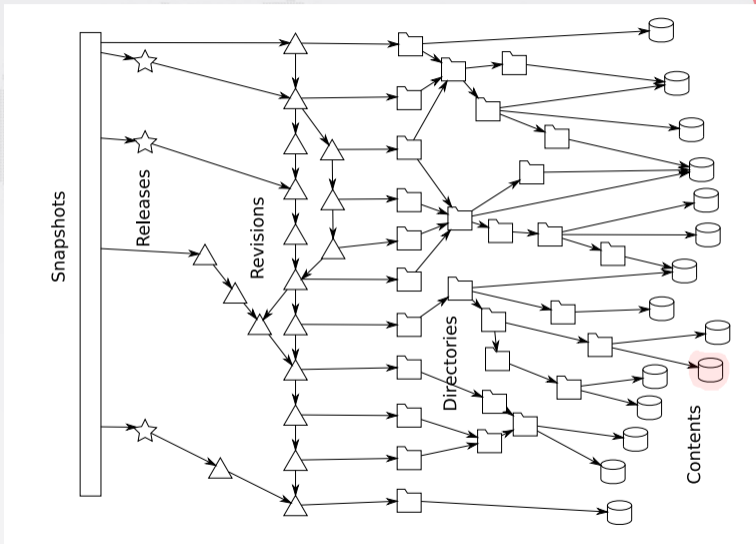


Contents



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

A worked example



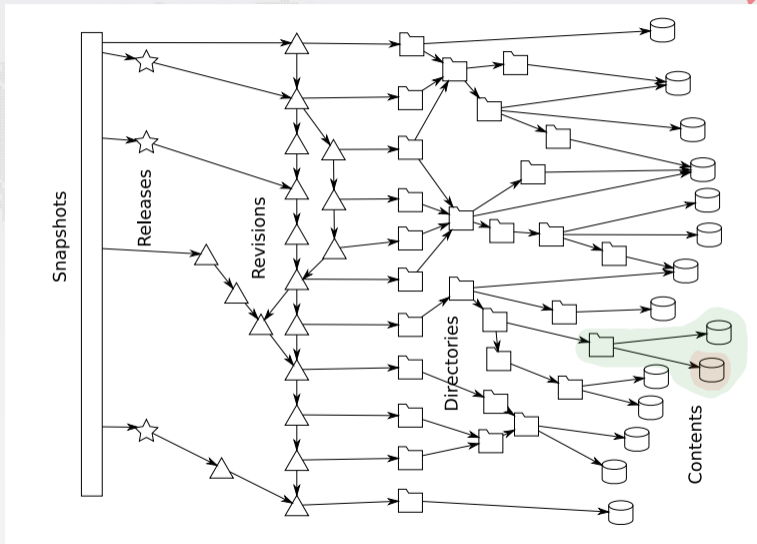


Directories


```
100644 blob c5baade4c44766042186ef858c0fd63d587ebf09 .gitignore
100644 blob 2d0a34af6f52cf3cf6b0c2f7bd0648fbd255e77f AUTHORS
100644 blob 94a9ed024d3859793618152ea559a168bbcbb5e2 LICENSE
100644 blob d9b2665a435a43f8a79a84e0867751dfb095c7bb MANIFEST.in
100644 blob 524175c2bad0b35b975f79284c2f5a6d5eaf2eb4 Makefile
100644 blob 5c7e3a5bbddb038682ba7793f440492ed9678bb3 Makefile.local
100644 blob 8617980629cd24e6080404f09aa749b085b3e07b README.db_testing
100644 blob 76b29f94cf815e0869c414d38d78d7ce08ec514e README.dev
040000 tree e1e10ece948af0b93adb0372afcf89f12e92618a bin
040000 tree 83e56d0beaf7793c77a45a345c80fcb8af503013 debian
040000 tree a34c9c4ba213f0cedc67f9816348d27955577af5 docs
100644 blob f2a6d32c6135aa7287bbd76167b01df2ae4f1539 requirements.txt
100755 blob eee147c36caf1bbc2d820da8dc026cb5b68180bc setup.py
040000 tree 224bb4c1f4c67fca1d160bffdd2d06094e7e1abf3 sql
040000 tree 8631c9cd77bbe993168107ab5baf51f40c6300be swl
040000 tree 8fb905b56ba8ed692f1209b2773b474c6c1d66c1 utils
```

id: 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d

A worked example



Revisions

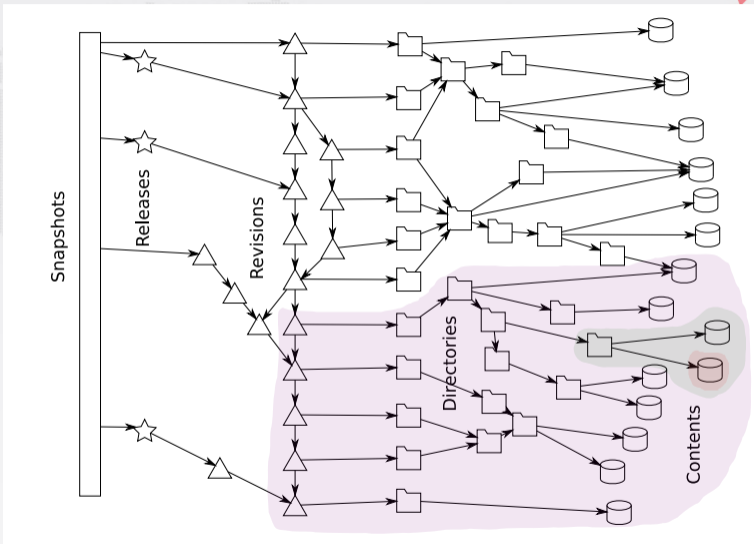
Details	Changes	Files
SHA: 963634dca6ba5dc37e3ee426ba091092c267f9f6		
Author: Nicolas Dandrimont <nicolas@dandrimont.eu> (Thu Sep 1 14:26:13 2016)		
Committer: Nicolas Dandrimont <nicolas@dandrimont.eu> (Thu Sep 1 14:26:13 2016)		
Subject: provenance.tasks: add the revision -> origin cache task		
Parent: fc3a8b59ca1df424d860f2c29ab07fee4dc35d10 : test...storage: properly pipeline origin and cont...		
provenance.tasks: add the revision -> origin cache task		
sw/wh/storage/provenance/tasks.py  77		

tree 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d
parent fc3a8b59ca1df424d860f2c29ab07fee4dc35d10
author Nicolas Dandrimont <nicolas@dandrimont.eu> 1472732773 +0200
committer Nicolas Dandrimont <nicolas@dandrimont.eu> 1472732773 +0200

provenance.tasks: add the revision -> origin cache task

id: 963634dca6ba5dc37e3ee426ba091092c267f9f6

A worked example



Releases

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

```
Release sw.h.storage v0.0.51
```

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

```
commit c0c9f16b1e134f593e7567570a1761b156e6eb1d
```

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

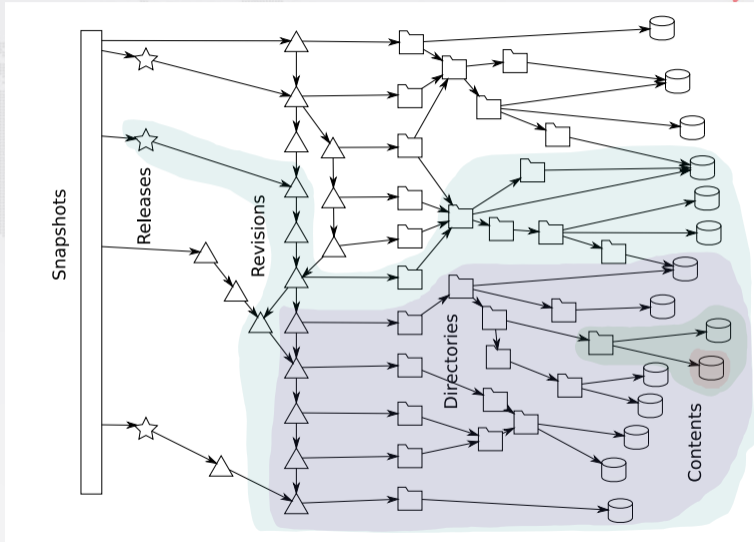
```
Release sw.h.storage v0.0.51
```

```
- Add new metadata column to origin_visit
- Update sw.h-add-directory script for updated API
---BEGIN PGP SIGNATURE---
```

```
iQIzBAABCAAdBQJXvZTNFhxuaWNvbGFzQGRhbmRyaW1vbnQuZXUACgkQ7AWLMo2+
neqorw//aq6SOB5DijzEa+kWN3rXgVS+1K1vEVh1wNKAw8eKJ7aX2kEILDtt7uf
ahpZ6pz3q8nqs6aC1+YrxBfcih3L2YtrdZeWXWqr8xWNMaEoYDb8qaphwh8AD5t2
ICBlit2ujXuCrDt93eKKPwvzZXg+h80sMwy35Dr6jW7Z7K4Mu/PgGlyIHPY55yo
IGEndWno7VFH1Vm6t1n5qB7I5mXRaqA+becqddubTZ2xjj+jpUqC8cyqN3hm/fL
qsJ2mu8kyz3t8tG/H1/pV+I5OwBlNpO5STH0tujoEvgPK/dHSP79QuHDHZFkCao
klj6kAWyU80Mxb+nKV/jeLbrR3+yWBFJ3Qp5a1/V8oOTh6E1dALcNmPeaKCoKtMt
d/gMRax111/g0EDfnsW67G6sDwKPKPHngfVLQ3nV3GaQQTnu1RpMz006H9/tAwzC
Gg/K1PdHT4hz0I46wYPZyje0U2VXGFu6vVU9vFQ4ZR/Wjn+0zZdcRdrJJSUOMn
RpTTfusbXUeXHGOpkXhSYTnvp1gdPc76U5TsK0aGe84AZm1Ik0mGrwXCvFPqYo
nhhibB5HBNMoqyF6yTSOpUbYK70tpYRRUGKwDeRk0wKSxkWKUZGtKzy6jYqJjo29
guLwZQif5qWQC80ontAL2+HvFfaVyckMejUhg62cP/+EHlvUk=
=kOxP
---END PGP SIGNATURE---
```

id: 85083a5cc14a441c89dea73f5bdf67c3f9c6afdb

A worked example



Snapshots

git show-refs

```
commit 08ffeb25770109525eb3ce21691466c53a1d9158 refs/heads/atime
commit ba5443a24e3f9fe323a46c292cec4fcbe61c67eb refs/heads/directory-listing-arrays
commit d69e0dbf892383ff6589b27f1c05d27238d9c5 refs/heads/foo
commit cf7ff9eea0eb22f8946908f5a8019f67de468e08 refs/heads/master
commit 7eca197fc66d2024047e54b1ed9e8b44361a0fc2 refs/heads/tmp-directory-add
commit 642a205f37de85005a85d427b53ee4fb2252e82e refs/heads/tmp/generic-releases
tag 20f043b1379c7f68d966597799fd4907c757f755 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 586fba4e580b4f5fab05f599367643cbb1a9c7f refs/tags/v0.0.15
tag 8cd8b885f4098bf363177742bd289f660e5be51c refs/tags/v0.0.16
tag a542444ee3f0fbcd35efb202fee035c809abc7d6 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 3fb168c2072a5d6252124257a1e5dfc0f5ffa1df refs/tags/v0.0.22
tag 8cddb0e8da4d731c5d262789e460a16ac3c72aba4 refs/tags/v0.0.23
...
```

id: b464cad1b66fff266a37b46ea6e7a04b545e904b

Let's look at some famous excerpts of source code

Apollo 11 source code (excerpt)

```
P63SP0T3      CA      BIT6          # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND      CHAN33
              EXTEND
              BZF      P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF      CODE500      # ASTRONAUT: PLEASE CRANK THE
              TC      BANKCALL      # SILLY THING AROUND
              CADR      GOPERF1
              TCF      GOTOP00H      # TERMINATE
              TCF      P63SP0T3      # PROCEED SEE IF HE'S LYING

P63SP0T4      TC      BANKCALL      # ENTER INITIALIZE LANDING RADAR
              CADR      SETPOS1

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

Quake III 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;
}
```

It works!

we have *intrinsic* identifiers for all 20+ billion objects in the archive

The Software Heritage ID schema (see <http://bit.ly/swhpids>)

`swh:1:cnt:94a9ed024d3859793618152ea559a168bbcbb5e2` full text of the GPL3 license

`swh:1:dir:d198bc9d7a6bcf6db04f476d29314f157507d505` Darktable source code

`swh:1:rev:309cf2674ee7a0749978cf8265ab91a60aea0f7d`

a **revision** in the development history of Darktable

`swh:1:rel:22ece559cc7cc2364edc5e5593d63ae8bd229f9f`

release 2.3.0 of Darktable, dated 24 December 2016

`swh:1:snp:c7c108084bc0bf3d81436bf980b46e98bd338453`

a **snapshot** of the entire Darktable repository (4 May 2017, GitHub)

Current resolvers: archive.softwareheritage.org and n2t.org