

# The SWH-ID

A digital fingerprint identifying software source code

Roberto Di Cosmo, Morane Gruenpeter

January 29th, 2020



# Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 
- 1 Software is our heritage
  - 2 Preserving all software source code
  - 3 The SWH-ID: the source code fingerprint
  - 4 Conclusion

# Source Code: *executable* and *human readable* knowledge



*"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");
}
```

Len Shustek, CHM

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

# The Paris call: Software Source Code is part of our Heritage

November, 2018 at the **UNESCO** headquarters experts signed *the engagement*



- **Recognise** software source code as a precious asset of humankind
- **Support** the development of shared infrastructures
- **Foster** international collaboration to build a common framework

*see full text*

- 
- 1 Software is our heritage
  - 2 Preserving all software source code
  - 3 The SWH-ID: the source code fingerprint
  - 4 Conclusion



## Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

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

Preserving our heritage, enabling better software and better science for all

### Reference catalog



find and reference **all** the source code

### Universal archive



preserve **all** the source code

### Research infrastructure



enable analysis of **all** the source code

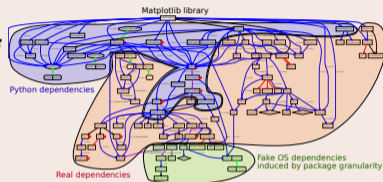
# Source code is *special*

Software **evolves** over time

- projects may last decades
- the *development history* is key to its *understanding*

Layers of **complexity**

- *millions* of lines of code
- large *web of dependencies*
- sophisticated *developer communities*



## Bottomline

- we must archive *all* the source code
- we must preserve *all* the history of its development
- we must **identify** *all* the archived software artifacts (more than 20 billions today!)

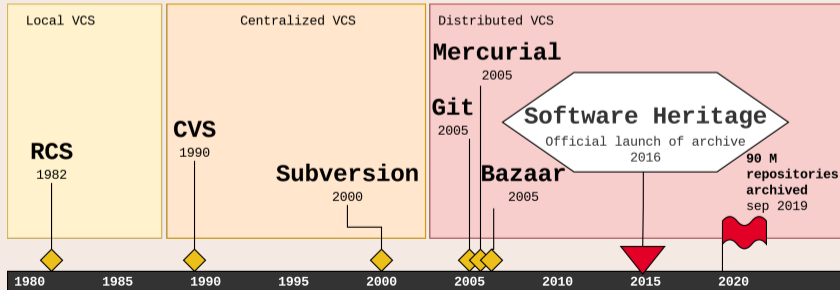
how can we do this?

# Evolution of software development

## Version control system (VCS)

- records changes made to a (set of) *source code file* (s)
- allows to operate on versions: diff/merge/fork/recover etc.
- **essential** tool for software development

## Three decades of evolution

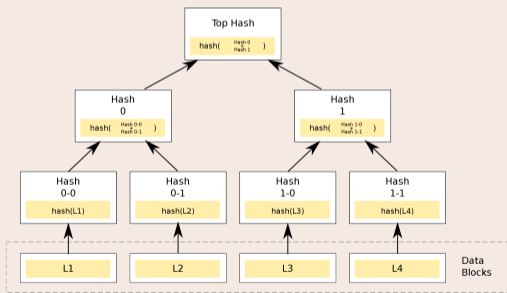




## Requirements for the D in DVCS

- **intrinsic** unique identifiers... (here: *cryptographic signature*, aka "hash")
- ... that work for **tree structures** (software directories)

## Merkle tree to the rescue (R. C. Merkle, Crypto 1979)



Combination of

- tree
- hash function

# A massive adoption

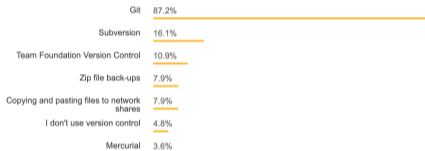
## Stack Overflow

[Survey 2018]

### Version Control

All Respondents

Professional Developers



74,298 responses; select all that apply

Git is the dominant choice for version control for developers today, with almost 90% of developers checking in their code via Git.

## In numbers

GitHub [Octoverse 2017] [Blog 2018]

- 100.000.000+ repositories
- 40.000.000+ developers worldwide

Bitbucket [Blog 2019]

- 28.000.000+ repositories
- 10.000.000+ developers worldwide

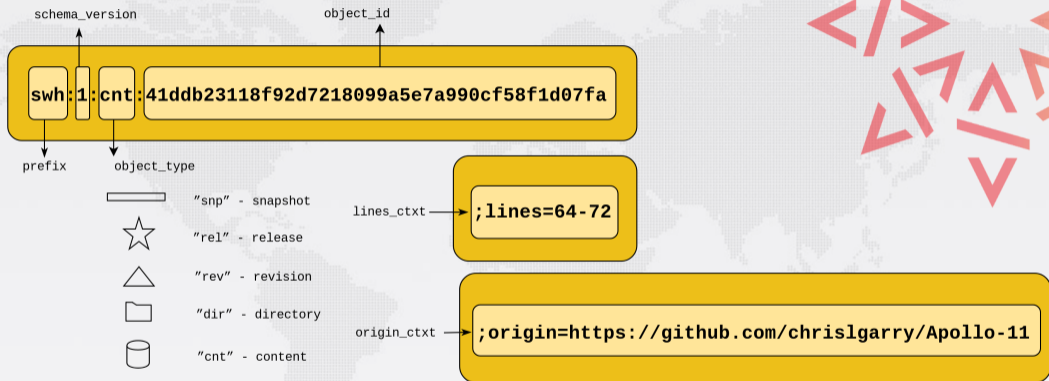
GitLab [Blog 2019]

- 1.000.000 MRs March 19'

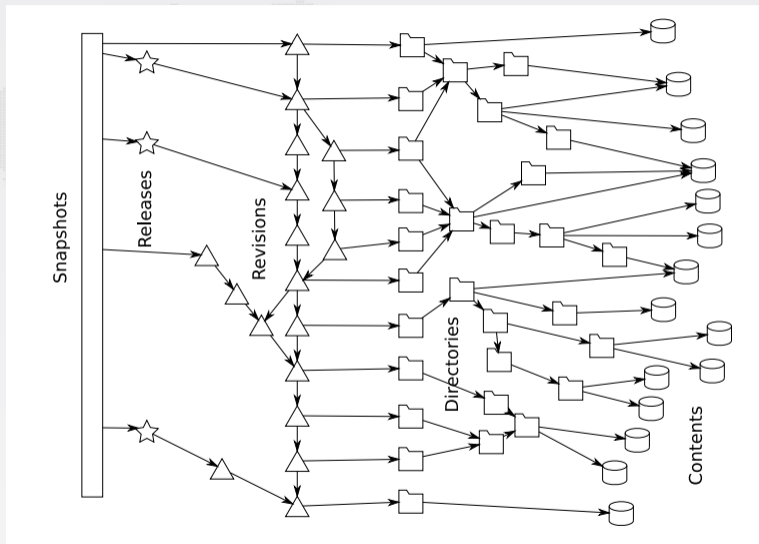
Let's use it!

- 
- 1 Software is our heritage
  - 2 Preserving all software source code
  - 3 The SWH-ID: the source code fingerprint
  - 4 Conclusion

# The SWH-ID schema



# 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.

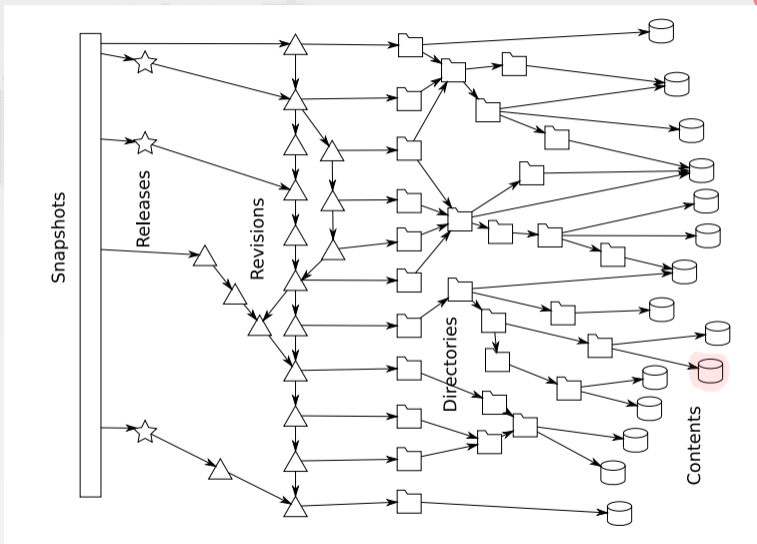
The licenses for most software and other practical works are designed
to take away your freedom to share and change the works. By contrast,
the GNU General Public License is intended to guarantee your freedom to
share and change all versions of a program—to make sure it remains free
software for all 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 prevent anyone from denying you
```

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

# A worked example





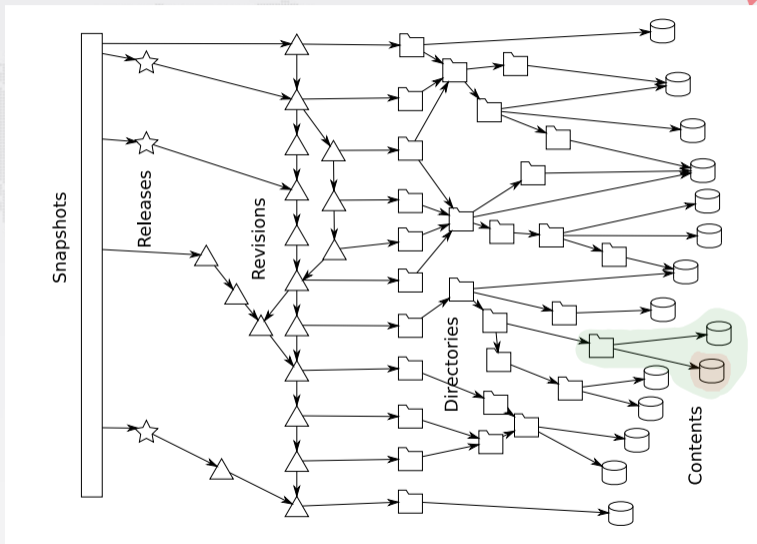
## Directories

```
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 e1e10ecef948af0b93adb0372afc89f12e92618a 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 swh
040000 tree 8fb905b56ba8ed692f1209b2773b474c6c1d66c1 utils
```


id: 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d



# A worked example



## Revisions

Details	Changes	Files
SHA: 963634dca6ba5dc37e3ee426ba091092c267f9f6		
Author: <a href="mailto:nicolas@dandrimont.eu">Nicolas Dandrimont &lt;nicolas@dandrimont.eu&gt;</a> (Thu Sep 1 14:26:13 2016)		
Committer: <a href="mailto:nicolas@dandrimont.eu">Nicolas Dandrimont &lt;nicolas@dandrimont.eu&gt;</a> (Thu Sep 1 14:26:13 2016)		
Subject: provenance.tasks: add the revision -> origin cache task		
Parent: <a href="#">fc3a8b59ca1df424d860f2c29ab07fee4dc35d10</a> : test...storage: properly pipeline origin and cont...		
provenance.tasks: add the revision -> origin cache task		
<a href="#">swb/storage/provenance/tasks.py</a>  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



## 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
```

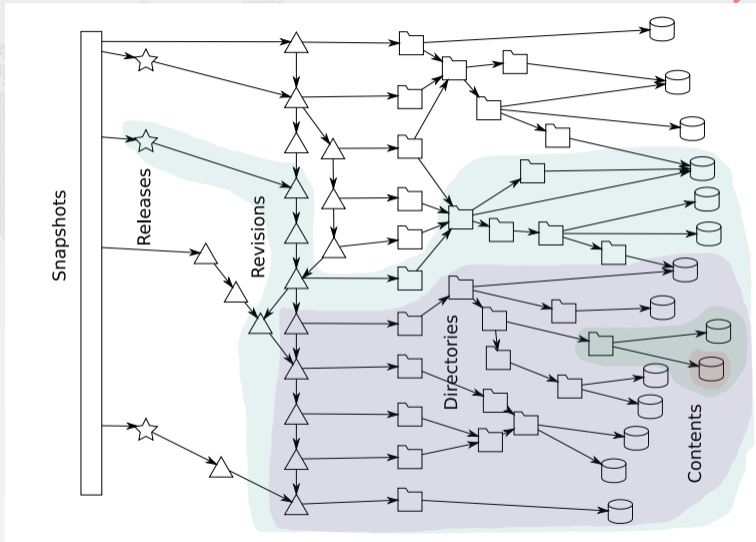
```
Release sw.h.storage v0.0.51
```

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

```
iQIzBAABCAAdBQJXvZTNFhXuaWNvbGFzQGRhbmRyaW1vbnQuZXUACgkQ7AWLMO2+
neqorw/aa65Ob5DjzEa+kWN3rXgV5+1K1vEVh1wNKAw8eKJ7aX2kEiLDtt7uf
ahpZ6pz3q8nqs6aC1+YrxBfcih3L2YtrdZeWXXWqr8xWNMaEoYDb8qaaphwh8AD5t2
ICBhit2ujtXuCrDt93eKKPwvzXG+h80sMWy35Dr6jW7Z7K4MuPGglyIHPY55yo
IGEndWno7Vfh1Vm6t1n5qB7I5mXRaqA+becqddubTZ2xij+jpUqC8cyqN3hm/fL
qsj2mu8kyz3t8tG/H1/pV+15OwBlNpO5TH0tujojEVgPK/dH5P79QuHDHZFkCao
klj6kAWyU80Mxb+nKVjeLbrR3+yWBFj3Qp5a1/V8o0Th6E1dALcNmPEaKCoKtMt
d/gMRax11I/g0EDfnsW67G6sDwKPKPHngfVLQ3nV3GaQQTnu1RpMz006H9/tAwzC
Gg/K1PdHT4hz0iI46wYPZje0U2VXGFu6vVU9vFQ4ZR/Wjn+0zZdcRdrjSUOMn
RpTTRfUbsXUeXHGOpkgXhSYTnvp1gdPc76U5TsK0aGe84A2m1lk0mGrwXCvFPqYo
nhhibB5HBNMoqyF6yTSOpUbyK70tpYRRUGKWDeRK0wKSxkWKUZGtKzy6jYqJzo29
gulwgZQif5qWQCB0oontAL2+HvPfaVyckMejUhg62cP/+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 d69e0dbf892383ff6589b27fbc1c05d27238d9c5 refs/heads/foo
commit cf7ff9eea0eb22f8946908f5a8019f67de468e08 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 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 215ea50daball1e082e0b72e76eb4b6073a87908 refs/tags/v0.0.21
tag 3fb168c2072a5d6252124257a1e5dfc0f5ffa1df refs/tags/v0.0.22
tag 8cdbee8da4d73fc5d262789e460a16ac3c72aba4 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      #
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

- 
- 1 Software is our heritage
  - 2 Preserving all software source code
  - 3 The SWH-ID: the source code fingerprint
  - 4 Conclusion



## Intrinsic identifiers...

- can be extracted from the **object itself**, hence:
  - no need for a *central authority*, nor maintenance
  - any modification to the object changes the identifier
- identifies the *object*, not the *metadata* !



## ... *for source code*

- Distributed Version Control Systems made them popular
- massively used every day by millions of software developers
- Software Heritage provides **SWH-IDs** for billions of software artifacts

[www.softwareheritage.org](http://www.softwareheritage.org) – learn more  
[save.softwareheritage.org](http://save.softwareheritage.org) – save code now  
[www.softwareheritage.org/swhap](http://www.softwareheritage.org/swhap) – legacy software acquisition process  
[forge.softwareheritage.org](http://forge.softwareheritage.org) – our own code

## Questions?

### References

-  **Jean-François Abramatic, Roberto Di Cosmo, Stefano Zacchioli**  
*Building the Universal Archive of Source Code*,  
Communications of the ACM, October 2018 ([10.1145/3183558](https://doi.org/10.1145/3183558))
-  **Roberto Di Cosmo, Morane Gruenpeter, Stefano Zacchioli**  
*Referencing Source Code Artifacts: a Separate Concern in Software Citation*,  
Computing in Science and Engineering, IEEE, pp.1-9. ([10.1109/MCSE.2019.2963148](https://doi.org/10.1109/MCSE.2019.2963148)) ([hal-02446202](https://hal.archives-ouvertes.fr/hal-02446202))