

# Software (source code) and Open Science

## Challenges and Opportunities

Roberto Di Cosmo

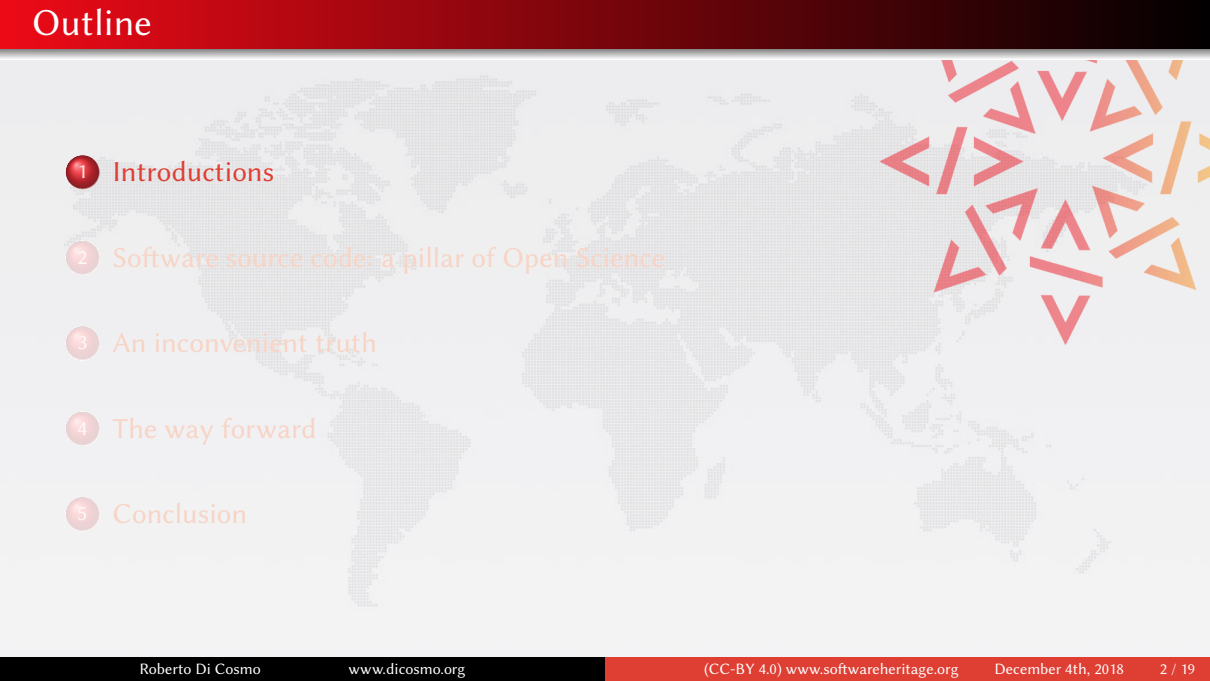
[roberto@dicosmo.org](mailto:roberto@dicosmo.org)

December 4th, 2018



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 
- 1 Introductions
  - 2 Software source code: a pillar of Open Science
  - 3 An inconvenient truth
  - 4 The way forward
  - 5 Conclusion

# Short Bio: Roberto Di Cosmo

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

2008 *Mancoosi project* [www.mancoosi.org](http://www.mancoosi.org)

2010 *IRILL* [www.irill.org](http://www.irill.org)

2015 *Software Heritage* at INRIA

2018 *National Committee for Open Science*, France

*Key mediator* for accessing *all* information (c) Banski

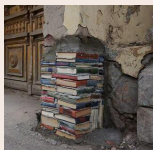


Information is a main pillar of our modern societies.

*Absent an ability to correctly interpret digital information, we are left with [...] "rotting bits" [...] of no value.*

*Vinton G. Cerf IEEE 2011*

## Key mediator for accessing *all* information (c) Banski



Information is a main pillar of our modern societies.

*Absent an ability to correctly interpret digital information, we are left with [...] "rotting bits" [...] of no value.*

Vinton G. Cerf IEEE 2011

## Software is *an essential component* of modern scientific research

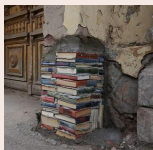


*[...] the vast majority describe experimental methods or software that have become essential in their fields.*

Top 100 papers (Nature, October 2014)

# Software is knowledge

*Key mediator for accessing all information* (c) Banski



Information is a main pillar of our modern societies.

*Absent an ability to correctly interpret digital information, we are left with [...] "rotting bits" [...] of no value.*

*Vinton G. Cerf IEEE 2011*

Software is an essential component of modern scientific research



*[...] the vast majority describe experimental methods or software that have become essential in their fields.*

Top 100 papers (Nature, October 2014)

Bottomline: Software embodies our *Knowledge* and *Cultural Heritage*

*It must be collected, referenced and made accessible!*

# The knowledge is in the source code!



*"The source code for a work means the preferred form of the work for making modifications to it."*

GPL Licence



# The knowledge is in the source code!



*"The source code for a work means the preferred form of the work for making modifications to it."*

GPL Licence

Hello World



# The knowledge is in the source code!



*"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
```

# The knowledge is in the source code!



*"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");
}
```

# 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 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;
}
```

## 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 algo 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.”*

## Apollo 11 Guidance Computer (~60.000 lines), 1969



"When I first got into it, nobody knew what it was that we were doing. It was like the Wild West."

Margaret Hamilton

# ~ 50 years, a lightning fast growth

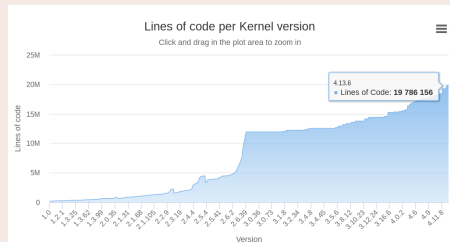
## Apollo 11 Guidance Computer (~60.000 lines), 1969



"When I first got into it, nobody knew what it was that we were doing. It was like the Wild West."

Margaret Hamilton

## Linux Kernel



... now in your pockets!

- 
- 1 Introductions
  - 2 Software source code: a pillar of Open Science
  - 3 An inconvenient truth
  - 4 The way forward
  - 5 Conclusion

# The scientific method...

## The experimental method



- make an *observation*
- formulate an *hypothesis*
- set up an **experiment**
- elaborate a *theory*

And then we **reproduce** and **verify**.

# The scientific method...

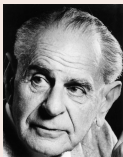
## The experimental method



- make an *observation*
- formulate an *hypothesis*
- set up an **experiment**
- elaborate a *theory*

And then we **reproduce** and **verify**.

## Reproducibility is the key



*non-reproducible single occurrences are of no significance to science*

*Karl Popper, The Logic of Scientific Discovery, 1934*



# ... evolves in the digital age!

For an experiment involving software, we need

- open access** to the scientific article describing it
- open data sets** used in the experiment
- source code** of all the components
- environment** of execution
- stable references** between all this



## ... evolves in the digital age!

For an experiment involving software, we need

- open access** to the scientific article describing it
- open data sets** used in the experiment
- source code** of all the components
- environment** of execution
- stable references** between all this

### Remark

The first two items are already widely discussed!

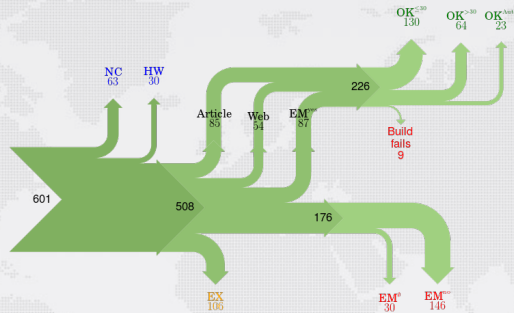
... what about *software*?

- 8 ACM conferences: ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12

- all very practical oriented

can we get the code to build and run?

[illegible]



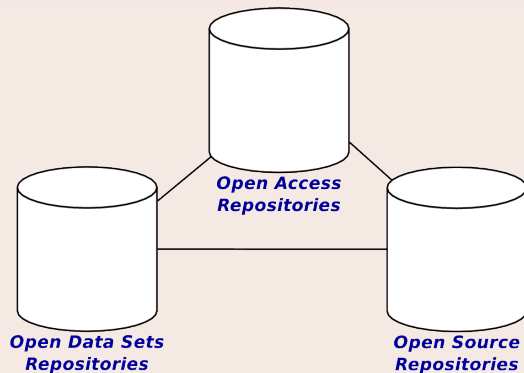
... that's a whopping 40% of **non reproducible** works!

## The main reasons

source code (*or the right version of it*) cannot be found

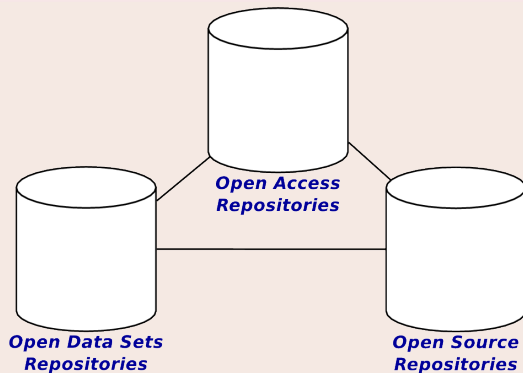
# Software Source code is an important pillar

## The Magic Triangle of Scientific Knowledge



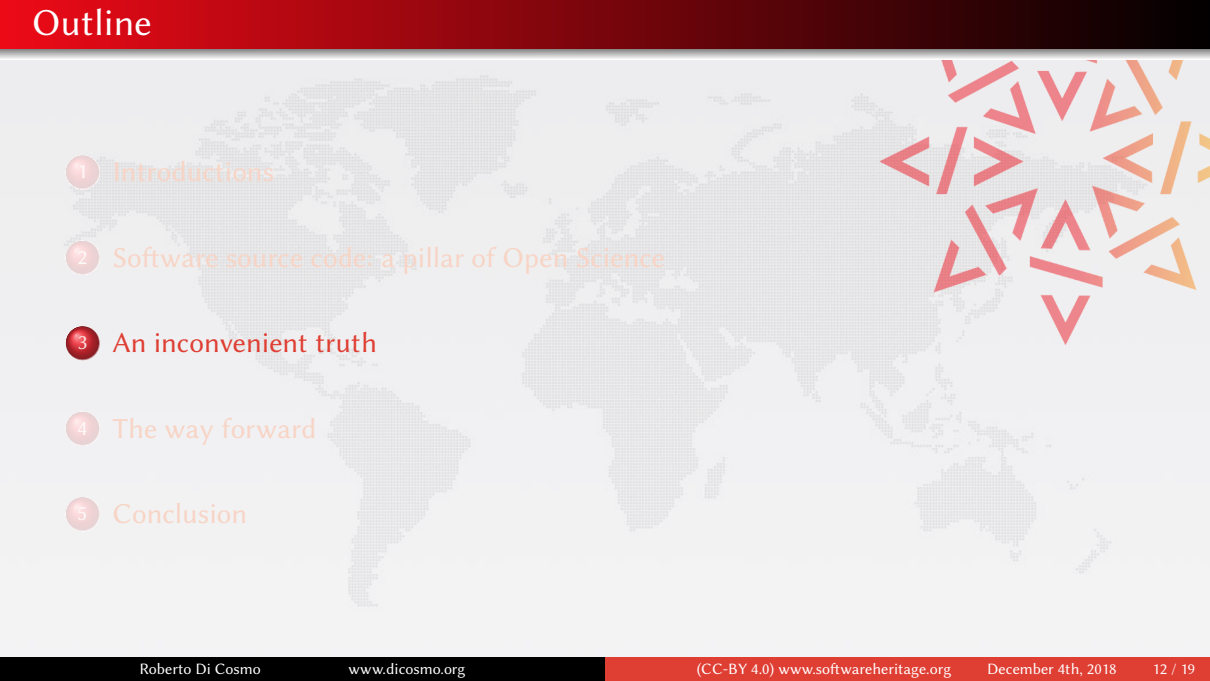
# Software Source code is an important pillar

## The Magic Triangle of Scientific Knowledge



Nota bene

The links in the picture are **essential**

- 
- 1 Introductions
  - 2 Software source code: a pillar of Open Science
  - 3 An inconvenient truth**
  - 4 The way forward
  - 5 Conclusion

# A forgotten pillar of Open Science

No reference catalog



to find and reference **all**  
the source code





# A forgotten pillar of Open Science

## No reference catalog



to find and reference **all**  
the source code

## No universal archive



to preserve **all** the source  
code



# A forgotten pillar of Open Science

## No reference catalog



to find and reference **all**  
the source code

## No universal archive



to preserve **all** the source  
code

## No research infrastructure



to enable analysis of **all**  
the source code

# A forgotten pillar of Open Science

## No reference catalog



to find and reference **all** the source code

## No universal archive



to preserve **all** the source code

## No research infrastructure



to enable analysis of **all** the source code

## Lack of recognition

not (yet) a first class citizen

- in the EOSC plan
- in the EU copyright reform
- in the scholarly works

# A forgotten pillar of Open Science

## No reference catalog



to find and reference **all** the source code

## No universal archive



to preserve **all** the source code

## No research infrastructure



to enable analysis of **all** the source code

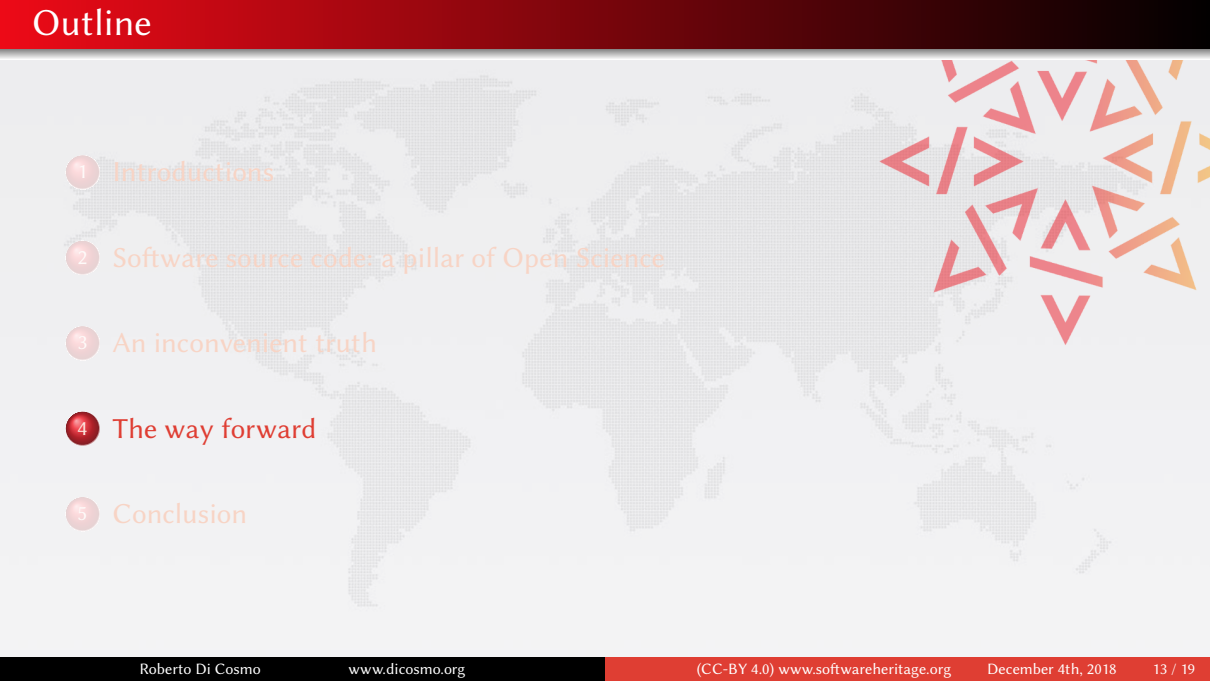
## Lack of recognition

not (yet) a first class citizen

- in the EOSC plan
- in the EU copyright reform
- in the scholarly works

## Lack of guidance on how to

- choose a license
- cite a software project
- relate to industry best practices
- make source code FAIR(\*)

- 
- 1 Introductions
  - 2 Software source code: a pillar of Open Science
  - 3 An inconvenient truth
  - 4 The way forward**
  - 5 Conclusion

# Raising Awareness

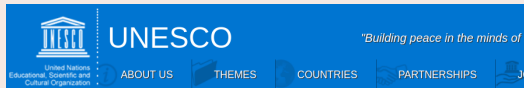
Inria Unesco agreement, April 3rd, 2017



Roberto Di Cosmo

[www.dicosmo.org](http://www.dicosmo.org)

Unesco Inria expert group, November 2018



Home > All News > Experts call for greater recognition of software source code as heritage for sustainable development

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

16 November 2018



(CC-BY 4.0) [www.softwareheritage.org](http://www.softwareheritage.org)

December 4th, 2018

13 / 19

## In the Research Data Alliance

Collaboration with a variety of international partners

- Source Code Interest Group
- Source Code Identification Working Group

## In the French Open Science Plan

- the GPLO group
  - software citation, reference, archival
  - software licensing
  - best practices
- support for Software Heritage



# Software Heritage



## Mission

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





# Software Heritage

## Mission

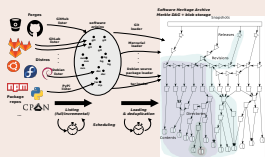
Collect, preserve and share the *source code* of *all the software* that is available

## Partners

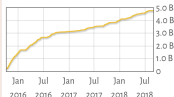
Initiator Inria

Industry Microsoft, Intel, Société Générale, Google, GitHub, FOSSID

Public sector UNESCO, DINSIC, DANS, UQAM, Bologna University

The largest software source code archive *ever*

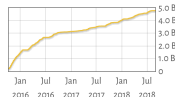
5,335,334,261



1,191,059,588

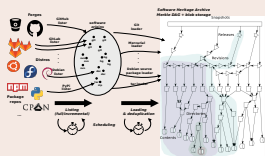


86,555,394



Must read: conceptual framework for DIOs and IDOs

[bit.ly/swhpidpaper](https://bit.ly/swhpidpaper)

The largest software source code archive *ever*

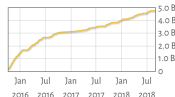
## Source files

5,335,334,261



## Commits

1,191,059,588



## Projects

86,555,394

Over 10 billions *intrinsic* identifiers (IDOs) for reproducibility

Must read: conceptual framework for DIOs and IDOs

bit.ly/swhidpaper

## Research software deposit

- moderated via **HAL**  
*open since 9/2018*

## Reference archive

- See for example

swmath.org

## Collaboration HUB

- industry, research
- digital preservation

Now part of the *French National Plan for Open Science*

# Reduce risk, avoid fragmentation

**Cultural Heritage**



**Industry**



**Research**



**Education**



Software Heritage



# Reduce risk, avoid fragmentation



Thomas Jefferson, February 18, 1791

*...let us save what remains: not by vaults and locks which fence them from the public eye and use in consigning them to the waste of time, but by such a multiplication of copies, as shall place them beyond the reach of accident.*

# Reduce risk, avoid fragmentation



Thomas Jefferson, February 18, 1791

*...let us save what remains: not by vaults and locks which fence them from the public eye and use in consigning them to the waste of time, but by such a multiplication of copies, as shall place them beyond the reach of accident.*

A *common* infrastructure

- **mutualisation** for sustainability
- open source, **non for profit**
- mirror network **open to all**
- **may** prevent a useless diaspora

# A word on FAIR

Research Source Code ... is just Source Code!

FAIR for Research Software Source Code is *different*



Research Source Code ... is just Source Code!

FAIR for Research Software Source Code is *different*

For Software Source Code, FAIR has a *different meaning*:

reFerenced with **intrinsic**, **verifiable** identifiers

- see the iPres 2018 article [bit.ly/swhpdpaper](http://bit.ly/swhpdpaper)

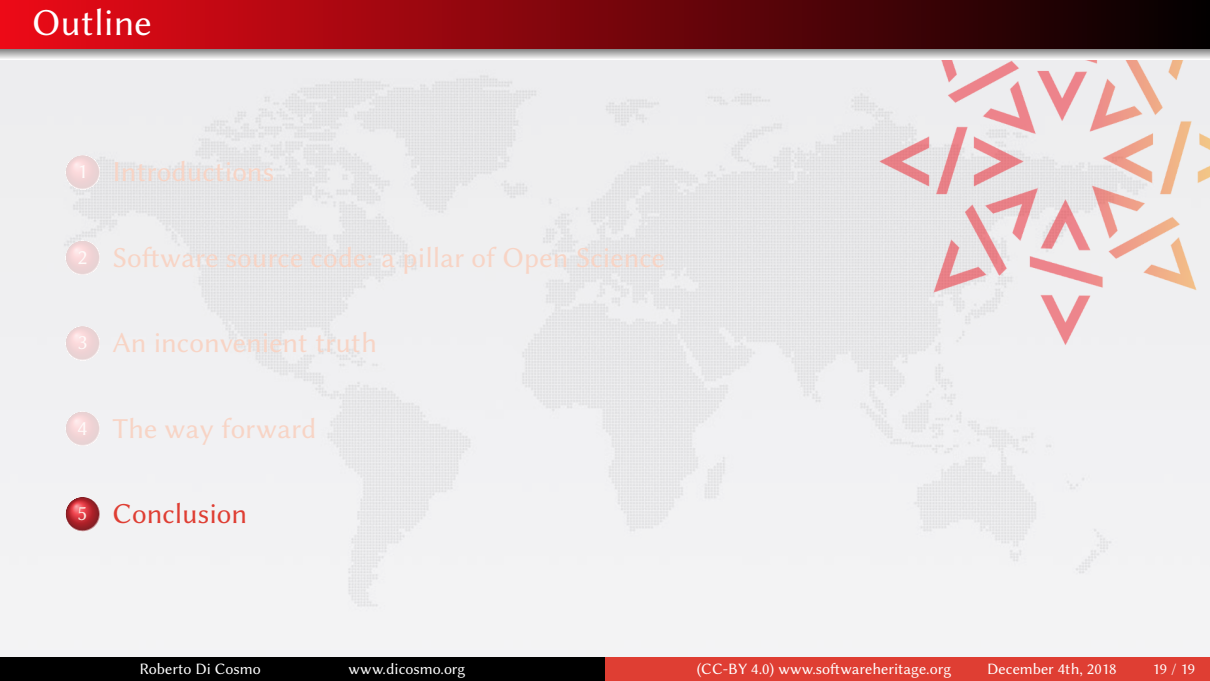
- example:

swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa;lines=53-82

Accessible in an **archive** that holds it for the **long term**

clted to **credit authors**, like all other scientific outputs

Reusable equipped with a proper **Open Source license**

- 
- 1 Introductions
  - 2 Software source code: a pillar of Open Science
  - 3 An inconvenient truth
  - 4 The way forward
  - 5 Conclusion**

## Challenges

### Software Source code:

- (forgotten) **pillar** of Open Science
- (undervalued) **key** to reproducibility
- (underrated) **scholarly production**



## Challenges

Software Source code:

- (forgotten) **pillar** of Open Science
- (undervalued) **key** to reproducibility
- (underrated) **scholarly production**

## Opportunities

Shared with Open Source communities

- **learn** from software development
- **adopt** proven approaches
- **avoid** dispersion of efforts

## Challenges

Software Source code:

- (forgotten) **pillar** of Open Science
- (undervalued) **key** to reproducibility
- (underrated) **scholarly production**

## Opportunities

Shared with Open Source communities

- **learn** from software development
- **adopt** proven approaches
- **avoid** dispersion of efforts



Jean-François Abramatic, Roberto Di Cosmo, Stefano Zacchioli

Building the Universal Archive of Source Code

Communication of the ACM, October 2018



Roberto Di Cosmo, Morane Gruenpeter, Stefano Zacchioli

Identifiers for Digital Objects: the Case of Software Source Code Preservation


iPRES 2018: Intl. Conf. on Digital Preservation



Roberto Di Cosmo, Publication scientifique: le rôle des États dans l'ère des TIC.

Upgrade, Vol. VII, No. 3, June 2006,

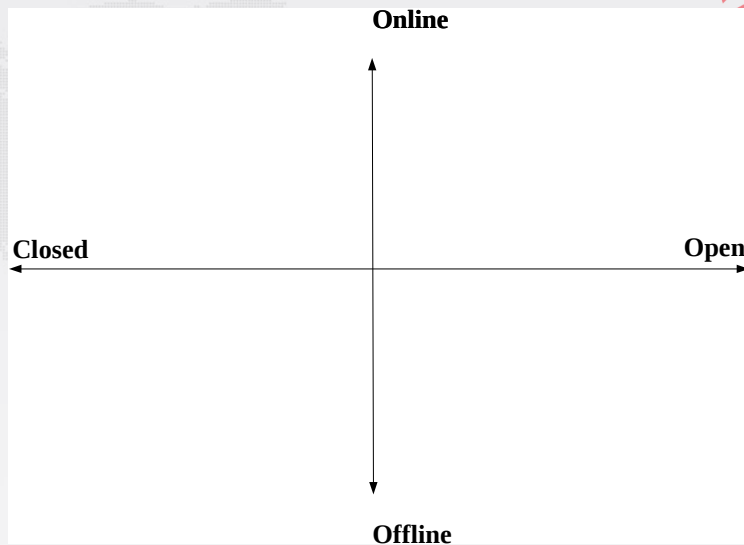
<http://www.dicosmo.org/FreeAccessToScience.pdf>

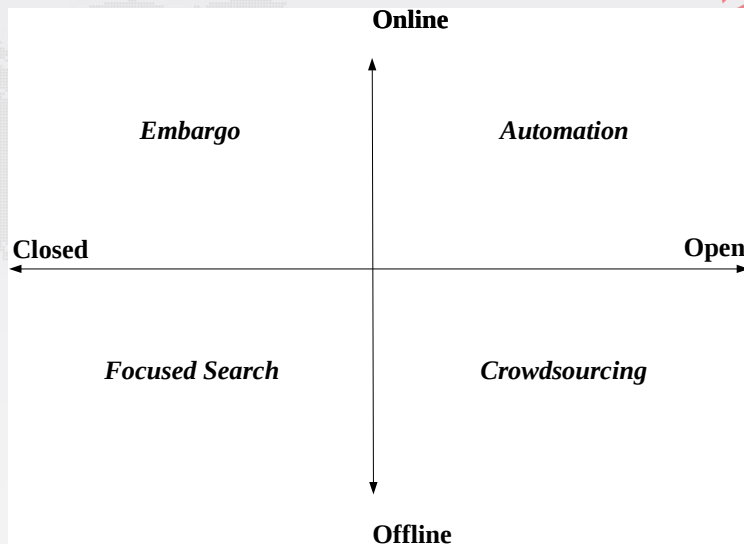


6 Strategy

7 Identifiers are not easy

8 Looking for the right identifiers









6 Strategy

7 Identifiers are not easy

8 Looking for the right identifiers

# URL decay disrupts the *web of reference*

Web links *are not* permanent (even *permalinks*)

*there is no general guarantee that a URL... which at one time points to a given object continues to do so*

*T. Berners-Lee et al. Uniform Resource Locators. RFC 1738.*

# URL decay disrupts the *web of reference*

Web links *are not* permanent (even *permalinks*)

*there is no general guarantee that a URL... which at one time points to a given object continues to do so*

*T. Berners-Lee et al. Uniform Resource Locators. RFC 1738.*

404

URLs used in articles *decay*!

Analysis of *IEEE Computer* (Computer), and the *Communications of the ACM* (CACM): 1995-1999

- the *half-life* of a referenced URL *is approximately 4 years* from its publication date  
D. Spinellis. The Decay and Failures of URL References.

Communications of the ACM, 46(1):71-77, January 2003.

# URL decay disrupts the *web of reference*

Web links *are not* permanent (even *permalinks*)

*there is no general guarantee that a URL... which at one time points to a given object continues to do so*

*T. Berners-Lee et al. Uniform Resource Locators. RFC 1738.*

404

URLs used in articles *decay*!

Analysis of *IEEE Computer* (Computer), and the *Communications of the ACM* (CACM): 1995-1999

- the *half-life* of a referenced URL *is approximately 4 years* from its publication date  
D. Spinellis. The Decay and Failures of URL References.

Communications of the ACM, 46(1):71-77, January 2003.

Similar findings in Lawrence, S. et al. *Persistence of Web References in Scientific Research*, IEEE Computer, 34(2), pp. 26-31, 2001.

## An example from Astronomy

Domain	links (broken)	.html	.txt	.dat	.gz	.tar	.fits	tilde
ooc.harvard.edu	802 (110)	336 (70)	0	0	4 (2)	5 (4)	1	0
heasarc.gsfc.nasa.gov	640 (33)	423 (27)	1	0	0	0	0	0
www.stsci.edu	498 (61)	205 (29)	3	0	0	0	0	15 (10)
asc.harvard.edu	471 (152)	212 (99)	0	0	0	0	0	1 (1)
ssc.spitzer.caltech.edu	427 (194)	125 (76)	3 (3)	0	0	0	0	0
cfa-www.harvard.edu	352 (68)	277 (52)	1	0	0	0	0	54 (17)
archive.stsci.edu	308 (58)	57 (9)	2	1 (0)	0	0	0	0
www.ipac.caltech.edu	285 (14)	209 (12)	0	0	0	0	0	0
www.atnf.csiro.au	211 (21)	12 (6)	0	0	0	0	0	7 (5)
space.mit.edu	193 (10)	58 (5)	1	0	0	0	0	2 (1)
www.astro.psu.edu	186 (4)	103 (1)	1	10	1	1	0	2
www.eso.org	186 (58)	54 (22)	1 (1)	0	0	0	0	4 (1)
irsa.ipac.caltech.edu	163 (5)	38	0	0	1	0	0	0
www.sdss.org	156 (2)	106 (1)	0	0	0	0	0	0
hea-www.harvard.edu	125 (37)	42 (17)	1	0	0	1	0	26 (16)
physics.nist.gov	125 (3)	63 (2)	0	0	0	0	0	0
www.noao.edu	120 (3)	50 (2)	0	0	0	0	0	0
xmm.vilspa.esa.es	118 (35)	23 (19)	0	0	8 (1)	0	0	1 (1)
www.astro.princeton.edu	115 (31)	43 (14)	0	0	0	0	0	53 (12)
sd.usno.navy.mil	110 (27)	98 (22)	3 (3)	0	0	0	0	1 (1)

This table lists total number of links and broken links (HTTP status codes 3xx, 4xx, and 5xx) to top domains (domains with over 100 links) found within articles published in the four main astronomy journals between 1997 and 2008. The table also shows, for each domain, the portion of links to common filename extensions, as well as links that contain the tilde character.

doi:10.1371/journal.pone.0104798.t001

## How Do Astronomers Share Data?

Pepe, Goodman, Muench, Crosas, Erdmann

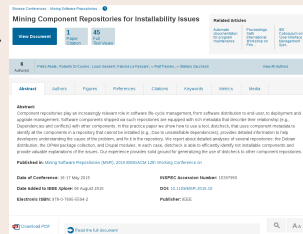
[dx.doi.org/10.1371/journal.pone.0104798](https://doi.org/10.1371/journal.pone.0104798)

PLOS August 28, 2014

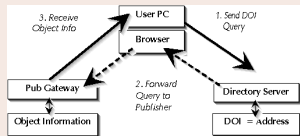
# DOI limitations

Example: doi:10.1109/MSR.2015.10

- to find what 10.1109/MSR.2015.10 is, go to a *resolver* (e.g. doi.org)
- this returns <http://ieeexplore.ieee.org/document/7180064/>
- at this URL we find ...



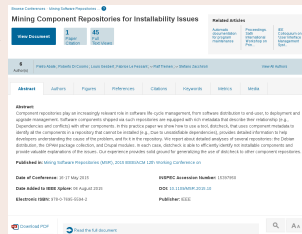
## Architecture of the DOI infrastructure



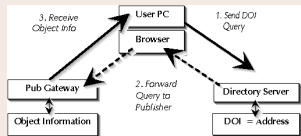
# DOI limitations

Example: doi:10.1109/MSR.2015.10

- to find what 10.1109/MSR.2015.10 is, go to a *resolver* (e.g. doi.org)
- this returns <http://ieeexplore.ieee.org/document/7180064/>
- at this URL we find ...



## Architecture of the DOI infrastructure



- DOI resolution *can change*
- content at URL *can change*
- no *intrinsic* way of noticing
- persistence based on *good will* of *multiple parties*



6 Strategy

7 Identifiers are not easy

8 Looking for the right identifiers



A *system of identifiers* is


- a set of labels (the identifiers)
- mechanisms to perform :

<i>Generation (minting)</i>	create a new label
<i>Assignment</i>	associate label to object
<i>Retrieval</i>	get object from a label

- optionally, mechanisms to perform:

<i>Verification</i>	check label and object
<i>Reverse Lookup</i>	get label from an object
<i>Description</i>	get metadata of an object

# Mechanisms offered in some systems of identifiers



Mech. / System	Handle	DOI	Ark	PURL
Generation	Yes	Yes	Yes	Yes
Assignment	Yes	Yes	Yes	Yes
Retrieval	Yes	Yes	Yes	Yes
Verification	N.A.	N.A.	N.A.	N.A.
Reverse Lookup	N.A.	N.A.	N.A.	N.A.
Description	Yes	Yes	Yes	N.A.

# Our challenges in the PID landscape

Typical properties of systems of identifiers

uniqueness, non ambiguity, persistence, abstraction (opacity)

# Our challenges in the PID landscape

## Typical properties of systems of identifiers

uniqueness, non ambiguity, persistence, abstraction (opacity)

## Key needed properties from our use cases

**gratis** identifiers are free (billions of objects)

**integrity** the associated object cannot be changed (sw dev, *reproducibility*)

**no middle man** no central authority is needed (sw dev, *reproducibility*)

# Our challenges in the PID landscape

## Typical properties of systems of identifiers

uniqueness, non ambiguity, persistence, abstraction (opacity)

## Key needed properties from our use cases

**gratis** identifiers are free (billions of objects)

**integrity** the associated object cannot be changed (sw dev, *reproducibility*)

**no middle man** no central authority is needed (sw dev, *reproducibility*)

we could not find systems with both **integrity** and **no middle man** !

# An important distinction: DIOs vs. IDOs

*The term “Digital Object Identifier” is construed as “digital identifier of an object,” rather than “identifier of a digital object”*

*Norman Paskin. 2010*



# An important distinction: DIOs vs. IDOs

*The term “Digital Object Identifier” is construed as “digital identifier of an object,” rather than “identifier of a digital object”*  
Norman Paskin. 2010

## DIO (Digital Identifier of an Object)

digital identifiers for (potentially) **non digital objects**

- epistemic complexity (manifestations, versions, locations, etc.)
- need an authority to ensure persistence and uniqueness

# An important distinction: DIOs vs. IDOs

*The term “Digital Object Identifier” is construed as “digital identifier of an object,” rather than “identifier of a digital object”*  
Norman Paskin. 2010

## DIO (Digital Identifier of an Object)

digital identifiers for (potentially) **non digital objects**

- epistemic complexity (manifestations, versions, locations, etc.)
- need an authority to ensure persistence and uniqueness

## IDO (Identifier of a Digital Object)

digital identifiers (only) for **digital objects**

- can provide both **integrity** and **no middle man**
- broadly used in modern software development (git, etc.)



# An important distinction: DIOs vs. IDOs

*The term “Digital Object Identifier” is construed as “digital identifier of an object,” rather than “identifier of a digital object”*  
Norman Paskin. 2010

## DIO (Digital Identifier of an Object)

digital identifiers for (potentially) **non digital objects**

- epistemic complexity (manifestations, versions, locations, etc.)
- need an authority to ensure persistence and uniqueness

## IDO (Identifier of a Digital Object)

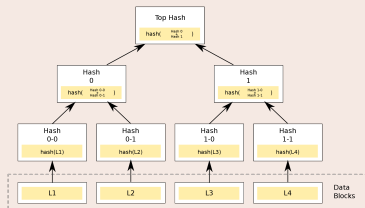
digital identifiers (only) for **digital objects**

- can provide both **integrity** and **no middle man**
- broadly used in modern software development (git, etc.)

for the core Software Heritage archive, **IDOs are enough**

# IDO in Software Development: the origins

## Merkle tree (R. C. Merkle, Crypto 1979)



Combination of

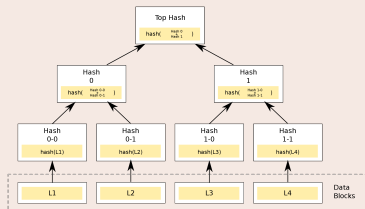
- tree
- hash function

## Classical cryptographic construction

fast, parallel signature of large data structures, built-in deduplication

# IDO in Software Development: the origins

## Merkle tree (R. C. Merkle, Crypto 1979)



Combination of

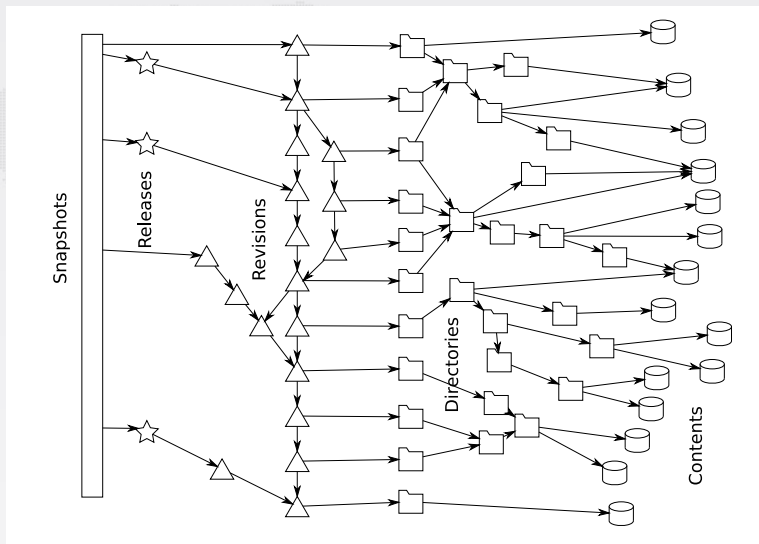
- tree
- hash function

## Classical cryptographic construction

fast, parallel signature of large data structures, built-in deduplication

- satisfies all three criteria: **gratis, integrity, no middle man!**
- widely used in industry (e.g., Git, nix, blockchains, IPFS, ...)

# IDO in Software Heritage: 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.

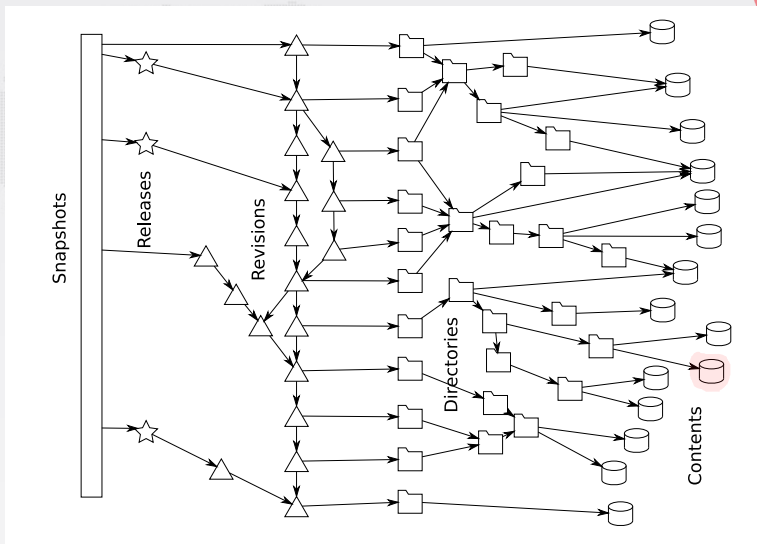
The licenses for most software and other practical works are designed  
to take away your freedom to share and change the works. By contrast,  
the GNU General Public License is intended to guarantee your freedom to  
share and change all versions of a program—to make sure it remains free  
software 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 others from denying you  
these rights by patenting or otherwise restricting the use of the software.

sha1: 8624bcdade55baeef...  
sha256: 8ceb4b9ee5aded...  
sha1\_git: **94a9ed024d385...**  
length: 35147

# IDO in Software Heritage: a worked example



# IDO in Software Heritage: a worked example

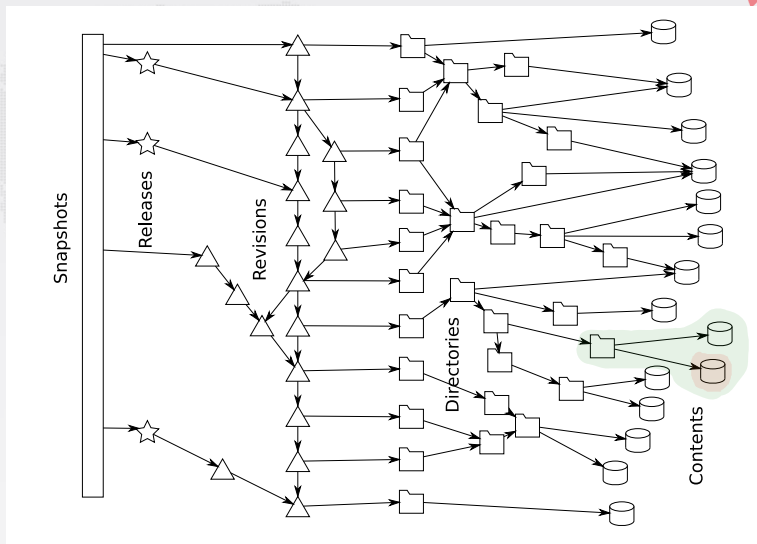


## Directories

```
100644 blob c5baade4c44766042186ef858c0fd63d587ebf09 .gitignore
100644 blob 2d0a34af6f52cf3cf6b0c2f7bd0648fbd255e77f AUTHORS
100644 blob 94a9ed024d3859793618152ea559a168bbcb5e2 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 e1e10ecf948af0b93adb0372afcb89f12e92618a bin
040000 tree 83e56d0beaf7793c77a45a345c80fcb8af503013 debian
040000 tree a34c9c4ba213f0cedc67f9816348d27955577af5 docs
100644 blob f2a6d32c6135aa7287bbd76167b01df2ae4f1539 requirements.txt
100755 blob eee147c36caflbbcd2d820da8dc026cb5b68180bc setup.py
040000 tree 224bb4c1f4c67fca1d160bffdd2d06094e7e1abf3 sql
040000 tree 8631c9cd77bbe993168107ab5baf51f40c6300be swl
040000 tree 8fb905b56ba8ed692f1209b2773b474c6c1d66c1 utils
```

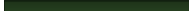
id: 515f00d44e92c65322aaa9bf3fa097c00ddb9c7d

# IDO in Software Heritage: a worked example





## Revisions

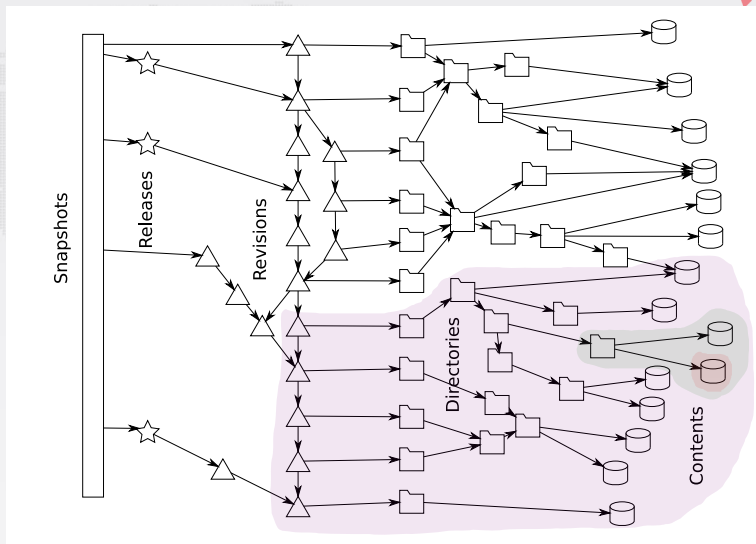
Details	Changes	Files
<p>SHA: 963634dca6ba5dc37e3ee426ba091092c267f9f6</p> <p>Author: <a href="mailto:nicolas@dandrimont.eu">Nicolas Dandrimont &lt;nicolas@dandrimont.eu&gt;</a> (Thu Sep 1 14:26:13 2016)</p> <p>Committer: <a href="mailto:nicolas@dandrimont.eu">Nicolas Dandrimont &lt;nicolas@dandrimont.eu&gt;</a> (Thu Sep 1 14:26:13 2016)</p> <p>Subject: <b>provenance.tasks: add the revision -&gt; origin cache task</b></p> <p>Parent: <a href="#">fc3a8b59ca1df424d860f2c29ab07fee4dc35d10</a> : test...storage: properly pipeline origin and cont...</p> <p>provenance.tasks: add the revision -&gt; origin cache task</p> <p><a href="#">sw/h/storage/provenance/tasks.py</a>  77</p>		

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](#)

# IDO in Software Heritage: a worked example



## 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

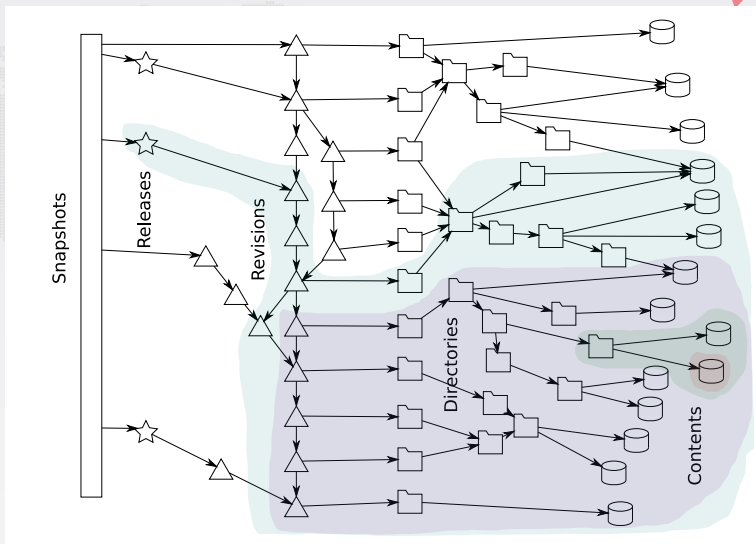
- Add new metadata column to origin\_visit
- Update swh-add-directory script for updated API

-----BEGIN PGP SIGNATURE-----

iQIzBAABCAAdBQJXvZTNFhXuaWNvbGFzQGRhbmRyaW1vbnQuZXUACgkQ7AWLMo2+  
neqorw/aaq6SOb5DjjzEa+kWN3rXgV5+1K1vEVh1wNKAwx8eKJ7aX2kEiLDtt7uf  
ahpZ6pz3q8nqs6aC1+YrxBfcih3l2YtrdZeWXXWqr8xWNMaFoYDb8qaphwh8AD5t2  
ICBli2ujtXuCrDt93eKkPwvzXg+h80sMWy35Dr6jWZ7K4Mu/PgGlyIHPY55yo  
IGEndWno7VfH1Vm6t1n5qB7l5mXRaqa+becqddubTZ2xij+jpLUqC8cyqN3hm/fL  
qsJ2mu8kyz3t8tG/H1/pV+I5OwBlNpO5STH0tuojoEVgPK/dH5P79QuHDHZFkCao  
klj6kAWyU80Mxb+nKVjjeLbrR3+yWBFJ3Qp5a1/V8oOTn6E1dAlcNMpEaKCoKtMt  
d/gMRax11l/g0EDfnsW67G6sDwKPKPHngfVLQ3nV3GaQQTnu1RpMz006H9/tAwzC  
Gg/K1PdHT4hz0I46wYPZye0U2VXGFu6vVU9vFQ4ZR/Wjn+0zMzdcRdrJlSUOMn  
RpTTfUsbXUeXHGOpkgXhSYTnvp1gdPc76USTsK0aGe84A2m1lk0mGrwXCvFPqIYo  
nhhibBSHBNMoqyF6yTSOpUbyK70tpYRRUGKwDeRK0wKSkxWKUZGtKzy6jYqJjo29  
gulwqZQif5qWQCB0ontAL2+HvPfFaVyckMejUhg62cP/+EHlvUk=  
=kOxP  
-----END PGP SIGNATURE-----

id: 85083a5cc14a441c89dea73f5bdf67c3f9c6afdb

# IDO in Software Heritage: a worked example



git show-refs

## Snapshots

```
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 3147c3d31ec46cf6492f881e908b1237ebdff72c7 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

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

swh:1:**cnt**:94a9ed024d3859793618152ea559a168bbcbb5e2

full text of the GPL3 license

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

swh:1:**cnt**:94a9ed024d3859793618152ea559a168bbcbb5e2

full text of the GPL3 license

swh:1:**dir**:d198bc9d7a6bcf6db04f476d29314f157507d505

Darktable source code

# The Software Heritage IDO 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



# The Software Heritage IDO 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

# The Software Heritage IDO 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)

# The Software Heritage IDO 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](http://archive.softwareheritage.org) and [n2t.org](http://n2t.org)