

Software Heritage

A revolutionary infrastructure for Open Science

Roberto Di Cosmo

July 10th, 2019
Singapore



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion



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

2005 *Open Access* debate

2007 *Free Software Thematic Group*
150 members 40 projects 200Me

2015 *Software Heritage* at INRIA

2018 *National Committee for Open Science*, France

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion



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 also 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.”

~ 50 years, a lightning fast growth

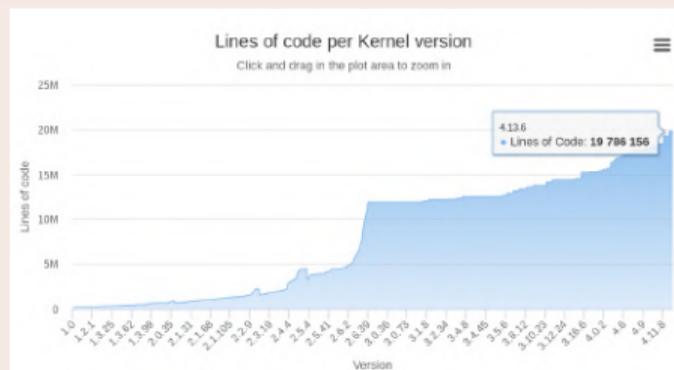
Apollo 11 Guidance Computer (~60.000 lines), 1969



"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!

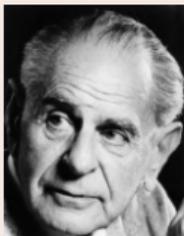
The experimental method



- 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

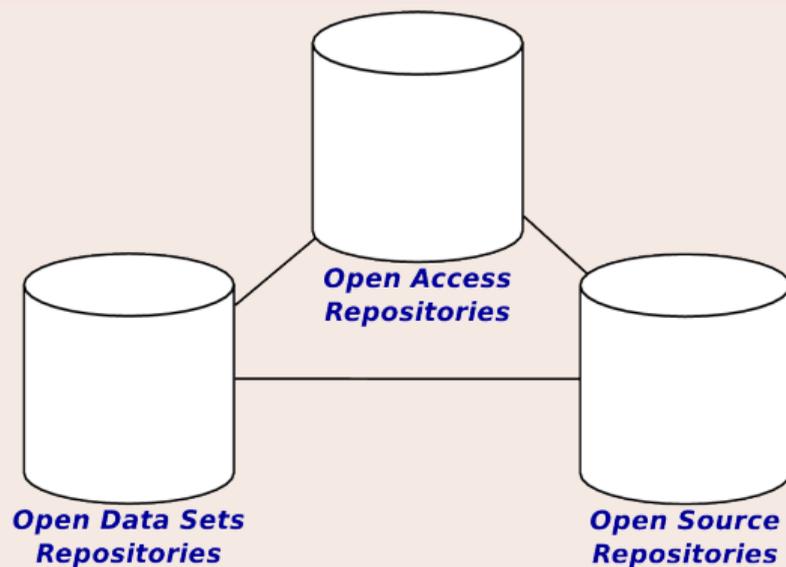
Remark

The first two items are already widely discussed!

... what about *software*?

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 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion



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 infrastructure plans
- in the copyright reform (EU)
- in the scholarly works

Lack of established guidance

- choose a license
- make source code available
- relate to industry best practices
- *cite* a software project

No catalog, no archive, no references: we are at a turning point

Looking at the past

- a lot of old software misplaced, lost, or behind barriers, but...
- most founding fathers are still here, and willing to share
- **urgent** to collect their knowledge

Only a few years left.

Looking at the future

- software development and use skyrockets: more programmers, and more code!
- **essential** to provide a **universal** platform for all the future software source code

Every year that goes by makes the problem worse.

it is **urgent** to take action!

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage**
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion





Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Our mission

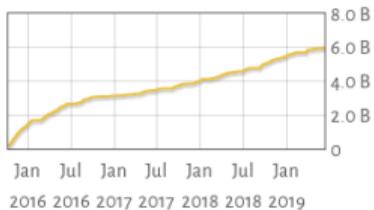
Collect, **preserve** and **share** the *source code* of *all the software* that is available

Past, present and future

Preserving the past, *enhancing* the present, *preparing* the future

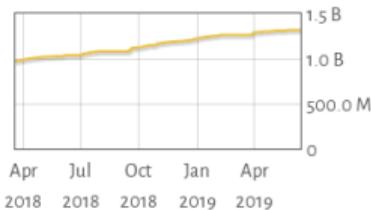
Source files

6,006,503,960



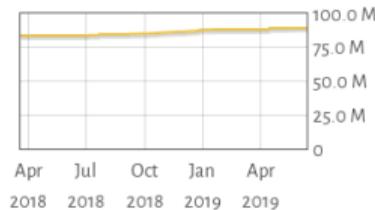
Commits

1,326,776,432



Projects

89,301,694



GitHub

debian



GitLab

npm

Google code

GITORIOUS

GNU

HAL
archives-ouvertes.fr

Inria
inventeurs du monde numérique

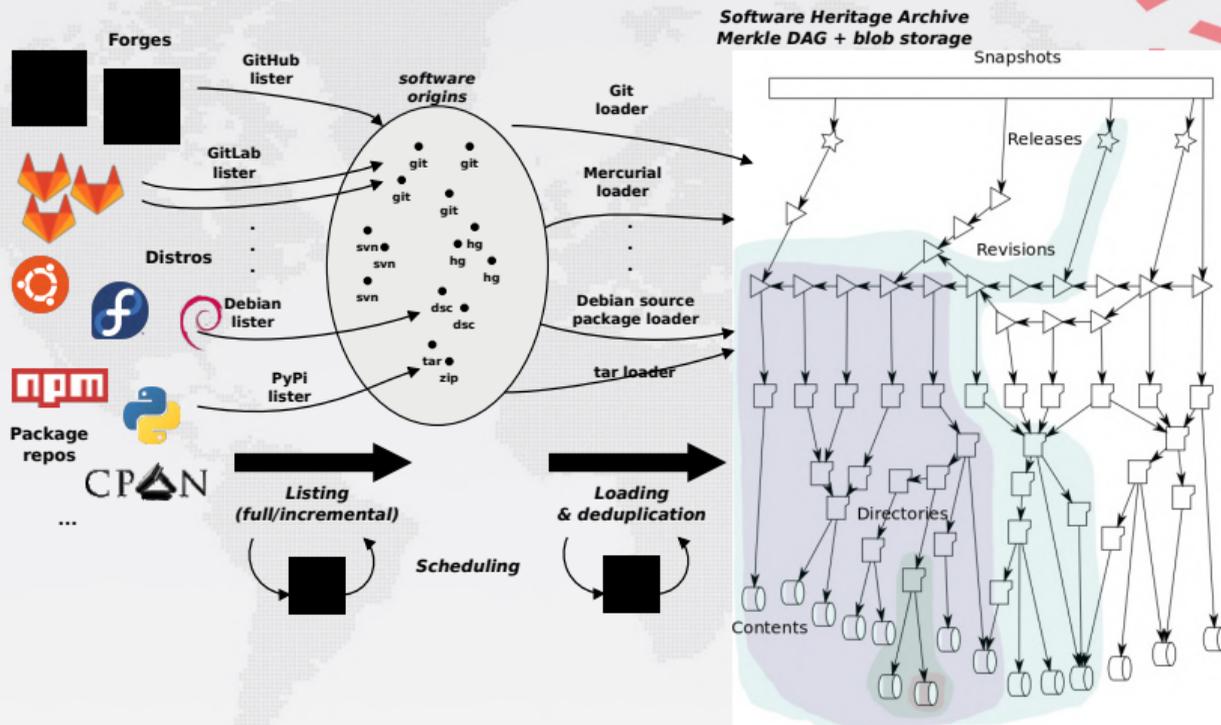
python
Package Index

- ~400 TB (uncompressed) blobs, ~20 B nodes, ~280 B edges
- The *richest* public source code archive, ... and growing daily!

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure**
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion



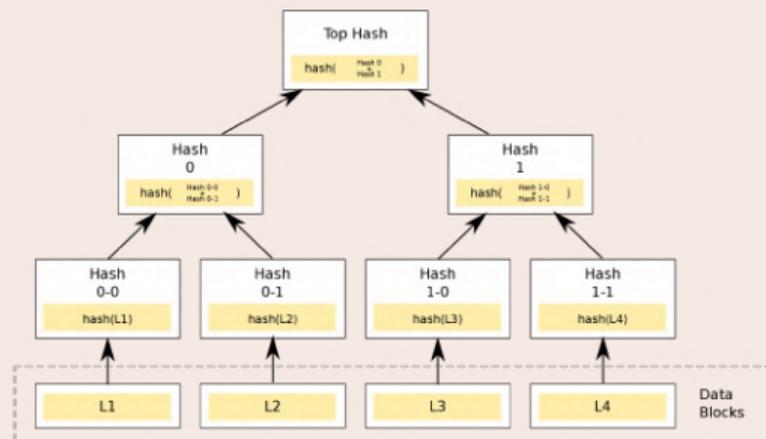
Automation, and storage



Full development history **permanently** archived in a **uniform data model**.

Much more than an archive!

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



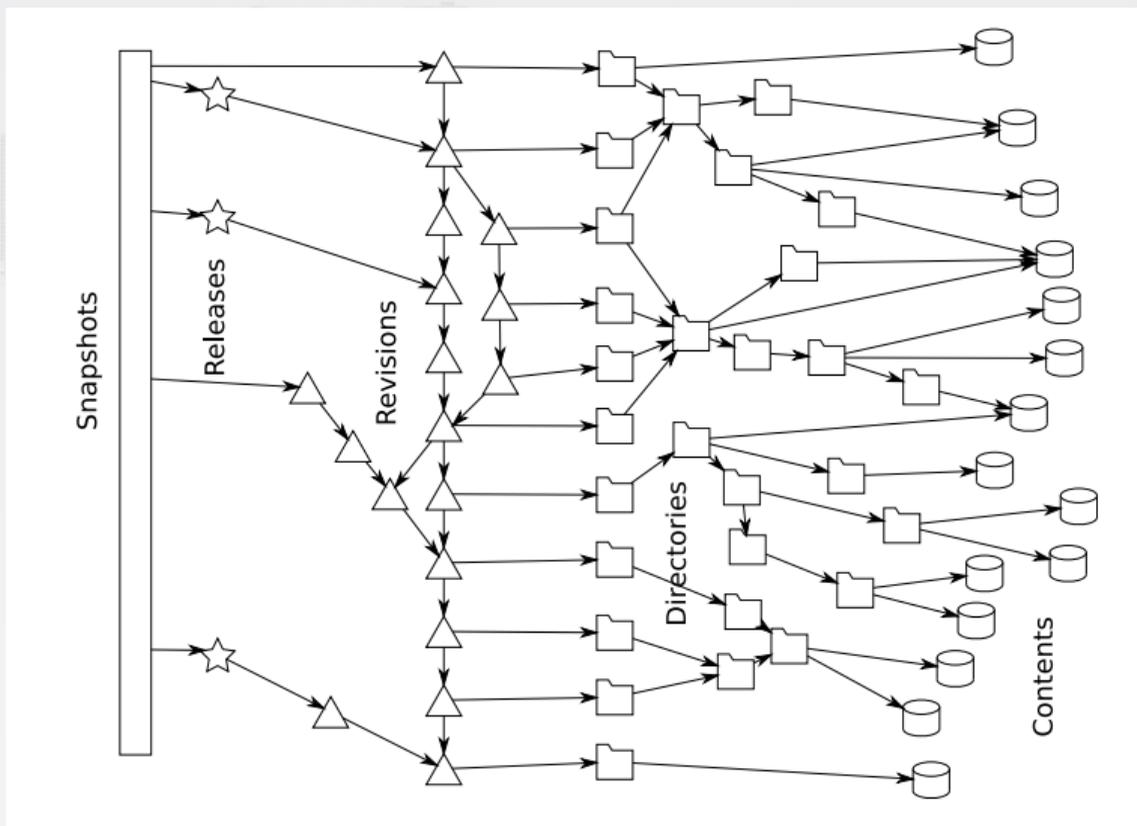
Combination of

- tree
- hash function

Classical cryptographic construction

- fast, parallel signature of large data structures
- widely used (e.g., Git, blockchains, IPFS, ...)
- **built-in deduplication**

The archive in pictures



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.

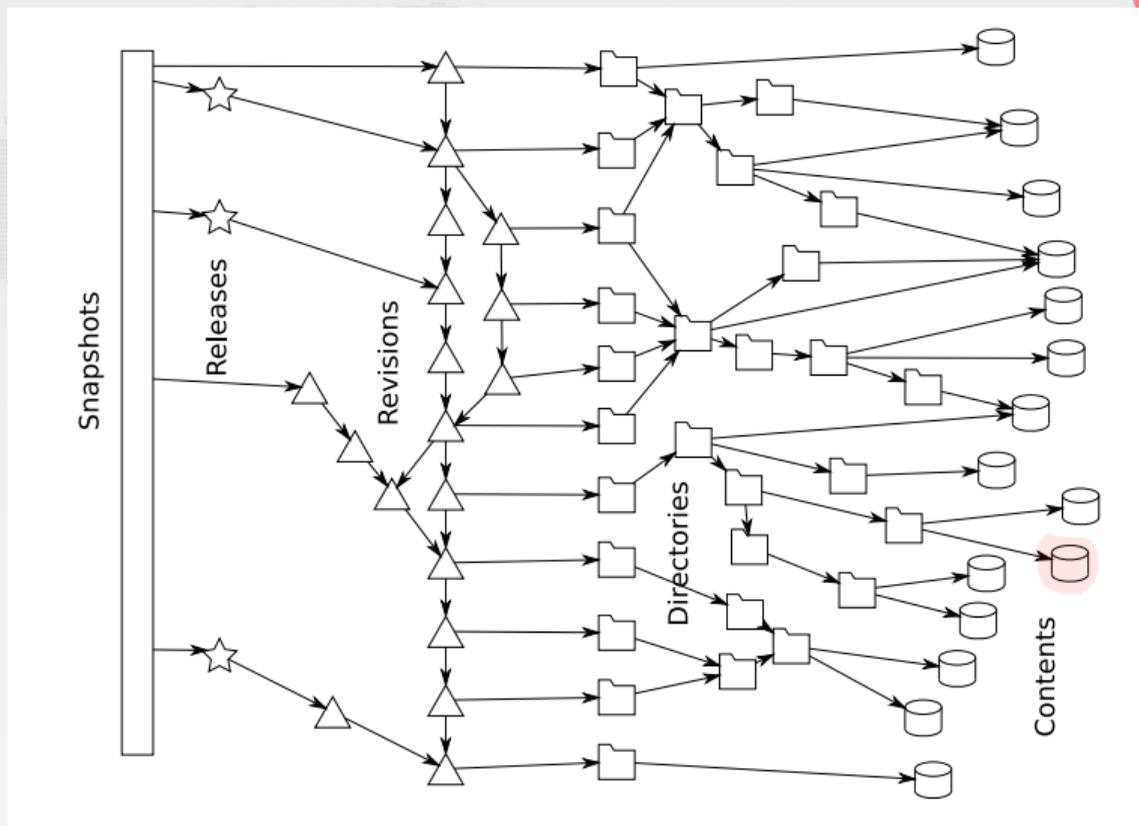
The licenses for most software and other practical works are designed
to take away your freedom to share and change the works. By contrast,
the GNU General Public License is intended to guarantee your freedom to
share and change all versions of a program--to make sure it remains free
software for all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

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

The archive in pictures



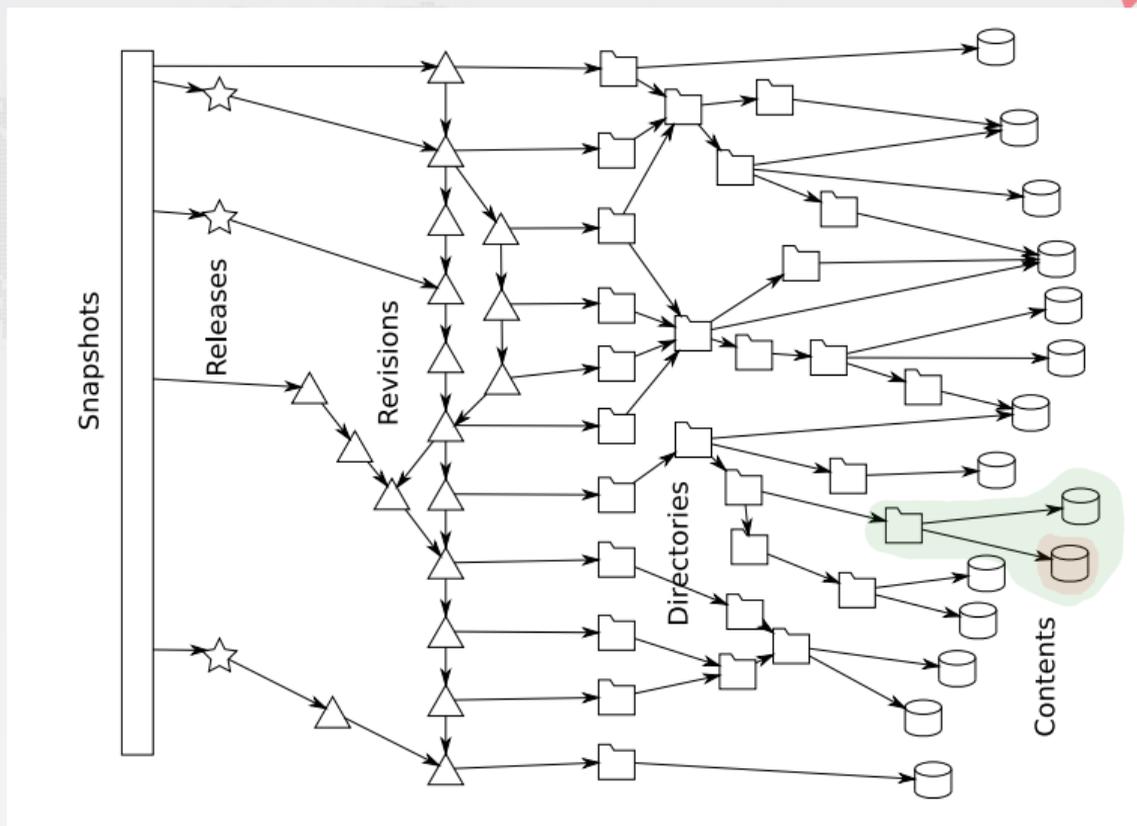


Directories

```
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 e1e10ecfe948af0b93adb0372afc89f12e92618a 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

The archive in pictures



Revisions

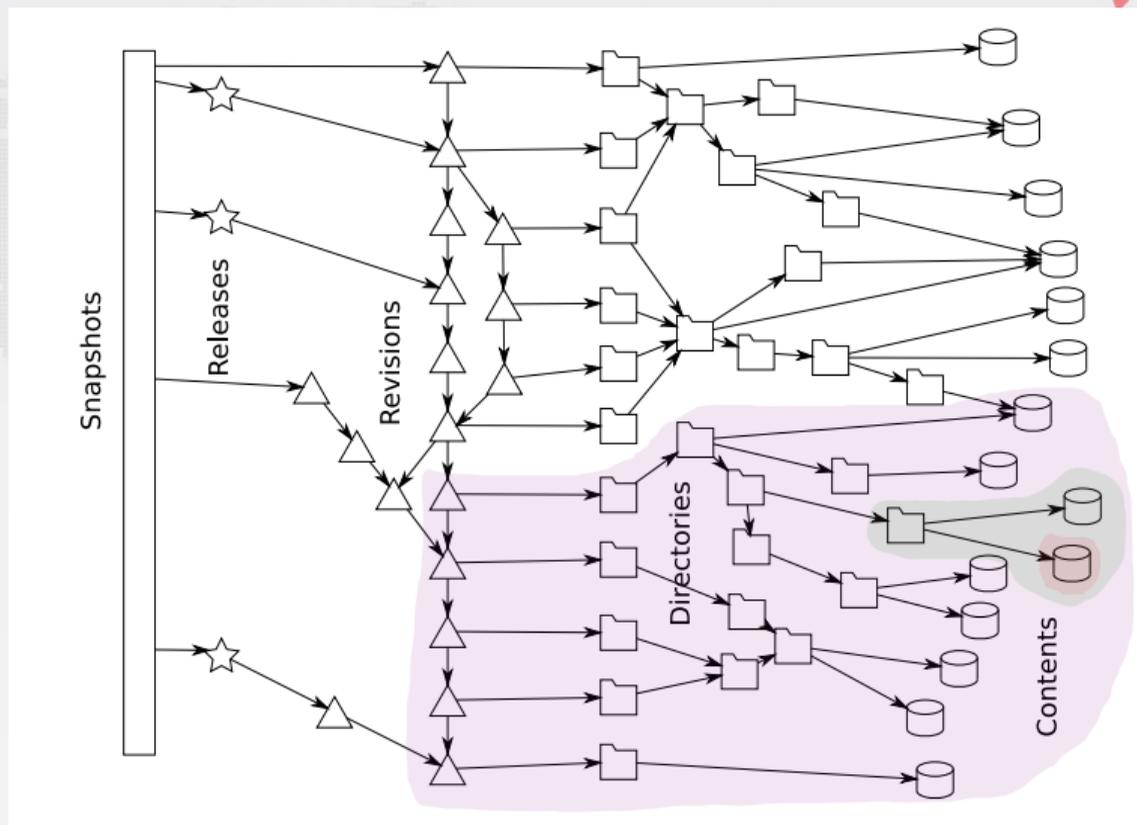
Details	Changes	Files
SHA: 963634dca6ba5dc37e3ee426ba091092c267f9f6		
Author: Nicolas Dandrimont <nicolas@dandrimont.eu> (Thu Sep 1 14:26:13 2016)		
Committer: Nicolas Dandrimont <nicolas@dandrimont.eu> (Thu Sep 1 14:26:13 2016)		
Subject: provenance.tasks: add the revision -> origin cache task		
Parent: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

The archive in pictures



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 c0c9f16b1e134f593e7567570a1761b156e6b1d

object c0c9f16b1e134f593e7567570a1761b156e6b1d
type commit
tag v0.0.51
tagger Nicolas Dandrimont <nicolas@dandrimont.eu> 1472042163 +0200

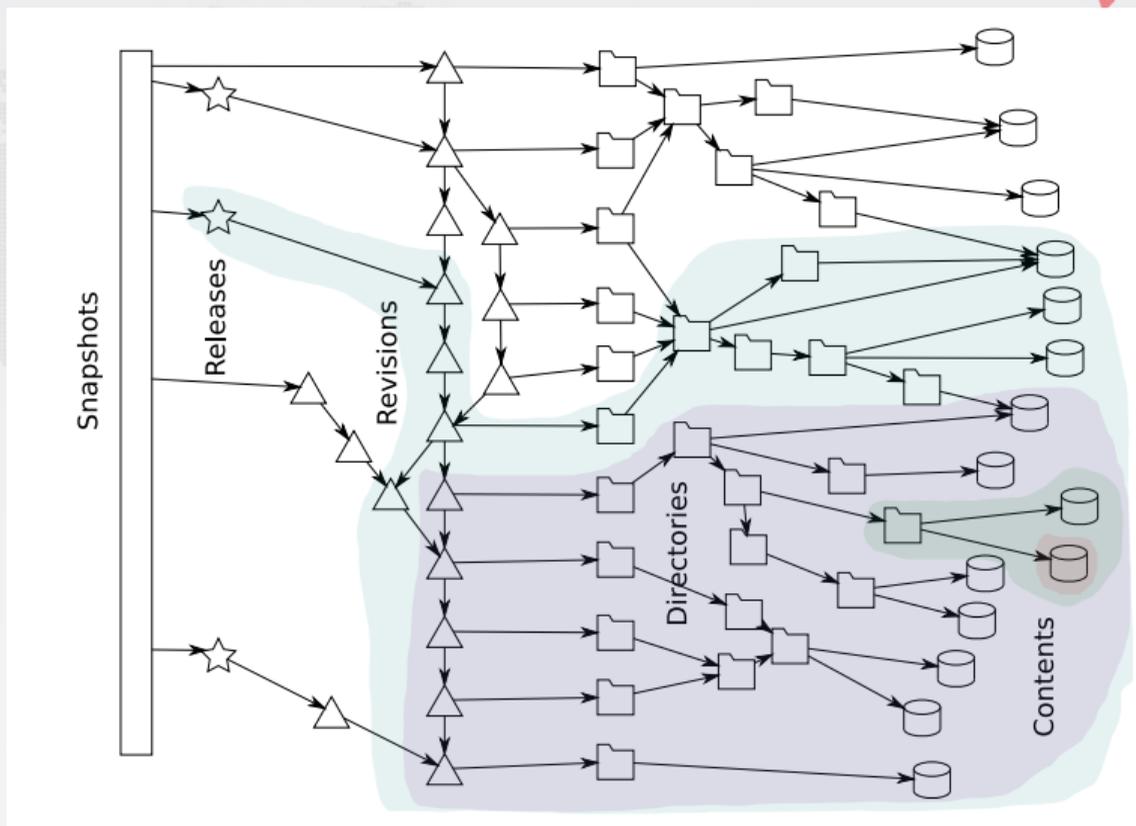
Release sw.h.storage v0.0.51

- 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+YrxBfcih3L2YtrdZeWwXqr8xWNMaEoYDb8aaphwh8AD5t2
ICBii2ujtXuCrDt93eKkPwvzZxg+h80sMWy35Dr6jWZ7K4MuPGglyIHPY55yo
IGEndWno7Vfh1Vm6t1n5qB7I5mXRaqA+becqddubTZ2xij+jpLlUqC8cyqN3hm/fl
qsj2mu8kyz3t8tG/H1/pV+150wBlNp05TH0tujojEVgPK/dH5P79QuHDHZFkCao
klj6kAWyU80Mxb+nKVjeLbrR3+yWBFj3Qp5a1/V8oOTh6E1dAlcNmPkaKCoKtMt
d/gMRax11I/g0EDfnsW67G6sDwKPKPHngfVLQ3nV3GaQQTnu1RpMz006H9/tAwzC
Gg/K1PdHT4hz0iI46wYPZje0U2VXGFu6vVU9vFQ4ZR/Wjn+0zZdcRdrjSUOMn
RpTTRfUbsXUeXHGOpkgXhSYTnvp1gdPc76U5TsK0aGe84AZm1lk0mGrwXCvFPqYo
nhhibB5HBNMoqyF6yTSOpUbyK70tpYRRUGKwDeRK0wKSxkWKUZGtKzy6jYqJ0z9
gulwgZQif5qWQCB0oontAL2+HvPfaVyckMejUhg62cP/+EHlvUk=
=kOxP
-----END PGP SIGNATURE-----

id: 85083a5cc14a441c89dea73f5bdf67c3f9c6afdb

The archive in pictures



Snapshots

git show-refs

```
commit 08ffeb25770109525eb3ce21691466c53a1d9158 refs/heads/atime
commit ba5443a24e3f9fe323a46c292cec4fcbe61c67eb refs/heads/directory-listing-arrays
commit d69e0dbf892383ff6589b27fbc1c05d27238d9c5 refs/heads/foo
commit cf7ff9eea0eb22f758946908f5a8019f67de468e08 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 228a2f1650dd1222e556559462e1e06fc4993d9 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

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 and n2t.org

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time**
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion



Demo time: let's highlight some features...

A "wayback machine" for software source code

- <http://archive.softwareheritage.org/browse>

Identification and sharing of billions of software artifacts

- <http://bit.ly/swhpids> for persistent identifiers

Depositing research software

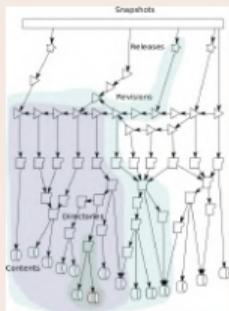
- <http://bit.ly/swdepositblog>

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure**
- 8 Building for the long term
- 9 Conclusion



A revolutionary infrastructure for industry

The *graph* of Software Development



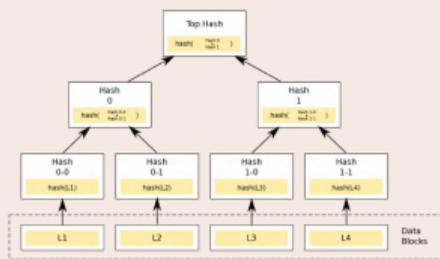
All of the software development in a **single graph**!

- **lookup** by content hash
- **wayback machine** for software development
 - <http://archive.softwareheritage.org/>
- ... and much more

The *blockchain* of Software Development

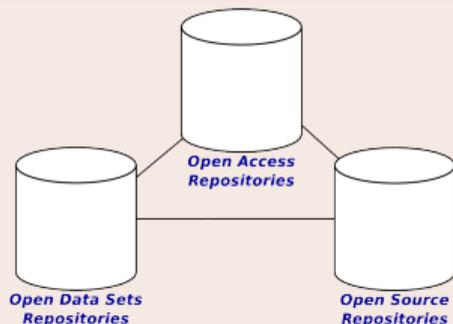
All of a software development... in a single **Merkle** graph!

Widely used crypto (e.g., Git, blockchains, IPFS, ...)



- built-in **deduplication**
- intrinsic, **unforgeable identifiers** at all levels
- simplifies **traceability** (licensing, supply chain management)

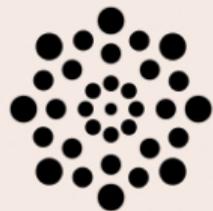
A pillar of Open Science



The *reference archive* of Research Software for **Open Science**

- **curated deposit** of research software
 - in collaboration with **HAL**, **CCSD** and **Inria IES**
 - now open *to all researchers!*
- **intrinsic** identifiers for **reproducibility**

Reference platform for *Big Code*



- unique **observatory** of all software development
- **big data, machine learning** paradise: classification, trends, coding patterns, code completion...

Save code now

two simple steps to get your repository archived:

- prepare your source code
- go to `https://archive.softwareheritage.org/browse/origin/save/`

that's it!

Reference code: precisely, and forever!

three simple steps:

- find your code in `https://archive.softwareheritage.org`
- select your code fragment (optional)
- get the link from the red permalink tab, and use it!

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term**
- 9 Conclusion



Cultural Heritage



Industry



Research



Education



Software Heritage

Technology

- transparency and FOSS
- replicas all the way down

Content

- intrinsic identifiers
- facts and provenance

Organization

- non-profit
- **mirror network**

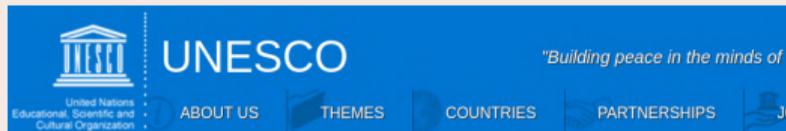
April 3rd 2017, Unesco Inria agreement

Inria
INVENTEURS DU MONDE NUMÉRIQUE



Roberto Di Cosmo

November 2018, Unesco Inria expert call



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



Growing Support

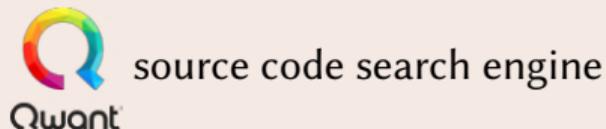
Sharing the vision



Donors, members, sponsors



Research collaboration



Global network

- first independent mirror
- increased reliability

Towards Reproducible Open Science

archive research software in SWH

reference it using *intrinsic identifiers*

build on top of SWH, *do not try to rebuild SWH!*

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**
- let's prevent a useless diaspora

Scientific challenges in building Software Heritage

graph queries, efficient storage, distributed archival, classification, search, ...

Using Software Heritage for research

the *Software Heritage graph dataset* is now available!

- AWS: <https://registry.opendata.aws/software-heritage/>
- guidelines: <https://upsilon.cc/~zack/research/publications/msr-2019-swh.pdf>

- 1 Introductions
- 2 Software source code: a pillar of Open Science
- 3 An inconvenient truth
- 4 Software Heritage
- 5 Under the hood: architecture and data structure
- 6 Demo time
- 7 A revolutionary infrastructure
- 8 Building for the long term
- 9 Conclusion



www.softwareheritage.org

@swheritage

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



Jean-François Abramatic, Roberto Di Cosmo, Stefano Zacchioli

Building the Universal Archive of Source Code

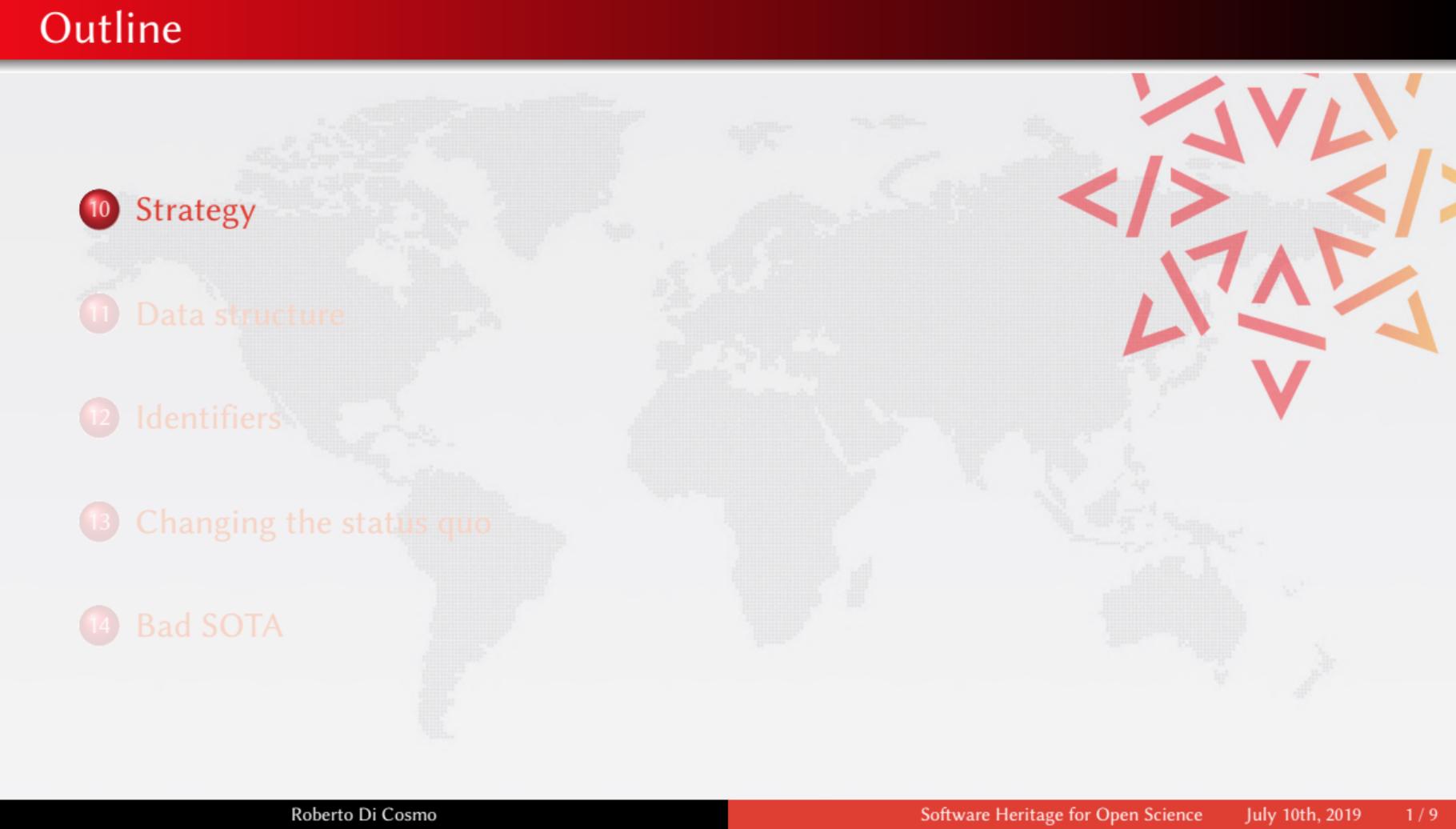
Communication of the ACM, October 2018

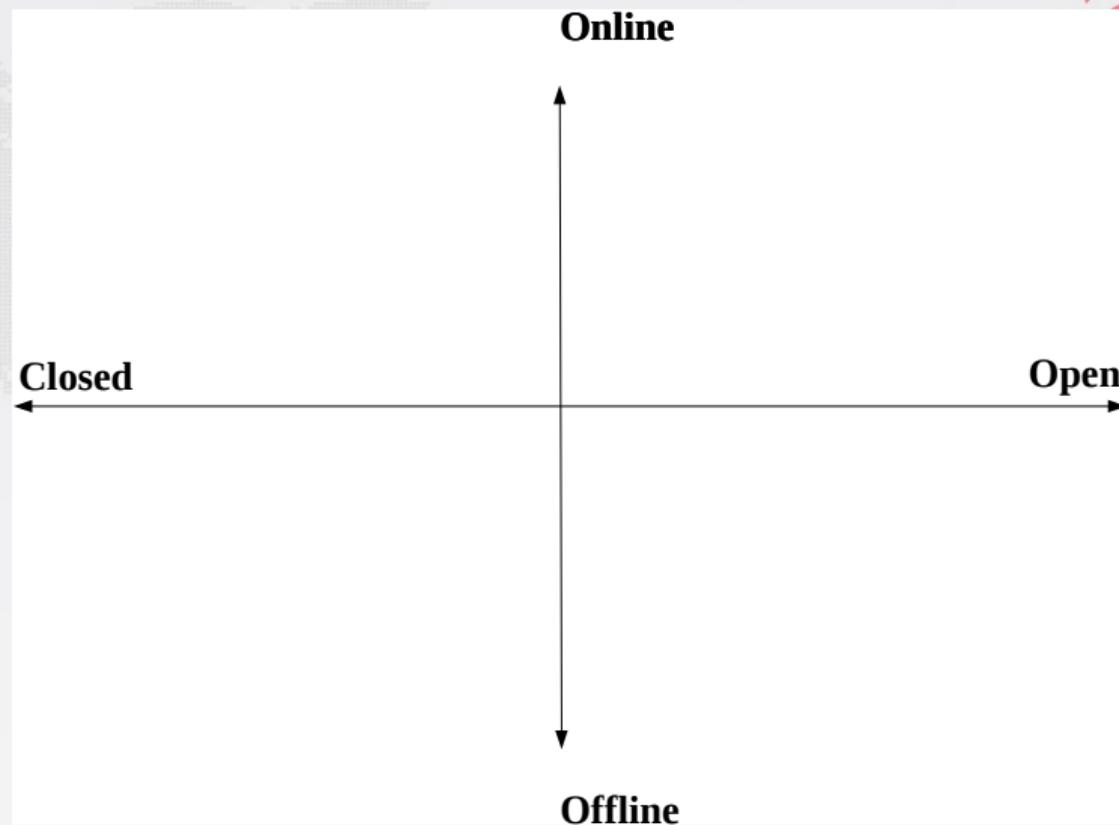


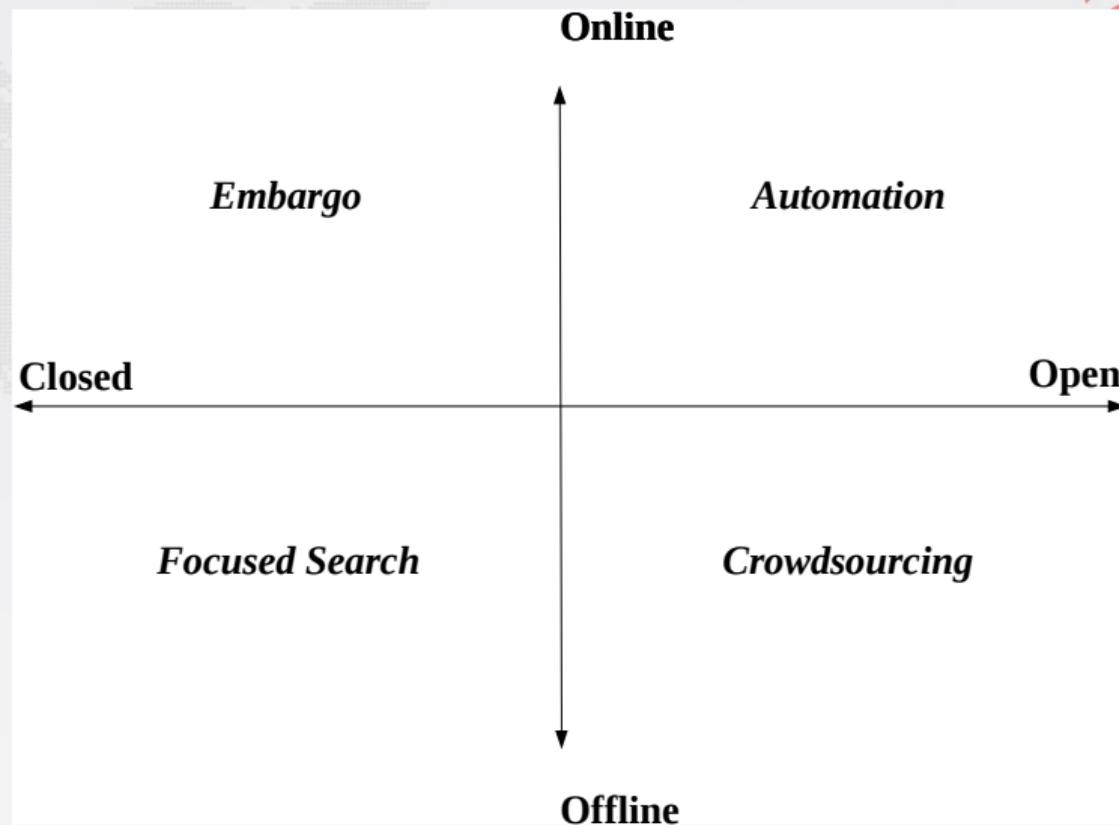
Roberto Di Cosmo, Morane Gruenpeter, Stefano Zacchioli

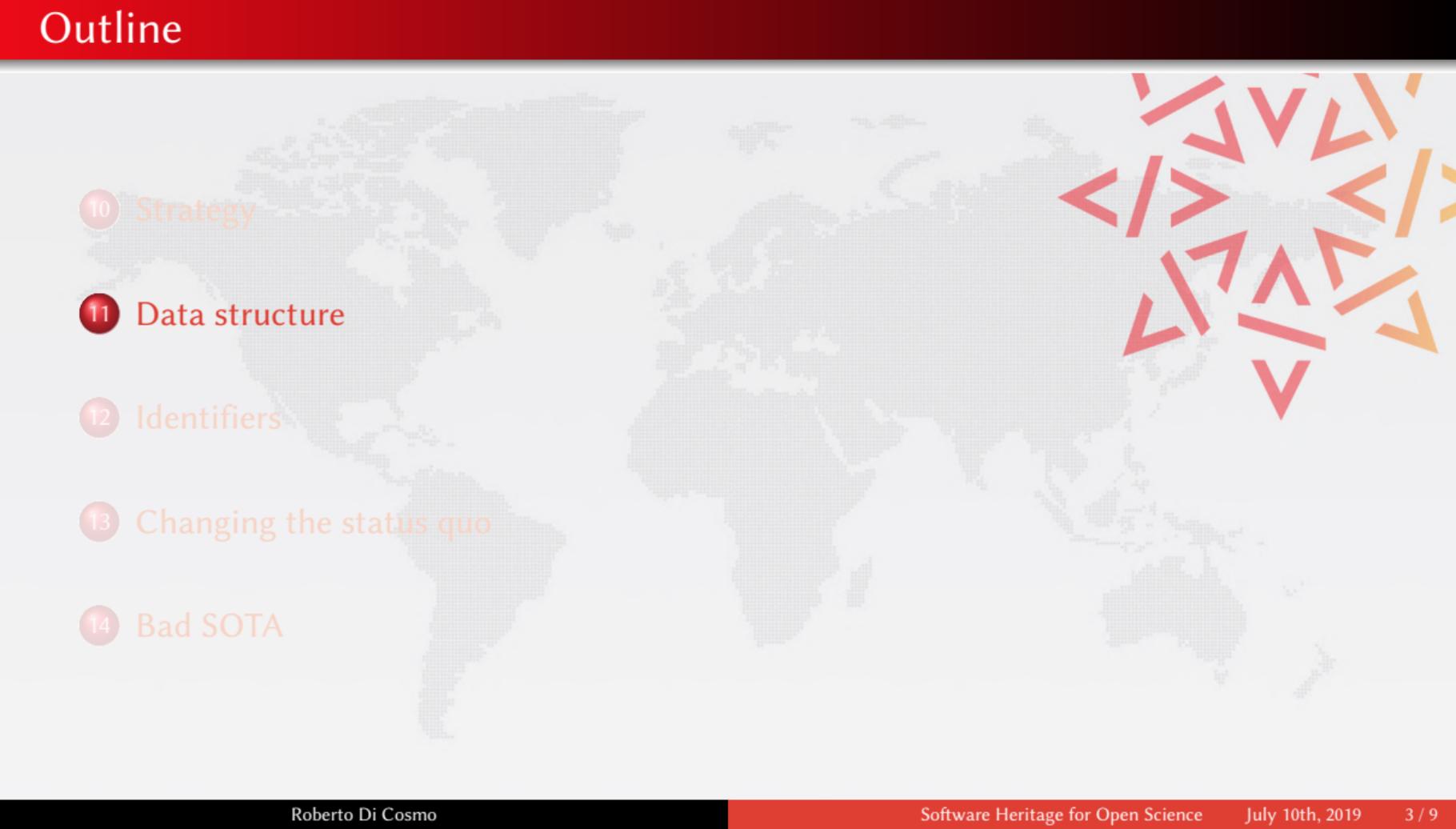
Identifiers for Digital Objects: the Case of Software Source Code Preservation

iPRES 2018: Intl. Conf. on Digital Preservation

- 
- 10 Strategy
 - 11 Data structure
 - 12 Identifiers
 - 13 Changing the status quo
 - 14 Bad SOTA





- 
- 10 Strategy
 - 11 Data structure
 - 12 Identifiers
 - 13 Changing the status quo
 - 14 Bad SOTA

A bird's eye view

```
origin https://forge.softwareheritage.org/source/helloworld.git
visit 1
timestamp Fri Feb 9 12:38:45 2018 +0100
snapshot 0861db5e...
```

```
origin https://forge.softwareheritage.org/...
visit 2
timestamp Fri Feb 9 13:29:00 2018 +0100
snapshot 0861db5e...
```

```
origin https://forge.softwareheritage.org/...
visit 3
timestamp Fri Feb 9 15:52:50 2018 +0100
snapshot 510aa88b...
```

```
<<Snapshot>>
510aa88b...
+branches
HEAD
refs/heads/master
refs/heads/doc
refs/tags/1.0
```

```
<<Snapshot>>
0861db5e...
+branches
HEAD
refs/heads/master
refs/tags/1.0
```

```
<<Release>>
edf82f21...
+author = "Foo Bar <foo@...>"
+name = "1.0"
+message = "1.0 release"
+timestamp = Thu Feb 8 15:51:00 2018 +0100
+target
```

```
<<Revision>>
1a99a56b...
+author = "Foo Bar <foo@...>"
+message = "Merge branch 'doc'"
+timestamp = Fri Feb 9 15:44:45 2018 +0100
+directory: Directory
+parents: Revision list
```

```
<<Revision>>
3d515253...
+author = "Foo Bar <foo@...>"
+message = "README: add homepage link"
+timestamp = Fri Feb 9 15:44:30 2018 +0100
+directory: Directory
+parents: Revision list
```

```
<<Revision>>
c7640e8d...
+author = "Foo Bar <foo@...>"
+message = "move source code to src/\n..."
+timestamp = Thu Feb 8 15:26:00 2018 +0100
+directory: Directory
+parents: Revision list
```

```
<<Revision>>
43ef7dcd...
+author = "Foo Bar <foo@...>"
+message = "add licensing information and README"
+timestamp = Thu Feb 8 10:54:09 2018 +0100
+directory: Directory
+parents: Revision list
```

```
<<Revision>>
a3ee21ad...
+author = "Foo Bar <foo@...>"
+message = "add buildtoolchain ..."
+timestamp = Thu Feb 8 10:49:29 2018 +0100
+directory: Directory
+parents: Revision list
```

```
<<Revision>>
1886826f...
+author = "Foo Bar <foo@...>"
+message = "implement a trivial ..."
+timestamp = Thu Feb 8 10:44:35 2018 +0100
+directory: Directory
+parents: Revision list = None
```

```
<<Directory>>
ded70c63...
+entries
"COPYING"
"Makefile"
"README.md"
"src"
```

```
<<Directory>>
ded70c63...
+entries
".gitignore"
"COPYING"
"Makefile"
"README.md"
"hello.c"
```

```
<<Directory>>
45f0c078...
+entries
"COPYING"
"Makefile"
"README.md"
"src"
```

```
<<Directory>>
fab8c908...
+entries
".gitignore"
"COPYING"
"Makefile"
"README.md"
"hello.c"
```

```
<<Directory>>
b94a90cd...
+entries
".gitignore"
"Makefile"
"hello.c"
```

```
<<Directory>>
6ca2e444...
+entries
"hello.c"
```

```
<<Content>>
4ec6b1e9...
+data = "...For more info..."
```

```
<<Content>>
a8becc46...
+data = "SRC_DIR = _"
```

```
<<Content>>
a1af0006...
+data = "...Yet another..."
```

```
<<Content>>
94a9e0d2...
+data = "...GNU GENERAL..."
```

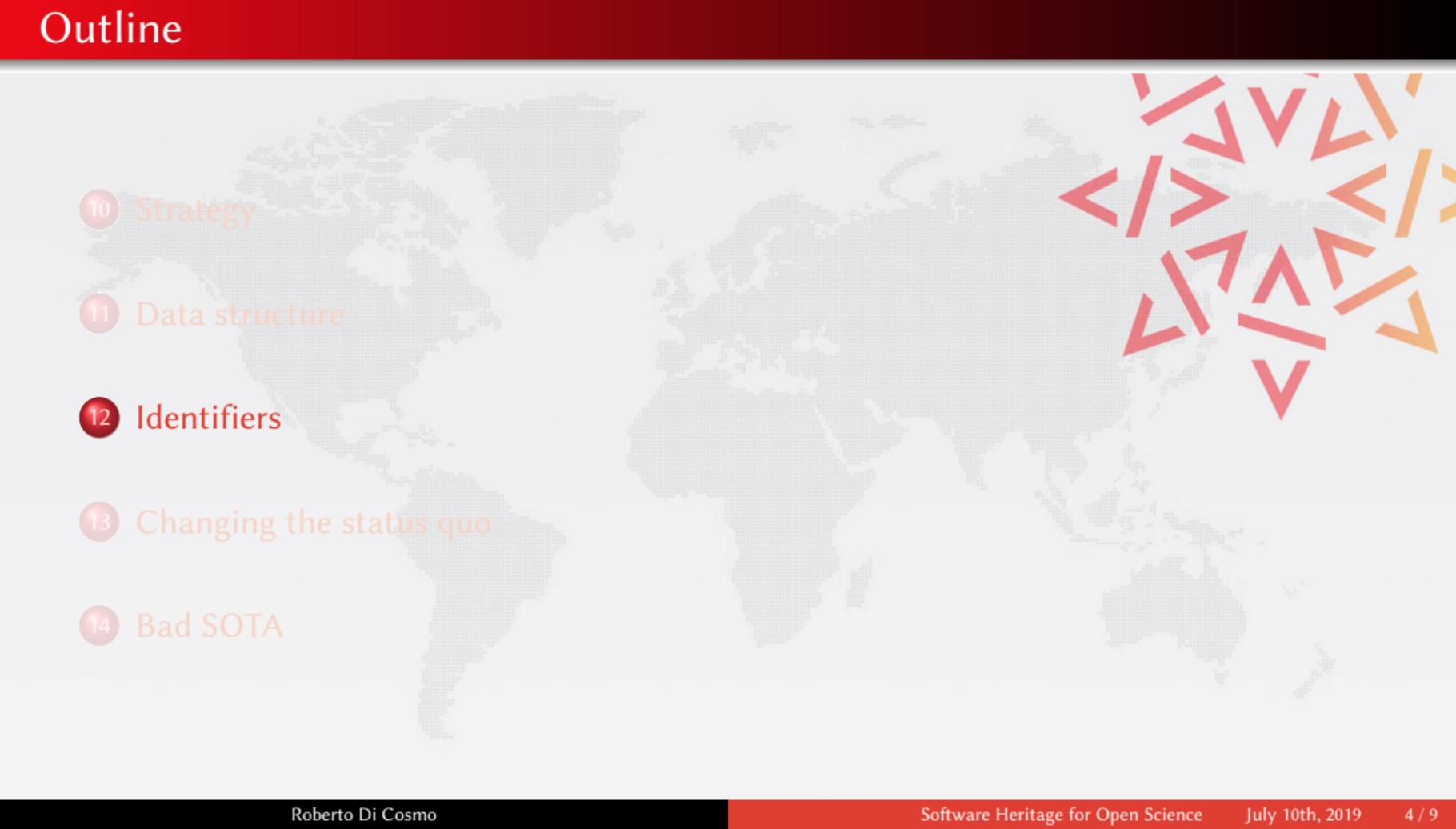
```
<<Content>>
59b32b2f...
+data = "...\nhello"
```

```
<<Content>>
225ae01b...
+data = "all: hello\n\n..."
```

```
<<Content>>
c839dea9...
+data = "#include _"
```

Archive content after visits 1, 2 and 3

Archive content after visits 1 and 2

- 
- 10 Strategy
 - 11 Data structure
 - 12 Identifiers**
 - 13 Changing the status quo
 - 14 Bad SOTA

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

DIO (Digital Identifier of an Object) 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) identifiers (only) for digital objects

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

IDOs and DIOs adress different needs

- for the core Software Heritage **IDOs are enough**
- we **must not** use DIOs for reproducibility

Limitations of DIOs

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 ...

Mining Component Repositories for Instability Issues

1 Type 45 Total Items

Abstract Authors Figures References Citations Requests Metrics Links

Abstract

Component repositories play an increasingly relevant role in software lifecycle management, from software distribution and asset to development and updates management. Software components stored in such repositories are organized with rich metadata that describe their relationship to other components and combined with other components. In this practice paper we show how to use a novel, effective, but more complex metadata to identify all the components in a repository that cannot be installed (e.g., missing unresolvable dependencies), provide detailed information to help developers understanding the cause of the problem, specific to the repository, this report about detailed analysis of several repositories. The authors distribute the other package information, and the paper includes, in each case, also how to use it effectively to identify the available components and provide reliable information of the assets. Our experience provides background for generating the case of distribution to other component repositories.

Published by Mining Software Repositories (MSR), 2015 IEEE/ACM 28th Working Conference on

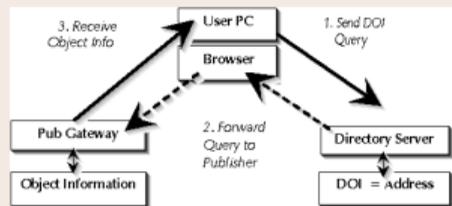
Date of Conference: 11-17 May 2015 IEEE Xplore Article Number: 7327100

Date Added to IEEE Xplore: 10 August 2015 DOI: 10.1109/MSR.2015.10

Electronic ISBN: 978-1-7651-0294-2 Publisher: IEEE

Download PDF Find this full document

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*

- 
- 10 Strategy
 - 11 Data structure
 - 12 Identifiers
 - 13 Changing the status quo
 - 14 Bad SOTA

Revisiting the "Software Citation Principles"

Recommendation for identifiers

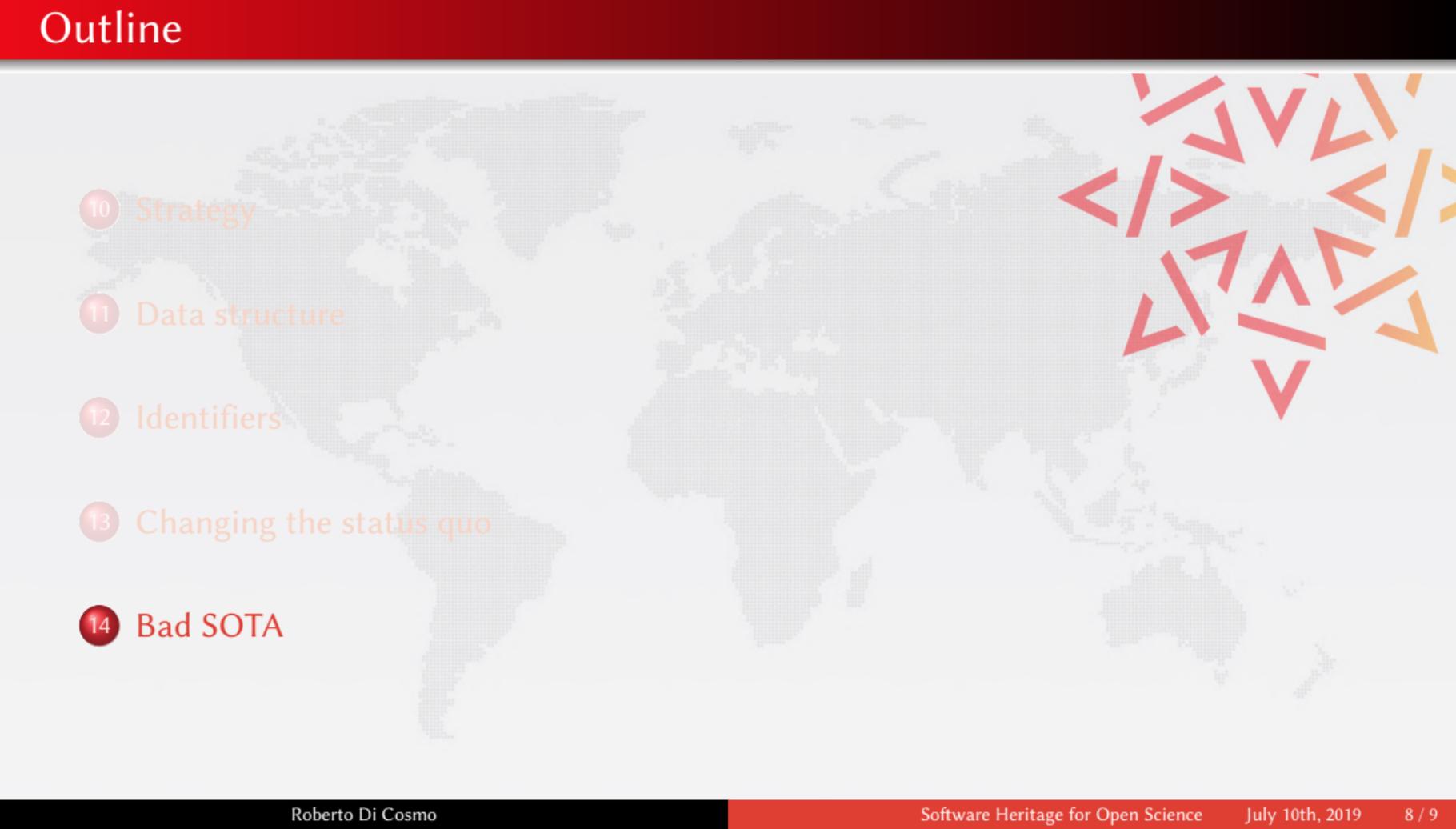
Use DIOs (*specifically, DOIs*) instead of IDOs (*specifically, git commit hashes*)

Original reasons to "avoid commit hashes"

- 1 Version numbers/commit references *are not guaranteed to be permanent.*
- 2 A repository address and version number *does not guarantee that the software is available [...]*
- 3 A particular version number/commit reference *may not represent a "preferred" point at which to cite the software [...]*

Software Heritage changes all this

- 1 SWH-IDs are permanent (and *do not depend on any particular VCS*)
- 2 Software Heritage is *an archive, and guarantees availability*
- 3 At Software Heritage, we separate citation from reference

- 
- 10 Strategy
 - 11 Data structure
 - 12 Identifiers
 - 13 Changing the status quo
 - 14 Bad SOTA**

An example from my research field, Computer Science

Analysis of 613 papers

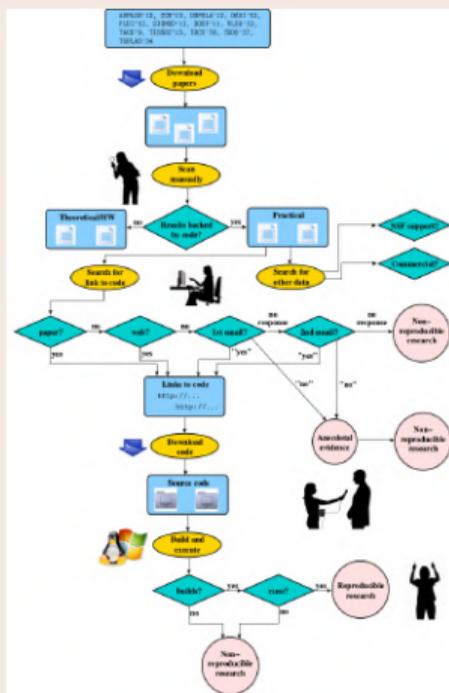
- 8 ACM conferences: ASPLOS'12, CCS'12, OOPSLA'12, OSDI'12, PLDI'12, SIGMOD'12, SOSP'11, VLDB'12
- 5 journals: TACO'9, TISSEC'15, TOCS'30, TODS'37, TOPLAS'34

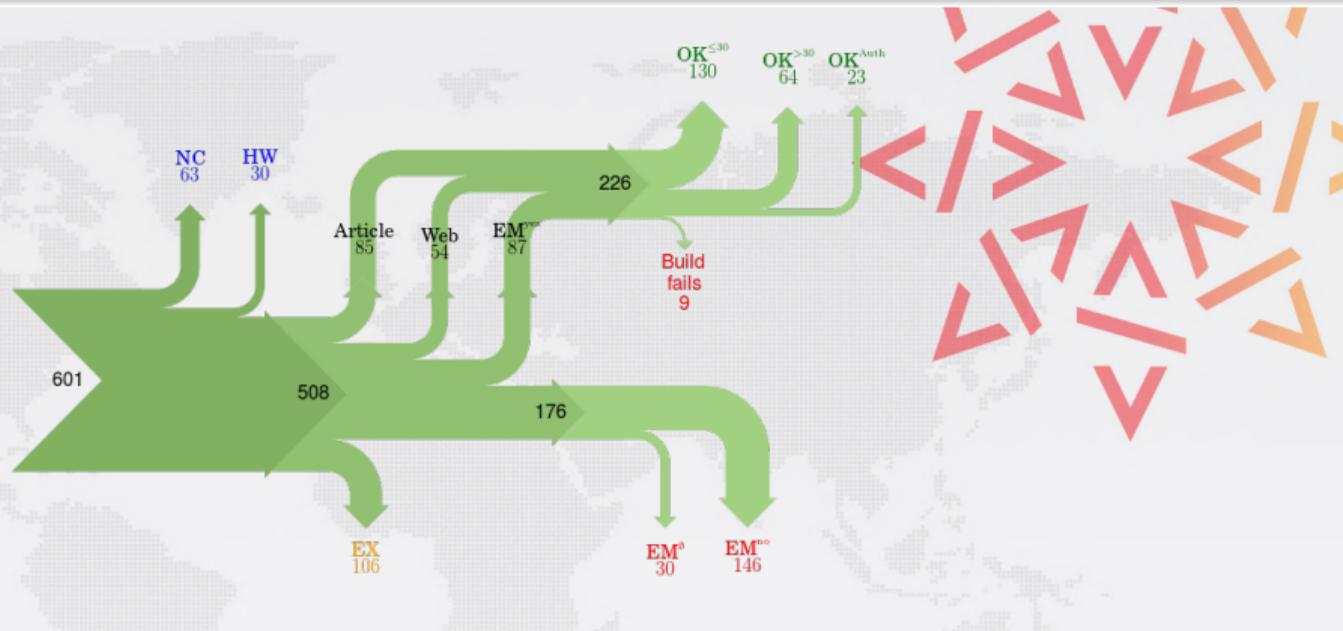
all very practical oriented

The basic question

can we get the code to build and run?

The workflow





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

The main reasons

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