# Software Heritage

#### Collecting, Preserving and Sharing all the Source Code



roberto@dicosmo.org

May 7th, 2019

# Software Heritage

THE GREAT LIBRARY OF SOURCE CODE



## 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
2015 Software Heritage at INRIA
2018 National Committee for Open Science, France

# Source code: *executable* and *human readable*

"The source code for a work making modifications to it."	neans the preferred form of the work for GPL Licence
Hello World	
Program (excerpt of binary)	Program (source code)
4004e6: 55	/* Hello World program */
4004e7: 48 89 e5	
4004ea: bf 84 05 40 00	<pre>#include<stdio.h></stdio.h></pre>
4004ef: b8 00 00 00 00	
4004f4: e8 c7 fe ff ff	<pre>void main()</pre>
4004f9: 90	{
4004fa: 5d	<pre>printf("Hello World");</pre>
4004fb: c3	}

# An essential part of our knowledge

#### Harold Abelson, Structure and Interpretation of Computer Programs

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

Len Shustek, Computer History Museum

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

#### Apollo 11 (~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 (in your pockets!)



#### Roberto Di Cosmo

(1985)

(2006)

# We are at a turning point

#### Loosing our heritage

- the key people are passing away
- we miss a catalog
- we miss an archive

#### Research lagging behind

- research software is undervalued
- software research is under-recognised
- reproducibility is not guaranteed

#### Generalised lack of awareness

- policy: EU copyright reform, Open Science
- general public: insufficient education

#### Diversity far from ideal

- "geek" / "silicon valley" culture
- gender and geography



## Software Heritage, in a nutshell



Roberto Di Cosmo

Software Heritage CC-BY 4.0 May 7th 2019 6 / 1

# All the source code: strategy



# Pull and push





# Using the archive

#### Intrinsic identifiers

- see the iPres 2018 article bit.ly/swhpidpaper
- example: swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa;lines=53-82

#### Wayback machine

- find lost source code
- example: Parmap in a 2012 research article (and how it should be)

#### Research software

moderated deposit via HAL: guidelines for authors and guidelines for moderators

Saving code explicitly

see Save code now

Roberto Di Cosmo



# **Raising Awareness**

#### April 3rd 2017, Unesco Inria agreement







#### 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



Roberto Di Cosmo

# **Growing Support**

#### Sharing the vision Donors, members, sponsors Söftware GitLab eclipse CW2 🕏 fsfe 🌺 SIG OFT SOCIETE GENERALE >= 100Ke/year Microsoft intel Liberti • Égalité • Fraters RÉPUBLIQUE FRANCAISE ADULLACT Computer >= 50Ke/year History Google CAST Museum RÉPUBLIQUE FRANCAISE >= 25Ke/year To software freedom conservancy Software Freedom Adacore & gandi.net DANS NOKIA Bell Labs FREE SOFTWARE SIF © creative GitHub FQSSID UOÀM openinventionnetwork



# You can help!



- save source code now!
- help build a shared process
- bootstrap the *focused search*

#### Preserve

- host ancillary material
- join the SWH Foundation
- become a *mirror*

#### Share

• build on top of SWH

#### Connect

• emulation, binaries, metadata, etc...

#### Awareness

• leverage the Paris Call on Software Source Code (Unesco)

# Join the revolution!

# Software Heritage

### www.softwareheritage.org

#### Library of Alexandria of code

- Working together to
  - recover the past



- preserve our heritage
- share the knowledge
- prepare the future

#### Learn more

#### www.softwareheritage.org/publications

@swheritage

- Jean-François Abramatic, Roberto Di Cosmo, Stefano Zacchiroli Building the Universal Archive of Source Code, Communication of the ACM, October 2018
- **Roberto Di Cosmo, Morane Gruenpeter, Stefano Zacchiroli** *Identifiers for Digital Objects: the Case of Software Source Code Preservation*, iPRES 2018: Intl. Conf. on Digital Preservation
  - Roberto Di Cosmo, Stefano Zacchiroli

Software Heritage: Why and How to Preserve Software Source Code, iPRES 2017: Intl. Conf. on Digital Preservation



# A principled infrastructure

# bit.ly/swhpaper



Roberto Di Cosmo

O Under the hood: architecture

#### 3 Under the hood: identifying billions of object

# Automation, and storage



- full development history permanently archived
- origins: GitHub (auto), Debian (auto), Gitlab.com, Gitorious, Google Code, GNU
- ~ 200Tb raw contents, ~ 10Tb graph (10Bn nodes, 100Bn edges)

# 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

Snapshots



# A bird's eye view



Roberto Di Cosmo

Under the hood: architecture

**8** Under the hood: identifying billions of objects

# 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)

- 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

identifiers (only) for digital objects

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



Roberto Di Cosmo

Software Heritage CC-BY 4.0

8/8