# Towards a Software Pillar for Open Science
## from policy to implementation

Roberto Di Cosmo

Chair, Software Chapter, National Committee for Open Science
Director, Software Heritage
Inria and Université de Paris Cité

July 11th 2023

**Software Heritage**
THE GREAT LIBRARY OF SOURCE CODE

# Outline

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 |
| 2010 | *IRILL* www.irill.org |
| 2015 | *Software Heritage* at INRIA |
| 2018 | *National Committee for Open Science*, France |
| 2021 | *EOSC Task Force on Infrastructures for Software*, European Union |

# Software *Source Code* is Precious Knowledge

## Harold Abelson, Structure and Interpretation of Computer Programs (1st ed.) 1985

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

## Apollo 11 source code (excerpt)

```
P63SPOT3    CA      BIT6        # IS THE LR ANTENNA IN POSITION 1 YET
            EXTEND
            RAND    CHAN33
            EXTEND
            BZF     P63SPOT4    # BRANCH IF ANTENNA ALREADY IN POSITION 1

            CAF     CODE500     # ASTRONAUT:    PLEASE CRANK THE
            TC      BANKCALL    #               SILLY THING AROUND
            CADR    GOPERF1
            TCF     GOTOPOOH    # TERMINATE
            TCF     P63SPOT3    # PROCEED    SEE IF HE'S LYING

P63SPOT4    TC      BANKCALL    # ENTER      INITIALIZE LANDING RADAR
            CADR    SETPOS1

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

## Quake III source code ( excerpt )

```c
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."*

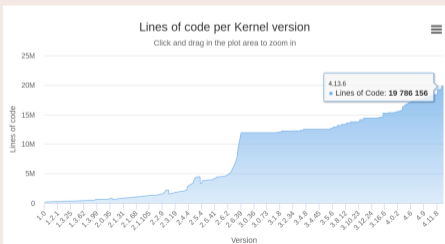# ~ 50 years, a lightning fast growth

## 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!

# Outline

# Software is eating the world…

## Business

THE WALL STREET JOURNAL.

Home   World   U.S.   Politics   Economy   Business   Tech   Markets   Opinion   Arts

ESSAY

## Why Software Is Eating The World

*By Marc Andreessen*
August 20, 2011

This week, Hewlett-Packard (where I am on the board) announced that it is exploring jettisoning its struggling PC business in favor of investing more heavily in software, where it sees better potential for growth. Meanwhile, Google plans to buy up the cellphone handset maker Motorola Mobility. Both moves surprised the tech world. But both moves are also in line with a trend I've observed, one that makes me optimistic about the future

Software companies

outperform or buy out
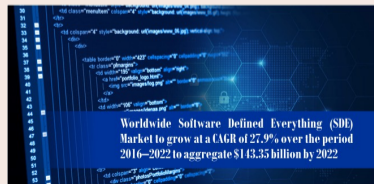
hardware companies
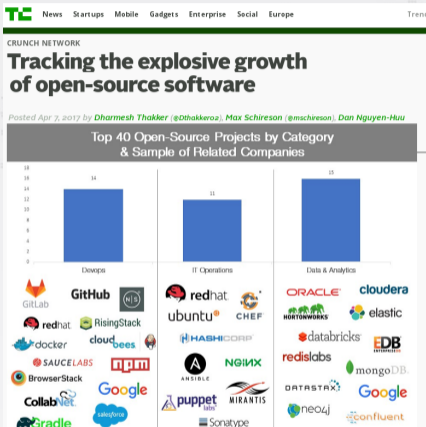
*Marc Andreesen, 2011*

## Technology

**Software Defined Everything**

Hardware gets commoditised

Software becomes the new value!

Worldwide Software Defined Everything (SDE) Market to grow at a CAGR of 27.9% over the period 2016–2022 to aggregate $143.35 billion by 2022

## Open Source Software

can be openly (re)used, modified, (re)distributed, *with full access to its source code!*

# Outline

# Free Software: 40 years, 4 layers, in a nutshell

## First 15 years: 1984-…                                    The early revolution

focus *freedom* for users and (especially) developers

keyword free software

## The second wave: 1999-…                            Progressive industry adoption

focus software quality and reduced cost

keyword open source (~25th anniversary!)

## The third wave: 2010-…                              Ecosystems, strategic alignment

focus community, organisation, foundations

keyword governance

## The fourth wave: 2015-…                                Industry consolidation

focus mergers and acquisitions

keyword control

# Key resources : competency, and adoption

## We really are in a *knowledge* economy!

- competencies
- talent
- network
- adoption
- mindshare

## Bottomline

*The infrastructure for (open) collaboration* is the new competitive advantage!
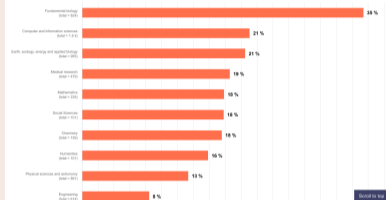
# Outline

# Software is a pillar of Open Science

## Software powers modern research



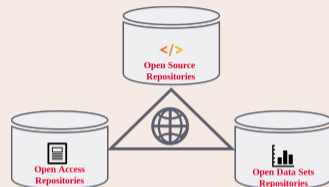Proportion of publications in France that mention code or software sharing by discipline

*Over 20% of articles across all disciplines share software*
*2023 French Open Science Monitor*

## Key pillar: software



Open Source Repositories

Open Access Repositories

Open Data Sets Repositories

Links are **important**

## Nota Bene

software may be a *tool*, a *research outcome* and a *research object*

access to the *source code* is essential!

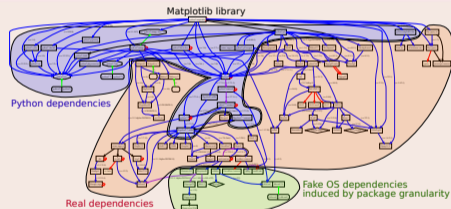Preserving (the history of) source code is necessary for *reproducibility*

# Source code is *special* (software is *not* data)

## Software *evolves* over time

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

## Complexity

- *millions* of lines of code
- large *web of dependencies*
  - easy to break, difficult to maintain
  - *research software* a thin top layer
- sophisticated *developer communities*



Matplotlib library

Python dependencies

Real dependencies

Fake OS dependencies induced by package granularity

## The human side

design, algorithm, code, test, documentation, community, funding

and so many more facets . . .

# How are we managing our software ?

## Reproducibility, maintenance in Academia



(articles: here, here, here and here)

## Security, integrity, traceability in Industry



Can they track the software that they

- ship, use, acquire
- has that bug or vulnerability

awareness is raising at the level of public policy

# Outline

# International highlights

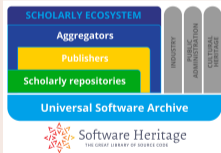## Paris Call on Software Source code (2019, UNESCO)

40 international experts call to *"promote software development as a valuable research activity, and research software as a key enabler for Open Science/Open Research, [...] recognising in the careers of academics their contributions to high quality software development, in all their forms"*

☞ Open Source in UNESCO recommendations for Open Science, 2021

## Software in the EOSC

2020 EOSC SIRS connect scholarly ecosystem via Software Heritage
2021 EOSC Task Force on Infrastructures for Research Software
2022 FAIRCORE4EOSC project WP6 implements SIRS report
2023 INFRAEOSC call on quality of scientific software

## And much more

Software track in OSEC 2022, Software working group launched in Science Europe, DFG adds software to model CV (9/22), NASA unveils Open Science policy (12/22), ...

# What is at stake

## ARDC

- Archive for retrieval (*reproducibility*)
- Reference for identification (*reproducibility*)
- Describe for discovery and reuse
- Cite/Credit for credit and evaluation

## Before ARDC

- Development practices and tools (VCS, build system, test suites, CI, code quality, . . . )
- Opening up towards a community (documentation, organization, communication)

Need training, tooling, infrastructures, best practices

## Beyond ARDC

- Policies (dissemination, reuse, careers, . . . )
- Sustainability (legal, financial, etc.)
- Technology transfer
- Advanced technologies and tools (quality, traceability, etc.)

a humbling challenge, and a complex one (we are not in a vacuum)

**MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR, DE LA RECHERCHE ET DE L'INNOVATION**
*Liberté Égalité Fraternité*

## SECOND FRENCH PLAN FOR OPEN SCIENCE

Generalising open science in France 2021-2024

*1*

---

Second French Plan for Open Science

**Launch on 6 July 2021** by Frédérique Vidal, Minister for Higher Education, Research and Innovation

- Multiplying the **levers for change** in order to **generalise open science practices**
- Structuring the **policy for opening up or sharing research data**
- New commitments to the **opening of source code produced by research**
- **European and international inclusion** in the context of the French Presidency of the European Union
- **Disciplinary and thematic variations**: open science policies must be adapted to disciplinary specificities

GENERALISING OPEN SCIENCE IN FRANCE 2021-2024

*2*

---

Path Three :
## Opening up and promoting source code produced by research

**7** Recognize and support the dissemination under an open source license of software produced by publicly funded research programmes

**8** Highlight the production of source code from higher education, research and innovation

**9** Define and promote an open source software policy

« The opening of software source code is a major challenge for the **reproducibility** of scientific results. »

« Distribution of software products under **open source licence** will be preferred. »

*3*

---

**Define and promote an open source software policy**

- Produce a **National Charter for Open Source Software** coming from higher education, research and innovation
- Develop the **link between data and software** through a network of **Chief Data Officers** in the various universities and research performing organisations
- Develop the **economic models of open source software** and make them known within commercialization services
- **Support Software Heritage** and recommend it for the archiving and referencing of source code

**Recognise source code as a contribution to research**

- Create an **open source research software prize**
- **Provide greater recognition** for software production in the career of researchers, research support staff

**Build an ecosystem that connects code, data and publications**

- Develop **proper coordination** between software forges, open publication archives, data repositories and the scientific publishing sector.

*4*

# Software College in the CoSO

## Five action lines

- **Identifying and highlighting** research software production
- Technical **tools** and best practices
- Translation and **sustainability**
- National, European, and International **coordination**
- Recognition, evaluation and **careers**

## Leveraging experience and connections

- Open Source thematic group in Systematic (since 2007, more on demand)
- Collaboration with DINUM, Eclipse Foundation, OW2, ...

## Coordination with other colleges

- The Open Science passport software booklet

**MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE**
Liberté
Égalité
Fraternité

## The first national Open Science award for Research Software

**2022 edition**

- 120+ high quality submissions
- 4 prizes
- 6 accessit
- 4 categories (inclusiveness)
- awarded by the Ministry of Research

## Institutionalised as an annual award

**2023 edition** now open, already inspired other countries (e.g. Australian award)

Detailed description and lessons learned **forthcoming**

## Twenty-three active members

Chairs: Roberto Di Cosmo and François Pellegrini

- Florent CHUFFART (Univ Grenoble Alpe)
- Mélanie CLÉMENT-FONTAINE (Univ Paris-Saclay - Versailles Saint-Quentin)
- Laurent COSTA (UMR 7041 ArScAn)
- Ludovic COURTÈS (Inria)
- Sébastien GÉRARD (Univ Paris-Saclay, CEA, List)
- Mathieu GIRAUD (CNRS, Univ Lille)
- Timothée GIRAUD (CNRS)
- Jean-Yves JEANNAS (Univ Lille, AFUL)
- Nicolas JULLIEN (IMT Atlantique)
- Daniel LE BERRE (Univ Artois, CNRS)
- Violaine LOUVET (CNRS / GRICAD - Univ Grenoble Alpes)
- Camille MAUMET (Inria, Univ Rennes, CNRS, Inserm)
- Clémentine MAURICE (CNRS)
- Grégory MIURA (Univ Bordeaux Montaigne)
- Raphaël MONAT (LIP6, Sorbonne Université)
- Patrick MOREAU (CNRS)
- Sophie RENAUDIN (AP-HP)
- Nicolas ROUGIER (Inria, Univ Bordeaux, CNRS)
- François SABOT (IRD)
- Sylvie TONDA-GOLDSTEIN (Inria)
- Samuel THIBAULT (Univ Bordeaux) (Univ Paris-Saclay)

# Outline

# How (not) to preserve and share research software

## A - Since the ~~1970's~~ 1990's

.zip or .tar file on:

- ~~ftp server~~ (e.g. gnu)
- web page (example)
- document archive (+ DOI sample)

## B - Since the 2000's

Rely on *software forges*

- institutional/project (e.g. example)
- free commercial ones: BitBucket, GitHub, GitLab, … (e.g. parmap)

## C: a mix of the two



Artifacts Available          Artifacts Evaluated & Functional

**Authors/Contributors:** Authors Info & Affiliations

**DOI:** https://doi.org/10.1145/████ **Version:** 1.0

**Description**

A source archive of ██████, and the version of █████ used in the paper eval. A more up-to-date version of █████ can be found at github.com/████/██████.git

**Assets**

Read Me ██████████████████

⬇ Download (3.5 KB)

## Can get no satisfaction…

    A  *Poor user experience*

    B  *No preservation guarantee*

    C  Can do *so much* better

# Forges are *not* archives!

## 2015: the first big bad news

Google Code and Gitorious.org shutdown: ~1M endangered repositories

- broken links in the web of knowledge (my papers too)

## Big bad news keep coming in

- summer 2019: BitBucket announces Mercurial VCS sunset
- july 2020: BitBucket erases *250.000+* repositories (including research software)
- summer 2022: GitLab.com considers erasing all projects that are inactive for a year

## In Academia too!

- 2021: Inria's old gforge is unplugged... breaks the Opam build chain for OCaml

We need a universal archive of software source code: now we have one!

# Software Heritage
## THE GREAT LIBRARY OF SOURCE CODE

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

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

---

**Reference catalog**



**find** and **reference** all
software source code

---

**Universal archive**



**preserve and share** all
software source code

---

**Research infrastructure**



**enable analysis** of all
software source code

---

# The largest software archive, a shared infrastructure

One infrastructure
open and shared

| Cultural Heritage | Industry | Research | Public Administration |

## Software Heritage

The largest archive ever built

| Source files | Commits | Projects |
|---|---|---|
| 13,974,813,954 | 2,912,845,019 | 207,160,527 |

| | | |
|---|---|---|
| Bitbucket 2,012,133 origins | git 19,494 origins | R 21,486 origins |
| debian 129,217 origins | 6,424 origins | GitHub 152,282,093 origins |
| GitLab 3,989,638 origins | Guix 12,451 origins | GNU 354 origins |
| heptapod 1,096 origins | launchpad 356,873 origins | Maven 93,710 origins |
| NixOS 12,451 origins | npm 1,799,296 origins | 4,083 origins |
| Phabricator 185 origins | python 427,135 origins | SOURCEFORGE 308,970 origins |

## Sharing the vision



And many more ...

## Donors, members, sponsors



Diamond sponsor

Platinum sponsors
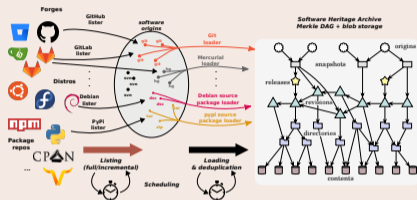
Gold sponsors

Silver sponsors

Bronze sponsors

# Addressing the four needs (see ICMS 2020 for details)

## Archive (15B+ files, 240M+ projects)



- save now, updateswh, webhooks
- deposit.softwareheritage.org

## Reference (30 billion SWHIDs)

### Intrinsic, cryptographically strong IDs



Now in SPDX 2.2, Wikidata
Specification: `https://swhid.org`

## Describe

- *Intrinsic metadata* from source code
- Contributed the Codemeta generator

## Cite/Credit

- Contributed *software citation* style
  biblatex-software, v 1.2-2 now on CTAN

# Mutualization and standardisation at work

## One archive, multiple infrastructures



universal software archive  *Software Heritage* connects with the global software ecosystem

scholarly repositories  institutional and disciplinary archives

publishers  journals, proceedings, preprints

aggregators  disciplinary catalogues, meta-portals, …

## Building interconnection and interoperability          FAIRCORE4EOSC HE (2022-2025)

Beta release: EOY 2023

*Interconnection* with SWH

repositories  HAL, InvenioRDM, Dataverse

publishers  Dagstuhl, episciences

agregators  swMath, OpenAire

*Interoperability*

metadata schema  *CodeMeta*

intrinsic identifier  *SWHID*

specifications  open/public

# Outline

# Call to action: best practices for ARDC are available... today!

## Archiving and referencing

For all source code used in research (*yes, even small scripts!*)

- ensure it is archived in Software Heritage (see save code now)
- get the proper SWHID for your software (see detailed HOWTO)
- add it to research articles for reproducibility (see detailed HOWTO)

## Describing and Citing/Crediting

For software you want to put forward (*mention in your CV, reports, etc., get citations and credit for it*), do the following extra steps:

- add codemeta.json with description (see the codemeta generator)
- reference in the HAL portal (french partners, see online HAL documentation)
- cite software using the biblatex-software package (in CTAN and TeXLive)

- train students and colleagues
- engage journals, conferences, learned societies

# Call to action: policy making

## A working agenda

- avoid proprietarisation: set the default to open
  - *publicly funded research software should be open source*, exceptions must be justified
  - set up institutional support
  - build common knowledge base for technology transfer offices
- establish intelligent and effective incentives
  - count quality software contributions in careers, avoid purely numerical indicators, keep the human in the loop (mind Goodhart's law)
- avoid balkanisation, support mutualised common infrastructures
  - build on common, shared, open, non profit infrastructures, like Software Heritage
  - acknowledge the predominant human component of digital infrastructures
    - recurrent funding of their cost
    - proper evaluation of their service

it's a long road, but together we can make it

# Questions?

## References

UNESCO, *Draft recommendations on Open Science*
2021, (online)

French Ministry of Research, *Second National Plan for Open Science*
2021, (online)

EOSC SIRS Task Force, *Scholarly Infrastructures for Research Software*
2020, Publications office of the European Commission, (10.2777/28598)

R. Di Cosmo, *Archiving and Referencing Source Code with Software Heritage*
International Conference on Mathematical Software 2020 (10.1007/978-3-030-52200-1_36)

J.F. Abramatic, R. Di Cosmo, S. Zacchiroli, *Building the Universal Archive of Source Code*
CACM, October 2018 (10.1145/3183558)