

Software Heritage

Does software preservation have an ethical impact on society?

Morane Gruenpeter

Software engineer and metadata specialist
Inria, Software Heritage

morane@softwareheritage.org

November 30th, 2020



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Introduction
- 2 The knowledge is in the source code !
- 3 Software Heritage: the universal source code archive
- 4 The Paris call: Software Source Code is our Heritage
- 5 Strategies for archiving
- 6 Software preservation and ethics



Goal: Building the Semantic Web of Free and Open Source Software



1999-2011 Harpist

2012-2015 Licence in Computer Science CNAM

2015-2017 Master STL - M2 R&D UPMC

2017 Internship *Software Heritage* (SWH)

2018-2019 European project EU2020 *CROSSMINER* (on SWH team)

2020-2022 European project *FAIRsFAIR* (on SWH team)

Working groups for Open Science and digital preservation

- the Research Data Alliance's **Software Source Code** Interest Group (SSC IG),
- the FORCE11's **Software Citation** Implementation Working Group (SCI WG),
- the joint RDA & FORCE11 **Software Identification** Working Group (SCID WG)
- WikiData for **Digital Preservation** initiative (WikiDigi).

Terminology

- understand terminology around software and source code

History

- a few elements to grasp recent and not so recent evolutions

Society

- software as a key to view society
- review initiatives working on software preservation

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

The World Wide Web, 1989, at CERN on a NeXT machine

"When somebody has learned how to program a computer ... You're joining a group of people who can do incredible things. They can make the computer do anything they can imagine."



From An Insight, An Idea with Tim Berners-Lee (2013)

What is software ?

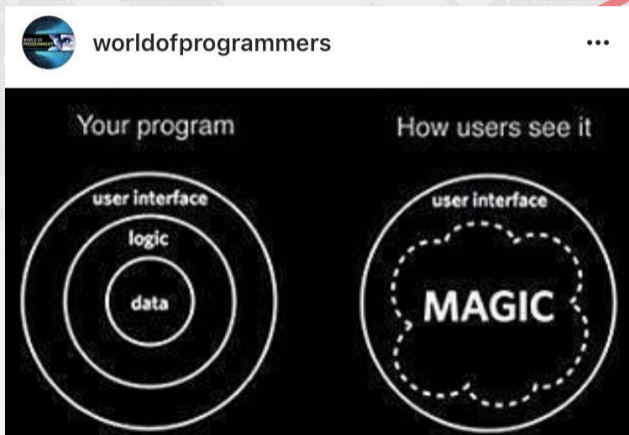


Image taken from [reddit - ProgrammerHumor](#)

Encyclopædia Britannica

“Software, instructions that tell a computer what to do. Software comprises the entire set of programs, procedures, and routines associated with the operation of a computer system. The term was coined to differentiate these instructions from hardware—i.e., the physical components of a computer system.”

[link](#)

Software as a concept

- software project / entity
- the creators and the community around it

Software artifact

- the binaries for different environments
- the **software source code** for each version

This is *software*?



Ceci n'est pas une pipe.

What about *software source code*?

- 1 Introduction
- 2 The knowledge is in the source code !
- 3 Software Heritage: the universal source code archive
- 4 The Paris call: Software Source Code is our Heritage
- 5 Strategies for archiving
- 6 Software preservation and ethics



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*

Executable and human readable knowledge

copyright law

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

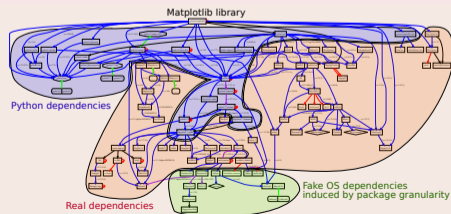
Harold Abelson

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
- sophisticated *developer communities*







Software Source Code human readable and executable knowledge


Full width Home Development Documentation Donate login

 **Software Heritage**
Archive

Features

-  Search
-  Downloads
-  Save code now
-  Help

```
52
53 # THE MASTER IGNITION ROUTINE IS DESIGNED FOR USE BY THE FOLLOWING LEM PROGRAMS: P12, P40, P42, P61, P63.
54 # IT PERFORMS ALL FUNCTIONS IMMEDIATELY ASSOCIATED WITH APS OR DPS IGNITION: IN PARTICULAR, EVERYTHING LYING
55 # BETWEEN THE PRE-IGNITION TIME CHECK -- ARE WE WITHIN 45 SECONDS OF TIG? -- AND TIG + 26 SECONDS, WHEN DPS
56 # PROGRAMS THROTTLE UP.
57 #
58 # VARIATIONS AMONG PROGRAMS ARE ACCOMODATED BY MEANS OF TABLES CONTAINING CONSTANTS (FOR AVEGEXIT, FOR
59 # WAITLIST, FOR PINBALL) AND TCF INSTRUCTIONS. USERS PLACE THE ADRES OF THE APPROPRIATE TABLE
60 # (OF P61TABLE FOR P61LM, FOR EXAMPLE) IN ERASABLE REGISTER 'WHICH' (E4). THE IGNITION ROUTINE THEN INDEXES BY
61 # WHICH TO OBTAIN OR EXECUTE THE PROPER TABLE ENTRY. THE IGNITION ROUTINE IS INITIATED BY A TCF BURNBABY,
62 # THROUGH BANKJUMP IF NECESSARY. THERE IS NO RETURN.
63 #
64 # THE MASTER IGNITION ROUTINE WAS CONCEIVED AND EXECUTED, AND (NOTA BENE) IS MAINTAINED BY ADLER AND EYLES.
65 #
66 #           HONI SOIT QUI MAL Y PENSE
67 #
68 #           *****
69 #           TABLES FOR THE IGNITION ROUTINE
70 #           *****
71 #
72 #           NOLI SE TANGERE
73
74 P12TABLE      VN      0674      # (0)
75              TCF      ULLGNOT   # (1)
76              TCF      COMFAIL3  # (2)
77              TCF      GOCUTOFF   # (3)
78              TCF      TASKOVER   # (4)
79              TCF      P12SPOT    # (5)
80              DEC      0          # (6)      NO ULLAGE
81              EBANK=  WHICH
82              2CADR  SERVEXIT    # (7)
83
84              TCF      DISPCHNG   # (11)
85              TCF      WAITABIT   # (12)
86              TCF      P12IGN     # (13)
87
88 P40TABLE      VN      0640      # (0)
```

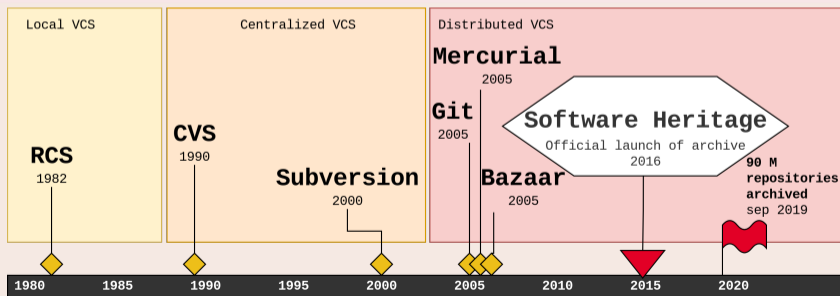
 Permalinks

Version Control System timeline

Version control system (VCS)

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

Three decades of evolution



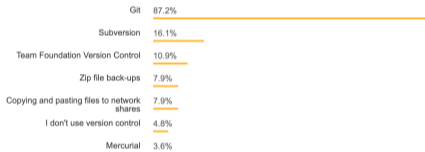
Stack Overflow

[Survey 2018]

Version Control

All Respondents

Professional Developers



74,298 responses; select all that apply

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

In numbers

GitHub [Octoverse 2017] [Blog 2018]

- 100.000.000+ repositories
- 40.000.000+ developers worldwide

Bitbucket [Blog 2019]

- 28.000.000+ repositories
- 10.000.000+ developers worldwide

GitLab [Blog 2019]

- 1.000.000 MRs March 19'

Let's use it!

- 1 Introduction
- 2 The knowledge is in the source code !
- 3 Software Heritage: the universal source code archive**
- 4 The Paris call: Software Source Code is our Heritage
- 5 Strategies for archiving
- 6 Software preservation and ethics





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 all software source code

Research infrastructure



enable analysis of all software source code

Cultural Heritage



Industry



Research



Education



Software Heritage

As of today the archive already contains and keeps safe for you the following amount of objects:

Source files

8,846,381,610

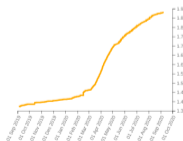


Directories

7,506,954,410

Commits

1,880,663,008

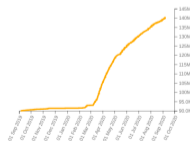


Authors

38,603,337

Projects

140,348,311



Releases

15,051,940

Raising awareness: landmark agreement, 3/4/2017; grand opening, 7/6/2018



Sharing the vision



Morane Gruenpeter

Sponsoring our work



Platinum sponsors



Gold sponsors



Silver sponsors



www.softwareheritage.org

November 30th, 2020

15 / 31

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 Introduction
- 2 The knowledge is in the source code !
- 3 Software Heritage: the universal source code archive
- 4 **The Paris call: Software Source Code is our Heritage**
- 5 Strategies for archiving
- 6 Software preservation and ethics



The Paris call: Software Source Code is our Heritage

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



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

see full text

Quotes

- "Preserving Software Source Code is **crucial** and **captures human civilization**"
- "It includes the need to **raising awareness** of the importance of SSC among decision makers, recognizing Software creators and the contributions of women and minorities to digital innovations".
- "Considering that documents produced and **preserved overtime**, in all their analog and digital forms through time and space, constitute the **primary means of knowledge creation and expression**, having an **impact on all areas of humanity's civilization** and its further progress".
- "Considering at the same time that the preservation of, and long term accessibility to documentary heritage underpins **fundamental freedoms** of opinion, expression and information as human rights".

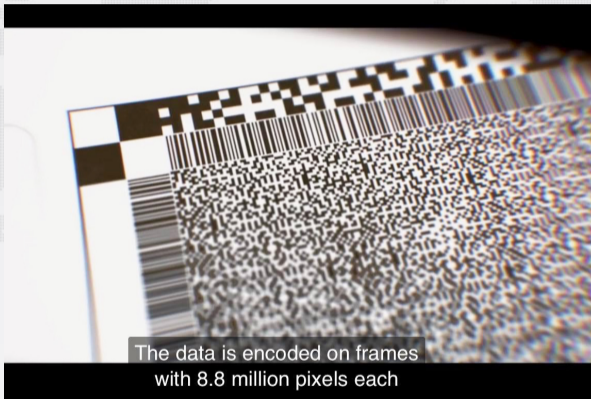
- 1 Introduction
- 2 The knowledge is in the source code !
- 3 Software Heritage: the universal source code archive
- 4 The Paris call: Software Source Code is our Heritage
- 5 Strategies for archiving
- 6 Software preservation and ethics



Joining forces in the urgent effort to preserve humankind's source code.



- A **testament to the importance** of software source code preservation
- A **multi-partners** strategy for archiving code
- A range of **storage solutions**, from real-time to long-term storage



The screenshot displays the Internet Archive homepage. At the top, a navigation bar includes links for ABOUT, HOME, PROJECTS, HELP, DONATE, CAREER, JOIN, VOLUNTEER, and PEOPLE. Below this is a search bar for the Wayback Machine. The main content area features a central banner with the Internet Archive logo (a classical building) and the text: "Internet Archive is a non-profit library of millions of free books, movies, software, music, websites, and more." To the right of the banner is an "Announcements" section with three items. Below the banner is a search bar with a "GO" button and a link to "Advanced Search". The "Top Collections at the Archive" section is a grid of ten collection cards, each with an icon, name, and item count.

Collection Name	Item Count
Community Audio	2,205,407 items
Community Video	1,016,433 items
Community Text	1,993,487 items
American Libraries	3,417,794 items
Community Data	254,669 items
The LibriVox Free Audiobook	14,790 items
Canadian Libraries	873,070 items
Electric Sheep	614 items
Live Music Archive	122,277 items
Community Images	261,094 items

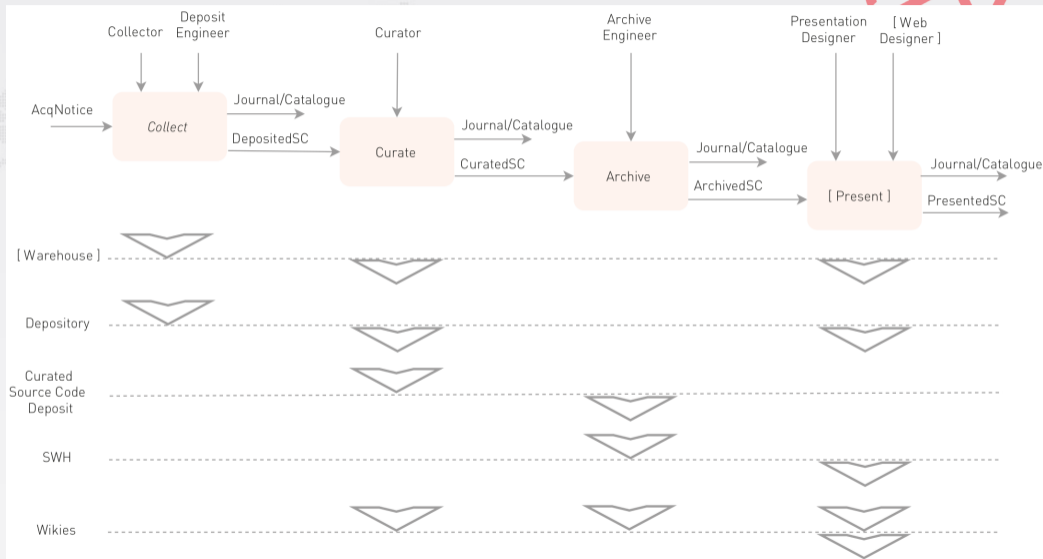
Paris Call on Software Source Code

“[We call to] support efforts to gather and preserve the artifacts and narratives of the history of computing, while the earlier creators are still alive”



- **Rescue** Legacy Software from different media
 - physical
 - digital
 - legacy / unsupported
 - recent / supported
- **Curate** the code
 - reconstructing the development history
 - collecting metadata
- And **illustrate** with dedicated presentations

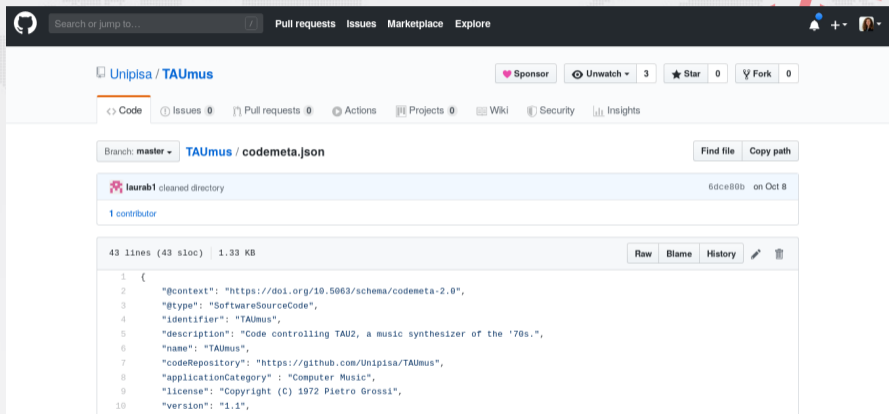
SWHAP: Four phases workflow to streamline the effort





- The control code for the music synthesizer TAU2
- FORTRAN II and TAUmus command language
- Istituto di Elaborazione dell'Informazione CNR
- Group led by the late M° P. Grossi
 - Le Sacre du Printemps (ABSTRACT)

SWHAP@PISA: Capturing metadata in branch master



Search or jump to... Pull requests Issues Marketplace Explore

Unipisa / TAUmus Sponsor Unwatch 3 Star 0 Fork 0

Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security Insights

Branch: master TAUmus / codemeta.json Find file Copy path

laurab1 cleaned directory 6dce89b on Oct 8

1 contributor

43 lines (43 sloc) 1.33 KB Raw Blame History

```
1 {
2   "@context": "https://doi.org/10.5063/schema/codemeta-2.0",
3   "@type": "SoftwareSourceCode",
4   "identifier": "TAUmus",
5   "description": "Code controlling TAU2, a music synthesizer of the '70s.",
6   "name": "TAUmus",
7   "codeRepository": "https://github.com/Unipisa/TAUmus",
8   "applicationCategory": "Computer Music",
9   "license": "Copyright (C) 1972 Pietro Grossi",
10  "version": "1.1",
```

SWHAP@PISA: Recreating development history in branch SourceCode

The screenshot shows the GitHub interface for the repository `Unipisa / TAUmus`. The page is viewed on the `SourceCode` branch. The commit history for October 8, 2019, is displayed, showing two commits:

- v1.1** - `be97ff8`: Pietro Grossi authored and `laurab1` committed on Oct 16, 1972
- v1.0** - `a99524e`: Pietro Grossi authored and `laurab1` committed on Sep 16, 1972

Navigation buttons for "Newer" and "Older" are visible at the bottom of the commit list.

Home [Archive](#) Development Documentation [Donate](#)

Software Heritage

Archive Access

- Browse
- Web API

Features

- Search
- Vault
- Save code now**

Miscellaneous

- Help

Save code now Beta version

Origin type: Origin url:

Submit

Software Heritage Archive

☰ Browse archived content for origin <https://github.com/Unipisa/TAUmus>

📅 Visits 📷 Snapshot date: 08 October 2019, 17:49 UTC 📁 Branches (3) 📦 Releases (0)

📄 Branch: refs/heads/SourceCode 8c34070 / SUBROUTINE_CALMUS.FOR

📄 Raw File 🗑 Select Language ⋮ Actions

SUBROUTINE_CALMUS.FOR

```
1 SUBROUTINE CALMUS
2 COMMON FR, T, R, S, INIZ, IFI
3 *ALR, DURA, NOA, IP, KONT
4 DIMENSION NNN(1700), I(10), NN(
5 1 FR(5000), T(5000) NPLLOD
6 REAL KFT(8)
7 READ(5,10)N, N1, N2
8 10 FORMAT(3I4)
9 N=4
10 LN=1
11 KK=0
12 KONT=0
13 KKONT=0
14 L=0
15 LL=0
16 KP=0
17 K=0
18 11 K=K+1
19 IF(K.LE.N)GO TO 12
20 LN=LN+1
21 IF(LN.GT.5)GO TO 50
22 GO TO 8
```

Permalinks

To reference or cite the objects present in the Software Heritage archive, permalinks based on persistent identifiers must be used instead of copying and pasting the url from the address bar of the browser (as there is no guarantee the current URI scheme will remain the same over time).

Select below a type of object currently browsed in order to display its associated persistent identifier and permalink.

📄 content 📁 directory ↶ revision 📷 snapshot

🌐 archived repository 🌐 archived swh:1:cnt:137968d1b65eabd6390647f95119f2eb24704962

swh:1:cnt:137968d1b65eabd6390647f95119f2eb24704962;origin=https://github.com/Unipisa/TAUmus

Add origin info Add selected lines info

📄 Copy identifier 📄 Copy permalink

- 1 Introduction
- 2 The knowledge is in the source code !
- 3 Software Heritage: the universal source code archive
- 4 The Paris call: Software Source Code is our Heritage
- 5 Strategies for archiving
- 6 Software preservation and ethics



Ethical Charter for using the archive data

<https://bit.ly/2K0yfeR>

- Avoid harm
- Protect Personal Data
- Avoid useless copies
- Care about derived data

Ethical Charter for Mirrors

<https://bit.ly/36q2nWx>

- Avoid harm
- Protect Personal Data
- Maintain coherent terms of use
- Ensure fair and non discriminatory access
- Foster Collaboration

Discussion topics

- 1 Transparency and conflict of interest
- 2 Digital divide
- 3 Ownership of the physical archive
- 4 Responsibility for the incorrect use of source codes
- 5 The positive potential of this archive
- 6 Possible role of philosophers in interdisciplinary teams with scientists



Software Heritage

Thank you! Any questions?

contact: morane@softwareheritage.org



Jean-François Abramatic, Roberto Di Cosmo, Stefano Zacchioli

Building the Universal Archive of Source Code, Communications 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

Acknowledgements

- Roberto Di Cosmo, Founder and Director of Software Heritage
- Leah Gruenpeter-Gold, PhD Philosophy Dept., Tel Aviv University