

# Software Heritage, l'archive mondiale du code source logiciel

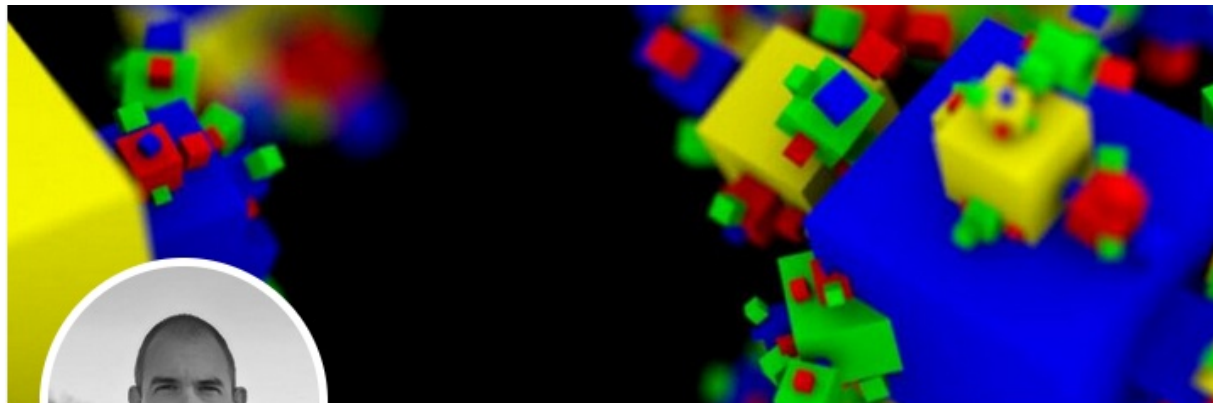
Pierre Poulain



Software Heritage  
THE GREAT LIBRARY OF SOURCE CODE

Les Rendez-vous du Centre des Humanités Numériques | 2022-04-21

# Bonjour 🖐️



**Pierre Poulain**

@pierrepo

Bioinformatics. Python. Unix. Associate prof at [@Univ\\_Paris](#) and [@IJMonod](#)  
[#proteomics](#) [#openscience](#) [#opensource](#) [@SWHeritage](#) ambassador  
[#pedagogy](#) [#Certifiens](#)

📍 Paris, France [🔗 cupnet.net](#)

✉️ [pierre.poulain@u-paris.fr](mailto:pierre.poulain@u-paris.fr)

🐦 [@pierrepo](#)

# Software Heritage ?



- Soutenir **Software Heritage** et recommander son adoption pour l'archivage et le référencement des codes sources.
- Proposer la standardisation du **Software Heritage Identifier (SWHID)**, qui complètera les DOI pour les logiciels.

Source : Ouvrir la science, 2021

# Software Heritage ?



## Feuille de route pour la Science ouverte à Université de Paris

### Université de Paris place la Science ouverte au cœur de son projet d'établissement.

De l'accès illimité et immédiat à l'information scientifique qu'autorise le support numérique dépendent la vitesse des échanges, la richesse des collaborations, la possibilité de reproduction des expériences et de réutilisation des données, la multiplication des innovations de transition voire de rupture, la qualité du dialogue entre les sciences et la société.

Issues d'initiatives pionnières et de déclarations historiques, l'ouverture et le partage des résultats publiés, des données sur lesquelles ils s'appuient et des codes et méthodes qui les ont produits, orientent désormais les politiques publiques aux échelles nationale ([Plan National pour la Science Ouverte 2021-2024](#)), européenne ([Plan S](#), [Horizon Europe 2021-2027](#)) et mondiale ([recommandation de l'UNESCO](#)).

### Université de Paris se donne trois objectifs principaux

**Faire de la Science ouverte un instrument de souveraineté scientifique** en se réappropriant les résultats de la recherche financée sur fonds publics, par le recensement et le signalement systématique des productions scientifiques d'Université de Paris, par la maîtrise des conditions de leur partage selon le principe « aussi ouvert que possible, aussi fermé que nécessaire » ;

**Faire de la Science ouverte un projet au service des personnels et des usagers** en levant les barrières économiques, techniques et juridiques à la circulation de l'information scientifique et technique, en facilitant la mise en conformité aux exigences des agences de financement en matière de Science ouverte, en simplifiant, en accompagnant, en encourageant les meilleures pratiques ;

**Faire de la Science ouverte un levier** pour l'accélération de l'innovation, l'intégrité académique et l'amélioration du dialogue science et société.

### Université de Paris définira les conditions et mettra en œuvre une triple ouverture

Des publications (*open access*).

3. Promouvoir les bonnes pratiques en lien avec l'utilisation de la plateforme Software Heritage
  - Encourager au dépôt dans Software Heritage, en particulier via HAL, des logiciels et des codes sources développés au sein d'Université de Paris
  - Intégrer le Software Heritage Identifier (SWHID) à la politique sur les identifiants ouverts.

Source : Université Paris Cité, 2021

# Au menu



1. Pourquoi l'archivage des codes sources des logiciels est important ?
2. Qu'est-ce que Software Heritage ?
3. Comment archiver son code source dans Software Heritage ?
4. Quelles sont les bonnes pratiques pour archiver un code source ?

*Cette présentation est garantie sans ligne de code*

# Les logiciels sont partout !



NEWS FEATURE

## THE TOP 100 PAPERS

Nature explores the most-cited research of all time.

BY RICHARD VAN NOORDEN,  
BRENDAN MAHER AND REGINA NUZZO

The discovery of high-temperature superconductors, the determination of DNA's double-helix structure, the first observations that the expansion of the Universe is accelerating — all of these breakthroughs won Nobel prizes and international acclaim. Yet none of the papers that announced them comes anywhere close to ranking among the 100 most highly cited papers of all time.

Citations, in which one paper refers to earlier works, are the standard means by which authors acknowledge the source of their methods, ideas and findings, and are often used as a rough measure of a paper's importance. Fifty years ago, Eugene Garfield published the Science Citation Index (SCI), the first systematic effort to track citations in the scientific literature. To mark the anniversary, Nature asked Thomson Reuters, which now owns the SCI, to list the 100 most highly cited papers of all time. (See the full list at [www.nature.com/top100](http://www.nature.com/top100).) The search covered all of Thomson Reuters' Web of Science, an online version of the SCI that also includes databases covering the social sciences, arts and humanities, conference proceedings and some books. It lists papers published from 1900 to the present day.

The exercise revealed some surprises, not least that it takes a staggering 12,119 citations to rank in the top 100 — and that many of the world's most famous papers do not make the cut. A few that do, such as the first observation<sup>1</sup>

of carbon nanotubes (number 36) are indeed classic discoveries. But the vast majority describe experimental methods or software that have become essential in their fields.

The most cited work in history, for example, is a 1951 paper<sup>2</sup> describing an assay to determine the amount of protein in a solution. It has now gathered more than 305,000 citations — a recognition that always puzzled its lead author, the late US biochemist Oliver Lowry. "Although I really know it is not a great paper ... I secretly get a kick out of the response," he wrote in 1977.

The colossal size of the scholarly literature means that the top-100 papers are extreme outliers. Thomson Reuters' Web of Science holds some 58 million items. If that corpus were scaled to Mount Kilimanjaro, then the 100 most-cited papers would represent just 1 centimetre at the peak. Only 14,499 papers — roughly a metre and a half's worth — have more than 1,000 citations (see "The paper mountain"). Meanwhile, the foothills comprise works that have been cited only once, if at all — a group that encompasses roughly half of the items.

Nobody fully understands what distinguishes the sliver at the top from papers that are merely very well known — but researchers' customs explain some of it. Paul Wouters, director of the Centre for Science and Technology Studies in Leiden, the Netherlands, says that many methods papers "become a standard reference that one cites in order to make clear

to other scientists what kind of work one is doing". Another common practice in science ensures that truly foundational discoveries — Einstein's special theory of relativity, for instance — get fewer citations than they might deserve: they are so important that they quickly enter the textbooks or are incorporated into the main text of papers as terms deemed so familiar that they do not need a citation.

Citation counts are riddled with other confounding factors. The volume of citations has increased, for example — yet older papers have had more time to accrue citations. Biologists tend to cite one another's work more frequently than, say, physicists. And not all fields produce the same number of publications. Modern bibliometricians therefore recoil from methods as crude as simply counting citations when they want to measure a paper's value. Instead, they prefer to compare counts for papers of similar age, and in comparable fields.

Nor is Thomson Reuters' list the only ranking system available. Google Scholar compiled its own top-100 list for Nature. It is based on many more citations because the search engine calls references from a much greater (although poorly characterized) literature base, including from a large range of books. In that list, available at [www.nature.com/top100](http://www.nature.com/top100), economics papers have more prominence. Google Scholar's list also features books, which Thomson Reuters did not analyse. But among the science papers, many of the same titles show up.

Yet even with all the caveats, the old-fashioned hall of fame still has value. If nothing else, it serves as a reminder of the nature of scientific knowledge. To make exciting advances, researchers rely on relatively unused papers to describe experimental methods, databases and software.

Here Nature tours some of the key methods that tens of thousands of citations have hoisted to the top of science's Kilimanjaro — essential, but rarely thrust into the limelight.

### BIOLOGICAL TECHNIQUES

For decades, the top-100 list has been dominated by protein biochemistry. The 1951 paper<sup>2</sup> describing the Lowry method for quantifying protein remains practically unreachable at number 1, even though many biochemists say that it and the competing Bradford assay<sup>3</sup> — described by paper number 3 on the list — are a tad outdated. In between, at number 2, is Laemmli buffer<sup>4</sup>, which is used in a different kind of protein analysis. The dominance of these techniques is attributable to the high volume of citations in cell and molecular biology, where they remain indispensable tools.

At least two of the biological techniques described by top-100 papers have resulted in Nobel prizes. Number 4 on the list describes the DNA-sequencing method<sup>5</sup> that earned the late Frederick Sanger his share of the 1980 Nobel Prize in Chemistry. Number 63 describes polymerase chain reaction ▶

PHOTO BY PAUL BEAN, DESIGN BY MICHELE FERRELLI/NATURE

« But the vast majority describe experimental methods or software that have become essential in their fields. »

Source: van Noorden et al, Nature, 2014  
DOI: 10.1038/514550a

# Science ouverte: montre-moi ton code !



**JCIM** JOURNAL OF CHEMICAL INFORMATION AND MODELING

pubs.acs.org/jcim Viewpoint

## Code Sharing in the Open Science Era

W. Patrick Walters\*

 Cite This: *J. Chem. Inf. Model.* 2020, 60, 4417–4420  Read Online

**ACCESS** |  Metrics & More |  Article Recommendations

**ABSTRACT:** Many high-profile scientific journals have established policies mandating the release of code accompanying papers that describe computational methods. Unfortunately, the majority of journals that publish papers in Computational Chemistry and Cheminformatics have yet to define such guidelines. This Viewpoint reviews the current state of reproducibility for the field and makes a case for the inclusion of code with computational papers.



Source: Walters, JCIM, 2020, DOI: 10.1021/acs.jcim.0c01000

# Les 3 piliers de la science ouverte



- Libre accès aux publications scientifiques (*open access*)
- Données ouvertes (*open data*)
- Codes logiciels ouverts (*open source*)



# L'essor des gestionnaires de versions



## "FINAL".doc



FINAL.doc!



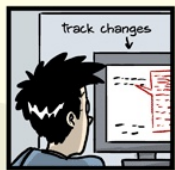
FINAL\_rev.2.doc



FINAL\_rev.6.COMMENTS.doc



FINAL\_rev.8.comments5.  
CORRECTIONS.doc



FINAL\_rev.18.comments7.  
corrections9.MORE.30.doc



FINAL\_rev.22.comments49.  
corrections.10.#@\$%WHYDID  
ICOMETOGRADSCHOOL?????.doc



Source : PhD Comics  
JORGE CHAM © 2012

WWW.PHDCOMICS.COM

OPINION ARTICLE



## Four simple recommendations to encourage best practices in research software [version 1; peer review: 3 approved]

✉ Rafael C. Jiménez [id](#)<sup>1</sup>, ✉ Mateusz Kuzak<sup>2</sup>, Monther Alhamdoosh [id](#)<sup>3</sup>, Michelle Barker<sup>4</sup>, Bérénice Batut [id](#)<sup>5</sup>, Mikael Borg<sup>6</sup>, Salvador Capella-Gutierrez [id](#)<sup>7</sup>, Neil Chue Hong<sup>8</sup>, Martin Cook<sup>1</sup>,

Source: Jiménez et al, F1000 Research, 2017  
DOI: 10.12688/f1000research.11407.1

## PLOS BIOLOGY

OPEN ACCESS

COMMUNITY PAGE

## Best Practices for Scientific Computing

Greg Wilson , D. A. Aruliah, C. Titus Brown, Neil P. Chue Hong, Matt Davis, Richard T. Guy, Steven H. D. Haddock, Kathryn D. Huff, Ian M. Mitchell, Mark D. Plumbley, Ben Waugh, Ethan P. White, Paul Wilson

Published: January 7, 2014 • <https://doi.org/10.1371/journal.pbio.1001745>

Source: Wilson et al, PLOS Biology, 2014  
DOI: 10.1371/journal.pbio.1001745

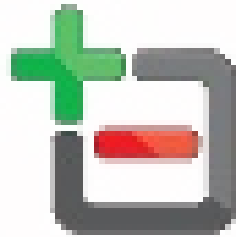
# et des plateformes de développement gratuites



**GitHub**



**GitLab**



**GITORIOUS**

**ATLASSIAN**



**Bitbucket**

**Google**<sup>™</sup>  
code

# et des plateformes de développement gratuites disparues



## Google Kills Off ... Code

Natasha Lomas @riptari / 10 ... ch 13, 2015

Source: TechCrunch

1,4 million de projects

## Sunsetting Mercurial support in Bitbucket

April 21, 2020 | 3 min read



Denise Chan

[Update Aug 26, 2020] ... have now been disabled and cannot be accessed.

[Update July 8, 2020] ... today, mercurial repositories, snippets, and wikis will turn to read-only mode. After July 8th, 2020 they will no longer be accessible.

250 000 dépôts

Source: BitBucket blog

# La reproductibilité en science nécessite un accès au code source sur le long terme



Héberger votre code sur une plateforme publique et gratuite est acceptable.  
Mais vous devez vous préparer à la fermeture de cette plateforme !



# Archiver du code comme des données ?



## Zenodo

OpenAIRE + CERN

## Figshare

Digital Science

Code citable avec un DOI  
Archivage automatisable depuis GitHub

# Archivez votre code sur



# Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Software Heritage archive tous les codes sources ouverts,  
pour toujours et gratuitement

- Organisation à but non lucratif
- Créée en 2016 à l'INRIA  
(Roberto Di Cosmo & Stefano Zacchiroli)
- Financée par l'UNESCO, le CNRS, UPC,  
Microsoft, Google, Huawei, Intel...



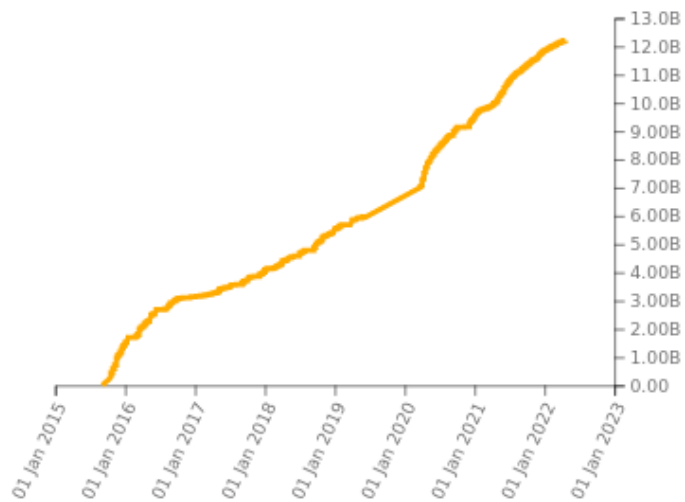
Source: Software Heritage 5 years anniversary (2021)

# Une archive conséquente



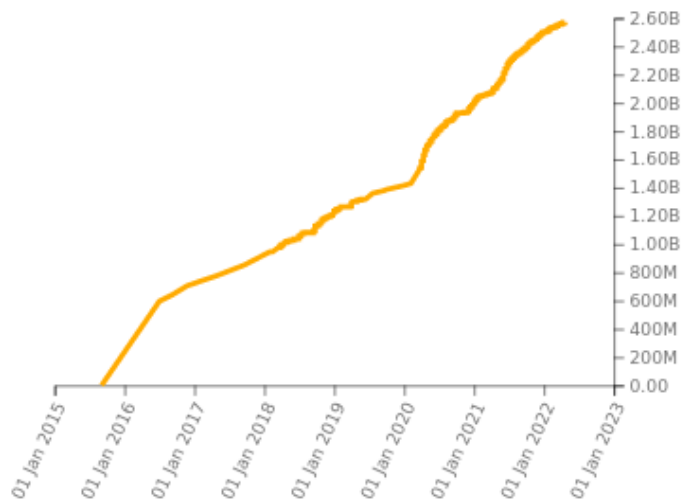
## Source files

12 204 306 258



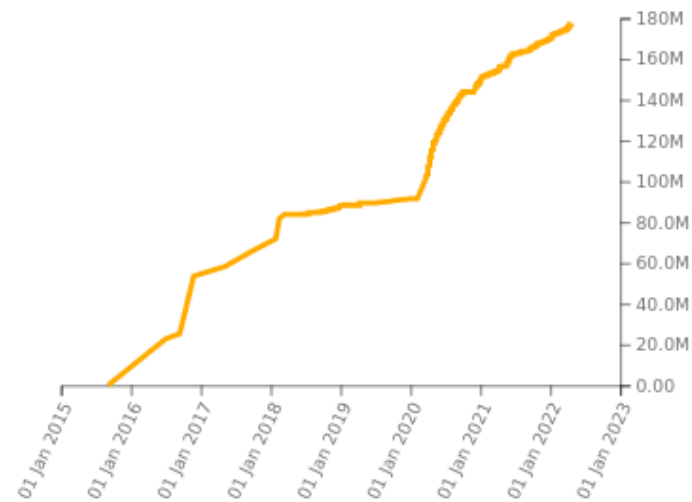
## Commits

2 570 865 678



## Projects

177 280 343



Source: Software Heritage

# qui intègre déjà de nombreuses « sources »



 Bitbucket

2,074,834 origins



git

22,235 origins





19,841 origins



 debian

126,944 origins



5,875 origins



GitHub

132,824,442 origins



 GitLab

4,082,399 origins



 Guix


11,821 origins



 GNU

354 origins



 heptapod

1,039 origins



 launchpad

20,417 origins



 NixOS

11,821 origins





1,802,916 origins



4,083 origins



462,245 origins



 SOURCEFORGE

313,582 origins



Source: Software Heritage



# Software Heritage sauve les codes sources



## Discontinued hosting

Discontinued hosting services. Those origins have been archived by Software Heritage.



122,014 origins



790,026 origins



336,795 origins



Source: Software Heritage


« Google Code content now safely collected », 2016


« Rescuing 250000+ endangered Mercurial repositories », 2020


# Sauvegardez votre code maintenant !



<https://archive.softwareheritage.org/save/>





 Software Heritage Archive

 **Save code now**




Enter a SWHID to resolve or keyword(s) to search for in origin URLs 


---

Features

-  Search
-  Downloads
-  **Save code now**
-  Help

You can contribute to extend the content of the Software Heritage archive by submitting an origin save request. To do so, fill the required info in the form below:

Origin type	Origin url	
git 	<input type="text" value=""/> 	Submit 

 [Browse save requests](#)

A "Save code now" request takes the following parameters:

# Sauvegarder mon code



patrickfuchs / buildH Public Unwatch 4 Fork 5 Starred 8

<> Code Issues 4 Pull requests Actions Projects Wiki Security Insights

master 1 branch 14 tags Go to file Add file Code

**patrickfuchs** Fix test UI ✓ 6429000 22 days ago 615 commits

.github/workflows	Simplify workflow file and add Python 3.9 for CI	8 months ago
binder	Fix environment.yml for pip	8 months ago
buildh	Make clearer that the topology is supplied in the json file	22 days ago
def_files	Merge pull request #139 from patrickfuchs/add_DPPC_CHARMM3...	10 months ago
devtools	Bump version: 1.6.0 → 1.6.1	3 months ago
docs	Add an example of ignore_CH3s to Notebook04	3 months ago
paper	Fix typos + add PS lipid	7 months ago
tests	Fix test UI	22 days ago
.gitignore	Update .gitignore	12 months ago
.readthedocs.yml	Update Read the Docs configuration (automatic)	3 months ago
.zenodo.json	Bump version: 1.6.0 → 1.6.1	3 months ago
AUTHORS	Update authors order	10 months ago
CHANGELOG.md	Bump version: 1.6.0 → 1.6.1	3 months ago
CODE_OF_CONDUCT.md	Update contact e-mail	10 months ago
CONTRIBUTING.md	Fix menu & add credit to Durand's blog post	8 months ago
LICENSE.txt	Update license text	11 months ago
Makefile	Update documentation and Makefile for PyPI package	12 months ago

**About**

Build hydrogen atoms from united-atom molecular dynamics of lipids and calculate the order parameters.

[buildh.readthedocs.io/](https://buildh.readthedocs.io/)

python molecular-dynamics-simulation order-parameters lipids united-atom

Readme BSD-3-Clause License Code of conduct 8 stars 4 watching 5 forks

**Releases** 8

v1.6.1 Latest on 20 Jan + 7 releases

**Contributors** 5

<https://github.com/patrickfuchs/buildH>

# Sauvegarder mon code



Save code now

Enter a SWHID to resolve or keyword(s) to search for in origin URLs



You can contribute to extend the content of the Software Heritage archive by submitting an origin save request. To do so, fill the required info in the form below:

Origin type

Origin url

git



https://github.com/patrickfuchs/buildH|

Submit

# Sauvegarder mon code



<https://github.com/patrickfuchs/buildH>

20 April 2022, 08:31 UTC

<> Code   Branches (2)   Releases (13)   Visits

★ Branch: HEAD   e433ce9 /   History   Download   Save again   Permalinks

Tip revision: 6429080a1c1deb608f8c6dacc5a83b1b72c58c77 authored by patrickfuchs on 29 March 2022, 14:19 UTC

Fix test UI

File	Mode	Size
└─ .github		
└─ binder		
└─ buildh		
└─ def_files		
└─ devtools		
└─ docs		
└─ paper		
└─ tests		
└─ .gitignore	-rw-r--r--	376 bytes
└─ .readthedocs.yml	-rw-r--r--	219 bytes
└─ .zenodo.json	-rw-r--r--	812 bytes
└─ AUTHORS	-rw-r--r--	177 bytes

[https://archive.softwareheritage.org/browse/origin/directory/?origin\\_url=https://github.com/patrickfuchs/buildH](https://archive.softwareheritage.org/browse/origin/directory/?origin_url=https://github.com/patrickfuchs/buildH)

**GOOD NEWS EVERYONE**



**SOFTWARE HERITAGE**

**ARCHIVES ANY OPEN-SOURCE CODE!**

# Bonnes pratiques : au-delà du code



## Metadonnées pour les humains

- **README**  
<https://readme.so/fr/editor>
- **AUTHORS**  
Ada Lovelance <ada@programming.org>  
Margaret Hamilton <margaret@nasa.com>
- **LICENSE**  
Licence ouverte compatible **SPDX**  
<https://choosealicense.com/>  
<https://reuse.software/>

## Metadonnées pour les machines

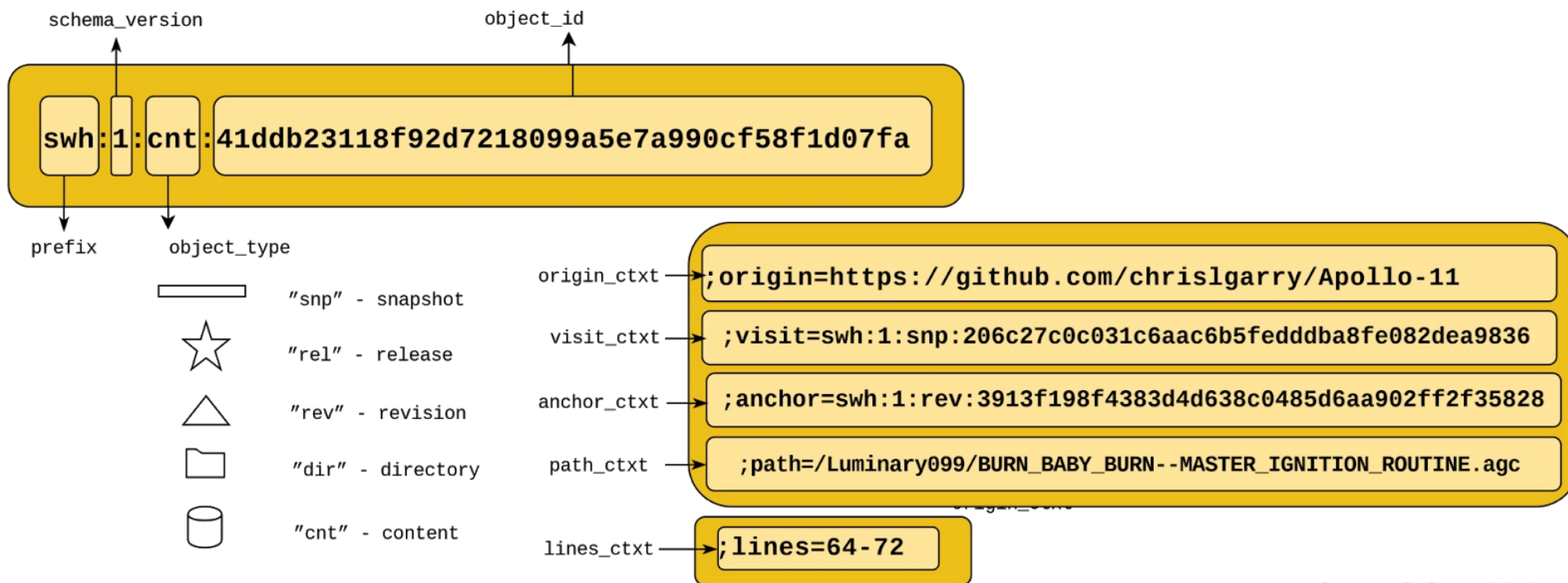
- **codemeta.json**  
avec un générateur pour  
les humains : CodeMeta Generator

Voir aussi « [HOWTO archive and reference your code](#) »

# Référencer son code source



Le ~~DOI~~ SWHID : un identifiant intrinsèque et persistant






Source: Software Heritage








# Référencer son code source




 <https://github.com/patrickfuchs/buildH> 

 20 April 2022, 08:31 UTC

 Code  Branches (2)  Releases (13)  Visits




 Branches



 Fix

**Permalinks**



To reference or cite the objects present in the Software Heritage archive, permalinks based on SoftWare Heritage persistent IDentifiers (SWHIDs) 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 SWHID and permalink.



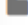





 directory  revision  snapshot

 archived  archived `swh:1:dir:e433ce9c1846a222446c8f22df7b1f4526a53033` Iframe embedding

```
swh:1:dir:e433ce9c1846a222446c8f22df7b1f4526a53033;
origin=https://github.com/patrickfuchs/buildH;
visit=swh:1:snp:6ea443e1de87de3d85cee62ba706a58b130dc00e;
anchor=swh:1:rev:6429080a1c1deb608f8c6dacc5a83b1b72c58c77
```

Add contextual information  Copy identifier  Copy permalink

 File

-  .github
-  binder
-  buildh
-  def\_files
-  devtools
-  docs
-  paper
-  tests

# Référencer son code source



L'adresse du dépôt sur GitHub :

<https://github.com/patrickfuchs/buildH>

La référence de l'archive dans SWH pour le README :



[https://archive.softwareheritage.org/browse/origin/directory/?origin\\_url=https://github.com/patrickfuchs/buildH](https://archive.softwareheritage.org/browse/origin/directory/?origin_url=https://github.com/patrickfuchs/buildH)

La référence de l'archive dans SWH pour un article scientifique, avec une version spécifique :

```
swh:1:dir:e433ce9c1846a222446c8f22df7b1f4526a53033;  
origin=https://github.com/patrickfuchs/buildH;  
visit=swh:1:snp:6ea443e1de87de3d85cee62ba706a58b130dc00e;  
anchor=swh:1:rev:6429080a1c1deb608f8c6dacc5a83b1b72c58c77
```

# Référencer son code source



## buildH: Build hydrogen atoms from united-atom molecular dynamics of lipids and calculate the order parameters

Python Submitted 08 July 2021 • Published 19 September 2021

Background  
Statement of Need  
Overview  
Acknowledgements  
References

Some notebooks are provided in the GitHub repository to explain how buildH works and how to analyze the data produced. In case of trouble, any user can post an issue on GitHub.

buildH is available in the Python Package Index (PyPI) as well as in the Bioconda repository. The current version 1.6.0 of buildH is archived in the Zenodo repository (<https://zenodo.org/record/5356246>) and in the Software Heritage archive ([swh:1:dir:4c63d5ca3497726a1e54ac152ce1667d7c004d2b](https://archive.softwareheritage.org/swh:1:dir:4c63d5ca3497726a1e54ac152ce1667d7c004d2b)).

### Acknowledgements

The authors thank the community of [NMRlipids](#) for useful discussions.

### References

Abraham, M. J., Murtola, T., Schulz, R., Páll, S., Smith, J. C., Hess, B., & Lindahl, E. (2015). GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX*, 1-2, 19–25. <https://doi.org/10.1016/j.softx.2015.06.001>

Antila, H., Buslaev, P., Favela-Rosales, F., Ferreira, T. M., Gushchin, I., Javanainen, M., Kav, B., Madsen, J. J., Melcr, J., Miettinen, M. S., Määttä, J., Nencini, R., Ollila, O. H. S., & Piggot, T. J. (2019). Headgroup structure and cation binding in phosphatidylserine lipid bilayers. *The Journal of Physical Chemistry B*, 123(43), 9066–9079. <https://doi.org/10.1021/acs.jpcc.9b06000>

<https://archive.softwareheritage.org/swh:1:dir:4c63d5ca3497726a1e54ac152ce1667d7c004d2b;origin=https://github.com/patrickfuchs/buildH;/visit=swh:1:snp:a63a8d07dbebeb442a06707be476817cec44ac72;anchor=swh:1:rev:9f05672515e1cdb0064eeb34f63844296193bc0d>

Software repository  
Paper review  
Download paper

Software archive  
Editor: [@richardjgowers](#) (all papers)  
Reviewers: [@lilyminium](#) (all reviews), [@blakeaw](#) (all reviews)

### Authors

Hubert Santuz (0000-0001-6149-9480)  
Bacle (0000-0002-3317-9110), Pierre Poulain (0000-0003-4177-3619), Patrick F. J. ...

<https://joss.theoj.org/papers/10.21105/joss.03521>

# Citer son code source



BibLaTeX style extension for software

[Software Release] B. Langmead and S. L. Salzberg, *Bowtie2* version 2.4.2, Oct. 2022. LIC: GPL. URL: <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>, VCS: <https://github.com/BenLangmead/bowtie2>, SWHID: `<swh:1:rel:97bacffea6e7c3f574ce5b566daba82aa18a11f;origin=https://github.com/BenLangmead/bowtie2;visit=swh:1:snp:c25778cfefc086c63c6f78eed230d0b9c88876ee>`.

[Software excerpt] MIT Instrumentation Laboratory, “AGC Luminary routine for changing LEM asset during landing”, from *Apollo 11 Guidance Computer (AGC) source code for the command and lunar module* 1967. VirtualAGC project. LIC: Public Domain. URL: <https://www.ibiblio.org/apollo>, VCS: <https://github.com/virtualagc/virtualagc>, SWHID: `<swh:1:cnt:64582b78792cd6c2d67d35da5a11bb80886a6409;origin=https://github.com/virtualagc/virtualagc;anchor=swh:1:rev:007c2b95f301f9438b8b74d7993b7a3b9a66255b;lines=245-261>`.

# Citer son code source



hal-03467422, version 1

## buildH

Hubert Santuz<sup>1,2</sup>, Amélie Bâcle<sup>3</sup>, Pierre Poulain<sup>4</sup>, Patrick Fuchs<sup>5,6</sup> [Détails](#)

- 1 LBT (UPR\_9080) - Laboratoire de biochimie théorique [Paris]
- 2 IBPC (FR\_550) - Institut de biologie physico-chimique
- 3 LitCh - Lipotoxicity and Channelopathies - ConicMeds
- 4 IJM (UMR\_7592) - Institut Jacques Monod
- 5 LBM UMR 7203 - Laboratoire des biomolécules
- 6 UPC - UFR SDV - Université Paris Cité - UFR Sciences du Vivant [Sciences]


**Abstract :** Build hydrogen atoms from united-atom molecular dynamics of lipids and calculate the order parameters

Type de document : [Logiciel](#)

Domaine : [Informatique \[cs\]](#) / [Bio-informatique \[q-bio.QM\]](#)

Liste complète des métadonnées [Voir](#)

### CONSULTER

 Software Heritage `swh:1:dir:4c63d5ca3497726a1e54ac152ce1667d7c004d2b;origin=https://github.com/patrickfuchs/buildH;/visit=swh:1:snp:a63a8d07d8eb442a06707be476817cec44ac72;anchor=swh:1:rev:9f05672515e1cdb0064eeb34f63844296193bc0d`

[Consulter](#)

<https://hal.archives-ouvertes.fr/hal-03467422>  
Contributeur : [Pierre Poulain](#) [Contacter le contributeur](#)  
Soumis le : lundi 6 décembre 2021 - 15:02:04  
Dernière modification le : mercredi 20 avril 2022 - 16:16:23

### MÉTADONNÉES

Keywords : [lipids](#) [order parameters](#)  
[molecular dynamics simulation](#) [united-atom](#)

version [1.6.0](#)

Licences <https://spdx.org/licenses/BSD-3-Clause>

Langage de programmation [Python](#)

Code Repository <https://github.com/patrickfuchs/buildH>

### COLLECTIONS

ESPCI | STIM | ENSCP | UP-SCIENCES | PSL | LBM | ENS-PARIS | UNIV-POITIERS | PARISTECH | SORBONNE-UNIVERSITE | UNIV-TOURS | SU-SCIENCES | LBT | ESPCI-PSL | IJM | CNRS | INC-CNRS | INSERM | UNIV-PARIS

### CITATION

Hubert Santuz, Amélie Bâcle, Pierre Poulain, Patrick Fuchs. buildH. 2019, (swh:1:dir:4c63d5ca3497726a1e54ac152ce1667d7c004d2b;origin=https://github.com/patrickfuchs/buildH;/visit=swh:1:snp:a63a8d07d8eb442a06707be476817cec44ac72;anchor=swh:1:rev:9f05672515e1cdb0064eeb34f63844296193bc0d). (hal-03467422)

### EXPORTER

[CodeMeta](#) [BibTeX](#) [TEI](#) [DC](#) [DCterms](#)  
[EndNote](#) [Datacite](#)

<https://hal.inria.fr/hal-03467422>

# En résumé



Archivez

<https://archive.softwareheritage.org/save/>



Décrivez avec des métadonnées

`README`, `LICENSE`, `AUTHORS`, `codemeta.json`



Référez

SWHID plutôt que DOI, contexte



Citez

Version, release, fichier, ligne



Merci !



The first 5 years of Software Heritage in 5 minutes!

Software Heritage FAQ

Open Science tutorial: source code deposit (SWH + HAL)