

Software Heritage

Large-Scale Research on Public Code Development

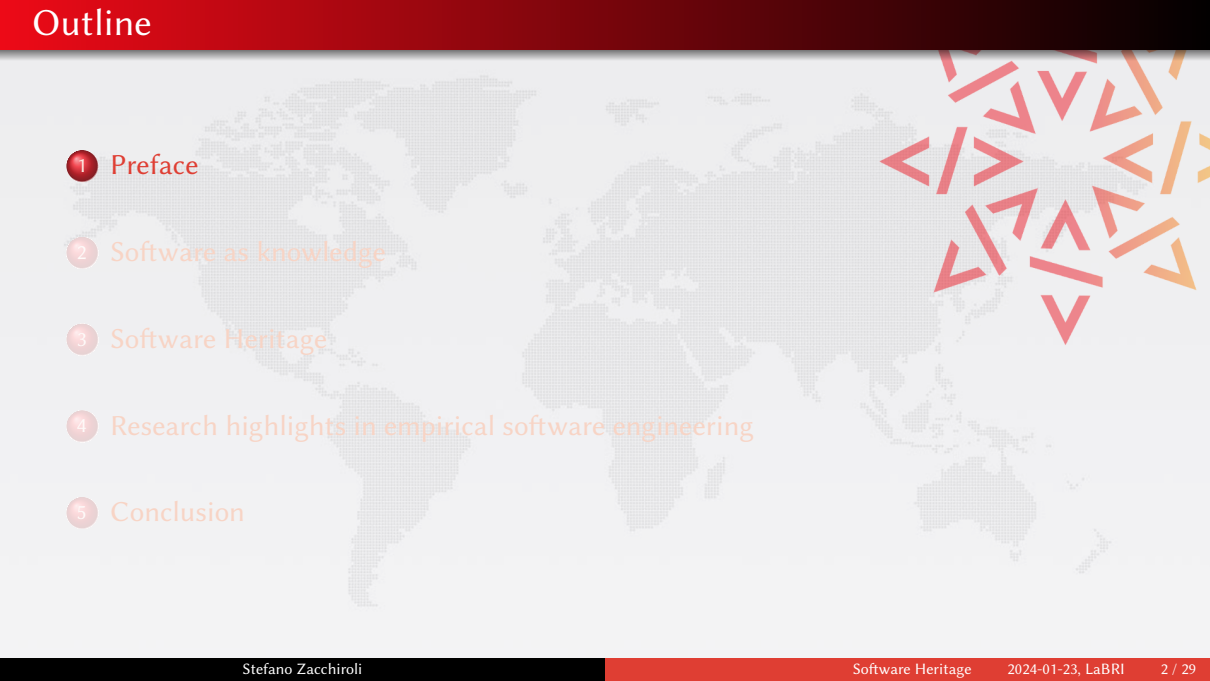
Stefano Zacchioli

Télécom Paris, Institut Polytechnique de Paris
`stefano.zacchioli@telecom-paris.fr`

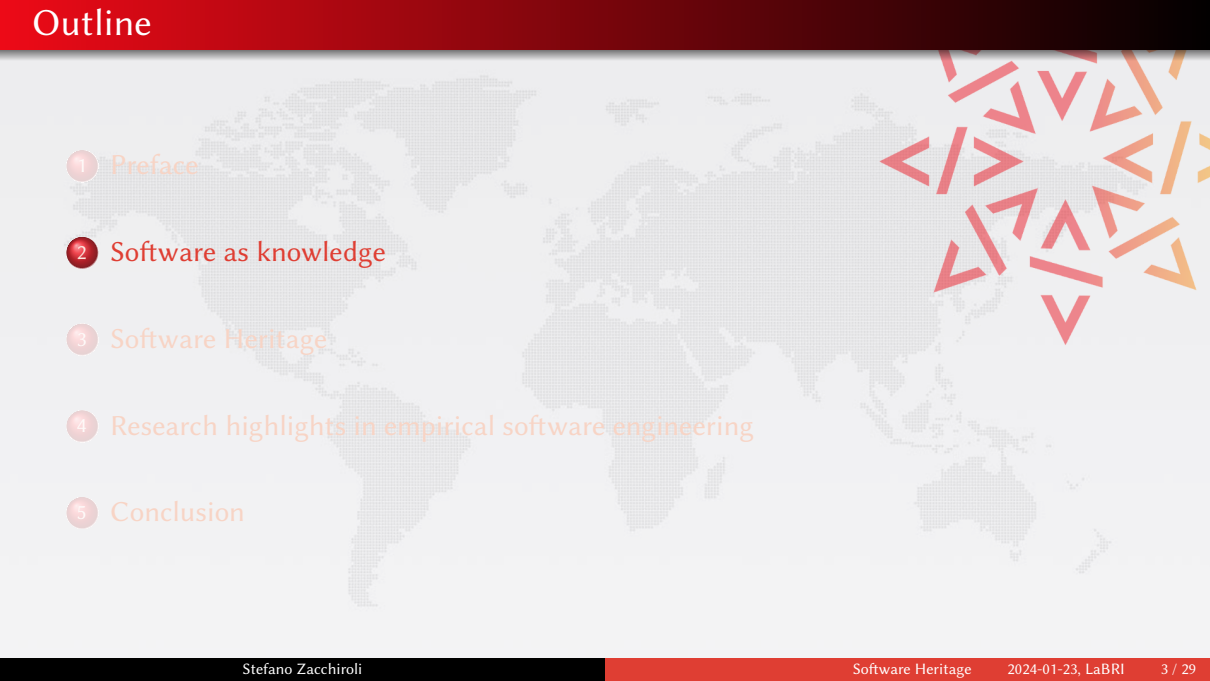
23 Jan 2024 — LaBRI, Bordeaux



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

- 
- 1 Preface
 - 2 Software as knowledge
 - 3 Software Heritage
 - 4 Research highlights in empirical software engineering
 - 5 Conclusion

- Professor of Computer Science, Télécom Paris, Institut Polytechnique de Paris
- Free/Open Source Software activist (20+ years)
- Debian Developer & Former 3x Debian Project Leader
- Former Open Source Initiative (OSI) director
- Software Heritage co-founder & CTO

- 
- 1 Preface
 - 2 Software as knowledge
 - 3 Software Heritage
 - 4 Research highlights in empirical software engineering
 - 5 Conclusion

Software is dual-form knowledge



"The source code for a work means the preferred form of the work for making modifications to it."

GPL Licence

Software is dual-form knowledge



"The source code for a work means the preferred form of the work for making modifications to it."

GPL Licence

Hello World

Software is dual-form knowledge



"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
```

Software is dual-form knowledge



"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");
}
```


Software *source code* is precious human 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.”

Software *source code* is precious human 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)

```
P63SP0T3      CA      BIT6          # IS THE LR ANTENNA IN POSITION 1 YET
EXTEND
RAND      CHAN33
EXTEND
BZF      P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

CAF      CODE500      # ASTRONAUT:   PLEASE CRANK THE
TC      BANKCALL      #                   SILLY THING AROUND
CADR      GOPERF1
TCF      GOTOP00H      # TERMINATE
TCF      P63SP0T3      # PROCEED     SEE IF HE'S LYING

P63SP0T4      TC      BANKCALL      # ENTER       INITIALIZE LANDING RADAR
CADR      SETPOS1

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

Software *source code* is precious human 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)

```
P63SP0T3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND   CHAN33
              EXTEND
              BZF    P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF    CODE500      # ASTRONAUT:  PLEASE CRANK THE
              TC     BANKCALL      #              SILLY THING AROUND
              CADR   GOPERF1
              TCF    GOTOP00H      # TERMINATE
              TCF    P63SP0T3      # PROCEED    SEE IF HE'S LYING

P63SP0T4      TC     BANKCALL      # ENTER      INITIALIZE LANDING RADAR
              CADR   SETPOS1

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

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;
}
```

Software *source code* is precious human 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)

```
P63SP0T3      CA      BIT6      # IS THE LR ANTENNA IN POSITION 1 YET
              EXTEND
              RAND    CHAN33
              EXTEND
              BZF     P63SP0T4      # BRANCH IF ANTENNA ALREADY IN POSITION 1

              CAF     CODE500      # ASTRONAUT:  PLEASE CRANK THE
              TC      BANKCALL     #              SILLY THING AROUND
              CADR    GOPERF1
              TCF     GOTOP00H     # TERMINATE
              TCF     P63SP0T3     # PROCEED   SEE IF HE'S LYING

P63SP0T4      TC      BANKCALL     # ENTER      INITIALIZE LANDING RADAR
              CADR    SETPOS1

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

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;
}
```

Len Shustek, *Computer History Museum*

2006

“Source code provides a view into the mind of the designer.”



A word cloud of terms related to software fragility, including: damage, disaster, malicious, attack, obsolete, dependencies, deletion, reference, storage, dangling, wear, corruption, encryption, and format. The words are arranged in a circular pattern with varying sizes and colors.

Like all digital information, FOSS is fragile

- link rot: projects are created, moved around, removed
- business-driven code loss (e.g., Gitorious, Google Code, Bitbucket)
- data rot: physical media with legacy software decay

Software source code is fragile



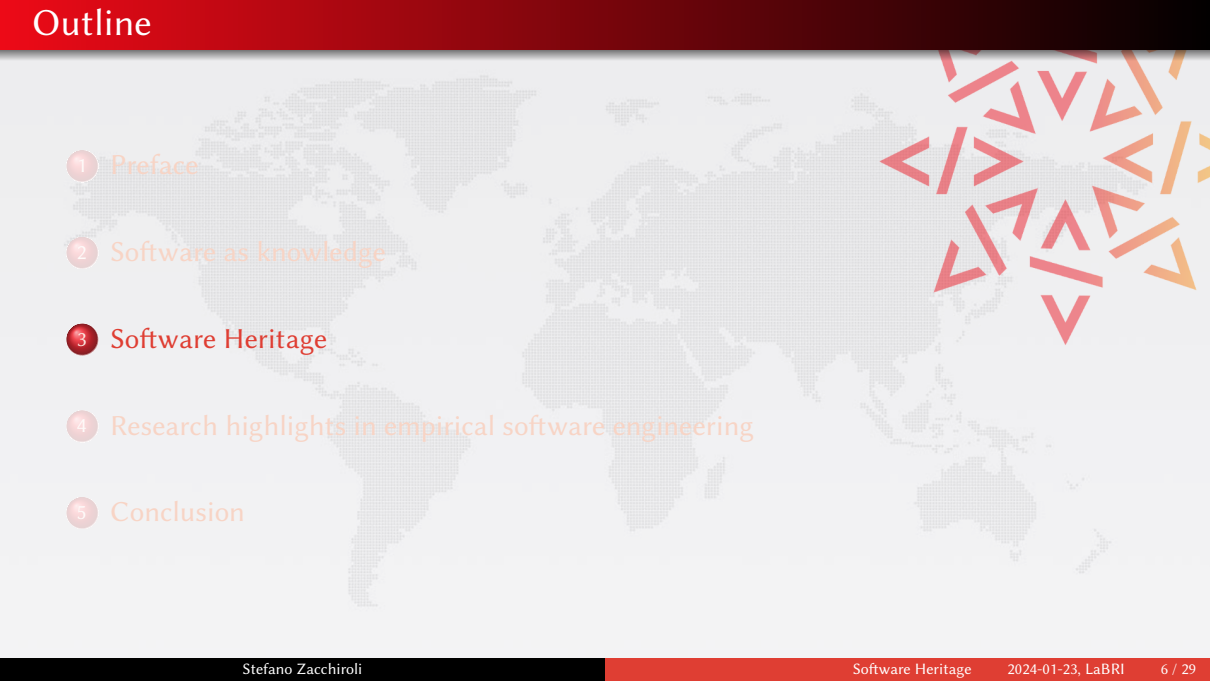
A word cloud of terms related to software fragility, including: damage, disaster, malicious, attack, obsolete, dependencies, deletion, reference, storage, dangling, wear, corruption, encryption, format, aging, media, and tear.

Like all digital information, FOSS is fragile

- link rot: projects are created, moved around, removed
- business-driven code loss (e.g., Gitorious, Google Code, Bitbucket)
- data rot: physical media with legacy software decay

If a website disappears you go to the Internet Archive...

where do you go if (a repository on) GitHub or GitLab goes away?

- 
- 1 Preface
 - 2 Software as knowledge
 - 3 Software Heritage
 - 4 Research highlights in empirical software engineering
 - 5 Conclusion



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



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



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



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

Archiving goals

Targets: VCS repositories & source code releases (e.g., tarballs, packages)

We DO archive

- file **content** (= blobs)
- **revisions** (= commits), with full metadata
- **releases** (= tags), ditto
- where (**origin**) & when (**visit**) we found any of the above

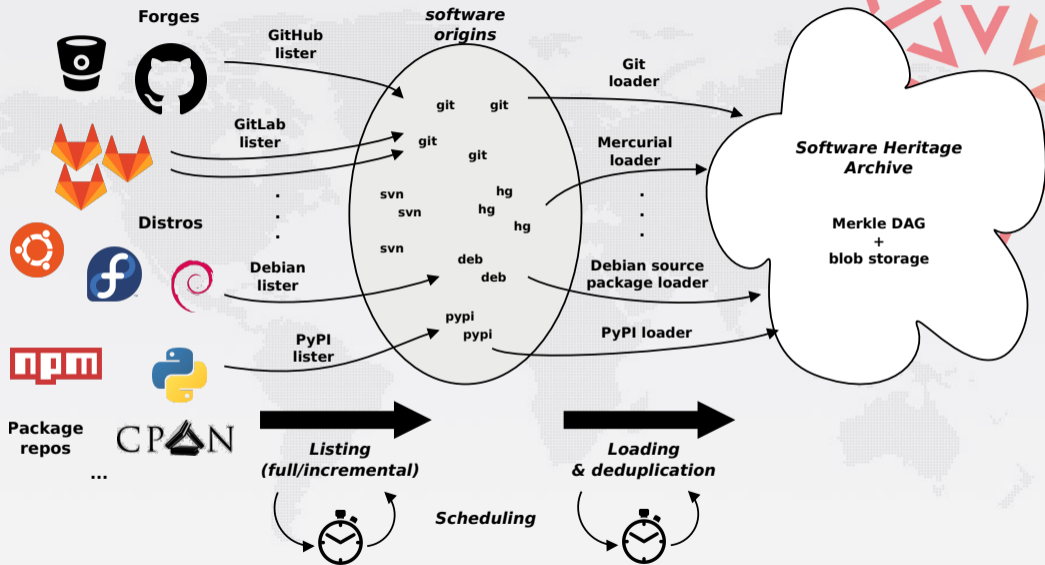
... in a VCS-/archive-agnostic **canonical data model**

We DON'T archive (yet)

- homepages, wikis
- BTS/issues/code reviews/etc.
- mailing lists

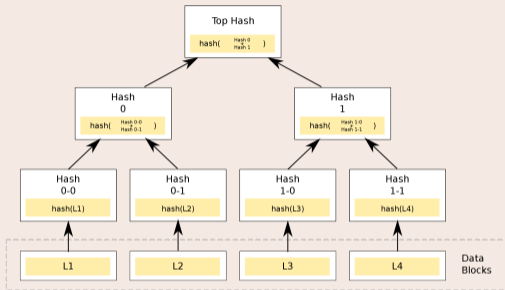
Long term vision: play our part in a *"semantic wikipedia of software"*

Data flow



Merkle trees

Merkle tree (R. C. Merkle, CRYPTO 1987)

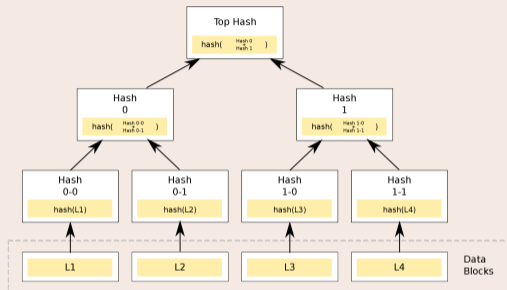


Combination of

- tree
- hash function

Merkle trees

Merkle tree (R. C. Merkle, CRYPTO 1987)

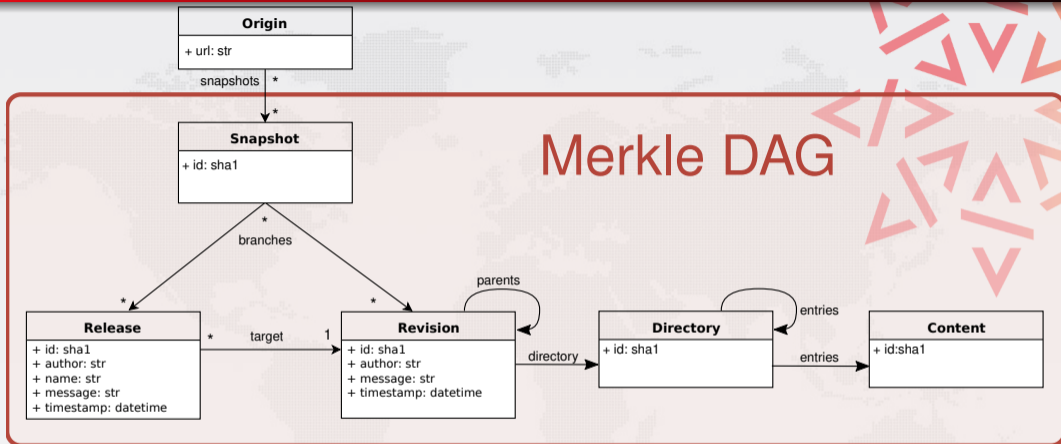


Combination of

- tree
- hash function

Classical cryptographic construction

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



A **global graph** linking together fully **deduplicated** source code artifact (files, commits, directories, releases, etc.) to the places that distribute them (e.g., Git repositories), providing a **unified view** on the entire *Software Commons*.

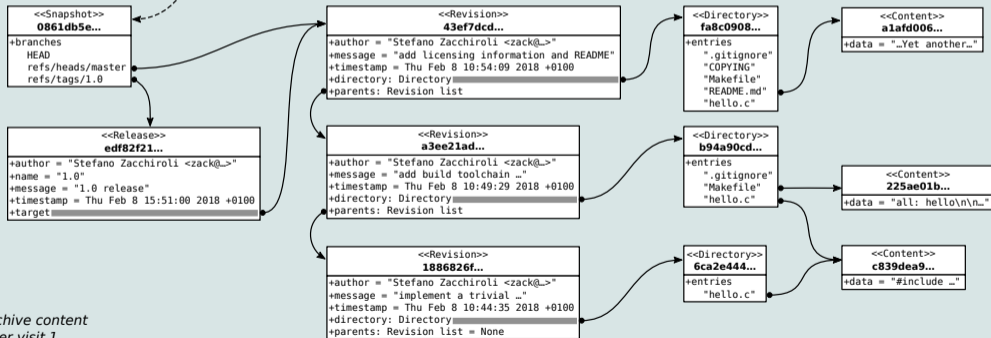
The archive: a (giant) Merkle DAG

origin
https://forge.softwareheritage.org/source/helloworld.git

visit
1

snapshot
0861db5e...

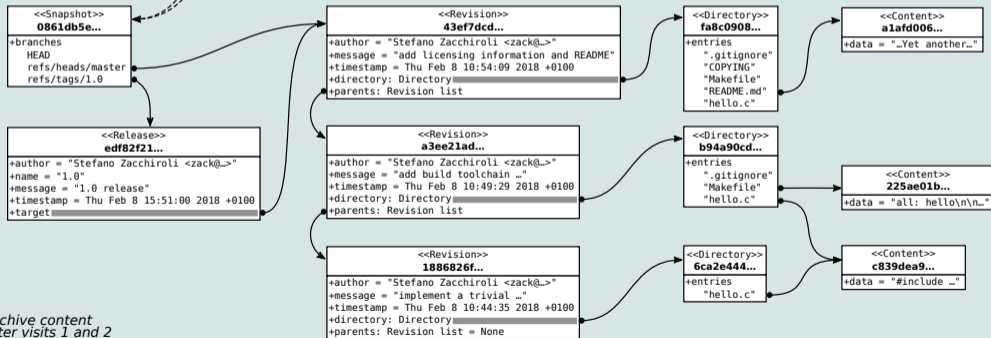
timestamp
Fri Feb 9 12:38:45 2018 +0100



Archive content
after visit 1

The archive: a (giant) Merkle DAG

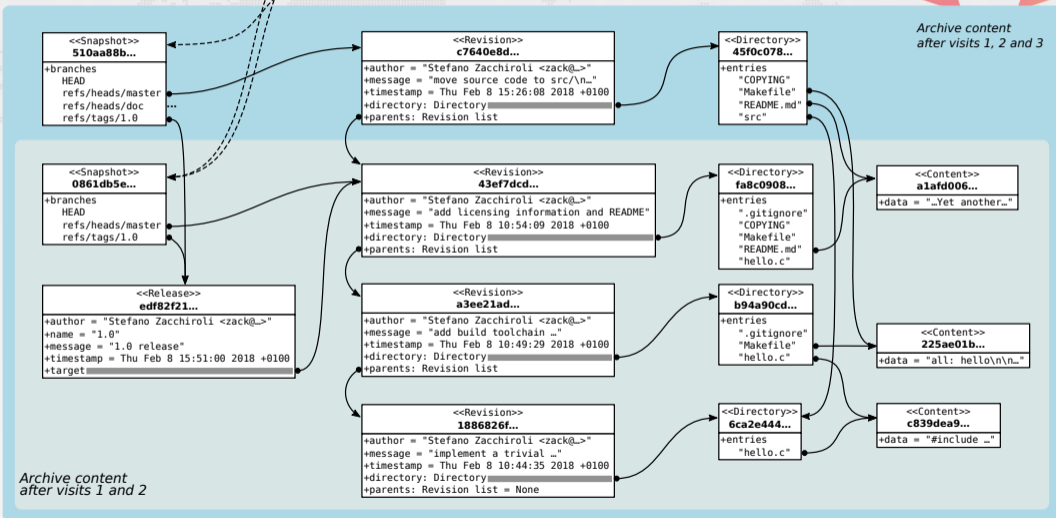
origin	visit	snapshot	timestamp
https://forge.softwareheritage.org/source/helloworld.git	1	0861db5e...	Fri Feb 9 12:38:45 2018 +0100
https://forge.softwareheritage.org/source/helloworld.git	2	0861db5e...	Fri Feb 9 13:29:00 2018 +0100

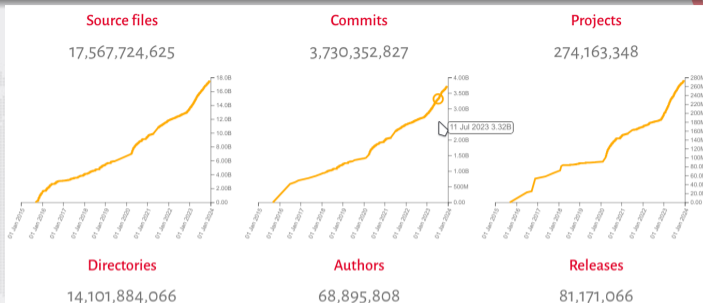


Archive content
after visits 1 and 2

The archive: a (giant) Merkle DAG

origin	visit	snapshot	timestamp
https://forge.softwareheritage.org/source/helloworld.git	1	0861db5e...	Fri Feb 9 12:38:45 2018 +0100
https://forge.softwareheritage.org/source/helloworld.git	2	0861db5e...	Fri Feb 9 13:29:00 2018 +0100
https://forge.softwareheritage.org/source/helloworld.git	3	510aa88b...	Fri Feb 9 15:52:50 2018 +0100








- on disk: ~1 PiB; as a graph ~35 B nodes, ~500 B edges
- the largest public source code archive in the world (and growing!)

- Browse [the archive](#)
- [Trigger archival](#) of your preferred software in a breeze
- Get and use SWHIDs ([full specification available online](#))
- The [Apollo 11 AGC source code example](#)
- Cite software [with the biblatex-software style](#) from CTAN
- Example use in a research article: compare Fig. 1 and conclusions
 - in [the 2012 version](#)
 - in [the updated version](#) using SWHIDs and Software Heritage
- Example in a journal: [an article from IPOL](#)
- [Curated deposit in SWH via HAL](#), see for example: [LinBox](#), [SLALOM](#), [Givaro](#), [NS2DDV](#), [SumGra](#), [Coq proof](#), ...
- Rescue landmark legacy software, see the [SWHAP process with UNESCO](#)

- 
- 1 Preface
 - 2 Software as knowledge
 - 3 Software Heritage
 - 4 Research highlights in empirical software engineering
 - 5 Conclusion

Graph dataset

Use case: large scale analyses of the most comprehensive corpus on the development history of free/open source software.



Antoine Pietri, Diomidis Spinellis, Stefano Zacchiroli

The Software Heritage Graph Dataset: Public software development under one roof

MSR 2019: 16th Intl. Conf. on Mining Software Repositories. IEEE

preprint: <http://deb.li/swhmsr19>

Dataset

- Relational representation of the full graph as a set of tables
- Available as open data: docs.softwareheritage.org/devel/swh-dataset/graph
- Chosen as subject for the **MSR 2020 Mining Challenge**

Formats

- Local use: set of Apache ORC files (10+ TiB in total)
- Live usage: Amazon Athena (SQL-queriable), Azure Data Lake


```
SELECT COUNT(*) AS c, word FROM (  
  SELECT LOWER(REGEXP_EXTRACT(FROM_UTF8(  
    message), '^w+')) AS word FROM revision)  
WHERE word != ''  
GROUP BY word ORDER BY COUNT(*) DESC LIMIT 5;
```

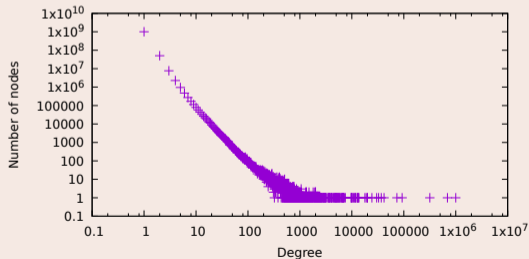
```
SELECT COUNT(*) AS c, word FROM (  
  SELECT LOWER(REGEXP_EXTRACT(FROM_UTF8(  
    message), '^\\w+')) AS word FROM revision)  
WHERE word != ''  
GROUP BY word ORDER BY COUNT(*) DESC LIMIT 5;
```

Count	Word
71 338 310	update
64 980 346	merge
56 854 372	add
44 971 954	added
33 222 056	fix

Fork arity

i.e., how often is a commit based upon?

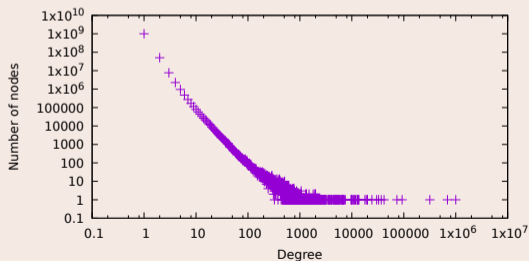
```
SELECT fork_deg, count(*) FROM (  
  SELECT id, count(*) AS fork_deg  
  FROM revision_history GROUP BY id) t  
GROUP BY fork_deg ORDER BY fork_deg;
```



Fork arity

i.e., how often is a commit based upon?

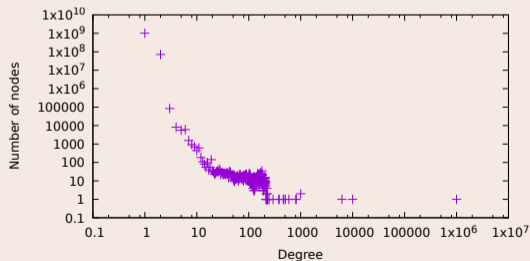
```
SELECT fork_deg, count(*) FROM (  
  SELECT id, count(*) AS fork_deg  
  FROM revision_history GROUP BY id) t  
GROUP BY fork_deg ORDER BY fork_deg;
```



Merge arity

i.e., how large are merges?

```
SELECT merge_deg, COUNT(*) FROM (  
  SELECT parent_id, COUNT(*) AS merge_deg  
  FROM revision_history GROUP BY parent_id)  
GROUP BY merge_deg ORDER BY merge_deg;
```





Stefano Zacchiroli

A Large-scale Dataset of (Open Source) License Text Variants

MSR 2022 (best dataset paper) + Empir. Soft. Eng. 28(6): 147 (2023)

preprint: <https://arxiv.org/abs/2308.11258>

Dataset

- 6.9 million unique full texts of FOSS license variants
- Detected using filename patterns across the entire SWH archive
 - LICENSE, COPYRIGHT, NOTICE, etc.
- Metadata: file lengths measures, detected MIME type, detected SPDX license (via ScanCode), example origin repository, oldest public commit of origin, ground truth


Use cases

- Empirical studies on FOSS licensing, including phylogenetics
- Training of automated license classifiers
- NLP analyses of legal texts

The Software Heritage Filesystem (SwhFS)

The **Software Heritage Filesystem (SwhFS)** is a user-space POSIX filesystem that enables browsing parts of the Software Heritage archive as if it were locally available.

- Code: forge.softwareheritage.org/source/swh-fuse
- Documentation: docs.softwareheritage.org/devel/swh-fuse

 **Thibault Allançon, Antoine Pietri, Stefano Zacchiroli**
The Software Heritage Filesystem (SwhFS): Integrating Source Code Archival with Development
ICSE 2021 (Tool track): The 43rd Intl. Conference on Software Engineering
<https://arxiv.org/abs/2102.06390>

The Software Heritage Filesystem (SwhFS) — example

```
$ mkdir swarfs
$ swarf fs mount swarfs/ # mount the archive
$ cd swarfs/

$ cat archive/swh:1:cnt:c839dea9e8e6f0528b468214348fee8669b305b2
#include <stdio.h>

int main(void) {
    printf("Hello, World!\n");
}

$ cd archive/swh:1:dir:1fee702c7e6d14395bbf5ac3598e73bcbf97b030
$ ls | wc -l
127
$ grep -i antenna THE_LUNAR_LANDING.s | cut -f 5
# IS THE LR ANTENNA IN POSITION 1 YET
# BRANCH IF ANTENNA ALREADY IN POSITION 1
```

The Software Heritage Filesystem (SwhFS) — example (cont.)

```
$ cd archive/swh:1:rev:9d76c0b163675505d1a901e5fe5249a2c55609bc

$ ls -F
history/  meta.json@  parent@  parents/  root@

$ jq '.author.name, .date, .message' meta.json
"Michal Golebiowski-Owczarek"
"2020-03-02T23:02:42+01:00"
"Data:Event:Manipulation: Prevent collisions with Object.prototype ..."

$ find root/src/ -type f -name '*.js' | xargs cat | wc -l
10136
```




Paolo Boldi, Antoine Pietri, Sebastiano Vigna, Stefano Zacchiroli

Ultra-Large-Scale Repository Analysis via Graph Compression

SANER 2020, 27th Intl. Conf. on Software Analysis, Evolution and Reengineering. IEEE

Research question

Is it possible to efficiently perform software development history analyses at the scale of Software Heritage archive on a single, relatively cheap machine?

Idea

Apply state-of-the-art graph compression techniques from the field of Web graph / social network analysis.

Results

The entire archive graph (35 B nodes, 500 B edges) can be loaded in 300 GiB and then traversed at the cost of tens of ns per edge (= a few hours for a full single-thread visit).

Java and gRPC APIs available: docs.softwareheritage.org/devel/swh-graph/grpc-api.html

Background — (Web) graph compression

Definition (The graph of the Web)

Directed graph that has Web pages as nodes and hyperlinks between them as edges.

Properties (1)

- **Locality:** pages link to pages whose URLs are lexicographically similar. URLs share long common prefixes.

→ use **D-gap compression**

Adjacency lists

Node	Outdegree	Successors
...
15	11	13,15,16,17,18,19,23,24,203,315,1034
16	10	15,16,17,22,23,24,315,316,317,3041
17	0	
18	5	13,15,16,17,50
...

D-gapped adjacency lists

Node	Outdegree	Successors
...
15	11	3,1,0,0,0,0,3,0,178,111,718
16	10	1,0,0,4,0,0,290,0,0,2723
17	0	
18	5	9,1,0,0,32
...

Background — (Web) graph compression (cont.)

Definition (The graph of the Web)

Directed graph that has Web pages as nodes and hyperlinks between them as edges.

Properties (2)

- **Similarity:** pages that are close together in lexicographic order tend to have many common successors.

→ use **reference compression**

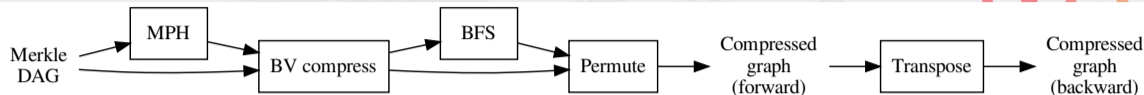
Adjacency lists

Node	Outd.	Successors
...
15	11	13,15,16,17,18,19,23,24,203,315,1034
16	10	15,16,17,22,23,24,315,316,317,3041
17	0	
18	5	13,15,16,17,50
...

Copy lists

Node	Ref.	Copy list	Extra nodes
...
15	0		13,15,16,17,18,19,23,24,203,315,1034
16	1	01110011010	22,316,317,3041
17			
18	3	11110000000	50
...	

Graph compression pipeline

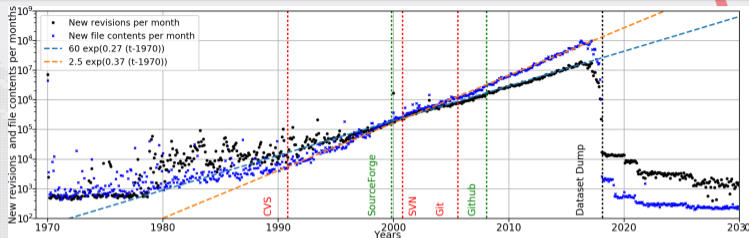


- **MPH**: minimal perfect hash, mapping Merkle IDs to 0..N-1 integers
- **BV compress**: Boldi-Vigna compression (based on MPH order)
- **BFS**: breadth-first visit to renumber
- **Permute**: update BV compression according to BFS order

(Re)establishing locality

- Key for good compression is a node ordering that ensures locality and similarity
- Which is very much *not* the case with Merkle IDs, ... but is the case *again* after BFS reordering

Software provenance and evolution



Key findings

- The amount of original commits in public code doubles every ~30 months and has been doing so for 20+ years; original source code files double every ~22 months
- It is possible to trace the provenance of source code artifacts at this scale in a compact relational model via the notion of isochrone graphs.



Rousseau, Di Cosmo, Zacchioli

Software Provenance Tracking at the Scale of Public Source Code

Empir. Softw. Eng. 25(4): 2930-2959 (2020)

Diversity, equity, and inclusion

Idea

Archived commit metadata contains public information that can be mined to study long-term trends of diversity, equity, and inclusion (DEI) traits of the global population of public code contributors.

Key findings on the gender gap

- Male authors contributed 92% of public code commits up to 2019.
- The ratio of female authors (and their contributions) has grown stably for 15 years reaching for the first time 10% of yearly contributions in 2019.
- The COVID-19 pandemic has reversed the trend (and it's not a coincidence!)

References

- Zacchiroli. *Gender differences in public code contributions: a 50-year perspective*. IEEE Software, 2021
- Rossi and Zacchiroli. *Worldwide gender differences in public code contributions (and how they have been affected by the COVID-19 pandemic)*. ICSE SEIS, 2022
- Casanueva, Rossi, Zimmermann, Zacchiroli. *The Impact of the COVID-19 Pandemic on Women's Contribution to Public Code*. Empir. Softw. Eng. Under review.

Key findings on the geographic gap

- Early decades of public code dominated by contributions from North America, followed by a period of alternating dominance between North America and Europe.
- Since then geographic diversity has increased constantly, with raising importance of contributions from Central and South America.
- The trend of increased female contributions is almost worldwide, with the notable exception of specific regions of Asia where it is either slower or flat.

References

- Rossi and Zacchiroli. *Geographic diversity in public code contributions*. MSR 2022

Ongoing work

Google AIR (Award for Inclusion Research) 2022, *What Causes the Lack of Diversity in Open Source?*

... and more (open research leads)

Cybersecurity

SWHSec (CampusCyber project) — main question: how can we leverage Software Heritage as a knowledge base to increase the security of open source software?

... and more (open research leads)

Cybersecurity

SWHSec (CampusCyber project) — main question: how can we leverage Software Heritage as a knowledge base to increase the security of open source software?

AI

- **Project classification** at the scale (*size and heterogeneity*) of Software Heritage (ongoing work with LabRI+UniBo)

... and more (open research leads)

Cybersecurity

SWHSec (CampusCyber project) — main question: how can we leverage Software Heritage as a knowledge base to increase the security of open source software?

AI

- **Project classification** at the scale (*size and heterogeneity*) of Software Heritage (ongoing work with LabRI+UniBo)
- **Code Commons**: producing research datasets for *ethical* training of LLMs on Software Heritage code

October 19, 2023

Software Heritage Statement on Large Language Models for Code



- 1 giving back to humanity
- 2 precisely identify training inputs
- 3 opt-out

- 1 Preface
- 2 Software as knowledge
- 3 Software Heritage
- 4 Research highlights in empirical software engineering
- 5 Conclusion

Conclusion

- Software Heritage archives public code and its history as a huge Merkle DAG
- Analyzing it at scale (35/500 B nodes/edges) is a significant big data undertaking
- Gold mine of research leads in and around empirical software engineering

Research on Software Heritage


www.softwareheritage.org/publications

Questions?

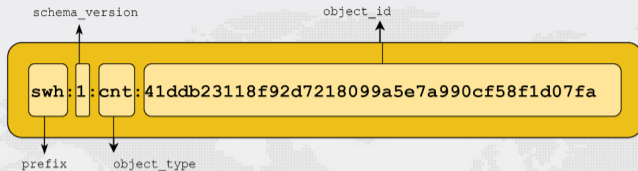
- Ask me!

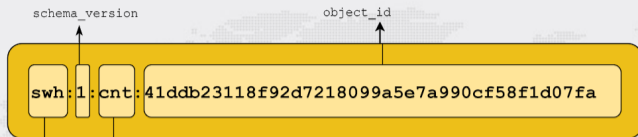
Contact

[Stefano Zacchioli](mailto:stefano.zacchioli@telecom-paris.fr) / stefano.zacchioli@telecom-paris.fr / [@zacchiro@mastodon.xyz](https://mstdn.org/@zacchiro)



Appendix





prefix

object_type

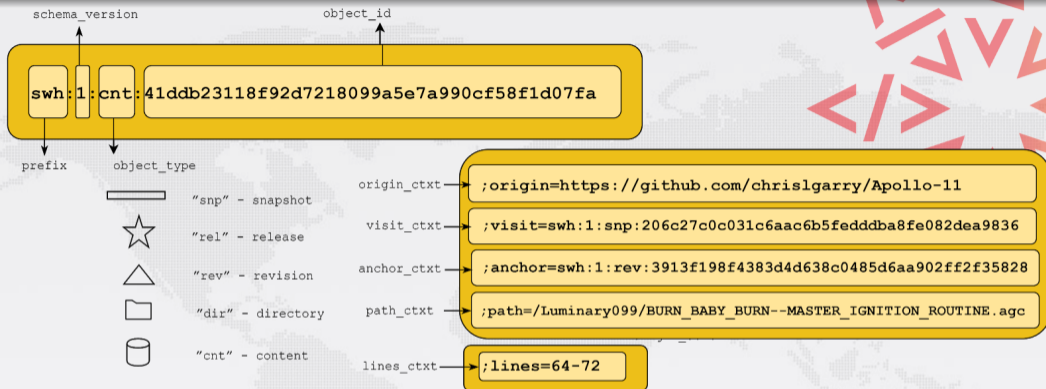
— "snp" - snapshot

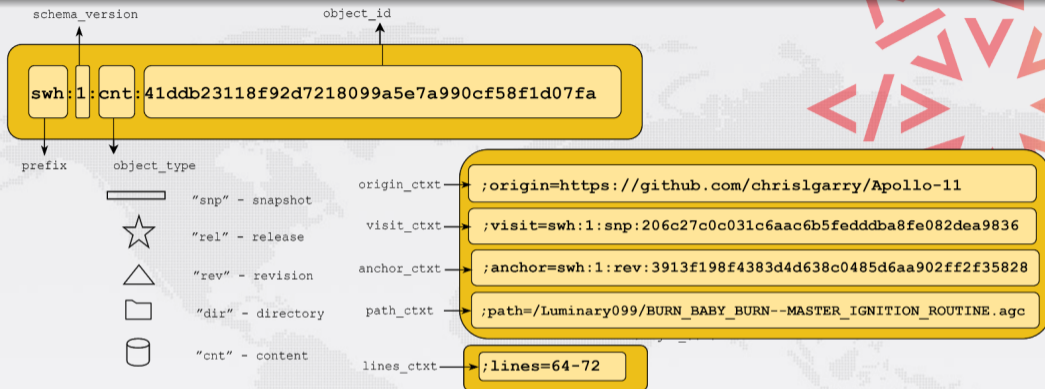
☆ "rel" - release

△ "rev" - revision

📁 "dir" - directory

🗄 "cnt" - content





An emerging standard

- in Linux Foundation's [SPDX 2.2](#)
- IANA-registered "swh:" URI prefix
- WikiData property [P6138](#)



An emerging standard

- in Linux Foundation's [SPDX 2.2](#)
- IANA-registered "swh:" URI prefix
- WikiData property [P6138](#)

Examples

- [Apollo 11 AGC excerpt](#)
- [Quake III rsqrt](#)

Sharing the vision



United Nations
Educational, Scientific and
Cultural Organization



www.softwareheritage.org/support/testimonials

Donors, members, sponsors



Diamond sponsor



Platinum sponsors



Gold sponsors



Silver sponsors



Bronze sponsors



www.softwareheritage.org/support/sponsors