

Software Heritage

A look at our team organization(s)

Nicolas Dandrimont
Tech Lead, Operations team
nicolas.dandrimont@inria.fr

April 22, 2026



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?
- 4 Tech team organization
- 5 Tech stack
- 6 Team's "Continuous Improvement"
- 7 Conclusions

Invisible fabric of digital society



Knowledge is in the source

```
/**
 * @brief The basic unit of the simulation and is associated to a geographical location.
 *
 * Interventions (e.g., school closures) are tracked at this level. It contains a list of its
 * members (people), places (schools, universities, workplaces etc.), road networks, links to
 * airports etc.
 */
struct Microcell
{
    /* Note use of short int here limits max run time to USHRT_MAX*ModelTimeStep - e.g. 65536*0.25=16384 days=
       Global search and replace of 'unsigned short int' with 'int' would remove this limit, but use more mem
    */

    int n; // Number of people in microcell
    int adunit; // admin unit microcell belongs to
    int* members; // array of members/hosts of microcell

    int* places[MAX_NUM_PLACE_TYPES]; // list of places (of various place types) within microcell
    unsigned short int NumPlacesByType[MAX_NUM_PLACE_TYPES]; // number of places (of various place types) with
    unsigned short int keyworkerproph, move_trig, place_trig, socdist_trig, keyworkerproph_trig;
    unsigned short int move_start_time, move_end_time;
    unsigned short int place_end_time, socdist_end_time, keyworkerproph_end_time;
    TreatStat moverest, treat, vacc, socdist, placeclose;
    unsigned short int treat_trig, vacc_trig;
    unsigned short int treat_start_time, treat_end_time;
    unsigned short int vacc_start_time;
    IndexList* AirportList;
};
```

Covid Sim ([excerpt](#))

Invisible fabric of digital society



Knowledge is in the source

```
/**
 * @brief The basic unit of the simulation and is associated to a geographical location.
 *
 * Interventions (e.g., school closures) are tracked at this level. It contains a list of its
 * members (people), places (schools, universities, workplaces etc.), road networks, links to
 * airports etc.
 */
struct Microcell
{
    /* Note use of short int here limits max run time to USHRT_MAX*ModelTimeStep - e.g. 65536*0.25=16384 days=
     * Global search and replace of 'unsigned short int' with 'int' would remove this limit, but use more mem
     */

    int n; // Number of people in microcell
    int adunit; // admin unit microcell belongs to
    int* members; // array of members/hosts of microcell

    int* places[MAX_NUM_PLACE_TYPES]; // list of places (of various place types) within microcell
    unsigned short int NumPlacesByType[MAX_NUM_PLACE_TYPES]; // number of places (of various place types) with
    unsigned short int keyworkerproph, move_trig, place_trig, socdist_trig, keyworkerproph_trig;
    unsigned short int move_start_time, move_end_time;
    unsigned short int place_end_time, socdist_end_time, keyworkerproph_end_time;
    TreatStat moverest, treat, vacc, socdist, placeclose;
    unsigned short int treat_trig, vacc_trig;
    unsigned short int treat_start_time, treat_end_time;
    unsigned short int vacc_start_time;
    IndexList* AirportList;
};
```

Covid Sim ([excerpt](#))

Len Shustek, Computer History Museum

2006

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

A Global Undertaking

From all continents

MSR '22, May 27-28, 2022, Pittsburgh, PA, USA

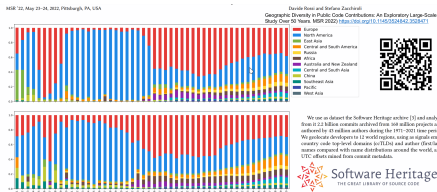
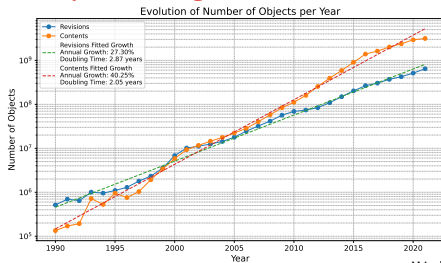
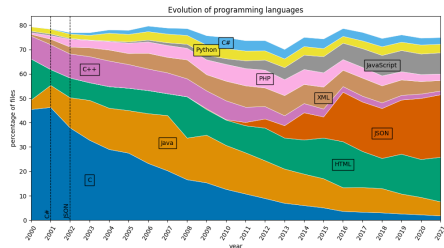


Figure 3: Ratio of commits (above) and active authors (below) by world zone over the 1971-2020 period.

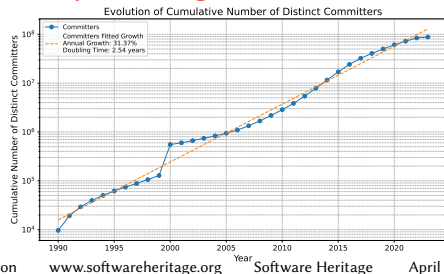
An exponential growth (code)



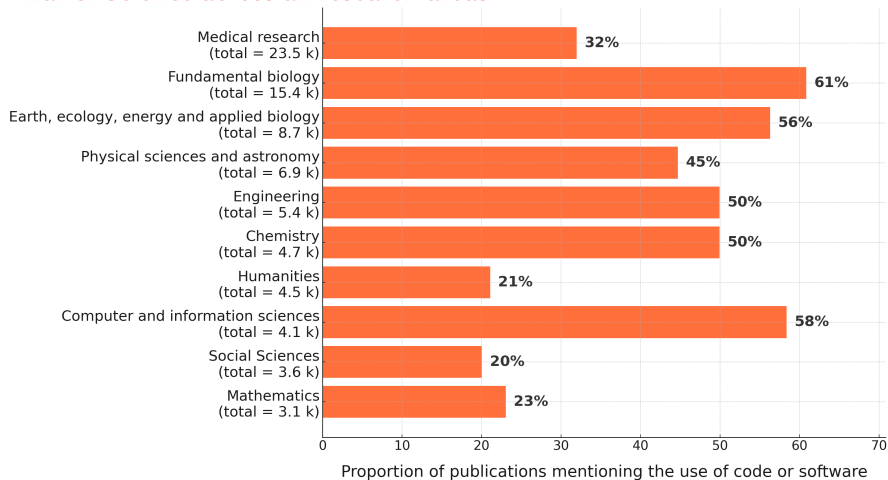
Many programming languages



An exponential growth (contributors)



Pillar of Science across all research areas



Source code is special

Software *evolves* over time

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

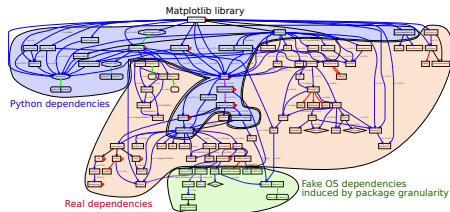
Source code is special

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*



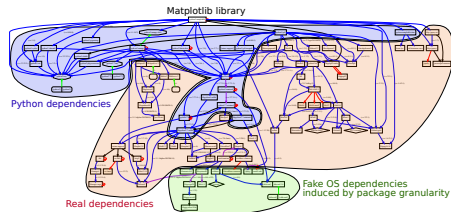
Source code is special

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*



Precious, endangered *executable* and *human readable* knowledge

key people **passing away**, platforms (GoogleCode, Gitorious, etc.) closing down ...
we need a **dedicated** infrastructure:

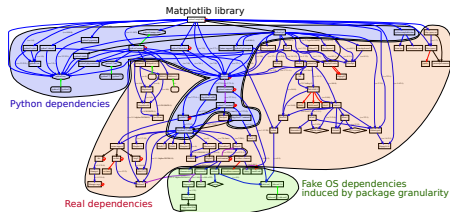
Source code is special

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*



Precious, endangered *executable* and *human readable* knowledge

key people **passing away**, platforms (GoogleCode, Gitorious, etc.) closing down ...

we need a **dedicated** infrastructure: now we have it!

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?
- 4 Tech team organization
- 5 Tech stack
- 6 Team's "Continuous Improvement"
- 7 Conclusions





Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

Inria

with



unesco



The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good *built in France since 2015*

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

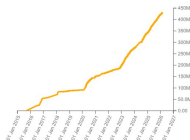
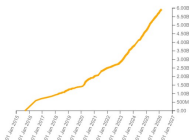
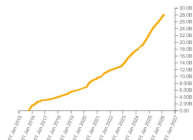


Public Administration



Projects

429,963,770





Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Inria with  **unesco**



The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good *built in France since 2015*

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

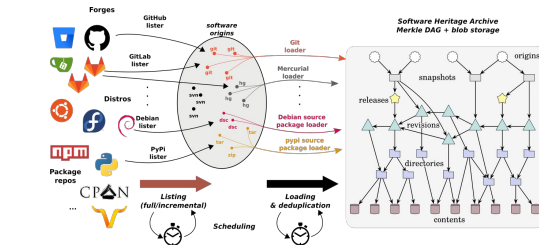


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

Inria with  unesco



The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good built in France since 2015

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

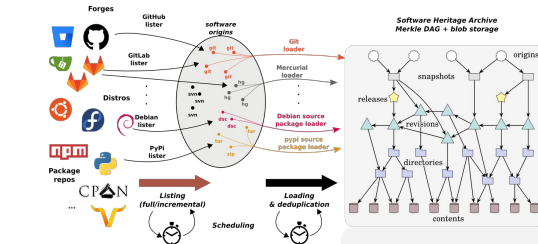


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph

- 50×10^9 nodes
- 1000×10^9 edges
~ 3 PB of storage



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

Inria with  unesco



The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good built in France since 2015

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

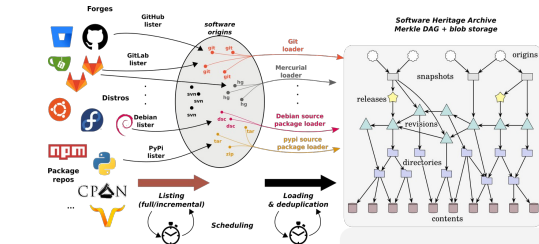


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph

- 50×10^9 nodes
- 1000×10^9 edges
~ 3 PB of storage

A revolutionary **infrastructure** ensures **availability** guarantees **integrity** enables **traceability**





Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Inria

with



unesco



The largest open source code archive: one infrastructure, open, shared, non profit

Unique digital common good built in France since 2015

Cultural Heritage



Source files

28,152,651,759

Industry



Commits

5,920,787,012

Research

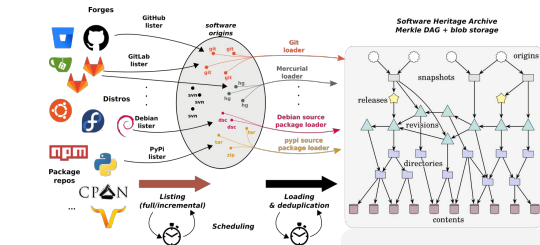


Public Administration



Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph

- 50 × 10⁹ nodes
- 1000 × 10⁹ edges
~ 3 PB of storage

A revolutionary **infrastructure** ensures **availability** guarantees **integrity** enables **traceability**



openinventionnetwork

AdaCore



GitHub



Unveiled in 2016



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

Unveiled in 2016



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

Unveiled in 2016



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

Unveiled in 2016



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

Sharing the vision



And many more ...

www.softwareheritage.org/support/testimonials

Sharing the vision



And many more ...

www.softwareheritage.org/support/testimonials

Members and sponsors



Diamond sponsors



Platinum sponsors



Gold sponsors



Silver sponsors



Bronze sponsors



*academia, industry, public sector, society
we are all concerned*

A *universal* software archive, as a shared infrastructure

One infrastructure
open and shared



Inria  UNESCO

A universal software archive, as a shared infrastructure

One infrastructure
open and shared



The largest archive ever built



A universal software archive, as a shared infrastructure

One infrastructure
open and shared

Cultural Heritage



Industry



Research



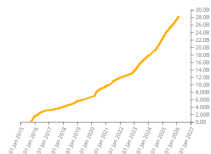
Public Administration



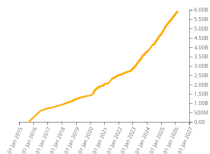
Software Heritage

The largest archive ever built

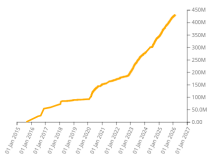
Source files
28,178,230,711



Commits
5,925,711,720



Projects
430,218,860



Directories
22,013,760,256

Authors
105,360,419

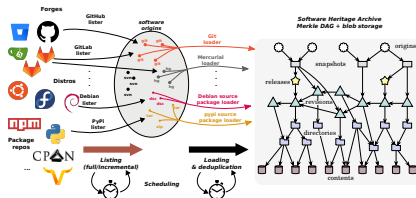
Releases
147,686,240



figures as of January 8 2026

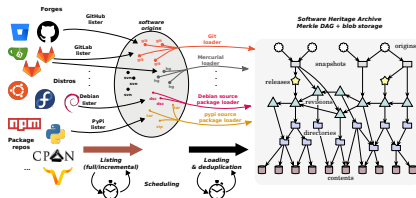
Addressing key needs in (Open) Science

Archive (28B+ files, 430M+ projects)



Addressing key needs in (Open) Science

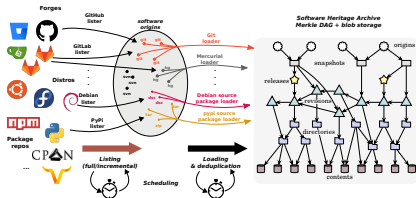
Archive (28B+ files, 430M+ projects)



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

Addressing key needs in (Open) Science

Archive (28B+ files, 430M+ projects)



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

Reference (50 billion SWHIDs)

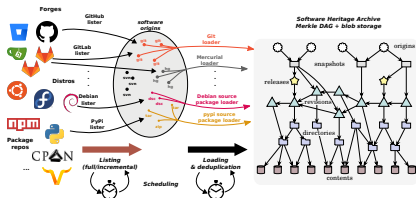
Intrinsic, cryptographically strong IDs



Now in SPDX 2.2, Wikidata
<https://swhid.org> - ISO/IEC 18670

Addressing key needs in (Open) Science

Archive (28B+ files, 430M+ projects)



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

Describe

- *Intrinsic metadata* from source code
- Contributed the [Codemeta generator](#)

Reference (50 billion SWHIDs)

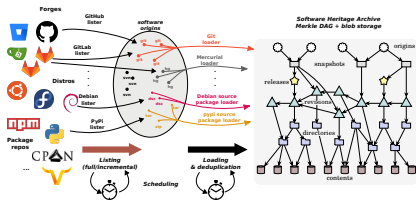
Intrinsic, cryptographically strong IDs



Now in SPDX 2.2, Wikidata
<https://swhid.org> - ISO/IEC 18670

Addressing key needs in (Open) Science

Archive (28B+ files, 430M+ projects)



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

Describe

- *Intrinsic metadata* from source code
- Contributed the [Codemeta generator](#)

Reference (50 billion SWHIDs)

Intrinsic, cryptographically strong IDs



Now in SPDX 2.2, Wikidata
<https://swhid.org> - ISO/IEC 18670

Cite/Credit

- Contributed [biblatex-software style](#)
- Software Citation from the archive!

An example is worth a thousand words

- Browse + Reference [ISO 18670] (Apollo 11 [excerpt], your work may be already there !)
- Trigger archival, use the updateswh browser extension, configure the webhooks
- Cite from the archive with biblatex-software (CTAN, ACMART)
- Describe with Codemeta (use codemeta generator)
- Curated deposit in SWH via HAL, see for example: LinBox, SLALOM, Givaro, NS2DDV, SumGra, Coq proof, ...
- Extracting all the software products for Inria, for CNRS, for CNES, for LIRMM or for Rémi Gribonval using HalTools
- Example with Parmap: devel on Github, archive in SWH, curated deposit in HAL
- Example research articles:
 - compare Fig. 1 and conclusions in the 2012 version and the updated version
 - SWHID in a replication experiment

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?**
- 4 Tech team organization
- 5 Tech stack
- 6 Team's "Continuous Improvement"
- 7 Conclusions

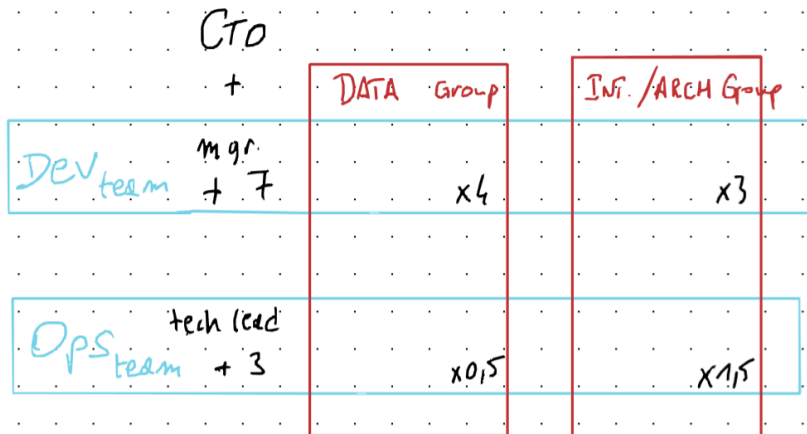
- Inria Rocquencourt
 - 5 racks - 70 machines (Prod and Staging infrastructures)
 - Proxmox cluster (VMs hosting + networked storage for workloads)
 - Kubernetes clusters (Bare metal + VMs)
 - Bare metal data silos: Cassandra clusters; Kafka clusters; Elasticsearch clusters
 - Snowflake servers (graph compression, etc.)
- CEA Saclay
 - 3 racks - 40 machines (Ceph cluster - 4PB - for mass storage)

- Inria Rocquencourt
 - 5 racks - 70 machines (Prod and Staging infrastructures)
 - Proxmox cluster (VMs hosting + networked storage for workloads)
 - Kubernetes clusters (Bare metal + VMs)
 - Bare metal data silos: Cassandra clusters; Kafka clusters; Elasticsearch clusters
 - Snowflake servers (graph compression, etc.)
- CEA Saclay
 - 3 racks - 40 machines (Ceph cluster - 4PB - for mass storage)
- Microsoft Azure
 - Object storage (copy of contents), hosted kubernetes clusters for ancillary infra
- Amazon S3
 - Object storage (copy of contents, dataset exports)



- 25 staff members (20 full time, 5 part time)
- 10 people in management / outreach / support
- 15 people in tech team (CTO + operations team + development team)

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?
- 4 Tech team organization**
- 5 Tech stack
- 6 Team's "Continuous Improvement"
- 7 Conclusions



Groups

- 2-week cycle for both groups (1h30 retrospective every 2 weeks)
- 15 mins "daily" checkpoints twice a week

Groups

- 2-week cycle for both groups (1h30 retrospective every 2 weeks)
- 15 mins "daily" checkpoints twice a week
- Communication toolkit:
 - synchronous: jitsi + matrix
 - asynchronous: GitLab issues / Epics; mailing lists
 - collaborative editing: HedgeDoc

Groups

- 2-week cycle for both groups (1h30 retrospective every 2 weeks)
- 15 mins "daily" checkpoints twice a week
- Communication toolkit:
 - synchronous: jitsi + matrix
 - asynchronous: GitLab issues / Epics; mailing lists
 - collaborative editing: HedgeDoc

Teams

- Weekly checkpoint for the dev & ops teams

Medium-term organization

Whole team

- Yearly (rolling) roadmap with high level priorities
- Monthly all hands meeting (management + outreach + tech team)

Medium-term organization

Whole team

- Yearly (rolling) roadmap with high level priorities
- Monthly all hands meeting (management + outreach + tech team)

Tech team

- Meeting every quarter (incl. retrospective and priorities for the next quarter)
- Regular individual meetings with CTO or team lead
- Regular (2-3 weekly) meeting for tech team leadership

Medium-term organization

Whole team

- Yearly (rolling) roadmap with high level priorities
- Monthly all hands meeting (management + outreach + tech team)

Tech team

- Meeting every quarter (incl. retrospective and priorities for the next quarter)
- Regular individual meetings with CTO or team lead
- Regular (2-3 weekly) meeting for tech team leadership

Management

- Weekly management meetings (Management team + Tech, Outreach, Comms team leaders)
- Follow up on external projects (e.g. CodeCommons projects)

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?
- 4 Tech team organization
- 5 Tech stack**
- 6 Team's "Continuous Improvement"
- 7 Conclusions

- Everything FLOSS (A)GPLv3+, developed [in the open](#)

The Software Heritage tech stack

- Everything FLOSS (A)GPLv3+, developed [in the open](#)
- [Internal stack](#) developed in house (Python + Rust) by the dev team

The Software Heritage tech stack

- Everything FLOSS (A)GPLv3+, developed [in the open](#)
- [Internal stack](#) developed in house (Python + Rust) by the dev team
- Deployment manifests mostly written by the ops team
 - [terraform](#) for low-level provisioning
 - [puppet](#) for system-level provisioning
 - [Helm](#) for kubernetes deployments

- Unit testing of each Python / Rust module on [Jenkins](#)
 - WIP: Gitlab CI integration

- Unit testing of each Python / Rust module on [Jenkins](#)
 - WIP: Gitlab CI integration
- [Integration tests](#) of the developed stack
 - Docker Compose / testinfra
 - Doubles as local development stack

- Unit testing of each Python / Rust module on [Jenkins](#)
 - WIP: Gitlab CI integration
- [Integration tests](#) of the developed stack
 - Docker Compose / testinfra
 - Doubles as local development stack
- Jenkins builds docker containers on releases
 - Pushes them to GitLab container registry
 - Creates a merge request to bump container versions in Helm chart

- [ArgoCD](#) to manage state of apps on kubernetes clusters

- [ArgoCD](#) to manage state of apps on kubernetes clusters
- 3 static Kubernetes clusters
 - "next-version": automatic deployment of latest container versions
 - staging: small-scale integration testing
 - production

- [ArgoCD](#) to manage state of apps on kubernetes clusters
- 3 static Kubernetes clusters
 - "next-version": automatic deployment of latest container versions
 - staging: small-scale integration testing
 - production
- WIP: autonomy for deployments in the two groups
 - Separate K8s namespaces
 - ...?

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?
- 4 Tech team organization
- 5 Tech stack
- 6 Team's "Continuous Improvement"**
- 7 Conclusions

Tech debt

- Organic growth over 10 years; little time to spend on iterative improvements of the whole stack
- Very complex stack, impossible to hold in one's head all at once
- Unexpected consequences of seemingly small changes

Tech debt

- Organic growth over 10 years; little time to spend on iterative improvements of the whole stack
- Very complex stack, impossible to hold in one's head all at once
- Unexpected consequences of seemingly small changes

Limited resources

- Very difficult to do production-level tests at the large scale of the Archive
- Difficult to iterate once a solution is in place, within limited resources
- Data/schema migrations take months to years

Tech debt

- Organic growth over 10 years; little time to spend on iterative improvements of the whole stack
- Very complex stack, impossible to hold in one's head all at once
- Unexpected consequences of seemingly small changes

Limited resources

- Very difficult to do production-level tests at the large scale of the Archive
- Difficult to iterate once a solution is in place, within limited resources
- Data/schema migrations take months to years

Communication, internal vs. external

- Free software-style open-by-default challenges, Especially 2nd language English
- Documentation often an afterthought

Over the past 5 years

- "Software Craftsmanship"
 - Agile methods / Pair/Mob programming / ...

Over the past 5 years

- "Software Craftsmanship"
 - Agile methods / Pair/Mob programming / ...
- ITIL Foundations
 - Common "language" for project/product/service driven management

Over the past 5 years

- "Software Craftsmanship"
 - Agile methods / Pair/Mob programming / ...
- ITIL Foundations
 - Common "language" for project/product/service driven management
- Kubernetes basics
 - Wanting to reduce the gap between dev and ops stack

Completed

- Recently recruited a tech writer to help keep our documentation useful (both for the team and for our users)

Completed

- Recently recruited a tech writer to help keep our documentation useful (both for the team and for our users)
- Introduction of SWHEP process (similar to RFC / PEP / ADR) for impactful changes in the tech stack

Completed

- Recently recruited a tech writer to help keep our documentation useful (both for the team and for our users)
- Introduction of SWHEP process (similar to RFC / PEP / ADR) for impactful changes in the tech stack

Ongoing

- Tech debt reduction
 - Automated dependency updates with [Mend Renovate](#)
 - GitLab CI in some projects
 - Migration to GitLab Enterprise for project management features

Completed

- Recently recruited a tech writer to help keep our documentation useful (both for the team and for our users)
- Introduction of SWHEP process (similar to RFC / PEP / ADR) for impactful changes in the tech stack

Ongoing

- Tech debt reduction
 - Automated dependency updates with [Mend Renovate](#)
 - GitLab CI in some projects
 - Migration to GitLab Enterprise for project management features
- Sharing ownership of our software deployment pipeline

Completed

- Recently recruited a tech writer to help keep our documentation useful (both for the team and for our users)
- Introduction of SWHEP process (similar to RFC / PEP / ADR) for impactful changes in the tech stack

Ongoing

- Tech debt reduction
 - Automated dependency updates with [Mend Renovate](#)
 - GitLab CI in some projects
 - Migration to GitLab Enterprise for project management features
- Sharing ownership of our software deployment pipeline
- Many iterative changes over time to make our organization more fluid

- 1 Software as a pillar of modern society
- 2 Meet Software Heritage: open, non profit, multistakeholder
- 3 What's under the hood?
- 4 Tech team organization
- 5 Tech stack
- 6 Team's "Continuous Improvement"
- 7 **Conclusions**

- Still a long road ahead!

- Still a long road ahead!
- Maybe early in building some structures (e.g. the Team/Group matrix)

- Still a long road ahead!
- Maybe early in building some structures (e.g. the Team/Group matrix)
- But we're looking to continue to grow the tech team

- Still a long road ahead!
- Maybe early in building some structures (e.g. the Team/Group matrix)
- But we're looking to continue to grow the tech team
- Some incremental improvements have felt slow but the results are starting to show

Questions?