# PIDs for software - shedding some light on a dark puzzle

**Daniel S. Katz**

Chief Scientist, NCSA, University of Illinois at Urbana-Champaign
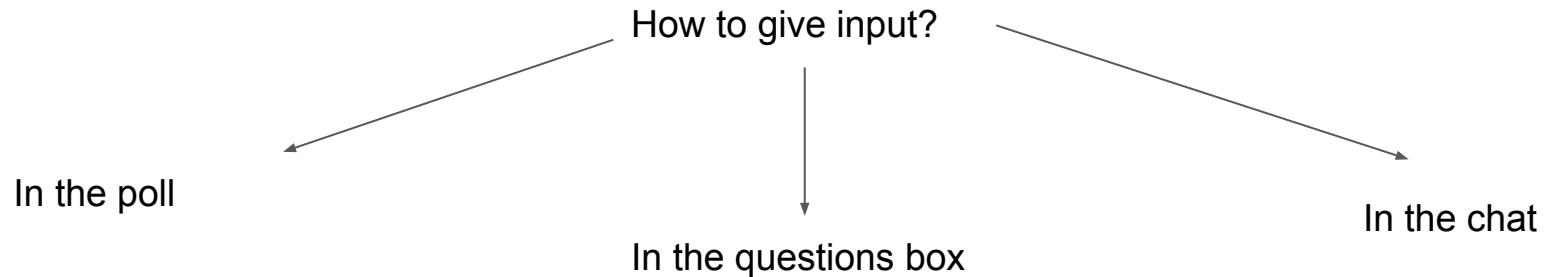
**Morane Gruenpeter**

Software Engineer

Metadata Specialist

Inria, Software Heritage

# This session's goals

1. **Share findings** from the Software Source Code Identification Working Group
   **DOI: 10.15497/RDA00053**
2. **Learn** about attendees' backgrounds
3. **Show the complex puzzle** of identifying software
4. Get attendee opinions about **where we are**
5. Get attendee opinions about **where we should go**

How to give input?

In the poll

In the questions box

In the chat

PIDapalooza

©FORCE11
The Future of Research Communications and e-Scholarship

RDA
RESEARCH DATA ALLIANCE

# Question 1: Are you **dealing** with software in your daily work?

- Yes, using,
- Yes, implementing
- Yes, developing
- Yes, designing
- Yes, testing
- Yes, all of the above
- No or Not yet

# Software is all around us...

**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

Tim Berners-Lee invented the **World Wide Web, 1989**, while working at CERN

"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 at 27:27 (25 January 2013)

**PIDapalooza**

© FORCE11
The Future of Research Communications and e-Scholarship

RDA
RESEARCH DATA ALLIANCE

# Question 2:  Have you **heard about** or **read** the RDA/FORCE11 Software Source Code Identification output?

- yes! I even read and commented during the community review,
- yes only heard about it,
- not yet :)

# Software Source Code Identification Working Group

The SCID WG Goal: **capture and analyze** the software identification state-of-the-art in the scholarly ecosystem

**Co-chairs**
- Roberto Di Cosmo
- Martin Fenner
- Daniel S. Katz

**Secretariat Liaison**
- Stefanie Kethers

RDA page
https://www.rd-alliance.org/groups/software-source-code-identification-wg

Repository
https://github.com/force11/force11-rda-scidwg

**Chronology...**
03/2018 Spawned at RDA **P11** in Berlin from the
- RDA Software Source Code IG &
- FORCE11 Software Citation Implementation WG

10/2018 - TAB endorsement

4/2019 - RDA **P13**, Philadelphia
- **WG kick-off**

10/2019 - **FORCE2019**, Edinburgh    Full day hackathon on research software

03/2020 - RDA VP15 session online

07/2020 - Output in community review **DOI:10.15497/RDA00053**

# Question 3: What kind of **stakeholder** are you in the scholarly ecosystem?

- a researcher, a research software engineer,
- a publisher, a citation manager,
- a funder, an evaluator,
- a registry admin, a curator
- or some other role

# Definition: Actors

# Question 4: Do you **need to identify** software in your work?

- Yes, in my articles, reports and other writings
- Yes, in my data
- Yes, in my software
- Yes, in different settings
- Not at the moment

# Question 5: What do you want to see when you click on a software PID?

- a project, a module,
- a release, a source code artifact, an executable,
- a docs page, a tutorial,
- a docker image
- or depending on the use-case at hand

# Identification target - what do we want to identify?

**Software concept / project / collection**
Description in registry, a homepage or any other form of metadata record
- Project versions (for example Python2 and Python3)
- Modules
- Sub-modules

**Software artifact**
- Executable (download link)
- Software source code
  - Dynamic artifact - current development code (on collaborative development platform)
  - Archived copy
    - Snapshot (all branches, all dev history)
    - Release / Package
    - Commit- a specific point in development history
    - Directory
    - File
    - Algorithm

**Software context**
- Complementary artifacts
  - the software environment, tutorial (Jupyter notebook), Data (input/output data), etc.
- Articles
- Documentation

Pyramid levels (top to bottom):
Project — GL1
Project versions — GL2
Modules — GL3
Sub-Modules — GL4
Snapshots — GL5
Releases — GL6
Commits — GL7
Directories — GL8
Files — GL9
Code fragments — GL10

GL= Granularity Level

# Question 6: What's the most important **use case** in your workflow?

- Get credit / Give credit
- To reproduce an experiment
- Curate software metadata
- Access the software source code to use it
- None of the above

# The use cases collection (a small excerpt)

| Actor | Use case description | Action | Identification target |
|---|---|---|---|
| Archive | Identify all the software artifacts I hold | Archiving, referencing | Release and smaller artifacts |
| Citation manager | Curate the software citation entries | Credit | Project, release |
| Curator / librarian / digital archivist | Catalog and browse the development history of legacy software source code for preservation purposes | Archiving | Project, release and smaller artifacts depending on the reference |
| Publisher | Create/retrieve identifiers quickly for use in the paper for all software including commercial packages. | Referencing, describing | Any item (all granularity levels) |
| Registry | Identify and curate the software entries I hold | Archiving, referencing, describing, credit | Project |
| Researcher as a software user (RSU) | Access and use SSC no longer available on a collaborative platform | Archiving | Snapshot, release, revision, directory |

PIDapalooza

FORCE11
The Future of Research Communications and e-Scholarship

RDA
RESEARCH DATA ALLIANCE

# Question 7: Are you using a **software identifier**?

- Yes, all the time
- Yes, sometimes
- Maybe, I'm not sure
- No, it is too difficult to do

# Landscape of Identifiers schemes (small excerpt)

HAL-ID

Digital Object Identifier

URI

WIKIDATA
Wiki Item identifier (Qxxx)

**ASCL.net**
Astrophysics Source Code Library

## Software Identification

**ARK**
Archival Resource Key

RRID

**Handle**
Handle System identifiers

git

swMATH
an information service for mathematical software

PIDapalooza

SWHID
Software Heritage identifiers

FORCE11
The Future of Research Communications and e-Scholarship

RDA
RESEARCH DATA ALLIANCE

# Intrinsic identifier: the Software Heritage ID ([SWHID](#))

- **Intrinsic**: compute a unique **digital fingerprint**
- **cryptographically strong** identifiers
- **decentralised**:
    - do not need a registry,
    - agreement on a standard

schema_version

object_id

```
swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa
```

prefix    object_type

```
"snp" - snapshot
"rel" - release
"rev" - revision
"dir" - directory
"cnt" - content
```

origin_ctxt — `;origin=https://github.com/chrislgarry/Apollo-11`

visit_ctxt — `;visit=swh:1:snp:206c27c0c031c6aac6b5fedddba8fe082dea9836`

anchor_ctxt — `;anchor=swh:1:rev:3913f198f4383d4d638c0485d6aa902ff2f35828`

path_ctxt — `;path=/Luminary099/BURN_BABY_BURN--MASTER_IGNITION_ROUTINE.agc`

lines_ctxt — `;lines=64-72`

[Documentation](#)

# Wikidata entities for identifying projects

Q1165184=SageMath.

A few examples of external identifiers properties of used on software entities:

- Arch package sagemath
- Debian stable package
- Fedora package
- Free Software Directory entry
- Freebase
- Gentoo package
- Open Hub ID sage
- Quora topic
- Ubuntu Package
- swMATH work ID 825
- SWHID snapshot (15.11.2020)
- and many more

# Identifiers from Registries

- To meet discipline-specific needs, discipline-specific registries have appeared
- **RRID** (life sciences), **ASCL** (astronomy), **swMath** (math)
- These let the community register different types of objects (including software) with an organization that curates the entries
- The entries then can be used in papers, including in citations
- But each is independent, and because software can cross disciplines, the registries can have overlapping contents

# DOI on JOSS articles

- Journal of Open Source Software (JOSS, https://joss.theoj.org)
  - A developer friendly, open access journal for research software packages
- Research software packages are identified with:
  - The article DOI: JOSS publishes a short paper with a Crossref DOI (https://doi.org/10.21105/joss.01686)
  - The software DOI: The author deposits the software in a repository that provides a DOI (e.g., Zenodo) (https://doi.org/10.5281/zenodo.3547813)
- As a result, there are 2 DOIs.
- Links from the paper and metadata DOI to:
  - the software deposit and its DOI,
  - the live version of the software (e.g., on GitHub)

# Summary

| Granularity level (GL) | ID target | Extrinsic identifiers | | | | | | | | | Intrinsic identifiers | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | *ASCL* | *ARK* | *DOI* | *HAL* | *URL* | *RRID* | *SwMath* | *Wikidata* | | *Hash* | *SWH* |
| | | | | | | | | | *entity* | *property* | | |
| GL1 | project | X | X | X | X | X | X | X | X | | | |
| GL2 | project version | | X | | | | | | X | | | |
| GL3 | module | | X | | | | | | X | | | |
| GL4 | repository | | X | | | X | | | | X | | |
| GL5 | repository snapshot | | X | | | | | | | X | | X |
| GL6 | release | | X | X** | | | | | | X | X | X |
| GL7 | commit | | X | | | | | | | X | X | X |
| GL8 | directory | | X | X** | X* | | | | | X | X | X |
| GL9 | file | | X | X** | | | | | | | X | X |
| GL10 | Code fragment | | X | | | | | | | | | X |

# Final question (8): Given this mess (dark puzzle), how can we clean it up?

## (considering communities, technologies, etc.)