



From Code Repositories to Knowledge Graphs of Research Software Metadata

**Daniel Garijo, Ontology Engineering Group,
Universidad Politécnica de Madrid, Spain**

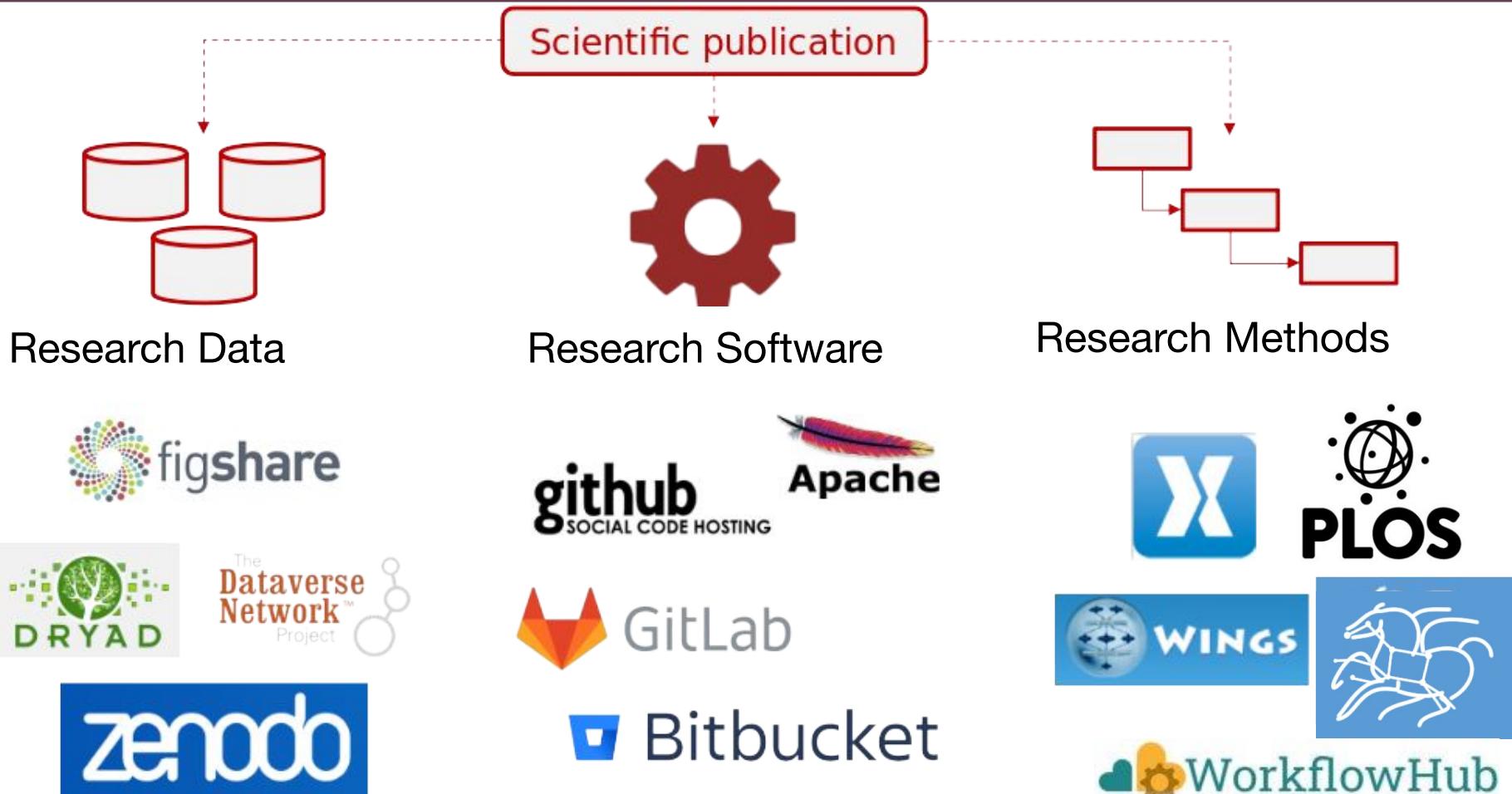
✉ daniel.garijo@upm.es

🐦 [@dgarijov](https://twitter.com/dgarijov)

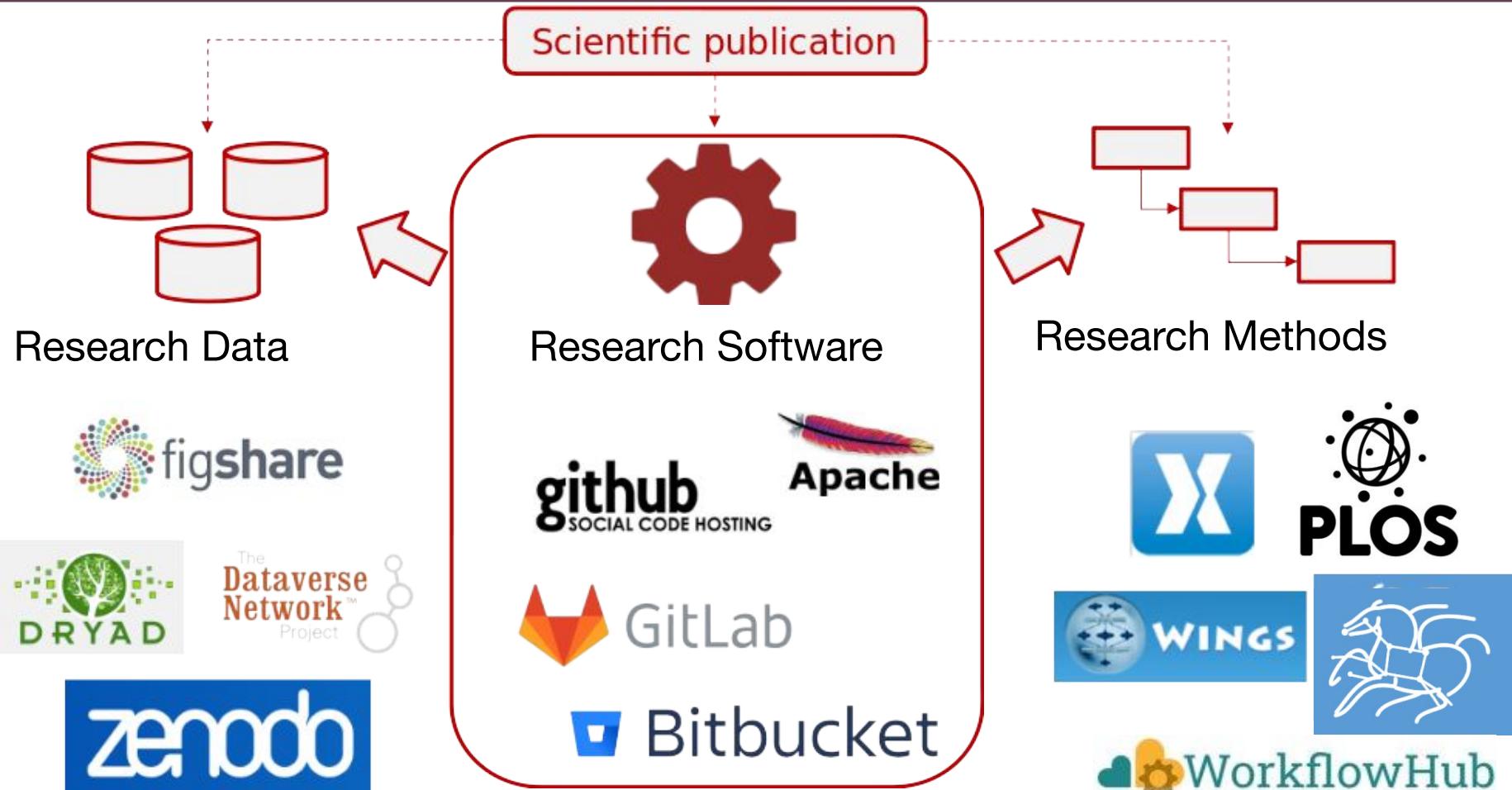
Research Software is one of the pillars of Open Science

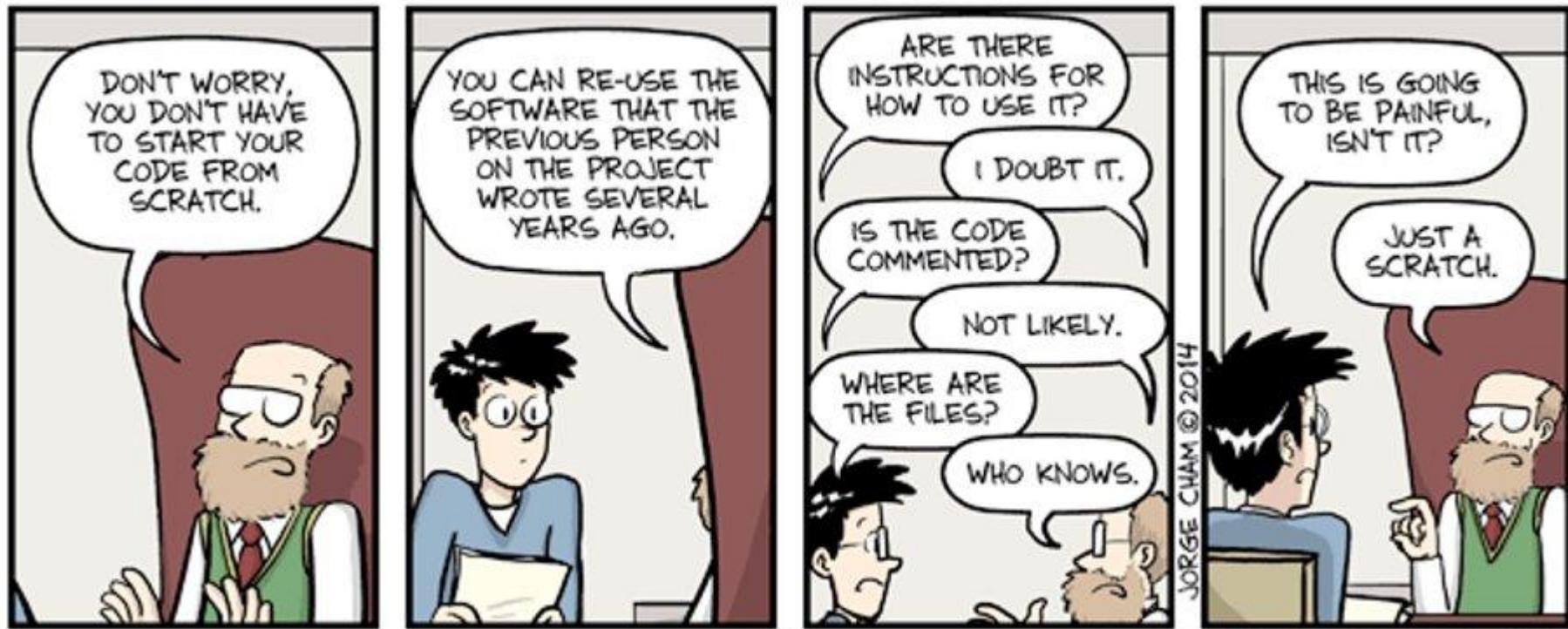


Research Software is one of the pillars of Open Science

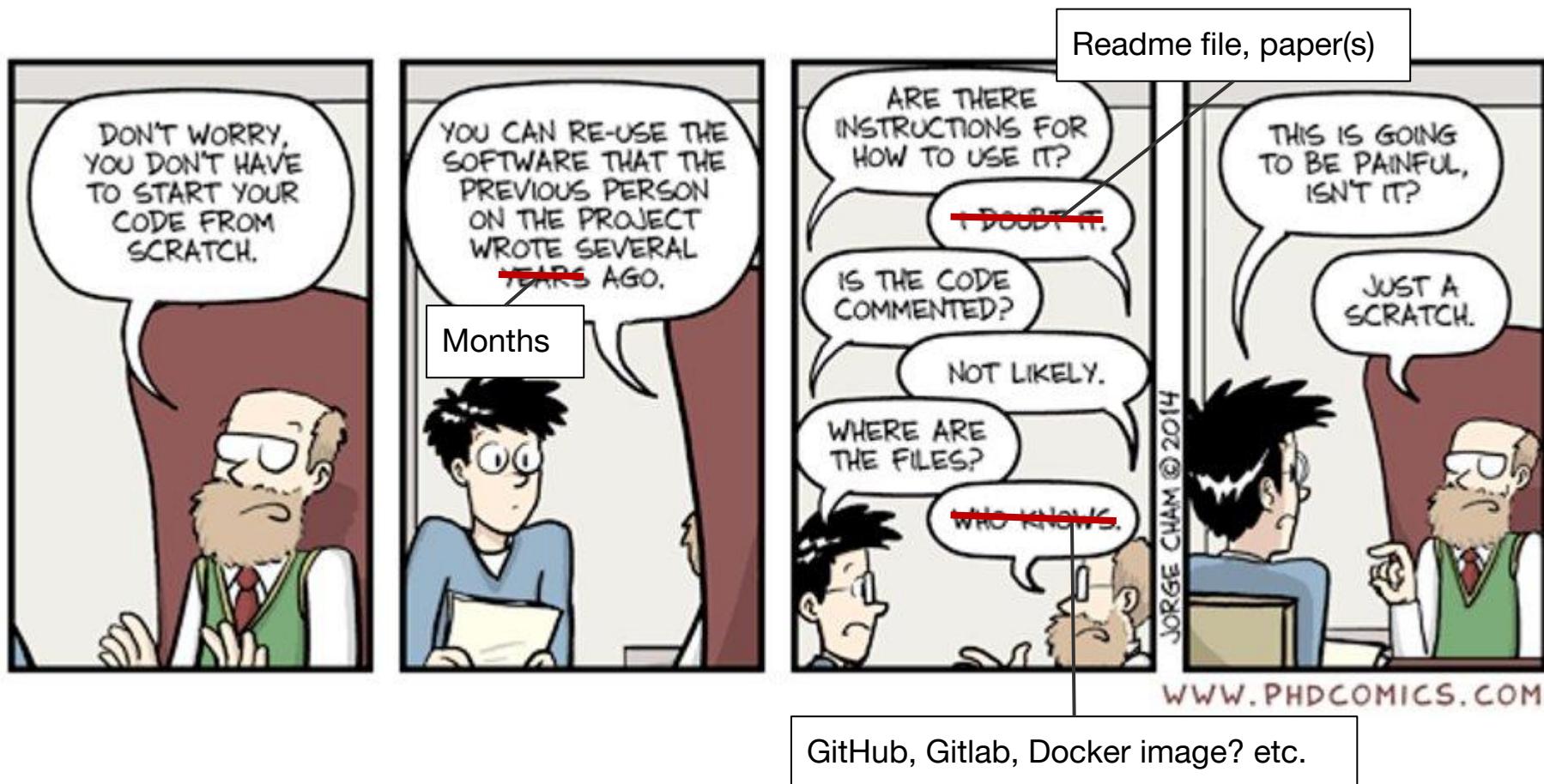


Research Software is one of the pillars of Open Science





Challenge 1: Lack of structured (machine-readable) documentation



In [1] we tried to reproduce an effort from **one** year before.

- All data were available online
- All tools were available online (except one, but authors had a replacement)
- > 250 hrs to full reproducibility
- > 100 hrs to get familiar with the tools and their I/O

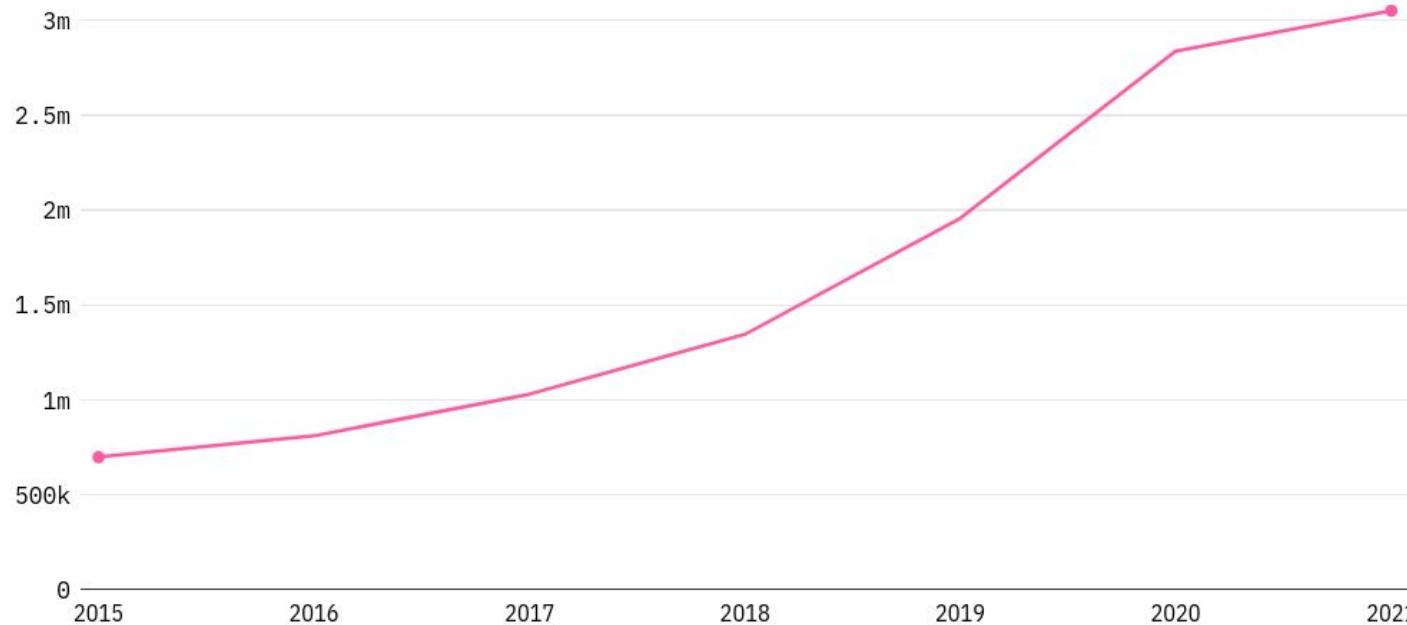
The screenshot shows a research article page. At the top, it says "OPEN ACCESS" and "PEER-REVIEWED". Below that is the title "RESEARCH ARTICLE". The title of the article is "The *Mycobacterium tuberculosis* Drugome and Its Polypharmacological Implications". Below the title, the authors listed are Sarah L. Kinnings, Li Xie, Kingston H. Fung, Richard M. Jackson, Lei Xie, and Philip E. Bourne. It was published on November 4, 2010, with the DOI <https://doi.org/10.1371/journal.pcbi.1000976>. To the right, there are metrics: 191 Save, 93 Citation, 20,259 View, and 1 Share. Below these are buttons for "Download PDF", "Print", and "Share". Further down, there's a "Check for updates" button and a "Subject Areas" section with several categories like "Drug discovery", "Protein structure", etc., each with a radio button next to it. The main content area has two columns: "Abstract" and "Abstract". The left column contains links to "Author Summary", "Introduction", "Results", "Discussion", "Methods", "Supporting Information", "Acknowledgments", "Author Contributions", and "References". The right column contains the abstract text, which discusses a computational approach to construct a drug-target network for *Mycobacterium tuberculosis*, identifying new targets and improving drug repositioning.

[1] Garijo, D., Kinnings, S., Xie, L., Xie, L., Zhang, Y., Bourne, P. E., & Gil, Y. (2013). Quantifying reproducibility in computational biology: the case of the tuberculosis drugome. *PLoS one*, 8(11), e80278.

Challenge 3: Comparing against existing tools

Millions of **open-source repositories** are updated/created every year

Number of first-time contributors for open source projects, by year



Source: GitHub Octoverse Report 2021

TECH MONITOR

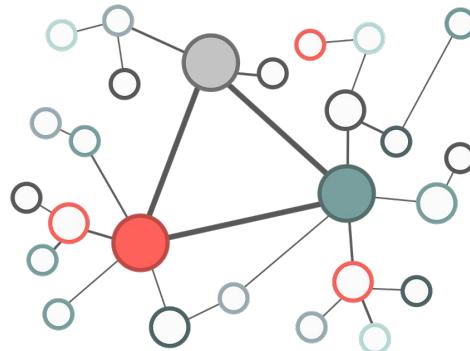
[Hucka et al]: Scientists still rely on three main methods for searching new software:

- Survey
- Recommendation from a colleague
- Search engine

1. Structured representation
2. Reuse
3. Compare
4. Search



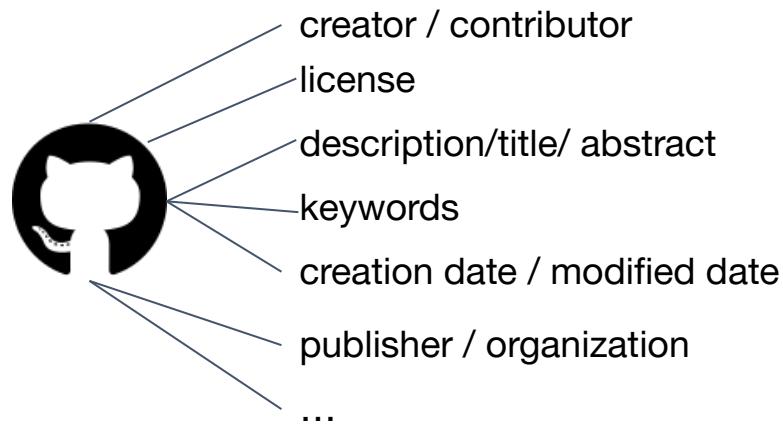
1. Structured representation
2. Reuse
3. Compare
4. Search



Knowledge Graphs

Creating KGs of Research Software metadata:

Representing RS at different levels of abstraction



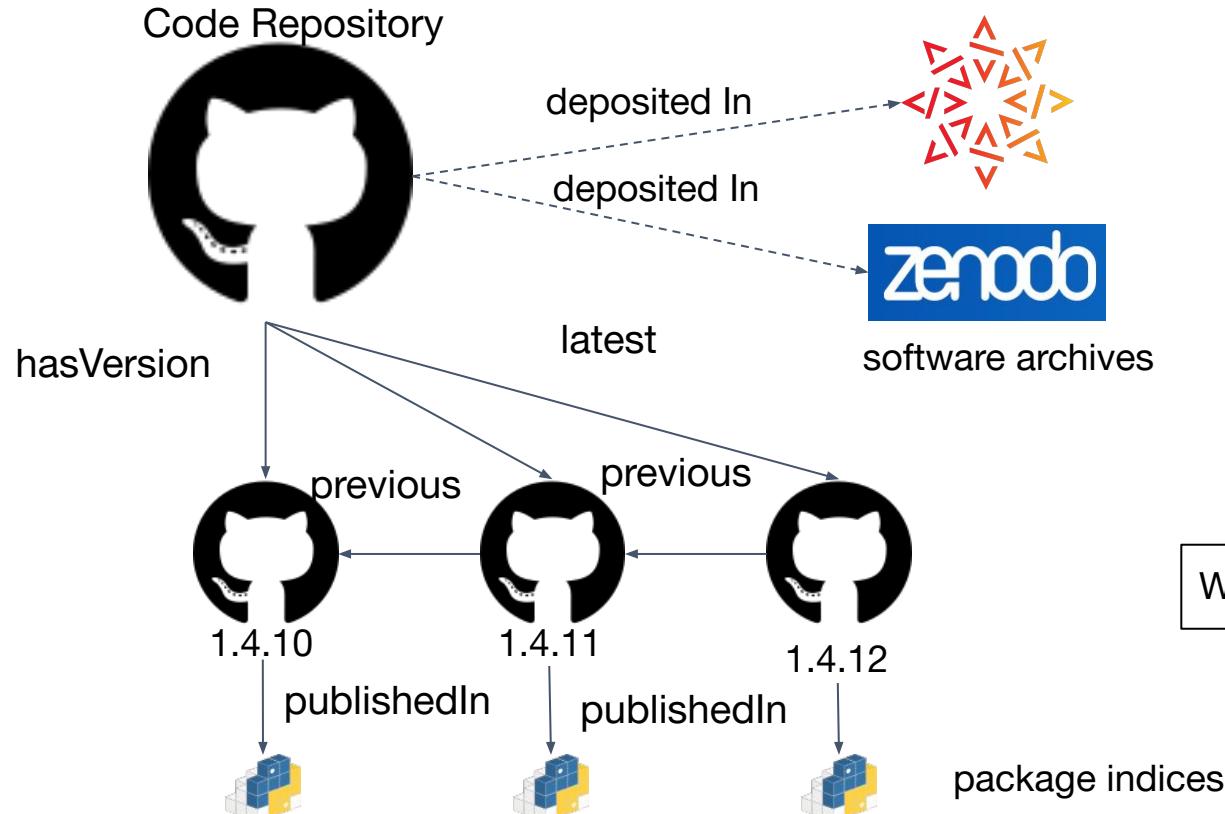
ORCIDs **not available**

Why is it needed?
- Search
- Compare

CodeMeta Schema.org

<https://w3id.org/okn/o/sd>

<https://w3id.org/software-types/>



Why is it needed?

- Search
- Compare
- Reuse

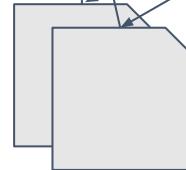
Which **identifiers** to use?

Code repository



**Citation File Format,
.bib**

resultOf



Research paper

refersTo

How to reconcile **code authors** with paper authors?

Why is it needed?

- Credit
- Trust

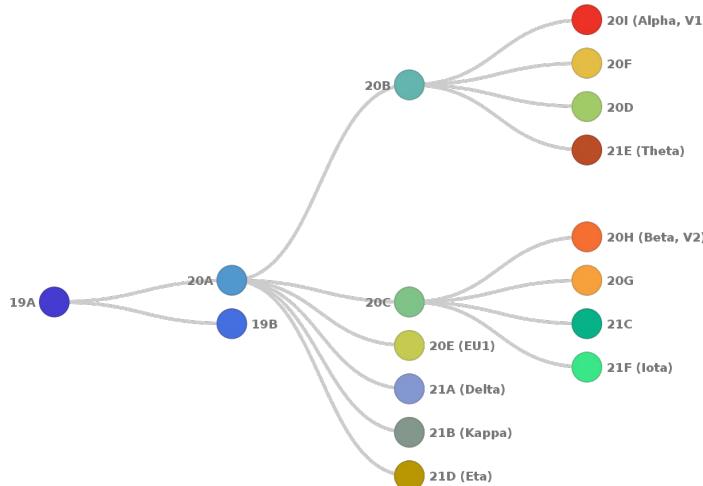
master ▾

- Commits on May 18, 2022
 - Merge pull request #536 from dgarijo/develop ...
dgarijo committed 21 days ago ✓
 - Merge branch 'master' into develop
dgarijo committed 21 days ago ✓
 - Merge pull request #531 from dgarijo/dependabot/github_actions/develop...
dgarijo committed 21 days ago ...
 - Merge pull request #532 from dgarijo/dependabot/github_actions/develop...
dgarijo committed 21 days ago ...
 - Merge pull request #533 from dgarijo/dependabot/github_actions/develop...
dgarijo committed 21 days ago ...

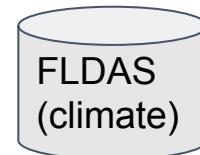
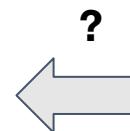
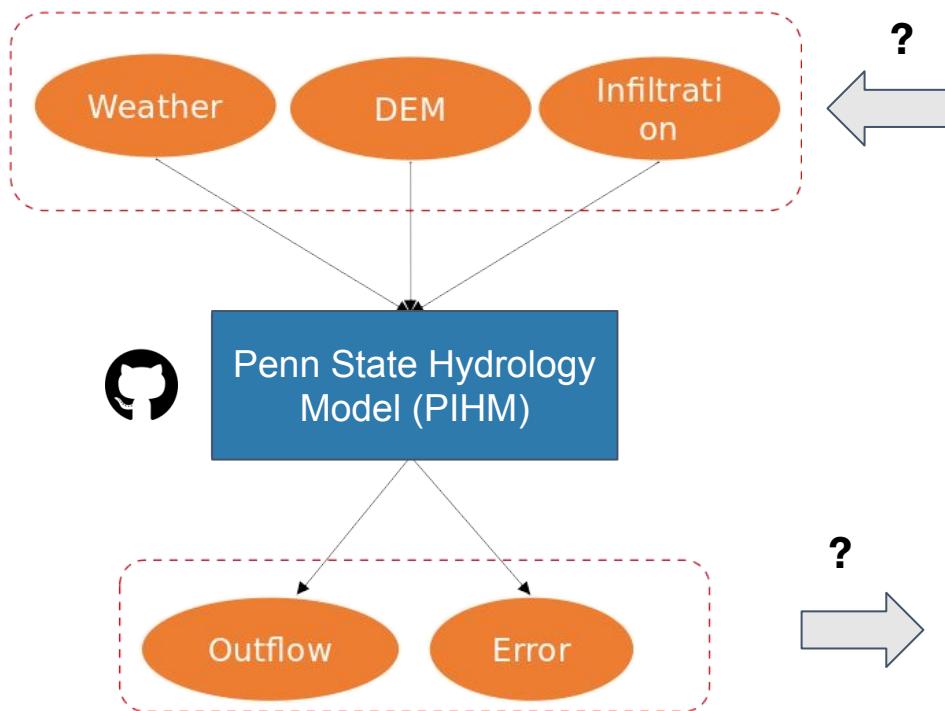
Contributions?
Development process?

Why is it needed?

- Credit
- Responsibility
- Accountability



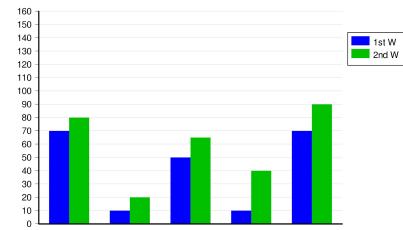
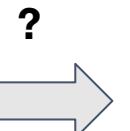
Describing inputs, outputs and their structure



Why is it needed?

- Reuse
- Reproducibility

Adapt new sources



Visualize result

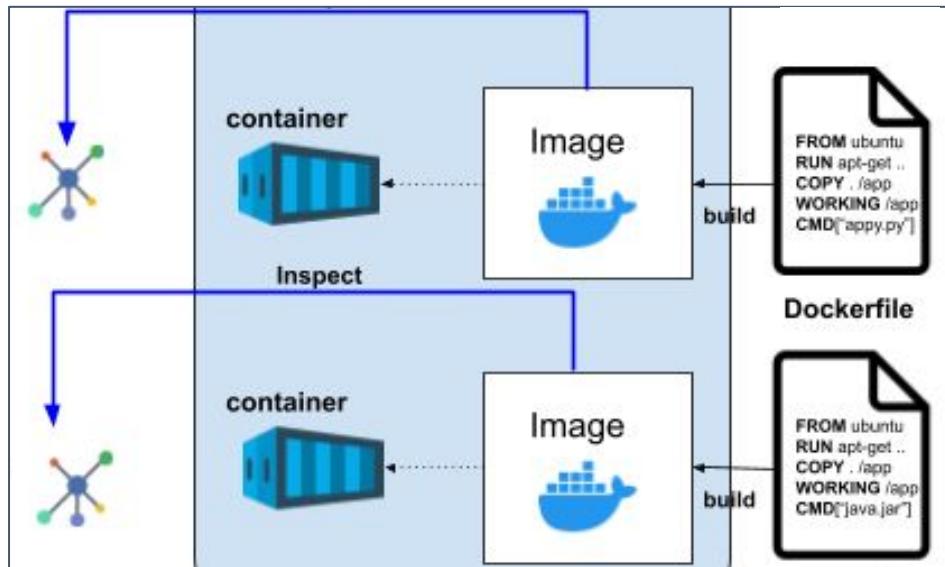


```
install_requires = [  
    "bs4==0.0.1",  
    "Click==7.0",  
    "click-option-group==0.5.3",  
    "markdown==3.3.6",  
    "matplotlib==3.5.0",  
    "nltk==3.6.6",  
    "numpy==1.22.0",  
    "pandas==1.3.4",  
    "rdflib==6.0.2",  
    "rdflib-jsonld==0.6.2",  
    "requests>=2.22.0",  
    "scikit-learn==1.0",  
    "textblob==0.17.1",  
    "validators==0.18.2",  
    "xgboost==1.5.0"  
]
```

Why is it needed?

- Reuse
- Security

Software images are created from configuration files (e.g., Dockerfiles)



Why is it needed?
- Reuse
- Security

Fig. by Jhon Toledo

Initial effort transforming part of DockerHub: <https://dockerpedia.inf.utfsm.cl/>

Osorio, M., Buil-Aranda, C., Santana-Perez, I., & Garijo, D. (2022). DockerPedia: A Knowledge Graph of Software Images and Their Metadata. *International Journal of Software Engineering and Knowledge Engineering*, 32(01), 71-89.

Creating KGs of Research Software metadata:

Knowledge extraction

Research problem: Harvesting Research Software metadata

Research Software metadata is not ~~abundant~~ machine readable

Can you please describe your software component with metadata?

I already did! Did you read the project readme?

Did you see the online documentation?

Perhaps the you saw the paper?



Many domain-specific registries are **curated by hand by experts**

- Documentation
 - Text classification
 - Named entity recognition and relation extraction

- Code
 - Static code analysis

 docs	update doc	13 days ago
 experiments	Added pipeline missed in previous version of create_models	8 months ago
 notebook	Fix #180	15 months ago
 src/somef	update version	13 days ago
 .gitignore	Fix test and added env to gitignore	29 days ago
 .readthedocs.yml	documentation	2 years ago
 CITATION.cff	Add citation file	4 months ago
 Dockerfile	updating Docker image	4 months ago
 LICENSE	initial cleanup	2 years ago
 README.md	update doc	13 days ago
 config.json	Created script to generate models and updated python version to 3.9	8 months ago
 mkdocs.yml	Fix #178	15 months ago
 pyproject.toml	minor package changes	4 months ago
 setup.py	Fix #437	28 days ago

Text classification: Software Metadata Extraction Framework

<https://github.com/KnowledgeCaptureAndDiscovery/somef/>



Repository



Results (Metadata)

dgarijo Merge pull request #174 from KnowledgeCaptureAndDiscovery/dev	
docs	Typos
experiments	Improved header analysis. Fix #166
notebook	Fix #96
src	Typos
.gitignore	Fix #147 and working towards automatic corpus va
.readthedocs.yml	documentation
Dockerfile	Fix #113 creating a Dockerfile
LICENSE	initial cleanup
README.md	Typos
config.json	Provide Fix for issues - 12, 35,36
mkdocs.yml	typos and reorganization
setup.py	Fix #113 creating a Dockerfile

- **Readme Analysis**
 - Supervised classification
 - Regular expressions
 - Header analysis
- **File exploration**
 - Notebooks
 - Dockerfiles
 - Documentation
- **GitHub API**



CodeMeta



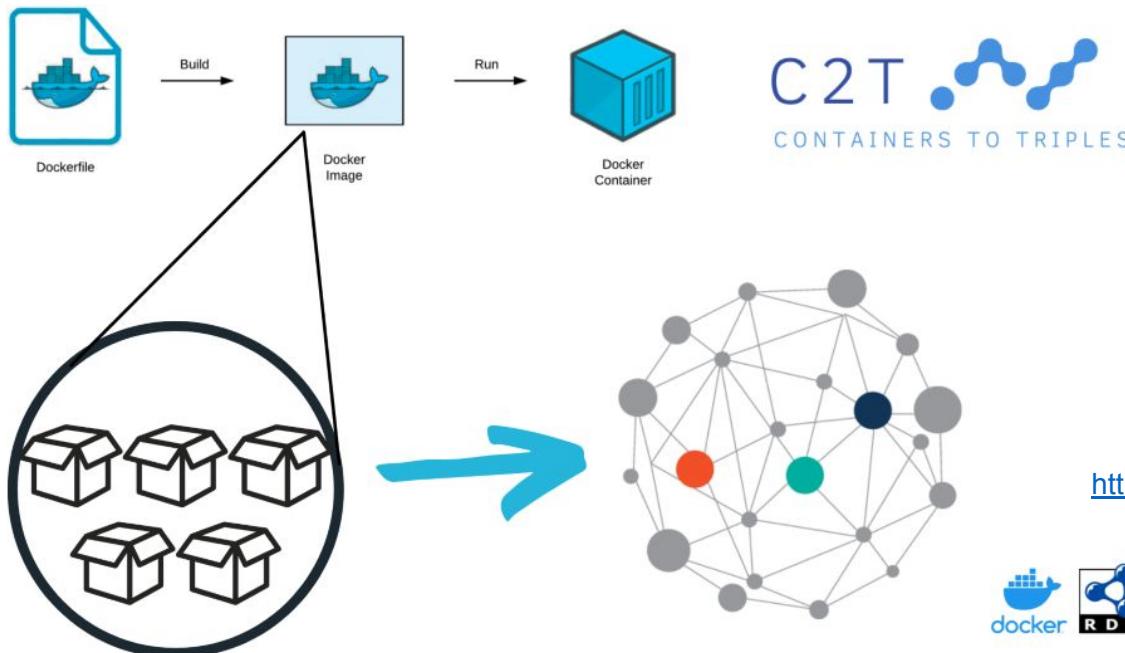
Kelley, A., & Garijo, D. (2021). A framework for creating knowledge graphs of scientific software metadata. *Quantitative Science Studies*, 1-37.

- Name (GA)
- Full title (RE)
- Description (SC, HA)
- Citation (SC, RE, HA)
- Installation instructions (SC, HA)
- Invocation (SC)
- Usage examples (HA)
- Documentation (HA, FE)
- Requirements (HA)
- Contributors (HA)
- FAQ (HA)
- Support (HA)
- License (GA, HA, FE)
- Stars (GA)
- Contact (HA)
- Download URL (HA, GA)
- DOI (RE)
- DockerFile (FE)
- Notebooks (FE)
- Executable notebooks (Binder, Collab) (RE)
- Owner: (GA)
- Keywords (GA)
- Source code (GA)
- Releases (GA)
- Changelog (GA)
- Issue tracker (GA)
- Programming languages (GA)
- Acknowledgements (HA)
- Logos (RE)
- Images (RE)
- Shell scripts (FE)
- Code of conduct (FE)
- Repository status (RE)
- Arxiv links (RE)
- Support channels (RE)
- Software category (SC) (Work in progress)
- ...

Method used (provenance):

- Supervised Classification (SC)
- Header Analysis and Synset comparison (HA)
- File Exploration (FE)
- Regular Expressions (RE)
- GitHub API (GA)

Converting containers and code into KGs



<https://osoc-es.github.io/c2t/website/>



<https://github.com/SoftwareUnderstanding/inspect4py>

Creating KGs of Research Software metadata:

Benefits

Early result: Automated software catalogs

 SOCA
SOFTWARE CATALOG
CREATOR

Software Catalog

Search for repositories...

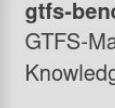
Title Stars Releases Last updated

Filter icons: 🔍 ⚖️ 📈 📄 “ doi 📈 📈 ? 📄 📈 📈 📈 📈

Morph-OME 
Online Mapping Editor

6 ★ v.2.1 3 ↗

📖 ⚖️ 📈 “ doi 📈

gtfs-bench 
GTFS-Madrid-Bench: A Benchmark for Knowledge Graph Construction Engines

11 ★ v1.2.2 5 ↗

📖 ⚖️ 📈 “ 📈 📈 📈

morph-csv 
Enhancing virtual KG access over tabular data with RML and CSVW

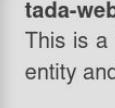
8 ★ v1.1.0 3 ↗

📖 ⚖️ 📈 “ 📈 📈 📈

pytada-hdt-entity 
A python library binding of the c++ library tada-hdt-entity

0 ★ v1.8 3 ↗

📖 ⚖️ 📈 “ doi 📈 📈

tada-web 
This is a web API project using tada-hdt-entity and the pytada-hdt-entity libraries

0 ★ v1.0 1 ↗

📖 ⚖️ “ doi 📈 📈

Widoco 
Wizard for documenting ontologies. WIDOCO is a step by step generator of HTML templates with the documentation of your ontology. It uses the LODE environment to create part of the template.

0 ★ 0 ↗

📖 ⚖️ 📈 “ doi 📈 📈

Useful for

- Comparison
- Exploring
- Reuse

Alpha available at: <https://software.oeg.fi.upm.es/> Github: <https://github.com/oeg-upm/soca>

morph-kgc

Powerful RDF Knowledge Graph Generation with [R2]RML Mappings

short description

notebooks



License

Apache License 2.0

Description:

A permissive license whose main conditions require preservation of copyright and license notices. Contributors provide an express grant of patent rights. Licensed works, modifications, and larger works may be distributed under different terms and without source code.

Permissions:

1. Commercial-use
2. Modifications
3. Distribution
4. Patent-use
5. Private-use

Usage

Learn quickly with the tutorial in [Google Colaboratory!](#)

PyPi is the fastest way to install Morph-KGC:

```
pip install morph-kgc
```

We recommend to use [virtual environments](#) to install Morph-KGC.
To run the engine via [command line](#) you just need to execute the following:

```
python3 -m morph_kgc config.ini
```

Check the [documentation](#) to can see how to generate the configuration INI file.
[Here](#) you can also see an example INI file.
It is also possible to run Morph-KGC as a [library](#) with [RDFLib](#) and [Oxigraph](#):

```
import morph_kgc
```

```
# generate the triples and load them to an RDFLib graph
g_rdflib = morph_kgc.materialize('/path/to/config.ini')
# work with the RDFLib graph
q_res = g_rdflib.query(' SELECT DISTINCT ?classes WHERE { ?s a ?classes } ')
```

```
# generate the triples and load them to Oxigraph
g_oxigraph = morph_kgc.materialize_oxigraph('/path/to/config.ini')
# work with Oxigraph
q_res = graph.query(' SELECT DISTINCT ?classes WHERE { ?s a ?classes } ')
```

```
# the methods above also accept the config as a string
config = """
[DataSource1]
mappings: /path/to/mapping/mapping_file.rml.ttl
db_url: mysql+pymysql://user:password@localhost
"""
g_rdflib = morph_kgc.materialize(config)
```

How to use it

```
python /morph-kgc/oeg-upm_morph-kgc/morph-kgc-main/src/morph_kgc/main.py
```

{Morph-KGC: Scalable Knowledge Graph Materialization with Mapping Partitions}

Citation

```
@article{arenas2022morph,
  title  = {{Morph-KGC: Scalable Knowledge Graph Materialization with I}},
  author = {Arenas-Guerrero, Julián and Chaves-Fraga, David and Toledo},
  journal = {Semantic Web},
  year   = {2022},
  url    = {http://www.semantic-web-journal.net/system/files/swj3135.p}
```

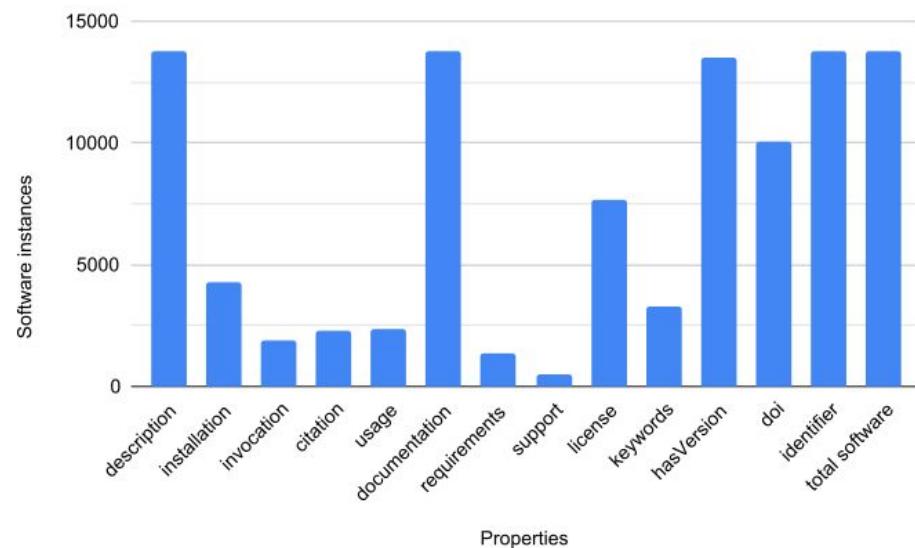
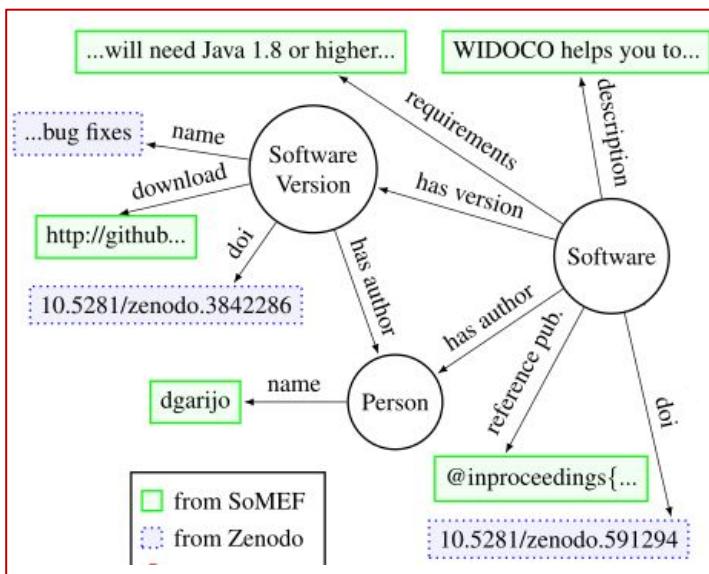
invocation

Dublin Core Metadata Initiative Conference, October, 2022

28

Extracting KGs from thousands of Open Source repositories

- Zenodo software (> 12000)
- Measuring best practices based on metadata



Summing up

Research software is a critical asset for Open Science

- Access information in structured, homogeneous manner
- Reusability
- Comparison
- Search



Pending Challenges:

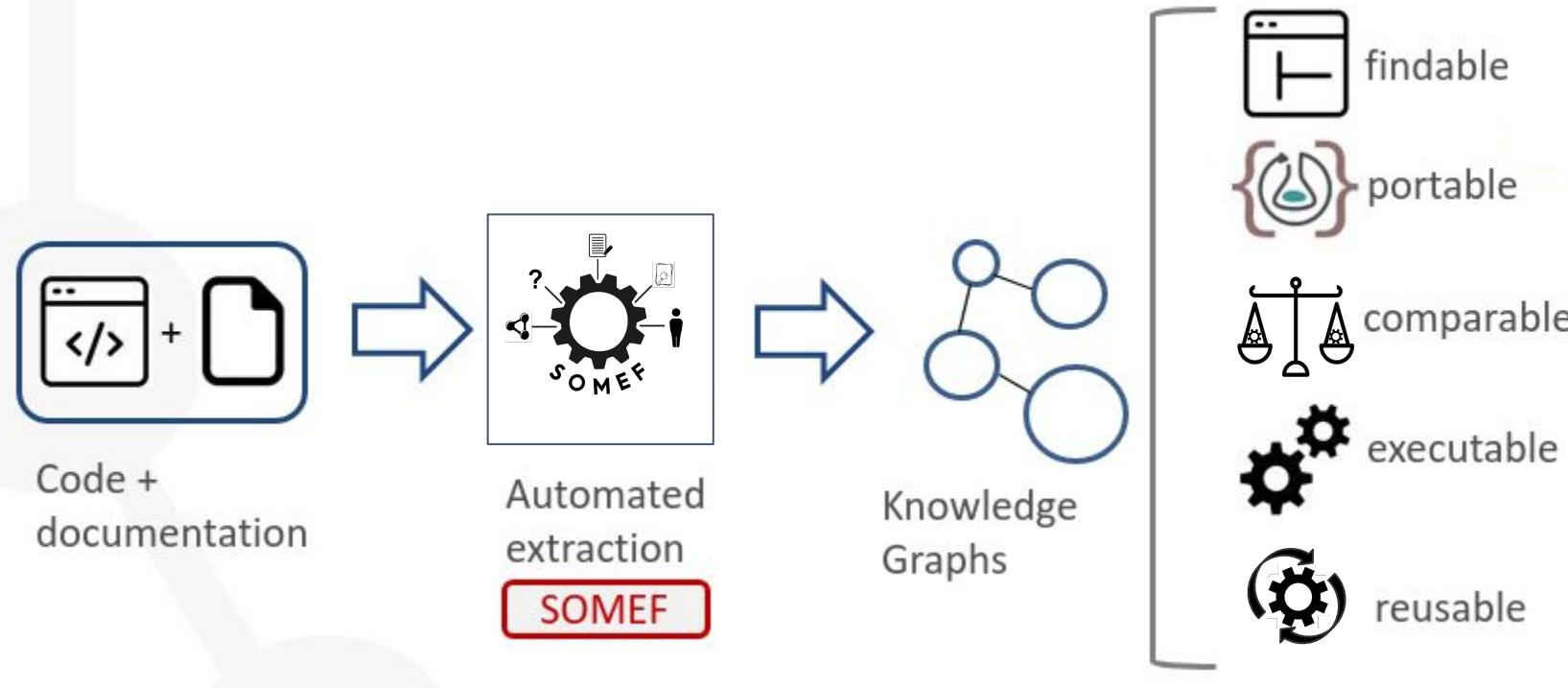
- Automated metadata extraction from existing sources
- Curation
- Reconciliation of entities (through KGs like Wikidata)
- Representing all levels of granularity

Acknowledgements



Thanks to Yolanda Gil, Varun Ratnakar, Maximiliano Osorio, Hernán Vargas, Deborah Khider, Allen Mao, Aidan Kelley, Haripriya Dharmala, Jiajing Wang, Rosa Filgueira, Pablo Calleja, Oscar Corcho, Laura Camacho, Jhon Toledo, Miguel Angel García, Esteban Gonzalez, Elena Montiel, Elvira Amador & all the students at UPM and USC who participated in the initiatives mentioned in this presentation

This work has been supported by the Madrid Government (Comunidad de Madrid-Spain) under the Multiannual Agreement with Universidad Politécnica de Madrid in the line Support for R&D projects for Beatriz Galindo researchers, in the context of the V PRICIT (Regional Programme of Research and Technological Innovation)



Let's create **machine-actionable** software metadata to promote Open Science!