



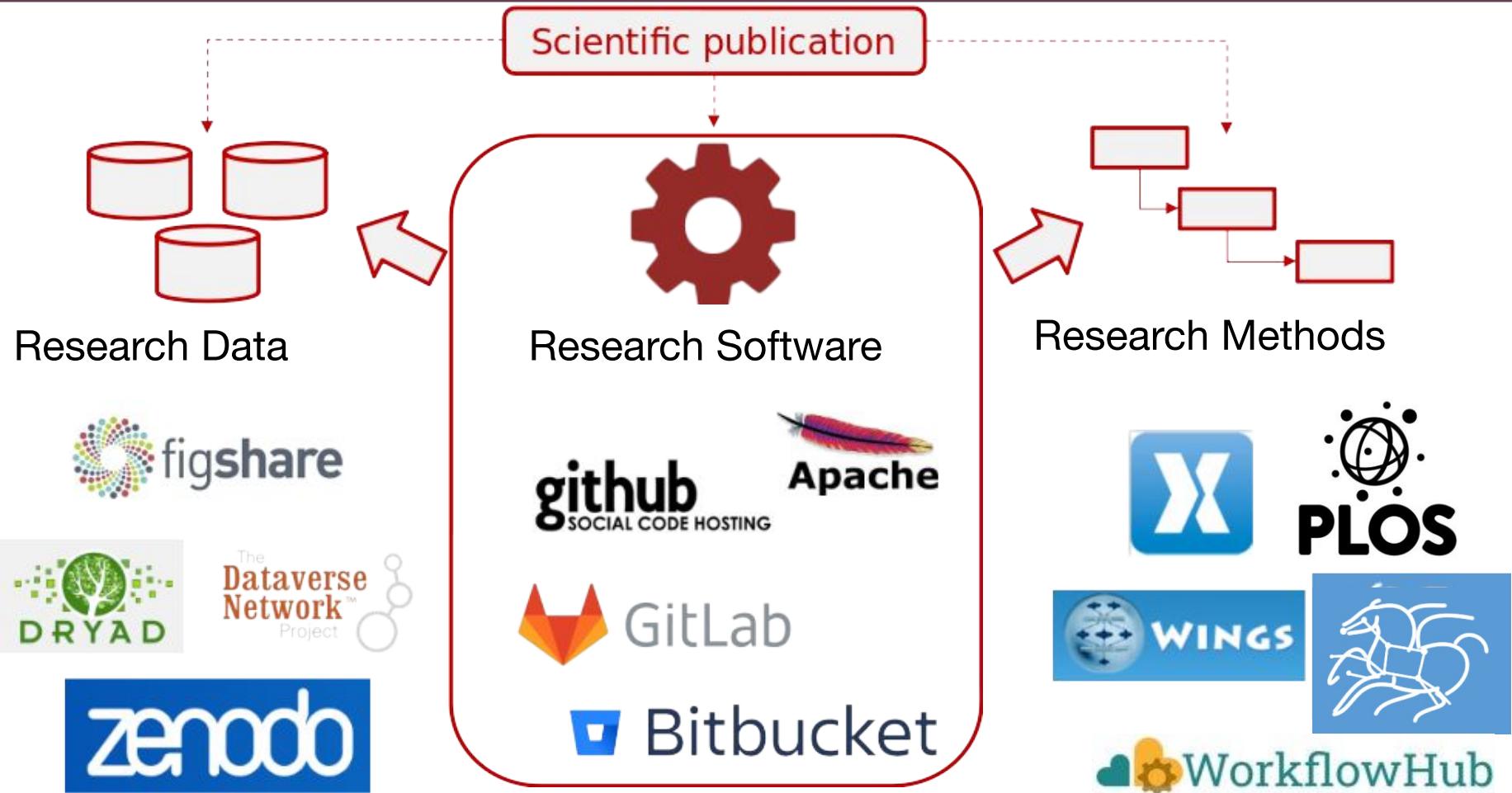
Open Science, Research Software and AI: A year at the Ontology Engineering Group

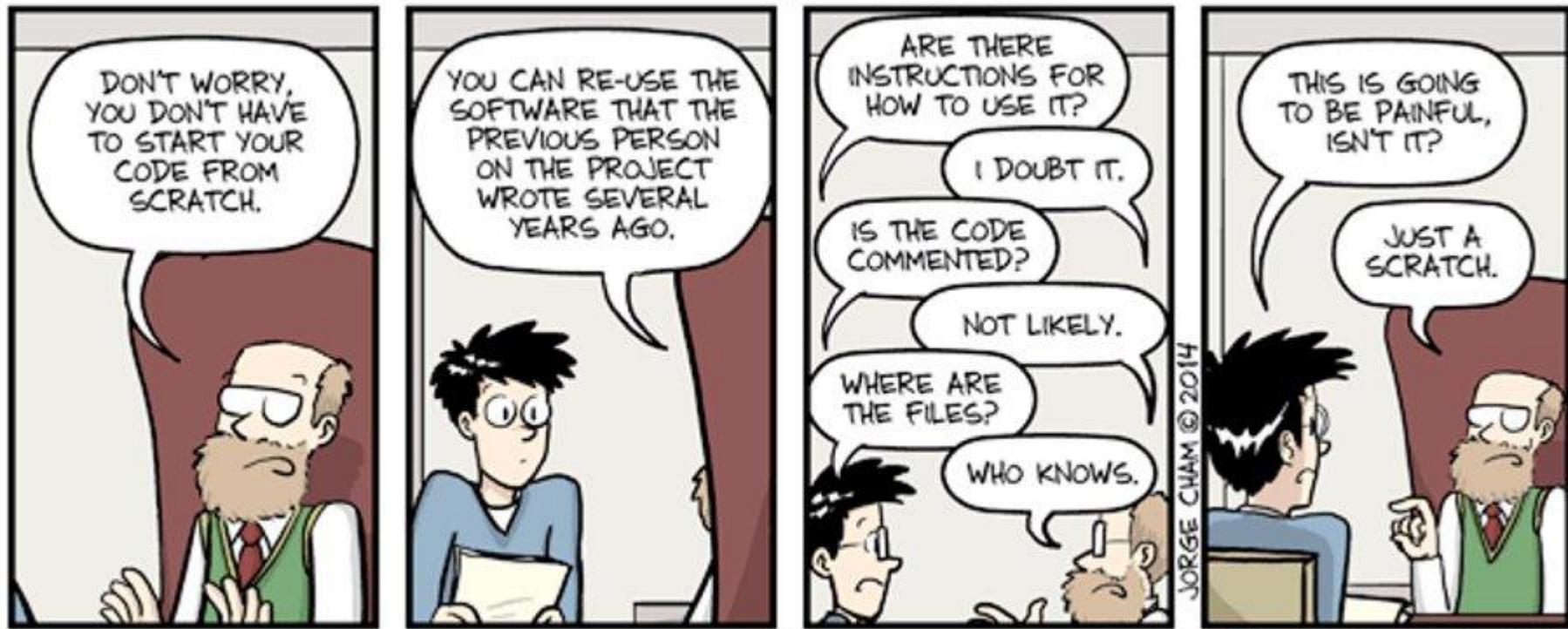
**Daniel Garijo, Ontology Engineering Group,
Universidad Politécnica de Madrid, Spain**

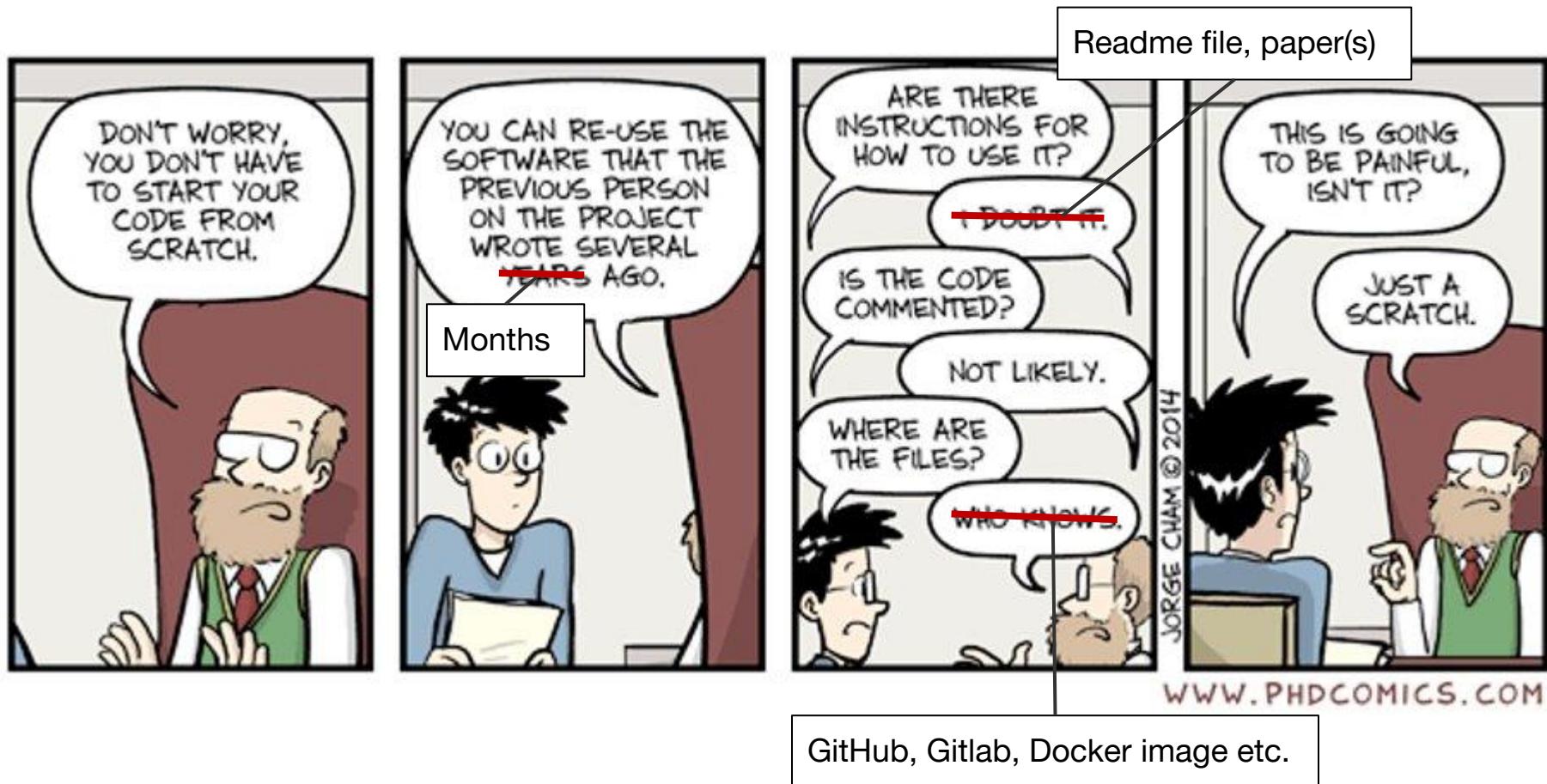
Slides based on my previous group report (July, 2021):
https://docs.google.com/presentation/d/14MExu9NXPMFhwsSDg2bd_uZAWfhuYsCNbmbfj1B42jQU/edit?usp=sharing

 daniel.garijo@upm.es
 @dgarijov

Research Software is one of the pillars of Open Science



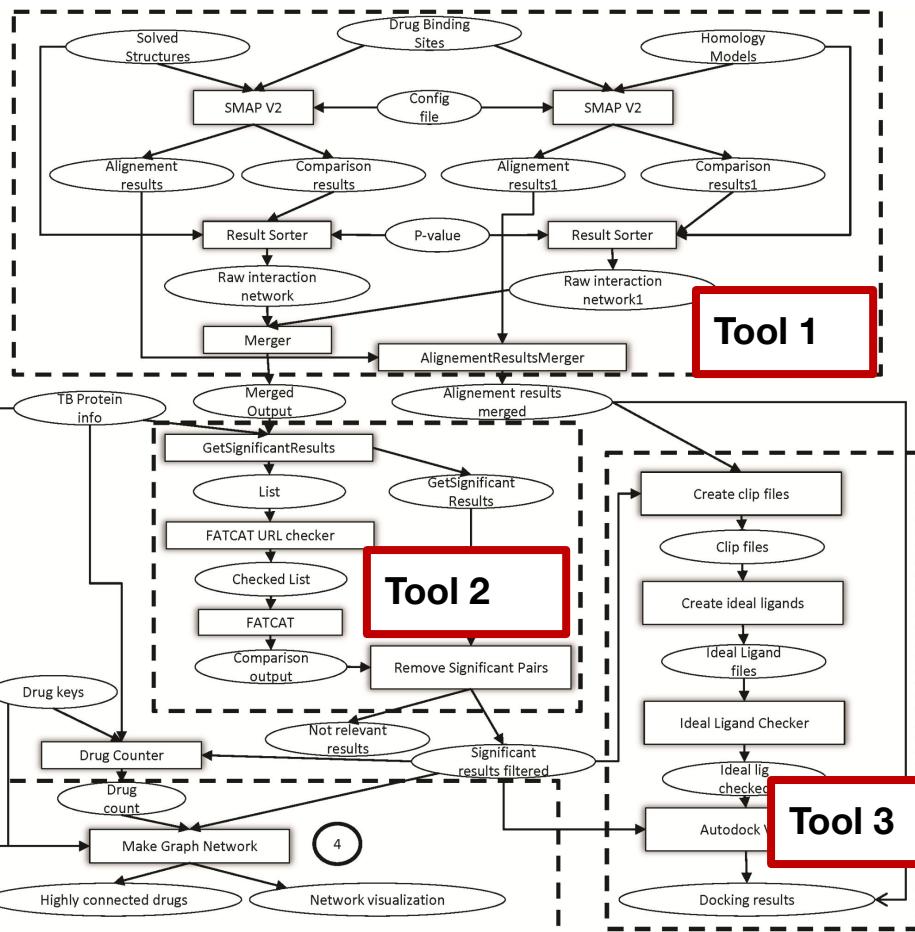




Reusability takes time, even when sources are available

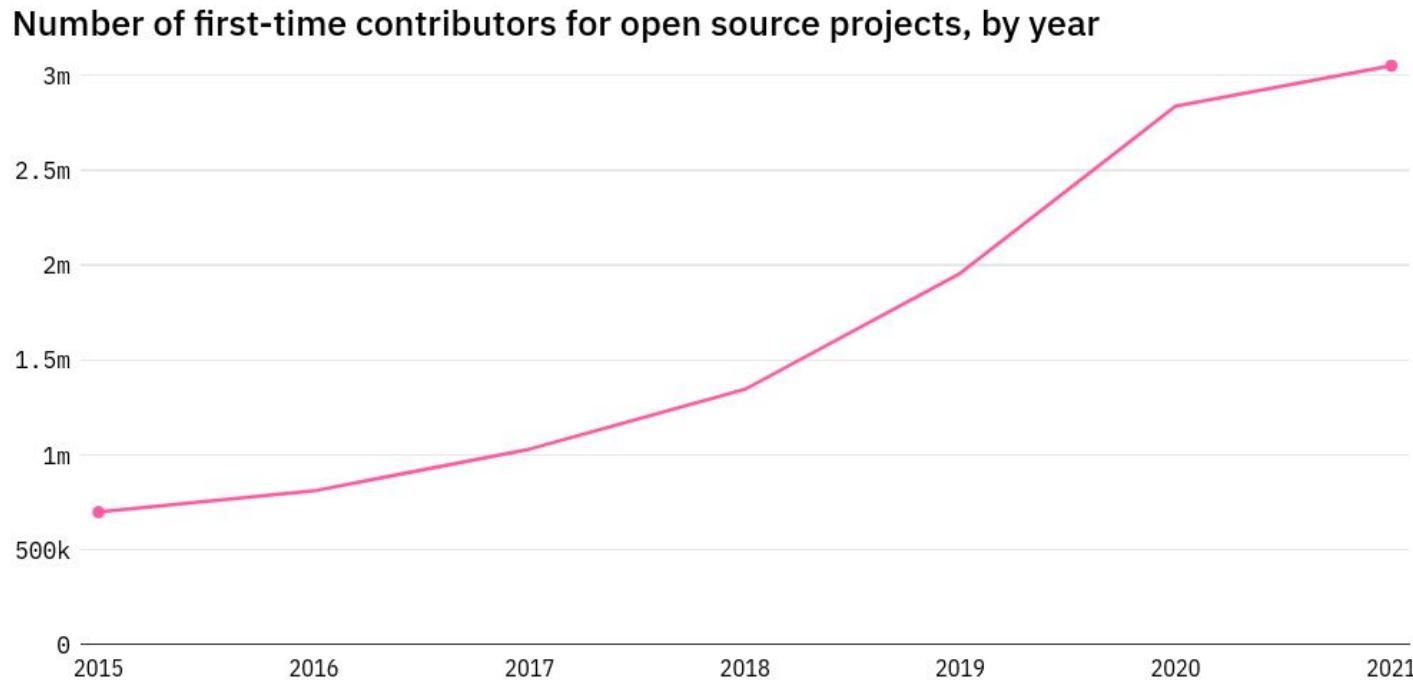
In [1] we tried to reproduce an effort from **one** year before.

- 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



[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.

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



Source: GitHub Octoverse Report 2021

TECH MONITOR

Can we **automatically** accelerate
software understanding?

Describe



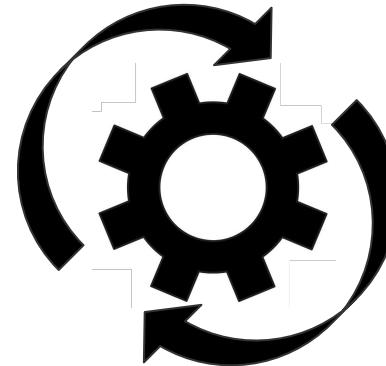
- what is in it?
- what's about?
- examples?
- relation to other resources?

Compare

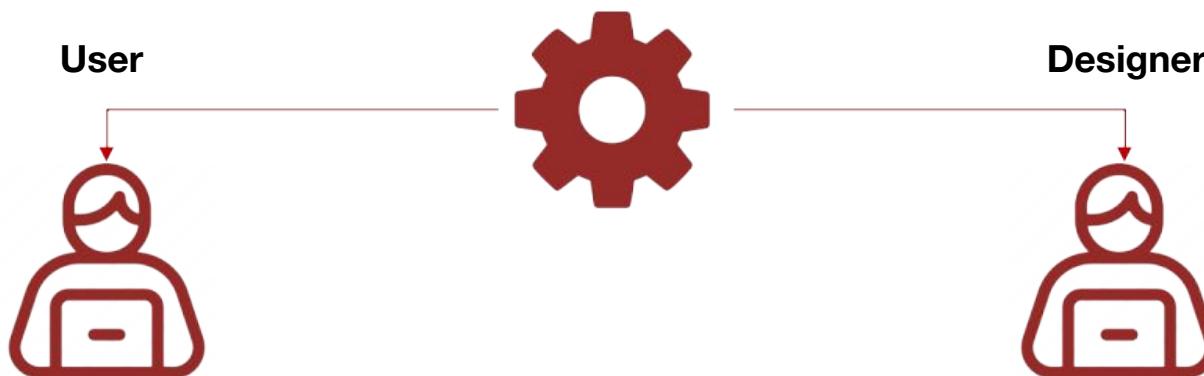


- similarities?
- differences?
- features?

Reuse



- run?
- repeat?
- reproduce?
- fix?
- combine?



- What does the **software component do**? Which of its methods should I use?
 - How **to transform my data** to use the software component?
 - How **to interpret the results** produced by the software component?
 - How **to invoke** the software component?
 - How **to configure** the software component with the right parameters?
 - How to **compare** against similar methods?
-
- How to ease **capturing the dependencies and installation instructions** of my software?
 - How to **encapsulate my software** so it can be used with other data?
 - How **to describe my software** so it can be used by others?
 - How **to test if my software** is ready to be used by others?

Before returning to the Ontology Engineering Group...



OntoSoft: Creating a Software registry for the Geosciences



<https://ontosoft.org/>

OntoSoft: Capturing Scientific Software Metadata. Eighth ACM International Conference on Knowledge Capture, Palisades, NY, 2015

Representing Software Metadata: OntoSoft

The screenshot shows two main views of the OntoSoft Software Repository. The top view is a 'Software List' showing five entries: DrEICH algorithm, PIHM, PIHMgls, TauDEM, and WBMsed. Each entry has a 'PUBLISH YOUR SOFTWARE' button and an 'EDIT' button. The bottom view is a 'Filter Software List' for the PIHM entry, displaying its metadata: Author (Tobias Körner), Keywords (Hydrological model OR Hydrology), Language (C++), and License (GNU General Public License v2.0). A link to the full software record is provided. An arrow points from the bottom view to the detailed software record below.

Software Repository
Describe your software so others can find and use it

PUBLISH YOUR SOFTWARE

Software List

COMPARE

Name

DrEICH algorithm

PIHM

PIHMgls

TauDEM

WBMsed

EDIT

Filter Software List

Search

Author

Keywords: Hydrological model
OR Hydrology

Language: C++

License: GNU General Public License v2.0

GNU General Public License v2.0

LOCATE

Importance Optional

What is the software called ? PIHM

What is a short description of this software ? PIHM is a multi-process, multi-scale hydrologic model where the major hydrological processes are fully coupled using the semi-discrete finite volume method. PIHM is a physical model for surface and groundwater, "tightly-coupled" to a D3D interface, PIHM-D3D, which is open source, platform independent and extensible. The tight coupling between D3D and the model is achieved by developing a shared data-model and hydrologic-model data structure.

Initial metadata was retrieved from http://cdsweb.cgd.ucar.edu/wiki/Model_PiHM

What are general categories (keywords, labels) for this software ? Hydrology, Basins, Continental

Is there a project website for this software ? <http://www.pihm.praxi.edu/githome.html>

Crowdsourced Software Metadata Registry

- Complements code repositories to make them **understandable**
- Software metadata **designed for scientists**
- Metadata is **curated by decentralized** communities of users
- **Training scientists** on best practices



<http://ontosoft.org>

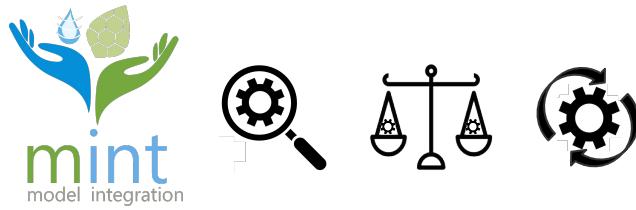
OntoSoft Software Community Training

Compare Software

DrEICH algorithm, PIHM, PIHMgis, TauDEM, WBMsed

| PIHM | PIHMgis | DrEICH | TauDEM | WBMsed |
|---|---|------------------------------------|---|---|
| | | | | |
| What are domain specific keywords for this software ? (eg: hydrology, climate) | | | | |
| Geomorphology, Hydrological, Bedrock channel ero- | Basins, Continental | Basins, GIS | Hydrologically corrected DEM, Watershed | Sediment flux, Global model, Hydrological model |
| What Operating Systems can the software run on ? | | | | |
| Unix Linux | Unix Windows Linux Mac OS | Unix Windows Linux Mac OS | Unix Windows Linux Mac OS | Unix Linux |
| Is there any test data available for the software ? | | | | |
| Test Data Location: http://onlinelibrary.wiley.com/doi/10.1002/2013WR015167/full | Test Data Location: http://source-forg.../pihmmodel/ | | Test Data Location: http://csdms.colorado.edu/wiki/Model:TauDEM#Testing | Test Data Location: http://csdms.colorado.edu/wiki/Mod...el:WBMsed#Testing |
| Test Data Description: Two test DEMs are included in the repository, | Test Data Description: Upper Juniata River 875 km ² : see: http://source-forg.../pihmmodel/ | | Test Data Description: The Logan River DEM is a small test dataset useful | Test Data Description: Extensive input dataset is available on the CSDMS |

MINT: Facilitating Model Integration in the Geosciences



<https://mint-project.info/>

Yolanda Gil, Daniel Garijo, Deborah Khider, et.al. 2021. Artificial Intelligence for Modeling Complex Systems: Taming the Complexity of Expert Models to Improve Decision Making. *ACM Trans. Interact. Intell. Syst.* 11, 2, Article 11 (June 2021), 49 pages. DOI:<https://doi.org/10.1145/3453172>

Adding structure to software Metadata: MINT (Model Integration)

Low grained machine-readable Software Metadata:

- (From OntoSoft) Attribution, license, funding, usage examples....
- Executable software components
- Software invocation
- Input & output files, variables and units
- Containers used to encapsulate and run software component

Input/output variables

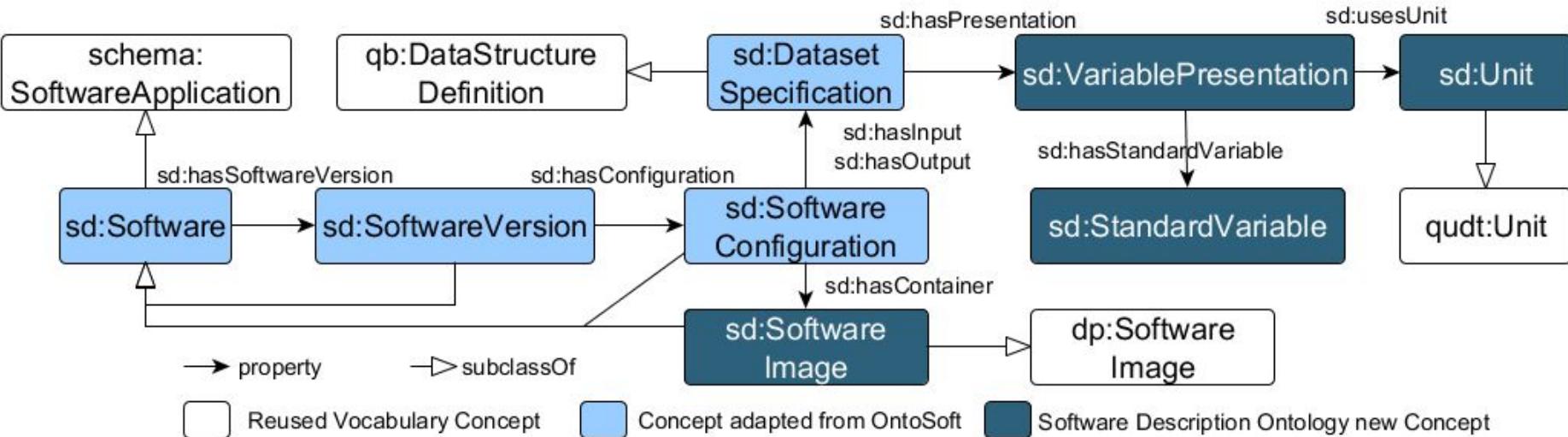
formation of river segments

IO Files:

| | Name | Description |
|--------|-------------------|---|
| INPUT | pihm-riv | Spatial geometry and material information of river segments |
| INPUT | pihm-geol | Geologic file |
| INPUT | pihm-ibc | Boundary condition information for elements |
| INPUT | pihm-modelinfo | PIHM model information aggregation file |
| INPUT | pihm-lc | Vegetation parameters of different land cover types |
| INPUT | pihm-base | Base file |
| INPUT | pihm-forc | PIHM forcing file with the majority of the relevant variables |
| INPUT | pihm-soil | Soil parameters for the soil types |
| INPUT | pihm-att | PIHM attribute file with index values of variables for timeseries |
| OUTPUT | pihm-et0 | Evaporation canopy file |
| OUTPUT | pihm-rivFlx9 | lateral outflux to the bed beneath river |
| OUTPUT | pihm-rivFlx4 | Baseflow to stream reach from aquifer on the left |
| OUTPUT | pihm-rech | Recharge Rate file |
| OUTPUT | pihm-rivFlx10 | lateral influx to the bed beneath river |
| OUTPUT | pihm-infiltration | Infiltration file |

pihm-riv

| Label | Long Name | Description | Standard Name | Units |
|-------------------|----------------------------|----------------------------|---|---------|
| Bed | Bed Depth | Bed Depth | channel_bed_thickness | m |
| KsatV | Bed Hydraulic Conductivity | Bed Hydraulic Conductivity | soil_water__vertical_saturated_hydraulic_conductivity | m day-1 |
| Water table value | Water table of the IC | Water table of the IC | | m |

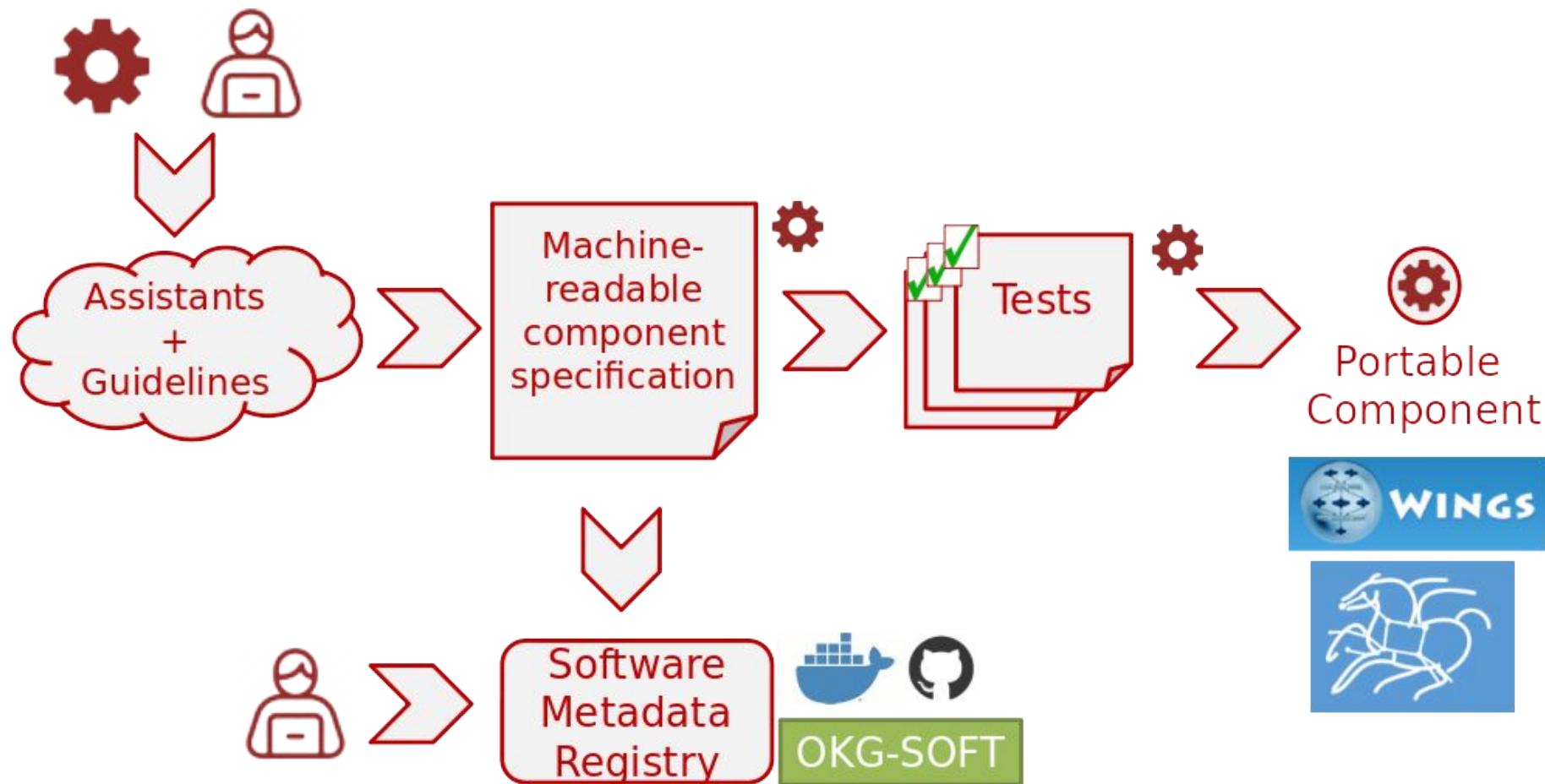


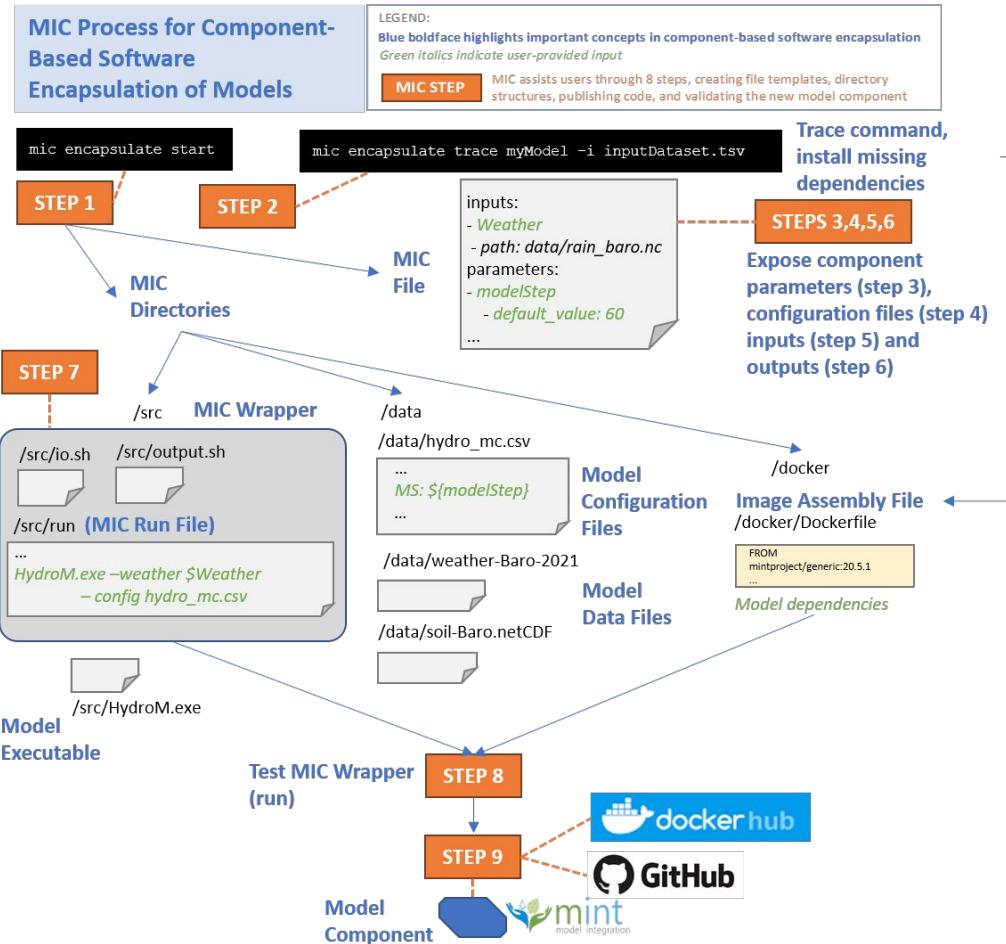
Extending:

- Schema.org/Codemeta (software metadata)
- W3C Data Cubes (Contents of inputs and outputs)
- NASA QUDT (Units)
- DockerPedia (Software images)
- Scientific Variables Ontology (Standard Variables)

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

Garijo, D., Osorio, M., Khider, D., Ratnakar, V., & Gil, Y. (2019, September). OKG-Soft: An open knowledge graph with machine readable scientific software metadata. In *2019 15th International Conference on eScience (eScience)* (pp. 349-358). IEEE.

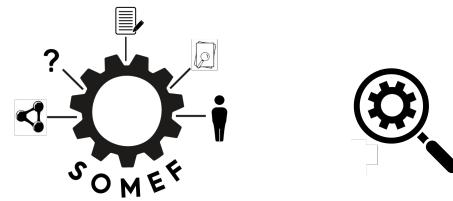




Software encapsulation methodology

- Input: **software component**
- Output:
 - Docker image
 - **Wrapper script** (GitHub)
 - **Metadata** (MINT model catalog)
- Powered by **ReproZip** (<https://www.reprozip.org/>) to automatically suggest I/O

SOMEF: Automating software metadata extraction



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

Problem: Harvesting Scientific Software Metadata

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?

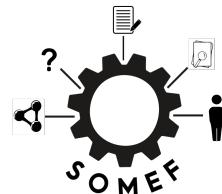
Did you see the online documentation?

Perhaps the you saw the paper?

Many domain-specific registries are curated by hand by experts



Repository



Extraction



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**



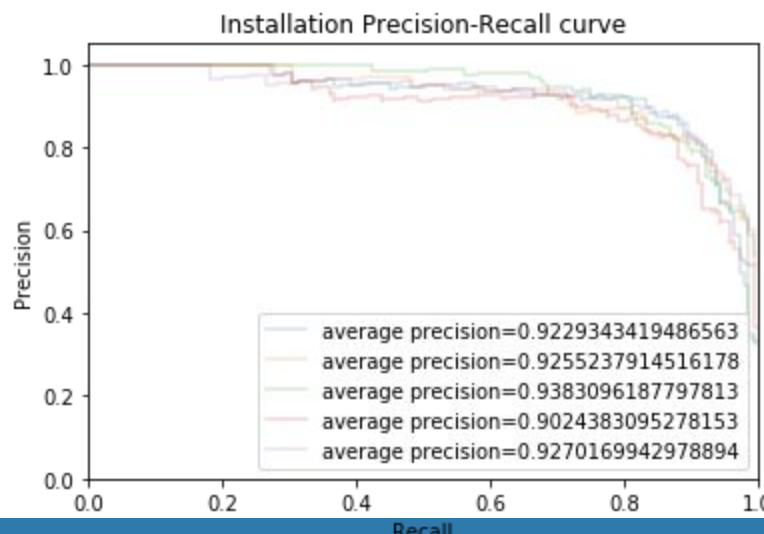
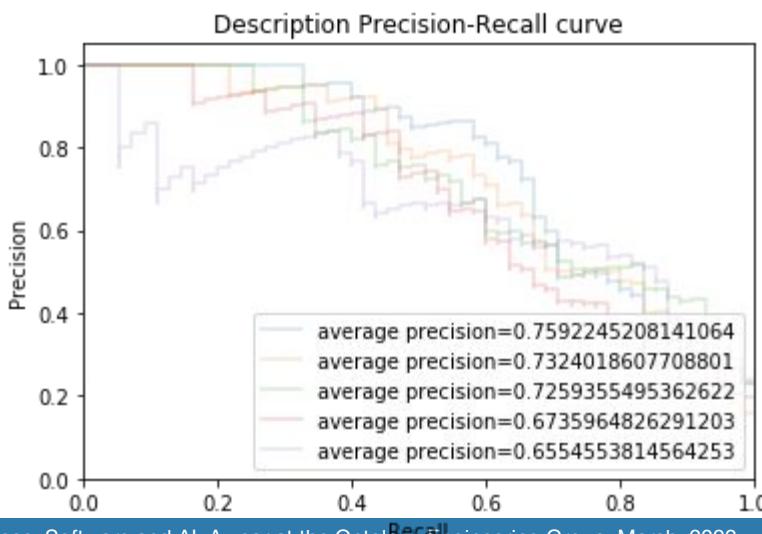
CodeMeta



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

- Paragraph-based text classification
- Four main categories:
 - Installation, citation, description, invocation.
- Binary classification problem

| Truth Value | Category | Apprx. Ratio | Count |
|-------------|--------------|--------------|-------|
| True | Description | 0.5 | 275 |
| False | Installation | 0.125 | 68 |
| | Invocation | 0.125 | 68 |
| | Citation | 0.125 | 68 |
| | Treebank | 0.125 | 68 |
| Total | | 1.0 | 547 |



- Extraction based on frequent header analysis
 - Fuzzy matching based on synsets

Installation



KGTK: Knowledge Graph Toolkit



Regular expressions, based on common practices (e.g., DOI, .bib, etc.)

The Knowledge Graph Toolkit (KGTK) is a comprehensive framework for the creation and exploitation of large hyper-relational knowledge graphs (KGs), designed for ease of use, scalability, and speed. KGTK represents KGs in tab-separated (TSV) files with four columns: edge-identifier, head, edge-label, and tail. All KGTK commands consume and produce KGs represented in this simple format, so they can be composed into pipelines to perform complex transformations on KGs. KGTK provides:

- 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)
- Stars (GA)
- Contact (HA)
- Download URL (HA, GA)
- DOI (RE)
- DockerFile (FE)
- Notebooks (FE)
- Executable notebooks (Binder) (RE)
- Owner: (GA)
- Keywords (GA)
- Source code (GA)
- Releases (GA)
- Changelog (GA)
- Issue tracker (GA)
- Programming languages (GA)
- Acknowledgements (HA)

Method used:

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



JSON

```
description:
  ▼ 0:
    ▼ excerpt: "WIDOCO helps you to publish and create an enriched and customized documentation of your classes, properties and data properties of the ontology, the OOPS! webservice by María being used. In addition, we use WebVowl to visualize the ontology and have extended Bub documentation of the terms in your ontology (based on [LODE](http://www.essepuntato.it/ annotation in JSON-LD snippets of the html produced.\n* Association of a provenance page means to complete it on the fly when generating your ontology. Check the [best practice WIDOCO].\n* Guidelines on the main sections that your document should have and how to create changelog of differences between the actual and the previous version of the ontology (both them independently and replace only those needed.\n* Content negotiation and serialization"
  ▼ confidence:
    0: 1
    technique: "wordnet"
  ▼ 1:
    ▼ excerpt: "For a complete list of the current improvements and next features, check the project website."
    ▼ confidence:
      0: 0.8231493588525339
      technique: "classifier"
  ▼ 2:
    ▼ excerpt: "Wizard for documenting ontologies. WIDOCO is a step by step generator of HTML template files for publishing ontologies on the web."
    ▼ confidence:
      0: 1
      technique: "metadata"
  citation:
    ▼ 0:
      ▼ excerpt: "@inproceedings{garijo2017widoco,\n  title={WIDOCO: a wizard for documenting ontologies},\n  organization={Springer, Cham},\n  doi = {10.1007/978-3-319-68204-4_9},\n  funding = {US National Science Foundation}\n}"
      ▼ confidence:
        0: 1
        technique: "classifier"
```





A **year** at the Ontology Engineering Group... **Research Projects**

Static code analysis in Python

- Extraction of available classes and functions
 - Documentation
- Requirements (reusing existing libraries)
- Call list
- File hierarchy
- Control flow (reusing existing libraries)
- Software invocation
 - Service? Package? Library? & invocation command
- Metadata export in JSON

Benefits

- Understanding, reuse, ML featurization, similarity, best practices



Rosa Filgueira



University of
St Andrews



Can we extract **additional context** from READMEs?

- E.g., installation instructions for Unix/Windows
- Hardware requirements and versions
- Library and OS versions
- Funding



Pablo Calleja
Marco Santucci (NER)
Andres Cardozo (RE)

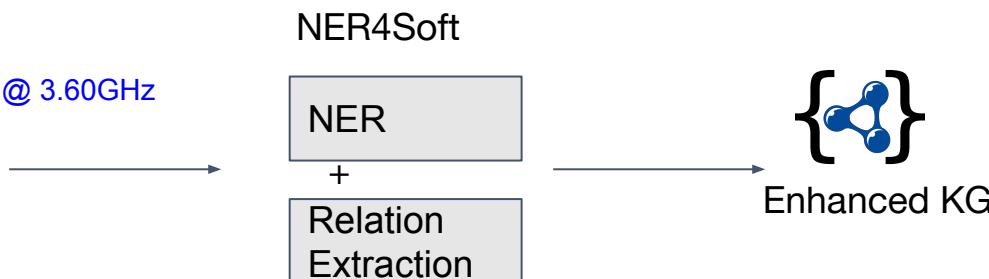
Requirements (Unix)

CPU: Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz

GPU: NVIDIA GTX1060

python 2.7:

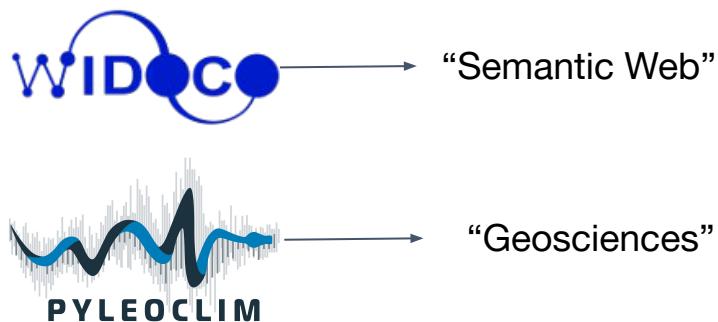
- pytorch == 0.3.1
- numpy
- opencv



<https://github.com/oeg-upm/ner4soft/>

Awesome lists capture common software projects under a common theme

- Can we use READMEs as labels for weak supervised classification?
- Text-based classification + some deep learning architectures (Bi-lstm)



Jenifer Ciuiu-Kiss

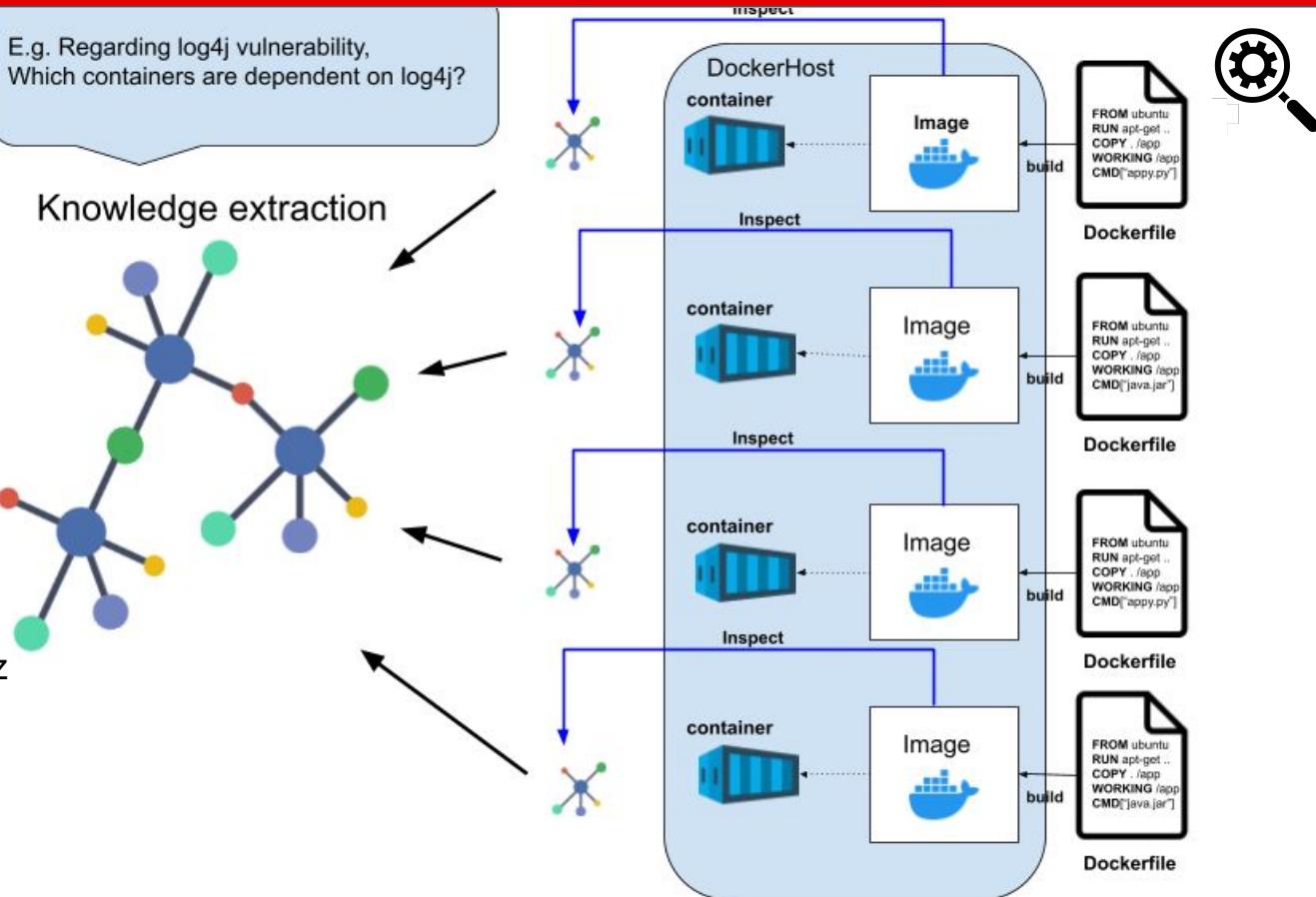


Extracting knowledge from other sources related to software (TFG)

Can we extract useful metadata from Docker images and Jupyter notebooks to gain additional insight in a repository?



Jhon Toledo
(Docker)
Daniel Rodriguez
(notebooks)



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.

How can we identify **similarities** between software components?

- Can we identify **differences**?
- What does similarity mean?
 - Metadata?
 - Description?
 - Documentation?
 - Functions?
 - KG diff?
 - Code?



Clark Wang



Topic modeling (TFM)

María Ayuso

- Are readme files appropriate for topic modeling-based similarity?
 - LDA
 - BERT-topics
 - Comparison against human-tagged categories



A **year** at the Ontology Engineering Group...
Engineering and demo projects

New (engineering) features in **SOMEF**

- Better parsing
- Pypi support
- Gitlab support
- New metadata fields
 - Colab notebooks
 - Status
 - sh scripts
 - logos
 - images
 - code of conduct
 - arxiv links
 - ...
- Additional types of files
- Better provenance (which technique was used and how)
- Robustness



Miguel Ángel García



SOMEF Vider

Victor Fernández

Description

Citation

APA Style
Mao, A., Garijo, D., & Fakhraei, S. (2019). SoMEF: A Framework for Capturing Scientific Software Metadata from its Documentation. 2019 IEEE International Conference on Big Data (Big Data), 3032–3037. <https://doi.org/10.1109/BigData47090.2019.9006447>

Bibtex
@INPROCEEDINGS{9006447, author={A. {Mao} and D. {Garijo} and S. {Fakhraei}}, booktitle={2019 IEEE International Conference on Big Data (Big Data)}, title={SoMEF: A Framework for Capturing Scientific Software Metadata from its Documentation}, year={2019}, doi={10.1109/BigData47090.2019.9006447}, url={http://dgarijo.com/papers/SoMEF.pdf}, pages={3032-3037} }

...

<https://somef.linkeddata.es/>

Software Catalog

Search for repositories...

Title Stars Releases Last updated



SOFTWARE CATALOG
CREATOR

Morph-OME

Online Mapping Editor



6 ★
v.2.1 3 ⚡



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 ⚡



tada-web

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



0 ★
v1.0 1 ⚡



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.



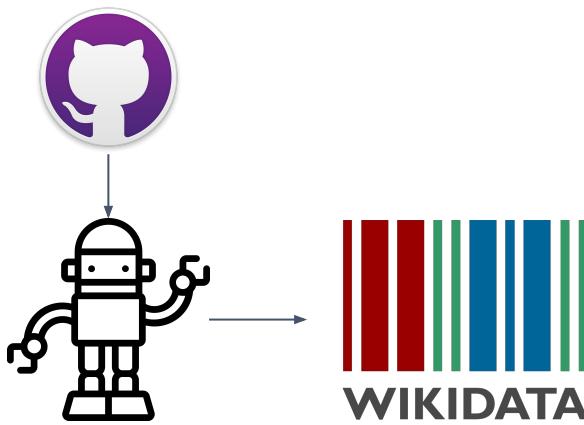
0 ★
0 ⚡



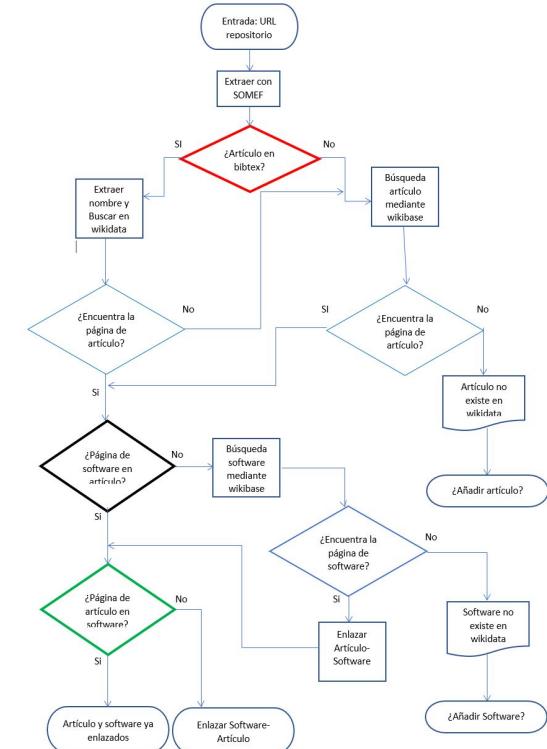
Daniel Rodriguez

alpha available at: <https://dakixr.github.io/scc/example/index.html>

Wikidata bot to link code repositories with Wikidata articles



Article exists in WD?
Software exists in WD?
Are they linked?



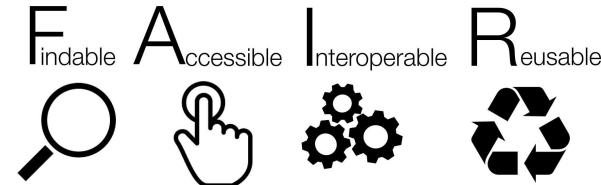
Jorge Bolinches

Summing up

Research software is a **critical asset for Open Science**

Accelerating **Software Understanding** requires:

- Automated description
- Assisted comparison
- Easy reuse



Additional challenges:

- Relation extraction between research outputs
- Create comparison benchmarks
- Representation of software metadata in RDF-star
- Can an assistant follow your **installation instructions?**
 - Tests?

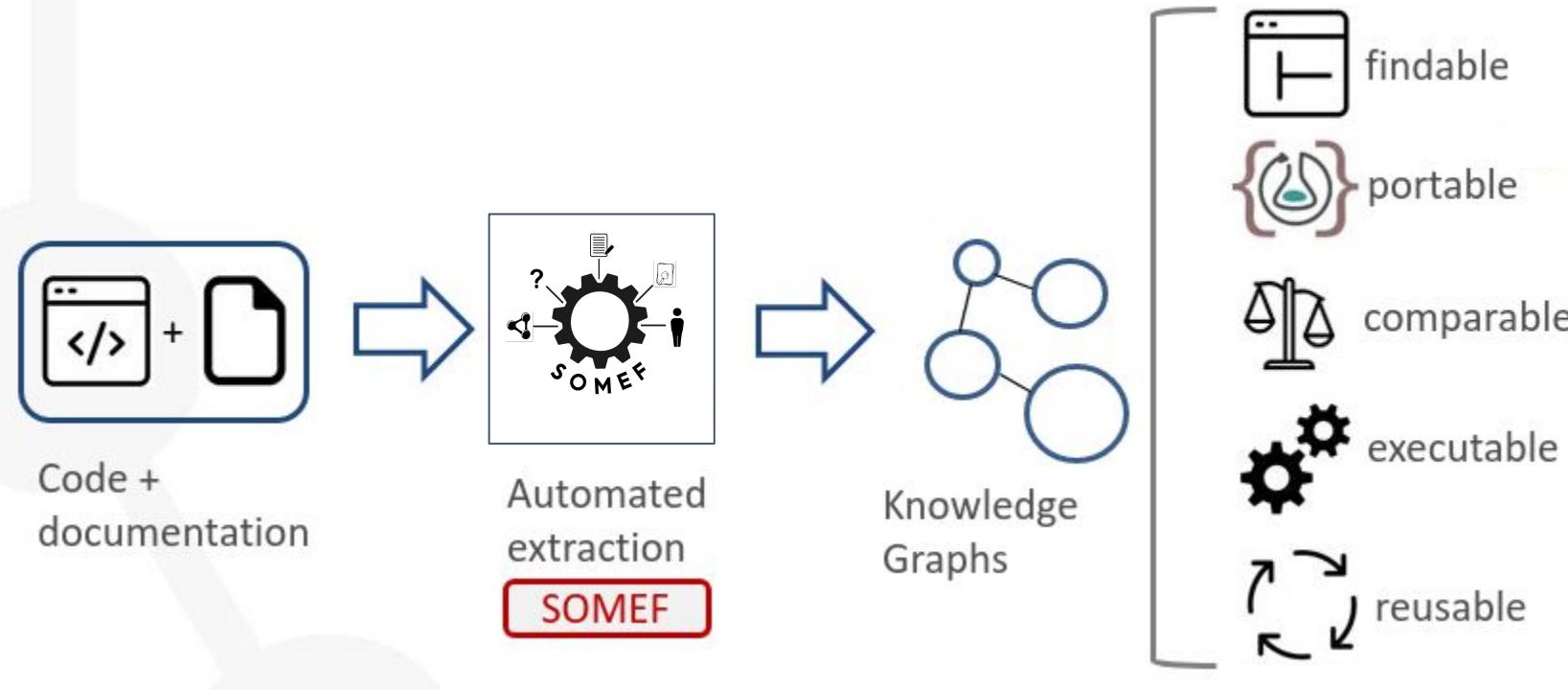


Acknowledgements



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 & 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