# Software Metadata Extraction from Code Repositories: An overview
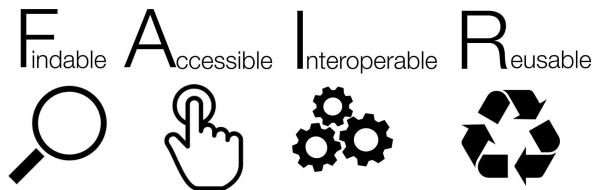
**Daniel Garijo, Ontology Engineering Group, Universidad Politécnica de Madrid, Spain**

✉ daniel.garijo@upm.es

🐦 @dgarijov

Now you are familiar with **Codemeta**

- Community standard extending Schema.org

- **JSON-LD**

- Metadata is key for **FAIR**



https://codemeta.github.io/terms/

## Describe



**Given a software project:**
- What is it about?
- Examples?
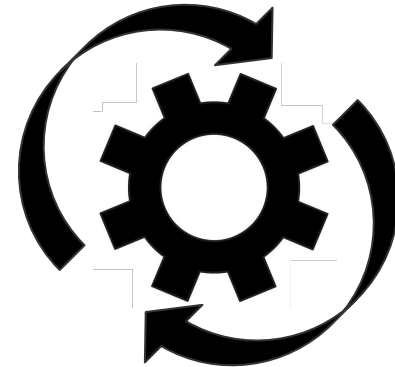- Relation to other resources (data, papers)?
- Metadata?

## Compare



**Given two or more tools:**
- What are their similarities?
- Differences?
- Main features?

## Reuse



**How to quickly:**
- run?
- repeat?
- reproduce?
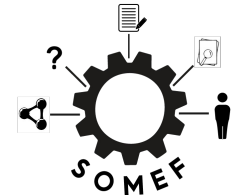- fix?
- combine?

## Manually

- Codemeta generator

https://codemeta.github.io/codemeta-generator/



## Automatically

- Software Metadata Extraction Framework (SoMEF)

Software metadata extraction for Software Management Plans. DeRSE 2024

4

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



Repository

Extraction

Results (Metadata)

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

CodeMeta

TURTLE

+ provenance

JSON

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

- Paragraph-based **text classification**
- Four main categories (binary classification):
  – Installation
  – Description
  – Invocation

| 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 |

| Classifier | Best pipeline | Precision | Recall | F-Measure |
|---|---|---|---|---|
| Description | CountVectorizer + LogisticRegression | 0.85 | 0.79 | 0.82 |
| Installation | TFIDFVectorizer + StochasticGradientDescent | 0.92 | 0.9 | 0.91 |
| Invocation | CountVectorizer + NaiveBayes | 0.88 | 0.9 | 0.89 |
| Citation | CountVectorizer + NaiveBayes | 0.89 | 0.98 | 0.93 |

Simple classification pipelines yield **adequate** results

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

**Installation**

**Installation through Docker**

Wordnet

```
docker pull uscisii2/kgtk
```

To run KGTK in the command line:

```
docker run -it --rm  --user root -e NB_GID=100 -e GEN_CERT=yes -e GRANT_SUDO=yes uscisii2/kgtk:latest
```

Installation instructions

# KGTK: Knowledge Graph Toolkit

DOI 10.5281/zenodo.3828068    build passing    coverage 33%

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:

Using READMEs to **categorize** software

- Creating a methodology to recognize categories based on awesome lists

… and recognize their **nature**:

- Workflow, website, library, tool, ontology…

Readme

**PYLEOCLIM**
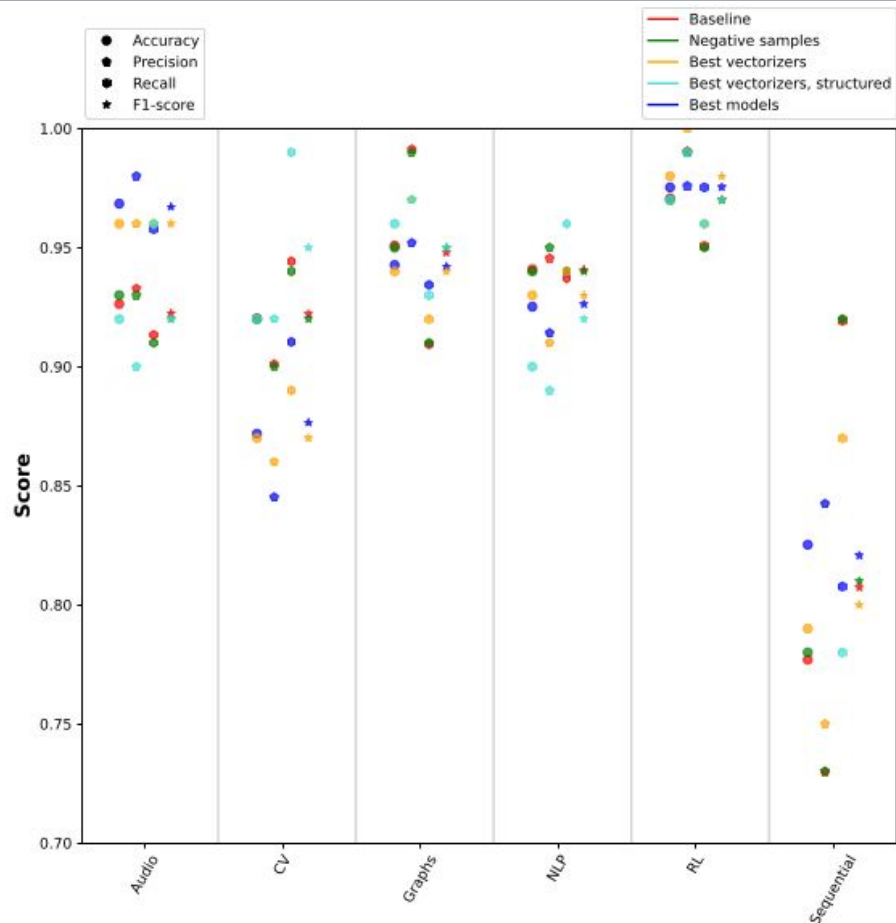
"Data Analysis"
"Package"

Fig. by Jenifer Tabitha

- **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)
- ...

## Method used (provenance):

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

So what can you achieve once you have **rich metadata**?

Search, compare, FAIR assessment…

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

Software metadata extraction for Software Management Plans. DeRSE 2024
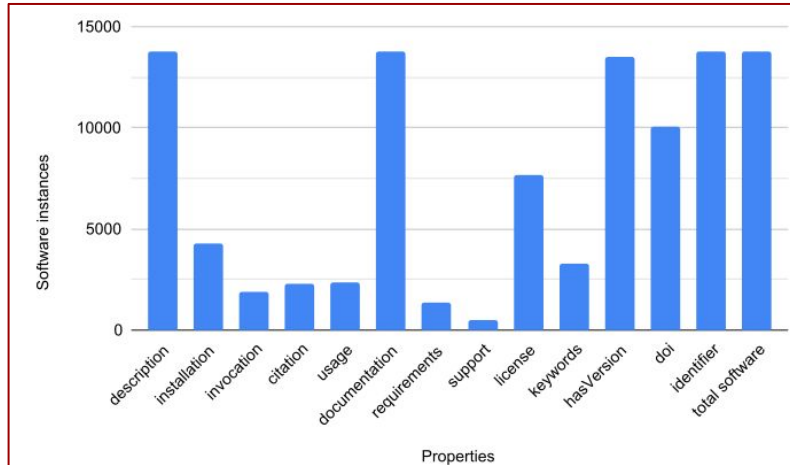
12

# A quick reality check…

We still have a **long way ahead of us**:
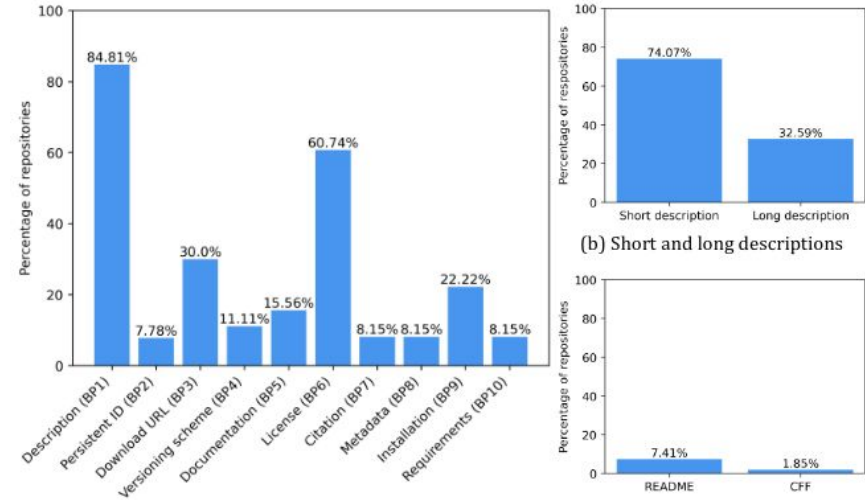- Best practices not widely adopted
- Citation practices are heterogeneous

## Software Engineering (Arxiv) [2]

| Citation practices | description | Bibtex | CFF | title |
|---|---|---|---|---|
| # Bidirectional papers | 759 | 307 | 49 | 353 |

## **Zenodo** [1]



## My **own lab** [3] :(



(b) Short and long descriptions

[1] Kelley, A., & Garijo, D. (2021). A framework for creating knowledge graphs of scientific software metadata. *Quantitative Science Studies*, 1-37.
[2] Garijo, D.; Arroyo, M.; Gonzalez, E.; Treude, C.; and Tarocco, N. Bidirectional Paper-Repository Tracing in Software Engineering. To appear in 21st International Conference on Mining Software Repositories, Cham, 2024. ACM
[3] Iglesias-Molina, A.; and Garijo, D. Towards Assessing FAIR Research Software Best Practices in an Organization Using RDF-star. Semantics 2023 Posters and Demos (CEUR), 3526. 2023.
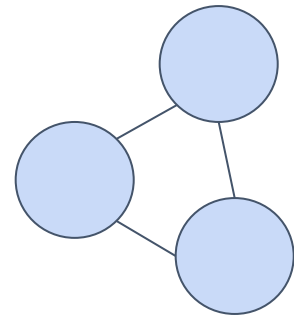
Research software **should become FAIR**!

**Automated** metadata extraction

- Scalable
- Helps fixing issues at the source (maintainable)
- Less effort
- But may contain **errors**

**Next steps**: beyond **Codemeta**

- inputs/ outputs
- models
- workflows
- containers
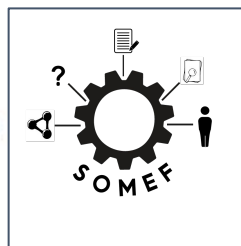- How to make the extracted metadata **more actionable**?

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

Code + documentation

Automated extraction

SOMEF

Knowledge Graphs

findable

portable

comparable

executable

reusable

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