# Don't mention it: An approach to assess challenges to using software mentions for citation and discoverability research

**Stephan Druskat**[1,2], **Neil P. Chue Hong**[3], **Sammie Buzzard**[4], **Olexandr Konovalov**[5], **and Patrick Kornek**[5]

[1]**German Aerospace Center (DLR), Institute for Software Technology, Berlin, Germany**
[2]**Humboldt-Universität zu Berlin, Department of Computer Science, Berlin, Germany**
[3]**EPCC, University of Edinburgh, Edinburgh, United Kingdom**
[4]**School of Earth and Environmental Sciences, Cardiff University, Cardiff, United Kingdom**
[5]**School of Computer Science, University of St Andrews, St Andrews, United Kingdom**

Corresponding author:
Stephan Druskat[1]

Email address: stephan.druskat@dlr.de

## ABSTRACT

Datasets collecting software mentions from scholarly publications can potentially be used for research into the software that has been used in the published research, as well as into the practice of software citation. Recently, new software mention datasets with different characteristics have been published. We present an approach to assess the usability of such datasets for research on research software. Our approach includes sampling and data preparation, manual annotation for quality and mention characteristics, and annotation analysis. We applied it to two software mention datasets for evaluation based on qualitative observation. Doing this, we were able to find challenges to working with the selected datasets to do research. Main issues refer to the structure of the dataset, the quality of the extracted mentions ($54\%$ and $23\%$ of mentions respectively are not to software), and software accessibility. While one dataset does not provide links to mentioned software at all, the other does so in a way that can impede quantitative research endeavors: (1) Links may come from different sources and each point to different software for the same mention. (2) The quality of the automatically retrieved links is generally poor (in our sample, $65.4\%$ link the wrong software). (3) Links exist only for a small subset (in our sample, $20.5\%$) of mentions, which may lead to skewed or disproportionate samples. However, the greatest challenge and underlying issue in working with software mention datasets is the still suboptimal practice of software citation: Software should not be mentioned, it should be cited following the software citation principles.

## INTRODUCTION

Until recently, a key challenge in trying to understand the research software landscape was a lack of knowledge of what software is being used. While there are some discipline-specific software catalogues and lists, and recent initiatives such as the Netherlands eScience Center Research Software Directory (Spaaks et al., 2020) and the Research Software Encyclopedia (The Research Software Encyclopedia project, 2021), there is no single, comprehensive directory or curated list of research software whose use is attested in the literature.

Approaches to build datasets that contain software references, e.g. for specific disciplines, often take publications as a starting point and identify software that was mentioned in the publication, either manually (Du et al., 2021) or through machine learning (). Domain registries such as *swMATH* for mathematical software (Dalitz et al., 2020) and *ASCL* for astrophysics software (Allen & Schmidt, 2015) use these techniques to identify possible candidates for inclusion. Other approaches rely on automatically mining code repositories, looking for key markers such as citation files or DOIs (Eitzen, 2020).

Researchers who want to do quantitative empirical research on research software and research software engineering (see Felderer et al. (2023)) have a need for datasets with specific features that allow them to access research software metadata and artifacts, including source code, i.e., datasets need to include URLs to repositories that store the relevant metadata or artifacts.

Examples for such research include

- accessing source code to study implementation details in research software;
- accessing source code repositories to study software metadata such as license information;
- accessing features of source code repositories such as the software's version history, issue trackers or integration processes (pull/merge requests) to study development processes and software provenance (Kurnatowski et al., 2021).

Recently, five datasets have been released that can potentially be used for the type of research described above. The *SoftCite* project (Du et al., 2021) is a human curated list of $4,093$ software mentions in the life and social sciences. This dataset, in turn, was used by a team at the Chan Zuckerberg Initiative to train a machine learning model that has been used to identify software references in the *CORD-19* collection of COVID-19-related research papers (Wang et al., 2020), which has been published as a raw dataset: *CORD-19 Software Mentions* (Wade & Williams, 2021) (*CSM*). The *CZ Software Mentions* dataset (Istrate, Veytsman, et al., 2022) (*CZI*) provides software mentions from the biomedical literature, their sources, textual context and metadata, extracted by a trained SciBERT model (). A subset of the complete corpus also disambiguates software entities, and provides links to software source code and/or artifact repositories. *SoftwareKG-PMC* (Schindler et al., 2021a) provides a knowledge graph of software mentions generated through a SciBert (Beltagy et al., 2019) model trained on a gold standard corpus of software mentions (). Escamilla et al. (2022) extracted URLs for git hosting platforms from the ArXiv and PMC corpora (*Extract-URLs*, Escamilla (2023)).

In this paper, we present our approach to assessing the usability of software mention datasets for (1) quantitative research on research software that requires access to software metadata or artifacts (i.e. for mining software repositories), and (2) research into the practice of software citation. We understand "usability" as the potential for a dataset to be used without further processing to access the source code of the mentioned software that is recorded in the dataset. An example for high usability would be the inclusion in a dataset of a correct URL to the source code repository for each software in the dataset. Our approach was originally developed to use CSM (Wade & Williams, 2021) in order to answer research questions (see *Research questions*) from the above types of research: (1) What is the impact of licenses on the type of software mention? (2) Has the practice of software citation improved as compared to Howison and Bullard (2015)? We consecutively applied our approach to assess the usability of CZI (Istrate, Veytsman, et al., 2022), which gave us opportunity to evaluate our approach, and compare the usability of two datasets that have been created in a similar way (see *Software mention datasets*). The approach consists mainly of the manual annotation of a stratified sample of the dataset with quality categories, categories of software mentions, and accessibility categories, and the subsequent analysis of the annotations.

In the following sections, we present related work (*Related work*) and describe our assessment approach and the methodology to apply this approach to the CSM and CZI datasets (*Methods*). We then present the results of the assessment and the results of two exploratory studies based on the assessed datasets (*Results*). Finally, we conclude with suggestions for future work (*Conclusion*).

## Related work

In 2015, Howison and Bullard described challenges with identifying and finding software that has been used for biology research, and crediting software authors (Howison & Bullard, 2015). These challenges relate to problematic practices of software citation: Howison and Bullard found that of the publications that mentioned software in any way, only 39% cited a formal publication relating to the software in the references, while informal mentions are included in 43% of investigated publications. Such informal mentions often fail to provide credit to software authors. Furthermore, only in 28% of cases were they able to identify the software versions that had been used for the research presented in the respective publications, and only 5% of the respective versions were referenced in a way that made it possible to find the actual software versions.

Since the publication of Howison and Bullard's paper, work has been done in the scholarly communications and research software communities to address the state of software citation practice. A. M. Smith et al. (2016) define the principles of software citation. These principles (italicised, see A. M. Smith et al. (2016) for their definitions) address the challenges to software identification, findability, and creditability mentioned by Howison and Bullard (2015), when applied in practice:

- *Importance* and *unique identification* disallow informal mentions that may obscure the identity of the software.
- *Persistence*, *accessibility* and *specificity* enable findability of and access to the specific version of referenced software.
- *Credit and attribution* enable credit for software authors.

The FORCE11 Software Citation Implementation Working Group (2017–2023)[1] worked with different stakeholders to endorse the software citation principles, helped implementing them, and developed guidelines for different stakeholder groups, e.g., for software developers (N. P. Chue Hong, Allen, Gonzalez-Beltran, et al., 2019), journal authors (N. P. Chue Hong, Allen,

---

[1] https://www.force11.org/group/software-citation-implementation-working-group. SD and OK were members of the working group that was co-chaired by NCH.

Gonzalez-Beltran, et al., 2019), publishers (Katz et al., 2021), software registries (Registries Task Force on Best Practices for Software et al., 2020), and libraries (Schmidt et al., 2021).

Software citation practices based on the software citation principles are also implemented in, or supported by: software metadata formats, that enable software authors to provide correct and complete citation metadata (e.g. Citation File Format (Druskat, Spaaks, et al., 2021), CodeMeta (Jones et al., 2017)); open access repositories and archives that provide DOIs or other unique identifiers for software (e.g. Zenodo (European Organization For Nuclear Research & OpenAIRE, 2013), Figshare, Software Heritage (Abramatic et al., 2018-10-01, 2018)); source code platforms that display citation information for software (versions) (e.g. GitHub (A. Smith, 2021)); reference managers that support import of correct and complete software citation metadata (e.g. Zotero, JabRef); publications that support and recommend citation of software versions (Katz et al., 2021).

Despite these activities towards better software citation, and an increase in the number of DOI registrations for software (Fenner et al., 2018), good software citation practice, it seems, has not yet permeated research culture: only "3.24% of all software DOIs registered by Zenodo are traceably cited at least once" (van de Sandt et al., 2019, p. 2). Instead, informal mentions of software in publications are still commonplace, which makes it hard to track software usage in research (see Du et al., 2021). This in turn is not only disadvantageous for the sustainability of the software in question (see Druskat, Katz, & Todorov, 2021, for a brief discussion of the relationship between citation and research software sustainability), it also impedes quantitative research on different aspects of research software: instead of being able to work with citation graphs that include software directly (), researchers are forced to trace software mentions. This is often a time-consuming and costly process: to build the *SoftCite* dataset, 36 student assistants were paid for 2 years to annotate software mentions in full-text PDF files (Du et al., 2021, p. 872). An alternative to the manual extraction of software mentions is the application of machine learning to datasets of literature. In this paper, we report on two datasets that were created using machine learning models.

## METHODS

The availability of new, large software mention datasets (see *Introduction*) promises improved access to research software whose use has been reported in the literature. If these datasets include direct links to software metadata and/or artifacts for research software, they can be used to build corpora of research software source code, or extract samples for empirical research. Additionally, these datasets may potentially be used to analyse the practice of software citation. We iteratively developed an approach to assess the usability of software mention datasets for these purposes. We then applied the approach to two of the datasets mentioned above. Our approach included three steps:

1. Take a stratified proportionate random sample from a dataset and prepare the sample for annotation.

2. Manually annotate the sample for mention extraction quality, mention categories and quality, following a set of annotation guidelines.

3. Analyse the sampling and preparation workflows, as well as the annotations, to assess the usability of the dataset, and report preliminary results for our research questions where possible.

Code, samples and annotated data are available as Druskat and Chue Hong (2023).

### Research questions
We developed and applied our approach to answer the following research questions.

**RQ1: Are software mention datasets usable as data sources for research on research software?** More precisely: **RQ1.1: Are software mention datasets usable as data sources for research on research software that requires access to software metadata or artifacts?** and **RQ1.2: Are software mention datasets usable as data sources for research into the practice of software citation?**

The goal of these research questions is to identify features of datasets that allow their use for research on research software, and vice versa, identify features that make datasets' use for this purpose harder or impossible. Intuitively, the inclusion in a dataset of links to software artifacts or metadata would enable research. Results would provide evidence of additional features whose presence, absence, or characteristics can enable or preclude reuse for research purposes.

**RQ2: Is open source software more cited in a way that allows credit for software authors than closed source software?**

The goal of this research question is to understand if the way that the software is licensed has an impact on the way that the software is cited or mentioned in publications. Howison and Bullard, 2015 define seven types of software mentions in publications, categorised as formal citations (to publications, user manuals or project websites), explicit mentions (like an instrument, of a URL, or just the software name) or implicit mentions (not even a name). We hypothesise that commercial software is cited more frequently using an in-text name mention or citation to project name or website, but that open source software is cited more frequently with a repository or associated research publication in the reference. The latter makes it easier

to credit the authors. Results would provide evidence concerning whether the increasing prevalence of Open Science / Open Research approaches could improve the quality of software citation.

**RQ3: Has the practice of software citation represented in software mention datasets improved in comparison to the practice described in Howison and Bullard (2015)?**

The goal of this research question is to find out if the practice of software citation has improved since the publication of Howison and Bullard (2015). Since 2015, a number of contributions have been made towards better software citation. In 2016 the software citation principles were published (A. M. Smith et al., 2016). They describe how software should be cited, and what metadata the respective references should ideally include: The software itself should be cited, rather than a substitute output, e.g., a paper describing software. The reference should also name the authors to enable academic credit, include a persistent identifier to the persisted artifact of the exact version of the software that was used, and allow access to the software itself. Following the principles paper, which was the main output of the FORCE11 Software Citation Working Group[2], the follow-up Software Citation Implementation WG[3] (2017–2023) created a number of outputs addressing different stakeholders to improve software citation practice: software developers (N. P. Chue Hong, Allen, Gonzalez-Beltran, et al., 2019), authors (N. P. Chue Hong, Allen, Gonzalez-Beltran, et al., 2019), publishers Jay et al. (2021) and libraries (N. Chue Hong et al., 2021). Additionally, the FAIR Principles for Research Software have drawn attention to the importance of software metadata that also "enables and encourages the citation of software" (N. P. Chue Hong et al., 2021, p. 8). And metadata formats were developed to make it easier to provide correct and complete software citation metadata in the first place ().

While we cannot prove causation between any of the above-mentioned activities and outputs and the data in software mention datasets, we hypothesise that software citation practice will overall have improved since the publication of Howison and Bullard (2015). Specifically, we expect that those of Howison and Bullard's (2015) categories of mentions that reflect the principles better are found relatively more often in mentions from publications in the datasets that were published in or after 2016.

### Software mention datasets

We applied our approach to two datasets:

- *CORD-19 Software Mentions* (Wade & Williams, 2021) (*CSM*)

- *CZ Software Mentions: A large dataset of software mentions in the biomedical literature* (Istrate, Veytsman, et al., 2022) (*CZI*)

CSM contains lists of software mentioned in 77,000 papers, and exceeds the size of the dataset used by Howison and Bullard (2015, random sample of 90 papers) by a factor of over 850. It was created by detecting software mentions in the CORD-19 (Wang et al., 2020) collection of full-text research papers related to the SARS-CoV-2 virus, using a machine learning model originally trained and evaluated on the *SoftCite* dataset (Du et al., 2021).

CZI contains software mentioned in $\approx$ 20.7 million papers (Istrate, Li, et al., 2022), thereby exceeding the size of CSM by a factor of $\approx$ 26,000. It was created by detecting software mentions in the NIH PMC-OA Commercial and Non-Commercial subsets, and a custom collection of papers provided by various publishers to the Chan Zuckerberg Initiative. The detection used a machine leraning model based on SciBERT(Beltagy et al., 2019), also trained on the *SoftCite* dataset (Du et al., 2021).

The datasets were chosen for different reasons. Firstly, both datasets have likely been prepared using the same or different versions of the same machine learning model[4], share authors (I. Williams), and have been prepared at the same institution (Chan Zuckerberg Initiative, Redwood City, CA, USA). This suggested that the methods applied in dataset creation would not have to be factored into a comparison of the two datasets.

Secondly, we were involved in preliminary work on CSM, which led to an earlier version of this paper. This preliminary work was conducted during a hack event at the *Software Sustainability Institute*'s *Collaborations Workshop 2021 Hack Day*[5] (Konovalov et al., 2021). During the hack event, a random sample of 100 mentions from CSM was taken and manually cleaned, and annotated with categories pertaining to the quality of the mention with respect to accessibility of the software and the quality of the mention extraction. Additionally, for each mention in the sample, identification of a source code repository was attempted, and the respective URL added to the dataset. The resulting dataset (`CORD19_software_popularity_sampled_QA_DOI.csv`) is available as part of Druskat and Chue Hong (2023).

---

[2]https://force11.org/group/software-citation-working-group/
[3]https://force11.org/groups/software-citation-implementation-working-group/
[4]See the code at Veytsman (2022) which uses the model that was used to create CSM (see The Software Mention Extraction authors (2022)).
[5]https://software.ac.uk/cw21/hack-day

## Sampling

Usually, sampling is a means to the end of obtaining data for a study, e.g. in sample studies along the lines of Stol and Fitzgerald (2018, p. 16f.). In our case, the sampling process was part of the study to answer RQ1, if software mention datasets are usable as data sources for research on research software. Here, we describe our methods for sampling CSM and CZI. We discuss results towards answering RQ1 (RQ1.1, RQ1.2) in *Results*.

To conduct our study towards answering RQ2 and RQ3, we took samples of both CSM and CZI. Sampling was done in Jupyter Notebooks () using NumPy (), Pandas (The pandas development team, 2023), Dask (The Dask 2023.3.1 developers, 2023), Matplotlib (), SciPy () and scikit-learn (Grisel et al., 2023).

To obtain a sample from CSM for annotation, we first downloaded the dataset (Wade & Williams, 2021). It consists of a CSV file which collects publications in rows. For each publication, there is some bibliometric information, such as DOI, title, source dataset, license, publication data, journal name and a list of URLs that resolve to a website providing a copy of the publication. Additionally, the software mentions that were extracted from the publication are given as a comma-separated list, e.g., `['GraphPad Prism', 'SigmaPlot', 'Systat']`. As we wanted to sample software mentions, not publications, we needed to explode the lists of software mentions in the dataset while preserving the bibliometric information, and clean up the exploded mention strings. The resulting data included $558,792$ software mentions. To better understand the data, we took counts of software mentions and plotted their distribution (Figure 1). The distribution shows extreme positive skew, around half the distinct mentinoned software has 1 mention, around 8% has more than 10 mentions, less than 1.5% have more than 50 mentions. We then took a simple random sample of $1,000$ rows from the dataset, exploded it to get one distinct software mention per row, took counts of mentions in the sample and plotted their distribution (Figure 2). The distribution of the sample also shows a strong positive skew. As was to be expected, Levene's test showed that the variances for number of mentions between the full dataset and our sample were not equal: $F(3296, 35) = 7.73, p = .0054$. As variance would not influence our study, we used this sample to take another random sample of 100 software mentions for developing and evaluating the annotation tagset.
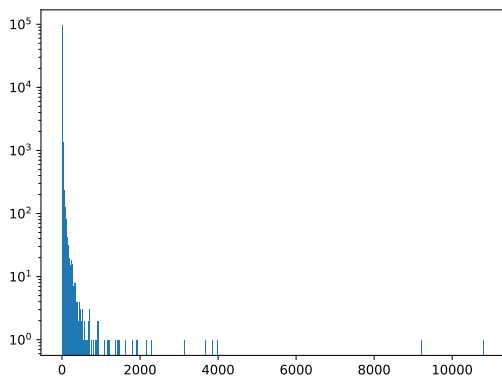


**Figure 1.** Distribution of mention counts over the complete exploded CSM dataset. *x*: distinct software mentions (*log*), *y*: sum of mentions for distinct software.
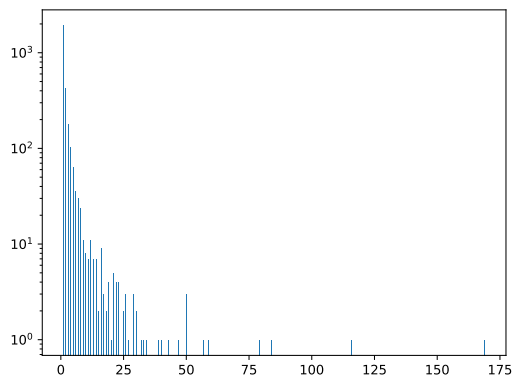


**Figure 2.** Distribution of mention counts over our sample from the CSM dataset. *x*: distinct software mentions (*log*), *y*: sum of mentions for distinct software.

CZI consists of several subsets (Istrate, Li, et al., 2022), representing the different steps in the dataset creation process: raw data, linked data, disambiguated data. As we were interested in evaluating the quality of original software mentions as found in the literature (RQ3), we used the *raw* and *linked* subsets. To obtain a sample from CZI, we used two consecutive Jupyter notebooks. In the first, we did the sampling. In the second, we merged the sampled mentions with the *linked* susbset of CZI to obtain repository data for the mentions where available.

As CZI is considerably larger than CSM, we used Dask (The Dask 2023.3.1 developers, 2023) instead of Pandas (The pandas development team, 2023) for working with the dataset, and ran the notebooks on a small CPU cluster at the German Aerospace Center's Institute for Software Technology using nbconvert (The nbconvert 7.2.10 developers, 2023). The sampling notebook downloaded and extracted the dataset archives. We then merged the raw subsets for each of the collections contained in the CZI dataset (commercial, non-commercial, publishers collection), and filtered those mentions that were curated as being to software (see Istrate, Li, et al. (2022)). The resulting dataset contained $20,792,352$ software mentions of $6,966$ distinct software. As for CSM, we plotted the distribution of software mention counts in the dataset (Figure 3) As the dataset at this point still used $\approx 1GB$ of disk space, we took a stratified proportionate sample of $\approx 100,000$ rows from the dataset after

brute-force deduplication of software names through capitalization, to avoid long computation times. The sample contained 99,973 software mentions of 6,966 distinct software. To visually check the stratification, we plotted the distribution of software mention counts in the sample (Figure 4). We then extracted one random row per distinct software name to avoid having duplicate instances of distinct software in our annotation sample, and sampled 100 rows randomly for annotation.
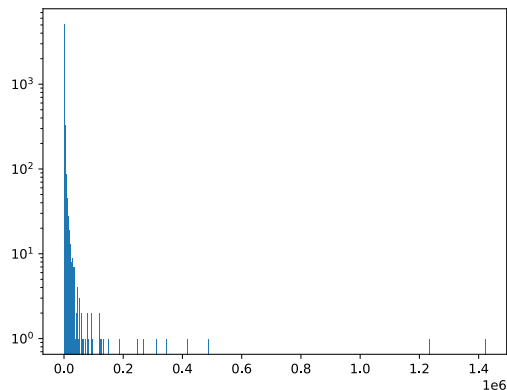


**Figure 3.** Distribution of mention counts over the complete filtered CZI dataset. *x*: distinct software mentions (*log*), *y*: sum of mentions for distinct software.
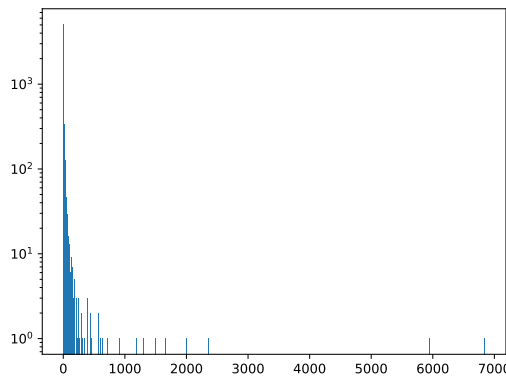


**Figure 4.** Distribution of mention counts over our 100k sample from the CZI dataset. *x*: distinct software mentions (*log*), *y*: sum of mentions for distinct software.

CZI contains a subset with URLs resulting from exact searches for software names in different software source code and artifact repositories (see Istrate, Li, et al. (2022, p. 8)). We wanted to leverage these data to add them to our assessment of the dataset quality (RQ1.1), and reuse the data to answer RQ2. To achieve this, we merged the relevant information from the subset of CZI containing repository data into our sample. Identification of the relevant information was possible through the unique mention IDs included in the raw and linked CZI subsets.

## Annotation

We manually annotated the samples from CSM and CZI on different layers, regarding the different research questions: (a) *quality of the extracted mention*, e.g., whether the extracted string included the complete software mention string (RQ1); (b) *quality of the mention in the publication*, i.e., whether and how the mention allows access to the software (RQ1); (c) for CSM, a URL where the software could be found (RQ1); (d) for CZI, *quality of the repository links*, i.e., whether links were extracted, and whether they referred to the correct software (RQ1); (e) *license and license type* for the mentioned software (RQ2); (f) *mention type* based on the mention types put forward in Howison and Bullard (2015, Table 6) (RQ3). Additionally, we added secondary annotations: whether the publication the mention was extracted from is a preprint; whether the publication the mention was extracted from is a paper describing the mentioned software; the confidence of the annotator regarding the correctness of the annotations.

Annotations were guided by the annotation guidelines summarized in the next section. The guidelines were developed iteratively by 1. annotating random samples, 2. analysing confidence annotations, 3. improving the annotation guidelines through discussion, 4. repeating from 1. until the guidelines were not further improved. Once no further improvements were made to the annotation guidelines, SD, NCH, SB and OK each annotated the same set of 50 random mentions from the CSM sample. We used these annotations to calculate inter-annotator agreement (Table 1). Finally, the annotation guidelines were applied to annotate software mentions in a second sample from CSM ($n = 100$) and a sample from CZI ($n = 100$), in addition to the already assessed annotations in the first CSM sample. The complete workflow is shown in Figure 5.
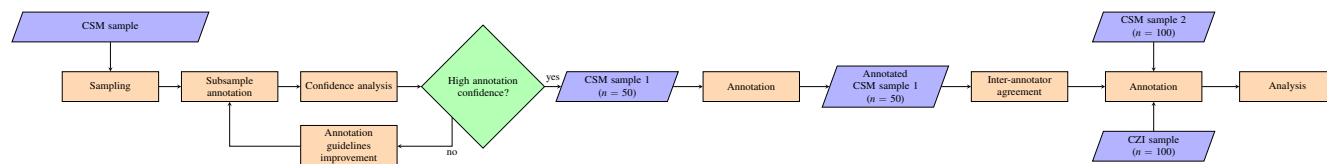


**Figure 5.** Visualization of the complete assessment workflow.

| Annotation layer | Krippendorff's $\alpha$ |
|---|---|
| Mention type | 0.55 |
| Quality of the mention in the publication | 0.72 |
| Quality of the mention extraction | 0.65 |
| Preprint | 0.80 |
| Software paper | 0.49 |
| All layers | 0.64 |

**Table 1.** Inter-annotator agreement for quality and mention annotations on a random sample of 50 mentions from the CSM sample.

### Annotation guidelines

For each software mention in the sample,

1. Resolve the first identifier for the publication in a web browser.

    1.1. If the publication is a preprint, use the next identifier if available.

    1.2. If the only available identifier is for a preprint, use the preprint.

2. Open the PDF for the publication.

    2.1. If you cannot access the PDF due to a paywall, use the next identifier.

    2.2. If there is no next identifier, use Unpaywall[6] to access an open version of the publication, or ask a co-author to retrieve the publication.

3. Search for the exact mention string in the PDF.

4. Verify for each search result that it is the exact search string. Note that:

    4.1. The mention string may be a substring of the complete software name (due to line breaks, composite names, etc.).

    4.2. There may be multiple software packages mentioned with similar names.

5. Annotate the quality of the mention retrieval according to Table 2.

6. Identify the best mention and annotate the mention type.

    6.1. Identify the best mention by adherence to the software citation principles.

    6.2. The *Order* column in Table 3 encodes the quality of the mention (from 1 = best to 6 = worst) by principles:

        • Importance is always the best. Citation of project name or website is better than citation of a publication. (Importance, Accessibility)

        • Citation of a publication is better than citation of a user manual. (Credit)

        • URLs in text are second best. (Accessibility)

        • Instrument-like citation is better than name-only mention. (Accessibility)

        • Name-only mentions are better than mention without name.

    6.3. Only use mentions matching the exact mention string, including capitalization.

    6.4. Only URLs found in the same paragraph as the mention, or in a footnote that is called from the same paragraph, shall be annotated with URL.

---

[6]https://unpaywall.org/

6.5. Citations to references must appear within the boundaries of the sentence that includes the mention.

6.5.1. Examples for citations to process:
- "We used SOFTWARE [1] for the analysis."
- "We used SOFTWARE for the analysis [1]."
- "We used SOFTWARE for the analysis. [1]"

6.5.2. Example for citations to ignore:
- "We used SOFTWARE and Otherthing for the analysis. We refuted the null hypothesis. The data provided evidence for something [1, 2]."

7. Annotate the quality of the mention (Table 4).

7.1. Differentiate between mention types `NA` and `SN`.

7.1.1. If it is clear that the authors considered the mentioned entity software, annotate as `SN`. Examples: listed as "computational method", compared with other software.

7.1.2. If still unclear, discuss with other annotators.

7.1.3. If still unclear, annotate as `UN`.

8. Annotate other layers.

| Code | Name |
|------|------|
| Y | Yes, name was correctly and completely retrieved from the publication for the dataset. |
| N | No, name was NOT correctly and completely retrieved from the publication for the dataset. |

**Table 2.** Annotations for quality of the mention extraction/retrieval.

| Code | Name | Definition | Order |
|------|------|-----------|-------|
| PUB | Cite to publication | Cites a paper/monograph primarily describing the mentioned software (NOT a review paper comparing different software), as it would for non-software cites. For non-software mentions, we don't judge the suitability of the referenced work. | 2 |
| PRO | Cite to project name or website | Cites the project name or website via a "fake" reference. | 1 |
| URL | URL in text | URL in text or in footnote | 4 |
| MAN | Cite to user manual | | 3 |
| INS | Instrument-like | Mention software in a manner similar to scientific instruments or materials, typically mentioning the name in text followed by the author or company and a location in parentheses. | 5 |
| NAM | In-text name mention only | | 6 |
| NOT | Not even name mentioned | | 7 |

**Table 3.** Annotations for mention types following Howison and Bullard (2015).

***License annotation***

Licenses for each software mention were annotated by NCH for the sample of CSM prepared during the initial hack event (see *Software mention datasets*), and by SD for the CZI sample. This was done by examining any associated code repository, website or documentation related to the mentioned software. For CSM, where a link to a repository or project website was identified as part of the preliminary work of the hack event, this was used and checked to see if a license was documented.

| Code | Name |
|------|------|
| SC | Software where a direct link to a code repository or distribution repository landing page (e.g., CRAN, PyPI) can be found in the mentioning paper, and the page includes author/version/license metadata. |
| SP | Software where a link to another website can be found in the mentioning paper and that website provides access to the source code, but the website does not provide author/version/license metadata. |
| SN | Software but no link to a code repository or website providing access to the source code can be found in the mentioning paper. Annotate as SN even if the reference is to a software paper that does include a link to a source code repository. |
| NA | Not software (only annotate this, retrieval quality and confidence) |
| UN | Other classification - unknown/needs further investigation, e.g., unclear from the information in the paper whether this is software or not. |

**Table 4.** Annotations for quality of the mention itself.

If a link was not present in the initial dataset, an additional attempt to find a source of documentation for the software was undertaken by NCH, and the license recorded if available. For CZI, any repositories for the software mention linked to from the *linked* subset of CZI were checked to see if they were for the mentioned software. If they were not, an attempt to find a source of documentation for the software was undertaken by SD and added to the dataset. Additionally, the quality of the repository URL extraction in CZI was annotated, and the license recorded. In a number of cases where the mention was to Software-as-a-Service and the license could not be identified, it was categorised in a subset of the "unknown" category.

### Analysis

Results pertaining to research question *RQ1: Are software mention datasets usable as data sources for research on research software?* and its subquestions were gained through qualitative observation during the sampling and annotation work described above.

In order to yield exploratory results for research question *RQ2: Is open source software more cited in a way that allows credit for software authors than closed source software?*, we grouped the licenses into the categories *closed* (closed source licenses), *academic* (academic use only, non-commercial licenses), *permissive* (minimally restrictive open source licenses), *copyleft* (open source licenses with reciprocal clauses) and *unknown* (license conditions could not be found, including Software-as-a-Service). We then clustered licenses into *open* (*permissive + copyleft*) and *closed* licenses (*closed + academic + unknown*). We also clustered mention types into quality categories *good* (PUB), *okay* (PRO + URL) and *poor* (INS + NAM).

To gain exploratory insights into the data with regard to *RQ3: Has the practice of software citation represented in software mention datasets improved in comparison to the practice described in Howison and Bullard (2015)?*, we grouped the mention type annotations by publication year of the mentioning publication, and compared the distribution over mention types with the results from Howison and Bullard (2015).

Analyses of the annotations were done in Jupyter Notebooks () using NumPy (), Pandas (The pandas development team, 2023) and Matplotlib (). Both our annotated data and the Jupyter notebooks with the analyses are available as part of Druskat and Chue Hong (2023).

## RESULTS

We developed an approach to assess the usability of software mention datasets for research on research software (RQ1). The approach includes taking a sample from the software mention dataset and preparing it for annotation, then annotating it manually for mention extraction quality and mention categories and quality, following a set of annotation guidelines. Finally, the annotations are analyzed to answer research questions.

We applied the approach to small samples (total $n = 250$) from two software mention datasets (CORD-19 Software Mentions (CSM), Istrate, Veytsman, et al. (2022); CZ Software Mentions (CZI), Wade and Williams (2021)), and assessed it through qualitative observation. We were particularly interested to see if software mention datasets can be used for quantitative research that requires access to software metadata or artifacts (RQ1.1). We also wanted to find out if they can be used for research into the practice of software citation (RQ1.2).

Through the application of our approach to assess the usability of software mention datasets for research on research software, we were able to uncover challenges to working with software mention datasets for the above-mentioned types of research.

We also found that practice of software citation must significantly improve in general to adhere to the software citation principles (A. M. Smith et al., 2016), thereby improving the quality of the source data used to create software mentions datasets.

## Exploratory studies

While it was our main goal to evaluate our approach to assessing the usability of software mention datasets for research on research software, some highly preliminary results were gained during the evaluation process.

Using our samples from CSM and CZI, we analyzed the software licenses for mentioned software. Figure 6 shows the distribution of license categories for closed and open licenses.
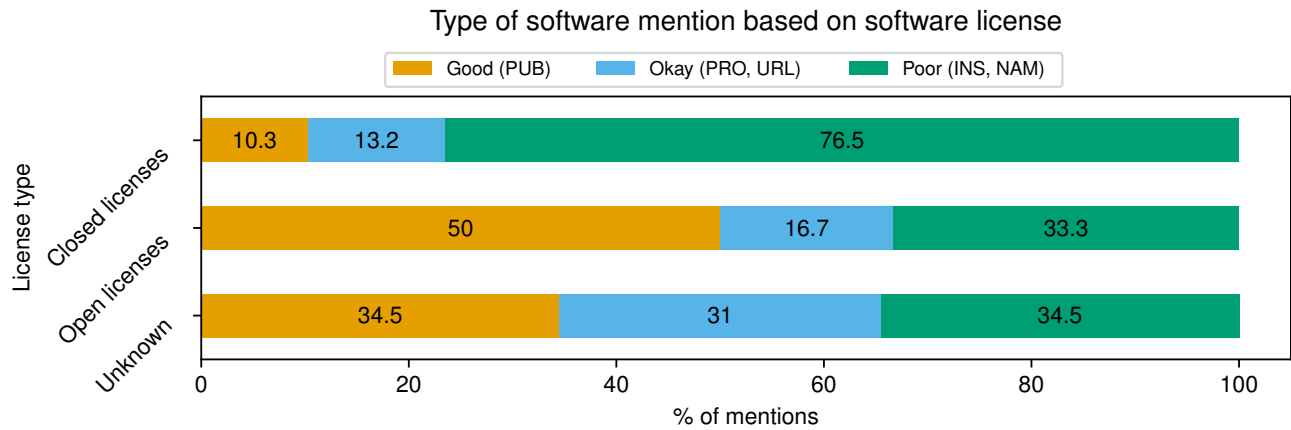


**Figure 6.** Percentages of mention types found in two samples from CSM and CZI, based on their software license, categorised by "quality".

While software available under a closed license (license types "closed" and "academic", see Table 5) is generally referenced using a poor quality mention (76%), half of the mentioned openly licensed software ("permissive", "copyleft") is referenced using a good quality mention, and another 16.7% using at least a medium quality mention. Table 6 shows the detailed distribution of mention types over license types.

| License category | License code | Description |
|---|---|---|
| Closed | Closed | Closed source licenses, generally commercial products |
| Academic | Academic | No cost for academic or non-commercial use |
| Permissive | Apache, Artistic, BSD, MIT, Unlimited | Minimally restrictive open source licenses |
| Copyleft | LGPL, GPL | Open source licenses with reciprocal clauses |
| Unknown | Unknown, Unknown (SaaS) | License conditions could not be found, with a subset for Software as a Service (SaaS) with no license for service or code |

**Table 5.** Categorisation of software licenses identified in our dataset.

Using our samples from CSM and CZI, we analyzed the quality of software mentions published after 2015, and compared them with the mention quality reported in (Howison & Bullard, 2015). Mentions in publications published in or after 2016 made up 75% of the CSM sample and 63% of the CZI sample. Table 7 shows the distribution of mention types over the samples from CSM, CZI and Howison and Bullard's data.

Figure 7 shows the same data. While the CSM data show an increase by 12.1% for URL type mentions as compared to Howison and Bullard's data, the other high and medium quality mentions have decreased (PUB −6.9%) or only slightly increased (PRO +0.8%). For our CZI sample, the quality of mentions seems to have decreased overall: PUB −7%, PRO −0.5%. There is only a slight increase in URL mentions (+1.7%), but a significant increase in name-only mentions (NAM: +18.9%), which represent more than half of the mentions in the sample published after 2015.

| | PUB | | PRO | | INS | | URL | | NAM | | Total | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Σ | % | Σ | % | Σ | % | Σ | % | Σ | % | Σ | % |
| License type | | | | | | | | | | | | |
| Closed | 3 | 1.84 | 1 | 0.61 | 23 | 14.11 | 6 | 3.68 | 24 | 14.72 | 57.0 | 34.97 |
| Academic | 4 | 2.45 | 0 | 0.00 | 1 | 0.61 | 2 | 1.23 | 4 | 2.45 | 11.0 | 6.75 |
| Permissive | 16 | 9.82 | 2 | 1.23 | 0 | 0.00 | 4 | 2.45 | 12 | 7.36 | 34.0 | 20.86 |
| Copyleft | 17 | 10.43 | 3 | 1.84 | 1 | 0.61 | 2 | 1.23 | 9 | 5.52 | 32.0 | 19.63 |
| Unknown | 10 | 6.13 | 4 | 2.45 | 0 | 0.00 | 5 | 3.07 | 10 | 6.13 | 29.0 | 17.79 |

**Table 6.** Distribution of mention types over license categories.

| | CZI | | CSM | | Howison & Bullard | |
|---|---|---|---|---|---|---|
| Mention type | Σ | % | Σ | % | Σ | % |
| PUB | 19 | 30.2 | 20 | 30.3 | 105 | 37.2 |
| MAN | 1 | 1.6 | 0 | 0.0 | 6 | 2.1 |
| PRO | 3 | 4.8 | 4 | 6.1 | 15 | 5.3 |
| INS | 4 | 6.3 | 11 | 16.7 | 53 | 18.8 |
| URL | 4 | 6.3 | 11 | 16.7 | 13 | 4.6 |
| NAM | 32 | 50.8 | 20 | 30.3 | 90 | 31.9 |

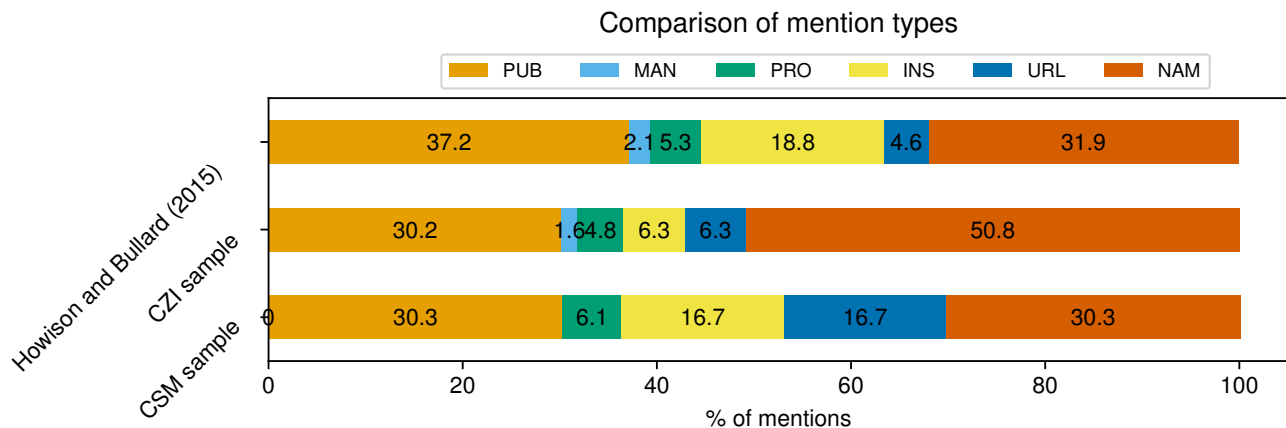**Table 7.** Distribution of mention types over samples.



**Figure 7.** Percentages of mention types found in the CSM and CZI samples compared to Howison and Bullard (2015). See Table 3 for the definition of mention types.

When we clustered the mention types into three more coarsely grained categories, we could observe the same trends. Table 8 shows the detailed distribution, Figure 8 provides an overview.

## DISCUSSION

The application of our approach to two software mentions datasets allowed us to define features that software mention datasets should include to enable research into research software and software citation practices (see *Results*).

During sampling, we found that the structure of the dataset has an impact on ease of sampling. CSM is based on publications: each row contains *all* mentions from a specific publication in a Python-like string list. This makes it more cumbersome and computationally expensive to extract a sample based on individual software mentions, as all rows must first be exploded, and data must be cleaned afterwards to exclude artifacts from list data. This structure also makes it harder to explore the dataset

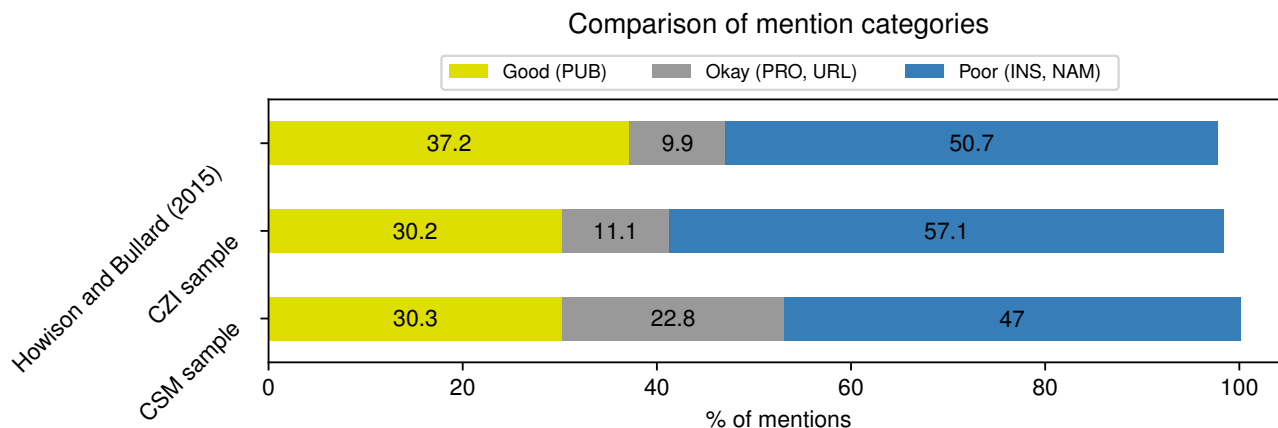|                             | Good (PUB) | Okay (PRO, URL) | Poor (INS, NAM) |
|-----------------------------|------------|-----------------|-----------------|
| CSM sample                  | 30.3       | 22.8            | 47.0            |
| CZI sample                  | 30.2       | 11.1            | 57.1            |
| Howison and Bullard (2015)  | 37.2       | 9.9             | 50.7            |

**Table 8.** Distribution of mention categories over samples.



**Figure 8.** Percentages of mention categories found in the CSM and CZI samples compared to Howison and Bullard (2015).

initially with regard to individual software mentions, as they are potentially distributed across many rows, which makes filtering difficult. This is much improved in CZI, which is based on individual mentions, starting with the initial raw data. Linking the different subsets of CZI is straightforward, as mentions have a unique ID. Therefore we reason that ordering data in software mention datasets by uniquely identified mention, not publication, is an important feature to make the dataset easily usable.

As both datasets are tabular data, manual annotation is technically straightforward, as respective columns could be added. Additionally, software for working with tabular data is readily available. We reason that being persisted as tabular data is a feature that improves usability of datasets for research.

While the annotation process for mention types and quality, and the quality of mention extraction, towards answering RQ3 was generally straightforward, the quality of the mention extraction has an impact on usability: in our sample of CSM, 19.3% of mentions were incorrectly extracted, i.e., the software name was not completely and correctly retrieved from the publication. Our sample data from CZI suggests that this may have improved here, with 7% of the mentions incorrectly extracted. Generally, an analysis of mention types could be made much more feasible if the dataset included all contexts for all mentions of a software across a publication. CZI includes context in its `text` column for raw data, which already makes it easier to find the respective mention. However, a single context does not necessarily represent the "best" mention, nor does it allow for an analysis of how a software is mentioned across a publication, where formal citations or references may have been made in another context. We therefore reason that ideally, software mention datasets should include contexts for all mentions of a software across a publication.

Our samples included mentions that were not to software, which may give a preliminary indication of the precision of the machine learning models that were used to extract the mentions. In our CSM sample, extended with the sample used for calculating inter-annotator agreement, 69 out of 150 mentions were not to software. This means that only little more than half (54%) of mentions were correctly identified as software. In our CZI sample, 23 out of 100 mentions were not to software, i.e. 77% were correctly identified. Assessing the precision of the models used for software mention identification in more detail should be part of future research. It seems necessary to improve these models to achieve higher precision.

The annotation of license information for software mentions based on the two dataset samples – as an example for quantitative research on research software – was difficult. CSM does not contain any links between the mentioned software and its documentation or software repository. This made it necessary to find a source of license documentation manually for each mentioned software. Such an approach does not scale and renders quantitative research infeasible. CZI, on the other hand, provide linking from different sources: GitHub, PyPI, CRAN, SciCrunch and Bioconductor (see Istrate, Li, et al. (2022)).

However, three aspects impede quantitative research based on these data.

1. CZI include links from different sources. In some cases, there were more than one link per software. This makes it difficult to leverage the links when they point to different targets. This is less problematic for cases where the targets represent the same software, although this information is hard to assess automatically. It becomes highly problematic, when the targets represent different software, which additionally may have the same name. In our sample of CZI, 7 out of 62 mentions (11.3%) had links to multiple different software.

2. Linking was created automatically in CZI, but the linking quality is generally poor: out of 55 mentions in our sample for which one, or multiple of the same link were provided, 36 (65.4%) gave a link to the wrong software. Of these, 9 were based on exact string matches between software name and mention string.

3. Not all software mentions have a link associated with them (16 out of 78 potentially linkable software mentions, 20.5% in our sample)

We are hopeful that the quality of machine learning models and algorithms for linking software mentions to repositories will be improved in the future. We nevertheless reason that links to repositories are not currently a dataset feature that improves the usability of mention datasets for quantitative research. Alternatively, datasets with repository URLs mined directly from publications (e.g., Escamilla et al. (2022)) could be used for quantitative research, although the semantics between the contents of the publication and the URL cannot be established to any satisfactory degree. One approach to improve the latter would be to combine URL mining and mention retrieval in the same dataset.

The underlying issue in terms of the challenges of using software mention datasets for research purposes as discussed here is the currently still suboptimal practice of software citation. If authors more strictly followed the software citation principles (A. M. Smith et al., 2016), software mentioned in the literature would be accessible per default, not via tedious manual search.

We believe that our approach to assess the usability of software mention datasets is generally valid, but it comes with its own shortcomings. Manual annotation of samples does not scale, and thus our approach can only be used for preliminary results from exploratory studies. The approach could potentially benefit not only from better linking models and algorithms as mentioned above, but also from a machine learning model to categorize mention types. Another shortcoming of our approach is the use of mention types that do not optimally reflect the software citation principles. Instead, we reused the types from Howison and Bullard (2015) to achieve comparability. Future work could therefore attempt the development of new mention type categories closed related to the software citation principles.

Based on the exploratory research presented here, we recommend features that software mention datasets can include to enable research on research software:

1. Software mentions datasets should be ordered by uniquely identified mention, not by mentioning publication.

2. Software mentions datasets should be made available as tabular data.

3. Software mentions datasets should include contexts for all mentions of a software across a publication.

4. Software mentions datasets should only include mentions that are to software, not to other entities.

5. URLs to software repositories in software mentions datasets should resolve to repositories containing the mentioned software.

**Exploratory studies**

The results from our exploratory studies (see *Exploratory studies*) are not representative. Their following discussion is therefore highly preliminary.

For *RQ2: Is open source software more cited in a way that allows credit for software authors than closed source software?* we hypothesized that commercial/close source software is cited more frequently using a lower-quality in-text name mention or citation to project name or website, and that open source software is cited more frequently with a repository or associated research publication in the reference. The results from our sample annotations seem to support this hypothesis. They showed that a third of openly licensed software still has a poor quality mention in our sample, suggesting that efforts towards better software citation are still necessary.

We were also interested in finding out the state of software citation – or mentioning – as compared to 2015. For *RQ3: Has the practice of software citation represented in software mention datasets improved in comparison to the practice described in Howison and Bullard (2015)?* we hypothesized that categories of mentions that reflect the principles better are found relatively more often in mentions from publications in the dataset samples that were published in or after 2016, than in the data

presented in Howison and Bullard (2015). Our results suggest that – at least in our data samples – the practice of software citation has not improved in the last 8 years.

The results from both exploratory studies would support previous work that describes challenges for software citation () and call for improved practice of software citation (Bouquin et al., 2023) and advocacy towards it (Du et al., 2022).

### Limitations

The positive evaluation of our assessment approach is threatened by medium value Krippendorff's $\alpha$ inter-annotator agreement scores (Table 1), where a value of $\alpha 1.0$ signifies complete agreement. Regarding the central "mention type" annotation, agreement showed only a slim positive trend ($\alpha 0.55$), and $\alpha 0.64$ across all layers. More confidence can be put into the assessment of the similarly important mention extraction quality annotations, at $\alpha 0.72$. Likewise, the qualitative observation towards answering RQ1 (RQ1.1, RQ1.2) is subjective by nature.

Our studies towards answering RQ2 and RQ3 were based on very small sample sizes. This would be a severe threat to validity of results that claim more than exploratory significance; note that we do not make such a claim here. Additionally, both sampled datasets include mostly biomedical software, and their data could not be used to make claims for research software in general.

## CONCLUSION

We presented an approach for the assessment of the usability of software mentions datasets for research on research software. Despite some small shortcomings, our approach is valid and applicable in the assessment of software mention datasets. It includes sampling and data preparation, manual annotation for quality and mention characteristics, and annotation analysis. We applied our approach to two software mention datasets (<**empty citation**>) to evaluate the approach. Through our approach, we were able to find challenges to working with the selected datasets. We were also able to define dataset features that would enable the use of software mention datasets in quantitative research on research software, and research into the practice of software citation. These features include: a dataset structure based on individual software mentions, persistence as tabular data, the inclusion of contexts for all mentions of individual software across an individual publication. We also found that automatically retrieved links to repositories for a software were of low quality in our sample; better retrieval algorithms and machine learning models are needed to improve the quality of links. Finally, the underlying issue with and challenge to working with software mention datasets is the suboptimal practice of software citation. We conclude that software should be *cited* using a formal citation and reference according to the software citation principles (A. M. Smith et al., 2016), and *not mentioned*.

## ACKNOWLEDGMENTS

## REFERENCES

Abramatic, J.-F., Cosmo, R. D., & Zacchiroli, S. (2018-10-01, 2018). Building the universal archive of source code (ACM, Ed.). *Communications of the ACM*, *61*(10), 29–31. https://doi.org/10.1145/3183558

Allen, A., & Schmidt, J. (2015). Looking before leaping: Creating a software registry. *Journal of Open Research Software*, *3*(1). https://doi.org/10.5334/jors.bv

Beltagy, I., Lo, K., & Cohan, A. (2019). SciBERT: A Pretrained Language Model for Scientific Text. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 3613–3618. https://doi.org/10.18653/v1/D19-1371

Bouquin, D., Trisovic, A., Bertuch, O., & Colón-Marrero, E. (2023). Advancing Software Citation Implementation (Software Citation Workshop 2022). (arXiv:2302.07500). https://doi.org/10.48550/arXiv.2302.07500

Chue Hong, N., Cope, J., Herterich, P., Katz, D. S., & Worthington, S. (2021). Recognising the value of software: How libraries can help the adoption of software citation. https://doi.org/10.6084/m9.figshare.14825268.v1

Chue Hong, N. P., Allen, A., Gonzalez-Beltran, A., de Waard, A., Smith, A. M., Robinson, C., Jones, C., Bouquin, D., Katz, D. S., Kennedy, D., Ryder, G., Hausman, J., Hwang, L., Jones, M. B., Harrison, M., Crosas, M., Wu, M., Löwe, P., Haines, R., Edmunds, S., Stall, S., Swaminathan, S., Druskat, S., Crick, T., Morrell, T., & Pollard, T. (2019). *Software Citation Checklist for Authors* (FORCE11 Software Citation Implementation Working Group Report). Zenodo. https://doi.org/10.5281/zenodo.3479199

Chue Hong, N. P., Allen, A., Gonzalez-Beltran, de Waard, A., Smith, A. M., Robinson, C., Jones, C., Bouquin, D., Katz, D. S., Kennedy, D., Ryder, G., Hausman, J., Hwang, L., Jones, M. B., Harrison, M., Crosas, M., Wu, M., Löwe, P., Haines, R., Edmunds, S., Stall, S., Swaminathan, S., Druskat, S., Crick, T., Morrell, T., & Pollard, T. (2019). *Software Citation Checklist for Developers* (FORCE11 Software Citation Implementation Working Group Report). Zenodo. https://doi.org/10.5281/zenodo.3482769

Chue Hong, N. P., Katz, D. S., Barker, M., Lamprecht, A.-L., Martinez, C., Psomopoulos, F. E., Harrow, J., Castro, L. J., Gruenpeter, M., Martinez, P. A., & al., e. (2021). FAIR principles for research software (FAIR4RS principles). *Research Data Alliance*. https://doi.org/10.15497/RDA00065

Dalitz, W., Sperber, W., & Chrapary, H. (2020). swMATH: A publication-based approach to mathematical software. *SIAM Newsletter*, *Volume 53*(Number 06 — July/August 2020). https://doi.org/10.12752/8009

Druskat, S. (2020). Software and Dependencies in Research Citation Graphs. *Computing in Science & Engineering*, *22*(2), 8–21. https://doi.org/10.1109/MCSE.2019.2952840

Druskat, S., & Chue Hong, N. (2023). Don't mention it: Challenges to using software mentions to investigate citation and discoverability - Data and Notebooks. https://doi.org/10.5281/zenodo.5518122

Druskat, S., Katz, D. S., & Todorov, I. T. (2021). Research Software Sustainability and Citation. *2021 IEEE/ACM International Workshop on Body of Knowledge for Software Sustainability (BoKSS)*, 1–2. https://doi.org/10.1109/BoKSS52540.2021.00008

Druskat, S., Spaaks, J. H., Chue Hong, N., Haines, R., Baker, J., Bliven, S., Willighagen, E., Pérez-Suárez, David, & Konovalov, O. (2021). Citation File Format. https://doi.org/10.5281/ZENODO.5171937

Du, C., Cohoon, J., Lopez, P., & Howison, J. (2021). Softcite dataset: A dataset of software mentions in biomedical and economic research publications. *Journal of the Association for Information Science and Technology*, *72*(7), 870–884. https://doi.org/10.1002/asi.24454

Du, C., Cohoon, J., Lopez, P., & Howison, J. (2022). Understanding progress in software citation: A study of software citation in the CORD-19 corpus. *PeerJ Computer Science*, *8*, e1022. https://doi.org/10.7717/peerj-cs.1022

Eitzen, C. (2020). *Research Software - Publication and Sustainability* [Master's thesis, Kiel University]. https://oceanrep.geomar.de/51354/

Escamilla, E. (2023). Extract-URLs. swh:1:snp:689cdf3440075d94e250b1da9f9e9d43c4efe675

Escamilla, E., Klein, M., Cooper, T., Rampin, V., Weigle, M. C., & Nelson, M. L. (2022). The Rise of GitHub in Scholarly Publications. In G. Silvello, O. Corcho, P. Manghi, G. M. Di Nunzio, K. Golub, N. Ferro, & A. Poggi (Eds.), *Linking Theory and Practice of Digital Libraries* (pp. 187–200, Vol. 13541). Springer International Publishing. https://doi.org/10.1007/978-3-031-16802-4_15

European Organization For Nuclear Research & OpenAIRE. (2013). Zenodo. https://doi.org/10.25495/7GXK-RD71

Felderer, M., Goedicke, M., Grunske, L., Hasselbring, W., Lamprecht, A.-L., & Rumpe, B. (2023). *Toward Research Software Engineering Research* (tech. rep.). Zenodo. https://doi.org/10.5281/zenodo.8020525

Fenner, M., Katz, D. S., Nielsen, L. H., & Smith, A. (2018). DOI Registrations for Software. https://doi.org/10.5438/1NMY-9902

Grisel, O., Mueller, A., Lars, Gramfort, A., Louppe, G., Prettenhofer, P., Fan, T. J., Blondel, M., Niculae, V., Nothman, J., Joly, A., Lemaitre, G., Vanderplas, J., Estève, L., kumar, m., du Boisberranger, J., Qin, H., Hug, N., Varoquaux, N., Layton, R., Metzen, J. H., Jalali, A., (Venkat) Raghav, R., Schönberger, J., Yurchak, R., Jerphanion, J., la Tour, T. D., Li, W., Marmo, C., & Woolam, C. (2023). Scikit-learn/scikit-learn: Scikit-learn 1.2.2. https://doi.org/10.5281/zenodo.7711792

Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., Fernández del Río, J., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., & Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, *585*, 357–362. https://doi.org/10.1038/s41586-020-2649-2

Howison, J., & Bullard, J. (2015). Software in the scientific literature: Problems with seeing, finding, and using software mentioned in the biology literature. *Journal of the Association for Information Science and Technology*, *67*(9), 2137–2155. https://doi.org/10.1002/asi.23538

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, *9*(3), 90–95. https://doi.org/10.1109/MCSE.2007.55

Istrate, A.-M., Li, D., Taraborelli, D., Torkar, M., Veytsman, B., & Williams, I. (2022). A large dataset of software mentions in the biomedical literature. https://doi.org/10.48550/arXiv.2209.00693

Istrate, A.-M., Veytsman, B., Li, D., Taraborelli, D., Torkar, M., & Williams, I. (2022). CZ Software Mentions: A large dataset of software mentions in the biomedical literature. https://doi.org/10.5061/DRYAD.6WWPZGN2C

Jay, C., Haines, R., & Katz, D. S. (2021). Software Must be Recognised as an Important Output of Scholarly Research. *International Journal of Digital Curation*, *16*(1), 6. https://doi.org/10.2218/ijdc.v16i1.745

Jones, M. B., Boettiger, C., Mayes, A. C., Smith, A., Slaughter, P., Niemeyer, K., Gil, Y., Fenner, M., Nowak, K., Hahnel, M., Coy, L., Allen, A., Crosas, M., Sands, A., Hong, N. C., Cruse, P., Katz, D., & Goble, C. (2017). CodeMeta: An exchange schema for software metadata. https://doi.org/10.5063/SCHEMA/CODEMETA-2.0

Katz, D. S., Bouquin, D., Hong, N. P. C., Hausman, J., Jones, C., Chivvis, D., Clark, T., Crosas, M., Druskat, S., Fenner, M., Gillespie, T., Gonzalez-Beltran, A., Gruenpeter, M., Habermann, T., Haines, R., Harrison, M., Henneken, E., Hwang, L., Jones, M. B., Kelly, A. A., Kennedy, D. N., Leinweber, K., Rios, F., Robinson, C. B., Todorov, I., Wu, M., & Zhang, Q. (2019). Software Citation Implementation Challenges. https://doi.org/10.48550/arXiv.1905.08674

Katz, D. S., Chue Hong, N. P., Clark, T., Muench, A., Stall, S., Bouquin, D., Cannon, M., Edmunds, S., Faez, T., Feeney, P., Fenner, M., Friedman, M., Grenier, G., Harrison, M., Heber, J., Leary, A., MacCallum, C., Murray, H., Pastrana, E., Perry, K., Schuster, D., Stockhause, M., & Yeston, J. (2021). Recognizing the value of software: A software citation guide. *F1000Research*, *9*, 1257. https://doi.org/10.12688/f1000research.26932.2

Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B., Bussonnier, M., Frederic, J., Kelley, K., Hamrick, J., Grout, J., Corlay, S., Ivanov, P., Avila, D., Abdalla, S., & Willing, C. (2016). Jupyter Notebooks – a publishing format for reproducible computational workflows. In F. Loizides & B. Schmidt (Eds.), *Positioning and power in academic publishing: Players, agents and agendas* (pp. 87–90). IOS Press.

Konovalov, O., Ye, H., Chisholm, L., Turner, M., Chue Hong, N. P., Buzzard, S., & Druskat, S. (2021). Introducing Habeas Corpus: A Collaborations Workshop 2021 HackDay Project. https://archive.softwareheritage.org/swh:1:cnt: d2e0fa61ff93704cbca9638d911d601498667019;origin=https://github.com/softwaresaved/habeas-corpus;visit=swh: 1:snp:17af7c58cac3146d7d6ab468175a10bf86b1126d;anchor=swh:1:rev:8cfd4ebb43dd6406fd413fd77be136899c72cc56; path=/docs/CW21%5C_Habeas%5C_Corpus%5C_Presentation.pdf

Kurnatowski, L., Stoffers, M., & Haupt, C. (2021). Provenance based software dashboards. *Workshop on the Science of Scientific-Software Development and Use*. https://elib.dlr.de/147617/

Peroni, S., & Shotton, D. (2020). OpenCitations, an infrastructure organization for open scholarship. *Quantitative Science Studies*, *1*(1), 428–444. https://doi.org/10.1162/qss_a_00023

Registries Task Force on Best Practices for Software, Monteil, A., Gonzalez-Beltran, A., Ioannidis, A., Allen, A., Lee, A., Bandrowski, A., Wilson, B. E., Mecum, B., Du, C. F., Robinson, C., Garijo, D., Katz, D. S., Long, D., Milliken, G., Ménager, H., Hausman, J., Spaaks, J. H., Fenlon, K., Vanderbilt, K., Hwang, L., Davis, L., Fenner, M., Crusoe, M. R., Hucka, M., Wu, M., Hong, N. C., Teuben, P., Stall, S., Druskat, S., Carnevale, T., & Morrell, T. (2020). Nine Best Practices for Research Software Registries and Repositories: A Concise Guide. *arXiv:2012.13117 [cs]*. Retrieved December 26, 2020, from http://arxiv.org/abs/2012.13117

Schindler, D., Bensmann, F., Dietze, S., & Krüger, F. (2021a). SoftwareKG-PMC. https://doi.org/10.5281/ZENODO.5713973

Schindler, D., Bensmann, F., Dietze, S., & Krüger, F. (2021b). SoMeSci- A 5 Star Open Data Gold Standard Knowledge Graph of Software Mentions in Scientific Articles. *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 4574–4583. https://doi.org/10.1145/3459637.3482017

Schindler, D., Bensmann, F., Dietze, S., & Krüger, F. (2022). The role of software in science: A knowledge graph-based analysis of software mentions in PubMed Central. *PeerJ Computer Science*, *8*, e835. https://doi.org/10.7717/peerj-cs.835

Schmidt, B., McGillivray, B., Larrousse, N., Broeder, D., Wilson, K., & Chue Hong, N. P. (2021). LIBER 2021 Session #3: Working with Software & Data. https://doi.org/10.5281/ZENODO.5036311

Smith, A. (2021). Enhanced support for citations on GitHub. https://github.blog/2021-08-19-enhanced-support-citations-github/

Smith, A. M., Katz, D. S., Niemeyer, K. E., & FORCE11 Software Citation Working Group. (2016). Software citation principles. *PeerJ Computer Science*, *2*(e86). https://doi.org/10.7717/peerj-cs.86

Spaaks, J. H., Klaver, T., Verhoeven, S., Diblen, F., Maassen, J., Sang Tjong Kim, E., Pawar, P., Meijer, C., Ridder, L., Kulik, L., Bakker, T., van Hees, V., Bogaardt, L., Mendrik, A., van Es, B., Attema, J., van Hage, W., Ranguelova, E., van Nieuwpoort, R., Gey, R., & Zach, H. (2020). Research Software Directory. https://doi.org/10.5281/ZENODO.1154130

Stol, K.-J., & Fitzgerald, B. (2018). The ABC of Software Engineering Research. *ACM Transactions on Software Engineering and Methodology*, *27*(3), 1–51. https://doi.org/10.1145/3241743

The Dask 2023.3.1 developers. (2023). Dask. https://pypi.org/project/dask/2023.3.1

The Jupyter Notebook 6.5.3 developers. (2023). Jupyter notebook (version 6.5.3-e3e14a1). https://pypi.org/project/notebook/6.5.3/

The Matplotlib 3.7.1 developers. (2023). Matplotlib. https://pypi.org/project/matplotlib/3.7.1/

The nbconvert 7.2.10 developers. (2023). Nbconvert. https://pypi.org/project/nbconvert/7.2.10

The NumPy 1.24.2 developers. (2023). NumPy. https://pypi.org/project/numpy/1.24.2/

The pandas development team. (2023). Pandas-dev/pandas: Pandas (version 1.5.3). https://doi.org/10.5281/zenodo.7549438

The Research Software Encyclopedia project. (2021). Research Software Encyclopedia. Retrieved August 24, 2021, from https://rseng.github.io/software/

The SciPy 1.10.1 developers. (2021). SciPy. https://pypi.org/project/scipy/1.10.1/

The Software Mention Extraction authors. (2022). Software mention extraction and linking from scientific articles. https://archive.softwareheritage.org/swh:1:cnt:149c13ef5b99c3ad52903305daee718bd09bfff7;origin=https://github.com/chanzuckerberg/software-mention-extraction;visit=swh:1:snp:592c9f600b721e8d56a59bcb424bf9a969a942fc;anchor=swh:1:rev:99296e04c9a982b5f31d4bf2dd33fec3894a385e;path=/README.md;lines=7-8

van de Sandt, S., Nielsen, L. H., Ioannidis, A., Muench, A., Henneken, E., Accomazzi, A., Bigarella, C., Lopez, J. B. G., & Dallmeier-Tiessen, S. (2019). Practice meets Principle: Tracking Software and Data Citations to Zenodo DOIs. *arXiv:1911.00295 [cs]*. Retrieved August 12, 2020, from http://arxiv.org/abs/1911.00295

Veytsman, B. (2022). Software mentions extractor (lines 65-69). https://archive.softwareheritage.org/swh:1:cnt:398bf69be085934e489d2a598ac;origin=https://github.com/chanzuckerberg/software-mentions;visit=swh:1:snp:699d5878f3cfc43c3fc57c98dd19455a370359d2;anchor=swh:1:rev:561ae152ff5e064fed85b9327afbeec36439fa63;path=/software-mentions-extractor/software-mentions-extractor.py;lines=65-69

Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Millman, K. J., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C. J., Polat, İ., Feng, Y., Moore, E. W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., & SciPy 1.0 Contributors. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in python. *Nature Methods*, *17*, 261–272. https://doi.org/10.1038/s41592-019-0686-2

Wade, A. D., & Williams, I. (2021). CORD-19 Software Mentions. https://doi.org/10.5061/DRYAD.VMCVDNCS0

Wang, L. L., Lo, K., Chandrasekhar, Y., Reas, R., Yang, J., Eide, D., Funk, K., Kinney, R. M., Liu, Z., Merrill, W., Mooney, P., Murdick, D., Rishi, D., Sheehan, J., Shen, Z., Stilson, B., Wade, A. D., Wang, K., Wilhelm, C., Xie, B., Raymond, D. A., Weld, D. S., Etzioni, O., & Kohlmeier, S. (2020). CORD-19: The COVID-19 Open Research Dataset. *NLPCOVID19*.