

Software and Dependencies in Research Citation Graphs

Stephan Druskat

Abstract—Following the widespread digitalization of scholarship, software has become essential for research, but the current sociotechnical system of citation does not reflect this sufficiently. Citation provides context for research, but the current model for the respective research citation graphs does not integrate software. In this paper, I develop a directed graph model to alleviate this, describe challenges for its instantiation, and give an outlook of useful applications of research citation graphs, including transitive credit.

Index Terms—software citation, citation graphs, transitive credit

1 INTRODUCTION

THE DIGITALIZATION of research changes research methods across disciplines, and produces new forms of research and knowledge [1]. In the process, research software has clearly become an integral part of digital research methodologies [2]. It embeds research knowledge, implements algorithms and models, and is a central component of digital scholarly integration and application [3]. It thus presents a significant, and increasingly vital, intellectual contribution to academic research.

Research software should therefore also be considered a legitimate research product [3], [4], [5]. Research products have the most value – and their outcomes can be understood most fully – when they are considered in their context [1, p. 10]. This context is standardly provided through citation. Therefore, one aspect of research software gaining the status of a research product is, that it must be integrated into the scholarly citation system.

The citation system in current digital scholarship is a sociotechnical system based on technical infrastructure, and involves different stakeholders. Stakeholders include domain-specific communities of researchers who use software, research software engineers (RSEs) and other software developers, research institutions, publishers, repository providers, index providers, and funders (cf. [6]). Technical infrastructure upon which the citation system is built most prominently include publication repositories; citation indices and aggregators; publishing services including websites; services provided by libraries; resolvers for digital identifiers; metadata formats; reference management, text processing, and other software.

Citation and the sociotechnical citation system broadly provide the following functions:

- **Context function:** The provision of context for research products by establishing a graph of research products with links between the identified citing and cited products, to enable traceability of outcomes, both over the past to understand how present knowledge was established, and into the future to understand how present knowledge is being used (cf. [1], [7], [8], [9])
- **Social functions:** The establishment of trust and authority [9], [10], [11]; the recognition of the value of a research product while providing credit to its authors [4], [7]; the potential for evaluation of individual researchers [7, ch. 6], individual research products [12], journals [13], and research groups, institutions, and countries [7, ch. 6]
- **Compliance function:** The assertion of compliance with good scholarly practice [7], [14]
- **Discursive function:** The organization and shaping of discourses of scholarly credibility, authority, and relevance through epistemic change via “dynamically rewriting the past” [7, p. xvi], [11], [15], [16])
- **Reproducibility function:** The enablement of research reproducibility through correct and complete citation [17], [18], [19], [20]

Software as a research product can be subject to all of the described functions – including the discursive function, albeit to a limited degree, see below – only if it is fully integrated in the citation system. While this is not currently the case [3], [18], [21], [22], [23], [24], progress is being made, driven by different stakeholders:

- Research software community initiatives such as the FORCE11 Software Citation Implementation Working Group (www.force11.org/group/software-citation-implementation-working-group) build on established community standards [4] and bring together stakeholders to shape technical infrastructure and policy, and develop guidance [6]; their activities

• Stephan Druskat is with the German Aerospace Center (DLR), Berlin, Germany, the Computer Science Department at Humboldt-Universität zu Berlin, Berlin, Germany, and the Department of English Studies at Friedrich Schiller University, Jena, Germany.
E-mail: stephan.druskat@dlr.de

© 2019 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.
DOI : <https://doi.org/10.1109/MCSE.2019.2952840>

concern the discursive and social functions directly, and the remaining functions indirectly.

- Domain, infrastructure and software communities develop software solutions for providing citation metadata [25], [26], create repositories, information services, and indices (e.g., [27], [28], [29], [30]), and develop metadata formats [31], [32]; their activities concern the context and social functions directly, and the remaining functions indirectly.
- Research policy researchers develop procedures for evaluating research software [33]; their activities concern the social and compliance functions directly, and the remaining functions indirectly.
- Funding agencies update funding policies (cf. [5]) and guidelines for scholarly practice [34] to incorporate citation of research software; their activities concern the compliance function directly, and the remaining functions indirectly.
- Publishers establish editorial policies that require the citation of software, sometimes as a subset of data [24], [35], or plan to do so [18]; their activities concern the context, social and reproducibility functions directly, and the remaining functions indirectly.

In this paper, I aim to contribute to the understanding of the requirements for the implementation of research software citation. To this end, I will investigate the output of the context function of citation, *research citation graphs* (RCGs), with the objective to answer the following research questions:

- **RQ1:** What are the necessary changes in the model of research citation graphs to allow for the integration of research software, and the adoption of the citation functions?
- **RQ2:** What are the requirements for the implementation of software citation based on an updated model of research citation graphs?
- **RQ3:** What are current challenges for the instantiation of research citation graphs?
- **RQ4:** What applications do research citation graphs enable?

2 RESEARCH CITATION GRAPHS

Research products and the references between them can be modeled as a directed graph $G_1 = (V, E)$ where V is a set of vertices (or “nodes”), and E is a set of ordered pairs of nodes (i.e., “directed edges”). The nodes in V represent research products, the edges represent reference relations (i.e., citation) between source nodes (the citing research product) and target nodes (the cited product). This most basic model of a *research citation graph* (RCG) enables the context function of citation: It helps understand what other research products a specific product relied on (“back-tracking”), or has led to (“forward-tracking”), in order to build on this understanding in research, or conduct evaluations and measurements. Both tracking methods can be implemented as graph traversal, where back-tracking follows outgoing edges, and forward-tracking follows incoming edges. The model also enables the social function of establishing trust and authority, where it is based on the citation of acknowledged trustworthy or authoritative

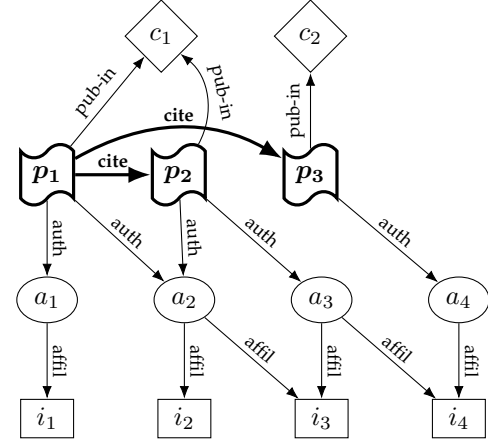


Fig. 1. Partial G_2 -type RCG for a research product which references two other research products. The graph includes the respective G_1 graph (in bold print). Nodes: p_n =research product, a_n =author, i_n =evaluable affiliation, c_n =evaluable publishing container; Edges: cite=“citation” relation, auth=“authored by” relation, affil=“affiliated with” relation, pub-in=“published in” relation.

research products rather than establishing references to their authors.

In order to fully exploit the social citation function, RCGs must model additional properties of research products:

- The provision of academic credit requires the inclusion of author nodes, and authorship relations between them and research products, as does the establishment of trust and authority if it is based on individuals.
- Evaluation requires the inclusion of two classes of entity nodes: affiliations (research groups, institutions, countries, etc.), and respective affiliation relations between them and authors; “product containers” (journals, edited volumes, repositories, archives, etc.), and respective published-in *part-of* relations between them and research products.

The model graph changes accordingly: Let P be the set of all vertices $\{p_1, \dots, p_n\}$ which represent research products, A the set of all vertices $\{a_1, \dots, a_n\}$ which represent authors, I the set of all vertices $\{i_1, \dots, i_n\}$ which represent evaluable author affiliations, and C the set of all vertices $\{c_1, \dots, c_n\}$ which represent entities which contain research products. Let \mathcal{V} be the set of disjoint sets $\{P, A, I, C\}$ of vertices in the RCG $G_2 = (V, E)$. Define $L : V \rightarrow \mathcal{V}$ to set

- $L(v) = P$ when $v \in P \in \mathcal{V}$,
- $L(v) = A$ when $v \in A \in \mathcal{V}$,
- $L(v) = I$ when $v \in I \in \mathcal{V}$,
- $L(v) = C$ when $v \in C \in \mathcal{V}$.

See Figure 1 for an exemplary visualization. In the past, the citation system, as the foundation of the academic credit and evaluation system, focused on journal articles, books, and conference papers [1], and G_2 -type RCGs are the output of the system’s context function at this stage. In addition to the more general reasons given above in section 1, research software must be integrated in the citation system also to realize the compliance and reproducibility functions of the current sociotechnical citation system. Definitions of and

guidelines for good scholarly practice change to reflect the digitalization of scholarship, and compliance with them requires research software to be cited. The recently updated guidelines for good scholarly practice by German funding agency Deutsche Forschungsgemeinschaft, for example, require that “provenance of data, organisms, materials and software used in the research process is indicated and reuse documented; the original sources are cited” [34, p. 14, my translation]. Computational reproducibility also naturally requires, amongst other information, e.g., about the computational environment, the correct identification of the software that has been used in published research [17], [36], and provision of this information through citation [4], [18], [37], [38].

In order to enable these functions in RCGs, they must model the citation-related specific properties which set software apart from other research products. These properties relate to the software citation principles of *specificity*, *unique identification*, and *attribution and credit* [4]. Software differs from textual research products in the form its artifacts can take, in its notion of finality and the relationships between its artifacts, in the citability of its concepts, in its dynamicity, in the containment relationships between a product and its contributions, and in the roles which contribute to it.

Take an academic paper as example: As a research product, it is perceived to be a single artifact, available in a “final” version. This finality is a function of peer review, editorial acceptance, and “adequate” publication. Final published papers cannot have a new version after publication. Instead, any changes must again pass peer review, acceptance and publication, at which stage the changed paper is considered to be a discrete new research product. Papers may also have non-final versions published as preprints.

In contrast, software development processes produce artifacts at different stages. Every commit (or “revision”) in a version control system produces an artifact, a collection of source code and other files, and build processes may additionally produce one or more binary artifacts. The notion of a “final” software product does not exist as such. This is due to the lack of a standardized publication process for software which is based on peer review. Instead, software can evolve over time, and any revision can be tagged as a version and “released”. None of the versions can be considered “final”, as at any time, any changes produce a newer version of the same software, not a discrete new research product. Also, while versions of papers may or may not differ in content – preprints may have the same content as a final publication – versions of software will usually represent differing source codes. When we talk about “software”, we usually mean a version of a software that has been released or used. Alternatively, “software” can also refer to the concept of a software (cf. [6]), rather than a version, revision, or artifact: “Microsoft Excel” refers to the concept of a spreadsheet application, of which its versions are realizations. Although a paper and its preprints are arguably also different realizations of the same concept, these paper concepts are never explicitly cited, whereas software concepts may be cited, e.g., in the case of pipelines or frameworks [6], or to understand the development of a software in computer science research, or in software comparison. In analogy, software concepts may be iden-

tifiable by unique identifiers as issued by repositories for digital research products such as Zenodo [39]. For papers, identifiers for concepts are not issued across repositories. Usually, preprints and final publications will be archived on separate platforms. Although preprint repositories such as arXiv do issue versioned identifiers, which point to the latest version when stripped of the version information, there are no cross-repository identifiers that can uniquely identify concepts of papers.

Defined as a “set of instructions that direct a computer to do a specific task” [40, p. 2], software is “functionally active” [41, p. 2] (“dynamic”), i.e., it performs functions on data, whereas papers are clearly static. A software may have different states, and execute along different paths at runtime. The final states and execution paths depend on configuration, interaction, and possibly the data that is being processed. States and execution paths define the actual “dynamic product” that is used to perform a specific software task.

A further difference between software and papers is the containment relationship between a product and the contributions to it. Contributions to a research product can be *active* or *passive*, and *direct* or *indirect*. In *direct active contributions* to papers, contributors influence the product directly through contributions of text, analyses, ideas, etc. Direct active contributions to software can take the form of source code, code comments, documentation, architectural design, API design, UI design, tests, code reviews, bug reports, etc. With *direct passive contributions*, a paper uses another product or parts thereof, by building on it, refuting it, refining its analyses, contextualizing its findings, etc. Direct passive contributions to software are mainly its “dependencies” – i.e., other software – but can also include other research products, such as publications that describe algorithms, models or methods implemented in the software. *Indirect contributions* to a research product are direct or indirect contributions to passive contributions to that product. Indirect software contributions to software are *transitive dependencies*.

“Dependencies” of a software S are software components to which S exhibits a degree of coupling. If S relies on functions of another software S_1 , without which S will not function as intended, S_1 is a dependency of S . This usually means that S calls functions from S_1 , or uses its API in another way, e.g., through inheritance. Dependencies can take different forms, as libraries, code fragments, or algorithms. The defining quality of a dependency is that it is not part of the original, directly contributed, source code of a software. Therefore, functions defined in a file X that are called from functions defined in another file Y , are part of a dependency iff file X is not part of the same codebase as file Y . This may include that file X has other authors than file Y . Original source code and dependency source code form the common codebase of a software. At runtime, direct passive contributions (dependencies) and indirect passive contributions (transitive dependencies, i.e., the dependencies of direct dependencies) become part of the same “software object”, as execution paths transcend boundaries between a software and its dependencies. In contrast, direct passive contributions to papers are part of the product in the form of references; indirect passive

contributions are not part of a paper at all, but must instead be retrieved via backtracking traversal of an RCG.

Direct active contributions to papers are recognized through authorship, and direct passive contributions through citation. Indirect contributions are not recognized, but can be discovered in RCGs. Similarly, direct active contributions to a software should be recognized through authorship – or acknowledged contributorship, see below – whereas reliance on dependencies as direct passive contributions should be recognized through citation. This holds despite the *part-of* relationship between dependencies, transitive dependencies, and the depending software at run-time, when dependencies arguably become direct active contributions to a software product. The software citation principles motivate this, by suggesting that software citation should generally address software source code [4], which makes dependencies passive contributions. Additionally, under a standard definition of authorship [42], the default recognition type for direct active contributions, authors of dependencies do not qualify for authorship of depending software: The contribution to a software through a dependency is substantial, but neither will they draft or revise the depending software, nor will they approve the version of the depending software to be published, or agree to be accountable for it. This categorically rules out authorship and motivates citation of dependencies instead.

Citation should attribute contributions to a research product to all contributors, and enable the provision of credit for a contribution. There is increasing acknowledgment of the fact that direct contributions to research products can take different forms than text production [43], and that metadata should represent different contribution types [44]. This is the case for all types of research product, and specifically for software, where creditable contributions greatly differ from those to papers, and include not just the writing of source code, but also contributions to the architecture, design, documentation, engineering, management, verification, validation, repair, maintenance, etc., of a software [45], [46], [47]. However, there is not yet a common understanding of which types of acknowledgeable and creditable contributions there are across different types of research software. It has also not yet been established whether there are qualitative differences between contribution types that may motivate a tiered concept of contributions, such as primary and secondary contributions. The schema.org ontology, for example, defines a term *contributor* as “a secondary contributor” [48], whereas community initiatives such as All Contributors [47] suggest to reward “every contribution, not just code” without defining differently valued contributions. It is therefore not possible to add contributor types to a model of RCGs at this stage in a meaningful and future-proof way. Instead, different stakeholders – software producers, researchers, publishers, institutions, policy makers, and others – should collaborate to develop and implement policies and ontologies that allow for a more differentiated model of author- and contributorship across all acknowledgeable research products. This model should replace the traditional authorship model and be reflected in the metadata for research products that is provided at publication, for purposes of citation, credit, etc.

In summary, the described specific properties of soft-

ware yield new requirements for a G_3 RCG model which integrates software research products, and supports the compliance and reproducibility functions:

- The versionability of software (and other) research products requires no new model elements, but the re-definition of the set P as the set of all vertices $\{p_1, \dots, p_n\}$ which represent *versions* of research products.
- The specific relationship between versions as realizations of concepts and the respective concepts requires the addition of concept nodes, i.e., a new set $O \in \mathcal{V}$, which is the set of all vertices $\{o_1, \dots, o_n\}$ which represent *concepts* of research products. Alternatively, P could be re-defined as the set of all vertices $\{p_1, \dots, p_n\}$ which represent *versions or concepts* of research products. For reasons of clarity, I will pursue the first option in the following.
- The relationships between versions and concepts also require the addition of realization relations, i.e., edges from nodes in P to nodes in O .
- For cases where concepts that overarch versions of a research product remain unidentifiable – as is usually the case with papers – the relationships between versions of a research product require the addition of order relations, i.e., edges from nodes in P to other nodes in P which define one version (the source node) as the predecessor of another version (the target node). These edges allow for the analysis of cumulative impact over versions of a research product.
- A differentiated model of contributorship to research products, and specifically software, requires the addition of edges representing different contribution types. As this is currently not possible in lieu of an agreed-upon model of acknowledgeable contributorship, the relations formerly denoting authorship have been labeled *contrib** in Figure 2, instead of *auth*, where the label is to be understood to reflect different types of contributions, including both traditional authorship and the more fine-grained contributions to digital research products.

These changes take into account almost all of the specific properties of software as a research product. “Dynamic products” have not been included here because they represent objects different from the principle target of software citation, i.e., source code (cf. [4]). They will be discussed further in section 3. See Figure 2 for a visualization showing an example of the updated G_3 model for RCGs. The figure shows a paper p_8 , which cites a software version p_1 . p_1 has been contributed to by two contributors, one of which (a_1) has made contributions of two different types, e.g., writing source code and writing documentation. p_1 also cites two dependencies p_3 and p_4 . p_3 cites its own dependency p_5 , making p_5 a transitive dependency of p_1 . p_4 is a realization of software concept o_1 , and precedes another version of the same software, p_6 , which is also a realization of o_1 . p_1 also cites a paper p_2 , e.g., describing an algorithm which p_1 implements. p_2 in turn has a predecessor in the preprint p_7 , with which it shares one contributor. For p_2 and p_7 , no concept has been published.

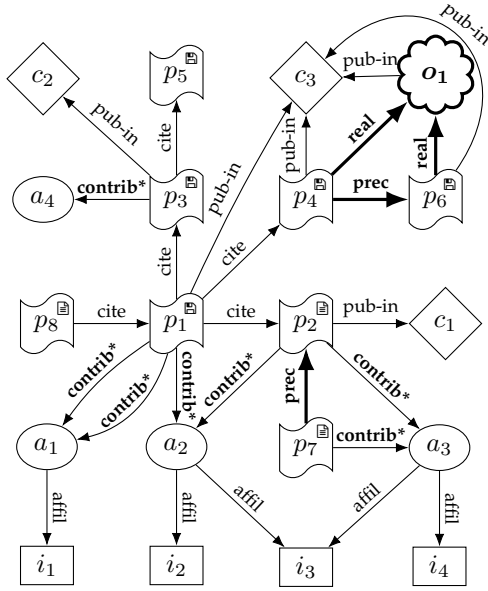


Fig. 2. Partial G_3 -type RCG for an exemplary textual research product (p_8) which cites an exemplary software (p_1) with two software dependencies (p_3, p_4), and a citation to a paper (p_2). Additional and changed elements in comparison to the G_2 model in bold print. Software products marked with \square , textual products marked with \diamond . Nodes: p_n =version of a research product, a_n =acknowledged contributor, i_n =evaluable affiliation, c_n =evaluable publishing container; Edges: *cite*="citation" relation, *contrib**="contributed to by" relation, *affil*="affiliated with" relation, *pub-in*="published in" relation, *prec*="precedes" relation, *real*="realizes" relation.

Taken together, the updated G_3 model for research citation graphs enables all functions of the current sociotechnical system of citation for research software (RQ1). With regards to the discursive function, some systemic limitations apply. First of all, the acknowledgement of direct passive software contributions to a research software, or failure to do so, can only be used to shape and organize discourses of scholarly credibility, authority, and relevance to a certain extent. As dependencies are hard-wired into the software product, only the preference of dependencies over functionally equivalent others before or while creating the product have the potential to affect a given discourse. To neglect, or favour, a specific dependency over another usually means a trade-off of functionality for a research software. These choices will be determined by the functional and/or engineering needs of any given software project, and are not likely to be used as discursive actions. A "re-writing of the past", by wilfully neglecting to *reference* specific dependencies and their import in the scholarly discourse, becomes impossible as soon as they are *used*, and thus become an inseparable part of, a software. Some build systems have a concept of optional dependencies. Optional dependencies, i.e., those that can optionally be used in a software which also defines a fallback mechanism for the case that the optional dependency cannot be resolved and used on a given system, may be more prone to allow discursive actions, but the discussion of this corner case is out of scope for this paper. Another case where the discursive function may be of interest is in references from a software to a text-based research product. This could be the case when a software

implements an algorithm described in a paper, and the paper should be cited to enable due credit for its contributors. Choosing to not cite the paper when it should be, is arguably discursive action with respect to creditability and relevance of research products. On the other hand, this goes against currently promoted good scholarly practice [34], and presents another corner case. Finally, and most importantly, the discursive function does not influence RCGs directly in terms of elements in a graph. Rather, its output must be studies on what is not in a given RCG instance when it should be, or what is in a graph when it should not be.

Based on the discussion and development of a useful model of research citation graphs that integrate software in this section, I will discuss current challenges for the instantiation of RCGs in the following section.

3 CHALLENGES FOR THE INSTANTIATION OF RCGs

Research citation graphs have a number of potential applications, which I will discuss in section 4. A prerequisite for these applications is the instantiation of research citation graphs for actual research products, including software. RCGs are built from distributed metadata, probably represented as linked data. An example of this is CodeMeta [32], a linked data exchange schema for software metadata which extends the schema.org vocabulary and is implemented in JSON-LD. CodeMeta files can, for example, represent software as well as other research products as references of a software. For journal papers, JATS [49] can be used to retrieve metadata and resolve references. Instantiated RCGs can, e.g., be stored in (graph) databases, represented in main memory, or visualized as graphs.

In theory, RCGs are instantiated by recursive resolution of references from a "root" research product to retrieve the set $P \in \mathcal{V}$ for a graph $G = (V, E)$; for each $p \in P$, retrieval of the product-specific metadata to retrieve the sets A, I, C , and O for G ; deduplication of the vertices in V for G .

This process yields requirements for an implementation of software citation (RQ2), which also reflect the software citation principles [4] (in parentheses):

- Publications and publication metadata – including references – are available digitally, and metadata are machine-actionable ("Persistence", "Accessibility");
- research products duly and correctly cite references including software in publications ("Importance", "Specificity", "Credit and attribution");
- publications are uniquely identifiable through a machine-actionable identifier ("Accessibility", "Unique identification");
- contributors, research institutions and other evaluable entities as well as publication platforms can be deduplicated, i.e., are uniquely identifiable.

To meet these requirements for research software, a number of challenges (RQ3) have to be overcome first. The publication process for textual research products is well-established and involves peer review, editorial acceptance and adequate publication together with curated and complete metadata and unique identification to enable the citation use case. No such process is yet in place for software,

although software journals such as JORS [50] or JOSS [51] aim to provide a similar workflow, but do not meet the “Importance” principle [4], as they publish metapapers, and not the software itself. Additionally, while references to text publications are recorded in their metapapers, software references (dependencies) are not. Alternatives include automated deposits from source code repositories such as GitHub to general purpose archives such as Zenodo, which provide unique identification, but do not require or curate citation-relevant metadata, including references. Similarly, Software Heritage [52] harvests source code and provides unique identification, but does not require citation-relevant metadata, including references.

In order to establish a publication process for software similar to that of papers, action is required from many different stakeholders. Research software creators will have to embed publication steps into their workflows, which also reflect short, iterative cycles of software development. These publication steps should be supported by repositories and archives, as well as publishers, which need to adopt and process metadata schemas that include references to other research products, and make reference metadata available in addition to product metadata [53]. Software publications must further be considered in evaluation processes by research institutions and funders, and thereby creating incentives for creators and institutions to publish research software. It will also have to be determined at which stage peer review and editorial acceptance – including the required curation of metadata – may be integrated. Peer review, editorial oversight, and metadata curation for software may be in the scope of publishers collaborating with archives. Instead of publishing metapapers, software journals could manage the peer review and acceptance process for a software deposit in archives. Suitable business models for this will have to be established, especially as software products can usually be published in cycles much shorter than those for papers. Another option for the curation of metadata would be for libraries to engage and adopt this task, again in collaboration with archives and research data management. Suitable metadata and exchange formats exist [31], [32].

Such adequate publication practices for software may also support the due citation of research software, which is still not widely established [3], [18], [21], [22], [23], [24]. In order to establish due citation of software, a culture change needs to take place across research, which can be driven from two directions: top-down by policy makers, institutions, funders, publishers and publishing platforms which can require the citation of software in research, and the provision of suitable metadata for software publications, and reward it based on adapted evaluation practices; bottom-up by educators, peer reviewers, editors, researchers, research software engineers, etc., which educate about, insist on, enable, and practice due citation of software and provision of citation metadata for software.

The unique identification of software products can be achieved through publication via archives that provide DOIs or similar persistent identifiers. While this solves the technical side of unique identification of software, the cultural challenge, i.e., the adoption of respective publication practices, remains to be solved as discussed above.

The unique identification of contributors and institutions

can be achieved through the use of identifiers such as ORCID. Again, this solves the technical side, but adoption remains a cultural challenge. This can be tackled through encouragement, request and requirement of adoption of identifiers from funders, publishers, and institutions, as well as through education and exemplary practice by researchers, software creators, and institutions.

While the requirements discussed above apply to all types of research products, there are two particular aspects of research software which further affect the instantiation of RCGs. As mentioned in section 2, software has two states: the static state of its source code, and the dynamic state at runtime. While strictly speaking, the dynamic state is irrelevant in a discussion of software citation as defined by the software citation principles [4], there are solutions that may enable the instantiation of RCGs constricted to a single software at runtime. Software that documents the executed paths taken in a complex software product to produce a research outcome at runtime can potentially produce an RCG for this execution. Duecredit [26] implements this concept for software written in Python. It registers references for portions of code at the module and function levels, and can inject references for dependencies. Its output is a list of references that represent the code that has been called during execution. This output can potentially be transformed to an RCG, albeit a local one which will usually contain only first and second level references of a software, depending on the downstream provision of the respective metadata by dependency projects. Additionally, the manually provided reference metadata cannot be verified.

While duecredit may not be suitable for the instantiation of non-local, larger-scale static RCGs, it brings into focus the fact that research software will often include dependencies or transitive dependencies that are not research products, and which will therefore not be published in an adequate way, and not come with citation-relevant metadata. In order to provide relevant context for research products, and enable all functions of citation for research software, these “hidden” contributions to research must be included in RCGs. This touches, in fact, the core of the discursive function of citation [16]. In lieu of publications, unique identifiers, and curated metadata for non-research software, this can be done by applying software engineering methods. Through static code analysis using manifests, build configurations, or import statements in conjunction with repository mining methods, a dependency graph can be retrieved for a given software, which can be transformed into a partial RCG for research software dependencies for which no machine-actionable metadata or unique identifiers are provided. I will investigate this method in future research.

In summary, the feasibility of instantiating RCGs that include research software is currently limited. This is due to unsatisfactory software publication practices, lack of provided correct and complete metadata, and insufficient software citation practices induced by lack of incentive and requirement to cite software. Solutions are being proposed and developed in technology [31], [32], policy [34], theory [4] and other areas [6]. These solutions support a culture change towards software citation implementation. Progress in software citation implementation will gradually unlock applications for RCGs. I will provide an outlook on

exemplary applications in the following section.

4 APPLICATIONS FOR RESEARCH CITATION GRAPHS

RCGs enable different analyses of the context of research products (RQ4). In this section, I will outline potentially useful analyses based on the visualization of the G_3 model for RCGs in Figure 2. The exemplary analyses also serve as indirect evaluation of the model.

Back-tracking exploration The context of research products can be explored using RCGs to find out which preceding research a research product builds on. This can be done by traversing the graph starting from p_8 and following outgoing edges of type *cite*. The analysis shows for example, that the paper p_8 indirectly builds on research published in p_2 . The implementation of software citation solicited above makes this insight possible, as without citation of p_1 in p_8 , and citation of p_2 in p_1 , the relation between p_8 and p_2 would remain hidden.

Citation tracking The context of research products can be explored using RCGs to find out which research builds upon a given research product. This can be done, e.g., by traversing the graph starting from a given research product node and following incoming edges of type *cite*. The analysis shows, for example, that paper p_2 has been cited by software p_1 , in addition to any other papers citing it (not visualized). Again, software citation makes this insight possible by providing not only references *to* software, but also references of other products *from* software.

Tracking of concept citation G_3 -type RCGs enable citation analyses for concepts in addition to products. This can be done by traversing the graph starting at o_1 , following incoming nodes of type *real* to products realizing the concept published in o_1 , and consecutively following outgoing *cite* relations from the realizing products. Given the citation of p_4 in p_1 , and assuming that p_6 was cited in another product p_9 (not visualized), this analysis yielded a citation count of 2 for the concept of a software which has been realized in two implementation versions of the software.

Contribution role analysis Traversing the graph starting from p_1 and following outgoing *contrib** relations allows an analysis of how roles are distributed over contributors to a research product. Once a sufficient model for contributions to research has been established, this also allows for a fine-grained evaluation of contributors with respect to their skill sets and creditable contributions.

Self-citation analyses RCGs enable self-citation analyses by finding nodes in A (for contributor-based analyses) or C (for, e.g., journal-based analyses) with more than one incoming relation of type *contrib**, and looking at whether their source nodes are connected directly with a *cite* relation.

Analysis of software development practices Traversing the graph from a software product node, e.g., p_6 , following incoming (or outgoing) *prec* relations, creates a timeline of versions of a software. Provided the respective publication dates deposited in the respective machine-readable metadata (not visualized), the common target node of the *real* relations can be taken into account to analyse the software development model employed for the implementation of

software concepts, e.g., o_1 . Given short timespans between versions, for example, an agile process could be assumed.

Credit for “hidden” contributions to research Assuming that p_5 is a commercially-developed software which was never intended to be published as a research product, and given that p_1 is research software, traversing the graph from p_1 following outgoing *cite* relations enables attribution of and credit for the contributors to p_5 (not visualized) for their contribution to the research published as p_1 . One obvious challenge here is the retrieval of complete and correct contributor information so that contributors to p_5 are correctly attributed and can receive credit.

Retrieval of transitive credit The last exemplary analysis (*Credit for “hidden” contributions to research*) already hints at the usability of RCGs for calculating fractional credit to be recorded in transitive credit maps for a research product. When fractional credit for a research product is not only distributed over authors – as is currently done implicitly through the order of authors lists for papers – but also over contributors and cited research products, a complete credit map for a research product is created. A credit map for a product A also feeds into the credit map of a product B that cites A. The main principle of transitive credit [9], [54] is that a contributor to product A can therefore receive credit for product B. If a fractional credit weight is determined for all contributors and references (i.e., built-upon research products), the credit weight for a single contributor or cited product can be determined transitively. If a software S is jointly developed by two people who share credit equally, both receive .5 fractional credit of the summary credit of 1 that can be distributed for any given research product. If S is cited in a paper P, and the fractional credit for S for its contribution to P is weighted to .2 (of 1), each contributor to S receives fractional credit of .1 transitively for P.

Different systems of weighting have been suggested for fractional credit, such as contribution taxonomies [43], [47] which could define default credit weights for different contribution types as templates for fractional credit. The same could be done for dependencies, which could be assigned a fixed (small) weight template for contributions to another software. However, this would disregard the fact that, for example, a software A has a different weight for a software B which provides a wrapper API for A, than for a software C which uses a single function of A and has a large number of additional dependencies. In future research, I will develop a weighting system of fractional credit for software dependencies which instead is based on software engineering metrics such as function call frequencies and complexity. In contrast to the retrieval of credit weights for other products than software, such a system does not rely on access to publication metadata, as it can use software engineering artifacts such as manifests and build configurations. It thus exploits the actual conditions found in the software and dependencies, rather than a meta-representation. The fractional credit weights will still need to be registered for publications, which requires solutions to challenges described in section 3.

Once fractional credit values are registered in metadata for research products during publication, RCGs make it easy to ascribe transitive credit to contributors to cited references, including dependencies. RCGs embed dependencies on a

par with other research products, following the software citation principles of “Importance” [4], while preserving the actual reference chains between software and its dependencies. In contrast, alternatives to the representation of dependencies for credit in RCGs (cf. [55]) may obfuscate these chains of references by citing software more granularly from a research product, i.e., software and dependencies are cited at one and the same depth rather than at the depth they have in the actual dependency graph; or, they place acknowledgeable contributions outside of citations altogether, and instead provide provenance information on software websites for example, thereby disregarding the “Importance” principle [24].

5 CONCLUSION

In this paper, I have introduced a directed graph model for research citation graphs that integrates software and dependencies (RQ1). This model can be used to determine requirements for the implementation of software citation based on established principles [4] (RQ2). These requirements are not currently met, and the current state of software citation poses challenges for the instantiation of the model (RQ3). These include: a lack of standardized publication practices for software; insufficient metadata provision and curation practices; a lack of incentives to cite software and give due credit to contributors to research software; insufficient use of unique identifiers for researchers and institutions. Some of these challenges can be tackled with software engineering methods and the application of good scholarly practice, others rely on a culture change concerning the attitude to software as a research product, and the implementation of respective practices. Once these challenges are overcome, research citation graphs based on the presented model enable a number of useful applications (RQ4), such as bibliometric and scientometric studies, analyses of software development workflows applied in research, and transitive credit.

ACKNOWLEDGMENTS

I would like to thank the discussion group on citation and rewarding systems at the Workshop on Sustainable Software Sustainability 2019 on 25 April 2019 in The Hague, Netherlands (www.software.ac.uk/wosss19). Discussion within the group has helped me to better understand the context for embedding software in the citation graph of research. The members of this group were: Neil Chue Hong, Gerard Coen, James Davenport, Leyla Garcia, Robert Haines, Catherine Jones, Adriaan Klinkenberg, Rachael Kotarski, Mateusz Kuzak, Brett Olivier, Esther Plomp, Shoaib Sufi, Stephanie van de Sandt, and Bettine van Willigen. I would also like to thank three anonymous reviewers for their very helpful suggestions.

REFERENCES

- [1] C. L. Borgman, *Scholarship in the Digital Age: Information, Infrastructure, and the Internet*. Cambridge, Mass: MIT Press, 2007, oCLC: ocm76794695.
- [2] C. Goble, “Better Software, Better Research,” *IEEE Internet Computing*, vol. 18, no. 5, pp. 4–8, Sep. 2014. [Online]. Available: <https://doi.org/10.1109/MIC.2014.88>
- [3] L. Hafer and A. E. Kirkpatrick, “Assessing Open Source Software As a Scholarly Contribution,” *Commun. ACM*, vol. 52, no. 12, pp. 126–129, Dec. 2009. [Online]. Available: <https://doi.org/10.1145/1610252.1610285>
- [4] A. M. Smith, D. S. Katz, K. E. Niemeyer, and FORCE11 Software Citation Working Group, “Software citation principles,” *PeerJ Computer Science*, vol. 2, no. e86, 2016. [Online]. Available: <https://doi.org/10.7717/peerj-cs.86>
- [5] H. Piwowar, “Altmetrics: Value all research products,” *Nature*, vol. 493, pp. 159–159, Jan. 2013. [Online]. Available: <https://doi.org/10.1038/493159a>
- [6] D. S. Katz, D. Bouquin, N. P. C. Hong, J. Hausman, C. Jones, D. Chivvis, T. Clark, M. Crosas, S. Druskat, M. Fenner, T. Gillespie, A. Gonzalez-Beltran, M. Gruenpeter, T. Habermann, R. Haines, M. Harrison, E. Henneken, L. Hwang, M. B. Jones, A. A. Kelly, D. N. Kennedy, K. Leinweber, F. Rios, C. B. Robinson, I. Todorov, M. Wu, and Q. Zhang, “Software Citation Implementation Challenges,” *arXiv:1905.08674 [cs]*, May 2019. [Online]. Available: <https://arxiv.org/abs/1905.08674>
- [7] N. De Bellis, *Bibliometrics and Citation Analysis: From the Science Citation Index to Cybermetrics*. Lanham, Md: Scarecrow Press, 2009, oCLC: ocn268952958.
- [8] E. Garfield, “Citation Indexes for Science: A New Dimension in Documentation through Association of Ideas,” *Science*, vol. 122, no. 3159, pp. 108–111, Jul. 1955. [Online]. Available: <https://doi.org/10.1126/science.122.3159.108>
- [9] D. Katz, “Transitive Credit as a Means to Address Social and Technological Concerns Stemming from Citation and Attribution of Digital Products,” *Journal of Open Research Software*, vol. 2, no. 1, p. e20, Jul. 2014. [Online]. Available: <https://doi.org/10.5334/jors.be>
- [10] D. Nicholas, A. Watkinson, R. Volentine, S. Allard, K. Levine, C. Tenopir, and E. Herman, “Trust and Authority in Scholarly Communications in the Light of the Digital Transition: Setting the scene for a major study,” *Learned Publishing*, vol. 27, no. 2, pp. 121–134, 2014. [Online]. Available: <https://doi.org/10.1087/20140206>
- [11] S. A. Greenberg, “How citation distortions create unfounded authority: Analysis of a citation network,” *BMJ*, vol. 339, p. b2680, Jul. 2009. [Online]. Available: <https://doi.org/10.1136/bmj.b2680>
- [12] C. Neylon and S. Wu, “Article-Level Metrics and the Evolution of Scientific Impact,” *PLOS Biology*, vol. 7, no. 11, p. e1000242, Nov. 2009. [Online]. Available: <https://doi.org/10.1371/journal.pbio.1000242>
- [13] E. Garfield, “Citation Analysis as a Tool in Journal Evaluation: Journals can be ranked by frequency and impact of citations for science policy studies,” *Science*, vol. 178, no. 4060, pp. 471–479, Nov. 1972. [Online]. Available: <https://doi.org/10.1126/science.178.4060.471>
- [14] A. Morin, J. Urban, P. D. Adams, I. Foster, A. Sali, D. Baker, and P. Sliz, “Shining Light into Black Boxes,” *Science*, vol. 336, no. 6078, pp. 159–160, Apr. 2012. [Online]. Available: <https://doi.org/10.1126/science.1218263>
- [15] M. Bunge, “Epistemic Change,” in *Epistemology & Methodology II: Understanding the World*, ser. Treatise on Basic Philosophy, M. Bunge, Ed. Dordrecht: Springer Netherlands, 1983, pp. 157–193. [Online]. Available: https://doi.org/10.1007/978-94-015-6921-7_4
- [16] M. Foucault, *The Archaeology of Knowledge*. New York, NY: Pantheon Books, 1982, oCLC: 254102097.
- [17] A. L. Berez-Kroeker, L. Gawne, S. S. Kung, B. F. Kelly, T. Heston, G. Holton, P. Pulsifer, D. I. Beaver, S. Chelliah, S. Dubinsky, R. P. Meier, N. Thieberger, K. Rice, and A. C. Woodbury, “Reproducible research in linguistics: A position statement on data citation and attribution in our field,” *Linguistics*, vol. 56, no. 1, pp. 1–18, 2018. [Online]. Available: <https://doi.org/10.1515/ling-2017-0032>
- [18] “Giving software its due,” *Nature Methods*, vol. 16, no. 3, pp. 207–207, Mar. 2019. [Online]. Available: <https://doi.org/10.1038/s41592-019-0350-x>
- [19] H. Cousijn, A. Kenall, E. Ganley, M. Harrison, D. Kernohan, T. Lemberger, F. Murphy, P. Polischuk, S. Taylor, M. Martone, and T. Clark, “A data citation roadmap for scientific publishers,” *Scientific Data*, vol. 5, p. 180259, Nov. 2018. [Online]. Available: <https://doi.org/10.1038/sdata.2018.259>
- [20] R. D. Peng, “Reproducible Research in Computational Science,” *Science*, vol. 334, no. 6060, pp. 1226–1227, Dec. 2011. [Online]. Available: <https://doi.org/10.1126/science.1213847>

- [21] J. Howison and J. Bullard, "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*, vol. 67, no. 9, pp. 2137–2155, 2016. [Online]. Available: <https://doi.org/10.1002/asi.23538>
- [22] K. Li, X. Lin, and J. Greenberg, "Software citation, reuse and metadata considerations: An exploratory study examining LAMMPS," *Proceedings of the Association for Information Science and Technology*, vol. 53, no. 1, pp. 1–10, 2016. [Online]. Available: <https://doi.org/10.1002/pra2.2016.14505301072>
- [23] K. Li, E. Yan, and Y. Feng, "How is R cited in research outputs? Structure, impacts, and citation standard," *Journal of Informetrics*, vol. 11, no. 4, pp. 989–1002, Nov. 2017. [Online]. Available: <https://doi.org/10.1016/j.joi.2017.08.003>
- [24] H. Park and D. Wolfram, "Research software citation in the Data Citation Index: Current practices and implications for research software sharing and reuse," *Journal of Informetrics*, vol. 13, no. 2, pp. 574–582, May 2019. [Online]. Available: <https://doi.org/10.1016/j.joi.2019.03.005>
- [25] C. Boettiger, "Citing R packages," Mar. 2012. [Online]. Available: <http://web.archive.org/web/20190823120400/https://www.carlboettiger.info/2012/03/20/citing-r-packages.html>
- [26] Y. Halchenko, Matteo Visconti di Oleggio Castello, Jason Gors, M. Szczepanik, P. R. Raamana, emirvine, C. Barnes, C. Markiewicz, J. Wilk, O. F. Gulban, O. Beckstein, L. Estève, K. Leinweber, and D. Völgyes, "Duecredit/duecredit 0.7.0," Zenodo, Aug. 2019. [Online]. Available: <https://doi.org/10.5281/zenodo.3376261>
- [27] GitHub, "Making Your Code Citable . GitHub Guides," Oct. 2016. [Online]. Available: <http://web.archive.org/web/20190823123222/https://guides.github.com/activities/citable-code/>
- [28] L. Shamir, J. F. Wallin, A. Allen, B. Berriman, P. Teuben, R. J. Nemiroff, J. Mink, R. J. Hanisch, and K. DuPrie, "Practices in source code sharing in astrophysics," *Astronomy and Computing*, vol. 1, pp. 54–58, Feb. 2013. [Online]. Available: <https://doi.org/10.1016/j.ascom.2013.04.001>
- [29] S. Bönisch, M. Brickenstein, H. Chrapary, G.-M. Greuel, and W. Sperber, "swMATH – A New Information Service for Mathematical Software," in *Intelligent Computer Mathematics*, ser. Lecture Notes in Computer Science, J. Carette, D. Aspinall, C. Lange, P. Sojka, and W. Windsteiger, Eds. Springer Berlin Heidelberg, 2013, pp. 369–373.
- [30] K. C. London, S. N. D. Service, F. S. S. D. Archive, University of Tartu, U. P. Fabra, and Centerdata, *TERESA: Tools E-Registry for E-Social Sciences And Humanities*. King's College London, Swedish National Data Service, Finnish Social Science Data Archive, University of Tartu, Universitat Pompeu Fabra and Centerdata, 2014. [Online]. Available: <https://github.com/DASISH/TERESA>
- [31] S. Druskat, J. H. Spaaks, N. Chue Hong, R. Haines, and J. Baker, "Citation File Format (CFF)," Aug. 2018. [Online]. Available: <https://doi.org/10.5281/zenodo.1003149>
- [32] M. B. Jones, C. Boettiger, A. C. Mayes, A. Smith, P. Slaughter, K. Niemeyer, Y. Gil, M. Fenner, K. Nowak, M. Hahnel, L. Coy, A. Allen, M. Crosas, A. Sands, N. C. Hong, P. Cruse, D. Katz, and C. Goble, *CodeMeta: An Exchange Schema for Software Metadata. Version 2.0*, 2017, published: KNB Data Repository. [Online]. Available: <https://doi.org/10.5063/schema/codemeta-2.0>
- [33] T. Gomez-Diaz and T. Recio, "On the evaluation of research software: The CDUR procedure," *F1000Research*, vol. 8, p. 1353, Aug. 2019. [Online]. Available: <https://doi.org/10.12688/f1000research.19994.1>
- [34] D. F. (DFG), "Leitlinien zur Sicherung guter wissenschaftlicher Praxis (Kodex) [Guidelines for Safeguarding Good Scientific Practice (Code)]," p. 32, Aug. 2019. [Online]. Available: http://web.archive.org/web/20190903173540/https://www.dfg.de/download/pdf/foerderung/rechtliche_rahmenbedingungen/gute_wissenschaftliche_praxis/kodex_gwp.pdf
- [35] B. Hanson, A. Sugden, and B. Alberts, "Making Data Maximally Available," *Science*, vol. 331, no. 6018, pp. 649–649, Feb. 2011. [Online]. Available: <https://doi.org/10.1126/science.1203354>
- [36] Y. AlNoamany and J. A. Borghi, "Towards computational reproducibility: Researcher perspectives on the use and sharing of software," *PeerJ Computer Science*, vol. 4, p. e163, Sep. 2018. [Online]. Available: <https://doi.org/10.7717/peerj-cs.163>
- [37] V. Stodden, M. McNutt, D. H. Bailey, E. Deelman, Y. Gil, B. Hanson, M. A. Heroux, J. P. A. Ioannidis, and M. Tauber, "Enhancing reproducibility for computational methods," *Science*, vol. 354, no. 6317, pp. 1240–1241, Dec. 2016. [Online]. Available: <https://doi.org/10.1126/science.aah6168>
- [38] The Yale Law School Round Table on Data and Core Sharing, "Reproducible Research," *Computing in Science Engineering*, vol. 12, no. 5, pp. 8–13, Sep. 2010. [Online]. Available: <https://doi.org/10.1109/MCSE.2010.113>
- [39] L. H. Nielsen, "Zenodo now supports DOI versioning!" May 2017. [Online]. Available: <http://web.archive.org/web/20190905092756/https://www.openaire.eu/blogs/zenodo-now-supports-doi-versioning-1>
- [40] W. H. K. Chun, "On Software, or the Persistence of Visual Knowledge," *Grey Room*, vol. 18, pp. 26–51, Jan. 2005. [Online]. Available: <https://doi.org/10.1162/1526381043320741>
- [41] D. S. Katz, K. E. Niemeyer, A. M. Smith, W. L. Anderson, C. Boettiger, K. Hinsén, R. Hooft, M. Hucka, A. Lee, F. Löffler, T. Pollard, and F. Rios, "Software vs. data in the context of citation," *PeerJ Inc.*, Tech. Rep. e2630v1, Dec. 2016. [Online]. Available: <https://doi.org/10.7287/peerj.preprints.2630v1>
- [42] International Committee of Medical Journal Editors, "ICMJE Recommendations: Defining the Role of Authors and Contributors." [Online]. Available: <http://web.archive.org/web/20190905080833/http://www.icmje.org/recommendations/browse/roles-and-responsibilities/defining-the-role-of-authors-and-contributors.html#two>
- [43] A. Brand, L. Allen, M. Altman, M. Hlava, and J. Scott, "Beyond authorship: Attribution, contribution, collaboration, and credit," *Learned Publishing*, vol. 28, no. 2, pp. 151–155, 2015. [Online]. Available: <https://doi.org/10.1087/20150211>
- [44] M. K. McNutt, M. Bradford, J. M. Drazen, B. Hanson, B. Howard, K. H. Jamieson, V. Kiermer, E. Marcus, B. K. Pope, R. Schekman, S. Swaminathan, P. J. Stang, and I. M. Verma, "Transparency in authors' contributions and responsibilities to promote integrity in scientific publication," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 115, no. 11, p. 2557, Mar. 2018. [Online]. Available: <https://doi.org/10.1073/pnas.1715374115>
- [45] J. Cheng and J. L. C. Guo, "Activity-Based Analysis of Open Source Software Contributors: Roles and Dynamics," *arXiv:1903.05277 [cs]*, Mar. 2019. [Online]. Available: <http://arxiv.org/abs/1903.05277>
- [46] E. Constantinou and G. M. Kapitsaki, "Developers Expertise and Roles on Software Technologies," in *2016 23rd Asia-Pacific Software Engineering Conference (APSEC)*, Dec. 2016, pp. 365–368. [Online]. Available: <https://doi.org/10.1109/APSEC.2016.061>
- [47] "Recognize all contributors." [Online]. Available: <http://web.archive.org/web/20190906132859/https://allcontributors.org/>
- [48] "Contributor - schema.org." [Online]. Available: <http://web.archive.org/web/20190906141516/https://schema.org/contributor>
- [49] J. Beck, "NISO Z39.96 The Journal Article Tag Suite (JATS): What Happened to the NLM DTDs?" *The Journal of Electronic Publishing*, vol. 14, no. 1, Aug. 2011. [Online]. Available: <https://doi.org/10.3998/3336451.0014.106>
- [50] "Journal of Open Research Software." [Online]. Available: <http://openresearchsoftware.metajnl.com/>
- [51] A. M. Smith, K. E. Niemeyer, D. S. Katz, L. A. Barba, G. Githinji, M. Gymrek, K. D. Huff, C. R. Madan, A. C. Mayes, K. M. Moerman, P. Prins, K. Ram, A. Rokem, T. K. Teal, R. V. Guimera, and J. T. Vanderplas, "Journal of Open Source Software (JOSS): Design and first-year review," *PeerJ Computer Science*, vol. 4, p. e147, Feb. 2018. [Online]. Available: <https://doi.org/10.7717/peerj-cs.147>
- [52] J.-F. Abramatic, R. Di Cosmo, and S. Zacchiroli, "Building the Universal Archive of Source Code," *Commun. ACM*, vol. 61, no. 10, pp. 29–31, Sep. 2018. [Online]. Available: <https://doi.org/10.1145/3183558>
- [53] D. Shotton, A. Dutton, S. Peroni, and T. Gray, "Setting our bibliographic references free: Towards open citation data," *Journal of Documentation*, vol. 71, no. 2, pp. 253–277, Feb. 2015. [Online]. Available: <https://doi.org/10.1108/JD-12-2013-0166>
- [54] D. S. Katz and A. M. Smith, "Transitive Credit and JSON-LD," *Journal of Open Research Software*, vol. 3, no. 1, p. e7, Nov. 2015. [Online]. Available: <https://doi.org/10.5334/jors.by>
- [55] S. Ahalt, T. Carsey, A. Couch, R. Hooper, L. Ibanez, R. Idaszak, M. B. Jones, J. Lin, and E. Robinson, "NSF Workshop on Supporting Scientific Discovery through Norms and Practices for

Software and Data Citation and Attribution," Tech. Rep. [Online]. Available: http://web.archive.org/web/20190911211016/https://softwaredatacitation.renci.org/Workshop%20Report/SoftwareDataCitation_workshop_report_2015_April_20_with_logo.pdf



Stephan Druskat holds an MA in English, Modern German Literature and Linguistics from the Free University of Berlin, Germany. He is a Research Software Engineer, working in linguistics, and a PhD candidate in Software Engineering at the German Aerospace Center (DLR) and the Computer Science Department at Humboldt-Universität zu Berlin in Berlin, Germany. In his work, he focuses on research software sustainability and software citation. He is a Special Collaborator of the Software Sustainability Institute (UK), and a board member of de-RSE e.V. - Society for Research Software (Germany).