

Data Citation and the Citation Graph

Peter Buneman¹, Dennis Dosso², Matteo Lissandrini³, and Gianmaria Silvello²

¹University of Edinburgh

²University of Padua

³Aalborg University

Abstract

The *citation graph* is a computational artifact that is widely used to represent the domain of published literature. It represents connections between published works, such as citations and authorship. Among other things, the graph supports the computation of bibliometric measures such as h-indexes and impact factors. There is now an increasing demand that we should treat the publication of data in the same way that we treat conventional publications. In particular, we should cite data for the same reasons that we cite other publications.

In this paper we discuss what is needed for the citation graph to represent data citation. We identify two challenges: (i) to model the evolution of credit appropriately (through references) over time and (ii) to model data citation not only to a dataset treated as a single object but also to parts of it. We describe an extension of the current citation graph model that addresses these challenges. It is built on two central concepts: citable units and reference subsumption. We discuss how this extension would enable data citation to be represented within the citation graph and how it allows for improvements in current practices for bibliometric computations both for scientific publications and for data.

Keywords: *data citation, bibliometrics, citation graph.*

1 Introduction

Citations and the Citation Graph

Citation is essential to the creation and propagation of knowledge and is a well-understood part of scholarship and scientific publishing. Citations allow us to identify the cited material, retrieve it, give credit to its creator, date it, and to provide partial knowledge of its subject and quality.

The *citation graph*, or citation network, is a model used to describe how citations link research entities, typically papers, journals, and books [Harzing and Van der Wal, 2008, Tang et al., 2008b]. It enables a number of important activities such as:

- *Exploration of the graph* to find publications of interest.
- *Tracking of authorship* of papers: citing and following citations is one way to attribute *credit* to authors and to *keep up-to-date* with the work of others.
- *Dissemination* of research findings: the exploration of citations and cited authors enables the dispersed communities of researchers to share their findings and engage in discussions.
- *Computation of bibliometrics* for the analysis of one researcher, venue, or publication impact in particular fields. The citation graph is the basis for nearly all the currently used bibliometrics, such as *impact factor* and *h-index*.

37 Throughout this paper, we refer to an idealized “citation graph” as though it were a real and unique dig-
38 ital artifact that represents papers and the citations between them. Of course, it is not unique: various
39 organizations have distinct implementations of it. Among these we count: Google Scholar, the Microsoft
40 Academic Graph (MAG), ¹ the Open Academic Graph (OAG) [Tang et al., 2008a], Semantic Scholar (SS) ²,
41 AMiner (AM) ³ and PubMed⁴ (this is a more a linked collection of documents than a full-fledged citation
42 graph), Scopus, ⁵ and the Web of Science⁶. These graphs differ in many aspects such as their coverage,
43 their being open- or closed-access, and their schema; but in all of these, the basic structure is a *directed*
44 graph, in which the vertices represent publications, and the edges represent citations from one publication
45 to another [Price, 1965].

46 Most of the information about papers is contained in annotations of the nodes. The edges are generally
47 typed but not annotated (an exception is MAG, which carries *context*, as we discuss later). While in early
48 models, nodes only represented papers and the only edges were "cites" edges, recently, citation graphs have
49 been extended with richer information [Peroni and Shotton, 2020]. These extensions may carry author nodes
50 with a "wrote" edge to papers, journal/conference nodes with a "part of" edge from papers, and subject
51 nodes with the corresponding edges. While representations differ, the purpose is similar: to provide the
52 services described above.

53 The need for data citation

54

55 Scientific publications increasingly rely on curated databases, which are numerous, “populated and updated
56 with a great deal of human effort” [Buneman et al., 2008], and at the core of current scientific research ⁷.
57 In this context, references to data are starting to be placed alongside traditional references. Hence, there
58 has been a strong demand [COD, 2013, FORCE-11, 2014] to give databases the same scholarly status of
59 traditional scientific works and to define a shared methodology to cite data. Scientific publishers (e.g.,
60 Elsevier, PLoS, Springer, Nature) took upon data citation by instituting policies to include data citations
61 in the reference lists.

62 The *open research culture* [Nosek et al., 2015] is based on methods and tools to share, discover, and access
63 experimental data. Moreover papers, journals, and articles should provide access to all the data that they
64 use [Cousijn et al., 2019]. Researchers and practitioners (e.g., journalists and data scientists) who make
65 use of electronic data should be able to cite the relevant data as they would cite a document from which
66 they had extracted information [Cousijn et al., 2017, Nature Physics Editorial, 2016]. As we shall see, the
67 citation graphs can become a fundamental tool in the pursuit of the goal of accessibility and networking
68 between papers and data.

69 We also observe that data occupy a crucial role today in research, emerging as a driving instrument in sci-
70 ence [Candela et al., 2015]. Data citations should be given the same scholarly status of traditional citations
71 and contribute to bibliometrics indicators [Belter, 2014, Peters et al., 2016]. Principles such as Findability,
72 Accessibility, Interoperability and Reusability (FAIR) [Wilkinson et al., 2016] require data to be easily find-
73 able and accessible, qualities that are more readily available once data can be appropriately cited. In this
74 sense, we can say that the FAIR principles encourage the adoption of data citation.

75 The reasons given for data citation are the same those given for a conventional citation [COD, 2013]: recog-
76 nition of the source (e.g., a title); credit for the author, curator, or agent; establishment of its currency
77 (when it was created); where it was located; and how it was extracted. The last three of these fall under the
78 general heading of provenance and are important when one wants to reproduce some analysis on the data
79 or establish the trustworthiness of a claim.

80 Data sets and databases are usually more complex and varied than textual documents, and they introduce

¹<https://www.microsoft.com/en-us/research/project/microsoft-academic-graph/>

²<https://www.semanticscholar.org/>

³<https://www.aminer.cn/>

⁴<https://www.ncbi.nlm.nih.gov/pubmed/>

⁵<https://www.scopus.com/home.uri>

⁶<https://clarivate.com/webofsciencegroup/solutions/web-of-science/>

⁷See <https://fairsharing.org/databases/> for a detailed list of curated scientific databases commonly used in research.

81 significant challenges for citation [Silvello, 2018]. Text publications have a fixed form, do not change over
82 time, are interpretable as independent units, share a standard format and representation model, and are
83 composed of predetermined, albeit domain-dependent, sets of elements that are considered as *citable*, e.g.,
84 the whole paper or book or a chapter. Scientific databases are structured according to diverse data models
85 and accessed with a variety of query languages. What can be cited may range from a single datum to data
86 subsets or aggregations specified by the person or agent that extracts the relevant data, and deciding *a priori*
87 what can and cannot be cited is rarely feasible. Data citation introduces multiple citation types, besides the
88 classical papers citing papers. These are (i) papers citing data; (ii) data citing papers; and (iii) data citing
89 data.

90 Data Citation in the Citation Graph

91 Our purpose in this paper is to discuss whether, in its current form, the current model of the citation
92 graph can properly accommodate data citation. We claim that, despite all the features and modifications
93 that have been added to various implementations of the citation graph, at least two significant features
94 are generally missing or poorly represented. These shortcomings already limit what we can represent with
95 existing implementations, and we argue that they make impossible the proper representation of data citations.

96 The first shortcoming concerns the *assignment of credit* when a referenced scientific work is corrected or
97 augmented with another version. A typical example case is that of a preprint paper that gets cited before
98 its peer-reviewed version is published. It is common for the authors to prefer that the preprint citations are
99 merged with those of the peer-reviewed version. Something similar happens also when an updated version
100 of a dataset is published.

101 In the case of data, we need to consider that a database may be composed of multiple independently citable
102 parts (e.g., a single record, a table, a view). Every single citable part can evolve and change over time and
103 obtain citations (also views or downloads, when monitoring other scientometrics signals) at a different point
104 in time. Therefore, it can be necessary to aggregate these statistics over all the versions of the same part to
105 measure its impact and that of the database. The MAG and S2ORC databases have also an explicit notion
106 of multiple versions of a paper, for example preprints and final published versions. It is however uncommon
107 to “move” citations from one version to another, following some criteria or algorithm to correctly allocate
108 citations. Yet, aggregating citation to a single version of a scientific work would have, among other things,
109 the desirable effect of allowing proper evaluation of the impact of the work.

110 The second feature is the representation of *context* of a citation. Context is required for various reasons.
111 It is typically used to describe the relevant part (e.g., page number) of a *cited* document. It may also
112 carry, as in MAG, the surrounding text within the *citing* document helping to understand the reason for the
113 citation; e.g., a simple mention, a confutation, or a validation, such as those described in the OpenCitations
114 ontology [Daquino et al., 2020]. In the case of data citations, the context can contain the query identifying
115 the cited data, expressed in different format (e.g. an URL, a file name, a SQL or SPARQL query, etc.).
116 Despite a great deal of attention dedicated to the citation context – see, for instance, the Citation Context
117 Analysis (CCA) discussed as early as the 1980s [Freeman et al., 2013] – there is no systematic approach to
118 representing it within citation graphs.

119 In fact, none of the largest citation-based systems, such as Scopus, MAG, and Google Scholar, properly
120 take into account scientific databases as objects for use in the research literature. Google Data Search⁸
121 allows us to search for indexed datasets, but it does not keep track of the citations to data or other types of
122 statistics, like clicks or downloads. Web of Science is one notable exception since it models data citations,
123 even though only at the database level, via the Data Citation Index (DCI) now maintained by Clarivate
124 Analytics [Force et al., 2016]. Note that DCI is not publicly available and the datasets are indexed after a
125 validation process.

126 Another effort is the Scholix framework [Burton et al., 2017], which can be regarded as a set of guidelines
127 and lightweight models that can be quickly adopted and expanded to facilitate interoperability among link
128 providers. Finally, an example of an initiative that includes data and databases among the entities of the
129 graph is the OpenAIRE Research Graph Data Model [Manghi et al., 2019], which leverages the OpenAIRE

⁸<https://datasetsearch.research.google.com/>

130 services to populate a *research graph* whose nodes include scientific results, organizations, funding agencies,
131 communities, and data sources.

132 The conventional approach is to treat a dataset as a single entity, in the same way, one would treat a scientific
133 publication. However, this is far from ideal since typically only a small part of the dataset or database is
134 cited, and the authorship – the people who have contributed to the database – can vary widely with the part
135 of the database being cited [Buneman et al., 2016].

136 In this paper, we discuss the extension of the current model to enable the proper inclusion of data citations
137 in the citation graph; and we discuss the *evolution* of a database: what happens to citations when new
138 versions of the database appear? For the versioning issue, we describe a relation between scientific works
139 (either papers or data) called *subsumption*. Through different policies, this relationship models effectively
140 how credit should be transferred through time when updated versions of data appear in the graph. Finally,
141 we discuss how to introduce data in the citation graph, considering the most common data citation strategies
142 currently used in the world of research. In particular, we take inspiration from one of the solutions proposed
143 by the *Research Data Alliance (RDA)*⁹. The RDA is a community-driven initiative launched in 2013 by
144 different commissions. One of its working group, the “Working Group on Data Citation: Making Dynamic
145 Data Citable” (WGDC), has as one of its goal the identification and citation of arbitrary views of data. As
146 potential solution, the WGDC recommends an identification method based on PIDs assigned to queries.

147 The focus of this work is on data citation; but to ease the comprehension of the paper, we first discuss the
148 limitations of the citation graph and the possible extensions we propose by focusing on textual documents
149 and then we extend the reasoning to data citation.

150 The paper is organized as follows: Section 2 describes some preliminary concepts and the limits of the citation
151 graphs; in Section 3 we discuss the proposed solutions for the first three issues; Section 4 presents the proposed
152 solution for the introduction of data in the citation graph; Section 5 sums up our main proposals and discusses
153 possible lines of research and development; Section 6 describes the related work; finally, Section 7 presents
154 conclusions and future work.

155 2 The Citation Graph: Concepts and Limits

156 2.1 Core concepts

157 **Citable Unit.** By citable unit (CU), we mean a published entity – be it a paper, a chapter, or portion of
158 data – which presents all the qualities necessary to be considered as a “citable work”. The characterization
159 of a CU that we use, given in [Wilke, 2015], requires that: (i) it has to be uniquely and unambiguously
160 identifiable and citable; (ii) it has to be *available* in perpetuity and in *unchanged* form; (iii) it has to be
161 *accessible*; (iv) it has to be *self-contained* and *complete*. Self-contained and complete means that whatever
162 new contribution is contained inside the piece of work, that contribution needs to be fully and clearly
163 explained. This is not always the case for certain publications. Consider in fact the slides of a scientific
164 presentation. As they are used merely as a support for the oral presentation, they often cannot be fully
165 understood without the corresponding talk. Also, the combination slides/ registration of the talk may be
166 incomplete, as many presenters tend to skip technical details during their presentations, referring to the
167 complete published work.

168 While some of these requirements are subjective, and not straightforward in databases, they still provide
169 a workable starting point. The requirement that is most problematic for databases is that the citable unit
170 has to be *unchanged*. Databases evolve rapidly, and creating a citable unit for each version may be counter-
171 productive. This is something we address in Section 4.2. Generally, what constitutes a citable unit is decided
172 by convention. We should also note that some citable units comprise other citable units. The proceedings of
173 a conference may be cited as may be a book on a topic whose chapters are written by different people and
174 may also be individually cited. There is thus a “part-of” relationship between CUs that we discuss later.

⁹<https://www.rd-alliance.org/>

175 In [Daquino et al., 2020] a similar concept, *bibliographic resource*, is defined as a resource that cites and can
176 be cited by other resources.

177 **Reference.** At the end of this paper, there is a list of references. Traditionally, a reference is a *pointer* to,
178 and a brief description of, another publication in the literature. It is a short text composed of fields such as
179 title, authors, year, venue and others, that enables us to identify and find the entity – i.e., a paper, a book,
180 or a survey – being referenced. Depending on aspects of the citing CU’s nature, like its field of research, the
181 publication venue, or even language, different attributes of the reference may vary such as the format or the
182 fields composing the reference. In physics, for example, titles are often omitted.

183 The important point is that, apart from the stylistic rendition of the reference, its contents are determined
184 by the cited CU, hence, to within stylistic variations, the reference to a CU will be the same in any paper.
185 In this paper, the reference determines the existence of a directed edge between two CUs, the citing and the
186 cited one.

187 **Citation.** There is no universal agreement on the distinction between *reference* and *citation*, and the two
188 terms are often used interchangeably [Price and Richardson, 2008, Altman and Crosas, 2014, Osareh, 1996,
189 Daquino et al., 2020].

190 One distinction proposed in [Gilbert and Woolgar, 1974] is that “reference” refers to the works mentioned in
191 the reference section or bibliography of a paper. A reference may be mentioned once or many times in an
192 article. Each of these mentions is considered a citation.

193 The distinction is crucial to our understanding of the citation graph. If we look at what goes in the body of a
194 paper, we may find, for example, “Austen, J. (2004). pp 101-104”. We note that this textual artifact contains
195 two parts. The first one is “Austen, J. (2004).”, which we call a **reference pointer**. A reference pointer is,
196 in general, a textual means that is used to denote a single bibliographic reference in the reference section
197 when mentioned in the body of a paper. The second part of the citation is composed of some additional
198 information, in this case “pp 101-104”, that may help the reader locate specific information within the cited
199 paper. Note that the same reference pointer can occur several times in a paper and may have differing
200 additional information, such as pp 10-25 and pp 110-120.

201 Therefore, we can say that a **citation** is composed of the combination of the reference pointer with the
202 (optional) information added to it in the paper’s body. The optional information in the paper’s body may
203 be referred to as a form of *context* for the citation. This implies that there is a many-one relationship between
204 citations and references, a fact that is supported by some discussions on the topic, for example “... every
205 citation should also have a corresponding entry in your reference list” [nzb, 2020].

206 **Reference Annotation.** We shall call this extra information such as “pp 101-104” *reference annotation*.
207 In this paper, the reference annotation consists of all the information added to a reference pointer to qualify
208 how it is used. This information is not part of the reference and can change depending on how that particular
209 resource is used.

210 The Citation Typing Ontology [Shotton, 2010] is replete with examples of other kinds of annotations such
211 as “refutes,” or “ridicules”, which are clearly about the relationship between the citing and cited docu-
212 ments. In the Microsoft Academic Graph [Sinha et al., 2015], the *context* – the text surrounding a citation
213 in the source document – may be recorded as another form of annotation. The OpenCitations ontol-
214 ogy [Daquino et al., 2020] contains a class called annotation¹⁰ attached to the in-text citation and to a
215 reference which has a similar role. Here, we do not need to distinguish between the context of a reference
216 pointer and its reference annotation – i.e., for our purposes these two concepts are the same, however it may
217 be that certain applications will require some finer distinctions.

218 These definitions differ slightly from those in [Daquino et al., 2020, Daquino et al., 2018], where a reference
219 (called bibliographic reference) and a reference pointer are *manifestations* of a citation. Moreover, in our

¹⁰<http://www.w3.org/ns/oa#Annotation>

220 example, the part “pp. 101-104” is a reference annotation, whereas in [Daquino et al., 2020] it is a *specializa-*
221 *tion* of the citation. We do not specifically model the concept of specialization, since it can be inferred from
222 the content of the reference annotation. Also, in [Daquino et al., 2020] the pointer may include additional
223 information, but the citation does not.

224 Summing up, we consider a reference annotation as a “box” that can contain information derived from the
225 context of a reference pointer.

226 Generally speaking, the Citation Context Analysis (CCA), whose basis was first developed in the early
227 1980s, is the syntactic and semantic analysis of citation content, used to analyze the context of research
228 behavior [Freeman et al., 2013]. CCA has been used as a promising addition to traditional quantitative
229 citation analysis methods. One of the main aspects of CCA is that it incorporates qualitative factors, such
230 as how one cites. In [Daquino et al., 2020] this idea is captured by the concept of **citation function**, i.e.
231 the function or purpose of the citation (e.g. to cite as background, extend, agree with the cited entity, etc.)
232 to which each in-text reference pointer relates. In our proposal, this qualitative factor, or citation function,
233 can be located in the reference annotation, and it could be inferred from the context of the reference pointer.

234 Even in a citation graph that represents conventional citations it is necessary to be able to attach information
235 to a reference to create proper citations. Yet, in some citation graph implementations, this is impossible,
236 because the reference relationship is represented as a directed, but unannotated edge. As noted above,
237 an exception is the Microsoft Academic Graph, which contains two kinds of edges between publications:
238 unannotated edges and edges annotated with context. The reason for this omission may be the difficulty
239 of collecting the relevant information; it may also be that it is not needed in the computation of most
240 bibliometrics.

241 **Part-of.** The *part-of* relationship exists between two citable units in the graph; it describes the situation
242 where one citation unit is somehow “contained” in the other. This is the case of papers published in an
243 instance of a venue (e.g., the 2020 version of the ACM SIGMOD), and these issues being part of the venues
244 themselves (e.g., ACM SIGMOD). This information is present for example in databases such as MAG and
245 AMiner.

246 In the case of data, the part-of relationship is particularly important. Many databases and datasets have a
247 hierarchical structure and may be cited at different levels of detail.

248 **Database categories and citation.** There is a broad spectrum of databases for which citation is appro-
249 priate. In discussing data citation it is helpful to divide them into three rough categories.

- 250 • *Static databases*, which are used to support claims in a publication. These are typically “one-off”
251 results of a set of experiments. For these databases, systems such as Mendeley ¹¹ store data alongside
252 the publication, so that a citation to the publication also serves as a citation to the data. Data
253 journals [Candela et al., 2015], i.e., journals publishing papers describing data sets, are also employed
254 as proxies to cite static data sets.
- 255 • *Evolving databases* of source data such as *weather data* [Philipp et al., 2010] or *satellite image data*
256 [Shanableh et al., 2019] that are collected for a wide range of purposes. Zenodo ¹², like Mendeley, stores
257 data together with its representative publication. However, a publication about a data set and the data
258 set itself can also have separate and unrelated DOIs. In this case the citation to the publication and to
259 the database are distinguished. Moreover, it allows to deposit multiple versions of the same database,
260 with new DOIs for each one, thus to keep track of usage stats like the number of downloads and views
261 on each version. A citation to the database, or even to a document that describes the whole database,
262 is generally regarded as inadequate. Usually, only a portion is used; hence, one needs to know the part
263 (the sensor, the location of the image, or the time range) from which the data was extracted.

¹¹<https://www.mendeley.com/>

¹²<https://zenodo.org/>

264 • Finally, we have *curated databases*. These have largely replaced conventional biological reference
265 works [Buneman et al., 2008], and like the works they replace, involve substantial human effort. One
266 advantage is that they are readily accessible and easy to search. Moreover, there are few limits on
267 their size and complexity, and they can evolve rapidly with the subject matter. For these, the citation
268 is a complex issue but is just as crucial for curated databases as it is for the reference works that they
269 replace.

270 The distinction between these three categories is not sharp, and there are many examples that lie in the
271 overlap. For example most source data databases involve a degree of curation.

272 2.2 Existing Limitations of the Citation Graph

273 While implementations of the citation graph differ, the basic model consists of a directed graph $\mathcal{G} = (V, E)$,
274 where V is the set of papers and $E \subseteq V \times V$ is the set of directed edges corresponding to the citations among
275 them: an edge $\langle p_1, p_2 \rangle$ connects the papers p_1 and p_2 , if p_1 cites p_2 . The following limitations of this simple
276 model are obstacles to the representation of data citation, but can already be seen in conventional citations
277 to papers.

278 **Lack of context.** While in the basic model of the citation graph the nodes often contain information such
279 as the *title*, the list of *authors* or the *venue* of publication, it is lacking the information about the *context*
280 of the citation, i.e. all that kind of information that could be inferred from the context of the reference
281 pointers, such as the specialization of the citation or the citation function. The only information provided
282 by the edge $\langle p_1, p_2 \rangle$ is that p_1 cites p_2 , but it does not specify the *why* or the *how* of this citation. In
283 the literature, we find the *contextual citation graphs*, which make apparent the textual contexts of each
284 citation [Lo et al., 2019, Bird et al., 2008, Daquino et al., 2020]. These graphs contain information about
285 reference annotations which is what, in this work, we consider as the citation context.

286 Note that a lack of citation context is not only an issue that is just related to data citation, but to the whole
287 scientific citation infrastructure and ecosystem. How one document is cited in another, whether cited as a
288 piece of evidence or a tool, could greatly influence how the scientific bibliographic universe is built and how
289 credit should be assigned between researchers.

290 **Versions.** Ideally, the papers in the citation graph should only cite papers in the past, i.e., papers that
291 already exist when the new paper is introduced in the graph [Lo et al., 2019]. If this is the case, the citation
292 graph is a DAG (Directed Acyclic Graph).

293 However, this often is not true since some of the papers in V go through revisions and modifications. This
294 happens for many reasons and with many variations. Among the possible cases: it may be that several
295 copies of one work are to be found on the internet; that one version is an “abstract” and is published in some
296 conference proceedings, and a “full version” is later published in some journal; that one version is published
297 in some archive online and then a fully-fledged paper is released in a conference or journal.

298 To receive credit, it is generally in the authors’ interest to have these documents seen as one. What appears
299 to happen in Google Scholar, for example, is that all versions are clustered together, and one of them, the
300 “main” version, is selected to be the recipient of all references.

301 Consider the following situation: document A is published, and a document P citing A is subsequently
302 published. Document B, a revision and possibly an extension of A, is then published, taking A’s place in
303 the graph. If this new version B contains new outgoing citations to P, then a cycle is created, and the graph
304 is no more a DAG ($P \rightarrow A \rightsquigarrow B \rightarrow P$). This problem may be solved by separating A and B.

305 Another source for cycles in citation graphs that cannot be avoided are papers by the same authors created
306 at the same time (e.g., a full paper written together with a demo paper or extended abstract). In this case,
307 the problem can be solved, for example, by conflating the papers.

308 Another problem arises when the system, for some reason, decides that B becomes the “main” representative
309 of the publication. In this case, what happens with services like Google Scholar, is that the references first
310 given to A are rerouted to B. This can be confusing as the reference annotation (e.g., the page number) may
311 no longer be valid.

312 **Citations to data.** One of the primary roles of data citation is to give credit and attribution to the work
313 of data creators and curators [COD, 2013]. If integrated into the citation graph, data citation represented
314 and analyzed in as are conventional citations, with data CUs and corresponding authors receiving citations
315 and thus credit for their work. However, services like Google Scholar or Scopus do not allow databases into
316 their citation graph.

317 Data journals [Candela et al., 2015] enable the publication of papers describing a database that works as a
318 proxy for it and its authors and receives its citations. This is a possible solution, but it is not complete since
319 it does not consider citations referring to *general* queries.

320 To give appropriate credit to the contributors to the various parts of a complex curated database, one
321 approach to data citation [Buneman et al., 2020] is to automatically create short papers, *citation summaries*,
322 for each citable part of the database and publish them in a dedicated on-line journal. This enables the
323 contributors to receive proper bibliometric credit for their contributions to the database. In this approach,
324 a new summary for a view is generated whenever that view changes substantially. This summary can then
325 be included in the current implementations of citation graphs and receive citations.

326 To conclude, unless there is some form of representation of the cited database or the cited query in the form
327 of a paper or journal, current citation graphs do not include databases as nodes and citations to data as
328 edges.

329 3 Extending the Citation Graph

330 We describe two key extensions to the citation graph needed to deal with both the structural complexity
331 and evolution of databases. These extensions already exist in a limited form in some implementations of the
332 citation graph. However, we need to specify them precisely and understand how they help with the limitations
333 described above and with data citations. What we propose is independent of any specific implementation
334 of the citation graph and, for the most part, it can be incorporated as extensions to those implementations
335 rather than requiring a completely new implementation of the supporting database.

336 3.1 Reference Annotation

337 As discussed above, a reference is represented by an edge in the citation graph. However, to represent a
338 citation accurately, we need to add *reference annotations*. That is, we need to annotate the edges. Unfor-
339 tunately, most data models currently implemented do not support data on edges¹³, so for consistency with
340 these models, our diagrams include a new kind of node rather than a new kind of edge.

341 Consider Figure 1. Two papers, P_1 and P_2 , are represented with circular nodes. We use these nodes to
342 represent citable units. They are annotated with all the information that usually constitutes one reference,
343 like title, authors, year of publication, journal name, and DOI.

344 In this example P_1 references P_2 . We can imagine the reference appearing in the “References” section of P_1
345 as something similar to Johansson, L. C. et al. (2019). XFEL structures of the human MT 2 melatonin receptor
346 reveal the basis of subtype selectivity. *Nature*, 569(7755), 289-292. doi: 10.1038/s41586-019-1144-0. The use
347 of this reference in the paper is reflected by the presence of the reference edge between P_1 and P_2 and the
348 reference node `reference_1`. This is a different kind of node, which contains information such as the edge
349 type (reference), the timestamp of when the citation was registered by the system and the type of reference

¹³Property graphs are an exception since they allow data to be assigned to edges.

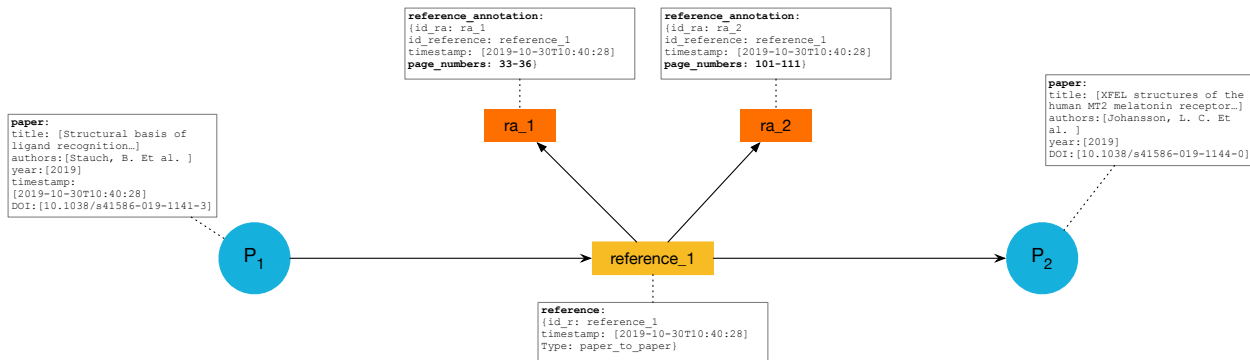


Figure 1: Use of references and reference annotations. Each reference is an edge connecting one citing unit to the cited one, and, if it exists, it is unique. One reference may have one or more reference annotations, each giving rise to a citation.

350 (in this case from a paper to another paper). The actual information contained by the node can be modeled
 351 according to whatever model we decide to follow, e.g., the aforementioned Open Citation ontology.

352 Suppose now that P_1 cites P_2 twice. Each time, it does not merely refer to the whole paper P_2 , but specific
 353 parts of it. The node $reference_1$ has two other neighbor nodes, called *reference annotation nodes*, ra_1
 354 and ra_2 . These two nodes contain the information describing the reference annotations found in P_1 used
 355 to cite P_2 , such as the context, references to particular tables or images, comment on the nature of the
 356 citation (e.g., that the authors of P_1 agree or disagree with P_2). In the example, these annotations carry
 357 page numbers. Hence, the combination of $reference_1$ with ra_1 makes one citation.

358 Reference and reference annotation nodes are the addition that we make to the citation graph to face the
 359 first problem.

360 3.2 Subsumption

361 Often new documents take the place of older versions, becoming also the recipients of both new and old
 362 citations. This behavior is handled behind the scenes by some existing implementations of the citation
 363 graph (notably Google Scholar). To deal with this phenomenon transparently, we propose the introduction
 364 in the citation graph of a new relation, called *subsumes*.

365 In Figure 2 we see a situation similar to the one of Figure 1, where P_1 is citing P_2 at time 1. Now, imagine
 366 that a new version of the same paper, P'_2 , is published and inserted in the citation graph at time 2. The
 367 reference for P'_2 should also have a version number or something that distinguishes it from P_2 . The relation
 368 *subsumes* between P'_2 and P_2 indicates that the former is a new version of the latter, and is, from now on,
 369 *the* paper to consult and reference.

370 In some scientific areas, a journal “paper” P'_2 may be treated as a version of an earlier conference “abstract”
 371 P_2 , even though the two differ substantially. Because of this we do not want to destroy the original link from
 372 P_1 to P_2 ; to do so would be to “rewrite history” and remove information from the graph and we strongly feel
 373 this should not be the case with the citation graph. The *subsumes* relation is present to indicate that one
 374 paper is a version of another and, crucially, that author credit can be transferred from the subsuming paper
 375 to the subsumed paper. On the one hand, the transfer of credit enables a more comprehensive measure
 376 of author contributions (e.g., increasing the number of citations on the latest version of the publication).
 377 On the other hand, credit transfer also transparently reflects the impact that the publication, seen as the
 378 aggregation of its different versions, has on the research community. Different types of subsumption can be
 379 defined. For example, the kind of subsumption that propagates the citations to the single papers to their
 380 journal, thus computing its impact factor.

381 It would be wrong to transfer the credit for writing a paper to more than one other paper, so the subsumption

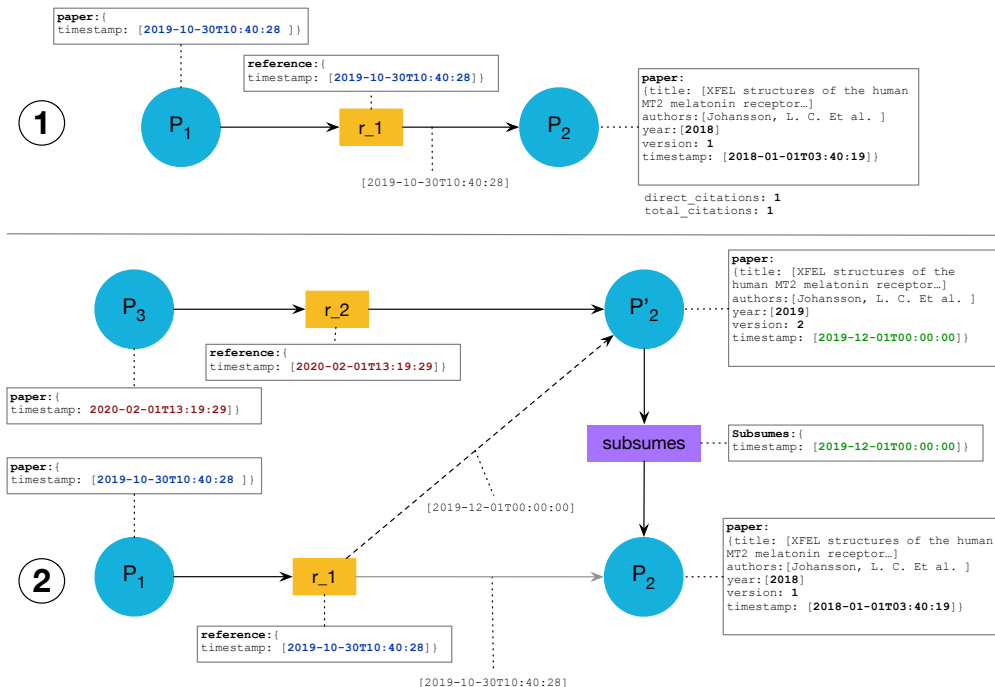


Figure 2: The subsumes relation between two CUs.

382 relation is many-one. It is necessarily acyclic, thus it is a *forest* with the roots of the trees in that forest
 383 being the papers that are designated to receive the credit. It may be useful to have a term for a root node
 384 on the subsumption graph, perhaps *primary citable units* (PCU). It is interesting to note a similar approach
 385 in the MAG¹⁴, which lists the CU P_2 under the PCU P'_2 , keeps the citation count for P_2 and P'_2 distinct,
 386 and reports, for example, “124 citations” for P_2 , “325 citations” for P'_2 but adds, to P'_2 , “449 citations for all”.

387 4 Data in the Citation Graph

388 Here we discuss how we place databases in the citation graph. We shall find that the two extensions we have
 389 discussed – edge annotation and subsumption – are essential to accommodating databases. In particular,
 390 they allow us to deal with databases, which tend to be updated and thus change much more frequently than
 391 papers. We could treat each version or instance of the database as a distinct document, but – at least for
 392 author credit – this would be a limitation, if not counter-productive.

393 First of all, we use the term “database” in the most general sense to refer to a conventional relational database,
 394 an ontology, some form of graph database, or a database that is a collection of files [Buneman et al., 2016].
 395 One might then say that anything one has termed a database is a citable unit. The problem is that *parts* of
 396 the database may also be citable units. The reason we need to discuss parts of the database is twofold: first,
 397 wherein the database one finds something is, like page numbers, a form of location or partial provenance;
 398 the second authorship may vary with what part of the database is being cited [Buneman et al., 2016].

399 With “part” of a database we intend a *view* [Buneman et al., 2016]. A view is a query which we again
 400 generalise to being anything from a relational query for a relational database, a directory path or URI for
 401 a collection of files, or some query in one of the several languages that have been developed for ontologies
 402 and graph databases. It is assumed that the database administrators will define these views and hence
 403 the citable units. MODIS [Justice et al., 1998] is an example of a large evolving database of earth images
 404 for which various subcollections have different authorship; and GtoPdb [Southan et al., 2015] is a complex

¹⁴<https://tinyurl.com/y9clyx8d>, retrieved 16 March 2020

405 curated relational database in which authorship is represented within the database and can be assigned to
406 views determined by the curators.

407 4.1 Part-of and reference annotation

408 Consider, for simplicity, the case in which the database is static, or that we are only interested in representing
409 citations for one version of the database (we address the more complex case of dynamic databases in the
410 next section). The first observation is that by defining the CUs as views, we immediately obtain a part-of
411 relationship: view V_1 is a part of view V_2 if V_1 can be answered from the result of V_2 . Formally, V_1 is part
412 of V_2 if there is a query Q such that for all possible instances of the database, $V_1(D) = Q(V_2(D))$.

413 We have already discussed reference annotations and the information they carry. Among the other things,
414 they contain information of the *where* in the cited document the relevant information being cited is to be
415 found. If we look at data citation, this notion of location has much greater importance. For example,
416 the DataCite schema [Group et al., 2016, Starr and Gastl, 2011] contains the support for the depiction of
417 geospatial data, with properties such as `GeoLocation` and in particular the sub-property `GeoLocationBox`,
418 which specifies a *bounding box*, that is the spatial limits of a box. Most generally we can describe the
419 “location” in the database as a *query* that extracted the relevant information. This is the approach taken
420 in systems that provide accurate provenance [Pröll and Rauber, 2013]. It meshes perfectly with what we
421 are suggesting: the query used to extract the data is a fundamental part of the data citation itself, and the
422 query – or possibly a URL which contains that query – is an essential part of the reference annotation in
423 the citation.

424 Many approaches can now be defined to decide how to introduce the CU corresponding to data in the citation
425 graph. Here we explore two possibilities, stemming from two of the most used strategies in the research world
426 today. We exemplified these two possible strategies in Figure 3.

427 In Figure 3.A we see that a database is represented with a node, DB_1 . A whole database is a citable unit, and
428 every time a paper wants to cite data in that database, it cites the entire database. The reference annotations
429 contain the queries to get the cited data. The paper P_1 presents two citations to DB_1 . Therefore, it has one
430 reference and two reference annotations containing the two different queries being used. P_2 is citing DB_1
431 only once. The total count of citations to DB_1 is two in this case.

432 With this solution DB_1 is the only recipient of citations. This means that its number of citations can
433 become very high. On the other hand, it may happen that the rightful authors and curators of the parts of
434 the database being actually cited do not receive any credit for their work.

435 In Figure 3.B we see the strategy adopted by the RDA. Every time a paper cites a data subset extracted
436 through a new query, a citable unit is created in the citation graph; we represent this CU as a view,
437 corresponding to that query. In this case, P_1 is citing DB_1 twice by using two different queries, thus there
438 are two distinct references, corresponding to the two views being cited, V_1 and V_2 . P_2 citing DB_1 with the
439 query Q_3 , generates another view, i.e., V_3 .

440 With this solution, new views are created every time it is necessary. This can produce an explosion in the
441 number of nodes in the citation graph, many of which receive only one citation. However, in this way it is
442 possible to cite the exact set of data extracted by the query and to give the credit of the citation to the
443 rightful authors of that data. Moreover, the three views of the example are connected to DB_1 by a “part of”
444 relationship. This means that DB_1 may inherit all their citations when needed.

445 We note that to assign only a single CU for the whole database and dispense with the part-of relationship
446 fails when, for example, authorship varies with the part of the database being cited. This is the case with
447 both MODIS and GtoPdb. In this case, we reiterate that it is up to the curators or database administrators
448 to determine the views that define the CUs.

449 In the case of GtoPdb, both the curators and the contributors agree that the PCU should be the data
450 summary [Buneman et al., 2020] for the most recent view of the database. Unfortunately, the PCU is not
451 determined by the curators but by the system that scans the dedicated journal and creates citation graphs.

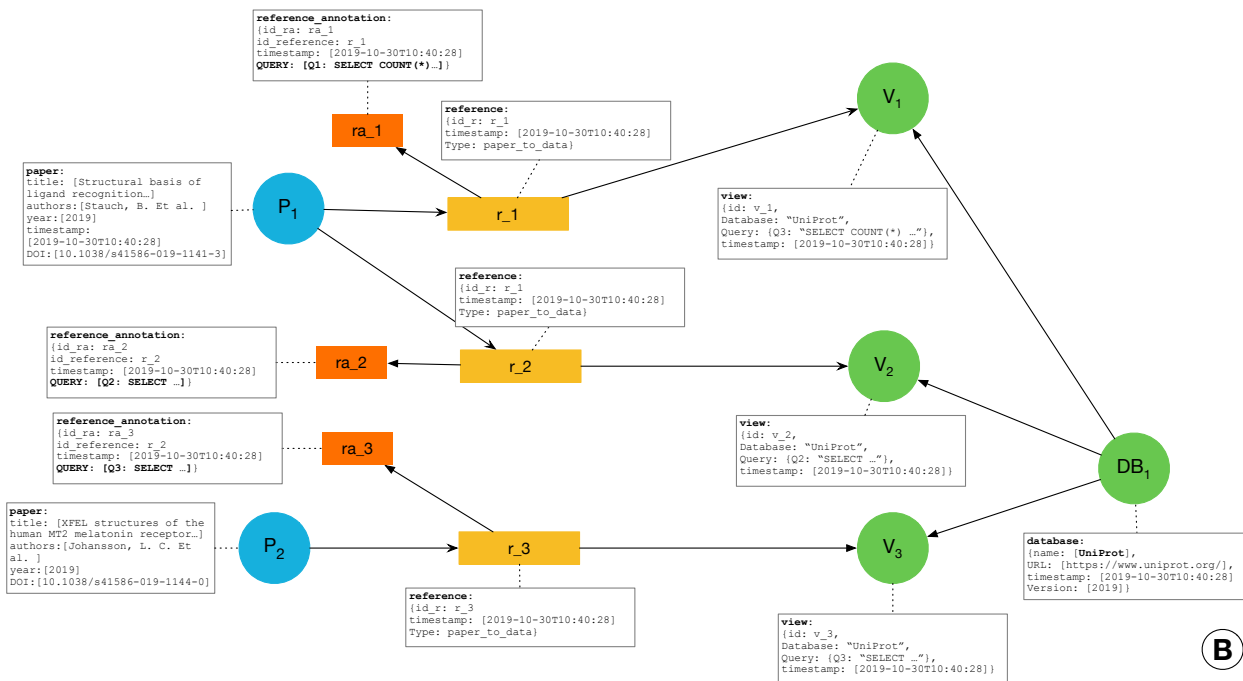
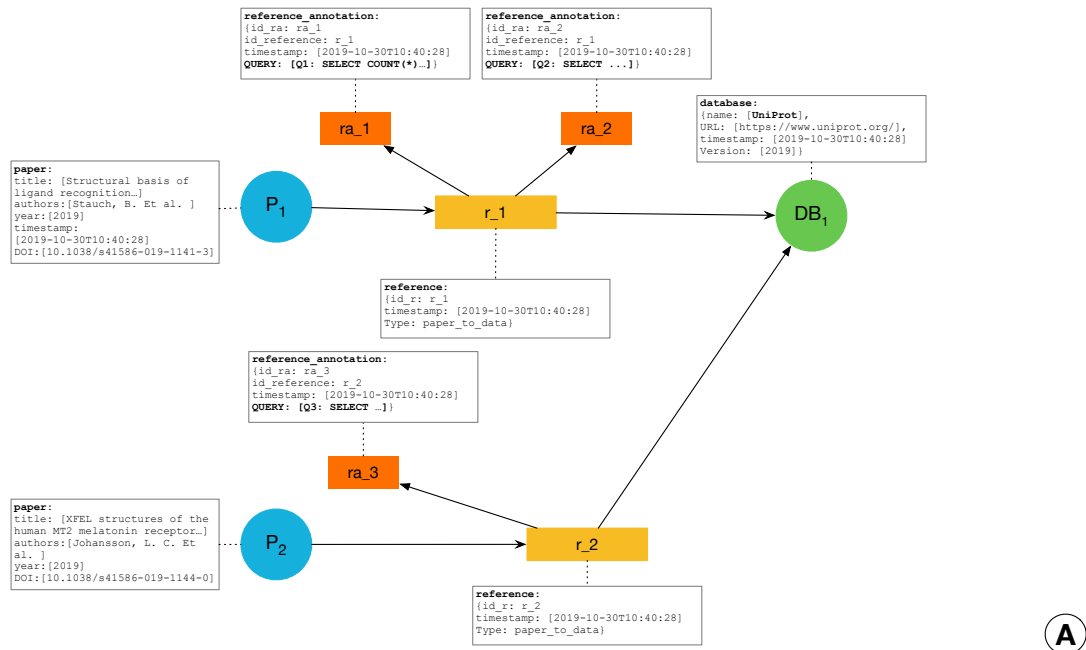


Figure 3: Two examples of possible strategies when citing data. A: always cite the database. B: create a view for every new query issued, if it does not already exist, and cite that view.

452 For a given database, it is the responsibility of the curators or administrators to determine the subsumption
453 relation. Even for conventional publications, we believe that the subsumption relation should be determined
454 by the authors and publishers.

455 4.2 Dealing with dynamic data: Subsumption for data

456 Most databases are not static. Unlike documents they are expected to evolve over time. If versions of a
457 database are released, say, every year, it might be appropriate to treat each version as a new CU. On the
458 other hand, as we discussed, a database in the Citation Graph can present a hierarchy of CUs connected
459 among them through the part-of property. Even though a database may change rapidly, the result of a view,
460 part-of a database, may remain unchanged. The lower a CU is in the part-of hierarchy, the less frequently it
461 will change. Also, even if a part-of CU does change, we may want to treat it as a new *version* of the previous
462 CU rather than an entirely, unrelated, new CU, just as we treat an extended or improved version of a paper.

463 Given these observations about the introduction of dynamic data, it is necessary to answer these questions:

- 464 • When is it necessary to introduce a new CU representing a view?
- 465 • In case a new CU has been introduced, when can it be considered a new version of the previous CU,
466 or an entirely new entity?
- 467 • In case a new CU has been introduced, how do we connect older CUs with the new ones, and still keep
468 track of their citation counts?

469 The answer to the first two questions can only be given by the database administrators. Every time a new
470 version of the database is released, the administrators go through the different CUs that compose the part-of
471 hierarchy of the database, and decide which ones need a new version. Recall that subsumption was needed
472 to transfer credit, in the case of papers, from one CU to another: the primary CU (PCU), i.e., the root of
473 the part-of hierarchy. The same can be done with data.

474 Since we have defined CUs by views when the database changes, we only need to consider creating a new
475 CU if the view changes. More precisely if D and D' are successive versions of the database and V is a
476 view, if $V(D) = V(D')$ the reference for $V(D)$ needs no change, and no new CU is necessary. However if
477 $V(D) \neq V(D')$, we may want to create a new CU.

478 Once it has been decided that a new CU needs to be created, it is necessary to determine whether the CU
479 associated with $V(D')$ is a new version of the CU for $V(D)$, or whether it is, instead, an entirely new CU.
480 The model we propose accommodates for both the possibilities; again, this is something that the database
481 administrators or curators can decide. If the content is different in the sense that there is some kind of
482 structural change, then an entirely different CU may be appropriate. Moreover, if the authorship changes,
483 then a different CU may be desirable since the two versions of the same CU are typically expected to have
484 the same authorship. These are only two examples of reasons why the DBAs may decide to consider the
485 new CU a new, independent, entity.

486 On the other hand, normally the change will be such that we want the CUs associated with $V(D)$ and $V(D')$
487 to be versions of each other, and the PCU can now become the later version $V(D')$. This preserves the
488 accuracy of the references and allows credit to accumulate on the latest version of the view.

489 In this second case, it is possible to connect the CU representing $V(D)$ to the one representing $V(D')$ through
490 the subsumption relationship. This new relationship has the precedence over the part-of relationship, and
491 thus new citations to the older version will be propagated to the new CU, and not upward to the older
492 hierarchy.

493 5 Discussion

494 Since citation graphs are currently unsuited for representing databases as first-class citizens, we have proposed
495 how to instead extend them to represent data citation in the citation graph. Among other things, this allows
496 us to capture the many citations given to databases and to give credit to the relevant authors or contributors.
497 The new model that we propose is based on a few adjustments, and builds on emerging practices in the world
498 of data citation. Above all, it has the goal of enabling an easy adoption since it is proposed as an extension
499 of existing models without requiring drastic changes. We argue that, with these extension of the current
500 model for citation graphs, we can fully achieve the goal of enabling data-citation without jeopardizing the
501 existing infrastructures.

502 The main limitation of existing models on which we focused are: (i) the lack of context on citations between
503 citable units; (ii) the inability to deal with different versions of the same CU; and consequently (iii) the
504 inability of introduce data, data evolution, and data citation (down to citable portion of a dataset) in the
505 citation graph. We showed how, by solving the first two problems through the introduction of reference
506 annotations and the subsumption property, we are also able to appropriately model data-citation in the
507 citation graph.

508 Unlike traditional scholarly publications, databases present a greater range of granularities and are subject to
509 more frequent change. Concerning the granularity of data, while it is possible to consider various scenarios,
510 we work with two main cases: either (1) only the whole database is treated as a node, or (2) each time a
511 new query is issued, a new node is added to the graph, connected to the whole database through a part-of
512 relationship.

513 The first solution is similar to what already happens with papers in data journals. In this case, the whole
514 database is represented through a single CU, i.e., one node in the graph. Every time a paper cites data in
515 the database, the citation goes to the database. Information such as the query and the rightful authors of
516 the citation may be inserted in the reference annotation of the citation. This solution is simple, but gives
517 all the citations to the whole database, thus without an explicit recognition for the rightful curators of the
518 cited data. Therefore, more computations are necessary to obtain the citation counts of the single queries
519 and the corresponding authors.

520 With the second solution, which follows the RDA specifications [Rauber et al., 2015], every time a new
521 query is issued to the database, a new CU (hence, a new node) is created. In this case, the graph represents
522 explicitly what is cited, and thus the rightful owners receive their citations without further computations.
523 However, this solution may result in an explosion of nodes. To mitigate this problem different techniques
524 could be deployed. For example, it could be possible to use algorithms of query containment to decide when
525 a query behind a citation can be answered from a CU already deployed. In this case, that CU could receive
526 that citation, instead of creating a new node. Of course, query containment is, in general, an NP-hard
527 problem, and it could become computationally prohibitive to exploit this solution, in particular in situations
528 where many nodes have already been created. Alternatively, the system could present to the interested user
529 a series of pre-computed queries, corresponding to already instantiated CUs, which may suit their citing
530 needs. In this way, the system already knows to which node to assign the citation.

531 We also observe that it could be possible to extend the proposed data model where, instead of nodes
532 presenting the metadata of the papers, the CUs are represented using or including the annotated full text of
533 a paper. In this way, annotations on the paper can be used to keep track of different types of information,
534 such as references and reference annotations. While this solution has a bigger expressive power, it also
535 increases the size and complexity of the model. As already discussed, the model proposed in this paper has
536 the advantage of being easy to implement on top of already existing systems. A new model, considering the
537 whole annotated text of a paper presents new implementation challenges, and thus requires the creation of
538 a new application from scratch.

539 It is important also to note how, as of today, there are many challenges to the implementation and proper
540 operation of data citation in general. Oftentimes the RDA guidelines for dynamic data citation are not
541 commonly implemented by many databases; it is often difficult to automatically produce context and thus
542 reference annotations that are machine-readable; and there are also many bad practices among researchers,

543 such as the one of depositing PDFs, images and tables of their papers in data repositories, calling them
544 research data. Although there are still many hindrances to the correct implementation of data citation,
545 the research community has still showed a great desire for the implementation of common techniques and
546 best practices for the correct application of these guidelines. Databases like Eagle-i ¹⁵ already provide
547 data citation snippets, others like GtoPdb automatically produce PDFs of their pages to allow an easier
548 citation of their data in form of CUs. Thus, it is our conviction that as data citation gets more traction and
549 get implemented appropriately, it would be crucial to account for it and integrate such information in the
550 common citation graph. In particular, a model like the one we propose in this paper will allow to achieve a
551 better and more fair implementation of data citation, will also benefit all researchers and become more and
552 more needed, as we transition toward the fourth paradigm of science. The more we will learn on the current
553 limits of data citation and how to address them, the faster we will come to the final goal of a correct system
554 for citing data.

555 Considering new possible research problems, we note that the citation graph in fact is, among other things,
556 a *historical record*, i.e., a record of how researchers interacted with information and other works to build
557 their expertise and new knowledge. Given this interpretation, then the graph should not be rewritable, that
558 is, it should *not* be possible to *rewrite history*. Therefore, the graph should be a timestamped “append-only”
559 structure in a way similar to the distributed ledgers. Thus, it should only be possible to insert data in it
560 without the possibility to overwrite or modify already existing information.

561 Among others, these requirements are necessary for the computation of impact factors [Garfield, 2006] where
562 it is necessary to know the number of citations received by a journal in the past two years. It is therefore
563 mandatory that this information is not modified over time. This is true also for other types of statistics that
564 researchers may be interested in.

565 In our examples, we have taken care to timestamp every element of information to make this possible.
566 The timestamps in particular indicate the moments the events “occurred” (e.g., when a citation happened),
567 not when they were inserted in the graph. However, there are several issues concerning the semantics and
568 representation of temporal information in the citation graph that requires further investigation.

569 If this property is correctly implemented, it should enable one to perform different types of query on the
570 graph. That is, past versions of the database should be accessible for accurate provenance. Ideally, given
571 the state of the graph in the present, it should be possible to rebuild a previous state at any given time in
572 the past. We call this property *history preservation*.

573 Several databases have this property. Weather data and geospatial data are generally accumulative [Justice et al., 1998].
574 Blockchains are also based on the idea that once added, a block cannot be removed or modified to guarantee
575 the preservation of the history of the transactions.

576 On the other hand, curated databases are not, in general, history preserving, in the sense that they are
577 updated and change with time. This is particularly problematic for data citation since one of its desiderata is
578 that a citation should always allow retrieving or at least knowing what was cited [Buneman, 2006]. Therefore,
579 we see the correct extension and implementation of history-preservation as an important future challenge to
580 be tackled in the implementation of a data-aware citation graph.

581 6 Related Work

582 Databases in relation to Data Citation

583 As we mentioned above, there are three main categories of databases that can be cited: (i) static databases;
584 (ii) evolving databases; and, (iii) curated databases. As a reasonable generalization, the problem of data
585 citation is easily solved for the first category since many systems and practices have been developed for
586 static databases. In this case, databases are treated as they were traditional publications since they are
587 never updated, the list of authors does not change, and even though only a portion of the database is cited,
588 the citation goes to the whole database. In this case, when we consider the citation graph, we have one

¹⁵<https://www.eagle-i.net/>

589 single node representing a database receiving all the citations from papers and data.

590 For the other two cases, data citation remains problematic. One relevant open issue is the citation of data
591 subsets generated on-the-fly by issuing general queries to the database. In this case, the main problems are
592 how to guarantee the persistence and accessibility of the data in the cited form and automatically provide a
593 complete and correct textual reference for the cited data.

594 The first problem is tackled by the RDA¹⁶. The RDA is a community-driven initiative launched in 2013
595 by different commissions, including the European Commission and the United States Government’s Na-
596 tional Science Foundation. Its goal is to build the social and technical infrastructures to enable open
597 sharing and re-use of data. The RDA “Working Group on Data Citation: Making Dynamic Data Citable”
598 (WGDC)¹⁷ [Rauber et al., 2016] has been working in the last years on large, dynamic and changing datasets.
599 While the WGDC first focused on RDBM as first forms of pilot solutions, many other types of databases
600 followed (XML, CSV, files, Git repositories, distributed databases such as VAMDC [Zwölf et al., 2019], mul-
601 tidimensional data cubes such as NetCDF/CCCA [Schubert, 2017]). The working group has finished the
602 development of its guidelines, and has now moved on into an adoption phase.

603 In particular, among the goals of the RDA WGDC [Rauber et al., 2015], there is the identification and
604 citation of arbitrary views of data. As potential solution, WGDC recommends an identification method
605 based on assigning PIDs to queries, that are then used as proxies for the data subset to be cited. The access
606 to a data subset is enabled by re-issuing the stored query and a citation is associated with the PID of the
607 query identifying the data [Rauber et al., 2016]. A PID is an identifier meant to uniquely and persistently
608 (i.e., continuously during the course of time) identify an object such as a publication, dataset and person,
609 usually in the context of digital objects that are accessible over the internet. Considering the citation graph,
610 this method based on PID adds a new citable unit every time a new query is cited and requires to check
611 query equivalence (and/or containment) to avoid the creation of a new citable unit for an already cited query.

612 The second aspect is characterized as a computational problem [Buneman et al., 2016] and some solu-
613 tions based on “query rewriting using views” [Davidson et al., 2017] have been proposed targeting general
614 queries citations for relational databases [Alawini et al., 2017b, Wu et al., 2018, Wu et al., 2019] and graph
615 databases [Alawini et al., 2017a].

616 Overall, most approaches do not consider the evolution of data and the fact that databases are not monolithic
617 objects. When those features are considered, some of the existing models propose the trivial solution of
618 treating databases and views as stand-alone objects. In our model instead, we explicitly model citable units
619 and their subsumption relationships, which allow to appropriately distribute credit.

620 Available Citation Graphs

621 The citation graph, or citation network, as a model of a graph where the vertices represent academic pa-
622 pers, has been long in use in the literature [Price, 1965] and has evolved considerably. There are different
623 implementations of citation graphs, which favor certain aspects of the information regarding publications,
624 citations, and authors, depending on the considered task. Some of them are provided explicitly for navi-
625 gational purposes, e.g., the Open Academic Graph (OAG). Others, instead, are the backbone of services
626 allowing search and exploration of scholarly works; these are the Microsoft Academic Graph (MAG), Google
627 Scholar, PubMed, Web of Science, Scopus, and Semantic Scholar.

628 The Microsoft Academic Graph (MAG) [Wang et al., 2019, Färber, 2019] is the backbone of the Microsoft
629 Academic Service (MAS), and its nodes represent five different entities: field of study, author, institution, pa-
630 per, venue, and event. An RDF version of MAG, called Microsoft Academic Knowledge Graph¹⁸ (MAKG) is
631 also available and connected to the Linked Open Data cloud.

632 The Open Academic Graph (OAG)¹⁹ is an open-source citation graph generated from the linking of two
633 other large academic graphs: MAG and ArnetMiner (or AMiner) [Wan et al., 2019], a free online service
634 used to index, search, and mine big scientific data, designed to search and perform data mining against

¹⁶<https://www.rd-alliance.org/>

¹⁷<https://www.rd-alliance.org/groups/data-citation-wg.html>

¹⁸<http://ma-graph.org>

¹⁹<https://aminer.org/open-academic-graph>

635 academic publications available on the Internet. This graph contains entities similar to the ones of MAG,
636 and it can be used as a unified sizable academic graph for the study of citation networks, paper content, and
637 the integration of multiple academic graphs with different fields and information.

638 The OpenAIRE Research Graph [Manghi et al., 2019] is the implementation of a full fledged Open Science
639 Graph. It is a collection of metadata and links connecting research entities, including articles, datasets,
640 software, etc., together with other elements such as organizations, funders, funding streams, projects, research
641 communities, and data sources²⁰. The graph today contains around 110M publications, 10M datasets, 180K
642 software research products, 7M other products with a total of 480M links between them. The aim of the
643 OpenAIRE RG is to bring discovery, monitoring, and assessment of science in the hands of the scientific
644 community [Fava, 2020].

645 The PID Graph [Fenner and Aryani, 2019, Fenner, 2020] is another example of implementation of Citation
646 Graph based around the concept of *PID* (Persistent IDentifier). The PID Graph targets citations aggregation:
647 (i) for all versions of a dataset or software source code; (ii) for all datasets hosted in a particular repository,
648 funded by a particular funder, or aggregated by a particular researcher; (iii) for a research object, such as a
649 publication, the data used in the paper, together with the software and samples used to create the dataset.
650 The PID graph adopts the outputs of the RDA WGDC. One peculiarity of the PID graph is that it does
651 not only include metadata about connections, but also metadata about the resources and implicit relations
652 about resources identified by the PIDs. This enables queries based on these metadata, making them more
653 expressive.

654 Google Scholar, PubMed, Web of Science, and Scopus are all relevant services providing citation graphs, but
655 their data is not directly accessible as a graph.

656 Google Scholar is an open general-purpose graph focusing on traditional publications and covering multiple
657 languages and publication venues. PubMed, instead, focuses on medicine and biomedical sciences [Roberts, 2001].
658 It covers medical bibliography from 1949 since today, with abstracts, review articles, and free full-text arti-
659 cles. Web of Science provides subscription-based access to multiple databases with comprehensive citation
660 data for many different academic disciplines [Falagas et al., 2008]. Finally, Scopus is Elsevier’s abstract and
661 indexing (closed) database featuring open access titles, indexes of web pages and patents, and links to both
662 citing and cited documents [Burnham, 2006]. While PubMed is an important resource for clinicians and
663 researchers, Scopus covers a wider journal range, offering also the capability for citation analysis, although
664 limited with respect of WoS, which covers articles published before 1995. Google Scholar, on the other hand,
665 presents all the pros and cons of a web search engines: it can help in the retrieval also of oblique information,
666 but it may present inadequate and less often updated citation information [Falagas et al., 2008].

667 Semantic Scholar is a project developed at the Allen Institute for Artificial Intelligence and is an AI-backed
668 search engine for scientific journal articles. It uses a combination of machine learning, natural language
669 processing, and machine vision to produce a semantic analysis of the papers of the network and to extract
670 figures, entities, and venues from the documents. It is designed to highlight the most important and influential
671 articles and to identify the connections between them [Fricke, 2018].

672 As we can see, many of these graphs and systems could work as good starting points for the implementation of
673 the proposed model. MAG and OAG already present the context, which can be used as reference annotation,
674 but lack the ability to accurately cite data and manage their versions. On the other hand, the OpenAIRE
675 graph is able to deal with granularity and different versions, but it still lacks the possibility to cite its data
676 with reference annotations, thus de facto it is still unable to deal with data citations. Nonetheless, many
677 of the systems implemented are close to the proposed model, and usually they lack one aspect (like the
678 versioning of the data or the presence of context). Therefore, we believe that a viable way forward would be
679 to implement the approach we propose on top of the already existing infrastructures.

680 Applications of the citation graphs are disparate. Some examples include: prediction of user queries over
681 the graph; recommendation systems for the generation of suggestions leveraging the relationships across the
682 different types of entities; exploration of papers, researchers, affiliations, and other entities; data integration;
683 data analysis and knowledge discovery of scholarly data through expert finding, geographic search, trend

²⁰<https://graph.openaire.eu>

684 analysis, review recommendation, association search, course search, academic performance evaluation, and
685 topic modeling [Wan et al., 2019].

686 Given the vital role of citation graphs and data citation, we argue that it is of crucial importance that existing
687 citation graphs be extended with the appropriate tools to model data citation in various forms. Most of the
688 existing citation graph do not expose their internal data model. Nonetheless, we can see they focus on the
689 same core assumption that citable objects are atomic elements with no citable portions and where evolution
690 through time is not considered. Hence, none of the models above tackle explicitly and directly the issues
691 linked to the task of modeling databases and subsets of databases, as well as the evolution of citable elements
692 through time, which is instead the goal of this work.

693 7 Conclusions and Future Work

694 Starting from a basic model of the citation graph in which the nodes are the papers, and the edges are the
695 citations between them, we highlighted three limitations of this model. They are: (i) the lack of context
696 for citations, that is, information about the *how* and *why* the citation is used along with which part of the
697 referenced object is used; (ii) the absence of a unified strategy of management of the versions of the papers
698 in the graph; and (iii) the difficulty of representing citations to databases and data generated by queries in
699 the graph.

700 To deal with these limitations, we proposed an implementation-agnostic model that includes reference an-
701 notations. These annotations contain the context of a citation (e.g., the page numbers of the citation, the
702 query issued to obtain the data or the considered bounding box).

703 We also discussed the subsumption property, which is used when a new version of a paper or a database is
704 introduced in the graph. This property indicates that the new version “takes the place” of the previous one
705 for the purpose of assigning credit. The old citations can be inherited from the new version or, depending
706 on the context, such as situations where authors have changed, different policies can be put in place.

707 Using these extensions to the basic model, we discussed how to represent data citation in the graph. While
708 important, this work is preliminary, and further work is needed if we are fully to incorporate citations or
709 other kinds of cross-reference between databases into the citation graph. Specifically:

- 710 • While we have used subsumption partly to deal with the evolution of citable units within databases, we
711 believe there is much more to be said about evolution in databases and in the citation graph itself. We
712 believe that all scientific databases should support “time travel”: it should be possible to ask queries
713 on some previous state of the database as easily as one asks queries on the current state. For many
714 databases, especially “source data”, it is important to support longitudinal queries, and this is true of
715 the citation graph itself.
- 716 • We have dealt with citations *to* databases, but what about citations *from* databases? If, as happens
717 in many curated databases, conventional citations are included within the database, then there should
718 be few problems, but what happens when a part of one database is created by a query from another
719 database? How is the citation represented; and how is it included in the citation graph?
- 720 • Finally, once we have properly supported databases within the citation graph, what kinds of biblio-
721 metric measures are possible? We have, for example, h-indexes and impact factors for conventional
722 publications. How can we appropriately measure the impact of databases?

723 We note that there is currently marginal interest to cite software and code even though interesting initiatives,
724 such as the FORCE 11 working group²¹, have been taking place and research groups are working on the
725 topic [Alliez et al., 2020, Katz et al., 2019, Katz et al., 2016]. This task presents a new set of problems, in
726 particular regarding authorship, since code is often copied or adapted from other repositories, passing from

²¹<https://www.force11.org/group/software-citation-working-group>

727 hand to hand, undergoing modifications. The characteristics of the lifecycle of software opens a whole new
728 set of problems and research questions about who is the righteous author of that piece of cited code and who
729 should receive credit from the citation.

730 Acknowledgement

731 The work was partially supported by the ExaMode project, as part of the European Union H2020 program
732 under Grant Agreement no. 825292. Matteo Lissandrini is supported by the European Union H2020 research
733 and innovation programme under the Marie Skłodowska-Curie grant agreement No 838216.

734 References

735 [COD, 2013] (2013). *Out of Cite, Out of Mind: The Current State of Practice, Polocy, and Technology for*
736 *the Citation of Data*, volume 12. CODATA-ICSTI Task Group on Data Citation Standards and
737 Practices.

738 [nzb, 2020] (2020). The difference between references and citations. [https://www.openpolytechnic.ac.
739 nz/current-students/study-tips-and-techniques/apa-referencing-and-avoiding-plagiarism/
740 the-difference-between-references-and-citations/](https://www.openpolytechnic.ac.nz/current-students/study-tips-and-techniques/apa-referencing-and-avoiding-plagiarism/the-difference-between-references-and-citations/), Retrieved March 2020.

741 [Alawini et al., 2017a] Alawini, A., Chen, L., Davidson, S. B., Portilho Da Silva, N., and Silvello, G.
742 (2017a). Automating Data Citation: The eagle-i Experience. In *2017 ACM/IEEE Joint Conference on*
743 *Digital Libraries, JCDL 2017*, pages 169–178. IEEE Computer Society.

744 [Alawini et al., 2017b] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017b). Automating data
745 citation in citedb. *Proc. VLDB Endow.*, 10(12):1881–1884.

746 [Alliez et al., 2020] Alliez, P., Di Cosmo, R., Guedj, B., Girault, A., Hacid, M.-S., Legrand, A., and
747 Rougier, N. P. (2020). Attributing and referencing (research) software: Best practices and outlook from
748 inria. *Computing in Science Engineering*, 22(1):39–52.

749 [Altman and Crosas, 2014] Altman, M. and Crosas, M. (2014). The evolution of data citation: from
750 principles to implementation. *IAssist quarterly*, 37(1-4):62–62.

751 [Belter, 2014] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation Analysis of
752 Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.

753 [Bird et al., 2008] Bird, S., Dale, R., Dorr, B. J., Gibson, B. R., Joseph, M. T., Kan, M., Lee, D., Powley,
754 B., Radev, D. R., and Tan, Y. F. (2008). The ACL anthology reference corpus: A reference dataset for
755 bibliographic research in computational linguistics. In *Proceedings of the International Conference on*
756 *Language Resources and Evaluation, LREC*. European Language Resources Association.

757 [Buneman, 2006] Buneman, P. (2006). How to cite curated databases and how to make them citable. In
758 *Scientific and Statistical Database Management, 2006. 18th International Conference on*, pages 195–203.
759 IEEE.

760 [Buneman et al., 2008] Buneman, P., Cheney, J., Tan, W.-C., and Vansummeren, S. (2008). Curated
761 databases. In *Proceedings of the twenty-seventh ACM SIGMOD-SIGACT-SIGART symposium on*
762 *Principles of database systems*, pages 1–12.

763 [Buneman et al., 2020] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D., Pawson,
764 A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn’t working, and what to do about it.
765 *Database*.

- 766 [Buneman et al., 2016] Buneman, P., Davidson, S., and Frew, J. (2016). Why data citation is a
767 computational problem. *Communications of the ACM*, 59(9):50–57.
- 768 [Burnham, 2006] Burnham, J. F. (2006). Scopus database: a review. *Biomedical digital libraries*, 3(1):1.
- 769 [Burton et al., 2017] Burton, A., Koers, H., Manghi, P., Stocker, M., Fenner, M., Aryani, A., La Bruzzo,
770 S., Diepenbroek, M., and Schindler, U. (2017). Scholix metadata schema for exchange of scholarly
771 communication links. *CERN: Geneva, Switzerland*.
- 772 [Candela et al., 2015] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Journals: A Survey.
773 *Journal of the Association for Information Science and Technology*, 66(9):1747–1762.
- 774 [Cousijn et al., 2019] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N. (2019). Bringing
775 citations and usage metrics together to make data count. *Data Science Journal*, 18(1).
- 776 [Cousijn et al., 2017] Cousijn, H., Kenall, A., Ganley, E., Harrison, M., Kernohan, D., Murphy, F.,
777 Polischuk, P., Martone, M., and Clark, T. (2017). A data citation roadmap for scientific publishers.
778 *bioRxiv*.
- 779 [Daquino et al., 2018] Daquino, M., Peroni, S., and Shotton, D. (2018). The OpenCitations data model.
780 Figshare. Online resource. <https://doi.org/10.6084/m9.figshare.3443876.v7>.
- 781 [Daquino et al., 2020] Daquino, M., Peroni, S., Shotton, D., Colavizza, G., Ghavimi, B., Lauscher, A.,
782 Mayr, P., Romanello, M., and Zumstein, P. (2020). The OpenCitations data model. In *International*
783 *Semantic Web Conference*, pages 447–463. Springer.
- 784 [Davidson et al., 2017] Davidson, S. B., Buneman, P., Deutch, D., Milo, T., and Silvello, G. (2017). Data
785 Citation: A Computational Challenge. In *Proc. of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium*
786 *on Principles of Database Systems, PODS 2017*, pages 1–4. ACM Press, New York, USA.
- 787 [Falagas et al., 2008] Falagas, M. E., Pitsouni, E. I., Malietzis, G. A., and Pappas, G. (2008). Comparison
788 of pubmed, scopus, web of science, and google scholar: strengths and weaknesses. *The FASEB journal*,
789 22(2):338–342.
- 790 [Färber, 2019] Färber, M. (2019). The microsoft academic knowledge graph: A linked data source with 8
791 billion triples of scholarly data. In *ISWC (2)*, volume 11779 of *Lecture Notes in Computer Science*, pages
792 113–129. Springer.
- 793 [Fava, 2020] Fava, I. (2020). Openaire research graph: Connecting open science – consultation phase.
794 <https://www.openaire.eu/openaire-research-graph-open-for-comments>, online September 2020.
- 795 [Fenner, 2020] Fenner, M. (2020). Powering the PID graph: announcing the DataCite GraphQL API.
796 <https://doi.org/10.5438/yfck-mv39>, online September 2020.
- 797 [Fenner and Aryani, 2019] Fenner, M. and Aryani, A. (2019). Introducing the PID graph.
798 <https://doi.org/10.5438/jwvf-8a66>, online September 2020.
- 799 [Force et al., 2016] Force, M., Robinson, N., Matthews, M., Auld, D., and Boletta, M. (2016). Research
800 data in journals and repositories in the web of science: Developments and recommendations. *Bulletin of*
801 *IEEE Technical Committee on Digital Libraries, Special Issue on Data Citation*, 12(1):27–30.
- 802 [FORCE-11, 2014] FORCE-11 (2014). *Data Citation Synthesis Group: Joint Declaration of Data Citation*
803 *Principles*. FORCE11, San Diego, CA, USA.
- 804 [Freeman et al., 2013] Freeman, G., Ding, Y., and Milojevic, S. (2013). Citation content analysis (cca): A
805 framework for syntactic and semantic analysis of citation content. *Journal of the American Society for*
806 *Information Science and Technology*, 64.
- 807 [Fricke, 2018] Fricke, S. (2018). Semantic scholar. *Journal of the Medical Library Association: JMLA*,
808 106(1):145.

- 809 [Garfield, 2006] Garfield, E. (2006). The history and meaning of the journal impact factor. *Jama*,
810 295(1):90–93.
- 811 [Gilbert and Woolgar, 1974] Gilbert, G. N. and Woolgar, S. (1974). Essay review: The quantitative study
812 of science: an examination of the literature. *Science studies*, 4(3):279–294.
- 813 [Group et al., 2016] Group, D. M. W. et al. (2016). Datacite metadata schema for the publication and
814 citation of research data.
- 815 [Harzing and Van der Wal, 2008] Harzing, A.-W. K. and Van der Wal, R. (2008). Google scholar as a new
816 source for citation analysis. *Ethics in science and environmental politics*, 8(1):61–73.
- 817 [Justice et al., 1998] Justice, C. O., Vermote, E., Townshend, J. R., Defries, R., Roy, D. P., Hall, D. K.,
818 Salomonson, V. V., Privette, J. L., Riggs, G., Strahler, A., et al. (1998). The moderate resolution
819 imaging spectroradiometer (modis): Land remote sensing for global change research. *IEEE transactions*
820 *on geoscience and remote sensing*, 36(4):1228–1249.
- 821 [Katz et al., 2019] Katz, D. S., Bouquin, D., Hong, N. P. C., Hausman, J., Jones, C., Chivvis, D., Clark,
822 T., Crosas, M., Druskat, S., Fenner, M., Gillespie, T., González-Beltrán, A. N., Gruenpeter, M.,
823 Habermann, T., Haines, R., Harrison, M., Henneken, E. A., Hwang, L. J., Jones, M. B., Kelly, A. A.,
824 Kennedy, D. N., Leinweber, K., Rios, F., Robinson, C., Todorov, I. T., Wu, M., and Zhang, Q. (2019).
825 Software citation implementation challenges. *CoRR*, abs/1905.08674.
- 826 [Katz et al., 2016] Katz, D. S., Niemeyer, K. E., Smith, A. M., Anderson, W. L., Boettiger, C., Hinsén, K.,
827 Hooft, R., Hucka, M., Lee, A., Löffler, F., Pollard, T., and Rios, F. (2016). Software vs. data in the
828 context of citation. *PeerJ Preprints*, (4):e2630v1.
- 829 [Lo et al., 2019] Lo, K., Wang, L. L., Neumann, M., Kinney, R., and Weld, D. S. (2019). GORC: A large
830 contextual citation graph of academic papers. *CoRR*, abs/1911.02782.
- 831 [Manghi et al., 2019] Manghi, P., Bardi, A., Atzori, C., Baglioni, M., Manola, N., Schirrwagen, J., and
832 Principe, P. (2019). The OpenAIRE research graph data model (version 1.3).
- 833 [Nature Physics Editorial, 2016] Nature Physics Editorial (2016). A statement about data. *Nature*
834 *Physics*, 12(10):889.
- 835 [Nosek et al., 2015] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J.,
836 Buck, S., Chambers, C. D., Chin, G., Christensen, G., Contestabile, M., Dafoe, A., Eich, E., Freese, J.,
837 Glennerster, R., Goroff, D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D., Kraut,
838 A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson, E., McNutt, M., Miguel, M., Paluck,
839 E. L., Simonsohn, U., Soderberg, C., Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S.,
840 Wagenmakers, E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research culture. *Science*,
841 348(6242):1422–1425.
- 842 [Osareh, 1996] Osareh, F. (1996). Bibliometrics, citation analysis and co-citation analysis: A review of
843 literature i. *Libri*, 46(3):149–158.
- 844 [Peroni and Shotton, 2020] Peroni, S. and Shotton, D. (2020). Opencitations, an infrastructure
845 organization for open scholarship. *Quantitative Science Studies*, 1(1):428–444.
- 846 [Peters et al., 2016] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J. (2016). Research
847 data explored: An extended analysis of citations and altmetrics. *Scientometrics*, 107(2):723–744.
- 848 [Philipp et al., 2010] Philipp, A., Bartholy, J., Beck, C., Erpicum, M., Esteban, P., Fettweis, X., Huth, R.,
849 James, P., Jourdain, S., Kreienkamp, F., Krennert, T., Lykoudis, S., Michalides, S. C.,
850 Pianko-Kluczynska, K., Post, P., Álvarez, D. R., Schiemann, R., Spekat, A., and Tymvios, F. S. (2010).
851 Cost733cat—a database of weather and circulation type classifications. *Physics and Chemistry of the*
852 *Earth, Parts A/B/C*, 35(9-12):360–373.

- 853 [Price, 1965] Price, D. J. D. S. (1965). Networks of scientific papers. *Science*, pages 510–515.
- 854 [Price and Richardson, 2008] Price, G. and Richardson, B. (2008). *MHRA style guide: a handbook for*
855 *authors, editors, and writers of theses*. MHRA.
- 856 [Pröll and Rauber, 2013] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic, large
857 databases: Model and reference implementation. In *Proceedings of the 2013 IEEE International*
858 *Conference on Big Data, 6-9 October 2013, Santa Clara, CA, USA*, pages 307–312.
- 859 [Rauber et al., 2016] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identification of
860 Reproducible Subsets for Data Citation, Sharing and Re-Use. *Bulletin of IEEE Technical Committee on*
861 *Digital Libraries, Special Issue on Data Citation*, 12(1):6–15.
- 862 [Rauber et al., 2015] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data citation of
863 evolving data: Recommendations of the working group on data citation (wgdc). *Result of the RDA Data*
864 *Citation WG*, 20.
- 865 [Roberts, 2001] Roberts, R. J. (2001). Pubmed central: The genbank of the published literature.
866 98(2):381–2.
- 867 [Schubert, 2017] Schubert, C. (2017). Implementing the rda data citation recommendations by the climate
868 change centre
869 austria (ccca) for a repository of netcdf files webinar. [https://www.rd-alliance.org/implementing%C2%](https://www.rd-alliance.org/implementing%C2%A0-rda-data-citation-recommendations-climate-change-centre-austria-ccca-repository-netcdf)
870 [A0-rda-data-citation-recommendations-climate-change-centre-austria-ccca-repository-netcdf](https://www.rd-alliance.org/implementing%C2%A0-rda-data-citation-recommendations-climate-change-centre-austria-ccca-repository-netcdf),
871 online, accessed December 2020.
- 872 [Shanableh et al., 2019] Shanableh, A., Al-Ruzouq, R., Gibril, M. B. A., Flesia, C., and Al-Mansoori, S.
873 (2019). Spatiotemporal mapping and monitoring of whiting in the semi-enclosed gulf using moderate
874 resolution imaging spectroradiometer (MODIS) time series images and a generic ensemble tree-based
875 model. *Remote Sensing*, 11(10):1193.
- 876 [Shotton, 2010] Shotton, D. (2010). CiTO, the Citation Typing Ontology. *Journal of Biomedical*
877 *Semantics*, 1(1):S6.
- 878 [Silvello, 2018] Silvello, G. (2018). Theory and Practice of Data Citation. *Journal of the American Society*
879 *for Information Science and Technology (JASIST)*, 69(1):6–20.
- 880 [Sinha et al., 2015] Sinha, S., Shen, Z., Song, Y., Ma, H., Eide, D., Hsu, B. J. P., and Wang, K. (2015). An
881 overview of microsoft academic service (MAS) and applications. In *Proceedings of the 24th International*
882 *Conference on World Wide Web (WWW ’15 Companion)*, pages 243–246, New York, NY, USA. ACM.
- 883 [Southan et al., 2015] Southan, C., Sharman, J. L., Benson, H. E., Faccenda, E., Pawson, A. J., Alexander,
884 S. P., Buneman, O. P., Davenport, A. P., McGrath, J. C., Peters, J. A., et al. (2015). The iuphar/bps
885 guide to pharmacology in 2016: towards curated quantitative interactions between 1300 protein targets
886 and 6000 ligands. *Nucleic acids research*, 44(D1):D1054–D1068.
- 887 [Starr and Gastl, 2011] Starr, J. and Gastl, A. (2011). isctdby: A metadata scheme for datacite. *D-Lib*
888 *Magazine*, 17(1/2).
- 889 [Tang et al., 2008a] Tang, J., Zhang, J., Yao, L., Li, J., Zhang, L., and Su, Z. (2008a). ArnetMiner:
890 extraction and mining of academic social networks. In *Proceedings of the 14th ACM SIGKDD*
891 *international conference on Knowledge discovery and data mining*, pages 990–998. ACM.
- 892 [Tang et al., 2008b] Tang, T., Zhang, J., Yao, L., Li, J., Zhang, L., and Su, Z. (2008b). Arnetminer:
893 Extraction and mining of academic social networks. In *KDD*, pages 990–998. ACM.
- 894 [Wan et al., 2019] Wan, H., Zhang, Y., Zhang, J., and Tang, J. (2019). AMiner: Search and mining of
895 academic social networks. *Data Intell.*, 1(1):58–76.

- 896 [Wang et al., 2019] Wang, K., Shen, Z., Huang, C., Wu, C., Eide, D., Dong, Y., Qian, J., Kanakia, A.,
897 Chen, A., and Rogahn, R. (2019). A review of microsoft academic services for science of science studies.
898 *Frontiers in Big Data*, 2:45.
- 899 [Wilke, 2015] Wilke, C. (2015). What constitutes a citable scientific work?
900 <https://serialmentor.com/blog/2015/1/2/what-constitutes-a-citable-scientific-work>.
- 901 [Wilkinson et al., 2016] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M.,
902 Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B., Bourne, P. E., et al. (2016). The fair guiding
903 principles for scientific data management and stewardship. *Scientific data*, 3.
- 904 [Wu et al., 2018] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data Citation: Giving
905 Credit Where Credit is Due. In *Proc. of the 2018 International Conference on Management of Data,*
906 *SIGMOD Conference 2018*, pages 99–114. ACM Press, New York, USA.
- 907 [Wu et al., 2019] Wu, Y., Alawini, A., Deutch, D., Milo, T., and Davidson, S. (2019). Provcite:
908 Provenance-based data citation. *Proc. VLDB Endow.*, 12(7):738–751.
- 909 [Zwölf et al., 2019] Zwölf, C. M., Moreau, N., Ba, Y. A., and Dubernet, M. L. (2019). Implementing in the
910 vamdc the new paradigms for data citation from the research data alliance. *Data Science Journal*, 18(1).