

# Credit Distribution through Data Provenance in Relational Scientific Databases

Dennis Dosso<sup>a</sup>, Susan B. Davidson<sup>b</sup>, Gianmaria Silvello<sup>a</sup>

<sup>a</sup>*Department of Information Engineering, University of Padua, Italy*

<sup>b</sup>*Department of Computer and Information Science, University of Pennsylvania, USA*

---

## Abstract

Digital data is an important form of research product for which citation, and the generation of credit or recognition for authors, is still not well understood. The notion of *data credit* has therefore recently emerged as a new metric, defined and based on data citation theory.

Data credit is a real value that represents the importance of data cited by a paper or by another research entity. Credit can be used to annotate data contained in a curated scientific database, and used as a measure for the importance and impact of that data in the research world. As such, it is a new method that, together with traditional citations, helps recognize the value of data and its creators.

In this paper we explore the problem of Data Credit Distribution, the process by which credit is distributed to the database parts responsible for the production of data being cited by a research entity.

We adopt as use case the IUPHAR/BPS Guide to Pharmacology (GtoPdb), a widely-used curated scientific relational database. We define two new distribution strategies based on two forms of data provenance, why-provenance and how-provenance.

Using different distribution strategies, we show how credit can highlight frequently used database areas and how it be used as a new bibliometric measure for data and their corresponding curators. In particular, credit rewards data and authors based on their research impact, not merely on the number of citations. We also show how different distribution strategies, based on different types of data provenance, can vary in their sensitivity to an input tuple in the generation of the output data and reward input tuples differently.

*Keywords:* Data Citation, Data Credit

---

## 1. Introduction

Citations are an essential component of scientific research, enabling research products to be found as well as the relationships between research products to be created and understood. They form a basis on which to give credit to authors, papers, and venues [19, 20, 54]. Citations are used, among other things, to decide on tenure, promotion, hiring, and funding of grants for researchers [21, 32, 37, 40].

Science and research are increasingly digital, and there are numerous curated databases that are at the core of scientific research efforts [12]. It is therefore generally accepted that data must be cited and citable [15, 38], and that data citations should contribute to the scientific reputation of researchers, scientists, data curators, and creators [4, 50]. It is also accepted that data citations should be counted alongside of traditional citations, and contribute to bibliometrics indicators [7, 44].

Many initiatives, at different levels, have been promoted to make data citation a reality. Scientific publishers, such as Elsevier, Springer and Nature, have been defining data policies and author guidelines to include data citations in the reference lists of published papers [19]. The European Commission has introduced the Open Research Data Pilot (ODP), whose aim is to improve and maximize the access and re-use of research data, together with an increase to the credit given to data creators and curators [48]. Initiatives such as FORCE11 and ESIP (Earth Science Information Partners) have collaborated on data and software citation principles and guidelines [26]. Other examples are the National Science Foundation (NSF), and the National Institute of Health (NIH) in the US [48].

Moreover, there are activities to promote and specify guidelines for data citations. A significant activity getting a broad adoption, is the Research Data Alliance (RDA), that produced a recommendation on citing specific subsets of dynamic data [47]. While this approach provides reference and access to a precise subset of data, it does not address specific credit concerns for that subset, such as when different authors contribute to a larger collection [43].

A central problem in the data citation process is how to attribute credit to data creators and curators [11]. How to handle and count the credit generated by data citation, and how it contributes to traditional and new

bibliometrics, are long-standing research issues [9, 28]. However, even when correctly applied, data citations and the bibliometrics computed using them do not always correctly or completely reward the creators of data used in a database. Data, in fact, is often cited at the “database level” or the “webpage level”. In the first case, the whole database is cited and therefore all credit goes to the key personnel of the database. In the second case, the database has a website with webpages that can be individually cited. The webpages use data extracted from the database, which is aggregated by topic and built to resemble a traditional research paper. Often the creators and curators of the webpage’s data are not credited or only marginally credited for their work [3].

Recently, the concepts of *data credit* and *Data Credit Distribution* (DCD) [27, 36, 53] have emerged, built on top of methodologies for data citation. Data credit is a value that is computed based on the importance of the data being cited in a paper, and represents the impact of the data on the citing paper. The DCD problem consists of distributing this credit to elements in the databases in the citation graph that are responsible for the generation of the data being cited. The goal of DCD is to improve and expand the reach of data citation, rather than being an alternative to it. This means that to employ DCD techniques, we need data citations in some form.

In this paper, we consider data credit as a measure of value for data in a (curated) scientific database. Credit is a real value that can be assigned to data of any kind and at any level of granularity. Therefore the concept of “data” is left intentionally vague, although in this paper we focus on relational databases. Credit is a positive *real* value, acting as a proxy for the value of data based on the measure of citations, accesses, clicks, downloads, or other surrogates for data use. We call DCD the process, method, or algorithm used to assign credit to a given datum or dataset.

The DCD problem differs from the traditional citation setting since:

1. When a paper  $p_1$  cites another paper  $p_2$ , a +1 citation “credit” is given to  $p_2$ , and to all its authors. It does not matter why or how paper  $p_1$  cites paper  $p_2$ ,<sup>1</sup> the result is always +1 to the citation count of  $p_2$  and of its authors. A different credit distribution strategy can assign a quantity of credit to  $p_2$ , and its authors, that is *proportional* to the

---

<sup>1</sup>Note that there is vast research on this topic and many alternative proposals, but none of them currently work at a large scale.

- 70 role played by  $p_2$  in  $p_1$ . Hence, we can weight the importance of the  
71 cited entities and assign credit according to their role.
- 72 2. Traditional citations are *atomic*: a citation from  $p_1$  to  $p_2$  can never  
73 be broken into pieces and assigned in part to  $p_2$  and in part to other  
74 papers or data that contributed to  $p_2$ . In contrast, with data credit,  
75 we use a *non-atomic* real value, which can be divided and distributed  
76 to multiple components of a database.
- 77 3. Credit can be *transitive*, that is, it can be propagated through one  
78 cited entity to other entities cited by it that contributed to its content.  
79 Citations, traditionally, are not.

80 We study the DCD problem in the context of relational databases (RDBs)  
81 since they are widely used <sup>2</sup> and are the main focus of current work in data  
82 citation methods [12, 14, 45]. RDBs are also frequently a test-bed for new  
83 methods that can be adapted to other databases, e.g., graphs or document  
84 databases. The “portions” of data in an RDB that can be credited can be  
85 defined at different levels of granularity, in particular: (i) the whole database,  
86 (ii) tables, (iii) tuples, and (iv) attributes. The ability to specify different  
87 levels of granularity in a relational database allows us to define the DCD  
88 problem at a particular level of granularity. In this paper, we focus on DCD  
89 at the tuple level.

90 The DCD process is summarized in Figure 1:

91 **Step 1** Scientists and experts contribute the curated information contained  
92 in a scientific database. These are called the “Data Curators”.

93 **Step 2** Other researchers use the data in their research, and when possible,  
94 cite them.

95 **Step 3** The citation to the data generates credit, that can be used as a  
96 proxy for the impact of the data on the citing paper. This credit is  
97 represented as a real value  $k \in \mathbb{R}_{>0}$ .

98 **Step 4** Given the database instance  $I$  and the query  $Q$ , it is possible to  
99 compute the *data provenance* of  $Q(I)$ . The provenance of  $Q(I)$  is a

---

<sup>2</sup>The “relational database market alone has revenue upwards of \$50B” [1].

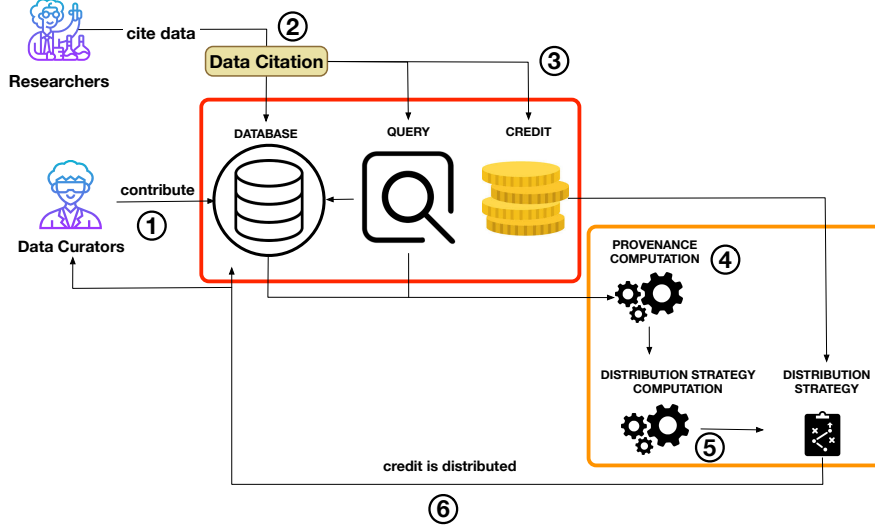


Figure 1: Overview of the credit distribution pipeline.

form of metadata that describes the generation process undertaken by  $Q$ , and the data used in  $I$  to generate the output [17]. Many different notions of provenance have been proposed in the literature for data in database management systems [13, 22, 30], describing different kinds of relationships between data in the input and the output of a query. As reported in [17], these provenances have been used in several applications beyond giving information on how queries work, for example, annotation propagation and the view update problem. In this paper, we consider three types of provenance: lineage, why-provenance, and how-provenance.

**Step 5** Provenance is input to the DCD problem, whose aim is to compute the *Credit Distribution Strategy* (CDS, also referred only as Distribution Strategy, DS). The CDS is a function that distributes  $k$  to the data in the input database  $I$ , and is defined on the basis of citation policies decided at the database administration level or at the domain community level. In this paper, since we base CDS on data provenance, we describe three CDS, each one based on a different form of provenance.

**Step 6** Once the CDS is computed, it is used to distribute the given credit  $k$  to the parts of the database that are responsible for the generation

119 of  $Q(I)$ . Transitively, this credit is also divided and given to the corre-  
120 sponding authors of those data.

121 This paper expands our recent work in [24], which addressed the problem  
122 of how to reward data and data curators who are typically overlooked in  
123 current citation systems. In that work, we first defined the problem of DCD  
124 in relational databases, and proposed a viable Distribution Strategy (DS)  
125 based on *lineage*, which is the simplest form of *data provenance*. The lineage  
126 of a tuple  $t$  in the output  $Q(I)$  is defined as the set of all and only the tuples  
127 in the database instance  $I$  that are “relevant” to the production of  $t$ , that  
128 is the tuple that are used by  $Q$  in the production of  $t$ . The lineage-based  
129 strategy equally redistributes the credit  $k$  to the tuples in the lineage set,  
130 thus each tuple receives credit  $k/|L_t|$ , where  $L_t$  is the lineage set of  $t$ .

131 One may argue that this DS is too simplistic, since lineage only tells  
132 the relevant tuple used to produce the output, and does not convey any  
133 information about their role or importance in the query. Therefore, one may  
134 desire to give more credit to the tuples that are more relevant or *essential*  
135 to the production of the output, i.e. those tuples that, if removed, would  
136 prevent the output tuple from appearing in the final result, or those tuples  
137 used more than once by the query.

138 Therefore, in this paper, we expand the ideas in [24] by proposing two  
139 new DSs based on other forms of data provenance: why-provenance [13]  
140 and how-provenance [30]. We compare them with the lineage-based solu-  
141 tion, and discuss why one may be preferred to another depending on the  
142 application and its goals. In particular, we show that why-provenance and  
143 how-provenance are more sensitive to the *role* of a tuple in a query, i.e. how  
144 many times the tuple is used and how it is used. The DS based on why-  
145 provenance gives more reward to tuples that are essential to the production  
146 of the result set, whereas the DS based on how-provenance also takes into  
147 consideration the different ways that a tuple is used.

148 For evaluation, we use a well-known curated database, the IUPHAR/BPS<sup>3</sup>  
149 Guide to Pharmacology [31], also known as GtoPdb<sup>4</sup>, which contains ex-  
150 pertly curated information about diseases, drugs, cellular drug targets, and

---

<sup>3</sup>International Union of Basic and Clinical Pharmacology/British Pharmacology Soci-  
ety

<sup>4</sup><https://www.guidetopharmacology.org/>

151 their mechanisms of action. We chose GtoPdb for two main reasons: (i) it  
152 is a widely-used and valuable curated relational database, (ii) many papers  
153 in the literature use, and cite its data (i.e., families, ligands, and receptors).  
154 Real queries used in papers can therefore be seen as data citations which, in  
155 turn, can be used to assign data credit.

156 We perform four sets of experiments. In the first one, real queries are ex-  
157 tracted from papers published in the British Journal of Pharmacology (BJP),  
158 that represent data citations to GtoPdb, and are used to distribute credit  
159 in the database using the three different provenance-based DSs. In the sec-  
160 ond and third experiment we analyse the behaviour of the different DS when  
161 complex citation queries are employed. In the fourth set of experiments we  
162 use both real and synthetic queries to assess the difference between tradi-  
163 tional citation and the notion of credit distribution in terms of rewarding  
164 those responsible for the data, e.g. data curators.

165 **Contributions** of this work include:

- 166 • The definition of new distribution strategies for the problem of Data  
167 Credit Distribution, based on why-provenance and how-provenance;
- 168 • An in-depth analysis of the effects of credit distribution on real-world  
169 curated data and of the differences between the three proposed Distri-  
170 bution Strategies.
- 171 • A comparison between the behavior of traditional citations and data  
172 credit in rewarding data curators.

173 **Outline.** The rest of the paper is organized as follows: Section 2 presents  
174 the background and related work. Section 3 describes the GtoPdb use case  
175 we adopted. Section 4 briefly presents the forms of provenance used in the  
176 paper. Section 5 describes the credit distribution problem and the proposed  
177 distribution strategies. In Section 6 we present the experimental evaluation.  
178 Finally, Section 8 draws some conclusions and outlines future work.

## 179 2. Background

180 *Data in Research.* The world of research is rapidly transitioning towards the  
181 *fourth paradigm of science* [33], that is, data-intensive scientific discovery,  
182 where data are important for scientific advances as well as for traditional  
183 publications [6].

184 The scientific community is promoting an *open research culture* [42],  
 185 founded on methods and tools to share, discover, and access experimental  
 186 data. The community has identified the FAIR principles (Findable, Acces-  
 187 sible, Interoperable, and Reusable) [51], that should be enforced by every  
 188 database. In particular, data should be accessible from the articles, journals,  
 189 and papers that cite or use them [19]. Aspects such as the need for the *repro-*  
 190 *ducibility* of experiments through the used data; the *availability* of scientific  
 191 data; the *connections* between data and the scientific results are all needed  
 192 aspects for the fourth paradigm, and are all relevant to the domain of *data*  
 193 *citation* [34].

194 *Data Citation: Principles and Motivations.* Data Citation principles were  
 195 proposed in [18], and later summarized and endorsed by the Joint Declaration  
 196 of Data Citation Principles (JDDCP) [39]. The principles are divided into  
 197 two groups [48]. The first one contains principles concerning the role of  
 198 data citation in scholarly and research activities such as the (i) *importance*  
 199 of data (why data citation is important and why data should be considered  
 200 as first-class citizens); (ii) *credit* and *attribution* to the creators and curators  
 201 of the data; (iii) *evidence*; (iv) *verifiability*; and *interoperability*, with these  
 202 last three requiring data citation methods to be flexible enough to operate  
 203 through different communities. The second group defines the main guidelines  
 204 to establish a data citation systems, and contains principles such as the (i)  
 205 *unique identification* of the data being cited; (ii) (*open*) *access* to data; (iii)  
 206 guarantee of *persistence* and *availability* of citations even after the lifespan  
 207 of the cited entity; the (iv) *specificity* of a citation, i.e. it must lead to the  
 208 data set originally cited.

209 It is possible to outline six main motivations for data citation [48]:

- 210 • *Data attribution*: identify the individuals that should be credited for  
 211 data with variable granularity.
- 212 • *Data connection*: connect papers to the data being used.
- 213 • *Data Discovery*: citations helps to find data records and subsets that  
 214 would be otherwise not findable via search engines.
- 215 • *Data Sharing*: share data obtained by researchers within the whole  
 216 community.



- *Data Impact*: highlight the results obtained in writing papers using specific data, the frequency and modality data were used.
- *Reproducibility*: data citation greatly impacts the reproducibility of science [5]. Many authoritative journals ask to share data and provide valid methodologies to reproduce experiments.

## 2.1. Data Citation in Relational Databases

In this paper, we develop our methods and experiments on relational databases. RDBs have been the main target of data citation methods since the surge of the data-centric research paradigm. The RDA “Working Group on Data Citation: Making Dynamic Data Citable”<sup>5</sup> [46] has been working in the last years on large, dynamic, and changing datasets. The working group has finished the development of its guidelines and has now moved on into an adoption phase. The datasets considered by the Working Group are often relational.

In one of its most recent sessions [47], the Working Group (WG) on Data Citation reported that there are various implementations of its guidelines for Data Citation on MySQL/Postgres relational databases. Some of these databases are: DEXHELPP<sup>6</sup> (Social Security Records); NERC (ARGO Global Array); EODC (Earth Observation Data Centre) [29]; LNEC (River dam monitoring); MDS (Million Song Database) [8]; CBMI<sup>7</sup> (Center for Biomedical Informatics); VMC (Vermont Monitoring Cooperative); CCA<sup>8</sup> (Climate Change Center Austria); VAMDC (Virtual Atomic and Molecular Data Center) [25, 55].

More examples of work on data citation in relational databases are [2, 12, 23, 52]. The website <https://fairsharing.org/> keeps a long updated list of curated and scientific databases (many of which are relational or graph-based) following FAIR guidelines. These databases are citable since they are compliant with the most recent guidelines, and they are in the vast majority of cases accessible via dynamically created Webpages. In all these databases is, therefore, possible to implement DCD on top of the existing infrastructures for citing data.

---

<sup>5</sup><https://www.rd-alliance.org/groups/data-citation-wg.html>

<sup>6</sup><http://www.dexhelpp.at/>

<sup>7</sup><https://medicine.missouri.edu/centers-institutes-labs/center-for-biomedical-informatics>

<sup>8</sup><https://ccca.ac.at/startseite>

248 Data citation techniques are primarily applied to relational databases  
249 because of their diffusion and also because the portions of data that are to  
250 be cited are easily identified: the whole database, a relation, a tuple, or  
251 even an attribute. Many papers [2, 10, 12] consider more complex citable  
252 units, recognizing that often the *views* of a database are the ones to be cited.  
253 Generally, a *view* is a query on the database. To this end, [52] suggested  
254 decomposing the database in a set of views, where each view is associated  
255 with its citation.

256 At present, the most common practices to cite databases include:

- 257 1. A database cited as a whole, even though only parts of the databases  
258 are used in the papers or datasets. Alternatively, the so-called “data pa-  
259 pers” can be cited, being traditional papers that describe a database [16].  
260 In this case, all the credit from the citations goes to the database ad-  
261 ministrators or to the authors of the data papers.
- 262 2. Subsets of data, obtained by issuing queries to a database, are individ-  
263 ually cited. This is the solution adopted by the *Resource Data Alliance*  
264 (RDA) working group on Data Citation [46]. In this case, the credit  
265 generated from citations can be distributed among the contributors of  
266 the portions of data being cited, and/or to the database administrators.
- 267 3. The database is accessible via a series of Webpages that arrange the  
268 content of the database by topic or theme. Examples in the life science  
269 domain include the Reactome Pathway database [35], the GtoPdb [31],  
270 and the VAMDC [55]. Every single Webpage is unequivocally identifi-  
271 able and can be individually cited.

## 272 2.2. Data Credit

273 Data credit is related to data citation: they both aim to recognize the  
274 work of data creators and curators. Data credit can therefore also be seen as  
275 a by-product of data citation, since credit attribution is impossible without  
276 the presence of data citations.

277 Katz [36] suggests the need for a *modified citation system* that includes  
278 the idea of *transient* and *fractional credit*, to be used by developers of research  
279 products as software and data. In the paper two considerations are made:  
280 (i) research objects such as data and software are currently not formally  
281 rewarded or recognized by the community; (ii) even in traditional papers,  
282 the contribution of each author to the work is hard to understand, unless  
283 explicitly specified in the paper. This is even more true for data, where  
284 different groups of people work on the same database.

285 In [36] credit is defined as a “quantity” that describes the importance of a  
 286 research entity, such as papers, software, or data, mentioned in a citation. It  
 287 also proposed the idea of a *distribution* of credit from research entities, such as  
 288 papers or data, to other research entities through citations. **Therefore, when**  
 289 **talking about data credit, here we are focusing on two aspects of the topic:**  
 290 *credit computation*, the process in which the quantity of credit generated by  
 291 the citation is computed, and *credit distribution*, the process by which credit  
 292 is distributed and assigned to the responsible entities that contributed to the  
 293 generation of the data being cited. In this paper we focus on the latter.

294 These two processes are done by exploiting the structure of the *citation*  
 295 *graph*, a directed graph whose nodes are publications and edges are citations.  
 296 This graph is the model at the core of systems such as Google Scholar and  
 297 the Web of Science. We add to this that the concept of credit can be built  
 298 on top of the existing infrastructure handling traditional and data citations.

299 Katz [36] further explores the idea of a *distribution* of credit from research  
 300 entities (i.e., papers and data) to other research entities through citations  
 301 that connect them. Thanks to traditional citations and now also to data  
 302 citations, this distribution is finally possible, at least between papers and  
 303 data. Some problems related to traditional citations can thus be solved by  
 304 citations:

- 305 1. Credit rewards research entities that to date are not (formally) recog-  
 306 nized (a goal shared with data citation).
- 307 2. Credit can reward authors *proportionally* to their role in generating the  
 308 entity. The more an author contributes to a paper, the more credit is  
 309 given to him. Zou and Peterson [54] work on something similar with  
 310 their zp-index, which includes in its formulation the position (and thus  
 311 the role) of a publication author to represent its impact in the work  
 312 itself.
- 313 3. Credit can be *transitively* channeled through a chain of papers citing  
 314 each other, thus enabling the rewarding of older papers that are no  
 315 more cited, since other papers summarize or report their content but  
 316 are nevertheless crucial in a research area for the influence of their  
 317 content.

318 Fang [27] presents a framework to distribute the credit generated by a  
 319 paper to its authors and to the papers in its reference list in a transitive way.  
 320 Let us consider the *citation graph* as the graph where the nodes are papers

and the links are the citations among them. In this graph, every paper is a source of credit, which is then transferred to the neighboring nodes. The quantity of credit received by each cited paper depends on its impact/role in the citing paper. So far, this theoretical framework is limited to papers, but it can be easily extended to a citation graph including both papers and data.

Zeng et al. [53] proposes the first method to compute credit within a network of papers citing data. Adopting a network flow algorithm, they simulate a random walker to estimate a score for each dataset, leveraging real-world usage data to compute the credit. This is the first step towards an automatic credit computation procedure. This proposal is, however, limited to assigning credit to whole datasets, and it does not deal with the granularity of data. It does not work to assign credit to a single research entity within a dataset. Differently from Zeng et al. [53], we do not treat the credit computation process, but we focus on the distribution process.

### 2.3. Data Provenance

To distribute credit, we base our methods on *data provenance*. Data provenance is information that describes the origin and the process of creation of data. It can also be seen as metadata pertaining to the derivation history of the data. It is particularly useful to help users to understand where data are coming from, and the process they went through. Data citation and data provenance are closely linked [3] since both are forms of annotations on data retrieved through queries. Data provenance has been widely studied in different areas of data management. In this paper, we focus on provenance for database management systems (DBMS). For further details on data provenance, please refer to surveys like [17] and [49].

Cheney et al. [17] presents four main types of data citation for DBMS: *lineage* [22], *why-provenance* [13], *how-provenance* [30] and *where-provenance* [13].

Let us start with the first three provenances. Given a database instance  $I$ , a query  $Q$ , and the result  $Q(D)$ , consider one tuple  $t$  of the output. Its provenance is information about its generation through the tuples of the input that are used by  $Q$ . Different types of provenance convey different levels of information. Since these three provenances are computed for each tuple of the output, they are also referred to as *tuple-based*.

Lineage is the simplest among the forms of provenance. It has been defined in different ways [17], but it can be thought of as the set of all the

357 tuples that are used in some way by the query to produce the output tuple,  
358 the ones that are somehow *relevant* to its generation.

359 The definition of why-provenance is based on the notion of *witness set*.  
360 A witness is a set of relevant tuples that guarantees the existence of  $t$  in  
361  $Q(D)$ . The lineage is therefore an example of a witness. The why-provenance  
362 of a tuple  $t$  is a peculiar set of witnesses – described in [13] – that are  
363 computed from the query, called *witness basis*. A witness basis may be  
364 composed of more than one witness. Therefore, the why-provenance contains  
365 more information than the lineage, since it describes *alternative* ways in  
366 which the same output may be generated.

367 The how-provenance takes the form of a polynomial, called *provenance*  
368 *polynomial*, where the variables are taken from the set of identifiers of the  
369 tuples (provided that each tuple in  $I$  has an identifier) and the coefficients are  
370 drew from  $\mathbb{N}$ . This provenance also contains information on *how* the input  
371 tuples are used. For example, when two tuples are combined by a join, they  
372 are also combined in the polynomial by the  $\cdot$  operator. When two or more  
373 tuples become equivalent due to a union or a projection, the corresponding  
374 monomials are combined by the  $+$  operator.

375 It has been shown in [17] that the how-provenance is the more general  
376 and informative of the three, containing the other two.

377 Where-provenance, differently from the other three, is *attribute-based*, so  
378 we do not take it into account in this work since we consider the tuple as the  
379 finest citable unit.

### 380 3. Use Case: GtoPdb

381 As use case we refer to the IUPHAR/BPS Guide to Pharmacology [31]  
382 or GtoPdb<sup>9</sup>. GtoPdb is a well-known and well structured scientific relational  
383 database that contains expertly curated information about diseases, drugs  
384 in clinical use, their cellular targets, and the mechanisms of action on the  
385 human body. It is curated and maintained by the GtoPdb Committee, and  
386 by 96 subcommittees, comprising 512 scientists collaborating with in-house  
387 curators who draw the information contained in the database from high-  
388 quality pharmacological and medicinal chemistry literature. Roughly 1000  
389 researchers from all over the world have contributed to the database, and the

---

<sup>9</sup><https://www.guidetopharmacology.org/>

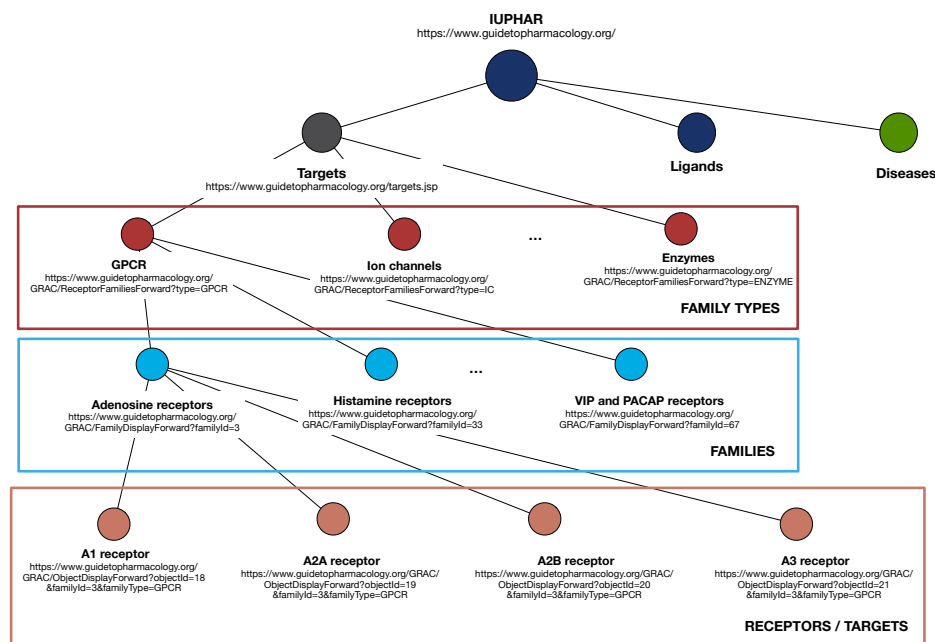


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

curators wanted to give recognition to these contributors. This led to some early work on data citation [10].

GtoPdb is relational, but its logical structure is hierarchical as shown in Figure 2. The information contained in the database is also organized into webpages focused on specific diseases, targets or ligands, and families for easier access by users. As depicted in Figure 2, the database can be thought of as a tree where the root is the database; the first level consists of all targets, ligands, and diseases; and the lower levels consists of specific targets, ligands and diseases. In this paper, we focus on targets; thus at the third level in the figure we show examples of family types, at the fourth level we show specific families of targets (a finer level of granularity), and finally, at the last level, the single targets (also known as receptors).

GtoPdb provides access to the webpages corresponding to all these nodes through URLs. The webpages corresponding to target families all present a similar structure, as shown in Figure 3 for the “Adenosine receptors” family. Each page has an *Overview*, a brief text describing the content of the page; a list of *Receptors* comprising the family; a section of *comments* about the

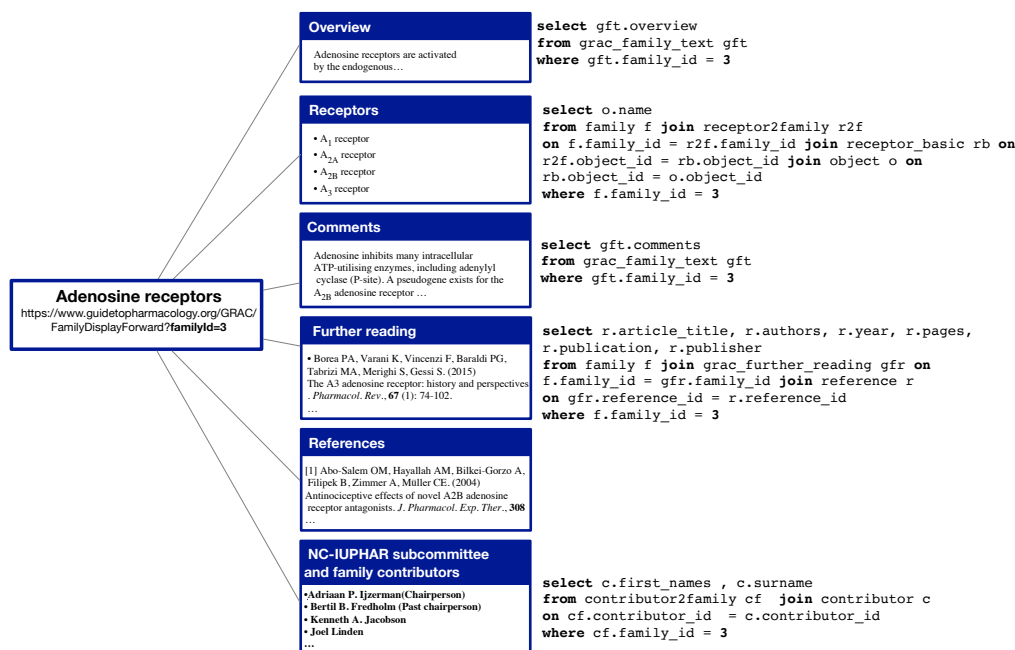


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

family; the *References*, a list of the papers consulted by the curators of the page, similar to a reference list of a paper; the *further reading* list, reporting papers that an interested reader may want to consult to obtain more insight on the family; and a final section called *How to cite this family page*, containing text snippets useful to cite the specific page or the whole database. Figure 3 shows the SQL code that retrieves the information used to build the corresponding sections (apart from the References section). Therefore, each family page can be considered a full-fledged traditional publication, consisting of title, authors, abstract (the overview), content, and references.

In practice, many papers in the literature only reference GtoPdb (the root) without including a reference to the specific page being cited. That is, they only cite a paper describing GtoPdb as a whole (e.g., [31]) and refer to targets, ligands, diseases, etc. only by name. Thus, citations to specific families are *de-facto* “hidden” to citation systems such as Google Scholar, and useless for the computation of bibliometrics.

In certain “lucky” cases, as with papers available in PDF and published

family			contributor2family		
id	name	type	id	family_id	contributor_id
$f_1$	Dopamine Receptors	gpcr	$c2f_1$	$f_1$	$c_1$
$f_2$	Bile Acid Receptor	gpcr	$c2f_2$	$f_1$	$c_2$
$f_3$	FAK Family	enzyme	$c2f_3$	$f_2$	$c_3$
$f_4$	YANK Family	enzyme	$c2f_4$	$f_4$	$c_1$

contributor		
id	Name	Country
$c_1$	John Smith	UK
$c_2$	Jim Doe	UK
$c_3$	Hans Zimmerman	Germany
$c_4$	Roberta Rossi	Italy

Table 1: Example of a database consisting of three tables. **family** includes some receptor families in the database; **contributor** contains the name and country of contributors; **contributor2family** connects contributors to the families they contributed to.

in the British Journal of Clinical Pharmacology <sup>10</sup> (BJCP), when a family, ligand, receptor name, etc. are used, they have a hyperlink pointing to the corresponding webpage in GtoPdb. Therefore, the citations to the families can be detected and counted using the URLs reported in the papers. However, these citations to GtoPdb webpages are not counted as such by citation systems, so they are not converted into credit for curators and collaborators.

For our running example, consider Table 1. This simplified version of GtoPdb illustrates three tables: **family**, **contributor** and **contributor2family**. The first table, **family**, has tuples representing families with three attributes: the id of the family, its name, and type. Table **contributor** consists of people who have helped generate the data of the database. The third table, **contributor2family**, serves as a link between the families and the people who contributed to them. For instance, “John Smith” ( $c_1$ ) contributed to “Dopamine Receptors” ( $f_1$ ) as well as to the “YANK Family” ( $f_4$ ). We use this example throughout the rest of the paper. In particular, we are using the id attribute of the tables as *provenance token* of its corresponding tuples, that is, as a symbol that serves to identify a tuple when talking about provenance.

<sup>10</sup><https://bpspubs.onlinelibrary.wiley.com/journal/13652125>



## 4. Data Provenances

In this section, we present the three types of provenance used in this paper: lineage, why-provenance, and how-provenance. We also discuss of Causality and Responsibility that, even though are not forms of data provenance *per se*, they are still used as basis to define a DS.

### 4.1. Lineage

Lineage was first introduced by Cui et al. [22]. Given a database instance  $I$  and query  $Q$ , lineage associates with each tuple  $o \in Q(I)$  the set of tuples in the input that contributed to its “production” [17]. As an example, consider the following SQL query Q1, applied to the database described in Table 1, that asks for the names of families curated by researchers based in the United Kingdom (UK):

```
Q1: SELECT DISTINCT f.name
FROM family AS f JOIN contributor2family AS c2f
ON f.id = c2f.family_id
JOIN contributor AS c ON c2f.contributor_id = c.id
WHERE c.country = 'UK'
```

id	name	lineage
$o_1$	Dopamine Receptors	$\{f_1, c2f_1, c_1, c2f_2, c_2\}$
$o_2$	YANK Family	$\{f_4, c2f_4, c_1\}$

Table 2: Result of an SQL query applied to the database instance in Table 1, which asks for the names of families curated by a researcher based in the UK. Attribute `id` is not part of the output and was added to succinctly identify each tuple as provenance token. Each tuple is also annotated with its lineage.

Table 2 shows the query result set, which consists of two tuples. We add an extra attribute `id` so that we can easily refer to each result tuple. The lineage for tuple  $o_1$  is the set  $\{f_1, c2f_1, c_1, c2f_2, c_2\}$ , since the tuple  $f_1$  was joined with  $c2f_1$  and then with  $c_1$ , and was also joined with  $c2f_2$  and  $c_2$ . No other tuple is used in the database to produce  $o_1$ . For tuple  $o_2$  the lineage is  $\{f_4, c2f_4, c_1\}$ . Lineage is defined for each tuple of the output, and can differ between tuples.

#### 4.2. Why-Provenance

Why-Provenance was first defined in terms of a deterministic semistructured data model and query language [13]. While why-provenance can be defined in many ways, we refer to [17], where it is expressed in terms of the relational model using the relational algebra.

In particular, while lineage aims to find all and only the tuples in the input relevant to the production of an output tuple, why-provenance aims to find sub-instances of the input that “witness” a part of the output. Given a tuple  $t$  in the query’s output, a *witness* is any sub-instance of the database that produces  $t$ . In particular, the whole database and the lineage of  $t$  are both witnesses of  $t$ . Since the definition of witness allows for the presence of “irrelevant” tuples, the set of all witnesses is finite (since the database instance  $I$  is finite), but it is potentially exponentially large [17].

Buneman et al. [13] defined the why-provenance of an output tuple  $t$  in the result  $Q(I)$  as a special *subset* of the set of witnesses called the *witness basis*. The witnesses of the basis depend on  $Q$ ; thus, each basis’s size is bounded by the size of  $Q$ . The witnesses of the basis exclude tuples that are irrelevant to  $t$  being produced by  $Q$ , and thus the basis tends to be very small compared to the set of all possible witnesses [17]. The witnesses are also *minimal*, in the sense that if one tuple is removed from one of these witnesses, it cannot produce the output.

id	name	why-provenance
$o_1$	Dopamine Receptors	$\{\{f_1, c2f_1, c_1\}, \{f_1, c2f_2, c_2\}\}$
$o_2$	YANK Family	$\{\{f_4, c2f_4, c_1\}\}$

Table 3: Result of a SQL query applied on the database of Table 1 with the why-provenance of the corresponding results.

In a sense, each witness in the witness basis captures one possible way in which the query can generate the output. To better understand this, consider the example in Table 3, where each tuple in the result of query Q1 is annotated with its why-provenance.

The why-provenance of output tuple  $o_2$  has only one witness, which coincides with its lineage. This happens because there is only one way this output tuple can be produced, i.e., for tuple  $f_4$  to be joined with  $c2f_4$  and  $c_1$ . On the other hand,  $o_1$  has a witness basis with of two witnesses, since there are two possible ways in which the query can generate  $o_1$ . One possibility is that

id	name	how-provenance
$o_1$	Dopamine Receptors	$f_1 \cdot c2f_1 \cdot c_1 + f_1 \cdot c2f_2 \cdot c_2$
$o_2$	YANK Family	$f_4 \cdot c2f_4 \cdot c_1$

Table 4: Result of the example SQL query Q1 with the corresponding how-provenances of the output tuples annotated.

495  $f_1$  is joined with  $c2f_1$  and  $c_1$  (the first witness), and the second possibility  
 496 is that  $f_1$  is joined with  $c2f_2$  and  $c_2$  (the second witness). This means that  
 497 to generate  $o_1$ , it is sufficient that only one of the two witnesses is present in  
 498 the input database.

#### 499 4.3. How-Provenance

500 While why-provenance describes the source tuples that witness an output  
 501 tuple in the result of the query, it leaves out information about how the source  
 502 tuples are used. How-provenance was therefore defined in [30] to capture this  
 503 information using a *semiring* algebraic structure, and is a form of provenance  
 504 that takes the form of a *polynomial*.

505 The key idea in Green et al. [30] is to use the two operators  $+$  and  $\cdot$  to  
 506 represent two basic transformations that source tuples undergo as a result  
 507 of applying a relational query to a database [17]. Two tuples may either be  
 508 joined together, as an effect of a join (represented with the  $\cdot$  operator) or  
 509 merged via union or projection (represented with the  $+$  operator).

510 Table 4 shows a simple example in which the two output tuples of our  
 511 running example are annotated with their respective how-provenances. Tuple  
 512  $o_2$  was produced through the join among the input tuples  $f_4$ ,  $c2f_4$ , and  $c_1$ .  
 513 The three provenance tokens are, therefore “multiplied” together. The case of  
 514  $o_1$  is slightly more complex. This tuple, as already discussed, can be obtained  
 515 through two different joins. The two monomials composing the polynomial  
 516 represent these two alternatives. They correspond, in a way, to the witnesses  
 517 of the why-provenance of  $o_1$ . The  $+$  operator represents the fact that the two  
 518 monomials describe alternative derivations. The output tuple is the result  
 519 of a merge of two distinct tuples after the projection on the attribute **name**.  
 520 This merge is due to the fact that the result of a relational algebra expression  
 521 is always a *set* of tuples, which corresponds to the presence of the **DISTINCT**  
 522 operator in an SQL query. This simple example gives the basic idea behind  
 523 how-provenance and how it allows us to track the operations that produced  
 524 an output tuple.

Provenance polynomials may also have monomials whose exponents and/or coefficients are greater than one, for example,  $3f_1 \cdot c2f_1 \cdot c_1 + f_1 \cdot c2f_2^3 \cdot c_2^3$ . This is a polynomial of a tuple produced by a query where the result of the join between the tuples  $f_1$ ,  $c2f_1$ , and  $c_1$  is produced three times and then merged (e.g. as the result of a union), and the tuples  $c2f_2$  and  $c_2$  are used three times in the operation described by the second monomial (e.g., with nested queries).

## 5. Credit Distribution and Distribution Strategies

We now give formal definitions of data credit and Data Credit Distribution (DCD), and present three different Distribution Strategies (DSs) based on the forms of provenance discussed earlier: Lineage-based DS, Why-Provenance-based DS, and How-Provenance-based DS. We also show how these strategies distribute credit in the IUPHAR example discussed earlier.

### 5.1. Data Credit and Data Credit Distribution

Given a database instance  $I$ , a *recipient of credit* is a unit of information within  $I$ . In the case of relational databases, recipients may be (i) the whole database; (ii) a table; (iii) a tuple; or (iv) an attribute.

*Data credit* is a value  $k \in \mathbb{R}_{>0}$ . Every recipient in a database is annotated with a quantity of credit as a proxy for its importance. In this paper, we focus on *tuples* as recipients of credit.

Given a *distribution strategy* (DS), *Data Credit Distribution* (DCD) takes a database instance  $I$ , quantity of credit  $k$ , and query  $Q$  over  $I$ , and splits  $k$  among the recipients of credit in  $I$ .

In the following, we use the notation in Cheney et al. [17]: Given an instance  $I$ , a *tuple location*  $(R, t)$  is a tuple  $t$  in relation  $R$ . With reference to the running example,  $(\text{family}, \langle f_1, \text{Dopamine Receptors}, \text{gpcr} \rangle)$  is the tuple location of the first tuple in the **family** relation. The set of all tuple locations in  $I$  is called *TupleLoc*. We use this to formally define DCD at the *tuple level*.

**Definition 5.1. Tuple Level Data Credit Distribution (DCD) [24]**  
 Given a query  $Q$  over  $I$  and  $k \in \mathbb{R}_{>0}$ , DCD is defined by the function  $f_{I,Q} : \text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$  such that  $f_{I,Q}(t, k) = h$  where  $0 \leq h \leq k$  and  $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$ . The function  $f_{I,Q}$  is the distribution strategy (DS).

As we can see, the DS is a function that annotates each tuple in the database with a real value, which is a fraction of the given quantity  $k$ . The only constraint is that the sum of the credit annotations on tuples must be  $k$ , i.e. that no credit is generated or destroyed during the distribution. Given  $I$  and  $Q$ , many different DSs may be defined as long as they sum up to  $k$ .

In what follows, we use information provided by data provenance to define distribution functions. For simplicity, we assume that the credit  $k$  is distributed equally across the set of output tuples (i.e. the result of a query), and discuss how the credit of one output tuple  $o$ ,  $k_o$ , is distributed across the instance  $I$ .

## 5.2. A Lineage-based Distribution Strategy

In the lineage-based distribution strategy, each tuple in the output of a query distributes credit equally to each input tuple that appears in its lineage. More formally:

**Definition 5.2.** *Lineage-based Distribution Strategy [24]*

Let  $I$  be a database instance,  $Q$  a query over  $I$ ,  $o \in Q(I)$  an output tuple and  $k_o$  the credit associated to  $o$ . Let  $L$  be the lineage of  $o$  and  $t$  be a tuple in  $I$ , then  $t$  receives credit equal to:

$$f_{I,Q}(t, k_o) = \begin{cases} 0 & \text{if } t \notin L \\ \frac{k_o}{|L|} & \text{if } t \in L \end{cases}$$

Note that lineage-based DS distributes credit only to input tuples that have a role in creating  $o$  by the query  $Q$ , and that each receives an equal share of credit via  $o$ . Thus, the more tuples in a lineage set, the less credit each tuple receives.

As an example, consider the output tuples of Table 2, and assume that each output tuple has credit  $k_o = 1$ . The lineage of the first tuple,  $o_1$ , is the set  $\{f_1, c2f_1, c_1, c2f_2, c_2\}$ . Therefore, each tuple in this set receives credit  $1/5$ . The other tuples of the database receive zero credit. The lineage of the second output tuple is  $\{f_4, c2f_4, c_1\}$ , therefore each of these tuples receives credit  $1/3$ .

At the end of the process, tuples  $f_1$ ,  $c2f_2$  and  $c_2$  each receive credit  $1/5$ , tuples  $f_4$  and  $c2f_4$  receive  $1/3$ , while tuple  $c_1$  receives  $8/15$ . Note that if a tuple appears in more than one lineage set, then it will accumulate credit from the distribution associated with each one of these sets, implying that

586 it has a more significant role in the context  $Q$ , as is the case with  $c_1$  in this  
 587 example.

588 Not all of the tuples in the lineage of an output tuple are necessary to be  
 589 present at the same time for the output tuple to appear in the query results.  
 590 For example, if the database only had the set of tuples  $\{f_1, c2f_1, c_1\}$  or the set  
 591  $\{f_1, c2f_2, c_2\}$ , the existence of  $o_1$  would still be guaranteed. In other words,  
 592 while  $f_1$  is always needed for  $o_1$  to appear in the output, only one of the sets  
 593 of tuples  $\{c2f_1, c_1\}$  and  $\{c2f_2, c_2\}$  is required. One could therefore argue that  
 594 it would be more fair for  $f_1$  to receive more credit than the other four tuples,  
 595 given its role in producing  $o_1$ .

596 This highlights one limitation of the lineage-based DS: while able to find  
 597 all and only the relevant tuples of the output, it does not distinguish the  
 598 *importance* of tuples in the query computations. We therefore present two  
 599 other, more sophisticated, forms of distribution strategies based on why- and  
 600 how-provenance.

### 601 5.3. A Why-Provenance-Based Distribution Strategy

602 The distribution strategy based on why-provenance first equally distributes  
 603 the credit  $k_o$  among the witnesses of the witness basis for  $o$ , and then equally  
 604 divides the credit of a witness among the tuples in the witness. Since a tuple  
 605 may appear in more than one witness, it will receive more than one portion  
 606 of credit from the same distribution. More formally:

607 **Definition 5.3.** *Why-Provenance-based Distribution Strategy*

608 *Let  $I$  be a database instance,  $Q$  a query over  $I$ ,  $o \in Q(I)$  an output tuple  
 609 and  $k_o$  the total credit associated to  $o$ . Let  $\mathcal{W} = \text{Why}(Q, I, o)$  be the witness  
 610 basis of  $o$  according to  $Q$  and  $I$ , and  $W \in \mathcal{W}$  be a witness.*

*Then tuple  $t$  in  $I$  receives credit equal to:*

$$f_{I,Q}(t, k_o) = \frac{k_o}{|\mathcal{W}|} \sum_{W \in \gamma(\mathcal{W}, t)} \frac{1}{|W|}$$

where  $\gamma$  is a function which returns all witnesses  $W$  in which  $t$  appears:

$$\gamma(\mathcal{W}, t) = \{W \in \mathcal{W} : t \in W\}$$

611 Figure 4 shows the distribution of credit with why-provenance-based DS  
 612 for tuple  $o_1$ . The credit is first equally divided between the two witnesses, so  
 613 that both receive credit  $1/2$ . The credit is then further divided among the

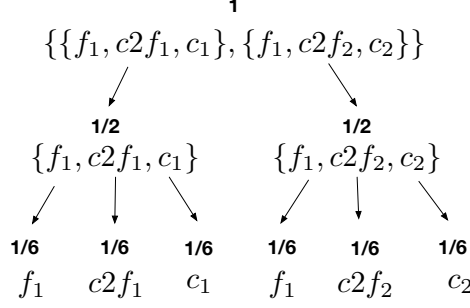


Figure 4: Distribution of credit using why-provenance-based DS for tuple  $o_1$ .

614 tuples in each witness. Since each witness has three tuples, each tuple in a  
 615 witness receives  $1/6$  of credit. At the end of the distribution,  $f_1$  receives a  
 616 total credit of  $1/3$ , and the other tuples receive  $1/6$  each. This distribution  
 617 better reflects the role of  $f_1$  in the generation of  $o_1$  since, as discussed earlier,  
 618 it is the only mandatory tuple for  $o_1$  to appear in the output; only one of the  
 619 two other pairs of tuples are necessary for  $o_1$  to appear in the result.

620 This example illustrates that why-provenance can better reward input  
 621 tuples depending on their role. Tuples that appear in more than one witness  
 622 are rewarded more than others.

#### 623 5.4. A How-Provenance Based Distribution Strategy

624 How-provenance conveys more information than why-provenance since  
 625 it not only captures what tuples are relevant to the output and in which  
 626 combination, but also how they are used. The “how” is captured through  
 627 the provenance polynomials.

628 The how-provenance-based DS therefore first distributes the credit to the  
 629 monomials of the polynomial accordingly to the weight represented by their  
 630 coefficients, then to the tuples of each monomial accordingly to the weights  
 631 represented by their exponents.

632 To define the DS more formally, we introduce some notation and illustrate  
 633 it using the provenance polynomial  $\mathcal{H}$  shown in Figure 5. This notation is  
 634 also reported for easy reference in Table 5.

635 We call  $c$  the function that, given a polynomial, returns the sum of the  
 636 coefficients of the polynomial; thus  $c(\mathcal{H}) = 3 + 1 = 4$ . We call  $e$  the function  
 637 that, given a monomial, returns the sum of its exponents, thus  $c(M_2) =$   
 638  $1 + 3 + 3 = 7$   $mc$  is the function that takes as input a monomial and returns  
 639 its coefficient.  $te$  is a function that takes as input a tuple and a monomial,

Table 5: Notations used in Definition 5.4.

$\mathcal{H}$	provenance polynomial
$M_i$	a monomial in $\mathcal{H}$
$t_j$	a tuple in $M_i$
$c(\mathcal{H})$	sum of $\mathcal{H}$ 's coefficients
$e(M_i)$	sum of $M_i$ 's exponents
$mc(M_i)$	$M_i$ 's coefficient
$te(t_j, M_i)$	exponent of $t_j$ in $M_i$
$\gamma(t_j, \mathcal{H})$	set of monomials in $\mathcal{H}$ containing $t_j$

$$\begin{aligned}
\mathcal{H} &= \underbrace{3f_1 \cdot c2f_1 \cdot c_1}_{M_1} + \underbrace{f_1 \cdot c2f_2^3 \cdot c_2^3}_{M_2} \\
c(\mathcal{H}) &= 4 & e(M_2) &= 7 \\
mc(M_1) &= 3 & mc(M_2) &= 1 \\
te(c_2, M_2) &= 3 & \gamma(c_1, \mathcal{H}) &= \{M_1\} \\
\gamma(f_1, \mathcal{H}) &= \{M_1, M_2\}
\end{aligned}$$

Figure 5: Illustration of notation used to define the how-provenance based DS in Definition 5.4.

640 and returns the exponent of the tuple in the monomial, if present; thus  
641  $te(c_2, M_2) = 3$ . Finally,  $\gamma$  takes as input a tuple and the whole polynomial,  
642 and returns a set containing the monomials containing that tuple, if present  
643 in the polynomial; thus  $\gamma(f_1, \mathcal{H}) = \{M_1, M_2\}$ .

644 **Definition 5.4.** *How-Provenance-Based Distribution Strategy*

645 *Let  $I$  be a database instance,  $Q$  a query over  $I$ ,  $o \in Q(I)$  an output tuple,  $\mathcal{H}$*   
646 *be the provenance polynomial for  $o$ , and  $k_o$  the credit given to  $o$ . The credit*  
647 *given to tuple  $t$  in  $I$  is:*

$$f_{I,Q}(t, k_o) = \frac{k_o}{c(\mathcal{H})} \sum_{M \in \gamma(t, \mathcal{H})} mc(M) \frac{te(t, M)}{e(M)}$$

648 Going back to the example of Table 4, consider  $o_1$  with provenance poly-  
649 nomial  $f_1c2f_1c_1 + f_1c2f_2c_2$ . The how-provenance-based DS firstly divides the  
650 credit between the two monomials. Since the coefficients of each monomial  
651 are 1, the credit is split in half. If they were, for example, 1 and 2 respectively,  
652 1/3 of the credit would go to the first monomial, and 2/3 to the second. Since



id	name
$oxs_1$	Dopamine Receptors

lineage	why-provenance	how-provenance
$\{f_1, c2f_1, c_1, c2f_2, c_2\}$	$\{\{f_1, c2f_1, c_1\}, \{f_1, c2f_2, c_2\}\}$	$f_1^2 c2f_1 c_1 + f_1^2 c2f_2 c_2$

Table 6: Result of query Q2 applied on the database of Table 1 and its different provenances. The reported numbers are the credit distributed through the process.

in our example each variable has exponent 1, the credit is further divided  
equally among the three variables. Thus, at the end of the computation,  $f_1$   
receives 1/3, and the other tuples receive 1/6. Consider instead the example  
where the polynomial is  $f_1 c2f_1 c_1 + f_1 c2f_2 c_2$ . The first monomial receives 1/2  
of the credit, then  $f_1$  receives 1/4 of credit, and the other two tuples receive  
1/8.

In this specific example, the how-provenance-based DS has the same out-  
come as the one based on why-provenance. We therefore consider another  
query over GtoPdb, Q2, that asks for the families of type `gpcr` that have as  
contributor a researcher located in the UK:

```

Q2: SELECT DISTINCT F.name
FROM family as F JOIN
(SELECT DISTINCT f.name AS name
FROM family AS f JOIN contributor2family AS c2f ON f.id = c2f.family_id
JOIN contributor AS c ON c2f.contributor_id = c.id
WHERE c.country = "UK") AS R ON F.name = R.name
WHERE F.type = "gpcr"

```

The result of Q2 is shown in Table 6, and consists of one tuple, anno-  
tated with each of the three provenances. As can be seen, lineage and why-  
provenance are identical to those of the tuple  $o_1$  in the previous example.  
The how-provenance, however, is different since tuple  $f_1$  is used twice: first  
in the join of the inner query, and second in the join of the outer query. This  
information is lost in the first two forms of provenances since they are sets,  
but it is captured in how-provenance through the use of the operator ‘.’.

Figure 6 shows the differences between the three DS for the tuple  $o_1$  of  
Table 6. Subfigure 6.a uses lineage, sub-figure 6.b uses why-provenance, and  
sub-figure 6.c uses how-provenance. The DS based on the provenance poly-  
nomial gives credit 1/2 to  $f_1$ , and 1/8 to the other tuples. This is reasonable  
since Q2 relies on  $f_1$  even more than Q1 does. The distribution based on

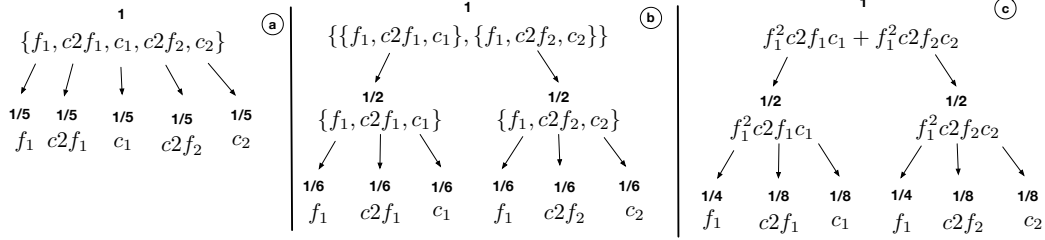


Figure 6: Comparison of different distributions strategies for tuple  $o_1$  produced by query Q2.

how-provenance can reward  $f_1$  more, showing that how-provenance is even more sensitive to the tuples' role in a query than why-provenance. This is a direct consequence of the fact that, as proven in [30], how-provenance is more general than why-provenance and lineage, in the sense that it contains more information.

## 6. Experimental Evaluation

To understand the trade-offs between these Distribution Strategies (DSs), we perform four sets of experiments using queries over target families presented on the GtoPdb website. The first set of experiments use real queries extracted from citations to GtoPdb published in the British Journal of Pharmacology. The second set uses synthetically produced provenance polynomials, corresponding to more complex queries, in order to better highlight the differences between the DSs. The third set of experiments considers the accrual of credit over time by the three strategies, again using synthetic queries. The fourth set of experiments shows how the DSs compare to traditional citations in giving credit to data curators using both real and synthetic queries.

All experiments were carried out on a MacBook Pro with a 2.4 GHz processor Intel Core i5 quad-core and 8 GB of memory at 2133 MHz. Code was written in Java, supported by a PostgreSQL database.<sup>11</sup>

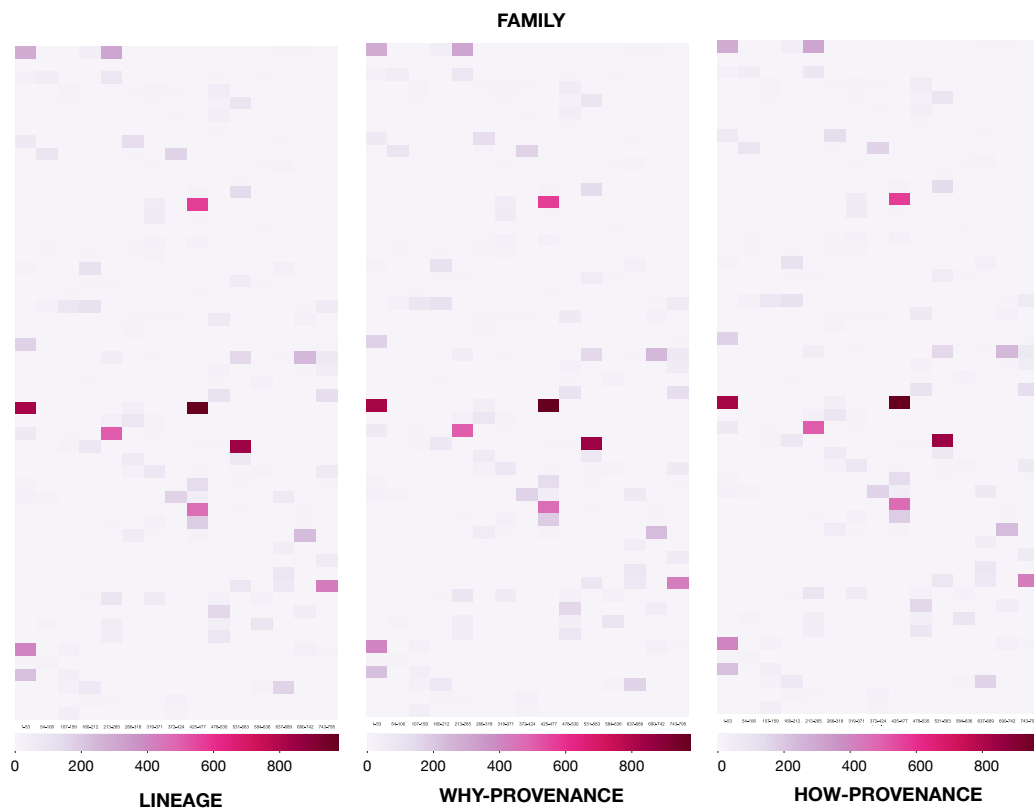


Figure 7: Comparison of three DS on the same table `family` using the distribution given by the queries retrieved from papers. Each cell is a tuple.

## 6.1. Real-world queries

Examples of real queries are drawn from papers published in the British Journal of Pharmacology (BJP).<sup>12</sup> Each time a paper in this journal cites a webpage from GtoPdb, it reports the URL of the page. From this URL, the query used to obtain the webpage data can be determined. We considered all 889 papers in BJCP citing the IUPHAR/BPS Guide to pharmacology [31] as of October 2020, and extracted all webpage URLs to GtoPdb contained

<sup>11</sup>For purposes of reproducibility, the code we used for our experiments and all queries are available here: [https://bitbucket.org/dennis\\_dosso/credit\\_distribution\\_project](https://bitbucket.org/dennis_dosso/credit_distribution_project).

<sup>12</sup><https://bpspubs.onlinelibrary.wiley.com>

709 within the paper.<sup>13</sup>

710 The queries that we inferred are those used to build target family web-  
711 pages within GtoPdb. An example was given in Figure 3, where we show  
712 how the structure of the “Adenosine receptors” family can be mapped into  
713 queries over the underlying database. In GtoPdb, all target family pages  
714 share a similar structure; the only difference is that individual sections, such  
715 as “contributors” or “further readings”, may be absent. Therefore, the same  
716 queries can be used to build all of the target family pages by changing the  
717 family id used in the query (for example, in Figure 3, it is 3). Note that  
718 the queries are fairly simple SQL queries, and fall into a class called “select-  
719 project-join” or “SPJ” queries. A total of more than 12K different queries  
720 were built in this way. Without loss of generality, we give each tuple in the  
721 output of a query a credit of 1.

722 *Results.* Figure 7 shows the heat-maps obtained by the distribution of credit  
723 according to the three different DS on one of the tables in the underlying  
724 database, **family**, which is often joined with other tables in the database to  
725 build the webpages. Each cell in a heat-map represents a tuple of the **family**  
726 table and the color indicates the amount of credit attributed to such tuple.  
727 It can be seen that the result of credit distribution over **family** is the same  
728 for all three strategies. The same result is also obtained with the other tables  
729 of the database used by the queries shown in Figure 3.

730 The reason why credit distribution is the same for all three strategies  
731 is that the queries are all simple SPJ queries, which use each table only  
732 once and do joins on key attributes. Under these conditions, each tuple of  
733 the output presents: (i) a how-provenance that is a single monomial with  
734 coefficient 1 and exponent 1 in each variable; (ii) a why-provenance with  
735 only one witness; and (iii) a lineage that coincides with the witness in the  
736 basis. Hence, for these queries, the three DSs behave in the same way: credit  
737 is uniformly distributed among the tuples present in each provenance.

738 To illustrate this, consider one of the queries in Figure 3 which is used to  
739 build the output webpage:

```
740 Q3: SELECT c.first_names, c.surname  
741 FROM contributor2family AS cf JOIN contributor AS c ON
```

---

<sup>13</sup>The IUPHAR/BPS Guide is a journal that describes the structure and evolution of GtoPdb. At the time of writing, it had received more than 1200 citations on Google Scholar.

```

742     cf.contributor_id = c.contributor_id
743     WHERE f.family_id = 3

```

744 Q3 returned 10 tuples from the version of GtoPdb used. The first tu-  
745 ple, <Bertil B., Fredholm>, has  $c_{939} \cdot c_{2f_{496}}$  as its provenance polynomial.  
746  $c_{939}$  represents the provenance token of a tuple in `contributor`, and  $c_{2f_{496}}$   
747 the provenance token of a tuple in table `contributor2family`. The why-  
748 provenance of this tuple is  $\{\{c_{939}, c_{2f_{496}}\}\}$  and its lineage is  $\{c_{939}, c_{2f_{496}}\}$ .  
749 Therefore, the credit assigned to these tuples is 1/2 using all three DS. This  
750 happens for all the tuples in the output of each query of GtoPdb, thus making  
751 the distributions equivalent over all outputs.

752 However, this is not the case with more complex queries. As we showed  
753 in the previous section, when two or more tuples are merged as a result of  
754 a projection or union, the credit distributions will differ between the three  
755 strategies.

## 756 6.2. Synthetic queries

757 To simulate synthetic queries, we randomly generated provenance poly-  
758 nomials in which the coefficients and exponents could be greater than 1. The  
759 queries involve three GtoPdb tables: `family`, `contributor2family`, and  
760 `contributor`. The polynomials were generated as follows (in particular, ev-  
761 ery time we write “randomly”, we mean using a uniform distribution): first,  
762 the number of monomials composing the polynomial is decided choosing ran-  
763 domly a number between 1 and 6. Then, we randomly choose a tuple from  
764 the tables `family`, one from the table `contributor2family` and one from  
765 table `contributor`, that are used as the monomial’s variables. Again, ran-  
766 domly, we choose a coefficient for this monomial (between 1 and 3) and an  
767 exponent for each tuple (between 1 and 4). For the next monomial, then, we  
768 decide if we want to keep the same tuple from the table `family` as first tuple  
769 of the new monomial. To do so, we generate a random number between 0 and  
770 1. If the number is above 0.2, we change the family tuple. An example can  
771 be found in Figure 8, which shows a sample synthetic provenance polynomial  
772 (the how-provenance) and the corresponding why-provenance and lineage ex-  
773 pressions. The resulting credit distribution for each DS is shown after the  
774 provenance expression.

775 As an example of how the distribution strategies behave with these syn-  
776 thetic queries, consider tuple  $f_5$  in Figure 8. This tuple receives the highest  
777 quantity of credit using lineage-based distribution, and less credit using why-

**How-provenance:**  $3f_1^3c_2f_1^2c_1^2 + 2f_1c_2f_2^3c_2^3 + 4f_5c_2f_{17}^4c_{18}^3$

**Credit distribution:**

$$f_1 = \frac{59}{315}, f_5 = \frac{1}{18}, c_2f_1 = \frac{2}{21}, c_2f_2 = \frac{2}{15}, c_2f_{17} = \frac{2}{9}, c_1 = \frac{2}{21}, c_2 = \frac{2}{15}, c_{18} = \frac{1}{6}$$

**Why-provenance:**  $\{\{f_1, c_2f_1, c_1\}, \{f_1, c_2f_2, c_2\}, \{f_5, c_2f_{17}, c_{18}\}\}$

**Credit distribution:**

$$f_1 = \frac{2}{9}, f_5 = \frac{1}{9}, c_2f_1 = \frac{1}{9}, c_2f_2 = \frac{1}{9}, c_2f_{17} = \frac{1}{9}, c_1 = \frac{1}{9}, c_2 = \frac{1}{9}, c_{18} = \frac{1}{9}$$

**Lineage:**  $\{f_1, f_5, c_2f_1, c_1, c_2f_1, c_2f_2, c_2f_{17}, c_1, c_2, c_{18}\}$

**Credit distribution:**

$$f_1 = \frac{1}{8}, f_5 = \frac{1}{8}, c_2f_1 = \frac{1}{8}, c_2f_2 = \frac{1}{8}, c_2f_{17} = \frac{1}{8}, c_1 = \frac{1}{8}, c_2 = \frac{1}{8}, c_{18} = \frac{1}{8}$$

Figure 8: Sample synthetic provenance polynomial (how-provenance) and corresponding why-provenance and lineage expressions with credit distributions.

778 and how-provenance because more information is available about the role of  
 779 the tuple in the overall computation. Generally speaking, the more complex  
 780 the distribution (the most complex being how-provenance), the more credit  
 781 is given to tuples which are more frequently used, and thus have a higher  
 782 impact in producing the output tuple.

783 Despite being synthetic, these provenance polynomials represent realistic  
 784 queries. The polynomials can be obtained by any nested query with join and  
 785 union operations that use the same tuple multiple times (in which case the  
 786 exponents are bigger than 1), and the same combination of operations more  
 787 than once (in which case the coefficients of monomials are bigger than 1).

788 *Results.* The results of credit distribution on the **family** table using 10K  
 789 randomly generated synthetic provenance polynomials are shown in Figure  
 790 9. We set the maximum value in the heat maps to the highest value reached  
 791 by a tuple in all three distributions (i.e., 9.4).

792 As can be seen, the three strategies generate significantly different credit  
 793 distributions indicated by the varying hues. However, there is a certain  
 794 amount of consistency between them in that tuples which are highly rewarded  
 795 by one strategy are also highly rewarded by the others. This shows that the  
 796 three DSs consistently reward certain tuples more than others.

797 Note that lineage-based DS gives the least credit to tuples in the **family**  
 798 table, indicated by an overall lighter hue. This is because the DS distributes  
 799 credit equally to all tuples appearing in the lineage. Since these queries also  
 800 use two other tables, credit is distributed to tuples in those tables.

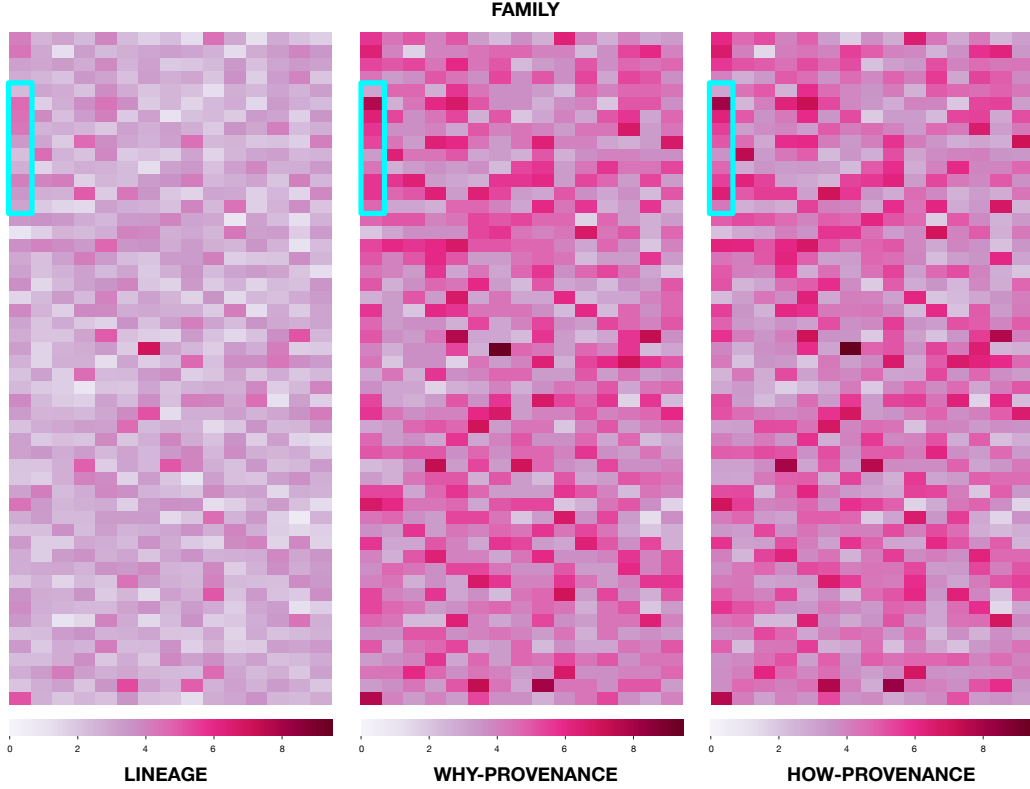


Figure 9: Comparison of three DS on the same table **family** after the distribution computed using 10K synthetic and randomly generated provenance polynomials. The tuples in the blue rectangles are used as example in the discussion connected to Figure 10.

801 Moving to why-provenance-based DS, we see that more credit is given to  
 802 tuples in the **family** table than with the previous strategy. This is because  
 803 the DS considers the different ways that a tuple is used, e.g. in joins with  
 804 other tuples. If the same tuple is present in more than one witness, it will  
 805 draw more credit and take it from other tuples in the witness basis. In this  
 806 case, tuples in **family** drew more credit, taking it from tuples in the other  
 807 two tables, due to the role that **family** tuples played in the queries that were  
 808 executed. We note that the lineage-based DS gives an average credit of 2.79  
 809 to each tuple in the table, while the DS based on why-provenance assigns  
 810 4.18. Moreover, lineage distributed a total of about 2200 units of credit to  
 811 the table, while the other DS assigned more than 3300 units. That is, the DS  
 812 based on why-provenance assigns on average 50% more credit to the **family**

813 table than the strategy based on lineage.

814 Finally, consider the how-provenance-based DS heat-map. As with why-  
815 provenance, more credit is typically given to tuples in **family** compared to  
816 lineage-based DS since it recognizes the role of these tuples in the queries, and  
817 the overall hue is deeper. The two distributions appear similar, although on  
818 closer inspection, slight differences between the two distributions can be seen.  
819 This is because how-provenance also considers the frequency with which tu-  
820 ples are used, not only the ways in which they are used. Therefore, although  
821 the overall distribution is similar, there are small differences due to the pres-  
822 ence of exponents and coefficients in the provenance polynomials, influencing  
823 the distribution of credit.

824 To better understand this difference, in the next subsection we consider  
825 the accrual of credit over time. In doing so, we will focus on the ten tuples  
826 shown within the large light blue rectangles in Figure 10. Each small rect-  
827 angle within a large blue rectangle is a tuple, and we number them from 1  
828 (top) to ten (bottom).

### 829 6.3. Credit accrual over time

830 Since credit accrues over time, we simulate the passage of time by varying  
831 the number of queries executed, and look at the “snapshots” of credit for each  
832 of the strategies using synthetic queries. The results are shown in Figure 10.

833 In this figure, four groups of heat-maps are shown. Each group represents  
834 a “snapshot” taken after 1K, 2K, 5K and 10K provenance polynomials have  
835 been considered for credit distribution. The ten tuples in each heat-map are  
836 those highlighted in the light blue boxes of Figure 9 from the **family** table.

837 The queries used are the same as the experiment of the previous section.  
838 The range of credit in each map goes from 0 (no credit) to 8 (the maximum  
839 quantity of credit reached on one of the tuples of the considered window at  
840 the “snapshot” with 10K queries). The color hue of the legend, as can be  
841 seen, still ranges from 0 to 9.5.

842 By the end of 1K queries, credit differentials between tuples as well as  
843 between strategies can be seen. For example, tuple 4 is usually rewarded the  
844 most credit by all three strategies. However, it receives the highest quantity of  
845 credit from the why-provenance-based strategy. Tuple 3 receives the highest  
846 quantity of credit overall with how-provenance. This trend continues to the  
847 end of 2k queries. By the end of 5k queries, tuple 2 emerges with the highest  
848 value of credit for why- and how-provenance, a position which is strengthened  
849 by the end of 10k queries. This is because tuple 2 is used several times





Figure 10: Comparison of the distribution of credit performed by the three DSs on a subset of 10 tuples taken from the **family** table, simulating the passing of time. The number at the top of each group of heat-maps represents the number of queries.

850 within queries being executed, which is rewarded strongly by why- and how-  
851 provenance but not taken into account in lineage.

852 While the relative value of credit “positions” of tuples within a DS strat-  
853 egy depends on what queries are being executed, the important thing to  
854 notice is the difference between the DSs over time: Overall, lineage gives far  
855 less credit to tuples in the **family** table than the other two strategies since  
856 credit is shared with tuples in other tables. However, the why- and how-  
857 provenance-based strategies recognize the more important role being played  
858 by the **Family** tuples than those in the other tables. The differences between  
859 the why- and how-provenance-based DSs are also relatively minor (about  
860 plus or minus 0.2 out of 9.5) in most cases. However, there are certain situ-

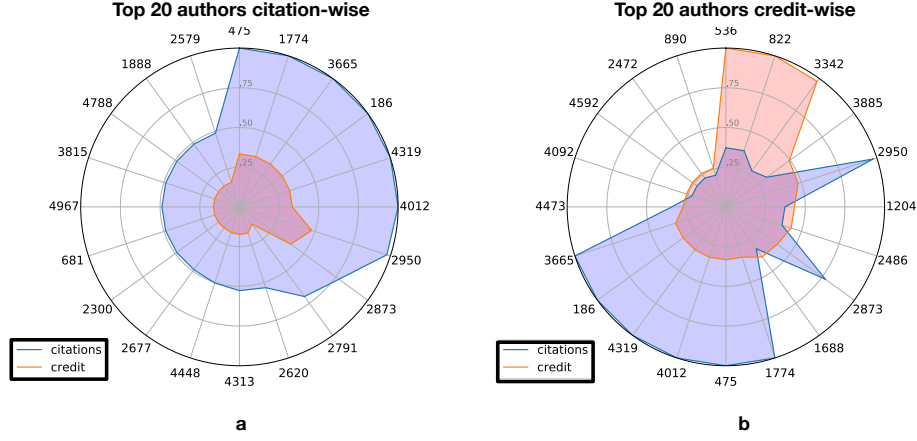


Figure 11: Radars presenting the top 20 authors citation-wise and credit wise, together with their (normalized between 0 and 1) values of citations and credit.

ations in which the role of a tuple is particularly critical in a query, and in this case the difference in the value of credit assigned is notably higher for how-provenance. An example of this can be seen in tuple 9 of the 10k group of Figure 10.

To sum up, the DS based on lineage is sufficient to highlight which tuples in the database are used by a query, and distributes credit equally to these tuples. The resulting distribution rewards tuples that are used by more queries, but does not reward how many times tuples are used in the same query. However, a DS based on why- or how-provenance may be better if the queries are complex, since they reward more tuples that have a critical role in generating the output. In particular, these two DSs may be useful for finding “hotspots” in the database based on the role of tuples, with the how-provenance-based DS being preferable if a higher sensitivity to the role of a tuple in queries is required.

#### 6.4. Credit vs Citations

In the last set of experiments, we compare traditional citations to the proposed credit distribution strategies to see the difference in reward for data authors and curators. Using both real-world and synthetic queries, we distribute credit to the authors responsible for the data under the different strategies. Our results show that credit rewards authors of data that is cited fewer times, but that has a higher impact on the query results.

882 To do so, we need to identify a set of authors and queries that cite data  
883 curated by them. Considering GtoPdb, each target family page has a list  
884 of curators, representing the people who are co-creators and curators of the  
885 data comprising the page. This list can be obtained using the last query  
886 shown in Figure 3. Each time a target family page is cited, we assign one  
887 *citation* to each author associated with the page. The authors also receive  
888 *credit* in the amount assigned to the data used by the query to construct the  
889 webpage, equally divided between the authors of the webpage.

890 *Results: Real-world queries.* As described in Section 6.1, we consider real-  
891 world queries taken from papers published in the BJP which reference web-  
892 pages in GtoPdb. Since for these queries there is no difference in the distri-  
893 bution of credit between the three DS, only one value for credit is used.

894 The results are shown in the radar plots of Figure 11, in which each  
895 number on the outer circle (e.g. 475, 1774 and 3665) represents an author  
896 (id) and the blue (red) line represents the normalized value of credit generated  
897 by citations (credit), respectively. The first radar plot, Figure 11.a, shows the  
898 top 20 authors in terms of *citations*, ordered in a clockwise direction, whereas  
899 Figure 11.b orders the authors based on *credit*. Comparing the author ids  
900 used in the outer circles of these two plots, it can immediately be seen that  
901 the “top authors” are very different using these two metrics, although there  
902 is some overlap (for example, authors 1774, 475, and 4012).

903 Diving a bit deeper to focus on the red and blue areas in each of the plots  
904 reveals that there is a significance difference between citations and credit:  
905 The top 20 authors in terms of citations do not have the highest values  
906 of credit (Figure 11.a). Conversely, the authors with the highest values of  
907 credit do not necessarily have a large number of citations (Figure 11.b). For  
908 example, author 536 has the highest value of credit, but is not even in the  
909 top 20 authors in terms of citations. This means that authors like 536, 822,  
910 and 3342 in Figure 11.b receive much more credit from their relatively few  
911 citations than authors like 475, who receives the largest number of citations.  
912 That is, the data underlying certain webpages is more “valuable” in terms  
913 of credit than a citation to the webpage.

914 The reason for the difference between citation and credit is partly due to  
915 the experimental setup: Each output tuple carries a credit of 1, and there  
916 can be many tuples used to generate a webpage. Thus a webpage that is  
917 created from more tuples will have a higher credit value than one created  
918 from fewer tuples. Furthermore, authors who collaborated with fewer people

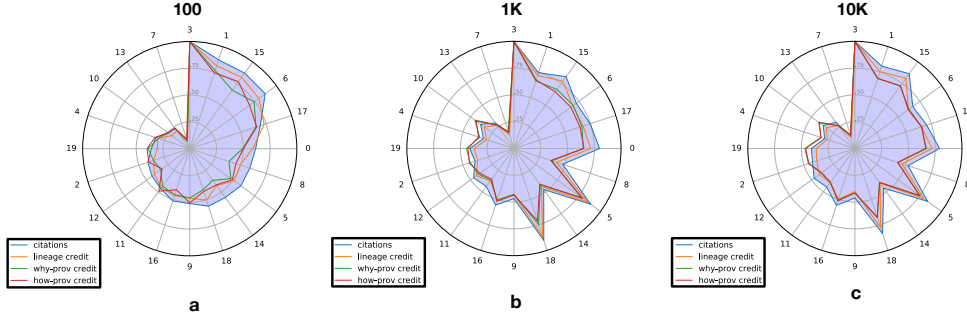


Figure 12: Radars presenting the 20 synthetic authors with corresponding citation and quantities of credit distributed through the 3 DS (all values normalized between 0 and 1) through different numbers of polynomials (respectively, 100, 1K and 10K). The order is the one defined by figure 1, i.e. descending order of citations obtained from 100 polynomials.

will receive a biggest share of the equally divided credit. However, all authors will receive a citation of one.

Credit distribution therefore rewards authors differently than traditional citations: An author who has curated larger quantities of cited data and collaborated with fewer co-authors, will receive larger quantities of credit. Thus, credit rewards them for their larger contribution to the database.

*Results: Synthetic queries.* We produced 100, 1K, and 10K batches of synthetic polynomials, as described in Section 6.2, and distributed credit through them to data. Since these polynomials are created by randomly selecting tuples from three tables, they usually correspond to a large set of authors who in reality did not collaborate. To make the size of the author set more realistic, we therefore created 20 synthetic authors, and randomly assigned one author to blocks of consecutive tuples in the database, with the size of each block varying between 10 and 40, to simulate different quantities of work performed by an author. Every time an author appears as curator of one or more tuples used in a polynomial, we assign them one citation. They also receive three kinds of credit, each one using a different DS.

Figure 12 shows three radar plots, one for each batch of synthetic polynomials. Each plot shows the top 20 authors in terms of citations (hence the authors and clockwise ordering is the same in each of the plots), and additionally shows the the normalized values of citation (blue line), lineage-based credit (yellow line), why-provenance-based credit (green line) and how-provenance-based (red line). As can be seen, given the synthetic nature of

the queries, the correlation between the number of citations and the quantity of credit assigned to the authors appears to be a much stronger than with the real-world queries of Figure 11. In fact, for Figure 12.a the linear correlation between the citation number and all three types of credit is always above 0.95 with p values in the order of  $1e-11$ . The credit distributed via lineage is closest to the number of citations (a linear correlation of 0.98, p value of  $6.15e-16$  in Figure 12.a), while the other two types of credit behave slightly differently (a linear correlation of around 0.95 in both cases in Figure 12.a). Similar observations can be made for Figure 12.b and 12.c.

What these figures show is that, in certain cases, authors who do not have a large number of citations receive more credit than others, as for example author 11 in Figure 12.a or author 19 in Figures 12.b and 12.c, especially when credit is distributed using how-provenance. This again shows how credit gives a different perspective on the role of data and authors by going beyond the limitations of traditional citations.

It is worth noting that, when scaling up to  $1K$  and  $10K$  polynomials, the credit distributions via why-provenance and how-provenance become almost identical (the linear correlation for the values of Figure 12.c is more than 0.99 with a p-value of  $1.32e-32$ ). This is consistent with what we observed in Figure 9.

## 7. Discussion

We note that, in our experiments, we always assumed that the credit carried by an output tuple is 1. Thus, each tuple in the output has equal importance. This in general may not be true, since different tuples in the output may have different weight, depending on the context of the citation. For example, data that is fundamental for the results of a paper may have more credit than data being cited as a reference. *Credit generation*, i.e. the process by which the credit of the output tuples is decided, is research problem with its own dignity and complexities, and we did not face it in this paper.

From the point of view of the model, even when the credit of the output tuples is different than 1, nothing needs to change in the models presented here, since they were defined for a generic value  $k$ . We note that, if the quantity of credit carried by an output tuple changes, as a consequence the final distribution will change, since certain tuples will be more “impactful” (i.e., distribute more credit) than others. With different quantities of credit,

978 therefore, new results, different from the ones obtained in the previous sec-  
979 tions, may be found. These results will depend on the nature of the context  
980 and the quantity of credit being considered.

## 981 8. Conclusions and Future Work

982 This paper defines two new distribution strategies based on why- and  
983 how-provenance, and compares them against the lineage-based distribution  
984 strategy defined in [24]. The first, why-provenance-based DS, uses the con-  
985 cept of a witness, and gives more credit to tuples that appear in more than  
986 one witness. In this way, tuples that are more important to the query and are  
987 used in different ways are rewarded more. The second, how-provenance-based  
988 DS, considers the frequency with which a tuple or combination of tuples is  
989 used in the query through the information contained in a provenance poly-  
990 nomial. In this case, the how-provenance-based DS is more sensitive than  
991 the why-provenance-based DS to the role and importance of tuples.

992 To show the differences between the three DSs, we performed extensive  
993 experiments based on GtoPdb, a curated scientific relational database, using  
994 both real and synthetic queries. In the first set of experiments, we used select-  
995 project-join (SPJ) queries extracted from citations to webpages in GtoPdb  
996 found in papers published in the British Journal of Pharmacology. Using  
997 these “real” queries, we distributed credit to tuples in different tables of the  
998 database, highlighting tuples that were more frequently used. We showed  
999 that, with these queries, the three strategies produce the same distribution.  
1000 This is because the SPJ queries were fairly simple, and did not use self-joins.  
1001 Therefore the formulas underlying the different DSs had the same output.

1002 In the second set of experiments, we synthetically produced more com-  
1003 plex provenance polynomials, corresponding to more complex queries, that  
1004 resulted in exponents and coefficients in the provenance polynomials that  
1005 were greater than (or equal to) 1. These experiments highlighted the differ-  
1006 ences between the three DSs. While the DS based on lineage rewards all the  
1007 tuples used by a query equally, the strategy based on why-provenance gives  
1008 more credit to tuples that are more critical to the query. In particular, why-  
1009 provenance consider the different ways in which a tuple is used in a query.  
1010 How-provenance is even more sensitive to the tuple’s role: it also considers  
1011 the frequency with which a tuple or a set of tuples is used.

1012 In the third set of experiments, we showed how the differences between  
1013 the DS are compounded over time, i.e. when more and more queries are

processed by the system.

In the fourth set of experiments we compared traditional citations to authors to the credit accrued to them via the DSs. We showed how, in both real-world and synthetic scenarios, credit rewards authors who contribute/curate data that has the highest impact, and therefore receives the biggest quantity of credit, and not necessarily the data with the highest citation count. In this sense, credit appears to be an useful new measure to discover data and their corresponding curators that have a high impact in the research world, even when they are cited few times or do not appear at all in the data that are cited (i.e. the case of data used to build the output of a query but that is not visualized in the output itself).

In future work, we plan to explore different strategies to generate and distribute credit. In this paper we assumed that each output tuple carries credit 1. In more sophisticated scenarios we can employ different strategies to compute credit, that reflect the importance of cited data. Also, other, and more sophisticated strategies could also be used to decide how credit is distributed between the authors, beyond the uniform distribution used here, in a way to reflect the work performed by them on the cited data.

We will also explore new applications for credit over relational databases. One example is *data pricing*, which gives a price to a query submitted by a user who wants to buy the produced information. Currently, a commonly strategy used for data pricing is based on query rewriting: A database stores a set of views with their price. When a new query arrives, the system rewrites it using the stored views to obtain a query price, a process that can be computationally expensive. We plan to distribute credit through carefully planned and representative queries, and use credit information to define a new, faster, and potentially more flexible pricing function.

Another application is *data reduction* [41], which addresses the problem of reducing the vast – and rapidly expanding – amount of data that is being produced.

Data credit can also address this problem, by helping find “hotspots” and “coldspots” of data. A hotspot is data in a database (e.g. a tuple) with a high quantity of credit, which is therefore valuable for the set of queries that execute frequently over the data and distribute the credit. On the other hand, a coldspot is data with a low quantity of credit, which is therefore considered less important and could be deleted or moved to cheaper and/or less efficient memory.

## Acknowledgement

The work was partially supported by the ExaMode project, as part of the European Union H2020 program under Grant Agreement no. 825292.

## References

- [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bernstein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L., Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Kossmann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo, T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo, A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suciu, D. (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–53.
- [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating data citation in citedb. *PVLDB*, 10(12):1881–1884.
- [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y. (2018). Data citation: A new provenance challenge. *IEEE Data Eng. Bull.*, 41(1):27–38.
- [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An Introduction to the Joint Principles for Data Citation. *Bulletin of the Association for Information Science and Technology*, 41(3):43–45.
- [5] Baggerly, K. (2010). Disclose all data in publications. *Nature*, 467(7314):401–401.
- [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D., Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I., Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and Goble, C. A. (2013). Why linked data is not enough for scientists. *Future Gener. Comput. Syst.*, 29(2):599–611.
- [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.



- 1080 [8] Bertin-Mahieux, T., Ellis, D., Whitman, B., and Lamere, P. (2011). The  
1081 million song dataset. In *Proceedings of the 12th International Conference*  
1082 *on Music Information Retrieval (ISMIR 2011)*, pages 591–596.
- 1083 [9] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In  
1084 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*  
1085 *nication*, pages 93–116. De Gruyter Mouton.
- 1086 [10] Buneman, P. (2006). How to cite curated databases and how to make  
1087 them citable. In *18th International Conference on Scientific and Statistical*  
1088 *Database Management, SSDBM*, pages 195–203. IEEE Computer Society.
- 1089 [11] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D.,  
1090 Pawson, A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn’t  
1091 working, and what to do about it. *Database J. Biol. Databases Curation*,  
1092 2020.
- 1093 [12] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation  
1094 is a computational problem. *Commun. ACM*, 59(9):50–57.
- 1095 [13] Buneman, P., Khanna, S., and Tan, W. C. (2001). Why and where: A  
1096 characterization of data provenance. In *Database Theory - ICDT 2001,*  
1097 *8th International Conference*, pages 316–330.
- 1098 [14] Buneman, P. and Silvello, G. (2010). A rule-based citation system for  
1099 structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.
- 1100 [15] Callaghan, S., Donegan, S., Pepler, S., Thorley, M., Cunningham, N.,  
1101 Kirsch, P., Ault, L., Bell, P., Bowie, R., Leadbetter, A. M., Lowry,  
1102 R. K., Moncoiffé, G., Harrison, K., Smith-Haddon, B., Weatherby, a.,  
1103 and Wright, D. (2012). Making Data a First Class Scientific Output:  
1104 Data Citation and Publication by NERC’s Environmental Data Centres.  
1105 *International Journal of Digital Curation*, 7(1):107–113.
- 1106 [16] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Jour-  
1107 nals: A Survey. *Journal of the Association for Information Science and*  
1108 *Technology*, 66(9):1747–1762.
- 1109 [17] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases:  
1110 Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–  
1111 474.

- 1112 [18] CODATA-ICSTI Task Group on Data Citation Standards and Practices  
1113 (2013). *Out of Cite, Out of Mind: The Current State of Practice, Policy,*  
1114 *and Technology for the Citation of Data*, volume 12.
- 1115 [19] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N.  
1116 (2019). Bringing citations and usage metrics together to make data count.  
1117 *Data Science Journal*, 18(1).
- 1118 [20] Cronin, B. (1984). *The Citation Process. The Role and Significance of*  
1119 *Citations in Scientific Communication*. London: Taylor Graham.
- 1120 [21] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evi-  
1121 dence of a structural shift in scholarly communication practices? *JASIST*,  
1122 52(7):558–569.
- 1123 [22] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of  
1124 view data in a warehousing environment. *ACM Trans. Database Syst.*,  
1125 25(2):179–227.
- 1126 [23] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model  
1127 for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on*  
1128 *Innovative Data Systems Research*. [www.cidrdb.org](http://www.cidrdb.org).
- 1129 [24] Dosso, D. and Silvello, G. (2020). Data credit distribution: A  
1130 new method to estimate databases impact. *Journal of Informetrics*,  
1131 14(4):101080.
- 1132 [25] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The vir-  
1133 tual atomic and molecular data centre (VAMDC) consortium. *Journal of*  
1134 *Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- 1135 [26] ESIP Data Preservation and Stewardship Committee (EDPSC) (2019).  
1136 Data citation guidelines for earth science data, version 2. Version 2, Earth  
1137 Science Information Partners.
- 1138 [27] Fang, H. (2018). A discussion of citations from the perspective of the  
1139 contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–  
1140 1520.
- 1141 [28] Garfield, E. (1999). Journal impact factor: a brief review. *Can. Med.*  
1142 *Assoc.*, 979-980.

- 1143 [29] Gößwein, B., Miksa, T., Rauber, A., and Wagner, W. (2019). Data  
1144 identification and process monitoring for reproducible earth observation  
1145 research. In *2019 15th International Conference on eScience (eScience)*,  
1146 pages 28–38. IEEE.
- 1147 [30] Green, T. J., Karvounarakis, G., and Tannen, V. (2007). Provenance  
1148 semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-*  
1149 *SIGART symposium on Principles of database systems*, pages 31–40. ACM.
- 1150 [31] Harding, S. D., Sharman, J. L., Faccenda, E., Southan, C., Pawson,  
1151 A. J., Ireland, S., Gray, A. J. G., Bruce, L., Alexander, S. P. H., Anderton,  
1152 S., Bryant, C., Davenport, A. P., Doerig, C., Fabbro, D., Levi-Schaffer, F.,  
1153 Spedding, M., Davies, J. A., and Nc-Iuphar (2018). The IUPHAR/BPS  
1154 guide to PHARMACOLOGY in 2018: updates and expansion to encom-  
1155 pass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Re-*  
1156 *search*, 46(Database-Issue):D1091–D1106.
- 1157 [32] Hartley, J. (2017). Authors and their citations: a point of view. *Scien-*  
1158 *tometrics*, 110(2):1081–1084.
- 1159 [33] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a  
1160 transformed scientific method.
- 1161 [34] Honor, L. B., Haselgrove, C., Frazier, J. A., and Kennedy, D. N. (2016).  
1162 Data citation in neuroimaging: proposed best practices for data identifi-  
1163 cation and attribution. *Frontiers in neuroinformatics*, 10:34.
- 1164 [35] Joshi-Tope, G., Gillespie, M., Vastrik, I., D’Eustachio, P., Schmidt, E.,  
1165 de Bono, B., Jassal, B., Gopinath, G. R., Wu, G. R., Matthews, L., Lewis,  
1166 S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of bio-  
1167 logical pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- 1168 [36] Katz, D. (2014). Transitive credit as a means to address social and  
1169 technological concerns stemming from citation and attribution of digital  
1170 products. *Journal of Open Research Software*, 2(1).
- 1171 [37] Kosten, J. (2016). A classification of the use of research indicators.  
1172 *Scientometrics*, 108(1):457–464.

- 1173 [38] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S.  
1174 (2011). Citation and Peer Review of Data: Moving Towards Formal Data  
1175 Publication. *International Journal of Digital Curation*, 6(2):4–37.
- 1176 [39] Martone, M. (2014). Joint declaration of data citation principles.  
1177 *FORCE11. San Diego CA. Data Citation Synthesis Group*. [https://www.](https://www.force11.org/datacitationprinciples)  
1178 [force11.org/datacitationprinciples](https://www.force11.org/datacitationprinciples), online September 2020.
- 1179 [40] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation  
1180 counts and rankings of LIS faculty: Web of science versus scopus and  
1181 google scholar. *Journal of the american society for information science*  
1182 *and technology*, 58(13):2105–2125.
- 1183 [41] Milo, T. (2019). Getting rid of data. *Journal of Data and Information*  
1184 *Quality (JDIQ)*, 12(1):1–7.
- 1185 [42] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D.,  
1186 Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G.,  
1187 Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff,  
1188 D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,  
1189 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson,  
1190 E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C.,  
1191 Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers,  
1192 E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research  
1193 culture. *Science*, 348(6242):1422–1425.
- 1194 [43] Parsons, M. A., Duerr, R. E., and Jones, M. B. (2019). The history and  
1195 future of data citation in practice. *Data Science Journal*, 18(1).
- 1196 [44] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J.  
1197 (2016). Research data explored: An extended analysis of citations and  
1198 altmetrics. *Scientometrics*, 107(2):723–744.
- 1199 [45] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic,  
1200 large databases: Model and reference implementation. In *Proceedings of*  
1201 *the 2013 IEEE International Conference on Big Data, 6-9 October 2013,*  
1202 *Santa Clara, CA, USA*, pages 307–312.
- 1203 [46] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identi-  
1204 fication of Reproducible Subsets for Data Citation, Sharing and Re-Use.

- 1205 *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue*  
1206 *on Data Citation*, 12(1):6–15.
- 1207 [47] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data  
1208 citation of evolving data: Recommendations of the working group on data  
1209 citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- 1210 [48] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*  
1211 *Sci. Technol.*, 69(1):6–20.
- 1212 [49] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data  
1213 provenance in e-science. *SIGMOD Record*, 34(3):31–36.
- 1214 [50] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-  
1215 spective. In of Sciences’ Board on Research Data, N. A. and Information,  
1216 editors, *Report from Developing Data Attribution and Citation Practices*  
1217 *and Standards: An International Symposium and Workshop*, pages 177–  
1218 178. National Academies Press: Washington DC.
- 1219 [51] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,  
1220 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,  
1221 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data  
1222 management and stewardship. *Scientific data*, 3.
- 1223 [52] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data  
1224 citation: Giving credit where credit is due. In *Proceedings of the 2018*  
1225 *International Conference on Management of Data, SIGMOD*, pages 99–  
1226 114.
- 1227 [53] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to  
1228 scientific datasets using article citation networks. *Journal of Informetrics*,  
1229 14(2).
- 1230 [54] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of  
1231 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.
- 1232 [55] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for  
1233 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*  
1234 *Molecular Spectroscopy*, 327:122–137.