

# 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 three new distribution strategies, the first two based on two forms of data provenance, why-provenance and how-provenance, and the third based on the concept of responsibility.

Using these distribution strategies we show how credit can highlight frequently used database areas and how it can 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.

## 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 [21, 22, 60]. Citations are used, among other things, to decide on tenure, promotion, hiring, and funding of grants for researchers [23, 36, 41, 44].

Science and research are increasingly digital, and there are numerous curated databases that are at the core of scientific research efforts [13]. It is therefore generally accepted that data must be cited and citable [16, 42], and that data citations should contribute to the scientific reputation of researchers, scientists, data curators, and creators [4, 55]. It is also accepted that data citations should be counted alongside of traditional citations, and contribute to bibliometrics indicators [7, 49].

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 [21]. 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 [53]. Initiatives such as FORCE11 and ESIP (Earth Science Information Partners) have collaborated on data and software citation principles and guidelines [29]. Other examples are the National Science Foundation (NSF), and the National Institute of Health (NIH) in the US [53].

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 [52]. 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 [48].

A central problem in the data citation process is how to attribute credit to data creators and curators [12]. 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 [10, 31]. 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) [30, 40, 59] 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

---

<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.

- 69 a quantity of credit to  $p_2$ , and its authors, that is *proportional* to the  
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 [13, 15, 50]. 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].

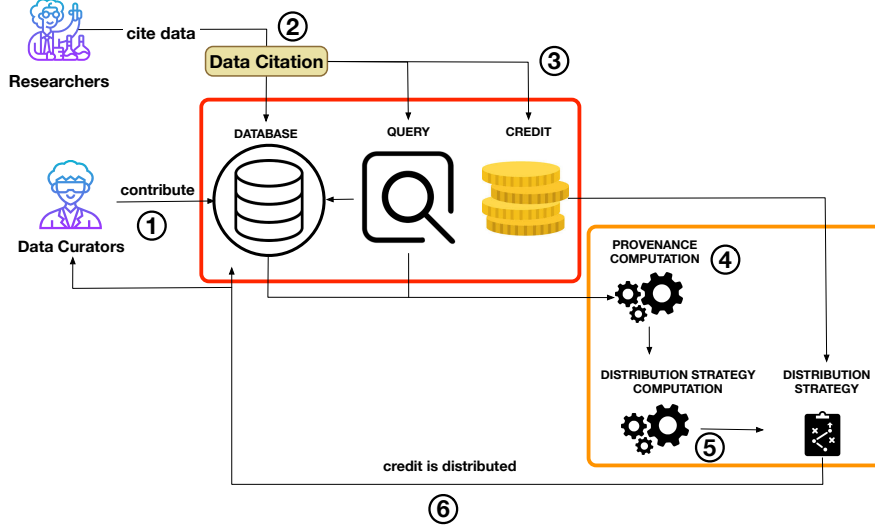


Figure 1: Overview of the credit distribution pipeline.

form of metadata that describes the generation process undertaken by  $Q$ , and the data used in  $I$  to generate the output [18]. Many different notions of provenance have been proposed in the literature for data in database management systems [14, 24, 33], describing different kinds of relationships between data in the input and the output of a query. As reported in [18], 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. Also, we consider the notions of causality and responsibility, that are built on top of provenance. In the following, for simplicity of exposition, when we say data provenance we also consider responsibility even though it is more of a form of enrichment of the information carried with lineage.

**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

120 describe four CDS, each one based on a different form of provenance.

121 **Step 6** Once the CDS is computed, it is used to distribute the given credit  
122  $k$  to the parts of the database that are responsible for the generation  
123 of  $Q(I)$ . Transitively, this credit is also divided and given to the corre-  
124 sponding authors of those data.

125 This paper expands our recent work in [26], which addressed the problem  
126 of how to reward data and data curators who are typically overlooked in  
127 current citation systems. In that work, we first defined the problem of DCD  
128 in relational databases, and proposed a viable Distribution Strategy (DS)  
129 based on *lineage*, which is the simplest form of *data provenance*. The lineage  
130 of a tuple  $t$  in the output  $Q(I)$  is defined as the set of all and only the tuples  
131 in the database instance  $I$  that are “relevant” to the production of  $t$ . The  
132 lineage-based strategy equally redistributes the credit  $k$  to the tuples in the  
133 lineage set, thus each tuple receives credit  $k/|L_t|$ , where  $L_t$  is the lineage set  
134 of  $t$ .

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

141 Therefore, in this paper, we expand the ideas in [26] by proposing three  
142 new DSs based on other forms of data provenance: why-provenance [14], how-  
143 provenance [33], and *responsibility* [45]. We use these form of provenances  
144 and the information that they carry to define new Distribution Strategies  
145 that have a different behavior with respect to the one of the lineage-based  
146 DS. We will show in the paper the formulas being define taking inspiration  
147 from the provenances and discuss their characteristics.

148 We compare them with the lineage-based solution, and discuss why one  
149 may be preferred to another depending on the application and its goals. In  
150 particular, we show that why-provenance, *responsibility* and how-provenance  
151 are more sensitive to the *role* of a tuple in a query, i.e. how many times the  
152 tuple is used and how it is used. The DSs based on why-provenance and  
153 *responsibility* give more reward to tuples that are essential to the production  
154 of the result set, whereas the DS based on how-provenance also takes into  
155 consideration the different ways that a tuple is used.

156 For evaluation, we use a well-known curated database, the IUPHAR/BPS<sup>3</sup>  
157 Guide to Pharmacology [35], also known as GtoPdb<sup>4</sup>, which contains ex-  
158 pertly curated information about diseases, drugs, cellular drug targets, and  
159 their mechanisms of action. We chose GtoPdb for two main reasons: (i) it  
160 is a widely-used and valuable curated relational database, (ii) many papers  
161 in the literature use, and cite its data (i.e., families, ligands, and receptors).  
162 Real queries used in papers can therefore be seen as data citations which, in  
163 turn, can be used to assign data credit.

164 We perform four sets of experiments. In the first one, real queries are ex-  
165 tracted from papers published in the British Journal of Pharmacology (BJP),  
166 that represent data citations to GtoPdb, and are used to distribute credit in  
167 the database using the three different provenance-based DSs. In the second  
168 and third experiment we analyze the behavior of the different DS when com-  
169 plex citation queries are employed. In the fourth set of experiments we use  
170 both real and synthetic queries to assess the difference between traditional  
171 citation and the notion of credit distribution in terms of rewarding those  
172 responsible for the data, e.g. data curators.

173 **Contributions** of this work include:

- 174 • Three new Distribution Strategies based on why- and how-provenance  
175 and on responsibility.
- 176 • An in-depth analysis of the effects of credit distribution on real-world  
177 curated data and of the differences between the three proposed Distri-  
178 bution Strategies.
- 179 • A comparison between the behavior of traditional citations and data  
180 credit in rewarding data curators.

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

---

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

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

## 187 2. Background

188 *Data in Research.* The world of research is rapidly transitioning towards the  
189 *fourth paradigm of science* [37], that is, data-intensive scientific discovery,  
190 where data are important for scientific advances as well as for traditional  
191 publications [6].

192 The scientific community is promoting an *open research culture* [47],  
193 founded on methods and tools to share, discover, and access experimental  
194 data. The community has identified the FAIR principles (Findable, Acces-  
195 sible, Interoperable, and Reusable) [57], that should be enforced by every  
196 database. In particular, data should be accessible from the articles, journals,  
197 and papers that cite or use them [21]. Aspects such as the need for the *repro-*  
198 *ducibility* of experiments through the used data; the *availability* of scientific  
199 data; the *connections* between data and the scientific results are all needed  
200 aspects for the fourth paradigm, and are all relevant to the domain of *data*  
201 *citation* [38].

202 *Data Citation: Principles and Motivations.* Data Citation principles were  
203 proposed in [20], and later summarized and endorsed by the Joint Declaration  
204 of Data Citation Principles (JDDCP) [43]. The principles are divided into  
205 two groups [53]. The first one contains principles concerning the role of  
206 data citation in scholarly and research activities such as the (i) *importance*  
207 of data (why data citation is important and why data should be considered  
208 as first-class citizens); (ii) *credit* and *attribution* to the creators and curators  
209 of the data; (iii) *evidence*; (iv) *verifiability*; and *interoperability*, with these  
210 last three requiring data citation methods to be flexible enough to operate  
211 through different communities. The second group defines the main guidelines  
212 to establish a data citation systems, and contains principles such as the (i)  
213 *unique identification* of the data being cited; (ii) *(open) access* to data; (iii)  
214 guarantee of *persistence* and *availability* of citations even after the lifespan  
215 of the cited entity; the (iv) *specificity* of a citation, i.e. it must lead to the  
216 data set originally cited.

217 It is possible to outline six main motivations for data citation [53]:

- 218 • *Data attribution:* identify the individuals that should be credited for  
219 data with variable granularity.
- 220 • *Data connection:* connect papers to the data being used.



- 221 • *Data Discovery*: citations helps to find data records and subsets that  
222 would be otherwise not findable via search engines.
- 223 • *Data Sharing*: share data obtained by researchers within the whole  
224 community.
- 225 • *Data Impact*: highlight the results obtained in writing papers using  
226 specific data, the frequency and modality data were used.
- 227 • *Reproducibility*: data citation greatly impacts the reproducibility of  
228 science [5]. Many authoritative journals ask to share data and provide  
229 valid methodologies to reproduce experiments.

### 230 2.1. Data Citation in Relational Databases

231 In this paper, we develop our methods and experiments on relational  
232 databases. RDBs have been the main target of data citation methods since  
233 the surge of the data-centric research paradigm. The RDA “Working Group  
234 on Data Citation: Making Dynamic Data Citable”<sup>5</sup> [51] has been working in  
235 the last years on large, dynamic, and changing datasets. The working group  
236 has finished the development of its guidelines and has now moved on into an  
237 adoption phase. The datasets considered by the Working Group are often  
238 relational.

239 In one of its most recent sessions [52], the Working Group (WG) on  
240 Data Citation reported that there are various implementations of its guide-  
241 lines for Data Citation on MySQL/Postgres relational databases. Some of  
242 these databases are: DEXHELPP<sup>6</sup> (Social Security Records); NERC (ARGO  
243 Global Array); EODC (Earth Observation Data Centre) [32]; LNEC (River  
244 dam monitoring); MDS (Million Song Database) [8]; CBMI<sup>7</sup> (Center for  
245 Biomedical Informatics); VMC (Vermont Monitoring Cooperative); CCA<sup>8</sup>  
246 (Climate Change Center Austria); VAMDC (Virtual Atomic and Molecular  
247 Data Center) [27, 61].

248 More examples of work on data citation in relational databases are [2, 13,  
249 25, 58]. The website <https://fairsharing.org/> keeps a long updated list

---

<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>

250 of curated and scientific databases (many of which are relational or graph-  
251 based) following FAIR guidelines. These databases are citable since they are  
252 compliant with the most recent guidelines, and they are in the vast majority  
253 of cases accessible via dynamically created Webpages. In all these databases  
254 is, therefore, possible to implement DCD on top of the existing infrastructures  
255 for citing data.

256 Data citation techniques are primarily applied to relational databases  
257 because of their diffusion and also because the portions of data that are to  
258 be cited are easily identified: the whole database, a relation, a tuple, or  
259 even an attribute. Many papers [2, 11, 13] consider more complex citable  
260 units, recognizing that often the *views* of a database are the ones to be cited.  
261 Generally, a *view* is a query on the database. To this end, [58] suggested  
262 decomposing the database in a set of views, where each view is associated  
263 with its citation.

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

- 265 1. A database cited as a whole, even though only parts of the databases  
266 are used in the papers or datasets. Alternatively, the so-called “data pa-  
267 pers” can be cited, being traditional papers that describe a database [17].  
268 In this case, all the credit from the citations goes to the database ad-  
269 ministrators or to the authors of the data papers.
- 270 2. Subsets of data, obtained by issuing queries to a database, are individ-  
271 ually cited. This is the solution adopted by the *Resource Data Alliance*  
272 (RDA) working group on Data Citation [51]. In this case, the credit  
273 generated from citations can be distributed among the contributors of  
274 the portions of data being cited, and/or to the database administrators.
- 275 3. The database is accessible via a series of Webpages that arrange the  
276 content of the database by topic or theme. Examples in the life science  
277 domain include the Reactome Pathway database [39], the GtoPdb [35],  
278 and the VAMDC [61]. Every single Webpage is unequivocally identifi-  
279 able and can be individually cited.

## 280 2.2. Data Credit

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

285 Katz [40] suggests the need for a *modified citation system* that includes  
286 the idea of *transient* and *fractional credit*, to be used by developers of research

287 products as software and data. In the paper two considerations are made:  
288 (i) research objects such as data and software are currently not formally  
289 rewarded or recognized by the community; (ii) even in traditional papers,  
290 the contribution of each author to the work is hard to understand, unless  
291 explicitly specified in the paper. This is even more true for data, where  
292 different groups of people work on the same database.

293 In [40] credit is defined as a “quantity” that describes the importance of a  
294 research entity, such as papers, software, or data, mentioned in a citation. It  
295 also proposed the idea of a *distribution* of credit from research entities, such as  
296 papers or data, to other research entities through citations. *Therefore, when*  
297 *talking about data credit, here we are focusing on two aspects of the topic:*  
298 *credit computation*, the process in which the quantity of credit generated by  
299 the citation is computed, and *credit distribution*, the process by which credit  
300 is distributed and assigned to the responsible entities that contributed to the  
301 generation of the data being cited. *In this paper we focus on the latter.*

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

307 Katz [40] further explores the idea of a *distribution* of credit from research  
308 entities (i.e., papers and data) to other research entities through citations  
309 that connect them. Thanks to traditional citations and now also to data  
310 citations, this distribution is finally possible, at least between papers and  
311 data. Some problems related to traditional citations can thus be solved by  
312 citations:

- 313 1. Credit rewards research entities that to date are not (formally) recog-  
314 nized (a goal shared with data citation).
- 315 2. Credit can reward authors *proportionally* to their role in generating the  
316 entity. The more an author contributes to a paper, the more credit is  
317 given to him. Zou and Peterson [60] work on something similar with  
318 their zp-index, which includes in its formulation the position (and thus  
319 the role) of a publication author to represent its impact in the work  
320 itself.
- 321 3. Credit can be *transitively* channeled through a chain of papers citing  
322 each other, thus enabling the rewarding of older papers that are no  
323 more cited, since other papers summarize or report their content but

324 are nevertheless crucial in a research area for the influence of their  
325 content.

326 Fang [30] presents a framework to distribute the credit generated by a  
327 paper to its authors and to the papers in its reference list in a transitive way.  
328 Let us consider the *citation graph* as the graph where the nodes are papers  
329 and the links are the citations among them. In this graph, every paper is  
330 a source of credit, which is then transferred to the neighboring nodes. The  
331 quantity of credit received by each cited paper depends on its impact/role  
332 in the citing paper. So far, this theoretical framework is limited to papers,  
333 but it can be easily extended to a citation graph including both papers and  
334 data.

335 Zeng et al. [59] proposes the first method to compute credit within a net-  
336 work of papers citing data. Adopting a network flow algorithm, they simulate  
337 a random walker to estimate a score for each dataset, leveraging real-world  
338 usage data to compute the credit. This is the first step towards an automatic  
339 credit computation procedure. This proposal is, however, limited to assign-  
340 ing credit to whole datasets, and it does not deal with the granularity of data.  
341 It does not work to assign credit to a single research entity within a dataset.  
342 Differently from Zeng et al. [59], we do not treat the credit computation  
343 process, but we focus on the distribution process.

### 344 2.3. Data Provenance

345 To distribute credit, we base our methods on *data provenance*. Data  
346 provenance is information that describes the origin and the process of cre-  
347 ation of data. It can also be seen as metadata pertaining to the derivation  
348 history of the data. It is particularly useful to help users to understand  
349 where data are coming from, and the process they went through. Data ci-  
350 tation and data provenance are closely linked [3] since both are forms of  
351 annotations on data retrieved through queries. Data provenance has been  
352 widely studied in different areas of data management. In this paper, we fo-  
353 cus on provenance for database management systems (DBMS). For further  
354 details on data provenance, please refer to surveys like [18] and [54].

355 Cheney et al. [18] presents four main types of data citation for DBMS: *lin-*  
356 *age* [24], *why-provenance* [14], *how-provenance* [33] and *where-provenance* [14].

357 Let us start with the first three provenances. Given a database instance  
358  $I$ , a query  $Q$ , and the result  $Q(D)$ , consider one tuple  $t$  of the output. Its  
359 provenance is information about its generation through the tuples of the

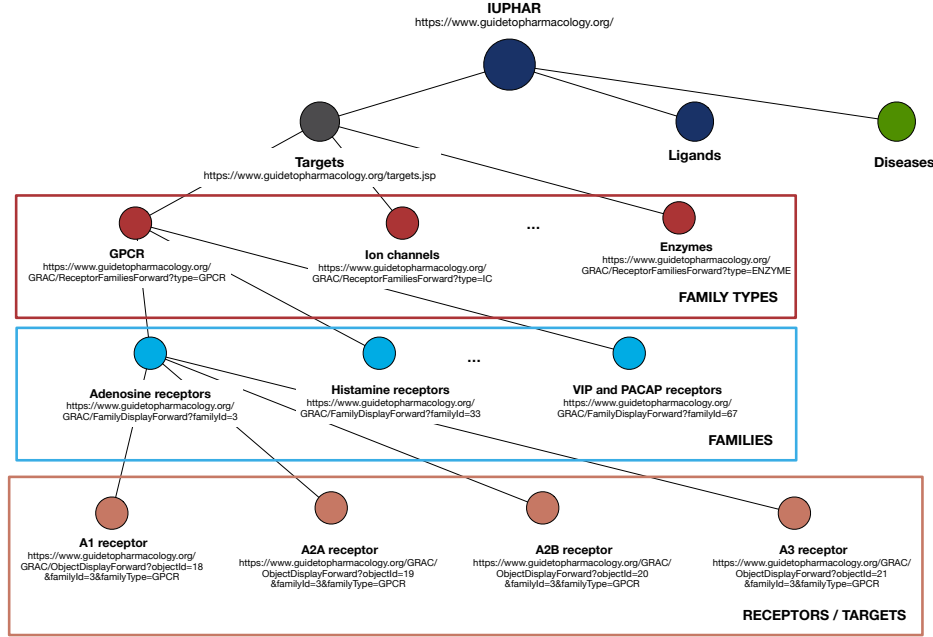


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

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*.

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

Also, here we consider the notions of causality and responsibility, as discussed in [45]. Causality is an enrichment of lineage, and is the attribution of a certain degree of importance to the tuples based on their role in the generation of the output. Responsibility is a value given to the tuples of the lineage to rank them based on their degree of causality (the more important the role of a tuple in generating the output, the higher its responsibility).

### 372 3. Use Case: GtoPdb

373 As use case we refer to the IUPHAR/BPS Guide to Pharmacology [35]  
374 or GtoPdb<sup>9</sup>. GtoPdb is a well-known and well structured scientific relational  
375 database that contains expertly curated information about diseases, drugs  
376 in clinical use, their cellular targets, and the mechanisms of action on the  
377 human body. It is curated and maintained by the GtoPdb Committee, and  
378 by 96 subcommittees, comprising 512 scientists collaborating with in-house  
379 curators who draw the information contained in the database from high-  
380 quality pharmacological and medicinal chemistry literature. Roughly 1000  
381 researchers from all over the world have contributed to the database, and the  
382 curators wanted to give recognition to these contributors. This led to some  
383 early work on data citation [11].

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

394 GtoPdb provides access to the webpages corresponding to all these nodes  
395 through URLs. The webpages corresponding to target families all present a  
396 similar structure, as shown in Figure 3 for the “Adenosine receptors” family.  
397 Each page has an *Overview*, a brief text describing the content of the page;  
398 a list of *Receptors* comprising the family; a section of *comments* about the  
399 family; the *References*, a list of the papers consulted by the curators of the  
400 page, similar to a reference list of a paper; the *further reading* list, reporting  
401 papers that an interested reader may want to consult to obtain more insight  
402 on the family; and a final section called *How to cite this family page*, con-  
403 taining text snippets useful to cite the specific page or the whole database.  
404 Figure 3 shows the SQL code that retrieves the information used to build the  
405 corresponding sections (apart from the References section). Therefore, each  
406 family page can be considered a full-fledged traditional publication, consist-

---

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

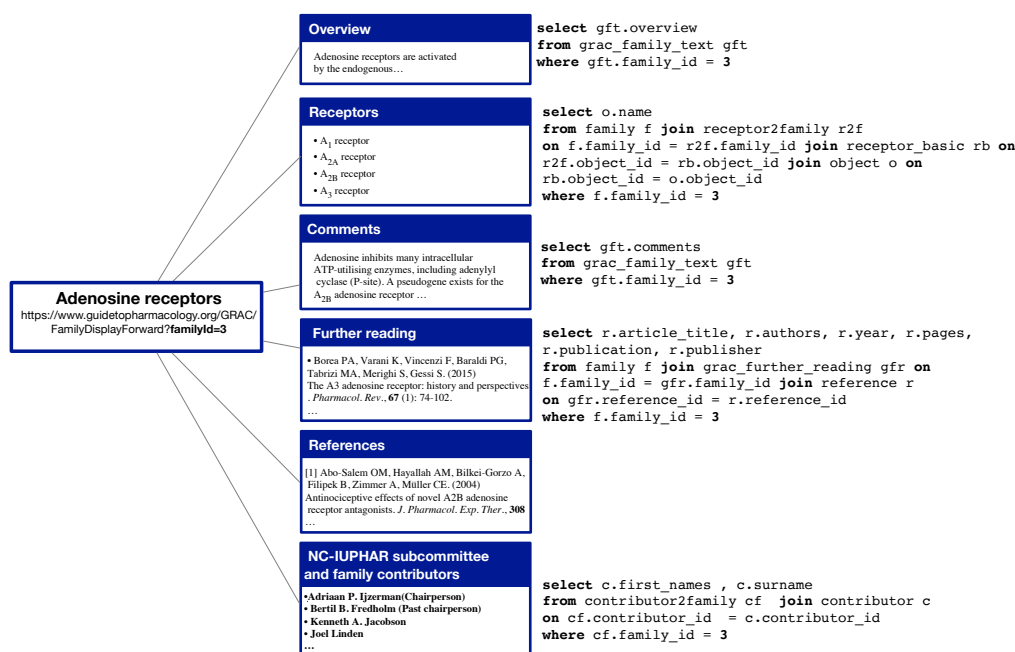


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

ing 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., [35]) 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 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.

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

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.

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.

## 4. Data Provenances

In this section, we present the three types of provenance used in this paper: lineage, why-provenance, and how-provenance. Also, we present the notions of Causality and Responsibility.

### 4.1. Lineage

Lineage was first introduced by Cui et al. [24]. Here we follow its definition as given by Cheney et al. [18]. Given a database instance  $I$ , the query  $Q$ , and



the result  $Q(D)$ , consider on tuple  $t$  in this output. Lineage is the simplest among the forms of provenance. It has been defined in different ways [18], but it can be thought as the set of all the tuples are used in some way by the query to produce the output tuple, i.e., the ones that are somehow *relevant* to its generation.

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 [14]. While why-provenance can be defined in many ways, we refer to [18], 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

468  $t$  in the query’s output, a *witness* is any sub-instance of the database that  
 469 produces  $t$ , i.e., a set that guarantees the existence of  $t$  in  $Q(D)$ . In particular,  
 470 the whole database and the lineage of  $t$  are both examples of witnesses of  $t$ .  
 471 Since the definition of witness allows for the presence of “irrelevant” tuples,  
 472 the set of all witnesses is finite (since the database instance  $I$  is finite), but  
 473 it is potentially exponentially large [18].

474 Buneman et al. [14] defined the why-provenance of an output tuple  $t$  in  
 475 the result  $Q(I)$  as a special *subset* of the set of witnesses called the *witness*  
 476 *basis*. The witnesses of the basis depend on  $Q$ ; thus, each basis’s size is  
 477 bounded by the size of  $Q$ . The witnesses of the basis exclude tuples that  
 478 are irrelevant to  $t$  being produced by  $Q$ , and thus the basis tends to be very  
 479 small compared to the set of all possible witnesses [18].

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.

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

484 The why-provenance of output tuple  $o_2$  has only one witness, which coin-  
 485 cides with its lineage. This happens because there is only one way this output  
 486 tuple can be produced, i.e., for tuple  $f_4$  to be joined with  $c2f_4$  and  $c_1$ . On  
 487 the other hand,  $o_1$  has a witness basis with of two witnesses, since there are  
 488 two possible ways in which the query can generate  $o_1$ . One possibility is that  
 489  $f_1$  is joined with  $c2f_1$  and  $c_1$  (the first witness), and the second possibility  
 490 is that  $f_1$  is joined with  $c2f_2$  and  $c_2$  (the second witness). This means that  
 491 to generate  $o_1$ , it is sufficient that only one of the two witnesses is present in  
 492 the input database.

### 493 4.3. How-Provenance

494 While why-provenance describes the source tuples that witness an output  
 495 tuple in the result of the query, it leaves out information about how the source  
 496 tuples are used. How-provenance was therefore defined in [33] to capture  
 497 this information using a *semiring* algebraic structure. It takes the form of

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.

498 a polynomial, called *provenance polynomial*, where the variables are taken  
 499 from the set  $X$  of identifiers of the tuples (provided that each tuple in  $I$  has  
 500 an identifier) and the coefficients are drawn from the set of natural numbers  
 501  $\mathbb{N}$ . This semiring therefore is commonly referred as  $\mathbb{N}(X)$  in the literature.

502 The key idea in Green et al. [33] is to use the two operators  $+$  and  $\cdot$  to  
 503 represent two basic transformations that source tuples undergo as a result  
 504 of applying a relational query to a database [18]. Two tuples may either be  
 505 joined together, as an effect of a join (represented with the  $\cdot$  operator) or  
 506 merged via union or projection (represented with the  $+$  operator).

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

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

#### 530 4.4. Causality and Responsibility

531 A formal study of causality was initiated in [19, 34] and later expanded  
 532 by Meliou et al. [45] to define the causes of answers and non-answers to  
 533 queries. Causality is, more precisely, related to the provenance of a query  
 534 result such as lineage and adds information to the one already provided by  
 535 the provenance.

536 In the following we define causality and responsibility as done in [45].  
 537 Differently from [45], we only focus on answers of a query, and not on  
 538 non-answers, since they are not relevant in the context of this paper. Let  
 539  $R_1, \dots, R_k$  be the relation names of a standard relational schema,  $D$  be a  
 540 database instance and  $q$  a conjunctive query. We also call  $D^n \subseteq D$  the set  
 541 of *endogenous tuples*, i.e. the tuples being actually considered to be possible  
 542 causes of a query output; while  $D^x = D - D^n$  is the set of *exogenous tuples*,  
 543 the tuples being considered external, unconcerned factors, thus deemed not  
 544 to be possible causes. This distinction between endogenous and exogenous  
 545 tuple is application dependent, and it can be done by the user at query time.  
 546 One example is with probabilistic databases with uncertain tuples, where  
 547 erroneous data may be contained. By considering these uncertain tuples  
 548 as part of the exogenous tuples dataset, we are factoring them out of the  
 549 computation of causality.

550 Then, given a tuple  $\bar{a}$  with the same arity as the query's answer, we write  
 551  $D \models q(\bar{a})$  when  $\bar{a}$  is an answer to  $q$  on  $D$ , and write  $D \not\models q(\bar{a})$  when  $\bar{a}$  is a  
 552 non-answer to  $q$  on  $D$ . Causality is defined as follows:

553 **Definition 4.1.** *Causality [45]*

554 *Let  $t \in D^n$  be an endogenous tuple, and  $\bar{a}$  a possible answer for  $q$ . Then:*

- 555 1.  *$t$  is called a counterfactual cause for  $\bar{a}$  in  $D$  if  $D \models q(\bar{a})$  and  $D - \{t\} \not\models$   
 556  $q(\bar{a})$*
- 557 2.  *$t \in D$  is called an actual cause for  $\bar{a}$  if there exists a set  $\Gamma \subseteq D^n$ , called  
 558 contingency for  $t$ , such that  $t$  is a counterfactual cause for  $\bar{a}$  in  $D - \Gamma$ .*

559  $t$  is a *counterfactual cause* if, by removing it from the database, we remove  
 560  $\bar{a}$  from the answer. Therefore, it can be fought as a tuple of the lineage which  
 561 is fundamental for the presence of  $\bar{a}$  in the answer. Vice-versa,  $t$  is an actual  
 562 cause if it is possible to find a contingency set of tuples such that, if that  
 563 set is removed, only then  $t$  becomes fundamental. In other words, when  $t$   
 564 is an actual cause, even if it was removed from the database,  $\bar{a}$  would still

id	name	responsibility
$o_1$	Dopamine Receptors	$f_1 : 1, c_2f_1 : 0.5, c_2f_2 : 0.5, c_1 : 0.5, c_2 : 0.5$
$o_2$	YANK Family	$f_4 : 1, c_2f_4 : 1, c_1 : 1$

Table 5: Result of the example SQL query Q1 with the corresponding responsibilities of the lineage tuples.

be present in the result set thanks to the contingency set. Checking the causality degree of tuples is NP-complete in general [28], but Meliou et al. [45] proved that the causality of conjunctive queries may be determined in PTIME.

The notion of *responsibility* was first defined in [19], and it measure the degree of causality as a function of the size of the smallest contingency set. It allows to rank the tuples in a lineage based on their degree of causality in generating the output.

**Definition 4.2.** *Responsibility [45] Let  $\bar{a}$  be an answer to a query  $q$ , and let  $t$  be a cause. The responsibility of  $t$  for the answer  $\bar{a}$  is:*

$$\rho_t = \frac{1}{1 + \min_{\Gamma} |\Gamma|}$$

where  $\Gamma$  ranges over all contingency sets for  $t$ .

As can be seen, a counterfactual cause will have the maximum responsibility of 1, while the bigger the minimum contingency of an actual cause, the smaller its responsibility since more tuples can still guarantee the presence of the answer  $\bar{a}$ .

While in general computing the responsibility is hard [19], Meliou et al. [45] showed that for each query without self-joins the responsibility is either computed in PTIME in the size of the database or checking if it has a responsibility below a given value is NP-hard.

As an example, consider Table 4, where we reported the tuples result of query Q1 together with the tuples of their lineage accompanied with their responsibility values. With output tuple  $o_1$ , the tuple  $f_1$  of the lineage is a counterfactual cause, since its contingency set is empty (when removed from the database,  $o_1$  disappears from the result set). Consequently, its responsibility is 1. On the other hand, the other tuples of the lineage are all actual causes.  $c_1$ , for example, has as minimal contingency set  $\{c_2f_2\}$ , and thus its responsibility is 0.5. For the output tuple  $o_2$ , all the tuples of the lineage are counterfactual causes, and thus they all have responsibility 1.

## 593 5. Credit Distribution and Distribution Strategies

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

### 599 5.1. Data Credit and Data Credit Distribution

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

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

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

609 In the following, we use the notation in Cheney et al. [18]: Given an  
 610 instance  $I$ , a *tuple location*  $(R, t)$  is a tuple  $t$  in relation  $R$ . With reference to  
 611 the running example,  $(\mathbf{family}, \langle f_1, \mathbf{Dopamine Receptors}, \mathbf{gpcr} \rangle)$  is the  
 612 tuple location of the first tuple in the **family** relation. The set of all tuple  
 613 locations in  $I$  is called *TupleLoc*. We use this to formally define DCD at the  
 614 *tuple level*.

#### 615 **Definition 5.1. Tuple Level Data Credit Distribution (DCD) [26]**

616 *Given a query  $Q$  over  $I$  and  $k \in \mathbb{R}_{>0}$ , DCD is defined by the function  $f_{I,Q} :$   
 617  $\text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$  such that  $f_{I,Q}(t, k) = h$  where  $0 \leq h \leq k$  and  
 618  $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$ . The function  $f_{I,Q}$  is the distribution strategy (DS).*

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

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

629 *5.2. A Lineage-based Distribution Strategy*

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

**Definition 5.2.** *Lineage-based Distribution Strategy [26]*

*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:*

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

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

637 As an example, consider the output tuples of Table 2, and assume that  
 638 each output tuple has credit  $k_o = 1$ . The lineage of the first tuple,  $o_1$ , is  
 639 the set  $\{f_1, c2f_1, c_1, c2f_2, c_2\}$ . Therefore, each tuple in this set receives credit  
 640  $1/5$ . The other tuples of the database receive zero credit. The lineage of the  
 641 second output tuple is  $\{f_4, c2f_4, c_1\}$ , therefore each of these tuples receives  
 642 credit  $1/3$ .

643 At the end of the process, tuples  $f_1$ ,  $c2f_2$  and  $c_2$  each receive credit  $1/5$ ,  
 644 tuples  $f_4$  and  $c2f_4$  receive  $1/3$ , while tuple  $c_1$  receives  $8/15$ . Note that if a  
 645 tuple appears in more than one lineage set, then it will accumulate credit  
 646 from the distribution associated with each one of these sets, implying that  
 647 it has a more significant role in the context  $Q$ , as is the case with  $c_1$  in this  
 648 example.

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

657 This highlights one limitation of the lineage-based DS: while able to find  
 658 all and only the relevant tuples of the output, it does not distinguish the

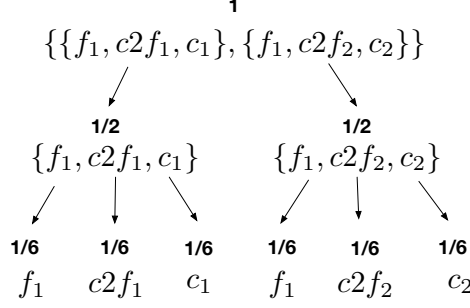


Figure 4: Distribution of credit using why-provenance-based DS for tuple  $o_1$ .

importance of tuples in the query computations. We therefore present two other, more sophisticated, forms of distribution strategies based on why- and how-provenance.

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

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

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

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

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

Figure 4 shows the distribution of credit with why-provenance-based DS for tuple  $o_1$ . The credit is first equally divided between the two witnesses, so that both receive credit  $1/2$ . The credit is then further divided among the tuples in each witness. Since each witness has three tuples, each tuple in a



Table 6: 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.

676 witness receives 1/6 of credit. At the end of the distribution,  $f_1$  receives a  
677 total credit of 1/3, and the other tuples receive 1/6 each. This distribution  
678 better reflects the role of  $f_1$  in the generation of  $o_1$  since, as discussed earlier,  
679 it is the only mandatory tuple for  $o_1$  to appear in the output; only one of the  
680 two other pairs of tuples are necessary for  $o_1$  to appear in the result.

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

#### 684 5.4. A How-Provenance Based Distribution Strategy

685 The how-provenance-based DS first distributes the credit to the mono-  
686 mials of the polynomial accordingly to the weight represented by their co-  
687 efficients, then to the tuples of each monomial accordingly to the weights  
688 represented by their exponents.

689 To define the DS more formally, we introduce some notation and illustrate  
690 it using the provenance polynomial  $\mathcal{H}$  shown in Figure 5. This notation is  
691 also reported for easy reference in Table 6.

692 We call  $c$  the function that, given a polynomial, returns the sum of its  
 693 coefficients; thus  $c(\mathcal{H}) = 3 + 1 = 4$ . We call  $e$  the function that, given a  
 694 monomial, returns the sum of its exponents, thus  $c(M_2) = 1 + 3 + 3 = 7$ .  
 695  $mc$  is the function that takes as input a monomial and returns its coefficient;  
 696 thus  $mc(M_1) = 3$ .  $te$  is a function that takes as input a tuple and a  
 697 monomial, and returns the exponent of the tuple in the monomial, if present;  
 698 thus  $te(c_2, M_2) = 3$ . Finally,  $\gamma$  takes as input a tuple and the whole polynomial,  
 699 and returns a set of monomials containing that tuple, if present in the  
 700 polynomial; thus  $\gamma(f_1, \mathcal{H}) = \{M_1, M_2\}$ .

701 **Definition 5.4.** *How-Provenance-Based Distribution Strategy*  
 702 Let  $I$  be a database instance,  $Q$  a query over  $I$ ,  $o \in Q(I)$  an output tuple,  $\mathcal{H}$   
 703 be the provenance polynomial for  $o$ , and  $k_o$  the credit given to  $o$ . The credit  
 704 given to tuple  $t$  in  $I$  is:

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

705 Going back to the example of Table 4, consider  $o_1$  with provenance polynomial  
 706  $f_1 c_2 f_1 c_1 + f_1 c_2 f_2 c_2$ . The how-provenance-based DS firstly divides the  
 707 credit between the two monomials. Since the coefficients of each monomial  
 708 are 1, the credit is split in half. If they were, for example, 1 and 2 respectively,  
 709 1/3 of the credit would go to the first monomial, and 2/3 to the second. Since  
 710 in our example each variable has exponent 1, the credit is further divided  
 711 equally among the three variables. Thus, at the end of the computation,  $f_1$   
 712 receives 1/3, and the other tuples receive 1/6. Consider instead the example  
 713 where the polynomial is  $f_1^2 c_2 f_1 c_1 + f_1^2 c_2 f_2 c_2$  and let us focus on the first  
 714 monomial. It receives 1/2 of the total credit, then  $f_1$  receives a portion of  
 715 credit equal to 1/4, while the other two tuples receive 1/8.

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

```
720 Q2: SELECT DISTINCT F.name
721 FROM family as F JOIN
722 (SELECT DISTINCT f.name AS name
723 FROM family AS f JOIN contributor2family AS c2f ON f.id = c2f.family_id
724 JOIN contributor AS c ON c2f.contributor_id = c.id
725 WHERE c.country = "UK") AS R ON F.name = R.name
726 WHERE F.type = "gpcr"
```

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 7: 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.

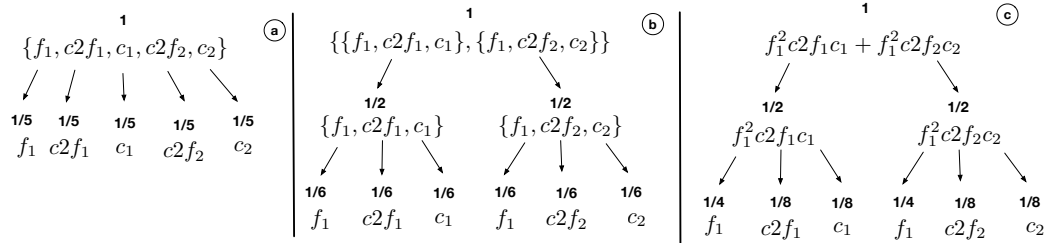


Figure 6: Comparison of different distributions strategies for tuple  $o_1$  produced by query Q2.

The result of Q2 is shown in Table 7, and consists of one tuple, annotated 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 7. Subfigure 7.a uses lineage, sub-figure 7.b uses why-provenance, and sub-figure 7.c uses how-provenance. The DS based on the provenance polynomial 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 how-provenance rewards  $f_1$  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 [33], how-provenance is more general than why-provenance and lineage, in the sense that it contains more information.

744

### 745 5.5. Responsibility-based Distribution Strategy

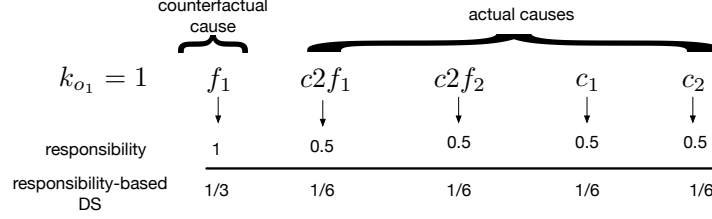


Figure 7: Example of distribution of credit using responsibility and normalized responsibility and the responsibility-based DS, assuming  $k_o = 1$ .

As we described in Section 4.3, causality and responsibility are not new forms of data provenance, but rather new information that is added to the already available lineage. Given the lineage of an output tuple  $o$ , it is first possible to distinguish the type of causality of each of its tuples, distinguishing between counterfactual and actual causes. Successively, it is possible to compute their responsibility, which, by itself, can be envisioned as a form of credit and assigned to the corresponding tuples.

One first option to define a distribution strategy using responsibility is to simply assign the responsibility as credit of the single tuple. Using the example of Table 5, in the case of output tuple  $o_1$ ,  $f_1$  receives credit 1, the other tuples credit 0.5. This strategy both generates the credit and gives it to the tuples.

However, we want a DS that is also a function of the input credit value  $k$  in order to be comparable with the other three strategies proposed so far. Therefore, we define a new DS based on responsibility that is a function of the quantity of credit  $k_o$  that assigns to each tuple of the lineage a portion of this credit weighted by its normalized quantity of responsibility. This function will give a bigger portion of credit to tuples that are higher in the responsibility ranking. Formally:

**Definition 5.5.** *Responsibility-based Distribution Strategy*

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

$$f_{I,Q}(t, k_o) = k_o \frac{\rho_t}{\sum_{t' \in L} \rho_{t'}}$$

Note that only the tuples that belong to the lineage will receive a quantity of credit  $> 0$ . The more important the tuple, i.e., the higher its responsibility,

772 the bigger the quantity of credit received.

773 Figure 7 shows the responsibility and the credit assigned to the tuples of  
774 the lineage of the output tuple  $o_1$  of Table 5. The only counterfactual tuple  
775  $f_1$  has responsibility 1 while the other have responsibility 0.5, as already  
776 discusses. Using the DS instead, and assuming that  $k_{o_1} = 1$ ,  $f_1$  receives  
777 credit  $1/3$ , while the others receive credit  $1/6$ . As we see, the DS in this case  
778 returns the same distribution obtained with why-provenance that was shown  
779 in Figure 6. This is not always the case though, as we show in the example  
780 of Section 6.2.

## 781 6. Experimental Evaluation

782 To understand the trade-offs between these Distribution Strategies (DSs),  
783 we perform four sets of experiments using queries over target families pre-  
784 sented on the GtoPdb website. The first set of experiments use real queries  
785 extracted from citations to GtoPdb published in the British Journal of Phar-  
786 macology. The second set uses synthetically produced provenance polyno-  
787 mials, corresponding to more complex queries, in order to better highlight  
788 the differences between the DSs. The third set of experiments considers  
789 the accrual of credit over time by the three strategies, again using synthetic  
790 queries. The fourth set of experiments shows how the DSs compare to tradi-  
791 tional citations in giving credit to data curators using both real and synthetic  
792 queries.

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

### 796 6.1. Real-world queries

797 Examples of real queries are drawn from papers published in the British  
798 Journal of Pharmacology (BJP).<sup>12</sup> Each time a paper in this journal cites a  
799 webpage from GtoPdb, it reports the URL of the page. From this URL, the  
800 query used to obtain the webpage data can be determined. We considered all  
801 889 papers in BJCP citing the IUPHAR/BPS Guide to pharmacology [35]

---

<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>

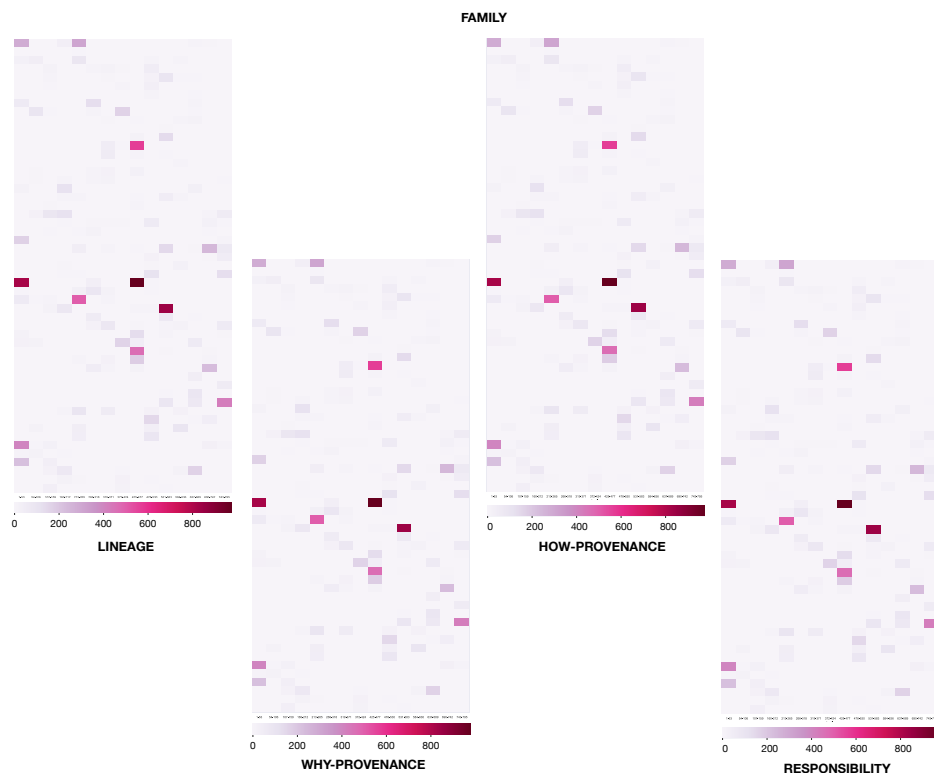


Figure 8: Comparison of four DS on the same table **family** using the distribution given by the queries retrieved from papers. Each cell is a tuple.

as of October 2020, and extracted all webpage URLs to GtoPdb contained within the paper.<sup>13</sup>

The queries that we inferred are those used to build target family webpages within GtoPdb. An example was given in Figure 3, where we show how the structure of the “Adenosine receptors” family can be mapped into queries over the underlying database. In GtoPdb, all target family pages share a similar structure; the only difference is that individual sections, such as “contributors” or “further readings”, may be absent. Therefore, the same queries can be used to build all of the target family pages by changing the family id used in the query (for example, in Figure 3, it is 3). Note that

<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.

the queries are fairly simple SQL queries, and fall into a class called “select-project-join” or “SPJ” queries. A total of more than 12K different queries were built in this way. Without loss of generality, we give each tuple in the output of a query a credit of 1.

*Results.* Figure 8 shows the heat-maps obtained by the distribution of credit according to the **four** different DS on one of the tables in the underlying database, **family**, which is often joined with other tables in the database to build the webpages. Each cell in a heat-map represents a tuple of the **family** table and the color indicates the amount of credit attributed to such tuple. It can be seen that the result of credit distribution over **family** is the same for all **four** strategies. The same result is also obtained with the other tables of the database used by the queries shown in Figure 3.

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

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

```
Q3: SELECT c.first_names, c.surname
FROM contributor2family AS cf JOIN contributor AS c ON
cf.contributor_id = c.contributor_id
WHERE f.family_id = 3
```

Q3 returned 10 tuples from the version of GtoPdb used. The first tuple, <Bertil B., Fredholm>, has  $c_{939} \cdot c_{2f_{496}}$  as its provenance polynomial.  $c_{939}$  represents the provenance token of a tuple in **contributor**, and  $c_{2f_{496}}$  the provenance token of a tuple in table **contributor2family**. The why-provenance of this tuple is  $\{\{c_{939}, c_{f_{496}}\}\}$ , its lineage is  $\{c_{939}, c_{2f_{496}}\}$ , both these tuples are counterfactual causes and have responsibility 1. Therefore, the credit assigned to these tuples is 1/2 using all four DS. This happens for all the tuples in the output of each query of GtoPdb, thus making the distributions equivalent over all outputs.

848 However, this is not the case with more complex queries. As we showed  
849 in the previous section, when two or more tuples are merged as a result of a  
850 projection or union, the credit distributions will differ between the first three  
851 strategies and often times also with the fourth DS.

## 852 6.2. Synthetic queries

853 To simulate synthetic queries, we randomly generated provenance poly-  
854 nomials in which the coefficients and exponents could be greater than 1.  
855 The queries involve three GtoPdb tables: `family`, `contributor2family`,  
856 and `contributor`. The polynomials were generated as follows (in particu-  
857 lar, every time we write “randomly”, we mean using a uniform distribution):  
858 first, the number of monomials composing the polynomial is decided choos-  
859 ing randomly a number between 1 and 6. Then, we randomly choose a tuple  
860 from the tables `family`, one from the table `contributor2family` and one  
861 from table `contributor`, that are used as the monomial’s variables. Again,  
862 randomly, we choose a coefficient for this monomial (between 1 and 3) and  
863 an exponent for each tuple (between 1 and 4). For the next monomial, then,  
864 we decide if we want to keep the same tuple from the table `family` as first  
865 tuple of the new monomial. To do so, we generate a random number between  
866 0 and 1. If the number is above 0.2, we change the family tuple.

867 An example can be found in Figure 9, which shows a sample synthetic  
868 provenance polynomial (the how-provenance), the corresponding why-provenance  
869 and lineage expressions, and the causality of the tuples of the lineage, to-  
870 gether with their responsibility. The resulting credit distribution for each  
871 DS is shown after the provenance expression.

872 As an example of how the distribution strategies behave with these syn-  
873 thetic queries, consider tuple  $f_5$  in Figure 9. This tuple receives the high-  
874 est quantity of credit using responsibility-based distribution, and less credit  
875 using, in order, lineage, why- and how-provenance. This is because more  
876 information is available about the role of the tuple in the overall compu-  
877 tation. Generally speaking, the more complex the distribution (the most  
878 complex being how-provenance), the more credit is given to tuples which  
879 are more frequently used, and thus have a higher impact in producing the  
880 output tuple. Responsibility, on its part, can be seen as an enrichment of  
881 the information brought by lineage. It enriches the tuples of the lineage with  
882 a value providing us with a ranking describing the importance of tuples in  
883 generating the output. As such, the responsibility-based DS moves part of



**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:**

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

**Causality:** counterfactual causes:  $\emptyset$ ,

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

**Responsibility:**

$$f_1 = \frac{1}{2}, f_5 = \frac{1}{2}, c_2f_1 = \frac{1}{3}, c_2f_2 = \frac{1}{3}, c_2f_{17} = \frac{1}{2}, c_1 = \frac{1}{3}, c_2 = \frac{1}{3}, c_{18} = \frac{1}{2}$$

**Credit distribution:**

$$f_1 = \frac{3}{20}, f_5 = \frac{3}{20}, c_2f_1 = \frac{1}{10}, c_2f_2 = \frac{1}{10}, c_2f_{17} = \frac{3}{20}, c_1 = \frac{1}{10}, c_2 = \frac{1}{10}, c_{18} = \frac{3}{20}$$

Figure 9: Sample synthetic provenance polynomial (how-provenance) and corresponding why-provenance, lineage, causality and responsibility values, together with the corresponding credit distributions.

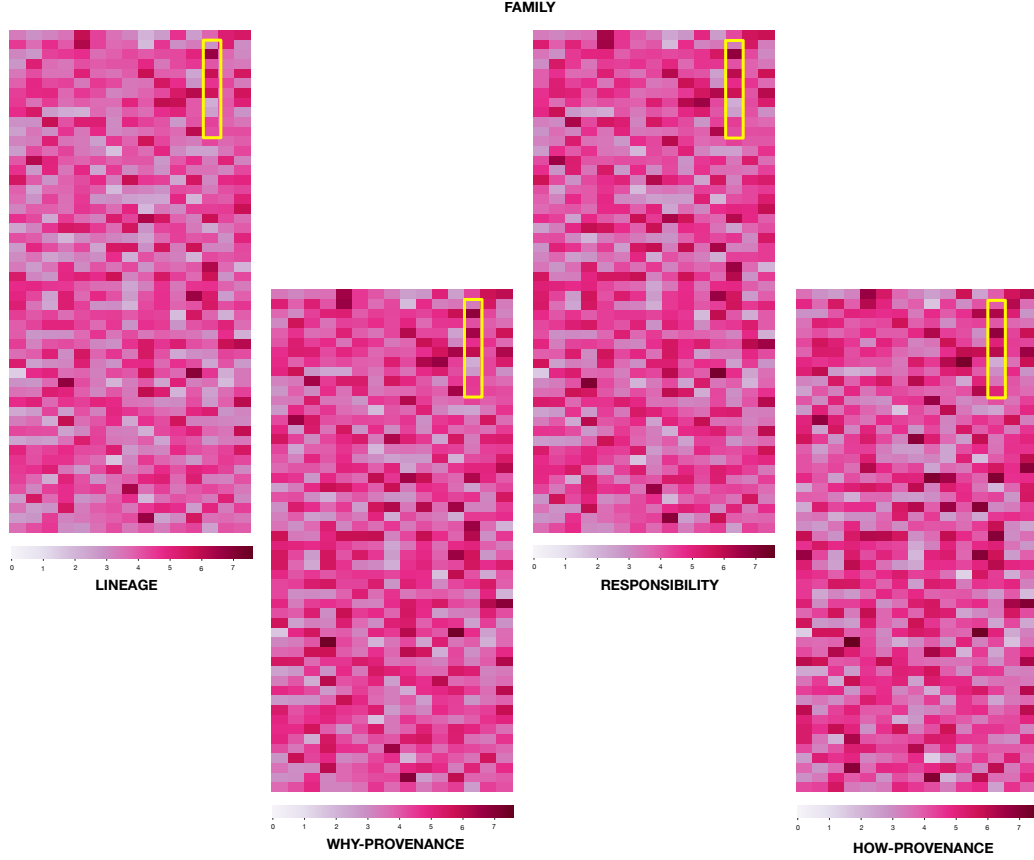


Figure 10: Comparison of three DS on the same table `family` after the distribution computed using 10K synthetic and randomly generated provenance polynomials. The tuples in the blue rectangles are used as example in the discussion connected to Figure 11.

the credit to  $f_1$ ,  $f_5$ ,  $c_2f_17$  and  $c_18$ , since they are tuples that are more important than the others in generating the outputs. This notion of “importance” is connected to their corresponding minimal contingency sets. For example,  $f_1$  has as minimal contingency set (one of the many)  $\{f_5\}$ , with cardinality 1. On the other hand,  $c_1$  has, as minimal contingency set (one of the many)  $\{f_5, c_2\}$ , with cardinality 2. This means that  $c_1$  is “less important” of the tuples with minimal contingency sets of lower cardinality, and this is reflected on the different quantity of credit being distributed.

Despite being synthetic, these provenance polynomials represent realistic queries. The polynomials can be obtained by any nested query with join and

union operations that use the same tuple multiple times (in which case the exponents are bigger than 1), and the same combination of operations more than once (in which case the coefficients of monomials are bigger than 1).

*Results.* The results of credit distribution on the **family** table using 10K randomly generated synthetic provenance polynomials are shown in Figure 10. We set the maximum value in the heat maps to the highest value reached by a tuple in all three distributions (i.e., 7.5).

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

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

Moving to why-provenance-based DS, we see that more credit is given to tuples in the **family** table than with the previous strategy. This is because the DS considers the different ways that a tuple is used, e.g. in joins with other tuples. If the same tuple is present in more than one witness, it will draw more credit and take it from other tuples in the witness basis. In this case, tuples in **family** drew more credit, taking it from tuples in the other two tables, due to the role that **family** tuples played in the queries that were executed. We also notice that the responsibility-based distribution strategy has a distribution that is quite similar to the one provided by why-provenance. It is often the case, for example when the witnesses of the why provenance share one common tuple, that the two distributions behave similarly. As a consequence, at times the generated polynomials are such that the two distributions behave in the same way, or very similarly.

We note that the lineage-based DS gives an average credit of 3.82 to each tuple in the table, while the DS based on why-provenance assigns 4.18 and the one based on responsibility 4.13. Moreover, lineage distributed a total of about 3121 units of credit to the **family** table, while responsibility assigned 3290 and why-provenance 3333.

Finally, consider the how-provenance-based DS heat-map. As with why-provenance, more credit is typically given to tuples in **family** compared to lineage-based DS, since it recognizes the role of these tuples in the queries,

931 and the overall hue is deeper. The two distributions appear similar, although  
 932 on closer inspection, slight differences can be seen. This is because how-  
 933 provenance also considers the frequency with which tuples are used, not only  
 934 the ways in which they are used. Therefore, although the overall distribution  
 935 is similar, there are small differences due to the presence of exponents and  
 936 coefficients in the provenance polynomials, influencing the distribution of  
 937 credit.

938 To better understand this difference, in the next subsection we consider  
 939 the accrual of credit over time. In doing so, we will focus on the ten tuples  
 940 shown within the large yellow rectangles in Figure 11. Each small rectangle  
 941 within a large blue rectangle is a tuple, and we number them from 1 (top) to  
 942 ten (bottom). *These ten tuples were selected specifically because they allow*  
 943 *us to see the evolution of the distribution of credit through time.*

### 944 6.3. Credit accrual over time

945 Since credit accrues over time, we simulate the passage of time by varying  
 946 the number of queries executed, and look at the “snapshots” of credit for each  
 947 of the strategies using synthetic queries. The results are shown in Figure 11.

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

952 The polynomials used are the same as the experiment of the previous  
 953 section. The range of credit in each map goes from 0 (no credit) to 7 (the  
 954 maximum quantity of credit reached – using how-provenance – on one of the  
 955 tuples of the considered window at the “snapshot” with 10K queries). The  
 956 color hue of the legend, as can be seen, still ranges from 0 to 7.5.

957 By the end of 1K queries, credit differentials between tuples as well as  
 958 between strategies can be seen. For example, tuple 3 is usually rewarded the  
 959 most credit by all three strategies. However, it receives the highest quantity of  
 960 credit from the why-provenance-based strategy. Tuple 3 receives the highest  
 961 quantity of credit overall with how-provenance. *Moreover, it can be seen*  
 962 *that tuples 1 and 7 increase their quantity of credit when how-provenance is*  
 963 *exploited.* Moving to 2K queries, it is possible to see that tuple 3 and 7 are  
 964 still the most rewarded by the strategies. This trend continues to the end of  
 965 2k queries.

966 By the end of 5k queries, tuple 7 emerges with the highest value of credit  
 967 for why- and how-provenance, a position which is strengthened by the end

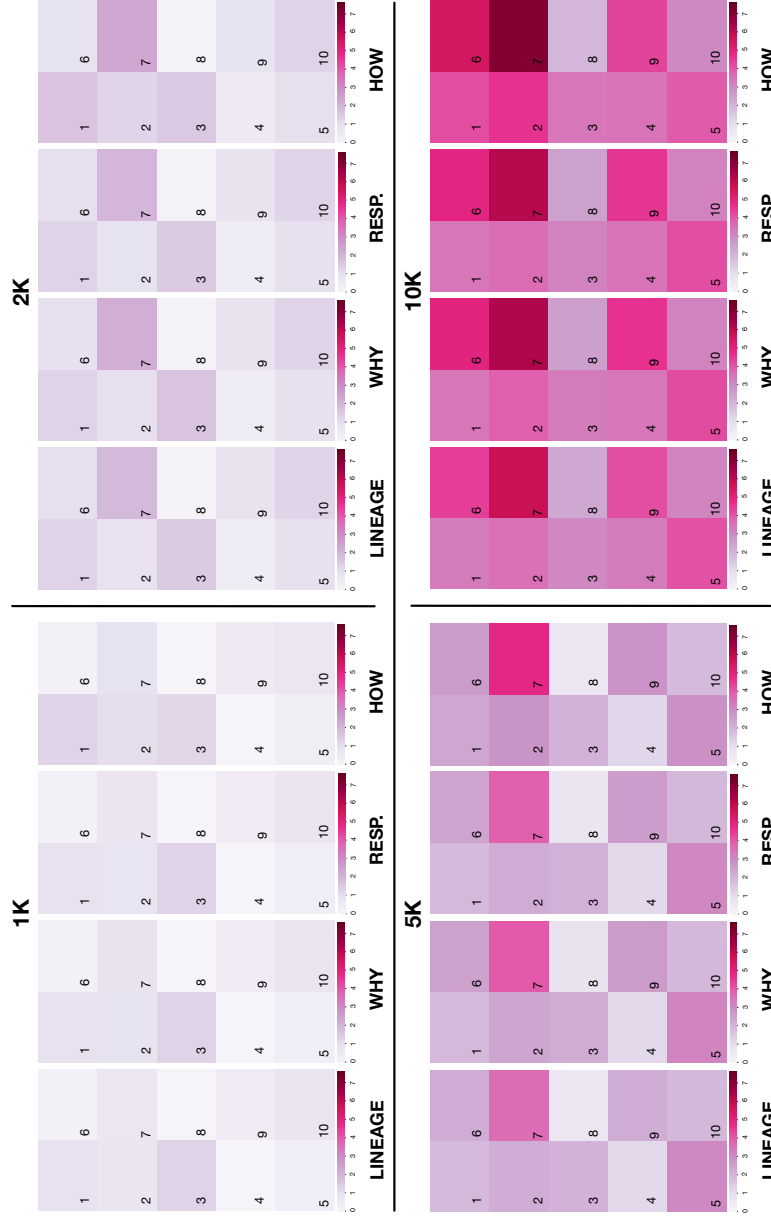


Figure 11: Comparison of the distribution of credit performed by the four DSs on a subset of 10 tuples taken from the `family` table, simulating the passing of time. The number at the top of each group of heat-maps represents the number of polynomials whose credit has been distributed.

968 of 10k queries. Moreover, with the passing of time, tuple 3 ceases to be one  
 969 of the most rewarded ones and new tuples, such as 6 and 9, emerge as being  
 970 particularly rewarded at 5K, while at 10K tuples 6 and 7 are the most re-  
 971 warded from the distributions. This is because tuple 7 is used several times  
 972 within queries being executed, which is rewarded strongly by why- and how-  
 973 provenance. We also note that the responsibility-based distribution confirms  
 974 its trend of being similar to why-provenance, although not completely identi-  
 975 cal. This is more evident at step10K, where tuple 7 is slightly less rewarded  
 976 using responsibility (6.12) with respect to why-provenance (6.24). This is  
 977 due to the fact that, among the polynomials being used for the experiments,  
 978 in some of them tuple 7 had a high responsibility but did not appear in al  
 979 witnesses, thus changing slightly the distribution.

980 While the relative value of credit “positions” of tuples within a DS strategy  
 981 depends on what queries are being executed, the important thing to notice  
 982 is the difference between the DSs over time: overall, lineage gives less credit  
 983 to tuples in the **family** table than the other two strategies since credit is  
 984 shared with tuples in other tables. However, the why-, **responsibility**- and  
 985 how-provenance-based strategies recognize the more important role being  
 986 played by the **family** tuples than those in the other tables. The differences  
 987 between why- and responsibility-based DS are, for the most times, negligible.  
 988 The differences between the why- and how-provenance-based DSs are also  
 989 relatively minor (about plus or minus 0.2 out of 9.5) in most cases. However,  
 990 there are certain situations in which the role of a tuple is particularly critical  
 991 in a query, and in this case the difference in the value of credit assigned is  
 992 notably higher for how-provenance, as we saw with tuple 7 in the example  
 993 of Figure 11.

994 To sum up, the DS based on lineage is sufficient to highlight which tuples  
 995 in the database are used by a query, and distributes credit equally to these  
 996 tuples. The resulting distribution rewards tuples that are used by more  
 997 queries, but does not reward how many times tuples are used in the same  
 998 query. However, a DS based on why-, **responsibility**- or how-provenance may  
 999 be better if the queries are complex, since they reward more tuples that have  
 1000 a critical role in generating the output. In particular, these **three** DSs may  
 1001 be useful for finding “hotspots” in the database based on the role of tuples,  
 1002 with the how-provenance-based DS being preferable if a higher sensitivity to  
 1003 the role of a tuple in queries is required.

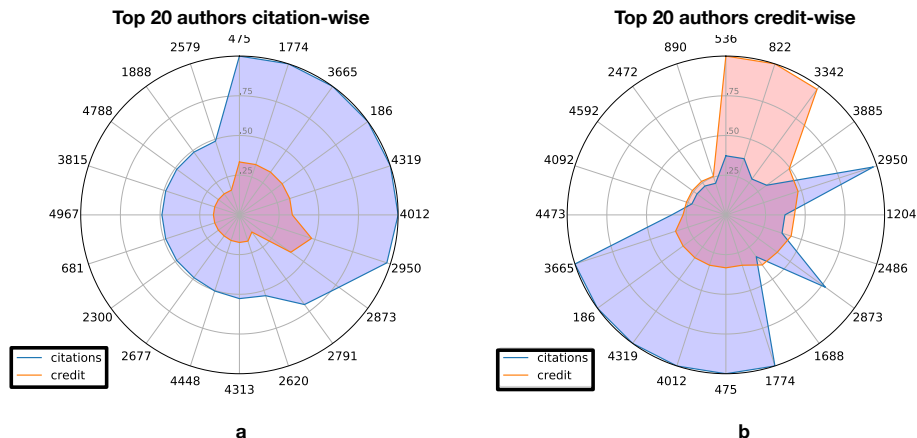


Figure 12: Radars presenting the top 20 authors citation-wise and credit wise, together with their (normalized between 0 and 1) values of citations and credit.

#### 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.

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

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

The results are shown in the radar plots of Figure 12, in which each number on the outer circle (e.g. 475, 1774 and 3665) represents an author

(id) and the blue (red) line represents the normalized value of credit generated by citations (credit), respectively. The first radar plot, Figure 12.a, shows the top 20 authors in terms of *citations*, ordered in a clockwise direction, whereas Figure 12.b orders the authors based on *credit*. Comparing the author ids used in the outer circles of these two plots, it can immediately be seen that the “top authors” are very different using these two metrics, although there is some overlap (for example, authors 1774, 475, and 4012).

Diving a bit deeper to focus on the red and blue areas in each of the plots reveals that there is a significance difference between citations and credit: The top 20 authors in terms of citations do not have the highest values of credit (Figure 12.a). Conversely, the authors with the highest values of credit do not necessarily have a large number of citations (Figure 12.b). For example, author 536 has the highest value of credit, but is not even in the top 20 authors in terms of citations. This means that authors like 536, 822, and 3342 in Figure 12.b receive much more credit from their relatively few citations than authors like 475, who receives the largest number of citations. That is, the data underlying certain webpages is more “valuable” in terms of credit than a citation to the webpage.

The reason for the difference between citation and credit is partly due to the experimental setup: Each output tuple carries a credit of 1, and there can be many tuples used to generate a webpage. Thus a webpage that is created from more tuples will have a higher credit value than one created from fewer tuples. Furthermore, authors who collaborated with fewer people 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 used the same synthetic polynomials described in Section 6.2, and we distributed credit with the first 100, 1K, and 10K of them. 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



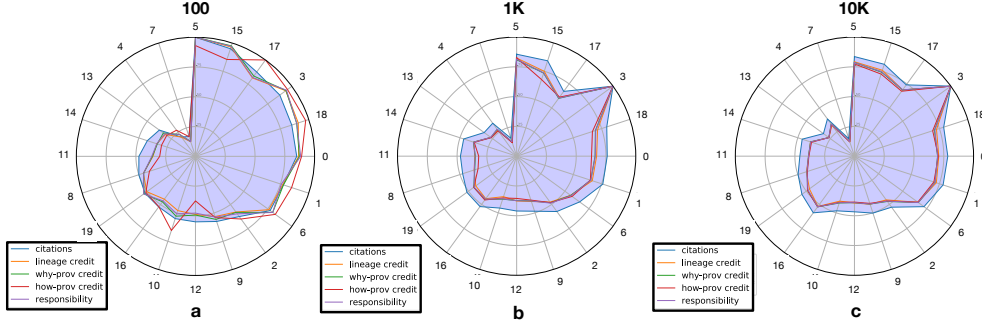


Figure 13: Radars presenting the 20 synthetic authors with corresponding citation and quantities of credit distributed through the 4 DSs (all values normalized between 0 and 1) through different numbers of polynomials (respectively, 100, 1K and 10K). The order is the one defined by figure a, i.e. descending order of citations obtained from 100 polynomials.

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 four kinds of credit, each one using a different DS.

Figure 13 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), how-provenance-based credit (red line), and responsibility-based credit (violet 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 12. In fact, for Figure 13.a the linear correlation between the citation number and all four types of credit is always above 0.94 with p values in the order of  $3e-8$ . The credit distributed via lineage is closest to the number of citations (a linear correlation of 0.99, p value of  $2e-16$  in Figure 13.a), while the other three types of credit behave slightly differently (a linear correlation of around 0.95 in all other three cases in Figure 13.a). Similar observations can be made for Figure 13.b and 13.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 authors 17 and 10 in Figure 13.a, and especially when credit is distributed using how-provenance. This again shows how credit gives a different per-

spective 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 become almost identical (the linear correlation for the values of Figure 13.c is more than 0.99 with a p-value of  $1.32e-32$ ). This is consistent with what we observed in Figure 10.

## 7. Discussion

*Credit Generation.* In this paper we focused on Credit Distribution, the problem of distributing credit generated by a citation to the parts of the database being used by the query subsumed by that citation. A different problem is credit generation, the task of generating credit *before* its distribution. Credit generation presents, in itself, a series of new problems. Among them, we count here:

1. *The correct generation of credit* Different types of citations may generate different quantities of credit. Data being cited in the related work may generate less credit than a result set of data that are extensively used throughout the paper. Different techniques may be employed to correctly compute the credit, such as the manual annotation by the authors of the data that are more relevant in their own assessment to the economy of the paper, or computations performed through NLP techniques to infer the importance of a citation.
2. *Credit produced by self-citations* Data credit, being built on top of traditional citations, inherits some of its problems. Authors, using self-citations, may generate and distribute credit to themselves, making their work appear much more impactful than it really is in reality. Different strategies may be exploited in this scenario, ranging from ignoring completely the credit generated from self-citations to applying a discount factor.
3. *Generic citations* As we discussed, citations may go to the whole database, or to large views in the database itself. In this case, credit may be assigned indiscriminately to large portions of data, losing the ability to accurately identify parts of the database that have high impact, and highlighting the whole database as being important, without really identifying any interesting part of it. This problem may also have

different solutions, such as ignoring queries that are too “general” and considering only queries that are discriminative.

4. *Different types of credit* In the real world, there are different types of research communities interested in different information in a database. Doctors’ interests and queries may differ from the interests and queries of ophthalmologists or pharmacists. For this reason, only distributing one generic credit generated from all possible queries may simply highlight data that are important in general, without taking into consideration the specific needs of communities. One possibility is to distinguish the type of credit, e.g., have one credit generated from queries coming from doctors, another type of credit generated from queries submitted by ophthalmologists, etc. In this way, it will be possible to accurately tailor the process of credit distribution around the information need of different categories of users.

*Credit Generation vs Credit Distribution.* 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, therefore, new results, different from the ones obtained in the previous sections, may be found. These results will depend on the nature of the context and the quantity of credit being considered.

*On the choice of the DS.* Depending on the type of task at hand, a different choice may be made for the DS to use. When the user only wants to highlight

the tuples being used in the database by a workload, the lineage-based DS is sufficient. When the user wants to know also the relative impact of tuples in the context of the query, the other DSs may be used. This may be true for applications such as data pricing, where we want to give a price to the parts of a database and credit may become a criterion to decide this price. In this context, other forms of provenance may be preferred since they allow to better understand the actual importance of data. While the real-world example that we used showed that the four DSs behave the same, this was due to the specific nature of the data and the queries being used.

In reality, the why-provenance of a query differs from the lineage of the same query whenever the output tuples can be computed in more than one way by the query, i.e., if there is more than one witness. While at the best of our knowledge there isn't any work that explores SQL query logs to validate the presence of this diversity, we still think this to be true in many case. To support this opinion, the work by Bonifati et al. [9] showed that in the context of SPARQL query logs submitted to various databases such as DBpedia and Wikidata, more than 90% of these queries are of type select, and more than 30% perform join operations through the and operator. These queries moreover contain triple patterns with cardinalities that range from 1 to 11 triples, highlighting their big complexity in certain cases. These queries, that many times are converted in their SQL versions, are composed by join operations that may result in why-provenances with cardinality bigger than 1. Other works, such as [56], showed that operations such as Inner joins can be found in at least 4.5% of queries in the considered workload, with a maximum number of times that operator is used in the same query equal to 164. Outer joins were found in 1% of the queries, and used up to 247 times in the same query. This is another evidence of the potentiality of the fact the why-provenances may become quite complex.

## 8. Conclusions and Future Work

This paper defines three new distribution strategies based on why- and how-provenance and on responsibility, and compares them against the lineage-based distribution strategy defined in [26]. The first, why-provenance-based DS, uses the concept of a witness, and gives more credit to tuples that appear in more than one witness. In this way, tuples that are more important to the query and are used in different ways are rewarded more. The second, how-provenance-based DS, considers the frequency with which a tuple

or combination of tuples is used in the query through the information contained in a provenance polynomial. In this case, the how-provenance-based DS is more sensitive than the why-provenance-based DS to the role and importance of tuples.

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

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

In the third set of experiments, we showed how the differences between 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

1230 distribute credit. In this paper we assumed that each output tuple carries  
1231 credit 1. In more sophisticated scenarios we can employ different strategies  
1232 to compute credit, that reflect the importance of cited data. Also, other,  
1233 and more sophisticated strategies could also be used to decide how credit is  
1234 distributed between the authors, beyond the uniform distribution used here,  
1235 in a way to reflect the work performed by them on the cited data.

1236 We will also explore new applications for credit over relational databases.  
1237 One example is *data pricing*, which gives a price to a query submitted by a  
1238 user who wants to buy the produced information. Currently, a commonly  
1239 strategy used for data pricing is based on query rewriting: A database stores a  
1240 set of views with their price. When a new query arrives, the system rewrites  
1241 it using the stored views to obtain a query price, a process that can be  
1242 computationally expensive. We plan to distribute credit through carefully  
1243 planned and representative queries, and use credit information to define a  
1244 new, faster, and potentially more flexible pricing function.

1245 Another application is *data reduction* [46], which addresses the problem  
1246 of reducing the vast – and rapidly expanding – amount of data that is being  
1247 produced.

1248 Data credit can also address this problem, by helping find “hotspots”  
1249 and “coldspots” of data. A hotspot is data in a database (e.g. a tuple) with  
1250 a high quantity of credit, which is therefore valuable for the set of queries  
1251 that execute frequently over the data and distribute the credit. On the other  
1252 hand, a coldspot is data with a low quantity of credit, which is therefore  
1253 considered less important and could be deleted or moved to cheaper and/or  
1254 less efficient memory.

## 1255 Acknowledgement

1256 The work was partially supported by the ExaMode project, as part of the  
1257 European Union H2020 program under Grant Agreement no. 825292.

## 1258 References

- 1259 [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bern-  
1260 stein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L.,  
1261 Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Koss-  
1262 mann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo,  
1263 T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo,

- 1264 A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suciu, D.  
1265 (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–  
1266 53.
- 1267 [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating  
1268 data citation in citedb. *PVLDB*, 10(12):1881–1884.
- 1269 [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y.  
1270 (2018). Data citation: A new provenance challenge. *IEEE Data Eng.*  
1271 *Bull.*, 41(1):27–38.
- 1272 [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An  
1273 Introduction to the Joint Principles for Data Citation. *Bulletin of the*  
1274 *Association for Information Science and Technology*, 41(3):43–45.
- 1275 [5] Baggerly, K. (2010). Disclose all data in publications. *Nature*,  
1276 467(7314):401–401.
- 1277 [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D.,  
1278 Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I.,  
1279 Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and  
1280 Goble, C. A. (2013). Why linked data is not enough for scientists. *Future*  
1281 *Gener. Comput. Syst.*, 29(2):599–611.
- 1282 [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation  
1283 Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.
- 1284 [8] Bertin-Mahieux, T., Ellis, D., Whitman, B., and Lamere, P. (2011). The  
1285 million song dataset. In *Proceedings of the 12th International Conference*  
1286 *on Music Information Retrieval (ISMIR 2011)*, pages 591–596.
- 1287 [9] Bonifati, A., Martens, W., and Timm, T. (2017). An analytical study of  
1288 large SPARQL query logs. *PVLDB*, 11(2):149–161.
- 1289 [10] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In  
1290 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*  
1291 *nication*, pages 93–116. De Gruyter Mouton.
- 1292 [11] Buneman, P. (2006). How to cite curated databases and how to make  
1293 them citable. In *18th International Conference on Scientific and Statistical*  
1294 *Database Management, SSDBM*, pages 195–203. IEEE Computer Society.

- [12] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D., Pawson, A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn't working, and what to do about it. *Database J. Biol. Databases Curation*, 2020.
- [13] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation is a computational problem. *Commun. ACM*, 59(9):50–57.
- [14] Buneman, P., Khanna, S., and Tan, W. C. (2001). Why and where: A characterization of data provenance. In *Database Theory - ICDT 2001, 8th International Conference*, pages 316–330.
- [15] Buneman, P. and Silvello, G. (2010). A rule-based citation system for structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.
- [16] Callaghan, S., Donegan, S., Pepler, S., Thorley, M., Cunningham, N., Kirsch, P., Ault, L., Bell, P., Bowie, R., Leadbetter, A. M., Lowry, R. K., Moncoiffé, G., Harrison, K., Smith-Haddon, B., Weatherby, a., and Wright, D. (2012). Making Data a First Class Scientific Output: Data Citation and Publication by NERC's Environmental Data Centres. *International Journal of Digital Curation*, 7(1):107–113.
- [17] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Journals: A Survey. *Journal of the Association for Information Science and Technology*, 66(9):1747–1762.
- [18] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases: Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–474.
- [19] Chockler, H. and Halpern, J. Y. (2004). Responsibility and blame: A structural-model approach. *J. Artif. Intell. Res.*, 22:93–115.
- [20] CODATA-ICSTI Task Group on Data Citation Standards and Practices (2013). *Out of Cite, Out of Mind: The Current State of Practice, Policy, and Technology for the Citation of Data*, volume 12.
- [21] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N. (2019). Bringing citations and usage metrics together to make data count. *Data Science Journal*, 18(1).



- [22] Cronin, B. (1984). *The Citation Process. The Role and Significance of Citations in Scientific Communication*. London: Taylor Graham.
- [23] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evidence of a structural shift in scholarly communication practices? *JASIST*, 52(7):558–569.
- [24] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of view data in a warehousing environment. *ACM Trans. Database Syst.*, 25(2):179–227.
- [25] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on Innovative Data Systems Research*. [www.cidrdb.org](http://www.cidrdb.org).
- [26] Dosso, D. and Silvello, G. (2020). Data credit distribution: A new method to estimate databases impact. *Journal of Informetrics*, 14(4):101080.
- [27] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The virtual atomic and molecular data centre (VAMDC) consortium. *Journal of Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- [28] Eiter, T. and Lukasiewicz, T. (2002). Complexity results for structure-based causality. *Artif. Intell.*, 142(1):53–89.
- [29] ESIP Data Preservation and Stewardship Committee (EDPSC) (2019). Data citation guidelines for earth science data, version 2. Version 2, Earth Science Information Partners.
- [30] Fang, H. (2018). A discussion of citations from the perspective of the contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–1520.
- [31] Garfield, E. (1999). Journal impact factor: a brief review. *Can. Med. Assoc.*, 979-980.
- [32] Gößwein, B., Miksa, T., Rauber, A., and Wagner, W. (2019). Data identification and process monitoring for reproducible earth observation research. In *2019 15th International Conference on eScience (eScience)*, pages 28–38. IEEE.

- [33] Green, T. J., Karvounarakis, G., and Tannen, V. (2007). Provenance semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 31–40. ACM.
- [34] Halpern, J. Y. and Pearl, J. (2013). Causes and explanations: A structural-model approach — part 1: Causes. *CoRR*, abs/1301.2275.
- [35] Harding, S. D., Sharman, J. L., Faccenda, E., Southan, C., Pawson, A. J., Ireland, S., Gray, A. J. G., Bruce, L., Alexander, S. P. H., Anderton, S., Bryant, C., Davenport, A. P., Doerig, C., Fabbro, D., Levi-Schaffer, F., Spedding, M., Davies, J. A., and Nc-Iuphar (2018). The IUPHAR/BPS guide to PHARMACOLOGY in 2018: updates and expansion to encompass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Research*, 46(Database-Issue):D1091–D1106.
- [36] Hartley, J. (2017). Authors and their citations: a point of view. *Scientometrics*, 110(2):1081–1084.
- [37] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a transformed scientific method.
- [38] Honor, L. B., Haselgrove, C., Frazier, J. A., and Kennedy, D. N. (2016). Data citation in neuroimaging: proposed best practices for data identification and attribution. *Frontiers in neuroinformatics*, 10:34.
- [39] Joshi-Tope, G., Gillespie, M., Vastrik, I., D’Eustachio, P., Schmidt, E., de Bono, B., Jassal, B., Gopinath, G. R., Wu, G. R., Matthews, L., Lewis, S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of biological pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- [40] Katz, D. (2014). Transitive credit as a means to address social and technological concerns stemming from citation and attribution of digital products. *Journal of Open Research Software*, 2(1).
- [41] Kosten, J. (2016). A classification of the use of research indicators. *Scientometrics*, 108(1):457–464.
- [42] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S. (2011). Citation and Peer Review of Data: Moving Towards Formal Data Publication. *International Journal of Digital Curation*, 6(2):4–37.

- 1388 [43] Martone, M. (2014). Joint declaration of data citation principles.  
 1389 *FORCE11. San Diego CA. Data Citation Synthesis Group*. [https://www.](https://www.force11.org/datacitationprinciples)  
 1390 [force11.org/datacitationprinciples](https://www.force11.org/datacitationprinciples), online September 2020.
- 1391 [44] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation  
 1392 counts and rankings of LIS faculty: Web of science versus scopus and  
 1393 google scholar. *Journal of the american society for information science*  
 1394 *and technology*, 58(13):2105–2125.
- 1395 [45] Meliou, A., Gatterbauer, W., Moore, K. F., and Suciu, D. (2010). The  
 1396 complexity of causality and responsibility for query answers and non-  
 1397 answers. *Proc. VLDB Endow.*, 4(1):34–45.
- 1398 [46] Milo, T. (2019). Getting rid of data. *Journal of Data and Information*  
 1399 *Quality (JDIQ)*, 12(1):1–7.
- 1400 [47] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D.,  
 1401 Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G.,  
 1402 Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff,  
 1403 D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,  
 1404 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson,  
 1405 E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C.,  
 1406 Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers,  
 1407 E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research  
 1408 culture. *Science*, 348(6242):1422–1425.
- 1409 [48] Parsons, M. A., Duerr, R. E., and Jones, M. B. (2019). The history and  
 1410 future of data citation in practice. *Data Science Journal*, 18(1).
- 1411 [49] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J.  
 1412 (2016). Research data explored: An extended analysis of citations and  
 1413 altmetrics. *Scientometrics*, 107(2):723–744.
- 1414 [50] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic,  
 1415 large databases: Model and reference implementation. In *Proceedings of*  
 1416 *the 2013 IEEE International Conference on Big Data, 6-9 October 2013,*  
 1417 *Santa Clara, CA, USA*, pages 307–312.
- 1418 [51] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identi-  
 1419 fication of Reproducible Subsets for Data Citation, Sharing and Re-Use.

- 1420 *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue*  
 1421 *on Data Citation*, 12(1):6–15.
- 1422 [52] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data  
 1423 citation of evolving data: Recommendations of the working group on data  
 1424 citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- 1425 [53] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*  
 1426 *Sci. Technol.*, 69(1):6–20.
- 1427 [54] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data  
 1428 provenance in e-science. *SIGMOD Record*, 34(3):31–36.
- 1429 [55] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-  
 1430 spective. In of Sciences’ Board on Research Data, N. A. and Information,  
 1431 editors, *Report from Developing Data Attribution and Citation Practices*  
 1432 *and Standards: An International Symposium and Workshop*, pages 177–  
 1433 178. National Academies Press: Washington DC.
- 1434 [56] Vogelsgesang, A., Haubenschild, M., Finis, J., Kemper, A., Leis, V.,  
 1435 Mühlbauer, T., Neumann, T., and Then, M. (2018). Get real: How bench-  
 1436 marks fail to represent the real world. In *Proceedings of the Workshop on*  
 1437 *Testing Database Systems*, pages 1–6.
- 1438 [57] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,  
 1439 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,  
 1440 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data  
 1441 management and stewardship. *Scientific data*, 3.
- 1442 [58] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data  
 1443 citation: Giving credit where credit is due. In *Proceedings of the 2018*  
 1444 *International Conference on Management of Data, SIGMOD*, pages 99–  
 1445 114.
- 1446 [59] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to  
 1447 scientific datasets using article citation networks. *Journal of Informetrics*,  
 1448 14(2).
- 1449 [60] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of  
 1450 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.

- 1451 [61] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for  
1452 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*  
1453 *Molecular Spectroscopy*, 327:122–137.