

Credit Distribution in Relational Scientific Databases

Dennis Dosso^a, Susan B. Davidson^b, Gianmaria Silvello^a

^a*Department of Information Engineering, University of Padua, Italy*

^b*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 then 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 based on how-provenance, the second based on the concept of responsibility, and the third on the Shapley value.

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 these distribution strategies vary in their sensitivity to the role of an input tuple in the generation of the output data, and reward input tuples differently.

Keywords: Data Citation, Data Credit, Provenance, Causality and Responsibility, Shapley value

1. Introduction

Citations are an essential component of scientific research, enabling research products to be found as well as the relationships between them to be created and understood. They form a basis on which to give credit to authors, papers, and venues [20, 21, 64]. Citations are used, among other things, to decide on tenure, promotion, hiring, and funding of grants for researchers [22, 35, 43, 47].

Science and research are increasingly digital, and there are numerous curated databases that are at the core of scientific research efforts [12]. It is therefore generally accepted that data must be cited and citable [15, 44], and that data citations should contribute to the scientific reputation of researchers, scientists, data curators, and creators [4, 59]. It is also accepted that data citations should be counted alongside of traditional citations, and contribute to bibliometrics indicators [7, 51].

A central problem in the data citation process is how to attribute credit to data creators and curators [11]. How to handle and count the credit generated by data citation, and how it contributes to traditional and new bibliometrics, are long-standing research issues [9, 30]. However, even when correctly applied, data citations and the bibliometrics computed using them do not always fully 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 idea of *Data Credit Distribution* (DCD) [29, 41, 63] has 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 is a proxy for 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.

In this paper, we consider data credit as a measure of value for data in

38 a (curated) scientific database. Credit is a real value that can be assigned
 39 to data of any kind and at any level of granularity. Therefore the concept
 40 of “data” is left intentionally vague, although in this paper we focus on
 41 relational databases. Credit acts as a proxy for the value of data based on
 42 the measure of citations, accesses, clicks, downloads, or other surrogates for
 43 data use.

44 We define DCD as *the process, method, or algorithm used to assign credit*
 45 *to a given datum or dataset*. It differs from the traditional citation setting
 46 since:

- 47 1. When a paper p_1 cites another paper p_2 , a +1 citation “credit” is given
 48 to p_2 , and to all its authors. It does not matter why or how paper
 49 p_1 cites paper p_2 ¹, the result is always +1 to the citation count of p_2
 50 and of its authors. A different credit distribution strategy can assign a
 51 quantity of credit to p_2 and its authors that is *proportional* to the role
 52 played by p_2 in p_1 . Hence, we can weight the importance of the cited
 53 entities and assign credit according to their role.
- 54 2. Traditional citations are *atomic*: a citation from p_1 to p_2 can never
 55 be broken into pieces and assigned in part to p_2 and in part to other
 56 papers or data that contributed to p_2 . In contrast, with data credit,
 57 we use a *non-atomic* real value, which can be divided and distributed
 58 to multiple components of a database.
- 59 3. Credit can be *transitive*, that is, it can be propagated through one
 60 cited entity to other entities cited by it that contributed to its content.
 61 Citations, traditionally, are not.

62 We study the DCD problem in the context of relational databases (RDBs)
 63 since they are widely used² and are the main focus of current work in data
 64 citation methods [12, 14, 52]. RDBs are also frequently a test-bed for new
 65 methods that can be adapted to other databases, e.g., graphs or document
 66 databases. The “portions” of data in an RDB that can be credited can be
 67 defined at different levels of granularity, in particular: (i) the whole database,
 68 (ii) tables, (iii) tuples, and (iv) attributes. The ability to specify different
 69 levels of granularity in a relational database allows us to define the DCD

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

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



Figure 1: Overview of the credit distribution pipeline.

problem at a particular level of granularity. In this paper, we focus on DCD at the tuple level.

The DCD process that we use is summarized in Figure 1:

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

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

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

Step 4 Given the database instance I and the query Q , the *data provenance* of $Q(I)$ is computed. The data provenance of $Q(I)$ is a form of metadata that captures how Q used I to generate the output [17].

Step 5 Provenance is input to the *Credit Distribution Strategy* (CDS, also referred only as Distribution Strategy, DS). CDS is a function f that takes as input the credit value k , divides it and distributes it to the

86 data in the input database I , and is defined on the basis of citation
 87 policies decided at the database administration level or at the domain
 88 community level.

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

93 This paper expands the work in [26] where we first defined the problem
 94 of DCD in relational databases, and proposed a viable Distribution Strategy
 95 (DS) based on *lineage* – the simplest form of *data provenance*. The lineage
 96 of a tuple t in the output $Q(I)$ is defined as the set of all and only the tuples
 97 in the database instance I that are “relevant” to the production of t . The
 98 corresponding strategy equally redistributes the credit k to the tuples in the
 99 lineage set, thus each tuple receives credit $k/|L_t|$, where L_t is the lineage set
 100 of t .

101 One may argue that this DS is too simplistic, since lineage does not convey
 102 any information about the role or importance of input tuples in the query.
 103 Therefore, one may desire to give more credit to the tuples that are more
 104 *important* to the production of the output, i.e. those tuples that, if removed,
 105 would prevent the output tuple from appearing in the final result, or those
 106 tuples used more than once by the query.

107 Therefore, in this paper, we expand the ideas in [26] by proposing new DSs
 108 based on another, more general, form of data provenance: how-provenance [32].
 109 We also propose other two DS based on the concepts of responsibility [48] and
 110 the Shapley value [25, 45]. We show how these DS differ from each other and
 111 discuss why one may be preferred to another depending on the application
 112 and its goals. We also show that the DSs based on responsibility give more
 113 credit to tuples that are essential to the production of the result set, whereas
 114 the how-provenance-based DS takes into consideration the different ways in
 115 which a tuple is used. Finally, the DS based on the Shapley value sees the
 116 process of distribution as a competitive game in which tuples that contribute
 117 more to the generation of the output are correspondingly rewarded more.

118 We use a well-known curated database called the IUPHAR/BPS³ Guide

³International Union of Basic and Clinical Pharmacology/British Pharmacology Society

119 to Pharmacology [34], also known as GtoPdb⁴, to evaluate the DSs. GtoPdb
120 contains expertly curated information about diseases, drugs, cellular drug
121 targets, and their mechanisms of action. We chose GtoPdb for two main
122 reasons: (i) it is a widely-used and valuable curated relational database, (ii)
123 many papers in the literature use, and cite, its data (i.e., families, ligands,
124 and receptors). Real queries used in papers can therefore be seen as data
125 citations which, in turn, can be used to assign data credit.

126 We perform four sets of experiments. In the first, real queries are ex-
127 tracted from papers published in the British Journal of Pharmacology (BJP),
128 that represent data citations to GtoPdb, and are used to distribute credit in
129 the database using the three different provenance-based DSs. In the second
130 and third experiment we analyze the behavior of the different DS when com-
131 plex citation queries are employed. In the fourth set of experiments we use
132 both real and synthetic queries to assess the difference between traditional
133 citation and the notion of credit distribution in terms of rewarding those
134 responsible for the data, e.g. data curators.

135 **Contributions** of this work include:

- 136 • Four new Distribution Strategies based on why-provenance, how-provenance,
137 responsibility and the Shapley value.
- 138 • An in-depth analysis of the effects of credit distribution on real-world
139 curated data and of the differences between the five proposed Distri-
140 bution Strategies.
- 141 • A comparison between the behavior of traditional citations and data
142 credit in rewarding data curators.

143 **Outline.** The rest of the paper is organized as follows: Section 2 presents
144 background material and related work. Section 3 describes the GtoPdb use
145 case. Section 4 briefly presents the forms of provenance used in the paper.
146 Section 5 describes the credit distribution problem and the proposed dis-
147 tribution strategies. In Section 6 we present the experimental evaluation,
148 followed by a discussion of our design decisions in Section 7. Section 8 draws
149 some conclusions and outlines future work.

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

150 2. Background

151 *Data in Research.* The world of research is rapidly transitioning towards the
152 *fourth paradigm of science* [36], that is, data-intensive scientific discovery,
153 where data are important for scientific advances as well as for traditional
154 publications [6].

155 The scientific community is promoting an *open research culture* [50],
156 founded on methods and tools to share, discover, and access experimental
157 data. The community has identified the FAIR principles (Findable, Acces-
158 sible, Interoperable, and Reusable) [61], that should be enforced by every
159 database. In particular, data should be accessible from the articles, journals,
160 and papers that cite or use them [20]. Aspects such as the need for the *repro-*
161 *ducibility* of experiments through the used data; the *availability* of scientific
162 data; the *connections* between data and the scientific results are all needed
163 aspects for the fourth paradigm, and are all relevant to the domain of *data*
164 *citation* [37].

165 *Data Citation: Principles and Motivations.* Data Citation principles were
166 proposed in [19], and later summarized and endorsed by the Joint Declaration
167 of Data Citation Principles (JDDCP) [46]. The principles are divided into
168 two groups [57]. The first group contains principles concerning the role of
169 data citation in scholarly and research activities such as the (i) *importance*
170 of data (why data citation is important and why data should be considered
171 as first-class citizens); (ii) *credit* and *attribution* to the creators and curators
172 of the data; (iii) *evidence*; (iv) *verifiability*; and *interoperability*, with these
173 last three requiring data citation methods to be flexible enough to operate
174 through different communities. The second group defines the main guidelines
175 to establish a data citation systems, and contains principles such as the (i)
176 *unique identification* of the data being cited; (ii) *(open) access* to data; (iii)
177 guarantee of *persistence* and *availability* of citations even after the lifespan
178 of the cited entity; the (iv) *specificity* of a citation, i.e. it must lead to the
179 data set originally cited.

180 The main motivations for data citation are outlined in [57] and range from
181 data attribution and connection to data sharing, impact and reproducibility.

182 2.1. Data Citation in Relational Databases

183 Relational databases have been the target of data citation methods since
184 the surge of the data-centric research paradigm. The RDA “Working Group

185 on Data Citation: Making Dynamic Data Citable”⁵ [53] has developed guide-
186 lines for citing large, dynamic, and changing datasets which have now moved
187 on into adoption phase. The datasets considered by the Working Group are
188 often relational.

189 In one of its most recent sessions [54], the Working Group (WG) on
190 Data Citation reported that there are various implementations of its guide-
191 lines for Data Citation on MySQL/Postgres relational databases. Some of
192 these databases are: DEXHELPP⁶ (Social Security Records); NERC (ARGO
193 Global Array); EODC (Earth Observation Data Centre) [31]; LNEC (River
194 dam monitoring); MDS (Million Song Database) [8]; CBMI⁷ (Center for
195 Biomedical Informatics); VMC (Vermont Monitoring Cooperative); CCA⁸
196 (Climate Change Center Austria); VAMDC (Virtual Atomic and Molecular
197 Data Center) [27, 65].

198 More examples of work on data citation in relational databases are [2,
199 12, 24, 62]. The website <https://fairsharing.org/> keeps an updated list
200 of curated and scientific databases (many of which are relational or graph-
201 based) following FAIR guidelines. These databases are citable since they are
202 compliant with the most recent guidelines, and they are in the vast majority
203 of cases accessible via dynamically created Webpages. In all these databases
204 it is, therefore, possible to implement DCD on top of the existing infrastruc-
205 tures for citing data.

206 Data citation techniques are primarily applied to relational databases
207 because of their pervasiveness as well as the “identifiability” of the portions
208 of data that are to be cited: the whole database, a relation, a tuple, or
209 even an attribute. Many papers [2, 10, 12] consider more complex citable
210 units, recognizing that often the *views* of a database are the ones to be cited.
211 Generally, a *view* is a query on the database. To this end, [62] suggested
212 decomposing the database into a set of views, where each view is associated
213 with its citation.

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

- 215 1. A database cited as a whole, even though only parts of the databases
216 are used in the papers or datasets. Alternatively, the so-called “data pa-

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

⁶<http://www.dexhelpp.at/>

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

⁸<https://ccca.ac.at/startseite>

pers” are cited, being traditional papers that describe a database [16]. In this case, all the credit from the citations goes to the database administrators or to the authors of the data papers.

2. Subsets of data, obtained by issuing queries to a database, are individually cited. This is the solution adopted by the *Resource Data Alliance* (RDA) working group on Data Citation [53]. In this case, the credit generated from citations is distributed among the contributors of the portions of data being cited, and/or to the database administrators.
3. The database is accessible via a series of Webpages that arrange the content of the database by topic or theme. Examples in the life science domain include the Reactome Pathway database [40], the GtoPdb [34], and the VAMDC [65]. Every single Webpage is unequivocally identifiable and can be individually cited.

2.2. Data Credit

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

Katz [41] suggests the need for a *modified citation system* that includes the idea of *transient* and *fractional credit*, to be used by developers of research products as software and data. Two considerations are made: (i) research objects such as data and software are currently not formally rewarded or recognized by the community; (ii) even in traditional papers, the contribution of each author to the work is hard to understand, unless explicitly specified in the paper. This is even more true for data, where different groups of people work on the same database.

In [41] credit is defined as a “quantity” that describes the importance of a research entity, such as papers, software, or data, mentioned in a citation. It also proposed the idea of a *distribution* of credit from research entities, such as papers or data, to other research entities through citations. Therefore, when discussing data credit, we need to consider *credit computation* – i.e., the process to compute the quantity of credit generated by the citation – and *credit distribution* – i.e., the process to distribute credit and to assign it to the entities that contributed to the creation/curation of the cited data. In this paper we focus on the latter.

These two processes are done by exploiting the structure of the *citation graph*, a directed graph whose nodes are publications and edges are citations.

This graph is the model at the core of systems such as Google Scholar and the Web of Science. We add to this that the concept of credit can be built on top of the existing infrastructure handling traditional and data citations.

Katz [41] further explores the idea of a *distribution* of credit from research entities (i.e., papers and data) to other research entities through citations that connect them. Thanks to traditional citations and now also to data citations, this distribution is finally possible, at least between papers and data. Some problems related to traditional citations can thus be solved by citations:

1. Credit rewards research entities that to date are not (formally) recognized (a goal shared with data citation).
2. Credit can reward authors *proportionally* to their role in generating the entity. The more an author contributes to a paper, the more credit is given to him. Zou and Peterson [64] work on something similar with their zp-index, which includes in its formulation the position (and thus the role) of a publication author to represent its impact in the work itself.
3. Credit can be *transitively* channeled through a chain of papers citing each other, thus enabling the rewarding of older papers that are no more cited, since other papers summarize or report their content but are nevertheless crucial in a research area for the influence of their content.

Fang [29] presents a framework to distribute the credit generated by a paper to its authors and to the papers in its reference list in a transitive way. Let us consider the *citation graph* as the graph where the nodes are papers and the links are the citations among them. In this graph, every paper is a source of credit, which is then transferred to the neighboring nodes. The quantity of credit received by each cited paper depends on its impact/role in the citing paper. So far, this theoretical framework is limited to papers, but it can be easily extended to a citation graph including both papers and data.

Zeng et al. [63] proposes the first method to compute credit within a network of papers citing data. Adopting a network flow algorithm, they simulate a random walker to estimate a score for each dataset, leveraging real-world usage data to compute the credit. This is the first step towards an automatic credit computation procedure. This proposal is, however, limited to assigning credit to whole datasets, and it does not deal with the granularity of data.

291 It does not work to assign credit to a single research entity within a dataset.
 292 Differently from Zeng et al. [63], we do not treat the credit computation
 293 process, but we focus on the distribution process.

294 2.3. Data Provenance

295 To distribute credit, we base two of our methods on *data provenance*.
 296 Data provenance is information that describes the origin and the process of
 297 creation of data. It can also be seen as metadata pertaining to the derivation
 298 history of the data. It is particularly useful to help users to understand
 299 where data are coming from, and the process they went through. Data
 300 citation and data provenance are closely linked [3] since both are forms of
 301 annotations on data retrieved through queries. Data provenance has been
 302 widely studied in different areas of data management. In this paper, we
 303 focus on provenance for database management systems (DBMS). For further
 304 details on data provenance, please refer to surveys like [17] and [58].

305 Cheney et al. [17] presents four main types of data citation for DBMS: *lin-*
 306 *age* [23], *why-provenance* [13], *how-provenance* [32] and *where-provenance* [13].

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

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

316 2.4. Causality and Responsibility

317 We also consider the notions of causality and responsibility, as defined
 318 in [48]. Causality is an enrichment of lineage, and it is the attribution of
 319 a certain degree of importance to the tuples of the lineage based on their
 320 role in the generation of the output. Responsibility is a value given to the
 321 tuples of the lineage to rank them based on their degree of causality (the
 322 more important the role of a tuple in generating the output, the higher its
 323 responsibility).

324 While computing responsibility for general queries is hard [18], Meliou
 325 et al. [48] proved a dichotomy result for conjunctive queries: for each query
 326 without self-joins, either its responsibility can be computed in PTIME in the

327 size of the database or checking if it has a responsibility below a given value
328 is NP-hard.

329 2.5. Shapley value

330 The Shapley value was introduced in 1952 [56], framed as a *cooperative*
331 *game* played by a set A of players, and defined by a *wealth function* v that
332 assigns to each coalition set $B \subseteq A$ the wealth $v(B)$. The question behind the
333 Shapley Value is how to quantify the contribution of each player to the overall
334 wealth. Informally, the Shapley value is defined as follows [45]: assume that
335 we select players randomly one by one and without replacement, starting
336 with the empty set. Every time a player a is selected, its addition to the
337 coalition B produces a change in the wealth of the coalition from $v(B)$ to
338 $v(B \cup \{a\})$. The Shapley value of a is the expectation of change that a causes
339 in this probabilistic process.

340 The Shapley value has been widely used, e.g. in economics, law, envi-
341 ronmental science, and network analysis, and has strong theoretical justifica-
342 tions. However, its use in databases as a metric for quantifying the influence
343 of a tuple on the output of a query (thereby presenting an alternative to
344 responsibility) has only recently been considered [45]. The initial theoretical
345 analysis in [45] showed mainly lower bounds on the complexity of the prob-
346 lem, and did not suggest a feasible implementation. However, very recently,
347 an efficient implementation for Boolean queries has been provided [25], both
348 in terms of an exact computation (which in practice works well for most
349 queries) and in inexact one (which is extremely fast and provides the same
350 ranking of tuples as the exact computation, but not necessarily the same
351 values).

352 3. Use Case: GtoPdb

353 The IUPHAR/BPS Guide to Pharmacology [34] (GtoPdb⁹) is a well-
354 known and well structured scientific relational database that contains ex-
355 pertly curated information about diseases, drugs in clinical use, their cellular
356 targets, and the mechanisms of action on the human body. It is curated and
357 maintained by the GtoPdb Committee and 96 subcommittees, comprising
358 512 scientists collaborating with in-house curators who draw the information

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

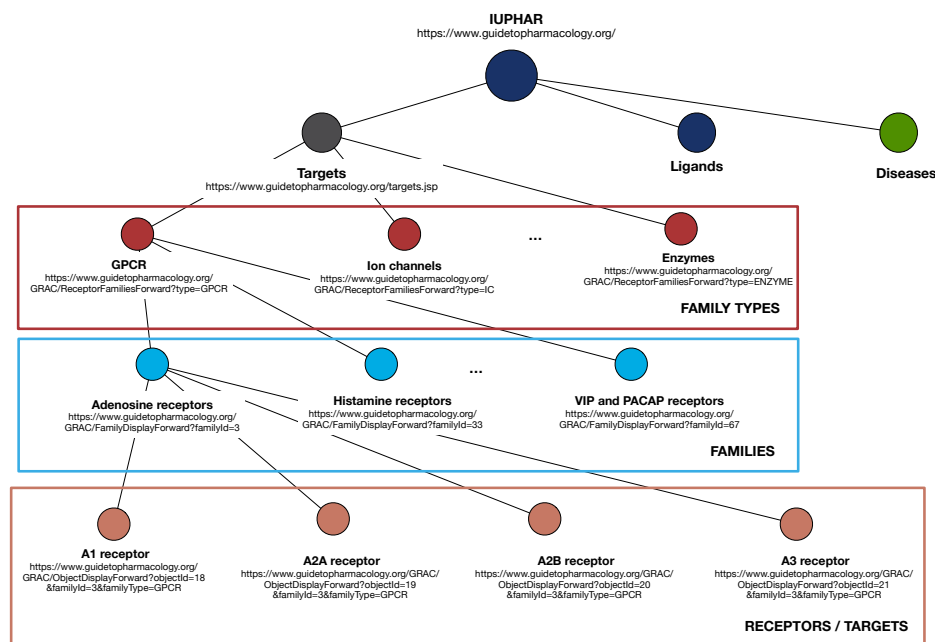


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

359 contained in the database from high-quality pharmacological and medicinal
 360 chemistry literature. Roughly 1000 researchers from all over the world have
 361 contributed to the database, and the curators wanted to give recognition to
 362 these contributors. This led to some early work on data citation [10].

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

373 GtoPdb provides access to the webpages corresponding to all these nodes
 374 through URLs. The webpages corresponding to target families all present a
 375 similar structure, as shown in Figure 3 for the “Adenosine receptors” family.

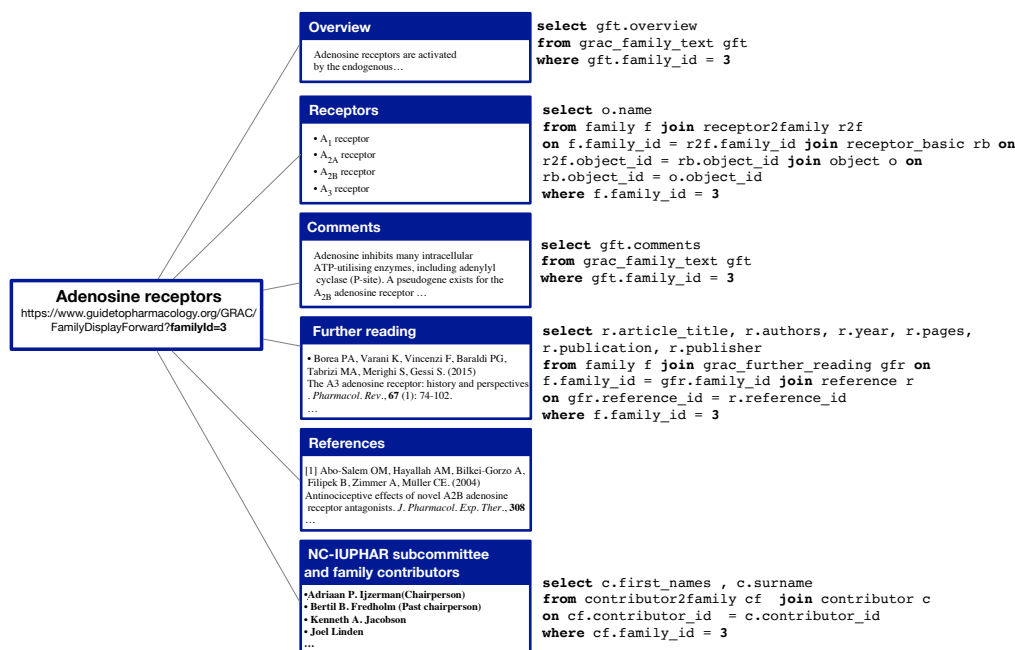


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

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

In practice, many papers in the literature only reference GtoPdb (the root) without including a reference to the specific page being cited. That is, they only cite a paper describing GtoPdb as a whole (e.g., [34]) 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.

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** contains receptor families; **contributor** contains the name and country of contributors; **contributor2family** connects contributors to the families they contributed to.

393 In certain “lucky” cases, as with papers available in PDF and published
394 in the British Journal of Clinical Pharmacology ¹⁰ (BJCP), when a family,
395 ligand, receptor name, etc. are used, they have a hyperlink pointing to the
396 corresponding webpage in GtoPdb. Therefore, the citations to the families
397 can be detected and counted using the URLs reported in the papers. How-
398 ever, these citations to GtoPdb webpages are not counted as such by citation
399 systems, so they are not converted into credit for curators and collaborators.

400 For our running example, consider Table 1. This simplified version of
401 GtoPdb contains three tables: **family**, **contributor** and **contributor2family**.
402 The first table, **family**, has tuples representing families with three attributes:
403 the id of the family, its name, and type. Table **contributor** contains peo-
404 ple who have helped generate the data in the database. The third table,
405 **contributor2family**, serves as a link between the families and the people
406 who contributed to them. For instance, “John Smith” (c_1) contributed to
407 “Dopamine Receptors” (f_1) as well as to the “YANK Family” (f_4). Through-
408 out the rest of the paper, we will use the **id** attribute of these tables as the
409 *provenance token* of its corresponding tuples, that is, as a symbol that serves
410 to identify a tuple when talking about provenance.

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

I	database instance
L_t	lineage set of an output tuple t
Γ	contingency set
ρ_t	responsibility of tuple t
Q	a query
\bar{Q}_o	Boolean query such that $\bar{Q}_o(I) = 1$ if o is present in $Q(I)$
\mathcal{W}	witness basis
W	a witness set
$\gamma(\mathcal{W}, t)$	set of witnesses in \mathcal{W} containing t
\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

Table 2: Notations used in this paper.

4. Provenance, Responsibility, and Shapley

We now describe how-provenance, the notions of causality, responsibility, and the Shapley value function.

4.1. How-Provenance

While why-provenance describes the source tuples that witness an output tuple in the result of the query, it leaves out information about how the source tuples are used. How-provenance was defined in [32] to capture the information about *how* the source tuples are used to build the output using a *semiring* algebraic structure. It takes the form of a polynomial, called *provenance polynomial*, where the variables are taken from the set X of identifiers of the tuples (provided that each tuple in I has an identifier) and the coefficients are drawn from the set of natural numbers \mathbb{N} .¹¹

1. The set K is a *commutative monoid* for the operator $+$ with a neutral element 0. Therefore, it has these properties:

¹¹This semiring is commonly referred as $\mathbb{N}[X]$ in the literature.

- 425 (a) $(a + b) + c = a + (b + c)$ (associative property)
- 426 (b) $0 + a = a + 0 = a$ (0 is the neutral element)
- 427 (c) $a + b = b + a$ (commutative property)
- 428 2. The set K is a *monoid* with identity element 1. Therefore, it has these
- 429 properties:
- 430 (a) $(a \cdot b) \cdot c = a \cdot (b \cdot c)$ (associative property)
- 431 (b) $1 \cdot a = a \cdot 1 = a$ (1 is the neutral element)
- 432 3. Multiplication is distributive on addition, i.e.:
- 433 (a) $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$
- 434 (b) $(a + b) \cdot c = (a \cdot c) + (b \cdot c)$
- 435 4. Multiplication by 0 annihilates K , i.e. $\forall x \in K, 0 \cdot x = x \cdot 0 = 0$

436 The key idea in Green et al. [32] is to use the two operators $+$ and \cdot to
 437 represent two basic transformations that source tuples undergo as a result
 438 of applying a relational query to a database [17]. Two tuples may either
 439 be joined together (a join is represented with the \cdot operator) or merged via
 440 union or projection (represented with the $+$ operator).

441 Formally, let K be a set containing an element 0. A K -relation is a
 442 function $R : U\text{-Tuples} \mapsto K$ which maps every U -tuple in an element in K
 443 such that its support, defined as $\text{supp}(R) = \{t | R(t) \neq 0\}$, is finite. U -tuples
 444 is the set of tuples with attributes in the set U . The K -relation is a finite
 445 function which models a relation R , tagging each tuple in R with an element
 446 of K and each tuple that is not in R with 0.

447 **Definition 4.1.** *Operations on the algebraic structure $(K, 0, 1, +, \cdot)$ [32]*
 448 *Let $(K, 0, 1, +, \cdot)$ be an algebraic structure with two binary operations $+$ and*
 449 *\cdot and two distinguished elements 0 and 1. The operations of the positive*
 450 *K -relational algebra are defined as follows:*

- 451 1. **Empty relation.** *For any set of attributes U , $\exists \emptyset : U\text{-Tuples} \mapsto K | \emptyset(t) =$*
 452 *0.*
- 453 2. **Selection** *Let $R : U\text{-Tuples} \mapsto K$ and σ be a selection predicate that*
 454 *maps each U -Tuple to either 0 or 1. Then $\sigma_\theta(R) : U\text{-Tuples} \mapsto K$ is*
 455 *defined by $(\sigma_\theta(R))(t) = R(t) \cdot \sigma(t)$.*
- 456 3. **Projection** *Let $R : U\text{-Tuples} \mapsto K$ and $V \subseteq U$. Then $\pi_V(R) : V\text{-Tuples}$*
 457 *$\mapsto K$ is defined by $(\pi_V(R))(t) = \sum_{t=t'[V] \vee R(t') \neq 0} R(t')$.*
- 458 4. **Union** *Let $R_1, R_2 : U\text{-Tuples} \mapsto K$. Then $R_1 \cup R_2 : U\text{-Tuples} \mapsto K$ is*
 459 *defined by $(R_1 \cup R_2)(t) = R_1(t) + R_2(t)$.*

460 5. *Natural join* Let $R_1 : U_1\text{-Tuples} \mapsto K$ and $R_2 : U_2\text{-Tuples} \mapsto K$. Then
 461 $R_1 \bowtie R_2 : U_1 \cup U_2\text{-Tuples} \mapsto K$ is defined by $(R_1 \bowtie R_2)(t) = R_1(t_1) \cdot$
 462 $R_2(t_2)$, where $t_1 = t[U_1]$ and $t_2 = t[U_2]$.

463 In the selection rule, θ is seen as a function $\theta : U\text{-Tuples} \mapsto \{0, 1\}$. It
 464 is observed in [17, 32] that if the K -relational semantics satisfies the same
 465 equivalence laws as positive relational algebra operators over bags, i.e. union
 466 $(+)$ is associative, commutative and has identity \emptyset and join (\cdot) is associa-
 467 tive, commutative and distributive over union, and projection and selection
 468 commute with each other, as well as with union and join, then $(K, 0, 1, +, \cdot)$
 469 must be a commutative semiring.

470 The semiring operations document *how* each output tuple is produced
 471 from source tuples. If each source tuple in a database D is tagged with a
 472 distinct id, the semiring gives the how-provenance for each output tuple in
 473 the form of a polynomial with coefficient from the set \mathbb{N} of natural numbers
 474 and variables from the set of source tuples id.

475 To define how-provenance in a compositional manner, as done in [17], let
 476 us consider the algebraic structure $(\mathbb{N}(\text{TupleLoc}), 0, 1, +, \cdot)$, where $\mathbb{N}(\text{TupleLoc})$
 477 is the set of polynomials whose coefficients are the natural numbers and the
 478 variable are from the set TupleLoc . The how-provenance of an output tuple
 479 is a function $\mathcal{H} = \text{How}(Q, I, o)$ that returns a polynomial in $\mathbb{N}(\text{TupleLoc})$
 480 called *provenance polynomial*.

Definition 4.2. *How-Provenance*

Let Q be a (complex) SPJRU query. Let I be a database instance, and t be a tuple in $Q(I)$. Then, the how-provenance of t according to Q and I , denoted as $\text{How}(Q, I, t)$, is an element of the set $\mathbb{N}(\text{TupleLoc})$ defined as follows:

$$\begin{aligned}
 \text{How}(\{u\}, I, t) &= \begin{cases} 1, & \text{if } t = u, \\ 0 & \text{otherwise.} \end{cases} \\
 \text{How}(R, I, t) &= \begin{cases} (R, t), & \text{if } t \in R, \\ 0 & \text{otherwise.} \end{cases} \\
 \text{How}(\sigma_\theta(Q), I, t) &= \theta(t) \cdot \text{How}(Q, I, t) \\
 \text{How}(\rho_{A \mapsto B}(Q), I, t) &= \text{How}(Q, I, t[B \mapsto A]) \\
 \text{How}(\pi_V(Q), I, t) &= \sum_{u \in \text{supp}(Q), u[V]=t} \text{How}(Q, I, u) \\
 \text{How}(Q_1 \bowtie Q_2, I, t) &= \text{How}(Q_1, I, t[U_1]) \cdot \text{How}(Q_2, I, t[U_2]) \\
 \text{How}(Q_1 \cup Q_2, I, t) &= \text{How}(Q_1, I, t) + \text{How}(Q_2, I, t)
 \end{aligned}$$

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 3: Result of Q1 over the database instance in Table 1 with the how-provenance polynomial of each output tuple.

481 $\{u\}$ is a query expression describing a constant, singleton relation, not a
 482 relation value per se. These constants correspond to K -relations that assign
 483 1 to u and 0 to all other tuples.

484 Table 3 shows the two output tuples of our running example annotated
 485 with their respective how-provenances. Tuple o_2 was produced by a join of
 486 the input tuples $f_4, c2f_4$, and c_1 . The three provenance tokens are therefore
 487 “multiplied” together. The case of o_1 is slightly more complex, as already
 488 discussed. It can be obtained by the joins of two different sets of tuples,
 489 so there are two monomials combined by $+$ representing these alternative
 490 derivations. Each monomial corresponds, in a way, to the witnesses of the
 491 why-provenance of o_1 .

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

499

500 4.2. Causality and Responsibility

501 A formal study of causality was introduced in [18, 33] and later expanded
 502 by Meliou et al. [48] to explain the causes of answers and non-answers to
 503 queries. In the following, we refer to the definition of causality and respon-
 504 sibility provided in [48]. In particular, we only focus on answers to a query
 505 since non-answers are not relevant in our context.

506 There are two types of “cause” tuples: counterfactual and actual. Let o
 507 be a tuple in the result of query q on the database instance I , and t a tuple
 508 in its lineage. We call t a *counterfactual cause* if, by removing t from I , o is
 509 also removed from the output (i.e., t is essential for the generation of t). We
 510 call t an *actual cause* if there is a set of tuples $\Gamma \subseteq I$ called a *contingency*

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 4: Result of Q1 over the database instance in Table 1 with the responsibilities of lineage tuples.

511 *set*, such that t is a counterfactual cause in $I - \Gamma$. In other words, t is an
512 actual cause if, even when removed from I , there is another set of tuples of
513 the lineage that guarantees the presence of o .

514 Computing the causality of tuples is NP-complete for general queries [28],
515 but for conjunctive queries can be computed in PTIME, as showed by Meliou
516 et al. [48].

517 The notion of *responsibility* measures the degree of causality as a function
518 of the size of the smallest contingency set [18]. This allows us to rank lineage
519 tuples based on their degree of causality in generating the output.

Definition 4.3. *Responsibility* [48]

Let o be an output tuple in the result of query Q on I , and let t be a cause for o . The responsibility of t for the answer o is:

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

520 where Γ ranges over all contingency sets for t .

521 Note that a counterfactual cause will have the maximum responsibility
522 of 1, and that the larger the minimum contingency of an actual cause is, the
523 smaller its responsibility will be since there are alternatives to guarantee the
524 presence of the answer o .

525 As an example, consider Table 3, where we reported the result set of Q1
526 and the tuples of the lineages with their responsibility values. Focusing on
527 o_1 : the lineage tuple f_1 is a counterfactual cause, since its contingency set is
528 empty (when removed from the database, o_1 disappears from the result set).
529 Consequently, its responsibility is 1. All the other tuples of the lineage are
530 actual causes. c_1 , for example, has as minimal contingency set $\{c_2f_2\}$, thus
531 its responsibility is 0.5. For the output tuple o_2 , all the tuples of the lineage
532 are counterfactual causes, thus their responsibility is 1.

533

id	name	Shapley value
o_1	Dopamine Receptors	$f_1 = \frac{7}{15}, c_2 f_1 = \frac{2}{15}, c_2 f_2 = \frac{2}{15}, c_1 = \frac{2}{15}, c_2 = \frac{2}{15}$
o_2	YANK Family	$f_4 = \frac{1}{3}, c_2 f_4 = \frac{1}{3}, c_1 = \frac{1}{3}$

Table 5: Result of **Q1** over the database instance in Table 1 with the Shapley values of the tuples of the lineage. In this case D^n corresponds to the lineage.

534 4.3. Shapley value

535 We use the definitions provided in [25]: Let q be a Boolean query and $f \in$
536 D be a fact, the Shapley value of f in D intuitively represents the contribution
537 of f to the query result.¹² The higher the value, the more f helps in satisfying
538 q . Formally, the Shapley value is defined as follows:

$$Shapley(q, D, f) = \sum_{E \subseteq D \setminus \{f\}} \frac{|E|! (|D| - |E| - 1)!}{|D|!} \left(q(E \cup \{f\}) - q(E) \right)$$

539 The sum in this value is performed on all possible subsets of D that do
540 not contain f . The value $(q(E \cup \{f\}) - q(E))$ is the “wealth” brought by f
541 when added to E . Thus, the Boolean query is used as a wealth function v :
542 its value is 1 only when the set $E \cup \{f\}$ makes the query true, and the set E
543 makes it false, i.e., when the addition of the fact f is determinant to making
544 the Boolean query true. The value $|E|! (|D| - |E| - 1)!$ is the number of all
545 the possible permutations over D where the facts in E come first, then f is
546 added, and then all the remaining facts. Thus, the value $\frac{|E|! (|D| - |E| - 1)!}{|D|!}$ can
547 be thought as a weight for the wealth brought by f when added to E .

548 To extend this definition to non-Boolean queries, we adopt the approach
549 in Deutch et al. [25]: the Shapley value of the fact f for the answer \bar{t} to
550 $Q(\bar{x})$ is the value $Shapley(Q[\bar{x}/\bar{t}], D, f)$, where $Q[\bar{x}/\bar{t}]$ is the Boolean query
551 defined by $Q[\bar{x}/\bar{t}](D) = 1$ if and only if \bar{t} is in the output of $Q(\bar{x})$ on D , and
552 0 otherwise. In other words, the definition of $Shapley(q, D, f)$ is extended
553 to queries $Q(\bar{x})$ with free variables by considering the Boolean query $Q[\bar{x}/\bar{t}]$
554 as a value function. This query can be seen as a function that takes as input
555 a set of facts and returns 1 if this set is a witness for \bar{t} , and 0 otherwise.

556 As an example, consider Table 5, that shows the Shapley values for the
557 lineage’s tuples of o_1 and o_2 , results of query **Q1**. We note that, to compute

¹²We ignore the distinction between endogenous and exogenous facts, since in our setting they are all assumed to be endogenous.

the Shapley value of an input tuple f it is sufficient to compute and sum the values $\frac{|E|!(|D|-|E|-1)!}{|D|!}$ for all the possible sets E such that $E \cup \{f\}$ is a witness and E is not. Thus, suppose we want to compute the Shapley value of the tuple f_1 . Let us call \bar{Q}_{1,o_1} the Boolean query such that $\bar{Q}_{1,o_1}(D) = 1$ if and only if o_1 is in the output of Q1 on D , and L_{o_1} is the lineage of o_1 . Then the Shapley value of f_1 with respect of o_1 is given by:

$$\begin{aligned} \text{Shapley}(\bar{Q}_{1,o_1}, L_{o_1}, f_1) &= \frac{2!2!}{5!} + \frac{2!2!}{5!} + \frac{3!}{5!} + \frac{3!}{5!} + \frac{3!}{5!} + \frac{3!}{5!} + \frac{4!}{5!} \\ &= \frac{7}{15} \end{aligned}$$

where for the first element of the sum the corresponding E is $\{c2f_1, c_1\}$, for the second element it is $\{c2f_2, c_2\}$, for the third $\{c2f_1, c2f_2, c_1\}$, for the fourth $\{c2f_1, c_1, c_2\}$, for the fifth $\{c2f_2, c_2, c_1\}$, for the sixth $\{c2f_1, c2f_2, c_2\}$, and for the seventh $\{c2f_1, c2f_2, c_1, c_2\}$. Every other possible subset E would make the factor equal to 0. Note that in this case we consider $D = L_{o_1}$, the lineage of o_1 , since these are the only facts in all the database that contribute to the generation of o_1 .

Similarly, for tuple c_1 (and the other tuples of the lineage), the computation is:

$$\begin{aligned} \text{Shapley}(\bar{Q}_{1,o_1}, L_{o_1}, c_1) &= \frac{2!2!}{5!} + \frac{3!}{5!} + \frac{3!}{5!} \\ &= \frac{2}{15} \end{aligned}$$

It can be seen that for all the tuples of o_2 's lineage the corresponding Shapley values are equal to $1/3$, since they are all equally responsible for the generation of the output. Thus the sum of the Shapley values of all the tuples in an output tuple's lineage is always equal to 1 when using a Boolean query as wealth function.

5. Credit Distribution and Distribution Strategies

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

5.1. Data Credit and Data Credit Distribution

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

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

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

593 In the following, we use the notation in Cheney et al. [17]: Given a
 594 database instance I , a *tuple location* (R, t) is a tuple t in relation R . With
 595 reference to the running example, $(\text{family}, \langle f_1, \text{Dopamine Receptors},$
 596 $\text{gpcr} \rangle)$ is the tuple location of the first tuple in the `family` relation. The set
 597 of all tuple locations in I is called *TupleLoc*. We use this to formally define
 598 DCD at the *tuple level*.

599 **Definition 5.1. Tuple Level Data Credit Distribution (DCD) [26]**
 600 *Given a query Q over I and $k \in \mathbb{R}_{>0}$, DCD is defined by the function $f_{I,Q} :$
 601 $\text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ such that $f_{I,Q}(t, k) = h$ where $0 \leq h \leq k$ and
 602 $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$. The function $f_{I,Q}$ is the distribution strategy (DS).*

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

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

613 5.2. A How-Provenance Based Distribution Strategy

614 The how-provenance-based DS first distributes the credit to the mono-
 615 mials of the polynomial accordingly to the weight represented by their co-
 616 efficients, then to the tuples of each monomial accordingly to the weights
 617 represented by their exponents.

618 To define the DS more formally, we introduce some notation and illustrate
 619 it using the provenance polynomial \mathcal{H} shown in Figure 4. This notation is
 620 also shown in Table 2 for easy reference.

621 We call c the function that, given a polynomial, returns the sum of its
 622 coefficients; thus $c(\mathcal{H}) = 3 + 1 = 4$. We call e the function that, given a

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

Figure 4: Illustration of notation used to define the how-provenance based DS

monomial, returns the sum of its exponents, thus $e(M_2) = 1 + 3 + 3 = 7$.
 mc is the function that takes as input a monomial and returns its coefficient; thus $mc(M_1) = 3$. te is a function that takes as input a tuple and a monomial, and returns the exponent of the tuple in the monomial, if present; thus $te(c_2, M_2) = 3$. Finally, γ takes as input a tuple and the whole polynomial, and returns a set of monomials containing that tuple, if present in the polynomial; thus $\gamma(f_1, \mathcal{H}) = \{M_1, M_2\}$, $\gamma(c_2, \mathcal{H}) = \{M_2\}$.

Definition 5.2. *How-Provenance-Based Distribution Strategy*

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

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

Going back to the example of Table 3, consider o_1 with provenance polynomial $f_1c2f_1c_1 + f_1c2f_2c_2$. The how-provenance-based DS firstly divides the credit between the two monomials. Since the coefficients of each monomial are 1, the credit is split in half. If they were, for example, 1 and 2 respectively, 1/3 of the credit would go to the first monomial, and 2/3 to the second. Since in our example each variable has exponent 1, the credit is further divided equally among the three variables. Thus, at the end of the computation, f_1 receives 1/3, and the other tuples receive 1/6.

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

id	name
oxs_1	Dopamine Receptors

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

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

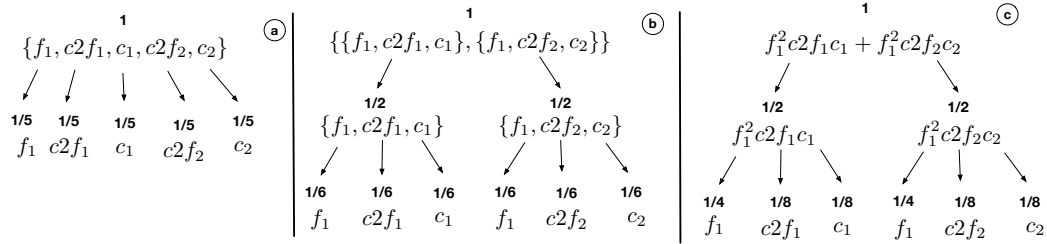


Figure 5: Comparison of different distributions strategies for tuple o_1 produced by query Q2.

```

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

```

649 The result of Q2 is shown in Table 6, and consists of one tuple, oxs_1 ,
650 annotated with each of the three provenances. As can be seen, lineage and
651 why-provenance are identical to those of the tuple o_1 in the previous example.
652 The how-provenance, however, is different since tuple f_1 is used twice: first
653 in the join of the inner query, and second in the join of the outer query. This
654 information is lost in the first two forms of provenances since they are sets,
655 but it is captured in how-provenance through the use of the operator ‘.’.

656 Figure 5 shows the differences between the three DS for the tuple o_1
657 of Table 6. Subfigure 6.a uses lineage, sub-figure 6.b uses why-provenance,
658 and sub-figure 6.c uses how-provenance. The DS based on the provenance
659 polynomial gives credit 1/2 to f_1 , and 1/8 to the other tuples. This is
660 reasonable since Q2 relies on f_1 even more than Q1 does. The distribution
661 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 [32], how-provenance is more general than why-provenance and lineage, in the sense that it contains more information.

5.3. Responsibility-based Distribution Strategy

As described in Section 4.1, causality and responsibility is new information that is added to lineage. One possible option for defining a distribution strategy using responsibility is to simply assign the responsibility of each tuple in the lineage of an output tuple as its credit. In this way, responsibility is both a way to compute credit and to distribute it. Using the example of Table 4, in the case of output tuple o_1 , f_1 receives credit 1 and the other tuples receive credit 0.5.

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. 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 will give a bigger portion of credit to tuples that are higher in the responsibility ranking. Formally:

Definition 5.3. *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:

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

where ρ_j is the responsibility of tuple j as in Definition 4.3.

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

Figure 6 shows the responsibility and credit assigned to the tuples of the lineage of the output tuple o_1 of Table 4. Assuming that $k_{o_1} = 1$, f_1 receives credit $1/3$, while the others receive credit $1/6$. As we see, the DS in this case returns the same distribution as that obtained using why-provenance as shown in Figure 5. This is not always the case though, as we show in the example of Section 6.2.

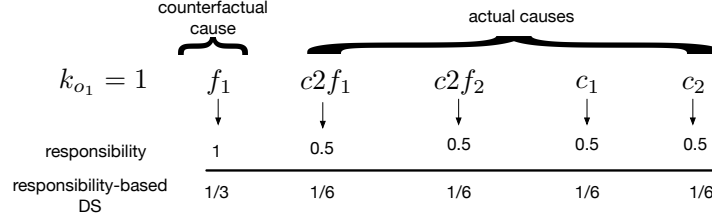


Figure 6: Example of distribution of credit using the responsibility-based DS, assuming $k_o = 1$.

693

694 5.4. Shapley value-based Distribution Strategy

695 As with responsibility, the Shapley value can be seen both as a method
 696 to generate and distribute credit. Moreover, it can be seen that, using the
 697 definition of Shapley value for Boolean queries given in Section 4.1, the sum
 698 of the Shapley values of all the tuples of the lineage L of an output tuple o
 699 is 1. Thus, the definition of a Shapley value-based distribution strategy is
 700 straightforward:

Definition 5.4. *Shapley Value-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 credit given to o . The credit given to tuple t in I is:

$$f_{I,Q}(t, k_o) = k_o \cdot \text{Shapley}(\bar{Q}_o, I, t)$$

701 where \bar{Q}_o is the Boolean query such that, given the set of facts D , $\bar{Q}_o(D) = 1$
 702 if and only if o is in the output of Q on D .

703 As shown in Table 5, tuple f_1 in o_1 's lineage takes credit 7/15 when
 704 $k_{o_1} = 1$, while the other tuples of the lineage take credit 2/15. This DS still
 705 rewards f_1 more than the other tuples, since it is more important than the
 706 other tuples of the lineage. This DS thus behaves differently from all the
 707 other four previous strategies. In particular, f_1 is rewarded more with this
 708 DS than with the others.

709 In the case of o_2 there is only one witness set, thus this DS behaves like
 710 all the others, distributing 1/3 of credit to each tuple in the lineage.

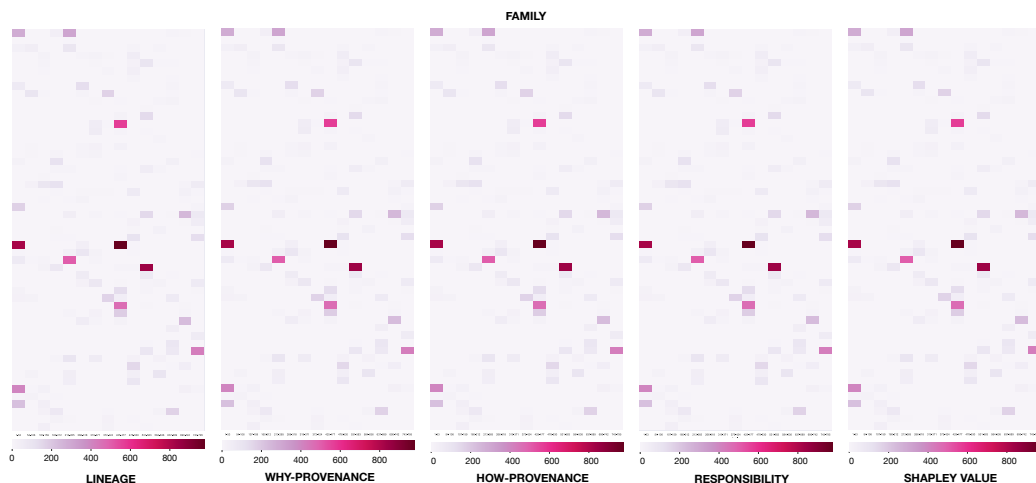


Figure 7: Comparison of four DS on the same table `family` using the distribution given by the queries retrieved from papers. Each cell is a tuple.

6. Experimental Evaluation

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

The source code for the experiments is written in Java and supported by a PostgreSQL database. 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.

727 6.1. Real-world queries

728 Examples of real queries are drawn from papers published in the British
 729 Journal of Pharmacology (BJP) ¹³. Each time a paper in this journal cites a
 730 webpage from GtoPdb, it reports the URL of the page. From this URL, the
 731 query used to obtain the webpage data can be determined. We considered all
 732 889 papers in BJCP citing the IUPHAR/BPS Guide to pharmacology [34]
 733 as of October 2020, and extracted all webpage URLs to GtoPdb contained
 734 within the paper.¹⁴

735 The queries that we inferred are those used to build target family web-
 736 pages within GtoPdb. An example was given in Figure 3, where we show
 737 how the structure of the “Adenosine receptors” family can be mapped into
 738 queries over the underlying database. In GtoPdb, all target family pages
 739 share a similar structure; the only difference is that individual sections, such
 740 as “contributors” or “further readings”, may be missing. Therefore, the same
 741 queries can be used to build all of the target family pages by changing the
 742 family id used in the query (for example, in Figure 3, it is 3). Note that
 743 the queries are fairly simple SQL queries, and fall into a class called “select-
 744 project-join” or “SPJ” queries. A total of more than 12K different queries
 745 were built in this way. Without loss of generality, we give each tuple in the
 746 output of a query a credit of 1.

747 *Results.* Figure 7 shows the heat-maps obtained by the distribution of credit
 748 according to the five DS on one of the tables in the underlying database,
 749 **family**, which is often joined with other tables in the database to build the
 750 webpages. Each cell in a heat-map represents a tuple of the **family** table
 751 and the color indicates the amount of credit attributed to such tuple. It can
 752 be seen that the result of credit distribution over **family** is the same for all
 753 five strategies. The same result is also obtained with the other tables of the
 754 database used by the queries shown in Figure 3.

755 The reason why credit distribution is the same for all five strategies is that
 756 the queries are all simple SPJ queries, which use each table only once and
 757 do joins on key attributes. Under these conditions, each tuple of the output
 758 presents: (i) a how-provenance that is a single monomial with coefficient

¹³<https://bpspubs.onlinelibrary.wiley.com>

¹⁴The 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.

one and exponent one in each variable; (ii) a why-provenance with only one witness; (iii) a lineage that is the same of the witness in the basis, (iv) all tuples are counterfactual causes when considering responsibility, and (v) all tuple have the same importance in the production of the output tuples according to their Shapley value. Hence, for these queries, the five DSs behave in the same way: credit is uniformly distributed among the tuples of the lineage.

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 $\{\{c_{939}, c_{2f_{496}}\}\}$, its lineage is $\{c_{939}, c_{2f_{496}}\}$, both these tuples are counterfactual causes and have a responsibility of one. Therefore, the credit assigned to these tuples is 1/2 using all five DS. This happens for all the tuples in the output of each query of GtoPdb, thus making the distributions equivalent over all outputs.

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

6.2. Synthetic queries

To see what happens with more complex queries, we synthetically generated provenance polynomials in which the coefficients and exponents could be greater than one, and picked them at random from a uniform distribution. The queries involve three GtoPdb tables: `family`, `contributor2family`, and `contributor`. The polynomials were generated as follows: first, the number of monomials was decided by randomly choosing a number between one and six. Then, we randomly chose a tuple from the `family` table, one from the `contributor2family` table and one from the `contributor` table; these are the variables of the monomial. We then chose a coefficient for the monomial (between one and three) and an exponent for each tuple (between one and

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

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

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:

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

Shapley value:

$$f_1 = 0.258\bar{3}, f_5 = \frac{1}{8}, c_2f_1 = 0.091\bar{6}, c_2f_2 = 0.091\bar{6}, c_2f_{17} = \frac{1}{8}, c_1 = 0.091\bar{6}, c_2 = 0.091\bar{6}, c_{18} = \frac{1}{8}$$

Figure 8: Sample synthetic provenance polynomial (how-provenance) and corresponding why-provenance, lineage, responsibility, and Shapley values, together with the corresponding credit distributions. The sum of Shapley values is equivalent to the quantity of credit being distributed (assuming that the input credit is equal to 1).

795 four). For the next monomial, we decided if we wanted to keep the same
 796 tuple from the table family as first tuple of the new monomial. To do so, we
 797 generated a random float number between zero and one. If the number was
 798 above 0.2, we changed the family tuple.

799 An example can be found in Figure 8, which shows a sample synthetic
 800 provenance polynomial (the how-provenance), the corresponding why-provenance,
 801 lineage, the causality of the tuples of the lineage, together with their respon-
 802 sibility, and, finally, the Shapley values of the lineage tuples. The resulting
 803 credit distribution for each DS is also shown.

804 As an example of how the distribution strategies behave with these syn-

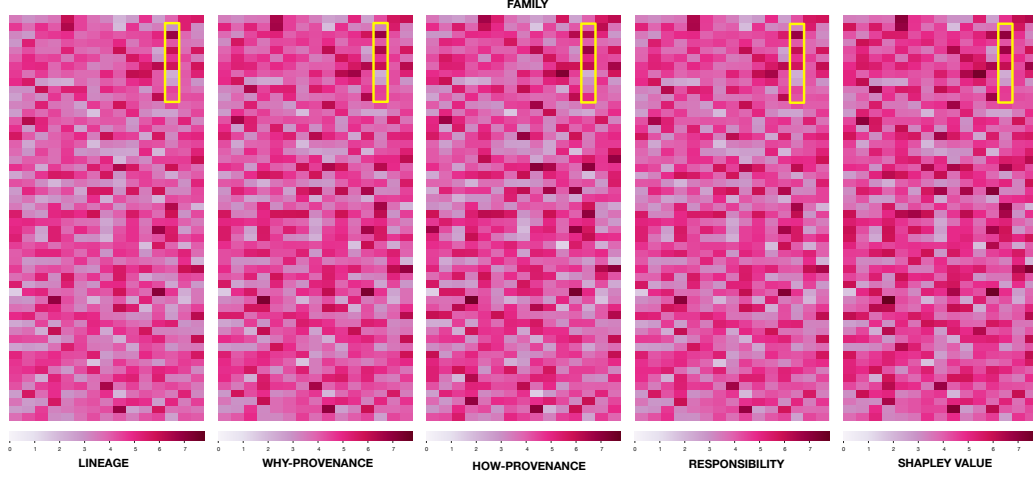


Figure 9: Comparison of three DS on the same table `family` after the distribution computed using 10K synthetic and randomly generated provenance polynomials. The tuples in the blue rectangles are used as example in the discussion connected to Figure 10.

805 thetic queries, consider tuple f_5 in Figure 8. This tuple receives the highest
 806 quantity of credit using responsibility-based distribution and less credit us-
 807 ing, in order, lineage, the Shapley value, why- and how-provenance. On the
 808 other hand, tuple f_1 is rewarded more by the Shapley value, then, in order,
 809 by why-provenance, how-provenance, responsibility, and finally lineage. This
 810 difference is explained considering the different role of the tuples in the gen-
 811 eration of the output and the characteristics of the distributions. Generally
 812 speaking, the more complex the distribution (e.g., the how-provenance), the
 813 more credit is given to tuples that are more frequently used or more critical in
 814 the production of the output. Depending on the situation, i.e. on the syntax
 815 of the query, the distributions may differ among them. Responsibility creates
 816 a ranking among lineage’s tuples describing the importance of their role in
 817 generating the output. As such, the responsibility-based DS gives more credit
 818 to f_1 , f_5 , c_2f_{17} and c_{18} due to their higher responsibility values. “Importance”
 819 is connected to their corresponding minimal contingency sets. For example,
 820 f_1 has a minimal contingency set (one of the many) $\{f_5\}$, with cardinality
 821 1. On the other hand, c_1 has, as minimal contingency set (one of the many)
 822 $\{f_5, c_2\}$, with cardinality two. This means that c_1 is the “least important”
 823 amongst the tuples with minimal contingency sets of lower cardinality, and
 824 this is reflected in the different quantities of credit being distributed.

Table 7: Results of the pairwise Kendall Tau confidence value on all the DSs on the **family** table (the p-values are all below 0.05).

	lineage	why	how	resp.	Shapley
lineage	1.0	0.88	0.73	0.91	0.81
why	0.88	1.0	0.75	0.93	0.92
how	0.73	0.75	1.0	0.74	0.74
resp.	0.91	0.93	0.74	1.0	0.89
Shapley	0.81	0.91	0.74	0.89	1.0

825 The Shapley value behaves similarly, but it rewards tuple f_1 the most and
826 then f_5 , $c2f_{17}$, c_{18} , and last all the other tuples of the lineage. Although both
827 Responsibility and the Shapley value create a ranking of the tuples based on
828 their role in the generation of the output, the corresponding functions behave
829 differently due to the syntax of the query; for this reason each different distri-
830 bution strategy highlights a slightly different aspect that can be considered
831 as “important” when distributing the credit.

832 Despite being synthetic, these provenance polynomials are realistic: they
833 can be obtained by any nested query with join and union operations that use
834 the same tuple multiple times (in which case the exponents are larger than
835 one), and the same combination of operations more than once (in which case
836 the coefficients of monomials are larger than one).

837 *Results.* The results of credit distribution on the **family** table using 10K
838 randomly generated synthetic provenance polynomials are shown in Figure
839 9. We set the maximum value in the heat maps to the highest value reached
840 by a tuple in all five distributions (i.e., 7.7, with the Shapley value-based
841 DS).

842 There is a certain amount of consistency between the strategies in that
843 tuples which are highly rewarded by one strategy are also highly rewarded by
844 the others. This shows that the four DSs consistently reward certain tuples
845 more than others.

846 Table 7 reports the pairwise Kendall τ correlation values¹⁵ for the five
847 DSs computed on the **family** table. As we see, there are certain DS that

¹⁵The Kendall’s τ coefficient is a statistic used to measure the ordinal association between two measured quantities [42]. Intuitively, it is high between two variables when observation have a similar rank.

848 are correlated to others, such as lineage with why-provenance, responsibility
849 and lineage, or responsibility and why-provenance. The others are mildly
850 correlated, such as the Shapley value with how-provenance, responsibility
851 and how-provenance, or why-provenance and lineage with how-provenance.
852 We see, therefore, that the DS based on how-provenance is the one that
853 correlates the least with the other DSs.

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

858 Moving to why-provenance-based DS, we see that more credit is given to
859 tuples in the **family** table than with the previous strategy. This is because
860 the DS considers the different ways that a tuple is used, e.g. in joins with
861 other tuples. If the same tuple is present in more than one witness, it will
862 draw more credit and take it from other tuples in the witness basis. In this
863 case, tuples in **family** drew more credit, taking it from tuples in the other
864 two tables, due to the role that **family** tuples played in the queries that were
865 executed.

866 Consider the how-provenance-based DS heat-map. As with why-provenance,
867 more credit is typically given to tuples in **family** compared to lineage-based
868 DS, since it recognizes the role of these tuples in the queries, and the over-
869 all hue is deeper. The two distributions appear similar, although on closer
870 inspection, slight differences can be seen. This is because how-provenance
871 also considers the frequency with which tuples are used, not only the ways in
872 which they are used. Therefore, although the overall distribution is similar,
873 there are small differences due to the presence of exponents and coefficients
874 in the provenance polynomials, influencing the distribution of credit.

875 The responsibility-based distribution strategy has a distribution that is
876 also quite similar to the one provided by why-provenance (which is also visible
877 from Table 7). It is often the case, for example when the witnesses of the
878 why provenance share one common tuple, that the two distributions behave
879 similarly.

880 Finally, the heat-map reporting the distribution produced by the Shapley
881 value is the one that, at a closer inspection, shows many differences. Although
882 the tuples that receive the biggest quantities of credit are the same, the hue
883 of this tuple is different. The Shapley value in certain circumstances differs
884 greatly from the other DSs, thus showing its ability to weight differently the
885 roles of the tuples.

886 We note that the lineage-based DS gives an average credit of 3.92 to each
887 tuple in the table, while the DS based on why-provenance assigns 4.19, how-
888 provenance 4.18, the one based on responsibility 4.13, and the one based on
889 the Shapley value 4.40. Moreover, lineage distributed a total of about 3121
890 units of credit to the **family** table, why-provenance 3333, how-provenance
891 3331, while responsibility assigned 3290, and the Shapley value 3505. Thus,
892 the Shapley value is the method that accumulates the highest quantity of
893 credit in this table.

894 To better understand the differences between DSs, in the next subsection
895 we consider the accrual of credit over time. In doing so, we will focus on the
896 ten tuples shown within the large yellow rectangles in Figure 10. Each small
897 rectangle within a large yellow rectangle is a tuple, and we number them
898 from 1 (top) to 10 (bottom). These ten tuples were cherry-picked because
899 they allow us to see the evolution of the distribution of credit through time.
900 There are other tuple sets that could have been selected driving us to the
901 same considerations.

902 6.3. Credit accrual over time

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

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

910 The polynomials used are the same as the experiment of the previous
911 section. The range of credit in each map goes from 0 (no credit) to 7 (the
912 maximum quantity of credit reached – using how-provenance – on one of the
913 tuples of the considered window at the “snapshot” with 10K queries). The
914 color hue of the legend, as can be seen, still ranges from 0 to 7.7.

915 By the end of 1K queries, credit differentials between tuples as well as
916 between strategies can be seen. For example, tuple 3 is usually rewarded
917 the most credit by all five strategies. Moreover, it can be seen that tuples 1
918 receives a higher quantity of credit when how-provenance is adopted, show-
919 ing how this form of provenance behaves differently from the others in this
920 context. Moving to 2K queries, it is possible to see that tuple 3 and 7 are
921 still the most rewarded by the strategies.

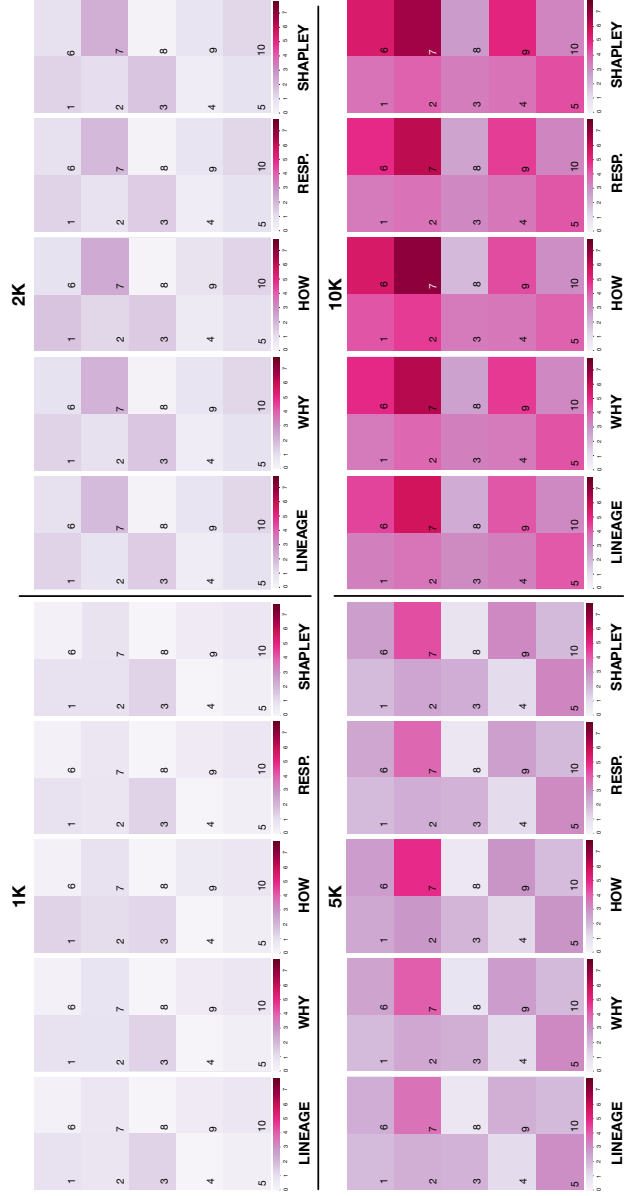


Figure 10: Comparison of the distribution of credit performed by the five DSs on a subset of 10 tuples taken from the `family` table, simulating the passing of time. The number at the top of each group of heat-maps represents the number of polynomials whose credit has been distributed.

922 By the end of 5K queries, tuple 7 emerges with the highest value of
 923 credit with all five DSs, a position which is strengthened with 10K queries.
 924 Moreover, with the passing of time, tuple 3 ceases to be one of the most
 925 rewarded ones and new tuples, such as 6 and 9, emerge as being particularly
 926 rewarded at 5K, while at 10K tuples 6 and 7 are the most rewarded from
 927 the distributions. This is because tuple 7 is used several times within queries
 928 being executed, which is rewarded strongly by why- and how-provenance. We
 929 also note that the responsibility-based distribution confirms its trend of being
 930 similar to why-provenance, although not identical. This is more evident at
 931 step 10K, where tuple 7 is slightly less rewarded using responsibility (6.12)
 932 with respect to why-provenance (6.24). The responsibility that rewards the
 933 more tuple 7 is the one based on how-provenance (credit 7.03), followed by the
 934 Shapley value (credit 6.64). This is due to the fact that tuple 7 had, among
 935 some of the polynomials being used for the experiments, a high responsibility
 936 but it did not appear in all witnesses. This changed slightly the distribution.

937 While the relative value of credit “positions” of tuples within a DS strat-
 938 egy depends on what queries are being executed, the important thing to
 939 notice is the difference between the DSs over time: overall, lineage gives less
 940 credit to tuples in the **family** table than the other strategies since credit
 941 is shared with tuples in other tables. The other strategies recognize the
 942 more important role being played by the **family** tuples than those in the
 943 other tables. The differences between why- and responsibility-based DS are,
 944 for the most times, negligible. The differences between the why- and how-
 945 provenance-based DSs are also relatively minor in most cases. However, there
 946 are certain situations in which the role of a tuple is particularly critical in a
 947 query, and in this case the difference in the value of credit assigned is notably
 948 higher for how-provenance and the Shapley value, as we saw with tuple 7 in
 949 the example of Figure 10.

950 To sum up, the DS based on lineage is sufficient to highlight which tu-
 951 ples in the database are used by a query, and distributes credit equally to
 952 these tuples. The resulting distribution rewards tuples that are used by
 953 more queries, but does not reward how many times tuples are used in the
 954 same query. However, a DS based on why-provenance, responsibility, Shap-
 955 ley value or how-provenance may be better if the queries are complex, since
 956 they reward more tuples that have a critical role in generating the output.
 957 In particular, these four DSs may be useful for finding “hotspots” in the
 958 database based on the role of tuples, with the how-provenance-based and
 959 Shapley value-based DSs being preferable if a higher sensitivity to the role

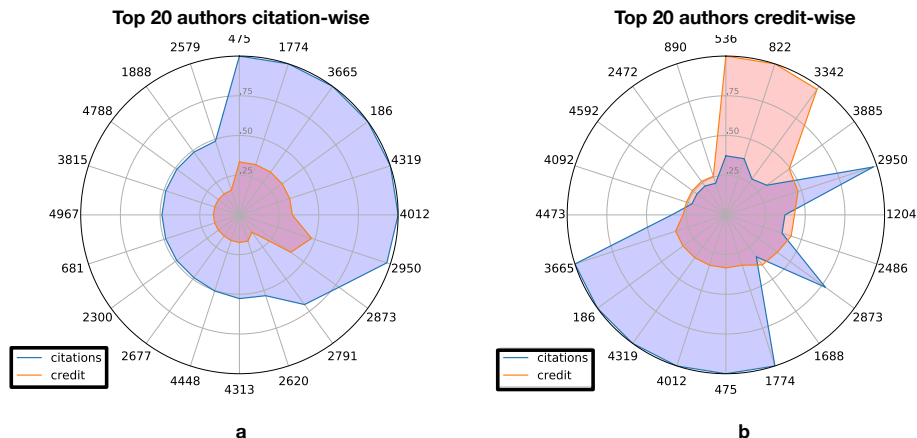


Figure 11: Radars presenting the top 20 authors citation-wise and credit wise, together with their (normalized between 0 and 1) values of citations and credit.

of a tuple in queries is required.

6.4. Credit vs Citations

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

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 webpages 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.

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

989 Diving a bit deeper to focus on the red and blue areas in each of the plots
 990 reveals that there is a significance difference between citations and credit:
 991 The top 20 authors in terms of citations do not have the highest values
 992 of credit (Figure 11.a). Conversely, the authors with the highest values of
 993 credit do not necessarily have a large number of citations (Figure 11.b). For
 994 example, author 536 has the highest value of credit, but is not even in the
 995 top 20 authors in terms of citations. This means that authors like 536, 822,
 996 and 3342 in Figure 11.b receive much more credit from their relatively few
 997 citations than authors like 475, who receives the largest number of citations.
 998 That is, the data underlying certain webpages is more “valuable” in terms
 999 of credit than a citation to the webpage.

1000 The reason for the difference between citation and credit is partly due to
 1001 the experimental setup: each output tuple carries a credit of 1, and there can
 1002 be many tuples used to generate a webpage. Thus a webpage that is created
 1003 from more tuples will have a higher credit value than one created from fewer
 1004 tuples. Furthermore, authors who collaborated with fewer people will receive
 1005 a biggest share of the equally divided credit. However, all authors will receive
 1006 a citation of one.

1007 Credit distribution therefore rewards authors differently than traditional
 1008 citations: an author who has curated larger quantities of cited data and
 1009 collaborated with fewer co-authors, will receive larger quantities of credit.
 1010 Thus, credit rewards them for their larger contribution to the database.

1011 *Results: Synthetic queries.* We used the same synthetic polynomials de-
 1012 scribed in Section 6.2, and we distributed credit with the first 100, 1K, and
 1013 10K of them. Since these polynomials are created by randomly selecting
 1014 tuples from three tables, they usually correspond to a set of data curated by
 1015 authors who, in reality, did not collaborate. To make the size of the author
 1016 set more realistic, we therefore created 20 synthetic authors, and randomly

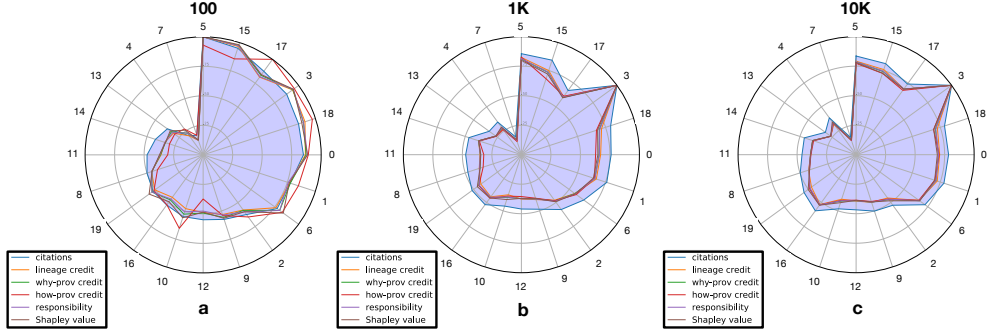


Figure 12: Radars presenting the 20 synthetic authors with corresponding citation and quantities of credit distributed through the 4 DSs (all values normalized between 0 and 1) through different numbers of polynomials (respectively, 100, 1K and 10K). The order is the one defined by figure a, i.e. descending order of citations obtained from 100 polynomials.

assigned one author to blocks of consecutive tuples in the database, with the size of each block varying between 10 and 40, to simulate different quantities of work performed by an author. Every time an author appears as curator of one or more tuples used in a polynomial, we assign them one citation. They also receive four kinds of credit, each one using a different DS.

Figure 12 shows three radar plots, one for the results obtained with 100 polynomials, one with 1K polynomials, one with 10K 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), responsibility-based credit (violet line), and the Shapley value-based credit (brown line).

As can be seen, given the synthetic nature of the queries, the correlation between the number of citations and the quantity of credit assigned to the authors appears to be a much stronger than with the real-world queries of Figure 11. In fact, for Figure 12.a the linear correlation between the citation number and all 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 12.a), while the other three types of credit behave slightly differently (a linear correlation of around 0.95 or above in all other four cases in Figure 12.a). Similar observations can be made for Figure 12.b and 12.c.

What these figures show is that, in certain cases, authors who do not have a large number of citations receive more credit than others, as for example authors 17, 18 and 10 in Figure 12.a, and especially when credit is distributed using how-provenance. This again shows how credit gives a different perspective on the role of data and authors by going beyond the limitations of traditional citations.

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

7. Discussion

Before concluding, we discuss some design decisions: the focus on Credit Distribution (as opposed to Credit Generation), and the choice of Distribution Strategies.

7.1. 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 referenced by the query. A different problem is Credit Generation, the task of generating credit which is then distributed. Credit Generation presents a series of issues which are shared by traditional citation practices. For instance, defining the quantity of credit to be generated for a given citation is still an open problem. Different types of citations may generate different quantities of credit. Data cited as previous work or as useful for previous work may generate less credit than other data extensively used to produce the results presented in a paper. The computation of credit could be done manually (although we must consider the complexity of the task, human biases and the resources required to carry it out) or automatically, but it must be based on a shared definition of impact which is still not agreed upon for data or for traditional citation. For this reason, we used a uniform credit assignment.

There is also the problem of *transitive credit distribution*, i.e., how to transitively propagate credit from one cited unit to another unit that was used to produce the one being cited. For this, a graph of cited units that propagate credit between the units depending on influence could be used. How to propagate credit is an open and non-trivial problem that needs to

consider the importance and impact of a citation in a work, be it a paper or data, and how to eventually compute the quantity of credit to be propagated.

Finally, in our experiments we assumed that the credit carried by an output tuple is one. Thus, each tuple in the output has equal importance. As described above, this assumption may be revised and different credit to different output tuples could be assigned. Note that from the distribution model viewpoint no change is required since the DCD is defined for a generic value k .

7.2. Choice of Distribution Strategies

In this paper we presented four different DSs, so the natural question is which one to use. This depends on the task at hand. When we want to highlight the tuples being used in the database by a workload, the lineage-based DS may be sufficient. When we also want to know the relative impact of tuples in the context of the query, the other DSs should be used since they give a better understanding of the importance of data.

In the real-world based experiments, the four DSs behaved the same, which was due to the specific nature of the data and the queries being used. However, the why-provenance of a query will differ 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. This is usually true when join and projection operators are used in the query.

To address the question of what types of queries are likely to extract cited data, we turn to the results of published studies on the characteristics of query workloads and the complexity of their queries [38, 55, 60]. These studies show that operations such as inner-/outer-joins and projections occur in a significant number of queries. Therefore why- and how-provenances may become quite complex in certain cases and provide a distribution of credit that is significantly different from the one obtained with lineage.

From a computational complexity standpoint, all five DSs are similar since we focused on SPJ queries. Going beyond SPJ queries, Green et al. [32] proposed the provenance semiring framework for SPJRU (Select, Project, Join, Rename, and Union queries), and Amsterdamer et al. [5] showed how to extend the framework to aggregate queries. Since lineage and why-provenance can be computed starting from how-provenance, it is possible to apply the first three DSs proposed in this paper to SPJRU and aggregation queries. Causality and subsequently Responsibility are harder to compute

1110 (NP-complete [48]) for general queries. Credit Distribution is more con-
 1111 cerned with Responsibility, which is in general hard to compute [18]. Meliou
 1112 et al. [48] proved a dichotomy result for conjunctive queries: for each query
 1113 without self-joins, either its responsibility can be computed in PTIME in the
 1114 size of the database, or checking if it has a responsibility below a given value is
 1115 NP-hard. Queries with self-joins are NP-hard in general. This makes respon-
 1116 sibility harder to be utilized for credit distribution in a real-world application,
 1117 since for this problem it is necessary to actually know the responsibility value,
 1118 not simply the ranking amongst tuples.

1119 As for the Shapley Value, Livshits et al. [45] studied the computational
 1120 complexity of calculating the Shapley values in query answering. They origi-
 1121 nally showed mainly lower bounds on the complexity of the problem, with
 1122 the exception of the sub-class of self-join free SPJ queries called *hierarchical*
 1123 queries, where they gave a polynomial-time algorithm. Very recently, Deutch
 1124 et al. [25] proved that the Shapley value can be efficiently (polynomial-time)
 1125 reduced to probabilistic query answering. This not only applies to hierarchi-
 1126 cal queries, but to general SPJ queries. This means that one can compute
 1127 Shapley values using a query engine for probabilistic databases, for exam-
 1128 ple, the practically effective *Knowledge Compilation* [39], making it a viable
 1129 solution for Credit Distribution via SPJ queries.

1130 8. Conclusions and Future Work

1131 This paper defines four new distribution strategies based on why-provenance,
 1132 how-provenance, responsibility, and the Shapley Value, and it compares them
 1133 against the lineage-based distribution strategy defined in [26]. The first, why-
 1134 provenance-based DS, uses the concept of a witness, and gives more credit
 1135 to tuples that appear in more than one witness. In this way, tuples that are
 1136 more important to the query and are used in different ways are rewarded
 1137 more. The second, how-provenance-based DS, considers the frequency with
 1138 which a tuple or combination of tuples is used in the query through the
 1139 information contained in a provenance polynomial. In this case, the how-
 1140 provenance-based DS is more sensitive than the why-provenance-based DS
 1141 to the role and importance of tuples. The third DS exploits the notion of
 1142 responsibility, a real value that ranks the lineage tuples based on their de-
 1143 gree of causality in generating the output. The responsibility-based DS was
 1144 shown to behave similarly to the why-provenance based DS. The fourth DS
 1145 uses the Shapley value function, used to rank the facts of the database, seen

1146 as players, in producing the required result. To do so, the wealth function in
1147 the Shapley value’s definition was adapted for general free-variable queries
1148 on the database.

1149 To show the differences between the five DSs, we performed extensive
1150 experiments based on GtoPdb, a curated scientific relational database, using
1151 both real and synthetic queries. In the first set of experiments, we used select-
1152 project-join (SPJ) queries extracted from citations to webpages in GtoPdb
1153 found in papers published in the British Journal of Pharmacology. Using
1154 these “real” queries, we distributed credit to tuples in different tables of the
1155 database, highlighting tuples that were more frequently used. We showed
1156 that, with these queries, the five strategies produce the same distribution.
1157 This is because the SPJ queries were fairly simple, and did not use self-joins.
1158 Therefore the formulas underlying the different DSs had the same output.

1159 In the second set of experiments, we synthetically produced more com-
1160 plex provenance polynomials, corresponding to more complex queries, that
1161 resulted in exponents and coefficients in the provenance polynomials that
1162 were greater than (or equal to) 1. These experiments highlighted the differ-
1163 ences between the five DSs. While the DS based on lineage rewards all the
1164 tuples used by a query equally, the strategies based on why-provenance and
1165 responsibility give more credit to tuples that are more critical to the query.
1166 In particular, why-provenance considers the different ways in which a tuple
1167 is used in a query, while responsibility considers the relative importance of
1168 a tuple in the generation of the output. The DS based on the Shapley value
1169 similarly rewards the tuples based on their participation. The more impactful
1170 the role of a tuple, the higher its reward in credit. This distribution proved
1171 to be different from the previous two and to reward even more tuples that
1172 are used in more than one witness. How-provenance is even more sensitive
1173 to the tuple’s role: it also considers the frequency with which a tuple or a
1174 set of tuples is used.

1175 In the third set of experiments, we showed how the differences between
1176 the DS are compounded over time, i.e. when more and more queries are
1177 processed by the system.

1178 In the fourth set of experiments we compared traditional citations to
1179 authors to the credit accrued to them via the DSs. We showed how, in
1180 both real-world and synthetic scenarios, credit rewards authors who con-
1181 tribute/curate data that has the highest impact, and therefore receives the
1182 biggest quantity of credit, and not necessarily the data with the highest ci-
1183 tation count. In this sense, credit appears to be an useful new measure to

discover data and their corresponding curators that have a high impact in the research world, even when they are cited few times or do not appear at all in the data that are cited (i.e. the case of data used to build the output of a query but that is not visualized in the output itself).

In future work, we plan to explore different strategies to generate and distribute credit. In this paper we assumed that each output tuple carries credit 1. In more sophisticated scenarios we can employ different strategies to compute credit, that reflect the importance of cited data. Other, more sophisticated, strategies could also be used to decide how credit is distributed between the authors, beyond the uniform distribution used here, in a way to reflect the work performed by them on the cited data. There are also a number of other intriguing applications for credit over relational databases. One such application is *data pricing*, which gives a price to a query submitted by a user who wants to buy the produced information. Currently, a common strategy used for data pricing is based on query rewriting: A database stores a set of views with their price. When a new query arrives, the system rewrites it using the stored views to obtain a query price, a process that can be computationally expensive. We plan to distribute credit through carefully planned and representative queries, and use credit information to define a new, faster, and potentially more flexible pricing function.

Another application is *data reduction* [49], which addresses the problem of reducing the vast – and rapidly expanding – amount of data that is being produced. Data credit can help address this problem by identifying “hotspots” and “coldspots” of data. A hotspot is data in a database (e.g. a tuple) with a high quantity of credit, which is therefore valuable for the set of queries that execute frequently over the data and distribute the credit. A coldspot is data with a low quantity of credit which can therefore be considered as less important, and could be deleted, summarized, or moved to cheaper and/or less efficient memory.

Acknowledgement

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

References

- [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bernstein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L.,

- 1219 Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Koss-
1220 mann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo,
1221 T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo,
1222 A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suciu, D.
1223 (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–
1224 53.
- 1225 [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating
1226 data citation in citedb. *PVLDB*, 10(12):1881–1884.
- 1227 [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y.
1228 (2018). Data citation: A new provenance challenge. *IEEE Data Eng.*
1229 *Bull.*, 41(1):27–38.
- 1230 [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An
1231 Introduction to the Joint Principles for Data Citation. *Bulletin of the*
1232 *Association for Information Science and Technology*, 41(3):43–45.
- 1233 [5] Amsterdamer, Y., Deutch, D., and Tannen, V. (2011). Provenance for ag-
1234 gregate queries. In Lenzerini, M. and Schwentick, T., editors, *Proceedings*
1235 *of the 30th ACM SIGMOD-SIGACT-SIGART Symposium on Principles*
1236 *of Database Systems, PODS 2011*, pages 153–164. ACM.
- 1237 [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D.,
1238 Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I.,
1239 Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and
1240 Goble, C. A. (2013). Why linked data is not enough for scientists. *Future*
1241 *Gener. Comput. Syst.*, 29(2):599–611.
- 1242 [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation
1243 Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.
- 1244 [8] Bertin-Mahieux, T., Ellis, D., Whitman, B., and Lamere, P. (2011). The
1245 million song dataset. In *Proceedings of the 12th International Conference*
1246 *on Music Information Retrieval (ISMIR 2011)*, pages 591–596.
- 1247 [9] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In
1248 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*
1249 *nication*, pages 93–116. De Gruyter Mouton.

- [10] Buneman, P. (2006). How to cite curated databases and how to make them citable. In *18th International Conference on Scientific and Statistical Database Management, SSDBM*, pages 195–203. IEEE Computer Society.
- [11] 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.
- [12] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation is a computational problem. *Commun. ACM*, 59(9):50–57.
- [13] 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.
- [14] Buneman, P. and Silvello, G. (2010). A rule-based citation system for structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.
- [15] 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.
- [16] 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.
- [17] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases: Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–474.
- [18] Chockler, H. and Halpern, J. Y. (2004). Responsibility and blame: A structural-model approach. *J. Artif. Intell. Res.*, 22:93–115.
- [19] 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.

- 1281 [20] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N.
 1282 (2019). Bringing citations and usage metrics together to make data count.
 1283 *Data Science Journal*, 18(1).
- 1284 [21] Cronin, B. (1984). *The Citation Process. The Role and Significance of*
 1285 *Citations in Scientific Communication*. London: Taylor Graham.
- 1286 [22] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evi-
 1287 dence of a structural shift in scholarly communication practices? *JASIST*,
 1288 52(7):558–569.
- 1289 [23] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of
 1290 view data in a warehousing environment. *ACM Trans. Database Syst.*,
 1291 25(2):179–227.
- 1292 [24] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model
 1293 for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on*
 1294 *Innovative Data Systems Research*. www.cidrdb.org.
- 1295 [25] Deutch, D., Frost, N., Kimelfeld, B., and Monet, M. (2021). Computing
 1296 the shapley value of facts in query answering.
- 1297 [26] Dosso, D. and Silvello, G. (2020). Data credit distribution: A
 1298 new method to estimate databases impact. *Journal of Informetrics*,
 1299 14(4):101080.
- 1300 [27] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The vir-
 1301 tual atomic and molecular data centre (VAMDC) consortium. *Journal of*
 1302 *Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- 1303 [28] Eiter, T. and Lukasiewicz, T. (2002). Complexity results for structure-
 1304 based causality. *Artif. Intell.*, 142(1):53–89.
- 1305 [29] Fang, H. (2018). A discussion of citations from the perspective of the
 1306 contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–
 1307 1520.
- 1308 [30] Garfield, E. (1999). Journal impact factor: a brief review. *Can. Med.*
 1309 *Assoc.*, 979-980.

- [31] 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.
- [32] 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.
- [33] Halpern, J. Y. and Pearl, J. (2013). Causes and explanations: A structural-model approach — part 1: Causes. *CoRR*, abs/1301.2275.
- [34] 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.
- [35] Hartley, J. (2017). Authors and their citations: a point of view. *Scientometrics*, 110(2):1081–1084.
- [36] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a transformed scientific method.
- [37] 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.
- [38] Jain, S., Moritz, D., Halperin, D., Howe, B., and Lazowska, E. (2016). Sqlshare: Results from a multi-year sql-as-a-service experiment. In *Proceedings of the 2016 International Conference on Management of Data*, pages 281–293.
- [39] Jha, A. K. and Suciu, D. (2013). Knowledge compilation meets database theory: Compiling queries to decision diagrams. *Theory Comput. Syst.*, 52(3):403–440.
- [40] 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,

- 1342 S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of bio-
1343 logical pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- 1344 [41] Katz, D. (2014). Transitive credit as a means to address social and
1345 technological concerns stemming from citation and attribution of digital
1346 products. *Journal of Open Research Software*, 2(1).
- 1347 [42] Kendall, M. G. (1938). A new measure of rank correlation. *Biometrika*,
1348 30(1/2):81–93.
- 1349 [43] Kosten, J. (2016). A classification of the use of research indicators.
1350 *Scientometrics*, 108(1):457–464.
- 1351 [44] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S.
1352 (2011). Citation and Peer Review of Data: Moving Towards Formal Data
1353 Publication. *International Journal of Digital Curation*, 6(2):4–37.
- 1354 [45] Livshits, E., Bertossi, L. E., Kimelfeld, B., and Sebag, M. (2020). The
1355 shapley value of tuples in query answering. In Lutz, C. and Jung, J. C.,
1356 editors, *23rd International Conference on Database Theory, ICDT 2020,*
1357 *March 30-April 2, 2020, Copenhagen, Denmark*, volume 155 of *LIPIcs*,
1358 pages 20:1–20:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.
- 1359 [46] Martone, M. (2014). Joint declaration of data citation principles.
1360 *FORCE11. San Diego CA. Data Citation Synthesis Group*. [https://www.](https://www.force11.org/datacitationprinciples)
1361 [force11.org/datacitationprinciples](https://www.force11.org/datacitationprinciples), online September 2020.
- 1362 [47] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation
1363 counts and rankings of LIS faculty: Web of science versus scopus and
1364 google scholar. *Journal of the american society for information science*
1365 *and technology*, 58(13):2105–2125.
- 1366 [48] Meliou, A., Gatterbauer, W., Moore, K. F., and Suciu, D. (2010). The
1367 complexity of causality and responsibility for query answers and non-
1368 answers. *Proc. VLDB Endow.*, 4(1):34–45.
- 1369 [49] Milo, T. (2019). Getting rid of data. *Journal of Data and Information*
1370 *Quality (JDIQ)*, 12(1):1–7.
- 1371 [50] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D.,
1372 Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G.,

- 1373 Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff,
1374 D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,
1375 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson,
1376 E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C.,
1377 Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers,
1378 E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research
1379 culture. *Science*, 348(6242):1422–1425.
- 1380 [51] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J.
1381 (2016). Research data explored: An extended analysis of citations and
1382 altmetrics. *Scientometrics*, 107(2):723–744.
- 1383 [52] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic,
1384 large databases: Model and reference implementation. In *Proceedings of*
1385 *the 2013 IEEE International Conference on Big Data, 6-9 October 2013,*
1386 *Santa Clara, CA, USA*, pages 307–312.
- 1387 [53] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identifi-
1388 cation of Reproducible Subsets for Data Citation, Sharing and Re-Use.
1389 *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue*
1390 *on Data Citation*, 12(1):6–15.
- 1391 [54] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data
1392 citation of evolving data: Recommendations of the working group on data
1393 citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- 1394 [55] Remil, Y., Bendimerad, A., Mathonat, R., Chaleat, P., and Kaytoue,
1395 M. (2021). ” what makes my queries slow?”: Subgroup discovery for sql
1396 workload analysis. *arXiv preprint arXiv:2108.03906*.
- 1397 [56] Shapley, L. S. (1954). A value for n-person games. In Kuhn, H. W. and
1398 Tucker, A. W., editors, *Contributions to the Theory of Games II*, pages
1399 307–317. Princeton University Press, Princeton.
- 1400 [57] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*
1401 *Sci. Technol.*, 69(1):6–20.
- 1402 [58] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data
1403 provenance in e-science. *SIGMOD Record*, 34(3):31–36.

- 1404 [59] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-
 1405 spective. In of Sciences’ Board on Research Data, N. A. and Information,
 1406 editors, *Report from Developing Data Attribution and Citation Practices*
 1407 *and Standards: An International Symposium and Workshop*, pages 177–
 1408 178. National Academies Press: Washington DC.
- 1409 [60] Vogelsgesang, A., Haubenschild, M., Finis, J., Kemper, A., Leis, V.,
 1410 Mühlbauer, T., Neumann, T., and Then, M. (2018). Get real: How bench-
 1411 marks fail to represent the real world. In *Proceedings of the Workshop on*
 1412 *Testing Database Systems*, pages 1–6.
- 1413 [61] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,
 1414 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,
 1415 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data
 1416 management and stewardship. *Scientific data*, 3.
- 1417 [62] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data
 1418 citation: Giving credit where credit is due. In *Proceedings of the 2018*
 1419 *International Conference on Management of Data, SIGMOD*, pages 99–
 1420 114.
- 1421 [63] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to
 1422 scientific datasets using article citation networks. *Journal of Informetrics*,
 1423 14(2).
- 1424 [64] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of
 1425 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.
- 1426 [65] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for
 1427 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*
 1428 *Molecular Spectroscopy*, 327:122–137.