

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 a basic form of research product for which citation, and the generation of credit or recognition for authors, are still not well understood. The notion of *data credit* has therefore recently emerged as a new measure, defined and based on data citation groundwork.

Data credit is a real value representing the importance of data cited by a research entity. We can use credit to annotate data contained in a curated scientific database and then as a proxy of the significance and impact of that data in the research world. It is a method that, together with 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 producing 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 focus on Select-Project-Join (SPJ) queries under bag semantics, and we define three distribution strategies based on how-provenance, responsibility, and 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 curators. In particular, credit rewards data and authors based on their research impact, not only on the citation count. 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 that allow us to find research products and create and understand their relationships. They form a basis to give credit to authors, papers, and venues [21, 22, 68]. Citations are used, among other things, to decide on tenure, promotion, hiring, and funding of grants for researchers [23, 36, 45, 50].

Science and research are increasingly digital, and numerous curated databases are at the core of scientific research efforts [13]. It is therefore generally accepted that data must be cited and citable [16, 46], and that data citations should contribute to the scientific reputation of researchers, scientists, data curators, and creators [4, 62]. It is also accepted that data citations should be counted alongside traditional citations and contribute to bibliometrics indicators [7, 54].

A central problem with data citation is how to attribute credit to data creators and curators [12]. How to handle and count the credit generated by data citation and how it contributes to traditional and new bibliometrics are long-standing research issues [10, 31]. However, data citations and their related bibliometrics do not always fully reward the creators of data used in a database, even when correctly applied. Data is often cited at the “database level” or the “webpage level”. In the first case, even though only a data subset was used, the whole database ends up being cited, and therefore all credit goes only 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 are built using data extracted from the database, which is aggregated by topic and layout 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) [30, 42, 67] 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 research entity (typically a paper), and is a proxy for the impact of the data on the citing entity. The DCD problem consists of distributing this credit to elements in the databases 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 a (curated) scientific database. Credit is a real value that can be assigned to

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

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

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

61 We study the DCD problem in the context of relational databases (RDBs)
62 since they are widely used² and are the main focus of current work in data
63 citation methods [13, 15, 55]. RDBs are also frequently a test-bed for new
64 methods that can be adapted to other databases, e.g., graphs or document
65 databases. The “portions” of data in an RDB that can be credited can be
66 defined at different levels of granularity, in particular: (i) the whole database,
67 (ii) tables, (iii) tuples, and (iv) attributes. The ability to specify different
68 levels of granularity in a relational database allows us to define the DCD
69 problem at a particular level of granularity. In this paper, we focus on DCD
70 at the tuple level.

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].

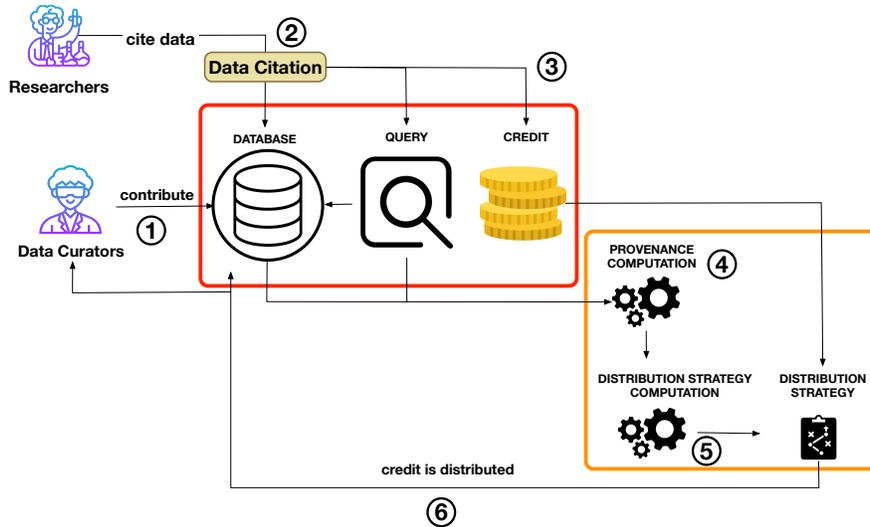


Figure 1: Overview of the credit distribution pipeline.

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

72 **Step 1** Scientists and experts create and curate the information contained
 73 in a scientific database. These are called the “Data Curators”.

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

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

79 **Step 4** Given the database instance I and the query Q , the *data provenance*
 80 of $Q(I)$ is computed as a form of metadata that captures how Q used
 81 I to generate the output [18].

82 **Step 5** Provenance is input to the *Credit Distribution Strategy* (CDS, also
 83 referred only as Distribution Strategy, DS). CDS is a function f that
 84 takes as input the credit k , distributes it to the data in the input
 85 database I , and is defined on the basis of citation policies decided at
 86 the database administration level or at the domain community level.

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

91 This paper expands the work in [27] where we first defined the problem
92 of DCD in relational databases, and proposed a viable Distribution Strategy
93 (DS) based on *lineage* – the simplest form of *data provenance*. The lineage
94 of a tuple t in the output $Q(I)$ is defined as the set of all and only the tuples
95 in the database instance I that are “relevant” to the production of t and
96 indicated as L_t . The corresponding strategy equally redistributes the credit
97 k to the tuples in the lineage set, thus each tuple receives credit $k/|L_t|$.

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

104 Therefore, in this paper, we expand the ideas in [27] by proposing new
105 DSs based on another form of data provenance: how-provenance [33]. We
106 also propose other two DS based on the concepts of responsibility [51] and
107 the Shapley value [26, 48]. We focus on SQL queries under the bag semantics
108 assumption.

109 We discuss why one provenance form may be preferred to another de-
110 pending on the application and its goals. We show that the DS based on
111 responsibility gives more credit to tuples that are essential to the production
112 of the result set. In contrast, the how-provenance-based DS considers the
113 different ways in which a tuple is used. Finally, we present an alternative
114 take on the problem with the Shapley-based DS that models the distribution
115 process as a competitive game in which tuples that contribute more to the
116 generation of the output are correspondingly rewarded more.

117 We use a well-known curated database called the IUPHAR/BPS³ Guide
118 to Pharmacology [35] – GtoPdb⁴ – to evaluate the DSs. GtoPdb contains
119 expertly curated information about diseases, drugs, cellular drug targets, and

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

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

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

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

134 **Contributions** of this work include:

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

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

149 2. Background

150 *Data in Research.* Research transitioned to the *fourth paradigm of science* [37],
151 that is, data-intensive scientific discovery, where data are essential for scien-
152 tific advances as well as for traditional publications [6].

153 The scientific community is promoting an *open research culture* [53],
154 founded on methods and tools to share, discover, and access experimental
155 data. A striking example is the FAIR principles (Findable, Accessible, Inter-
156 operable, and Reusable) [64], which every database should enforce. In par-
157 ticular, data should be accessible from the articles, journals, and papers that
158 cite or use them [21]. The need for *reproducibility* of experiments through the
159 used data; the *availability* of scientific data; and, the *connections* between
160 data and the scientific results are all needed aspects to operationalize the
161 fourth paradigm, and relevant for *data citation* [38].

162 *Data Citation: Principles and Motivations.* Data Citation principles were
163 proposed in [20], and later summarized and endorsed by the Joint Declara-
164 tion of Data Citation Principles (JDDCP) [49]. The principles are divided
165 into two groups [60]. The first group is about the role of data citation in
166 scholarly and research activities such as the (i) *importance* of data (why data
167 citation is important and why data should be considered as first-class citi-
168 zens); (ii) *credit* and *attribution* to the creators and curators of the data;
169 (iii) *evidence*; (iv) *verifiability*; and *interoperability*, with these last three re-
170 quiring data citation methods to be flexible enough to operate through differ-
171 ent communities. The second group defines the main guidelines to establish
172 a data citation systems, and contains principles such as the (i) *unique iden-*
173 *tification* of the data being cited; (ii) (*open*) *access* to data; (iii) guarantee
174 of *persistence* and *availability* of citations even after the lifespan of the cited
175 entity; the (iv) *specificity* of a citation, i.e. it must lead to the data set
176 originally cited.

177 The main motivations for data citation are outlined in [60] and range from
178 data attribution and connection to data sharing, impact and reproducibility.

179 2.1. Data Citation in Relational Databases

180 Relational databases have been the target of data citation methods since
181 the surge of the data-centric research paradigm. The RDA “Working Group
182 on Data Citation: Making Dynamic Data Citable”⁵ [56] (hereafter, RDA-
183 WGDC) has developed guidelines for citing large, dynamic, and changing
184 datasets which have now moved on into adoption phase. The datasets con-
185 sidered by the Working Group are often relational.

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

186 The RDA-WGDC [57] reported that there are various implementations
187 of its guidelines for Data Citation on MySQL/Postgres relational databases.
188 Some of these databases are: DEXHELPP⁶ (Social Security Records); NERC
189 (ARGO Global Array); EODC (Earth Observation Data Centre) [32]; LNEC
190 (River dam monitoring); MDS (Million Song Database) [9]; CBMI⁷ (Center
191 for Biomedical Informatics); VMC (Vermont Monitoring Cooperative); CCA⁸
192 (Climate Change Center Austria); VAMDC (Virtual Atomic and Molecular
193 Data Center) [28, 69].

194 More examples of work on data citation in relational databases are [2,
195 13, 25, 65]. The website <https://fairsharing.org/> keeps an updated list
196 of curated and scientific databases (many of which are relational or graph-
197 based) following FAIR guidelines. These databases are citable since they are
198 compliant with the most recent guidelines, and they are in the vast majority
199 of cases accessible via dynamically created webpages. In all these databases it
200 is, therefore, possible to implement DCD on top of the existing infrastructures
201 for citing data.

202 Data citation techniques are primarily applied to relational databases
203 because of their pervasiveness as well as the “identifiability” of the portions
204 of data that are to be cited: the whole database, a relation, a tuple, or
205 even an attribute. Many papers [2, 11, 13] consider more complex citable
206 units, recognizing that often the *views* of a database are the ones to be cited.
207 Generally, a *view* is a query on the database. To this end, [65] suggested
208 decomposing the database into a set of views, where each view is associated
209 with its citation.

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

- 211 1. A database cited as a whole, even though only parts of the databases
212 are used in the papers or datasets. Alternatively, the so-called “data pa-
213 pers” are cited, being traditional papers that describe a database [17].
214 In this case, all the credit from the citations goes to the database ad-
215 ministrators or to the authors of the data papers.
- 216 2. Subsets of data, obtained by issuing queries to a database, are individ-
217 ually cited. This is the solution adopted by the RDA-WGDC [56]. In

⁶<http://www.dexhelpp.at/>

⁷[https://medicine.missouri.edu/centers-institutes-labs/
center-for-biomedical-informatics](https://medicine.missouri.edu/centers-institutes-labs/center-for-biomedical-informatics)

⁸<https://ccca.ac.at/startseite>

218 this case, the credit generated from citations is distributed among the
219 contributors of the portions of data being cited, and/or to the database
220 administrators.

221 3. The database is accessible via a series of webpages that arrange the
222 content of the database by topic or theme. Examples in the life science
223 domain include the Reactome Pathway database [41], the GtoPdb [35],
224 and the VAMDC [69]. Every single Webpage is unequivocally identifi-
225 able and can be individually cited.

226 2.2. Data Credit

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

231 In this framework, Katz [42] suggests the need for a *modified citation*
232 *system* that includes the idea of *transient* and *fractional credit*, to be used
233 by developers of research products as software and data. Two considerations
234 are made: (i) research objects such as data and software are currently not
235 formally rewarded or recognized by the community; (ii) even in traditional
236 papers, the contribution of each author to the work is hard to understand,
237 unless explicitly specified in the paper. This is even more true for data, where
238 different groups of people work on the same database.

239 In [42] credit is defined as a “quantity” that describes the importance of a
240 research entity, such as papers, software, or data, mentioned in a citation. It
241 also proposed the idea of a *distribution* of credit from research entities, such
242 as papers or data, to other research entities through citations. Therefore,
243 when discussing data credit, we need to consider *credit computation* – i.e.,
244 the process to compute the quantity of credit generated by the citation – and
245 *credit distribution* – i.e., the process to distribute credit and to assign it to
246 the entities that contributed to the creation/curation of the cited data. In
247 this paper we focus on the latter.

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

253 Katz [42] further explores the idea of a *distribution* of credit from re-
254 search entities (i.e., papers and data) to other research entities through cita-

255 tions that connect them. Thanks to the idea of “credit distribution”, some
256 problems related to traditional citations can be addressed:

- 257 1. Credit rewards research entities that to date are not (formally) recog-
258 nized (a goal shared with data citation).
- 259 2. Credit can reward authors *proportionally* to their role in generating the
260 entity. The more an author contributes to a paper, the more credit is
261 given to him. Zou and Peterson [68] work on something similar with
262 their zp-index, which includes in its formulation the position (and thus
263 the role) of a publication author to represent its impact in the work
264 itself.
- 265 3. Credit can be *transitively* channeled through a chain of papers citing
266 each other, thus enabling the rewarding of older papers that are no
267 more cited, since other papers summarize or report their content but
268 are nevertheless crucial in a research area for the influence of their
269 content.

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

279 Zeng et al. [67] proposes the first method to compute credit within a
280 network of papers citing data. Adopting a network flow algorithm, they
281 simulate a random walk to estimate a score for each dataset, leveraging real-
282 world usage data to compute the credit. This is the first step towards an
283 automatic credit computation procedure. This proposal is, however, limited
284 to assigning credit to whole datasets, and it does not deal with the granu-
285 larity of data. It does not work to assign credit to a single research entity
286 within a dataset. Differently from Zeng et al. [67], we do not treat the credit
287 computation process, but we focus on the distribution process.

288 2.3. Data Provenance

289 To distribute credit, we base our methods on the *data provenance* body-
290 work. Data provenance is information that describes the origin and the

291 process of creation of data. It can also be seen as metadata pertaining to
292 the derivation history of the data. It is particularly useful to help users to
293 understand where data are coming from, and the process they went through.
294 Data citation and data provenance are closely linked [3] since both are forms
295 of annotations on data retrieved through queries. Data provenance has been
296 widely studied in different areas of data management. In this paper, we fo-
297 cus on provenance for database management systems (DBMS). For further
298 details on data provenance, please refer to surveys like [18] and [61].

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

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

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

310 Green et al. [33] defined the semiring model which captures all of the
311 above provenance models – lineage, why-provenance, how-provenance and
312 where-provenance – and expresses set semantics, bag semantics and some
313 extensions of the relational model. For data credit distribution, the results
314 achieved with lineage and why-provenance are subsumed by those obtained
315 using how-provenance, which we focus on in this work.

316 2.4. Causality and Responsibility

317 We also consider the notions of causality and responsibility, as defined
318 in [51]. 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 [19], Meliou
325 et al. [51] 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 [59], 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 [48]: 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 [48]. The initial theoretic-
345 al analysis in [48] showed lower bounds on the complexity of the problem,
346 but did not suggest a feasible implementation. However, very recently, an
347 efficient implementation for Boolean queries has been provided [26], both in
348 terms of an exact computation (it works well for most queries) and in inexact
349 one (it is extremely fast and provides the same ranking of tuples as the exact
350 computation, but not necessarily the same values).

351 3. Use Case: GtoPdb

352 The IUPHAR/BPS Guide to Pharmacology [35] (GtoPdb⁹) is a well-
353 known scientific relational database that contains expertly curated informa-
354 tion about diseases, drugs in clinical use, their cellular targets, and the mech-
355 anisms of action on the human body. It is curated and maintained by the
356 GtoPdb Committee and 96 subcommittees, comprising 512 scientists collab-
357 orating with in-house curators who draw the information contained in the
358 database from high-quality pharmacological and medicinal chemistry litera-
359 ture. Roughly 1000 researchers from all over the world have contributed to

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

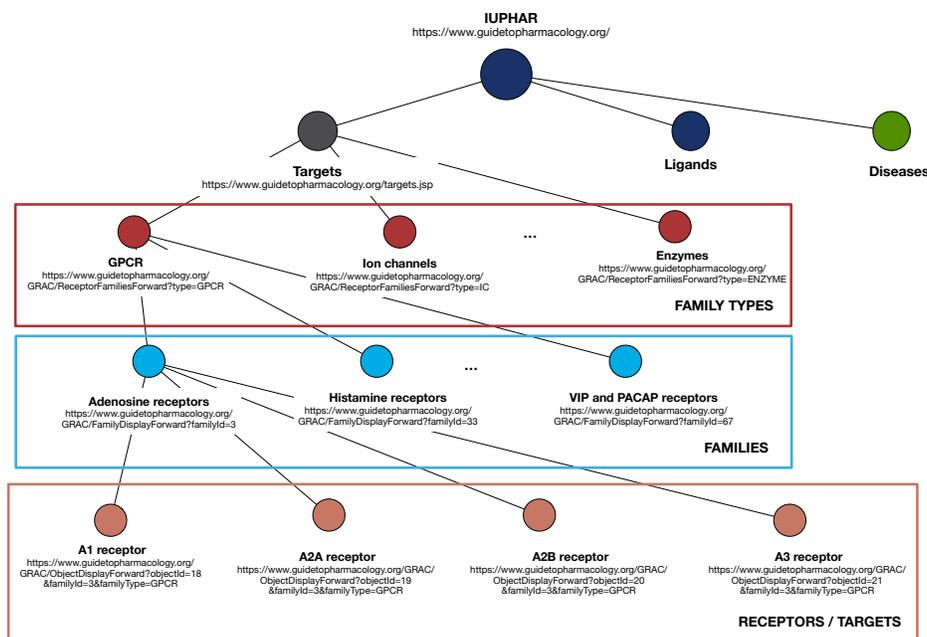


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

360 the database, and the curators wanted to give recognition to these contribu-
 361 tors. This led to some early work on data citation [11].

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

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

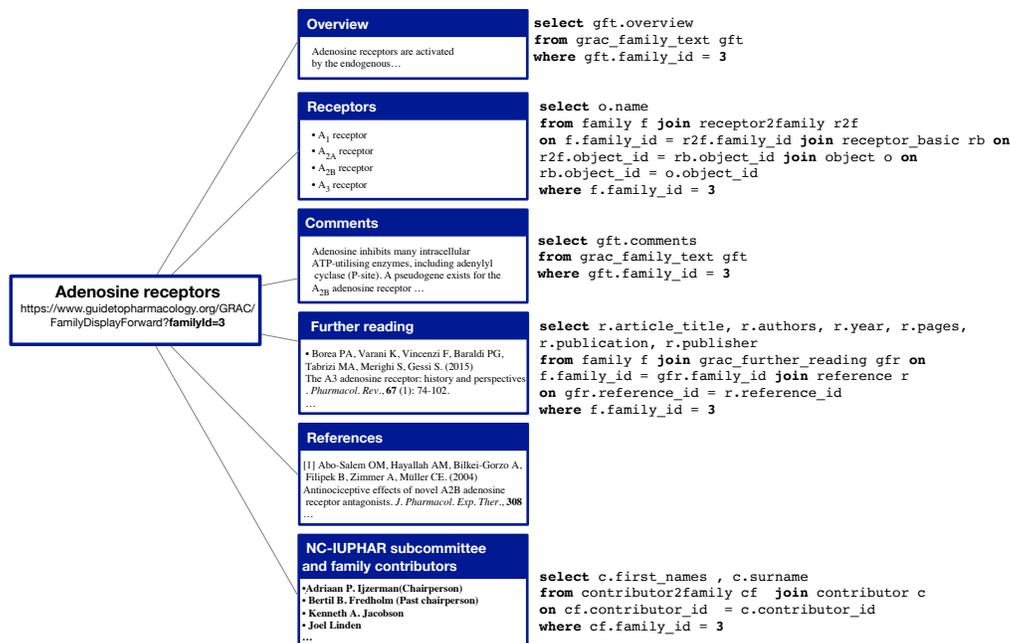


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

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

386 In practice, many papers in the literature only reference GtoPdb (the
 387 root) without including a reference to the specific page being cited. That is,
 388 they only cite a paper describing GtoPdb as a whole (e.g., [35]) and refer
 389 to targets, ligands, diseases, etc. only by name. Thus, citations to specific
 390 families are *de-facto* “hidden” to citation systems such as Google Scholar,
 391 and useless for the computation of bibliometrics.

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

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

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

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

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

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

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

I	database instance
L, 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{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.

410 4. Provenance, Responsibility, and Shapley Value

411 We now introduce how-provenance, causality, responsibility, and the Shap-
412 ley value function. In the following we use the notion of the *lineage* of an
413 output tuple [18, 24]. The *lineage set* L of a tuple $o \in Q(I)$ is the set of *all*
414 and *only* the tuples in the database instance I that are used by query Q to
415 produce the output tuple o .

416 4.1. How-Provenance

417 How-provenance was first defined in [33] to capture the information about
418 how the source tuples are used exploiting a *semiring* algebraic structure. It
419 takes the form of a polynomial, called *provenance polynomial*, where the
420 variables are taken from the set X of identifiers of the tuples (provided that
421 each tuple in I has an identifier) and the coefficients are drawn from the set
422 of natural numbers \mathbb{N} .¹¹

423 In the following, we rely on the commonly-adopted notation of [18]. Let
424 \mathbf{D} be a finite domain of data values $\{d_1, \dots, d_n\}$ and \mathcal{U} a collection of *field*
425 *names* (also attribute names). We use U, V to denote finite subset of \mathcal{U} .

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

426 A tuple t is a function $U \mapsto \mathbf{D}$, from the attributes $\{A_1, \dots, A_n\} \in U$ to
 427 the data values in \mathbf{D} , written as $(A_1 : d_1, \dots, A_n : d_n)$. A tuple assigning
 428 values to each field name in U is called U -tuple. We write $Tuple$ for the set
 429 of all tuples, U - $Tuple$ for the set of all U -tuples. We write $t.A$ or $t \bullet A$ for
 430 the value of the A -field of t and $t[U]$ for the restriction of tuple t over $U \subseteq V$
 431 to field names in U . We write $t[A \mapsto B]$ for the result of renaming field A to
 432 B in t (assuming B is not already present in t).

433 A *relation* or table $r : U$ is a finite set of tuples over U . We call \mathcal{R} a finite
 434 collection of *relation names*. A *schema* \mathbf{R} is a mapping $(R_1 : U_1, \dots, R_n : U_n)$
 435 from \mathcal{R} to a finite subsets of \mathcal{U} (assigning to a every relation name a set of
 436 attributes). A *database* or *instance* I is a function $I : (R_1 : U_1, \dots, R_n : U_n)$
 437 mapping each $R_i : U_i \in \mathbf{R}$ to a relation r_i over U_i .

438 A *tuple location* is defined as a tuple in one relation of the database
 439 tagged with its name. A tuple location is indicated with (R, t) , where R is
 440 the relation in the database, and t is the tuple in R . With reference to the
 441 running example of Table 1, $(\mathbf{family}, \langle f_1, \text{Dopamine Receptors}, \text{gpcr} \rangle)$
 442 is the tuple location of the first tuple in the \mathbf{family} relation. The set of all
 443 the tuple locations in I is called *TupleLoc*.

444 A semiring K is a *set* equipped with two operations, typically denoted
 445 with the symbols $+$ and \cdot , satisfying the following axioms [8, pg. 26]:

- 446 1. The set K is a *commutative monoid* for the operator $+$ with a neutral
 447 element 0 . Therefore, it has these properties:
 - 448 (a) $(a + b) + c = a + (b + c)$ (associative property)
 - 449 (b) $0 + a = a + 0 = a$
 - 450 (c) $a + b = b + a$ (commutative property)
- 451 2. The set K is a *monoid* with identity element 1 . Therefore, it has these
 452 properties:
 - 453 (a) $(a \cdot b) \cdot c = a \cdot (b \cdot c)$ (associative property)
 - 454 (b) $1 \cdot a = a \cdot 1 = a$ (1 is the neutral element)
- 455 3. Multiplication is distributive on addition, i.e.:
 - 456 (a) $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$
 - 457 (b) $(a + b) \cdot c = (a \cdot c) + (b \cdot c)$
- 458 4. Multiplication by 0 annihilates K , i.e. $\forall x \in K, 0 \cdot x = x \cdot 0 = 0$.

459 The key idea in Green et al. [33] is to use the two operators $+$ and \cdot to
 460 represent two basic transformations that source tuples undergo as a result
 461 of applying a relational query to a database [18]. Two tuples may either

462 be joined together (a join is represented with the \cdot operator) or merged via
 463 union or projection (represented with the $+$ operator).

464 Now we formally introduce the mathematical framework behind how-
 465 provenance [33]. Let K be a set containing an element 0 , U a set of attributes
 466 and U -Tuples the set of tuples with attributes in the set U (each such tuple is
 467 called, for brevity, U -tuple). A K -relation is a function $R : U\text{-Tuples} \mapsto K$
 468 which maps every U -tuple in an element in K such that its support, defined
 469 as $\text{supp}(R) = \{t \mid R(t) \neq 0\}$, is finite. Thus, it is possible to see the K -
 470 relation as a finite function that models a relation R , tagging each tuple in
 471 R with an element of K and each tuple that is not in R with 0 .

472 **Definition 4.1.** *Operations on the algebraic structure $(K, 0, 1, +, \cdot)$ [33]*

473 *Let $(K, 0, 1, +, \cdot)$ be an algebraic structure with two binary operations $+$ and*
 474 *\cdot and two distinguished elements 0 and 1 . The operations of the positive*
 475 *K -relational algebra are defined as follows:*

- 476 1. *Empty relation.* For any set of attributes U , $\exists \emptyset : U\text{-Tuples} \mapsto K \mid \emptyset(t) =$
 477 0 .
- 478 2. *Selection* Let $R : U\text{-Tuples} \mapsto K$ and σ be a selection predicate that
 479 maps each U -Tuple to either 0 or 1 . Then $\sigma_\theta(R) : U\text{-Tuples} \mapsto K$ is
 480 defined by $(\sigma_\theta(R))(t) = R(t) \cdot \sigma(t)$.
- 481 3. *Projection* Let $R : U\text{-Tuples} \mapsto K$ and $V \subseteq U$. Then $\pi_V(R) : V\text{-Tuples}$
 482 $\mapsto K$ is defined by $(\pi_V(R))(t) = \sum_{t=t'[V] \vee R(t') \neq 0} R(t')$.
- 483 4. *Union* Let $R_1, R_2 : U\text{-Tuples} \mapsto K$. Then $R_1 \cup R_2 : U\text{-Tuples} \mapsto K$ is
 484 defined by $(R_1 \cup R_2)(t) = R_1(t) + R_2(t)$.
- 485 5. *Natural join* Let $R_1 : U_1\text{-Tuples} \mapsto K$ and $R_2 : U_2\text{-Tuples} \mapsto K$. Then
 486 $R_1 \bowtie R_2 : U_1 \cup U_2\text{-Tuples} \mapsto K$ is defined by $(R_1 \bowtie R_2)(t) = R_1(t_1) \cdot$
 487 $R_2(t_2)$, where $t_1 = t[U_1]$ and $t_2 = t[U_2]$.

488 It is observed in [33] that if the K -relational semantics satisfies the same
 489 equivalence laws as positive relational algebra operators over bags, i.e. union
 490 ($+$) is associative, commutative and has identity \emptyset and join (\cdot) is associa-
 491 tive, commutative and distributive over union, and projection and selection
 492 commute with each other, as well as with union and join, then $(K, 0, 1, +, \cdot)$
 493 must be a commutative semiring.

494 Let us consider the algebraic structure $(\mathbb{N}(\text{TupleLoc}), 0, 1, +, \cdot)$, where
 495 $\mathbb{N}(\text{TupleLoc})$ is the set of polynomials whose coefficients are the natural
 496 numbers and the variable are from the set TupleLoc . The how-provenance
 497 of an output tuple is a function $\mathcal{H} = \text{How}(Q, I, o)$ that returns a polynomial

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.

498 in $\mathbb{N}(TupleLoc)$ called *provenance polynomial*. The following definition is
 499 adapted from [18] by considering the case applying to our work, i.e., $Q^K(I)$
 500 with $K = 1$.

Definition 4.2. *How-Provenance*

Let Q be an 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 $How(Q, I, t)$, is an element of the set $\mathbb{N}(TupleLoc)$ defined as follows:

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

501 Here $\{u\}$ is a query expression describing a constant, singleton relation,
 502 not a relation value per se. These constants correspond to K -relations that
 503 assign 1 to u and 0 to all other tuples. The summation in the projection
 504 case is finite since the support of a K -relation is assumed to be finite. In the
 505 selection rule, θ is seen as a function $\theta : U\text{-Tuples} \mapsto \{0, 1\}$.

506 *Example.* Let us consider the following SQL query Q1, applied to the database
 507 described in Table 1, asking for the names of families curated by researchers
 508 based in the United Kingdom (UK):

```

509     Q1: SELECT DISTINCT f.name
510         FROM family AS f JOIN contributor2family AS c2f
511         ON f.id = c2f.family_id

```

```

512 JOIN contributor AS c ON c2f.contributor_id = c.id
513 WHERE c.country = 'UK'

```

514 Table 3 shows the two output tuples of query Q1 annotated with their
515 respective how-provenances. Tuple o_2 was produced by a join of the input
516 tuples f_4 , $c2f_4$, and c_1 . The three provenance tokens are therefore “multi-
517 plied” together. The case of o_1 is slightly more complex. It can be obtained
518 by the joins of two different sets of tuples, so there are two monomials com-
519 bined by + representing these alternative derivations.

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

527 4.2. Causality and Responsibility

528 A formal study of causality was introduced in [19, 34] and later expanded
529 by Meliou et al. [51] to explain the causes of answers and non-answers to
530 queries. In the following, we refer to the definition of causality and respon-
531 sibility provided in [51]. In particular, we only focus on answers to a query
532 since non-answers are not relevant in our context.

533 There are two types of “cause” tuples: counterfactual and actual. Let o
534 be a tuple in the result of query Q on the database instance I , and t a tuple
535 in its lineage. We call t a *counterfactual cause* if, by removing t from I , o is
536 also removed from the output (i.e., t is essential for the generation of o).

537 We call t an *actual cause* if there is a set of tuples $\Gamma \subseteq I$ called a *contingency set*,
538 such that t is a counterfactual cause in $I - \Gamma$. In other words, t is
539 an actual cause if, even when removed from I , there is another set of tuples
540 of the lineage that guarantees the presence of o .

541 Computing the causality of tuples is NP-complete for general queries [29],
542 but for conjunctive queries it can be computed in PTIME, as showed by
543 Meliou et al. [51].

544 The notion of *responsibility* measures the degree of causality as a function
545 of the size of the smallest contingency set [19]. This allows us to rank lineage
546 tuples based on their degree of causality in generating the output.

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.

Definition 4.3. *Responsibility [51]*

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:

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

547 where Γ ranges over all contingency sets for t .

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

552 *Example.* Let us consider Table 4, where we reported the result set of Q1
553 and the tuples of the lineages with their responsibility values. Focusing on
554 o_1 : the lineage tuple f_1 is a counterfactual cause, since its contingency set is
555 empty (when removed from the database, o_1 disappears from the result set).
556 Consequently, its responsibility is 1. All the other tuples of the lineage are
557 actual causes. c_1 , for example, has as minimal contingency set $\{c_2f_2\}$, thus
558 its responsibility is 0.5. For the output tuple o_2 , all the tuples of the lineage
559 are counterfactual causes, thus their responsibility is 1.

560 *4.3. Shapley value*

561 We rely on the definition provided in [26]. Let Q be a Boolean query
562 and $f \in D$ be a fact, the Shapley value of f in D intuitively represents the
563 contribution of f to the query result.¹² The higher the value, the more f
564 helps in satisfying Q . Formally, the Shapley value is defined as follows:

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

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.

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

565 The sum is performed on all possible subsets of D that do not contain f . The
566 value $(Q(E \cup \{f\}) - Q(E))$ is the “wealth” brought by f when added to E .
567 Thus, the Boolean query is used as a wealth function v : its value is 1 only
568 when the set $E \cup \{f\}$ makes the query true, and the set E makes it false,
569 i.e., when the addition of the fact f is determinant to making the Boolean
570 query true. The value $|E|! (|D| - |E| - 1)!$ is the number of all the possible
571 permutations over D where the facts in E come first, then f is added, and
572 then all the remaining facts. Thus, the value $\frac{|E|! (|D| - |E| - 1)!}{|D|!}$ can be thought
573 as a weight for the wealth brought by f when added to E .

574 To extend this definition to non-Boolean queries, we adopt the approach
575 in Deutch et al. [26]: the Shapley value of the fact f for the answer \bar{t} to
576 $Q(\bar{x})$ is the value $Shapley(Q[\bar{x}/\bar{t}], D, f)$, where $Q[\bar{x}/\bar{t}]$ is the Boolean query
577 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
578 0 otherwise. In other words, the definition of $Shapley(q, D, f)$ is extended
579 to queries $Q(\bar{x})$ with free variables by considering the Boolean query $Q[\bar{x}/\bar{t}]$
580 as a value function. This query can be seen as a function that takes as input
581 a set of facts and returns 1 if this set is a witness for \bar{t} , and 0 otherwise.

582 *Example.* Let us consider Table 5, that shows the Shapley values for the
583 lineage’s tuples of o_1 and o_2 , results of query Q1. We note that, to compute
584 the Shapley value of an input tuple f it is sufficient to compute and sum the
585 values $\frac{|E|! (|D| - |E| - 1)!}{|D|!}$ for all the possible sets E such that $E \cup \{f\}$ is a witness
586 and E is not. Thus, suppose we want to compute the Shapley value of the
587 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
588 only if o_1 is in the output of Q1 on D , and L_{o_1} is the lineage of o_1 . Then the
589 Shapley value of f_1 with respect of o_1 is given by:

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

590 where for the first element of the sum the corresponding E is $\{c2f_1, c_1\}$, for
 591 the second element it is $\{c2f_2, c_2\}$, for the third $\{c2f_1, c2f_2, c_1\}$, for the fourth
 592 $\{c2f_1, c_1, c_2\}$, for the fifth $\{c2f_2, c_2, c_1\}$, for the sixth $\{c2f_1, c2f_2, c_2\}$, and for
 593 the seventh $\{c2f_1, c2f_2, c_1, c_2\}$. Every other possible subset E would make
 594 the factor equal to 0. Note that in this case we consider $D = L_{o_1}$, the lineage
 595 of o_1 , since these are the only facts in all the database that contribute to the
 596 generation of o_1 .

Similarly, for tuple c_1 (and the other tuples of the lineage), the computa-
 tion is:

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

597 We can see that for all the tuples of o_2 's lineage the corresponding Shapley
 598 values are equal to $1/3$, since they are all equally responsible for the gener-
 599 ation of the output. Thus the sum of the Shapley values of all the tuples in
 600 an output tuple's lineage is always equal to 1 when using a Boolean query
 601 as wealth function.

602 5. Credit Distribution and Distribution Strategies

603 We now give formal definitions of data credit and Data Credit Distri-
 604 bution (DCD), and present the three different Distribution Strategies (DSs)
 605 base on how-provenance, responsibility, and Shapley value. We also show
 606 how these strategies distribute credit in the IUPHAR example presented
 607 above.

608 5.1. Data Credit and Data Credit Distribution

609 Given a database instance I , a *recipient of credit* is a unit of information
 610 within I ; in this work, we focus on tuples as recipients. *Data credit* is a value
 611 $k \in \mathbb{R}_{>0}$. Every recipient in a database is annotated with a quantity of credit
 612 as a proxy for its importance.

613 Given a DS, DCD takes a database instance I , a quantity of credit k ,
 614 query $Q(I)$, and it divides k among the tuples in I .

$$\begin{aligned}
\mathcal{H} &= \underbrace{3f_1 \cdot c_2 f_1 \cdot c_1}_{M_1} + \underbrace{f_1 \cdot c_2 f_2^3 \cdot c_2^3}_{M_2} \\
c(\mathcal{H}) &= 4 & e(M_2) &= 7 \\
mc(M_1) &= 3 & mc(M_2) &= 1 \\
te(c_2, M_2) &= 3 & \gamma(c_1, \mathcal{H}) &= \{M_1\} \\
\gamma(f_1, \mathcal{H}) &= \{M_1, M_2\}
\end{aligned}$$

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

615 **Definition 5.1.** *Tuple Level Data Credit Distribution (DCD) [27]*
616 *Given a query Q over I and $k \in \mathbb{R}_{>0}$, the tuple level DCD is defined by*
617 *the function $f_{I,Q} : \text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ such that $f_{I,Q}(t, k) = h$ where*
618 *$0 \leq h \leq k$ and $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$. The function $f_{I,Q}$ is the distribution*
619 *strategy (DS).*

620 As we can see, the DS is a function that annotates each tuple in the
621 database with a real value, which is a fraction of the given quantity k . The
622 only constraint is that the sum of the credit annotations on tuples is k .

623 In the following, we use information provided by data provenance to de-
624 fine distribution functions. For simplicity, we assume that the credit k is
625 distributed equally across the set of output tuples, and discuss how the credit
626 k_o of one output tuple o , is distributed across the instance I .

627 5.2. A How-Provenance Based Distribution Strategy

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

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

635 We call c the function that, given a polynomial, returns the sum of its
636 coefficients; e.g., $c(\mathcal{H}) = 3 + 1 = 4$. We call e the function that, given a
637 monomial, returns the sum of its exponents, e.g., $e(M_2) = 1 + 3 + 3 = 7$.
638 mc is the function that takes a monomial as input and returns its coef-
639 ficient; e.g., $mc(M_1) = 3$. te is a function that takes as input a tuple
640 and a monomial, and returns the exponent of the tuple in the monomial,
641 if present; e.g., $te(c_2, M_2) = 3$. Finally, γ takes as input a tuple and the

id	name	how-provenance
<i>oxs₁</i>	Dopamine Receptors	$f_1^2 c_2 f_1 c_1 + f_1^2 c_2 f_2 c_2$

Figure 5: 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.

642 whole polynomial, and returns a set of monomials containing that tuple;
643 e.g., $\gamma(f_1, \mathcal{H}) = \{M_1, M_2\}$, $\gamma(c_2, \mathcal{H}) = \{M_2\}$.

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:

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

644 Going back to the example of Table 3, consider o_1 with provenance poly-
645 nomial $f_1 c_2 f_1 c_1 + f_1 c_2 f_2 c_2$. The how-provenance-based DS firstly divides
646 the credit between the two monomials. Since the coefficients of each mono-
647 mial are 1, the credit is split in half. If they were, for example, 1 and 2
648 respectively, 1/3 of the credit would go to the first monomial, and 2/3
649 to the second. Since in our example each variable has exponent 1, the credit is
650 further divided equally among the three variables. Thus, at the end of the
651 computation, f_1 receives 1/3, and the other tuples receive 1/6.

652 As a further example, let us consider a query Q2 over GtoPdb, asking
653 for the families of type `gpcr` that have researchers located in the UK as
654 contributors.

```

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

```

662 The result of Q2 is shown in Figure 5, and consists of one tuple, *oxs₁*,
663 annotated with its how-provenance. As we can see, the how-provenance
664 shows that f_1 is used twice: first in the join of the inner query, and second
665 in the join of the outer query.

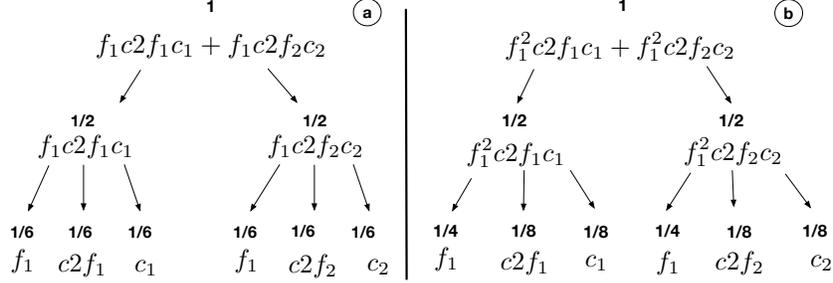


Figure 6: Comparison of different distributions obtained with the how-provenance-based DS with queries Q1 and Q2.

666 Figure 6 shows how the DS based on how-provenance behaves on the
 667 polynomial from query Q1 (Figure 6.a) and that from query Q2 (Figure 6.b).

668 In Figure 6.a, tuple f_1 receives credit $1/3$ and the other tuples receive
 669 $1/6$, while in Figure 6.b tuple f_1 receives credit $1/2$ and the others receive
 670 $1/8$. This is reasonable since Q2 relies on f_1 more than Q1, and it shows how
 671 how-provenance is sensitive to the tuples' role in a query.

672

673 5.3. Responsibility-based Distribution Strategy

674 As described in Section 4.2, causality and responsibility is new information
 675 that is added to lineage. One option for a responsibility-based DS is to
 676 assign the responsibility of each tuple in the lineage of an output tuple as
 677 its credit. In this way, responsibility is both a way to compute credit and to
 678 distribute it. Referring to the example of Table 4, in the case of output tuple
 679 o_1 , f_1 receives credit 1 and the other tuples receive credit 0.5.

680 However, we want a DS that is also a function of the input credit value
 681 k . So, we define a new DS that is a function of the quantity of credit k_o that
 682 assigns to each tuple of the lineage a portion of this credit weighted by its
 683 normalized quantity of responsibility. This function gives a bigger portion of
 684 credit to tuples that are higher in the responsibility ranking.

Definition 5.3. *Responsibility-based Distribution Strategy*

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

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

	$\underbrace{f_1}_{\text{counterfactual cause}}$	$\underbrace{c_2 f_1 \quad c_2 f_2 \quad c_1 \quad c_2}_{\text{actual causes}}$			
$k_{o_1} = 1$	f_1	$c_2 f_1$	$c_2 f_2$	c_1	c_2
	↓	↓	↓	↓	↓
responsibility	1	0.5	0.5	0.5	0.5
responsibility-based DS	1/3	1/6	1/6	1/6	1/6

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

685 Figure 7 shows the responsibility and credit assigned to the tuples of the
686 lineage of the output tuple o_1 of Table 4. Assuming that $k_{o_1} = 1$, f_1 receives
687 credit $1/3$, while the others receive credit $1/6$.

688

689 5.4. Shapley value-based Distribution Strategy

690 As with responsibility, the Shapley value can be seen both as a method
691 to generate and distribute credit. Moreover, it can be seen that, using the
692 definition of Shapley value for Boolean queries given in Section 4.3, the sum
693 of the Shapley values of all the tuples of the lineage L of an output tuple o
694 is 1.

Definition 5.4. Shapley Value-Based Distribution Strategy

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

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

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

697 As shown in Table 5, tuple f_1 in o_1 's lineage takes credit $7/15$ when
698 $k_{o_1} = 1$, while the other tuples of the lineage take credit $2/15$. This DS still
699 rewards f_1 more than the other tuples, since it is more important than the
700 other tuples of the lineage. However, this DS behaves differently from the
701 other two previous strategies. In particular, f_1 is rewarded more with this
702 DS than with the others.

703 In the case of o_2 there is only one monomial in the provenance polyno-
704 mial and all the three tuples appearing in it are counterfactual causes. The
705 consequence, in this type of cases, is that the three distributions behave in
706 the same way. Here, all three tuples of o_2 's lineage receive credit $1/3$.

707 6. Experimental Evaluation

708 To understand the trade-offs between these Distribution Strategies (DSs),
709 we perform five sets of experiments using queries over target families pre-
710 sented on the GtoPdb website. The first set of experiments uses real queries
711 extracted from citations to GtoPdb published in the British Journal of Phar-
712 macology. The second set uses synthetically produced provenance polyno-
713 mials, corresponding to more complex queries, in order to better highlight
714 the differences between the DSs. The third set of experiments considers
715 the accrual of credit over time by the three strategies, again using synthetic
716 queries. The fourth set of experiments shows how the DSs compare to tradi-
717 tional citations in giving credit to data curators using both real and synthetic
718 queries. In the last set of experiments we report the execution time required
719 to compute how-provenance, responsibility and Shapley values of the output
720 tuples.

721 The source code for the experiments is written in Java and supported by
722 a PostgreSQL database. For purposes of reproducibility, the source code
723 and all queries are available at [https://bitbucket.org/dennis_dosso/
724 credit_distribution_project](https://bitbucket.org/dennis_dosso/credit_distribution_project).

725 6.1. Real-world queries

726 Examples of real queries are drawn from papers published in the British
727 Journal of Pharmacology (BJP)¹³. Each time a paper in this journal cites a
728 webpage from GtoPdb, it reports the URL of the page. From this URL, the
729 query used to obtain the webpage data can be determined. We considered all
730 889 papers in BJCP citing the IUPHAR/BPS Guide to pharmacology [35]
731 as of October 2020, and extracted all webpage URLs to GtoPdb contained
732 within the paper.¹⁴

733 The queries that we inferred are those used to build target family web-
734 pages within GtoPdb. An example was given in Figure 3, where we show
735 how the structure of the “Adenosine receptors” family can be mapped into
736 queries over the underlying database. In GtoPdb, all target family pages
737 share a similar structure; the only difference is that individual sections, such

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.

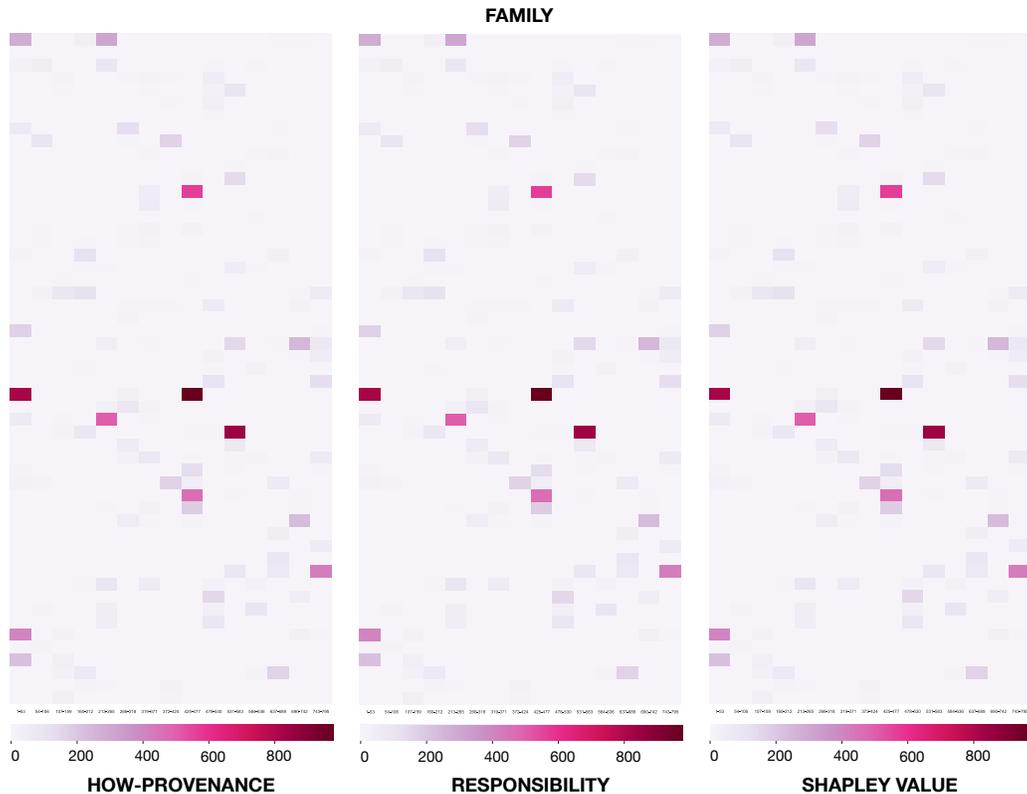


Figure 8: Comparison of four DS on the same table `family` using the distribution given by the queries retrieved from papers. Each cell is a tuple.

738 as “contributors” or “further readings”, may be missing. Therefore, the same
 739 queries can be used to build all of the target family pages by changing the
 740 family id used in the query (for example, in Figure 3, it is 3). Note that
 741 the queries are fairly simple SPJ SQL queries. A total of more than 12K
 742 different queries were built in this way. Without loss of generality, we give
 743 each tuple in the output of a query a credit of 1.

744 *Results.* Figure 8 shows the heat-maps obtained by the distribution of credit
 745 according to the three DS on one of the tables in the underlying database,
 746 `family`, which is often joined with other tables in the database to build the
 747 webpages. Each cell in a heat-map represents a tuple of the `family` table
 748 and the color indicates the amount of credit attributed to such tuple. It can
 749 be seen that the result of credit distribution over `family` is the same for all

750 three strategies. The same result is also obtained with the other tables of
751 the database used by the queries shown in Figure 3.

752 The reason why credit distribution is the same for all strategies is that the
753 queries are all simple SPJ queries, which use one tuple from each table only
754 once and do joins on key attributes (these are always 1-to-1 joins). Under
755 these conditions, each output tuple presents: (i) a how-provenance that is
756 a single monomial with coefficient one and exponent one in each variable;
757 (ii) all tuples are counterfactual causes when considering responsibility, thus
758 they have responsibility 1, and (iii) all tuples have the same importance in the
759 production of the output tuples according to their Shapley value. Hence, for
760 these queries, the DSs behave in the same way: credit is uniformly distributed
761 among the tuples of the lineage.

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

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

768 Q3 returned 10 tuples from the version of GtoPdb used. The first tuple,
769 <Bertil B., Fredholm>, has $c_{939} \cdot c2f_{496}$ as its provenance polynomial. c_{939}
770 represents the provenance token of a tuple in `contributor`, and $c2f_{496}$ the
771 provenance token of a tuple in table `contributor2family`. Also, both these
772 tuples are counterfactual causes and have a responsibility of one. Therefore,
773 the credit assigned to these tuples is 1/2 using all five DS. This happens
774 for all the tuples in the output of each query of GtoPdb, thus making the
775 distributions equivalent over all outputs.

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

779 6.2. Synthetic queries

780 To see what happens with more complex queries, we synthetically gener-
781 ated provenance polynomials in which the coefficients and exponents could
782 be greater than one, and picked them at random from a uniform distribution.
783 The queries involve three GtoPdb tables: `family`, `contributor2family`, and
784 `contributor`. The polynomials were generated as follows: first, the number
785 of monomials was decided by randomly choosing a number 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:

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

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 9: Sample synthetic provenance polynomial (how-provenance) and corresponding 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).

786 six. Then, we randomly chose a tuple from the `family` table, one from the
787 `contributor2family` table and one from the `contributor` table; these are
788 the variables of the monomial. We then chose a coefficient for the monomial
789 (between one and three) and an exponent for each tuple (between one and
790 four). For the next monomial, we decided if we wanted to keep the same
791 tuple from the table `family` as first tuple of the new monomial. To do so,
792 we generated a random float number between zero and one. If the number
793 was above 0.2, we changed the family tuple. This number was chosen ar-
794 bitrarily to obtain polynomials that presented a certain “variation” in their
795 monomials, i.e., to make sure that not all monomials started with the same
796 tuple.

797 An example can be seen in Figure 9, which shows a sample synthetic
798 provenance polynomial (the how-provenance), the causality of the tuples of
799 the lineage, together with their responsibility, and, finally, the Shapley values
800 of the lineage tuples. The resulting credit distribution for each DS is also
801 shown.

802 As an example of how the distribution strategies behave with these syn-
803 thetic queries, consider tuple f_5 in Figure 9. This tuple receives the highest

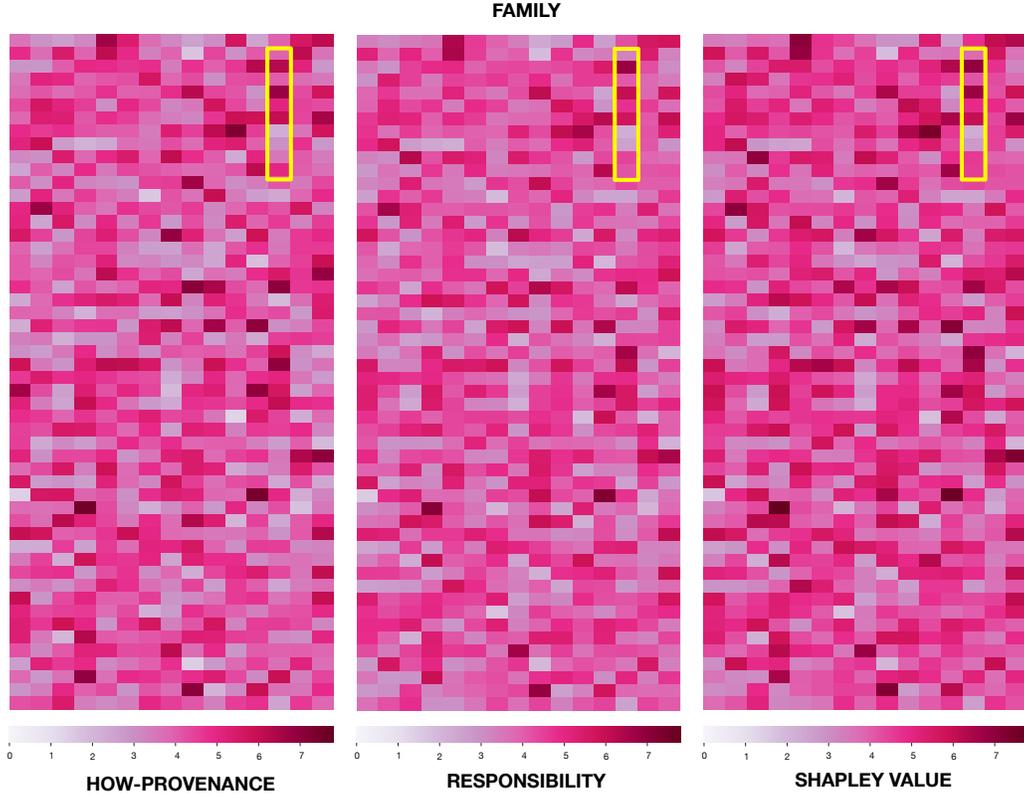


Figure 10: Comparison of three DS on the same table `family` after the distribution computed using 10K synthetic and randomly generated provenance polynomials. The tuples in the blue rectangles are used as example in the discussion connected to Figure 11.

804 quantity of credit using responsibility-based distribution and less credit using,
 805 in order, the Shapley value and how-provenance. On the other hand, tuple
 806 f_1 is rewarded more by the Shapley value, then, in order, how-provenance
 807 and responsibility. This difference is explained considering the different role
 808 of the tuples in the generation of the output and the characteristics of the
 809 distributions.

810 Responsibility creates a ranking among lineage’s tuples describing the im-
 811 portance of their role in generating the output. As such, the responsibility-
 812 based DS gives more credit to f_1 , f_5 , c_2f_{17} and c_{18} due to their higher respon-
 813 sibility values. “Importance” is connected to their corresponding minimal
 814 contingency sets. For example, f_1 has a minimal contingency set (one of
 815 the many) $\{f_5\}$, with cardinality 1. On the other hand, c_1 has, as minimal

Table 6: Results of the pairwise Kendall Tau confidence value on all the DSs on the `family` table (the p-values are all below 0.05).

	how	resp.	Shapley
how	1.0	0.74	0.74
resp.	0.74	1.0	0.89
Shapley	0.74	0.89	1.0

816 contingency set (one of the many) $\{f_5, c_2\}$, with cardinality two. This means
 817 that c_1 is the “least important” amongst the tuples with minimal contin-
 818 gency sets of lower cardinality, and this is reflected in the different quantities
 819 of credit being distributed.

820 The Shapley value behaves similarly, but it rewards tuple f_1 the most
 821 and then f_5, c_2, f_{17}, c_{18} , and last all the other tuples of the lineage. Although
 822 both Responsibility and the Shapley value create a ranking of the tuples
 823 based on their role in the generation of the output, the corresponding func-
 824 tions behave differently due to the syntax of the query. For this reason each
 825 different distribution strategy highlights a slightly different aspect that can
 826 be considered as “important” when distributing the credit.

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

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

836 There is consistency between the strategies in that tuples which are highly
 837 rewarded by one strategy are also highly rewarded by the others. This shows
 838 that the four DSs consistently reward certain tuples more than others.

839 Table 6 reports the pairwise Kendall τ correlation values¹⁵ for the three
 840 DSs computed on the `family` table. As we see, the distribution based on

¹⁵The Kendall’s τ coefficient is a statistic used to measure the ordinal association between two measured quantities [43]. Intuitively, it is high between two variables when observation have a similar rank.

841 how-provenance is the one that correlates less with the other two strategies,
842 while it seems that the DSs based on responsibility and the Shapley value are
843 more correlated one with the other. This may be explained because, while
844 how-provenance captures how the tuples are used, the other two strategies
845 are concerned with the importance of the tuples in the lineage of the query
846 (responsibility) and the role that the tuples have in the query seen as a
847 coalition game (Shapley value). Hence, the three DSs represent different
848 viewpoints about the “importance” of a tuple, and this reflects on their
849 distributions. Moreover, we have to consider that how-provenance is a
850 *provenance*, and our approach uses its information to obtain a metric, while
851 Responsibility and Shapley value are metrics. The main difference between
852 the three resides in the definition of the metric itself. The definition of
853 Shapley value resides on the concept of coalition and in the different possible
854 combinations in which a coalition is built. Responsibility, on the other hand,
855 is based on the concept of minimal contingency. The metric that we derived
856 in this paper from how-provenance, instead, exploits the information in the
857 polynomial to obtain a value metric that is not based on the concept of a set
858 (respectively, coalition and contingency). This may be a further explanation
859 of why how-provenance correlates the least with the other twos.

860 Considering the three heat-maps reported in Figure 10, it is evident that
861 there are many similarities. However, upon closer inspection, it is possible
862 to see that they are behaving differently, with certain tuples rewarded more
863 with one strategy than with the others.

864 The heat-map reporting the distribution produced by the Shapley value is
865 the one that, at a closer inspection, shows more evident differences. Although
866 the tuples that receive the biggest quantities of credit are the same, the hue
867 of these tuple is different.

868 We note that the how-provenance-based DS gives an average credit of
869 4.18 to each tuple in the table, while the responsibility-based 4.13, and the
870 Shapley-based 4.40. Moreover, how-provenance distributed a total of about
871 3331 units of credit to the **family** table, while responsibility assigned 3290,
872 and the Shapley value 3505 (the difference of credit is due to the fact that,
873 depending on the DS, other tables used in the joins are rewarded more).

874 To better understand the differences between DSs, in the next subsection
875 we consider the accrual of credit over time. In doing so, we will focus on the
876 ten tuples shown within the large yellow rectangles in Figure 11. Each small
877 rectangle within a large yellow rectangle is a tuple, and we number them
878 from 1 (top) to 10 (bottom). These ten tuples were cherry-picked because

879 they allow us to see the evolution of the distribution of credit through time.
 880 There are other tuple sets that could have been selected driving us to the
 881 same considerations.

882 6.3. Credit accrual over time

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

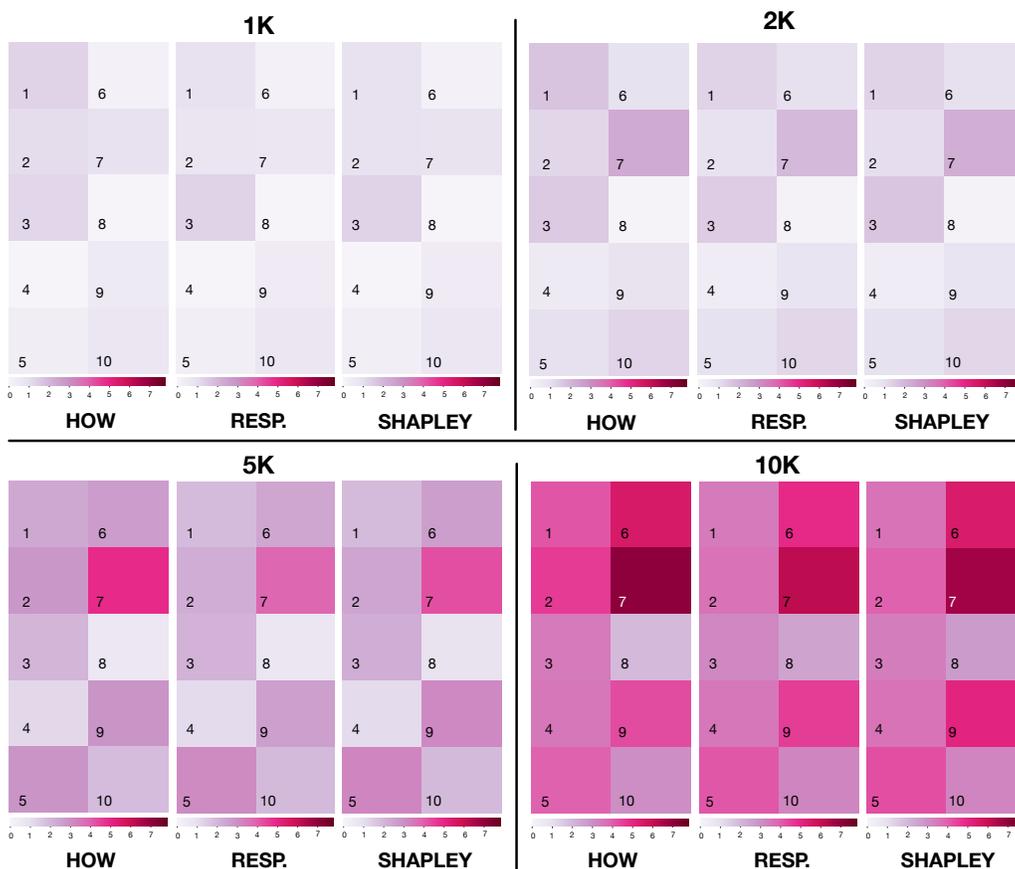


Figure 11: Comparison of the distribution of credit performed by the five DSs on a subset of 10 tuples taken from the `family` table, simulating the passing of time. The number at the top of each group of heat-maps represents the number of polynomials whose credit has been distributed.

886 In this figure, four groups of heat-maps are shown. Each group represents
 887 a “snapshot” taken after 1K, 2K, 5K and 10K provenance polynomials have

888 been considered for credit distribution. The ten tuples in each heat-map are
889 those highlighted in the yellow boxes of Figure 10 from the `family` table.

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

895 By the end of 1K queries, credit differentials between tuples as well as
896 between strategies can be seen. For example, tuple 3 is usually rewarded the
897 most credit by all three strategies. Moreover, it can be seen that tuple 1
898 receives a higher quantity of credit when how-provenance is adopted, show-
899 ing how this form of provenance behaves differently from the others in this
900 context. Moving to 2K queries, it is possible to see that tuples 3 and 7 are
901 still the most rewarded by the strategies.

902 By the end of 5K queries, tuple 7 emerges with the highest value of
903 credit with all three DSs, a position which is strengthened with 10K queries.
904 Moreover, with the passing of time, tuple 3 ceases to be amongst the most
905 rewarded ones and new tuples, such as 6 and 9, emerge as being particularly
906 rewarded at 5K, while at 10K tuples 6 and 7 are the most rewarded. The DS
907 that rewards the more tuple 7 is the one based on how-provenance (credit
908 7.03), followed by the Shapley value (credit 6.64). This is due to the fact
909 that tuple 7 had, among some of the polynomials being used for the exper-
910 iments, a high responsibility but it did not appear in all the monomials of
911 the provenance polynomials. This changed slightly the distribution.

912 To sum up, the DS based on how-provenance highlights which tuples in
913 the database are used by a query. It distributes credit to the tuples based on
914 their role in the queries. In particular, tuples that were used more frequently
915 and in many different ways receive more credit. The distributions based on
916 responsibility and the Shapley value are more concerned with the importance
917 of individual tuples in generating the output. Responsibility, in particular, is
918 concerned in the role of the tuple as an actual or counterfactual cause, and
919 will reward tuples that are more “fundamental” for the output. On the other
920 hand, the Shapley value sees tuples as players in a coalition game where all
921 the tuples of the lineage “work” toward the production of the output. The
922 tuples whose role is more important in the game defined by the query are
923 rewarded with higher quantities of credit.

924 These three DSs may be useful for finding “hotspots” in the database
925 based on the role of tuples. The preference of one over the others depends

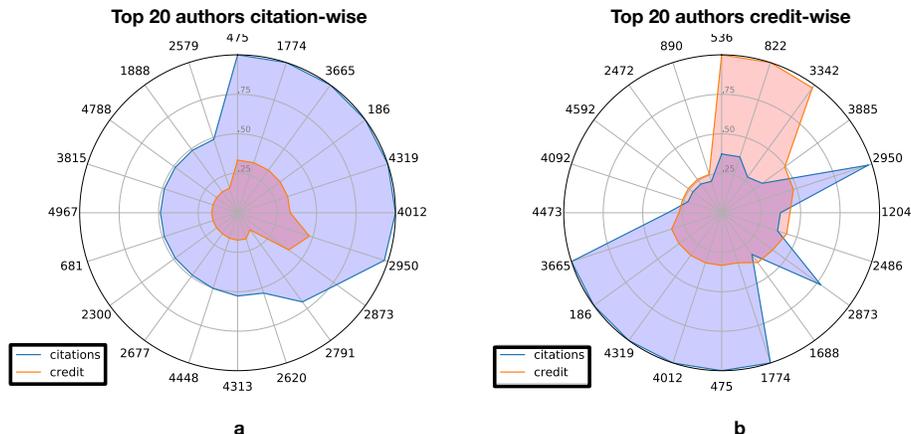


Figure 12: Radars presenting the top 20 authors citation-wise and credit wise, together with their (normalized between 0 and 1) values of citations and credit.

926 on the type of sensitivity to the role of a tuple in queries that is required by
 927 the context as dictated by the preferences of the users or the peculiarities of
 928 the application at hand.

929 6.4. Credit vs Citations

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

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

944 *Results: Real-world queries.* As described in Section 6.1, we consider real-
 945 world queries taken from papers published in the BJP which reference web-

946 pages in GtoPdb. Since for these queries there is no difference in the distri-
947 bution of credit between the DSs, only one value for credit is used.

948 The results are shown in the radar plots of Figure 12, in which each
949 number on the outer circle (e.g. 475, 1774 and 3665) represents an author
950 (id) and the blue (red) line represents the normalized value of credit generated
951 by citations (credit), respectively. The first radar plot, Figure 12.a, shows the
952 top 20 authors in terms of *citations*, ordered in a clockwise direction, whereas
953 Figure 12.b orders the authors based on *credit*. Comparing the author ids
954 used in the outer circles of these two plots, it can immediately be seen that
955 the “top authors” are very different using these two metrics, although there
956 is some overlap (for example, authors 1774, 475, and 4012).

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

968 The reason for the difference between citations and credit is partly due
969 to the experimental setup: each output tuple carries a credit of 1, and there
970 can be many tuples used to generate a webpage. Thus a webpage that is
971 created from more tuples will have a higher credit value than one created
972 from fewer tuples. Furthermore, authors who collaborated with fewer people
973 will receive a biggest share of the equally divided credit. However, all authors
974 will receive a citation of one.

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

979 *Results: Synthetic queries.* We used the same synthetic polynomials de-
980 scribed in Section 6.2, and we distributed credit with the first 100, 1K, and
981 10K of them. Since these polynomials are created by randomly selecting
982 tuples from three tables, they usually correspond to a set of data curated by

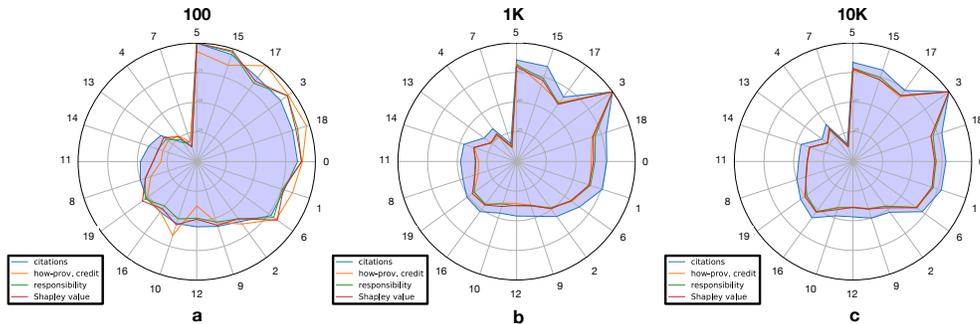


Figure 13: Radars presenting the 20 synthetic authors with corresponding citation and quantities of credit distributed through the 3 DSs (all values normalized between 0 and 1) through different numbers of polynomials (respectively, 100, 1K and 10K). The order is the one defined by figure a, i.e. descending order of citations obtained from 100 polynomials.

983 authors who, in reality, did not collaborate. To make the size of the author
 984 set more realistic, we therefore created 20 synthetic authors, and randomly
 985 assigned one author to blocks of consecutive tuples in the database, with the
 986 size of each block varying between 10 and 40, to simulate different quantities
 987 of work performed by an author. Every time an author appears as curator of
 988 one or more tuples used in a polynomial, we assign them one citation. They
 989 also receive three kinds of credit, each one using a different DS.

990 Figure 13 shows three radar plots, one for the results obtained with 100
 991 polynomials, one with 1K polynomials, one with 10K polynomials. Each
 992 plot shows the top 20 authors in terms of citations (hence the authors and
 993 clockwise ordering is the same in each of the plots), and additionally shows
 994 the normalized values of citation (blue line), how-provenance-based credit
 995 (yellow line), responsibility-based credit (green line), and the Shapley value-
 996 based credit (red line).

997 As can be seen, given the synthetic nature of the queries, the correlation
 998 between the number of citations and the quantity of credit assigned to the
 999 authors appears to be a much stronger than with the real-world queries of
 1000 Figure 12. In fact, for Figure 13.a the linear correlation between the citation
 1001 number and all three types of credit is always above 0.94 with p -values in
 1002 the order of $3e - 8$.

1003 What these figures show is that, in certain cases, authors who do not have
 1004 a large number of citations receive more credit than others, as for example au-
 1005 thors 17, 18 and 10 in Figure 13.a, and especially when credit is distributed

Table 7: Average execution time (ms) to compute how-provenance, responsibility and Shapley values of one output tuple. The accompanying z-values were computed with confidence of 95% and $\alpha = 5$.

	how-provenance	responsibility	Shapley
real queries	57.29 ± 0.25	58.16 ± 0.02	85.18 ± 0.24
synthetic queries	-	1.48 ± 0.05	39.79 ± 2.87

1006 using how-provenance. This again shows how credit gives a different per-
 1007 spective on the role of data and authors by going beyond the limitations of
 1008 traditional citations.

1009 It is worth noting that, when scaling up to $1K$ and $10K$ polynomials, the
 1010 credit distributions become almost identical (the linear correlation for the
 1011 values of Figure 13.c is more than 0.99 with a p -value of $1.32e - 32$). This is
 1012 consistent with what we observed in Figure 10.

1013 6.5. Execution Time

1014 We studied the time required to compute the how-provenance, responsi-
 1015 bility and Shapley value of the output tuples used in the previous exper-
 1016 iments on GtoPdb, for both real and synthetic queries. All experiments were carried
 1017 out on a MacBook Pro with a 2.4 GHz processor Intel Core i5 quad-core and
 1018 8 GB of memory at 2133 MHz.

1019 Recall that we first compute the how-provenance of real queries, obtaining
 1020 a total of 58,037 polynomials. For synthetic queries, we directly produce the
 1021 polynomials so it was not necessary to compute the how-provenance, whereas
 1022 responsibility and Shapley values of the output tuples were computed starting
 1023 from these polynomials.

1024 Table 7 reports the average time required to compute how-provenance,
 1025 responsibility, and Shapley values of one output tuple, both in the case of
 1026 real and synthetic queries (here, we consider all 10,000 produced synthetic
 1027 polynomials when computing the average). All times are reported in millisec-
 1028 onds. The time reported in the table to compute how-provenance is obtained
 1029 using the code provided in [66], while the responsibility and Shapley value
 1030 times are the result of the sum of this time with the time required to compute
 1031 them starting from how-provenance.

1032 From this table, we can see that the overhead required to compute respon-
 1033 sibility is small, while the overhead for the Shapley value is larger, primarily
 1034 due to the need to compute the power set of the lineage. We also note that

1035 the execution times for a single tuple are relatively small, but become size-
1036 able when the queries present a large result set and, in particular, for tuples
1037 with big lineage sets.

1038 What we can see from these results is that how-provenance is efficient and
1039 gives an informative distribution of credit for SPJ queries. Responsibility is
1040 still efficient, and gives a slightly different perspective on credit distribution.
1041 The Shapley value adds significant computational overhead, but is still fea-
1042 sible for small/medium databases and SPJ queries. Moreover, recent work is
1043 investigating new efficient and approximated ways to compute the Shapley
1044 value.

1045 In the following Sections we provide a bigger picture of computing how-
1046 provenance, responsibility and Shapley value for queries beyond SPJ, based
1047 on the latest findings in the literature.

1048 7. Discussion

1049 Before concluding, we discuss some design decisions: the focus on Credit
1050 Distribution (as opposed to Credit Generation), the choice of Distribution
1051 Strategies and, finally, how the concept of Game Provenance can open up
1052 new possibilities for Credit Distribution in new contexts and for new classes
1053 of queries.

1054 7.1. Credit Generation

1055 Credit Generation is the task of generating the credit to be distributed
1056 by a DS. Credit Generation presents a series of issues shared by traditional
1057 citation practices. For instance, defining the quantity of credit generated
1058 for a given citation is still an open problem. Different types of citations may
1059 generate different amounts of credit. Data cited as previous work or as useful
1060 for previous work may generate less credit than other data extensively used
1061 to produce the results presented in a paper. The computation of credit could
1062 be done manually (although we must consider the complexity of the task,
1063 human biases, and the resources required to carry it out) or automatically,
1064 but it must be based on a shared definition of impact, which is still not
1065 agreed upon for data or traditional citation. For this reason, we used a
1066 uniform credit assignment function.

1067 There is also the problem of *transitive credit distribution*, i.e., how to
1068 transitively propagate credit from one cited unit to another unit that was
1069 used to produce the one being cited. For this, a graph of cited units that

1070 propagate credit between the units depending on influence could be used.
1071 How to propagate credit is an open and non-trivial problem that needs to
1072 consider the importance and impact of a citation in a work, be it a paper or
1073 data, and how to eventually compute the quantity of credit to be propagated.

1074 Finally, in our experiments we assumed that the credit carried by an
1075 output tuple is one. Thus, each tuple in the output has equal importance.
1076 As described above, this assumption may be revised and different credit to
1077 different output tuples could be assigned. Note that from the distribution
1078 model viewpoint no change is required since the DCD is defined for a generic
1079 value k .

1080 7.2. Choice of Distribution Strategies

1081 In this paper we presented three different DSs, so the natural question
1082 is which one to use. This depends on the task at hand. When we want to
1083 highlight the tuples being used in the database by a workload, the lineage-
1084 based DS proposed in [27] may be sufficient. When we also want to know the
1085 relative impact of tuples in the context of the query, the other DSs should
1086 be used since they give a better understanding of the importance of data.

1087 In the real-world-based experiments presented in the paper, the three DSs
1088 behaved the same, which was due to the specific nature of the data and the
1089 queries being used. However, the how-provenance of a query will differ from
1090 the lineage of the same query whenever the output tuples can be computed
1091 in more than one way by the query. This is usually true when join and
1092 projection operators are used in the query. This means that how-provenance
1093 DS may be preferred to the simple lineage-based one when more complex
1094 provenance polynomials may be expected.

1095 To address the question of what types of queries are likely to extract
1096 cited data, we turn to the results of published studies on the characteristics
1097 of query workloads and the complexity of their queries [39, 58, 63]. These
1098 studies show that operations such as inner-/outer-joins and projections occur
1099 in many queries. Therefore how-provenances may become quite complex
1100 in some instances and provide a distribution of credit that is significantly
1101 different from the one obtained with lineage.

1102 From the perspective of computational complexity, all three DSs are sim-
1103 ilar since we focused on SPJ queries, although there is a slightly larger over-
1104 head with the Shapley value (see Section 6.5). However, the tests were con-
1105 ducted on a relatively small database using a rather naïve algorithm to com-
1106 pute responsibility and the Shapley values. Hence, on a big database, the

1107 Shapley value might become prohibitively expensive to use. On the other
1108 hand, faster algorithms to calculate the Shapley value are being investigated
1109 and might speed up the process at least for a specific class of queries (e.g.,
1110 SPJ) [26].

1111 Going beyond SPJ queries, Green et al. [33] proposed the provenance
1112 semiring framework for SPJRU (Select, Project, Join, Rename, and Union
1113 queries), and Amsterdamer et al. [5] showed how to extend the framework to
1114 aggregate queries. This makes the DS based on how-provenance also suited
1115 for these important types of queries.

1116 Responsibility is harder to compute for *general queries* (NP-complete).
1117 Meliou et al. [51] proved a dichotomy result for conjunctive queries: for each
1118 query without self-joins, either its responsibility can be computed in PTIME
1119 in the size of the database, or checking if it has a responsibility below a
1120 given value is NP-hard. Queries with self-joins are NP-hard in general. This
1121 makes responsibility harder to be utilized for credit distribution in a real-
1122 world application, since for this problem it is necessary to actually know the
1123 responsibility value, not simply the ranking amongst tuples.

1124 The Shapley Value has (at least) four properties that are widely believed
1125 to be important:

- 1126 1. *Efficiency*: The sum of the Shapley values of all agents equals the
1127 value of the *grand coalition*, so that all the gain is distributed among
1128 the agents.
- 1129 2. *Symmetry*: If i and j are two actors who are equivalent in the sense
1130 that $v(S \cup \{i\}) = v(S \cup \{j\})$ for every subset $S \subseteq N$ that contained
1131 neither i nor j , then their Shapley values are the same.
- 1132 3. *Null player*: The Shapley value of a *null* player i in a game v is zero.
1133 A player i is null if $v(S \cup \{i\}) = v(S)$ for all coalitions S such that
1134 $i \notin S$.
- 1135 4. *Linearity*: If two coalition games described by gain functions v and w
1136 are combined, then the distributed gains should correspond to the gains
1137 derived from v and the gains derived from w , that is: $Shap_i(v + w) =$
1138 $Shap_i(v) + Shap_i(w)$ for every $i \in N$. Also, for any real number a ,
1139 $Shap_i(a \cdot v) = a \cdot Shap_i(v)$.

1140 Livshits et al. [48] studied the computational complexity of calculating
1141 the Shapley values in query answering. They showed lower bounds on the
1142 complexity of the problem, with the exception of the sub-class of self-join free
1143 SPJ queries called *hierarchical* queries, where they gave a polynomial-time

1144 algorithm (which, however, do not appear to be useful for real world sce-
1145 narios [26]). Very recently, Deutch et al. [26] proved that the Shapley value
1146 can be efficiently (polynomial-time) reduced to probabilistic query answering.
1147 This not only applies to hierarchical queries, but to general SPJ queries. This
1148 means that one can compute Shapley values using a query engine for prob-
1149 abilistic databases, for example, the practically effective *Knowledge Compi-*
1150 *lation* [40]. More precisely, the approach in [26] shows that their approach
1151 can exactly compute the Shapley value quickly in most cases while, in other
1152 cases, the relative order given by the Shapley value may be obtained. This
1153 new work makes the Shapley value a viable solution for Credit Distribution
1154 for many queries.

1155 We can conclude that, given the current state-of-the-art in computing
1156 provenances, the how-provenance-based DS is, at the moment, one of the
1157 most informative and cost-efficient type of provenance that can be used.
1158 The other forms of information such as responsibility and Shapley may still
1159 be used in the majority of cases, that may incur in computational problems,
1160 in particular with large databases and query logs.

1161

1162 7.3. The case of Game Provenance and Query Evaluation Games

1163

1164 *Game Provenance.* Köhler et al. [44] described the notion of *game prove-*
1165 *nance*, i.e. a form of provenance in the context of games.

1166 A generic game is modeled as a graph $G = (V, M)$, where the set of nodes
1167 V represents the possible *positions* in the game, while the set $M \subseteq V \times V$
1168 represents the possible moves from one position to another. A *play* π is a
1169 sequence (finite or infinite) of edges M that describes the subsequent moves
1170 performed by two players, I and II, that play one after the other. The player
1171 that finds themselves in a position where no move is possible loses (π is
1172 lost by that player), at which point the other player wins (π is won by that
1173 player).

1174 Since any First-Order (FO) query $\varphi(\bar{x})$ on an input database instance
1175 D can be expressed as a non-recursive Datalog⁻ (Datalog with negation)
1176 program Q_φ , Köhler et al. [44] observe that the evaluation of $Q(D)$ can
1177 be seen as a game between players I and II who argue whether an atom
1178 $A \in Q(D)$.

1179 [44] also shows that game provenance coincides with semiring provenance
1180 (i.e., how provenance) for positive queries but that, unlike semiring prove-

1181 nance, it naturally extends to full FO queries with negation. This provenance
1182 can be represented as a particular type of tree, called *operator tree*.

1183 Therefore, game provenance opens up new possibilities for credit distribu-
1184 tion. First of all, new DSs based on the information provided by the operator
1185 trees of queries can be devised. These new DSs can be based on the operator
1186 tree topology, propagating the credit as a flux through its nodes and edges,
1187 devising new methods and dynamic for the distribution. Second, new DSs
1188 for the class of FO queries with negation may be devised. In particular, as
1189 shown in [47], these operator trees can also be used for *why-not provenance*,
1190 i.e., to explain the *absence* of a fact from the query output. In this case,
1191 new strategies may produce credit corresponding to “missing” facts in the
1192 query output. This, in turn, may allow to assign credit to “missing” facts in
1193 the database instance whose absence is critical for the missing output fact.
1194 This information can be useful for the database administrators to under-
1195 stand if some valuable information is missing, and help them decide whether
1196 and where to allocate the necessary resources to create/add those data if
1197 possible/sensible.

1198 8. Conclusions and Future Work

1199 This paper expanded on our previous work on data credit and data credit
1200 distribution based on the notion of lineage in [27] by defining three new dis-
1201 tribution strategies based on how-provenance, responsibility, and the Shapley
1202 Value. The how-provenance-based DS considers the frequency with which a
1203 tuple or combination of tuples is used in the query through the information
1204 contained in a provenance polynomial. In this case, the how-provenance-
1205 based DS is more sensitive than the lineage-based DS to the role and im-
1206 portance of tuples. The second DS exploits the notion of responsibility, a
1207 real value that ranks the lineage tuples based on their degree of causality in
1208 generating the output. The third DS is based on the Shapley value function,
1209 used to rank the facts of the database, seen as players, in producing the re-
1210 quired result. To do so, the wealth function in the Shapley value’s definition
1211 was adapted for general free-variable queries on the database.

1212 To show the differences between the three new DSs, we performed exten-
1213 sive experiments based on GtoPdb, a curated scientific relational database,
1214 using both real and synthetic queries. In the first set of experiments, we
1215 used select-project-join (SPJ) queries extracted from citations to webpages
1216 in GtoPdb found in papers published in the British Journal of Pharmacol-

1217 ogy. Using these “real” queries, we distributed credit to tuples in different
1218 tables of the database, highlighting tuples that were more frequently used.
1219 We showed that, with these queries, the three strategies produce the same
1220 distribution. This is because the SPJ queries were fairly simple, and did not
1221 use self-joins. Therefore the formulas underlying the different DSs had the
1222 same output.

1223 In the second set of experiments, we synthetically produced more com-
1224 plex provenance polynomials, corresponding to more complex queries, that
1225 resulted in exponents and coefficients in the provenance polynomials that
1226 were greater than (or equal to) 1. These experiments highlighted the differ-
1227 ences between the three DSs. While the DS based on lineage presented in
1228 [27] rewards all the tuples used by a query equally, the strategy based on
1229 responsibility gives more credit to tuples that are more critical to the query.
1230 Responsibility considers the relative importance of a tuple in the generation
1231 of the output. The DS based on the Shapley value similarly rewards the
1232 tuples based on their participation. The more impactful the role of a tuple,
1233 the higher its reward in credit. This distribution proved to be different from
1234 the previous two and to reward even more tuples that are used in more than
1235 one monomial. How-provenance is even more sensitive to the tuple’s role: it
1236 also considers the frequency with which a tuple or a set of tuples is used.

1237 In the third set of experiments, we showed how the differences between
1238 the DS are compounded over time, i.e. when more and more queries are
1239 processed by the system.

1240 In the fourth set of experiments we compared traditional citations to
1241 authors to the credit accrued to them via the DSs. We showed how, in
1242 both real-world and synthetic scenarios, credit rewards authors who con-
1243 tribute/curate data that has the highest impact, and therefore receives the
1244 biggest quantity of credit, and not necessarily the data with the highest ci-
1245 tation count. In this sense, credit appears to be an useful new measure to
1246 discover data and their corresponding curators that have a high impact in
1247 the research world, even when they are cited few times or do not appear at
1248 all in the data that are cited (i.e., the case of data used to build the output
1249 of a query but that is not visualized in the output itself).

1250 In the last set of experiments we showed how, on GtoPdb, all the ap-
1251 proaches present reasonable execution times, but we noted how the compu-
1252 tation of Shapley value may become unfeasible on bigger databases and with
1253 bigger queries. Very recent works such as [26] showed that it is still possible
1254 to compute the Shapley value in polynomial time in many cases.

1255 In future work, we plan to explore different strategies to generate and
1256 distribute credit. In this paper we assumed that each output tuple carries
1257 credit 1. In more sophisticated scenarios we can employ different strategies
1258 to compute credit, that reflect the importance of cited data. Other, more
1259 sophisticated, strategies could also be used to decide how credit is distributed
1260 between the authors, beyond the uniform distribution used here, in a way
1261 to reflect the work performed by them on the cited data. There are also a
1262 number of other intriguing applications for credit over relational databases.
1263 One such application is *data pricing*, which gives a price to a query submitted
1264 by a user who wants to buy the produced information. Currently, a common
1265 strategy used for data pricing is based on query rewriting: A database stores a
1266 set of views with their price. When a new query arrives, the system rewrites
1267 it using the stored views to obtain a query price, a process that can be
1268 computationally expensive. We plan to distribute credit through carefully
1269 planned and representative queries, and use credit information to define a
1270 new, faster, and potentially more flexible pricing function.

1271 Another application is *data reduction* [52], which addresses the problem of
1272 reducing the vast – and rapidly expanding – amount of data that is being pro-
1273 duced. Data credit can help address this problem by identifying “hotspots”
1274 and “coldspots” of data. A hot spot is data in a database (e.g. a tuple) with
1275 a high quantity of credit, which is therefore valuable for the set of queries
1276 that execute frequently over the data and distribute the credit. A cold spot
1277 is data with a low quantity of credit which can therefore be considered as less
1278 important, and could be deleted, summarized, or moved to cheaper and/or
1279 less efficient memory.

1280 Acknowledgement

1281 The work was partially supported by the ExaMode project, as part of the
1282 European Union H2020 program under Grant Agreement no. 825292.

- 1283 [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bern-
1284 stein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L.,
1285 Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Koss-
1286 mann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo,
1287 T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo,
1288 A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suci, D.
1289 (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–
1290 53.

- 1291 [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating
1292 data citation in citedb. *PVLDB*, 10(12):1881–1884.
- 1293 [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y.
1294 (2018). Data citation: A new provenance challenge. *IEEE Data Eng.*
1295 *Bull.*, 41(1):27–38.
- 1296 [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An
1297 Introduction to the Joint Principles for Data Citation. *Bulletin of the*
1298 *Association for Information Science and Technology*, 41(3):43–45.
- 1299 [5] Amsterdamer, Y., Deutch, D., and Tannen, V. (2011). Provenance for ag-
1300 gregate queries. In Lenzerini, M. and Schwentick, T., editors, *Proceedings*
1301 *of the 30th ACM SIGMOD-SIGACT-SIGART Symposium on Principles*
1302 *of Database Systems, PODS 2011*, pages 153–164. ACM.
- 1303 [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D.,
1304 Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I.,
1305 Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and
1306 Goble, C. A. (2013). Why linked data is not enough for scientists. *Future*
1307 *Gener. Comput. Syst.*, 29(2):599–611.
- 1308 [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation
1309 Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.
- 1310 [8] Berstel, J. and Perrin, D. (1985). *Theory of codes*. Academic Press.
- 1311 [9] Bertin-Mahieux, T., Ellis, D., Whitman, B., and Lamere, P. (2011). The
1312 million song dataset. In *Proceedings of the 12th International Conference*
1313 *on Music Information Retrieval (ISMIR 2011)*, pages 591–596.
- 1314 [10] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In
1315 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*
1316 *nication*, pages 93–116. De Gruyter Mouton.
- 1317 [11] Buneman, P. (2006). How to cite curated databases and how to make
1318 them citable. In *18th International Conference on Scientific and Statistical*
1319 *Database Management, SSDBM*, pages 195–203. IEEE Computer Society.
- 1320 [12] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D.,
1321 Pawson, A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn’t

- 1322 working, and what to do about it. *Database J. Biol. Databases Curation*,
1323 2020.
- 1324 [13] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation
1325 is a computational problem. *Commun. ACM*, 59(9):50–57.
- 1326 [14] Buneman, P., Khanna, S., and Tan, W. C. (2001). Why and where: A
1327 characterization of data provenance. In *Database Theory - ICDT 2001*,
1328 *8th International Conference*, pages 316–330.
- 1329 [15] Buneman, P. and Silvello, G. (2010). A rule-based citation system for
1330 structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.
- 1331 [16] Callaghan, S., Donegan, S., Pepler, S., Thorley, M., Cunningham, N.,
1332 Kirsch, P., Ault, L., Bell, P., Bowie, R., Leadbetter, A. M., Lowry,
1333 R. K., Moncoiffé, G., Harrison, K., Smith-Haddon, B., Weatherby, a.,
1334 and Wright, D. (2012). Making Data a First Class Scientific Output:
1335 Data Citation and Publication by NERC’s Environmental Data Centres.
1336 *International Journal of Digital Curation*, 7(1):107–113.
- 1337 [17] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Jour-
1338 nals: A Survey. *Journal of the Association for Information Science and*
1339 *Technology*, 66(9):1747–1762.
- 1340 [18] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases:
1341 Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–
1342 474.
- 1343 [19] Chockler, H. and Halpern, J. Y. (2004). Responsibility and blame: A
1344 structural-model approach. *J. Artif. Intell. Res.*, 22:93–115.
- 1345 [20] CODATA-ICSTI Task Group on Data Citation Standards and Practices
1346 (2013). *Out of Cite, Out of Mind: The Current State of Practice, Policy,*
1347 *and Technology for the Citation of Data*, volume 12.
- 1348 [21] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N.
1349 (2019). Bringing citations and usage metrics together to make data count.
1350 *Data Science Journal*, 18(1).
- 1351 [22] Cronin, B. (1984). *The Citation Process. The Role and Significance of*
1352 *Citations in Scientific Communication*. London: Taylor Graham.

- 1353 [23] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evi-
1354 dence of a structural shift in scholarly communication practices? *JASIST*,
1355 52(7):558–569.
- 1356 [24] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of
1357 view data in a warehousing environment. *ACM Trans. Database Syst.*,
1358 25(2):179–227.
- 1359 [25] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model
1360 for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on*
1361 *Innovative Data Systems Research*. www.cidrdb.org.
- 1362 [26] Deutch, D., Frost, N., Kimelfeld, B., and Monet, M. (2021). Computing
1363 the shapley value of facts in query answering.
- 1364 [27] Dosso, D. and Silvello, G. (2020). Data credit distribution: A
1365 new method to estimate databases impact. *Journal of Informetrics*,
1366 14(4):101080.
- 1367 [28] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The vir-
1368 tual atomic and molecular data centre (VAMDC) consortium. *Journal of*
1369 *Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- 1370 [29] Eiter, T. and Lukasiewicz, T. (2002). Complexity results for structure-
1371 based causality. *Artif. Intell.*, 142(1):53–89.
- 1372 [30] Fang, H. (2018). A discussion of citations from the perspective of the
1373 contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–
1374 1520.
- 1375 [31] Garfield, E. (1999). Journal impact factor: a brief review. *Can. Med.*
1376 *Assoc.*, 979-980.
- 1377 [32] Gößwein, B., Miksa, T., Rauber, A., and Wagner, W. (2019). Data
1378 identification and process monitoring for reproducible earth observation
1379 research. In *2019 15th International Conference on eScience (eScience)*,
1380 pages 28–38. IEEE.
- 1381 [33] Green, T. J., Karvounarakis, G., and Tannen, V. (2007). Provenance
1382 semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-*
1383 *SIGART symposium on Principles of database systems*, pages 31–40. ACM.

- 1384 [34] Halpern, J. Y. and Pearl, J. (2013). Causes and explanations: A
1385 structural-model approach — part 1: Causes. *CoRR*, abs/1301.2275.
- 1386 [35] Harding, S. D., Sharman, J. L., Faccenda, E., Southan, C., Pawson,
1387 A. J., Ireland, S., Gray, A. J. G., Bruce, L., Alexander, S. P. H., Anderton,
1388 S., Bryant, C., Davenport, A. P., Doerig, C., Fabbro, D., Levi-Schaffer, F.,
1389 Spedding, M., Davies, J. A., and Nc-Iuphar (2018). The IUPHAR/BPS
1390 guide to PHARMACOLOGY in 2018: updates and expansion to encom-
1391 pass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Re-*
1392 *search*, 46(Database-Issue):D1091–D1106.
- 1393 [36] Hartley, J. (2017). Authors and their citations: a point of view. *Scien-*
1394 *tometrics*, 110(2):1081–1084.
- 1395 [37] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a
1396 transformed scientific method.
- 1397 [38] Honor, L. B., Haselgrove, C., Frazier, J. A., and Kennedy, D. N. (2016).
1398 Data citation in neuroimaging: proposed best practices for data identifi-
1399 cation and attribution. *Frontiers in neuroinformatics*, 10:34.
- 1400 [39] Jain, S., Moritz, D., Halperin, D., Howe, B., and Lazowska, E. (2016).
1401 Sqlshare: Results from a multi-year sql-as-a-service experiment. In *Pro-*
1402 *ceedings of the 2016 International Conference on Management of Data*,
1403 pages 281–293.
- 1404 [40] Jha, A. K. and Suci, D. (2013). Knowledge compilation meets database
1405 theory: Compiling queries to decision diagrams. *Theory Comput. Syst.*,
1406 52(3):403–440.
- 1407 [41] Joshi-Tope, G., Gillespie, M., Vastrik, I., D’Eustachio, P., Schmidt, E.,
1408 de Bono, B., Jassal, B., Gopinath, G. R., Wu, G. R., Matthews, L., Lewis,
1409 S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of bio-
1410 logical pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- 1411 [42] Katz, D. (2014). Transitive credit as a means to address social and
1412 technological concerns stemming from citation and attribution of digital
1413 products. *Journal of Open Research Software*, 2(1).
- 1414 [43] Kendall, M. G. (1938). A new measure of rank correlation. *Biometrika*,
1415 30(1/2):81–93.

- 1416 [44] Köhler, S., Ludäscher, B., and Zinn, D. (2013). First-order provenance
1417 games. In *In Search of Elegance in the Theory and Practice of Computation*
1418 - *Essays Dedicated to Peter Buneman*, volume 8000 of *Lecture Notes in*
1419 *Computer Science*, pages 382–399. Springer.
- 1420 [45] Kosten, J. (2016). A classification of the use of research indicators.
1421 *Scientometrics*, 108(1):457–464.
- 1422 [46] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S.
1423 (2011). Citation and Peer Review of Data: Moving Towards Formal Data
1424 Publication. *International Journal of Digital Curation*, 6(2):4–37.
- 1425 [47] Lee, S., Ludäscher, B., and Glavic, B. (2018). PUG: A framework and
1426 practical implementation for why & why-not provenance (extended ver-
1427 sion). *CoRR*, abs/1808.05752.
- 1428 [48] Livshits, E., Bertossi, L. E., Kimelfeld, B., and Sebag, M. (2020). The
1429 shapley value of tuples in query answering. In Lutz, C. and Jung, J. C.,
1430 editors, *23rd International Conference on Database Theory, ICDT 2020,*
1431 *March 30-April 2, 2020, Copenhagen, Denmark*, volume 155 of *LIPIcs*,
1432 pages 20:1–20:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.
- 1433 [49] Martone, M. (2014). Joint declaration of data citation principles.
1434 *FORCE11. San Diego CA. Data Citation Synthesis Group*. [https://www.](https://www.force11.org/datacitationprinciples)
1435 [force11.org/datacitationprinciples](https://www.force11.org/datacitationprinciples), online September 2020.
- 1436 [50] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation
1437 counts and rankings of LIS faculty: Web of science versus scopus and
1438 google scholar. *Journal of the american society for information science*
1439 *and technology*, 58(13):2105–2125.
- 1440 [51] Meliou, A., Gatterbauer, W., Moore, K. F., and Suciu, D. (2010). The
1441 complexity of causality and responsibility for query answers and non-
1442 answers. *Proc. VLDB Endow.*, 4(1):34–45.
- 1443 [52] Milo, T. (2019). Getting rid of data. *Journal of Data and Information*
1444 *Quality (JDIQ)*, 12(1):1–7.
- 1445 [53] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D.,
1446 Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G.,
1447 Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff,

- 1448 D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,
1449 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson,
1450 E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C.,
1451 Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers,
1452 E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research
1453 culture. *Science*, 348(6242):1422–1425.
- 1454 [54] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J.
1455 (2016). Research data explored: An extended analysis of citations and
1456 altmetrics. *Scientometrics*, 107(2):723–744.
- 1457 [55] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic,
1458 large databases: Model and reference implementation. In *Proceedings of*
1459 *the 2013 IEEE International Conference on Big Data, 6-9 October 2013,*
1460 *Santa Clara, CA, USA*, pages 307–312.
- 1461 [56] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identi-
1462 fication of Reproducible Subsets for Data Citation, Sharing and Re-Use.
1463 *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue*
1464 *on Data Citation*, 12(1):6–15.
- 1465 [57] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data
1466 citation of evolving data: Recommendations of the working group on data
1467 citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- 1468 [58] Remil, Y., Bendimerad, A., Mathonat, R., Chaleat, P., and Kaytoue,
1469 M. (2021). ” what makes my queries slow?”: Subgroup discovery for sql
1470 workload analysis. *arXiv preprint arXiv:2108.03906*.
- 1471 [59] Shapley, L. S. (1954). A value for n-person games. In Kuhn, H. W. and
1472 Tucker, A. W., editors, *Contributions to the Theory of Games II*, pages
1473 307–317. Princeton University Press, Princeton.
- 1474 [60] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*
1475 *Sci. Technol.*, 69(1):6–20.
- 1476 [61] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data
1477 provenance in e-science. *SIGMOD Record*, 34(3):31–36.
- 1478 [62] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-
1479 spective. In of Sciences’ Board on Research Data, N. A. and Information,

- 1480 editors, *Report from Developing Data Attribution and Citation Practices*
1481 *and Standards: An International Symposium and Workshop*, pages 177–
1482 178. National Academies Press: Washington DC.
- 1483 [63] Vogelsgesang, A., Haubenschild, M., Finis, J., Kemper, A., Leis, V.,
1484 Mühlbauer, T., Neumann, T., and Then, M. (2018). Get real: How bench-
1485 marks fail to represent the real world. In *Proceedings of the Workshop on*
1486 *Testing Database Systems*, pages 1–6.
- 1487 [64] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,
1488 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,
1489 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data
1490 management and stewardship. *Scientific data*, 3.
- 1491 [65] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data
1492 citation: Giving credit where credit is due. In *Proceedings of the 2018*
1493 *International Conference on Management of Data, SIGMOD*, pages 99–
1494 114.
- 1495 [66] Wu, Y., Alawini, A., Deutch, D., Milo, T., and Davidson, S. B. (2019).
1496 ProvCite: provenance-based data citation. *Proceedings of the VLDB En-*
1497 *dowment*, 12(7):738–751.
- 1498 [67] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to
1499 scientific datasets using article citation networks. *Journal of Informetrics*,
1500 14(2).
- 1501 [68] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of
1502 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.
- 1503 [69] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for
1504 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*
1505 *Molecular Spectroscopy*, 327:122–137.