

Credit Distribution through Data Provenance in Relational Scientific Databases

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

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

^b*Department of Computer and Information Science, University of Pennsylvania, USA*

Abstract

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

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

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

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

Using these distribution strategies we show how credit can highlight frequently used database areas and how it can be used as a new bibliometric measure for data and their corresponding curators. In particular, credit rewards data and authors based on their research impact, not merely on the number of citations. We also show how 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

1. Introduction

Citations are an essential component of scientific research, enabling research products to be found as well as the relationships between them to be created and understood. They form a basis on which to give credit to authors, papers, and venues [20, 21, 63]. Citations are used, among other things, to decide on tenure, promotion, hiring, and funding of grants for researchers [22, 36, 42, 46].

Science and research are increasingly digital, and there are numerous curated databases that are at the core of scientific research efforts [12]. It is therefore generally accepted that data must be cited and citable [15, 43], and that data citations should contribute to the scientific reputation of researchers, scientists, data curators, and creators [4, 58]. It is also accepted that data citations should be counted alongside of traditional citations, and contribute to bibliometrics indicators [7, 50].

A central problem in the data citation process is how to attribute credit to data creators and curators [11]. How to handle and count the credit generated by data citation, and how it contributes to traditional and new bibliometrics, are long-standing research issues [9, 31]. However, even when correctly applied, data citations and the bibliometrics computed using them do not always fully reward the creators of data used in a database. Data, in fact, is often cited at the “database level” or the “webpage level”. In the first case, the whole database is cited and therefore all credit goes to the key personnel of the database. In the second case, the database has a website with webpages that can be individually cited. The webpages use data extracted from the database, which is aggregated by topic and built to resemble a traditional research paper. Often the creators and curators of the webpage’s data are not credited or only marginally credited for their work [3].

Recently, the idea of *Data Credit Distribution* (DCD) [30, 41, 62] has emerged, built on top of methodologies for data citation. Data credit is a value that is computed based on the importance of the data being cited in a paper, and is a proxy for the impact of the data on the citing paper. The DCD problem consists of distributing this credit to elements in the databases in the citation graph that are responsible for the generation of the data being cited. The goal of DCD is to improve and expand the reach of data citation,

36 rather than being an alternative to it.

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

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

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

62 We study the DCD problem in the context of relational databases (RDBs)
63 since they are widely used² and are the main focus of current work in data
64 citation methods [12, 14, 51]. RDBs are also frequently a test-bed for new
65 methods that can be adapted to other databases, e.g., graphs or document
66 databases. The “portions” of data in an RDB that can be credited can be
67 defined at different levels of granularity, in particular: (i) the whole database,
68 (ii) tables, (iii) tuples, and (iv) attributes. The ability to specify different

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.

69 levels of granularity in a relational database allows us to define the DCD
70 problem at a particular level of granularity. In this paper, we focus on DCD
71 at the tuple level.

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

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

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

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

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

83 **Step 5** Provenance is input to the *Credit Distribution Strategy* (CDS, also
84 referred only as *Distribution Strategy*, DS). CDS is a function f that

85 takes as input the credit value k , divides it and distributes it to the
86 data in the input database I , and is defined on the basis of citation
87 policies decided at the database administration level or at the domain
88 community level.

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

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

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

107 Therefore, in this paper, we expand the ideas in [27] by proposing new DSs
108 based on three other forms of data provenance: why-provenance [13], how-
109 provenance [33], and responsibility [47].³ We show how these DS differ from
110 each other as well as the one based on lineage, and discuss why one may be
111 preferred to another depending on the application and its goals. In particular,
112 we show that why-provenance-, how-provenance- and responsibility- based
113 DSs are more sensitive than one based on lineage to the *role* of a tuple in a
114 query, i.e. how many times the tuple is used and how it is used. We also
115 show that why-provenance- and responsibility-based DSs give more credit
116 to tuples that are essential to the production of the result set, whereas the
117 how-provenance-based DS also takes into consideration the different ways in

³A discussion of basing a DS on the recently proposed notion of Shapley value[25, 44] can be found in Section 7.2.

which a tuple is used. The newly defined DS strategies exploit the additional information provided by these forms of provenance to weight the contribution of the tuples on the basis of their role in producing the output.

The evaluation is based on a well-known curated database, the IUPHAR/BPS⁴ Guide to Pharmacology [35], also known as GtoPdb⁵, which contains expertly curated information about diseases, drugs, cellular drug targets, and their mechanisms of action. We chose GtoPdb for two main reasons: (i) it is a widely-used and valuable curated relational database, (ii) many papers in the literature use, and cite, its data (i.e., families, ligands, and receptors). Real queries used in papers can therefore be seen as data citations which, in turn, can be used to assign data credit.

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

Contributions of this work include:

- Three new Distribution Strategies based on why-provenance, how-provenance, and responsibility.
- An in-depth analysis of the effects of credit distribution on real-world curated data and of the differences between the three proposed Distribution Strategies.
- A comparison between the behavior of traditional citations and data credit in rewarding data curators.

Outline. The rest of the paper is organized as follows: Section 2 presents background material and related work. Section 3 describes the GtoPdb use case. Section 4 briefly presents the forms of provenance used in the paper.

⁴International Union of Basic and Clinical Pharmacology/British Pharmacology Society

⁵<https://www.guidetopharmacology.org/>

149 Section 5 describes the credit distribution problem and the proposed dis-
150 tribution strategies. In Section 6 we present the experimental evaluation,
151 followed by a discussion of our design decisions in Section 7. Section 8 draws
152 some conclusions and outlines future work.

153 2. Background

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

158 The scientific community is promoting an *open research culture* [49],
159 founded on methods and tools to share, discover, and access experimental
160 data. The community has identified the FAIR principles (Findable, Acces-
161 sible, Interoperable, and Reusable) [60], that should be enforced by every
162 database. In particular, data should be accessible from the articles, journals,
163 and papers that cite or use them [20]. Aspects such as the need for the *repro-*
164 *ducibility* of experiments through the used data; the *availability* of scientific
165 data; the *connections* between data and the scientific results are all needed
166 aspects for the fourth paradigm, and are all relevant to the domain of *data*
167 *citation* [38].

168 *Data Citation: Principles and Motivations.* Data Citation principles were
169 proposed in [19], and later summarized and endorsed by the Joint Declaration
170 of Data Citation Principles (JDDCP) [45]. The principles are divided into
171 two groups [56]. The first group contains principles concerning the role of
172 data citation in scholarly and research activities such as the (i) *importance*
173 of data (why data citation is important and why data should be considered
174 as first-class citizens); (ii) *credit* and *attribution* to the creators and curators
175 of the data; (iii) *evidence*; (iv) *verifiability*; and *interoperability*, with these
176 last three requiring data citation methods to be flexible enough to operate
177 through different communities. The second group defines the main guidelines
178 to establish a data citation systems, and contains principles such as the (i)
179 *unique identification* of the data being cited; (ii) *(open) access* to data; (iii)
180 guarantee of *persistence* and *availability* of citations even after the lifespan
181 of the cited entity; the (iv) *specificity* of a citation, i.e. it must lead to the
182 data set originally cited.

183 * SBD: Is the next paragraph necessary? Could we just say
184 "The main motivations for data citation are outlined in [56]." *

185 It is possible to outline six main motivations for data citation [56]:

- 186 • *Data attribution*: identify the individuals that should be credited for
187 data with variable granularity.
- 188 • *Data connection*: connect papers to the data being used.
- 189 • *Data Discovery*: citations helps to find data records and subsets that
190 would be otherwise not findable via search engines.
- 191 • *Data Sharing*: share data obtained by researchers within the whole
192 community.
- 193 • *Data Impact*: highlight the results obtained in writing papers using
194 specific data, the frequency and modality data were used.
- 195 • *Reproducibility*: data citation greatly impacts the reproducibility of
196 science [5]. Many authoritative journals ask to share data and provide
197 valid methodologies to reproduce experiments.

198 2.1. Data Citation in Relational Databases

199 Relational databases have been the target of data citation methods since
200 the surge of the data-centric research paradigm. The RDA “Working Group
201 on Data Citation: Making Dynamic Data Citable”⁶ [52] has developed guide-
202 lines for citing large, dynamic, and changing datasets which have now moved
203 on into adoption phase. The datasets considered by the Working Group are
204 often relational.

205 In one of its most recent sessions [53], the Working Group (WG) on
206 Data Citation reported that there are various implementations of its guide-
207 lines for Data Citation on MySQL/Postgres relational databases. Some of
208 these databases are: DEXHELPP⁷ (Social Security Records); NERC (ARGO
209 Global Array); EODC (Earth Observation Data Centre) [32]; LNEC (River
210 dam monitoring); MDS (Million Song Database) [8]; CBMI⁸ (Center for
211 Biomedical Informatics); VMC (Vermont Monitoring Cooperative); CCA⁹

⁶<https://www.rd-alliance.org/groups/data-citation-wg.html>

⁷<http://www.dexhelpp.at/>

⁸<https://medicine.missouri.edu/centers-institutes-labs/center-for-biomedical-informatics>

⁹<https://ccca.ac.at/startseite>

212 (Climate Change Center Austria); VAMDC (Virtual Atomic and Molecular
213 Data Center) [28, 64].

214 More examples of work on data citation in relational databases are [2,
215 12, 24, 61]. The website <https://fairsharing.org/> keeps an updated list
216 of curated and scientific databases (many of which are relational or graph-
217 based) following FAIR guidelines. These databases are citable since they are
218 compliant with the most recent guidelines, and they are in the vast majority
219 of cases accessible via dynamically created Webpages. In all these databases
220 it is, therefore, possible to implement DCD on top of the existing infrastruc-
221 tures for citing data.

222 Data citation techniques are primarily applied to relational databases
223 because of their pervasiveness as well as the “identifiability” of the portions
224 of data that are to be cited: the whole database, a relation, a tuple, or
225 even an attribute. Many papers [2, 10, 12] consider more complex citable
226 units, recognizing that often the *views* of a database are the ones to be cited.
227 Generally, a *view* is a query on the database. To this end, [61] suggested
228 decomposing the database into a set of views, where each view is associated
229 with its citation.

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

- 231 1. A database cited as a whole, even though only parts of the databases
232 are used in the papers or datasets. Alternatively, the so-called “data pa-
233 pers” are cited, being traditional papers that describe a database [16].
234 In this case, all the credit from the citations goes to the database ad-
235 ministrators or to the authors of the data papers.
- 236 2. Subsets of data, obtained by issuing queries to a database, are individ-
237 ually cited. This is the solution adopted by the *Resource Data Alliance*
238 (RDA) working group on Data Citation [52]. In this case, the credit
239 generated from citations is distributed among the contributors of the
240 portions of data being cited, and/or to the database administrators.
- 241 3. The database is accessible via a series of Webpages that arrange the
242 content of the database by topic or theme. Examples in the life science
243 domain include the Reactome Pathway database [40], the GtoPdb [35],
244 and the VAMDC [64]. Every single Webpage is unequivocally identifi-
245 able and can be individually cited.

246 2.2. Data Credit

247 Data credit is related to data citation: they both aim to recognize the
248 work of data creators and curators. Data credit can therefore also be seen as

249 a by-product of data citation, since credit attribution is impossible without
250 the presence of data citations.

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

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

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

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

- 279 1. Credit rewards research entities that to date are not (formally) recog-
280 nized (a goal shared with data citation).
- 281 2. Credit can reward authors *proportionally* to their role in generating the
282 entity. The more an author contributes to a paper, the more credit is
283 given to him. Zou and Peterson [63] work on something similar with
284 their zp-index, which includes in its formulation the position (and thus
285 the role) of a publication author to represent its impact in the work

286 itself.
287 3. Credit can be *transitively* channeled through a chain of papers citing
288 each other, thus enabling the rewarding of older papers that are no
289 more cited, since other papers summarize or report their content but
290 are nevertheless crucial in a research area for the influence of their
291 content.

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

301 Zeng et al. [62] proposes the first method to compute credit within a net-
302 work of papers citing data. Adopting a network flow algorithm, they simulate
303 a random walker to estimate a score for each dataset, leveraging real-world
304 usage data to compute the credit. This is the first step towards an automatic
305 credit computation procedure. This proposal is, however, limited to assign-
306 ing credit to whole datasets, and it does not deal with the granularity of data.
307 It does not work to assign credit to a single research entity within a dataset.
308 Differently from Zeng et al. [62], we do not treat the credit computation
309 process, but we focus on the distribution process.

310 2.3. Data Provenance

311 To distribute credit, we base our methods on *data provenance*. Data
312 provenance is information that describes the origin and the process of cre-
313 ation of data. It can also be seen as metadata pertaining to the derivation
314 history of the data. It is particularly useful to help users to understand
315 where data are coming from, and the process they went through. Data ci-
316 tation and data provenance are closely linked [3] since both are forms of
317 annotations on data retrieved through queries. Data provenance has been
318 widely studied in different areas of data management. In this paper, we fo-
319 cus on provenance for database management systems (DBMS). For further
320 details on data provenance, please refer to surveys like [17] and [57].

321 Cheney et al. [17] presents four main types of data citation for DBMS: *lin-*
322 *eage* [23], *why-provenance* [13], *how-provenance* [33] and *where-provenance* [13].

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

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

332 We also consider the notions of causality and responsibility, as defined
333 in [47]. Causality is an enrichment of lineage, and it is the attribution of
334 a certain degree of importance to the tuples of the lineage based on their
335 role in the generation of the output. Responsibility is a value given to the
336 tuples of the lineage to rank them based on their degree of causality (the
337 more important the role of a tuple in generating the output, the higher its
338 responsibility).

339 3. Use Case: GtoPdb

340 The IUPHAR/BPS Guide to Pharmacology [35] (GtoPdb¹⁰) is a well-
341 known and well structured scientific relational database that contains ex-
342 pertly curated information about diseases, drugs in clinical use, their cellular
343 targets, and the mechanisms of action on the human body. It is curated and
344 maintained by the GtoPdb Committee and 96 subcommittees, comprising
345 512 scientists collaborating with in-house curators who draw the information
346 contained in the database from high-quality pharmacological and medicinal
347 chemistry literature. Roughly 1000 researchers from all over the world have
348 contributed to the database, and the curators wanted to give recognition to
349 these contributors. This led to some early work on data citation [10].

350 GtoPdb is relational, but its logical structure is hierarchical as shown
351 in Figure 2. The information contained in the database is also organized
352 into webpages focused on specific diseases, targets or ligands, and families
353 for easier access by users. As depicted in Figure 2, the database can be

¹⁰<https://www.guidetopharmacology.org/>

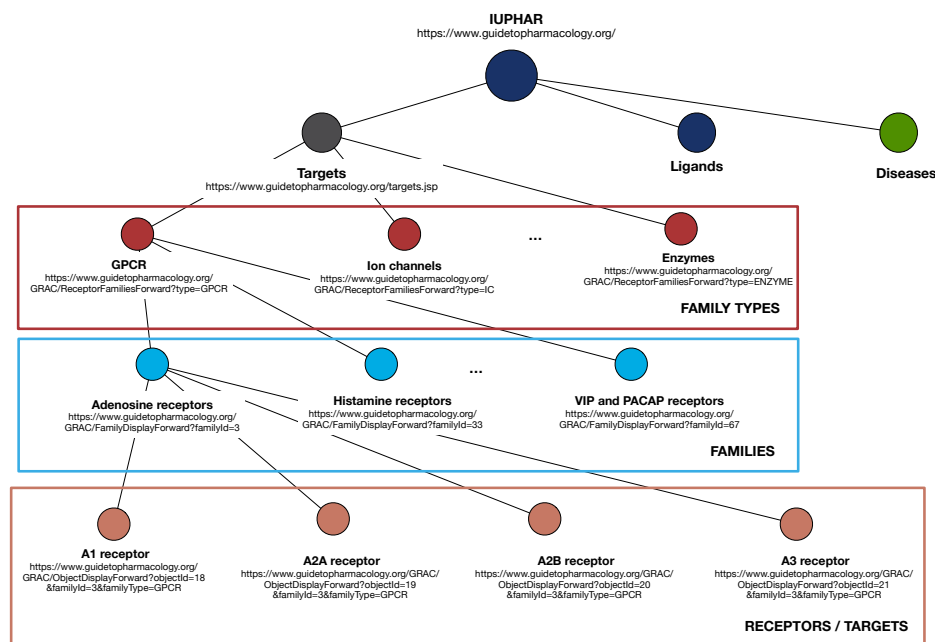


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

354 thought of as a tree where the root is the database; the first level consists
 355 of all targets, ligands, and diseases; and the lower levels consists of specific
 356 targets, ligands and diseases. In this paper, we focus on targets; thus the
 357 figure at the third level shows examples of family types, at the fourth level
 358 of specific families of targets (a finer level of granularity), and finally, at the
 359 last level, the single targets (also known as receptors).

360 GtoPdb provides access to the webpages corresponding to all these nodes
 361 through URLs. The webpages corresponding to target families all present a
 362 similar structure, as shown in Figure 3 for the “Adenosine receptors” family.
 363 Each page has an *Overview*, a brief text describing the content of the page;
 364 a list of *Receptors* comprising the family; a section of *comments* about the
 365 family; the *References*, a list of the papers consulted by the curators of the
 366 page, similar to a reference list of a paper; the *further reading* list, reporting
 367 papers that an interested reader may want to consult to obtain more insight
 368 on the family; and a final section called *How to cite this family page*, con-
 369 taining text snippets useful to cite the specific page or the whole database.
 370 Figure 3 shows the SQL code that retrieves the information used to build the

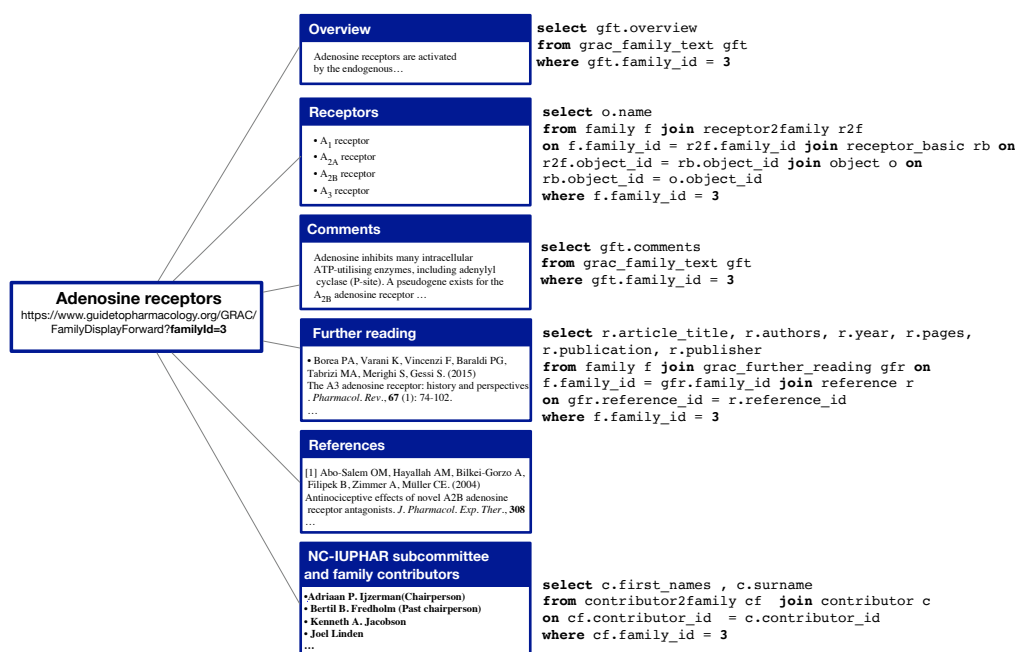


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

371 corresponding sections (apart from the References section). Therefore, each
 372 family page can be considered a full-fledged traditional publication, consist-
 373 ing of title, authors, abstract (the overview), content, and references.

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

380 In certain “lucky” cases, as with papers available in PDF and published
 381 in the British Journal of Clinical Pharmacology ¹¹ (BJCP), when a family,
 382 ligand, receptor name, etc. are used, they have a hyperlink pointing to the
 383 corresponding webpage in GtoPdb. Therefore, the citations to the families
 384 can be detected and counted using the URLs reported in the papers. How-

¹¹<https://bpspubs.onlinelibrary.wiley.com/journal/13652125>

family			contributor2family		
id	name	type	id	family_id	contributor_id
f_1	Dopamine Receptors	gpcr	c_2f_1	f_1	c_1
f_2	Bile Acid Receptor	gpcr	c_2f_2	f_1	c_2
f_3	FAK Family	enzyme	c_2f_3	f_2	c_3
f_4	YANK Family	enzyme	c_2f_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.

385 ever, these citations to GtoPdb webpages are not counted as such by citation
 386 systems, so they are not converted into credit for curators and collaborators.
 387 For our running example, consider Table 1. This simplified version of
 388 GtoPdb contains three tables: **family**, **contributor** and **contributor2family**.
 389 The first table, **family**, has tuples representing families with three attributes:
 390 the id of the family, its name, and type. Table **contributor** contains peo-
 391 ple who have helped generate the data in the database. The third table,
 392 **contributor2family**, serves as a link between the families and the people
 393 who contributed to them. For instance, “John Smith” (c_1) contributed to
 394 “Dopamine Receptors” (f_1) as well as to the “YANK Family” (f_4). Through-
 395 out the rest of the paper, we will use the **id** attribute of these tables as the
 396 *provenance token* of its corresponding tuples, that is, as a symbol that serves
 397 to identify a tuple when talking about provenance.

398 4. Data Provenances

399 We now describe the three types of provenance used in this paper – lin-
 400 eage, why-provenance, and how-provenance – as well as the notion of Causal-
 401 ity and Responsibility.

4.1. Lineage

Lineage is the simplest form of provenance. It was first introduced by Cui et al. [23], and can be thought of as the set of all tuples that are used by the query to generate the output [17].

As an example, consider the following SQL query **Q1**, applied to the database described in Table 1, asking for the names of families curated by researchers based in the United Kingdom (UK):

```

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

```

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

Table 2: Result of **Q1** over the database instance in Table 1 with the lineage of each output tuple. Attribute **id** is not part of the output, and was added to identify each tuple.

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

4.2. Why-Provenance

Why-Provenance was first defined in terms of a deterministic semistructured data model and query language [13]. We use here its definition in terms of the relational model [17].

While lineage aims to find all and only the tuples in the input relevant to the production of an output tuple, why-provenance aims to find sub-instances of the input that “witness” a part of the output. Given a tuple t in the query’s output $Q(D)$, a *witness* is any sub-instance of the database that produces t , i.e., a set that guarantees the existence of t in $Q(D)$. In particular, the whole database and the lineage of t are both examples of witnesses of t . Since the

431 definition of witness allows for the presence of “irrelevant” tuples, the set
 432 of all witnesses is finite (since the database instance I is finite), but it is
 433 potentially exponentially large [17].

434 Buneman et al. [13] defined the why-provenance of an output tuple t in
 435 the result $Q(I)$ as a special *subset* of the set of witnesses called the *witness*
 436 *basis*. The witnesses of the basis exclude tuples that are irrelevant to t being
 437 produced by Q , and thus the basis tends to be very small compared to the
 438 set of all possible witnesses [17].

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

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

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

443 The why-provenance of output tuple o_2 has only one witness, which co-
 444 incides with its lineage. This happens because there is only one way this
 445 output tuple can be produced, i.e., for tuple f_4 to be joined with $c2f_4$ and c_1 .
 446 On the other hand, o_1 has a witness basis of two witnesses, since there are
 447 two possible ways in which the query can generate o_1 . One possibility is that
 448 f_1 is joined with $c2f_1$ and c_1 (the first witness), and the second possibility
 449 is that f_1 is joined with $c2f_2$ and c_2 (the second witness). This means that
 450 to generate o_1 , it is sufficient that only one of the two witnesses is present in
 451 the input database.

452 4.3. How-Provenance

453 While why-provenance describes the source tuples that witness an output
 454 tuple in the result of the query, it leaves out information about how the source
 455 tuples are used. How-provenance was therefore defined in [33] to capture
 456 this information using a *semiring* algebraic structure. It takes the form of
 457 a polynomial, called *provenance polynomial*, where the variables are taken
 458 from the set X of identifiers of the tuples (provided that each tuple in I has

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

Table 4: Result of Q1 over the database instance in Table 1 with the how-provenance polynomial of each output tuple.

an identifier) and the coefficients are drawn from the set of natural numbers \mathbb{N} .¹²

The key idea in Green et al. [33] is to use the two operators $+$ and \cdot to represent two basic transformations that source tuples undergo as a result of applying a relational query to a database [17]. Two tuples may either be joined together (a join is represented with the \cdot operator) or merged via union or projection (represented with the $+$ operator).

Table 4 shows the two output tuples of our running example annotated with their respective how-provenances. Tuple o_2 was produced by a join of the input tuples $f_4, c2f_4$, and c_1 . The three provenance tokens are therefore “multiplied” together. The case of o_1 is slightly more complex, as already discussed. It can be obtained by the joins of two different sets of tuples, so there are two monomials combined by $+$ representing these alternative derivations. Each monomial corresponds, in a way, to the witnesses of the why-provenance of o_1 .

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

4.4. Causality and Responsibility

A formal study of causality was introduced in [18, 34] and later expanded by Meliou et al. [47] to explain the causes of answers and non-answers to queries. In the following, we refer to the definition of causality and respon-

¹²This semiring is commonly referred as $\mathbb{N}[X]$ in the literature.

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

Table 5: Result of Q1 over the database instance in Table 1 with the responsibilities of lineage tuples.

486 sibility provided in [47]. In particular, we only focus on answers to a query
 487 since non-answers are not relevant in our context.

488 There are two types of “cause” tuples: counterfactual and actual. Let t
 489 be a tuple in the query’s output $Q(I)$, and t' a tuple in its lineage. We call
 490 t' a *counterfactual cause* if, by removing t' from I , t is also removed from
 491 the output (i.e., t' is essential for the generation of t). We call t' an *actual*
 492 *cause* if there is a set of tuples $\Gamma \subseteq I$ called a *contingency set*, such that t'
 493 is a counterfactual cause in $I - \Gamma$. In other words, t' is an actual cause if,
 494 even when removed from I , there is another set of tuples of the lineage that
 495 guarantees the presence of t .

496 Computing the causality of tuples is NP-complete for general queries [29],
 497 but for conjunctive queries can be computed in PTIME Meliou et al. [47].

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

Definition 4.1. *Responsibility* [47]

Let \bar{a} be an answer to a query q , and let t be a cause. The responsibility of t for the answer \bar{a} is:

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

501 where Γ ranges over all contingency sets for t .

502 Note that a counterfactual cause will have the maximum responsibility
 503 of 1, and that the larger the minimum contingency of an actual cause is, the
 504 smaller its responsibility will be since there are alternatives to guarantee the
 505 presence of the answer \bar{a} .

506 As an example, consider Table 4, where we reported the result set of Q1
 507 and the tuples of the lineages with their responsibility values. Focusing on
 508 o_1 : the lineage tuple f_1 is a counterfactual cause, since its contingency set is
 509 empty (when removed from the database, o_1 disappears from the result set).
 510 Consequently, its responsibility is 1. All the other tuples of the lineage are
 511 actual causes. c_1 , for example, has as minimal contingency set $\{c_2f_2\}$, thus

its responsibility is 0.5. For the output tuple o_2 , all the tuples of the lineage are counterfactual causes, thus their responsibility is 1.

While computing responsibility for general queries is hard [18], Meliou et al. [47] proved a dichotomy result for conjunctive queries: for each query without self-joins, either its responsibility can be computed in PTIME in the size of the database or checking if it has a responsibility below a given value is NP-hard.

519

520 4.5. Shapley Value

The Shapley value is named after Lloyd Shapley, who introduced it for the first time in his 1952 work [55]. He considered a *cooperative game* played by a set A of players, defined by a *wealth function* v that assigns to each coalition set $B \subseteq A$ the wealth $v(B)$. The question behind the Shapley Value is how to quantify the contribution of each player to the overall wealth. Informally, the Shapley value is defined as follows [44]: assume that we select players randomly one by one and without replacement, starting with the empty set. Every time a player a is selected, its addition to the coalition B produces a change in the wealth of the coalition from $v(B)$ to $v(B \cup \{a\})$. The Shapley value of a is the expectation of change that a causes in this probabilistic process.

The Shapley value can be used in different research areas beyond cooperative games, such as economics, law, environmental science, and network analysis. Here, as recently done by Livshits et al. [44], we use it as a way of quantifying the contribution of input facts (tuples) to query answers. As defined in [26], in the context of relational databases, given a query $q(\bar{x})$, a database D , an input fact $f \in D$ (here seen as a player) and a tuple \bar{t} of same arity as \bar{x} , the Shapley value of f in D intuitively represents the contribution of f to the presence (or absence) of \bar{t} in the query result. Formally, the Shapley value is defined as follows:

Definition 4.2. *Shapley value [26]*

Let the database D be partitioned into two sets of facts: a set D^x of exogenous facts, and a set D^n of endogenous facts. Let q be a Boolean query and $f \in D^n$ be an endogenous fact. The Shapley value of f in D for query q is defined

id	name	responsibility
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 6: 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.

as:

$$Shapley(q, D^n, D^x, f) = \sum_{B \subseteq D^n \setminus \{f\}} \frac{|B|! (|D^n| - |B| - 1)!}{|D^n|!} (q(D^x \cup B \cup \{f\}) - q(D^x \cup B))$$

541 The sum is performed on all possible coalitions of fact B that do not
542 contain the player f . Thus, the value $(q(D^x \cup B \cup \{f\}) - q(D^x \cup B))$ is the
543 wealth brought by f when added to B . As we see, the Boolean query is used
544 as wealth function v , thus this value will be 1 only when the set $D^x \cup B \cup \{f\}$
545 makes the query true, and the set $D^x \cup B$ still makes it false, i.e., when
546 the addition of the fact f is determinant to make the Boolean query true.
547 The value $|B|! (|D^n| - |B| - 1)!$ is the number of all the possible permutations
548 over D^n where the facts in B come first, then f is added, and then all the
549 remaining facts. Thus, the value $\frac{|B|! (|D^n| - |B| - 1)!}{|D^n|!}$ can be thought as a weight
550 for the wealth brought by the addition of f to the coalition B .

551 To extend this definition to non-Boolean queries, we use the same straight-
552 forward approach used in Deutch et al. [26]: the Shapley value of the fact f
553 for the answer \bar{t} to $q(\bar{x})$ is the value $Shapley(q[\bar{x}/\bar{t}], D^n, D^x, f)$, where $q[\bar{x}/\bar{t}]$
554 is the Boolean query defined by $q[\bar{x}/\bar{t}](D) = 1$ if and only if \bar{t} is in the output
555 of $q(\bar{x})$ on D , and 0 otherwise. *** DD: I added this paragraph down here**
556 **hoping to make things clearer. If you think it fails to do so, feel**
557 **free to delete it.** * In other words, the definition of $Shapley(q, D^n, D^x, f)$
558 is extended to such queries $q(\bar{x})$ with free variables by considering the query
559 $q[\bar{x}/\bar{t}]$ instead as value function. This query can be seen as a function that
560 takes as input a set of facts and returns 1 if this set is a witness for \bar{t} , and 0
561 otherwise.

562 As an example, consider table 6, that shows the Shapley values for the
563 lineage's tuples of o_1 and o_2 . In fact, it is only necessary to consider as endoge-
564 nous tuples the ones of the lineage, since they are the only ones contributing
565 to the generation of the output. In this case, to compute the Shapley value

566 of an input tuple f it is sufficient to compute and the values $\frac{|B|!(|D^n|-|B|-1)!}{|D^n|!}$
 567 for all the possible sets B such that $B \cup \{f\}$ is a witness and B alone instead
 568 is not. Thus, suppose we want to compute the Shapley value of the tuple
 569 f_1 . Let us call \bar{Q}_{1,o_1} the Boolean query such that $\bar{Q}_{1,o_1}(D) = 1$ if and only if
 570 o_1 is in the output of $\mathbf{Q1}$, and L_{o_1} the lineage of o_1 . The computation is the
 571 following:

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

572 Where, for the first element of the sum the corresponding B is $\{c2f_1, c_1\}$,
 573 for the second element it is $\{c2f_2, c_2\}$, for the third it is $\{c2f_1, c2f_2, c_1\}$, for
 574 the fourth it is $\{c2f_1, c_1, c_2\}$, for the fifth it is $\{c2f_2, c_2, c_1\}$, for the sixth it is
 575 $\{c2f_1, c2f_2, c_2\}$, and for the seventh $\{c2f_1, c2f_2, c_1, c_2\}$. Every other possible
 576 coalition B would make the factor equal to 0.

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

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

577 As we can see, the sum of the Shapley values of all the tuples in an output
 578 tuple's lineage is always equal to 1 when using a Boolean query as wealth
 579 function.

580 5. Credit Distribution and Distribution Strategies

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

587 5.1. Data Credit and Data Credit Distribution

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

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

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

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

603 **Definition 5.1. Tuple Level Data Credit Distribution (DCD) [27]**
 604 *Given a query Q over I and $k \in \mathbb{R}_{>0}$, DCD is defined by the function $f_{I,Q} :$
 605 $\text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ such that $f_{I,Q}(t, k) = h$ where $0 \leq h \leq k$ and
 606 $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$. The function $f_{I,Q}$ is the distribution strategy (DS).*

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

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

617 5.2. A Lineage-based Distribution Strategy

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

Definition 5.2. Lineage-based Distribution Strategy [27]

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

then t receives credit equal to:

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

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

625 As an example, consider the output tuples of Table 2, and assume that
 626 each output tuple has credit $k_o = 1$. The lineage of the first tuple, o_1 , is
 627 the set $\{f_1, c2f_1, c_1, c2f_2, c_2\}$. Therefore, each tuple in this set receives credit
 628 $1/5$. The other tuples of the database receive zero credit. The lineage of the
 629 second output tuple is $\{f_4, c2f_4, c_1\}$, therefore each of these tuples receives
 630 credit $1/3$.

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

637 Not all of the tuples in the lineage of an output tuple are necessary to be
 638 present at the same time for the output tuple to appear in the query results.
 639 For example, if the database only had the set of tuples $\{f_1, c2f_1, c_1\}$ or the set
 640 $\{f_1, c2f_2, c_2\}$, the existence of o_1 would still be guaranteed. In other words,
 641 while f_1 is always needed for o_1 to appear in the output, only one of the sets
 642 of tuples $\{c2f_1, c_1\}$ and $\{c2f_2, c_2\}$ is required. One could therefore argue that
 643 it would be more fair for f_1 to receive more credit than the other four tuples,
 644 given its role in producing o_1 .

645 This highlights one limitation of the lineage-based DS: while able to find
 646 all and only the relevant tuples of the output, it does not distinguish the
 647 *importance* of tuples in the query computations. We therefore present three
 648 other, more sophisticated, forms of distribution strategies based on why-
 649 provenance, how-provenance, and responsibility.

650 5.3. A Why-Provenance-Based Distribution Strategy

651 The distribution strategy based on why-provenance first equally distributes
 652 the credit k_o among the witnesses of the witness basis for o , and then equally

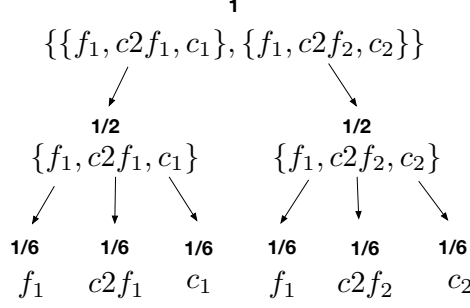


Figure 4: Distribution of credit using why-provenance-based DS for tuple o_1 .

divides the credit of a witness among the tuples in the witness. Since a tuple may appear in more than one witness, it will receive more than one portion of credit from the same distribution. More formally:

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

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

Then tuple t in I receives credit equal to:

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

where γ is a function which returns all witnesses W in which t appears:

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

Figure 4 shows the distribution of credit with why-provenance-based DS for tuple o_1 . The credit is first equally divided between the two witnesses, so that both receive credit $1/2$. The credit is then further divided among the tuples in each witness. Since each witness has three tuples, each tuple in a witness receives $1/6$ of credit. At the end of the distribution, f_1 receives a total credit of $1/3$, and the other tuples receive $1/6$ each. This distribution better reflects the role of f_1 in the generation of o_1 since, as discussed earlier, it is the only mandatory tuple for o_1 to appear in the output; only one of the two other pairs of tuples are necessary for o_1 to appear in the result.

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

\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 7: Notation used in Definition 5.4.

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

Figure 5: Illustration of notation used to define the how-provenance based DS

5.4. A How-Provenance Based Distribution Strategy

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

To define the DS more formally, we introduce some notation and illustrate it using the provenance polynomial \mathcal{H} shown in Figure 5. This notation is also shown in Table 7 for easy reference.

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

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

GMS:
move
nota-
tion
table at
the be-
ginning
and
cover
all the
nota-
tion
of the
section.

id	name
oxs_1	Dopamine Receptors

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

Table 8: 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.

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:

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

689 Going back to the example of Table 4, consider o_1 with provenance poly-
690 nomial $f_1 c2f_1 c_1 + f_1 c2f_2 c_2$. The how-provenance-based DS firstly divides
691 the credit between the two monomials. Since the coefficients of each mono-
692 mial are 1, the credit is split in half. If they were, for example, 1 and 2
693 respectively, 1/3 of the credit would go to the first monomial, and 2/3 to
694 the second. Since in our example each variable has exponent 1, the credit is
695 further divided equally among the three variables. Thus, at the end of the
696 computation, f_1 receives 1/3, and the other tuples receive 1/6.

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

```

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

```

708 The result of Q2 is shown in Table 8, and consists of one tuple, oxs_1 ,
709 annotated with each of the three provenances. As can be seen, lineage and
710 why-provenance are identical to those of the tuple o_1 in the previous example.
711 The how-provenance, however, is different since tuple f_1 is used twice: first

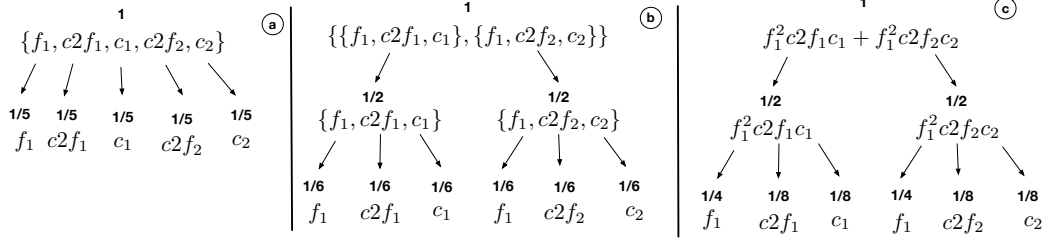


Figure 6: Comparison of different distributions strategies for tuple o_1 produced by query Q2.

in the join of the inner query, and second in the join of the outer query. This information is lost in the first two forms of provenances since they are sets, but it is captured in how-provenance through the use of the operator ‘.’.

Figure 6 shows the differences between the three DS for the tuple o_1 of Table 8. Subfigure 8.a uses lineage, sub-figure 8.b uses why-provenance, and sub-figure 8.c uses how-provenance. The DS based on the provenance polynomial gives credit $1/2$ to f_1 , and $1/8$ to the other tuples. This is reasonable since Q2 relies on f_1 even more than Q1 does. The distribution based on how-provenance rewards f_1 more, showing that how-provenance is even more sensitive to the tuples’ role in a query than why-provenance. This is a direct consequence of the fact that, as proven in [33], how-provenance is more general than why-provenance and lineage, in the sense that it contains more information.

725

726 5.5. Responsibility-based Distribution Strategy

As described in Section 4.3, causality and responsibility is new information that is added to lineage. One possible option for defining a distribution strategy using responsibility is to simply assign the responsibility of each tuple in the lineage of an output tuple as its credit. In this way, responsibility is both a way to compute credit and to distribute it. Using the example of Table 5, in the case of output tuple o_1 , f_1 receives credit 1 and the other tuples receive credit 0.5.

However, we want a DS that is also a function of the input credit value k in order to be comparable with the other three strategies proposed so far. We define a new DS based on responsibility that is a function of the quantity of credit k_o that assigns to each tuple of the lineage a portion of this credit

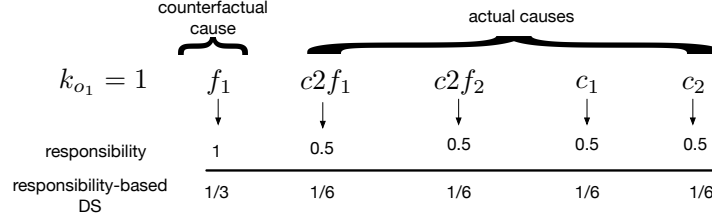


Figure 7: Example of distribution of credit using the responsibility-based DS, assuming $k_o = 1$.

738 weighted by its normalized quantity of responsibility. This will give a bigger
 739 portion of credit to tuples that are higher in the responsibility ranking. For-
 740 mally:

741

Definition 5.5. *Responsibility-based Distribution Strategy*

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

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

742 where ρ_j is the responsibility of tuple j as in Definition 4.1.

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

746 Figure 7 shows the responsibility and credit assigned to the tuples of the
 747 lineage of the output tuple o_1 of Table 5. Assuming that $k_{o_1} = 1$, f_1 receives
 748 credit $1/3$, while the others receive credit $1/6$. As we see, the DS in this
 749 case returns the same distribution as that obtained using why-provenance as
 750 shown in Figure 6. This is not always the case though, as we show in the
 751 example of Section 6.2.

752

753 *5.6. Shapley value-based Distribution Strategy*

754 Similarly to Responsibility, the Shapley value can be seen both as a
 755 method to generate and distribute credit. Moreover, it can be seen that,
 756 using the definition of Shapley value for Boolean queries given in Section 4.3,

757 the sum of the Shapley values of all the tuples of the lineage L of an out-
 758 put tuple o is 1. Thus, the definition of a Shapley value-based distribution
 759 strategy is straightforward:

Definition 5.6. *Shapley Value-Based Distribution Strategy*

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

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

760 Where \bar{Q}_o is the Boolean query such that $\bar{Q}_o(I) = 1$ if and only if o is in the
 761 output of Q on I .

762 As shown in Table 6, tuple f_1 in o_1 's lineage takes credit 7/15 when
 763 $k_{o_1} = 1$, while the other tuples of the lineage take credit 2/15. This DS still
 764 rewards f_1 more than the other tuples, since it is more important than the
 765 other tuples of the lineage. This DS thus behaves differently from all the
 766 other four previous strategies. In particular, f_1 is rewarded more with this
 767 DS than with the others.

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

770 6. Experimental Evaluation

771 To understand the trade-offs between these Distribution Strategies (DSs),
 772 we perform four sets of experiments using queries over target families pre-
 773 sented on the GtoPdb website. The first set of experiments use real queries
 774 extracted from citations to GtoPdb published in the British Journal of Phar-
 775 macology. The second set uses synthetically produced provenance polyno-
 776 mials, corresponding to more complex queries, in order to better highlight
 777 the differences between the DSs. The third set of experiments considers
 778 the accrual of credit over time by the three strategies, again using synthetic
 779 queries. The fourth set of experiments shows how the DSs compare to tradi-
 780 tional citations in giving credit to data curators using both real and synthetic
 781 queries.

782 The source code for the experiments is written in Java and supported by
 783 a PostgreSQL database. For purposes of reproducibility, the code we used
 784 for our experiments and all queries are available here: [https://bitbucket.](https://bitbucket.org/dennis_dosso/credit_distribution_project)
 785 [org/dennis_dosso/credit_distribution_project](https://bitbucket.org/dennis_dosso/credit_distribution_project).

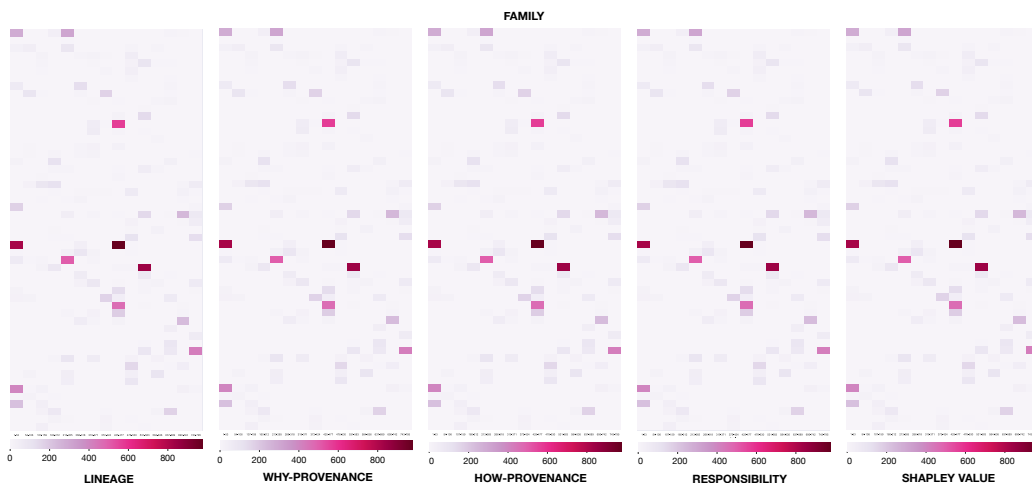


Figure 8: Comparison of four DS on the same table **family** using the distribution given by the queries retrieved from papers. Each cell is a tuple.

6.1. Real-world queries

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

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

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.

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

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

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

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

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

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

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

841 6.2. Synthetic queries

842 To see what happens with more complex queries, we synthetically gener-
843 ated provenance polynomials in which the coefficients and exponents could
844 be greater than one, and picked them at random from a uniform distribution.
845 The queries involve three GtoPdb tables: **family**, **contributor2family**, and
846 **contributor**. The polynomials were generated as follows: first, the number
847 of monomials was decided by randomly choosing a number between one and
848 six. Then, we randomly chose a tuple from the **family** table, one from the
849 **contributor2family** table and one from the **contributor** table; these are
850 the variables of the monomial. We then chose a coefficient for the monomial
851 (between one and three) and an exponent for each tuple (between one and
852 four). For the next monomial, we decided if we wanted to keep the same
853 tuple from the table **family** as first tuple of the new monomial. To do so, we
854 generated a random float number between zero and one. If the number was
855 above 0.2, we changed the family tuple.

856 An example can be found in Figure 9, which shows a sample synthetic
857 provenance polynomial (the how-provenance), the corresponding why-provenance,
858 lineage, and the causality of the tuples of the lineage, together with their re-
859 sponsibility. The resulting credit distribution for each DS is also shown.

860 As an example of how the distribution strategies behave with these syn-
861 thetic queries, consider tuple f_5 in Figure 9. This tuple receives the high-
862 est quantity of credit using responsibility-based distribution and less credit
863 using, in order, lineage, why- and how-provenance. This happens because
864 more information is available about the role of the tuple in the overall com-
865 putation. Generally speaking, the more complex the distribution (e.g., the
866 how-provenance), the more credit is given to tuples that are more frequently
867 used, thus having a higher impact in producing the output tuple. Respon-
868 sibility creates a ranking among lineage’s tuples describing the importance
869 of their role in generating the output. As such, the responsibility-based DS
870 gives more credit to f_1 , f_5 , c_2f_17 and c_18 due to their higher responsibility
871 values. “Importance” is connected to their corresponding minimal contin-
872 gency sets. For example, f_1 has a minimal contingency set (one of the many)
873 $\{f_5\}$, with cardinality 1. On the other hand, c_1 has, as minimal contingency
874 set (one of the many) $\{f_5, c_2\}$, with cardinality two. This means that c_1

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

Credit distribution:

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

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

Credit distribution:

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

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

Credit distribution:

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

Causality: counterfactual causes: \emptyset ,

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

Responsibility:

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

Credit distribution:

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

Shapley value:

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

Figure 9: Sample synthetic provenance polynomial (how-provenance) and corresponding why-provenance, lineage, responsibility, and Shapley values, together with the corresponding credit distributions. In the case of the Shapley value, the value and the distribution of credit are the same.

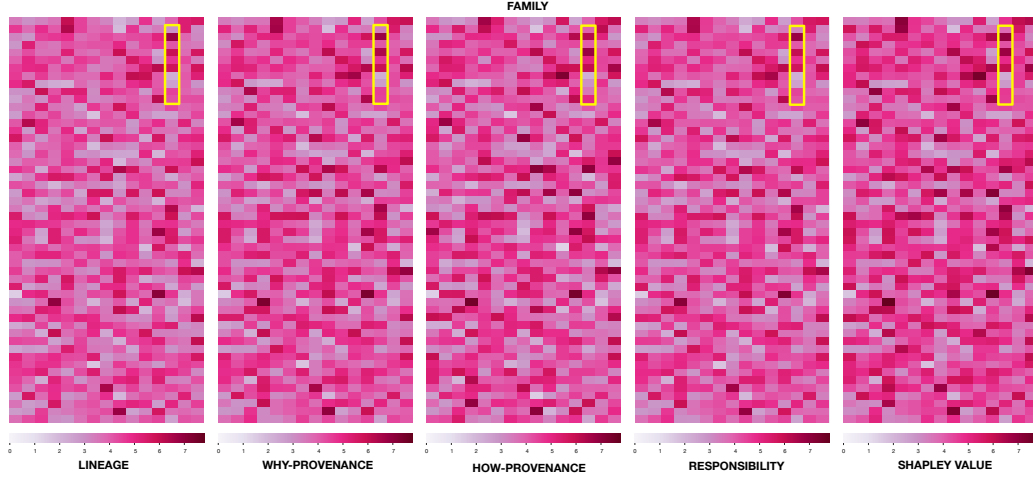


Figure 10: Comparison of three DS on the same table **family** after the distribution computed using 10K synthetic and randomly generated provenance polynomials. The tuples in the blue rectangles are used as example in the discussion connected to Figure 11.

875 is the “least important” amongst the tuples with minimal contingency sets
 876 of lower cardinality, and this is reflected in the different quantities of credit
 877 being distributed.

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

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

887 As can be seen, the four strategies generate different credit distributions,
 888 indicated by the varying hues. However, there is a certain amount of consis-
 889 tency between them in that tuples which are highly rewarded by one strategy
 890 are also highly rewarded by the others. This shows that the four DSs consis-
 891 tently reward certain tuples more than others.

892 Note that lineage-based DS gives the least credit to tuples in the **family**
 893 table, indicated by an overall lighter hue. This is because the DS distributes
 894 credit equally to all tuples appearing in the lineage. Since these queries also

895 use two other tables, credit is distributed to tuples in those tables.

896 Moving to why-provenance-based DS, we see that more credit is given to
897 tuples in the **family** table than with the previous strategy. This is because
898 the DS considers the different ways that a tuple is used, e.g. in joins with
899 other tuples. If the same tuple is present in more than one witness, it will
900 draw more credit and take it from other tuples in the witness basis. In
901 this case, tuples in **family** drew more credit, taking it from tuples in the
902 other two tables, due to the role that **family** tuples played in the queries
903 that were executed. We also notice that the responsibility-based distribution
904 strategy has a distribution that is quite similar to the one provided by why-
905 provenance. It is often the case, for example when the witnesses of the
906 why provenance share one common tuple, that the two distributions behave
907 similarly.

908 We note that the lineage-based DS gives an average credit of 3.92 to each
909 tuple in the table, while the DS based on why-provenance assigns 4.18, how-
910 provenance 4.18, and the one based on responsibility 4.13. Moreover, lineage
911 distributed a total of about 3121 units of credit to the **family** table, why-
912 provenance 3333, how-provenance 3331, while responsibility assigned 3290.

913 Finally, consider the how-provenance-based DS heat-map. As with why-
914 provenance, more credit is typically given to tuples in **family** compared to
915 lineage-based DS, since it recognizes the role of these tuples in the queries,
916 and the overall hue is deeper. The two distributions appear similar, although
917 on closer inspection, slight differences can be seen. This is because how-
918 provenance also considers the frequency with which tuples are used, not only
919 the ways in which they are used. Therefore, although the overall distribution
920 is similar, there are small differences due to the presence of exponents and
921 coefficients in the provenance polynomials, influencing the distribution of
922 credit.

923 To better understand this difference, in the next subsection we consider
924 the accrual of credit over time. In doing so, we will focus on the ten tuples
925 shown within the large yellow rectangles in Figure 11. Each small rectangle
926 within a large blue rectangle is a tuple, and we number them from 1 (top) to
927 10 (bottom). These ten tuples were cherry-picked because they allow us to
928 see the evolution of the distribution of credit through time. There are other
929 tuple sets that could have been selected driving us to the same considerations.

6.3. Credit accrual over time

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

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

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

By the end of 1K queries, credit differentials between tuples as well as between strategies can be seen. For example, tuple 3 is usually rewarded the most credit by all four strategies. However, it receives the highest quantity of credit from the why-provenance- and responsibility-based strategy (1.33). Moreover, it can be seen that tuples 1 and 7 receive a higher quantity of credit when how-provenance is adopted, showing how this form of provenance behaves differently from the others in this context. Moving to 2K queries, it is possible to see that tuple 3 and 7 are still the most rewarded by the strategies. This trend continues to the end of 2K queries.

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

While the relative value of credit “positions” of tuples within a DS strategy depends on what queries are being executed, the important thing to

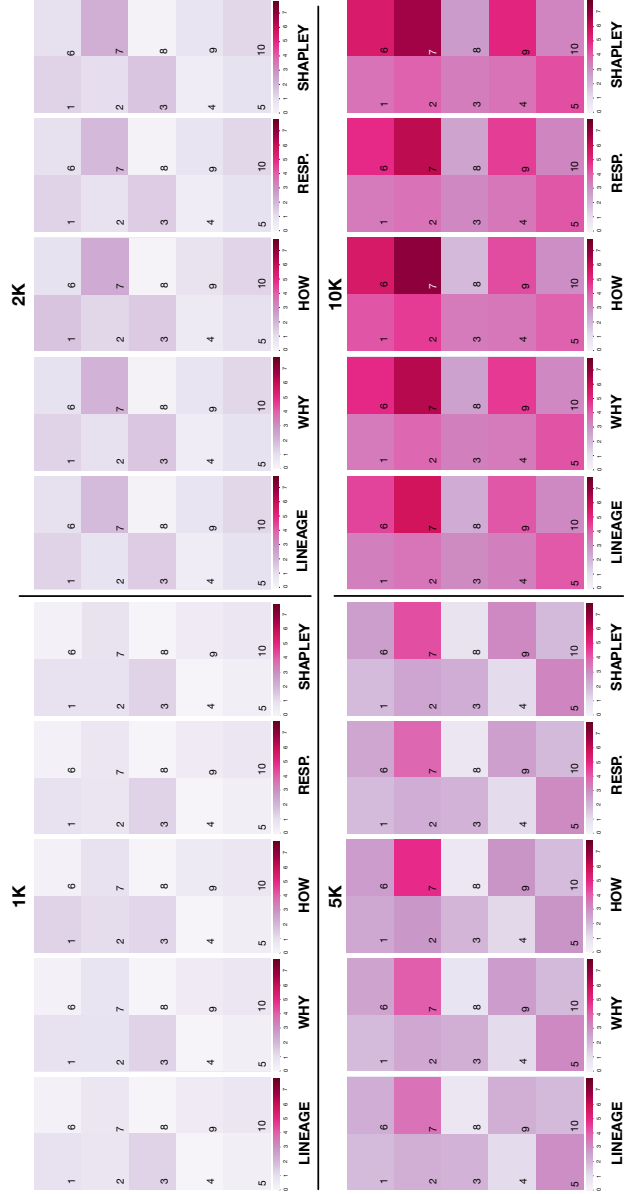


Figure 11: Comparison of the distribution of credit performed by the **five** DSs on a subset of 10 tuples taken from the **family** table, simulating the passing of time. The number at the top of each group of heat-maps represents the number of polynomials whose credit has been distributed.

notice is the difference between the DSs over time: overall, lineage gives less credit to tuples in the `family` table than the other two strategies since credit is shared with tuples in other tables. However, the why-, **responsibility**- and how-provenance-based strategies recognize the more important role being played by the `family` tuples than those in the other tables. The differences between why- and responsibility-based DS are, for the most times, negligible. The differences between the why- and how-provenance-based DSs are also relatively minor in most cases. However, there are certain situations in which the role of a tuple is particularly critical in a query, and in this case the difference in the value of credit assigned is notably higher for how-provenance, as we saw with tuple 7 in the example of Figure 11.

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

6.4. *Credit vs Citations*

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

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

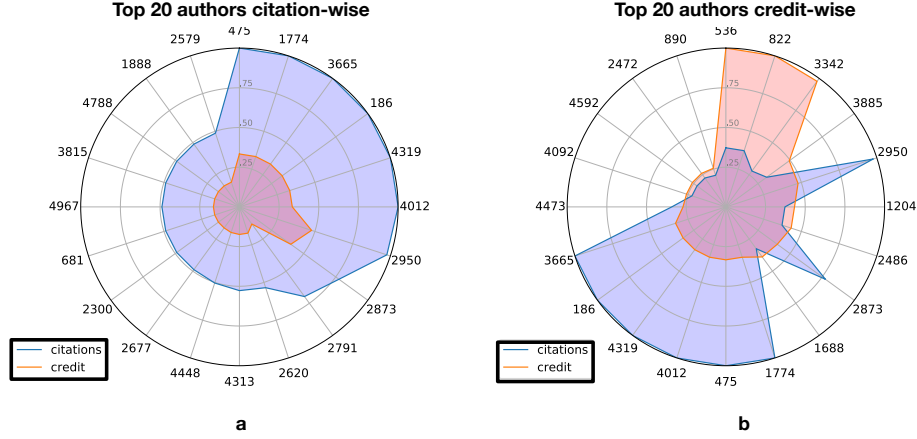


Figure 12: Radars presenting the top 20 authors citation-wise and credit wise, together with their (normalized between 0 and 1) values of citations and credit.

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

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

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

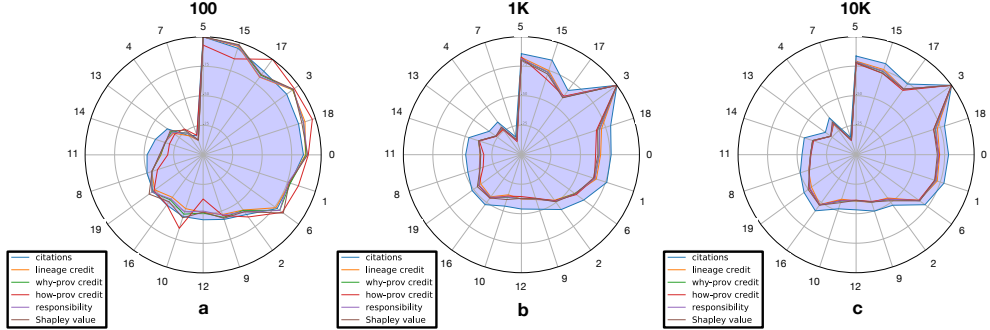


Figure 13: Radars presenting the 20 synthetic authors with corresponding citation and quantities of credit distributed through the 4 DSs (all values normalized between 0 and 1) through different numbers of polynomials (respectively, 100, 1K and 10K). The order is the one defined by figure a, i.e. descending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

That is, the data underlying certain webpages is more “valuable” in terms of credit than a citation to the webpage.

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

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

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

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

Figure 13 shows three radar plots, one for the results obtained with 100 polynomials, one with 1K polynomials, one with 10K polynomials. Each plot shows the top 20 authors in terms of citations (hence the authors and clockwise ordering is the same in each of the plots), and additionally shows the the normalized values of citation (blue line), lineage-based credit (yellow line), why-provenance-based credit (green line), how-provenance-based credit (red line), and responsibility-based credit (violet line). As can be seen, given the synthetic nature of the queries, the correlation between the number of citations and the quantity of credit assigned to the authors appears to be a much stronger than with the real-world queries of Figure 12. In fact, for Figure 13.a the linear correlation between the citation number and all four types of credit is always above 0.94 with p values in the order of $3e-8$. The credit distributed via lineage is closest to the number of citations (a linear correlation of 0.99, p value of $2e-16$ in Figure 13.a), while the other three types of credit behave slightly differently (a linear correlation of around 0.95 in all other three cases in Figure 13.a). Similar observations can be made for Figure 13.b and 13.c.

What these figures show is that, in certain cases, authors who do not have a large number of citations receive more credit than others, as for example authors 17 and 10 in Figure 13.a, and especially when credit is distributed using how-provenance. This again shows how credit gives a different perspective on the role of data and authors by going beyond the limitations of traditional citations.

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

7. Discussion

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

7.1. Credit Generation

In this paper we focused on Credit Distribution, the problem of distributing credit generated by a citation to the parts of the database referenced by

1084 the query. A different problem is Credit Generation, the task of generating
 1085 credit which is then distributed. Credit Generation presents a series of issues
 1086 which are shared by traditional citation practices. For instance, defining the
 1087 quantity of credit to be generated for a given citation is still an open prob-
 1088 lem. Different types of citations may generate different quantities of credit.
 1089 * **SBD: Not sure what the next sentence means: "related to the**
 1090 **results"? I get "previous work".** * Data cited as related to the results
 1091 or as useful for previous work may generate less credit than other data ex-
 1092 tensively used to produce the results presented in a paper. The computation
 1093 of credit could be done manually (although we must consider the complex-
 1094 ity of the task, human biases and the resources required to carry it out) or
 1095 automatically, but it must be based on a shared definition of impact which
 1096 is still not agreed upon for data or for traditional citation. For this reason,
 1097 we used a uniform credit assignment.

1098 There is also the problem of *transitive credit distribution*, i.e., how to
 1099 transitively propagate credit from one cited unit to another unit that was
 1100 used to produce the one being cited. For this, a graph of cited units that
 1101 propagate credit between the units depending on influence could be used.
 1102 How to propagate credit is an open and non-trivial problem that needs to
 1103 consider the importance and impact of a citation in a work, be it a paper or
 1104 data, and how to eventually compute the quantity of credit to be propagated.

1105 * **SBD: Revised below, make sure you are ok with it.** *

1106 Finally, in our experiments we assumed that the credit carried by an
 1107 output tuple is one. Thus, each tuple in the output has equal importance.
 1108 As described above, this assumption may be revised and different credit to
 1109 different output tuples could be assigned. Note that from the distribution
 1110 model viewpoint no change is required since the DCD is defined for a generic
 1111 value k .
 1112

1113 7.2. Choice of Distribution Strategies

1114 In this paper we presented four different DSs, so the natural question is
 1115 which one to use. This depends on the task at hand. When we want to
 1116 highlight the tuples being used in the database by a workload, the lineage-
 1117 based DS may be sufficient. When we also want to know the relative impact
 1118 of tuples in the context of the query, the other DSs should be used since they
 1119 give a better understanding of the importance of data.

1120 In the real-world based experiments, the four DSs behaved the same,
 1121 which was due to the specific nature of the data and the queries being used.

GMS:
 I com-
 mented
 out a list
 of items
 I did not
 agree
 much
 with.
 Please
 take a
 look and
 resume
 what you
 think
 is good
 for you
 (if any).
 – DD:
 Some
 of the
 points
 were
 raised
 because
 of Re-
 viewer 2.
 Since we
 discuss
 about
 them in
 the first
 para-
 graph,
 we can
 keep the
 list out
 of the
 paper.

However, the why-provenance of a query will differ from the lineage of the same query whenever the output tuples can be computed in more than one way by the query, i.e., if there is more than one witness. This is usually true when join and projection operators are used in the query.

*** SBD: this paragraph went all over the place and was confusing. I eliminated a lot of details that you might think are important (original is in an eat environment. ***

To address the question of what types of queries are likely to extract cited data, we turn to the results of published studies on the characteristics of query workloads and the complexity of their queries [39, 54, 59]. These studies show that operations such as inner-/outer-joins and projections occur in a significant number of queries. Therefore why- and how-provenances may become quite complex in certain cases and provide a distribution of credit that is significantly different from the one obtained with lineage.

*** Is there more to say here? What are the general queries for which responsibility is hard to compute, and can the various provenances handle them at all? I know that provenance semi-rings has been extended to SPJU and aggregate queries, so imagine this means the others can be extended since it is a general framework. *** From a complexity standpoint, all four DS are similar since we focused on SPJ queries. However, responsibility is hard to compute for general queries. In terms of implementation, lineage is the simplest to compute since it only cares about a tuple being used, while the other provenances also need additional information to be taken into consideration.

Another promising DS that could be developed is one based on Shapley values. This function has been widely used in knowledge representation and machine learning, and has strong theoretical justifications. However, its use in databases as a metric for quantifying the influence of a tuple on the output of a query (thereby presenting an alternative to responsibility) has only recently been proposed [44]. Furthermore, the initial theoretical analysis in [44] showed mainly lower bounds on the complexity of the problem, and did not suggest a feasible implementation. However, very recently, an efficient implementation for boolean queries (queries that output true or false) has been provided [25], both in terms of an exact computation (which in practice works well for most queries) and an inexact one (which is extremely fast and provides the same ranking of tuples as the exact computation, but not necessarily the same values). In future work, we will explore a Shapley-based DS and test its performance in Credit Distribution.

8. Conclusions and Future Work

This paper defines three new distribution strategies based on why-provenance, how-provenance, and responsibility, and compares them against the lineage-based distribution strategy defined in [27]. The first, why-provenance-based DS, uses the concept of a witness, and gives more credit to tuples that appear in more than one witness. In this way, tuples that are more important to the query and are used in different ways are rewarded more. The second, how-provenance-based DS, considers the frequency with which a tuple or combination of tuples is used in the query through the information contained in a provenance polynomial. In this case, the how-provenance-based DS is more sensitive than the why-provenance-based DS to the role and importance of tuples. The third DS exploits the notion of responsibility, a real value which ranks the lineage tuples based on their degree of causality in generating the output. The responsibility-based DS was shown to behave similarly to the why-provenance based DS.

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

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

1197 In the third set of experiments, we showed how the differences between
1198 the DS are compounded over time, i.e. when more and more queries are
1199 processed by the system.

1200 In the fourth set of experiments we compared traditional citations to
1201 authors to the credit accrued to them via the DSs. We showed how, in
1202 both real-world and synthetic scenarios, credit rewards authors who con-
1203 tribute/curate data that has the highest impact, and therefore receives the
1204 biggest quantity of credit, and not necessarily the data with the highest ci-
1205 tation count. In this sense, credit appears to be an useful new measure to
1206 discover data and their corresponding curators that have a high impact in
1207 the research world, even when they are cited few times or do not appear at
1208 all in the data that are cited (i.e. the case of data used to build the output
1209 of a query but that is not visualized in the output itself).

1210 In future work, we plan to explore different strategies to generate and
1211 distribute credit. In this paper we assumed that each output tuple carries
1212 credit 1. In more sophisticated scenarios we can employ different strategies
1213 to compute credit, that reflect the importance of cited data. Other, more
1214 sophisticated, strategies could also be used to decide how credit is distributed
1215 between the authors, beyond the uniform distribution used here, in a way to
1216 reflect the work performed by them on the cited data. *We also plan to explore*
1217 *a DS based on Shapley values [44], exploiting the practical implementation*
1218 *recently suggested in [25].*

1219 There are also a number of other intriguing applications for credit over
1220 relational databases. One such application is *data pricing*, which gives a
1221 price to a query submitted by a user who wants to buy the produced in-
1222 formation. Currently, a common strategy used for data pricing is based on
1223 query rewriting: A database stores a set of views with their price. When a
1224 new query arrives, the system rewrites it using the stored views to obtain a
1225 query price, a process that can be computationally expensive. We plan to
1226 distribute credit through carefully planned and representative queries, and
1227 use credit information to define a new, faster, and potentially more flexible
1228 pricing function.

1229 Another application is *data reduction* [48], which addresses the problem of
1230 reducing the vast – and rapidly expanding – amount of data that is being pro-
1231 duced. Data credit can help address this problem by identifying “hotspots”
1232 and “coldspots” of data. A hotspot is data in a database (e.g. a tuple) with
1233 a high quantity of credit, which is therefore valuable for the set of queries
1234 that execute frequently over the data and distribute the credit. A coldspot is

1235 data with a low quantity of credit which can therefore be considered as less
1236 important, and could be deleted, summarized, or moved to cheaper and/or
1237 less efficient memory.

1238 Acknowledgement

1239 The work was partially supported by the ExaMode project, as part of the
1240 European Union H2020 program under Grant Agreement no. 825292.

1241 References

- 1242 [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bern-
1243 stein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L.,
1244 Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Koss-
1245 mann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo,
1246 T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo,
1247 A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suciu, D.
1248 (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–
1249 53.
- 1250 [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating
1251 data citation in citedb. *PVLDB*, 10(12):1881–1884.
- 1252 [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y.
1253 (2018). Data citation: A new provenance challenge. *IEEE Data Eng.*
1254 *Bull.*, 41(1):27–38.
- 1255 [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An
1256 Introduction to the Joint Principles for Data Citation. *Bulletin of the*
1257 *Association for Information Science and Technology*, 41(3):43–45.
- 1258 [5] Baggerly, K. (2010). Disclose all data in publications. *Nature*,
1259 467(7314):401–401.
- 1260 [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D.,
1261 Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I.,
1262 Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and
1263 Goble, C. A. (2013). Why linked data is not enough for scientists. *Future*
1264 *Gener. Comput. Syst.*, 29(2):599–611.

- 1265 [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation
1266 Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.
- 1267 [8] Bertin-Mahieux, T., Ellis, D., Whitman, B., and Lamere, P. (2011). The
1268 million song dataset. In *Proceedings of the 12th International Conference*
1269 *on Music Information Retrieval (ISMIR 2011)*, pages 591–596.
- 1270 [9] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In
1271 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*
1272 *nication*, pages 93–116. De Gruyter Mouton.
- 1273 [10] Buneman, P. (2006). How to cite curated databases and how to make
1274 them citable. In *18th International Conference on Scientific and Statistical*
1275 *Database Management, SSDBM*, pages 195–203. IEEE Computer Society.
- 1276 [11] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D.,
1277 Pawson, A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn’t
1278 working, and what to do about it. *Database J. Biol. Databases Curation*,
1279 2020.
- 1280 [12] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation
1281 is a computational problem. *Commun. ACM*, 59(9):50–57.
- 1282 [13] Buneman, P., Khanna, S., and Tan, W. C. (2001). Why and where: A
1283 characterization of data provenance. In *Database Theory - ICDT 2001,*
1284 *8th International Conference*, pages 316–330.
- 1285 [14] Buneman, P. and Silvello, G. (2010). A rule-based citation system for
1286 structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.
- 1287 [15] Callaghan, S., Donegan, S., Pepler, S., Thorley, M., Cunningham, N.,
1288 Kirsch, P., Ault, L., Bell, P., Bowie, R., Leadbetter, A. M., Lowry,
1289 R. K., Moncoiffé, G., Harrison, K., Smith-Haddon, B., Weatherby, a.,
1290 and Wright, D. (2012). Making Data a First Class Scientific Output:
1291 Data Citation and Publication by NERC’s Environmental Data Centres.
1292 *International Journal of Digital Curation*, 7(1):107–113.
- 1293 [16] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Jour-
1294 nals: A Survey. *Journal of the Association for Information Science and*
1295 *Technology*, 66(9):1747–1762.

- 1296 [17] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases:
1297 Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–
1298 474.
- 1299 [18] Chockler, H. and Halpern, J. Y. (2004). Responsibility and blame: A
1300 structural-model approach. *J. Artif. Intell. Res.*, 22:93–115.
- 1301 [19] CODATA-ICSTI Task Group on Data Citation Standards and Practices
1302 (2013). *Out of Cite, Out of Mind: The Current State of Practice, Policy,*
1303 *and Technology for the Citation of Data*, volume 12.
- 1304 [20] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N.
1305 (2019). Bringing citations and usage metrics together to make data count.
1306 *Data Science Journal*, 18(1).
- 1307 [21] Cronin, B. (1984). *The Citation Process. The Role and Significance of*
1308 *Citations in Scientific Communication*. London: Taylor Graham.
- 1309 [22] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evi-
1310 dence of a structural shift in scholarly communication practices? *JASIST*,
1311 52(7):558–569.
- 1312 [23] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of
1313 view data in a warehousing environment. *ACM Trans. Database Syst.*,
1314 25(2):179–227.
- 1315 [24] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model
1316 for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on*
1317 *Innovative Data Systems Research*. www.cidrdb.org.
- 1318 [25] Deutch, D., Frost, N., Kimelfeld, B., and Monet, M. (2022a). Com-
1319 puting the Shapley Value of Facts in Query Answering. In Bonifati, A.
1320 and Abbadi, A. E., editors, *SIGMOD ’22: International Conference on*
1321 *Management of Data, Philadelphia, June 12-17, 2022*. ACM.
- 1322 [26] Deutch, D., Frost, N., Kimelfeld, B., and Monet, M. (2022b). Computing
1323 the shapley value of facts in query answering. *the 2022 Special Interest*
1324 *Group on Management of Data conference (SIGMOD)*. in print.
- 1325 [27] Dosso, D. and Silvello, G. (2020). Data credit distribution: A
1326 new method to estimate databases impact. *Journal of Informetrics*,
1327 14(4):101080.

- [28] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The virtual atomic and molecular data centre (VAMDC) consortium. *Journal of Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- [29] Eiter, T. and Lukasiewicz, T. (2002). Complexity results for structure-based causality. *Artif. Intell.*, 142(1):53–89.
- [30] Fang, H. (2018). A discussion of citations from the perspective of the contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–1520.
- [31] Garfield, E. (1999). Journal impact factor: a brief review. *Can. Med. Assoc.*, 979-980.
- [32] Gößwein, B., Miksa, T., Rauber, A., and Wagner, W. (2019). Data identification and process monitoring for reproducible earth observation research. In *2019 15th International Conference on eScience (eScience)*, pages 28–38. IEEE.
- [33] Green, T. J., Karvounarakis, G., and Tannen, V. (2007). Provenance semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 31–40. ACM.
- [34] Halpern, J. Y. and Pearl, J. (2013). Causes and explanations: A structural-model approach — part 1: Causes. *CoRR*, abs/1301.2275.
- [35] Harding, S. D., Sharman, J. L., Faccenda, E., Southan, C., Pawson, A. J., Ireland, S., Gray, A. J. G., Bruce, L., Alexander, S. P. H., Anderton, S., Bryant, C., Davenport, A. P., Doerig, C., Fabbro, D., Levi-Schaffer, F., Spedding, M., Davies, J. A., and Nc-Iuphar (2018). The IUPHAR/BPS guide to PHARMACOLOGY in 2018: updates and expansion to encompass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Research*, 46(Database-Issue):D1091–D1106.
- [36] Hartley, J. (2017). Authors and their citations: a point of view. *Scien-tometrics*, 110(2):1081–1084.
- [37] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a transformed scientific method.

- [38] Honor, L. B., Haselgrove, C., Frazier, J. A., and Kennedy, D. N. (2016). Data citation in neuroimaging: proposed best practices for data identification and attribution. *Frontiers in neuroinformatics*, 10:34.
- [39] Jain, S., Moritz, D., Halperin, D., Howe, B., and Lazowska, E. (2016). Sqlshare: Results from a multi-year sql-as-a-service experiment. In *Proceedings of the 2016 International Conference on Management of Data*, pages 281–293.
- [40] Joshi-Tope, G., Gillespie, M., Vastrik, I., D’Eustachio, P., Schmidt, E., de Bono, B., Jassal, B., Gopinath, G. R., Wu, G. R., Matthews, L., Lewis, S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of biological pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- [41] Katz, D. (2014). Transitive credit as a means to address social and technological concerns stemming from citation and attribution of digital products. *Journal of Open Research Software*, 2(1).
- [42] Kosten, J. (2016). A classification of the use of research indicators. *Scientometrics*, 108(1):457–464.
- [43] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S. (2011). Citation and Peer Review of Data: Moving Towards Formal Data Publication. *International Journal of Digital Curation*, 6(2):4–37.
- [44] Livshits, E., Bertossi, L. E., Kimelfeld, B., and Sebag, M. (2020). The shapley value of tuples in query answering. In Lutz, C. and Jung, J. C., editors, *23rd International Conference on Database Theory, ICDT 2020, March 30–April 2, 2020, Copenhagen, Denmark*, volume 155 of *LIPICs*, pages 20:1–20:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.
- [45] Martone, M. (2014). Joint declaration of data citation principles. *FORCE11. San Diego CA. Data Citation Synthesis Group*. <https://www.force11.org/datacitationprinciples>, online September 2020.
- [46] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation counts and rankings of LIS faculty: Web of science versus scopus and google scholar. *Journal of the american society for information science and technology*, 58(13):2105–2125.

- [47] Meliou, A., Gatterbauer, W., Moore, K. F., and Suciu, D. (2010). The complexity of causality and responsibility for query answers and non-answers. *Proc. VLDB Endow.*, 4(1):34–45.
- [48] Milo, T. (2019). Getting rid of data. *Journal of Data and Information Quality (JDIQ)*, 12(1):1–7.
- [49] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D., Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G., Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff, D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D., Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson, E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C., Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers, E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research culture. *Science*, 348(6242):1422–1425.
- [50] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J. (2016). Research data explored: An extended analysis of citations and altmetrics. *Scientometrics*, 107(2):723–744.
- [51] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic, large databases: Model and reference implementation. In *Proceedings of the 2013 IEEE International Conference on Big Data, 6-9 October 2013, Santa Clara, CA, USA*, pages 307–312.
- [52] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identification of Reproducible Subsets for Data Citation, Sharing and Re-Use. *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue on Data Citation*, 12(1):6–15.
- [53] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data citation of evolving data: Recommendations of the working group on data citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- [54] Remil, Y., Bendimerad, A., Mathonat, R., Chaleat, P., and Kaytoue, M. (2021). ” what makes my queries slow?”: Subgroup discovery for sql workload analysis. *arXiv preprint arXiv:2108.03906*.

- 1420 [55] Shapley, L. S. (1954). A value for n-person games. In Kuhn, H. W. and
1421 Tucker, A. W., editors, *Contributions to the Theory of Games II*, pages
1422 307–317. Princeton University Press, Princeton.
- 1423 [56] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*
1424 *Sci. Technol.*, 69(1):6–20.
- 1425 [57] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data
1426 provenance in e-science. *SIGMOD Record*, 34(3):31–36.
- 1427 [58] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-
1428 spective. In of Sciences’ Board on Research Data, N. A. and Information,
1429 editors, *Report from Developing Data Attribution and Citation Practices*
1430 *and Standards: An International Symposium and Workshop*, pages 177–
1431 178. National Academies Press: Washington DC.
- 1432 [59] Vogelsgesang, A., Haubenschild, M., Finis, J., Kemper, A., Leis, V.,
1433 Mühlbauer, T., Neumann, T., and Then, M. (2018). Get real: How bench-
1434 marks fail to represent the real world. In *Proceedings of the Workshop on*
1435 *Testing Database Systems*, pages 1–6.
- 1436 [60] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,
1437 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,
1438 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data
1439 management and stewardship. *Scientific data*, 3.
- 1440 [61] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data
1441 citation: Giving credit where credit is due. In *Proceedings of the 2018*
1442 *International Conference on Management of Data, SIGMOD*, pages 99–
1443 114.
- 1444 [62] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to
1445 scientific datasets using article citation networks. *Journal of Informetrics*,
1446 14(2).
- 1447 [63] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of
1448 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.
- 1449 [64] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for
1450 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*
1451 *Molecular Spectroscopy*, 327:122–137.