

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, United States*

Abstract

In the current world of research data is a fundamental method to disseminate scientific knowledge, to determine scholarship, and to provide credit and recognition to the authors of research endeavors. However, issues like data citation, handling and counting the credit generated by such citations are still open research questions.

In this context, data credit has recently emerged as a new measure of value, defined and built on top of the data citation theory. Data credit is a real value that represents the importance of data cited by a paper, or by another research entity. As such, credit can be used to annotate data contained in curated scientific databases, and it can be considered as a measure for their importance and impact in the research world. As such, it is a new method that, together with traditional citations, helps to recognize the value of data and its creators in a world more and more dependent on data.

In this paper we explore the problem of Data Credit Distribution, the process by which credit is divided and assigned to the data in a database that are 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 curated and well-known scientific relational database. We define two new distribution strategies, functions that perform this task, based on two form of data provenance, why-provenance, and how-provenance.

Using different distribution strategies, we show how credit can highlight areas of a database that are frequently used, and how it can work as a new bibliometric measure for data and their corresponding curators. Credit in particular rewards data and authors based on their research impact, and not

merely on the number of citations. Also, we show how different distribution strategies, based on different types of data provenance, can be more sensible to the role of an input tuple in the generation of the output, and thus rewarding it differently.

Keywords: Data Citation, Data Credit

1 Introduction

Citations are an essential component of scientific research, enabling research products to be found as well as the relationships between research products to be understood. They form a basis on which to give credit to authors, papers, and venues [55, 19, 20]. Citations are used, among other things, to decide on tenure, promotion, hiring, and funding of grants for researchers [41, 21, 32, 38].

Nowadays, science and research are increasingly digital. 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 [39, 15], and that data citations should contribute to the scientific reputation of researchers, scientists, data curators, and creators [4, 50]. It is also accepted that data citations should be counted alongside of traditional citations, and contribute to bibliometrics indicators [7, 44].

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

Recently, the concepts of *data credit* and *Data Credit Distribution* (DCD) [26, 36, 54] have emerged, built on top of methodologies for data citation. Data

credit is a value that is computed based on the importance of the data being cited in a paper, and represents the impact of the data on the citing paper. The Data Credit Distribution problem consists of distributing this credit to elements in the databases in the citation graph that are responsible for the generation of the data being cited. The goal of DCD is to improve and expand the reach of data citation, rather than being an alternative to it. This means that to employ DCD techniques, we need data citations in some form.

[37] defined credit as a “quantity” that describes the importance of a research entity, such as papers or data mentioned in a citation, and proposed the idea of a *distribution* of credit from research entities, such as papers or data, to other research entities through citations. This can be done by exploiting the structure of the *citation graph*, a directed graph whose nodes are publications and edges are citations. This graph is the model at the core of systems such as Google Scholar and the Web of Science. Zeng et al. [54] and Fang [26] further explored this concept by defining frameworks for the computation and distribution of credit between papers, authors, and data used by papers in the citation graph.

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

The DCD problem differs from the traditional citation setting since:

1. In a traditional setting, when a paper cites another paper, a +1 “credit” is given to the cited paper (and to its authors). It does not matter why or how paper p_1 cites paper p_2 ¹, the result is always +1 from p_1 to p_2 and thus a +1 to the citation count of the authors of p_2 . With a different credit distribution strategy, the “value” given to the cited entity can be *proportional* to the role played in the citing entity. Hence, we can weigh the importance of the cited entities and assign credit according to their role.

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



Figure 1: Overview of the credit distribution pipeline.

2. Traditional citations are considered to be *atomic*. A citation from p_1 to p_2 can never be broken into pieces and assigned in part to p_2 and in part to other papers or data that contributed to p_2 . This is due to the intrinsic difficulty in grasping the role and “weight” of the other papers and data, and in automating the credit assignment process. In contrast, we consider data credit to be a *non-atomic* real value, which can be divided and distributed to multiple components of a database.
3. Credit can be *transitive*, that is, it can be propagated through one cited entity to other entities cited by it that contributed to its content.

We study the DCD problem in the context of relational databases (RDBs) since they are widely used² and are the main focus of current work in data citation methods [14, 12, 45]. RDBs are also frequently a test-bed for new methods that can be adapted to other databases, e.g., graphs or document databases. Furthermore, the “portions” of data in an RDB that can be credited can be defined at different levels of granularity, in particular: (i) the whole database, (ii) tables, and (iii) tuples.

The DCD process is summarized in Figure 1:

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

- 82 **Step 1** Scientists and experts contribute the curated information contained
83 in a scientific database. These are called the “Data Curators”.
- 84 **Step 2** Other researchers use the data in their research, and when possible,
85 cite them.
- 86 **Step 3** The citation to the data generates credit, that can be used as a
87 proxy for the impact of the data on the citing paper. This credit is
88 represented as a real value $k \in \mathbb{R}_{>0}$.
- 89 **Step 4** Given the database instance I and the query Q , it is possible to
90 compute the *data provenance* of $Q(I)$. The provenance of $Q(I)$ is a
91 form of metadata that describes the generation process undertaken by
92 Q , and the data used in I to generate the output [17]. Many different
93 notions of provenance have been proposed in the literature for data in
94 database management systems [22, 13, 30], describing different kinds
95 of relationships between data in the input and the output of a query.
96 As reported in [17], these provenances have been used in several appli-
97 cations beyond giving information on how queries work, for example,
98 annotation propagation and the view update problem. In this paper,
99 we consider three types of provenance: lineage, why-provenance, and
100 how-provenance.
- 101 **Step 5** Provenance is input to the CDC problem, whose aim is to compute
102 the *Credit Distribution Strategy* (CDS, also referred only as Distribu-
103 tion Strategy, DS). The CDS is a function that distributes k to the data
104 in the input database I , and is defined on the basis of citation policies
105 decided at the database administration level or at the domain commu-
106 nity level. In this paper, since we base CDS on data provenance, we
107 describe three CDS, each one based on a different form of provenance.
- 108 **Step 6** Once the CDS is computed, it is used to distribute the given credit
109 k to the parts of the database that are responsible for the generation
110 of $Q(I)$. Transitively, this credit is also divided and given to the corre-
111 sponding authors of those data.

112 This paper expands our recent work in [24], which addressed the problem
113 of how to reward data and data curators who are typically overlooked in
114 current citation systems. In that work, we first defined the problem of DCD

115 in relational databases, and proposed a viable Distribution Strategy (DS)
 116 based on *lineage*, which is the simplest form of *data provenance*. The lineage
 117 of a tuple t in the output $Q(I)$ is defined as the set of all and only the tuples
 118 in the database instance I that are “relevant” to the production of t , that
 119 is the tuple that are used by Q in the production of t . The lineage-based
 120 strategy equally redistributes the credit k to the tuples in the lineage set,
 121 thus each tuple receives credit $k/|L_t|$, where L_t is the lineage set of t .

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

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

139 For evaluation, we use a well-known curated database, the IUPHAR/BPS³
 140 Guide to Pharmacology [31], also known as GtoPdb⁴, which contains ex-
 141 pertly curated information about diseases, drugs, cellular drug targets, and
 142 their mechanisms of action. We chose GtoPdb for two main reasons: (i) it
 143 is a widely-used and valuable curated relational database, (ii) many papers
 144 in the literature use, and cite its data (i.e., families, ligands, and receptors).
 145 Real queries used in papers can therefore be seen as data citations which, in
 146 turn, can be used to assign data credit.

147 We perform three sets of experiments. In the first one, real queries are ex-

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

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

148 tracted from papers published in the British Journal of Pharmacology (BJP),
149 that represent data citations to GtoPdb, and are used to distribute credit
150 in the database using the three different provenance-based DSs. In the sec-
151 ond and third experiment we analyse the behaviour of the different DS when
152 complex citation queries are employed.

153 **Contributions.** Contributions of this work include:

- 154 • The definition of new distribution strategies for the problem of Data
155 Credit Distribution, based on why-provenance and how-provenance;
- 156 • An in-depth analysis of the effects of credit distribution on real-world
157 curated data and of the differences between the three proposed Distri-
158 bution Strategies.

159 **Outline.** The rest of the paper is organized as follows: Section 2 presents the
160 background and related work. Section 3 describes the use case we adopted.
161 Section 4 briefly presents the forms of provenance used in the paper. Section
162 5 describes the problem of DCD and the proposed DS. In Section 6 we present
163 the experimental evaluation. Finally, Section 7 draws some conclusions and
164 outlines future work.

165 2. Background

166 *Data in Research.* As described by Jim Gray in his last talk [33], the world of
167 research is rapidly transitioning towards the *fourth paradigm of science*, that
168 is, data-intensive scientific discovery, where data are important for scientific
169 advances as well as for traditional publications [6].

170 The scientific community is promoting an *open research culture* [43],
171 founded on methods and tools to share, discover, and access experimental
172 data. The community has identified the FAIR principles (Findable, Acces-
173 sible, Interoperable, and Reusable) [52], that should be enforced by every
174 database. In particular, data should be accessible from the articles, journals,
175 and papers that cite or use them [19]. Aspects such as the need for the *repro-*
176 *ducibility* of experiments through the used data; the *availability* of scientific
177 data; the *connections* between data and the scientific results are all needed
178 aspects for the fourth paradigm, and are all relevant to the domain of *data*
179 *citation* [34].

180 *Data Citation: Principles and Motivations.* Data Citation principles were
 181 first described in detail in [18], and later summarized and endorsed by the
 182 Joint Declaration of Data Citation Principles (JDDCP) [40]. The principles
 183 are divided into two groups [48]. The first one contains principles concerning
 184 the role of data citation in scholarly and research activities such as the (i)
 185 *importance* of data (why data citation is important and why data should be
 186 considered as first-class citizens); (ii) *credit* and *attribution* to the creators
 187 and curators of the data; (iii) *evidence*; (iv) *verifiability*; and *interoperability*,
 188 with these last three requiring data citation methods to be flexible enough to
 189 operate through different communities. The second group defines the main
 190 guidelines to establish a data citation systems, and contains principles such
 191 as the (i) *unique identification* of the data being cited; (ii) (*open*) *access* to
 192 data; (iii) guarantee of *persistence* and *availability* of citations even after the
 193 lifespan of the cited entity; the (iv) *specificity* of a citation, i.e. it must lead
 194 to the data set originally cited.

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

- 196 • *Data attribution*: identify the individuals that should be credited for
 197 data with variable granularity.
- 198 • *Data connection*: connect papers to the data being used.
- 199 • *Data Discovery*: citations helps to find data records and subsets that
 200 would be otherwise not findable via search engines.
- 201 • *Data Sharing*: share data obtained by researchers within the whole
 202 community.
- 203 • *Data Impact*: highlight the results obtained in writing papers using
 204 specific data, the frequency and modality data were used.
- 205 • *Reproducibility*: data citation greatly impacts the reproducibility of
 206 science [5]. Many authoritative journals ask to share data and provide
 207 valid methodologies to reproduce experiments.

208 2.1. Data Citation in Relational Databases

209 In this paper, we develop our methods and experiments on relational
 210 databases. RDBs have been the main target of data citation methods since
 211 the surge of the data-centric research paradigm. The RDA “Working Group

212 on Data Citation: Making Dynamic Data Citable”⁵ [46] has been working in
213 the last years on large, dynamic, and changing datasets. The working group
214 has finished the development of its guidelines and has now moved on into an
215 adoption phase. The datasets considered by the WG are often relational.

216 In one of its most recent sessions [47], the Working Group (WG) on
217 Data Citation reported that there are various implementations of its guide-
218 lines for Data Citation on MySQL/Postgres relational databases. Some of
219 these databases are: DEXHELPP⁶ (Social Security Records); NERC (ARGO
220 Global Array); EODC (Earth Observation Data Centre) [29]; LNEC (River
221 dam monitoring); MDS (Million Song Database) [8]; CBMI⁷ (Center for
222 Biomedical Informatics); VMC (Vermont Monitoring Cooperative); CCA⁸
223 (Climate Change Center Austria); VAMDC (Virtual Atomic and Molecular
224 Data Center) [25, 56].

225 More examples of work on data citation in relational databases are [12,
226 53, 2, 23]. The website <https://fairsharing.org/> keeps a long updated
227 list of curated and scientific databases (many of which are relational or graph-
228 based) following FAIR guidelines. These databases are citable since they are
229 compliant with the most recent guidelines, and they are in the vast majority
230 of cases accessible via dynamically created Webpages. In all these databases
231 is, therefore, possible to implement DCD on top of the existing infrastructures
232 for citing data.

233 Data citation techniques are primarily applied to relational databases
234 because of their diffusion and also because the portions of data that are to
235 be cited are easily identified: the whole database, a relation, a tuple, or
236 even an attribute. Many papers [10, 12, 2] consider more complex citable
237 units, recognizing that often the *views* of a database are the ones to be cited.
238 Generally, a *view* is a query on the database. To this end, [53] suggested
239 decomposing the database in a set of views, where each view is associated
240 with its citation.

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

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

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

⁶<http://www.dexhelpp.at/>

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

⁸<https://ccca.ac.at/startseite>

- pers” can be cited, being traditional papers that describe a database [16].
 In this case, all the credit from the citations goes to the database administrators or to the authors of the data papers.
2. Subsets of data, obtained by issuing queries to a database, are individually cited. This is the solution adopted by the *Resource Data Alliance* (RDA) working group on Data Citation [46]. In this case, the credit generated from citations can be distributed among the contributors of the portions of data being cited, and/or to the database administrators.
 3. The database is accessible via a series of Webpages that arrange the content of the database by topic or theme. Examples in the life science domain include the Reactome Pathway database [35], the GtoPdb [31], and the VAMDC [56]. Every single Webpage is unequivocally identifiable and can be individually cited.

Despite all the research efforts dedicated to the study and promotion of data citation, none of the largest citation-based systems, such as Elsevier Scopus, Web of Science, Microsoft Academia, or Google Scholar, consider scientific datasets as citable objects in academic work. Clarivate Analytics Data Citation Index (DCI) [27] is an exception, since its infrastructure tracks data usage in scientific domains and provides the technical means to connect datasets and repositories to scientific papers. However, DCI considers only citations to (previously registered and approved) databases as a whole and does not count citations to database portions such as views, tables, or tuples.

2.2. Data Credit

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

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

In [36] credit is defined as a “quantity” that describes the importance of a research entity, such as papers, software, or data, mentioned in a citation.

281 We add that the concept of credit can be built on top of the existing infras-
 282 tructure handling traditional and data citations. [36] further explores the
 283 idea of a *distribution* of credit from research entities (i.e., papers and data)
 284 to other research entities through citations that connect them. Thanks to
 285 traditional citations and now also to data citations, this distribution is fi-
 286 nally possible, at least between papers and data. Some problems related to
 287 traditional citations can thus be solved by citations:

- 288 1. Credit rewards research entities that to date are not (formally) recog-
 289 nized (a goal shared with data citation).
- 290 2. Credit can reward authors *proportionally* to their role in generating
 291 the entity. The more an author contributes to a paper, the more credit
 292 is given to him. [55] work on something similar with their zp-index,
 293 which includes in its formulation the position (and thus the role) of a
 294 publication author to represent its impact in the work itself.
- 295 3. Credit can be *transitively* channeled through a chain of papers citing
 296 each other, thus enabling the rewarding of older papers **that are no**
 297 **more cited, since other papers summarize or report their con-**
 298 **tent. Gianmaria: I do not understand this token, what do you**
 299 **mean with: papers that are no more cited?** but are nevertheless
 300 crucial in a research area for the influence of their content.

301 [26] presents a framework to distribute the credit generated by a paper to
 302 its authors and to the papers in its reference list in a transitive way. Let us
 303 consider the *citation graph* as the graph where the nodes are papers and the
 304 links are the citations among them. In this graph, every paper is a source of
 305 credit, which is then transferred to the neighboring nodes. The quantity of
 306 credit received by each cited paper depends on its impact/role in the citing
 307 paper. So far, this theoretical framework is limited to papers, but it can be
 308 easily extended to a citation graph including both papers and data.

309 [54] proposes the first method to compute credit within a network of
 310 papers citing data. Adopting a network flow algorithm, they simulate a
 311 random walker to estimate a score for each dataset, leveraging real-world
 312 usage data to compute the credit. This is the first step towards an automatic
 313 credit computation procedure. This proposal is, however, limited to assigning
 314 credit to whole datasets, and it does not deal with the granularity of data.
 315 It does not work to assign credit to a single research entity within a dataset.
 316 Differently from [54], we do not treat the credit computation process, but we
 317 focus on the distribution process.

318 2.3. Data Provenance

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

329 [17] presents four main types of data citation for DBMS: *lineage* [22],
 330 *why-provenance* [13], *how-provenance* [30] and *where-provenance* [13].

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

337 Lineage is somehow the simplest among the forms of provenance. It has
 338 been defined in different ways [17], but it can be thought of as the set of all
 339 the tuples that are used in some way by the query to produce the output
 340 tuple, the ones that are somehow *relevant* to its generation.

341 The definition of why-provenance is based on the notion of *witness set*.
 342 A witness is a set of relevant tuples that guarantees the existence of t in
 343 $Q(D)$. The lineage is therefore an example of a witness. The why-provenance
 344 of a tuple t is a peculiar set of witnesses – described in [13] – that are
 345 computed from the query, called *witness basis*. A witness basis may be
 346 composed of more than one witness. Therefore, the why-provenance contains
 347 more information than the lineage, since it describes *alternative* ways in
 348 which the same output may be generated.

349 The how-provenance takes the form of a polynomial, called *provenance*
 350 *polynomial*, where the variables are taken from the set of identifiers of the
 351 tuples (provided that each tuple in I has an identifier) and the coefficients are
 352 taken from \mathbb{N} . This provenance also contains information on *how* the input
 353 tuples are used. For example, when two tuples are combined by a join, they
 354 are also combined in the polynomial by the \cdot operator. When two or more

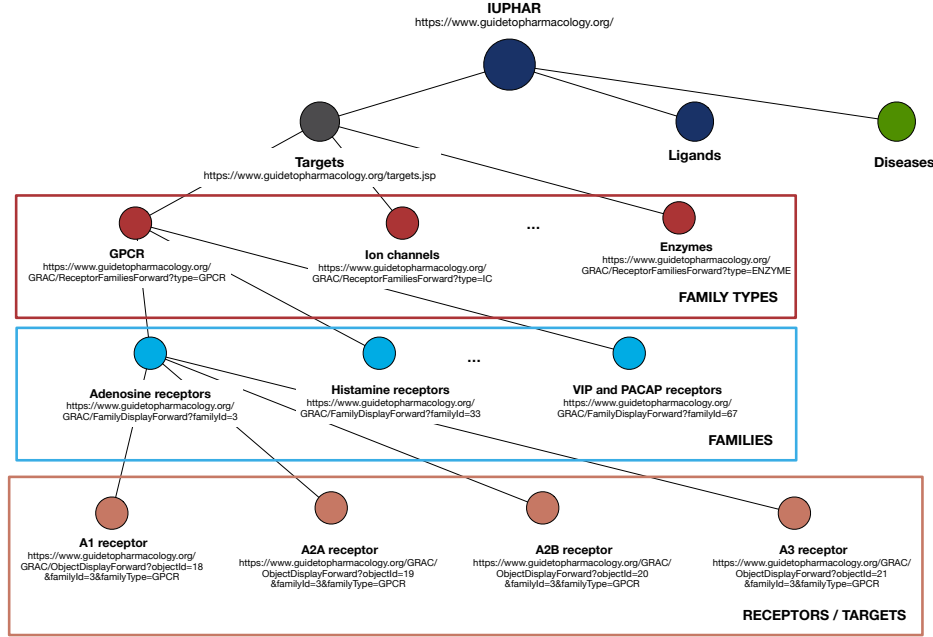


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

355 tuples become equivalent due to a union or a projection, the corresponding
 356 monomials are combined by the $+$ operator.

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

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

362 3. Use Case: GtoPdb

363 As use case we refer to the IUPHAR/BPS Guide to Pharmacology [31]
 364 or GtoPdb⁹. GtoPdb is a well-known and well structured scientific relational
 365 database that contains expertly curated information about diseases, drugs
 366 in clinical use, their cellular targets, and the mechanisms of action on the
 367 human body. It is curated and maintained by the GtoPdb Committee, and

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

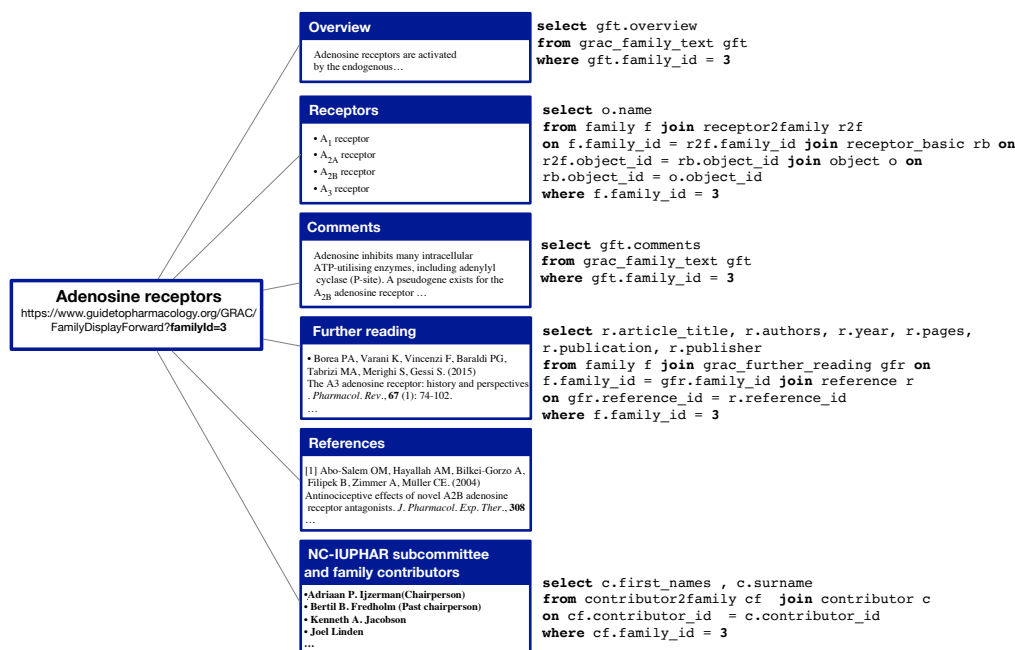


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

by 96 subcommittees, comprising 512 scientists collaborating with in-house curators who draw the information contained in the database from high-quality pharmacological and medicinal chemistry literature. Roughly 1000 researchers from all over the world have contributed to the database, and the curators wanted to give recognition to these contributors. This led to some early work on data citation [10].

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

GtoPdb provides access to the webpages corresponding to all these nodes

385 through URLs. The webpages corresponding to target families all present a
386 similar structure, as shown in Figure 3 for the “Adenosine receptors” family.
387 Each page has an *Overview*, a brief text describing the content of the page;
388 a list of *Receptors* comprising the family; a section of *comments* about the
389 family; the *References*, a list of the papers consulted by the curators of the
390 page, similar to a reference list of a paper; the *further reading* list, reporting
391 papers that an interested reader may want to consult to obtain more insight
392 on the family; and a final section called *How to cite this family page*, con-
393 taining text snippets useful to cite the specific page or the whole database.
394 Figure 3 shows the SQL code that retrieves the information used to build the
395 corresponding sections (apart from the References section). Therefore, each
396 family page can be considered a full-fledged traditional publication, consist-
397 ing of title, authors, abstract (the overview), content, and references.

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

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

411 For our running example, consider Table 1. This simplified version of
412 GtoPdb illustrates three tables: **family**, **contributor** and **contributor2family**.
413 The first table, **family**, has tuples representing families with three attributes:
414 the id of the family, its name, and type. Table **contributor** consists of peo-
415 ple who have helped generate the data of the database. The third table,
416 **contributor2family**, serves as a link between the families and the people
417 who contributed to them. For instance, “John Smith” (c_1) contributed to
418 “Dopamine Receptors” (f_1) as well as to the “YANK Family” (f_4). We use
419 this example throughout the rest of the paper. In particular, we are using

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

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

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

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

the *id* attribute of the tables as *provenance token* of its corresponding tuples, that is, as a symbol that serves to identify a tuple when talking about provenance.

4. Data Provenances

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

4.1. Lineage

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

```

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

```


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

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

438 Table 2 shows the query result, which consists of two tuples. We add
439 an extra attribute `id` so that we can easily refer to each result tuple. The
440 lineage for tuple o_1 is the set $\{f_1, c2f_1, c_1, c2f_2, c_2\}$, since the tuple f_1 was
441 joined with $c2f_1$ and then with c_1 , and was also joined with $c2f_2$ and c_2 . No
442 other tuple is used in the database to produce o_1 . For tuple o_2 the lineage is
443 $\{f_4, c2f_4, c_1\}$. Lineage is defined for each tuple of the output, and can differ
444 between tuples.

445 4.2. Why-Provenance

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

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

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

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

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

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

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

479 4.3. How-Provenance

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

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

490 Table 4 shows a simple example in which the two output tuples of our
 491 running example are annotated with their respective how-provenances. Tuple
 492 o_2 was produced through the join among the input tuples $f_4, c2f_4$, and c_1 .
 493 The three provenance tokens are, therefore “multiplied” together. The case of
 494 o_1 is slightly more complex. This tuple, as already discussed, can be obtained
 495 through two different joins. The two monomials composing the polynomial

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

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

represent these two alternatives. They correspond, in a way, to the witnesses of the why-provenance of o_1 . The $+$ operator represents the fact that the two monomials describe alternative derivations. The output tuple is the result of a merge of two distinct tuples after the projection on the attribute **name**. This merge is due to the fact that the result of a relational algebra expression is always a *set* of tuples, which corresponds to the presence of the **DISTINCT** operator in an SQL query. This simple example gives the basic idea behind how-provenance and how it allows us to track the operations that produced an output tuple.

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

5. Credit Distribution and Distribution Strategies

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

5.1. Data Credit and Data Credit Distribution

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

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

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

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

534 **Definition 5.1. Tuple Level Data Credit Distribution (DCD) [24]**
 535 Given a query Q over I and $k \in \mathbb{R}_{>0}$, DCD is defined by the function $f_{I,Q} :$
 536 $\text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ such that $f_{I,Q}(t, k) = h$ where $0 \leq h \leq k$ and
 537 $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$. The function $f_{I,Q}$ is the *distribution strategy* (DS).

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

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

548 5.2. A Lineage-based Distribution Strategy

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

Definition 5.2. Lineage-based Distribution Strategy [24]

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

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

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

556 As an example, consider the output tuples of Table 2, and assume that
 557 each output tuple has credit $k_o = 1$. The lineage of the first tuple, o_1 , is
 558 the set $\{f_1, c2f_1, c_1, c2f_2, c_2\}$. Therefore, each tuple in this set receives credit
 559 $1/5$. The other tuples of the database receive zero credit. The lineage of the
 560 second output tuple is $\{f_4, c2f_4, c_1\}$, therefore each of these tuples receives
 561 credit $1/3$.

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

568 Not all of the tuples in the lineage of an output tuple are necessary to be
 569 present at the same time for the output tuple to appear in the query results.
 570 For example, if the database only had the set of tuples $\{f_1, c2f_1, c_1\}$ or the set
 571 $\{f_1, c2f_2, c_2\}$, the existence of o_1 would still be guaranteed. In other words,
 572 while f_1 is always needed for o_1 to appear in the output, only one of the sets
 573 of tuples $\{c2f_1, c_1\}$ and $\{c2f_2, c_2\}$ is required. One could therefore argue that
 574 it would be more fair for f_1 to receive more credit than the other four tuples,
 575 given its role in producing o_1 .

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

581 5.3. A Why-Provenance-Based Distribution Strategy

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

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

588 *Let I be a database instance, Q a query over I , $o \in Q(I)$ an output tuple*

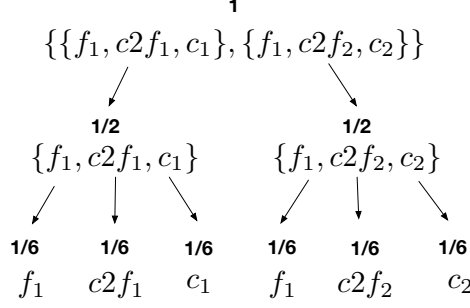


Figure 4: Distribution of credit using why-provenance-based DS for tuple o_1 .

589 and k_o the total credit associated to o . Let $\mathcal{W} = \text{Why}(Q, I, o)$ be the witness
 590 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:

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

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

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

603 5.4. A How-Provenance Based Distribution Strategy

604 How-provenance conveys more information than why-provenance since
 605 it not only captures what tuples are relevant to the output and in which

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

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

606 combination, but also how they are used. The “how” is captured through
607 the provenance polynomials.

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

612 To define the DS more formally, we introduce some notation and illustrate
613 it using the provenance polynomial \mathcal{H} shown in Figure 5.

614 We call c the function that, given a polynomial, returns the sum of the
615 coefficients of the polynomial; thus $c(\mathcal{H}) = 3 + 1 = 4$. We use the same name
616 for the function that, given a monomial, returns the sum of its exponents;
617 thus $c(M_2) = 1 + 3 + 3 = 7$. mc is the function that takes as input a monomial
618 and returns its coefficient. e is a function that takes as input a tuple and a
619 monomial, and returns the exponent of the tuple in the monomial, if present;
620 thus $e(c_2, M_2) = 3$. γ takes as input a tuple and the whole polynomial, and
621 returns a set containing the monomials containing that tuple, if present in
622 the polynomial; thus $\gamma(f_1, \mathcal{H}) = \{M_1, M_2\}$.

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

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

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

627 Going back to the example of Table 4, consider o_1 with provenance poly-
628 nomial $f_1 c_2 f_1 c_1 + f_1 c_2 f_2 c_2$. The how-provenance-based DS firstly divides

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

the credit between the two monomials. Since the coefficients of each monomial are 1, the credit is split in half. If they were, for example, 1 and 2 respectively, 1/3 of the credit would go to the first monomial, and 2/3 to the second. Since in our example each variable has exponent 1, the credit is further divided equally among the three variables. Thus, at the end of the computation, f_1 receives 1/3, and the other tuples receive 1/6. If, for example, the first monomial was $f_1^2 c2f_1 c_1$, then the portion of credit of this monomial would be divided in this way: 1/2 to f_1 and 1/4 to each of the other two tuples.

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

```

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

```

The result of Q2 is shown in Table 5, and consists of one tuple, annotated with each of the three provenances. As can be seen, lineage and why-provenance are identical to those of the tuple o_1 in the previous example. The how-provenance, however, is different since tuple f_1 is used twice: first in the join of the inner query, and second in the join of the outer query. This information is lost in the first two forms of provenances since they are sets, but it is captured in how-provenance through the use of the operator ‘.’.

Figure 6 shows the differences between the three DS for the tuple o_1 of Table 5. Subfigure 5.a uses lineage, sub-figure 5.b uses why-provenance, and

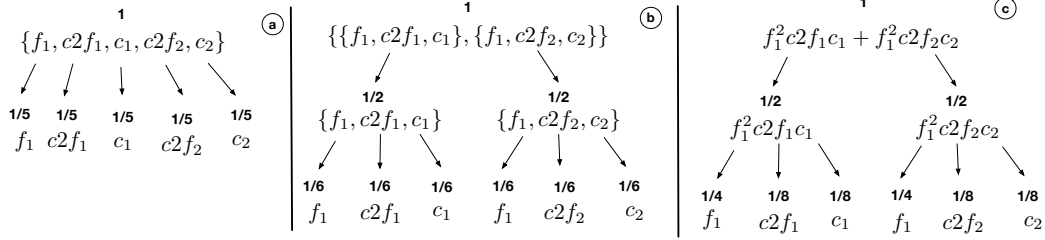


Figure 6: Comparison of different distributions strategies for tuple o_1 produced by query Q2.

sub-figure 5.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 can reward f_1 more, showing that how-provenance is even more sensitive to the tuples' role in a query than why-provenance. This is a direct consequence of the fact that, as proven in [30], how-provenance is more general than why-provenance and lineage, in the sense that it contains more information.

6. Experimental Evaluation

To understand the trade-off between these Distribution Strategies (DS), we perform three sets of experiments using queries over GtoPdb. The first set of experiments used real queries extracted from citations to GtoPdb published in the British Journal of Pharmacology. The second set uses different sets of synthetically produced provenance polynomials, corresponding to more complex queries, highlighting the differences between the different DS employed. In the third set of experiments, we compare traditional citations and credit in rewarding data curators.

6.1. Real-world queries

We evaluate the proposed distribution strategies on GtoPdb, and in particular, we focus on target families described on the GtoPdb website. There are eight family types: *GPCR*, *Ion channels*, *NHRs*, *Kinases*, *Catalytic receptors*, *Transporters*, *Enzymes* and *Other protein targets*.

When a paper uses data from GtoPdb, it can cite the full database, the webpage of interest, or a subset of data extracted with a query. We

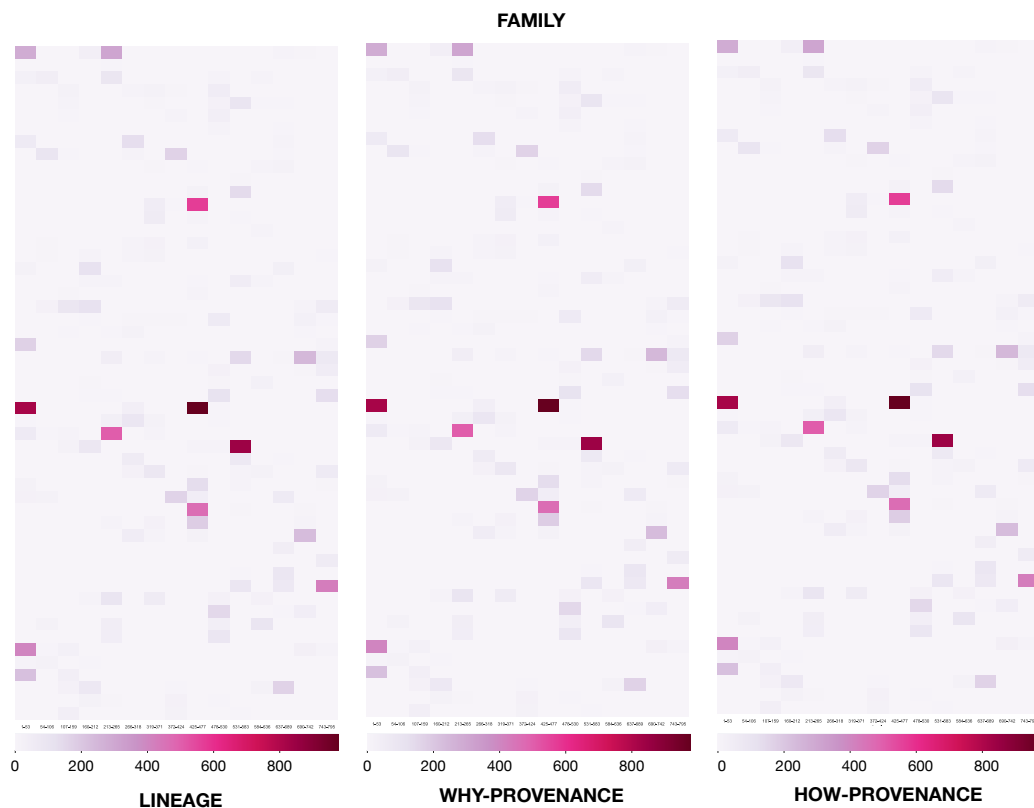


Figure 7: Comparison of three DS on the same table `family` using the distribution given by the queries retrieved from papers.

682 consider as sources of citations the papers published in the British Jour-
683 nal of Pharmacology (BJP) ¹¹, since each time they cite a webpage from
684 GtoPdb, they report the URL of that page. From that URL, it is possible
685 to reverse-engineer the queries used to obtain the pages' data. In particular,
686 we considered all the 889 papers in BJCP citing the IUPHAR/BPS Guide to
687 pharmacology [31] as of October 2020. The IUPHAR/BPS guide is a data
688 journal that describes the structure and evolution of GtoPdb. Every two
689 years, the GtoPdb consortium releases such a journal to describe the evolu-
690 tion of the databases. At the time of writing, [31] received more than 1200
691 citations on Google Scholar.

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

692 The queries that we inferred are those used to build a target family web-
693 page that we reported in Figure 3, where we see how the structure of the
694 “Adenosine receptors” family is mapped into the queries to get the informa-
695 tion reported in the corresponding webpage. In GtoPdb, all target family
696 pages share a similar structure (the only difference is that individual sections,
697 such as “contributors” or “further readings”, may be absent). Therefore, the
698 same queries can build all the target family pages by simply changing the
699 family id used in the query (in Figure 3, it is 3). All these queries are SPJ.
700 A total of more than 12K different queries were built in this way¹². Without
701 any loss of generality, we decided that each tuple in these queries’ output
702 carries a default credit of 1.

703 Figure 7 shows the heat-maps obtained by the distribution of credit per-
704 formed by the three different DS on the `family` table of GtoPdb. `family` is a
705 table describing the characteristics and necessary information of the receptor
706 families and, as can be seen in Figure 3, it is often used in join with other
707 tables to get the data to build a webpage.

708 The result of the distribution is the same using the three strategies. The
709 same effect is also obtained with the other tables of the database used by
710 the queries shown in Figure 3. This is because of the conditions in which we
711 produced this experiment. Indeed, the considered queries are all SPJ using
712 each table only once in the join condition and joins are on key attributes.
713 With these specific conditions, each tuple of the output presents: (i) a how-
714 provenance that is a single monomial with coefficient 1 and exponent 1 in
715 each variable; (ii) a why-provenance that is composed of only one witness;
716 (iii) a lineage that coincides with the only witness in the basis. Hence, given
717 these queries, the three distributions act in the same way. The credit is
718 always uniformly distributed among the tuples present in each provenance.

719 To better clarify what is happening, let us consider one of the types of
720 queries used to build the output webpage, as shown in Figure 3:

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

¹²For reproducibility purposes, the code we used for our experiments and all the produced queries are available here: https://bitbucket.org/dennis_dosso/credit_distribution_project.

Q3 returns a series of 10 tuples from the considered GtoPdb version. The first tuple produced by this query, <Bertil B., Fredholm>, has $c_{939} \cdot c2f_{496}$ as provenance polynomial. c_{939} represents the provenance token of a tuple in **contributor**, the same for $c2f_{496}$ in table **contributor2family**. The why-provenance of this tuple is $\{\{c_{939}, c2f_{496}\}\}$ and its lineage is $\{c_{939}, c2f_{496}\}$. Therefore, the credit assigned to these tuples is 1/2 using all three DS. This happens for all the tuples in the output of each query of GtoPdb, thus making the distributions equivalent to their output.

This is not always the case with general queries and other databases. As we showed in the examples in the previous section, when two or more tuples are merged by the effect of a projection or union, we see sensible differences between the three distribution strategies.

6.2. Synthetic queries

To better show the differences between the three DS, let us consider the case reported in Figure 8. The figure reports a distribution of credit performed on the table **family** through the generation of 10K *synthetic* polynomials. We randomly generated provenance polynomials that might be the how-provenance of randomly generated synthetic queries, using the three GtoPdb tables **family**, **contributor2family**, and **contributor**. An example of such synthetic polynomial is:

$$3f_1^3c2f_1^2c_1^2 + 2f_1c2f_2^3c_2^3 + 4f_5c2f_{17}^4c_{18}^3$$

As can be seen, we made sure to also include coefficients and exponents that differ from 1. Its corresponding why-provenance is:

$$\{\{f_1, c2f_1, c_1\}, \{f_1, c2f_2, cf_2\}, \{f_5, c2f_{17}, c_{18}\}\}$$

its lineage is:

$$\{f_1, f_5, c2f_1, c_1, c2f_1, c2f_2, c2f_{17}, c_1, c_2, c_{18}\}$$

These types of polynomials are not impossible to obtain in real applications. They can be obtained by any nested queries with join and union operations that use multiple times the same tuples (e.g., the presence of exponents bigger than 1) and the same combination of operations more than

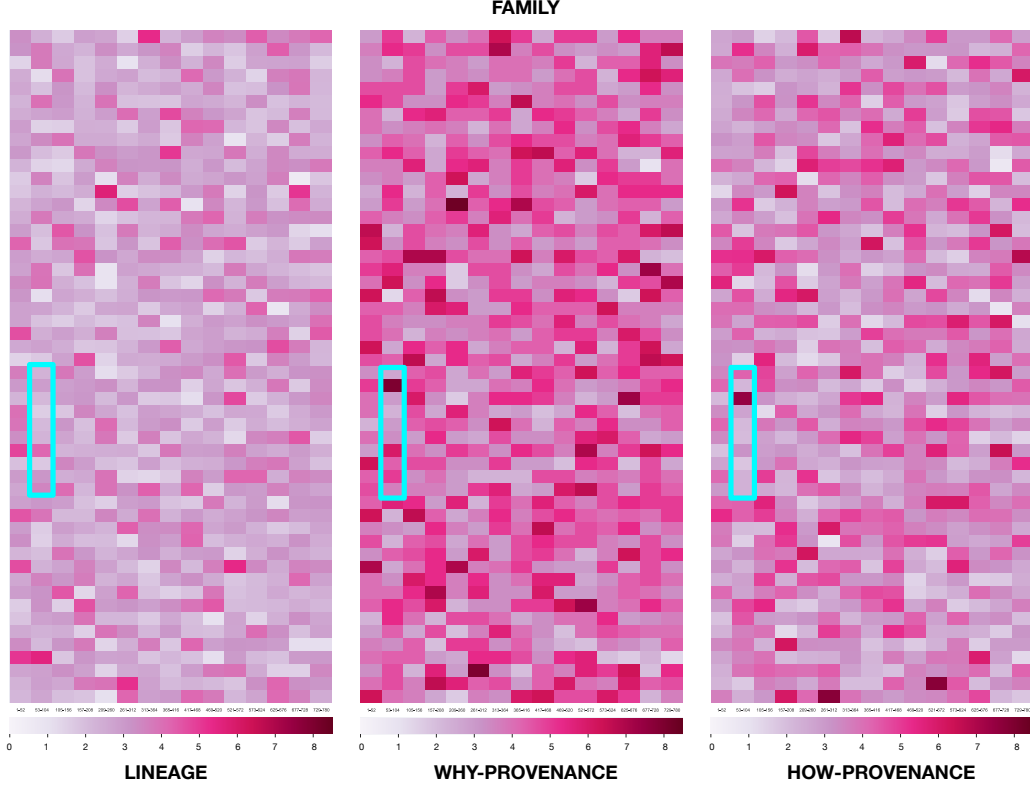


Figure 8: Comparison of three DS on the same table **family** after the distribution computed using 10K synthetic and randomly generated provenance polynomials. The tuples in the blue rectangles are used as example in the discussion connected to Figure 9.

755 once (e.g., the presence of coefficients for monomials bigger than 1). We
 756 randomly generated a set of 10K synthetic polynomials.

757 Using the how-provenance, the distribution obtained from the example
 758 polynomial we are considering is the following:

$$f_1 = \frac{59}{315}, f_5 = \frac{1}{18}, c_2 f_1 = \frac{2}{21}, c_2 f_2 = \frac{2}{15}, c_2 f_{17} = \frac{2}{9}, c_1 = \frac{2}{21}, c_2 = \frac{2}{15}, c_{17} = \frac{1}{6}$$

759

760 Using the why-provenance, the output is:

$$f_1 = \frac{2}{9}, f_5 = \frac{1}{9}, c_2 f_1 = \frac{1}{9}, c_2 f_2 = \frac{1}{9}, c_2 f_{17} = \frac{1}{9}, c_1 = \frac{1}{9}, c_2 = \frac{1}{9}, c_{17} = \frac{1}{9}$$

761

762 Finally, with the lineage, the distribution is:

$$f_1 = \frac{1}{8}, f_5 = \frac{1}{8}, c2f_1 = \frac{1}{8}, c2f_2 = \frac{1}{8}, c2f_{17} = \frac{1}{8}, c_1 = \frac{1}{8}, c_2 = \frac{1}{8}, c_{17} = \frac{1}{8}$$

763

764 To highlight how the distributions behave differently with these poly-
 765 nomials, consider tuple f_5 . f_5 receives the highest quantity of credit when
 766 we use the lineage-based distribution. Why-provenance and how-provenance
 767 distribute less credit to that tuple because more information is available for
 768 the computation and the algorithms weigh less and less its role.

769 Generally speaking, the more complex the distribution, the more polar-
 770 ized the credit is toward the tuples that are more frequently used or with a
 771 higher impact in producing the output tuple.

772 Going back to Figure 8, we can see how the three provenances behaved
 773 differently. We set the maximum value for the heat-maps to the highest
 774 value reached by a tuple in all three distributions (i.e., 8.33). Note that
 775 lineage is the form of provenance giving less credit to the tuples of the **family**
 776 table. This is because this DS equally distributes the credit to all the tuples
 777 appearing in the lineage. Since these queries use other two tables, the credit
 778 is also given to those tables' tuples.

779 Moving to the heat-map reporting the distribution performed by the DS
 780 based on why-provenance, we see that this time more credit is given overall
 781 to the tuples of the table. This DS is the one that distributes more credit to
 782 the **family** table, among the three strategies. This is because the DS based
 783 on why-provenance also considers the different ways a tuple is used, e.g., in
 784 other joins. If the same tuple is present in more than one witness, it is more
 785 probable that it will attract more credit, withdrawing it from the other tuples
 786 in the witness basis. In this case, **family** drew more credit, taking it from
 787 the other two tables, due to the role of its tuples in the queries that were
 788 executed.

789 Let us consider the heat-map produced by the DS based on the how-
 790 provenance. We can see how, although it presents more credit in its tuples
 791 that the one present in the lineage heat-map, it does not reaches the levels
 792 of the why-provenance heat-map. This is due to the fact that this DS is
 793 even more sophisticated, weighting even more than the previous DS the role
 794 of tuples in the production of the output. The result is a distribution that
 795 still rewards the tuples of this table more than lineage, but not in the same
 796 measure as the DS based on why-provenance, since the other tuples in the



Figure 9: Comparison of the distribution of credit performed by the three DSs on a subset of 10 tuples taken from table `family` simulating the passing of time. The number on top of each group of heat-maps represent the number of queries computed.

797 other tables are able to attract more credit due to their roles in the queries.
798 ← **Gianmaria: This is not clear and I suggest a rewriting of this**
799 **paragraph.**

800 To show how the DS based on different provenances may differ in their
801 behavior also through the course of time, let us consider Figure 9.

802 In this figure, we report four groups of heat-maps. Each group presents
803 three maps obtained by selecting the same ten tuples from the GtoPdb
804 `family` table after an incremental distribution of credit (the tuples of ranks
805 ranging from 79 to 89). These are the same tuples highlighted in the blue
806 boxes in Figure 8. In particular, the four groups represents “snapshots”
807 taken during an incremental accumulation of credit on the database, at dif-

808 ferent moments chosen when a certain number of executed queries is reached
809 (specifically, 1K, 2K, 5K and 10K). Figure 8 represents the end of the process.

810 In this way, we simulate the passing of time on a database where credit
811 distribution is performed. Each group of heat-maps can be thought of as a
812 snapshot of that set of tuples at a certain moment. The queries utilized are
813 the same as the experiment reported in the previous section. The range of
814 credit in each map goes from 0 (no credit) to 6 (maximum quantity of credit
815 reached on a tuple at the “snapshot” with 10K queries).

816 Focusing on the 1K and 2K groups, we see that the tuples highlighted
817 by the three DS are almost the same. Still, there are small differences, in
818 particular in tuple 5.

819 The first interesting differences come to light with 5K queries. In particu-
820 lar, we note how tuple 7 is rewarded poorly by the DS based on lineage, while
821 it is rewarded more by why-provenance-based DS and most of all by the DS
822 based on how-provenance. This is because tuple 7 appears in a relatively low
823 number of lineages, but its role is critical to these queries; thus, the other
824 DS reward it more. On the other hand, a tuple 5 is highly rewarded by the
825 DS based on lineage and why-provenance, and less by how-provenance. Al-
826 though tuple 5 appears in many queries and used in different combinations,
827 its exponents in the provenance polynomials must be low, therefore giving it
828 low credit with how-provenance. It is also interesting to note how other tu-
829 ples like tuple 2 now surpass certain tuples, like tuple 1 that up to 2K queries
830 presented the highest values of credit. This shows how credit can keep track
831 of the “hotspots” in a database over time. The presence of new queries and
832 new credit distributions can change the hotspots in a table, showing how the
833 research community’s interests may change during time.

834 Finally, the highest differences are shown in the 10K group. In this case,
835 we see a situation similar to the one with 5K queries. Like 8 or 10, specific
836 tuples receive more credit with why-provenance and how-provenance, rather
837 than with lineage. This is still due to the critical role of the tuple in the
838 queries where it appears.

839 From this progression, we see how, given the peculiar synthetic prove-
840 nance polynomials that we presented, we can see the differences between the
841 three distributions. These differences become more evident with time, i.e.,
842 the more credit is distributed to the tuples.

843 The DS based on lineage is sufficient when a user only wants to highlight
844 the tuples of the database used by a query (and not only visualized in the
845 output). However, it equally distributes the credit to the tuples of the lineage,

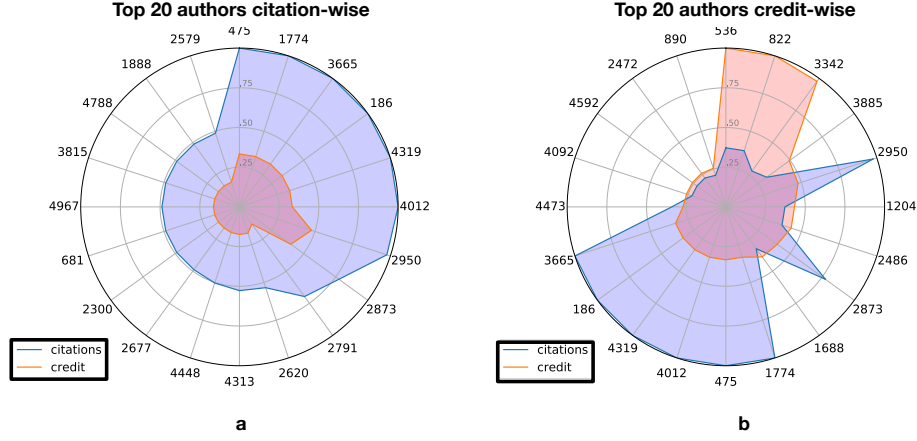


Figure 10: Radars presenting the top 20 authors citation-wise and credit wise, together with their (normalized between 0 and 1) values of citations and credit.

therefore not considering the information on the tuples' role in the production of the output.

For this reason, a user may want (depending on the nature of the queries) to use DS based on why-provenance and how-provenance. Using the why-provenance and how-provenance DS, it is possible to change the distribution of credit to the tuple, rewarding more the tuples that have a more critical role in generating the output. Therefore, these two DS can be preferred when the user aims to find “hotspots” in the database based on the tuples' role.

6.3. Credit vs Citations

We compare traditional citations and credit for the last set of experiments to check their behavior difference when rewarding data curators. Consider the two radar plots in Figure 10. Figure 10.a reports the top 20 author (we identify the authors with their ID instead of their name), ordered based on the normalized value of citations distributed by the queries taken from the papers published in BJP as described in Section 6.1, together with their normalized value of credit. An author transitively receives credit from the data s/he created or curated. The credit assigned to data is then split equally to the authors of those tuples. As shown in Section 6.1, there is no difference for these queries in the distribution of credit between the three DS. Thus these values are equal for the three distributions. The second plot is similar to the first one, but the authors are ordered based on the received credit. As

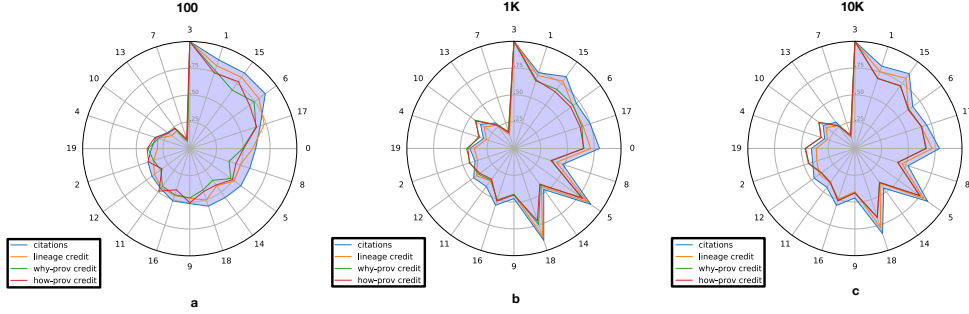


Figure 11: Radars presenting the 20 synthetic authors with corresponding citation and quantities of credit distributed through the 3 DS (all values normalized between 0 and 1) through different numbers of queries (respectively, 100, 1K and 10K). The order is the descending one of the citations of the authors with 100 queries.

we see, the quantity of credit and the number of citations differ sizeably; i.e., an author with the highest number of citations does not necessarily have the highest credit value. As shown in Figure 10.b, the authors with the highest value of credit do not also have the highest number of citations. This means that there are citations that are more “valuable” for an author regarding credit. This is because the quantity of credit assigned by these citations is very high, i.e., the impact of those cited data is high.

This shows how credit can reward authors that are cited less than others but, nonetheless, have a high impact on the research community.

Let us consider Figure 11. In this case, we produced 100, 1K, and 10K synthetic polynomials, and we distributed credit through them. Since these polynomials correspond to queries whose authors are not easily identifiable, we created 20 “synthetic” authors, and we randomly assigned one author to each tuple in the database. The authors receive “blocks” of consecutive tuples, with each block of the size varying between 10 and 40 to simulate different quantities of “work” performed by an author. Every time a tuple was used in a provenance polynomial, we assigned one citation to the author corresponding to the tuple. The same author also receives the three different credits given to the tuple at the end of the distribution process using the three DS.

The three radar plots in the upper row presents 20 authors are sorted based on the normalized number of received citations, together with the corresponding normalized quantities of credits assigned using the 3 different strategies. The ones in the bottom row present the authors ordered based on

the quantity of credit assigned by the DS based on how-provenance. In this way we can see the behavior of the different DS in rewarding the authors with the highest levels of citations. Given the synthetic nature of these queries, the correlation between the number of citations and the quantity of credit assigned to the authors appears to be a much stronger with respect to the case with the real-world queries.

Nonetheless, it is possible to note that credit still behaves differently,. In particular we still see that certain authors that are not in the top 10 positions citation-wise are still rewarded with high quantities of credit, showing their importance with respect to their impact. Interestingly, scaling up to $1K$ and $10K$ polynomials, it appears that the distribution performed via why-provenance and how-provenance become equivalent for the authors. We can note that, although not exactly equal, the values of credit assigned to the authors by those DS become quite similar with these higher quantities of polynomials.

6.4. Execution times

# of polynomials	lineage	why-prov.	how-prov.
100	226.6 ms	192.0 ms	185.5 ms
200	431.2 ms	392.2 ms	403.2 ms
500	1.013 s	934.2 ms	881.8 ms
1K	2.041 s	1.934 s	1.744 s
2K	3.773 s	3.491 s	3.510 s
5K	8.992 s	8.653 s	8.889 s
10K	17.10 s	16.84 s	16.84 s
20K	34.59 s	35.30 s	39.70 s
100K	3.289 min	3.442 min	3.652 min
1M	35.91 min	34.87 min	37.91 min

Table 6: The times required to perform the three DS for different number of synthetic polynomials.

In Table 5 we report the time required to compute the distribution using the DS based on the three provenances. As we see, the execution time grows linearly with the number of polynomials that are submitted to the system. With a high number of polynomials (1M), the time required by the DS based on lineage and why-provenance is lower than the time needed for the DS based on how-provenance. This is due to the more significant number of

913 operations required to calculate the how-provenance DS and distribute the
 914 portions of credit to be assigned to the different tuples. We note that, since
 915 we created these polynomials on-the-fly, these values do not include the time
 916 required to compute the provenances. Therefore, limited to the time required
 917 to distribute credit, the three DS are equivalent in terms of performances.
 918 The first differences can be seen only with high number of polynomials, when
 919 lineage and why-provenance may be preferred if there are no requirements
 920 to assign credit with the strategy implemented by the how-provenance-based
 921 DS.

922 All the experiments were carried on a MacBook Pro 13-inch, 2019 with
 923 2.4 GHz processor Intel Core i5 quad-core, 8 GB of memory at 2133 MHz
 924 with code written in Java and the support of a PostgreSQL database.

925 7. Conclusions

926 This paper expanded on our previous work on data credit and data credit
 927 distribution in [24] by defining two new distribution strategies, based on
 928 why- and how-provenance. The first distribution is based on the concept of
 929 witness, and it can give more credit to tuples that appear in more than one
 930 witness. In other words, tuples that are more important to the query and are
 931 used in different ways are also rewarded more by the strategy. The second
 932 DS, based on how-provenance, considers the frequency in which a tuple or a
 933 combination of tuples is used in the query through the information contained
 934 in the provenance polynomial. In this case, the distribution is even more
 935 sensitive than the first one to the role and importance of tuples.

936 To show the differences between the three DS (also considering the one
 937 based on lineage, defined in our previous work), we performed different ex-
 938 periments on GtoPdb, a curated scientific relational database, with the use
 939 of both real and synthetic queries. In the first set of experiments, we used
 940 SPJ queries extracted by data citations present in papers published in the
 941 British Journal of Pharmacology. Employing these queries, we were able to
 942 distribute the credit to the tuples in different tables of the database, high-
 943 lighting the tuples used more than others. We showed that with these queries,
 944 the three strategies produce the same distribution. These are SPJ queries
 945 that do not present self-joins, and therefore the formulas at the base of the
 946 DS have the same output.

947 In the second set of experiments, we synthetically produced more complex
 948 provenance polynomials, corresponding to more complex synthetic queries,

949 that present exponents and coefficients different than 1. In this way, we
950 showed that, even though all three DS can highlight all the tuples used by
951 the queries in the database, the three have different behaviors. While the DS
952 based on lineage rewards all the tuples used by a query in equal measure, the
953 strategy based on why-provenance tends to reward the tuples more critical
954 to the query. In particular, why-provenance can consider the different ways
955 in which one tuple is used in a query. How-provenance is even more sensitive
956 to the tuples' role: it can also consider the frequency by which a tuple or a
957 set of tuples is used in the case of more complex queries. Depending on the
958 goal of a user, one provenance may be preferred to another.

959 We also showed how the differences between the DS become more and
960 more evident with the passing of time, i.e. when more and more polynomials
961 are processed by the system.

962 In the third set of experiments we compared the citations to the authors
963 to the credit brought to them. We showed how, both in the real-world and
964 synthetic scenarios the credit rewards more the authors that have a higher
965 impact, i.e. the authors connected to the data that produce the highest
966 quantities of credit, and not necessarily the data with the highest citation
967 count. In this sense, credit appears to be an useful new measure to discover
968 data and their corresponding curators that have a high impact in the research
969 world, even when they are cited few times or do not appear at all in the data
970 that are cited (i.e. the case of data used to build the output of a query but
971 that is not visualized in the output itself).

972 In future work, we plan to explore the different potential applications of
973 credit on relational databases. One example is the so-called *data pricing*.
974 Data pricing consists of giving a price to a query submitted by a user who
975 wants to buy the produced information. Currently, a commonly used strategy
976 to face data pricing is based on query rewriting. A database stores a set of
977 views correlated with their price. When a new query arrives, the system tries
978 to rewrite it using the stored views and obtain a query price. This process
979 is computationally expensive. We plan to distribute credit through carefully
980 planned and representative queries and use it as information to define a new,
981 faster, and potentially more flexible pricing function.

982 Another application is *data reduction* [42], concerned with reducing the
983 vast mole of data that is produced in the evolving world of research and
984 information technology. Data reduction deals with different aspects of dealing
985 with huge amounts of data, such as finding reduced and relevant data streams
986 from the multiple gigabytes of data produced by big data systems every

987 second or dealing with the curse of dimensionality which requires unbounded
988 computational resources to uncover actionable knowledge patters [51].

989 Data credit can also help to find “hotspots” and “coldspots”. A hotspot
990 is data in a database (a tuple or a single attribute, for example) that presents
991 a high quantity of credit and is therefore valuable for the set of queries that
992 distributed that credit. On the other hand, a coldspot is data that present
993 low quantities of credit and can be considered useless or less relevant and can
994 therefore be removed or moved in another cheaper and less efficient memory
995 location.

996 References

- 997 [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bern-
998 stein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L.,
999 Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Koss-
1000 mann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo,
1001 T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo,
1002 A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suciu, D.
1003 (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–
1004 53.
- 1005 [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating
1006 data citation in citedb. *PVLDB*, 10(12):1881–1884.
- 1007 [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y.
1008 (2018). Data citation: A new provenance challenge. *IEEE Data Eng.*
1009 *Bull.*, 41(1):27–38.
- 1010 [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An
1011 Introduction to the Joint Principles for Data Citation. *Bulletin of the*
1012 *Association for Information Science and Technology*, 41(3):43–45.
- 1013 [5] Baggerly, K. (2010). Disclose all data in publications. *Nature*,
1014 467(7314):401–401.
- 1015 [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D.,
1016 Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I.,
1017 Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and
1018 Goble, C. A. (2013). Why linked data is not enough for scientists. *Future*
1019 *Gener. Comput. Syst.*, 29(2):599–611.

- 1020 [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation
1021 Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.
- 1022 [8] Bertin-Mahieux, T., Ellis, D., Whitman, B., and Lamere, P. (2011). The
1023 million song dataset. In *Proceedings of the 12th International Conference*
1024 *on Music Information Retrieval (ISMIR 2011)*, pages 591–596.
- 1025 [9] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In
1026 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*
1027 *nication*, pages 93–116. De Gruyter Mouton.
- 1028 [10] Buneman, P. (2006). How to cite curated databases and how to make
1029 them citable. In *18th International Conference on Scientific and Statistical*
1030 *Database Management, SSDBM*, pages 195–203. IEEE Computer Society.
- 1031 [11] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D.,
1032 Pawson, A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn’t
1033 working, and what to do about it. *Database J. Biol. Databases Curation*,
1034 2020.
- 1035 [12] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation
1036 is a computational problem. *Commun. ACM*, 59(9):50–57.
- 1037 [13] Buneman, P., Khanna, S., and Tan, W. C. (2001). Why and where: A
1038 characterization of data provenance. In *Database Theory - ICDT 2001,*
1039 *8th International Conference*, pages 316–330.
- 1040 [14] Buneman, P. and Silvello, G. (2010). A rule-based citation system for
1041 structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.
- 1042 [15] Callaghan, S., Donegan, S., Pepler, S., Thorley, M., Cunningham, N.,
1043 Kirsch, P., Ault, L., Bell, P., Bowie, R., Leadbetter, A. M., Lowry,
1044 R. K., Moncoiffé, G., Harrison, K., Smith-Haddon, B., Weatherby, a.,
1045 and Wright, D. (2012). Making Data a First Class Scientific Output:
1046 Data Citation and Publication by NERC’s Environmental Data Centres.
1047 *International Journal of Digital Curation*, 7(1):107–113.
- 1048 [16] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Jour-
1049 nals: A Survey. *Journal of the Association for Information Science and*
1050 *Technology*, 66(9):1747–1762.

- 1051 [17] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases:
1052 Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–
1053 474.
- 1054 [18] CODATA-ICSTI Task Group on Data Citation Standards and Practices
1055 (2013). *Out of Cite, Out of Mind: The Current State of Practice, Policy,*
1056 *and Technology for the Citation of Data*, volume 12.
- 1057 [19] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N.
1058 (2019). Bringing citations and usage metrics together to make data count.
1059 *Data Science Journal*, 18(1).
- 1060 [20] Cronin, B. (1984). *The citation process. The role and significance of*
1061 *citations in scientific communication*. London: Taylor Graham.
- 1062 [21] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evi-
1063 dence of a structural shift in scholarly communication practices? *JASIST*,
1064 52(7):558–569.
- 1065 [22] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of
1066 view data in a warehousing environment. *ACM Trans. Database Syst.*,
1067 25(2):179–227.
- 1068 [23] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model
1069 for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on*
1070 *Innovative Data Systems Research*. www.cidrdb.org.
- 1071 [24] Dosso, D. and Silvello, G. (2020). Data credit distribution: A
1072 new method to estimate databases impact. *Journal of Informetrics*,
1073 14(4):101080.
- 1074 [25] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The vir-
1075 tual atomic and molecular data centre (VAMDC) consortium. *Journal of*
1076 *Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- 1077 [26] Fang, H. (2018). A discussion of citations from the perspective of the
1078 contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–
1079 1520.
- 1080 [27] Force, M., Robinson, N., Matthews, M., Auld, D., and Boletta, M.
1081 (2016). Research data in journals and repositories in the web of science:

- 1082 Developments and recommendations. *Bulletin of IEEE Technical Com-*
1083 *mittee on Digital Libraries, Special Issue on Data Citation*, 12(1):27–30.
- 1084 [28] Garfield, E. (1999). Journal impact factor: a brief review. *Can. Med.*
1085 *Assoc.*, 979-980.
- 1086 [29] Gößwein, B., Miksa, T., Rauber, A., and Wagner, W. (2019). Data
1087 identification and process monitoring for reproducible earth observation
1088 research. In *2019 15th International Conference on eScience (eScience)*,
1089 pages 28–38. IEEE.
- 1090 [30] Green, T. J., Karvounarakis, G., and Tannen, V. (2007). Provenance
1091 semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-*
1092 *SIGART symposium on Principles of database systems*, pages 31–40. ACM.
- 1093 [31] Harding, S. D., Sharman, J. L., Faccenda, E., Southan, C., Pawson,
1094 A. J., Ireland, S., Gray, A. J. G., Bruce, L., Alexander, S. P. H., Anderton,
1095 S., Bryant, C., Davenport, A. P., Doerig, C., Fabbro, D., Levi-Schaffer, F.,
1096 Spedding, M., Davies, J. A., and Nc-Iuphar (2018). The IUPHAR/BPS
1097 guide to PHARMACOLOGY in 2018: updates and expansion to encom-
1098 pass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Re-*
1099 *search*, 46(Database-Issue):D1091–D1106.
- 1100 [32] Hartley, J. (2017). Authors and their citations: a point of view. *Scien-*
1101 *tometrics*, 110(2):1081–1084.
- 1102 [33] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a
1103 transformed scientific method.
- 1104 [34] Honor, L. B., Haselgrove, C., Frazier, J. A., and Kennedy, D. N. (2016).
1105 Data citation in neuroimaging: proposed best practices for data identifi-
1106 cation and attribution. *Frontiers in neuroinformatics*, 10:34.
- 1107 [35] Joshi-Tope, G., Gillespie, M., Vastrik, I., D’Eustachio, P., Schmidt, E.,
1108 de Bono, B., Jassal, B., Gopinath, G. R., Wu, G. R., Matthews, L., Lewis,
1109 S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of bio-
1110 logical pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- 1111 [36] Katz, D. (2014). Transitive credit as a means to address social and
1112 technological concerns stemming from citation and attribution of digital
1113 products. *Journal of Open Research Software*, 2(1).

- 1114 [37] Katz, D. S., Hong, N., Clark, T., Fenner, M., and Martone, M. (2020).
 1115 Software and data citation. *Computing in Science & Engineering*, 22 (2):4–
 1116 7.
- 1117 [38] Kosten, J. (2016). A classification of the use of research indicators.
 1118 *Scientometrics*, 108(1):457–464.
- 1119 [39] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S.
 1120 (2011). Citation and Peer Review of Data: Moving Towards Formal Data
 1121 Publication. *International Journal of Digital Curation*, 6(2):4–37.
- 1122 [40] Martone, M. (2014). Joint declaration of data citation principles.
 1123 *FORCE11. San Diego CA. Data Citation Synthesis Group*. [https://www.](https://www.force11.org/datacitationprinciples)
 1124 [force11.org/datacitationprinciples](https://www.force11.org/datacitationprinciples), online September 2020.
- 1125 [41] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation
 1126 counts and rankings of LIS faculty: Web of science versus scopus and
 1127 google scholar. *Journal of the american society for information science*
 1128 *and technology*, 58(13):2105–2125.
- 1129 [42] Milo, T. (2019). Getting rid of data. *Journal of Data and Information*
 1130 *Quality (JDIQ)*, 12(1):1–7.
- 1131 [43] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D.,
 1132 Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G.,
 1133 Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff,
 1134 D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,
 1135 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson,
 1136 E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C.,
 1137 Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers,
 1138 E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research
 1139 culture. *Science*, 348(6242):1422–1425.
- 1140 [44] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J.
 1141 (2016). Research data explored: An extended analysis of citations and
 1142 altmetrics. *Scientometrics*, 107(2):723–744.
- 1143 [45] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic, large
 1144 databases: Model and reference implementation. In *Proceedings of the*
 1145 *2013 IEEE International Conference on Big Data*, pages 307–312. IEEE.

- 1146 [46] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identifi-
 1147 cation of Reproducible Subsets for Data Citation, Sharing and Re-Use.
 1148 *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue*
 1149 *on Data Citation*, 12(1):6–15.
- 1150 [47] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data
 1151 citation of evolving data: Recommendations of the working group on data
 1152 citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- 1153 [48] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*
 1154 *Sci. Technol.*, 69(1):6–20.
- 1155 [49] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data
 1156 provenance in e-science. *SIGMOD Record*, 34(3):31–36.
- 1157 [50] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-
 1158 spective. In of Sciences’ Board on Research Data, N. A. and Information,
 1159 editors, *Report from Developing Data Attribution and Citation Practices*
 1160 *and Standards: An International Symposium and Workshop*, pages 177–
 1161 178. National Academies Press: Washington DC.
- 1162 [51] Ur Rehman, M. H., Liew, C. S., Abbas, A., Jayaraman, P. P., Wah,
 1163 T. V., and Khan, S. U. (2016). Big data reduction methods: a survey.
 1164 *Data Science and Engineering*, 1(4):265–284.
- 1165 [52] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,
 1166 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,
 1167 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data
 1168 management and stewardship. *Scientific data*, 3.
- 1169 [53] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data
 1170 citation: Giving credit where credit is due. In *Proceedings of the 2018*
 1171 *International Conference on Management of Data, SIGMOD*, pages 99–
 1172 114.
- 1173 [54] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to
 1174 scientific datasets using article citation networks. *Journal of Informetrics*,
 1175 14(2).
- 1176 [55] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of
 1177 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.

- 1178 [56] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for
1179 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*
1180 *Molecular Spectroscopy*, 327:122–137.