

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

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

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

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

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

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

1. A database cited as a whole, even though only parts of the databases 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

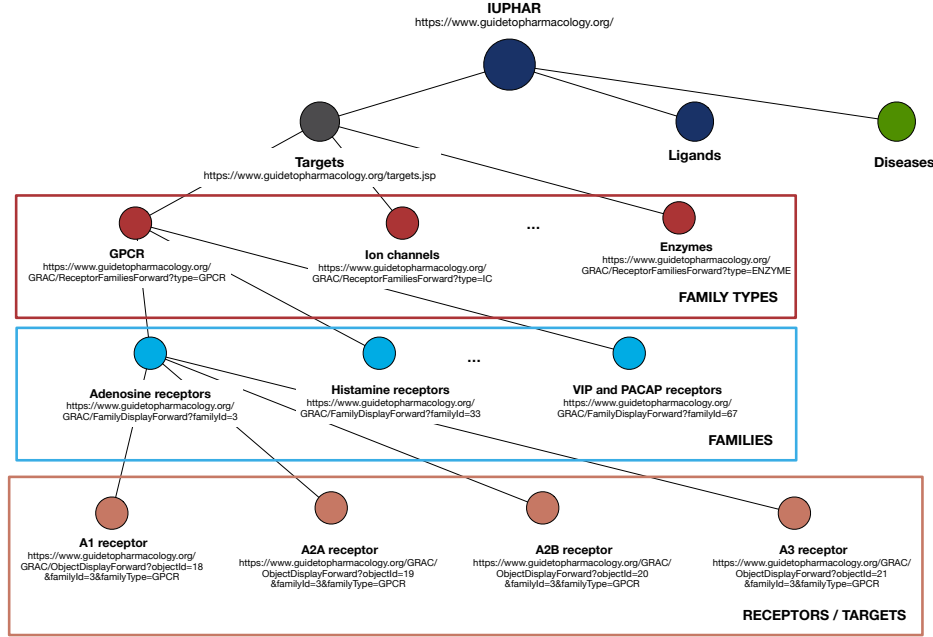


Figure 2: Partial map of the GtoPdb hierarchical structure grouping the targets into families and family types.

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/>

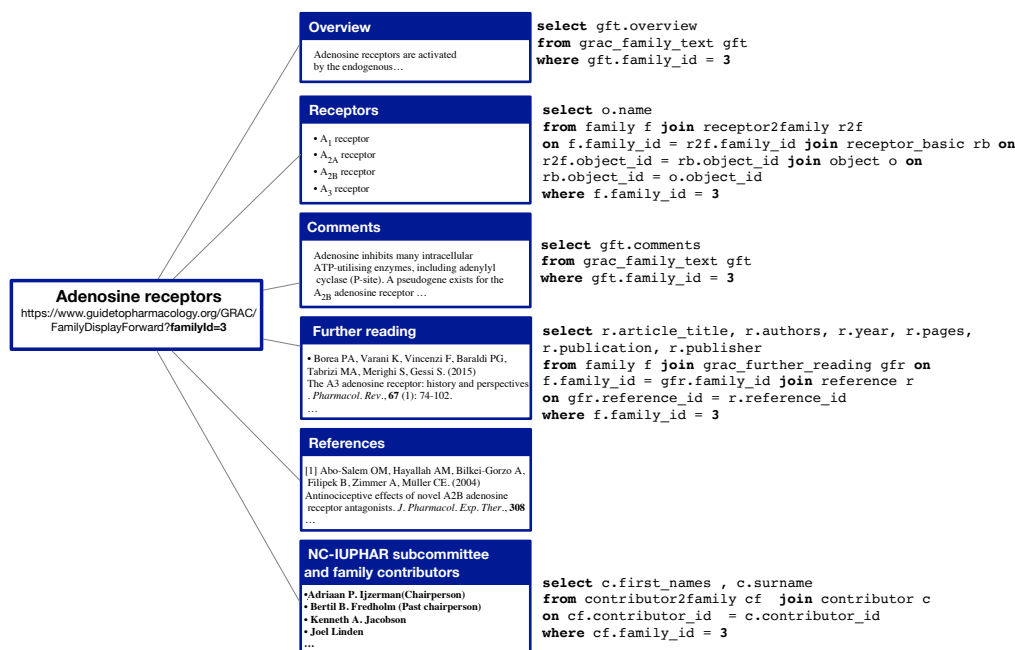


Figure 3: Basic web-page structure of “Adenosine receptors” family (ID 3), with queries used to retrieve the information contained in every section, except references.

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	$f_1 \cdot c_2 f_1 \cdot c_1 + f_1 \cdot c_2 f_2 \cdot c_2$
o_2	YANK Family	$f_4 \cdot c_2 f_4 \cdot c_1$

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.

A provenance polynomial may also present monomials that have exponents and coefficients different from one. One provenance polynomial can also be, for example, $3f_1 \cdot c_2 f_1 \cdot c_1 + f_1 \cdot c_2 f_2^3 \cdot c_2^3$. This is a polynomial of a tuple produced by a query where the result of the join between the tuples f_1 , $c_2 f_1$, and c_1 is produced three times and then merged (e.g. as the result of a projection), and the tuples $c_2 f_2$ 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

*** This whole section is heavily modified. *** 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.

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

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

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

535 **Definition 5.1. Tuple Level Data Credit Distribution (DCD) [24]**
 536 *Given a query Q over I and $k \in \mathbb{R}_{>0}$, DCD is defined as the computation*
 537 *of the function $f_{I,Q} : \text{TupleLoc} \times \mathbb{R}_{>0} \rightarrow \mathbb{R}_{\geq 0}$ such that $f_{I,Q}(t, k) = h$ where*
 538 *$0 \leq h \leq k$ and $\sum_{t \in \text{TupleLoc}} f_{I,Q}(t, k) = k$. The function $f_{I,Q}$ is the distribution*
 539 *strategy (DS).*

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

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

550 5.2. A Lineage-based Distribution Strategy

551 In the lineage-based distribution strategy, each tuple in the output of
 552 a query distributes credit equally to each input tuple that appears in its
 553 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:

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

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

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

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

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

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

583 5.3. A Why-Provenance-Based Distribution Strategy

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

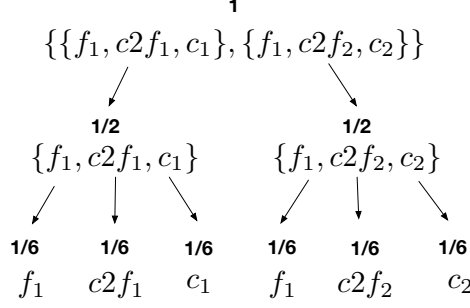


Figure 4: Distribution of credit using why-provenance-based DS for tuple o_1 .

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

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

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

Then tuple t in I receives credit equal to:

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

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

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

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

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

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

Figure 5: Example of provenance polynomial and the different notations used to define the how-provenance based distribution of Definition 5.4.

5.4. A How-Provenance Based Distribution Strategy

How-provenance conveys more information than why-provenance since it not only captures what tuples are relevant to the output and in which combination, but also how they are used. The “how” is captured through the provenance polynomials.

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

To define the distribution strategy based on the how-provenance, we introduce some preliminaries. Consider, as example, the provenance polynomial \mathcal{H} presented in Figure 5.

In this figure we show the notation that we use to refer to different information taken from the provenance polynomial. We call c the function that returns the sum of the coefficients of the polynomial. We use the same name for the function that, taken in input one monomial, in the example M_1 , outputs the sum of its exponents. mc is the function that takes in input a monomial and returns its coefficient. e is a function with parameters a monomial and a tuple, that returns the exponent of that tuple in the monomial, if present. γ takes in input a tuple and the whole polynomial, and returns a set containing the monomials containing that tuple, if present in the polynomial.

More formally, consider the provenance polynomial $\mathcal{H} = H(Q, I, o)$ of a tuple o . We define:

1. $c(\mathcal{H}) = n$ the function $c : \mathbb{N}[TupleLoc] \mapsto \mathbb{N}$ that, given a polynomial,

- 630 returns the sum of its coefficients;
- 631 2. $c(M)$ the function $c : \mathcal{M} \mapsto \mathbb{N}$ that, given a monomial M , returns the
632 sum of its exponents (with $\mathcal{M} \subset \mathbb{N}[TupleLoc]$ such that \mathcal{M} is made
633 only by the monomials M in $\mathbb{N}[TupleLoc]$);
- 634 3. $e(t, M)$ the function $e : TupleLoc \times \mathcal{M} \mapsto \mathbb{N}$ that, given in input a
635 tagged tuple and a monomial, returns the exponent of that tuple inside
636 the monomial;
- 637 4. $mc(M)$ the function $mc : \mathcal{M} \mapsto \mathbb{N}$ that, given in input one monomial,
638 returns its coefficient;
- 639 5. $\gamma(t, \mathcal{H})$ the function $\gamma : TupleLoc \times \mathbb{N}[TupleLoc] \mapsto \mathcal{M}$ that, given a
640 tuple t and a provenance polyomial \mathcal{H} , returns the (possibly empty)
641 set of monomials M in \mathcal{H} such that t appears in M .

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

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

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

645 Going back to the example of Table 4, consider o_1 with provenance poly-
646 nomial $f_1c_2f_1c_1 + f_1c_2f_2c_2$. The DS firstly divides the credit between the
647 two monomials. Since the coefficients of each monomial are 1, the credit is
648 split in half. If they were, for example, 1 and 2 respectively, 1/3 of the credit
649 would go to the first monomial, and 2/3 to the second. Since in our example
650 each variable has exponent 1, the credit is further divided equally among the
651 three variables. Thus, at the end of the computation, f_1 receives 1/3, and the
652 other tuples receive 1/6. If, for example, the first monomial was $f_1^2c_2f_1c_1$,
653 then the portion of credit of this monomial would be divided in this way:
654 1/2 to f_1 and 1/4 to each of the other two tuples.

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

659 Q2: SELECT DISTINCT F.name
660 FROM family as F JOIN
661 (SELECT DISTINCT f.name AS name

id	name
oxs_1	Dopamine Receptors

lineage	why-provenance	how-provenance
$\{f_1, c2f_1, c_1, c2f_2, c_2\}$	$\{\{f_1, c2f_1, c_1\}, \{f_1, c2f_2, c_2\}\}$	$f_1(f_1c2f_1c_1 + f_1c2f_2c_2)$

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.

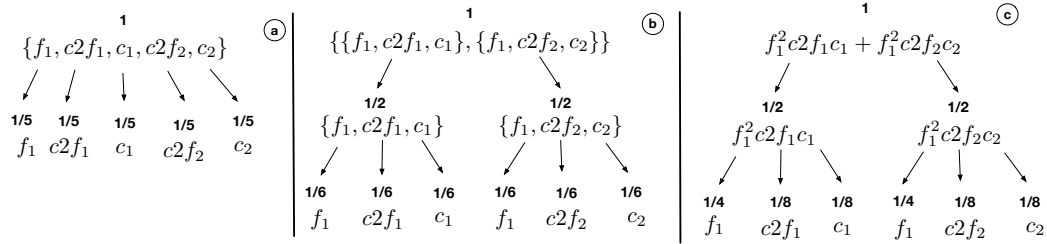


Figure 6: Comparison of different distributions strategies for tuple o_1 produced by query Q2.

```

662 FROM family AS f JOIN contributor2family AS c2f ON f.id = c2f.family_id
663 JOIN contributor AS c ON c2f.contributor_id = c.id
664 WHERE c.country = 'UK') AS R ON F.name = R.name
665 WHERE F.type = 'gpcr'

```

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

673 *** This polynomial still doesn’t have coefficients other than 1,**
674 **but ok. Why don’t you rewrite the polynomial to $f_1^2c2f_1c_1 + f_1^2c2f_2c_2$**
675 **to make the exponents clearer? ***

676 Figure 6 shows the differences between the three DS for the tuple o_1 of
677 Table 5. Subfigure 5.a uses lineage, sub-figure 5.b uses why-provenance, and
678 sub-figure 5.c uses how-provenance. The DS based on the provenance poly-
679 nomial gives credit 1/2 to f_1 , and 1/8 to the other tuples. This is reasonable
680 since Q2 relies on f_1 even more than Q1 does. The distribution based on

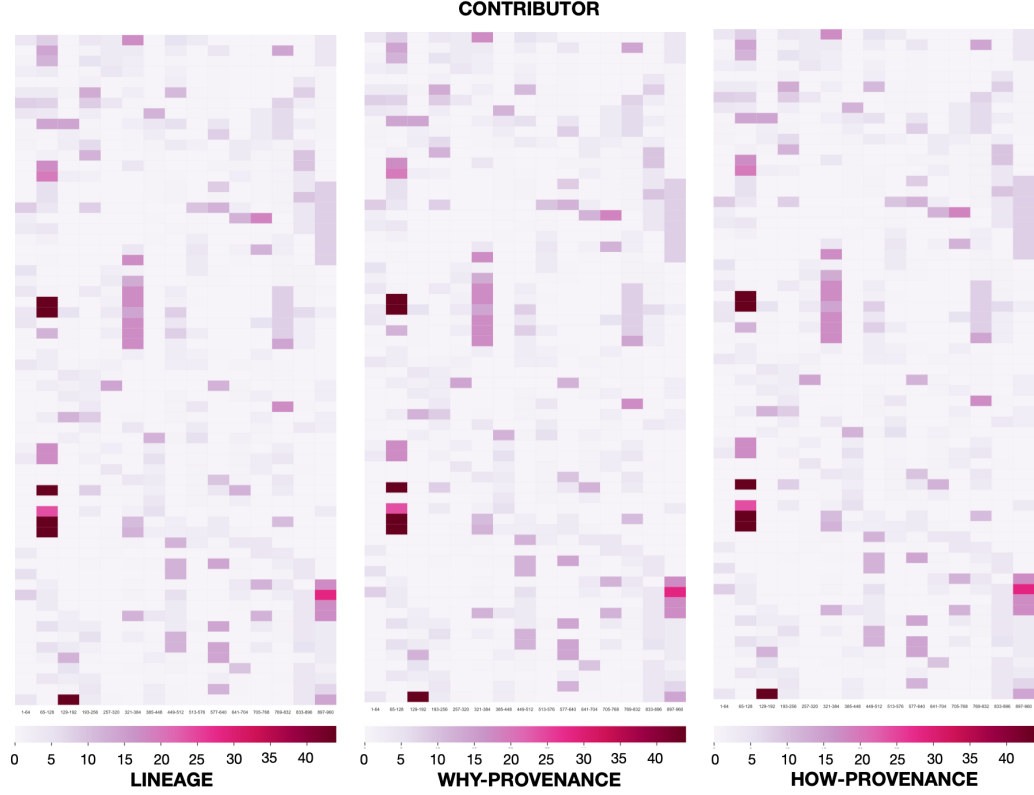


Figure 7: Comparison of three DS on the same table `contributor` using the distribution given by the queries retrieved from papers.

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: comparing provenances

We evaluate the proposed distribution strategies on GtoPdb, and in we focus on target families, all of those are described in webpages. GtoPdb particular identifies 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

692 family webpage of interest, or a subset of data extracted with a query. In this
 693 work, we consider a full-fledged data citation context in which papers cite
 694 the specific *data* subset of interest and not the webpage or the full database
 695 acting as data proxies. Therefore, when a paper cites family data, it is citing
 696 a set of queries needed to retrieve all the information provided by the family
 697 webpage, i.e., one query for each section composing a page, as depicted in
 698 Figure 3. In the figure, we can see how the structure of one family, “Adenosine
 699 receptors”, is mapped into several queries to obtain the information to build
 700 the corresponding webpage. In GtoPdb, all family pages share a similar
 701 structure (the only differences may be the presence/absence and length of
 702 the receptors lists, further readings, and contributors sections). Therefore,
 703 the same queries are used to build all other pages by simply changing the
 704 family id (which, in our example, is 3). All these queries are SPJ.

705 As already stated, many papers that draw information from the GtoPdb
 706 website¹¹ cite papers published every two years by the GtoPdb Committee on
 707 Receptor Nomenclature and Drug Classification (NC-IUPHAR). To obtain
 708 a set of citations capable of representing what happens, we consider a paper
 709 subset citing the 2018 GtoPdb [31] data paper. At the time of writing, this
 710 paper received more than 1200 citations.

711 As explained in Section 3, in the papers published in the British Journal of
 712 Clinical Pharmacology, that cite GtoPdb, the name of families are hyperlinks
 713 that point to the corresponding webpages. We considered all the 460 papers
 714 in BJCP citing [31] as of February 2020. We automatically extracted the
 715 URL references to family pages were automatically extracted to guide in
 716 building the queries to produce corresponding webpages. A total of 5,945
 717 different queries were built in this way.¹²

718 Figure 7 shows the heat-maps obtained by three different DS on the table
 719 **contributor**. It is immediately evident that the result of the distribution is
 720 the same with the three strategies. The same effect is also obtained in the
 721 other tables of the database used by the considered queries. Why is that? It
 722 is the case that the conditions in which we produced this experiment are quite
 723 peculiar. The queries that we used share similar characteristics. They are all
 724 SPJ queries, each of them utilizes each table only once in the join condition

¹¹<https://www.guidetopharmacology.org>

¹²For reproducibility purposes, the code we used for our experiments and all the produced queries can be found at the following link: https://bitbucket.org/dennis_dosso/credit_distribution_project.

(there are no self-joins), and all the joins are made using key attributes. In this particular condition, each tuple of the output presents: (i) a how-provenance that is a single monomial with coefficient 1 and exponent 1 in each variable; (ii) a why-provenance that is composed by only one witness; (iii) a lineage that coincides with the only witness in the witness basis. It is easy to see how, given these queries, the three distributions act in the same way. The credit is always uniformly distributed among the tuples appearing in each provenance.

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

```

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

```

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

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.

To give an example of how the CDS can differ from one another in their behavior, let us consider a different query:

```

753 Q4: SELECT f.name AS name
754 FROM family AS F JOIN
755 (SELECT DISTINCT f.family_id, f.name
756 FROM "family" AS f JOIN contributor2family AS cf ON
757 f.family_id = cf.family_id
758 JOIN contributor c ON
759 cf.contributor_id = c.contributor_id
760 WHERE c.country = 'UK') AS R

```

761 ON F.name = R.name

762 Here the innermost query retrieves all the names and ids of the families
 763 written by an author from the UK producing a relation called *R*. This
 764 relation is then joined with the table **family** on the attribute **name**.

765 One output tuple of this query is <Histamine receptors>, that has the
 766 following provenance polynomial:

$$f_{625}(f_{625}c_2f_{656}c_{184} + f_{625}c_2f_{113}c_{180} + f_{625}c_2f_{283}c_{198} + \\ + f_{625}c_2f_{550}c_{865} + f_{625}c_2f_{573}c_{101} + f_{625}c_2f_{95}c_{109})$$

767 As already discussed, the different monomials represent possible *alternatives*
 768 of combinations of tuples that produce the considered output tuple.
 769 Tuple f_{625} is used each time with different joins, thus it appears in each
 770 monomial. The last join, performed in the outmost query, is responsible
 771 for the final multiplication of f_{625} with the rest of the polynomial between
 772 parenthesis.

773 From this polynomial we compute the why-provenance as a set of six
 774 different witnesses:

$$\{\{f_{625}, c_2f_{656}, c_{184}\}, \\ \{f_{625}, c_2f_{113}, c_{180}\} \\ \{f_{625}, c_2f_{283}, c_{198}\}, \\ \{f_{625}, c_2f_{550}, c_{865}\}, \\ \{f_{625}, c_2f_{573}, c_{101}\}, \\ \{f_{625}, c_2f_{95}, c_{109}\}\}$$

775 And corresponding lineage:

$$\{f_{625}, c_2f_{656}, c_{184}, c_2f_{113}, c_{180}, c_2f_{283}, c_{198}, c_2f_{550}, c_{865}, c_2f_{573}, c_{101}, c_2f_{95}, c_{109}\}$$

776 This was only one tuple among the 86 obtained from this query. If we
 777 assign credit 1 to all these tuples and distribute it with the different strategies,
 778 we obtain the result shown in Figure 8 for the table **contributor**. At first
 779 sight, it may appear that the three distributions produce the same result.
 780 This is only partially true: the heat maps appear equal, but the absolute
 781 values assigned to each tuple are different. This is more evident if we look
 782 at the legend of each heat-map, where the maximum quantity of credit is
 783 different for each distribution. The one performed through lineage is around
 784 1.8, the why-provenance's one is around 1.4, and the one based on how-
 785 provenance is around 1.1.

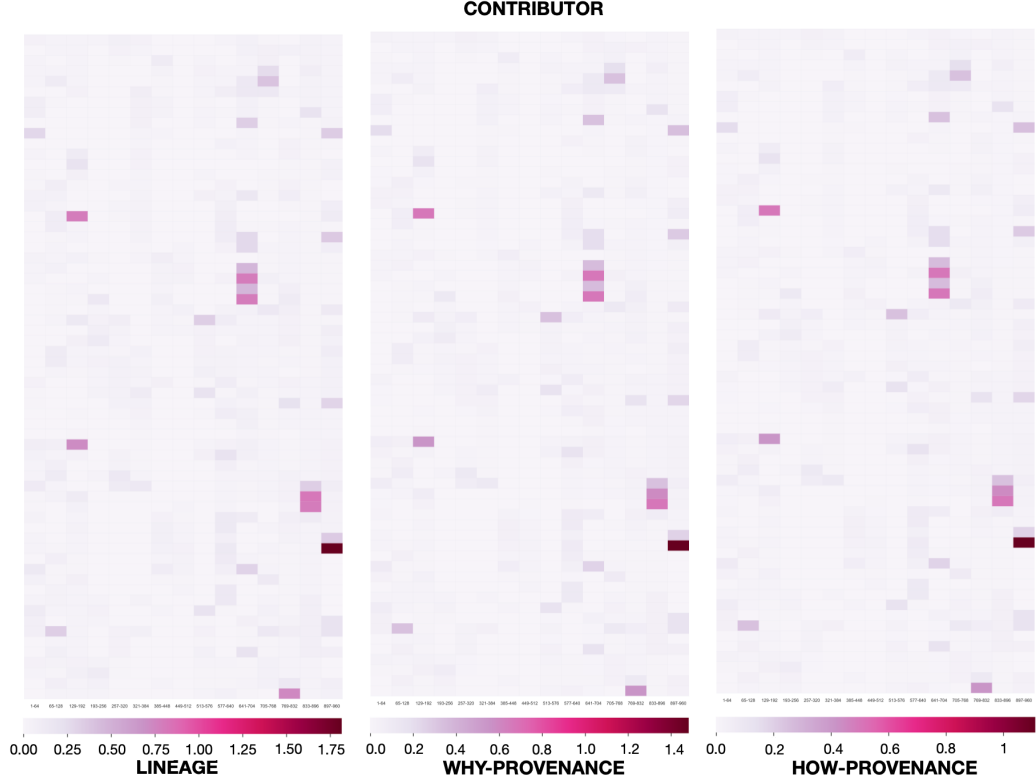


Figure 8: Comparison of three DS on the same table `family` after the distribution of the credit connected to query Q4.

786 To understand what is happening with this query in this specific ta-
 787 ble, consider the output tuple `<Histamine receptors>` and its provenances,
 788 as discussed above. Let us focus on its lineage. There are a total of six
 789 authors for this family and 13 tuples in total in the lineage. Thus, using
 790 the lineage-based DS, each tuple belonging to the `contributor` table (i.e.
 791 $c_{184}, c_{180}, c_{198}, c_{865}, c_{101}, c_{109}$) receives credit equal to $1/13$. Tuple f_{625} too re-
 792 ceives a portion of credit equal to $1/13$.

793 Let us consider now why-provenance. Tuple f_{625} appears six times in
 794 six different witnesses composed of 3 elements each. From each witness it
 795 receives a portion of credit equal to $1/18$, thus its total credit is $1/3$. On the
 796 other hand, all the authors appear only once in each witness, thus each of
 797 them receives credit $1/18$. In this case, why-provenance is recognizing more
 798 credit to tuple f_{625} , since it appears in each witness. The consequence is

799 that this distribution is equally *subtracting* credit from the other tuples in
800 the witnesses and giving it to f_{625} . In Figure 8 we are only looking at table
801 **contributor**. This same effect is reproduced for each tuple of the output of
802 query Q4, thus the *absolute* credit values on the tuples vary depending on the
803 deployed strategy. What happens is that the tuples in table **contributor**
804 receive less credit than the one received using lineage, but in the same pro-
805 portions. The heat map appears thus equal to the one obtained with lineage.
806 This same effect is also present with the how-provenance-based CDS. In this
807 case, tuple f_{625} is rewarded even more, since it appears with an exponent 2
808 in each monomial, thus attracting even more credit.

809 This is also why when we look at the legend for each part of Figure 8,
810 the maximum value reached with the lineage-based DS is higher than the
811 one reached with the why-provenance-based DS, which in turn is higher than
812 the one obtained with the how-provenance. This is because the different
813 strategies reward less and less the tuples of table **contributor** and more the
814 ones in table **family**.

815 This clearly shows the ability of the different strategies to adapt to sit-
816 uations. All three of them can highlight the relevant tuples in the table.
817 However, they differ in the way they reward the tuples. Depending on the
818 task, one provenance can be preferred to the other. If the only interest is
819 to highlight the relevant tuples, lineage is sufficient. If the interest is also
820 to reward more the tuples that are fundamental to the output, one can also
821 choose why- or how-provenance, knowing that how-provenance rewards even
822 more than why-provenance the relevant tuples that are indispensable for the
823 output.

824 Let us consider another interesting case we show in Figure 9. The figure
825 reports a distribution of credit performed on **family** through the generation
826 of *synthetic* polynomials. In this last case, we did not produce full-fledged
827 queries. Rather, we randomly generated provenance polynomials that might
828 be the how-provenance of randomly generated synthetic queries. An example
829 of such synthetic polynomial is:

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

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

$$\{\{f_1, c_2f_1, c_1\}, \{f_1, c_2f_2, cf_2\}, \{f_5, c_2f_{17}, c_{18}\}\}$$

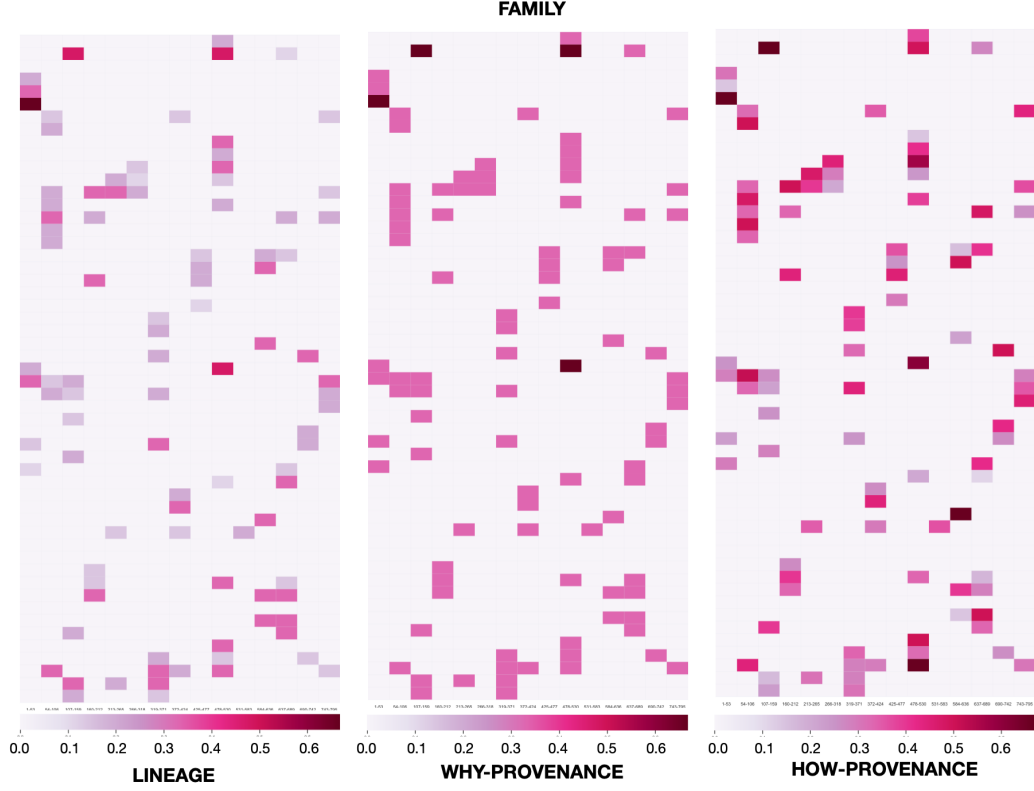


Figure 9: Comparison of three DS on the same table `family` after the distribution computed on provenances randomly generated.

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. They can be obtained by writing nested queries with join and union operations that use multiple times the same tuples (thus the presence of exponents bigger than 1) and that use the same combination of operations more than once (thus the presence of coefficients for monomials bigger than 1). We randomly generated a set of 100 such polynomials.

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

$$f_1 = \frac{59}{315}, f_5 = \frac{1}{18}, c2f_1 = \frac{2}{21}, c2f_2 = \frac{2}{15}, c2f_{17} = \frac{2}{9}, c_1 = \frac{2}{21}, c_2 = \frac{2}{15}, c_{17} = \frac{1}{6}$$

841 Using why-provenance, this is the output:

$$f_1 = \frac{2}{9}, f_5 = \frac{1}{9}, c2f_1 = \frac{1}{9}, c2f_2 = \frac{1}{9}, c2f_{17} = \frac{1}{9}, c_1 = \frac{1}{9}, c_2 = \frac{1}{9}, c_{17} = \frac{1}{9}$$

842 Finally, with lineage, this is the distribution:

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

843 To highlight how the distributions behave differently with these polynomi-
 844 als, consider tuple f_5 . f_5 receives the highest quantity of credit when we use
 845 the lineage-based distribution. Why-provenance and how-provenance reduce
 846 its quantity of credit since more information is available for the computation
 847 and the algorithms weigh less and less its role.

848 Generally speaking, the more complex the distribution, the more polar-
 849 ized the credit is toward the tuples that are used more frequently or with a
 850 higher impact in the production of the output tuple. Looking at the heat-
 851 maps of Figure 9, it appears that lineage tends to distribute credit more
 852 “equally” among the tuples, with only one or two tuples receiving higher
 853 quantities of credit, primarily because they are used in many different queries.

854 Why-provenance produces more tuples that are rewarded with high values
 855 of credit. Moreover, it appears that the other tuples that are not on the top
 856 of the spectrum are rewarded even more evenly compared to the DS based on
 857 lineage. That is, why-provenance, in this case, rewarded many tuples with
 858 roughly the same quantity of credit, and few tuples (but more compared to
 859 the DS based on lineage) with higher quantities of credit. This is due to
 860 the fact that why-provenance not only rewards the presence of a tuple in the
 861 computation but also the ways in which it is used.

862 How-provenance, finally, produces the distribution more sensible to the
 863 way a tuple is used in a query. Compared to the previous two DS, it also takes
 864 into consideration how many times a tuple is used, and weighs this factor
 865 in the distribution. It is interesting to see how certain tuples that received
 866 the lowest values of credit with lineage are now rewarded with higher values,
 867 showing that their fundamental role in certain queries outshines the fact that
 868 other tuples were used more frequently in the set of queries.

869 For our last set of experiments, consider Figure 10. We still use the 100
 870 polynomials described above and the credit distributed through them. Since
 871 these polynomials correspond to queries whose corresponding authors are not
 872 easily identifiable, we considered 20 “synthetic” authors, and we randomly

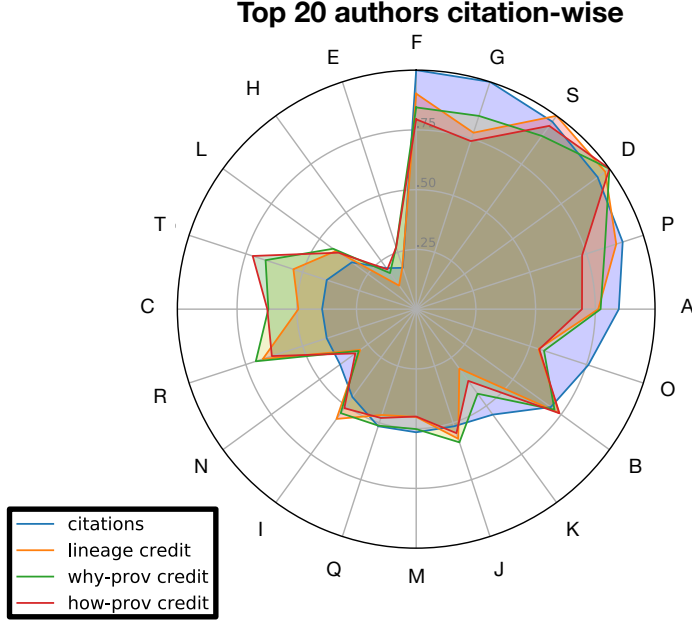


Figure 10: Top 20 authors by number of citations and their credit given through the three different DS.

873 assigned one author to each tuple in the database. The authors receive
874 “blocks” of consecutive tuples, with each block of the size varying between
875 10 and 40. Every time a tuple was used in a provenance polynomial, we
876 assigned one citation to the author corresponding to the tuple. The same
877 author also receives the three different credits assigned to the tuple at the
878 end of the distribution process using the three DS.

879 Figure 10 presents the radar plot where the 20 authors are sorted based on
880 the normalized number of received citations, together with the corresponding
881 normalized quantities of credits. Credit presents a different behavior from
882 one of the citations, and each form of credit, i.e., the credit obtained from
883 the different DS, behaves differently from the others. For example, it appears
884 that authors T, C, and R that are low in the number of citations are still
885 rewarded more than other more cited authors in terms of credit. Even if
886 the tuples of these authors received fewer citations, they still received more
887 credit than other more cited tuples. This shows how credit can be an effective
888 new method to use together with traditional citations to reward curators,
889 highlighting aspects lost using the traditional bibliometrics.

890 The three DS are all effective ways to distribute credit, and there is not
891 one distribution that is preferable to the other all the time. It all depends on
892 the needs of the users. Lineage is to be preferred when users only want to see
893 the tuples used in queries and reward more the tuples used in many queries.
894 It only rewards based on the *presence* of the tuples. Why-provenance is more
895 versatile when users also want to consider how many ways a tuple is used;
896 thus, in a way, its *versatility* inside the queries that used it. Finally, how-
897 provenance also counts how many times a tuple is used, its *frequency* in the
898 computation of a query.

899 7. Conclusions

900 This paper expanded on our previous work on data credit and data credit
901 distribution by defining two new distribution strategies, based on the why-
902 and how-provenance. The first distribution is based on the concept of witness,
903 and it can give more credit to tuples that appear in more than one witness.
904 In other words, tuples that are more important to the query and are used in
905 different ways by a query are also rewarded more by the distribution. The
906 second distribution, based on how-provenance, considers the frequency in
907 which a tuple or a combination of tuples is used in the query through the
908 provenance polynomial information. In this sense, it is even more sensitive
909 than the first one.

910 To show the differences between the three DS (also considering the one
911 based on lineage, defined in our previous work), we performed different ex-
912 periments on GtoPdb, a curated scientific relational database. In the first set
913 of experiments, we used SPJ queries extracted by data citations present in
914 papers published in the British Journal of Pharmacology. Employing these
915 queries, we were able to distribute the credit to the tuples in different tables
916 of the database, highlighting the tuples used more than others. We showed
917 that with these queries, the three strategies produce the same distribution.
918 With the specific type of queries that do not present self-joins, the formulas
919 at the base of the strategies have the same output. In this particular case,
920 the tuples are used in the same way by the queries; thus, the DSs do not
921 register any particular difference in the tuples' role.

922 In the second and third sets of experiments, we synthetically produced
923 more complex queries, i.e., nested queries whose provenance polynomials
924 presents coefficients and exponents bigger than 1. In this way, we showed
925 that, even though all three DS can highlight all the tuples used by the queries

926 in the database, the three have different behaviors. While the DS based on
927 lineage rewards all the tuples used by a query in equal measure, the strategy
928 based on why-provenance tends to reward the tuples more critical to the
929 query. In particular, why-provenance can consider the different ways in which
930 one tuple is used in a query. How-provenance is even more sensitive to the
931 tuples' role: it can also consider the frequency by which a tuple or a set of
932 tuples is used in the case of more complex queries. Depending on the goal of
933 a user, one provenance may be preferred to another.

934 In the fourth set of experiments, we showed how, compared with tra-
935 ditional citations, the credit distributed with the three strategies works as
936 a new tool highlighting different aspects of an author's role in the research
937 context identified by queries. Authors with a limited number of citations
938 can still have a high quantity of credit due to the importance of the data to
939 which they contributed to the queries.

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

950 Another application is *data reduction* [42], concerned with reducing the
951 vast mole of data that is produced in the evolving world of research and
952 information technology. Data reduction deals with different aspects of dealing
953 with huge amounts of data, such as finding reduced and relevant data streams
954 from the multiple gigabytes of data produced by big data systems every
955 second or dealing with the curse of dimensionality which requires unbounded
956 computational resources to uncover actionable knowledge patterns [51].

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

964 References

- 965 [1] Abadi, D., Ailamaki, A., Andersen, D., Bailis, P., Balazinska, M., Bern-
 966 stein, P., Boncz, P., Chaudhuri, S., Cheung, A., Doan, A., Dong, L.,
 967 Franklin, M. J., Freire, J., Halevy, A., Hellerstein, J. M., Idreos, S., Koss-
 968 mann, D., Kraska, T., Krishnamurthy, S., Markl, V., Melnik, S., Milo,
 969 T., Mohan, C., Neumann, T., Chin Ooi, B., Ozcan, F., Patel, J., Pavlo,
 970 A., Popa, R., Ramakrishnan, R., Ré, C., Stonebraker, M., and Suciu, D.
 971 (2020). The seattle report on database research. *SIGMOD Rec.*, 48(4):44–
 972 53.
- 973 [2] Alawini, A., Davidson, S. B., Hu, W., and Wu, Y. (2017). Automating
 974 data citation in citedb. *PVLDB*, 10(12):1881–1884.
- 975 [3] Alawini, A., Davidson, S. B., Silvello, G., Tannen, V., and Wu, Y.
 976 (2018). Data citation: A new provenance challenge. *IEEE Data Eng.*
 977 *Bull.*, 41(1):27–38.
- 978 [4] Altman, M., Borgman, C. L., Crosas, M., and Martone, M. (2015). An
 979 Introduction to the Joint Principles for Data Citation. *Bulletin of the*
 980 *Association for Information Science and Technology*, 41(3):43–45.
- 981 [5] Baggerly, K. (2010). Disclose all data in publications. *Nature*,
 982 467(7314):401–401.
- 983 [6] Bechhofer, S., Buchan, I. E., De Roure, D., Missier, P., Ainsworth, J. D.,
 984 Bhagat, J., Couch, P. A., Cruickshank, D., Delderfield, M., Dunlop, I.,
 985 Gamble, M., Michaelides, D. T., Owen, S., Newman, D. R., Sufi, S., and
 986 Goble, C. A. (2013). Why linked data is not enough for scientists. *Future*
 987 *Gener. Comput. Syst.*, 29(2):599–611.
- 988 [7] Belter, C. W. (2014). Measuring the Value of Research Data: A Citation
 989 Analysis of Oceanographic Data Sets. *PLoS ONE*, 9(3):e92590.
- 990 [8] Bertin-Mahieux, T., Ellis, D. P. W., Whitman, B., and Lamere, P. (2011).
 991 The million song dataset.
- 992 [9] Borgman, C. L. (2016). Data Citation as a Bibliometric Oxymoron. In
 993 Sugimoto, C. R., editor, *Theories of Informetrics and Scholarly Commu-*
 994 *nication*, pages 93–116. De Gruyter Mouton.

995 [10] Buneman, P. (2006). How to cite curated databases and how to make
 996 them citable. In *18th International Conference on Scientific and Statistical*
 997 *Database Management, SSDBM*, pages 195–203. IEEE Computer Society.

998 [11] Buneman, P., Christie, G., Davies, J. A., Dimitrellou, R., Harding, S. D.,
 999 Pawson, A. J., Sharman, J. L., and Wu, Y. (2020). Why data citation isn’t
 1000 working, and what to do about it. *Database*.

1001 [12] Buneman, P., Davidson, S. B., and Frew, J. (2016). Why data citation
 1002 is a computational problem. *Commun. ACM*, 59(9):50–57.

1003 [13] Buneman, P., Khanna, S., and Tan, W. C. (2001). Why and where: A
 1004 characterization of data provenance. In *Database Theory - ICDT 2001,*
 1005 *8th International Conference*, pages 316–330.

1006 [14] Buneman, P. and Silvello, G. (2010). A rule-based citation system for
 1007 structured and evolving datasets. *IEEE Data Eng. Bull.*, 33(3):33–41.

1008 [15] Callaghan, S., Donegan, S., Pepler, S., Thorley, M., Cunningham, N.,
 1009 Kirsch, P., Ault, L., Bell, P., Bowie, R., Leadbetter, A. M., Lowry,
 1010 R. K., Moncoiffé, G., Harrison, K., Smith-Haddon, B., Weatherby, a.,
 1011 and Wright, D. (2012). Making Data a First Class Scientific Output:
 1012 Data Citation and Publication by NERC’s Environmental Data Centres.
 1013 *International Journal of Digital Curation*, 7(1):107–113.

1014 [16] Candela, L., Castelli, D., Manghi, P., and Tani, A. (2015). Data Jour-
 1015 nals: A Survey. *Journal of the Association for Information Science and*
 1016 *Technology*, 66(9):1747–1762.

1017 [17] Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases:
 1018 Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–
 1019 474.

1020 [18] CODATA-ICSTI Task Group on Data Citation Standards and Practices
 1021 (2013). *Out of Cite, Out of Mind: The Current State of Practice, Policy,*
 1022 *and Technology for the Citation of Data*, volume 12.

1023 [19] Cousijn, H., Feeney, P., Lowenberg, D., Presani, E., and Simons, N.
 1024 (2019). Bringing citations and usage metrics together to make data count.
 1025 *Data Science Journal*, 18(1).

- 1026 [20] Cronin, B. (1984). *The citation process. The role and significance of*
1027 *citations in scientific communication*. London: Taylor Graham.
- 1028 [21] Cronin, B. (2001). Hyperauthorship: A postmodern perversion or evi-
1029 dence of a structural shift in scholarly communication practices? *JASIST*,
1030 52(7):558–569.
- 1031 [22] Cui, Y., Widom, J., and Wiener, J. L. (2000). Tracing the lineage of
1032 view data in a warehousing environment. *ACM Trans. Database Syst.*,
1033 25(2):179–227.
- 1034 [23] Davidson, S. B., Deutch, D., Milo, T., and Silvello, G. (2017). A model
1035 for fine-grained data citation. In *CIDR 2017, 8th Biennial Conference on*
1036 *Innovative Data Systems Research*. www.cidrdb.org.
- 1037 [24] Dosso, D. and Silvello, G. (2020). Data credit distribution: A
1038 new method to estimate databases impact. *Journal of Informetrics*,
1039 14(4):101080.
- 1040 [25] Dubernet, M. L., Antony, B. K., Ba, Y. A., et al. (2016). The vir-
1041 tual atomic and molecular data centre (VAMDC) consortium. *Journal of*
1042 *Physics B: Atomic, Molecular and Optical Physics*, 49(7):074003.
- 1043 [26] Fang, H. (2018). A discussion of citations from the perspective of the
1044 contribution of the cited paper to the citing paper. *JASIST*, 69(12):1513–
1045 1520.
- 1046 [27] Force, M., Robinson, N., Matthews, M., Auld, D., and Boletta, M.
1047 (2016). Research data in journals and repositories in the web of science:
1048 Developments and recommendations. *Bulletin of IEEE Technical Com-*
1049 *mittee on Digital Libraries, Special Issue on Data Citation*, 12(1):27–30.
- 1050 [28] Garfield, E. (1999). Journal impact factor: a brief review.
- 1051 [29] Gößwein, B., Miksa, T., Rauber, A., and Wagner, W. (2019). Data
1052 identification and process monitoring for reproducible earth observation
1053 research. In *2019 15th International Conference on eScience (eScience)*,
1054 pages 28–38. IEEE.
- 1055 [30] Green, T. J., Karvounarakis, G., and Tannen, V. (2007). Provenance
1056 semirings. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-*
1057 *SIGART symposium on Principles of database systems*, pages 31–40. ACM.

- [31] Harding, S. D., Sharman, J. L., Faccenda, E., Southan, C., Pawson, A. J., Ireland, S., Gray, A. J. G., Bruce, L., Alexander, S. P. H., Anderton, S., Bryant, C., Davenport, A. P., Doerig, C., Fabbro, D., Levi-Schaffer, F., Spedding, M., Davies, J. A., and Nc-Iuphar (2018). The IUPHAR/BPS guide to PHARMACOLOGY in 2018: updates and expansion to encompass the new guide to IMMUNOPHARMACOLOGY. *Nucleic Acids Research*, 46(Database-Issue):D1091–D1106.
- [32] Hartley, J. (2017). Authors and their citations: a point of view. *Scientometrics*, 110(2):1081–1084.
- [33] Hey, T., Tansley, S., and Tolle, K. M. (2009). Jim Gray on eScience: a transformed scientific method.
- [34] Honor, L. B., Haselgrove, C., Frazier, J. A., and Kennedy, D. N. (2016). Data citation in neuroimaging: proposed best practices for data identification and attribution. *Frontiers in neuroinformatics*, 10:34.
- [35] Joshi-Tope, G., Gillespie, M., Vastrik, I., D’Eustachio, P., Schmidt, E., de Bono, B., Jassal, B., Gopinath, G. R., Wu, G. R., Matthews, L., Lewis, S., Birney, E., and Stein, L. (2005). Reactome: a knowledgebase of biological pathways. *Nucleic Acids Research*, 33(Database-Issue):428–432.
- [36] Katz, D. (2014). Transitive credit as a means to address social and technological concerns stemming from citation and attribution of digital products. *Journal of Open Research Software*, 2(1).
- [37] Katz, D. S., Hong, N., Clark, T., Fenner, M., and Martone, M. (2020). Software and data citation. *Computing in Science & Engineering*, 22 (2):4–7.
- [38] Kosten, J. (2016). A classification of the use of research indicators. *Scientometrics*, 108(1):457–464.
- [39] Lawrence, B., Jones, C., Matthews, B., Pepler, S., and Callaghan, S. (2011). Citation and Peer Review of Data: Moving Towards Formal Data Publication. *International Journal of Digital Curation*, 6(2):4–37.
- [40] Martone, M. (2014). Joint declaration of data citation principles. *FORCE11. San Diego CA. Data Citation Synthesis Group*. doi: <https://doi.org/10.25490/a97f-egykh>, url: <https://www.force11.org/datacitationprinciples> (visited on 2020/03/17).

- 1091 [41] Meho, L. I. and Yang, K. (2007). Impact of data sources on citation
1092 counts and rankings of LIS faculty: Web of science versus scopus and
1093 google scholar. *Journal of the american society for information science*
1094 *and technology*, 58(13):2105–2125.
- 1095 [42] Milo, T. (2019). Getting rid of data. *Journal of Data and Information*
1096 *Quality (JDIQ)*, 12(1):1–7.
- 1097 [43] Nosek, B. A., Alter, G., Banks, G. C., Borsboom, D., Bowman, S. D.,
1098 Breckler, S. J., Buck, S., Chambers, C. D., Chin, G., Christensen, G.,
1099 Contestabile, M., Dafoe, A., Eich, E., Freese, J., Glennerster, R., Goroff,
1100 D., Green, D. P., Hesse, B., Humphreys, M., Ishiyama, J., Karlan, D.,
1101 Kraut, A., Lupia, A., Mabry, P., Madon, T., Malhotra, N., Mayo-Wilson,
1102 E., McNutt, M., Miguel, M., Paluck, E. L., Simonsohn, U., Soderberg, C.,
1103 Spellman, B. A., Turitto, J., VandenBos, G., Vazire, S., Wagenmakers,
1104 E. J., Wilson, R., and Yarkoni, T. (2015). Promoting an open research
1105 culture. *Science*, 348(6242):1422–1425.
- 1106 [44] Peters, I., Kraker, P., Lex, E., Gumpenberger, C., and Gorraiz, J.
1107 (2016). Research data explored: An extended analysis of citations and
1108 altmetrics. *Scientometrics*, 107(2):723–744.
- 1109 [45] Pröll, S. and Rauber, A. (2013). Scalable data citation in dynamic, large
1110 databases: Model and reference implementation. In *Proceedings of the*
1111 *2013 IEEE International Conference on Big Data*, pages 307–312. IEEE.
- 1112 [46] Rauber, A., Ari, A., van Uytvanck, D., and Pröll, S. (2016). Identifi-
1113 cation of Reproducible Subsets for Data Citation, Sharing and Re-Use.
1114 *Bulletin of IEEE Technical Committee on Digital Libraries, Special Issue*
1115 *on Data Citation*, 12(1):6–15.
- 1116 [47] Rauber, A., Asmi, A., van Uytvanck, D., and Proell, S. (2015). Data
1117 citation of evolving data: Recommendations of the working group on data
1118 citation (wgdc). *Result of the RDA Data Citation WG*, 20.
- 1119 [48] Silvello, G. (2018). Theory and practice of data citation. *J. Assoc. Inf.*
1120 *Sci. Technol.*, 69(1):6–20.
- 1121 [49] Simmhan, Y., Plale, B., and Gannon, D. (2005). A survey of data
1122 provenance in e-science. *SIGMOD Record*, 34(3):31–36.

- 1123 [50] Spengler, S. (2012). Data Citation and Attribution: A Funder’s Per-
 1124 spective. In of Sciences’ Board on Research Data, N. A. and Information,
 1125 editors, *Report from Developing Data Attribution and Citation Practices*
 1126 *and Standards: An International Symposium and Workshop*, pages 177–
 1127 178. National Academies Press: Washington DC.
- 1128 [51] Ur Rehman, M. H., Liew, C. S., Abbas, A., Jayaraman, P. P., Wah,
 1129 T. V., and Khan, S. U. (2016). Big data reduction methods: a survey.
 1130 *Data Science and Engineering*, 1(4):265–284.
- 1131 [52] Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G.,
 1132 Axton, M., Baak, A., Blomberg, N., Boiten, J., da Silva Santos, L. B.,
 1133 Bourne, P. E., et al. (2016). The fair guiding principles for scientific data
 1134 management and stewardship. *Scientific data*, 3.
- 1135 [53] Wu, Y., Alawini, A., Davidson, S. B., and Silvello, G. (2018). Data
 1136 citation: Giving credit where credit is due. In *Proceedings of the 2018*
 1137 *International Conference on Management of Data, SIGMOD*, pages 99–
 1138 114.
- 1139 [54] Zeng, T., Wu, L., Bratt, S., and Acuna, D. E. (2020). Assigning credit to
 1140 scientific datasets using article citation networks. *Journal of Informetrics*,
 1141 14(2).
- 1142 [55] Zou, C. and Peterson, J. B. (2016). Quantifying the scientific output of
 1143 new researchers using the zp-index. *Scientometrics*, 106(3):901–916.
- 1144 [56] Zwölf, C. M., Moreau, N., and Dubernet, M.-L. (2016). New Model for
 1145 Datasets Citation and Extraction Reproducibility in VADMC. *Journal of*
 1146 *Molecular Spectroscopy*, 327:122–137.