# On The Reasonable Effectiveness of Relational Diagrams\*

Explaining Relational Query Patterns and the Pattern Expressiveness of Relational Languages

WOLFGANG GATTERBAUER, (1) Northeastern University, USA CODY DUNNE, (1) Northeastern University, USA

Comparing relational languages by their logical expressiveness is well understood. Less well understood is how to compare relational languages by their ability to represent *relational query patterns*. Indeed, what are query patterns other than "a certain way of writing a query"? And how can query patterns be defined across procedural and declarative languages, irrespective of their syntax? To the best of our knowledge, we provide the first semantic definition of relational query patterns by using a variant of structure-preserving mappings between the relational tables of queries. This formalism allows us to analyze the *relative pattern expressiveness* of relational language fragments and create a hierarchy of languages with equal logical expressiveness yet different pattern expressiveness. Notably, for the non-disjunctive language fragment, we show that relational calculus can express a larger class of patterns than the basic operators of relational algebra.

Our language-independent definition of query patterns opens novel paths for assisting database users. For example, these patterns could be leveraged to create visual query representations that faithfully represent query patterns, speed up interpretation, and provide visual feedback during query editing. As a concrete example, we propose Relational Diagrams, a complete and sound diagrammatic representation of safe relational calculus that is provably (i) unambiguous, (ii) relationally complete, and (iii) able to represent all query patterns for unions of non-disjunctive queries. Among all diagrammatic representations for relational queries that we are aware of, ours is the only one with these three properties. Furthermore, our anonymously preregistered user study shows that Relational Diagrams allow users to recognize patterns meaningfully faster and more accurately than SQL.

### **ACM Reference Format:**

arXiv:2401.04758v1 [cs.DB] 9 Jan 2024

Wolfgang Gatterbauer and Cody Dunne. 2024. On The Reasonable Effectiveness of Relational Diagrams: Explaining Relational Query Patterns and the Pattern Expressiveness of Relational Languages. *Proc. ACM Manag. Data* 2, 1 (SIGMOD), Article 61 (February 2024), 71 pages. https://doi.org/10.1145/3639316

## 1 INTRODUCTION

When designing and comparing query languages, we are usually concerned with *logical expressive-ness*: can a language express a particular query we want? For relational languages, questions of expressiveness have been studied for decades, and formalisms for comparing expressiveness are well-developed and understood.

We do not yet have a similarly developed machinery to reason about *relational query patterns* across languages. Intuitively, a query pattern should capture "a certain way of writing a query." To be universally applicable, a formalization would have to be applicable across the four major

\*The title is a reference to Wigner's 1960 article [78] in which he states that "the enormous usefulness of mathematics in the natural sciences is something bordering on the mysterious and that there is no rational explanation for it." While there have been similar observations of surprising effectiveness for both data [45] and logic [46] in our community, our strong experimental evidence of Relational Diagrams helping users understand relational patterns better is actually quite expected.

Authors' addresses: Wolfgang Gatterbauer Northeastern University, USA, w.gatterbauer northeastern.edu; Cody Dunne Northeastern University, USA, c.dunne northeastern.edu.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2024 Copyright held by the owner/author(s). 2836-6573/2024/2-ART61

https://doi.org/10.1145/3639316

languages—Datalog, Relational Algebra (RA), Relational Calculus (RC), and SQL—and thus be orthogonal to questions of syntax and procedural or declarative language design.

We posit that identifying patterns in queries could open novel paths for assisting users [42], especially learners who try to understand the structure behind relational queries written in different languages. It could help learners spot similarities in queries across different schemas and thus more easily separate intent (the logic) from the particular syntactic expression. On an even more fundamental level, establishing a separation on "pattern expressiveness" between relational languages could lead to new insights into the intrinsic properties of relational languages and algebraic limits of visualizations. An important insight that we establish in this paper is that visual languages which build upon the operators of RA cannot faithfully express all query patterns, and instead necessitate reformulating queries and thus changing their patterns.

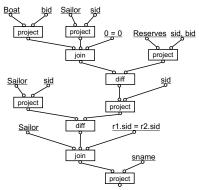


Fig. 1. DFQL [12] visualization of the TRC query from Example 1. Notice the 3 instances of the Sailor relation and thus a different "structure" of the visualization from the original query.

EXAMPLE 1 (UNDERSTANDING THE STRUCTURE OF A TRC QUERY). Imagine Kiyana, a theory-leaning undergraduate student, trying to understand relational query languages better. Kiyana has been reading the chapters on relational calculus across several books. In the textbook by Ramakrishnan and Gehrke [64] (page 121 of Sect. 4.3.1) she finds the query "(Q9) Find the names of sailors who have reserved all boats" written as follows:

$$\{P \mid \exists S \in Sailor \, \forall B \in Boat(\exists R \in Reserves \\ (S.sid = R.sid \land R.bid = B.bid \land P.sname = S.sname))\}$$
 (1)

She tries to understand "the structure" of the query and translates it first into RA, and then from there into DFQL (Dataflow Query Language) [12, 20, 44]. DFQL is a visual representation that is relationally complete by mapping its visual symbols to the operators of RA. Kiyana quickly notices that she cannot translate the query into RA without using additional Sailor relations.

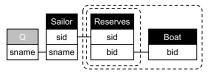
$$Q = \pi_{sname} (Sailor \bowtie (\pi_{sid}Sailor - \pi_{sid})((\pi_{sid}Sailor \times \pi_{bid}Boat) - \pi_{sid,bid}Reserves)))$$

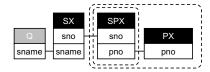
As a result, she does not find the resulting DFQL visualization (Fig. 1) very helpful because there is an obvious mismatch in "its structure" with 3 instances of Sailor relations. She wonders whether she is missing an obvious simpler translation into RA or whether there is none. As is, she does not find the query visualization helpful.

As a consequence, *no query visualization that relies on the operators of RA* could help Kiyana with what she would like to see: a simple visual representation that captures the structure of the query as written in its original logical form.

EXAMPLE 2 (COMPARING RC QUERIES FROM TEXTBOOKS). Kiyana continues looking through different textbooks and finds in Date's textbook [27] (page 224 of Sect. 8.3) the query "8.3.6 Get supplier names

<sup>&</sup>lt;sup>1</sup>DFQL is one of several visual query languages mentioned as relationally complete in an influential survey [20]. Kiyana found a detailed online documentation [44].





(a)  $q_1$ : "Find names of sailors who reserved all boats."

(b)  $q_2$ : "Find names of suppliers who supply all parts."

Fig. 2. Relational Diagrams representations of the two queries from Example 1 ([64]) and Example 2 ([27]). Notice the similar "relational query patterns."

for suppliers who supply all parts" written as follows:

From the natural language description, the query seems to follow a similar pattern as the earlier one ("Return X which have a relationship with all Y"). But that apparent similarity is difficult to see from the two expressions. She wonders whether there is a simple way to see that those two queries somehow follow a "similar structure."

In this paper, we show that there is indeed a simple and arguably-natural visualization that allows Kiyana to (*i*) represent her queries in a way that preserves their logical structure (or pattern), (*ii*) decide whether two logically-equivalent queries have the same pattern, and (*iii*) see whether any two queries, even across different schemas, use a "similar pattern." We call this visualization Relational Diagrams [67]. See Fig. 2 and notice how every relation from the two queries maps to exactly one relation in the Relational Diagrams. Also, notice how the similar structure of both queries becomes natural to see.

Our 1st contribution: query patterns. We develop a precise language-independent notion of relational query patterns that allows us to compare the patterns of two queries. Our definition is semantic (in the sense that the definition involves relations over sets of attributes) instead of syntactic (which would involve structural properties which are inherently language-dependent). The intuition behind our formalism is to reason about mappings between the (existentially or universally) quantified relations referenced in two queries. Yet it is not trivial to turn this intuition into a working definition that can be applied to any relational query and language (we include examples to show that seemingly easier mapping definitions would fail on queries). We believe that our notion is the "right" definition and show how to use it to compare relational query languages by their abilities to express query patterns present in other languages and thus compare their relative pattern expressiveness. In particular, we contribute a novel hierarchy of pattern-expressiveness among the non-disjunctive fragments of four relational query languages.

Our 2<sup>nd</sup> contribution: Relational Diagrams. We formalize an arguably simple and intuitive diagrammatic representation of relational queries called Relational Diagrams [67] and prove that (i) it is unambiguous (every diagram has a unique logical interpretation), (ii) it is relationally complete (every relational query can be expressed in a logically-equivalent Relational Diagram), and (iii) that it can express all query patterns in the non-disjunctive fragment of relational query languages and those with union at the root. In particular, we prove that no prior or future diagrammatic representation based on RA could represent all relational query patterns from RC. Our user study (Section 6.2) shows that our formalisms helps users recognize patterns faster than with SQL.

**Outline of the paper.** Section 2 defines the *non-disjunctive fragment* of relational query languages for Datalog, Relational Algebra (RA), Tuple Relational Calculus (TRC), and SQL, and proves

that they have equivalent logical expressiveness. Section 3 shows that the non-disjunctive fragment allows for an arguably natural diagrammatic representation system that we term Relational Diagrams\*. We give the formal translation from the non-disjunctive fragment of TRC to Relational Diagrams\* and back, and define their formal validity. We use Relational Diagrams\* for the remainder of the paper to illustrate "query patterns." Section 4 formalizes the notion of a relational query pattern and contributes a novel hierarchy of pattern expressiveness among the above four languages for the non-disjunctive fragment. We prove that Relational Diagrams have strong structure-preserving properties in that they can express all query patterns in this fragment. Section 4.4 formalizes "similar patterns" across different schemas. This extended notion allows us to see similarities across queries that use different relations and are thus not logically equivalent. Section 5 adds a single visual element (called a union cell) to Relational Diagrams\* to make the resulting Relational Diagrams relationally complete.<sup>2</sup> We also show that even without that element, Relational Diagrams\*can express all logical statements of first-order logic. This allows us to compare our diagrammatic formalism against a long history of diagrams for representing logical sentences. Section 6 includes two studies. The first shows that more logical queries across five popular textbooks have pattern-isomorphic representations in Relational Diagrams than either RA, Datalog, QBE, or QueryVis. The second controlled user experiment demonstrates that using Relational Diagrams instead of SQL helps users recognize patterns across different schemas faster and more often correctly. Section 7 contrasts our formalism with selected related work. In particular, we discuss the connection to Peirce's existential graphs [62, 69, 71] and show that our formalism is more general and solves interpretational problems of Peirce's graphs, which have been the focus of intense research for over a century.

Due to space constraints, we had to move proofs, several intuitive illustrating examples, all study details, and more detailed comparison against related work to the appendix.

## 2 THE NON-DISJUNCTIVE FRAGMENT OF RELATIONAL QUERY LANGUAGES

This section defines the *non-disjunctive fragment* of relational query languages. Throughout, we assume a linear order over the active domain and thus explicitly allow built-in predicates using ordered operators such as <, in addition to equality = and disequality  $\neq$ .

We assume the reader to be familiar with Datalog (non-recursive Datalog with negation), RA (Relational Algebra), TRC (safe Tuple Relational Calculus), SQL (Structured Query Language), and the necessary safety conditions for TRC and Datalog to be equivalent in logical expressiveness to RA. We also assume familiarity with concepts such as relations, predicates, atoms, and the named and unnamed perspective of relational algebra. The most comprehensive exposition of these topics we know of is Ullman's 1988 textbook [77], together with resources for translating between SQL and relational calculus [14, 29]. These connections are also discussed in most database textbooks [34, 37, 64, 72], though in less detail. We only cover TRC and not Domain Relational Calculus (DRC) as the 1-to-1 correspondence between DRC and TRC is straight-forward [34, 72], and—as we will discuss in Section 7.1—TRC has a more natural translation into diagrams than DRC.

## 2.1 Non-recursive Datalog with negation

We start with Datalog since the definition is most straightforward. Datalog expresses disjunction (or union) by repeating an Intensional Database (IDB) predicate in the head of multiple rules. For

<sup>&</sup>lt;sup>2</sup>Although disjunctions can be composed of conjunction and negation using De Morgan's law  $(A \lor B = \neg (\neg A \land \neg B))$ , this additional visual symbol is necessary: for safe relational queries, DeMorgan is not enough, as there is no way to write a safe Tuple Relational Calculus (TRC) expression "*Return all entries that appear in either R or S*" that avoids a union operator. This is part of the textbook argument for the union operator being an essential, non-redundant operator for relational algebra.

example, consider the following query in Datalog:

$$Q(x) := R(x, y), S(x), T(\_), y > 5.$$

$$Q(x) := R(x, y), S(\_), T(x), y > 5.$$
(3)

The underscore stands for a variable that appears only once [37]. This query cannot be expressed without defining at least one IDB at least twice, in our case the result table Q(x). This leads to a natural definition of the non-disjunctive fragment of Datalog  $\overline{}$ :

Definition 1 (Datalog\*). Non-disjunctive non-recursive Datalog with negation (Datalog\*) is the non-recursive fragment of Datalog¬ with built-in predicates where every IDB appears in the head of exactly one rule and can be used maximally once in any body.

Notice that Datalog\* inherits all restrictions from non-recursive Datalog¬ with built-in predicates [77], and thus rules out the existence of an IDB in both the head and the body of the same rule. The restriction of IDB's being used maximally once rules out views to be used multiple times (including simple copies of input tables).

## 2.2 Relational Algebra (RA)

We focus on the subfragment of basic RA  $(\times, \sigma, \bowtie_c, \pi, -)$  that contains no union operator  $\cup$  and in which all selection conditions are simple (i.e. they do not use the disjunction operator  $\vee$ ). A simple condition is  $C = (X\theta Y)$  where X is an attribute, Y is either an attribute or a constant, and  $\theta$  is a comparison operator from  $\{=, \neq, <, \leq, >, \geq, \}$ . Notice that conjunctions of selections can be modeled as concatenation of selections, e.g.,  $\sigma_{C_1 \wedge C_2}(R)$  is the same as  $\sigma_{C_1}(\sigma_{C_2}(R))$ . The Datalog query from (3) cannot be expressed in that fragment and requires either the union operator  $\cup$  as in:

$$\pi_A(\sigma_{B>5}(R) \bowtie S \times \rho_{A\to C}(T)) \cup \pi_A(\sigma_{B>5}(R) \bowtie T \times \rho_{A\to D}(S))$$

or the disjunction operator  $\vee$  as in:

$$\pi_A(\sigma_{A=D\vee A=C}(\sigma_{B>5}(R)\times\rho_{A\to D}(S)\times\rho_{A\to C}(T)))$$

Definition 2 ( $RA^*$ ). The non-disjunctive fragment of basic Relational Algebra ( $RA^*$ ) results from disallowing the union operators  $\cup$  and by restricting selections to conjunctions of simple predicates.

## 2.3 Tuple Relational Calculus (TRC)

Recall that safe TRC only allows existential quantification (and not universal quantification) [77]. Predicates are either join predicates " $r.A \theta s.B$ " or selection predicates " $r.A \theta v$ ", with r,s being table variables and v a domain value. WLOG, every existential quantifier can be pulled out as early as to either be at the start of the query, or directly following a negation operator. For example, instead of  $\neg(\exists r \in R[r.A = 0 \land \exists s \in S[s.B = r.B]])$  we rather write this sentence canonically as  $\neg(\exists r \in R, s \in S[r.A = 0 \land s.B = r.B])$ . This canonical representation implies that a set of existential quantifiers is always preceded by the negation operator, except for the table variables outside any scope of negation operators. Also, WLOG, we only allow equality conditions with the result table. For example, instead of  $\{q(A) \mid \exists r \in R, s \in S[q.A = r.A \land s.A > q.A])\}$  we rather write  $\{q(A) \mid \exists r \in R, s \in S[q.A = r.A \land s.A > r.A])\}$ . Recall that at least one equality predicate for each output attribute is required due to standard safety conditions [77].

We will define an additional requirement that each predicate contains a local (or what we refer to as *guarded*) attribute whose table is quantified within the scope of the last negation. For example, we do not allow  $\neg(\exists r \in R[\neg(r.A=0)])$  because the table variable r is defined outside the scope of the most inner negation around the predicate r.A=0. However, we allow the logically-equivalent  $\neg(\exists r \in R[r.A \neq 0])$  where the table variable r is existentially quantified within the same scope as the attribute  $r.A \neq 0$ .

*Definition 3 (Guarded predicate).* A predicate is *guarded* if it contains at least one attribute of a table that is existentially quantified inside the same negation scope as that predicate.

Intuitively, guarding a predicate guarantees that the predicates can be applied in the same logical scope where a table is defined. This requirement also avoids a hidden disjunction. To illustrate, consider the following TRC query:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg(\exists s \in S[r.A = 0 \land s.B = r.B])]\}$$

This query contains no apparent disjunction, however the predicate "r.A = 0" could be pulled outside the negation, and after applying De Morgan's law on the expression we get a disjunction:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land (r.A \neq 0 \lor \neg (\exists s \in S[s.B = r.B]))]\}$$

To avoid both disjunctions and "hidden disjunctions", the non-disjunctive fragment *only allows conjunctions of guarded predicates*:

Definition 4 (TRC\*). The non-disjunctive fragment of safe TRC (TRC\*) restricts predicates to conjunctions of guarded predicates.

In order to express the Datalog query from (3) we need the disjunction operator. A possible translation is:

$$\{q(A) \mid \exists r \in R, s \in S, t \in T[q.A = r.A \land r.B > 5 \land (r.A = s.A \lor r.A = t.A)]\}$$

### 2.4 SQL under set semantics

Structured Query Language (SQL) uses bag instead of set semantics and uses a ternary logic with NULL values. In order to treat SQL as a logical query language, we assume binary logic and no NULL values in the input database. It has been pointed out that "SQL's logic of nulls confuses people" and even programmers tend to think in terms of the familiar two-valued logic [76]. Our focus here is devising a general formalism to capture logical query patterns across relational languages, not on devising a visual representation of SQL's idiosyncrasies. To emphasize the set semantic interpretation, we write the DISTINCT operator in all our SQL statements.

We define the non-disjunctive fragment of SQL as the Extended Backus–Naur form (EBNF) [60] grammar shown in Fig. 3, interpreted under set semantics (no duplicates by using DISTINCT) and under binary logic (no null values allowed in the input tables). We also require the same syntactic restriction as for TRC\*: *every predicate needs to be guarded* (Definition 3), i.e., every predicate must reference at least one table within the scope of the last NOT. This restriction excludes hidden disjunctions, such as "NOT(NOT(P1) and NOT(P2))" which is equivalent to "P1 or P2".

Definition 5. SQL\*: Non-disjunctive SQL under set semantics (SQL\*) is the syntactic restriction of SQL under binary logic (no NULL values in the input tables) to the grammar defined in Fig. 3, and additionally requiring every predicate to be guarded.

Every SQL\* query can be brought into a canonical form that maintains a straightforward one-to-one correspondence with TRC\*. The idea is to replace membership and quantified subqueries with existential subqueries (see grammar in Fig. 3) and then unnest any existential quantifiers, i.e., to only use "not exists". This pulling up quantification as early as possible is identical to the way we defined the canonical form of TRC\*.

The Datalog query from (3) cannot be expressed in SQL\* and requires either a UNION operator or disjunction as in Fig. 4.

```
Q:= SELECT [DISTINCT] (C \{, C\} \mid *)
                                         Non-Boolean query
     FROM R {, R}
     [WHERE P]
     SELECT NOT (P)
                                         Boolean query
   | SELECT [NOT] EXISTS (Q)
                                         Boolean query
                                         column or attribute
C::= [T.]A
R := T [[AS] T]
                                         table (table alias)
P::= P {AND P}
                                         conjunction of predicates
     COC
                                            join predicate
     COV
                                            selection predicate
     NOT '('P')'
                                            negation
     [NOT] EXISTS '('Q')'
                                            existential subquery
     C [NOT] IN '('Q')'
                                            membership subquery
     C O (ALL '('Q')' | ANY '('Q')')
                                            quantified subquery
O::= = | <> | < | \le | \ge | >
                                         comparison operator
                                         table identifier
T::=
A::=
                                         attribute identifier
V::=
                                         string or number
```

Fig. 3. EBNF Grammar of  $SQL^*$ : Statements enclosed in [ ] are optional; statements separated by | indicate a choice between alternatives; parentheses without quotation marks () group alternative choices; parentheses with quotation marks '(' ')' form part of the test. Additionally, the main query requires the DISTINCT keyword (if non-Boolean), and all join and selection predicates need to be *guarded* (Definition 3), i.e., reference at least one table within the scope of the last NOT.

```
SELECT DISTINCT R.A
FROM R, S, T
WHERE R.B > 5
AND (R.A = S.A OR R.A = T.A)
```

Fig. 4. Example SQL with disjunction.

### 2.5 Logical expressiveness of the fragment

We show that the four languages restricted to the non-disjunctive fragment are equivalent in their logical expressiveness. The proof is available in Appendix C and is an adaptation of the standard proofs of equal expressiveness as found, for example, in [77]. However, the translations also need to pay attention to the restricted fragment (e.g. we cannot use union to define an active domain) and attempt to keep the numbers of extensional database atoms the same, if possible. This detail will be important later in Section 4, where we show that those four fragments differ in the types of query patterns they can express.

```
THEOREM 6. [Logical expressiveness] Datalog*, RA^*, TRC^*, and SQL^* have the same logical expressiveness.
```

## 3 RELATIONAL DIAGRAMS\*

This section introduces our diagrammatic representation of relational queries. It details the basic visual elements of Relational Diagrams\*(Section 3.1), gives the formal translation from TRC\* (Section 3.2) and back (Section 3.3), and shows that there is a one-to-one correspondence between TRC\* expressions and Relational Diagrams\*, thereby proving their validity (Section 3.4).

### 3.1 Visual elements

In designing our diagrammatic representation, we started from existing widely-used visual metaphors and then added the minimum necessary visual elements to obtain expressiveness for full  $TRC^*$ . In the following five points, we discuss both (i) necessary specifications for Relational

Diagrams\* and (*ii*) concrete design choices that are not formally required but justified based on best practices from HCI and visualization guidelines. We use the term *canvas* to refer to the plane in which a Relational Diagram\* is displayed. We illustrate with Fig. 5.

- (1) Tables and attributes: We use the set-of-mappings definition of relations [77] in which a tuple is a mapping from attributes' names to values, in contrast to the set-of-lists representation in which order of presentation matters and which more closely matches the typical vector representation. Thus a table is represented by any visual grouping of its attributes. We use the typical UML convention of representing tables as rectangular boxes with table names on top and attribute names below in separate rows. Table names are shown with white text on a black background and, to differentiate them, attributes use black text on a white background. For example, table with attribute A. Similar to Datalog and RA (and different from SQL and TRC), we do not use table aliases. Such table aliases create extra cognitive burden and are only needed in languages where references to repeated table instances cannot be otherwise disambiguated. We also reduce visual complexity by only showing attributes that are used in the query, similar to SQL and TRC (and different from Datalog). Database users are commonly familiar with relational schema diagrams. Thus, we believe that a simple conjunctive query should be visualized similarly to a typical database schema representation, as used prominently in standard introductory database textbooks [34, 72].
- (2) Selection predicates: Selection predicates are filters and are shown "in place." For example, an attribute " $r_2.C > 1$ " is shown as  $\boxed{C>1}$  in the corresponding instance of table R. An attribute participating in multiple selection predicates is repeated at least as many times as there are selections (e.g. to display " $r_2.C > 1 \land r_2.C < 3$ ", we would repeat R.C twice as  $\boxed{C>1}$  and  $\boxed{C<3}$ ). An attribute participating in k selection predicates is repeated k times.
- (3) Join predicates: Equi-join predicates (e.g. " $s_2.A = t_2.A$ "), which arguably are the most common type of join in practice, are represented by lines connecting joined attributes. For other less-frequent theta join operators  $\{\neq, <, \leq, \geq, >\}$ , we add the operator as a label on the line and use an arrowhead to indicate the reading order and correct application of the operator *in the direction of the arrow*. For example, for a predicate " $r_1.A > r_2.B$ ", the label is > and the arrow points from attribute A of the first R occurrence to B of the second:  $A \stackrel{>}{\rightarrow} B$ . Notice that the direction of arrows can be flipped, along with flipping the operator, while maintaining the identical meaning:  $A \stackrel{<}{\leftarrow} B$ . To avoid ambiguity with the standard left-to-right reading convention for operators, we normalize arrows to never point from right to left. An attribute participating in multiple join predicates needs to be shown only once and has several lines connecting it to other attributes. An attribute participating in one or more join predicates and k selection predicates is shown k+1 times.<sup>3</sup>
- (4) Negation boxes: In TRC\*, negations are either avoided (e.g.  $\neg(R.A = S.B)$ ) is identical to  $R.A \neq S.B$ ) or placed before the existential quantifiers. We represent a negation with a closed line that partitions the canvas into a subcanvas that is negated (inside the bounding box) and everything else that is not (outside of the bounding box). As a convention, we use dashed rounded rectangles.<sup>4</sup> Recursive partitioning of the canvas allows us to represent a tree-based nesting order that corresponds to the nested scopes of quantified tuple variables in TRC (and also the nesting order of subqueries in SQL). We call *the main canvas* the root of that nesting hierarchy and each node a *partition of the canvas*.

<sup>&</sup>lt;sup>3</sup>In practice, one can reduce the size of a Relational Diagram\* by reusing an existing selection predicate for joins. This comes at the conceptual complication that the exact graph topology of the Relational Diagram\* (which attributes are connected) is not uniquely determined (though it still allows only one correct interpretation). In our example Fig. 5d, one could remove the attribute R.C of  $r_2$  and connect Q.D to either C>1 or C<3 instead.

<sup>&</sup>lt;sup>4</sup>Rectangles allow better use of space than ellipses, and rounded corners together with dashed lines distinguish those negation boxes clearly from the rectangles with solid edges and right angles used for tables and attributes.

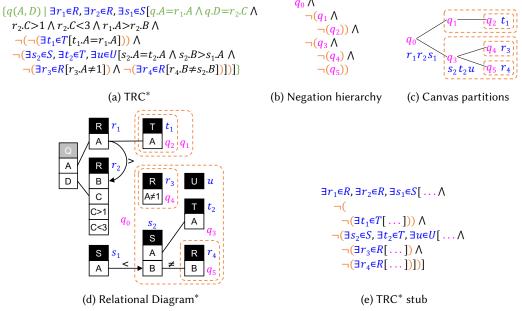


Fig. 5. Section 3.2: Example TRC\* expression (a), derivation of the negation hierarchy (b, c), and corresponding Relational Diagram\* (d). Colored partitions  $q_i$  (purple) and table variables  $r_i$  (blue) are not part of the diagrams and shown only to discuss the correspondence. Section 3.3: TRC\* stub after step 2 of the translation (e).

(5) Output table: We display an output table to emphasize the compositional nature of relational queries: a relational query uses several tables as input, and returns one new table as output. We use the same symbol for that output table as the TRC expression, for which we most commonly use Q. We use a gray background  $\bigcirc$  to make this table visually distinct from input tables.

## 3.2 From TRC\* to Relational Diagrams\*

We next describe the 5-step translation from any valid TRC\* expression to a Relational Diagrams\*. We illustrate by translating a TRC\* expression (Fig. 5a) into a Relational Diagrams\* (Fig. 5d). Notice that the translation critically leverages three conditions fulfilled by the input: (1) Safe TRC (and thus also TRC\*) only allows existential and not universal quantification [77], (2) TRC\* only allows conjunction between predicates, and (3) all predicates in TRC\* are guarded (recall Definition 3).

- (1) Creating canvas partitions: The scopes of the negations in a TRC are nested by definition. We translate this hierarchy of the scopes for each negation (the *negation hierarchy*) into a nested partition of the canvas. Fig. 5c illustrates the nested partitions as derived from the negation hierarchy Fig. 5b of the original TRC\* expression. Notice that the double negation " $\neg(\neg(...))$ " results in the scope  $q_1$  of the negation hierarchy to be empty.
- (2) Placing tables: Each table variable defines a table that gets placed into the canvas partition that corresponds to the respective negation scope. For example, the tables corresponding to the table variables  $r_1$ ,  $r_2$ , and  $s_1$  are outside any negation scope and thus placed in the root partition  $q_0$ . Similar to Datalog and RA (and in contrast to TRC and SQL) Relational Diagrams\* do not need table aliases.
- (3) *Placing selection predicates*: The predicates within each scope are combined via conjunction and are thus added one after the other. Since all selection predicates are guarded, the selection predicates

can be placed in the same partition as their respective table, which allows correct interpretation (see Section 3.3). For example, for  $\neg(\exists r_3 \in R[r_3.A \neq 1])$ , the predicate " $A \neq 1$ " is placed directly below R in  $q_4$ . An example of a predicate that is not guarded would be  $\exists r_3 \in R[\neg(r_3.A = 1)]$ : the scope of the negation contains a predicate of a table that is not existentially quantified in that scope.

- (4) Placing join predicates: For each join predicate, we add the two attributes (if not already present) and connect them via an edge with any comparison operator drawn in the middle. An attribute participating in multiple join predicates needs to be shown only once. Equi-joins are the standard and no operator is shown. Asymmetric joins include an arrowhead at one end of the edge (see Section 3.1). As for guarded join predicates one or both attributes are in the partition of a local table, the negation can be correctly interpreted. An example of an unguarded predicate would be  $\neg(r_4.B = s_2.B)$ . What is possible is the logically-equivalent  $r_4.B \neq s_2.B$  (as long as one of the two attributes is in the local scope of the last negation. In our example, this is the case in  $\neg(\exists r_4 \in R[r_4.B \neq s_2.B])$ .
- (5) Place and connect output table: The safety conditions for TRC [77] imply that output predicates can only be chosen from tables outside of all negations, thus in the root scope or partition  $q_0$ . If the query is non-Boolean, we add a new table named Q for query with a unique gray background to imply the difference from table references. If the query is Boolean, there is no output table and the query represents a logical sentence that is true or false.

**Completeness.** This five-step translation guarantees uniqueness of the following aspects: (1) nesting hierarchy (corresponding to the negation hierarchy), (2) where tables are placed (canvas partitions corresponding to the negation scope), (3) which attributes have selection predicates, and (4) which attributes participate in joins and how. The following aspects are not uniquely defined (without impact on the later interpretation): (1) the order of attributes below each table; (2) the direction of arrows can be flipped with a simultaneous label flip e.g.,  $s_1.A \leftarrow s_2.B$  and  $s_1.A \rightarrow s_2.B$  are identical (by convention we avoid arrows from right-to-left, but allow them up-to-down and down-to-up); (3) the size of visual elements and their relative arrangement; and (4) any optional changes in style (e.g. other than dashed negation boxes, distinct visual appearance between tables and attributes).

## 3.3 From Relational Diagram\*to TRC

We next describe the reverse five-step translation from any valid Relational Diagram\* to a valid and unique TRC\* expression. At the end, we summarize the conditions of a Relational Diagram\* to be valid, which are the set of requirements listed for each of the five steps. We again illustrate with the examples from Fig. 5.

- (1) Determine the nested scopes of negation: From the nested canvas partitions (Fig. 5c), create the nested scopes of the negation operators of the later TRC\* expression (Fig. 5b).
- (2) Quantification of table variables: Each table in a partition corresponds to an existentially-quantified table variable. WLOG, we use a small letter indexed by the number of occurrence for repeated tables. We add those quantified table variables in the respective scope of the negation hierarchy (Fig. 5c). For example, table T in  $q_2$  becomes  $\exists t_1 \in T[\ldots]$  and replaces  $q_2$  in Fig. 5e. Notice that partition  $q_1$  is empty and the resulting negation scope does not contain any expression other than another negation scope. We require that the leaves of the partition are not empty and contain at least one table. Otherwise, expressions  $\land \neg()$  and  $\land \neg(\neg())$  would both have to be true, leaving the meaning of an empty leaf partition ambiguous. This also implies that an empty canvas (there is only one partition, in which root and leaf are empty) is not a valid Relational Diagram\*.
- (3) Selection predicates: Selection attributes are placed into the scope in which its table is defined. For example, the predicate  $R.A \neq 1$  in partition  $q_4$  leads to  $\neg(\exists r_3 \in R[r_3.A \neq 1])$ .

(4) Join predicates: For join predicates (lines connecting attributes in Relational Diagrams\* with optional direction and operator), we have a validity condition that they can only connect attributes of tables that are in the same partition or different partitions that are in a direct-descendant relationship. In our example, T.A in  $q_2$  connects to R.A in  $q_0$  (here  $q_0$  is the root and grandparent of  $q_2$ .) However, we could not connect any attribute in  $q_5$  with any attribute in  $q_4$  (which are siblings in the nesting hierarchy). This requirement is the topological equivalent of scopes for quantified variables in TRC and guarantees that only already-defined table variables are referenced. Each such predicate is placed in the scope of the lower of the two partitions in the hierarchy, which guarantees the predicate to be guarded. For example, the inequality join connecting S.B in  $q_3$  and R.B in  $q_5$  is placed in the scope of  $q_5$ .

(5) Output table: The validity condition for the output table is that each of its one or more attributes is connected to exactly one attribute from a table in the root partition  $q_0$ . This corresponds to the standard safety condition of safe TRC. This step adds the set parentheses, the output tables, and its attribute and output predicates shown in green in Fig. 5a for non-Boolean queries.

**Soundness.** Notice that this five-step translation guarantees that the resulting TRC\* is uniquely determined up to (1) renaming of the tuple variables; (2) reordering the predicates in conjunctions, and (3) flipping the left/right positions of attributes in each predicate. It follows that Relational Diagrams\* are sound, and their logical interpretation is *unambiguous*.

## 3.4 Valid Relational Diagrams\*

In order for a Relational Diagrams\* to be valid we require that each of the conditions for the five-step translation process is fulfilled.

Definition 7 (Validity). A Relational Diagram\* is valid iff:

- (1) The nested hierarchy of optional negation boxes partitions the canvas (any two dashed boxes are either disjoint or one is completely contained within the other).
- (2) Each table, its attributes, and its selection predicates are discernible and reside in exactly one canvas partition.
- (3) Each leaf in the canvas partition contains at least one table.
- (4) Joins only happen between attributes of tables in partitions that are descendants (not siblings or their descendants). Join predicates with asymmetric operators such as < and > require a line with directionality (e.g. an arrowhead).
- (5) If there is an output table, then it has at least one attribute, and each attribute connects to exactly one attribute in the root partition  $q_0$  (safety condition of TRC).

THEOREM 8 (UNAMBIGUOUS RELATIONAL DIAGRAMS\*). Every valid Relational Diagram\* has an unambiguous interpretation in TRC\*.

The constructive translations from Sections 3.2 and 3.3 form the proof. Also notice that there is an additional validity condition we will add later in Definition 16 that will extend the logical expressiveness to include disjunction and go beyond TRC\*.

### 3.5 Logical statements

Boolean queries (or logical sentences) are formulas without free variables. Being able to express *relational sentences* (or constraints) allows us to compare our formalism against a long history of formalisms for logical statements [40]. An additional freedom with sentences is that the otherwise important safety conditions of relational calculus vanish. Thus, we need to be able to express statements that do not have any existentially-quantified relations in the main canvas. We next give an intuitive example, with more example given in Appendix E.

```
SELECT not exists

(SELECT *
FROM Sailor s
WHERE not exists
(SELECT b.bid
FROM Boat b, Reserves r
WHERE b.color = 'red'
AND r.bid = b.bid
AND r.sid = s.sid))

(a)

(Beserves Boat
bid
bid
bid
color = 'red'
```

Fig. 6. Example 3: All sailors reserve some red boat.

```
Example 3. Consider the statement: "All sailors reserve a red boat." \neg (\exists s \in Sailor[\neg (\exists b \in Boat, r \in Reserves[b.color = 'red' \land r.bid = b.bid \land r.sid = s.sid])]) \tag{4}
```

The first 4 steps of the translation in Section 3.2 still work: the root canvas  $q_0$  does not contain any relation (Fig. 6b). Similarly, the equivalent canonical SQL\* statement contains no FROM clause before the first NOT. Notice that Definition 12 of query pattern isomorphism still works as it is defined based on the relational tables.

## 3.6 A note on implementation

Creating valid Relational Diagrams\* programmatically requires a spatial layout algorithm that ensures that tables, predicates, and nested multi-layer canvas partitions are drawn unambiguously. To improve readability, it should also reduce edge crossings and edge bendiness. For initial work in that direction, please see our optimization model approach called STRATISFIMAL LAYOUT [30].

## **4 RELATIONAL QUERY PATTERNS**

Example 1 illustrated that—while two languages may well have the same logical expressiveness—one of them may have more ways to represent "logical patterns" than the other. We are interested in making this intuition more precise and *establishing a language-independent formalism* that captures the so-far vague notion of a relational query pattern. The formalism should allow us to study the "relative pattern expressiveness" of languages, i.e.: *Can languages*  $\mathcal{L}_2$  *express all patterns that language*  $\mathcal{L}_1$  *can?* We will then apply this formalism to the non-disjunctive fragment and compare the four previously-defined relational query languages by their relative abilities to represent "the same set of patterns" as other languages.

In the following, we often need to distinguish between a query as the *query expression* (the actual syntax in a particular query language) and a query as a *logical function* that maps a set of input tables to an output table. If we need to be precise, we refer to the function implied by a query as the *query semantics* and the actual syntax as the *query expression*. We use the word *signature* to refer to an ordered argument list as the input to a function, and use bracket notation for indexing. For example, the signature of f(x, y) is S = (x, y), and the first element is S[1] = x.

## 4.1 Defining relational query patterns

**Intuition.** Our goal is to define relational patterns in a way that allows us to analyze and compare *any relational query languages* irrespective of their syntax. Our idea is to formalize patterns based on the only common symbols in queries across languages: references to the *input relations* from the database. Since every relational query language needs to use input tables, the resulting formalism

generalizes. Intuitively, we will define two queries to be *pattern-isomorphic*<sup>5</sup> if there is a one-to-one correspondence that pairs each table in one query with a table in the other query that "plays the same semantic role." This means that when applying identical changes to these paired input tables (e.g. inserting a tuple), both queries will make identical changes to their outputs. However, for queries with multiple occurrences of the same input table (also called self-joins), we need to treat such repeated occurrences of the same input table as if they were *independent tables*. We will refer to such repeated table occurrences as "*table references*."

For example, consider  $q = R - (\pi_A R \times S)$ . The semantics of this query expression is a function q(R,S) that maps input relations R and S to an output table, and its signature would be just its relational input (R,S). However, we will not be interested in the signature of a query semantics, but rather the *signature of a query expression*, since we need to capture that two occurrences of R play different semantic roles in the query. In order to capture these different roles, we assign unique names to each table reference, resulting in what we call the *dissociated query*  $q' = R_1 - (\pi_A R_2 \times S)$ . We then formally define the *relational query pattern* of q as the logical function  $q'(R_1, R_2, S)$  expressed by the dissociated query q'. Notice that our definition of "dissociation" is inspired by, yet slightly different from its original use in the context of probabilistic inference [43] and the complexity of resilience and causal responsibility [35].

**Formalization.** To make our definitions precise across relational query languages with different syntax, we need to unambiguously refer to the individual occurrences of relational input tables in a given query expression, irrespective of the language used.

Definition 9 (Query signature). A table reference in a query expression q is any existentially or universally quantified reference to an input table. The signature S of q is the ordered list of its table references.

For example, the symbol "R" is a table reference in the SQL fragment "FROM R as R1", the TRC fragment " $\exists r_1 \in R$ ", the RA fragment " $\pi_A R$ ", and the Datalog fragment " $R(x, _)$ ".6 In contrast, the symbol "R" is not a table reference in the SQL fragment "WHERE R=1" as it is part of a reference to an attribute of a previously-defined table variable and not part of an existentially-quantified statement. The signature of a conjunctive SQL query with FROM clause "FROM R as R1, R as R2, S" is S = (R, R, S).

Definition 10 (Dissociated query). A dissociation of a query expression q with signature S is a modified query q' with S being replaced with a table signature S' of same size (i.e. |S'| = |S|), where every table in S' has a different name, and every table S'[i] has the same schema as table S[i] for all  $i \in [|S|]$ .

We call S' the dissociated signature of q. It is easy to dissociate a query by simply replacing duplicate names in S with fresh names. For simplicity, we will use subscripts when dissociating tables.

Example 4 (Dissociation). The RA query  $q = R - (\pi_A R \times S)$  has signature S = (R, R, S) with two of the three table references referring to the same input table R. Replacing the signature S with a dissociated signature  $S' = (R_1, R_2, S)$  leads to the dissociated query  $q' = R_1 - (\pi_A R_2 \times S)$ . Since the dissociated tables  $R_1, R_2$  inherit the schema information from table R, the dissociated query is still a

<sup>&</sup>lt;sup>5</sup>Recall that an *isomorphism* is a reversible *structure-preserving* mapping between two structures. For it to be reversible, it needs to be surjective (each element in the target is mapped to) and injective (different elements in the source map to different elements in the target) [36]. We use the term *pattern-isomorphic* (instead of structure-isomorphic) since our focus is on particular relational structures we refer to as patterns.

<sup>&</sup>lt;sup>6</sup>Existential quantification either happens explicitly as in TRC, or implicitly as in RA.

valid query and represents a new relational function  $q'(R_1, R_2, S)$  that maps three different input tables to an output.

The intuition behind this formalism is that the dissociated query defines a function that maps a set of *table references* (not just a set of input tables) to an output table. Thus, *the dissociated query is a semantic definition* of a relational query pattern across different relational query languages. Two queries use the same query pattern if their dissociated queries are logically equivalent, up to renaming and reordering of the input tables.

Definition 11 (Relational pattern). Given a query expression q with signature S. The relational pattern of q is the logical function defined by its dissociated query q'(S').

Definition 12 (Pattern isomorphism). Given two logically-equivalent queries  $q_1$  and  $q_2$  with signatures  $S_1$  and  $S_2$ , and dissociated queries  $q_1'(S_1')$  and  $q_2'(S_2')$ , respectively. The queries are pattern-isomorphic iff  $q_1'(S_1') = q_2'(\pi(S_1'))$  for some permutation  $\pi$ . In that case, we call the bijection  $S_1[i] \mapsto S_2[\pi(i)]$  between the query signatures a pattern-preserving mapping.

Example 5 (Patterns). Next, consider the Datalog query

$$I(x, y) := R(x, \_), S(y).$$
  
 $Q_1(x, y) := R(x, z), \neg I(x, y).$ 

with signature  $S_1 = (R, S, R)$ . Then its dissociated query is

$$I(x, y) := R_1(x, \_), S_2(y).$$
  
 $Q'(x, y) := R_3(x, z), \neg I(x, y).$ 

with signature  $S'_1 = (R_1, S_2, R_3)$ . Notice that Q' defines a logical function  $Q'(R_1, S_2, R_3)$  mapping two input tables with the same schema as R and an input table S to a binary output table.

Next, consider the RA query q from Example 4 with signature  $S_2 = (R, R, S)$ . Notice that above Datalog query Q and this RA query q are pattern-isomorphic since their dissociated queries define the same logical function up to permutation in the signatures:  $Q'_1(R_1, S_2, R_3) = q'(R_3, R_1, S_2) = q'(\pi(R_1, S_2, R_3))$  for permutation  $\pi = (2, 3, 1)$ . Thus the mapping  $(S_1[1], S_1[2], S_1[3]) \mapsto (S_2[2], S_2[3], S_2[1])$  is a pattern-preserving mapping.

Complexity of deciding pattern isomorphism. Deciding whether two relational queries are pattern-isomorphic is undecidable, in general (we need to determine whether two queries are equivalent, both before and after dissociation). This follows from Trakhtenbrot's theorem stating that the problem of validity in first-order logic on finite models is undecidable, and thus also the logical equivalence of relational queries (see, e.g., the reduction in [7]). However, we get a one-sided guarantee: if we can determine whether two queries are logically equivalent, then we can also determine whether they are pattern-isomorphic. In practice, equivalence of relational queries can often be determined, even for sophisticated SQL queries with grouping and aggregation evaluated over bags or sets [19].

## 4.2 Discussion with illustrating example

We next give an example of queries that have different query patterns although they are logically equivalent and have the same table signature. This detailed example motivates to a large extent why we define query patterns based on the dissociated signature.

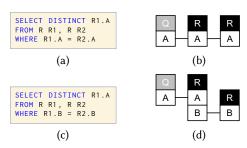


Fig. 7. Example 6: Two queries (a) and (c) with identical signatures that are logically equivalent but *not* pattern-isomorphic. Their associated Relational Diagrams are shown in (b) and (d), respectively.

EXAMPLE 6 (DIFFERENT PATTERNS). Consider table R(A, B) and the two Datalog queries  $Q_1(R)$  and  $Q_2(R)$  with

$$Q_1(x) := R(x, \_), R(x, \_).$$
  
 $Q_2(x) := R(x, y), R(\_, y).$ 

Both queries are logically equivalent to  $Q(x) := R(x, \_)$ , and thus also logically equivalent to each other. However,  $Q_1$  and  $Q_2$  represent different patterns:  $Q_1$  never uses the second attribute of R whereas  $Q_2$  uses that attribute to join both occurrences of R. This difference becomes even more apparent in SQL: Fig. 7a would even work if R was unary, whereas Fig. 7c requires R to be at least binary.

We next show that table dissociation allows us to formally distinguish the two patterns. First, notice that both queries have two occurrences of R as table references, and hence we need to associate each individual appearance to the "role" it plays semantically in the query. We achieve this by first dissociating the two occurrences of R into two fresh input tables (with the same schema). The two dissociated queries are  $Q_1'(R_1, R_2)$  and  $Q_2'(R_3, R_4)$  with

$$Q'_1(x) := R_1(x, \_), R_2(x, \_).$$
  
 $Q'_2(x) := R_3(x, y), R_4(\_, y).$ 

It is easy to verify that neither of the two possible mappings  $h_1 = \{(R_1, R_3), (R_2, R_4)\}$  and  $h_2 = \{(R_1, R_4), (R_2, R_3)\}$ , preserves logical equivalence for the dissociated queries. For example,  $R_1(1, 2), R_2(1, 3)$  returns an answer for  $Q_1'$  but not for  $Q_2'$ , under neither  $h_1$  nor  $h_2$ .

However,  $Q_1$  is pattern-isomorphic to the TRC query  $q_3(R)$  with

$$\{q_3(A) \mid \exists r_1 \in \mathbb{R}, r_2 \in \mathbb{R} [q.A = r_1.A \land r_1.A = r_2.A]\}$$

To see that, notice that its dissociated query  $q_3'(R_5, R_6)$  with

$$\{q_3'(A) \mid \exists r_1 \in \mathbb{R}_5, r_2 \in \mathbb{R}_6 [q.A = r_1.A \land r_1.A = r_2.A]\}$$

allows the isomorphism  $h_3 = \{(R_1, R_5), (R_2, R_6)\}$  from  $Q_1'$  to  $q_3'$  that preserves logical equivalence. By the same arguments,  $Q_1$  is pattern-isomorphic to the SQL query in Fig. 7a, and  $Q_2$  is pattern-isomorphic to the SQL query in Fig. 7c.

By design, our definition excludes views and intermediate tables such as Intensional Database Predicates in Datalog from the definition of table references. To see why, consider a query returning nodes that form the starting point of a length-3 directed path:

$$Q_1(x) := E(x, y), E(y, z), E(z, w).$$

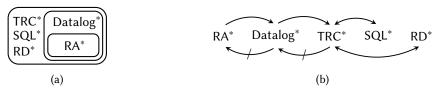


Fig. 8. Theorem 14: (a) A diagram summarizing the representation hierarchy between the non-disjunctive fragments of 4 query languages and Relational Diagrams\* (shown as RD\*). (b) Directions of pattern-preservation (and non-preservation) used in the proofs.

The following Datalog query uses the same logical pattern (find three edges that join and keep the starting node), even though it defines the intermediate intensional database predicate *I*:

$$I(y) := \underline{E}(y, z), \underline{E}(z, w).$$
  
$$Q'_1(x) := \underline{E}(x, y), I(y).$$

## 4.3 Comparing the "pattern-expressiveness" of relational languages

We next add the final definition needed to formally compare relational query languages based on their relative abilities to represent relational query patterns.

Definition 13 (Representation equivalence). We say that a query language  $\mathcal{L}_2$  can pattern-represent a query language  $\mathcal{L}_1$  (written as  $\mathcal{L}_1 \subseteq^{\text{rep}} \mathcal{L}_2$ ) iff for every legal query expression  $q_1 \in \mathcal{L}_1$  there is a pattern-isomorphic query  $q_2 \in \mathcal{L}_2$ . We call a query languages  $\mathcal{L}_2$  pattern-dominating another language  $\mathcal{L}_1$  (written as  $\mathcal{L}_1 \subseteq^{\text{rep}} \mathcal{L}_2$ ) iff  $\mathcal{L}_1 \subseteq^{\text{rep}} \mathcal{L}_2$  but  $\mathcal{L}_1 \not\supseteq^{\text{rep}} \mathcal{L}_2$ . We call  $\mathcal{L}_1, \mathcal{L}_2$  representation equivalent (written as  $\mathcal{L}_1 \equiv^{\text{rep}} \mathcal{L}_2$ ) iff  $\mathcal{L}_1 \subseteq^{\text{rep}} \mathcal{L}_2$  and  $\mathcal{L}_1 \supseteq^{\text{rep}} \mathcal{L}_2$ , i.e., both language can represent the same set of relational patterns.

We are now ready to state our result on the hierarchy of pattern expressiveness of the nondisjunctive fragment of the four languages defined earlier (Section 2) and our proposed relational diagrammatic representation Relational Diagrams\* (Section 3):

THEOREM 14 (REPRESENTATION HIERARCHY). 
$$RA^* \subseteq Pep$$
 Datalog $^* \subseteq Pep$   $TRC^* \subseteq Pep$   $SQL^* \subseteq Pep$  Relational Diagrams $^*$  (see Fig. 8).

The proofs is provided in Appendix F. In addition, Appendix G.1 shows how the separation between RA\* and Datalog\* disappears after adding the antijoin operator to the basic operators of RA\*, while the separation from TRC\* remains. The proof demonstrates that relational calculus has relational patterns that cannot be expressed in relational algebra. The important consequence is that RA\*, Datalog\* or any diagrammatic language modeled after them would not be a suitable target language for helping users understand all existing relational query patterns (including those used by SQL\*). Our related work (Section 7) shows that most existing visual query representations are modeled after relational algebra in that they model data flowing between relational operators, which implies they cannot faithfully represent all relational query patterns from TRC\* or SQL\*.

## 4.4 Similar patterns across schemas

We next extend the notion of pattern equivalence to allow comparing queries across different schemas. We call this concept "pattern similarity" and define it as a Boolean condition: two queries either have a similar pattern or not. The intuition is simple and best illustrated with the two queries from Fig. 2: As written those queries are not logically equivalent and thus they can't be pattern-isomorphic. However, if we first replace the table and attribute names from  $q_1$  with table

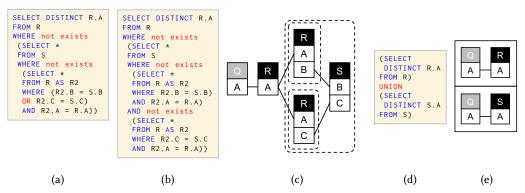


Fig. 9. Illustrations for Example 8 on replacing disjunctions: (a) SQL with disjunctions, (b) logically-equivalent (yet not representation-equivalent) SQL\* statement, and (c) Relational Diagrams. Illustrations for Example 9 on creating the union of queries: (d) union of SQL\* statements, and (e) Relational Diagrams with union cells.

names from  $q_2$  in a reversible (thus bijective) way, then the thus modified query  $q'_1$  would be pattern-isomorphic to  $q_2$ .

More formally, call a *schema mapping*  $\lambda$  from query  $q_1$  to  $q_2$ , a bijective mapping that replaces table names, attribute names, constants, and attribute order appearing in  $q_1$  with those from  $q_2$ .

Definition 15 (Similar Patterns). Given two queries  $q_1$  and  $q_2$ . The queries use a similar pattern iff there is a schema mapping  $\lambda$  from  $q_1$  to  $q_2$  s.t.  $\lambda(q_1)$  and  $q_2$  are pattern-isomorphic.

Example 7. For our example from Fig. 2, consider a mapping  $\lambda$  that replaces 'Sailor' with 'SX', 'Reserves' with 'SPX', 'Boat' with 'PX', 'sname' with 'sname', 'sid' with 'sno', and 'bid' with 'pno'. Then the thus modified query  $\lambda(q_1)$  is pattern-isomorphic with  $q_2$ .

#### 5 RELATIONAL COMPLETENESS

To make Relational Diagrams relationally complete, we now remove the non-disjunction restrictions and allow disjunctions and unions in all four relational query languages (Section 2). This means we must also add a corresponding syntactic device to Relational Diagrams that achieves logical equivalence to the other relational query languages. Unfortunately, this means that Relational Diagrams are no longer representation-equivalent to TRC. Can this be addressed in the future by a better diagram design? Based on the current understanding of the inherent limits of diagrams to express disjunctive information [70, 71] (see also the colored car example in Appendix K), such an extension would require adding *non-diagrammatic abstractions* (i.e. "syntactic devices").

The syntactic device that makes Relational Diagrams relationally complete is inspired by the representation of disjunction in Datalog. It was also proposed by Peirce in his discussion of Euler diagrams [62, 4.366] (see also [70, sect. 2.3.1]): we allow placing several Relational Diagrams\* on the same canvas, each in a separate *union cell*. Each cell of the canvas then represents only conjunctive information, yet the relation among the different cells is disjunctive (a union of the outputs).

We next illustrate with two examples logical transformations that are not pattern-preserving but that guarantee relational completeness. These transformations, together with union cells, make Relational Diagrams *relationally complete*: every query expressible in full RA, safe TRC, Datalog¬, or our prior SQL\* fragment extended by union and disjunctions of predicates<sup>7</sup> can then be represented

<sup>&</sup>lt;sup>7</sup>Extend the grammar from Fig. 3 with one additional rule: P::= '('P OR P')' and making adjustments for allowing the UNION clause between non-Boolean queries.

as a logically-equivalent Relational Diagram. The first example shows how to avoid disjunctions if they are not at the root level. The second shows how to replace disjunctions in the root by unions of queries.

Example 8 (Replacing disjunctions). Consider the SQL query from Fig. 9a which contains a disjunction and is not in SQL\*. Using De Morgan's Law with quantifiers  $\neg \exists r \in R[A \lor B] = \neg (\exists r \in R[A] \lor \exists r \in R[B]) = \neg \exists r \in R[A] \land \neg \exists r \in R[B]$ , we can first reformulate the conditions including disjunction as DNF, and then distribute the quantifier over the conjuncts. This leads to a disjunction-free query, yet leads to an increased number of table references:

```
 \{q(A) \mid \exists r \in R[q.A = r.A \land \neg(\exists s \in S \\ [\neg(\exists r_2 \in R[(r_2.B = s.B \lor r_2.C = s.C) \land r_2.A = r.A])])]\} 
 = \{q(A) \mid \exists r \in R[q.A = r.A \land \neg(\exists s \in S \\ [\neg(\exists r_2 \in R[r_2.B = s.B \land r_2.A = r.A]) \land \\ \neg(\exists r_3 \in R[r_3.C = s.C \land r_3.A = r.A])]\}
```

Fig. 9b shows this query as representation-equivalent SQL\* query, and Fig. 9c as Relational Diagram.

Example 9 (Union of queries). Consider two unary tables R(A) and S(A) and the TRC query

$$\{q(A) \mid \exists r \in R[q.A = r.A] \lor \exists s \in S[q.A = s.A]\}$$

We can write this query as a union of disjunction-free TRC\* queries:

$$\{q(A) \mid \exists r \in R[q.A = r.A]\} \cup \{q(A) \mid \exists s \in S[q.A = s.A]\}$$

Figure 9d shows a pattern-isomorphic SQL query, and Fig. 9e shows it as Relational Diagrams with two separate Relational Diagrams\* queries, each in a separate union cell, and each with the same attribute signature in the output table. This query cannot be rewritten without the union operator in RA, nor Relational Diagrams\* without union cells.

The additional validity criterion for multiple union cells follows the conditions of union or disjunction in the *named perspective* [1] of query languages: for disjunction in TRC, each operand needs to have the same arity, and the mapping between them is achieved by reusing the same variables.

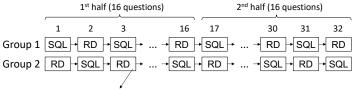
Definition 16 (Validity—extending Definition 7)).

(6) The output tables in multiple cells for the same query need to have the same name and same set of attributes.

THEOREM 17 (COMPLETENESS). Relational Diagrams (Relational Diagrams\* extended with union cells) are relationally complete.

The proof is in Appendix L. It uses the earlier proven logical expressiveness of Relational Diagrams\* and the fact that disjunctions can either be rewritten with DeMorgan or be pushed to the root. It also immediately follows that Relational Diagrams\* (without union cells) can already express any logical statement.

COROLLARY 18 (COMPLETENESS). Any logical statement in first-order logic can be expressed by a logically-equivalent Relational Diagrams\*.



- All participants see the j<sup>th</sup> question on the j<sup>th</sup> schema.
- The actual question they see is chosen from 4 patterns and 2 conditions.

Fig. 11. Illustration of the randomization and counterbalancing in our within-subjects study design.

#### 6 TWO APPLICABILITY STUDIES

## 6.1 Relational Diagrams for Textbook Queries

We analyzed the proportion of relational calculus queries encountered in learning scenarios that have pattern-isomorphic representations in either Relational Diagrams, RA, Datalog, QueryVis [26], or QBE [81]. For that purpose, we identify 59 queries across 5 popular textbooks with sections on relational calculus [22, 27, 34, 64, 72].



Fig. 10. Section 6.1: Fraction among 59 queries from 5 textbooks with pattern-isomorphic representations in listed languages.

Among those 59 queries, the number of queries that have pattern-isomorphic representations are 56 (94.9%) for Relational Diagrams, 53 (89.8%) for QueryVis, 49 (87.5%) for QBE, 48 (85.7%) for RA, and 47 (79.7%) for Datalog. The fraction for QueryVis happens to be identical to the *non-disjunctive* fragment. Standard Datalog cannot express disjunctions in the body of a query and thus performs worse than RA.<sup>8</sup> For QBE, notice that QBE 1) can express disjunctions within the same relations, yet 2) also requires the same safety conditions as Datalog. Furthermore, theta joins require the use of a non-diagrammatic conditions box [34, Appendix C]. RA extended with a primitive antijoin operator covers the same fraction as QBE. More details and all queries are given in Appendix N.

### 6.2 Controlled user study

We conducted a controlled experiment on Amazon Mechanical Turk (MTurk) [6] to evaluate the utility of Relational Diagrams for recognizing patterns. Our study investigates 3 main questions: (1) Can SQL users *identify common relational query patterns faster* using Relational Diagrams than SQL? (2) Can participants identify patterns faster over time, thus can users learn the patterns under repeated exposure to the same patterns? (3) Do participants have a similar accuracy (i.e. a comparable numbers of correct responses) using Relational Diagrams or SQL? We chose SQL as a baseline for comparison because we expect that fewer workers on MTurk understand TRC.

**Open practices.** Following best practices in user design, we preregistered the study design on OSF before collecting the data [41]. All code for generating the stimuli, the stimuli, the tutorial provided to participants, the resulting data (pilot n = 13, study n = 50), the analysis code, and changes from the preregistration are available on OSF [41]. More details on the study design and procedure are provided in Appendix O.

<sup>&</sup>lt;sup>8</sup>Modern variants of Datalog exist that can express disjunctions in the body, such as Souffle [74], but those are not standard.

Counterbalanced within-subjects study with randomization. We asked participants 32 questions, each having them identify which of 4 relational query patterns was presented. The exposure for each participant alternates between two conditions (Relational Diagrams and formatted SQL text). Each question uses a different relational schema found or adapted from common textbooks. We used counterbalancing and randomization to reduce ordering effects. Concretely, we start half the participants using SQL (group 1) and the rest using Relational Diagrams (group 2), after which participants alternate conditions with each question (see Figure 11). Moreover, we randomize the order in which patterns are presented such that each pattern-condition combination appeared twice in the first 16 questions ( $1^{st}$  half) and twice in the second 16 questions ( $2^{nd}$  half). The result is that each participant sees each of 4 patterns, in each of 2 conditions, exactly 2 times ( $16 = 4 \times 2 \times 2$ ), in each of the 2 halves (first and second 16 questions). This randomization helps reduce cheating as well as order effects. This setup necessitated creating 256 different stimuli (32 schemas  $\times$  4 patterns  $\times$  2 conditions), creation of which we semi-automated. The 4 patterns, illustrated with the sailors-reserve-boats schema, were:

- (1) Find sailors who have reserved some boat.
- (2) Find sailors who have not reserved any boat.
- (3) Find sailors who have not reserved all boats.
- (4) Find sailors who have reserved all boats.

We chose these patterns because we are interested in how Relational Diagrams can be used in educational settings, and they represent 4 different query structures. In particular, double-negation from pattern (4) is challenging for novice users to understand and is easy to misinterpret [57]. Moreover, pattern (4) does not have a pattern-isomorphic representation in RA.

Several prior user studies [55, 65, 66] have shown that diagrammatic representations of queries can help users understand queries faster. Key innovations in our design are: (1) We repeatedly expose participants to 4 identical patterns across 32 different schemas, allowing us to semi-automatically design 256 questions instead of a few hand-curated ones (each participant saw only 32, one for each schema). (2) Our questions are balanced across the first 16 and second 16 questions, which allows us to track learning over time. We are not aware of prior study design that allowed studying learning in an online study. (3) Our setup is randomized and parameterized, which creates a space of  $2 \cdot 2540^4$  possible treatments, i.e., participants are unlikely to see the same question sequence, reducing the chance of cheating. (4) We used monetary incentives for both time and accuracy, inspired by our recent work on QueryVis [55],

**Participants.** We conducted an n = 13 pilot study in the lab. After registering our study on OSF [41] and receiving approval from our Institutional Review Board (IRB), we began recruiting participants on MTurk [6]. Participants needed to have at least 500 completed tasks approved by requesters and at least 97% of their completed tasks approved. For us to approve their task, they needed to have at least 50% accuracy (i.e. answer at least 16 of our 32 questions correctly). Thus a participant who answered every question in SQL correctly and every question in Relational Diagrams incorrectly (or v.v.) would be included. Among the 120 task submissions, only 58 were approved. Our preregistration specified 50 participants, so we dropped 8 records and kept the counterbalancing intact by using the data from the first 25 participants who started in each condition.

**Tutorial.** Participants were given a self-paced 8-page tutorial on Relational Diagrams. The tutorial introduced our basic visual notations by showing SQL examples and their diagrams. The mean (resp. median) time spent on the tutorial and consent form was approximately 6.33 (respectively 3.5) minutes. The tutorial is available in our supplemental material [41].

**Analysis.** (1) As a within-subjects study, we first determined the *per-participant median time* for each condition. We used the median (instead of mean) for time because it is robust to outliers,

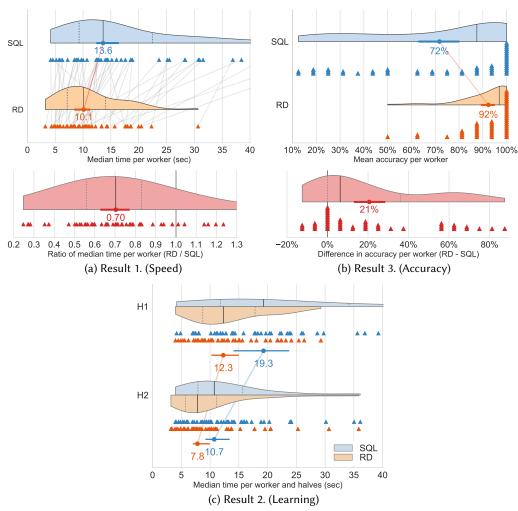


Fig. 12. User study: Triangles show median times per condition or mean accuracy for each of the n=50 successful participants. Violin plots show the data distribution, the median with a solid line, and the 25% and 75% quantiles with dashed lines. Error bars show the 95% BCa bootstrapped confidence intervals (CI) around the mean or median. Lines connect related marks. Relational Diagram is abbreviated here by RD.

although the median often requires more participants for the same statistical power. We computed the relative time Relational Diagrams/SQL needed per participant and again calculated the *median across* all participants. Here we used the median of the ratios as the median creates an unbiased estimator (in contrast to the mean of ratios, see Appendix O for details). (2) Likewise, we computed the median per-participant time for each condition spent on the  $1^{st}$  half (16 questions), on the  $2^{nd}$  half (16 questions), and the median ratio of the  $2^{nd}/1^{st}$  times. (3) We also computed the perparticipant accuracy for each condition and their difference. Then, across all participants, we calculated the mean of the differences in accuracy. Here we used the mean (instead of median) since the values are bounded within [0,1] (i.e. there are no outliers) and mean is more appropriate for discrete values like accuracies (i.e. 16/16, 15/16, ...). We analyzed these mean/median effect

sizes [25, 32] and used bias-corrected and accelerated (BCa) 95% confidence intervals (CIs) to show their range of plausible values [31, 33].

Results. We summarize 3 key takeaways. The executed Python notebook has more details [41].

**Result 1.** (Speed) We have strong evidence that participants were meaningfully faster at identifying patterns using Relational Diagrams than SQL: median ratio Relational Diagrams/SQL = 0.70, 95% CI [0.63, 0.77].

Our choice of visualization is a variant of Raincloud plots [5] and is inspired from recent work in the visualization literature [24] discussing various ways to juxtapose multiple visualizations ("clouds + rain + lightning") in the same chart for increasing information content. In that framework, each of our charts consists of (i) density plots that show an overview of the shape of the distribution (the "cloud"), (ii) unjittered dot plots that show the raw data (the "rain": here we deviate from [24] in using triangles instead of circles which, in our opinion, are more easily countable due to their visible vertices), and (iii) 95% confidence intervals that provide summary statistics (the "lightning"). Furthermore, whenever we compare alternative modalities ("repeated measures"), we also use (4) paired plots with lines connecting summary statistics and/or raw data.

Figure 12a (top) uses a paired plot to show the individual median per-participant times (and overall median across participants together with confidence interval) for both SQL (13.61, 95% CI [12.37, 16.43] in blue on the top) and Relational Diagrams (10.11 95% CI [8.38, 11.26] in orange on the bottom). Fig. 12a (bottom) shows the per-participant ratios between median times. Notice that the 95% confidence interval of the overall median [0.63, 0.77] does not overlap 1.00, which gives strong evidence for our conclusions.

Result 2. (Learning) Participants got meaningfully faster during the study in both conditions.

Figure 12c shows the individual times for H1 (1<sup>st</sup> half, i.e. the first 16 questions) and H2 (2<sup>nd</sup> half, i.e. the last 16 questions), for both SQL (in blue on the top) and Relational Diagrams (in orange on the bottom), together with medians and CIs. We see that the overall trend (Relational Diagrams being faster than SQL) is repeated across both halfs, and additionally that learning is taking place (participants need less time in H2 than in H1 in both conditions). The median ratios H1/H2 we used for inference (not shown but in our supplemental material) are 0.71, 95% CI [0.63, 0.79] for Relational Diagrams, and 0.70, 95% CI [0.51, 0.79] for SQL.

**Result 3.** (Accuracy) Participants were considerably more often correct with Relational Diagrams than with SQL: mean difference in accuracy Relational Diagrams – SQL = 21%, 95% CI [13%, 29%].

Figure 12b (top) shows that the per-participant accuracies and the overall mean accuracies were meaningfully higher with Relational Diagrams than with SQL. Notice that each participant answered 16 questions in each condition, thus possible scores are 16/16 = 1,  $15/16 \approx 0.94$ , etc. Thus accuracy per user and modality is discretized in multiples of 1/16 (in contrast to completion time, which is a continuous value and differs, even if slightly, between any two users). We thus use stacked triangles akin to a Wilkinson dot plot [79] to avoid overplotting and show individual data points. Figure 12b (bottom) shows the per-participant difference in mean accuracy. As the 95% CI of the overall mean [13%, 29%] does not overlap 0, we have strong evidence for our conclusion.

**Participant comments.** Participants could optionally write feedback at the end of the study. Few participants did, but those who did were encouraging, such as, "I found your diagrams very helpful in understanding the queries. At first I didn't get it, but after staring at the diagrams for a

few minutes it clicked and everything became super simple. I saw the patterns and it became just looking for the correct pattern to know which query was being used."

### 7 RELATED WORK

## 7.1 Peirce's beta Existential Graphs

Relational Diagrams represent nested quantifiers similarly as the influential and widely-studied *Existential Graphs* by Charles Sanders Peirce [62, 69, 71] for expressing logical statements (i.e. Boolean queries). Peirce's graphs come in two variants called alpha and beta. Alpha graphs represent propositional logic, and beta graphs represent first-order logic (FOL). Both variants use so-called *cuts* to express negation (similar to our negation boxes), and beta graphs use a syntactical element called the *Line of Identity* (LI) to denote both *the existence of objects* and *the identity between objects*.

**Differences.** The four key differences of beta graphs vs. Relational Diagrams are: (1) beta graphs can only represent sentences and not queries; (2) beta graphs cannot represent constants, thus selections cannot be modeled and instead require dedicated predicates; (3) beta graphs can only represent identity predicates (and no comparisons); and (4) Lines of Identity (LIs) in beta graphs have multiple meanings (existential quantification and identity between objects) and are a primary symbol. This *function overload of LIs* can make reading the graphs ambiguous. We, in contrast, have predicates inspired by TRC. Lines only connect two attributes and have no loose ends. Interpreting a graph as a TRC formula is straightforward and can be summarized in a simple set of rules (recall Section 3). We discuss this important conceptual difference in detail in Appendix P.1.

## 7.2 QueryVis

Some of our design decisions are similar to an earlier query representation called QueryVis [26, 38, 55]. In QueryVis diagrams, grouping boxes are used to group all tables within a local scope, i.e., for each individual query block. Those boxes thus cannot show their respective nesting, and an additional symbol of directed arrows is needed to "encode" the nesting. The high-level consequence of those design decisions is that (1) QueryVis does not guarantee to unambiguously visualize nested queries with nesting depth  $\geq$  4 (please see Appendix P.2. for a minimum example), (2) each grouping box needs to contain at least one relation (thus QueryVis cannot represent the query in Fig. 5), and (3) QueryVis cannot represent general Boolean sentences (e.g. the sentence "All sailors have reserved some red boat"). Thus QueryVis is not sound and not relationally complete, even for the disjunctive fragment.

## 7.3 Other relationally-complete formalisms

Appendix P.4 compares Relational Diagrams to other related visualizations like DFQL (Dataflow Query Language) [12, 20]. On a high level, all visual formalisms that we are aware of and that were proven to be relationally complete (including those listed in [12]) are at their core visualizations of relational algebra operators. This applies even to the more abstract graph data structures (GDS) from [11] and the later graph model (GM) from [13], which are related to our concept of query representation. The key difference is that GDS and GM are formulated inductively based on mappings onto operators of relational algebra. They thus mirror dataflow-type languages where visual symbols (directed hyperedges) represent operators like set difference connecting two relational symbols, leading to a new third symbol as output. Even QBE [81] uses the query pattern from RA and Datalog¬ of implementing relational division (or universal quantification) in a dataflow-type, sequential manner. Similarly, SIEUFERD [9], a direct manipulation spreadsheet-like interface, uses direct translation of relational algebra operators to prove SQL-92 completeness. This translation

<sup>&</sup>lt;sup>9</sup>Every beta graph has lines, and graphs with lines but no predicates have meanings. See, e.g., the definition in [71, p. 41].

involves expressing set difference with outer joins and "IS NULL" conditions. We have proved that there are simple queries in relational calculus that cannot be represented in relational algebra with the same number of relational symbols. Thus any visual formalism based on relational algebra cannot represent the full range of relational query patterns.

## 7.4 Other diagrammatic and non-diagrammatic query representations

Visual SQL [50] is a visual query language that also supports query visualization. With its focus on query specification, it maintains the one-to-one correspondence to SQL, and syntactic variants of the same query lead to different representations. SQLVis [59] shares motivation with QueryVis. Similar to Visual SQL, it places a stronger focus on the actual syntax of a SQL query and syntactic variants like nested EXISTS queries change the visualization, and join conditions are expressed as text. StreamTrace [10] focuses on visualizing temporal queries with workflow diagrams and a timeline. It is an example of visualizations for spatiotemporal domains and not the logic behind general relational queries. DataPlay [2, 3] allows a user to specify their query by interactively modifying a query tree with quantifiers and observing changes in the matching/non-matching data. It does not have a union operator and is thus not relationally complete. For a more detailed discussion we refer to two recent tutorials on visual representations of relational queries [39, 40].

### 8 CONCLUSIONS AND FUTURE WORK

We motivated a criterion called *pattern-isomorphism* that captures the patterns across relational languages and gave evidence for its importance in designing diagrammatic representations. We formulated the non-disjunctive fragments of Datalog¬, RA, safe TRC, and corresponding SQL (interpreted under set semantics) that naturally generalize conjunctive queries to nested queries with negation. We prove that this important fragment allows a rather intuitive and, in hindsight, natural diagrammatic representation that can preserve the query pattern used across all four languages. We further prove that this formalism, extended with a representation of union, is complete for full safe relational calculus (though not pattern-preserving) and showed via user studies strong evidence that this diagrammatic representation allows users to understand query patterns faster and more accurately than SQL, even with minimal training.

Finding a pattern-preserving diagrammatic representation for disjunction and even more general features of SQL (such as grouping and aggregates) is an open problem. For example, it is not clear how to achieve an intuitive and principled diagrammatic representation for arbitrary nestings of disjunctions, such as " $R.A < S.E \land (R.B < S.F \lor R.C < S.G)$ " or " $(R.A > 0 \land R.A < 10) \lor (R.A > 20 \land R.A < 30)$ " with minimal additional notations. Grounded in a long history of diagrammatic representations of logic, we gave intuitive arguments for why visualizing disjunctions is inherently more difficult than conjunctions, with some experts believing it is not possible [70, 71] unless one adds non-diagrammatic abstractions.

### **ACKNOWLEDGMENTS**

This work was supported in part by a Khoury seed grant program, and the National Science Foundation (NSF) under award numbers IIS-1762268, IIS-2145382, and IIS-1956096. It was conducted in part while WG was visiting the Simons Institute for the Theory of Computing. We like to thank Mirek Riedewald for helpful comments on an early version of this paper, and Jan Van den Bussche for insightful comments about an earlier version of the proof of separation lemma 20.

#### REFERENCES

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. Foundations of Databases. Addison-Wesley. http://webdam.inria.fr/Alice/
- [2] Azza Abouzied, Joseph M. Hellerstein, and Avi Silberschatz. 2012. DataPlay: interactive tweaking and example-driven correction of graphical database queries. In UIST. ACM, 207–218. https://doi.org/10.1145/2380116.2380144
- [3] Azza Abouzied, Joseph M. Hellerstein, and Avi Silberschatz. 2012. Playful Query Specification with DataPlay. PVLDB 5, 12 (2012), 1938–1941. https://doi.org/10.14778/2367502.2367542
- [4] Javad Akbarnejad, Gloria Chatzopoulou, Magdalini Eirinaki, Suju Koshy, Sarika Mittal, Duc On, Neoklis Polyzotis, and Jothi S. Vindhiya Varman. 2010. SQL QueRIE Recommendations. PVLDB 3, 1 (2010), 1597–1600. https://doi.org/10.1145/2839509.2844640
- [5] Micah Allen, Davide Poggiali, Kirstie Whitaker, Tom Rhys Marshall, Jordy van Langen, and Rogier A. Kievit. 2019. Raincloud Plots: a multi-platform tool for robust data visualization. Wellcome Open Research 4 (2019). https://doi.org/10.12688/wellcomeopenres.15191.2
- [6] Amazon Mechanical Turk (MTurk). 2023. https://www.mturk.com.
- [7] Marcelo Arenas, Pablo Barcelo, Leonid Libkin, Wim Martens, and Andreas Pieris. 2022. *Database Theory: Querying Data.* Open source at https://github.com/pdm-book/community.
- [8] Peter C. Austin and Janet E. Hux. 2002. A brief note on overlapping confidence intervals. *Journal of Vascular Surgery* 36, 1 (2002), 194–195. https://doi.org/10.1067/mva.2002.125015
- [9] Eirik Bakke and David R. Karger. 2016. Expressive Query Construction through Direct Manipulation of Nested Relational Results. In SIGMOD. ACM, 1377–1392. https://doi.org/10.1145/2882903.2915210
- [10] Leilani Battle, Danyel Fisher, Robert DeLine, Mike Barnett, Badrish Chandramouli, and Jonathan Goldstein. 2016. Making Sense of Temporal Queries with Interactive Visualization. In CHI. ACM, 5433-5443. https://doi.org/10.1145/2858036.2858408
- [11] Tiziana Catarci. 1991. On the Expressive Power of Graphical Query Languages. In Visual Database Systems, II.

  Proceedings of the IFIP TC2/WG 2.6 Second Working Conference on Visual Database Systems. (IFIP Transactions, Vol. A-7).

  North-Holland, 411–421. https://dblp.org/rec/conf/vdb/Catarci91
- [12] Tiziana Catarci, Maria Francesca Costabile, Stefano Levialdi, and Carlo Batini. 1997. Visual Query Systems for Databases: A Survey. J. Vis. Lang. Comput. 8, 2 (1997), 215–260. https://doi.org/10.1006/jvlc.1997.0037
- [13] Tiziana Catarci, Giuseppe Santucci, and Michele Angelaccio. 1993. Fundamental Graphical Primitives for Visual Query Languages. Inf. Syst. 18, 2 (1993), 75–98. https://doi.org/10.1016/0306-4379(93)90006-M
- [14] Stefano Ceri and Georg Gottlob. 1985. Translating SQL Into Relational Algebra: Optimization, Semantics, and Equivalence of SQL Queries. IEEE Trans. Software Eng. 11, 4 (1985), 324–345. https://doi.org/10.1109/TSE.1985.232223
- [15] Stefano Ceri, Georg Gottlob, and Letizia Tanca. 1989. What you Always Wanted to Know About Datalog (And Never Dared to Ask). IEEE Trans. Knowl. Data Eng. 1, 1 (1989), 146–166. https://doi.org/10.1109/69.43410
- [16] Claudio Cerullo and Marco Porta. 2007. A System for Database Visual Querying and Query Visualization: Complementing Text and Graphics to Increase Expressiveness. In DEXA. IEEE, 109–113. https://doi.org/10.1109/DEXA.2007.91
- [17] Ashok K. Chandra and Philip M. Merlin. 1977. Optimal Implementation of Conjunctive Queries in Relational Data Bases. In Proceedings of the Ninth Annual ACM Symposium on Theory of Computing (Boulder, Colorado, USA) (STOC '77). ACM, New York, NY, USA, 77–90. https://doi.org/10.1145/800105.803397
- [18] Gloria Chatzopoulou, Magdalini Eirinaki, and Neoklis Polyzotis. 2009. Query Recommendations for Interactive Database Exploration. In SSDBM (LNCS, Vol. 5566). Springer, 3–18. https://doi.org/10.1007/978-3-642-02279-1\_2
- [19] Shumo Chu, Brendan Murphy, Jared Roesch, Alvin Cheung, and Dan Suciu. 2018. Axiomatic Foundations and Algorithms for Deciding Semantic Equivalences of SQL Queries. PVLDB 11, 11 (2018), 1482–1495. https://doi.org/10.1 4778/3236187.3236200
- [20] Gard J. Clark and C. Thomas Wu. 1994. DFQL: Dataflow query language for relational databases. *Inf. Manag.* 27, 1 (1994), 1–15. https://doi.org/10.1016/0378-7206(94)90098-1
- [21] Edgar F. Codd. 1970. A Relational Model of Data for Large Shared Data Banks. Commun. ACM 13, 6 (1970), 377–387. https://doi.org/10.1145/362384.362685
- [22] Thomas M. Connolly and Carolyn E. Begg. 2015. Database Systems: A Practical Approach to Design, Implementation and Management, Global Edition (5 ed.). Pearson Addison Wesley. https://www.pearson.com/en-gb/subject-catalog/p/data base-systems-a-practical-approach-to-design-implementation-and-management-global-edition/P200000003964/
- [23] Wikipedia contributors. 2021. Existential graph Wikipedia, The Free Encyclopedia. https://en.wikipedia.org/wiki/Existential\_graph [Online; accessed June-2021].
- [24] Michael Correll. 2023. Teru Teru Bozu: Defensive Raincloud Plots. Computer Graphics Forum (EuroVis) 42, 3 (2023), 235–246. https://doi.org/10.1111/cgf.14826
- [25] Geoff Cumming. 2013. Understanding the new statistics: Effect sizes, confidence intervals, and meta-analysis. Routledge. https://doi.org/10.4324/9780203807002

- [26] Jonathan Danaparamita and Wolfgang Gatterbauer. 2011. QueryViz: Helping Users Understand SQL queries and their patterns. In EDBT. ACM, 558–561. https://doi.org/10.1145/1951365.1951440
- [27] Christopher J. Date. 2003. An introduction to database systems (8 ed.). Pearson/Addison Wesley Longman. https://dl.acm.org/doi/10.5555/861613
- [28] Frithjof Dau. 2006. Fixing Shin's Reading Algorithm for Peirce's Existential Graphs. In Diagrams (International Conference on Theory and Application of Diagrams) (LNCS, Vol. 4045). Springer, 88–92. https://doi.org/10.1007/117831 83 10
- [29] Jan Van den Bussche and Stijn Vansummeren. 2009. Translating SQL into the relational algebra. Course notes, Hasselt University and Université Libre de Bruxelles. https://dipot.ulb.ac.be/dspace/bitstream/2013/198813/1/sql2alg\_eng.pdf
- [30] Sara Di Bartolomeo, Mirek Riedewald, Wolfgang Gatterbauer, and Cody Dunne. 2021. STRATISFIMAL LAYOUT: A modular optimization model for laying out layered node-link network visualizations. *IEEE Transactions on Visualization and Computer Graphics (VIS'21)* 28, 1 (2021), 324–334. https://doi.org/10.1109/TVCG.2021.3114756 Preprint & Supplemental Material: https://osf.io/qdyt9.
- [31] Pierre Dragicevic. 2016. Fair Statistical Communication in HCI. Springer International Publishing, Cham, 291–330. https://doi.org/10.1007/978-3-319-26633-6\_13
- [32] Pierre Dragicevic. 2018. Can we call mean differences "effect sizes"? https://transparentstatistics.org/2018/07/05/meanings-effect-size/
- [33] Bradley Efron. 1987. Better Bootstrap Confidence Intervals. J. Amer. Statist. Assoc. 82, 397 (1987), 171–185. https://doi.org/10.1080/01621459.1987.10478410
- [34] Ramez Elmasri and Sham Navathe. 2015. Fundamentals of database systems (7 ed.). Addison Wesley. https://dl.acm.org/doi/book/10.5555/2842853
- [35] Cibele Freire, Wolfgang Gatterbauer, Neil Immerman, and Alexandra Meliou. 2015. The Complexity of Resilience and Responsibility for Self-Join-Free Conjunctive Queries. PVLDB 9, 3 (2015), 180–191. https://doi.org/10.14778/2850583.2 850592
- [36] Jean H Gallier. 2011. Discrete mathematics. Springer. https://doi.org/10.1007/978-1-4419-8047-2
- [37] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. 2008. Database systems: The complete book (2 ed.). Prentice Hall Press. https://dl.acm.org/doi/book/10.5555/1450931
- [38] Wolfgang Gatterbauer. 2011. Databases will Visualize Queries too. PVLDB 4, 12 (2011), 1498–1501. https://doi.org/10.1 4778/3402755.3402805
- [39] Wolfgang Gatterbauer. 2023. A Tutorial on Visual Representations of Relational Queries. *PVLDB* 16, 12 (2023), 3890–3893. https://doi.org/10.14778/3611540.3611578. Tutorial page: https://northeastern-datalab.github.io/visual-query-representation-tutorial/, Slides: https://northeastern-datalab.github.io/visual-query-representation-tutorial/slides/VLDB\_2023-Visual\_Representations\_of\_Relational\_Queries.pdf.
- [40] Wolfgang Gatterbauer. 2024. A Comprehensive Tutorial on over 100 Years of Diagrammatic Representations of Logical Statements and Relational Queries. In ICDE. IEEE. Tutorial page: https://northeastern-datalab.github.io/diagrammatic-representation-tutorial/.
- [41] Wolfgang Gatterbauer and Cody Dunne. 2023. Supplemental material for "On the reasonable effectiveness of Relational Diagrams". Homepage: https://relationaldiagrams.com/. Main suplemental material folder on OSF: https://osf.io/q9g6u/. Textbook analysis: https://osf.io/u7c4z. Study tutorial: https://osf.io/mruzw. Stimuli-generating code: https://osf.io/kgx4y. The stimuli: https://osf.io/d5qaj. Stimuli/schema index CSV: https://osf.io/u8bf9. Stimuli/schema index JSON: https://osf.io/sn83j. Server code for hosting the study: https://osf.io/suj4a. Collected data: https://osf.io/8vm42. The final analysis code: https://osf.io/f2xe3. Preregistered analysis code: https://osf.io/4zpsk/.
- [42] Wolfgang Gatterbauer, Cody Dunne, H. V. Jagadish, and Mirek Riedewald. 2022. Principles of Query Visualization. IEEE Data Eng. Bull. 45, 3 (2022), 47–67. http://sites.computer.org/debull/A22sept/p47.pdf
- [43] Wolfgang Gatterbauer and Dan Suciu. 2014. Oblivious bounds on the probability of Boolean functions. *TODS* 39, 1 (2014), 5:1–5:34. https://doi.org/10.1145/2532641
- [44] Paruntungan Girsang. 1994. The comparison of SQL, QBE, and DFQL as query languages for relational databases. Master's thesis. Naval Postgraduate School, Monterey, California. https://core.ac.uk/download/pdf/36723678.pdf
- [45] Alon Y. Halevy, Peter Norvig, and Fernando Pereira. 2009. The Unreasonable Effectiveness of Data. *IEEE Intell. Syst.* 24, 2 (2009), 8–12. https://doi.org/10.1109/MIS.2009.36
- [46] Joseph Y. Halpern, Robert Harper, Neil Immerman, Phokion G. Kolaitis, Moshe Y. Vardi, and Victor Vianu. 2001. On the unusual effectiveness of logic in computer science. *Bulletin of Symbolic Logic* 7, 2 (2001), 213–236. https://doi.org/10.2307/2687775
- [47] Jayant R. Haritsa. 2010. The Picasso Database Query Optimizer Visualizer. PVLDB 3, 2 (2010), 1517–1520. https://doi.org/10.14778/1920841.1921027
- [48] David Hibert and Wilhelm Ackermann. 1928. *Grundzüge der theoretischen Logik. By.* Berlin, J. Springer. https://doi.org/10.2307/2018808

- [49] Bill Howe and Garret Cole. 2010. SQL is Dead; Long Live SQL: Lightweight Query Services for Ad Hoc Research Data. In 4th Microsoft eScience Workshop. https://homes.cs.washington.edu/~billhowe/projects/2014/03/22/SQLShare.html
- [50] Hannu Jaakkola and Bernhard Thalheim. 2003. Visual SQL High-Quality ER-Based Query Treatment. In ER (Workshops) (LNCS). Springer, 129–139. https://doi.org/10.1007/978-3-540-39597-3\_13
- [51] H. V. Jagadish, Adriane Chapman, Aaron Elkiss, Magesh Jayapandian, Yunyao Li, Arnab Nandi, and Cong Yu. 2007. Making database systems usable. In SIGMOD. 13–24. https://doi.org/10.1145/1247480.1247483
- [52] Nodira Khoussainova, Magdalena Balazinska, Wolfgang Gatterbauer, YongChul Kwon, and Dan Suciu. 2009. A Case for A Collaborative Query Management System. In 4th Biennial Conference on Innovative Data Systems Research (CIDR). www.cidrdb.org. https://arxiv.org/abs/0909.1778
- [53] Nodira Khoussainova, YongChul Kwon, Magdalena Balazinska, and Dan Suciu. 2010. SnipSuggest: A Context-Aware SQL-Autocomplete System. *PVLDB* 4, 1 (2010), 22–33. https://doi.org/10.14778/1880172.1880175
- [54] Mirjam J Knol, Wiebe R Pestman, and Diederick E Grobbee. 2011. The (mis) use of overlap of confidence intervals to assess effect modification. European Journal of Epidemiology 26 (2011), 253–254. https://doi.org/10.1007%2Fs10654-011-9563-8
- [55] Aristotelis Leventidis, Jiahui Zhang, Cody Dunne, Wolfgang Gatterbauer, H. V. Jagadish, and Mirek Riedewald. 2020. QueryVis: Logic-based Diagrams help Users Understand Complicated SQL Queries Faster. In SIGMOD. ACM, 2303–2318. https://doi.org/10.1145/3318464.3389767
- [56] Guoliang Li, Ju Fan, Hao Wu, Jiannan Wang, and Jianhua Feng. 2011. DBease: Making Databases User-Friendly and Easily Accessible. In 5th Biennial Conference on Innovative Data Systems Research (CIDR). www.cidrdb.org, 45–56. http://cidrdb.org/cidr2011/Papers/CIDR11\_Paper6.pdf
- [57] Wo-Shun Luk and Steve Kloster. 1986. ELFS: English Language from SQL. ACM Trans. Database Syst. 11, 4 (dec 1986), 447–472. https://doi.org/10.1145/7239.384276
- [58] David Maier. 1983. The Theory of Relational Databases. Computer Science Press. https://dl.acm.org/doi/book/10.5555/ 1097039
- [59] Daphne Miedema and George Fletcher. 2021. SQLVis: Visual Query Representations for Supporting SQL Learners. In Symposium on Visual Languages and Human-Centric Computing (VL/HCC). IEEE, 1–9. https://doi.org/10.1109/VL/HCC 51201.2021.9576431
- [60] Richard E. Pattis. 2013. EBNF: A Notation to Describe Syntax. https://ics.uci.edu/~pattis/misc/ebnf2.pdf. (accessed on September 21, 2021).
- [61] Mark E. Payton, Matthew H. Greenstone, and Nathaniel Schenker. 2003. Overlapping confidence intervals or standard error intervals: What do they mean in terms of statistical significance? *Journal of Insect Science* 3, 1 (01 2003), 34. https://doi.org/10.1093/jis/3.1.34
- [62] Charles Sanders Peirce. 1933. Collected Papers. Vol. 4. Harvard University Press. https://doi.org/10.1177/000271623417 400185
- [63] PostgreSQL. 2022. https://www.postgresql.org.
- [64] Raghu Ramakrishnan and Johannes Gehrke. 2002. Database Management Systems (3 ed.). McGraw-Hill, Inc., USA. https://dl.acm.org/doi/book/10.5555/560733
- [65] Phyllis Reisner. 1981. Human Factors Studies of Database Query Languages: A Survey and Assessment. ACM Comput. Surv. 13, 1 (1981), 13–31. https://doi.org/10.1145/356835.356837
- [66] Phyllis Reisner, Raymond F. Boyce, and Donald D. Chamberlin. 1975. Human Factors Evaluation of Two Data Base Query Languages: Square and Sequel. In AFIPS (AFIPS '75). ACM, 447–452. https://doi.org/10.1145/1499949.1500036
- [67] Relational Diagrams. 2023. https://www.relationaldiagrams.com.
- [68] Don D. Roberts. 1973. The Existential Graphs of Charles S. Peirce. The Hague: Mouton. https://doi.org/10.1515/978311 0226225
- [69] Don D. Roberts. 1992. The existential graphs. Computers & Mathematics with Applications 23, 6 (1992), 639–663. https://doi.org/10.1016/0898-1221(92)90127-4
- [70] Sun-Joo Shin. 1995. The Logical Status of Diagrams. Cambridge University Press. https://doi.org/10.1017/CBO9780511 574696
- [71] Sun-Joo Shin. 2002. The Iconic Logic of Peirce's Graphs. The MIT Press. https://doi.org/10.7551/mitpress/3633.001.0001
- [72] Avi Silberschatz, Henry F. Korth, and S. Sudarshan. 2020. Database System Concepts (7 ed.). McGraw-Hill Book Company. https://www.db-book.com/db7/index.html
- [73] Alkis Simitsis, Kevin Wilkinson, Jason Blais, and Joe Walsh. 2014. VQA: vertica query analyzer. In SIGMOD. 701–704. https://doi.org/10.1145/2588555.2594531
- [74] Soufflé. 2023. https://souffle-lang.github.io/rules.
- [75] Alexandru Toth. 2011. Snowflake joins. https://sourceforge.net/projects/revj/files/ [Online; accessed June-2021].
- [76] Etienne Toussaint, Paolo Guagliardo, Leonid Libkin, and Juan Sequeda. 2022. Troubles with Nulls, Views from the Users. *PVLDB* 15, 11 (2022), 2613–2625. https://www.vldb.org/pvldb/vol15/p2613-guagliardo.pdf

- [77] Jeffrey D. Ullman. 1988. Principles of Database and Knowledge-base Systems, Vol. I. Computer Science Press, Inc. https://dl.acm.org/doi/book/10.5555/42790
- [78] Eugene Wigner. 1960. The Unreasonable Effectiveness of Mathematics in the Natural Sciences. Communications in Pure and Applied Mathematics 13, 1 (1960), 1–14. https://doi.org/10.1002/cpa.3160130102
- [79] Leland Wilkinson. 1999. Dot Plots. The American Statistician 53, 3 (1999), 276–281. https://doi.org/10.1080/00031305.1 999.10474474
- [80] Joseph Jay Zeman. 1964. *The graphical logic of C.S. Peirce.* Ph.D. Dissertation. University of Chicago, Dept. of Philosophy.
- [81] Moshé M. Zloof. 1977. Query-by-Example: A Data Base Language. IBM Systems Journal 16, 4 (1977), 324–343. https://doi.org/10.1147/sj.164.0324

#### A NOMENCLATURE

Symbol	Definition
$\overline{\mathcal{S}}$	Table signature or ordered list of its table references of a query
q	Dissociated query starting from query q
$\mathcal{L}_1 \subseteq^{\text{rep}} \mathcal{L}_2$	Language $\mathcal{L}_2$ can pattern-represent language $\mathcal{L}_1$ , i.e., $\mathcal{L}_2$ can represent all relational query
	patterns of language $\mathcal{L}_1$ .
$\mathcal{L}_1 \not\subseteq^{\text{rep}} \mathcal{L}_2$	Language $\mathcal{L}_2$ can not pattern-represent language $\mathcal{L}_1$ , i.e., there are query patterns in language
	$\mathcal{L}_1$ that $\mathcal{L}_2$ cannot represent
$\mathcal{L}_1 \subsetneq^{\text{rep}} \mathcal{L}_2$	Language $\mathcal{L}_2$ pattern-dominates language $\mathcal{L}_1$ , i.e., $\mathcal{L}_2$ can represent all relational query patterns
	of language $\mathcal{L}_1$ and $\mathcal{L}_2$ can represent relational query patterns that $\mathcal{L}_1$ cannot.
$\mathcal{L}_1 \equiv^{\text{rep}} \mathcal{L}_2$	Languages $\mathcal{L}_1$ and $\mathcal{L}_2$ are representation equivalent, i.e., they can express the identical set of
	relational query patterns.

### B ADDITIONAL EXAMPLE FOR Section 1: LIMITS OF RELATIONAL ALGEBRA

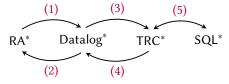
We give an alternative minimal example that illustrates the limits of relational algebra in expressing even simple relational query patterns.

EXAMPLE 10 (RA vs. Datalog). Consider the Datalog  $\$  (non-recursive Datalog with negation) query in Fig. 13g, which returns all tuples in R(A,B) whose attribute B does not appear in the unary table S(B). The query references each of the input tables R and S exactly once. As the proof of our later Theorem 14 shows, basic Relational Algebra (RA) cannot express this query by referencing each of the tables R and S only once. Figures 13b and 13e show two logically-equivalent queries in RA, each of which references the table R(A,B) twice. We also added equivalent Datalog  $\$  queries, which for those two RA expressions use the same "query pattern" (a concept we will formalize later). Intuitively (and we prove this later more formally), Datalog  $\$  can express strictly more query patterns than RA; it has a higher "pattern-expressiveness" despite having the same logical expressiveness. We believe that any diagrammatic language for illustrating and reasoning about query patterns used in queries should be able to express the full range of possible patterns across existing relational query languages (such as the one in Fig. 13g). It follows that any diagrammatic representation of relational queries that relies on a one-to-one mapping with the operators of RA cannot represent the full spectrum of query patterns of relational queries.

Example 10 illustrated that query languages with equal expressiveness may not necessarily be equally able to express the same range of logical patterns—they are not *representation-equivalent*. In other words, there may be queries have have no patterns-preserving translations into the other language.

## C PROOF FOR Section 2, Theorem 6

PROOF OF THEOREM 6. We prove each of the directions in turn:



(1) RA\*  $\rightarrow$  Datalog\*: The proof for this direction is an easy induction on the size of the algebraic expression. It is a minor adaptation of the translation from RA to Datalog $^{\neg}$  proposed by Ullman [77],

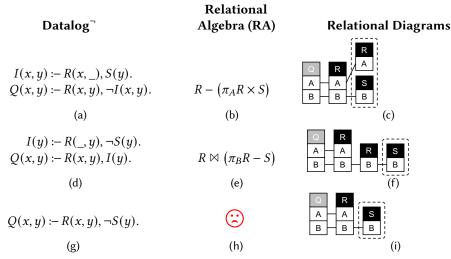


Fig. 13. Given tables R(A,B) and S(B). The first column (a, d, g) shows three logically-equivalent Datalog queries that use three different "query patterns." The first two (a, d) can also be expressed in Relational Algebra (RA) (b, e), whereas the third pattern (g) cannot be expressed in RA, i.e. it is not possible to write a logically-equivalent query in basic RA that references only one occurrence of R and S each. Our paper makes these notions of query patterns precise. The third column (c, f, i) shows Relational Diagrams that use the same "query patterns" as the Datalog "queries in the first column. To the best of our knowledge, our proposed Relational Diagrams are the first diagrammatic representation of relational queries that are (i) unambiguous, (ii) relationally-complete and (iii) able to represent the full range of query patterns for union of non-disjunctive relational calculus.

yet it also pays attention to the restricted fragment and keeps the numbers of atoms constant during the translation. Formally, we show that if an RA\* expression has i occurrences of operators, then there is a Datalog\* program that produces the value of the expression as the relation for one of its Intensional Database (IDB) predicates.

The basis is i=0, that is, a single operand. If this operand is a given relation R, then R is an Extensional Database (EDB) relation and thus "available" without the need for any rules. For the induction, consider an expression whose outermost operator is one of 5 operators: Cartesian product  $\times$ , selection  $\sigma$ , theta join  $\bowtie_c$ , projection  $\pi$ , or difference -. Notice that union  $\cup$  is missing. Also, the rename operator is trivial since Datalog uses position information instead of names.

Case 1:  $Q = E_1 \times E_2$ : Let RA\* expressions  $E_1$  and  $E_2$  have Datalog\* predicates  $e_1$  and  $e_2$  whose rules define their relations, and assume their relations are of arities d and m, respectively. Then define q, the predicate for Q, by:

$$q(x_1, \ldots, x_{d+m}) := e_1(x_1, \ldots, x_d), e_2(x_{d+1}, \ldots, x_{d+m}).$$

Case 2:  $Q = \sigma_c E$ : By restricting our language from RA to RA\*, we only allow selections  $\sigma_c(\varphi)$  where the condition c is a conjunction of simple selections  $c = c_1 \wedge c_2 \wedge \cdots$ , i.e., each selection  $c_i$  is of the form  $\sigma_{A_{i1}\theta A_{i2}}$  (join predicate) or  $\sigma_{A_{i1}\theta v}$  (selection predicate). Let e be a Datalog\* predicate whose relation is the same as the relation for E, and suppose e has arity e. Then the rule for e0 is:

$$q(x_1,\cdots,x_d):=e_1(x_1,\cdots,x_d),c_{\theta}.$$

where  $c_{\theta}$  is a conjunction of join predicates  $x_i \theta x_j$  and selection predicates  $x_i \theta v$  where  $x_i$  and  $x_j$  index the attributes  $A_i$  and  $A_j$ .

Case 3:  $Q = E_1 \bowtie_c E_2$ : While the join operator is not a basic operator of relational algebra, built-in predicates are, in practice, commonly expressed directly through join conditions. This case follows immediately from cases 1 and 2, and the definition of joins as  $Q = E_1 \bowtie_c E_2 = \sigma_c(E_1 \times E_2)$ .

Case 4:  $Q = \pi_{A_{i1},...,A_{id}}(E)$ : Let E 's relation have arity m, and let e be a predicate of arity m whose rules produce the relation for E. Then the rule for the predicate q corresponding to expression Q is:

$$q(x_{i_1},\ldots,x_{i_d}) := e(x_1,\ldots,x_m).$$

Case 5:  $Q = E_1 - E_2$ : We know by definition of the set difference that  $E_1$  and  $E_2$  must have the same arities. Assume those to be d, and that there are predicates  $e_1$  and  $e_2$  whose rules define their relations to be the same as the relations for  $E_1$  and  $E_2$ , respectively. Then we use rule:

$$q(x_1, \dots, x_d) := e_1(x_1, \dots, x_d), \neg e_2(x_1, \dots, x_d).$$

to define a predicate q whose relation is the same as the relation for Q. We can easily check that safety for this rule is fulfilled as all variables appearing in the negated  $e_2$  also appear in the positive  $e_1$ .

(2) Datalog\*  $\rightarrow$  RA\*: Typical textbook translations from Datalog to RA, such as the one by Ullman [77, Th 3.8, Alg 3.2, Alg 3.6] need to compute the active domain by projecting all EDB relations onto each of their components and then taking the union of these projections and the set of constants appearing in the rules. Since RA\* lacks the union operator, we cannot create the active domain from a union of all constants used in the database.

Let  $Datalog^*$  program  $\mathcal{P}$  be a collection of safe, nonrecursive Datalog rules, possibly with negated subgoals. By the safety condition, every variable that appears anywhere in the rule must appear in some nonnegated, relational subgoal of the body, or must be bound by an equality (or a sequence of equalities) to a variable of such an ordinary predicate or to a constant [15]. From the definition of  $Datalog^*$ , each DB appears in exactly one rule as the head. Then, for each DB predicate Q of P, there is an expression Q of relational algebra that computes the relation for Q. Since P is nonrecursive, we can order the predicates according to a topological sort of the dependency graph; that is, if Q appears as a subgoal in a rule for Q, then Q precedes Q in the order.

If a rule contains built-in predicates in the body (join predicates  $x_i\theta x_j$  or selection predicates  $x_i\theta v$ ), the translation first focuses on the body without predicates and then applies a selection  $\sigma_c$  where the selection condition c is a conjunction of the built-in predicates.

To express negated subgoals in the body, we need to use the set difference, and this requires us to complement negated subgoals with additional attributes. Concretely, take a general Datalog rule with built-in predicates:

$$q(\mathbf{x}) := p_1(\mathbf{x}_1), \dots, p_k(\mathbf{x}_k), \neg n_1(\mathbf{y}_1), \dots, \neg n_m(\mathbf{y}_m), c_{\theta}.$$

From the safety conditions of this rule, we know that all variables in the built-in predicates  $c_{\theta}$  need to appear in positive atoms. Similarly, all variables in negated atoms also need to appear in positive atoms:  $\bigcup_{i=1}^{k} \mathbf{x}_{i} \supseteq \bigcup_{i=1}^{m} \mathbf{y}_{i}$ . Let  $\mathbf{z}$  be the set of complementing attributes, i.e., the attributes that only appear in positive atoms:  $\mathbf{z} = \bigcup_{i=1}^{k} \mathbf{x}_{i} - \bigcup_{i=1}^{m} \mathbf{y}_{i}$ .

Let  $P_i$  and  $N_i$  be the RA\* expressions corresponding to Datalog\* predicates  $p_i$  and  $n_i$ . If  $\mathbf{z} = \emptyset$ , then define Q' as

$$(P_1 \bowtie \ldots \bowtie P_k) - (N_1 \bowtie \ldots \bowtie N_m)$$

Otherwise define Q' as

$$(P_1 \bowtie \ldots \bowtie P_k) - ((N_1 \bowtie \ldots \bowtie N_m) \times \pi_{\mathbf{z}}(P_1 \bowtie \ldots \bowtie P_k))$$
(5)

Finally define  $Q = \pi_A(\sigma_\theta(Q'))$  with  $A = (A_{i1}, \dots, A_{id})$  representing the set of attributes indexed by x. Then this expression Q translates one single rule with built-in predicates into a valid relational algebra expression Q without union or disjunction. Since every IDB predicate q appears in only one rule, we do not need union nor disjunctions even if multiple rules are translated. Since all variables in the built-in predicate  $c_\theta$  need to appear in Q', the selection  $\sigma_\theta$  can be correctly applied on Q'. It then follows by induction, on the order in which the IDB predicates are considered, that each has a relation defined by some expression in RA\*.

(3) Datalog\*  $\rightarrow$  TRC\*: We consider a general Datalog rule:

$$q(\mathbf{x}) := p_1(\mathbf{x}_1), \dots, p_k(\mathbf{x}_k), \neg n_1(\mathbf{y}_1), \dots, \neg n_m(\mathbf{y}_m), c_{\theta}.$$

Here  $c_{\theta}$  is a conjunction of built-in predicates that adhere to the standard safety conditions [15]. From those safety conditions, we know that all variables in the built-in predicates  $c_{\theta}$  need to appear in positive atoms, and all variables in negated atoms also need to appear in positive atoms.

The rule then translates into a TRC fragment

$$\{q(\mathbf{A}) \mid p_1 \in P_1, \dots, p_k \in P_k [c_{\text{out}} \land c_p \\ \land \neg (\exists n_1 \in N_1, \dots, n_m \in N_m [c_{\text{in}}])]\}$$

Here A is a set of attributes that correspond to the variables returned by the Datalog rule (from safety conditions, only attributes from the positive relations can be returned),  $c_{\text{out}}$  is a conjunction of equality predicates linking attributes from the output table q to attributes from the input tables  $P_i$ ,  $c_p$  is a conjunction of equality predicates and comparison predicates specified by  $c_\theta$  between the positive relations or constants, and  $c_{\text{in}}$  is a conjunction of equality predicates between exactly one negative relation and either a positive relation or a constant. Notice that this translation guarantees that all predicates (including those in  $c_{\text{in}}$ ) are *guarded*.

(4) TRC\*  $\rightarrow$  Datalog\*: In this translation, we start from the canonical representation of TRC\* (Section 2.3) where a set of existential quantifiers is always preceded by the negation operator (except for the table variables at the root of the query). This implies that we can decompose any query in TRC\* and write it as nested query components, each delimited by the scope of one negation operator. Each query component is then of the form:

$$\{q(\mathbf{A}) \mid p_1 \in P_1, \dots, p_k \in P_k [c_{\text{out}} \land c_p \land \neg q_1(\mathbf{A}_1) \land \dots \land \neg q_m(\mathbf{A}_m)]\}$$

Here **A** is a set of attributes that correspond to the variables returned by the query (or, equivalently, variables that are passed to a nested query that determine whether that nested query is true or false),  $c_{\text{out}}$  is a conjunction of equality predicates linking attributes from the output table q to attributes from the local tables  $P_i$ ,  $c_p$  is a conjunction of equality and comparison predicates between the positive relations or constants, and  $\mathbf{A}_j$  are attributes from the input tables  $P_1, \ldots, P_k$  used in the nested query  $q_j$ .

Notice that for safe queries, only attributes from the positive relations can be returned, i.e., the output attributes need to be connected via equality predicates specified in  $c_{\text{out}}$  to attributes from  $P_1, \ldots, P_k$ . However, nested queries do not need to be safe, and output attributes can be (i) connected directly to further nested queries, or (ii) connected to input tables  $P_1, \ldots, P_k$  via built-in instead of equality predicates. This last point is the main complication we need to take care of during the translation. We proceed in two steps:

(1) First assume that each query is safe. Then each subquery can be immediately translated into a separate rule by induction on the nesting hierarchy from the inside out. The basis of the induction

is the leaf queries which can't contain nested queries and thus are of the form:

$$\{q(\mathbf{A}) \mid p_1 \in P_1, \dots, p_k \in P_k[c_{\text{out}} \land c_p]\}$$

A leaf query is translated into

$$q(\mathbf{x}) := P_1(\mathbf{x}_1), \dots, P_k(\mathbf{x}_k), c_{\theta}.$$

where **x** are attributes chosen from the relations  $P_i$  as specified in  $c_{\text{out}}$ , and  $c_{\theta}$  is the conjunction of comparison predicates between the positive relations  $P_i$ .

For the induction step, assume that each nested  $q_i(\mathbf{A}_i)$  is safe and translated into a rule  $q_i$  Then the safe query q

$$\{q(\mathbf{A}) \mid p_1 \in P_1, \dots, p_k \in P_k [c_{\text{out}} \land c_p \\ \land \neg q_1(\mathbf{A}_1) \land \dots \land \neg q_m(\mathbf{A}_m)]\}$$

is translated into a rule

$$q(\mathbf{x}) := P_1(\mathbf{x}_1), \dots, P_k(\mathbf{x}_k), \neg N_1(\mathbf{y}_1), \dots, \neg N_m(\mathbf{y}_m), c_{\theta}.$$

where  $\mathbf{x}$  are attributes chosen from the positive relations  $P_i$  as specified in  $c_{\text{out}}$ ,  $c_{\theta}$  is a conjunction of comparison predicates between the positive relations  $P_i$ , and  $\mathbf{y}_j$  are chosen from the variables  $\mathbf{x} = \bigcup_{i=1}^{k} \mathbf{x}_i$  used in the positive relations  $P_1, \ldots, P_k$ .

(2) Next assume that a nested query is valid yet not safe for either of the two previously-stated reasons: (i) either some  $\neg q_j(\mathbf{A}_j)$  uses an attribute  $q.A_{it}$  from the output  $q(\mathbf{A})$  directly; or (ii) some predicate in  $c_{\text{out}}$  connects an attribute  $p_j.A_{is}$  from  $P_1,\ldots,P_k$  to an output attribute  $q.A_{it}$  with a built-in predicate  $p_j.A_{is}\theta q.A_{it}$  instead of an equality predicate. In both cases, we can make this query safe by first adding an additional existentially-quantified table  $p_{k+1} \in P_{k+1}$  to the query that has an attribute  $A_i'$  containing the domain of  $q.A_{it}$ . In our translation, we use the same table that is used in the outer query to connect to and therefore to "bound"  $q.A_{it}$  to the active domain. We then have to add appropriate predicates: for case (i), we add the equality predicate  $p_{k+1}.A_i' = q.A_{it}$  to  $c_{\text{out}}$  and replace the attribute  $q.A_{it}$  previously used in  $q_j(A_j)$  with instead  $p_{k+1}.A_i'$ ; for case (ii), we also add the equality predicate  $p_{k+1}.A_i' = q.A_{it}$  to  $c_{\text{out}}$  and replace the previously used predicate  $p_j.A_{is}\theta q.A_{it}$  with instead  $p_j.A_{is}\theta p_{k+1}.A_i'$ . After thus making the subquery safe, we can use the translation described earlier.

It follows that every query in TRC\* can be translated into a logically equivalent query in Datalog\*. We next illustrate both cases for when we need to add existentially-quantified tables in turn.

EXAMPLE 11 (ALL QUANTIFICATION IN DATALOG\*). We illustrate the translation for case (i) above with the following query:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg (\exists s \in S[\neg (\exists r_2 \in R[r_2.B = s.B \land r_2.A = r.A])])]\}$$

It represents relational division and is shown as Relational Diagrams,  $SQL^*$ ,  $RA^*$ , and Datalog\* in Figs. 24 and 25, and will also be re-used in Example 18. Based on our extended safety condition for  $TRC^*$  Definition 3, all predicates are guarded, i.e., they contain at least one attribute of a table that is existentially quantified inside the same negation scope as that predicate. Those guarded attributes are: r.A in r.A=q.A, r.B in r.B=s.B, and r.B=s.B.

Rewriting the query based on its recursive nested negation hierarchy gives us 3 query components:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg (q_1(r.A))] \}$$

$$\{q_1(A) \mid \exists s \in S[\neg (q_2(q_1.A, s.B))] \}$$

$$\{q_2(A, B) \mid \exists r_2 \in R[q_2.B = r_2.B \land r_2.A = q_2.A] \}$$

Now notice that  $q_1$  is not safe because  $q_1.A$  is used within the negated scope  $\neg(q_2(q_1.A, s.B))$  without being existentially quantified (or somehow "bound" to an element in the active domain) within  $q_1$ . In other words, the recursive call  $q_2(q_1.A, s.B)$  passes a predicate through the call hierarchy from  $q_2$  directly to q without being "bound" while passing through  $q_1$ .

We can make  $q_1$  safe (or equivalently "bound" the predicate using  $q_1.A$ ) by adding another table  $r_3 \in R$  in  $q_1$  that accepts and hands over that attribute in the call hierarchy:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg (q_1(r.A))] \}$$

$$\{q_1(A) \mid \exists s \in S, \exists r_3 \in R[q_1.A = r_3.A \land \neg q_2(r_3.A, s.B)] \}$$

$$\{q_2(A,B) \mid \exists r_2 \in R[q_2.B = r_2.B \land r_2.A = q_2.A]] \}$$

This rewritten query now allows a direct translation into Datalog\* from the inside out:

$$Q_2(x, y) := R(x, y).$$
  
 $Q_1(x) := R(x, \_), S(y), \neg Q_2(x, y).$   
 $Q(x) := R(x, \_), \neg Q_1(x).$ 

EXAMPLE 12 (BUILT-IN PREDICATES IN DATALOG\*). We next illustrate the translation for case (ii) of a built-in predicate with the following query asking for values from R for which no smaller value appears in S:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg(\exists s \in S[s.A < r.A])]\}$$

This query is additionally used in Example 21 as  $Q_3$  and illustrated in Figs. 26c, 26g, 26k, 26o and 26s. Rewriting the query based on its recursive nested negation hierarchy gives us 2 query components:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg (q_1(r.A))]\}$$
$$\{q_1(A) \mid \exists s \in S[s.A < q_1.A]\}$$

Now notice that  $q_1$  is not safe because  $q_1.A$  is connected to an existentially-quantified table  $s \in S$  only via a built-in predicate  $s.A < q_1.A$  instead of an equality predicate.

We can make  $q_1$  safe by adding another table  $r_2 \in R$  in  $q_1$ :

$$\begin{aligned} & \{q(A) \mid \exists r \in R[q.A = r.A \land \neg (q_1(r.A))]\} \\ & \{q_1(A) \mid \exists s \in S, r_2 \in R[s.A < r_2.A \land r_2.A = q_1.A]\} \end{aligned}$$

This rewritten query (also illustrated in Figs. 26d, 26h, 26l, 26p and 26t) now allows a direct translation into Datalog\* from the inside out:

$$Q_1(x) := R(x), S(y), x > y.$$
  
 $Q(x) := R(x), \neg Q_1(x).$ 

- (5) TRC\*  $\leftrightarrow$  SQL\*: We prove equivalence in three steps: We first reduce the syntactic variety of SQL\*, then define a canonical form, and finally prove a one-to-one mapping between that canonical SQL\* and canonical TRC\*.
- 1. Starting from the grammar in Fig. 3, we first transform "membership subqueries" (Fig. 14a) and "quantified subqueries" (Figs. 14b and 14c) into equivalent "existential subqueries." Here O' is the complement operator of O (for example "<" for ">=") and C1 and C2 represent different columns or attributes.
- 2. Analogous to the canonical form of TRC\*, we pull existential quantifiers of tables (table variables defined in FROM clauses) as early as possible such that they either appear in the root

```
WHERE {P AND} C1 [not] IN
                                     WHERE {P AND} [not] exists
                                       (SFLECT *
(SELECT C2
FROM R {, R}
[WHERE P {AND P}])
                                      FROM R {, R}
                                      WHERE [P {AND P} AND]
          (a) Transforming membership subqueries
                                      WHERE {P AND} not exists
WHERE {P AND} C1 O ALL
 (SELECT C2
                                       (SELECT *
                                      FROM R {, R}
WHERE [P {AND P} AND]
 FROM R {, R}
 [WHERE P {AND P}])
                                      C1 0' C2)
        (b) Transforming ALL quantified subqueries
                                     WHERE {P AND} exists
WHERE {P AND} C1 O ANY
(SELECT C2
                                       (SFLECT *
FROM R {, R}
[WHERE P {AND P}])
                                       FROM R {, R}
                                       WHERE [P {AND P} AND]
                                      C1 0 C2)
        (c) Transforming ANY quantified subqueries
FROM R1 {, R2}
                                     FROM R1 {, R2}, R3 {,
WHERE {P1 AND} exists
                                     WHERE {P1 AND} P2 {AND P3}
 (SELECT *
FROM R3 {, R4}
WHERE P2 {AND P3})
```

Fig. 14. Part (5) in the proof of Theorem 6: SQL\* queries can be brought into a canonical form by replacing membership subqueries (a) and quantified subqueries (b, c) with existential subqueries that push join predicates into the local scope of the nested query, and then unnesting non-negated subqueries (d).

(d) Unnesting non-negated subqueries

query, or directly following a not exists (Fig. 14d). This reduction is deterministic, and every valid TRC\* query is equivalent to exactly one canonical TRC\* query.

3. The resulting canonical SQL\* is now in a direct 1-to-1 correspondence to TRC\*, and the translation between SQL\* and TRC\* is then the matter of translating the different syntactic expressions between the two languages: From our grammar (Fig. 3), SELECT DISTINCT C {, C} is equivalent to the output definition in TRC\*  $\{q(A) \mid \ldots\}$ , each FROM R {, R} defines the existentially-quantified tuple variables  $\exists r \in R[\ldots]$ , each not exists(SELECT \* FROM R {, R} ...) corresponds to negated existentially-quantified tuple variables,  $\neg(\exists r \in R[\ldots])$ , and the syntax of predicates is identical. The Boolean variants are similar yet the resulting TRC\* is a logical statement and thus without curly braces for set identifiers.

EXAMPLE 13 (SQL\* vs. TRC\*). Figure 15 shows three different non-disjunctive queries in  $TRC^*$ , various syntactic variants of  $SQL^*$ , and Relational Diagrams.  $SQL^*$  queries (b, h, m) are canonical and in a direct 1-to-1 relationship with  $TRC^*$ .

#### D MORE EXPLANATIONS FOR Section 3

A table can be represented by any visual grouping of its attributes (see Fig. 16 for examples). Our choice in Relational Diagrams is to use the typical UML convention of representing tables as rectangular boxes with the table name on top and attribute names below in separate rows (Fig. 16a). Any alternative choice may affect the readability and usability of Relational Diagrams, yet does not

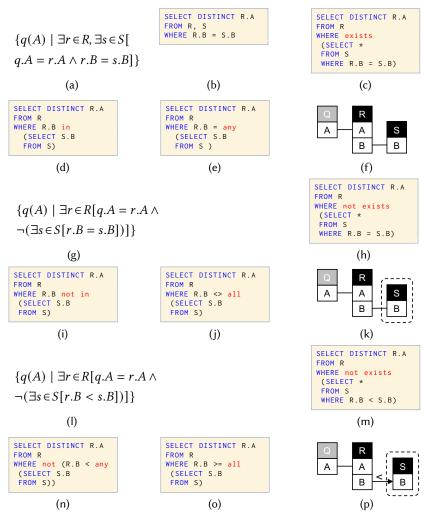


Fig. 15. Example 13: SQL has a redundant syntax if interpreted under set semantics ("SELECT DISTINCT"), binary logic (tables contain no null values) and compared with TRC. Here, queries (a)–(e), queries (g)–(j), and queries (l)–(n) are equivalent. On the right, (f), (k), and (o) show the three corresponding Relational Diagrams (Section 3) that abstract away the syntactic variants and focus on the logical patterns of the queries. SQL queries (b), (h), (m) are canonical and isomorphic to the TRC queries.

affect their semantics and pattern expressiveness. Our focus in this paper is pattern expressiveness, not usability.

# E MORE EXAMPLES FOR Section 3.5: LOGICAL SENTENCES (OR BOOLEAN QUERIES)

Example 14 (Sailors reserving all red boats). Consider the sailor database [64] that models sailors reserving boats: Sailor(sid, sname, rating, age), Reserves(sid, bid, day), Boat(bid,bname,color),

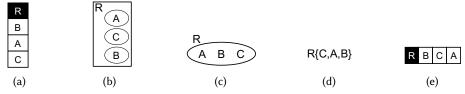


Fig. 16. Appendix D: A few alternative ways to visualize a table and its set of attributes as a group of nodes. Inspired by a familiar UML convention for ER diagrams, we chose (a).

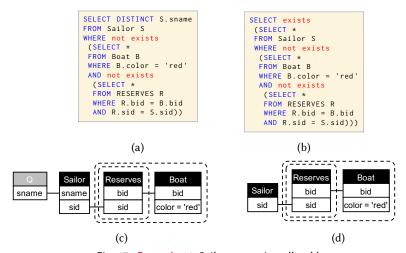


Fig. 17. Example 14: Sailors reserving all red boats.

and the query "Find sailors who reserved all red boats:"

Contrast it with the logical statement "There is a sailor who reserved all red boats." In TRC\*, the difference is achieved by leaving away curly brackets and any mentions of the output table (highlighted for a different example in green color in Fig. 5a):

$$\exists s \in Sailor[ \\ \neg (\exists b \in Boat[b.color = 'red' \land \\ \neg (\exists r \in Reserves[r.bid = b.bid \land r.sid = s.sid])])$$
 (7)

Similarly, Relational Diagrams loses the output table (contrast Fig. 17c with Fig. 17d and their respective SQL statements).

We give an example that shows that in order to express sentences (instead of queries), and to be relationally complete (in that we would like to be able to express all logical sentences), we actually would not have to introduce the visual union. This is in stark contrast to the union at the root being *necessary* for queries.

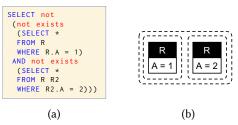


Fig. 18. Example 15:  $\exists r \in R[R.A = 1 \lor R.A = 2]$ .

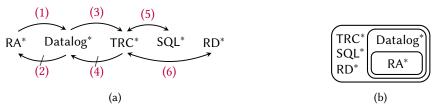


Fig. 19. Appendix F: Directions used in proof for Theorem 14 (a) and resulting representation hierarchy (b). We use two separations lemmas for (2) and (4) in Appendix F.1.

Example 15 (Disjunctions). Consider the simplest disjunction

$$\exists r \in R[R.A = 1 \lor R.A = 2]$$

We can remove the disjunction with a double negation:

$$\exists r \in R[R.A = 1] \lor \exists r \in R[R.A = 2]$$

$$\neg(\neg(\exists r \in R[R.A = 1] \lor \exists r \in R[R.A = 2]))$$

$$\neg(\neg(\exists r \in R[R.A = 1]) \land \neg(\exists r \in R[R.A = 2]))$$

The first 4 steps of the translation in Section 3.2 still work and leads to Fig. 18b. For SQL, the query uses the second new rule to express double negation before the first FROM clause.

# F PROOF FOR Section 4, Theorem 14

The most interesting parts of this proof are two separation results (Fig. 19). We bring those in two separate lemmas first.

# F.1 Two Separation Lemmas

LEMMA 19 (RA\*  $\not\supseteq$ <sup>REP</sup> DATALOG\*). The following Datalog\* query over schema R(A, B), S(B) has no pattern-isomorphic query in RA\*:

$$Q(x,y) := R(x,y), \neg S(y) \tag{8}$$

PROOF LEMMA 19. We show that query (8) (also used in our earlier Example 10) has no pattern-isomorphic query in RA\*. The simple intuition is that the binary *minus operator from RA requires* the same arity of the two input relations. Thus one cannot apply the minus operator directly to combine *R* and *S* as in Datalog\*. Any possible sequence that includes a minus thus either uses the minus on 1 attribute, or 2 attributes (or more attributes). We show that any such option requires at least 3 table instances.

Case 1: Minus on 2 (or more) attributes: Having 2 (or more) attributes for the minus requires us to increase the arity of the right side and thus *S*. This in turn requires a cross product with the

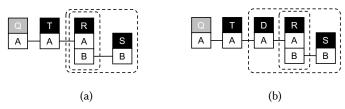


Fig. 20. Illustration for Lemma 20 to show that Datalog\* ⊉<sup>rep</sup> TRC\*, which forms part (4) of the proof of Theorem 14.

domain from *R*.*A* before the minus as in the following translation:

$$R - (\pi_A R \times S)$$

This in turn increases the number of input table instances used from 2 to at least 3, which prevents a pattern-preserving representation.

Case 2: Minus on 1 attribute: Having 1 attribute on the minus requires us to increase the arity *after applying the minus* (because our output has arity 2). This in turn unavoidably increases the number of input table instances to at least 3, which again prevents a pattern-preserving representation. An example translation first uses a projection on R before the minus (to reduce the left input to arity 1) and then a subsequent join again with R after the minus:

$$R \bowtie_B (\pi_B R - S)$$

It follows that any RA\* expression representing the Datalog\* expression (8) needs at least 3 references to input tables and thus cannot preserve the representation from (8).  $\Box$ 

LEMMA 20 (DATALOG\*  $\not\supseteq^{\text{REP}}$  TRC\*). The following TRC\* query over a schema T(A), R(A, B), S(B) has no pattern-isomorphic query in Datalog\*:

$$\{q(A) \mid \exists t \in T[q.A = t.A \land \neg (\exists s \in S[ \neg (\exists r \in R[r.B = s.B \land r.A = t.A])])]\}$$

$$(9)$$

PROOF LEMMA 20. The TRC\* query (9) returns attribute values from T.A that co-occur in R with all attribute values from S.B. We will prove that Datalog\* cannot isomorphically represent this query (which corresponds to the dissociated query (14) for relational division from Example 18).

The intuitive reason is that the predicate "t.A = r.A" represents a join across two different negation scopes (Figure 20a shows that query as Relational Diagram where that join predicate "t.A = r.A" crosses two negation boxes). The safety condition of Datalog\* requires that each variable occurring in a negated atom also occurs in at least one non-negated atom of the same rule [15]. As such, it can express negation only one rule at a time (each rule only allows the application of one negation), and it cannot represent the join predicate r.A = t.A across two negations without an additional domain table to fulfill the safety condition. The standard translation thus uses two separate rules, each of which can express maximally one negation scope, and each such rule needs to fulfill the safety conditions. The first rule finds all the A values that do not co-occur with all S.B values. The second rule then finds the complement against the domain from T.A:

$$I(x) := D(x), S(y), \neg R(x, y)$$
$$Q(x) := T(x), \neg I(x)$$

Notice that the first "extra" atom D(x) is needed for the safety condition of Datalog\* that requires that each variable occurring in a negated atom also occur in at least one non-negated atom of the same rule [15]. Here D.A represents any domain that includes all values that appear in T.A (it can

be a subset or even empty only if S.B is also empty). Thus in general, we require  $D.A \supseteq T.A$  for this Datalog query to be logically equivalent with our TRC\* query (9). Figure 20b shows this last query as Relational Diagram. This domain table cannot be R or S (which may contain different domains) and one occurrence of T needs to be already used outside any negation). We thus need to use an additional reference to T as "guard" for that join predicate. Thus Datalog\* requires adding an additional table reference. Thus Datalog\* cannot preserve the pattern from (9).

This was an intuition, and we make this more precise: We next show that the TRC\* query (9) cannot be expressed in Datalog\* without some additional domain table D(x), which would have to be T(x) if T, R, S were the only tables available in a database. In other words, it cannot be expressed with each of the 3 tables T, R, S appearing only once.

First, observe that q is positive monotone in T and R and negative monotone in S, i.e. adding tuples to T or R can never remove a tuple from the output, and adding tuples to S can never add a tuple to the output.

Next consider a general Datalog program without built-in predicates. Here  $p_i$  are positive atoms, and  $n_i$  negative atoms:

$$\begin{aligned} q_d(\mathbf{x}_d) &:= p_{d1}(\mathbf{x}_{d1}), \dots, p_{dk_d}(\mathbf{x}_{dk_d}), \neg n_{d1}(\mathbf{y}_{d1}), \dots, \neg n_{dm_d}(\mathbf{y}_{dm_d}).\\ & \dots \\ q_1(\mathbf{x}_1) &:= p_{11}(\mathbf{x}_{11}), \dots, p_{1k_1}(\mathbf{x}_{1k_1}), \neg n_{11}(\mathbf{y}_{11}), \dots, \neg n_{1m_1}(\mathbf{y}_{1m_1}).\\ q_0(\mathbf{x}_0) &:= p_{10}(\mathbf{x}_{10}), \dots, p_{0k_0}(\mathbf{x}_{0k_0}), \neg n_{01}(\mathbf{y}_{01}), \dots, \neg n_{0m_0}(\mathbf{y}_{0m_0}). \end{aligned}$$

Recall from Definition 1 of Datalog\* that (i) every IDB appears in the head of exactly one rule (no disjunction), and (ii) every IDB can be used maximally once in any body (no re-use of IDBs). Furthermore (iii) we will assume every relation T, R, S appears exactly once in the body and show this leads to a contradiction.

WLOG we only focus on "canonical" Datalog programs in which IDBs in the body can only appear in negated atoms. This is WLOG, because an IDB that appears in a positive atom (recall it can appear only once in the body) can always be replaced with the body of its defining rule and give an equivalent Datalog program with the same number of table occurrences.

Next, define the "nesting depth"  $ND(q_i)$  of  $q_i$  recursively as the number of rules to traverse to reach  $q_0$ , which has by definition depth 0. In other words, treat each rule as a hyperedge with atoms in head and body as vertices. Then  $ND(q_i)$  is the length of the path (the number of hyperedges or rules to traverse) to reach  $q_0$  starting from  $q_i$ .

Next, define the "sign" of an EDB in  $q_i$  as positive if is it appears as positive atom in  $q_i$  and  $(ND(q_i) \mod 2) = 0$ , or as negative atom and  $(ND(q_i) \mod 2) = 1$ . Analogously for negative.

Recall that every relation and every IDB is used only once, and thus appears either as positive or negative. It follows from an induction argument on the nesting depth that the query  $q_0$  is positive monotone in relation R iff it appears with positive sign in the Datalog program. Analogously for negative monotone.

Our proof now proceeds by simple enumeration over all possible canonical Datalog\* programs that use T, R, S only once, consistent with their defined "signs" (T and R positive, and S negative). Case 1. Consider one rule  $q_0(x) : -T(x)$ , R(z),  $\neg S(y)$ . A single rule cannot express q', thus we

need at least 2 rules.

 $<sup>\</sup>overline{^{10}}$ Assume for a moment that we did not require  $D.A \supseteq T.A$  and instead used D.A = R.A. Then the query would incorrectly return Q(0) for the example database T(0), D(1), S(5), R(1, 5).

Case 2. The only way to have canonical rules up to nesting depth 2 with T in  $q_0$  and R and S appearing with a sign consistent with their expected monotonicity in  $q_0$  is

$$q_2(y) := R(x, y)$$
  
 $q_1(y) := S(y), \neg q_2(y).$   
 $q_0(x) := T(x), \neg q_1(y).$ 

which is not equal to q'. Thus, we can only have a Datalog rule with nesting depth 1.

Case 3. There is no way to have two canonical rules with nesting depth 1 with T in  $q_0$  and R and S appearing with a sign consistent with their expected monotonicity in  $q_0$  since one rule needs to contain S as a positive atom, and the other one would have to have R as a negative atom, which is not possible without an additional atom to make this safe. Thus we can only have two rules, one of which is nesting depth 1.

Case 4. Assume we have two canonical rules, one of nesting level 1. Since *R* is positive, it can either appear in level 0 as positive or in level 1 negated. *R* appearing in level 0 leads to a contradiction:

$$q_1(y) := S(y).$$
  
 $q_0(x) := T(x), R(\mathbf{z}), \neg q_1(y).$ 

Case 5. The last option is of the form:

$$q_1(x) := S(y), \neg R(x, y).$$
  
 $q_0(x) := T(x), \neg q_1(x).$ 

For that program to be safe, both x and y in R need to be bound to an element in the active domain, which can only happen if unary S is accompanied with another relation. Contradiction.

# F.2 Proof Representation Hierarchy

PROOF OF THEOREM 14. We prove each of the directions in turn (Fig. 19). The logical equivalences already follow from the proof of Theorem 6 in Appendix C. We need to prove that certain directions are guaranteed to be pattern-preserving, and for the other directions that do not preserve the structure in general, we give minimum counterexamples.

- (1) RA\*  $\subseteq$  rep Datalog\*: This direction follows immediately from the proof part (1) of Theorem 6 in Appendix C by observing each of the mappings in the 5 cases to be pattern-preserving.
- (2) RA\* ⊉<sup>rep</sup> Datalog\*: Lemma 19 showed that the set difference (or minus –) from RA\* cannot isomorphically represent negation from Datalog\* if the complementing set of attributes is non-empty (see (5)).
- (3) Datalog\*  $\subseteq$  rep TRC\*: This direction follows immediately from the proof of Theorem 6 by observing the mappings of each Datalog rule to be pattern-preserving.
- (4) Datalog\* ≱<sup>rep</sup> TRC\*: Lemma 20 showed that Datalog\* cannot isomorphically represent the TRC\* query (9) (which corresponds to the dissociated query (14) for relational division from Example 18).
- (5) TRC\*  $\equiv$  rep SQL\*: This also follows immediately from the proof part (5) in Appendix C which preserves 1-to-1 correspondences of the mappings in either direction.
- (6) TRC\*  $\equiv$  rep RD\*: This follows immediately from the step-by-step translations in Sections 3.2 and 3.3 which keeps a 1-to-1 correspondence between table references and thus form the proof. □



Fig. 21. Appendix G.1: Directions used in proof for Theorem 21 (a) and resulting representation hierarchy (b).

# G EXTENSION FOR Section 4: PATTERN EXPRESSIVENESS OF RA\* WITH ADDITIONAL OPERATORS

# G.1 RA\* with antijoins (RA\*<sup>▶</sup>)

We have so far focused on the pattern expressiveness of RA\* using only the basic algebraic operators except for the union. We will now show that adding the antijoin operator can extend the expressiveness to the same as Datalog\*.

Given two relations R and S, and a conjunction c of equality predicates between attributes from R and S, the antijoin  $R \triangleright_c S$  (sometimes written as  $R \bowtie_c S$ ) returns all tuples in R that do not have any tuple in S that joins with R based on the equality predicates in C [72]. For example,  $R \triangleright_{R.A=S.B} S$  outputs all tuples in R that do not have any tuple in S whose S. B attribute value matches that tuples R. A attribute value. It is formally defined as  $R \triangleright_c S = R - \pi_{schema(R)}(R \bowtie_c S)$ . It is customary to leave away the conditions C if they are equijoins on identically named attributes and thus corresponding to a natural join. Thus,  $R \triangleright_S = R - \pi_{schema(R)}(R \bowtie_S)$ .

THEOREM 21 (RA\*\*).  $RA^{**}$  ( $RA^*$  extended with the antijoin operator  $\triangleright$ ) and Datalog\* are representation-equivalent.

THEOREM 21. We again prove each of the two directions in turn (Fig. 21) and build upon the proof of Theorem 6 in Appendix C.

(1)  $RA^{**} \subseteq ^{rep} Datalog^*$ : Building upon Appendix C, we need to show that the translation of the antijoin is pattern-preserving.

Case 6:  $Q = E_1 \triangleright_c E_2$ : Assume that there are predicates  $e_1$  and  $e_2$  whose rules define their relations to be the same as the relations for  $E_1$  and  $E_2$ . We partition the attributes of  $E_1$  and  $E_2$  into three sets based on the equality conditions c: those that appear in  $E_1$  but not in  $E_2$  (indexed by  $\mathbf{x}$ ), those that appear in both (indexed by  $\mathbf{y}$ ), and those that appear only in  $E_2$  (indexed by  $\mathbf{z}$ ).

Then we use two rules:

$$q'(\mathbf{y}) := e_2(\mathbf{y}, \mathbf{z}).$$
  
 $q(\mathbf{x}, \mathbf{y}) := e_1(\mathbf{x}, \mathbf{y}), \neg q'(\mathbf{y}).$ 

to define a first predicate q' whose relation contains all the attribute values from  $E_2$  that could join with  $E_1$  and q a second predicate that contains the result of the antijoin. We can easily see that the two rules preserve the referenced input tables, and that the safety for these rules is fulfilled as all variables appearing in the negated q' also appear in the positive  $e_1$ .

(2) Datalog\* → RA\*\*: Building upon Appendix C, we show that the translation of rules with negated subgoals in the body can be achieved in a pattern-preserving way by using the antijoin. Concretely, take a general Datalog rule with built-in predicates:

$$q(\mathbf{x}) := p_1(\mathbf{x}_1), \dots, p_k(\mathbf{x}_k), \neg n_1(\mathbf{y}_1), \dots, \neg n_m(\mathbf{y}_m), c_{\theta}.$$

```
FROM R LEFT JOIN

(SELECT X.A

FROM (SELECT R.A, S.B

FROM R, S) AS X

LEFT JOIN R

ON (R.B = X.B

AND R.A = X.A)

WHERE R.A IS NULL) AS Y

ON R.A = Y.A

WHERE Y.A IS NULL
```

Fig. 22. Example 17: Relational division in SQL with antijoins (syntactically expressed and LEFT JOINS with "IS NULL" selections).

From the safety conditions of this rule, we know that all variables in the built-in predicates  $c_{\theta}$  need to appear in positive atoms. Similarly, all variables in negated atoms also need to appear in positive atoms:  $\bigcup_{1}^{k} \mathbf{x}_{i} \supseteq \mathbf{y}_{i}, \forall i \in [m]$ . Let  $P_{i}$  and  $N_{i}$  be the RA\* expressions corresponding to Datalog\* predicates  $p_{i}$  and  $n_{i}$ , A represent the set of attributes indexed by  $\mathbf{x}$ , and  $c_{i}$  represent the conjunction of equality predicates implied by the set of re-used variables  $\mathbf{y}_{i}$  for each negated predicate  $n_{i}, i \in [m]$ . Then the following expression translates one single rule with built-in predicates into a valid RA\*\* expression in a pattern-preserving way:

$$\pi_{\mathbf{A}}(\sigma_{\theta}(\dots(((P_1 \bowtie \dots \bowtie P_k) \triangleright_{c1} N_1) \triangleright_{c2} N_2) \dots \triangleright_{c_m} N_m))$$
(10)

Since all variables in the built-in predicate  $c_{\theta}$  need to appear in  $\bigcup_{i=1}^{k} \mathbf{x}_{i}$ , the selection  $\sigma_{\theta}$  can be correctly applied on Q'. It then follows by induction on the order in which the IDB predicates are considered that each has a relation defined by some expression in RA\*.

Example 16. Using the antijoin operator, the Datalog\* query Q(x,y) := R(x,y),  $\neg S(y)$ , from Example 10 can be translated in a pattern-preserving way into the  $RA^{**}$  expression  $R \triangleright S$  (or  $R \triangleright_{R.B=S.B} S$  with explicit join conditions)

Example 17 (Relational division with antijoins). We can use the antijoin operator to also express relational division

$$I(x) := R(x, \_), S(y), \neg R(x, y).$$
  
 $Q(x) := R(x, \_), \neg I(x).$ 

in RA\*<sup>▶</sup> as

$$\pi_A R \triangleright \pi_A ((\pi_A R \times S) \triangleright R)$$

While standard SQL does not have a dedicated antijoin operator, it can model the operator via a left join and "IS NULL" selection as shown in Fig. 22 Notice that the syntax is quite different from the SQL expressions shown elsewhere throughout this paper (e.g. we require subqueries in the FROM clause because we require two "IS NULL" selections) and, for equivalence, still requires that the input tables have no NULL values (relations with NULLs are only an intermediate representation to express the antijoin). However, this syntax is then still pattern-isomorphic with relational division expressed above in Datalog and in Fig. 25b in various languages. This example illustrates the powerful abilities of reasoning in terms of patterns across languages and various syntactic constructs.

# G.2 Adding relational division to Datalog\*

We saw adding the antijoin operator to RA\* makes resulting language RA\* equally pattern-expressive as Datalog\*. We next show that adding the relational division to Datalog\* does not make



Fig. 23. Illustration for Appendix G.2.

it as expressive as TRC\*. We achieve this by giving a TRC\* query  $q_{S3}$  that requires two additional tables in the translation to Datalog\* (and thus also to RA\* since Datalog\* can pattern represent RA\*).

$$\{q_{S3}(A) \mid \exists r_0 \in R_0[q_{S3}.A_0 = r_0.A_0 \land \\ \neg \exists r_1 \in R_1[\\ \neg \exists r_2 \in R_2[r_2.A_1 = r_1.A_1 \land \\ \neg \exists r_3 \in R_3[r_3.A_2 = r_2.A_2 \land r_3.A_3 = r_0.A_0])])\}\}$$

$$(11)$$

We assume a schema of only binary relations:  $R_0(A_0, B)$ ,  $R_1(A_0, A_1)$ ,  $R_2(A_1, A_2)$ ,  $R_3(A_2, A_3)$ . Figure 23b shows  $q_{S3}$  in Relational Diagrams. Now assume that we add an additional "operator" to Datalog\* that captures the semantics of relational division (shown in Figure 23a):

$$I(x) := R_0(x, \_), R_1(\_, y), \neg R_0(x, y).$$

$$Q_D(x) := R_0(x, \_), \neg I(x).$$
(12)

We represent this "operator" as a function  $f_D$  that we add to the constructs of the resulting language, and show that the resulting language has not pattern-isomorphic representation of  $q_{S3}$ .

$$Q_D(x) := f_D[R_0(x, y), R_1(z, y)]. \tag{13}$$

First, notice that  $f_D$  is monotone positive in  $R_0$  and monotone negative in  $R_1$ .  $q_{S3}$  is monotone positive in  $R_0$  and  $R_2$  and monotone negative in  $R_1$ .

Next notice that the output schema of  $f_D$  is unary, while the inputs to  $f_D$  are binary. Thus  $f_D$  cannot be composed. This prevents that the resulting query has two occurrences of  $f_D$  as otherwise all 4 input tables  $R_0$ ,  $R_1$ ,  $R_2$ ,  $R_3$  would be used as inputs to the two occurrences of  $f_D$  while we cannot express the query. It follows  $f_D$  must occur exactly once (it must occur at least once since  $q_{S3}$  cannot be expressed without  $f_D$ ).

The proof now succeeds be checking all finitely many ways in which  $q_{S3}$  could be expressed with one occurrence of  $f_D$  while respecting the monotonicity constraints and realizing that none of these queries is equivalent to  $q_{S3}$ .

# G.3 Conjecture about limits of adding operators

We leave it open whether RA can be extended to express the full range of relational patterns available in logical languages such as TRC\*. We hypothesize that this gap between procedural languages (either via an algebraic language like RA or declarative languages that are restricted to rule-by-rule evaluation such as Datalog) and logical languages cannot be breached with any natural operator that can be readily implemented in a procedural way.

# **H MORE EXAMPLES FOR Section 4**

We next illustrate, with the help of the more complicated example of relational division, that there is a pattern-preserving mapping from RA to TRC, but not in the other direction.

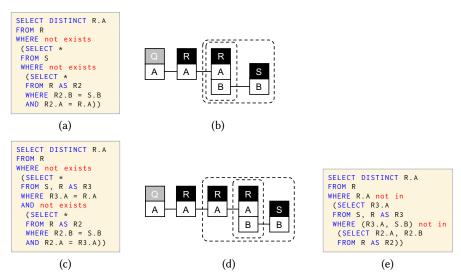


Fig. 24. Example 18: Relational division in SQL (a)(c)(e) and as Relational Diagrams (b)(d). All 5 representations are *logically equivalent*, but only the partitions  $\{(a), (b)\}$  and  $\{(c), (d), (e)\}$  are also *pattern-isomorphic* (which is what we expect). We illustrate in Fig. 25 how the tables in this query correspond to the equivalent queries in Relational Diagrams, SQL\*, RA\*, Datalog\*, TRC\*.

EXAMPLE 18 (TRC AND RA ARE NOT REPRESENTATION-EQUIVALENT). Assume a schema R(A, B), S(B). Consider the relational division asking for attribute values from R. At that co-occur in R with all attribute values from S. B. The translation into TRC is

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg(\exists s \in S[ \neg(\exists r_2 \in R[r_2.B = s.B \land r_2.A = r.A])])]\}$$

$$(14)$$

The corresponding canonical SQL statement is shown in Fig. 24a. Relational division expressed in primitive RA is

$$\pi_A R - \pi_A ((\pi_A R \times S) - R) \tag{15}$$

The translation into Datalog uses two rules:

$$I(x) := R(x, \_), S(y), \neg R(x, y).$$
  
 $Q(x) := R(x, \_), \neg I(x).$  (16)

The atoms  $R(x, \_)$  are needed for the safety condition of Datalog $^{\neg}$ . This translation is part of a standard proof for equivalence of expressiveness between RA and safe TRC in textbooks such as [1, 58, 77].

Now notice an arguably important difference between the three expressions: TRC (14) uses the atom R two times, whereas RA (15) and Datalog (16) use R three times. It turns out that there is no way to represent relational division in primitive RA or Datalog with only two occurrences of the R symbol (see Theorem 14).

There is, however, an alternative representation in TRC that preserves the RA structure with three occurrences of R:

$$\{q(A) \mid \exists r \in R[q.A = r.A \land \neg(\exists s \in S, \exists r_3 \in R[r_3.A = r.A \land \neg(\exists r_2 \in R[r_2.B = s.B \land r_2.A = r_3.A])])\}\}$$

$$(17)$$

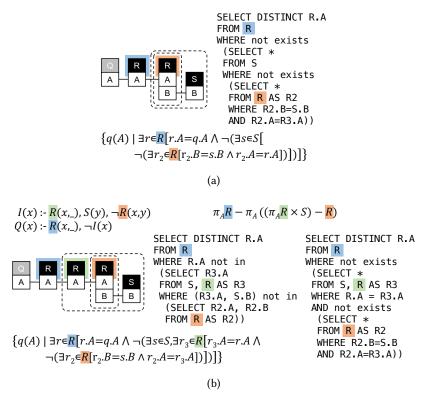


Fig. 25. Example 18 and Fig. 24 continued: Two logically-equivalent sets (a) and (b) of relational division in 5 query languages (Relational Diagrams, SQL\*, RA\*, Datalog\*, TRC\*). The queries in (a) use 2 occurrences of R, whereas the ones in (b) use 3 occurrences of R. We highlight the two or three occurrences of R across the different languages that can be mapped to each other according to our pattern isomorphism defined in Definition 12. We prove in Appendix F that for (a), there is no pattern-isomorphic representation in Datalog¬ (thus neither in RA).

Notice that now there is a natural 1-to-1 correspondence between the atoms in TRC (17) and the atoms in RA (15). This correspondence is even more intuitive by mapping the correspondence between two logically-equivalent SQL statements (Figs. 24c and 24e) and RA (15): for example, lines 4–8 in Fig. 24e translate into the RA fragment " $\pi_A((\pi_A R \times S) - R)$ ", which corresponds to the IDB predicate Temp(x) in Datalog (16).

In other words, while all of these 7 queries are logically equivalent, they partition into two disjoint sets that are "pattern-isomorphic":

This suggests that TRC and RA are not representation equivalent.

We next prove the pattern-isomorphism between RA query (15) and TRC query (17) with our formalism. First, write their dissociated queries  $q'_{RA}(R_1, R_2, S, R_3)$  and  $q'_{TRC}(R_1, S, R_2, R_3)$  with

$$q'_{RA} = \pi_A R_1 - \pi_A ((\pi_A R_2 \times S) - R_3)$$

$$\{q'_{TRC}(A) \mid \exists r \in R_1[q.A = r.A \land \neg(\exists s \in S, \exists r_3 \in R_2[r_3.A = r.A \land \neg(\exists r_2 \in R_3[r_2.B = s.B \land r_2.A = r_3.A])])\}\}$$

Second, define the homomorphism  $h(R_1, R_2, S, R_3) = (R_1, S, R_2, R_3)$ , which is injective and surjective and thus an isomorphism between the signature of the dissociated queries. We can now easily verify that the dissociated queries are logically equivalent after composition with  $h: q'_{RA} \equiv q'_{TRC} \circ h$ . In other words:

$$q'_{RA}(R_1, R_2, S, R_3) \equiv q'_{TRC}(h(R_1, R_2, S, R_3)) = q'_{TRC}(R_1, S, R_2, R_3)$$

Figure 25 illustrates the pattern isomorphism within two sets of queries with the color-highlighted R atoms. Notice in Fig. 25a the correspondences between the blue and orange highlighted tables R across the SQL and the TRC statements and the Relational Diagram. Notice that the constraints between their "A" attributes ("R2.A=R.A") make use of the nesting hierarchy: two levels of the "not exists" nesting hierarchy in SQL and two levels in the negation hierarchy in TRC. Similarly, in Fig. 25b the two SQL variants use different syntactic constructs to represent the single negation hierarchy between R3 and R2. Then, see how Datalog represents this constraint without referring to the explicit attributes "A" but by positional reference and using the repeated variable x together with "not" to represent the same logical constraint. RA represents the same logical constraint by projecting attribute "A" from the green instance R on the left side of a cross product, before the set difference with the yellow instance R. Despite the extreme syntactic variants and logical equivalence of all these queries, by defining individual pairwise isomorphisms between extensional tables, our formalisms allow us to partition these queries into two sets within which the queries are pattern-isomorphic.

Example 19 (Example 3 Continued). For completeness, we also translate here the statement (or Boolean query) "All sailors reserve a red boat" into Datalog\* and RA\*. Boolean queries are usually not shown in RA but are straightforward: The queries project away all attributes as last operation:

$$\pi_{\emptyset}(\pi_{sid}Sailor - \pi_{sid}(Reserves \bowtie (\sigma_{color='red'}Boat)))$$

The translation into Datalog uses the same pattern:

$$I(s) := Reserves(s, b, \_), Boat(b, \_, 'red').$$
  
 $Q := Sailor(s, \_, \_), \neg I(s).$ 

Notice that all 4 languages use the same relational pattern for this query.

#### I MORE DISCUSSIONS FOR Section 4

Our formalism is similar in spirit to edge-preserving graph homomorphisms that map two nodes in graph  $G_1$  linked by an edge to two nodes in graph  $G_2$  that are also linked by an edge. In our pattern-preserving isomorphisms between queries, the role of nodes is played by the table signatures in the queries and the queries themselves play the role of the edges. Notice also the difference in homomorphisms between conjunctive queries for determining query containment [17]: in that formalism, the role of nodes is played by variables (and constants) and the relational atoms play the role of edges. Also, notice from Example 6 that a simpler mapping between the set of tables used (instead of the repeated table references) in two queries alone would not work.

Our definition—by design—does not distinguish query patterns based on operator execution orders. Also by design, our definition does not include any notion of views or intermediate tables. This is achieved by excluding Intensional Database Predicates (as in Datalog) from the definition of table references. We again illustrate the intuition for that design with examples.

EXAMPLE 20 (JOIN ORDERS AND VIEWS DO NOT AFFECT PATTERNS). Assume that the edges of a directed graph are stored in a binary relation E(A, B). Consider a query returning nodes that form the starting point of a length-3 directed path. In unnamed RA where indices replace attribute names [1], we can write the query in two different ways, one applying projections early and the other late:

$$q_1(E) = \pi_1 \sigma_{2=3 \land 4=5} (\mathbf{E} \times \mathbf{E} \times \mathbf{E})$$
  
$$q_2(E) = \pi_1 \sigma_{2=3} (\mathbf{E} \times \pi_1 \sigma_{2=3} (\mathbf{E} \times \pi_1 \mathbf{E}))$$

We can also write these two queries in the more familiar named perspective of RA. These queries encode the same algebraic operations but are more verbose, since the named perspective of RA requires a rename operator  $\rho$  to express the identical sequence of operations:

$$q_1(E) = \pi_{E.A}\sigma_{E.B=F.A \land F.B=G.A}(\mathbf{E} \times \rho_{E \to F}\mathbf{E} \times \rho_{E \to G}\mathbf{E})$$
  
$$q_2(E) = \pi_{E.A}\sigma_{E.B=F.A}(\mathbf{E} \times \pi_{F.A}\sigma_{F.B=G.A}(\rho_{E \to F}\mathbf{E} \times \rho_{E \to G}\pi_A\mathbf{E}))$$

All four RA query expressions use the same relational pattern according to our definition. We believe it is essential to separate the notion of a relational query pattern from concerns regarding operator execution order, because the latter is not meaningful for queries written in declarative relational languages. To see that, consider the logically-equivalent query in Datalog:

$$Q_3(x) := \underline{E}(x, y), \underline{E}(y, z), \underline{E}(z, w).$$

It declaratively specifies what attributes the three tables need to be joined on, but it does not specify any order or joins nor when projections happen. Furthermore, notice that RA query  $q_1$  does not even specify a concrete join order between the three table references and the query can be seen as a shorthand for either of two sequences, more precisely written as  $(E \times E) \times E$  or  $E \times (E \times E)$  depending on the preferred definition of the shorthand.

For a similar reason, temporary tables such as Intensional Database Predicates in Datalog do not count as table references. Thus, the following Datalog query uses the same logical pattern (find three edges that join and keep the starting node), even though it defines the intermediate intensional database predicate I:

$$I(y) := E(y, z), E(z, w).$$

$$Q_4(x) := E(x, y), I(y).$$

Notice that it follows immediately that two pattern-isomorphic queries need to have the same number of table references. We believe that such a pattern-preserving mapping between queries is important if we want to help readers understand the *exact logic* (the exact *logical pattern*) behind a relational query, irrespective of the language it is written in. In particular, if we want to help users understand the relational pattern of an existing relational query with diagrams (recall that we focus on set semantics and binary logic), we must be able to create this 1-to-1 correspondence between the query and its diagrammatic representation.

# J DRILL-DOWN FOR Section 4: LIMITS OF DATALOG FOR REPRESENTING PATTERNS

We have shown earlier that Datalog\* cannot represent all Query patterns from TRC\*. We next use another example to illustrate that this limit of Datalog does not only appear with deeply nested queries; it already appears for simply nested queries with inequality conditions and is an immediate consequence of Datalog's safety conditions for built-in predicates.

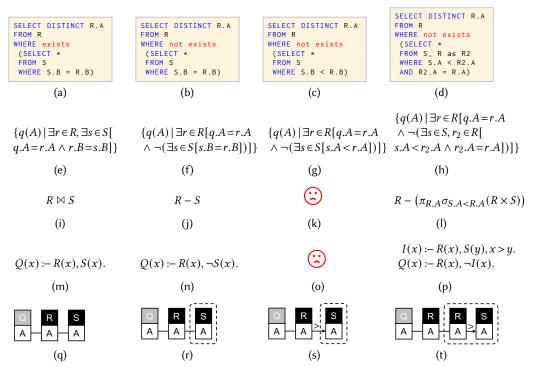


Fig. 26. Example 21: The 5 rows show  $SQL^*$ ,  $TRC^*$ ,  $RA^*$ ,  $Datalog^*$ , and Relational Diagrams, respectively. The first two column show  $Q_1$  and  $Q_2$ , respectively. The last two columns show  $Q_3$ . Notice that the  $SQL^*$ ,  $TRC^*$ , and Relational Diagram queries from the 3rd column have no pattern-isomorphic query in  $RA^*$  or  $Datalog^*$ . Instead,  $RA^*$  and  $Datalog^*$  require an additional cross-join with RA, which is shown in the 4th column.

EXAMPLE 21 (LIMITS OF DATALOG). Consider two unary tables R(A) and S(A) and three questions:

 $Q_1$ : Find values from R that also appear in S.

 $Q_2$ : Find values from R that do not appear in S

 $Q_3$ : Find values from R for which no smaller value appears in S.

These queries are shown in  $SQL^*$ ,  $TRC^*$ ,  $RA^*$ ,  $Datalog^*$ , and as Relational Diagram in Fig. 26, with  $Q_1$  in the 1st column (a, e, i, m, q),  $Q_2$  in the 2nd (b, f, j, n, r), and  $Q_3$  in the 3rd (c, g, k, o, s).

Notice in the third column  $(Q_3)$  that the  $SQL^*$  (c),  $TRC^*$  (g), and Relational Diagram (s) queries have no pattern-isomorphic query in  $RA^*$  (k) or Datalog\* (o). The safety condition of Datalog requires each variable to appear in a non-negated atom. This criterion requires a cross-join with the domain of RA in a separate rule before the negation can be applied on an equality predicate. For the same reason,  $RA^*$  cannot apply the set difference directly and also requires an additional cross-join with RA. Even extending the available operators with an antijoin does not change this since antijoins are only defined for equality conditions [72].

The 4th column (d, h, l, p, t) shows the resulting resulting RA\* and Datalog\* queries together with their pattern-isomorphic queries in  $SQL^*$ ,  $TRC^*$ , and as Relational Diagrams.

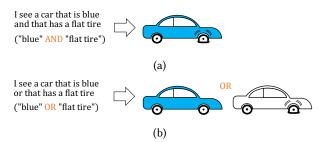


Fig. 27. A single situation (concrete arrangement of items) can only show conjunctive information [71], i.e., the car is blue. To show that the car is blue OR it has a flat tire, you need two situations and a visual construct to show disjunction (here, just the text 'OR').

# K Section 5: WHY DISJUNCTIONS ARE HARDER TO REPRESENT

It is useful to understand why disjunctions are more difficult to represent in a diagrammatic representation. As a motivating example, assume Alice calls Bob and tells him "I see a blue car that has a flat tire." What is the mental image that Bob has from this information? It is a car with two conditions: it is blue, and it has a flat tire as in Fig. 27a. Next, assume that Alice instead tells Bob "I see a car that is either blue or that has a flat tire." What is the mental image that Bob has from that information?

There is no single mental image that could capture that situation. Bob needs to imagine two *different* images. If Bob sees one image with two different cars (one blue, the other with a flat tire), then he actually sees *two separate cars*. Bob needs to add some additional visual symbol representing the logic that those are two *different overlapping possible worlds*. In the words of Shin [71], "any situation" (think of a concrete arrangement of items) "can only display conjunctive information." This diagrammatic representation problem is not as apparent in text: "Car.color = 'blue' OR Car.tire = 'flat'". Shin goes further and claims that "Any diagrammatic system that seeks to represent disjunctive information needs to bring in an artificial syntactic device with its own convention" [70].

#### L PROOFS FOR Section 5

PROOF THEOREM 17. Given a safe TRC expression, we pull any existential quantifier as early as possible: either to be at the start of the query or directly following a negation operator.

First, consider a nested query with disjunctions in the WHERE conditions, possibly nested with conjunctions. Rewrite the conditions as DNF, i.e. as

$$\neg (\exists \mathbf{r} \in \mathbf{R}[c_1 \lor c_2 \cdots \lor c_k])$$

Here,  $\mathbf{r}$  is a set of table variables,  $\mathbf{R}$  a set of tables, and each  $c_i$  is a conjunction of guarded predicates. Next, rewrite it as:

$$\neg (\exists \mathbf{r}_1 \in \mathbf{R}[c_1]) \land \neg (\exists \mathbf{r}_2 \in \mathbf{R}[c_2]) \cdots \land \neg (\exists \mathbf{r}_k \in \mathbf{R}[c_k])$$

This fragment is in TRC\* and can be visualized by Relational Diagrams\*.

 $<sup>^{11}</sup>$ To provide some additional intuition, recall that *conjunctions* of selections can be simply modeled as a concatenation of simple selections, e.g.  $\sigma_{C_1 \wedge C_2}(R)$  is the same as  $\sigma_{C_1}(\sigma_{C_2}(R))$ . Thus conjunctions are an inherently more natural logical connective than disjunctions; disjunctions cannot be represented without additional visual symbols.

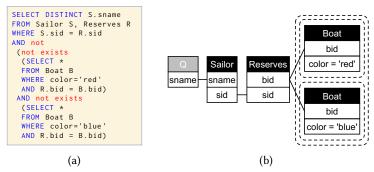


Fig. 28. Example 22: The Query "Find sailors who reserve a red or blue boat" can be represented with double negation in the non-disjunctive fragments of SQL\* (a) and Relational Diagram\* (b).

Second, for remaining disjunctions in the top query  $q_0$ , rewrite the query as a union over queries without disjunction:

$$\{q(\mathbf{A}) \mid \exists \mathbf{r}_1 \in \mathbf{R}_1[c_1] \lor \exists \mathbf{r}_2 \in \mathbf{R}_2[c_2] \cdots \lor \exists \mathbf{r}_k \in \mathbf{R}_k[c_k] \}$$

$$= \{q(\mathbf{A}) \mid \exists \mathbf{r} \in \mathbf{R}[c_1] \} \cup \{q(\mathbf{A}) \mid \exists \mathbf{r} \in \mathbf{R}[c_1] \} \cdots \cup \{q(\mathbf{A}) \mid \exists \mathbf{r} \in \mathbf{R}[c_k] \}$$

Since each individual query is in  $TRC^*$ , the whole query can be represented by Relational Diagrams by representing each individual  $TRC^*$  in a union cell.

#### M MORE EXAMPLES FOR Section 5

Example 22 (Red or blue). Consider the following  $TRC^*$  query asking for sailors who have reserved a red or a blue boat:

```
\{q(sname) \mid \exists s \in Sailor, r \in Reserves[q.sname = s.sname \land s.sid = r.sid \land \exists b \in Boat[
b.bid = r.bid \land (b.color = 'red' \lor b2.color = 'blue')]]\}
```

Using De Morgan's Law  $(A \vee B) = \neg(\neg A \wedge \neg B)$  applied to quantifiers, we can transform the disjunction into double-negation with conjunction. This transformation comes at the cost of repeated uses of extensional tables and is thus not pattern-preserving:

```
\{q(sname) \mid \exists s \in Sailor, r \in Reserves[q.sname = s.sname \land s.sid = r.sid \land \neg (
\neg (\exists b1 \in Boat[b1.bid = r.bid \land b1.color = `red']) \land
\neg (\exists b2 \in Boat[b2.bid = r.bid \land b2.color = `blue']))]\}
(18)
```

Figure 28a shows (18) translated into canonical SQL\* and Fig. 28b its translation into a Relational Diagram. Notice how the non-disjunctive fragment repeats the boats table twice.

EXAMPLE 23 (UNION OF QUERIES). Consider three tables R(A, B), S(A), and T(A) and the running query from Section 2. We can replace disjunction by pulling it to the root and replacing the query

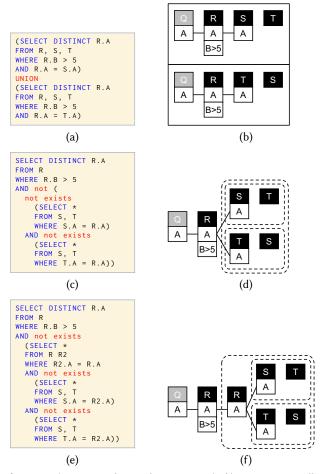


Fig. 29. Illustrations for Example 23 on replacing disjunctions: (a, b) using union cells to replace disjunctions; (c, d) using DeMorgan to replace disjunctions by double-negation and conjunction; (e, f) further changing the table signatures to allow pattern-isomorphic representations in Datalog\* and RA\*\*.

```
with a union of non-disjunctive TRC^* queries:  \{q(A) \mid \exists r \in R, \exists s \in S, \exists t \in T[q.A = r.A \land r.B > 5 \land (r.A = s.A \lor r.A = t.A)]\} 
 = \{q(A) \mid \exists r \in R, \exists s \in S, \exists t \in T[q.A = r.A \land r.B > 5 \land r.A = s.A]\} \cup \{q(A) \mid \exists r \in R, \exists s \in S, \exists t \in T[q.A = r.A \land r.B > 5 \land r.A = t.A]\}
```

Fig. 29a shows the representation-equivalent SQL query. Fig. 29b shows this union of two separate Relational Diagrams\* queries, each in a separate union cell, and each with the same set attributes in the output table. Notice that we cannot leave away the non-connected tables in the individual queries; if any of the tables are empty, then the query needs to return an empty result.

We can also rewrite this query directly as a non-disjunctive TRC\* query by using DeMorgan and repeating table references:

```
 \{q(A) \mid \exists r \in R, \exists s \in S, \exists t \in T[q.A = r.A \land r.B > 5 \land \\ (r.A = s.A \lor r.A = t.A)]\} 
 = \{q(A) \mid \exists r \in R[q.A = r.A \land r.B > 5 \land \\ \exists s \in S, \exists t \in T[r.A = s.A \lor r.A = t.A]]\} 
 = \{q(A) \mid \exists r \in R[q.A = r.A \land r.B > 5 \land \\ (\exists s \in S, \exists t \in T[r.A = s.A] \lor \exists s \in S, \exists t \in T[r.A = t.A])]\} 
 = \{q(A) \mid \exists r \in R[q.A = r.A \land r.B > 5 \land \\ \neg (\neg \{\exists s \in S, \exists t \in T[r.A = s.A]\} \land \neg \{\exists s \in S, \exists t \in T[r.A = t.A]\})\}\}
```

Figure 29a shows the pattern-isomorphic SQL\* query, and Fig. 29b shows the pattern-isomorphic Relational Diagrams\* query.

Based on our results, this query cannot be represented with a pattern-isomorphic RA\*\* nor Datalog\* queries. But we can represent it by increasing the table signature: following the steps from part (4) in the proof of Theorem 6:

```
 \{q(A) \mid \exists r \in R[q.A = r.A \land r.B > 5 \land \\ \neg (\exists r_2 \in R[r_2.A = r.A \land \\ \neg (\exists s \in S, \exists t \in T[r_2.A = s.A]) \land \neg (\exists s \in S, \exists t \in T[r_2.A = t.A])])] \}
```

This query now can be represented in pattern-isomorphic RA\* and Datalog\* queries as follows:

```
I1(x) := S(x), T(\_).
I2(x) := S(\_), T(x).
I3(x) := R(x,\_), \neg I1(x), \neg I2(x).
Q(x) := R(x,y), \neg I3(x), y > 5.
\pi_A(\sigma_{B>5}R) - ((R \triangleright_{R.A=S.A} \pi_{S.A}(S \times T)) \triangleright_{R.A=T.A} \pi_{T.A}(S \times T))
```

Figures 29e and 29f show their pattern-isomorphic representations in SQL\* and Relational Diagrams\*.

# N TEXTBOOK ANALYSIS (Section 6.1)

We give here more details on the 59 queries analyzed in Section 6.1. These queries and our analysis is available on OSF. 12

**Approach.** We identified 5 popular database textbooks [22, 27, 34, 64, 72] that include chapters on relational calculus. Textbooks were considered which the authors of this study had previously used examples from in their own database classes. A 6th popular textbook ("the complete book" from Stanford [37]) is not included since it does not contain a section on relational calculus.

In each of the 5 books, we identified chapters that illustrate relational calculus queries with examples and extracted all example queries from these chapters. In case a query is given multiple times in pattern-isomorphic variants (e.g. first in Tuple Relational Calculus and then again in Domain Relational Calculus), we list the query only once in its earlier representation. In addition, we also identified chapters on SQL queries and extracted those queries that have a logically-equivalent representation in relational calculus (thus no aggregates, arithmetic attributes, or outer joins).

<sup>12</sup> Queries: https://osf.io/u7c4z/

In total we identified 59 queries across those 5 textbooks. For each of the 5 languages Relational Diagrams, RA, Datalog, QueryVis [26], and QBE [81], we analyzed which of those queries can be expressed *in a pattern-isomorphic representation.*<sup>13</sup> For the more interesting queries (e.g. those that have pattern-isomorphic representations in Relational Diagrams but not in RA) we included the Relational Diagrams representation in the documentation.

**Detailed analysis.** (1) From the "cow book" by Ramakrishnan and Gehrke [64], we extracted 25 queries from chapters 4.3.1 on TRC, 4.3.2 on DRC, and 5.2-5.4 on SQL. The number of queries that have pattern-isomorphic representations are: 24 for Relational Diagrams, 22 for QueryVis (it has no union operator and can only visualize a strict subset of the non-disjunctive fragment), 19 in QBE, and 18 in RA or Datalog. The number for RA increases to 19 if the antijoin is added as an additional primitive relational operator (Appendix G.1).

- (2) From the "sailboat book" by Silberschatz, Korth, and Sudarshan [72], we extracted 8 queries from chapters 27.1 (on TRC) and 27.2 (on DRC). The number of queries that have pattern-isomorphic representations are: 8 for Relational Diagrams, and 7 for QueryVis, QBE, RA and Datalog.
- (3) From the "stone formation book" by Elmasri and Navathe [34], we extracted 9 unique queries from chapter 8.6 (on TRC) and found only pattern-isomorphic queries in chapter 8.7 (on DRC). The number of queries that have pattern-isomorphic representations are: 8 for Relational Diagrams, QueryVis, RA, QBE, and 7 for Datalog.
- (4) From Date's book [27], we extracted 9 unique queries from chapter 8.3 (on TRC). The number of queries that have pattern-isomorphic representations are: 8 for Relational Diagrams and QueryVis, 7 for RA and QBE, and 6 for Datalog.
- (5) From the "bookshelf book" by Connolly and Begg [22], we extracted 8 unique queries from section 5.2 (on TRC). All queries have pattern-isomorphic representations in all languages.

#### O CONTROLLED USER EXPERIMENT (Section 6.2)

This section gives additional details on the user study (Appendix O.1) and provides links to the supplemental materials including preregistration, tutorial, actual stimuli, stimuli-creation code, collected data, and analysis code and results (Appendix O.2).

# O.1 Additional details on experimental design

**Procedure.** After a short tutorial (see Fig. 30), participants are asked to answer 32 questions, each with a different schema found or adapted from common textbooks. Each question asks the participant to select which of 4 plain-English patterns is similar to the ones listed in Section 6.2 correctly matches the shown query (but each question applied to a different schema, such as *Students taking classes, Actors playing in movies, Suppliers supplying parts*, etc.). The correct answer is provided after each question so participants *can learn from mistakes*. We analyze the time a participant spends on the question itself (not learning from the answer) as well as whether their chosen answer is correct.

Randomization and Counterbalancing. To reduce ordering effects, we start half the participants using SQL (group 1) and the other using Relational Diagrams (group 2), after which participants alternate conditions with each question (see Section 6.2 and Figure 11).

The total number of possible sequences for each condition and for each half is the multiset permutation  $\binom{8}{2,2,2,2} = 2520.$  The total number of treatments is thus 2520<sup>4</sup> (each of 2 conditions

<sup>&</sup>lt;sup>13</sup>Recall that since all considered query languages (except QueryVis) are relationally complete, all 59 queries have logically-equivalent representations in each language. The question here is whether the languages have *pattern-isomorphic* representations, not merely logically-equivalent representations.

 $<sup>^{14}\</sup>mbox{https://en.wikipedia.org/wiki/Permutation#Permutations_of_multisets:}$  Each condition is shown 8 times per half. Each of the 4 patterns is shown 2 times.

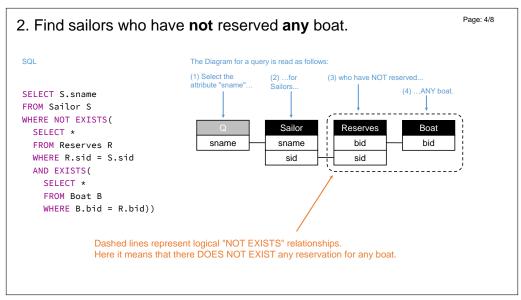


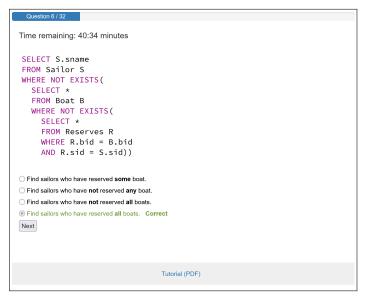
Fig. 30. One page of the 8-page tutorial for user study participants.

and each of 2 halves chosen independently from the 2520 possible sequences, irrespective of which conditions are used). We *sample a new treatment (group and sequence) for each participant*. I.e. recording the correct answers to one's treatment, and then sharing that information with another worker will not be very helpful to that worker. The chance that, among n participants, any two share the same questions (irrespective of condition) is approximately  $1 - e^{\frac{n^2/2}{2520^4}}$ , which is around 1/5000 for n = 50. In our study, none of the 171 workers who viewed the consent form were randomly assigned to the same treatment.

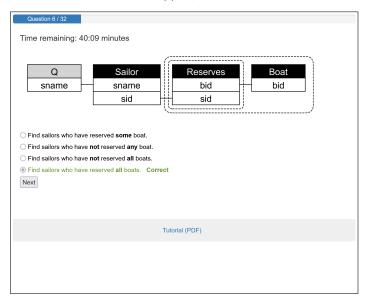
MTurk inclusion/rejection criteria. We recruited participants through Amazon Mechanical Turk (AMT), limiting participation to adult workers who live in the United States, have at least 500 completed tasks approved by requesters, have at least 97% of their completed tasks approved overall, and who self-determine themselves as experienced SQL users based on the following prompt in the HIT description: "Workers should be familiar with SQL at the level of an advanced undergraduate database class, in particular with nested SQL queries." The HIT is accepted if participants correctly answer at least 16 of the questions and submit within 50 minutes of starting. We purposefully do not use a qualification exam on SQL as time-consuming qualification exams discourage many workers from participating and are prone to cheating (standardized questions and answers tend to be shared on AMT-specific forums). Instead, we rely on self-determination of SQL experience. Workers know that failing to submit 16/32 correct questions will get their HIT rejected, meaning they do not get paid, and their overall approval rate goes down-possibly precluding them from doing other HITs. However, this approach can lead to many workers who skim or skip directions being rejected, lowering the requester's HIT approval rate. The requester approval rate is visible in the MTurk user interface, and workers may be wary of accepting HITs from a requester with a low rate of approving submissions.

**Payment for time and accuracy.** To make sure participants have a stake in the variables we measure (accuracy and speed), we motivate correct and fast answers by paying participants

<sup>&</sup>lt;sup>15</sup>https://en.wikipedia.org/wiki/Birthday problem#Approximations



#### (a) SQL



#### (b) Relational Diagrams

Fig. 31. The user interface a participant would see for answering questions. The shown stimuli are pattern (4) on the sailors-reserve-boats schema for both conditions. This particular schema was used for the tutorial, so we excluded it from the 32 questions on the actual study.

according to their answering correctly and fast (largely inspired by our earlier user study on QueryVis [55]). The base pay for an accepted submission is \$6.00 USD. (1) *Accuracy bonus*: For every correctly answered question after the 16<sup>th</sup>, the participant receives a bonus payment of \$0.20 USD for a maximum pay of \$9.20 USD. (2) *Speed bonus*: Based on total test completion time, the

participant will receive a percentage bonus on total pay (including the accuracy bonus). Completion within 11 minutes awards a 5% bonus for a total maximum pay of \$9.66. Each minute faster gets the participant an additional 5% bonus up to 40% for completing within 4 minutes, with a maximum pay of \$12.88. These values were determined using pilot data to target workers receiving a U.S. average living wage of  $$25.02 \text{ USD/hr}^{16}$  and average \$8.76 in total for our study.

Choice of 50 participants. Our goal was to recruit 50 successful participants. We made the task (HIT) available in two batches of 60. 171 MTurk workers viewed the consent form, 162 accepted the task and began the study, 146 answered at least one question, 133 completed all questions, and 120 submitted the task. 42 workers accepted the HIT and began the study but returned it without completing it (162 - 120). Only 58 of the 120 submissions had above the 50% accuracy threshold. Among those 58 participants, we used the data of the first 25 who started the task for each condition to get our planned 50 participants, balanced across two groups (Fig. 11).

**Stimuli and user interface.** Figure 31 demonstrates the interface a participant would see for answering questions during the study. (a) shows an example of a question with pattern (4) and the SQL condition, (b) shows the same question posed with Relational Diagrams. Note that we showed participants each schema only once—this figure is simply for illustrative purposes. A progress bar and countdown timer at the top inform participants of their progress. Upon selecting an answer and pressing 'Next', the chosen and correct answers are highlighted. A link at the bottom provides quick access to a PDF of the tutorial.

**Median vs. mean.** We carefully considered using the mean or the median as summary statistics. Each has pros and cons, and decided to use the median for time and the mean for accuracy. We considered 4 properties: (1) The median is more robust against outliers. A good example of this is how median income is more representative than mean income when the focus is on the experience of individual citizens instead of the overall performance of a country. The completion times for each question in our study are very imbalanced and have heavier tails than a Gaussian distribution. (2) On the other hand, the median has a slower convergence behavior. For example, when sampling from a Gaussian distribution, the 95% confidence interval (CI) for the median will be larger than the one for the mean. Thus to get the same statistical power, more data points are needed. (3) The mean is more appropriate for summarizing discrete choices. As an extreme example, consider a repeated Boolean event (how often does a coin come up head first). The median only summarizes which of the two events occurred more times, whereas the mean gives a better estimate of the probability (and even a lossless one together with the number of tosses). (4) When summarizing ratios (e.g., the time improvements in our case), the mean gives a biased estimate, even in the limit of infinitely many data points. For example, consider two conditions A and B, and three users, U1, U2, and U3, with recorded times for both conditions being:

User		A	В	Ratio B/A		
J	J1	200	100	0.5		
J	J2	150	150	1.0		
τ	J3	100	200	2.0		

Thus one user is faster using B, one user is slower, and one is the same in both conditions. However, the mean ratio is 1.17, whereas the median ratio is 1.0.

For these reasons, we are using the median for summarizing individual per-participant times and ratios of time improvement. However, we use the mean to summarize relative error rates and their differences.

 $<sup>^{16}</sup> https://livingwage.mit.edu/articles/103-new-data-posted-2023-living-wage-calculator, which is a superconductive of the control of the$ 

	Pattern	Median time or ratio	95% CI
RD	P1	9.25	[6.67, 10.47]
SQL	P1	11.62	[10.26, 13.35]
ratio RD/SQL	P1	.64	[.49, .78]
RD	P2	11.62	[10.26, 13.35]
SQL	P2	13.32	[11.95, 17.53]
ratio RD/SQL	P2	.83	[.70, .97]
RD	P3	10.89	[9.51, 13.22]
SQL	P3	14.13	[11.19, 21.01]
ratio RD/SQL	P3	.66	[.53, .77]
RD	P4	8.14	[7.02, 11.58]
SQL	P4	11.95	[10.60, 14.92]
ratio RD/SQL	P4	.71	[.60, .86]

Table 1. Median times and ratios along with 95% BCa confidence intervals for Fig. 32.

# O.2 Supplemental Material

All supplemental material required to reproduce our results from the data we collected or replicate our study with additional participants is available on the Open Science Framework (OSF). Here we provide links to key materials:

- The main OSF folder: https://osf.io/q9g6u/
- Study tutorial: https://osf.io/mruzw
- Stimuli-generating code: https://osf.io/kgx4y
- The stimuli: https://osf.io/d5qaj
- Stimuli/schema index CSV: https://osf.io/u8bf9
- Stimuli/schema index JSON: https://osf.io/sn83j
- Server code for hosting the study: https://osf.io/suj4a.
- Collected data: https://osf.io/8vm42
- The final analysis code: https://osf.io/f2xe3
- The analysis can be compared with our planned and preregistered analysis code: https://osf.io/4 zpsk/

#### O.3 Additional results (preregistered)

Here we describe an additional result included as part of the preregistration but not listed under its hypotheses. This result was demoted to the appendix for lack of space.

**Result 4.** (**Patterns and Speed**) For all 4 different patterns, we have strong evidence that participants were meaningfully faster at identifying each pattern using Relational Diagrams than SQL.

We also calculated the median per-participant time for identifying each of the 4 patterns in each of the two conditions. Figure 32 shows the interesting observation that pattern 4 with double negation ("Find *sailors* who have **not** *reserved* **all** *boats*") was not the one that took participants the longest to recognize. The median improvement of Relational Diagrams over SQL per pattern and participant is still statistically significant *for each pattern individually* (i.e. the 95% CIs were fully below the ratio 1.00). Please see Table 1 for details.

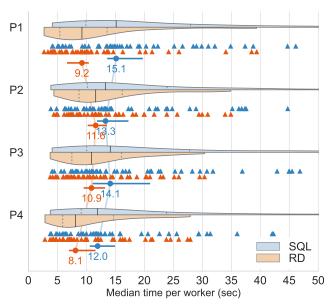


Fig. 32. Result (4): P1–P4 stand for patterns 1–4, respectively. The median time per worker for each pattern and condition is shown with bootstrapped 95% BCa confidence intervals and two views of the data distribution. It becomes clear that different patterns take different times to recognize, and Relational Diagrams allow participants to do this faster than SQL. Although confidence intervals per pattern overlap here, the analysis of the 95% confidence intervals of the ratios per participant and pattern show statistical significance. For a discussion of why these CIs shouldn't be directly compared, see [8, 54, 61].

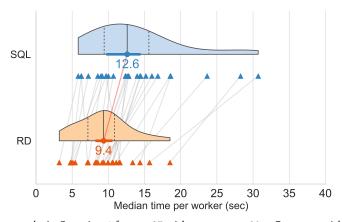


Fig. 33. Exploratory analysis: Question 1 for n = 27 with accuracy > 90%. Compare with Figure 12a (top).

## O.4 Exploratory analyses (not preregistered)

This section details an exploratory analysis of a subset of the accepted participants, which we conducted after collecting the data. We did not preregister these analyses. The results should be taken as preliminary findings that need to be backed up by summative studies.

(5) **Results with stricter accuracy requirement.** Our threshold for acceptance was 50% accuracy, thus a participant needs to correctly answer at least k = 16 among the n = 32 questions. What is the expected fraction of users who finish the HIT by completely randomly choosing from the 4

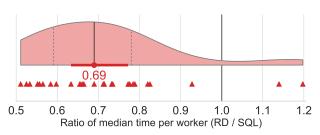


Fig. 34. Exploratory analysis: Question 1 for n = 27 with accuracy > 90%. Compare with Figure 12a (bottom).

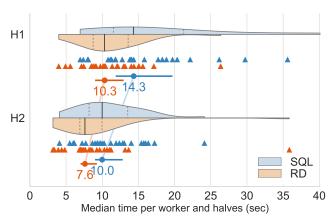


Fig. 35. Exploratory analysis: Question 2 for n = 27 with accuracy > 90%. Compare with Figure 12c.

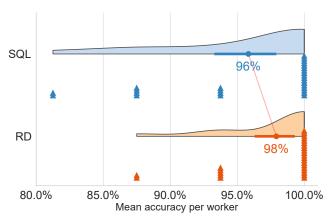


Fig. 36. Exploratory analysis: Question 3 for n = 27 with accuracy > 90%. Compare with Figure 12b (top).

available choices? Modeling the number of successful answers X as a random variable following a binomial distribution B(n,p) with p=1/4, the probability of getting exactly k right is  $\mathbb{P}[X=k]=\binom{n}{k}p^k(1-p)^{n-k}$ . We are interested in  $\mathbb{P}[X\geq k]=1-\mathbb{P}[X\leq k-1]$  which is 0.2% for n=32, k=16, p=1/4. Given that 162 MTurk users attempted the HIT, only 133 answered all the questions, of those only 58 got at least half correct, and some successful participants had low accuracy on the 16

 $<sup>^{17}</sup>$ This can be calculated explicitly e.g. with scipy.stats.binom (https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.binom.html).

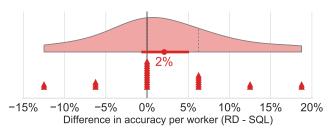


Fig. 37. Exploratory analysis: Question 3 for n = 27 with accuracy > 90%. Compare with Figure 12b (bottom).

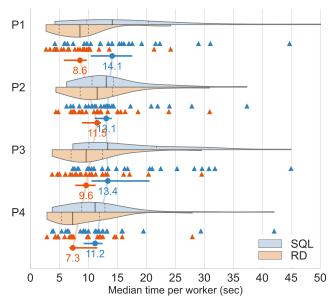


Fig. 38. Question 4 for n = 27 with accuracy > 90%. Compare with Figure 32.

questions shown in SQL (Fig. 12b), we asked how the results would change by only considering those among the first 50 users with an accuracy of 90% (thus maximally 3 incorrect answers).

Figures 33 to 35 and 38 show our preliminary finding that filtering to users with at least 90% correct answers has results similar to our planned analysis. The median time, confidence intervals, and data distribution are surprisingly similar despite the smaller sample size. We believe this is because the higher accuracy threshold reduces a lot of the randomness from users who barely passed (see a few participants with very low accuracy for SQL in Figure 12b (top)).

Regarding accuracy, Figures 36 and 37 show the expected difference caused by the increased threshold: there is now much less evidence for a difference in accuracy between conditions (which can partially be explained by the fact that most participants in that group actually perfect accuracies).

Notice that filtering the 50 users down to those with an accuracy of 90% or more led to a slight imbalance with 13 users in group 1 and 14 users in group 2.

# O.5 Results from pilot study with PhD students

Before preregistering our study design, we conducted a pilot study with n = 13 PhD students studying databases at a U.S. University. This helped us refine our hypotheses; test our website,

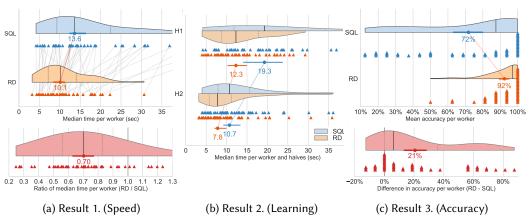


Fig. 39. This figure duplicates Fig. 12 to support direct comparison with data from the pilot testing in Fig. 40. See the caption of Fig. 12 and surrounding text for a description of the encodings used.

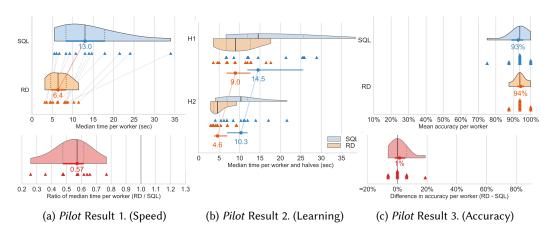


Fig. 40. Results of analyzing data from the n=13 pilot participants using the approach we preregistered for the full study. These pilot participants were PhD students studying databases. Compare with Fig. 39 (Fig. 12) where we present the main study results for the Amazon Mechanical Turk participants. Results 1 (Speed) and 2 (Learning) were consistent between the PhD students and MTurk workers. For Result 3 (Accuracy), the PhD students had much higher accuracy using SQL than the MTurk workers, and their accuracy was comparable when using Relational Diagrams.

database, and MTurk implementation code; and test our analysis code. Pilot testing also helped us to identify any confusing or incorrect stimuli.

To provide an exploratory comparison of our main results versus using student participants, we analyzed the pilot data using our preregistered analysis code and present the results in Fig. 40. Compare these with the main results in Fig. 12, duplicated here in Fig. 39 to support direct comparison. Note that these results should be considered *exploratory* rather than *summative* as the study design was not preregistered at that point. Moreover, the results are limited by our using an incorrect counterbalancing approach that we have since rectified. We do not believe the introduced ordering effects would dramatically affect the results, but we cannot rule out a possible effect.

Our Pilot Results 1 (Speed) and 2 (Learning) were consistent with what we found in the eventual study. However, Pilot Result 3 (Accuracy) stands in contrast.

**Pilot Result 1.** (Speed) We have some exploratory evidence that PhD database students were meaningfully faster at identifying patterns using Relational Diagrams than SQL: median ratio Relational Diagrams/SQL = 0.57, 95% CI [0.48, 0.62].

Figure 40a (top) shows the median per-participant times per condition (and overall median across participants), and Fig. 40a (bottom) the per-participant ratios between median times.

**Pilot Result 2.** (Learning) We have some exploratory evidence that PhD students got meaningfully faster during the study in both conditions.

Figure 40b shows the individual times for H1 ( $1^{st}$  half) and H2 ( $2^{nd}$  half), together with medians and CIs. Not shown are the median ratios H1/H2 we used for inference, which were 0.54, 95% CI [0.72, 0.80] for Relational Diagrams = and 0.59, 95% CI [0.43, 0.70] for SQL.

**Pilot Result 3.** (Accuracy) We have some exploratory evidence that PhD database students have comparable correctness with Relational Diagrams and SQL: mean difference in accuracy Relational Diagrams – SQL = 1%, 95% CI [-2%, 5%].

Figure 40c (top) shows the per-participant accuracy in each condition and the overall mean accuracies. Figure 40c (bottom) shows the per-participant difference in accuracy and overall mean. This last pilot finding is what led to our preregistered hypothesis that participants make a comparable number of correct responses using SQL and Relational Diagrams. However, it turns out that MTurk workers were considerably more often correct using Relational Diagrams than SQL. An additional summative study is necessary to say whether there is actually a difference between the proportion of correct responses provided by MTurk workers (who are presumably less skilled with SQL) and PhD students studying databases (who are presumably more skilled). It appears that may be true though, in which case a subsequent study could better answer the question: "Do Relational Diagrams help inexperienced SQL users identify common relational query patterns more accurately but only have comparable accuracy for experienced users?" Future studies could also better answer the question: "How comparable in SQL proficiency are self-described SQL experts on MTurk to students who get good grades in university database courses?"

# P MORE RELATED WORK FOR Section 7

#### P.1 Peirce's existential graphs (Section 7.1)

We discussed earlier that Lines of Identity (LIs) in beta graphs have multiple meanings (existential quantification and identity), and this "function overload" can make the graphs ambiguous. We now discuss this important point in more detail.

**Problems from abusing lines in beta graphs.** While over 100 years old, Peirce's beta system has led over the years to multiple misinterpretations and an ongoing discussions about how to interpret a valid beta graph correctly. The literature contains many attempts to provide formal "interpretations" and provide consistent readings of these graphs. How can it be that something supposed to be formal still allows so much ambiguities? In our opinion, beta graphs have one important design problem leading to those misunderstandings: it is the *overloading of the meaning of the Lines of Identity (LI)*, and thus one instances of an abuse of lines as symbols. As mentioned before, LIs are used to denote two different concepts: (i) the *existence* of objects (intuitively an existential quantification of a variable in DRC such as  $\exists x$ ), and (ii) the *identity* between objects (intuitively, R.A=S.A in TRC). This non-separation of concerns and function overload leads to unfixable ambiguities. We illustrate with a few examples.

— is a red boat color = 'red'

(a) Beta graph (b) Relational Diagram

Fig. 41. Example 24 illustrates that lines in beta graphs (called "Lines of Identities" or LIs) suffer from "function overload": LIs are used as a symbol for both equivalence between objects, and *existential quantification*. Relational Diagrams avoid this function overload by using lines only to represent built-in (comparison) predicates, and assuming proper predicates (relations with attributes) to be existentially quantified.

EXAMPLE 24 (A RED BOAT). Consider the sentence "There is a red boat" shown in Fig. 41a as a beta graph. As beta graphs cannot represent constants, the graph requires a special unary predicate "red boat." The LI represents both "there exists something" and "that something is equal to a red boat." Thus a line (which arguably suggests two items being connected or joined) is meant as a quantified variable, and the beta graph can be interpreted in DRC as:

$$\exists x [RedBoat(x)]$$

Figure 41b shows the same sentence as a Relational Diagram. Notice that "there exists something" is represented by just placing this something (a predicate) on the canvas. There is no need for an existential line. Also notice how the modern UML diagram allows predicates (relational atoms) with several attributes, and one of those attributes can be set equal to a constant (here  $\boxed{\text{color} = \text{`red'}}$ ). It can be read rather naturally like a TRC statement:

$$\exists b \in Boat[B.color=`red']$$

The interpretation of beta graphs where one LI represents one existentially-quantified variable can at times be intuitive and simply correspond to a modern DRC interpretation (see also Example 26). However, such a simple interpretation is not always possible.

Example 25 (Exactly one red boat). Consider the sentence "There exists exactly one red boat." Figures 42a and 42b show two beta graphs with different cut nestings that can both be read as

$$\exists x [RedBoat(x) \land \neg (\exists y [RedBoat(y) \land x \neq y])]$$

Now a single LI needs to represent two existentially-quantified variables, and two different nestings of the cuts can represent the same statement. Contrast this with Fig. 42c, read in TRC as:

$$\exists b \in Boat[B.color = `red' \land \neg (\exists b_2 \in Boat \\ [b_2.color = `red' \land b.bid \neq b_2.bid])]$$

Notice here that the inequality is simply represented by a label of a join between two predicates. Two tuple variables are represented by two different atoms and the interpretation is unambiguous: There exists a boat whose color is red, and there does not exist another boat whose color is red and whose bid is different.

The fact that one LI can branch into multiple endings (also called *ligatures*), may have *loose endings*, and may represent multiple existentially-quantified variables, together with cuts being applied to such LI's can quickly lead to hard-to-interpret diagrams (see e.g., the increasingly-unreadable figures in [71, pp. 42-49]). This led to several attempts in the literature to provide "reading algorithms" of those graphs (e.g., [68, 71, 80]) and rather complicated proofs of the expressiveness of beta graphs [80], assuming a correct reading. For example, the paper by Dau [28] points out an error in Shin's reading algorithm [71]. However, Dau's correction to Shin [28] itself also has errors. For example, the interpretation of the right-most diagram in [28, Fig 2] (reproduced as Fig. 43b) is

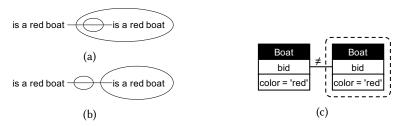


Fig. 42. Example 25: The combination of Ll's and nesting symbols (called "cuts" in beta graphs) provide ambiguous ways to nest cuts.

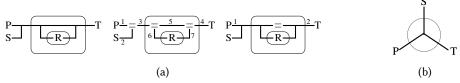


Fig. 43. (a): Figure copied from Dau [28] discussing an example beta graph whose interpretation provided by Shin [71] is incorrect, together with two alternative ways of splitting the LIs in order to interpret the graph correctly. The details of the arguments are intricate and not important here. What matters is that a lot of disagreement exists as to how to interpret LI's correctly. Relational Diagrams avoid this problem entirely by using lines only as comparison predicates. (b): Right-most diagram of Figure 2 in Dau [28]: a difficult-to-interpret Peirce beta graph. The reading algorithm presented is wrong and misses one equality.

wrong and misses one equality. The given interpretation is

$$\exists x. \exists y. \exists z [S(x) \land P(y) \land T(z) \land \neg (x = y \land y = z)]$$

whereas it should be

$$\exists x. \exists y. \exists z [S(x) \land P(y) \land T(z) \land \neg (x = y \land y = z \land x = z)]$$

This is just an intuitive example how difficult beta graphs are in practice to interpret, even by the experts, and even by experts pointing out errors from other experts.

Why Relational Diagrams avoid the problem. Relational Diagrams use the line only for connecting two attributes. The type of connection is unambiguously represented by a label. Quantification is represented by predicates themselves. Thus, on a more philosophical level, we think that our visual formalism solves those problems based on a more modern interpretation of first order logic: TRC was created by Edgar Codd in the 1960s and 1970s in order to provide a declarative database-query language for data manipulation in the relational data model [21]. In contrast, beta graphs were proposed even before first-order logic, which was only clearly articulated some years after Peirce's death in the 1928 first edition of David Hilbert and Wilhelm Ackermann's "Grundzüge der theoretischen Logik" [48]. Zeman, in his 1964 PhD thesis [80], was the first to note that beta graphs are isomorphic to first-order logic with equality. However, the secondary literature, especially Roberts [68] and Shin [71], does not agree on just how this is so [23]. We did not start from Peirce's beta graphs and attempt to fix issues that have been occupying a whole community for years. Rather, we started from the modern UML reading of relational schemas and an understanding of TRC, and tried to achieve a minimal visual extension to provide relational completeness and pattern-isomorphism to TRC. This happens to provide a natural solution for the interpretation problems of beta graphs. We believe that Relational Diagrams provide a clean, unambiguous, and, in hindsight, simple abstraction of query patterns.

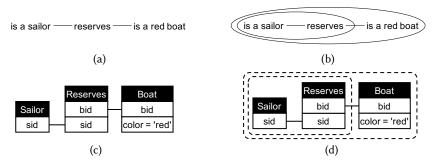


Fig. 44. Example 26: Diagrams comparing the representations of negation in beta graphs (top) and Relational Diagrams (bottom).

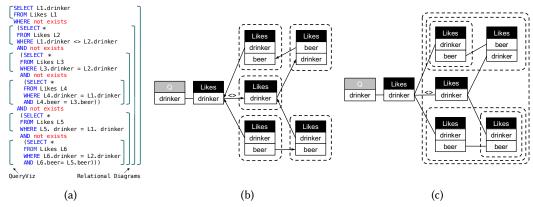


Fig. 45. Illustrations for Example 27: (a) Unique-set-query "Find drinkers with a unique beer taste" used by [55], (b) QueryVis diagram with reading order encoded by arrows (redrawn according to [55]), (c) Relational Diagrams with a nested scoping and no need for arrows.

Example 26 (Nested negation). Figure 44 shows 2 beta graphs:

a: There exists a sailor who reserved a red boat.

**b**: All red boats were reserved by some sailor.

Beta graphs cannot represent constants and thus need to replace a selection of boats that are red with a dedicated new predicate "is a red boat." Their respective translations into DRC are:

a: 
$$\exists x, y[Sailor(x) \land RedBoat(y) \land Reserves(x, y)]$$
  
b:  $\neg(\exists y[RedBoat(y) \land \neg(\exists x[Sailor(x) \land Reserves(x, y)])])$ 

Contrast the beta graphs with their respective Relational Diagrams in Fig. 44 and TRC:

c: 
$$\exists s \in Sailor, b \in Boat, r \in Reserves[r.bid = b.bid \land r.sid = s.sid \land b.color = 'red']$$
d:  $\neg(\exists b \in Boat[B.color = 'red' \land \neg(\exists s \in Sailor, r \in Reserves])$ 
(19)

$$[r.bid = b.bid \land r.sid = s.sid])]) \tag{20}$$

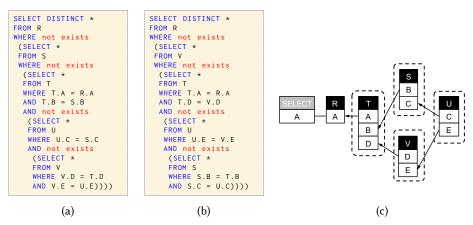


Fig. 46. Example 28: Minimal example showing that QueryVis is not sound for nested queries with 4 levels: Two different queries (a), (b) that are translated into the same QueryVis diagram (c).

## P.2 QueryVis (Section 7.2)

We use the "unique beer taste" query that was proposed in the QueryVis paper [55] to show the difference in design decisions.

Example 27 (Unique-set-query). Consider the SQL query from Fig. 45a asking to find "drinkers who like a unique set of beers," i.e. no other drinker likes the exact same set of beers. The scoping brackets to the left of the query in Fig. 45a show the content of boxes used by QueryVis, which include all tables from each individual query block. Without the additional visual symbol of arrows, this diagram becomes ambiguous to interpret. To mitigate this problem, the design of QueryVis [26, 55] uses directed arrows with an implied reading order (Fig. 45b).

The scoping brackets to the right in Fig. 45a show the nesting of the variables scopes in queries, which are also reflected in the dashed bounding boxes in Relational Diagrams (Fig. 45c).

The design decision in QueryVis [55] are justified in terms of usability (for "most" queries the diagrams are not ambiguous and the reduction in nesting simplifies their interpretation), yet requires overloading of the meaning of arrows. Two conceptual problems with these diagrams are: (1) QueryVis requires each partition of the canvas to contain a relation from the relational schema. Our earlier examples from Figs. 18 and 28 show examples that can thus not be handled. (2) QueryVis does not guarantee unambiguous visualizations for nested queries with nesting depth  $\geq$  4. This was alluded to already in [55], and we next give an example to illustrate:

Example 28 (Ambiguous QueryVis). We next give a minimum example for when QueryVis becomes ambiguous. Consider the two different SQL queries Figs. 46a and 46b. Following the algorithm given in [55], both lead to the same visual representation Fig. 46c. In other words, it is not possible to uniquely interpret the diagram in Fig. 46c.

**Problems from abusing lines in QueryVis.** Similar to beta graphs, we think, that the design of QueryVis abuses the line symbol by using it for two purposes: for (*i*) joining atoms and (*ii*) for representing the negation hierarchy. In contrast, Relational Diagrams use the line only for connecting two attributes and represent the negation hierarchy explicitly by nesting negation boxes. Relational Diagrams fix the completeness and soundness issues, and, in addition, can show logical sentences and queries or sentences lacking tables in one or more of the negation scopes of nested queries.

_ld	
Boat bid bname color BadSids sid	
_B red IId	

(a) Temporary relation BadSids(sid) ("I." stands for insert)

Sailor	sid	sname	rating	age	BadSids	sid
	_ld	PS			7	_ld

(b) Actual answer Q(sid) ("P." stands for print)

Fig. 47. Example 29: QBE needs to create a temporary relation BadSids in order to express relational division. It does follow the query pattern of relational algebra and not relational calculus.

# P.3 Query-By-Example (QBE)

The development of QBE [81] was strongly influenced by DRC. However, QBE can express relational division only by using COUNT or by breaking the query into two logical steps and using a temporary relation [64, online appendix]. But in doing so, QBE uses the query pattern from RA and Datalog of implementing relational division (or universal quantification) in a dataflow-type, sequential manner, requiring two occurrences of the Sailor table.

EXAMPLE 29 (SAILORS RESERVING ALL RED BOATS IN QBE). Consider the query "Find sailors who have reserved all red boats" QBE needs to first create an intermediate relation "BadSids" that stores all Sailors for whom there is a red boat that is not reserved by the sailor (Fig. 47a), and then finds all the other Sailors (Fig. 47b). The pattern of this query in QBE thus matches exactly the one of Datalog (it requires two occurrences of the relational Sailor instead of one as in calculus), which is arguably a more dataflow (one relation accessed after the other) than logical query language pattern:

$$I_{1}(x, y) := Reserves(x, y, \_).$$

$$I_{2}(x) := Sailor(x, \_, \_), Boat(y, \_, 'red'), \neg I_{1}(x, y).$$

$$Q(y) := Sailor(x, y, \_, \_), \neg I_{2}(x).$$
(21)

More formally, the QBE query from Fig. 47 is pattern-isomorphic to the Datalog query in (21). Furthermore, the following logically-equivalent TRC query (22) has no pattern-isomorphic representation in QBE: (i.e. with one single occurrence of the Sailor relation).

$$\{q(\text{sname}) \mid \exists s \in \text{Sailor}[q.\text{sname} = s.\text{sname} \land \\ \neg(\exists b \in \text{Boat}[b.\text{color} = \text{`red'} \land \\ \neg(\exists r \in \text{Reserves}[r.\text{bid} = b.\text{bid} \land r.\text{sid} = s.\text{sid}])]) \}$$
 (22)

#### P.4 DFQL (Section 7.3)

DFQL (Dataflow Query Language) is an example visual representation that is relationally complete [12, 20, 44] by mapping its visual symbols to the operators of relational algebra. Aside from providing a basic set of operators derived from the requirements for being as expressive as first-order predicate calculus, DFQL also provides a diagrammatic representation of grouping operators in both comparison functions and aggregations. In contrast to several other similarly-motivated visual query languages that represent operators of relational algebra with different visual symbols, a detailed Master's thesis [44] lists the language and its constructs in enough detail to allow us to create visualizations of new queries in DFQL. Following the same procedurality as RA, DFQL expresses the dataflow in a top-down tree-like structure. However, since DFQL focuses on the 1-to-1

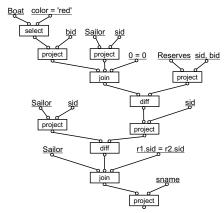


Fig. 48. Example 30: DFQL visualization of the query "Find sailors who have reserved all red boats" that is structurally equivalent to the pattern used by relational algebra. The *diff* operator is equivalent to binary – (minus) in RA and the tautology "0 = 0" in *join* operator is required to create a Cartesian Join in DFQL [20]. Compare the difficulty in perceiving a logical pattern in this visualization against the one from Relational Diagrams in Fig. 17c.

correspondence to relational algebra, it also can not generate a pattern-isomorphic diagram for query Fig. 24a which has no pattern-isomorphic representation in RA. See the following example for details:

Example 30 (Sailors reserving all red boats" (recall Example 14 and Fig. 17c) in the formalism of DFQL, the entire query can be visualized in one single connected tree-like diagram (unlike QBE which needs to visualize a temporary table to hold the intermediary values). However, since the language is based on the operators of RA, there is no pattern-isomorphic expression of the query (Fig. 17c) in relational algebra. Instead, the logically-equivalent representation in RA is as follows:

$$Q = \pi_{sname} \left( Sailor \bowtie \left( \pi_{sid} Sailor - \pi_{sid} \left( (\pi_{sid} Sailor + \pi_{bid} \sigma_{color='red'} Boat - \pi_{sid,bid} Reserves \right) \right) \right)$$
(23)

The join between Sailor S and Sailor S2 is necessary to project column sname from the table. This later query can be visualized by DFQL in a pattern-preserving way as Fig. 48. One can easily find a 1-to-1 mapping between DFQL operators and this RA expression.

Notice that for the same arguments, there is also no pattern-isomorphic expression of the query shown in Fig. 13g and DFQL needs two extensional tables for input table R to represent that query.

# P.5 Tools for query visualizations

The four projects that we know of that focus on the problem of visualizing existing relational queries are QueryVis [26, 38, 55] (which we showed is not relationally complete, yet which inspired a lot of our work), GraphSQL [16], Visual SQL [50], (both of which maintain the 1-to-1 correspondence to SQL, and syntactic variants of the same query like Fig. 15 lead to different representations), and Snowflake join [75] (which is a pure query visualization approach that focuses on join queries with optional grouping, but does not support any nested queries with negation). Compared to all these visual representations, ours is the only one that is relationally complete and that can preserve and represent all logical patterns in the non-disjunctive fragment of relational query languages.

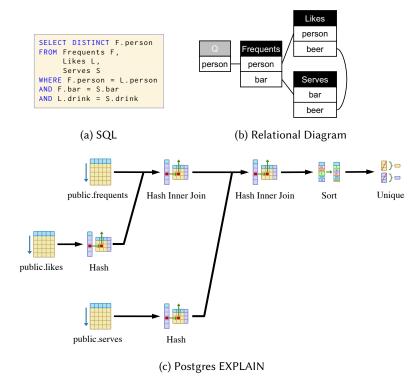


Fig. 49. Example 31: Find persons who frequent some bar that serves some drink they like.

# P.6 Tools for query plan visualizations

Relational Diagrams are a query visualization and thus conceptually different from a *query plan visualization* that many SQL users will be familiar from the EXPLAIN command in PostgreSQL [63]. We explain with an example used with kind permission from the authors of [42].

Example 31 (Cyclic query). Consider the following cyclic query over the beer drinkers database introduced by Ullman [77]:

```
\{q.person \mid \exists f \in Frequents, \exists l \in Likes, \exists s \in Serves \\ [q.person = f.person \land f.person = l.person \land \\ l.drink = s.drink, \land s.bar = f.bar]\}
```

The query asks for drinkers who frequent bars that serve some beer they like. Figure 49a shows the query in SQL and Fig. 49b as Relational Diagram.

Figure 49c shows a query plan chosen by PostgreSQL [63]. Notice that the produced query plans does not captures the cyclic nature of the join of the query and instead shows the query as a tree.

A *query plan visualization* targets the physical query execution and represents HOW a query is executed and often helps reason the user about ways to make the query run faster. In contrast, a query visualization, such as Relational Diagrams, attempts to represent WHAT a query does (i.e. its intent) and possibly the relational pattern it uses. Contrast the query visualization in Fig. 49b, which shows the join pattern and that this query is cyclic.

Similarly, query visualizations are also different from *query dashboards*, i.e. tools that help monitor key characteristics of queries (such as Vertica analyzer [73]) or visualize and compare the cost or speed of execution plans (such as Picasso [47]).

# P.7 Applications for query interpretation

Query Interpretation is the problem of reading and understanding an existing query. It is often as hard as query composition, i.e., creating a new query [65]. In the past, several projects have focused on building Query Management Systems that help users issue queries by leveraging an existing log of queries. Known systems to date include CQMS [52, 53], SQL QuerIE [4, 18], DBease [56], and SQLShare [49]. All of those are motivated by making SQL composition easier and thus databases more usable [51], especially for non-sophisticated database users. An essential ingredient of such systems is a query browse facility, i.e., a way that allows the user to browse and quickly choose between several queries proposed by the system. This, in turn, requires a user to quickly understand existing queries.

Whereas visual systems for *specifying* queries have been studied extensively (a 1997 survey by Catarci et al. [12] cites over 150 references), the explicit reverse problem of visualizing and thereby helping *interpret a relational query that has already been written* has not drawn much attention, despite very early [65, 66] and very recent work [55] repeatedly showing that visualizations of relational queries can help users understand them faster than SQL text.