

---

# EFFICIENT EVALUATION OF ARBITRARY RELATIONAL CALCULUS QUERIES

MARTIN RASZYK <sup>a</sup>, DAVID BASIN <sup>b</sup>, SRĐAN KRSTIĆ <sup>b</sup>, AND DMITRIY TRAYTEL <sup>c</sup>

<sup>a</sup> DFINITY, Zurich, Switzerland  
*e-mail address:* martin.raszyk@dfinity.org

<sup>b</sup> Department of Computer Science, ETH Zürich, Zurich, Switzerland  
*e-mail address:* {basin, srdan.krstic}@inf.ethz.ch

<sup>c</sup> Department of Computer Science, University of Copenhagen, Copenhagen, Denmark  
*e-mail address:* traytel@di.ku.dk

---

**ABSTRACT.** The relational calculus (RC) is a concise, declarative query language. However, existing RC query evaluation approaches are inefficient and often deviate from established algorithms based on finite tables used in database management systems. We devise a new translation of an arbitrary RC query into two safe-range queries, for which the finiteness of the query’s evaluation result is guaranteed. Assuming an infinite domain, the two queries have the following meaning: The first is closed and characterizes the original query’s relative safety, i.e., whether given a fixed database, the original query evaluates to a finite relation. The second safe-range query is equivalent to the original query, if the latter is relatively safe. We compose our translation with other, more standard ones to ultimately obtain two SQL queries. This allows us to use standard database management systems to evaluate arbitrary RC queries. We show that our translation improves the time complexity over existing approaches, which we also empirically confirm in both realistic and synthetic experiments.

## 1. INTRODUCTION

Codd’s theorem states that all domain-independent queries of the relational calculus (RC) can be expressed in relational algebra (RA) [Cod72]. A popular interpretation of this result is that RA suffices to express all interesting queries. This interpretation justifies why SQL evolved as the practical database query language with the RA as its mathematical foundation. SQL is declarative and abstracts over the actual RA expression used to evaluate a query. Yet, SQL’s syntax inherits RA’s deliberate syntactic limitations, such as union-compatibility, which ensure domain independence. RC does not have such syntactic limitations, which arguably makes it a more attractive declarative query language than both RA and SQL. The main problem of RC is that it is not immediately clear how to evaluate even domain-independent queries, much less how to handle the domain-dependent (i.e., not domain-independent) ones.

As a running example, consider a shop in which brands (unary finite relation  $B$  of brands) sell products (binary finite relation  $P$  relating brands and products) and products are reviewed by users with a score (ternary finite relation  $S$  relating products, users, and scores). We

---

*Key words and phrases:* Relational calculus, relative safety, safe range, query translation.

consider a brand *suspicious* if there is a user and a score such that all the brand’s products were reviewed by that user with that score. An RC query computing suspicious brands is

$$Q^{susp} := B(b) \wedge \exists u, s. \forall p. P(b, p) \longrightarrow S(p, u, s).$$

This query is domain independent and follows closely our informal description. It is not, however, clear how to evaluate it because its second conjunct is domain dependent as it is satisfied for every brand that does not occur in  $P$ . Finding suspicious brands using RA or SQL is a challenge, which only the best students from an undergraduate database course will accomplish. We give away an RA answer next (where  $-$  is the set difference operator and  $\triangleright$  is the anti-join, also known as the *generalized* difference operator [AHV95]):

$$\pi_{brand}((\pi_{user, score}(S) \times B) - \pi_{brand, user, score}((\pi_{user, score}(S) \times P) \triangleright S)) \cup (B - \pi_{brand}(P))).$$

The highlighted expressions  $\pi_{user, score}(S)$  are called *generators*. They ensure that the left operands of the anti-join and set difference operators include or have the same columns (i.e., are union-compatible) as the corresponding right operands. (Following Codd [Cod72], one could also use the active domain to obtain canonical, but far less efficient, generators.)

Van Gelder and Topor [GT87, GT91] present a translation from a decidable class of domain-independent RC queries, called *evaluable*, to RA expressions. Their translation of the evaluable  $Q^{susp}$  query would yield different generators, replacing both highlighted parts by  $\pi_{user}(S) \times \pi_{score}(S)$ . That one can avoid this Cartesian product as shown above is subtle: Replacing only the first highlighted generator with the product results in an inequivalent RA expression.

Once we have identified suspicious brands, we may want to obtain the users whose scoring made the brands suspicious. In RC, omitting  $u$ ’s quantifier from  $Q^{susp}$  achieves just that:

$$Q_{user}^{susp} := B(b) \wedge \exists s. \forall p. P(b, p) \longrightarrow S(p, u, s).$$

In contrast, RA cannot express the same property as it is domain dependent (hence also not evaluable and thus out of scope for Van Gelder and Topor’s translation):  $Q_{user}^{susp}$  is satisfied for every user if a brand has no products, i.e., it does not occur in  $P$ . Yet,  $Q_{user}^{susp}$  is satisfied for finitely many users on every database instance where  $P$  contains at least one row for every brand from the relation  $B$ , in other words  $Q_{user}^{susp}$  is *relatively safe* on such database instances.

How does one evaluate queries that are not evaluable or even domain dependent? The main approaches from the literature (Section 2) are either to use variants of the active domain semantics [BL00, HS94, AGSS86] or to abandon finite relations entirely and evaluate queries using finite representations of infinite (but well-behaved) relations such as systems of constraints [Rev02] or automatic structures [BG04]. These approaches favor expressiveness over efficiency. But unlike query translations, they cannot benefit from decades of practical database research and engineering.

In this work, we translate arbitrary RC queries to RA expressions under the assumption of an infinite domain. To deal with queries that are domain dependent, our translation produces two RA expressions, instead of a single equivalent one. The first RA expression characterizes the original RC query’s relative safety, the decidable question of whether the query evaluates to a finite relation for a given database, which can be the case even for a domain-dependent query, e.g.,  $Q_{user}^{susp}$ . If the original query is relatively safe on a given database, i.e., produces some finite result, then the second RA expression evaluates to the same finite result. Taken together, the two RA expressions solve the *query capturability* problem [AH91]: they allow us to enumerate the original RC query’s finite evaluation result, or to learn that it would be infinite using RA operations on the unmodified database.

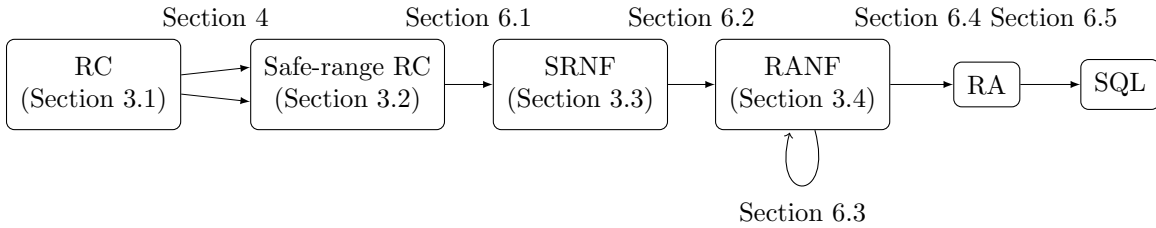


Figure 1: Overview of our translation.

Figure 1 summarizes our translation’s steps and the sections where they are presented. Starting from an RC query, it produces two SQL queries via transformations to safe-range queries, the safe-range normal form (SRNF), the relational algebra normal form (RANF), and RA, respectively (Section 3). This article’s main contribution is the first step: translating an RC query into two safe-range RC queries (Section 4), which fundamentally differs from Van Gelder and Topor’s approach and produces better generators, like  $\pi_{user,score}(S)$  above. Our generators strictly improve the time complexity of query evaluation (Section 5).

After the standard transformations from safe-range to RANF queries and from there to RA expressions, we translate the RA expressions into SQL using the `radb` tool [Yan19] (Section 6). We leverage various ideas from the literature to optimize the overall result. For example, we generalize Claußen et al. [CKMP97]’s approach to avoid evaluating Cartesian products like  $\pi_{user,score}(S) \times P$  in RANF queries by using count aggregations (Section 6.3).

The translation to SQL enables any standard database management system (DBMS) to evaluate RC queries. We implement our translation and then use either PostgreSQL or MySQL for query evaluation. Using a real Amazon review dataset [NLM19] and our synthetic benchmark that generates hard database instances for random RC queries (Section 7), we evaluate our translation’s performance (Section 8). The evaluation shows that our approach outperforms Van Gelder and Topor’s translation (which also uses a standard DBMS for evaluation) and other RC evaluation approaches based on constraint databases and structure reduction.

In summary, our three main contributions are as follows:

- We devise a translation of an arbitrary RC query into a pair of RA expressions as described above. The time complexity of evaluating our translation’s results improves upon Van Gelder and Topor’s approach [GT91].
- We implement our translation and extend it to produce SQL queries. The resulting tool RC2SQL makes RC a viable input language for any standard DBMS. We evaluate our tool on synthetic and real data and confirm that our translation’s improved asymptotic time complexity carries over into practice.
- To challenge RC2SQL (and its competitors) in our evaluation, we devise the *Data Golf* benchmark that generates hard database instances for randomly generated RC queries.

This article extends our ICDT 2022 conference paper [RBKT22b] with a more complete description of the translation. In particular, it describes the steps that follow our main contribution – the translation of RC queries into two safe-range queries. In addition, we formally verify the functional correctness (but not the complexity analysis) of the main contribution using the Isabelle/HOL proof assistant [RT22]. The theorems and examples that have been verified in Isabelle are marked with a special symbol ( $\mathfrak{S}$ ). The formalization helped us identify and correct a technical oversight in the algorithm from the conference paper (even though the problem was compensated for by the subsequent steps of the translation in our implementation).

## 2. RELATED WORK

We recall Trakhtenbrot’s theorem and the fundamental notions of *capturability* and *data complexity*. Given an RC query over a *finite* domain, Trakhtenbrot [Tra50] showed that it is undecidable whether there exists a (finite) structure and a variable assignment satisfying the query. In contrast, the question of whether a *fixed* structure and a *fixed* variable assignment satisfies the given RC query is decidable [AGSS86].

Kifer [Kif88] calls a query class *capturable* if there is an algorithm that, given a query in the class and a database instance, enumerates the query’s evaluation result, i.e., all tuples satisfying the query. Avron and Hirshfeld [AH91] observe that Kifer’s notion is restricted because it requires every query in a capturable class to be domain independent. Hence, they propose an alternative definition that we also use: A query class is *capturable* if there is an algorithm that, given a query in the class, a (finite or infinite) domain, and a database instance, determines whether the query’s evaluation result on the database instance over the domain is finite and enumerates the result in this case. Our work solves Avron and Hirshfeld’s *capturability* problem additionally assuming an infinite domain.

Data complexity [Var82] is the complexity of recognizing if a tuple satisfies a fixed query over a database, as a function of the database size. Our *capturability* algorithm provides an upper bound on the data complexity for RC queries over an infinite domain that have a finite evaluation result (but it cannot decide if a tuple belongs to a query’s result if the result is infinite).

Next, we group related approaches to evaluating RC queries into three categories.

**Structure reduction.** The classical approach to handling arbitrary RC queries is to evaluate them under a finite structure [Lib04]. The core question here is whether the evaluation produces the same result as defined by the natural semantics, which typically considers infinite domains. Codd’s theorem [Cod72] affirmatively answers this question for domain-independent queries, restricting the structure to the *active domain*. Ailamazyan et al. [AGSS86] show that RC is a *capturable* query class by extending the active domain with a few additional elements, whose number depends only on the query, and evaluating the query over this finite domain. *Natural-active collapse* results [BL00] generalize Ailamazyan et al.’s [AGSS86] result to extensions of RC (e.g., with order relations) by combining the structure reduction with a translation-based approach. Hull and Su [HS94] study several semantics of RC that guarantee the finiteness of the query’s evaluation result. In particular, the “output-restricted unlimited interpretation” only restricts the query’s evaluation result to tuples that only contain elements in the active domain, but the quantified variables still range over the (finite or infinite) underlying domain. Our work is inspired by all these theoretical landmarks, in particular Hull and Su’s work (Section 4.1). Yet we avoid using (extended) active domains, which make query evaluation impractical.

**Query translation.** Another strategy is to translate a given query into one that can be evaluated efficiently, for example as a sequence of RA operations on finite tables. Van Gelder and Topor pioneered this approach [GT87, GT91] for RC. A core component of their translation is the choice of generators, which replace the active domain restrictions from structure reduction approaches and thereby improve the time complexity. Extensions to scalar and complex function symbols have also been studied [EHJ93, LYL08]. All these approaches focus on syntactic classes of RC, for which domain independence is given, e.g., the *evaluable* queries of Van Gelder and Topor (Appendix A). Our approach is inspired by Van Gelder and Topor’s work but generalizes it to handle arbitrary RC queries at the cost

of assuming an infinite domain. Also, we further improve the time complexity of Van Gelder and Topor’s approach by choosing better generators.

**Evaluation with infinite relations.** Constraint databases [Rev02] obviate the need for using RA operations on finite tables. This yields significant expressiveness gains as domain independence need not be assumed. Yet the efficiency of the quantifier elimination procedures employed cannot compare with the simple evaluation of the RA’s projection operation. Similarly, automatic structures [BG04] can represent the results of arbitrary RC queries finitely, but struggle with large quantities of data. We demonstrate this in our evaluation where we compare our translation to several modern incarnations of the above approaches, all based on binary decision diagrams [MLAH99, Møl02, CGS09, KM01, BKMZ15].

### 3. PRELIMINARIES

We introduce the RC syntax and semantics and define relevant classes of RC queries.

**3.1. Relational Calculus.** A signature  $\sigma$  is a triple  $(\mathcal{C}, \mathcal{R}, \iota)$ , where  $\mathcal{C}$  and  $\mathcal{R}$  are disjoint finite sets of constant and predicate symbols, and the function  $\iota : \mathcal{R} \rightarrow \mathbb{N}$  maps each predicate symbol  $r \in \mathcal{R}$  to its arity  $\iota(r)$ . Let  $\sigma = (\mathcal{C}, \mathcal{R}, \iota)$  be a signature and  $\mathcal{V}$  a countably infinite set of variables disjoint from  $\mathcal{C} \cup \mathcal{R}$ . The following grammar defines the syntax of RC queries:

$$Q ::= \perp \mid \top \mid x \approx t \mid r(t_1, \dots, t_{\iota(r)}) \mid \neg Q \mid Q \vee Q \mid Q \wedge Q \mid \exists x. Q.$$

Here,  $r \in \mathcal{R}$  is a predicate symbol,  $t, t_1, \dots, t_{\iota(r)} \in \mathcal{V} \cup \mathcal{C}$  are terms, and  $x \in \mathcal{V}$  is a variable. We write  $\exists \vec{v}. Q$  for  $\exists v_1. \dots \exists v_k. Q$  and  $\forall \vec{v}. Q$  for  $\neg \exists \vec{v}. \neg Q$ , where  $\vec{v}$  is a variable sequence  $v_1, \dots, v_k$ . If  $k = 0$ , then both  $\exists \vec{v}. Q$  and  $\forall \vec{v}. Q$  denote just  $Q$ . Quantifiers have lower precedence than conjunctions and disjunctions, e.g.,  $\exists x. Q_1 \wedge Q_2$  means  $\exists x. (Q_1 \wedge Q_2)$ . We use  $\approx$  to denote the equality of terms in RC to distinguish it from  $=$ , which denotes syntactic object identity. We also write  $Q_1 \longrightarrow Q_2$  for  $\neg Q_1 \vee Q_2$ . However, writing  $Q_1 \vee Q_2$  for  $\neg(\neg Q_1 \wedge \neg Q_2)$  would complicate later definitions, e.g., the safe-range queries (Section 3.2).

We define the subquery partial order  $\sqsubseteq$  on queries as the (reflexive and transitive) subterm relation on the datatype of RC queries. For example,  $Q_1$  is a subquery of the query  $Q_1 \wedge \neg \exists y. Q_2$ . We denote by  $\text{sub}(Q)$  the set of subqueries of a query  $Q$ , by  $\text{fv}(Q)$  the set of *free variables* in  $Q$ , and by  $\text{av}(Q)$  be the set of all (free and bound) variables in a query  $Q$ . Furthermore, we denote by  $\vec{\text{fv}}(Q)$  the sequence of free variables in  $Q$  based on some fixed ordering of variables. We lift this notation to sets of queries in the standard way. A query  $Q$  with no free variables, i.e.,  $\text{fv}(Q) = \emptyset$ , is called *closed*. Queries of the form  $r(t_1, \dots, t_{\iota(r)})$  and  $x \approx c$  are called *atomic predicates*. We define the predicate  $\text{ap}(\cdot)$  characterizing atomic predicates, i.e.,  $\text{ap}(Q)$  is true iff  $Q$  is an atomic predicate. Queries of the form  $\exists \vec{v}. r(t_1, \dots, t_{\iota(r)})$  and  $\exists \vec{v}. x \approx c$  are called *quantified predicates*. We denote by  $\tilde{\exists} x. Q$  the query obtained by existentially quantifying a variable  $x$  from a query  $Q$  if  $x$  is free in  $Q$ , i.e.,  $\tilde{\exists} x. Q := \exists x. Q$  if  $x \in \text{fv}(Q)$  and  $\tilde{\exists} x. Q := Q$  otherwise. We lift this notation to sets of queries in the standard way. We use  $\tilde{\exists} x. Q$  (instead of  $\exists x. Q$ ) when constructing a query to avoid introducing bound variables that never occur in  $Q$ .

A structure  $\mathcal{S}$  over a signature  $(\mathcal{C}, \mathcal{R}, \iota)$  consists of a non-empty domain  $\mathcal{D}$  and interpretations  $c^{\mathcal{S}} \in \mathcal{D}$  and  $r^{\mathcal{S}} \subseteq \mathcal{D}^{\iota(r)}$ , for each  $c \in \mathcal{C}$  and  $r \in \mathcal{R}$ . We assume that all the relations  $r^{\mathcal{S}}$  are *finite*. Note that this assumption does *not* yield a finite structure (as defined in finite model theory [Lib04]) since the domain  $\mathcal{D}$  can still be infinite. A (*variable*) *assignment* is a mapping  $\alpha : \mathcal{V} \rightarrow \mathcal{D}$ . We extend  $\alpha$  to constant symbols  $c \in \mathcal{C}$  with  $\alpha(c) = c^{\mathcal{S}}$ . We write  $\alpha[x \mapsto d]$  for the assignment that maps  $x$  to  $d \in \mathcal{D}$  and is otherwise identical to  $\alpha$ . We lift this

$$\begin{array}{l}
(\mathcal{S}, \alpha) \not\models \perp; (\mathcal{S}, \alpha) \models \top; \\
(\mathcal{S}, \alpha) \models r(t_1, \dots, t_{\iota(r)}) \text{ iff } (\alpha(t_1), \dots, \alpha(t_{\iota(r)})) \in r^{\mathcal{S}}; \\
(\mathcal{S}, \alpha) \models (Q_1 \vee Q_2) \text{ iff } (\mathcal{S}, \alpha) \models Q_1 \text{ or } (\mathcal{S}, \alpha) \models Q_2; \\
(\mathcal{S}, \alpha) \models (Q_1 \wedge Q_2) \text{ iff } (\mathcal{S}, \alpha) \models Q_1 \text{ and } (\mathcal{S}, \alpha) \models Q_2;
\end{array}
\left|
\begin{array}{l}
(\mathcal{S}, \alpha) \models (x \approx t) \text{ iff } \alpha(x) = \alpha(t); \\
(\mathcal{S}, \alpha) \models (\neg Q) \text{ iff } (\mathcal{S}, \alpha) \not\models Q; \\
(\mathcal{S}, \alpha) \models (\exists x. Q) \text{ iff } (\mathcal{S}, \alpha[x \mapsto d]) \models Q, \\
\text{for some } d \in \mathcal{D}.
\end{array}
\right.$$

Figure 2: The semantics of RC.

$$\begin{array}{l}
x \approx x \equiv \top, \quad \neg \perp \equiv \top, \quad \neg \top \equiv \perp, \quad \exists x. \perp \equiv \perp, \quad \exists x. \top \equiv \top, \\
Q \wedge \perp \equiv \perp, \quad \perp \wedge Q \equiv \perp, \quad Q \wedge \top \equiv Q, \quad \top \wedge Q \equiv Q, \\
Q \vee \perp \equiv Q, \quad \perp \vee Q \equiv Q, \quad Q \vee \top \equiv \top, \quad \top \vee Q \equiv \top.
\end{array}$$

Figure 3: Constant propagation rules.

notation to sequences  $\vec{x}$  and  $\vec{d}$  of pairwise distinct variables and arbitrary domain elements of the same length. The semantics of RC queries for a structure  $\mathcal{S}$  and an assignment  $\alpha$  is defined in Figure 2. We write  $\alpha \models Q$  for  $(\mathcal{S}, \alpha) \models Q$  if the structure  $\mathcal{S}$  is fixed in the given context. For a fixed  $\mathcal{S}$ , only the assignments to  $Q$ 's free variables influence  $\alpha \models Q$ , i.e.,  $\alpha \models Q$  is equivalent to  $\alpha' \models Q$ , for every variable assignment  $\alpha'$  that agrees with  $\alpha$  on  $\text{fv}(Q)$ . For closed queries  $Q$ , we write  $\models Q$  and say that  $Q$  holds, since closed queries either hold for all variable assignments or for none of them. We call a finite sequence  $\vec{d}$  of domain elements  $d_1, \dots, d_k \in \mathcal{D}$  a *tuple*. Given a query  $Q$  and a structure  $\mathcal{S}$ , we denote the set of satisfying tuples for  $Q$  by

$$[[Q]]^{\mathcal{S}} = \{\vec{d} \in \mathcal{D}^{|\text{fv}(Q)|} \mid \text{there exists an assignment } \alpha \text{ such that } (\mathcal{S}, \alpha[\text{fv}(Q) \mapsto \vec{d}]) \models Q\}.$$

We omit  $\mathcal{S}$  from  $[[Q]]^{\mathcal{S}}$  if  $\mathcal{S}$  is fixed. We call the values from  $[[Q]]$  assigned to  $x \in \text{fv}(Q)$  *column*  $x$ .

The *active domain*  $\text{adom}^{\mathcal{S}}(Q)$  of a query  $Q$  and a structure  $\mathcal{S}$  is a subset of the domain  $\mathcal{D}$  containing the interpretations  $c^{\mathcal{S}}$  of all constant symbols that occur in  $Q$  and the values in the relations  $r^{\mathcal{S}}$  interpreting all predicate symbols that occur in  $Q$ . Since  $\mathcal{C}$  and  $\mathcal{R}$  are finite and all  $r^{\mathcal{S}}$  are finite relations of a finite arity  $\iota(r)$ , the active domain  $\text{adom}^{\mathcal{S}}(Q)$  is also a finite set. We omit  $\mathcal{S}$  from  $\text{adom}^{\mathcal{S}}(Q)$  if  $\mathcal{S}$  is fixed in the given context.

Queries  $Q_1$  and  $Q_2$  over the same signature are *equivalent*, written  $Q_1 \equiv Q_2$ , if  $(\mathcal{S}, \alpha) \models Q_1 \iff (\mathcal{S}, \alpha) \models Q_2$ , for every  $\mathcal{S}$  and  $\alpha$ . Queries  $Q_1$  and  $Q_2$  over the same signature are *inf-equivalent*, written  $Q_1 \overset{\infty}{\equiv} Q_2$ , if  $(\mathcal{S}, \alpha) \models Q_1 \iff (\mathcal{S}, \alpha) \models Q_2$ , for every  $\mathcal{S}$  with an *infinite* domain  $\mathcal{D}$  and every  $\alpha$ . Clearly, equivalent queries are also inf-equivalent.

A query  $Q$  is *domain-independent* if  $[[Q]]^{\mathcal{S}_1} = [[Q]]^{\mathcal{S}_2}$  holds for every two structures  $\mathcal{S}_1$  and  $\mathcal{S}_2$  that agree on the interpretations of constants ( $c^{\mathcal{S}_1} = c^{\mathcal{S}_2}$ ) and predicates ( $r^{\mathcal{S}_1} = r^{\mathcal{S}_2}$ ), while their domains  $\mathcal{D}_1$  and  $\mathcal{D}_2$  may differ. Agreement on the interpretations implies  $\text{adom}^{\mathcal{S}_1}(Q) = \text{adom}^{\mathcal{S}_2}(Q) \subseteq \mathcal{D}_1 \cap \mathcal{D}_2$ . It is undecidable whether an RC query is domain-independent [Pao69, Var81].

We denote by  $\text{cp}(Q)$  the query obtained from a query  $Q$  by exhaustively applying the rules in Figure 3. Note that  $\text{cp}(Q)$  is either of the form  $\perp$  or  $\top$  or contains no  $\perp$  or  $\top$  subqueries.

**Definition 3.1.** *The substitution of the form  $Q[x \mapsto y]$  is the query  $\text{cp}(Q')$ , where  $Q'$  is obtained from a query  $Q$  by replacing all occurrences of the free variable  $x$  by the variable  $y$ , potentially also renaming bound variables to avoid capture.*

**Definition 3.2.** *The substitution of the form  $Q[x/\perp]$  is the query  $\text{cp}(Q')$ , where  $Q'$  is obtained from a query  $Q$  by replacing with  $\perp$  every atomic predicate or equality containing the free variable  $x$ , except for  $(x \approx x)$  which is replaced by  $\top$ .*

We lift the substitution notation to sets of queries in the standard way.

$\text{gen}(x, \perp, \emptyset);$	$\text{cov}(x, x \approx x, \emptyset);$	
$\text{gen}(x, Q, \{Q\})$ if $\text{ap}(Q)$ and $x \in \text{fv}(Q);$	$\text{cov}(x, Q, \emptyset)$	if $x \notin \text{fv}(Q);$
$\text{gen}(x, \neg Q, \mathcal{G})$ if $\text{gen}(x, Q, \mathcal{G});$	$\text{cov}(x, x \approx y, \{x \approx y\})$	if $x \neq y;$
$\text{gen}(x, \neg(Q_1 \vee Q_2), \mathcal{G})$	$\text{cov}(x, y \approx x, \{x \approx y\})$	if $x \neq y;$
if $\text{gen}(x, (\neg Q_1) \wedge (\neg Q_2), \mathcal{G});$	$\text{cov}(x, Q, \{Q\})$	if $\text{ap}(Q)$ and $x \in \text{fv}(Q);$
$\text{gen}(x, \neg(Q_1 \wedge Q_2), \mathcal{G})$	$\text{cov}(x, \neg Q, \mathcal{G})$	if $\text{cov}(x, Q, \mathcal{G});$
if $\text{gen}(x, (\neg Q_1) \vee (\neg Q_2), \mathcal{G});$	$\text{cov}(x, Q_1 \vee Q_2, \mathcal{G}_1 \cup \mathcal{G}_2)$	if $\text{cov}(x, Q_1, \mathcal{G}_1)$ and $\text{cov}(x, Q_2, \mathcal{G}_2);$
$\text{gen}(x, Q_1 \vee Q_2, \mathcal{G}_1 \cup \mathcal{G}_2)$	$\text{cov}(x, Q_1 \vee Q_2, \mathcal{G})$	if $\text{cov}(x, Q_1, \mathcal{G})$ and $Q_1[x/\perp] = \top;$
if $\text{gen}(x, Q_1, \mathcal{G}_1)$ and $\text{gen}(x, Q_2, \mathcal{G}_2);$	$\text{cov}(x, Q_1 \vee Q_2, \mathcal{G})$	if $\text{cov}(x, Q_2, \mathcal{G})$ and $Q_2[x/\perp] = \top;$
$\text{gen}(x, Q_1 \wedge Q_2, \mathcal{G})$	$\text{cov}(x, Q_1 \wedge Q_2, \mathcal{G}_1 \cup \mathcal{G}_2)$	if $\text{cov}(x, Q_1, \mathcal{G}_1)$ and $\text{cov}(x, Q_2, \mathcal{G}_2);$
if $\text{gen}(x, Q_1, \mathcal{G})$ or $\text{gen}(x, Q_2, \mathcal{G});$	$\text{cov}(x, Q_1 \wedge Q_2, \mathcal{G})$	if $\text{cov}(x, Q_1, \mathcal{G})$ and $Q_1[x/\perp] = \perp;$
$\text{gen}(x, Q \wedge x \approx y, \mathcal{G}[y \mapsto x])$	$\text{cov}(x, Q_1 \wedge Q_2, \mathcal{G})$	if $\text{cov}(x, Q_2, \mathcal{G})$ and $Q_2[x/\perp] = \perp;$
if $\text{gen}(y, Q, \mathcal{G});$	$\text{cov}(x, \exists y. Q_y, \exists y. \mathcal{G})$	
$\text{gen}(x, Q \wedge y \approx x, \mathcal{G}[y \mapsto x])$		if $x \neq y$ and $\text{cov}(x, Q_y, \mathcal{G})$ and $(x \approx y) \notin \mathcal{G};$
if $\text{gen}(y, Q, \mathcal{G});$	$\text{cov}(x, \exists y. Q_y, \exists y. (\mathcal{G} \setminus \{x \approx y\}) \cup \mathcal{G}_y[y \mapsto x])$	
$\text{gen}(x, \exists y. Q_y, \exists y. \mathcal{G})$		if $x \neq y$ and $\text{cov}(x, Q_y, \mathcal{G})$ and $\text{gen}(y, Q_y, \mathcal{G}_y).$
if $x \neq y$ and $\text{gen}(x, Q_y, \mathcal{G}).$		

Figure 4: The *generated* relation.Figure 5: The *covered* relation.

The function  $\text{flat}^\oplus(Q)$ , where  $\oplus \in \{\vee, \wedge\}$ , computes a set of queries by “flattening” the operator  $\oplus$ :  $\text{flat}^\oplus(Q) := \text{flat}^\oplus(Q_1) \cup \text{flat}^\oplus(Q_2)$  if  $Q = Q_1 \oplus Q_2$  and  $\text{flat}^\oplus(Q) := \{Q\}$  otherwise.

**3.2. Safe-Range Queries.** The class of *safe-range* queries [AHV95] is a decidable subset of domain-independent RC queries. Its definition is based on the notion of the range-restricted variables of a query. A variable is called *range restricted* if “its possible values all lie within the active domain of the query” [AHV95]. Intuitively, atomic predicates restrict the possible values of a variable that occurs in them as a term. An equality  $x \approx y$  can extend the set of range-restricted variables in a conjunction  $Q \wedge x \approx y$ : If  $x$  or  $y$  is range restricted in  $Q$ , then both  $x$  and  $y$  are range restricted in  $Q \wedge x \approx y$ .

We formalize range-restricted variables using the *generated* relation  $\text{gen}(x, Q, \mathcal{G})$ , defined in Figure 4. Specifically,  $\text{gen}(x, Q, \mathcal{G})$  holds if  $x$  is a range-restricted variable in  $Q$  and every satisfying assignment for  $Q$  satisfies some quantified predicate, referred to as *generator*, from  $\mathcal{G}$ . A similar definition by Van Gelder and Topor [GT91, Figure 5] uses a set of atomic (not quantified) predicates  $\mathcal{A}$  as generators and defines the rule  $\text{gen}_{\text{vgt}}(x, \exists y. Q_y, \mathcal{A})$  if  $x \neq y$  and  $\text{gen}_{\text{vgt}}(x, Q_y, \mathcal{A})$  (Appendix A, Figure 17). In contrast, we modify the rule’s conclusion to existentially quantify the variable  $y$  in all queries in  $\mathcal{G}$  where  $y$  is free:  $\text{gen}(x, \exists y. Q_y, \exists y. \mathcal{G})$ . Hence,  $\text{gen}(x, Q, \mathcal{G})$  implies  $\text{fv}(\mathcal{G}) \subseteq \text{fv}(Q)$ . We now formalize these relationships.

**Lemma 3.3.**  $\clubsuit$  Let  $Q$  be a query,  $x \in \text{fv}(Q)$ , and  $\mathcal{G}$  be a set of quantified predicates such that  $\text{gen}(x, Q, \mathcal{G})$ . Then (i) for every  $Q_{qp} \in \mathcal{G}$ , we have  $x \in \text{fv}(Q_{qp})$  and  $\text{fv}(Q_{qp}) \subseteq \text{fv}(Q)$ , (ii) for every  $\alpha$  such that  $\alpha \models Q$ , there exists a  $Q_{qp} \in \mathcal{G}$  such that  $\alpha \models Q_{qp}$ , and (iii)  $Q[x/\perp] = \perp$ .

**Definition 3.4.** Let  $\text{gen}(x, Q)$  hold iff  $\text{gen}(x, Q, \mathcal{G})$  holds for some  $\mathcal{G}$ . Let  $\text{nongens}(Q) := \{x \in \text{fv}(Q) \mid \text{gen}(x, Q) \text{ does not hold}\}$  be the set of free variables in a query  $Q$  that are not range restricted. A query  $Q$  has range-restricted free variables if every free variable of  $Q$  is range restricted, i.e.,  $\text{nongens}(Q) = \emptyset$ . A query  $Q$  has range-restricted bound variables if the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$  is range restricted, i.e.,  $\text{gen}(y, Q_y)$  holds. A query is safe range if it has range-restricted free and range-restricted bound variables.

$\text{ranf}(\perp)$	;	$\text{ranf}(\top)$	;
$\text{ranf}(Q)$		if $\text{ap}(Q)$	;
$\text{ranf}(\neg Q)$		if $\text{ranf}(Q)$ and $\text{fv}(Q) = \emptyset$	;
$\text{ranf}(Q_1 \vee Q_2)$		if $\text{ranf}(Q_1)$ and $\text{ranf}(Q_2)$ and $\text{fv}(Q_1) = \text{fv}(Q_2)$	;
$\text{ranf}(Q_1 \wedge Q_2)$		if $\text{ranf}(Q_1)$ and $\text{ranf}(Q_2)$	;
$\text{ranf}(Q_1 \wedge \neg Q_2)$		if $\text{ranf}(Q_1)$ and $\text{ranf}(Q_2)$ and $\text{fv}(Q_2) \subseteq \text{fv}(Q_1)$	;
$\text{ranf}(Q \wedge (x \approx y))$		if $\text{ranf}(Q)$ and $\{x, y\} \cap \text{fv}(Q) \neq \emptyset$	;
$\text{ranf}(Q \wedge \neg(x \approx y))$		if $\text{ranf}(Q)$ and $\{x, y\} \subseteq \text{fv}(Q)$	;
$\text{ranf}(\exists x. Q_x)$		if $\text{ranf}(Q_x)$ and $x \in \text{fv}(Q_x)$	.

Figure 6: Characterization of RANF queries.

**3.3. Safe-Range Normal Form.** A query  $Q$  is in safe-range normal form (SRNF) if the query  $Q'$  in every subquery  $\neg Q'$  of  $Q$  is an atomic predicate, equality, or an existentially quantified query [AHV95]. In Section 6.1 we define function  $\text{srnf}(Q)$  that returns a SRNF query equivalent to a query  $Q$ . Intuitively, the function  $\text{srnf}(Q)$  proceeds by pushing negations downwards [AHV95, Section 5.4], distributing existential quantifiers over disjunction [GT91, Rule (T9)], and dropping bound variables that never occur [GT91, Definition 9.2]. We include the last two rules to optimize the time complexity of evaluating the resulting query.

If a query  $Q$  is safe range, then  $\text{srnf}(Q)$  is also safe range.

**3.4. Relational Algebra Normal Form.** Relation algebra normal form (RANF) is a class of safe-range queries that can be easily mapped to RA [AHV95] and evaluated using the RA operations for projection, column duplication, selection, set union, binary join, and anti-join.

Figure 6 defines the predicate  $\text{ranf}(\cdot)$  characterizing RANF queries. The translation of safe-range queries (Section 3.2) to equivalent RANF queries proceeds via SRNF (Section 3.3). A safe-range query in SRNF can be translated to an equivalent RANF query by subquery rewriting using the following rules [AHV95, Algorithm 5.4.7]:

$$\begin{aligned} Q \wedge (Q_1 \vee Q_2) &\equiv (Q \wedge Q_1) \vee (Q \wedge Q_2), & (R1) \\ Q \wedge (\exists x. Q_x) &\equiv (\exists x. Q \wedge Q_x), & (R2) \\ Q \wedge \neg Q' &\equiv Q \wedge \neg(Q \wedge Q'). & (R3) \end{aligned}$$

Subquery rewriting is a nondeterministic process (as the rewrite rules can be applied in an arbitrary order) that impacts the performance of evaluating the resulting RANF query. We translate a safe-range query in SRNF to an equivalent RANF query by a recursive function  $\text{sr2ranf}(\cdot)$  inspired by the rules (R1)–(R3) and fully specified in Figure 12 in Section 6.2.

**3.5. Query Cost.** To assess the time complexity of evaluating a RANF query  $Q$ , we define the *cost* of  $Q$  over a structure  $\mathcal{S}$ , denoted  $\text{cost}^{\mathcal{S}}(Q)$ , to be the sum of intermediate result sizes over all RANF subqueries of  $Q$ . Formally,  $\text{cost}^{\mathcal{S}}(Q) := \sum_{Q' \sqsubseteq Q, \text{ranf}(Q')} \left| \llbracket Q' \rrbracket^{\mathcal{S}} \right| \cdot |\text{fv}(Q')|$ . This corresponds to evaluating  $Q$  following its RANF structure (Section 3.4, Figure 6) using the RA operations. The complexity of these operations is linear in the combined input and output size (ignoring logarithmic factors due to set operations). The output size (the number of tuples times the number of variables) is counted in  $\left| \llbracket Q' \rrbracket^{\mathcal{S}} \right| \cdot |\text{fv}(Q')|$  and the input size is counted as the output size for the input subqueries. Repeated subqueries are only considered once, which does not affect the asymptotics of query cost. In practice, the evaluation results for common subqueries can be reused.



## 4. QUERY TRANSLATION

Our approach to evaluating an arbitrary RC query  $Q$  over a fixed structure  $\mathcal{S}$  with an infinite domain  $\mathcal{D}$  proceeds by translating  $Q$  into a pair of safe-range queries  $(Q_{fin}, Q_{inf})$  such that

- (FV)  $\text{fv}(Q_{fin}) = \text{fv}(Q)$  unless  $Q_{fin}$  is syntactically equal to  $\perp$ ;  $\text{fv}(Q_{inf}) = \emptyset$ ;
- (EVAL)  $\llbracket Q \rrbracket$  is an infinite set if  $Q_{inf}$  holds; otherwise  $\llbracket Q \rrbracket = \llbracket Q_{fin} \rrbracket$  is a finite set.

Since the queries  $Q_{fin}$  and  $Q_{inf}$  are safe range, they are domain-independent and thus  $\llbracket Q_{fin} \rrbracket$  is a finite set. In particular,  $\llbracket Q \rrbracket$  is a finite set if  $Q_{inf}$  does not hold. Our translation generalizes Hull and Su's case distinction that restricts bound variables [HS94] to restrict all variables. Moreover, we use Van Gelder and Topor's idea to replace the active domain by a smaller set (generator) specific to each variable [GT91] while further improving the generators. Unless explicitly noted, in the rest of the article we assume a fixed structure  $\mathcal{S}$ .

**4.1. Restricting One Variable.** Let  $x$  be a free variable in a query  $\tilde{Q}$  with range-restricted bound variables. This assumption on  $\tilde{Q}$  will be established by translating an arbitrary query  $Q$  bottom-up (Section 4.2). In this section, we develop a translation of  $\tilde{Q}$  into an equivalent query  $\tilde{Q}'$  that satisfies the following:

- $\tilde{Q}'$  has range-restricted bound variables;
- $\tilde{Q}'$  is a disjunction;  $x$  is range restricted in the first disjunct; the remaining disjuncts are all binary conjunctions of a query not containing  $x$  with a query of a special form containing  $x$ . The special form, central to our translation, is either an equality  $x \approx y$  or a query satisfied by infinitely many values of  $x$  for all values of the remaining free variables.

We now restate Hull and Su's [HS94] and Van Gelder and Topor's [GT91] approaches using our notation in order to clarify how we generalize both approaches. In particular, Hull and Su's approach is already stated in a generalized way that restricts a *free* variable.

**Hull and Su.** Let  $\tilde{Q}$  be a query with range-restricted bound variables and  $x \in \text{fv}(\tilde{Q})$ . Then

$$\tilde{Q} \equiv \left( \tilde{Q} \wedge \text{AD}(x, \tilde{Q}) \right) \vee \left( \bigvee_{y \in \text{fv}(\tilde{Q}) \setminus \{x\}} (\tilde{Q}[x \mapsto y] \wedge x \approx y) \right) \vee \left( \tilde{Q}[x/\perp] \wedge \neg(\text{AD}(x, \tilde{Q}) \vee \bigvee_{y \in \text{fv}(\tilde{Q}) \setminus \{x\}} x \approx y) \right).$$

Here  $\text{AD}(x, \tilde{Q})$  stands for an RC query with a single free variable  $x$  that is satisfied by an assignment  $\alpha$  if and only if  $\alpha(x) \in \text{adom}(\tilde{Q})$ . Hull and Su's translation distinguishes the following three cases for a fixed assignment  $\alpha$  (each corresponding to a top-level disjunct above):

- if  $\text{AD}(x, \tilde{Q})$  holds (and hence  $\alpha(x) \in \text{adom}(\tilde{Q})$ ), then we do not alter the query  $\tilde{Q}$ ;
- if  $x \approx y$  holds for some free variable  $y \in \text{fv}(\tilde{Q}) \setminus \{x\}$ , then  $x$  can be replaced by  $y$  in  $\tilde{Q}$ ;
- otherwise,  $\tilde{Q}$  is equivalent to  $\tilde{Q}[x/\perp]$ . Specifically, all atomic predicates having  $x$  free can be replaced by  $\perp$  (as  $\alpha(x) \notin \text{adom}(\tilde{Q})$ ), all equalities  $x \approx y$  and  $y \approx x$  for  $y \in \text{fv}(\tilde{Q}) \setminus \{x\}$  can be replaced by  $\perp$  (as  $\alpha(x) \neq \alpha(y)$ ), and all equalities  $x \approx z$  for a bound variable  $z$  can be replaced by  $\perp$  (as  $\alpha(x) \notin \text{adom}(\tilde{Q})$  and  $z$  is range restricted in its subquery  $\exists z. Q_z$ , by assumption). In the last case,  $\text{gen}(z, Q_z)$  holds and thus, for all  $\alpha'$  extending  $\alpha$ , we have  $\alpha' \models \exists z. Q_z$  if and only if there exists a  $d \in \text{adom}(Q_z) \subseteq \text{adom}(\tilde{Q})$  such that  $\alpha'[z \mapsto d] \models Q_z$ .

**Example 4.1.** Consider the query  $Q := \text{B}(y) \vee x \approx y$ . Then  $\text{AD}(x, Q) = \text{B}(x)$  and following Hull and Su we obtain that  $Q$  is equivalent to the disjunction of the following three queries:

- $(\text{B}(y) \vee x \approx y) \wedge \text{B}(x)$ , which is equivalent to  $\text{B}(x) \wedge \text{B}(y)$ ;

- $(\mathbf{B}(y) \vee x \approx y)[x \mapsto y] \wedge x \approx y$ , which is syntactically equal to  $x \approx y$  due to constant propagation that is part of the substitution operator;
- $(\mathbf{B}(y) \vee x \approx y)[x/\perp] \wedge \neg(\mathbf{B}(x) \vee x \approx y)$ , which is equivalent to  $\mathbf{B}(y) \wedge \neg\mathbf{B}(x) \wedge \neg x \approx y$ .

Note that in this example, each disjunct covers a different subset of  $Q$ 's satisfying assignments and all three disjuncts are necessary to cover all of  $Q$ 's satisfying assignments.

**Van Gelder and Topor.** Let  $\tilde{Q}$  be an *evaluable* query with range-restricted bound variables,  $x \in \text{fv}(\tilde{Q})$ . Then there exists a set  $\mathcal{A}$  of atomic predicates such that

$$\tilde{Q} \equiv \left( \tilde{Q} \wedge \bigvee_{Q_{ap} \in \mathcal{A}} \exists \vec{\text{fv}}(Q_{ap}) \setminus \{x\}. Q_{ap} \right) \vee \left( \tilde{Q}[x/\perp] \right).$$

Note that  $\exists \vec{\text{fv}}(Q) \setminus \{x\}. Q$  is the query in which all free variables of  $Q$  except for  $x$  are existentially quantified. Van Gelder and Topor restrict their attention to evaluable queries, which do not contain equalities between variables. (They only discuss an incomplete approach to supporting such equalities [GT91, Appendix A].) Thus, their translation lacks the corresponding disjuncts that Hull and Su have.

To avoid enumerating the entire active domain  $\text{adom}(\tilde{Q})$ , Van Gelder and Topor replace the query  $\text{AD}(x, \tilde{Q})$  used by Hull and Su by the query  $\bigvee_{Q_{ap} \in \mathcal{A}} \exists \vec{\text{fv}}(Q_{ap}) \setminus \{x\}. Q_{ap}$  constructed from the atomic predicates from  $\mathcal{A}$ . Because their translation must yield an equivalent query (for every finite or infinite domain),  $\mathcal{A}$  and  $\tilde{Q}$  must satisfy, for all  $\alpha$ ,

$$\begin{aligned} \alpha \models \neg \bigvee_{Q_{ap} \in \mathcal{A}} \exists \vec{\text{fv}}(Q_{ap}) \setminus \{x\}. Q_{ap} &\implies (\alpha \models \tilde{Q} \iff \alpha \models \tilde{Q}[x/\perp]) & (\text{VGT}_1) & \quad \text{and} \\ \alpha \models \tilde{Q}[x/\perp] &\implies \alpha \models \forall x. \tilde{Q} & (\text{VGT}_2). \end{aligned}$$

Note that (VGT<sub>2</sub>) does not hold for the query  $\tilde{Q} := \neg\mathbf{B}(x)$  and thus Van Gelder and Topor only consider a proper subset of all RC queries, called *evaluable*. For evaluable queries, Van Gelder and Topor use the *constrained* relation  $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$ , defined in Appendix A, Figure 17, to construct a set of atomic predicates  $\mathcal{A}$  that satisfies (VGT<sub>1</sub>).

**Our Translation.** Let  $\tilde{Q}$  be a query with range-restricted bound variables,  $x \in \text{fv}(\tilde{Q})$ . Then there exists a set  $\mathcal{A}$  of atomic predicates and a set of equalities  $\mathcal{E}$  such that

$$\begin{aligned} \tilde{Q} \equiv & \left( \tilde{Q} \wedge \bigvee_{Q_{ap} \in \mathcal{A}} \exists \vec{\text{fv}}(Q_{ap}) \setminus \text{fv}(\tilde{Q}). Q_{ap} \right) \vee \left( \bigvee_{x \approx y \in \mathcal{E}} (\tilde{Q}[x \mapsto y] \wedge x \approx y) \right) \vee \\ & \left( \tilde{Q}[x/\perp] \wedge \neg \left( \bigvee_{Q_{ap} \in \mathcal{A}} \exists \vec{\text{fv}}(Q_{ap}) \setminus \text{fv}(\tilde{Q}). Q_{ap} \vee \bigvee_{x \approx y \in \mathcal{E}} x \approx y \right) \right). \end{aligned}$$

In contrast to Van Gelder and Topor, we only require that  $\mathcal{A}$  satisfies (VGT<sub>1</sub>) in our translation, which also allows us to translate non-evaluable queries, such as  $\tilde{Q} := \neg\mathbf{B}(x)$  above. Note that we also existentially quantify only these variables that are not free in  $\tilde{Q}$ , whereas Van Gelder and Topor quantify all variables except  $x$ . For our introductory example  $Q^{\text{sup}}$ , this modification allows our translation to use the quantified predicate  $\exists p. \mathbf{S}(p, u, s)$  to restrict both  $u$  and  $p$  simultaneously. In contrast, Van Gelder and Topor's approach restricts them separately using  $\exists p, u. \mathbf{S}(p, u, s)$  and  $\exists p, s. \mathbf{S}(p, u, s)$ , so that the Cartesian product of these quantified predicates may need to be computed in their translated queries.

In contrast to Hull and Su, we do not consider the equalities of  $x$  with all other free variables in  $\tilde{Q}$ , but only such equalities  $\mathcal{E}$  that occur in  $\tilde{Q}$ . We jointly compute the sets  $\mathcal{A}$  and  $\mathcal{E}$  using the *covered* relation  $\text{cov}(x, Q, \mathcal{G})$  (in contrast to  $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$  relation). Figure 5 shows the definition of this relation. The set  $\mathcal{G}$  computed by the covered relation contains atomic predicates that satisfy (VGT<sub>1</sub>) and are already quantified as described above. The set

also contains the relevant equalities that can be used in our translation. For every variable  $x$  and query  $\tilde{Q}$  with range-restricted bound variables, there exists at least one set of quantified predicates and equalities  $\mathcal{G}$  such that  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  and  $(\text{VGT}_1)$  holds for the set of atomic predicate subqueries in  $\mathcal{G}$  (i.e., for  $\{Q' \mid \text{ap}(Q') \wedge \exists Q \in \mathcal{G}. Q' \sqsubseteq Q\}$ ). As the *cover* set  $\mathcal{G}$  in  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  may contain both quantified predicates and equalities between two variables, we define a function  $\text{qps}(\mathcal{G})$  that collects all *generators*, i.e., quantified predicates and a function  $\text{eqs}(x, \mathcal{G})$  that collects all *variables*  $y$  distinct from  $x$  occurring in equalities of the form  $x \approx y$ . We use  $\text{qps}^\vee(\mathcal{G})$  to denote the query  $\bigvee_{Q_{qp} \in \text{qps}(\mathcal{G})} Q_{qp}$ . We state the soundness and completeness of the relation  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  in the next lemma, which follows by induction on the derivation of  $\text{cov}(x, \tilde{Q}, \mathcal{G})$ .

**Lemma 4.2.**  $\clubsuit$  *Let  $\tilde{Q}$  be a query with range-restricted bound variables,  $x \in \text{fv}(\tilde{Q})$ .*

**Completeness:** *Then there exists a set  $\mathcal{G}$  of quantified predicates and equalities such that  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  holds and,*

**Soundness:** *for any  $\mathcal{G}$  satisfying  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  and all  $\alpha$ ,*

$$\alpha \models \neg(\text{qps}^\vee(\mathcal{G}) \vee \bigvee_{y \in \text{eqs}(x, \mathcal{G})} x \approx y) \implies (\alpha \models \tilde{Q} \iff \alpha \models \tilde{Q}[x/\perp]).$$

Finally, to preserve the dependencies between the variable  $x$  and the remaining free variables of  $Q$  occurring in the quantified predicates from  $\text{qps}(\mathcal{G})$ , we do not project  $\text{qps}(\mathcal{G})$  on the single variable  $x$ , i.e., we restrict  $x$  by  $\text{qps}^\vee(\mathcal{G})$  instead of  $\exists \text{fv}(Q) \setminus \{x\}. \text{qps}(\mathcal{G})$  as by Van Gelder and Topor. From Lemma 4.2, we derive our optimized translation characterized by the following lemma.

**Lemma 4.3.**  $\clubsuit$  *Let  $\tilde{Q}$  be a query with range-restricted bound variables,  $x \in \text{fv}(\tilde{Q})$ , and  $\mathcal{G}$  be such that  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  holds. Then  $x \in \text{fv}(Q_{qp})$  and  $\text{fv}(Q_{qp}) \subseteq \text{fv}(\tilde{Q})$ , for every  $Q_{qp} \in \text{qps}(\mathcal{G})$ , and*

$$\begin{aligned} \tilde{Q} \equiv & \left( \tilde{Q} \wedge \text{qps}^\vee(\mathcal{G}) \right) \vee \left( \bigvee_{y \in \text{eqs}(x, \mathcal{G})} (\tilde{Q}[x \mapsto y] \wedge x \approx y) \right) \vee \\ & \left( \tilde{Q}[x/\perp] \wedge \neg(\text{qps}^\vee(\mathcal{G}) \vee \bigvee_{y \in \text{eqs}(x, \mathcal{G})} x \approx y) \right). \end{aligned} \quad (\star)$$

Note that  $x$  is only guaranteed to be range restricted in  $(\star)$ 's first disjunct. However, it only occurs in the remaining disjuncts in subqueries of a special form that are conjoined at the top-level to the disjuncts. These subqueries of a special form are equalities of the form  $x \approx y$  or negations of a disjunction of quantified predicates with a free occurrence of  $x$  and equalities of the form  $x \approx y$ . We will show how to handle such occurrences in Section 4.2 and Section 4.3. Moreover, the negation of the disjunction can be omitted if  $(\text{VGT}_2)$  holds.

**4.2. Restricting Bound Variables.** Let  $x$  be a free variable in a query  $\tilde{Q}$  with range-restricted bound variables. Suppose that the variable  $x$  is not range restricted, i.e.,  $\text{gen}(x, \tilde{Q})$  does not hold. To translate  $\exists x. \tilde{Q}$  into an inf-equivalent query with range-restricted bound variables ( $\exists x. \tilde{Q}$  does not have range-restricted bound variables precisely because  $x$  is not range restricted in  $\tilde{Q}$ ), we first apply  $(\star)$  to  $\tilde{Q}$  and distribute the existential quantifier binding  $x$  over disjunction. Next we observe that

$$\exists x. (\tilde{Q}[x \mapsto y] \wedge x \approx y) \equiv \tilde{Q}[x \mapsto y] \wedge \exists x. (x \approx y) \equiv \tilde{Q}[x \mapsto y],$$

where the first equivalence follows because  $x$  does not occur free in  $\tilde{Q}[x \mapsto y]$  and the second equivalence follows from the straightforward validity of  $\exists x. (x \approx y)$ . Moreover, we observe

**input:** An RC query  $Q$ .  
**output:** A query  $\tilde{Q}$  with range-restricted bound variables such that  $Q \stackrel{\infty}{\equiv} \tilde{Q}$ .

```

1 auxiliary function fixbound( $Q, x$ ) =
2 |  $\{Q_{fix} \in Q \mid x \in \text{nongens}(Q_{fix})\}$ 
3 function rb( $Q$ ) =
4 | switch  $Q$  do
5 |   case  $\neg Q'$  do return  $\neg \text{rb}(Q')$ ;
6 |   case  $Q'_1 \vee Q'_2$  do return
7 |      $\text{rb}(Q'_1) \vee \text{rb}(Q'_2)$ ;
8 |   case  $Q'_1 \wedge Q'_2$  do return
9 |      $\text{rb}(Q'_1) \wedge \text{rb}(Q'_2)$ ;
10 |  case  $\exists x. Q_x$  do
11 |     $Q := \text{flat}^\vee(\text{rb}(Q_x))$ ;
12 |    while  $\text{fixbound}(Q, x) \neq \emptyset$  do
13 |       $Q_{fix} \leftarrow \text{fixbound}(Q, x)$ ;
14 |       $\mathcal{G} \leftarrow \{G \mid \text{cov}(x, Q_{fix}, G)\}$ ;
15 |       $Q := (Q \setminus \{Q_{fix}\}) \cup$ 
16 |         $\{Q_{fix} \wedge \text{qps}^\vee(G)\} \cup$ 
17 |         $\bigcup_{y \in \text{eqs}(x, G)} \{Q_{fix}[x \mapsto y]\} \cup$ 
18 |         $\{Q_{fix}[x/\perp]\}$ ;
19 |    return  $\bigvee_{\tilde{Q} \in Q} \exists x. \tilde{Q}$ ;
20 |  otherwise do return  $Q$ ;

```

Figure 7: Restricting bound variables.

**input:** An RC query  $Q$ .  
**output:** Safe-range query pair  $(Q_{fin}, Q_{inf})$  for which (FV) and (EVAL) hold.

```

1 auxiliary function fixfree( $Q_{fin}$ ) =
2 |  $\{(Q_{fix}, E) \in Q_{fin} \mid \text{nongens}(Q_{fix}) \neq \emptyset\}$ 
3 auxiliary function inf( $Q_{fin}, Q$ ) =
4 |  $\{(Q_f, E) \in Q_{fin} \mid \text{disjointvars}(Q_f, E) \neq \emptyset \vee$ 
5 |    $\text{fv}(Q_f) \cup \text{fv}(E) \neq \text{fv}(Q)\}$ 
6 function split( $Q$ ) =
7 |  $Q_{fin} := \{(rb(Q), \emptyset)\}$ ;  $Q_{inf} := \emptyset$ ;
8 | while  $\text{fixfree}(Q_{fin}) \neq \emptyset$  do
9 |    $(Q_{fix}, E) \leftarrow \text{fixfree}(Q_{fin})$ ;
10 |    $x \leftarrow \text{nongens}(Q_{fix})$ ;
11 |    $\mathcal{G} \leftarrow \{G \mid \text{cov}(x, Q_{fix}, G)\}$ ;
12 |    $Q_{fin} := (Q_{fin} \setminus \{(Q_{fix}, E)\}) \cup$ 
13 |      $\{(Q_{fix} \wedge \text{qps}^\vee(G), E)\} \cup$ 
14 |      $\bigcup_{y \in \text{eqs}(x, G)} \{(Q_{fix}[x \mapsto y], E \cup \{(x, y)\})\}$ ;
15 |    $Q_{inf} := Q_{inf} \cup \{Q_{fix}[x/\perp]\}$ ;
16 | while  $\text{inf}(Q_{fin}, Q) \neq \emptyset$  do
17 |    $(Q_f, E) \leftarrow \text{inf}(Q_{fin}, Q)$ ;
18 |    $Q_{fin} := Q_{fin} \setminus \{(Q_f, E)\}$ ;
19 |    $Q_{inf} := Q_{inf} \cup \{Q_f \wedge (\bigwedge_{Q \in \approx(E)} Q)\}$ ;
20 | return  $(\bigvee_{(Q_f, E) \in Q_{fin}} (\bigwedge^{\approx}(Q_f, E)),$ 
21 |    $\text{rb}(\bigvee_{Q_i \in Q_{inf}} \exists \vec{fv}(Q_i). Q_i))$ ;

```

Figure 8: Restricting free variables.

the following inf-equivalence (recall: an equivalence that holds for infinite domains only):

$$\exists x. (\tilde{Q}[x/\perp] \wedge \neg(\text{qps}^\vee(\mathcal{G}) \vee \bigvee_{y \in \text{eqs}(x, \mathcal{G})} x \approx y)) \stackrel{\infty}{\equiv} \tilde{Q}[x/\perp]$$

because  $x$  is not free in  $\tilde{Q}[x/\perp]$  and there exists a value  $d$  for  $x$  in the infinite domain  $\mathcal{D}$  such that  $x \neq y$  holds for all finitely many  $y \in \text{eqs}(x, \mathcal{G})$  and  $d$  is not among the finitely many values interpreting the quantified predicates in  $\text{qps}(\mathcal{G})$ . Altogether, we obtain the following lemma.

**Lemma 4.4.**  $\clubsuit$  *Let  $\tilde{Q}$  be a query with range-restricted bound variables,  $x \in \text{fv}(\tilde{Q})$ , and  $\mathcal{G}$  be a set of quantified predicates and equalities such that  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  holds. Then*

$$\exists x. \tilde{Q} \stackrel{\infty}{\equiv} \left( \exists x. \tilde{Q} \wedge \text{qps}^\vee(\mathcal{G}) \right) \vee \left( \bigvee_{y \in \text{eqs}(x, \mathcal{G})} (\tilde{Q}[x \mapsto y]) \right) \vee \left( \tilde{Q}[x/\perp] \right). \quad (\star\exists)$$

Our approach for restricting all bound variables recursively applies Lemma 4.4. Because the set  $\mathcal{G}$  such that  $\text{cov}(x, Q, \mathcal{G})$  holds is not necessarily unique, we introduce the following (general) notation. We denote the non-deterministic choice of an object  $X$  from a non-empty set  $\mathcal{X}$  as  $X \leftarrow \mathcal{X}$ . We define the recursive function  $\text{rb}(Q)$  in Figure 7, where  $\text{rb}$  stands for range-restrict bound variables. The function converts an arbitrary RC query  $Q$  into an inf-equivalent query with range-restricted bound variables. We proceed by describing the

case  $\exists x. Q_x$ . First,  $\text{rb}(Q_x)$  is recursively applied on Line 9 to establish the precondition of Lemma 4.4 that the translated query has range-restricted bound variables. Because existential quantification distributes over disjunction, we flatten disjunction in  $\text{rb}(Q_x)$  and process the individual disjuncts independently. We apply  $(\star\exists)$  to every disjunct  $Q_{fix}$  in which the variable  $x$  is not already range restricted. For every  $Q'_{fix}$  added to  $\mathcal{Q}$  after applying  $(\star\exists)$  to  $Q_{fix}$  the variable  $x$  is either range restricted or does not occur in  $Q'_{fix}$ , i.e.,  $x \notin \text{nongens}(Q'_{fix})$ . This entails the termination of the loop on Lines 10–13.

**Example 4.5.**  $\clubsuit$  Consider the query  $Q_{user}^{susp} := \text{B}(b) \wedge \exists s. \forall p. \text{P}(b, p) \longrightarrow \text{S}(p, u, s)$  from Section 1. Restricting its bound variables yields the query

$$\text{rb}(Q_{user}^{susp}) = \text{B}(b) \wedge \left( \left( \exists s. (\neg \exists p. \text{P}(b, p) \wedge \neg \text{S}(p, u, s)) \wedge (\exists p. \text{S}(p, u, s)) \right) \vee \left( \neg \exists p. \text{P}(b, p) \right) \right).$$

The bound variable  $p$  is already range restricted in  $Q_{user}^{susp}$  and thus only  $s$  must be restricted. Applying  $(\star)$  to restrict  $s$  in  $\neg \exists p. \text{P}(b, p) \wedge \neg \text{S}(p, u, s)$ , then existentially quantifying  $s$ , and distributing the existential quantifier over disjunction would yield the first disjunct in  $\text{rb}(Q_{user}^{susp})$  above and  $\exists s. (\neg \exists p. \text{P}(b, p)) \wedge \neg(\exists p. \text{S}(p, u, s))$  as the second disjunct. Because there exists some value in the infinite domain  $\mathcal{D}$  that does not belong to the finite interpretation of the atomic predicate  $\text{S}(p, u, s)$ , the query  $\exists s. \neg(\exists p. \text{S}(p, u, s))$  is a tautology over  $\mathcal{D}$ . Hence,  $\exists s. (\neg \exists p. \text{P}(b, p)) \wedge \neg(\exists p. \text{S}(p, u, s))$  is inf-equivalent to  $\neg \exists p. \text{P}(b, p)$ , i.e., the second disjunct in  $\text{rb}(Q_{user}^{susp})$ . This reasoning justifies that instead of  $(\star)$  our algorithm applies  $(\star\exists)$  to restrict  $s$  in  $\exists s. \neg \exists p. \text{P}(b, p) \wedge \neg \text{S}(p, u, s)$ .

**4.3. Restricting Free Variables.** Given an arbitrary query  $Q$ , we translate the inf-equivalent query  $\text{rb}(Q)$  with range-restricted bound variables into a pair of safe-range queries  $(Q_{fin}, Q_{inf})$  such that our translation's main properties (FV) and (EVAL) hold. Our translation is based on the following lemma.

**Lemma 4.6.**  $\clubsuit$  Let  $x$  be a free variable in a query  $\tilde{Q}$  with range-restricted bound variables and let  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  for a set of quantified predicates and equalities  $\mathcal{G}$ . If  $\tilde{Q}[x/\perp]$  is not satisfied by any tuple, then

$$\llbracket \tilde{Q} \rrbracket = \llbracket \left( \tilde{Q} \wedge \text{qps}^\vee(\mathcal{G}) \right) \vee \left( \bigvee_{y \in \text{eqs}(x, \mathcal{G})} (\tilde{Q}[x \mapsto y] \wedge x \approx y) \right) \rrbracket. \quad (\star)$$

If  $\tilde{Q}[x/\perp]$  is satisfied by some tuple, then  $\llbracket \tilde{Q} \rrbracket$  is an infinite set.

*Proof.* If  $\tilde{Q}[x/\perp]$  is not satisfied by any tuple, then  $(\star)$  follows from  $(\star)$ . If  $\tilde{Q}[x/\perp]$  is satisfied by some tuple, then the last disjunct in  $(\star)$  applied to  $\tilde{Q}$  is satisfied by infinitely many tuples obtained by assigning  $x$  some value from the infinite domain  $\mathcal{D}$  such that  $x \neq y$  holds for all finitely many  $y \in \text{eqs}(x, \mathcal{G})$  and  $x$  does not appear among the finitely many values interpreting the quantified predicates from  $\text{qps}(\mathcal{G})$ .  $\square$

We remark that  $\llbracket \tilde{Q} \rrbracket$  might be an infinite set of tuples even if  $\tilde{Q}[x/\perp]$  is never satisfied, for some  $x$ . This is because  $\tilde{Q}[y/\perp]$  might be satisfied by some tuple, for some  $y$ , in which case Lemma 4.6 (for  $y$ ) implies that  $\llbracket \tilde{Q} \rrbracket$  is an infinite set of tuples. Still,  $(\star)$  can be applied to  $\tilde{Q}$  for  $x$  resulting in a query satisfied by the same infinite set of tuples.

Our approach is implemented by the function  $\text{split}(Q)$  defined in Figure 8. In the following, we describe this function and justify its correctness, formalized by the input/output specification. In  $\text{split}(Q)$ , we represent the queries  $Q_{fin}$  and  $Q_{inf}$  using a set  $\mathcal{Q}_{fin}$  of pairs consisting of a query and a relation representing a set of equalities and a set  $\mathcal{Q}_{inf}$  of queries such that

$$Q_{fin} := \bigvee_{(Q_f, E) \in \mathcal{Q}_{fin}} (\bigwedge^{\approx} (Q_f, E)), \quad Q_{inf} := \bigvee_{Q_i \in \mathcal{Q}_{inf}} \exists \vec{fv}(Q_i). Q_i,$$

and, for every  $(Q_f, E) \in \mathcal{Q}_{fin}$ , the relation  $E$  represents a set of equalities between variables. Hereby,  $\bigwedge^{\approx}(Q_f, E)$  is a query that is equivalent to  $\bigwedge_{Q \in \{Q_f\} \cup \approx(E)} Q$  where  $\approx(E)$  abbreviates  $\{x \approx y \mid (x, y) \in E\}$ . However, the  $\bigwedge^{\approx}(Q_f, E)$  operator carefully assembles the conjunction to ensure that the resulting query is safe range (whenever possible). In particular, the operator must iteratively conjoin the equalities from  $\approx(E)$  to  $Q_f$  in a left-associative fashion and always pick next an equation for which one of the variables is free in  $Q_f$  or in the equalities conjoined so far, if such an equation exists. (If no such equation exists, the operator is free to conjoin the remaining equations in an arbitrary order.)

Our algorithm proceeds as follows. As long as there exists some  $(Q_{fix}, E) \in \mathcal{Q}_{fin}$  such that  $\text{nongens}(Q_{fix}) \neq \emptyset$ , we apply  $(\star)$  to  $Q_{fix}$  and add the query  $Q_{fix}[x/\perp]$  to  $\mathcal{Q}_{inf}$ . We remark that if we applied  $(\star)$  to the entire disjunct  $\bigwedge^{\approx}(Q_{fix}, E)$ , the loop on Lines 7–12 might not terminate. Note that, for every  $(Q'_{fix}, E')$  added to  $\mathcal{Q}_{fin}$  after applying  $(\star)$  to  $Q_{fix}$ ,  $\text{nongens}(Q'_{fix})$  is a proper subset of  $\text{nongens}(Q_{fix})$ . This entails the termination of the loop on Lines 7–12. Finally, if  $\llbracket Q_{fix} \rrbracket$  is an infinite set of tuples, then  $\llbracket \bigwedge^{\approx}(Q_{fix}, E) \rrbracket$  is an infinite set of tuples too. This is because the equalities in  $E$  merely duplicate columns of the query  $Q_{fix}$ . Hence, it indeed suffices to apply  $(\star)$  to  $Q_{fix}$  instead of  $\bigwedge^{\approx}(Q_{fix}, E)$ .

After the loop on Lines 7–12 in Figure 8 terminates, for every  $(Q_f, E) \in \mathcal{Q}_{fin}$ , the query  $Q_f$  is safe range and  $E$  is a conjunction of equalities such that  $\text{fv}(Q_f) \cup \text{fv}(E) \subseteq \text{fv}(Q)$ . However, the query  $\bigwedge^{\approx}(Q_f, E)$  does not have to be safe range, e.g., if  $Q_f := B(x)$  and  $E := \{(x, y), (u, v)\}$ . Given a relation  $E$ , let  $\text{classes}(E)$  be the set of equivalence classes of free variables  $\text{fv}(Q^{\approx})$  with respect to the (partial) equivalence closure of  $E$ , i.e., the smallest symmetric and transitive relation that contains  $E$ . For instance,  $\text{classes}(\{(x, y), (y, z), (u, v)\}) = \{\{x, y, z\}, \{u, v\}\}$ . Let  $\text{disjointvars}(Q_f, E) := \bigcup_{V \in \text{classes}(E), V \cap \text{fv}(Q_f) = \emptyset} V$  be the set of all variables in equivalence classes from  $\text{classes}(E)$  that are disjoint from  $Q_f$ 's free variables. Then,  $\bigwedge^{\approx}(Q_f, E)$  is safe range if and only if  $\text{disjointvars}(Q_f, E) = \emptyset$ .<sup>1</sup>

Now if  $\text{disjointvars}(Q_f, E) \neq \emptyset$  and the query  $\bigwedge^{\approx}(Q_f, E)$  is satisfied by some tuple, then  $\llbracket \bigwedge^{\approx}(Q_f, E) \rrbracket$  is an infinite set of tuples because all equivalence classes of variables in  $\text{disjointvars}(Q_f, E) \neq \emptyset$  can be assigned arbitrary values from the infinite domain  $\mathcal{D}$ . In our example with  $Q_f := B(x)$  and  $E := \{(x, y), (u, v)\}$ , we have  $\text{disjointvars}(Q_f, E) = \{u, v\} \neq \emptyset$ . Moreover, if  $\text{fv}(Q_f) \cup \text{fv}(E) \neq \text{fv}(Q)$  and  $\bigwedge^{\approx}(Q_f, E)$  is satisfied by some tuple, then this tuple can be extended to infinitely many tuples over  $\text{fv}(Q)$  by choosing arbitrary values from the infinite domain  $\mathcal{D}$  for the variables in the non-empty set  $\text{fv}(Q) \setminus (\text{fv}(Q_f) \cup \text{fv}(E))$ . Hence, for every  $(Q_f, E) \in \mathcal{Q}_{fin}$  with  $\text{disjointvars}(Q_f, E) \neq \emptyset$  or  $\text{fv}(Q_f) \cup \text{fv}(E) \neq \text{fv}(Q)$ , we remove  $(Q_f, E)$  from  $\mathcal{Q}_{fin}$  and add  $\bigwedge^{\approx}(Q_f, E)$  to  $\mathcal{Q}_{inf}$ . Note that we only remove pairs from  $\mathcal{Q}_{fin}$ , hence the loop on Lines 13–16 terminates. Afterwards, the query  $Q_{fin}$  is safe range. However, the query  $Q_{inf}$  does not have to be safe range. Indeed, every query  $Q_i \in \mathcal{Q}_{inf}$  has

<sup>1</sup>This statement contained the error we discovered while formalizing the result presented in our conference paper [RBKT22b]. There we had wrongly used the naive conjunction  $Q_f \wedge (\bigwedge_{Q \in \approx(E)} Q)$ , which will not be safe range whenever  $E$  has more than one element, instead of the more carefully constructed  $\bigwedge^{\approx}(Q_f, E)$ .

range-restricted bound variables, but not all the free variables of  $Q_i$  need be range restricted and thus the query  $\exists \vec{fv}(Q_i). Q_i$  does not have to be safe range. But the query  $Q_{inf}$  is closed and thus the inf-equivalent query  $rb(Q_{inf})$  with range-restricted bound variables is safe range.

**Lemma 4.7.**  $\clubsuit$  *Let  $Q$  be an RC query and  $\text{split}(Q) = (Q_{fin}, Q_{inf})$ . Then the queries  $Q_{fin}$  and  $Q_{inf}$  are safe range;  $fv(Q_{fin}) = fv(Q)$  unless  $Q_{fin}$  is syntactically equal to  $\perp$ ; and  $fv(Q_{inf}) = \emptyset$ .*

**Lemma 4.8.**  $\clubsuit$  *Let  $Q$  be an RC query and  $\text{split}(Q) = (Q_{fin}, Q_{inf})$ . If  $\models Q_{inf}$ , then  $\llbracket Q \rrbracket$  is an infinite set. Otherwise,  $\llbracket Q \rrbracket = \llbracket Q_{fin} \rrbracket$  is a finite set.*

By Lemma 4.7,  $Q_{fin}$  is a safe-range (and thus also domain-independent) query. Hence, for the fixed structure, the tuples in  $\llbracket Q_{fin} \rrbracket$  only contain elements in the active domain  $\text{adom}(Q_{fin})$ , i.e.,  $\llbracket Q_{fin} \rrbracket = \llbracket Q_{fin} \rrbracket \cap \text{adom}(Q_{fin})^{|\text{fv}(Q_{fin})|}$ . Our translation does not introduce new constants in  $Q_{fin}$  and thus  $\text{adom}(Q_{fin}) \subseteq \text{adom}(Q)$ . Hence, by Lemma 4.8, if  $\not\models Q_{inf}$ , then  $\llbracket Q_{fin} \rrbracket$  is equal to the “output-restricted unlimited interpretation” [HS94] of  $Q$ , i.e.,  $\llbracket Q_{fin} \rrbracket = \llbracket Q \rrbracket \cap \text{adom}(Q)^{|\text{fv}(Q)|}$ . In contrast, if  $\models Q_{inf}$ , then  $\llbracket Q_{fin} \rrbracket = \llbracket Q \rrbracket \cap \text{adom}(Q)^{|\text{fv}(Q)|}$  does not necessarily hold. For instance, for  $Q := \neg B(x)$ , our translation yields  $\text{split}(Q) = (\perp, \top)$ . In this case, we have  $Q_{inf} = \top$  and thus  $\models Q_{inf}$  because  $\neg B(x)$  is satisfied by infinitely many tuples over an infinite domain. However, if  $B(x)$  is never satisfied, then  $\llbracket Q_{fin} \rrbracket = \emptyset$  is not equal to  $\llbracket Q \rrbracket \cap \text{adom}(Q)^{|\text{fv}(Q)|}$ .

Next, we demonstrate different aspects of our translation on a few examples. Thereby, we use a mildly modified algorithm that performs constant propagation after all steps that could introduce constants  $\top$  or  $\perp$  in a subquery. This optimization keeps the queries small, but is not necessary for termination and correctness. (In contrast, the constant propagation that is part of the substitution operators  $Q[x \mapsto y]$  and  $Q[x/\perp]$  is necessary.) We have verified in Isabelle that our results hold for the modified algorithm. That is, for all above theorems, we proved two variants: one with and one without additional constant propagation steps.

**Example 4.9.**  $\clubsuit$  *Consider the query  $Q := B(x) \vee P(x, y)$ . The variable  $y$  is not range restricted in  $Q$  and thus  $\text{split}(Q)$  restricts  $y$  by a conjunction of  $Q$  with  $P(x, y)$ . However, if  $Q[y/\perp] = B(x)$  is satisfied by some tuple, then  $\llbracket Q \rrbracket$  contains infinitely many tuples. Hence,  $\text{split}(Q) = ((B(x) \vee P(x, y)) \wedge P(x, y), \exists x. B(x))$ . Because  $Q_{fin} = (B(x) \vee P(x, y)) \wedge P(x, y)$  is only used if  $\not\models Q_{inf}$ , i.e., if  $B(x)$  is never satisfied, we could simplify  $Q_{fin}$  to  $P(x, y)$ . However, our translation does not implement such heuristic simplifications.*

**Example 4.10.**  $\clubsuit$  *Consider the query  $Q := B(x) \wedge u \approx v$ . The variables  $u$  and  $v$  are not range restricted in  $Q$  and thus  $\text{split}(Q)$  chooses one of these variables (e.g.,  $u$ ) and restricts it by splitting  $Q$  into  $Q_f = B(x)$  and  $E = \{(u, v)\}$ . Now, all variables are range restricted in  $Q_f$ , but the variables in  $Q_f$  and  $E$  are disjoint. Hence,  $\llbracket Q \rrbracket$  contains infinitely many tuples whenever  $Q_f$  is satisfied by some tuple. In contrast,  $\llbracket Q \rrbracket = \emptyset$  if  $Q_f$  is never satisfied. Hence, we have  $\text{split}(Q) = (\perp, \exists x. B(x))$ .*

**Example 4.11.**  $\clubsuit$  *Consider the query  $Q_{user}^{susp} := B(b) \wedge \exists s. \forall p. P(b, p) \longrightarrow S(p, u, s)$  from Section 1. Restricting its bound variables yields the query  $rb(Q_{user}^{susp}) = B(b) \wedge ((\exists s. (\neg \exists p. P(b, p) \wedge \neg S(p, u, s)) \wedge (\exists p. S(p, u, s))) \vee (\neg \exists p. P(b, p)))$  derived in Example 4.5. Splitting  $Q_{user}^{susp}$  yields*

$$\text{split}(Q_{user}^{susp}) = \left( rb(Q_{user}^{susp}) \wedge (\exists s, p. S(p, u, s)), \exists b. B(b) \wedge \neg \exists p. P(b, p) \right).$$

To understand  $\text{split}(Q_{\text{user}}^{\text{susp}})$ , we apply  $(\star)$  to  $\text{rb}(Q_{\text{user}}^{\text{susp}})$  for the free variable  $u$ :

$$\text{rb}(Q_{\text{user}}^{\text{susp}}) \equiv \left( \text{rb}(Q_{\text{user}}^{\text{susp}}) \wedge (\exists s, p. \text{S}(p, u, s)) \right) \vee \left( \text{B}(b) \wedge (\neg \exists p. \text{P}(b, p)) \wedge \neg \exists s, p. \text{S}(p, u, s) \right).$$

If the subquery  $\text{B}(b) \wedge (\neg \exists p. \text{P}(b, p))$  from the second disjunct is satisfied for some  $b$ , then  $Q_{\text{user}}^{\text{susp}}$  is satisfied by infinitely many values for  $u$  from the infinite domain  $\mathcal{D}$  that do not belong to the finite interpretation of  $\text{S}(p, u, s)$  and thus satisfy the subquery  $\neg \exists s, p. \text{S}(p, u, s)$ . Hence,  $\llbracket Q_{\text{user}}^{\text{susp}} \rrbracket^{\mathcal{S}} = \llbracket \text{rb}(Q_{\text{user}}^{\text{susp}}) \rrbracket^{\mathcal{S}}$  is an infinite set of tuples whenever  $\text{B}(b) \wedge \neg \exists p. \text{P}(b, p)$  is satisfied for some  $b$ . In contrast, if  $\text{B}(b) \wedge \neg \exists p. \text{P}(b, p)$  is not satisfied for any  $b$ , then  $Q_{\text{user}}^{\text{susp}}$  is equivalent to  $\text{rb}(Q_{\text{user}}^{\text{susp}}) \wedge (\exists s, p. \text{S}(p, u, s))$  obtained also by applying  $(\star)$  to  $Q_{\text{user}}^{\text{susp}}$  for the free variable  $u$ .

## 5. COMPLEXITY ANALYSIS

We analyze the time complexity of capturing  $Q$ , i.e., checking if  $\llbracket Q \rrbracket$  is finite and enumerating  $\llbracket Q \rrbracket$  in this case. To bound the asymptotic time complexity of capturing a fixed  $Q$ , we need to apply an additional standard translation step to both queries produced by our translation to obtain two RANF queries. Query cost (Section 3.5) can then be applied to the resulting two queries to bound computation time based on the cardinalities of subquery evaluation results.

**Definition 5.1.** Let  $Q$  be an RC query and  $\text{split}(Q) = (Q_{\text{fin}}, Q_{\text{inf}})$ . Let  $\hat{Q}_{\text{fin}} := \text{sr2ranf}(Q_{\text{fin}})$  and  $\hat{Q}_{\text{inf}} := \text{sr2ranf}(Q_{\text{inf}})$  be the equivalent RANF queries. We define  $\text{rw}(Q) := (\hat{Q}_{\text{fin}}, \hat{Q}_{\text{inf}})$ .

Since function  $\text{sr2ranf}(\cdot)$  is a standard translation step, we present it in Section 6 (see Figure 12). Note that the proof of Lemma 5.6 relies on its algorithmic details.

We ignore the (constant) time complexity of computing  $\text{rw}(Q) = (\hat{Q}_{\text{fin}}, \hat{Q}_{\text{inf}})$  and focus on the time complexity of evaluating the RANF queries  $\hat{Q}_{\text{fin}}$  and  $\hat{Q}_{\text{inf}}$ , i.e., the query cost of  $\hat{Q}_{\text{fin}}$  and  $\hat{Q}_{\text{inf}}$ . Without loss of generality, we assume that the input query  $Q$  has pairwise distinct (free and bound) variables to derive a set of quantified predicates from  $Q$ 's atomic predicates and formulate our time complexity bound. Still, the RANF queries  $\hat{Q}_{\text{fin}}$  and  $\hat{Q}_{\text{inf}}$  computed by our translation need not have pairwise distinct (free and bound) variables.

We define the relation  $\lesssim_Q$  on  $\text{av}(Q)$  such that  $x \lesssim_Q y$  iff the scope of an occurrence of  $x \in \text{av}(Q)$  is contained in the scope of an occurrence of  $y \in \text{av}(Q)$ . Formally, we define  $x \lesssim_Q y$  iff  $y \in \text{fv}(Q)$  or  $\exists x. Q_x \sqsubseteq \exists y. Q_y \sqsubseteq Q$  for some  $Q_x$  and  $Q_y$ . Note that  $\lesssim_Q$  is a preorder on all variables and a partial order on the bound variables for every query with pairwise distinct (free and bound) variables.

Let  $\text{aps}(Q)$  be the set of all atomic predicates in a query  $Q$ . We denote by  $\overline{\text{aps}}(Q)$  the set of quantified predicates obtained from  $\text{aps}(Q)$  by performing the variable substitution  $x \mapsto y$ , where  $x$  and  $y$  are related by equalities in  $Q$  and  $x \lesssim_Q y$ , and existentially quantifying from a quantified predicate  $Q_{qp}$  the innermost bound variable  $x$  in  $Q$  that is free in  $Q_{qp}$ . Let  $\text{eqs}^*(Q)$  be the transitive closure of equalities occurring in  $Q$ . Formally, we define  $\overline{\text{aps}}(Q)$  by:

- $Q_{ap} \in \overline{\text{aps}}(Q)$  if  $Q_{ap} \in \text{aps}(Q)$ ;
- $Q_{qp}[x \mapsto y] \in \overline{\text{aps}}(Q)$  if  $Q_{qp} \in \overline{\text{aps}}(Q)$ ,  $(x, y) \in \text{eqs}^*(Q)$ , and  $x \lesssim_Q y$ ;
- $\exists x. Q_{qp} \in \overline{\text{aps}}(Q)$  if  $Q_{qp} \in \overline{\text{aps}}(Q)$ ,  $x \in \text{fv}(Q_{qp}) \setminus \text{fv}(Q)$ , and  $x \lesssim_Q y$  for all  $y \in \text{fv}(Q_{qp})$ .

We bound the time complexity of capturing  $Q$  by considering subsets  $\mathcal{Q}_{qps}$  of quantified predicates  $\overline{\text{aps}}(Q)$  that are *minimal* in the sense that every quantified predicate in  $\mathcal{Q}_{qps}$  contains a unique free variable that is not free in any other quantified predicate in  $\mathcal{Q}_{qps}$ . Formally, we define  $\text{minimal}(\mathcal{Q}_{qps}) := \forall Q_{qp} \in \mathcal{Q}_{qps}. \text{fv}(\mathcal{Q}_{qps} \setminus \{Q_{qp}\}) \neq \text{fv}(Q_{qp})$ . Every minimal



subset  $\mathcal{Q}_{qps}$  of quantified predicates  $\overline{\text{qps}}(Q)$  contributes the product of the numbers of tuples satisfying each quantified predicate  $Q_{qp} \in \mathcal{Q}_{qps}$  to the overall bound (that product is an upper bound on the number of tuples satisfying the join over all  $Q_{qp} \in \mathcal{Q}_{qps}$ ). Similarly to Ngo et al. [NRR13], we use the notation  $\tilde{O}(\cdot)$  to hide logarithmic factors incurred by set operations.

**Theorem 5.2.** *Let  $Q$  be a fixed RC query with pairwise distinct (free and bound) variables. The time complexity of capturing  $Q$ , i.e., checking if  $\llbracket Q \rrbracket$  is finite and enumerating  $\llbracket Q \rrbracket$  in this case, is in  $\tilde{O}\left(\sum_{\mathcal{Q}_{qps} \subseteq \overline{\text{qps}}(Q), \text{minimal}(\mathcal{Q}_{qps})} \prod_{Q_{qp} \in \mathcal{Q}_{qps}} \llbracket Q_{qp} \rrbracket\right)$ .*

Before we prove Theorem 5.2 we first provide some examples to reinforce the intuition behind our claim. Examples 5.3 and 5.4 show that the time complexity from Theorem 5.2 cannot be achieved by the translation of Van Gelder and Topor [GT91] or over finite domains. Example 5.5 shows how equalities affect the bound in Theorem 5.2.

**Example 5.3.** *Consider the query  $Q := \mathbf{B}(b) \wedge \exists u, s. \neg \exists p. \mathbf{P}(b, p) \wedge \neg \mathbf{S}(p, u, s)$ , equivalent to  $Q^{sup}$  from Section 1. Then  $\text{aps}(Q) = \{\mathbf{B}(b), \mathbf{P}(b, p), \mathbf{S}(p, u, s)\}$  and  $\overline{\text{qps}}(Q) = \{\mathbf{B}(b), \mathbf{P}(b, p), \exists p. \mathbf{P}(b, p), \mathbf{S}(p, u, s), \exists p. \mathbf{S}(p, u, s), \exists s, p. \mathbf{S}(p, u, s), \exists u, s, p. \mathbf{S}(p, u, s)\}$ . The translated query  $Q_{vgt}$  by Van Gelder and Topor [GT91]*

$$\begin{aligned} & \left( (\exists s, p. \mathbf{S}(p, u, s)) \wedge (\exists u, p. \mathbf{S}(p, u, s)) \wedge \mathbf{B}(b) \right) \wedge \\ & \neg \exists p. \left( (\exists s, p. \mathbf{S}(p, u, s)) \wedge (\exists u, p. \mathbf{S}(p, u, s)) \wedge \mathbf{P}(b, p) \right) \wedge \neg \mathbf{S}(p, u, s) \end{aligned}$$

restricts the variables  $u$  and  $s$  by  $\exists s, p. \mathbf{S}(p, u, s)$  and  $\exists u, p. \mathbf{S}(p, u, s)$ , respectively. Note that this corresponds to the RA expression shown in Section 1 with the highlighted generators replaced with  $\pi_{\text{user}}(\mathbf{S}) \times \pi_{\text{score}}(\mathbf{S})$ .

Consider an interpretation of  $\mathbf{B}$  by  $\{(c') \mid c' \in \{1, \dots, n\}\}$ ,  $\mathbf{P}$  by  $\{(c', c') \mid c' \in \{1, \dots, n\}\}$ , and  $\mathbf{S}$  by  $\{(c, c', c') \mid c \in \{1, \dots, n\}, c' \in \{1, \dots, m\}\}$ ,  $n, m \in \mathbb{N}$ . Computing the join of  $\mathbf{P}(b, p)$ ,  $\exists s, p. \mathbf{S}(p, u, s)$ , and  $\exists u, p. \mathbf{S}(p, u, s)$ , which is a Cartesian product, results in a time complexity in  $\Omega(n \cdot m^2)$  for  $Q_{vgt}$ . In contrast, Theorem 5.2 yields an asymptotically better time complexity in  $\tilde{O}(n + m + n \cdot m)$  for our translation, more precisely:

$$\tilde{O}(\llbracket \mathbf{B}(b) \rrbracket + \llbracket \mathbf{P}(b, p) \rrbracket + \llbracket \mathbf{S}(p, u, s) \rrbracket + (\llbracket \mathbf{B}(b) \rrbracket + \llbracket \mathbf{P}(b, p) \rrbracket) \cdot \llbracket \mathbf{S}(p, u, s) \rrbracket),$$

which corresponds to the complexity of evaluating the RA expression shown in Section 1.

**Example 5.4.** *The query  $\neg \mathbf{S}(x, y, z)$  is satisfied by a finite set of tuples over a finite domain  $\mathcal{D}$  (as is every query over a finite domain). For an interpretation of  $\mathbf{S}$  by  $\{(c, c, c) \mid c \in \mathcal{D}\}$ , the equality  $|\mathcal{D}| = \llbracket \mathbf{S}(x, y, z) \rrbracket$  holds and the number of satisfying tuples is*

$$\llbracket \neg \mathbf{S}(x, y, z) \rrbracket = |\mathcal{D}|^3 - \llbracket \mathbf{S}(x, y, z) \rrbracket = \llbracket \mathbf{S}(x, y, z) \rrbracket^3 - \llbracket \mathbf{S}(x, y, z) \rrbracket \in \Omega(\llbracket \mathbf{S}(x, y, z) \rrbracket^3),$$

which exceeds the bound  $\tilde{O}(\llbracket \mathbf{S}(x, y, z) \rrbracket)$  of Theorem 5.2. Hence, our infinite domain assumption is crucial for achieving the better complexity bound.

**Example 5.5.** *Consider the following query over the infinite domain  $\mathcal{D} = \mathbb{N}$  of natural numbers:*

$$\begin{aligned} Q := & \forall u. (u \approx 0 \vee u \approx 1 \vee u \approx 2) \longrightarrow \\ & (\exists v. \mathbf{B}(v) \wedge (u \approx 0 \longrightarrow x \approx v) \wedge (u \approx 1 \longrightarrow y \approx v) \wedge (u \approx 2 \longrightarrow z \approx v)). \end{aligned}$$

Note that this query is equivalent to  $Q \equiv \mathbf{B}(x) \wedge \mathbf{B}(y) \wedge \mathbf{B}(z)$  and thus it is satisfied by a finite set of tuples of size  $\llbracket \mathbf{B}(x) \rrbracket \cdot \llbracket \mathbf{B}(y) \rrbracket \cdot \llbracket \mathbf{B}(z) \rrbracket = \llbracket \mathbf{B}(x) \rrbracket^3$ . The set of atomic predicates of  $Q$  is  $\text{aps}(Q) = \{\mathbf{B}(v)\}$  and it must be closed under the equalities occurring in  $Q$  to yield a

valid bound in Theorem 5.2. In this case,  $\overline{\text{qps}}(Q) = \{\mathbf{B}(v), \exists v. \mathbf{B}(v), \mathbf{B}(x), \mathbf{B}(y), \mathbf{B}(z)\}$  and the bound in Theorem 5.2 is  $\|\mathbf{B}(v)\| \cdot \|\mathbf{B}(x)\| \cdot \|\mathbf{B}(y)\| \cdot \|\mathbf{B}(z)\| = \|\mathbf{B}(x)\|^4$ . In particular, this bound is not tight, but it still reflects the complexity of evaluating the RANF queries produced by our translation as it does not derive the equivalence  $Q \equiv \mathbf{B}(x) \wedge \mathbf{B}(y) \wedge \mathbf{B}(z)$ .

Now, to prove Theorem 5.2, we need to introduce guard queries and the set of quantified predicates of a query. Given a RANF query  $\hat{Q}$ , we define a *guard* query  $\text{guard}(\hat{Q})$  that is implied by  $\hat{Q}$ , i.e.,  $\text{guard}(\hat{Q})$  can be used to over-approximate the set of satisfying tuples for  $\hat{Q}$ . We use this over-approximation in our proof of Theorem 5.2. The guard query  $\text{guard}(\hat{Q})$  has a simple structure: it is the disjunction of conjunctions of quantified predicates and equalities.

We now define the set of quantified predicates  $\text{qps}(Q)$  occurring in the guard query  $\text{guard}(Q)$ . For an atomic predicate  $Q_{ap} \in \text{aps}(Q)$ , let  $\mathcal{B}_Q(Q_{ap})$  be the set of sequences of bound variables for all occurrences of  $Q_{ap}$  in  $Q$ . For example, consider a query  $Q_{ex} := ((\exists z. (\exists y, z. P_3(x, y, z)) \wedge P_2(y, z)) \wedge P_1(z)) \vee P_3(x, y, z)$ . Then  $\text{aps}(Q_{ex}) = \{P_1(z), P_2(y, z), P_3(x, y, z)\}$  and  $\mathcal{B}_{Q_{ex}}(P_3(x, y, z)) = \{\{y, z\}, \square\}$ , where  $\square$  denotes the empty sequence corresponding to the occurrence of  $P_3(x, y, z)$  in  $Q_{ex}$  for which the variables  $x, y, z$  are all free in  $Q_{ex}$ . Note that the variable  $z$  in the other occurrence of  $P_3(x, y, z)$  in  $Q_{ex}$  is bound to the innermost quantifier. Hence, neither  $[z, y]$  nor  $[z, y, z]$  are in  $\mathcal{B}_{Q_{ex}}(P_3(x, y, z))$ . Furthermore, let  $\text{qps}(Q)$  be the set of the quantified predicates obtained by existentially quantifying sequences of bound variables in  $\mathcal{B}_{Q'}(Q_{ap})$  from the atomic predicates  $Q_{ap} \in \text{aps}(Q')$  in all subqueries  $Q'$  of  $Q$ . Formally,  $\text{qps}(Q) := \bigcup_{Q' \sqsubseteq Q, Q_{ap} \in \text{aps}(Q')} \{\exists \vec{v}. Q_{ap} \mid \vec{v} \in \mathcal{B}_{Q'}(Q_{ap})\}$ . For instance,  $\text{qps}(Q_{ex}) = \{P_3(x, y, z), \exists z. P_3(x, y, z), \exists y, z. P_3(x, y, z), P_2(y, z), \exists z. P_2(y, z), P_1(z)\}$ .

A crucial property of our translation, which is central for the proof of Theorem 5.2, is the relationship between the quantified predicates  $\text{qps}(\hat{Q})$  for a RANF query  $\hat{Q}$  produced by our translation and the original query  $Q$ . The relationship is formalized in the following lemma.

**Lemma 5.6.** *Let  $Q$  be an RC query with pairwise distinct (free and bound) variables and let  $\text{rw}(Q) = (\hat{Q}_{fin}, \hat{Q}_{inf})$ . Let  $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$ . Then  $\text{qps}(\hat{Q}) \subseteq \overline{\text{qps}}(Q)$ .*

*Proof.* Let  $\text{split}(Q) = (Q_{fin}, Q_{inf})$ . We observe that  $\text{aps}(Q_{fin}) \subseteq \overline{\text{qps}}(Q)$ ,  $\text{eqs}^*(Q_{fin}) \subseteq \text{eqs}^*(Q)$ ,  $\lesssim_{Q_{fin}} \subseteq \lesssim_Q$ ,  $\text{aps}(Q_{inf}) \subseteq \overline{\text{qps}}(Q)$ ,  $\text{eqs}^*(Q_{inf}) \subseteq \text{eqs}^*(Q)$ , and  $\lesssim_{Q_{inf}} \subseteq \lesssim_Q$ . Hence,  $\overline{\text{qps}}(Q_{fin}) \subseteq \overline{\text{qps}}(Q)$  and  $\overline{\text{qps}}(Q_{inf}) \subseteq \overline{\text{qps}}(Q)$ .

Next we observe that  $\text{qps}(Q') \subseteq \overline{\text{qps}}(Q')$  for every query  $Q'$ . Finally, we show that  $\text{qps}(\hat{Q}_{fin}) \subseteq \text{qps}(Q_{fin})$  and  $\text{qps}(\hat{Q}_{inf}) \subseteq \text{qps}(Q_{inf})$ . We observe that  $\mathcal{B}_{\text{cp}(Q')}(Q_{ap}) \subseteq \mathcal{B}_{Q'}(Q_{ap})$ ,  $\mathcal{B}_{\text{srnf}(Q')}(Q_{ap}) \subseteq \mathcal{B}_{Q'}(Q_{ap})$ , and then  $\text{qps}(\text{cp}(Q')) \subseteq \text{qps}(Q')$ ,  $\text{qps}(\text{srnf}(Q')) \subseteq \text{qps}(Q')$ , for every query  $Q'$ .

Assume that  $Q' \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}$  is a safe-range query in which no variable occurs both free and bound, no bound variables shadow each other, i.e., there are no subqueries  $\exists x. Q_x \sqsubseteq Q'_x$  and  $\exists x. Q'_x \sqsubseteq Q' \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}$ , and every two subqueries  $\exists x. Q_x \sqsubseteq Q_1$  and  $\exists x. Q'_x \sqsubseteq Q_2$  such that  $Q_1 \wedge Q_2 \sqsubseteq Q' \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}$  have the property that  $\exists x. Q_x$  or  $\exists x. Q'_x$  is a quantified predicate. Then the free variables in  $\bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}$  never clash with the bound variables in  $Q'$ , i.e., Line 26 in Figure 12 is never executed. Next we observe that  $\mathcal{B}_{\text{sr2ranf}(Q', \mathcal{Q})}(Q_{ap}) \subseteq \mathcal{B}_{Q' \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}}(Q_{ap})$  (this subset relation only holds when considering queries modulo  $\alpha$ -equivalence, i.e., queries that have the same binding structure but differ in the used bound variable names are considered to be equal) and then  $\text{qps}(\text{sr2ranf}(Q', \mathcal{Q})) \subseteq \text{qps}(Q' \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q})$ . Because  $Q_{fin}, Q_{inf}$  have the assumed properties and  $\text{qps}(\text{srnf}(Q')) \subseteq \text{qps}(Q')$ , for every query  $Q'$ , we get  $\text{qps}(\hat{Q}_{fin}) = \text{qps}(\text{sr2ranf}(Q_{fin})) \subseteq \text{qps}(Q_{fin})$  and  $\text{qps}(\hat{Q}_{inf}) = \text{qps}(\text{sr2ranf}(Q_{inf})) \subseteq \text{qps}(Q_{inf})$ .  $\square$

Recall Example 5.3. The query  $\exists u, p. S(p, u, s)$  is in  $\mathbf{qps}(Q_{vgt})$ , but not in  $\overline{\mathbf{qps}}(Q)$ . Hence,  $\mathbf{qps}(Q_{vgt}) \subseteq \overline{\mathbf{qps}}(Q)$ , i.e., an analogue of Lemma 5.6 for Van Gelder and Topor's translation, does not hold.

Every tuple satisfying a RANF query  $\hat{Q}$  belongs to the set of tuples satisfying the join over some minimal subset  $\mathcal{Q}_{qps} \subseteq \mathbf{qps}(\hat{Q})$  of quantified predicates and satisfying equalities duplicating some of  $\mathcal{Q}_{qps}$ 's columns. Hence, we define the guard query  $\mathbf{guard}(\hat{Q})$  as follows:

$$\mathbf{guard}(\hat{Q}) := \bigvee_{\substack{\mathcal{Q}_{qps} \subseteq \mathbf{qps}(\hat{Q}), \text{minimal}(\mathcal{Q}_{qps}), \\ \mathcal{Q}^{\approx} \subseteq \{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in \text{fv}(\hat{Q})\}, \\ \text{fv}(\mathcal{Q}_{qps}) \cup \text{fv}(\mathcal{Q}^{\approx}) = \text{fv}(\hat{Q})}} \left( \bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \wedge \bigwedge_{Q^{\approx} \in \mathcal{Q}^{\approx}} Q^{\approx} \right).$$

Note that  $\{x \approx y \mid x \in V \wedge y \in V'\}$  denotes the set of all equalities  $x \approx y$  between variables  $x \in V$  and  $y \in V'$ . We express the correctness of the guard query in the following lemma.

**Lemma 5.7.** *Let  $\hat{Q}$  be a RANF query. Then, for all variable assignments  $\alpha$ ,*

$$\alpha \models \hat{Q} \implies \alpha \models \mathbf{guard}(\hat{Q}).$$

Moreover,  $\text{fv}(\mathbf{guard}(\hat{Q})) = \text{fv}(\hat{Q})$  unless  $\mathbf{guard}(\hat{Q}) = \perp$ . Hence,  $\llbracket \hat{Q} \rrbracket$  satisfies

$$\llbracket \hat{Q} \rrbracket \subseteq \bigcup_{\substack{\mathcal{Q}_{qps} \subseteq \mathbf{qps}(\hat{Q}), \text{minimal}(\mathcal{Q}_{qps}), \\ \mathcal{Q}^{\approx} \subseteq \{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in \text{fv}(\hat{Q})\}, \\ \text{fv}(\mathcal{Q}_{qps}) \cup \text{fv}(\mathcal{Q}^{\approx}) = \text{fv}(\hat{Q})}} \left[ \bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \wedge \bigwedge_{Q^{\approx} \in \mathcal{Q}^{\approx}} Q^{\approx} \right].$$

*Proof.* The statement follows by well-founded induction over the definition of  $\mathbf{ranf}(\hat{Q})$ .  $\square$

We now derive a bound on  $\llbracket \hat{Q}' \rrbracket$ , for an arbitrary RANF subquery  $\hat{Q}' \sqsubseteq \hat{Q}$ ,  $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$ .

**Lemma 5.8.** *Let  $Q$  be an RC query with pairwise distinct (free and bound) variables and let  $\text{rw}(Q) = (\hat{Q}_{fin}, \hat{Q}_{inf})$ . Let  $\hat{Q}' \sqsubseteq \hat{Q}$  be a RANF subquery of  $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$ . Then*

$$\left| \llbracket \hat{Q}' \rrbracket \right| \leq \sum_{\mathcal{Q}_{qps} \subseteq \overline{\mathbf{qps}}(Q), \text{minimal}(\mathcal{Q}_{qps})} 2^{|\text{av}(\hat{Q})|} \cdot \prod_{Q_{qp} \in \mathcal{Q}_{qps}} \llbracket Q_{qp} \rrbracket.$$

*Proof.* Applying Lemma 5.7 to the RANF query  $\hat{Q}'$  yields

$$\llbracket \hat{Q}' \rrbracket \subseteq \bigcup_{\substack{\mathcal{Q}_{qps} \subseteq \mathbf{qps}(\hat{Q}'), \text{minimal}(\mathcal{Q}_{qps}), \\ \mathcal{Q}^{\approx} \subseteq \{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in \text{fv}(\hat{Q}')\}, \\ \text{fv}(\mathcal{Q}_{qps}) \cup \text{fv}(\mathcal{Q}^{\approx}) = \text{fv}(\hat{Q}')}} \left[ \bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \wedge \bigwedge_{Q^{\approx} \in \mathcal{Q}^{\approx}} Q^{\approx} \right].$$

We observe that  $\llbracket \bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \wedge \bigwedge_{Q^{\approx} \in \mathcal{Q}^{\approx}} Q^{\approx} \rrbracket \leq \llbracket \bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \rrbracket \leq \prod_{Q_{qp} \in \mathcal{Q}_{qps}} \llbracket Q_{qp} \rrbracket$  where the first inequality follows from the fact that equalities  $Q^{\approx} \in \mathcal{Q}^{\approx}$  can only restrict a set of tuples and duplicate columns. Because  $\hat{Q}'$  is a subquery of  $\hat{Q}$ , it follows that  $\mathbf{qps}(\hat{Q}') \subseteq \mathbf{qps}(\hat{Q})$ . Lemma 5.6 yields  $\mathbf{qps}(\hat{Q}) \subseteq \overline{\mathbf{qps}}(Q)$ . Hence, we derive  $\mathbf{qps}(\hat{Q}') \subseteq \overline{\mathbf{qps}}(Q)$ .

The number of equalities in  $\{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in \text{fv}(\hat{Q}')\}$  is at most

$$|\text{fv}(\mathcal{Q}_{qps})| \cdot |\text{fv}(\hat{Q}')| \leq |\text{fv}(\hat{Q}')|^2 \leq |\text{av}(\hat{Q})|^2.$$

The first inequality holds because  $\text{fv}(\mathcal{Q}_{qps}) \cup \text{fv}(\mathcal{Q}^{\approx}) = \text{fv}(\hat{Q}')$  and thus  $\text{fv}(\mathcal{Q}_{qps}) \subseteq \text{fv}(\hat{Q}')$ . The second inequality holds because the variables in a subquery  $\hat{Q}'$  of  $\hat{Q}$  are in  $\text{av}(\hat{Q})$ . Hence, the number of subsets  $\mathcal{Q}^{\approx} \subseteq \{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in \text{fv}(\hat{Q}')\}$  is at most  $2^{|\text{av}(\hat{Q})|^2}$ .  $\square$

We now bound the query cost of a RANF query  $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$  over the fixed structure  $\mathcal{S}$ .

**Lemma 5.9.** *Let  $Q$  be an RC query with pairwise distinct (free and bound) variables and let  $\text{rw}(Q) = (\hat{Q}_{fin}, \hat{Q}_{inf})$ . Let  $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$ . Then*

$$\text{cost}^{\mathcal{S}}(\hat{Q}) \leq \left| \text{sub}(\hat{Q}) \right| \cdot \left| \text{av}(\hat{Q}) \right| \cdot 2^{|\text{av}(\hat{Q})|} \cdot \sum_{\mathcal{Q}_{qps} \subseteq \overline{\text{qps}}(Q), \text{minimal}(\mathcal{Q}_{qps})} \prod_{Q_{qp} \in \mathcal{Q}_{qps}} \llbracket Q_{qp} \rrbracket.$$

*Proof.* Recall that  $|\text{sub}(\hat{Q})|$  denotes the number of subqueries of the query  $\hat{Q}$  and thus bounds the number of RANF subqueries  $\hat{Q}'$  of the query  $\hat{Q}$ . For every subquery  $\hat{Q}'$  of  $\hat{Q}$ , we first use the fact that  $|\text{fv}(\hat{Q}')| \leq |\text{av}(\hat{Q})|$  to bound  $\llbracket \hat{Q}' \rrbracket \cdot |\text{fv}(\hat{Q}')| \leq \llbracket \hat{Q}' \rrbracket \cdot |\text{av}(\hat{Q})|$ . Then we use the estimation of  $\llbracket \hat{Q}' \rrbracket$  by Lemma 5.8.  $\square$

Finally, we prove Theorem 5.2.

*Proof of Theorem 5.2.* We derive Theorem 5.2 from Lemma 5.9 and the fact that the quantities  $|\text{sub}(\hat{Q})|$ ,  $|\text{av}(\hat{Q})|$ , and  $2^{|\text{av}(\hat{Q})|^2}$  only depend on the query  $Q$  and thus they do not contribute to the asymptotic time complexity of capturing a fixed query  $Q$ .  $\square$

## 6. IMPLEMENTATION

We have implemented our translation RC2SQL consisting of roughly 1000 lines of OCaml code [RBKT22a]. It consists of multiple translation steps that take an arbitrary relational calculus (RC) query and produce two SQL queries.

Figure 9 summarizes the order of the translation steps and the functions that implement them. The function `split`( $\cdot$ ) (Section 4.3), applied in the first step, is the main part of our translation. Recall that it takes an arbitrary RC query and returns two safe-range RC queries. Next, the function `srnf`( $\cdot$ ) (Section 6.1) converts both queries to safe-range normal form (SRNF), followed by the function `sr2ranf`( $\cdot, \cdot$ ) (Section 6.2) that converts SRNF queries into relation algebra normal form (RANF). Both normal forms were defined in Section 3. For simplicity, we define a function `sr2ranf`( $\cdot$ ) that combines the previous two functions and can be applied to any safe-range RC query. In addition to the worst-case complexity, we further improve our translation’s average-case complexity by implementing the optimizations inspired by Claußen et al. [CKMP97]. The function `optcnt`( $\cdot$ ) (Section 6.3) implements these optimizations on the RANF queries. Finally, to derive SQL queries from the RANF queries we first obtain equivalent relational algebra (RA) expressions following a (slightly modified) standard approach [AHV95] implemented by the function `ranf2ra`( $\cdot$ ) (Section 6.4). To translate the RA expressions into SQL, we reuse a publicly available RA interpreter `radb` [Yan19] (Section 6.5). We name the composition of the last two steps `ranf2sql`( $\cdot$ ).

To resolve the nondeterministic choices present in our algorithms (Section 6.6) we always choose the alternative with the lowest query cost. The query cost is estimated by using a sample structure of constant size, called a *training database*. A good training database should preserve the relative ordering of queries by their cost over the actual database as much as possible. Nevertheless, our translation satisfies the correctness and worst-case complexity claims independently of the choice of the training database.

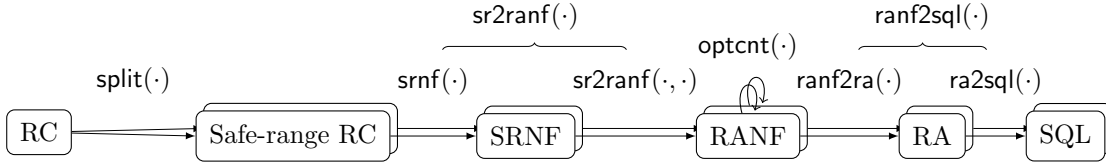


Figure 9: Overview of the functions used in our implementation.

**input:** An RC query  $Q$ .

**output:** A SRNF query  $Q_{srnf}$  such that  $Q \equiv Q_{srnf}$ ,  $\text{fv}(Q) = \text{fv}(Q_{srnf})$ .

```

1 function srnf( $Q$ ) =
2   switch  $Q$  do
3     case  $\neg Q'$  do
4       switch  $Q'$  do
5         case  $\neg Q''$  do return srnf( $Q''$ );
6         case  $Q_1 \vee Q_2$  do return srnf( $(\neg Q_1) \wedge (\neg Q_2)$ );
7         case  $Q_1 \wedge Q_2$  do return srnf( $(\neg Q_1) \vee (\neg Q_2)$ );
8         case  $\exists \vec{v}. Q_{\vec{v}}$  do
9           if  $\vec{v} \cap \text{fv}(Q_{\vec{v}}) = \emptyset$  then return srnf( $\neg Q_{\vec{v}}$ );
10          else
11            switch srnf( $Q_{\vec{v}}$ ) do
12              case  $Q_1 \vee Q_2$  do return srnf( $(\neg \exists \vec{v}. Q_1) \wedge (\neg \exists \vec{v}. Q_2)$ );
13              otherwise do return  $\neg \exists \vec{v} \cap \text{fv}(Q_{\vec{v}}). \text{srnf}(Q_{\vec{v}})$ ;
14            otherwise do return  $\neg \text{srnf}(Q')$ ;
15         case  $Q_1 \vee Q_2$  do return srnf( $Q_1$ )  $\vee$  srnf( $Q_2$ );
16         case  $Q_1 \wedge Q_2$  do return srnf( $Q_1$ )  $\wedge$  srnf( $Q_2$ );
17         case  $\exists \vec{v}. Q_{\vec{v}}$  do
18           switch srnf( $Q_{\vec{v}}$ ) do
19             case  $Q_1 \vee Q_2$  do return srnf( $(\exists \vec{v}. Q_1) \vee (\exists \vec{v}. Q_2)$ );
20             otherwise do return  $\exists \vec{v} \cap \text{fv}(Q_{\vec{v}}). \text{srnf}(Q_{\vec{v}})$ ;
21         otherwise do return  $Q$ ;

```

Figure 10: Translation to SRNF.

Overall, the translation is formally defined as

$$\text{RC2SQL}(Q) := (Q'_{fin}, Q'_{inf})$$

where  $Q'_{fin} := \text{ranf2sql}(\text{optcnt}(\text{sr2ranf}(Q_{fin})))$ ,  $Q'_{inf} := \text{ranf2sql}(\text{optcnt}(\text{sr2ranf}(Q_{inf})))$ , and  $(Q_{fin}, Q_{inf}) := \text{split}(Q)$ .

**6.1. Translation to SRNF.** Figure 10 defines the function  $\text{srnf}(Q)$  that yields a SRNF query equivalent to  $Q$ . The function  $\text{srnf}(Q)$  pushes negations downwards (Lines 6–7), eliminates double negations (Line 5), drops bound variables that do not occur in the query (Line 9), and distributes existential quantifiers over disjunction (Line 19). The termination of  $\text{srnf}(Q)$  follows using the measure  $\text{m}(Q)$ , shown in Figure 11, that decreases for proper subqueries, after pushing negations and distributing existential quantification over disjunction.

Next we prove a lemma that we use as a precondition for translating safe-range queries in SRNF to queries in RANF.

$$\begin{aligned}
m(\perp) = m(\top) &= m(x \approx t) = 1 \\
m(r(t_1, \dots, t_{\nu(r)})) &= 1 \\
m(\neg Q) &= 2 \cdot m(Q) \\
m(Q_1 \vee Q_2) &= 2 \cdot m(Q_1) + 2 \cdot m(Q_2) + 2 \\
m(Q_1 \wedge Q_2) &= m(Q_1) + m(Q_2) + 1 \\
m(\exists x. Q_x) &= 2 \cdot m(Q_x)
\end{aligned}$$

Figure 11: Measure on RC queries.

**Lemma 6.1.** *Let  $Q_{srnf}$  be a query in SRNF. Then  $\text{gen}(x, \neg Q')$  does not hold for any variable  $x$  and subquery  $\neg Q'$  of  $Q_{srnf}$ .*

*Proof.* Using Figure 4,  $\text{gen}(x, \neg Q')$  can only hold if  $\neg Q'$  has the form  $\neg \neg Q$ ,  $\neg(Q_1 \vee Q_2)$ , or  $\neg(Q_1 \wedge Q_2)$ . The SRNF query  $Q_{srnf}$  cannot have a subquery  $\neg Q'$  that has any such form.  $\square$

**6.2. Translation to RANF.** The function  $\text{sr2ranf}(Q, \mathcal{Q}) = (\hat{Q}, \overline{\mathcal{Q}})$ , defined in Figure 12, where  $\text{sr2ranf}$  stands for *safe range to relational algebra normal form*, takes a safe-range query  $Q \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}$  in SRNF, or in existential normal form (ENF) (see Appendix B) and returns a RANF query  $\hat{Q}$  such that  $Q \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q} \equiv \hat{Q} \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q}$ . To restrict variables in  $Q$ , the function  $\text{sr2ranf}(Q, \mathcal{Q})$  conjoins a subset of queries  $\overline{Q} \subseteq \mathcal{Q}$  to  $Q$ . Given a safe-range query  $Q$ , we first convert  $Q$  into SRNF and set  $\mathcal{Q} = \emptyset$ . Then we define  $\text{sr2ranf}(Q) := \hat{Q}$ , where  $(\hat{Q}, \_) := \text{sr2ranf}(\text{srnf}(Q), \emptyset)$ , to be a RANF query  $\hat{Q}$  equivalent to  $Q$ . The termination of  $\text{sr2ranf}(Q, \mathcal{Q})$  follows from the lexicographic measure  $(2 \cdot m(Q) + \text{eqneg}(Q) + 2 \cdot \sum_{\overline{Q} \in \mathcal{Q}} m(\overline{Q}) + 2 \cdot |\mathcal{Q}|, m(Q) + \sum_{\overline{Q} \in \mathcal{Q}} m(\overline{Q}))$ . Here  $m(Q)$  is defined in Figure 11,  $\text{eqneg}(Q) := 1$  if  $Q$  is an equality between two variables or the negation of a query, and  $\text{eqneg}(Q) := 0$  otherwise.

Next we describe the definition of  $\text{sr2ranf}(Q, \mathcal{Q})$  that follows [AHV95, Algorithm 5.4.7]. Note that no constant propagation (Figure 3) is needed in [AHV95, Algorithm 5.4.7], because the constants  $\perp$  and  $\top$  are not in the query syntax [AHV95, Section 5.3]. Because  $\text{gen}(x, \perp)$  holds and  $x \notin \text{fv}(\perp)$ , we need to perform constant propagation to guarantee that every disjunct has the same set of free variables (e.g., the query  $\perp \vee B(x)$  must be translated to  $B(x)$  to be in RANF). We flatten the disjunction and conjunction using  $\text{flat}^\vee(\cdot)$  and  $\text{flat}^\wedge(\cdot)$ , respectively. In the case of a conjunction  $Q^\wedge$ , we first split the queries from  $\text{flat}^\wedge(Q^\wedge)$  and  $\mathcal{Q}$  into queries  $\mathcal{Q}^+$  that do not have the form of a negation and queries  $\mathcal{Q}^-$  that do. Then we take out equalities between two variables and negations of equalities between two variables from the sets  $\mathcal{Q}^+$  and  $\mathcal{Q}^-$ , respectively. To partition  $\text{flat}^\wedge(Q^\wedge) \cup \mathcal{Q}$  this way, we define the predicates  $\text{neg}(Q)$  and  $\text{eq}(Q)$  characterizing equalities between two variables and negations, respectively, i.e.,  $\text{neg}(Q)$  is true iff  $Q$  has the form  $\neg Q'$  and  $\text{eq}(Q)$  is true iff  $Q$  has the form  $x \approx y$ . Finally, the function  $\text{sort}^\wedge(\mathcal{Q})$  converts a set of queries into a RANF conjunction, defined in Figure 6, i.e., a left-associative conjunction in RANF. Note that the function  $\text{sort}^\wedge(\mathcal{Q})$  must order the queries  $x \approx y$  so that either  $x$  or  $y$  is free in some preceding conjunct, e.g.,  $B(x) \wedge x \approx y \wedge y \approx z$  is in RANF, but  $B(x) \wedge y \approx z \wedge x \approx y$  is not. In the case of an existentially quantified query  $\exists \vec{v}. Q_{\vec{v}}$ , we rename the variables  $\vec{v}$  to avoid a clash of the free variables in the set of queries  $\mathcal{Q}$  with the bound variables  $\vec{v}$ .

Finally, we resolve the nondeterministic choices in  $\text{sr2ranf}(Q, \mathcal{Q})$  by minimizing the cost of the resulting RANF query with respect to a training database (Section 6.6).

**input:** A safe-range query  $Q \wedge \bigwedge_{\bar{Q} \in \bar{\mathcal{Q}}} \bar{Q}$  such that for all subqueries of the form  $\neg Q'$ ,  $\text{gen}(x, \neg Q')$  does not hold for any variable  $x$ .

**output:** A RANF query  $\hat{Q}$  and a subset of queries  $\bar{\mathcal{Q}} \subseteq \mathcal{Q}$  such that  $Q \wedge \bigwedge_{\bar{Q} \in \bar{\mathcal{Q}}} \bar{Q} \equiv \hat{Q} \wedge \bigwedge_{\bar{Q} \in \bar{\mathcal{Q}}} \bar{Q}$ ; for all  $\mathcal{S}$  and  $\alpha$ ,  $(\mathcal{S}, \alpha) \models \hat{Q} \implies (\mathcal{S}, \alpha) \models \bigwedge_{\bar{Q} \in \bar{\mathcal{Q}}} \bar{Q}$  holds;  $\hat{Q} = \text{cp}(\hat{Q})$ ; and  $\text{fv}(Q) \subseteq \text{fv}(\hat{Q}) \subseteq \text{fv}(Q) \cup \text{fv}(\bar{\mathcal{Q}})$ , unless  $\hat{Q} = \perp$ .

```

1 function sr2ranf(Q, Q) =
2   if ranf(Q) then
3     return (cp(Q), ∅)
4   switch Q do
5     case x ≈ y do
6       return sr2ranf(x ≈ y ∧ ⋀_{Q̄ ∈ Q̄} Q̄, ∅)
7     case ¬Q' do
8       Q̄ ← {Q̄ ⊆ Q | (¬Q') ∧ ⋀_{Q̄ ∈ Q̄} Q̄ is safe range};
9       if Q̄ = ∅ then
10        (Q̂', -) := sr2ranf(Q', ∅);
11        return (cp(¬Q̂'), ∅);
12      else return sr2ranf((¬Q') ∧ ⋀_{Q̄ ∈ Q̄} Q̄, ∅);
13     case Q₁ ∨ Q₂ do
14       Q̄ ← {Q̄ ⊆ Q | ⋁_{Q' ∈ flat^∨(Q)} (Q' ∧ ⋀_{Q̄ ∈ Q̄} Q̄) is safe range};
15       foreach Q' ∈ flat^∨(Q) do (Q̂', -) := sr2ranf(Q' ∧ ⋀_{Q̄ ∈ Q̄} Q̄, ∅);
16       return (cp(⋁_{Q' ∈ flat^∨(Q)} Q̂'), Q̄);
17     case Q₁ ∧ Q₂ do
18       Q⁻ := {Q' ∈ flat^∧(Q) ∪ Q | neg(Q')}; Q⁺ := (flat^∧(Q) ∪ Q) \ Q⁻;
19       Q̃ := {Q' ∈ Q⁺ | eq(Q')}; Q⁺ := Q⁺ \ Q̃;
20       Q̄ := {¬Q' ∈ Q⁻ | eq(Q')}; Q⁻ := Q⁻ \ Q̄;
21       foreach Q' ∈ Q⁺ do (Q̂', Q_{Q'}) := sr2ranf(Q', (Q⁺ ∪ Q̃) \ {Q'});
22       foreach ¬Q' ∈ Q⁻ do (Q̂', -) := sr2ranf(Q', Q⁺ ∪ Q̃);
23       Q̄ ← {Q̄ ⊆ Q⁺ | Q⁺ ⊆ ⋃_{Q' ∈ Q̄} (Q_{Q'} ∪ {Q'})};
24       return (cp(sort^∧(⋃_{Q' ∈ Q̄} {Q̂'}) ∪ Q̃ ∪ ⋃_{¬Q' ∈ Q⁻} {¬Q̂'} ∪ Q̄), ⋃_{Q' ∈ Q̄} (Q_{Q'} ∩ Q));
25     case ∃v̄. Q_{v̄} do
26       if fv(Q) ∩ v̄ ≠ ∅ then w̄ ← {w̄ | |w̄| = |v̄| and ((fv(Q_{v̄}) \ v̄) ∪ fv(Q)) ∩ w̄ = ∅};
27       else w̄ := v̄;
28       Q_{w̄} := Q_{v̄}[v̄ ↦ w̄];
29       Q̄ ← {Q̄ ⊆ Q | Q_{w̄} ∧ ⋀_{Q̄ ∈ Q̄} Q̄ is safe range};
30       (Q̂_{w̄}, -) := sr2ranf(Q_{w̄} ∧ ⋀_{Q̄ ∈ Q̄} Q̄, ∅);
31       return (cp(∃w̄. Q̂_{w̄}), Q̄);
32   otherwise do return (cp(Q), ∅);

```

Figure 12: Translation of a safe-range query in SRNF to RANF.

**6.3. Optimization using Count Aggregations.** In this section, we introduce count aggregations and describe a generalization of Claußen et al. [CKMP97]’s approach to evaluate RANF queries using count aggregations. Consider the query

$$Q_x \wedge \neg \exists y. (Q_x \wedge Q_y \wedge \neg Q_{xy}),$$

where  $\text{fv}(Q_x) = \{x\}$ ,  $\text{fv}(Q_y) = \{y\}$ , and  $\text{fv}(Q_{xy}) = \{x, y\}$ . This query is obtained by applying our translation to the query  $Q_x \wedge \forall y. (Q_y \longrightarrow Q_{xy})$ . The cost of the translated query is dominated by the cost of the Cartesian product  $Q_x \wedge Q_y$ . Consider the subquery  $Q' := \exists y. (Q_x \wedge Q_y \wedge \neg Q_{xy})$ . A assignment  $\alpha$  satisfies  $Q'$  iff  $\alpha$  satisfies  $Q_x$  and there exists a value  $d$  such that  $\alpha[y \mapsto d]$  satisfies  $Q_y$ , but not  $Q_{xy}$ , i.e., the number of values  $d$  such that  $\alpha[y \mapsto d]$  satisfies  $Q_y$  is not equal to the number of values  $d$  such that  $\alpha[y \mapsto d]$  satisfies both  $Q_y$  and  $Q_{xy}$ . An alternative evaluation of  $Q'$  evaluates the queries  $Q_x$ ,  $Q_y$ ,  $Q_y \wedge Q_{xy}$  and computes the numbers of values  $d$  such that  $\alpha[y \mapsto d]$  satisfies  $Q_y$  and  $Q_y \wedge Q_{xy}$ , respectively, i.e., computes count aggregations. These count aggregations are then used to filter assignments  $\alpha$  satisfying  $Q_x$  to get assignments  $\alpha$  satisfying  $Q'$ . The asymptotic time complexity of the alternative evaluation never exceeds that of the evaluation computing the Cartesian product  $Q_x \wedge Q_y$  and asymptotically improves it if  $\|Q_x\| + \|Q_y\| + \|Q_{xy}\| \ll \|Q_x \wedge Q_y\|$ . Furthermore, we observe that a assignment  $\alpha$  satisfies  $Q_x \wedge \neg Q'$  if  $\alpha$  satisfies  $Q_x$ , but not  $Q'$ , i.e., the number of values  $d$  such that  $\alpha[y \mapsto d]$  satisfies  $Q_y$  is equal to the number of values  $d$  such that  $\alpha[y \mapsto d]$  satisfies  $Q_y \wedge Q_{xy}$ .

Next we introduce the syntax and semantics of count aggregations. We extend RC’s syntax by  $[\text{CNT } \vec{v}. Q_{\vec{v}}](c)$ , where  $Q$  is a query,  $c$  is a variable representing the result of the count aggregation, and  $\vec{v}$  is a sequence of variables that are bound by the aggregation operator. The semantics of the count aggregation is defined as follows:

$$(\mathcal{S}, \alpha) \models [\text{CNT } \vec{v}. Q_{\vec{v}}](c) \text{ iff } (M = \emptyset \longrightarrow \text{fv}(Q) \subseteq \vec{v}) \text{ and } \alpha(c) = |M|,$$

where  $M = \{\vec{d} \in \mathcal{D}^{|\vec{v}|} \mid (\mathcal{S}, \alpha[\vec{v} \mapsto \vec{d}]) \models Q\}$ . We use the condition  $M = \emptyset \longrightarrow \text{fv}(Q) \subseteq \vec{v}$  instead of  $M \neq \emptyset$  to set  $c$  to a zero count if the group  $M$  is empty and there are no group-by variables (like in SQL). The set of free variables in a count aggregation is  $\text{fv}([\text{CNT } \vec{v}. Q_{\vec{v}}](c)) = (\text{fv}(Q) \setminus \vec{v}) \cup \{c\}$ . Finally, we extend the definition of  $\text{ranf}(Q)$  with the case of a count aggregation:

$$\text{ranf}([\text{CNT } \vec{v}. Q_{\vec{v}}](c)) \text{ iff } \text{ranf}(Q) \text{ and } \vec{v} \subseteq \text{fv}(Q) \text{ and } c \notin \text{fv}(Q).$$

We formulate translations introducing count aggregations in the following two lemmas.

**Lemma 6.2.** *Given  $Q \neq \emptyset$ , let  $\exists \vec{v}. Q_{\vec{v}} \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \neg \bar{Q}$  be a RANF query. Let  $c, c'$  be fresh variables that do not occur in  $\text{fv}(Q_{\vec{v}})$ . Then*

$$\begin{aligned} (\exists \vec{v}. Q_{\vec{v}} \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \neg \bar{Q}) &\equiv ((\exists \vec{v}. Q_{\vec{v}}) \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \neg (\exists \vec{v}. Q_{\vec{v}} \wedge \bar{Q})) \vee \\ &(\exists c, c'. [\text{CNT } \vec{v}. Q_{\vec{v}}](c) \wedge \\ &[\text{CNT } \vec{v}. \bigvee_{\bar{Q} \in \mathcal{Q}} (Q_{\vec{v}} \wedge \bar{Q})](c') \wedge \neg(c = c')). \end{aligned} \quad (\#)$$

Moreover, the right-hand side of (#) is in RANF.



**Lemma 6.3.** *Given  $\mathcal{Q} \neq \emptyset$ , let  $\hat{Q} \wedge \neg \exists \vec{v}. Q_{\vec{v}} \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \neg \overline{Q}$ , be a RANF query. Let  $c, c'$  be fresh variables that do not occur in  $\text{fv}(\hat{Q}) \cup \text{fv}(Q_{\vec{v}})$ . Then*

$$\begin{aligned} (\hat{Q} \wedge \neg \exists \vec{v}. Q_{\vec{v}} \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \neg \overline{Q}) &\equiv (\hat{Q} \wedge \neg(\exists \vec{v}. Q_{\vec{v}})) \vee \\ &(\exists c, c'. \hat{Q} \wedge [\text{CNT } \vec{v}. Q_{\vec{v}}](c) \wedge \\ &[\text{CNT } \vec{v}. \bigvee_{\overline{Q} \in \mathcal{Q}} (Q_{\vec{v}} \wedge \overline{Q})](c') \wedge (c = c')). \end{aligned} \quad (\#\#)$$

Moreover, the right-hand side of  $(\#\#)$  is in RANF.

Note that the query cost does not decrease after applying the translation  $(\#)$  or  $(\#\#)$  because of the subquery  $[\text{CNT } \vec{v}. Q_{\vec{v}}](c)$  in which  $Q_{\vec{v}}$  is evaluated before the count aggregation is computed. For the query  $\exists y. ((Q_x \wedge Q_y) \wedge \neg Q_{xy})$  from before, we would compute  $[\text{CNT } y. Q_x \wedge Q_y](c)$ , i.e., we would not (yet) avoid computing the Cartesian product  $Q_x \wedge Q_y$ . However, we could reduce the scope of the bound variable  $y$  by further translating

$$[\text{CNT } y. Q_x \wedge Q_y](c) \equiv Q_x \wedge [\text{CNT } y. Q_y](c).$$

This technique, called *mini-scoping*, can be applied to a count aggregation  $[\text{CNT } \vec{v}. Q_{\vec{v}}](c)$  if the aggregated query  $Q_{\vec{v}}$  is a conjunction that can be split into two RANF conjuncts and the variables  $\vec{v}$  do not occur free in one of the conjuncts (that conjunct can be pulled out of the count aggregation). Mini-scoping can be analogously applied to queries of the form  $\exists \vec{v}. Q_{\vec{v}}$ .

Moreover, we can split a count aggregation over a conjunction  $Q_{\vec{v}} \wedge Q'_{\vec{v}}$  into a product of count aggregations if the conjunction can be split into two RANF conjuncts with disjoint sets of bound variables, i.e.,  $\vec{v} \cap \text{fv}(Q_{\vec{v}}) \cap \text{fv}(Q'_{\vec{v}}) = \emptyset$ :

$$[\text{CNT } \vec{v}. Q_{\vec{v}} \wedge Q'_{\vec{v}}](c) \equiv (\exists c_1, c_2. [\text{CNT } \vec{v} \cap \text{fv}(Q_{\vec{v}}). Q_{\vec{v}}](c_1) \wedge [\text{CNT } \vec{v} \cap \text{fv}(Q'_{\vec{v}}). Q'_{\vec{v}}](c_2) \wedge c = c_1 \cdot c_2).$$

Here  $c_1$  and  $c_2$  are fresh variables that do not occur in  $\text{fv}(Q_{\vec{v}}) \cup \text{fv}(Q'_{\vec{v}}) \cup \{c\}$ . Note that mini-scoping is only a heuristic and it can both improve and harm the time complexity of query evaluation. We leave the application of other more general optimization algorithms [KNR16, OS16]) as future work.

We implement the translations from Lemmas 6.2 and 6.3 and mini-scoping in the function  $\text{optcnt}(\cdot)$ . Given a RANF query  $\hat{Q}$ ,  $\text{optcnt}(\hat{Q})$  is an equivalent RANF query after introducing count aggregations and performing mini-scoping. The function  $\text{optcnt}(\hat{Q})$  uses a training database to decide how to apply the translations from Lemmas 6.2 and 6.3 and mini-scoping. More specifically, the function  $\text{optcnt}(\hat{Q})$  tries several possibilities and chooses one that minimizes the query cost of the resulting RANF query.

**Example 6.4.** *We show how to introduce count aggregations into the RANF query*

$$\hat{Q} := Q_x \wedge \neg \exists y. (Q_x \wedge Q_y \wedge \neg Q_{xy}).$$

After applying the translation  $(\#\#)$  and mini-scoping to this query, we obtain the following equivalent RANF query:

$$\begin{aligned} \text{optcnt}(\hat{Q}) &:= (Q_x \wedge \neg(Q_x \wedge \exists y. Q_y)) \vee \\ &(\exists c, c'. Q_x \wedge [\text{CNT } y. Q_y](c) \wedge [\text{CNT } y. Q_y \wedge Q_{xy}](c') \wedge (c = c')). \end{aligned}$$

**6.4. Translating RANF to RA.** Our translation of a RANF query into SQL has two steps: we first translate the query to an equivalent RA expression, which we then translate to SQL using a publicly available RA interpreter `radb` [Yan19].

We define the function  $\text{ranf2ra}(\hat{Q})$  translating RANF queries  $\hat{Q}$  into equivalent RA expressions  $\text{ranf2ra}(\hat{Q})$ . The translation is based on Algorithm 5.4.8 by Abiteboul et al. [AHV95], which we modify as follows. We adjust the way closed RC queries are handled. Chomicki and Toman [CT95] observed that closed RC queries cannot be handled by SQL, since SQL allows neither empty projections nor 0-ary relations. They propose to use a unary auxiliary predicate  $A \in \mathcal{R}$  whose interpretation  $A^S = \{t\}$  always contains exactly one tuple  $t$ . Every closed query  $\exists x. Q_x$  is then translated into  $\exists x. A(t) \wedge Q_x$  with an auxiliary free variable  $t$ . Every other closed query  $\hat{Q}$  is translated into  $A(t) \wedge \hat{Q}$ , e.g.,  $B(42)$  is translated into  $A(t) \wedge B(42)$ . We also use the auxiliary predicate  $A$  to translate queries of the form  $x \approx c$  and  $c \approx x$  because the single tuple  $(t)$  in  $A^S$  can be mapped to any constant  $c$ . Finally, we extend [AHV95, Algorithm 5.4.8] with queries of the form [CNT  $\vec{v}. Q_{\vec{v}}(c)$ ].

**6.5. Translating RA to SQL.** The `radb` interpreter, abbreviated here by the function  $\text{ra2sql}(\cdot)$ , translates an RA expression into SQL by simply mapping the RA connectives into their SQL counterparts. The function  $\text{ra2sql}(\cdot)$  is primitive recursive on RA expressions. We modify `radb` to further improve performance of the query evaluation as follows.

A RANF query  $Q_1 \wedge \neg Q_2$ , where  $\text{ranf}(Q_1)$ ,  $\text{ranf}(Q_2)$ , and  $\text{fv}(Q_2) \subseteq \text{fv}(Q_1)$  is translated into RA expression  $\text{ranf2ra}(Q_1) \triangleright \text{ranf2ra}(Q_2)$ , where  $\triangleright$  denotes the anti-join operator and  $\text{ranf2ra}(Q_1)$ ,  $\text{ranf2ra}(Q_2)$  are the equivalent relational algebra expressions for  $Q_1$ ,  $Q_2$ , respectively. The `radb` interpreter only supports the anti-join operator  $\text{ranf2ra}(Q_1) \triangleright \text{ranf2ra}(Q_2)$  expressed as  $\text{ranf2ra}(Q_1) - (\text{ranf2ra}(Q_1) \bowtie \text{ranf2ra}(Q_2))$ , where  $-$  denotes the set difference operator and  $\bowtie$  denotes the natural join. Alternatively, the anti-join operator can be directly mapped to `LEFT JOIN` in SQL. We generalize `radb` to use `LEFT JOIN` since it performs better in our empirical evaluation [RBKT22a].

The `radb` interpreter introduces a separate SQL subquery in a `WITH` clause for every subexpression in the RA expression. We extend `radb` to additionally perform common subquery elimination, i.e., to merge syntactically equal subqueries. Common subquery elimination is also assumed in our query cost (Section 3.5).

Finally, the function  $\text{ranf2sql}(\hat{Q})$  is defined as  $\text{ranf2sql}(\hat{Q}) := \text{ra2sql}(\text{ranf2ra}(\hat{Q}))$ , i.e., as a composition of the two translations from RANF to RA and from RA to SQL.

**6.6. Resolving Nondeterministic Choices.** To resolve the nondeterministic choices in our algorithms, we suppose that the algorithms have access to a *training database*  $\mathcal{T}$  of constant size. The training database is used to compare the cost of queries over the actual database and thus it should preserve the relative ordering of queries by their cost over the actual database as much as possible. Still, our translation satisfies the correctness and worst-case complexity claims (Section 4.3 and 5) for every choice of the training database. The training databases used in our empirical evaluation are obtained using the function `dg` (Section 7) with  $|\mathcal{T}^+| = |\mathcal{T}^-| = 2$ . Because of its constant size, the complexity of evaluating a query over the training database is constant and does not impact the asymptotic time complexity of evaluating the query over the actual database using our translation. There are two types of nondeterministic choices in our algorithms:

**input:** An RC query  $Q$  satisfying CON, CST, VAR, REP,  $\gamma \in \{0, 1\}$ .  
**output:** Two sets of variables  $\mathcal{V}^+$  and  $\mathcal{V}^-$  whose values must be equal in every tuple in  $\mathcal{T}_{\bar{v}}^+$  and  $\mathcal{T}_{\bar{v}}^-$  when computing a Data Golf structure.

```

1 function  $\text{dg}^{\approx}(Q, \gamma) =$ 
2   switch  $Q$  do
3     case  $r(t_1, \dots, t_{l(r)})$  do return  $(\emptyset, \emptyset)$ ;
4     case  $x \approx y$  do return  $(\{x, y\}, \emptyset)$ ;
5     case  $\neg Q'$  do
6        $(\mathcal{V}^+, \mathcal{V}^-) := \text{dg}^{\approx}(Q', \gamma)$ ;
7       return  $(\mathcal{V}^-, \mathcal{V}^+)$ ;
8     case  $Q'_1 \vee Q'_2$  or  $Q'_1 \wedge Q'_2$  do
9        $(\mathcal{V}_1^+, \mathcal{V}_1^-) := \text{dg}^{\approx}(Q'_1, \gamma)$ ;
10       $(\mathcal{V}_2^+, \mathcal{V}_2^-) := \text{dg}^{\approx}(Q'_2, \gamma)$ ;
11      if  $\gamma = 0$  then return  $(\mathcal{V}_1^+ \cup \mathcal{V}_2^+, \mathcal{V}_1^- \cup \mathcal{V}_2^-)$ ;
12      else if  $Q = Q'_1 \vee Q'_2$  then return  $(\mathcal{V}_1^+ \cup \mathcal{V}_2^-, \mathcal{V}_1^- \cup \mathcal{V}_2^-)$ ;
13      else if  $Q = Q'_1 \wedge Q'_2$  then return  $(\mathcal{V}_1^+ \cup \mathcal{V}_2^+, \mathcal{V}_1^+ \cup \mathcal{V}_2^-)$ ;
14     case  $\exists y. Q_y$  do return  $\text{dg}^{\approx}(Q_y, \gamma)$ ;

```

Figure 13: Computing sets of variables  $\mathcal{V}^+$  and  $\mathcal{V}^-$  to reflect equalities in a query when computing a Data Golf structure.

- Choosing some  $X \leftarrow \mathcal{X}$  in a while-loop. As the while-loops always update  $\mathcal{X}$  with  $\mathcal{X} := (\mathcal{X} \setminus \{X\}) \cup f(X)$  for some  $f$ , the order in which the elements of  $\mathcal{X}$  are chosen does not matter.
- Choosing a subset of queries  $\bar{Q} \subseteq Q$  in the function  $\text{sr2ranf}(Q, Q)$ . Because  $\text{sr2ranf}(Q, Q)$  yields a RANF query, we enumerate all *minimal* subsets (a subset  $\bar{Q} \subseteq Q$  is minimal if there exists no proper subset  $\bar{Q}' \subsetneq \bar{Q}$  that could be used instead of  $\bar{Q}$ ) and choose one that minimizes the query cost of the RANF query.
- Choosing a variable  $x \in V$  and a set  $\mathcal{G}$  such that  $\text{cov}(x, \tilde{Q}, \mathcal{G})$ , where  $\tilde{Q}$  is a query with range-restricted bound variables and  $V \subseteq \text{fv}(\tilde{Q})$ . Observe that the measure  $\text{m}(Q)$  on queries, defined in Figure 11, decreases for the queries in the premises of the rules for  $\text{gen}(x, \tilde{Q}, \mathcal{G})$  and  $\text{cov}(x, \tilde{Q}, \mathcal{G})$ , defined in Figures 4 and 5. Hence, deriving  $\text{gen}(x, \tilde{Q}, \mathcal{G})$  and  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  either succeeds or gets stuck after at most  $\text{m}(\tilde{Q})$  steps. In particular, we can enumerate all sets  $\mathcal{G}$  such that  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  holds. Because we derive one additional query  $\tilde{Q}[x \mapsto y]$  for every  $y \in \text{eqs}(x, \mathcal{G})$  and a single query  $\tilde{Q} \wedge \text{qps}^{\vee}(\mathcal{G})$ , we choose  $x \in V$  and  $\mathcal{G}$  minimizing  $|\text{eqs}(x, \mathcal{G})|$  as the first objective and  $\sum_{Q_{qp} \in \text{qps}(\mathcal{G})} \text{cost}^T(Q_{qp})$  as the second objective. Our particular choice of  $\mathcal{G}$  with  $\text{cov}(x, \tilde{Q}, \mathcal{G})$  is merely a heuristic and does not provide any additional guarantees compared to every other choice of  $\mathcal{G}$  with  $\text{cov}(x, \tilde{Q}, \mathcal{G})$ . This process can be further improved by adopting query plan search heuristics as is done by most of the existing database management systems.

## 7. DATA GOLF BENCHMARK

In this section, we describe the *Data Golf* benchmark, which we use to generate structures (i.e., database instances) for our empirical evaluation. The technical description of

**input:** An RC query  $Q$  satisfying CON, CST, VAR, REP, a sequence of pairwise distinct variables  $\vec{v}$ ,  $\text{av}(Q) \subseteq \vec{v}$ , sets of tuples  $\mathcal{T}_{\vec{v}}^+$  and  $\mathcal{T}_{\vec{v}}^-$  over  $\vec{v}$  such that all values of variables from  $\text{av}(Q)$  in these tuples are pairwise distinct (also across tuples) except that, in every tuple in  $\mathcal{T}_{\vec{v}}^+$  ( $\mathcal{T}_{\vec{v}}^-$ ), the variables in  $\mathcal{V}^+$  ( $\mathcal{V}^-$ ) have the same value (which is different across tuples), where  $\text{dg}^{\approx}(Q, \gamma) = (\mathcal{V}^+, \mathcal{V}^-)$ ,  $\gamma \in \{0, 1\}$ .

**output:** A structure  $\mathcal{S}$  such that  $\mathcal{T}_{\vec{v}}^+[\vec{fv}(Q)] \subseteq \llbracket Q \rrbracket$ ,  $\mathcal{T}_{\vec{v}}^-[\vec{fv}(Q)] \cap \llbracket Q \rrbracket = \emptyset$ , and  $\llbracket Q' \rrbracket$  and  $\llbracket \neg Q' \rrbracket$  contain at least  $\min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$  tuples, for every  $Q' \sqsubseteq Q$ .

```

1 function dg( $Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma$ ) =
2   switch  $Q$  do
3     case  $r(t_1, \dots, t_{i(r)})$  do return  $\{r^{\mathcal{S}} \mapsto \mathcal{T}_{\vec{v}}^+[t_1, \dots, t_{i(r)}]\}$ ;
4     case  $x \approx y$  do return  $\emptyset$ ;
5     case  $\neg Q'$  do return  $\text{dg}(Q', \vec{v}, \mathcal{T}_{\vec{v}}^-, \mathcal{T}_{\vec{v}}^+, \gamma)$ ;
6     case  $Q'_1 \vee Q'_2$  or  $Q'_1 \wedge Q'_2$  do
7        $(\mathcal{V}_1^+, \mathcal{V}_1^-) := \text{dg}^{\approx}(Q'_1, \gamma)$ ;  $(\mathcal{V}_2^+, \mathcal{V}_2^-) := \text{dg}^{\approx}(Q'_2, \gamma)$ ;
8       if  $\gamma = 0$  then  $(\mathcal{V}^1, \mathcal{V}^2) := (\mathcal{V}_1^+ \cup \mathcal{V}_2^-, \mathcal{V}_1^- \cup \mathcal{V}_2^-)$ ;
9       else if  $Q = Q'_1 \wedge Q'_2$  then  $(\mathcal{V}^1, \mathcal{V}^2) := (\mathcal{V}_1^- \cup \mathcal{V}_2^-, \mathcal{V}_1^- \cup \mathcal{V}_2^-)$ ;
10      else if  $Q = Q'_1 \vee Q'_2$  then  $(\mathcal{V}^1, \mathcal{V}^2) := (\mathcal{V}_1^+ \cup \mathcal{V}_2^+, \mathcal{V}_1^- \cup \mathcal{V}_2^-)$ ;
11       $(\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2) \leftarrow \{(\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2) \mid |\mathcal{T}_{\vec{v}}^1| = |\mathcal{T}_{\vec{v}}^2| = \min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$ , all values in tuples in
         $\mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2$  are chosen to be pairwise distinct (also across tuples) except that,
        in every tuple in  $\mathcal{T}_{\vec{v}}^1$  ( $\mathcal{T}_{\vec{v}}^2$ ), the variables in  $\mathcal{V}^1$  ( $\mathcal{V}^2$ ) have the same value
        (which is different across tuples)\};
12      if  $\gamma = 0$  then return
         $\text{dg}(Q'_1, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^2, \gamma) \cup \text{dg}(Q'_2, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^2, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^1, \gamma)$ ;
13      else if  $Q = Q'_1 \vee Q'_2$  then return
         $\text{dg}(Q'_1, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^2, \gamma) \cup \text{dg}(Q'_2, \vec{v}, \mathcal{T}_{\vec{v}}^1 \cup \mathcal{T}_{\vec{v}}^2, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^+, \gamma)$ ;
14      else if  $Q = Q'_1 \wedge Q'_2$  then return
         $\text{dg}(Q'_1, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^-, \mathcal{T}_{\vec{v}}^1 \cup \mathcal{T}_{\vec{v}}^2, \gamma) \cup \text{dg}(Q'_2, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^2, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^1, \gamma)$ ;
15      case  $\exists y. Q_y$  do return  $\text{dg}(Q_y, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ ;

```

Figure 14: Computing the Data Golf structure.

this benchmark is only needed to fully understand the setup of our empirical evaluation (Section 8), but its details are independent of our query translation (Section 4–6).

Given an RC query, we seek a structure that yields a nontrivial evaluation result for the overall query and for all its subqueries. Intuitively, the structure makes query evaluation potentially more challenging compared to the case where some subquery evaluates to a trivial (e.g., empty) result. More specifically, Data Golf has two objectives. The first resembles the *regex golf* game’s objective [Ell13] (hence the name) and aims to find a structure on which the result of a given query contains a given *positive* set of tuples and does not contain any tuples from another given *negative* set. The second objective is to ensure that all the query’s subqueries evaluate to a non-trivial result.

Formally, given a query  $Q$  and two sets of tuples  $\mathcal{T}^+$  and  $\mathcal{T}^-$  over a fixed domain  $\mathcal{D}$ , representing assignments of  $\text{av}(Q)$  and satisfying further assumptions on their values, Data Golf

produces a structure  $\mathcal{S}$  (represented as a partial mapping from predicate symbols to their interpretations) such that the projections of tuples in  $\mathcal{T}^+$  ( $\mathcal{T}^-$ , respectively) to  $\vec{fv}(Q)$  are in  $\llbracket Q \rrbracket$  (disjoint from  $\llbracket Q \rrbracket$ , respectively) and  $\llbracket \llbracket Q' \rrbracket \rrbracket$  and  $\llbracket \llbracket \neg Q' \rrbracket \rrbracket$  are at least  $\min\{|\mathcal{T}^+|, |\mathcal{T}^-|\}$ , for every  $Q' \sqsubseteq Q$ . To be able to produce such a structure  $\mathcal{S}$ , we make the following assumptions on  $Q$ :

- (CON) for every subquery  $\exists y. Q_y$  of  $Q$  we have  $\text{con}_{\text{vgt}}(y, Q_y, \mathcal{A})$  (Figure 17) for some set of atomic predicates  $\mathcal{A}$  and, moreover,  $\{y\} \subsetneq \text{fv}(Q_{ap})$  holds for every  $Q_{ap} \in \mathcal{A}$ ; these conditions prevent subqueries like  $\exists y. \neg P_2(x, y)$  and  $\exists y. (P_2(x, y) \vee P_1(y))$ , respectively;
- (CST)  $Q$  contains no subquery of the form  $x \approx c$ , which is satisfied by exactly one tuple;
- (VAR)  $Q$  contains no closed subqueries, e.g.,  $P_1(42)$ , because a closed subquery is either satisfied by all possible tuples or no tuple at all; and
- (REP)  $Q$  contains no repeated predicate symbols and no equalities  $x \approx y$  in  $Q$  share a variable; this avoids subqueries like  $P_1(x) \wedge \neg P_1(x)$  and  $x \approx y \wedge \neg x \approx y$ .

Given a sequence of pairwise distinct variables  $\vec{v}$  and a tuple  $\vec{d}$  of the same length, we may interpret the tuple  $\vec{d}$  as a *tuple over*  $\vec{v}$ , denoted as  $\vec{d}(\vec{v})$ . Given a sequence  $t_1, \dots, t_k \in \vec{v} \cup \mathcal{C}$  of terms, we denote by  $\vec{d}(\vec{v})[t_1, \dots, t_k]$  the tuple obtained by evaluating the terms  $t_1, \dots, t_k$  over  $\vec{d}(\vec{v})$ . Formally, we define  $\vec{d}(\vec{v})[t_1, \dots, t_k] := (d'_i)_{i=1}^k$ , where  $d'_i = \vec{d}_j$  if  $t_i = \vec{v}_j$  and  $d'_i = t_i$  if  $t_i \in \mathcal{C}$ . We lift this notion to sets of tuples over  $\vec{v}$  in the standard way.

Data Golf is formalized by the function  $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ , defined in Figure 14, where  $\vec{v}$  is a sequence of pairwise distinct variables containing all variables in  $Q$ , i.e.,  $\text{av}(Q) \subseteq \vec{v}$ ,  $\mathcal{T}_{\vec{v}}^+$  and  $\mathcal{T}_{\vec{v}}^-$  are sets of tuples over  $\vec{v}$ , and  $\gamma \in \{0, 1\}$  is a *strategy*. To reflect  $Q$ 's equalities in the sets  $\mathcal{T}_{\vec{v}}^+$  and  $\mathcal{T}_{\vec{v}}^-$ , given a strategy  $\gamma$ , we define the function  $\text{dg}^{\approx}(Q, \gamma) = (\mathcal{V}^+, \mathcal{V}^-)$  (Figure 13) that computes two sets of variables  $\mathcal{V}^+$  and  $\mathcal{V}^-$  whose values must be equal in every tuple in  $\mathcal{T}_{\vec{v}}^+$  and  $\mathcal{T}_{\vec{v}}^-$ , respectively. The values of the remaining variables ( $\vec{v} \setminus \mathcal{V}^+$  and  $\vec{v} \setminus \mathcal{V}^-$ , respectively) must be pairwise distinct and also different from the value of the variables in  $\mathcal{V}^+$  and  $\mathcal{V}^-$ , respectively. In the case of a conjunction or a disjunction, we add disjoint sets  $\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2$  of tuples over  $\vec{v}$  to  $\mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-$  so that the intermediate results for the subqueries are neither equal nor disjoint. We implement two strategies (parameter  $\gamma$ ) to choose these sets  $\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2$ .

Let  $\mathcal{S}$  be a Data Golf structure computed by  $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ . We justify why  $\mathcal{S}$  satisfies  $\mathcal{T}_{\vec{v}}^+[\vec{fv}(Q)] \subseteq \llbracket Q \rrbracket$  and  $\mathcal{T}_{\vec{v}}^-[\vec{fv}(Q)] \cap \llbracket Q \rrbracket = \emptyset$ . We proceed by induction on the query  $Q$ . Because of (REP), the Data Golf structures for the subqueries  $Q_1, Q_2$  of a binary query  $Q_1 \vee Q_2$  or  $Q_1 \wedge Q_2$  can be combined using the union operator. The only case that does not follow immediately is that  $\mathcal{T}_{\vec{v}}^-[\vec{fv}(Q)] \cap \llbracket Q \rrbracket = \emptyset$  for a query  $Q$  of the form  $\exists y. Q_y$ . We prove this case by contradiction. Without loss of generality we assume that  $\vec{fv}(Q_y) = \vec{fv}(Q) \cdot y$ . Suppose that  $\vec{d} \in \mathcal{T}_{\vec{v}}^-[\vec{fv}(Q)]$  and  $\vec{d} \in \llbracket Q \rrbracket$ . Because  $\vec{d} \in \mathcal{T}_{\vec{v}}^-[\vec{fv}(Q)]$ , there exists some  $d$  such that  $\vec{d} \cdot d \in \mathcal{T}_{\vec{v}}^-[\vec{fv}(Q_y)]$ . Because  $\vec{d} \in \llbracket Q \rrbracket$ , there exists some  $d'$  such that  $\vec{d} \cdot d' \in \llbracket Q_y \rrbracket$ . By the induction hypothesis,  $\vec{d} \cdot d \notin \llbracket Q_y \rrbracket$  and  $\vec{d} \cdot d' \notin \mathcal{T}_{\vec{v}}^-[\vec{fv}(Q_y)]$ . Because  $\text{con}_{\text{vgt}}(y, Q_y, \mathcal{A})$  holds for some  $\mathcal{A}$  satisfying (CON), the query  $Q_y$  is equivalent to  $(Q_y \wedge \bigvee_{Q_{ap} \in \mathcal{A}} Q_{ap}) \vee Q_y[y/\perp]$ . We have  $\vec{d} \cdot d' \in \llbracket Q_y \rrbracket$ . If the tuple  $\vec{d} \cdot d'$  satisfies  $Q_y[y/\perp]$ , then  $\vec{d} \cdot d \in \llbracket Q_y \rrbracket$  (contradiction) because the variable  $y$  does not occur in the query  $Q_y[y/\perp]$  and thus its assignment in  $\vec{d} \cdot d'$  can be arbitrarily changed. Otherwise, the tuple  $\vec{d} \cdot d'$  satisfies some atomic predicate  $Q_{ap} \in \mathcal{A}$  and (CON) implies  $\{y\} \subsetneq \text{fv}(Q_{ap})$ . Hence, the tuples  $\vec{d} \cdot d$  and  $\vec{d} \cdot d'$  agree on the assignment of a variable  $x \in \text{fv}(Q_{ap}) \setminus \{y\}$ . Let  $\overline{\mathcal{T}}_{\vec{v}}^+$  and  $\overline{\mathcal{T}}_{\vec{v}}^-$  be the sets in the recursive call of  $\text{dg}$  on the atomic predicate from  $Q_{ap}$ . Because  $\vec{d} \cdot d \in \mathcal{T}_{\vec{v}}^-[\vec{fv}(Q_y)]$  and  $\mathcal{T}_{\vec{v}}^-[\vec{fv}(Q_y)] \subseteq \overline{\mathcal{T}}_{\vec{v}}^+[\vec{fv}(Q_y)] \cup \overline{\mathcal{T}}_{\vec{v}}^-[\vec{fv}(Q_y)]$ ,

the tuple  $\vec{d} \cdot d$  is in  $\overline{\mathcal{T}}_{\vec{v}}^+[\vec{\text{fv}}(Q_y)] \cup \overline{\mathcal{T}}_{\vec{v}}^-[\vec{\text{fv}}(Q_y)]$ . Because  $\vec{d} \cdot d'$  satisfies the quantified predicate  $Q_{ap}$ , the tuple  $\vec{d} \cdot d'$  is in  $\overline{\mathcal{T}}_{\vec{v}}^+[\vec{\text{fv}}(Q_y)]$ . Next we observe that the assignments of every variable (in particular,  $x$ ) in the tuples from the sets  $\overline{\mathcal{T}}_{\vec{v}}^+$ ,  $\overline{\mathcal{T}}_{\vec{v}}^-$  are pairwise distinct (there can only be equal values of variables within a single tuple). Because the tuples  $\vec{d} \cdot d$  and  $\vec{d} \cdot d'$  agree on the assignment of  $x$ , they must be equal, i.e.,  $\vec{d} \cdot d = \vec{d} \cdot d'$  (contradiction).

The sets  $\mathcal{T}_{\vec{v}}^+$ ,  $\mathcal{T}_{\vec{v}}^-$  only grow in  $\text{dg}$ 's recursion and the properties (CON), (CST), (VAR), and (REP) imply that  $Q$  has no closed subquery. Hence,  $\mathcal{T}_{\vec{v}}^+[\vec{\text{fv}}(Q)] \subseteq \llbracket Q \rrbracket$  and  $\mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$  imply that  $\llbracket Q' \rrbracket$  and  $\llbracket \neg Q' \rrbracket$  contain at least  $\min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$  tuples, for every  $Q' \sqsubseteq Q$ .

**Example 7.1.** Consider the query  $Q := \neg \exists y. P_2(x, y) \wedge \neg P_3(x, y, z)$ . This query  $Q$  satisfies (CON), (CST), (VAR), and (REP). In particular,  $\text{con}_{\text{vgt}}(y, P_2(x, y) \wedge \neg P_3(x, y, z), \mathcal{A})$  holds for  $\mathcal{A} = \{P_2(x, y)\}$  with  $\{y\} \subsetneq \text{fv}(P_2(x, y))$ . We choose  $\vec{v} = (x, z, y)$ ,  $\mathcal{T}_{\vec{v}}^+ = \{(0, 4, 8), (2, 6, 10)\}$ , and  $\mathcal{T}_{\vec{v}}^- = \{(12, 16, 20), (14, 18, 22)\}$ . The function  $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$  first flips  $\mathcal{T}_{\vec{v}}^+$  and  $\mathcal{T}_{\vec{v}}^-$  because  $Q$ 's main connective is negation.

For conjunction (a binary operator), two additional sets of tuples are computed:  $\mathcal{T}_{\vec{v}}^1 = \{(24, 28, 32), (26, 30, 34)\}$  and  $\mathcal{T}_{\vec{v}}^2 = \{(36, 40, 44), (38, 42, 46)\}$ . Depending on the strategy ( $\gamma = 0$  or  $\gamma = 1$ ), one of the following structures is computed:  $\mathcal{S}_0 = \{P_2 \mapsto \{(12, 20), (14, 22), (24, 32), (26, 34)\}, P_3 \mapsto \mathcal{T}_{xyz}^+\}$ , or  $\mathcal{S}_1 = \{P_2 \mapsto \{(12, 20), (14, 22), (0, 8), (2, 10)\}, P_3 \mapsto \mathcal{T}_{xyz}^+\}$ , where  $\mathcal{T}_{xyz}^+ = \{(0, 8, 4), (2, 10, 6), (24, 32, 28), (26, 34, 30)\}$ .

The query  $P_1(x) \wedge Q$  is satisfied by the finite set of tuples  $\mathcal{T}_{\vec{v}}^+[x, z]$  under the structure  $\mathcal{S}_1 \cup \{P_1 \mapsto \{(0), (2)\}\}$  obtained by extending  $\mathcal{S}_1$  ( $\gamma = 1$ ). In contrast, the same query  $P_1(x) \wedge Q$  is satisfied by an infinite set of tuples including  $\mathcal{T}_{\vec{v}}^+[x, z]$  and disjoint from  $\mathcal{T}_{\vec{v}}^-[x, z]$  under the structure  $\mathcal{S}_0 \cup \{P_1 \mapsto \{(0), (2)\}\}$  obtained by extending  $\mathcal{S}_0$  ( $\gamma = 0$ ).

## 8. EMPIRICAL EVALUATION

We empirically validate the evaluation performance of the queries output by RC2SQL. We also assess RC2SQL's translation time, the average-case time complexity of query evaluation, scalability to large databases, and DBMS interoperability. To this end, we answer the following research questions:

- RQ1 How does RC2SQL's query evaluation perform compared to the state-of-the-art tools on both domain-independent and domain-dependent queries?
- RQ2 How does RC2SQL's query evaluation scale on large synthetic databases?
- RQ3 How does RC2SQL's query evaluation perform on real-world databases?
- RQ4 How does the count aggregation optimization impact RC2SQL's performance?
- RQ5 Can RC2SQL use different DBMSs for query evaluation?
- RQ6 How long does RC2SQL take to translate different queries (without query evaluation)?

We organize our evaluation into five experiments. Four experiments (SMALL, MEDIUM, LARGE, and REAL) focus on the type and size of the structures we use for query evaluation. The fifth experiment (INFINITE) focuses on the evaluation of non-evaluable (i.e., domain-dependent) queries that may potentially produce infinite evaluation results.

To answer RQ1, we compare our tool with the translation-based approach by Van Gelder and Topor [GT91] (VGT), the structure reduction approach by Ailamazyan et al. [AGSS86], and the DDD [MLAH99, Møl02], LDD [CGS09], and MonPoly<sup>REG</sup> [BKMZ15] tools that evaluate RC queries directly using infinite relations encoded as binary decision diagrams.

Experiments:	SMALL	MEDIUM	LARGE	INFINITE	REAL
RC2SQL	✓	✓	✓	✓*	✓
VGT	✓	✓	TO	N/A	✓
DDD	✓	TO	TO	✓	✓
LDD	✓	TO	TO	✓	✓
MonPoly <sup>REG</sup>	✓	TO	TO	✓	✓
Ailamazyan et al.	TO	TO	TO	TO	TO

\* Only states that the result is infinite.

Table 1: Applicability and performance of all the tools on all the experiments. TO = Timeout of 300s on all experiment runs, N/A = Not applicable

We could not find a publicly available implementation of Van Gelder and Topor’s translation. Therefore, the tool VGT for evaluable RC queries is derived from our implementation by modifying the function  $\text{rb}(\cdot)$  in Figure 7 to use the relation  $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$  (Appendix A, Figure 17) instead of  $\text{cov}(x, Q, \mathcal{G})$  (Figure 5) and to use the generator  $\bigvee_{Q_{ap} \in \mathcal{A}} \exists \vec{fv}(Q) \setminus \{x\}. Q_{ap}$  instead of  $\text{qps}^{\vee}(\mathcal{G})$ . Evaluable queries  $Q$  are always translated into  $(Q_{fin}, \perp)$  by  $\text{rw}(\cdot)$  because all of  $Q$ ’s free variables are range restricted. We exclude VGT from the comparison on non-evaluable queries (experiment INFINITE). Similarly, the implementation of Ailamazyan et al.’s approach was not available; hence we used our formally-verified implementation [Ras22]. The implementations of the remaining tools were publicly available.

We use Data Golf structures of growing size (experiments SMALL, MEDIUM, and LARGE) to answer RQ2. In contrast, to answer RQ3, we use real-world structures obtained from the Amazon review dataset [NLM19] (experiment REAL).

To answer RQ4, we also consider variants of the translation-based approaches without the step that uses count aggregation optimization  $\text{optcnt}(\cdot)$ , superscripted with a minus ( $^-$ ).

SQL queries computed by the translations are evaluated using the PostgreSQL and MySQL DBMS (RQ5). We superscript the tool names with <sup>P</sup> and <sup>M</sup> accordingly. In the LARGE experiment, we only use PostgreSQL because it consistently performed better than MySQL in the MEDIUM experiment. In all our experiments, the translation-based tools used a Data Golf structure with  $|\mathcal{T}^+| = |\mathcal{T}^-| = 2$  as the training database. We run our experiments on an AMD Ryzen 7 PRO 4750U computer with 32 GB RAM. The relations in PostgreSQL and MySQL are recreated before each invocation to prevent optimizations based on caching recent query evaluation results. We measure the query evaluation times of all the tools and the translation time of our RC2SQL tool (RQ6). We provide all our experiments in an easily reproducible and publicly available artifact [RBKT22a].

In the SMALL, MEDIUM, and LARGE experiments, we generate ten pseudorandom queries (denoted as  $Q_i, 1 \leq i \leq 10$ , see Appendix C) with a fixed size 14 and Data Golf structures  $\mathcal{S}$  (strategy  $\gamma = 1$ ). The queries satisfy the Data Golf assumptions along with a few additional ones: the queries are not safe range, every bound variable actually occurs in its scope, disjunction only appears at the top-level, and only pairwise distinct variables appear as terms in predicates. The queries have 2 free variables and every subquery has at most 4 free variables. We control the size of the Data Golf structure  $\mathcal{S}$  in our experiments using a parameter  $n = |\mathcal{T}^+| = |\mathcal{T}^-|$ . Because the sets  $\mathcal{T}^+$  and  $\mathcal{T}^-$  grow in the recursion on subqueries, relations in a Data Golf structure typically have more than  $n$  tuples. The values of the parameter  $n$  for Data Golf structures are summarized in Figure 15.

The INFINITE experiment consists of five pseudorandom queries  $Q_i^I$ ,  $1 \leq i \leq 5$  (Appendix C) that are *not* evaluable and  $\text{rw}(Q_i^I) = (Q_{i,\text{fin}}, Q_{i,\text{inf}})$ , where  $Q_{i,\text{inf}} \neq \perp$ . Specifically, the queries are of the form  $Q_1 \wedge \forall x, y. Q_2 \longrightarrow Q_3$ , where  $Q_1, Q_2$ , and  $Q_3$  are either atomic predicates or equalities. We choose the queries so that the number of their satisfying tuples is not too high, e.g., quadratic in the parameter  $n$ , because no tool can possibly enumerate so many tuples within the timeout. For each  $1 \leq i \leq 5$ , we compare the performance of our tool to tools that directly evaluate  $Q_i^I$  on structures generated by the two Data Golf strategies (parameter  $\gamma$ ), which trigger infinite or finite evaluation results on the considered queries. For infinite results, our tool outputs this fact (by evaluating  $Q_{i,\text{inf}}$ ), whereas the other tools also output a finite representation of the infinite result. For finite results, all tools produce the same output.

Figure 15 shows the empirical evaluation results for the experiments SMALL, MEDIUM, LARGE, and INFINITE. All entries are execution times in seconds, TO is a timeout, and RE is a runtime error. In the experiments SMALL, MEDIUM, and LARGE, the columns correspond to ten unique pseudorandom queries (the same queries are used in all the three experiments). In the INFINITE experiment, we use five unique pseudorandom queries and two Data Golf strategies. The time it takes for our translation RC2SQL to translate each query is shown in the first line for the experiments SMALL and INFINITE because the queries in the experiments MEDIUM and LARGE are the same as in SMALL. The remaining lines show evaluation times with the lowest time for a query typeset in bold. We omit the rows for tools that time out or crash on all queries of an experiment, e.g., Ailamazyan et al. [AGSS86]. Table 1 summarizes which tools were used in which experiment. We conclude that our translation RC2SQL significantly outperforms all other tools on all queries (RQ1) (except VGT on the fourth query, but on the smallest structure) and scales well to higher values of  $n$ , i.e., larger relations in the Data Golf structures, on all queries (RQ2). The count aggregation optimization (RQ4) provides no conclusive benefit for our translation, while it consistently improves VGT’s performance, especially in the MEDIUM experiment.

We also evaluate the tools on the queries  $Q^{\text{ susp}}$  and  $Q_{\text{user}}^{\text{ susp}}$  from the introduction and on the more challenging query  $Q_{\text{text}}^{\text{ susp}} := \text{B}(b) \wedge \exists u, s, t. \forall p. \text{P}(b, p) \longrightarrow \text{S}(p, u, s) \vee \text{T}(p, u, t)$  with an additional relation T that relates user’s review text (variable  $t$ ) to a product. The query  $Q_{\text{text}}^{\text{ susp}}$  computes all brands for which there is a user, a score, and a review text such that all the brand’s products were reviewed by that user with that score or by that user with that text. We use both Data Golf structures (strategy  $\gamma = 1$ ) and real-world structures obtained from the Amazon review dataset [NLM19]. The real-world relations P, S, and T are obtained by projecting the respective tables from the Amazon review dataset for two chosen product categories: gift cards (abbreviated GC) consisting of 147 194 reviews of 1 548 products and musical instruments (MI) consisting of 1 512 530 reviews of 120 400 products. The relation B contains all brands from P that have at least three products. Because the tool by Ailamazyan et al., DDD, LDD, and MonPoly<sup>REG</sup> only support integer data, we injectively remap the string and floating-point values from the Amazon review dataset to integers.

Figure 16 shows the empirical evaluation results: the time it takes for our translation RC2SQL to translate each query is shown in the first line and the execution times on Data Golf structures (left) and on structures derived from the real-world dataset for two specific product categories (right) are shown in the remaining lines. We remark that VGT cannot handle the query  $Q_{\text{user}}^{\text{ susp}}$  as it is not evaluable [GT91], hence we mark the correspond cells in Figure 16 with  $-$ . Our translation RC2SQL significantly outperforms all other tools (except VGT on  $Q^{\text{ susp}}$ , where RC2SQL and VGT have similar performance) on both Data Golf and



Query	$Q_1$	$Q_2$	$Q_3$	$Q_4$	$Q_5$	$Q_6$	$Q_7$	$Q_8$	$Q_9$	$Q_{10}$
Experiment SMALL, Evaluable pseudorandom queries $Q$ , $ \text{sub}(Q)  = 14$ , $n = 500$ :										
Translation time	0.6	0.0	0.1	0.0	0.0	0.0	0.0	0.1	0.0	0.0
RC2SQL <sup>P</sup>	<b>0.2</b>	0.2	0.3	0.2	0.2	<b>0.2</b>	0.2	<b>0.1</b>	<b>0.2</b>	<b>0.1</b>
RC2SQL <sup>M</sup>	0.3	0.2	0.3	0.2	0.2	<b>0.2</b>	<b>0.1</b>	0.2	0.3	0.2
RC2SQL <sup>-P</sup>	<b>0.2</b>	<b>0.1</b>	<b>0.2</b>	0.2	<b>0.1</b>	<b>0.2</b>	<b>0.1</b>	<b>0.1</b>	<b>0.2</b>	<b>0.1</b>
RC2SQL <sup>-M</sup>	<b>0.2</b>	<b>0.1</b>	0.3	0.2	<b>0.1</b>	<b>0.2</b>	<b>0.1</b>	<b>0.1</b>	0.4	0.3
VGTP <sup>P</sup>	1.5	0.2	1.8	1.3	1.9	8.6	2.5	1.5	15.6	4.8
VGTM <sup>M</sup>	0.3	0.2	0.3	<b>0.1</b>	0.2	56.3	6.1	0.2	155.7	13.5
VGTP <sup>-P</sup>	16.4	6.0	10.7	6.3	TO	6.2	3.1	2.3	48.0	8.8
VGTM <sup>-M</sup>	129.1	73.7	97.3	66.8	TO	52.1	19.1	12.8	TO	61.1
DDD	3.8	2.4	5.5	RE	1.2	3.7	4.7	2.2	17.9	5.3
LDD	35.9	16.0	46.2	15.6	9.0	28.3	12.9	17.4	206.0	32.6
MonPoly <sup>REG</sup>	31.3	11.2	29.7	10.5	10.2	31.2	10.7	21.0	103.3	24.0
Experiment MEDIUM, Evaluable pseudorandom queries $Q$ , $ \text{sub}(Q)  = 14$ , $n = 20000$ :										
RC2SQL <sup>P</sup>	2.1	1.1	2.2	1.2	1.1	1.2	<b>0.5</b>	0.9	1.8	<b>1.0</b>
RC2SQL <sup>M</sup>	6.2	2.6	6.6	2.7	2.4	4.4	1.6	2.7	9.0	3.3
RC2SQL <sup>-P</sup>	<b>1.4</b>	<b>0.8</b>	<b>1.4</b>	<b>0.8</b>	<b>0.8</b>	<b>1.0</b>	<b>0.5</b>	<b>0.6</b>	<b>1.7</b>	<b>1.0</b>
RC2SQL <sup>-M</sup>	4.3	1.9	5.0	1.8	1.9	3.3	1.6	2.2	8.3	3.2
VGTP <sup>P</sup>	2.9	1.0	3.0	2.2	2.8	TO	TO	2.3	TO	TO
VGTM <sup>M</sup>	4.6	2.4	5.2	2.6	2.7	TO	TO	3.0	TO	TO
VGTP <sup>-P</sup>	TO	TO	TO	TO	TO	TO	TO	TO	TO	TO
VGTM <sup>-M</sup>	TO	TO	TO	TO	TO	TO	TO	TO	TO	TO
Experiment LARGE, Evaluable pseudorandom queries $Q$ , $ \text{sub}(Q)  = 14$ , tool = RC2SQL <sup>P</sup> :										
$n = 40000$	4.5	2.3	4.4	2.3	2.2	2.5	1.0	1.7	3.7	1.8
$n = 80000$	8.8	4.5	8.7	4.8	4.5	5.0	1.8	3.3	7.2	3.7
$n = 120000$	14.1	6.8	12.8	7.2	7.0	7.1	2.8	5.1	10.8	4.9
Query	$Q_1^I$	$Q_2^I$	$Q_3^I$	$Q_4^I$	$Q_5^I$	$Q_1^I$	$Q_2^I$	$Q_3^I$	$Q_4^I$	$Q_5^I$
Experiment INFINITE, Non-evaluable pseudorandom queries $Q$ , $ \text{sub}(Q)  = 7$ , $n = 4000$ :										
	Infinite results ( $\gamma = 0$ )					Finite results ( $\gamma = 1$ )				
Translation time	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
RC2SQL <sup>P</sup>	0.5	0.5	0.5	0.5	0.5	0.5	1.5	0.7	0.7	1.4
RC2SQL <sup>M</sup>	<b>0.3</b>	0.5	<b>0.3</b>	0.5	0.5	<b>0.4</b>	<b>0.7</b>	0.4	0.6	<b>0.6</b>
RC2SQL <sup>-P</sup>	<b>0.3</b>	<b>0.3</b>	<b>0.3</b>	<b>0.3</b>	<b>0.3</b>	<b>0.4</b>	TO	<b>0.3</b>	<b>0.5</b>	TO
RC2SQL <sup>-M</sup>	0.4	1.0	0.5	0.7	0.9	0.6	TO	0.6	0.6	TO
DDD	32.4	81.0	32.7	60.9	81.6	32.1	68.6	31.9	59.4	68.2
LDD	TO	TO	TO	TO	TO	288.5	TO	TO	TO	TO
MonPoly <sup>REG</sup>	175.0	TO	175.4	TO	TO	160.3	299.1	160.0	TO	TO

Figure 15: Experiments SMALL, MEDIUM, LARGE, and INFINITE. We use the following abbreviations: TO = Timeout of 300s, RE = Runtime Error. Reported translation time for RC2SQL.

Query Param. $n$	$Q^{susp}$		$Q_{user}^{susp}$		$Q_{text}^{susp}$		Dataset	$Q^{susp}$		$Q_{user}^{susp}$		$Q_{text}^{susp}$	
	$10^3$	$10^4$	$10^3$	$10^4$	$10^3$	$10^4$		GC	MI	GC	MI	GC	MI
Translation time	0.0		0.0		0.3			0.0		0.0		0.3	
RC2SQL <sup>P</sup>	1.2	1.4	1.8	2.2	3.5	4.1		<b>1.6</b>	<b>8.4</b>	2.5	<b>11.9</b>	4.9	<b>51.3</b>
RC2SQL <sup>M</sup>	<b>0.3</b>	<b>1.1</b>	<b>0.3</b>	<b>1.4</b>	<b>0.5</b>	<b>2.6</b>		1.9	54.2	<b>2.4</b>	71.0	<b>4.5</b>	TO
RC2SQL <sup>-P</sup>	28.1	TO	28.6	TO	228.6	TO		132.9	TO	131.8	TO	TO	TO
RC2SQL <sup>-M</sup>	TO	TO	TO	TO	TO	TO		TO	TO	TO	TO	TO	TO
VGT <sup>P</sup>	1.5	1.7	–	–	204.3	TO		1.9	10.0	–	–	TO	TO
VGT <sup>M</sup>	<b>0.3</b>	1.4	–	–	TO	TO		1.7	60.9	–	–	TO	TO
VGT <sup>-P</sup>	TO	TO	–	–	TO	TO		TO	TO	–	–	TO	TO
VGT <sup>-M</sup>	TO	TO	–	–	TO	TO		TO	TO	–	–	TO	TO
DDD	4.0	TO	4.1	TO	18.6	TO		61.3	TO	61.0	TO	126.7	TO
LDD	22.3	TO	22.2	TO	148.5	TO		TO	TO	TO	TO	TO	TO
MonPoly <sup>REG</sup>	22.4	TO	22.6	TO	84.1	TO		TO	TO	TO	TO	TO	TO

Figure 16: Experiment REAL with the queries  $Q^{susp}$ ,  $Q_{user}^{susp}$ ,  $Q_{text}^{susp}$ . We use the following abbreviations: GC = Gift Cards, MI = Musical Instruments, TO = Timeout of 300s. Reported translation time for RC2SQL.

real-world structures (RQ3). VGT<sup>-</sup> translates  $Q^{susp}$  into a RANF query with a higher query cost than RC2SQL<sup>-</sup>. However, the optimization `optcnt(·)` manages to rectify this inefficiency (RQ4) and thus VGT exhibits a comparable performance as RC2SQL. Specifically, the factor of  $80\times$  in query cost between VGT<sup>-</sup> and RC2SQL<sup>-</sup> improves to  $1.1\times$  in query cost between VGT and RC2SQL on a Data Golf structure with  $n = 20$  [RBKT22a]. Nevertheless, VGT does not finish evaluating the query  $Q_{text}^{susp}$  on GC and MI datasets within 5 minutes, unlike RC2SQL. Finally, RC2SQL’s translation took less than 1 second on all the queries (RQ6).

## 9. CONCLUSION

We presented a translation-based approach to evaluating arbitrary relational calculus queries over an infinite domain with improved time complexity over existing approaches. This contribution is an important milestone towards making the relational calculus a viable query language for practical databases. In future work, we plan to integrate into our base language features that database practitioners love, such as inequalities, bag semantics, and aggregations.

**Acknowledgments.** This research has been supported by the Swiss National Science Foundation grants “Big Data Monitoring“ (167162) and “Model-driven Security & Privacy” (204796) as well as by a Novo Nordisk Fonden start package grant (NNF20OC0063462).

## REFERENCES

- [AGSS86] Alfred K. Ailamazyan, Mikhail M. Gilula, Alexei P. Stolboushkin, and Grigorii F. Schwartz. Reduction of a relational model with infinite domains to the case of finite domains. *Doklady Akademii Nauk SSSR*, 286(2):308–311, 1986. URL: <http://mi.mathnet.ru/dan47310>.
- [AH91] Arnon Avron and Yoram Hirshfeld. On first order database query languages. In *LICS 1991*, pages 226–231. IEEE Computer Society, 1991. doi:10.1109/LICS.1991.151647.

- [AHV95] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995. URL: <http://webdam.inria.fr/Alice/>.
- [BG04] Achim Blumensath and Erich Grädel. Finite presentations of infinite structures: Automata and interpretations. *Theory Comput. Syst.*, 37(6):641–674, 2004. doi:10.1007/s00224-004-1133-y.
- [BKMZ15] David Basin, Felix Klaedtke, Samuel Müller, and Eugen Zălinescu. Monitoring metric first-order temporal properties. *J. ACM*, 62(2):15:1–15:45, 2015. doi:10.1145/2699444.
- [BL00] Michael Benedikt and Leonid Libkin. Relational queries over interpreted structures. *J. ACM*, 47(4):644–680, 2000. doi:10.1145/347476.347477.
- [CGS09] Sagar Chaki, Arie Gurfinkel, and Ofer Strichman. Decision diagrams for linear arithmetic. In *FMCAD 2009*, pages 53–60. IEEE, 2009. doi:10.1109/FMCAD.2009.5351143.
- [CKMP97] Jens Claußen, Alfons Kemper, Guido Moerkotte, and Klaus Peithner. Optimizing queries with universal quantification in object-oriented and object-relational databases. In Matthias Jarke, Michael J. Carey, Klaus R. Dittrich, Frederick H. Lochovsky, Pericles Loucopoulos, and Manfred A. Jeusfeld, editors, *VLDB 1997*, pages 286–295. Morgan Kaufmann, 1997. URL: <http://www.vldb.org/conf/1997/P286.PDF>.
- [Cod72] E. F. Codd. Relational completeness of data base sublanguages. *Research Report / RJ / IBM / San Jose, California*, RJ987, 1972.
- [CT95] Jan Chomicki and David Toman. Implementing temporal integrity constraints using an active DBMS. *IEEE Trans. Knowl. Data Eng.*, 7(4):566–582, 1995. doi:10.1109/69.404030.
- [EHJ93] Martha Escobar-Molano, Richard Hull, and Dean Jacobs. Safety and translation of calculus queries with scalar functions. In Catriel Beeri, editor, *PODS 1993*, pages 253–264. ACM Press, 1993. doi:10.1145/153850.153909.
- [Ell13] Erling Ellingsen. Regex golf, 2013. URL: <https://alf.nu/RegexGolf>.
- [GT87] Allen Van Gelder and Rodney W. Topor. Safety and correct translation of relational calculus formulas. In Moshe Y. Vardi, editor, *PODS 1987*, pages 313–327. ACM, 1987. doi:10.1145/28659.28693.
- [GT91] Allen Van Gelder and Rodney W. Topor. Safety and translation of relational calculus queries. *ACM Trans. Database Syst.*, 16(2):235–278, 1991. doi:10.1145/114325.103712.
- [HS94] Richard Hull and Jianwen Su. Domain independence and the relational calculus. *Acta Informatica*, 31(6):513–524, 1994. doi:10.1007/BF01213204.
- [Kif88] Michael Kifer. On safety, domain independence, and capturability of database queries (preliminary report). In Catriel Beeri, Joachim W. Schmidt, and Umeshwar Dayal, editors, *JCDKB 1988*, pages 405–415. Morgan Kaufmann, 1988. doi:10.1016/b978-1-4832-1313-2.50037-8.
- [KM01] Nils Klarlund and Anders Møller. *MONA v1.4 User Manual*. BRICS, Department of Computer Science, University of Aarhus, 2001. URL: <http://www.brics.dk/mona/>.
- [KNR16] Mahmoud Abo Khamis, Hung Q. Ngo, and Atri Rudra. FAQ: questions asked frequently. In Tova Milo and Wang-Chiew Tan, editors, *PODS 2016*, pages 13–28. ACM, 2016. doi:10.1145/2902251.2902280.
- [Lib04] Leonid Libkin. *Elements of Finite Model Theory*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2004. doi:10.1007/978-3-662-07003-1.
- [LYL08] Hong-Cheu Liu, Jeffrey Xu Yu, and Weifa Liang. Safety, domain independence and translation of complex value database queries. *Inf. Sci.*, 178(12):2507–2533, 2008. doi:10.1016/j.ins.2008.02.005.
- [MLAH99] Jesper B. Møller, Jakob Lichtenberg, Henrik Reif Andersen, and Henrik Hulgaard. Difference decision diagrams. In Jörg Flum and Mario Rodríguez-Artalejo, editors, *CSL 1999*, volume 1683 of *LNCS*, pages 111–125. Springer, 1999. doi:10.1007/3-540-48168-0\_9.
- [Mø02] Jesper B. Møller. DDDLIB: A library for solving quantified difference inequalities. In Andrei Voronkov, editor, *CADE 2002*, volume 2392 of *LNCS*, pages 129–133. Springer, 2002. doi:10.1007/3-540-45620-1\_9.
- [NLM19] Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *EMNLP 2019*, pages 188–197. Association for Computational Linguistics, 2019. doi:10.18653/v1/D19-1018.
- [NRR13] Hung Q. Ngo, Christopher Ré, and Atri Rudra. Skew strikes back: new developments in the theory of join algorithms. *SIGMOD Rec.*, 42(4):5–16, 2013. doi:10.1145/2590989.2590991.

- [OS16] Dan Olteanu and Maximilian Schleich. Factorized databases. *SIGMOD Rec.*, 45(2):5–16, 2016. doi:10.1145/3003665.3003667.
- [Pao69] Robert A. Di Paola. The recursive unsolvability of the decision problem for the class of definite formulas. *J. ACM*, 16(2):324–327, 1969. doi:10.1145/321510.321524.
- [Ras22] Martin Raszyk. First-order query evaluation. *Archive of Formal Proofs*, 2022. Formal proof development. URL: [https://www.isa-afp.org/entries/Eval\\_F0.html](https://www.isa-afp.org/entries/Eval_F0.html).
- [RBKT22a] Martin Raszyk, David Basin, Srđan Krstić, and Dmitriy Traytel. Implementation and empirical evaluation associated with this article, 2022. URL: <https://github.com/mraszyk/rc2sql/tree/lmcs22>.
- [RBKT22b] Martin Raszyk, David Basin, Srđan Krstić, and Dmitriy Traytel. Practical relational calculus query evaluation. In Dan Olteanu and Nils Vortmeier, editors, *ICDT 2022*, volume 220 of *LIPICs*, pages 11:1–11:21. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022. doi:10.4230/LIPICs.ICDT.2022.11.
- [Rev02] Peter Z. Revesz. *Introduction to Constraint Databases*. Texts in Computer Science. Springer, 2002. doi:10.1007/b97430.
- [RT22] Martin Raszyk and Dmitriy Traytel. Making arbitrary relational calculus queries safe-range. *Archive of Formal Proofs*, 2022. Formal proof development. URL: [https://www.isa-afp.org/entries/Safe\\_Range\\_RC.html](https://www.isa-afp.org/entries/Safe_Range_RC.html).
- [Tra50] Boris A. Trakhtenbrot. Impossibility of an algorithm for the decision problem in finite classes. *Doklady Akademii Nauk SSSR*, 70(4):569–572, 1950.
- [Var81] Moshe Y. Vardi. The decision problem for database dependencies. *Inf. Process. Lett.*, 12(5):251–254, 1981. doi:10.1016/0020-0190(81)90025-9.
- [Var82] Moshe Y. Vardi. The complexity of relational query languages (extended abstract). In Harry R. Lewis, Barbara B. Simons, Walter A. Burkhard, and Lawrence H. Landweber, editors, *STOC 1982*, pages 137–146. ACM, 1982. doi:10.1145/800070.802186.
- [Yan19] Jun Yang. radb, 2019. URL: <https://github.com/junyang/radb>.

$\text{gen}_{\text{vgt}}(x, Q, \{Q\})$	if $\text{ap}(Q)$ and $x \in \text{fv}(Q)$ ;
$\text{gen}_{\text{vgt}}(x, \neg\neg Q, \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, Q, \mathcal{A})$ ;
$\text{gen}_{\text{vgt}}(x, \neg(Q_1 \vee Q_2), \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, (\neg Q_1) \wedge (\neg Q_2), \mathcal{A})$ ;
$\text{gen}_{\text{vgt}}(x, \neg(Q_1 \wedge Q_2), \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, (\neg Q_1) \vee (\neg Q_2), \mathcal{A})$ ;
$\text{gen}_{\text{vgt}}(x, \neg\exists y. Q_y, \mathcal{A})$	if $x \neq y$ and $\text{gen}_{\text{vgt}}(x, \neg Q_y, \mathcal{A})$ ;
$\text{gen}_{\text{vgt}}(x, Q_1 \vee Q_2, \mathcal{A}_1 \cup \mathcal{A}_2)$	if $\text{gen}_{\text{vgt}}(x, Q_1, \mathcal{A}_1)$ and $\text{gen}_{\text{vgt}}(x, Q_2, \mathcal{A}_2)$ ;
$\text{gen}_{\text{vgt}}(x, Q_1 \wedge Q_2, \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, Q_1, \mathcal{A})$ ;
$\text{gen}_{\text{vgt}}(x, Q_1 \wedge Q_2, \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, Q_2, \mathcal{A})$ ;
$\text{gen}_{\text{vgt}}(x, \exists y. Q_y, \mathcal{A})$	if $x \neq y$ and $\text{gen}_{\text{vgt}}(x, Q_y, \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, Q, \emptyset)$	if $x \notin \text{fv}(Q)$ ;
$\text{con}_{\text{vgt}}(x, Q, \{Q\})$	if $\text{ap}(Q)$ and $x \in \text{fv}(Q)$ ;
$\text{con}_{\text{vgt}}(x, \neg\neg Q, \mathcal{A})$	if $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, \neg(Q_1 \vee Q_2), \mathcal{A})$	if $\text{con}_{\text{vgt}}(x, (\neg Q_1) \text{ and } (\neg Q_2), \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, \neg(Q_1 \wedge Q_2), \mathcal{A})$	if $\text{con}_{\text{vgt}}(x, (\neg Q_1) \vee (\neg Q_2), \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, \neg\exists y. Q_y, \mathcal{A})$	if $x \neq y$ and $\text{con}_{\text{vgt}}(x, \neg Q_y, \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, Q_1 \vee Q_2, \mathcal{A}_1 \cup \mathcal{A}_2)$	if $\text{con}_{\text{vgt}}(x, Q_1, \mathcal{A}_1)$ and $\text{con}_{\text{vgt}}(x, Q_2, \mathcal{A}_2)$ ;
$\text{con}_{\text{vgt}}(x, Q_1 \wedge Q_2, \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, Q_1, \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, Q_1 \wedge Q_2, \mathcal{A})$	if $\text{gen}_{\text{vgt}}(x, Q_2, \mathcal{A})$ ;
$\text{con}_{\text{vgt}}(x, Q_1 \wedge Q_2, \mathcal{A}_1 \cup \mathcal{A}_2)$	if $\text{con}_{\text{vgt}}(x, Q_1, \mathcal{A}_1)$ and $\text{con}_{\text{vgt}}(x, Q_2, \mathcal{A}_2)$ ;
$\text{con}_{\text{vgt}}(x, \exists y. Q_y, \mathcal{A})$	if $x \neq y$ and $\text{con}_{\text{vgt}}(x, Q_y, \mathcal{A})$ .

Figure 17: The relations  $\text{gen}_{\text{vgt}}(x, Q, \mathcal{A})$  and  $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$  [GT91].

## APPENDIX A. EVALUABLE QUERIES

The classes of *evaluable* queries [GT91, Definition 5.2] and *allowed* queries [GT91, Definition 5.3] are decidable subsets of domain-independent RC queries. The evaluable queries characterize exactly the domain-independent queries with no repeated predicate symbols [GT91, Theorem 10.5]. Every evaluable query can be translated to an equivalent allowed query [GT91, Theorem 8.6] and every allowed query can be translated to an equivalent RANF query [GT91, Theorem 9.6].

**Definition A.1.** A query  $Q$  is called *evaluable* if

- every variable  $x \in \text{fv}(Q)$  satisfies  $\text{gen}_{\text{vgt}}(x, Q)$  and
- the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$  satisfies  $\text{con}_{\text{vgt}}(y, Q_y)$ .

A query  $Q$  is called *allowed* if

- every variable  $x \in \text{fv}(Q)$  satisfies  $\text{gen}_{\text{vgt}}(x, Q)$  and
- the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$  satisfies  $\text{gen}_{\text{vgt}}(y, Q_y)$ ,

where the relation  $\text{gen}_{\text{vgt}}(x, Q)$  is defined to hold iff there exists a set of atomic predicates  $\mathcal{A}$  such that  $\text{gen}_{\text{vgt}}(x, Q, \mathcal{A})$  and the relation  $\text{con}_{\text{vgt}}(x, Q)$  is defined to hold iff there exists a set of atomic predicates  $\mathcal{A}$  such that  $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$ , respectively. The relations  $\text{gen}_{\text{vgt}}(x, Q, \mathcal{A})$  and  $\text{con}_{\text{vgt}}(x, Q, \mathcal{A})$  are defined in Figure 17.

The termination of the rules in Figure 17 follow using the measure  $m(Q)$  (Figure 11). We now relate the definitions from Figure 4 and Figure 17 with the following lemmas.

**Lemma A.2.** Let  $x$  and  $y$  be free variables in a query  $Q$  such that  $\text{gen}_{\text{vgt}}(x, \neg Q)$  and  $\text{gen}_{\text{vgt}}(y, Q)$  hold. Then we get a contradiction.

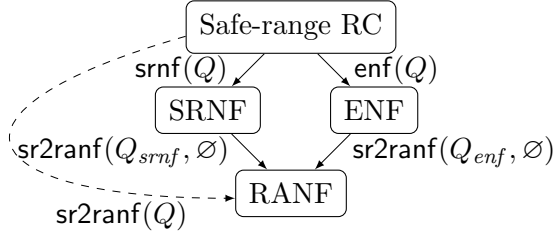


Figure 18: Alternative translation from safe-range RC to RANF query via ENF.

*Proof.* This is proved by induction on the query  $Q$  using the measure  $m(Q)$  on queries defined in Figure 11, which decreases in every case of the definition in Figure 17.  $\square$

**Lemma A.3.** *Let  $Q$  be a query such that  $\text{gen}_{\text{vgt}}(y, Q_y)$  holds for the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$ . Suppose that  $\text{gen}_{\text{vgt}}(x, Q)$  holds for a free variable  $x \in \text{fv}(Q)$ . Then  $\text{gen}(x, Q)$  holds.*

*Proof.* This is proved by induction on the query  $Q$  using the measure  $m(Q)$  on queries defined in Figure 11, which decreases in every case of the definition in Figure 17.

Lemma A.2 and the assumption that  $\text{gen}_{\text{vgt}}(y, Q_y)$  holds for the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$  imply that  $\text{gen}_{\text{vgt}}(x, Q)$  cannot be derived using the rule  $\text{gen}_{\text{vgt}}(x, \neg\exists y. Q_y)$ , i.e.,  $Q$  cannot be of the form  $\neg\exists y. Q_y$ . Every other case in the definition of  $\text{gen}_{\text{vgt}}(x, Q)$  has a corresponding case in the definition of  $\text{gen}(x, Q)$ .  $\square$

**Lemma A.4.** *Let  $Q$  be an allowed query, i.e.,  $\text{gen}_{\text{vgt}}(x, Q)$  holds for every free variable  $x \in \text{fv}(Q)$  and  $\text{gen}_{\text{vgt}}(y, Q_y)$  holds for the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$ . Then  $Q$  is a safe-range query, i.e.,  $\text{gen}(x, Q)$  holds for every free variable  $x \in \text{fv}(Q)$  and  $\text{gen}(y, Q_y)$  holds for the bound variable  $y$  in every subquery  $\exists y. Q_y$  of  $Q$ .*

*Proof.* The lemma is proved by applying Lemma A.3 to every free variable of  $Q$  and to the bound variable  $y$  in every subquery of  $Q$  of the form  $\exists y. Q_y$ .  $\square$

Lemma A.4 shows that every allowed query is safe range. But there exist safe-range queries that are not allowed, e.g.,  $B(x) \wedge x \approx y$ .

## APPENDIX B. EXISTENTIAL NORMAL FORM

Recall that our translation uses a standard approach to obtain RANF queries from safe-range queries via SRNF [AHV95]. In this section we introduce existential normal form (ENF), an alternative normal form to SRNF that can also be used to translate safe-range queries to RANF queries, and we discuss why we opt for using SRNF instead.

Figure 18 shows an overview of the RC fragments and query normal forms (nodes) and the functions we use to translate between them (edges). The dashed edge shows the translation of a safe-range query to RANF we opt for in this article. It is the composition of the two translations from safe-range RC to SRNF and from SRNF to RANF, respectively. In the rest of this section we introduce ENF and the corresponding translations to RANF.

ENF was introduced by Van Gelder and Topor [GT91] to translate an allowed query [GT91] into an equivalent RANF query. Given a safe-range query in ENF, the rules (R1)–(R3) from Section 3.4 can be applied to obtain an equivalent RANF query [EHJ93, Lemma 7.8].

We remark that the rules (R1)–(R3) are not sufficient to yield an equivalent RANF query for the original definition of ENF [GT91]. This issue has been identified and fixed by Escobar-Molano et al. [EHJ93]. Unlike SRNF, a query in ENF can have a subquery of the form  $\neg(Q_1 \wedge Q_2)$ , but no subquery of the form  $\neg Q_1 \vee Q_2$  or  $Q_1 \vee \neg Q_2$ . A function  $\text{enf}(Q)$  that yields an ENF query equivalent to  $Q$  can be defined in terms of subquery rewriting using the rules in [EHJ93, Figure 2].

Analogously to [EHJ93, Lemma 7.4], if a query  $Q$  is safe range, then  $\text{enf}(Q)$  is also safe range. Next we prove the following lemma that we could use as a precondition for translating safe-range queries in ENF to queries in RANF.

**Lemma B.1.** *Let  $Q_{\text{enf}}$  be a query in ENF. Then  $\text{gen}(x, \neg Q')$  does not hold for any variable  $x$  and subquery  $\neg Q'$  of  $Q_{\text{enf}}$ .*

*Proof.* Assume that  $\text{gen}(x, \neg Q')$  holds for a variable  $x$  in a subquery  $\neg Q'$  of  $Q_{\text{enf}}$ . We derive a contradiction by induction on  $\text{m}(Q_{\text{enf}})$ . According to Figure 4 and by definition of ENF,  $\text{gen}(x, \neg Q')$  can only hold if  $Q'$  is a conjunction. Then  $\text{gen}(x, \neg Q')$  implies  $\text{gen}(x, \neg Q_1)$  for some query  $Q_1 \in \text{flat}^\wedge(Q')$  that is not a negation (by definition of ENF) or conjunction (by definition of  $\text{flat}^\wedge(\cdot)$ ), i.e.,  $Q_1$  is a disjunction (according to Figure 4). Then  $\text{gen}(x, \neg Q_1)$  implies  $\text{gen}(x, \neg Q_2)$  for some query  $Q_2 \in \text{flat}^\vee(Q_1)$  that is not a negation (by definition of ENF) or disjunction (by definition of  $\text{flat}^\vee(\cdot)$ ), i.e.,  $Q_2$  is a conjunction (according to Figure 4). Next we observe that  $\neg Q_2$  is in ENF because  $Q_2$  is a subquery of the ENF query  $Q_{\text{enf}}$ ,  $Q_2$  is a conjunction, and  $Q_2$  is a subquery of a disjunction ( $Q_1$ ) in  $Q_{\text{enf}}$ . Moreover,  $\text{m}(\neg Q_2) < \text{m}(Q_1) < \text{m}(Q') < \text{m}(Q_{\text{enf}})$ . This allows us to apply the induction hypothesis to the ENF query  $\neg Q_2$  and its subquery  $\neg Q_2$  (note that a query is a subquery of itself) and derive that  $\text{gen}(x, \neg Q_2)$  does not hold, which is a contradiction.  $\square$

Although applying the rules (R1)–(R3) to  $\text{enf}(Q)$  instead of  $\text{srnf}(Q)$  may result in a RANF query with fewer subqueries, the query cost, i.e., the time complexity of query evaluation, can be arbitrarily larger. We illustrate this in the following example that is also included in our artifact [RBKT22a]. We thus opt for using SRNF instead of ENF for translating safe-range queries into RANF.

**Example B.2.** *The safe-range query  $Q_{\text{enf}} := P_2(x, y) \wedge \neg(P_1(x) \wedge P_1(y))$  is in ENF and RANF, but not SRNF. Applying the rule (R1) to  $\text{srnf}(Q_{\text{enf}})$  yields the RANF query  $Q_{\text{srnf}} := (P_2(x, y) \wedge \neg P_1(x)) \vee (P_2(x, y) \wedge \neg P_1(y))$  that is equivalent to  $Q_{\text{enf}}$ . The costs of the two queries over a structure  $\mathcal{S}$  are  $\text{cost}^{\mathcal{S}}(Q_{\text{enf}}) = 2 \cdot \|\llbracket P_2(x, y) \rrbracket\| + \|\llbracket P_1(x) \rrbracket\| + \|\llbracket P_1(y) \rrbracket\| + 2 \cdot \|\llbracket P_1(x) \wedge P_1(y) \rrbracket\| + 2 \cdot \|\llbracket Q_{\text{enf}} \rrbracket\|$  and  $\text{cost}^{\mathcal{S}}(Q_{\text{srnf}}) = 2 \cdot \|\llbracket P_2(x, y) \rrbracket\| + \|\llbracket P_1(x) \rrbracket\| + 2 \cdot \|\llbracket P_2(x, y) \rrbracket\| + \|\llbracket P_1(y) \rrbracket\| + 2 \cdot \|\llbracket P_2(x, y) \wedge \neg P_1(x) \rrbracket\| + 2 \cdot \|\llbracket P_2(x, y) \wedge \neg P_1(y) \rrbracket\| + 2 \cdot \|\llbracket Q_{\text{srnf}} \rrbracket\|$ , respectively. Note that the cost of  $Q_{\text{enf}}$  can be arbitrarily larger if  $P_1(x) \wedge P_1(y)$  evaluates to a large intermediate result, i.e.,  $\|\llbracket P_1(x) \wedge P_1(y) \rrbracket\| \gg \|\llbracket P_2(x, y) \rrbracket\|$ . In contrast, the cost of  $Q_{\text{srnf}}$  can only be larger by a constant factor.*

## APPENDIX C. QUERIES USED IN THE EVALUATION

Figure 19 shows the randomly generated queries that we use in our evaluation. The index of each predicate indicates the predicate’s arity.

$$\begin{aligned}
Q_1 &:= \neg(\exists x_2. \neg(\exists x_3. ((P_3(x_1, x_0, x_3)) \wedge (\neg(\exists x_4. P_4(x_1, x_3, x_4, x_2)))) \wedge (\exists x_4. P_2(x_1, x_4)))) \vee (Q_2(x_1, x_0)) \\
Q_2 &:= \neg(\exists x_2. \neg(\exists x_3. (\neg(P_4(x_1, x_0, x_2, x_3))) \wedge (P_3(x_1, x_3, x_0)))) \wedge (\exists x_2. \exists x_3. (x_1 = x_2) \wedge (Q_3(x_0, x_2, x_3))) \\
Q_3 &:= (\exists x_2. \neg(\exists x_3. ((P_3(x_1, x_0, x_3)) \wedge (P_1(x_0))) \wedge ((x_0 = x_1) \wedge (\neg(P_4(x_1, x_2, x_3, x_0))))) \wedge (\exists x_2. Q_3(x_1, x_2, x_0)) \\
Q_4 &:= (\exists x_2. P_3(x_1, x_2, x_0)) \wedge (\neg(x_0 = x_1)) \wedge (\neg(\exists x_2. \neg(\exists x_3. (P_2(x_0, x_3)) \wedge (\neg(P_4(x_1, x_3, x_2, x_0)))))) \\
Q_5 &:= (\exists x_2. \exists x_3. \neg(\exists x_4. (\exists x_5. P_3(x_0, x_5, x_4)) \wedge (\neg(P_4(x_0, x_4, x_2, x_3)))) \wedge ((\neg(x_0 = x_1)) \wedge (P_2(x_0, x_1))) \\
Q_6 &:= (\exists x_2. \neg(\exists x_3. ((\neg(\exists x_4. P_3(x_0, x_3, x_4))) \wedge (Q_3(x_0, x_1, x_3))) \wedge (\neg(R_3(x_0, x_2, x_1))))) \wedge (\exists x_2. S_3(x_0, x_1, x_2)) \\
Q_7 &:= (\exists x_2. (P_3(x_0, x_2, x_1)) \wedge (x_0 = x_2)) \wedge (\exists x_2. \neg(\exists x_3. (\neg(\exists x_4. \exists x_5. P_4(x_0, x_2, x_5, x_4))) \wedge (P_2(x_1, x_3)))) \\
Q_8 &:= \neg(\exists x_2. \neg(\exists x_3. (\exists x_4. (P_3(x_0, x_3, x_1)) \wedge ((P_4(x_0, x_3, x_1, x_4)) \wedge (x_3 = x_4))) \wedge (\neg(\exists x_4. Q_4(x_0, x_4, x_3, x_2)))))) \\
Q_9 &:= P_2(x_1, x_0) \wedge (\exists x_2. \neg(\exists x_3. ((\neg(P_4(x_1, x_3, x_2, x_0))) \wedge (P_3(x_0, x_3, x_1))) \wedge (\neg(Q_2(x_0, x_3))))) \wedge (x_0 = x_1) \\
Q_{10} &:= (\exists x_2. P_3(x_1, x_2, x_0)) \vee (\neg(\exists x_2. \neg(\exists x_3. (\neg(P_2(x_1, x_2))) \wedge (\exists x_4. (P_1(x_0)) \wedge (P_4(x_0, x_3, x_1, x_4))))) \\
Q_1^I &:= P_1(x_0) \wedge \neg(\exists x_2. \exists x_3. (P_3(x_0, x_2, x_3) \wedge \neg(x_0 = x_1))) \\
Q_2^I &:= P_1(x_0) \wedge \neg(\exists x_2. \exists x_3. (P_3(x_0, x_2, x_3) \wedge \neg P_2(x_1, x_3))) \\
Q_3^I &:= P_1(x_0) \wedge \neg(\exists x_2. \exists x_3. (P_3(x_0, x_3, x_2) \wedge \neg(x_1 = x_2))) \\
Q_4^I &:= P_1(x_1) \wedge \neg(\exists x_2. \exists x_3. (P_3(x_1, x_3, x_2) \wedge \neg P_4(x_1, x_0, x_3, x_2))) \\
Q_5^I &:= P_1(x_0) \wedge \neg(\exists x_2. \exists x_3. (P_3(x_0, x_2, x_3) \wedge \neg Q_3(x_1, x_2, x_3)))
\end{aligned}$$

Figure 19: Randomly generated queries.