

Practical Relational Calculus Query Evaluation

Martin Raszyk 


Department of Computer Science, ETH Zürich, Switzerland

David Basin 

Department of Computer Science, ETH Zürich, Switzerland

Srdan Krstić 

Department of Computer Science, ETH Zürich, Switzerland

Dmitriy Traytel 

Department of Computer Science, University of Copenhagen, Denmark

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.

2012 ACM Subject Classification Theory of computation → Database query languages (principles)

Keywords and phrases Relational calculus, relative safety, safe-range, query translation

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

1 Introduction

Codd’s theorem states that all domain-independent queries of the relational calculus (RC) can be expressed in relational algebra (RA) [10]. 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, the SQL 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 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 the informal description. It is not clear how to evaluate it though because its second conjunct is domain-dependent as it is satisfied



© Martin Raszyk, David Basin, Srdan Krstić, and Dmitriy Traytel;
licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:30

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

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):

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

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 [10], one could in principle also use the active domain to obtain canonical but far less efficient generators.)

Van Gelder and Topor [13, 14] 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 [2, 5, 15] or to abandon finite relations entirely and evaluate queries using finite representations of infinite (but well-behaved) relations such as systems of constraints [26] or automatic structures [6]. These approaches favor expressiveness over efficiency. 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 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 [3]: 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.

Our translation of the RC query to two RA expressions proceeds in several steps via safe-range queries and the relational algebra normal form (Section 3). We focus on the first step of translating RC to two safe-range RC queries (Section 4), which fundamentally differs from Van Gelder and Topor's approach and produces better generators like $\pi_{user,score}(S)$. The better generators strictly improve the time complexity of query evaluation (Section 4.4).

After the more standard transformations to relational algebra normal form and from there to RA expressions, we translate the resulting RA expressions into SQL using the `radb` tool [29]. Along the way to SQL, we leverage different ideas from the literature to optimize the overall result (Section 5). For example, we generalize Claußen et al. [9]'s approach to avoid evaluating Cartesian products like $\pi_{user,score}(S) \times P$ in the above translation by using count aggregations.

The overall translation allows us to use standard database management systems to evaluate RC queries. We implement our translation and use PostgreSQL to evaluate the resulting queries. With this setup, we empirically show that our approach outperforms Van Gelder and Topor’s translation (which also uses PostgreSQL for evaluation) and other approaches (Section 5). In our evaluation, we use the realistic Amazon review dataset [23] and design a synthetic benchmark that generates hard database instances for random RC queries.

In summary, the following are our three main contributions:

- 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 [14].
- We implement our translation and extend it to produce SQL queries. The resulting tool RC2SQL makes RC a viable input language for standard database management systems. We evaluate our tool on synthetic and realistic 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.

2 Related Work

We first recall the fundamental notions of *capturability* and *data complexity*. Kifer [16] 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 [3] 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: A query class is *capturable* if there is an algorithm that, given a query in the class, a 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 [28] 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 [18]. 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 [10] affirmatively answers this question for domain-independent queries, restricting the structure to the *active domain*. Ailamazyan et al. [2] show that by extending the active domain with a few additional elements, whose number depends only on the query, one can decide relative safety for RC. Further extensions are *natural-active collapse* results that combine the structure reduction with a translation-based approach [5, 15]. Our work is inspired by these theoretical landmarks (Section 4.1). Yet we avoid using (extended) active domains, which make query evaluation impractical due to a prohibitively high time complexity.

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. Van Gelder and Topor pioneered this approach [13, 14] 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. Subsequently, extensions to scalar and complex function symbols have been studied [12, 19]. 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 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 [26] obviate the need for using finite tables when evaluating RC queries. This yields significant expressiveness gains over RC. Yet the efficiency of the quantifier elimination procedures employed cannot compare with the simple evaluation of a projection operation in RA. Similarly, automatic structures [6] 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 [4, 7, 17, 20, 21].

3 Preliminaries

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

3.1 Relational Calculus

A signature σ 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 $\exists \vec{v}. Q$ as a shorthand 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 the query Q . We use \approx to denote the equality of terms in RC in order to distinguish it from $=$, which denotes the meta-level equality. We point out that defining $Q_1 \vee Q_2$ as a shorthand for $\neg(\neg Q_1 \wedge \neg Q_2)$ would complicate the definitions of the relevant classes of RC queries, e.g., the safe-range queries.

We define the subquery partial order \sqsubseteq on queries inductively on the structure of RC queries, e.g., Q_y is a subquery of the query $Q \wedge \neg \exists y. Q_y$. One can also view \sqsubseteq as the (transitive) subterm relation on the datatype of RC queries. We denote by $\text{sub}(Q)$ the set of subqueries of a query Q and by $\text{fv}(Q)$ the set of *free variables* in 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 $\exists x. Q$ the query obtained by existentially quantifying a variable x from a query Q if x is free in Q , i.e., $\exists x. Q := \exists x. Q$ if $x \in \text{fv}(Q)$ and $\exists x. Q := Q$ otherwise. We lift this notation to sets of queries in the standard way. We use $\exists x. Q$ (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 [18]) since the domain \mathcal{D} can still be infinite. A (variable) *assignment* is a mapping $\alpha : \mathcal{V} \rightarrow \mathcal{D}$. We additionally define α on constant symbols $c \in \mathcal{C}$ as $\alpha(c) = c^{\mathcal{S}}$. We write $\alpha[x \mapsto d]$ for the assignment that maps x to $d \in \mathcal{D}$ and is otherwise identical to α . We lift this notation to sequences \vec{x} and \vec{d} of pairwise distinct variables and arbitrary domain elements of

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

■ **Figure 1** The semantics of RC.

the same length. The semantics of RC queries for a structure \mathcal{S} and an assignment α is defined in Figure 1. 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 α' that agrees with α 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

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

We omit \mathcal{S} from $[[Q]]^{\mathcal{S}}$ if \mathcal{S} is fixed in the given context. We call values from $[[Q]]^{\mathcal{S}}$ assigned to $x \in \text{fv}(Q)$ as Q 's 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 a also 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 α . 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 α . 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 [24, 27].

We denote by $Q[x \mapsto y]$ the query obtained from the query Q after replacing each free occurrence of the variable x by the variable y (possibly renaming bound variables to avoid capture) and performing constant propagation (Appendix B), i.e., simplifications like $(x \approx x) \equiv \top$, $Q \wedge \perp \equiv \perp$, $Q \vee \perp \equiv Q$, etc. We lift this notation to sets of queries in the standard way. Finally, we denote by $Q[x/\perp]$ the query obtained from Q after replacing every atomic predicate or equality containing a free variable x by \perp (except for $x \approx x$) and performing constant propagation.

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 [1] is a decidable subset of domain-independent RC queries. Its definition is based on the notion of 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” [1]. Intuitively, atomic predicates restrict the possible values of a variable that occurs in them as a term. A similar idea applies to equalities between a variable and a constant. 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 a range-restricted variable in Q , then both x and y are range-restricted variables in $Q \wedge x \approx y$.

We formalize range-restricted variables using the *generated* relation $\text{gen}(x, Q, \mathcal{G})$, defined in Figure 2. Specifically, $\text{gen}(x, Q, \mathcal{G})$ holds if x is a range-restricted variable in Q and every

$\text{gen}(x, \perp, \emptyset);$	$\text{cov}(x, x \approx x, \emptyset);$	
$\text{gen}(x, Q, \{Q\})$ 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 2** The *generated* relation.■ **Figure 3** The *covered* relation.

224 satisfying assignment for Q satisfies some quantified predicate, referred to as *generator*, from
 225 \mathcal{G} . Note that, unlike in a similar definition by Van Gelder and Topor [14, Fig. 5], the rule
 226 $\text{gen}(x, \exists y. Q_y, \exists y. \mathcal{G})$ existentially quantifies the bound variable y from all queries in \mathcal{G} where
 227 y occurs. Hence, $\text{gen}(x, Q, \mathcal{G})$ implies $\text{fv}(\mathcal{G}) \subseteq \text{fv}(Q)$. Now we formalize these relationships.

228 ► **Lemma 1.** *Let Q be a query, $x \in \text{fv}(Q)$, and \mathcal{G} be a set of quantified predicates such that*
 229 *$\text{gen}(x, Q, \mathcal{G})$. Then (i) $x \in \text{fv}(Q_{qp})$ and $\text{fv}(Q_{qp}) \subseteq \text{fv}(Q)$ hold for every $Q_{qp} \in \mathcal{G}$, (ii) for every*
 230 *assignment α such that $\alpha \models Q$, there exists $Q_{qp} \in \mathcal{G}$ such that $\alpha \models Q_{qp}$, and (iii) $Q[x/\perp] = \perp$.*

231 ► **Definition 2.** *We define $\text{gen}(x, Q)$ to hold iff there exists a set \mathcal{G} such that $\text{gen}(x, Q, \mathcal{G})$. Let*
 232 *$\text{nongens}(Q) := \{x \in \text{fv}(Q) \mid \text{gen}(x, Q) \text{ does not hold}\}$ be the set of free variables in a query Q*
 233 *that are not range-restricted. A query Q has range-restricted free variables if every free variable*
 234 *of Q is range-restricted, i.e., $\text{nongens}(Q) = \emptyset$. A query Q has range-restricted bound variables*
 235 *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.*
 236 *A query is safe-range if it has range-restricted free and range-restricted bound variables.*

237 Relational algebra normal form (RANF) is a class of safe-range queries that can be easily
 238 mapped to RA [1]. In Appendix C, we define the predicate $\text{ranf}(\cdot)$ characterizing RANF
 239 queries and the translation $\text{sr2ranf}(\cdot)$ of a safe-range query into an equivalent RANF query.

240 3.3 Query Cost

241 To assess the time complexity of evaluating a RANF query Q , we define the *cost* of Q over
 242 a structure \mathcal{S} , denoted $\text{cost}^{\mathcal{S}}(Q)$, to be the sum of intermediate result sizes over all RANF
 243 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
 244 evaluating Q following its RANF structure (Appendix C, Figure 11) using RA operations pro-
 245 jection, column duplication, selection, set union, binary join, and anti-join. The complexity of
 246 these operations is linear in the combined input and output size (ignoring logarithmic factors
 247 due to set operations). The output size (the number of tuples times the number of variables)
 248 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 sub-
 249 queries. Repeated subqueries are only considered once, which does not affect the asymptotics
 250 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) $\llbracket Q \rrbracket = \llbracket Q_{fin} \rrbracket$ is a finite set if Q_{inf} does not hold; $\llbracket Q \rrbracket$ is an infinite set if Q_{inf} holds.

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 of tuples. In particular, $\llbracket Q \rrbracket$ is a finite set of tuples if Q_{inf} does not hold. In our translation, we combine Hull and Su's case distinction [15] to restrict variables and Van Gelder and Topor's idea [14] to replace the active domain by a smaller set specific to each variable.

4.1 Restricting One Variable

Let \tilde{Q} 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}' so that x is range-restricted except in a single subquery of \tilde{Q}' that is the negation of a disjunction of quantified predicates with a free occurrence of x and equalities of the form $x \approx y$ for some y . From the case distinction “for the corresponding variable: in or out of adom , and equality or inequality to other ‘previous’ variables if out of adom ” [15], we translate \tilde{Q} into the following equivalent query:

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

Here, $x \in \text{adom}(\tilde{Q})$ stands for an RC query expressing that the single free variable x is in $\text{adom}^S(\tilde{Q})$. The translation distinguishes the following three cases for a fixed valuation α :

- if $x \in \text{adom}(\tilde{Q})$ holds, 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 $\tilde{Q}[x/\perp]$, i.e., all atomic predicates with a free occurrence of x can be replaced by \perp (because $\neg x \in \text{adom}(\tilde{Q})$), all equalities $x \approx y$ and $y \approx x$ for $y \in \text{fv}(\tilde{Q}) \setminus \{x\}$ can be replaced by \perp (because $x \neq y$), and all equalities $x \approx z$ for a bound variable z can be replaced by \perp (because $\neg x \in \text{adom}(\tilde{Q})$ and z is range-restricted in its subquery $\exists z. Q_z$, by assumption, i.e., $\text{gen}(z, Q_z)$, thus $z \in \text{adom}(Q_z)$, and thus $z \in \text{adom}(\tilde{Q})$).

Note that $\exists \vec{\text{fv}}(\tilde{Q}) \setminus \{x\}. Q$ is the query in which all free variables of Q except x are existentially quantified. Given a set of quantified predicates \mathcal{G} , we write $\exists \vec{\alpha}. \mathcal{G}$ for $\bigvee_{Q_{qp} \in \mathcal{G}} \exists \vec{\alpha}. Q_{qp}$. To avoid enumerating the entire active domain $\text{adom}^S(Q)$ of the query Q and a structure \mathcal{S} , Van Gelder and Topor [14] replace the condition $x \in \text{adom}(Q)$ in their translation by $\exists \vec{\text{fv}}(\mathcal{G}) \setminus \{x\}. \mathcal{G}$, where generator set \mathcal{G} is a subset of atomic predicates. Because their translation [14] must yield an equivalent query (for every finite or infinite domain), \mathcal{G} must satisfy

$$\neg \exists \vec{\text{fv}}(\mathcal{G}) \setminus \{x\}. \mathcal{G} \implies Q \equiv Q[x/\perp] \quad (\text{VGT}_1) \quad \text{and} \quad Q[x/\perp] \implies \forall x. Q \quad (\text{VGT}_2).$$

Note that $Q[x/\perp] \implies \forall x. Q$ does not hold for the query $Q := \neg B(x)$ and thus a generator set \mathcal{G} of atomic predicates satisfying (VGT₂) only exists for a proper subset of all RC queries. In contrast, we only require that \mathcal{G} satisfies (VGT₁) in our translation. To this end, we define a *covered* relation $\text{cov}(x, Q, \mathcal{G})$ (in contrast to Van Gelder and Topor's *constrained* relation $\text{con}_{\text{vgt}}(x, Q, \mathcal{G})$ defined in Appendix A, Figure 8) such that, for every variable x and query \tilde{Q} with range-restricted bound variables, there exists at least one set \mathcal{G} such that $\text{cov}(x, \tilde{Q}, \mathcal{G})$ and (VGT₁) holds. Figure 3 shows the definition of this relation. Unlike the generator set \mathcal{G} in $\text{gen}(x, Q, \mathcal{G})$, the *cover* set \mathcal{G} in $\text{cov}(x, Q, \mathcal{G})$ may also contain equalities between two

variables. Hence, 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 $\text{cov}(x, Q, \mathcal{G})$ in the following lemma.

► **Lemma 3.** *Let \tilde{Q} be a query with range-restricted bound variables and $x \in \text{fv}(\tilde{Q})$. Then there exists a set \mathcal{G} of quantified predicates and equalities such that $\text{cov}(x, \tilde{Q}, \mathcal{G})$ and*

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

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 $\exists \vec{fv}(Q) \setminus \{x\}. \text{qps}(\mathcal{G})$. This yields our optimized translation characterized by the following lemma.

► **Lemma 4.** *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 $\text{cov}(x, \tilde{Q}, \mathcal{G})$. 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*

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

Note that x is not guaranteed to be range-restricted in (\star) 's last disjunct. However, it occurs only in the negation of a disjunction of quantified predicates with a free occurrence of x and equalities of the form $x \approx c$ or $x \approx y$. We will show how to handle such occurrences in Sections 4.2 and 4.3. Moreover, the negation of the disjunction can be omitted if (VGT₂) 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 that

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

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 5.** *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*

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

Our approach for restricting all bound variables recursively applies Lemma 5. 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 4, where 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

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

```

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

Figure 4 Restricting bound variables.

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

```

1 function fixfree( $\mathcal{Q}_{fin}$ ) =
  { $(Q_{fix}, Q^=) \in \mathcal{Q}_{fin} \mid \text{nongens}(Q_{fix}) \neq \emptyset$ };
2 function inf( $\mathcal{Q}_{fin}, Q$ ) = { $(Q_\phi, Q^=) \in$ 
   $\mathcal{Q}_{fin} \mid \text{disjointvars}(Q_\phi, Q^=) \neq \emptyset \vee$ 
   $\text{fv}(Q_\phi \wedge Q^=) \neq \text{fv}(Q)$ };
3 function split( $Q$ ) =
4    $\mathcal{Q}_{fin} := \{(\text{rb}(Q), \top)\}$ ,  $\mathcal{Q}_{inf} := \emptyset$ ;
5   while  $\text{fixfree}(\mathcal{Q}_{fin}) \neq \emptyset$  do
6      $(Q_{fix}, Q^=) \leftarrow \text{fixfree}(\mathcal{Q}_{fin})$ ;
7      $x \leftarrow \text{nongens}(Q_{fix})$ ;
8      $\mathcal{G} \leftarrow \{\mathcal{G} \mid \text{cov}(x, Q_{fix}, \mathcal{G})\}$ ;
9      $\mathcal{Q}_{fin} := (\mathcal{Q}_{fin} \setminus \{(Q_{fix}, Q^=)\}) \cup$ 
       $\{(Q_{fix} \wedge \text{qps}^\vee(\mathcal{G}), Q^=)\} \cup$ 
       $\bigcup_{y \in \text{eqs}(x, \mathcal{G})} \{(Q_{fix}[x \mapsto y], Q^= \wedge x \approx y)\}$ ;
10     $\mathcal{Q}_{inf} := \mathcal{Q}_{inf} \cup \{Q_{fix}[x/\perp]\}$ ;
11  while  $\text{inf}(\mathcal{Q}_{fin}, Q) \neq \emptyset$  do
12     $(Q_\phi, Q^=) \leftarrow \text{inf}(\mathcal{Q}_{fin}, Q)$ ;
13     $\mathcal{Q}_{fin} := \mathcal{Q}_{fin} \setminus \{(Q_\phi, Q^=)\}$ ;
14     $\mathcal{Q}_{inf} := \mathcal{Q}_{inf} \cup \{Q_\phi \wedge Q^=\}$ ;
15  return  $(\bigvee_{(Q_\phi, Q^=) \in \mathcal{Q}_{fin}} (Q_\phi \wedge Q^=),$ 
     $\text{rb}(\bigvee_{Q_\infty \in \mathcal{Q}_{inf}} \exists \text{fv}(Q_\infty). Q_\infty))$ ;
```

Figure 5 Restricting free variables.

case $\exists x. Q_x$. First, $\text{rb}(Q_x)$ is recursively applied on Line 8 to establish the precondition of Lemma 5 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 fact entails the termination of the loop on Lines 9–12.

► **Example 6.** 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

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

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 $\neg \exists p. P(b, p) \wedge \neg S(p, u, s)$, then existentially quantifying s , and distributing the existential over disjunction yields the first disjunct in $\text{rb}(Q_{user}^{susp})$ above and $\exists s. (\neg \exists p. P(b, p)) \wedge \neg (\exists p. S(p, u, s))$ 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 $S(p, u, s)$, the query $\exists s. \neg (\exists p. S(p, u, s))$ is a tautology over \mathcal{D} . Hence, $\exists s. (\neg \exists p. P(b, p)) \wedge \neg (\exists p. S(p, u, s))$ is inf-equivalent to $\neg \exists p. P(b, p)$, i.e., the second disjunct in $\text{rb}(Q_{user}^{susp})$. This reasoning is used by $(\star\exists)$ applied to restrict s in $\exists s. \neg \exists p. P(b, p) \wedge \neg S(p, u, s)$.

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 the main properties of our translation FV and EVAL hold. Our translation is based on the following lemma.

► **Lemma 7.** *Let a structure \mathcal{S} with an infinite domain \mathcal{D} be fixed. 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 never satisfied, then*

$$\tilde{Q} \equiv (\tilde{Q} \wedge \text{qps}^\vee(\mathcal{G})) \vee \bigvee_{y \in \text{eqs}(x, \mathcal{G})} (\tilde{Q}[x \mapsto y] \wedge x \approx y), \quad (\star)$$

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 never satisfied, 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 $\text{qps}(\mathcal{G})$. ◀

Our approach is implemented by the function $\text{split}(Q)$ defined in Figure 5. In the following, we describe this function and informally 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 query pairs and a set \mathcal{Q}_{inf} of queries such that

$$Q_{fin} := \bigvee_{(Q_\phi, Q^=) \in \mathcal{Q}_{fin}} (Q_\phi \wedge Q^=), \quad Q_{inf} := \bigvee_{Q_\infty \in \mathcal{Q}_{inf}} \exists \vec{fv}(Q_\infty) \cdot Q_\infty,$$

and, for every $(Q_\phi, Q^=) \in \mathcal{Q}_{fin}$, $Q^=$ is a conjunction of equalities. As long as there exists some $(Q_{fix}, Q^=) \in \mathcal{Q}_{fin}$ 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 $Q_{fix} \wedge Q^=$, the loop on Lines 5–10 might not terminate. Note that, for every $(Q'_{fix}, Q'^=)$ 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 5–10. Finally, if $\llbracket Q_{fix} \rrbracket$ is an infinite set of tuples, then $\llbracket Q_{fix} \wedge Q^= \rrbracket$ is an infinite set of tuples, too. This is because the equalities in $Q^=$ merely duplicate columns of the query Q_{fix} . Hence, it indeed suffices to apply (\star) to Q_{fix} instead of $Q_{fix} \wedge Q^=$.

After the loop on Lines 5–10 in Figure 5 terminates, for every $(Q_\phi, Q^=) \in \mathcal{Q}_{fin}$, Q_ϕ is a safe-range query and $Q^=$ is a conjunction of equalities such that $\text{fv}(Q_\phi \wedge Q^=) = \text{fv}(Q)$. However, the query $Q_\phi \wedge Q^=$ need not be safe-range, e.g., if $Q_\phi := B(x)$ and $Q^= := (x \approx y \wedge u \approx v)$. Given a set of equalities $Q^=$, let $\text{classes}(Q^=)$ be the set of equivalence classes of free variables $\text{fv}(Q^=)$ with respect to $Q^=$. For instance, $\text{classes}(\{x \approx y, y \approx z, u \approx v\}) = \{\{x, y, z\}, \{u, v\}\}$. Let $\text{disjointvars}(Q_\phi, Q^=) := \bigcup_{V \in \text{classes}(\text{flat}^\wedge(Q^=)), V \cap \text{fv}(Q_\phi) = \emptyset} V$ be the set of all variables in equivalence classes from $\text{classes}(\text{flat}^\wedge(Q^=))$ that are disjoint from Q_ϕ 's free variables. Then, $Q_\phi \wedge Q^=$ is safe-range if and only if $\text{disjointvars}(Q_\phi, Q^=) = \emptyset$ (recall the definition of safe-range).

Now if $\text{disjointvars}(Q_\phi, Q^=) \neq \emptyset$ and $Q_\phi \wedge Q^=$ is satisfied by some tuple, then $\llbracket Q_\phi \wedge Q^= \rrbracket$ is an infinite set of tuples because all equivalence classes of variables in $\text{disjointvars}(Q_\phi, Q^=) \neq \emptyset$ can be assigned arbitrary values from the infinite domain \mathcal{D} . In our example with $Q_\phi := B(x)$ and $Q^= := (x \approx y \wedge u \approx v)$, we have $\text{disjointvars}(Q_\phi, Q^=) = \{u, v\} \neq \emptyset$. Moreover, if $\text{fv}(Q_\phi \wedge Q^=) \neq \text{fv}(Q)$ and $Q_\phi \wedge Q^=$ 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 $\text{fv}(Q) \setminus \text{fv}(Q_\phi \wedge Q^=)$. Hence, for

every $(Q_\infty, Q^\infty) \in \mathcal{Q}_{fin}$ with $\text{disjointvars}(Q_\infty, Q^\infty) \neq \emptyset$ or $\text{fv}(Q_\infty \wedge Q^\infty) \neq \text{fv}(Q)$, we remove (Q_∞, Q^∞) from \mathcal{Q}_{fin} and add $Q_\infty \wedge Q^\infty$ to \mathcal{Q}_{inf} . Note that we only remove pairs from \mathcal{Q}_{fin} , which entails the termination of the loop on Lines 11–14. Afterwards, the query Q_{fin} is safe-range. However, the query Q_{inf} need not be safe-range. Indeed, every query $Q_\infty \in \mathcal{Q}_{inf}$ has range-restricted bound variables, but not all the free variables of Q_∞ need be range-restricted and thus the query $\exists \vec{v}(Q_\infty). Q_\infty$ need not be safe-range. But the query Q_{inf} is closed and thus the inf-equivalent query $\text{rb}(Q_{inf})$ with range-restricted bound variables is safe-range.

► **Lemma 8.** *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; $\text{fv}(Q_{fin}) = \text{fv}(Q)$ unless Q_{fin} is syntactically equal to \perp ; and $\text{fv}(Q_{inf}) = \emptyset$.*

► **Lemma 9.** *Let a structure \mathcal{S} with an infinite domain \mathcal{D} be fixed. Let Q be an RC query and $\text{split}(Q) = (Q_{fin}, Q_{inf})$. If $\not\models Q_{inf}$, then $\llbracket Q \rrbracket = \llbracket Q_{fin} \rrbracket$ is a finite set. If $\models Q_{fin}$, then $\llbracket Q \rrbracket$ is an infinite set.*

► **Example 10.** Consider the query $Q_{user}^{susp} := B(b) \wedge \exists s. \forall p. P(b, p) \rightarrow S(p, u, s)$ from Section 1. Restricting its bound variables yields the query $\text{rb}(Q_{user}^{susp}) = B(b) \wedge ((\exists s. (\neg \exists p. P(b, p) \wedge \neg S(p, u, s)) \wedge (\exists p. S(p, u, s))) \vee (\neg \exists p. P(b, p)))$ derived in Example 6. Splitting Q_{user}^{susp} yields

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

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

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

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

► **Definition 11.** *Let Q be an RC query and $\text{split}(Q) = (Q_{fin}, Q_{inf})$. Let $\hat{Q}_{fin} := \text{sr2ranf}(Q_{fin})$, $\hat{Q}_{inf} := \text{sr2ranf}(Q_{inf})$ be the equivalent RANF queries. We define $\text{rw}(Q) := (\hat{Q}_{fin}, \hat{Q}_{inf})$.*

4.4 Complexity Analysis

In this section, we analyze the time complexity of capturing Q , i.e., checking if $\llbracket Q \rrbracket$ is finite and enumerating $\llbracket Q \rrbracket$ if it is finite. To bound the asymptotic time complexity of capturing a fixed Q , we ignore the (constant) time complexity of computing $\text{rw}(Q) = (\hat{Q}_{fin}, \hat{Q}_{inf})$ and focus on the time complexity of evaluating the RANF queries \hat{Q}_{fin} and \hat{Q}_{inf} , i.e., the query cost of \hat{Q}_{fin} and \hat{Q}_{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. Nevertheless, the RANF queries \hat{Q}_{fin} and \hat{Q}_{inf} computed by our translation need not have pairwise distinct (free and bound) variables.

Let $\text{av}(Q)$ be the set of all (free and bound) variables in a query Q . 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 variable substitution $x \mapsto y$,

where x, 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 $\overline{\text{qps}}(Q)$ as follows:

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

When restricting a variable by $\text{qps}^\vee(\mathcal{G})$, we only introduce quantified predicates $Q_{qp} \in \overline{\text{qps}}(Q)$.

We bound the complexity of capturing Q by considering subsets \mathcal{Q}_{qps} of quantified predicates $\overline{\text{qps}}(Q)$ 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. [22], we use the notation $\tilde{O}(\cdot)$ to hide logarithmic factors incurred by set operations.

► **Theorem 12.** *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 $\llbracket Q \rrbracket$ if it is finite, is in $\tilde{O}\left(\sum_{\mathcal{Q}_{qps} \subseteq \overline{\text{qps}}(Q), \text{minimal}(\mathcal{Q}_{qps})} \prod_{Q_{qp} \in \mathcal{Q}_{qps}} |\llbracket Q_{qp} \rrbracket|\right)$.*

We prove Theorem 12 in Appendix D. Examples 13 and 14 show that the time complexity from Theorem 12 can be achieved neither for the translation of Van Gelder and Topor [14] nor for finite domains. Example 15 shows how equalities affect the bound in Theorem 12.

► **Example 13.** Consider the query $Q := \text{B}(b) \wedge \exists u, s. \neg \exists p. \text{P}(b, p) \wedge \neg \text{S}(p, u, s)$, equivalent to Q^{susp} from Section 1. Then $\text{aps}(Q) = \{\text{B}(b), \text{P}(b, p), \text{S}(p, u, s)\}$ and $\overline{\text{qps}}(Q) = \{\text{B}(b), \text{P}(b, p), \exists p. \text{P}(b, p), \text{S}(p, u, s), \exists p. \text{S}(p, u, s), \exists s, p. \text{S}(p, u, s), \exists u, s, p. \text{S}(p, u, s)\}$. The translated query Q_{vgt} by Van Gelder and Topor [14] restricts the variables r and s by $\exists s, p. \text{S}(p, u, s)$ and $\exists u, p. \text{S}(p, u, s)$, respectively. Computing the join of $\exists s, p. \text{S}(p, u, s)$ and $\exists u, p. \text{S}(p, u, s)$, which is a Cartesian product, results in a time complexity $\Omega(|\llbracket \text{S}(p, u, s) \rrbracket|^2)$ for Q_{vgt} . Finally, we observe that $|\llbracket \text{S}(p, u, s) \rrbracket|^2$ can be arbitrarily larger than the bound in Theorem 12:

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

► **Example 14.** The query $\neg \text{S}(x, y, z)$ is satisfied by a finite set of tuples over a finite domain \mathcal{D} (as is every other query over a finite domain). The number of satisfying tuples is

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

although the bound in Theorem 12 is $\tilde{O}(|\llbracket \text{S}(x, y, z) \rrbracket|)$.

► **Example 15.** Consider the following query over the domain $\mathcal{D} = \mathbb{N}$ of natural numbers:

$$Q := \forall u. (u = 0 \vee u = 1 \vee u = 2) \longrightarrow (\exists v. \text{B}(v) \wedge (u = 0 \longrightarrow x = v) \wedge (u = 1 \longrightarrow y = v) \wedge (u = 2 \longrightarrow z = v)).$$

Note that this query is equivalent to $Q \equiv \text{B}(x) \wedge \text{B}(y) \wedge \text{B}(z)$ and thus it is satisfied by a finite set of tuples of size $|\llbracket \text{B}(x) \rrbracket| \cdot |\llbracket \text{B}(y) \rrbracket| \cdot |\llbracket \text{B}(z) \rrbracket| = |\llbracket \text{B}(x) \rrbracket|^3$. The set of atomic predicates of Q is $\text{aps}(Q) = \{\text{B}(v)\}$ and it must be closed under the equalities occurring in Q to yield a valid bound in Theorem 12. In this case, $\overline{\text{qps}}(Q) = \{\text{B}(v), \exists v. \text{B}(v), \text{B}(x), \text{B}(y), \text{B}(z)\}$ and the bound in Theorem 12 is $|\llbracket \text{B}(v) \rrbracket| \cdot |\llbracket \text{B}(x) \rrbracket| \cdot |\llbracket \text{B}(y) \rrbracket| \cdot |\llbracket \text{B}(z) \rrbracket| = |\llbracket \text{B}(x) \rrbracket|^4$. In particular, the 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 $Q \equiv \text{B}(x) \wedge \text{B}(y) \wedge \text{B}(z)$.

5 Implementation and Empirical Evaluation

We have implemented our translation RC2SQL consisting of roughly 1000 lines of OCaml code [25]. Although our translation satisfies the worst-case complexity bound (Theorem 12), we further improve its average-case complexity by implementing the following optimizations.

- We use a sample structure of constant size, called a *training database*, to estimate the query cost when resolving the nondeterministic choices in our algorithms (Appendix E.1). 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 (Section 4.3 and 4.4) for every choice of the training database. Later we describe the training databases used in our empirical evaluation.
- We use the function `optcnt` optimizing RANF subqueries of the form $\exists \vec{y}. Q^+ \wedge \bigwedge_{i=1}^k \neg Q_i^-$ using the count aggregation operator (Appendix E.2). Inspired by Claußen et al. [9], we compare the number of valuations of \vec{y} that satisfy Q^+ and $\bigvee_{i=1}^k (Q^+ \wedge Q_i^-)$, respectively.
- To compute an SQL query from a RANF query, we define the function `ranf2sql`(\cdot) (Appendix E.3). We first obtain an equivalent RA expression using the standard approach [1] but adjusting the case of closed queries [8]. To translate RA expressions into SQL, we reuse a publicly available RA interpreter `radb` [29]. We modify its implementation to improve the performance of the resulting SQL query. We map the anti-join operator $\hat{Q}_1 \triangleright \hat{Q}_2$ to a more efficient LEFT JOIN, if $\text{fv}(\hat{Q}_2) \subsetneq \text{fv}(\hat{Q}_1)$, and we perform common subquery elimination.

Appendix E provides more details on the implementation of our translation and optimizations.

To validate our translation’s improved asymptotic time complexity, we compare it with the translation by Van Gelder and Topor [14] (VGT), an implementation of the algorithm by Ailamazyan et al. [2] that uses an extended active domain as the generators, and the DDD [20,21], LDD [7], and `MonPoly`^{REG} [4] tools that support direct RC query evaluation using binary decision diagrams. 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 `rb`(\cdot) in Figure 4 to use the relation `convgt`(x, Q, \mathcal{G}) (Appendix A, Figure 8) instead of `cov`(x, Q, \mathcal{G}) (Figure 3) and to use the generator $\exists \text{fv}(Q) \setminus \{x\}. \text{qps}^\vee(\mathcal{G})$ instead of $\text{qps}^\vee(\mathcal{G})$. Evaluable queries Q are always translated into (Q_{fin}, \perp) by `rw`(\cdot) because all of Q ’s free variables are range-restricted. We also consider translation variants that omit the count aggregation optimization `optcnt`(\cdot), marked with a minus ($-$).

SQL queries computed by the translations are evaluated using the PostgreSQL database engine. We have also used the MySQL database engine but decided to omit it from our evaluation after discovering that it computed incorrect results for some queries. This issue was reported and subsequently confirmed by MySQL developers. We run our experiments on an Intel Core i5-4200U CPU computer with 8 GB RAM. The relations in PostgreSQL are recreated before each invocation to prevent optimizations based on caching recent query evaluation results. We provide all our experiments in an easily reproducible artifact [25].

In our experiments, we use pseudorandom structures generated by our benchmark *Data Golf* (Appendix F), which has two objectives. The first resembles the *regex golf* game’s objective [11] (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}^- , representing valuations of $\text{fv}(Q)$, Data Golf produces a pseudorandom structure \mathcal{S} , such that $\mathcal{T}^+ \subseteq \llbracket Q \rrbracket$, $\mathcal{T}^- \cap \llbracket Q \rrbracket = \emptyset$, and $\llbracket Q' \rrbracket, \llbracket \neg Q' \rrbracket$ contain at least $\min\{|\mathcal{T}^+|, |\mathcal{T}^-|\}$ tuples each, for every $Q' \sqsubseteq Q$. Data Golf assumes that Q has only constrained [14] bound variables whose generators also contain other variables, i.e., the bound variable y in every subquery $\exists y. Q_y$

23:14 Practical Relational Calculus Query Evaluation

Experiment	SMALL, Evaluable pseudorandom queries Q , $ \text{sub}(Q) = 14$, $n = 500$:									
RC2SQL	0.3	0.2	0.3	0.3	0.2	0.3	0.3	0.2	0.2	0.3
RC2SQL ⁻	0.3	0.2	148.4	0.2	0.2	0.3	0.3	0.3	5.9	0.3
VGT	31.5	6.8	4.2	2.5	40.0	9.3	2.5	2.2	11.7	2.8
VGT ⁻	33.5	4.8	118.5	6.3	10.9	22.1	32.8	10.8	12.2	21.8
DDD	9.2	2.2	RE	6.4	5.8	RE	5.1	RE	2.2	5.1
LDD	58.8	23.7	166.6	39.1	56.8	37.3	62.5	TO	16.1	55.7
MonPoly ^{REG}	64.0	31.0	143.3	57.5	67.4	54.5	72.2	171.3	33.5	71.4
Experiment	MEDIUM, Evaluable pseudorandom queries Q , $ \text{sub}(Q) = 14$, $n = 20000$:									
RC2SQL	2.7	1.5	4.0	2.2	1.5	2.9	3.4	1.7	1.1	2.7
RC2SQL ⁻	2.0	1.0	TO	1.8	1.6	2.3	2.3	1.6	TO	1.9
VGT	TO	TO	7.9	4.0	TO	TO	5.0	4.5	TO	4.9
VGT ⁻	TO	TO	TO	TO	TO	TO	TO	TO	TO	TO
Experiment	LARGE, Evaluable pseudorandom queries Q , $ \text{sub}(Q) = 14$, tool = RC2SQL:									
$n = 40000$	3.5	2.9	8.2	4.2	3.0	5.6	6.6	3.5	1.9	5.8
$n = 80000$	7.6	5.9	17.9	8.5	6.3	12.0	14.1	7.9	4.0	12.1
$n = 120000$	10.6	8.9	25.4	12.1	8.8	16.7	21.1	10.7	7.4	17.1
Experiment	INFINITE, Non-evaluable pseudorandom queries Q , $ \text{sub}(Q) = 7$, $n = 4000$:									
	Infinite results ($\gamma = 0$)					Finite results ($\gamma = 1$)				
RC2SQL	0.1	0.1	0.8	0.8	0.1	1.1	1.1	0.9	2.4	0.9
RC2SQL ⁻	0.1	0.1	0.5	0.5	0.1	0.7	0.7	0.6	TO	2.0
DDD	40.4	37.3	44.9	115.6	76.1	90.7	49.2	49.2	99.5	90.6
LDD	TO	TO	TO	TO	TO	TO	TO	TO	TO	TO
MonPoly ^{REG}	TO	TO	TO	TO	TO	TO	TO	TO	TO	TO

TO = Timeout of 300s, RE = Runtime Error

■ **Figure 6** Experiments SMALL, MEDIUM, LARGE, and INFINITE.

satisfies $\text{con}_{\text{vgt}}(y, Q_y, \mathcal{G})$ for some set \mathcal{G} such that $\{y\} \subsetneq \text{fv}(Q_{qp})$ for every $Q_{qp} \in \mathcal{G}$, and further assumes that Q has no subquery of the form $x \approx c$, no closed atomic predicate, and no repeated predicate symbols. We control the size of the structure \mathcal{S} in our experiments using a parameter $n = |\mathcal{T}^+| = |\mathcal{T}^-|$. All our experiments also used a Data Golf structure with $n = 2$ as the training database for RC2SQL and VGT. We remark that relations in a Data Golf structure have typically more than n tuples because the sets \mathcal{T}^+ , \mathcal{T}^- grow in the recursion on subqueries. We use two strategies to obtain Data Golf structures (parameter γ in Appendix F).

In the SMALL, MEDIUM, and LARGE experiments, we generate ten pseudorandom queries with a fixed size 14. The queries satisfy the Data Golf assumptions along with a few additional ones: the queries are not safe-range, have no repeated equalities, 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. The values of the parameter n for Data Golf structures are summarized in Figure 6.

The INFINITE experiment considers five pseudorandom queries Q that are *not* evaluable and $\text{rw}(Q) = (Q_{\text{fin}}, Q_{\text{inf}})$, such that $Q_{\text{inf}} \neq \perp$. 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. For each query Q , we compare the performance of our tool to tools that directly evaluate Q on pseudorandom structures generated by the two Data Golf strategies (parameter γ), which trigger infinite or finite evaluation results on the considered queries. For infinite results, our tool outputs this fact (by evaluating Q_{inf}), whereas the other tools also output a finite representation of the infinite result. For finite results, all tools produce the same output.

Figure 6 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. Each column shows evaluation times for a unique pseudorandom query.

Query	Q^{susp}		Q_{user}^{susp}		Q_{text}^{susp}		Query	Q^{susp}		Q_{user}^{susp}		Q_{text}^{susp}	
Param. n	10^3	10^4	10^3	10^4	10^3	10^4	Dataset	GC	MI	GC	MI	GC	MI
RC2SQL	2.1	2.6	3.1	3.6	6.4	8.0	RC2SQL	3.5	22.5	4.4	23.1	10.9	120.0
RC2SQL ⁻	61.6	TO	62.7	TO	387.6	TO	RC2SQL ⁻	274.1	TO	270.6	TO	TO	TO
VGT	3.6	3.0	—	—	216.6	TO	VGT	3.7	20.1	—	—	TO	TO
VGT ⁻	TO	TO	—	—	TO	TO	VGT ⁻	TO	TO	—	—	TO	TO
DDD	6.6	TO	5.6	TO	25.8	TO	DDD	84.4	TO	84.3	TO	173.9	TO
LDD	33.2	TO	35.5	TO	214.6	TO	LDD	TO	TO	TO	TO	TO	TO
MonPoly ^{REG}	50.0	TO	49.1	TO	182.4	TO	MonPoly ^{REG}	TO	TO	TO	TO	TO	TO

GC = Gift Cards dataset, MI = Musical Instruments dataset, TO = Timeout of 600s

■ **Figure 7** Experiment with the queries Q^{susp} , Q_{user}^{susp} , and Q_{text}^{susp} .

The lowest time for a query is typeset in bold. We do not report the translation time because it does not contribute to the time complexity for a fixed query. Still, RC2SQL's translation time is below 0.6 seconds in all experiments. We also omit the rows for tools that time out or crash on all queries of an experiment, e.g., Ailamazyan et al. [2]. We conclude that our translation RC2SQL significantly outperforms all other tools on all queries and scales well to higher values of n , i.e., larger relations in the Data Golf structures, on all queries.

We also evaluate the tools on the queries Q^{susp} and Q_{user}^{susp} from introduction and on the more challenging query $Q_{text}^{susp} := B(b) \wedge \exists u, s, t. \forall p. P(b, p) \rightarrow S(p, u, s) \vee T(p, u, t)$ with an additional relation T that relates user's review text (variable t) to a product. The query Q_{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 pseudorandom Data Golf structures (strategy $\gamma = 1$) and real-world structures obtained from the Amazon review dataset [23]. The real-world relations P , S , and T are obtained by projecting the respective tables from the Amazon review dataset for some chosen product categories (abbreviated GC and MI in Figure 7) and 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^{REG} only support integer data, we injectively remap the string and floating-point values from the Amazon review dataset to integers.

Figure 7 shows the empirical evaluation results: execution times on Data Golf structures (left) and execution times on structures derived from the real-world dataset for two specific product categories (right). We remark that VGT cannot handle the query Q_{user}^{susp} as it is not evaluable [14]. Our translation RC2SQL significantly outperforms most other tools (except VGT on Q^{susp}) on both pseudorandom and real-world structures. VGT⁻ translates Q^{susp} into a RANF query with a higher query cost than RC2SQL⁻. However, the optimization `optcnt(·)` manages to rectify this inefficiency and thus VGT exhibits a comparable performance as RC2SQL. Specifically, the factor of $80\times$ in query cost between VGT⁻ and RC2SQL⁻ improves to $1.1\times$ in query cost between VGT and RC2SQL on a Data Golf structure with $n = 20$ [25]. The factor of $1.2\times$ in query cost between VGT and RC2SQL is also observed on the MI dataset although PostgreSQL evaluates the query produced by VGT $0.9\times$ faster. Nevertheless, VGT does not finish evaluating the query Q_{text}^{susp} within 10 minutes, unlike RC2SQL.

6 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.

References

- 1 Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995. URL: <http://webdam.inria.fr/Alice/>.
- 2 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. *Dokl. Akad. Nauk SSSR*, 286:308–311, 1986. URL: <http://mi.mathnet.ru/dan47310>.
- 3 Arnon Avron and Yoram Hirshfeld. On first order database query languages. In *LICS, July 15-18, 1991, Amsterdam, The Netherlands*, pages 226–231. IEEE Computer Society, 1991. doi:10.1109/LICS.1991.151647.
- 4 David A. Basin, Felix Klaedtke, Samuel Müller, and Eugen Zalinescu. Monitoring metric first-order temporal properties. *J. ACM*, 62(2):15:1–15:45, 2015. doi:10.1145/2699444.
- 5 Michael Benedikt and Leonid Libkin. Relational queries over interpreted structures. *J. ACM*, 47(4):644–680, 2000. doi:10.1145/347476.347477.
- 6 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.
- 7 Sagar Chaki, Arie Gurfinkel, and Ofer Strichman. Decision diagrams for linear arithmetic. In *FMCAD, 15-18 November 2009, Austin, Texas, USA*, pages 53–60. IEEE, 2009. doi:10.1109/FMCAD.2009.5351143.
- 8 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.
- 9 Jens Claßen, 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, August 25-29, 1997, Athens, Greece*, pages 286–295. Morgan Kaufmann, 1997. URL: <http://www.vldb.org/conf/1997/P286.PDF>.
- 10 E. F. Codd. Relational completeness of data base sublanguages. *Research Report / RJ / IBM / San Jose, California*, RJ987, 1972.
- 11 Erling Ellingsen. Regex golf, 2013. <https://alf.nu/RegexGolf>.
- 12 Martha Escobar-Molano, Richard Hull, and Dean Jacobs. Safety and translation of calculus queries with scalar functions. In Catriel Beeri, editor, *PODS, May 25-28, 1993, Washington, DC, USA*, pages 253–264. ACM Press, 1993. doi:10.1145/153850.153909.
- 13 Allen Van Gelder and Rodney W. Topor. Safety and correct translation of relational calculus formulas. In Moshe Y. Vardi, editor, *PODS, March 23-25, 1987, San Diego, California, USA*, pages 313–327. ACM, 1987. doi:10.1145/28659.28693.
- 14 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.
- 15 Richard Hull and Jianwen Su. Domain independence and the relational calculus. *Acta Informatica*, 31(6):513–524, 1994. doi:10.1007/BF01213204.
- 16 Michael Kifer. On safety, domain independence, and capturability of database queries (preliminary report). In Catriel Beeri, Joachim W. Schmidt, and Umeshwar Dayal, editors, *Proceedings of the Third International Conference on Data and Knowledge Bases: Improving Usability and Responsiveness, June 28-30, 1988, Jerusalem, Israel*, pages 405–415. Morgan Kaufmann, 1988. doi:10.1016/b978-1-4832-1313-2.50037-8.
- 17 Nils Klarlund and Anders Møller. *MONA v1.4 User Manual*. BRICS, Department of Computer Science, University of Aarhus, January 2001. URL: <http://www.brics.dk/mona/>.
- 18 Leonid Libkin. *Elements of Finite Model Theory*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2004. URL: <http://www.cs.toronto.edu/~7Elibkin/fmt>, doi:10.1007/978-3-662-07003-1.
- 19 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.

- 640 20 Jesper B. Møller. DDDLlib: A library for solving quantified difference inequalities. In Andrei
 641 Voronkov, editor, *CADE, July 27-30, 2002, Copenhagen, Denmark*, volume 2392 of *Lecture*
 642 *Notes in Computer Science*, pages 129–133. Springer, 2002. doi:10.1007/3-540-45620-1_9.
- 643 21 Jesper B. Møller, Jakob Lichtenberg, Henrik Reif Andersen, and Henrik Hulgaard. Difference
 644 decision diagrams. In Jörg Flum and Mario Rodríguez-Artalejo, editors, *CSL, September*
 645 *20-25, 1999, Madrid, Spain*, volume 1683 of *Lecture Notes in Computer Science*, pages 111–125.
 646 Springer, 1999. doi:10.1007/3-540-48168-0_9.
- 647 22 Hung Q. Ngo, Christopher Ré, and Atri Rudra. Skew strikes back: new developments in the
 648 theory of join algorithms. *SIGMOD Rec.*, 42(4):5–16, 2013. doi:10.1145/2590989.2590991.
- 649 23 Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying recommendations using distant-
 650 labeled reviews and fine-grained aspects. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun
 651 Wan, editors, *EMNLP, November 3-7, 2019, Hong Kong, China*, pages 188–197. Association
 652 for Computational Linguistics, 2019. doi:10.18653/v1/D19-1018.
- 653 24 Robert A. Di Paola. The recursive unsolvability of the decision problem for the class of definite
 654 formulas. *J. ACM*, 16(2):324–327, 1969. doi:10.1145/321510.321524.
- 655 25 Martin Raszyk, David Basin, Srđan Krstić, and Dmitriy Traytel. Implementation and
 656 evaluation associated with this paper, 2021. <https://github.com/rc2sql/rc2sql>.
- 657 26 Peter Z. Revesz. *Introduction to Constraint Databases*. Texts in Computer Science. Springer,
 658 2002. doi:10.1007/b97430.
- 659 27 Moshe Y. Vardi. The decision problem for database dependencies. *Inf. Process. Lett.*, 12(5):251–
 660 254, 1981. doi:10.1016/0020-0190(81)90025-9.
- 661 28 Moshe Y. Vardi. The complexity of relational query languages (extended abstract). In
 662 Harry R. Lewis, Barbara B. Simons, Walter A. Burkhard, and Lawrence H. Landweber,
 663 editors, *STOC, May 5-7, 1982, San Francisco, California, USA*, pages 137–146. ACM, 1982.
 664 doi:10.1145/800070.802186.
- 665 29 Jun Yang. radb, 2019. <https://github.com/junyang/radb>.

666 A Evaluable Queries

667 The classes of *evaluable* queries [14, Def. 5.2] and *allowed* queries [14, Def. 5.3] are decidable
 668 subsets of domain-independent RC queries. The evaluable queries characterize exactly the
 669 domain-independent queries with no repeated predicate symbols [14, Theorem 10.5]. Every
 670 evaluable query can be translated to an equivalent allowed query [14, Theorem 8.6]. Every
 671 allowed query can be translated to an equivalent RANF query [14, Theorem 9.6].

672 ► **Definition 16.** A query Q is called evaluable if
 673 ■ every variable $x \in \text{fv}(Q)$ satisfies $\text{gen}_{\text{vgt}}(x, Q)$ and
 674 ■ the bound variable y in every subquery $\exists y. Q_y$ of Q satisfies $\text{con}_{\text{vgt}}(y, Q_y)$.
 675 A query Q is called allowed if
 676 ■ every variable $x \in \text{fv}(Q)$ satisfies $\text{gen}_{\text{vgt}}(x, Q)$ and
 677 ■ the bound variable y in every subquery $\exists y. Q_y$ of Q satisfies $\text{gen}_{\text{vgt}}(y, Q_y)$,
 678 where the relation $\text{gen}_{\text{vgt}}(x, Q)$ is defined to hold iff there exists a set \mathcal{G} such that $\text{gen}_{\text{vgt}}(x, Q, \mathcal{G})$
 679 and the relation $\text{con}_{\text{vgt}}(x, Q)$ is defined to hold iff there exists a set \mathcal{G} such that $\text{con}_{\text{vgt}}(x, Q, \mathcal{G})$,
 680 respectively. The relations $\text{gen}_{\text{vgt}}(x, Q, \mathcal{G})$ and $\text{con}_{\text{vgt}}(x, Q, \mathcal{G})$ are defined in Figure 8.

681 In Figure 9 we introduce a measure $\text{measure}(Q)$ on queries, that decreases for proper
 682 subqueries, after pushing negation, and after distributing existential quantification over
 683 disjunction. Hence, the termination of the rules in Figures 2, 3, and 8 and the termination
 684 of the functions in Figures 13 and 14 follow using the measure $\text{measure}(Q)$.

685 We relate the definitions from Figure 2 and Figure 8 with the following lemmas.

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

■ **Figure 8** The relations $\text{gen}_{\text{vgt}}(x, Q, \mathcal{G})$ and $\text{con}_{\text{vgt}}(x, Q, \mathcal{G})$ [14].

686 ► **Lemma 17.** *Let x and y be free variables in a query Q such that $\text{gen}_{\text{vgt}}(x, \neg Q)$ and*
 687 *$\text{gen}_{\text{vgt}}(y, Q)$ hold. Then we get a contradiction.*

688 **Proof.** The lemma is proved by induction on the query Q using the measure $\text{measure}(Q)$ on
 689 queries defined in Figure 9, which decreases in every case of the definition in Figure 8. ◀

690 ► **Lemma 18.** *Let Q be a query such that $\text{gen}_{\text{vgt}}(y, Q_y)$ holds for the bound variable y in*
 691 *every subquery $\exists y. Q_y$ of Q . Suppose that $\text{gen}_{\text{vgt}}(x, Q)$ holds for a free variable $x \in \text{fv}(Q)$.*
 692 *Then $\text{gen}(x, Q)$ holds.*

693 **Proof.** The lemma is proved by induction on the query Q using the measure $\text{measure}(Q)$ on
 694 queries defined in Figure 9, which decreases in every case of the definition in Figure 8.

695 Lemma 17 and the assumption that $\text{gen}_{\text{vgt}}(y, Q_y)$ holds for the bound variable y in
 696 every subquery $\exists y. Q_y$ of Q imply that $\text{gen}_{\text{vgt}}(x, Q)$ cannot be derived using the rule
 697 $\text{gen}_{\text{vgt}}(x, \neg\exists y. Q_y)$, i.e., Q cannot be of the form $\neg\exists y. Q_y$. Every other case in the definition
 698 of $\text{gen}_{\text{vgt}}(x, Q)$ has a corresponding case in the definition of $\text{gen}(x, Q)$. ◀

699 ► **Lemma 19.** *Let Q be an allowed query, i.e., $\text{gen}_{\text{vgt}}(x, Q)$ holds for every free variable*
 700 *$x \in \text{fv}(Q)$ and $\text{gen}_{\text{vgt}}(y, Q_y)$ holds for the bound variable y in every subquery $\exists y. Q_y$ of Q .*
 701 *Then Q is a safe-range query, i.e., $\text{gen}(x, Q)$ holds for every free variable $x \in \text{fv}(Q)$ and*
 702 *$\text{gen}(y, Q_y)$ holds for the bound variable y in every subquery $\exists y. Q_y$ of Q .*

703 **Proof.** The lemma is proved by applying Lemma 18 to every free variable of Q and to the
 704 bound variable y in every subquery of Q of the form $\exists y. Q_y$. ◀

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

■ **Figure 9** The measure $\text{measure}(Q)$ on RC queries.

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

■ **Figure 10** Constant propagation rules.

705 Lemma 19 shows that every allowed query is safe-range. But there exist safe-range queries
 706 that are not allowed, e.g., $B(x) \wedge x \approx y$.

707 B Constant Propagation

708 We introduce constant propagation rules in Figure 10. We denote by $\text{cp}(Q)$ the query
 709 obtained from a query Q by exhaustively applying the rules in Figure 10. Note that $\text{cp}(Q)$ is
 710 either of the form \perp or \top or contains neither \perp nor \top as a subquery.

711 The following definitions introduce substitution of a variable by another variable and
 712 removing all free occurrences of a free variable.

713 ► **Definition 20.** *The substitution of the form $Q[x \mapsto y]$ is the query $\text{cp}(Q')$ where Q' is*
 714 *obtained from a query Q by replacing all occurrences of the free variable x by the variable y ,*
 715 *potentially also renaming bound variables to avoid capture.*

716 ► **Definition 21.** *The substitution of the form $Q[x/\perp]$ is the query $\text{cp}(Q')$ where Q' is*
 717 *obtained from a query Q by replacing with \perp every atomic predicate or equality containing*
 718 *the free variable x , except for $(x \approx x) \equiv \top$.*

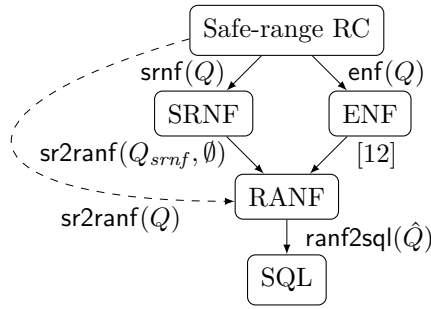
719 C Query Normal Forms

720 In this section, we introduce relational algebra normal form (RANF), which is a syntactic class
 721 of safe-range RC queries that admit a simple construction of an equivalent RA expression.
 722 We also introduce other normal forms useful for translating safe-range queries to RANF
 723 queries. Note that a query normal form concerns the structure of the query rather than
 724 functional dependencies between attributes in relations (e.g., as in 1NF, 2NF, 3NF).

725 Figure 11 defines the predicate $\text{ranf}(\cdot)$ characterizing RANF queries. The translation of
 726 safe-range queries (Section 3.2) to equivalent RANF queries proceeds via safe-range normal
 727 form (SRNF) introduced in [1, Section 5.4] and summarized here in Appendix C.1. A
 728 safe-range query in SRNF can be translated to an equivalent RANF query by subquery

$\text{ranf}(\perp)$; $\text{ranf}(\top)$;
 $\text{ranf}(Q)$ 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 11** Characterization of RANF queries.



■ **Figure 12** Overview of query normal forms.

rewriting using the following rules [1, 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 sr2ranf inspired by the rules (R1)–(R3).

Existential normal form (ENF) was introduced by Van Gelder and Topor [14] to translate an allowed query [14] into an equivalent RANF query. Given a safe-range query in ENF, the rules (R1)–(R3) can be applied to obtain an equivalent RANF query [12, 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 [14]. This issue has been identified and fixed by Escobar-Molano et al. [12]. 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 [12, Fig. 2]. 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 [25]. We thus opt for using SRNF instead of ENF for translating safe-range queries into RANF.

► **Example 22.** 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 $Q_{\text{srnf}} :=$

input : A 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  $Q$ ;
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 13 Translation to 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_{enf} . The costs of the two queries over a structure \mathcal{S} are $\text{cost}^S(Q_{enf}) = 2 \cdot \|\llbracket P_2(x, y) \rrbracket\| + \|\llbracket P_1(x) \rrbracket\| + \|\llbracket P_1(y) \rrbracket\| + 2 \cdot \|\llbracket P_1(x) \wedge P_1(y) \rrbracket\| + 2 \cdot \|\llbracket Q_{enf} \rrbracket\|$ and $\text{cost}^S(Q_{srnf}) = 2 \cdot \|\llbracket P_2(x, y) \rrbracket\| + \|\llbracket P_1(x) \rrbracket\| + 2 \cdot \|\llbracket P_2(x, y) \rrbracket\| + \|\llbracket P_1(y) \rrbracket\| + 2 \cdot \|\llbracket P_2(x, y) \wedge \neg P_1(x) \rrbracket\| + 2 \cdot \|\llbracket P_2(x, y) \wedge \neg P_1(y) \rrbracket\| + 2 \cdot \|\llbracket Q_{srnf} \rrbracket\|$, respectively. Note that the cost of Q_{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_{srnf} can only be larger by a constant factor.

Figure 12 shows an overview of the RC fragments and query normal forms (nodes) and the functions we use to translate between them and to SQL (edges). The dashed edge shows the translation of a safe-range query to RANF we opt for in this paper. 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 the normal forms and translations. Appendix E.3 shows how we translate RANF queries to SQL.

C.1 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 has the form of an atomic predicate, equality, or an existentially quantified query [1]. Figure 13 defines the function $\text{srnf}(Q)$ that yields a SRNF query equivalent to a query Q . The function $\text{srnf}(Q)$ proceeds by pushing negation [1, Section 5.4], distributing existential quantifiers over disjunction [14, Rule (T9)], and dropping bound variables that never oc-

cur [14, Definition 9.2]. We include the last two rules to optimize the time complexity of evaluating the resulting RANF query. The termination of the function $\text{srnf}(Q)$ follows using the measure $\text{measure}(Q)$, defined in Figure 9.

Using Figure 2, if a query Q is safe-range, then $\text{srnf}(Q)$ is also safe-range. Next we prove the following lemma that we use as a precondition for translating safe-range SRNF queries to RANF queries.

► **Lemma 23.** *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 $\neg Q_{\text{srnf}}$.*

Proof. Using Figure 2, $\text{gen}(x, \neg Q')$ 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. ◀

C.2 Relational Algebra Normal Form

The function $\text{sr2ranf}(Q, \mathcal{Q}) = (\hat{Q}, \overline{\mathcal{Q}})$, defined in Figure 14, where 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 and returns a RANF query \hat{Q} and a subset $\overline{\mathcal{Q}} \subseteq \mathcal{Q}$ such that $Q \wedge \bigwedge_{\overline{Q} \in \mathcal{Q}} \overline{Q} \equiv \hat{Q} \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}$. To restrict variables in Q , the function $\text{sr2ranf}(Q, \mathcal{Q})$ conjoins a subset of queries $\overline{\mathcal{Q}} \subseteq \mathcal{Q}$ 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 $(2 \cdot \text{measure}(Q) + \text{eqneg}(Q) + 2 \cdot \sum_{\overline{Q} \in \mathcal{Q}} \text{measure}(\overline{Q}) + 2 \cdot |\mathcal{Q}|, \text{measure}(Q))$, where $\text{measure}(Q)$ is defined in Figure 9, $\text{eqneg}(Q) := 1$ if Q has the form of an equality between two variables or a negation, and $\text{eqneg}(Q) := 0$ otherwise.

Next we describe the definition of $\text{sr2ranf}(Q, \mathcal{Q})$ that follows [1, Algorithm 5.4.7]. Note that no constant propagation (Appendix B) is needed in [1, Algorithm 5.4.7], because the constants \perp and \top are not in the query syntax [1, 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 \perp 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 $\text{flat}^\wedge(Q^\wedge)$ 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 Q^\wedge 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 11, i.e., a left-associative conjunction in RANF. Note that the function $\text{sort}^\wedge(\mathcal{Q})$ must place 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 clash of the free variables in the set of queries \mathcal{Q} with the bound variables \vec{v} .

Finally, the nondeterministic choices in the function $\text{sr2ranf}(Q, \mathcal{Q})$ are resolved by minimizing the cost of the resulting RANF query with respect to a training database (Appendix E.1).

D Complexity Analysis

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, let $Q_{ex} := ((\exists z. (\exists y, z. P_3(x, y, z)) \wedge P_2(y, z)) \wedge$

input : A safe-range query $Q \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}$ 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 $\overline{\mathcal{Q}} \subseteq \overline{\mathcal{Q}}$ such that $Q \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q} \equiv \hat{Q} \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}$, $\hat{Q} \implies \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}$, $\hat{Q} = \text{cp}(\hat{Q})$, and $\text{fv}(Q) \subseteq \text{fv}(\hat{Q}) \subseteq \text{fv}(Q) \cup \text{fv}(\overline{\mathcal{Q}})$ unless $\hat{Q} = \perp$.

```

1 function sr2ranf( $Q, \overline{\mathcal{Q}}$ ) =
2   if ranf( $Q$ ) then return ( $\text{cp}(Q), \emptyset$ );
3   switch  $Q$  do
4     case  $x \approx y$  do return sr2ranf( $x \approx y \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}$ ,  $\emptyset$ );
5     case  $\neg Q'$  do
6        $\overline{\mathcal{Q}} \leftarrow \{\overline{Q} \subseteq \overline{\mathcal{Q}} \mid (\neg Q') \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q} \text{ is safe-range}\}$ ;
7       if  $\overline{\mathcal{Q}} = \emptyset$  then
8          $(\hat{Q}', \_ ) := \text{sr2ranf}(Q', \emptyset)$ ;
9         return ( $\text{cp}(\neg \hat{Q}'), \emptyset$ );
10      else return sr2ranf( $(\neg Q') \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}$ ,  $\emptyset$ );
11     case  $Q_1 \vee Q_2$  do
12        $\overline{\mathcal{Q}} \leftarrow \{\overline{Q} \subseteq \overline{\mathcal{Q}} \mid \bigvee_{Q' \in \text{flat}^\vee(Q)} (Q' \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}) \text{ is safe-range}\}$ ;
13       foreach  $Q' \in \text{flat}^\vee(Q)$  do  $(\hat{Q}', \_ ) := \text{sr2ranf}(Q' \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}, \emptyset)$ ;
14       return ( $\text{cp}(\bigvee_{Q' \in \text{flat}^\vee(Q)} \hat{Q}'), \overline{\mathcal{Q}}$ );
15     case  $Q_1 \wedge Q_2$  do
16        $\mathcal{Q}^- := \{Q' \in \text{flat}^\wedge(Q) \mid \text{neg}(Q')\}$ ;  $\mathcal{Q}^+ := \text{flat}^\wedge(Q) \setminus \mathcal{Q}^-$ ;
17        $\mathcal{Q}^= := \{Q' \in \mathcal{Q}^+ \mid \text{eq}(Q')\}$ ;  $\mathcal{Q}^+ := \mathcal{Q}^+ \setminus \mathcal{Q}^=$ ;
18        $\mathcal{Q}^\neq := \{\neg Q' \in \mathcal{Q}^- \mid \text{eq}(Q')\}$ ;  $\mathcal{Q}^- := \mathcal{Q}^- \setminus \mathcal{Q}^\neq$ ;
19       foreach  $Q' \in \mathcal{Q}^+$  do  $(\hat{Q}', \mathcal{Q}_{Q'}) := \text{sr2ranf}(Q', \mathcal{Q} \cup ((\mathcal{Q}^+ \cup \mathcal{Q}^=) \setminus \{Q'\}))$ ;
20       foreach  $\neg Q' \in \mathcal{Q}^-$  do  $(\hat{Q}', \_ ) := \text{sr2ranf}(Q', \mathcal{Q} \cup (\mathcal{Q}^+ \cup \mathcal{Q}^=))$ ;
21        $\overline{\mathcal{Q}} \leftarrow \{\overline{Q} \subseteq \overline{\mathcal{Q}} \mid \mathcal{Q}^+ \subseteq \bigcup_{Q' \in \overline{\mathcal{Q}}} (\mathcal{Q}_{Q'} \cup \{Q'\})\}$ ;
22       return ( $\text{cp}(\text{sort}^\wedge(\bigcup_{Q' \in \overline{\mathcal{Q}}} \hat{Q}' \cup \mathcal{Q}^= \cup \bigcup_{\neg Q' \in \mathcal{Q}^-} \neg \hat{Q}' \cup \mathcal{Q}^\neq), \bigcup_{Q' \in \overline{\mathcal{Q}}} (\mathcal{Q}_{Q'} \cap \mathcal{Q}))$ );
23     case  $\exists \vec{v}. Q_{\vec{v}}$  do
24       if  $\text{fv}(\overline{\mathcal{Q}}) \cap \vec{v} \neq \emptyset$  then  $\vec{w} \leftarrow \{\vec{w} \mid |\vec{w}| = |\vec{v}| \text{ and } ((\text{fv}(Q_{\vec{v}}) \setminus \vec{v}) \cup \text{fv}(\overline{\mathcal{Q}})) \cap \vec{w} = \emptyset\}$ ;
25       else  $\vec{w} := \vec{v}$ ;
26        $Q_{\vec{w}} := Q_{\vec{v}}[\vec{v} \mapsto \vec{w}]$ ;
27        $\overline{\mathcal{Q}} \leftarrow \{\overline{Q} \subseteq \overline{\mathcal{Q}} \mid Q_{\vec{w}} \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q} \text{ is safe-range}\}$ ;
28        $(\hat{Q}_{\vec{w}}, \_ ) := \text{sr2ranf}(Q_{\vec{w}} \wedge \bigwedge_{\overline{Q} \in \overline{\mathcal{Q}}} \overline{Q}, \emptyset)$ ;
29       return ( $\text{cp}(\exists \vec{w}. \hat{Q}_{\vec{w}}), \overline{\mathcal{Q}}$ );
30   otherwise do return ( $\text{cp}(Q), \emptyset$ );

```

■ **Figure 14** Translation of safe-range SRNF to RANF.

811 $P_1(z)) \vee P_3(x, y, z)$. Then $\text{aps}(Q_{ex}) = \{P_1(z), P_2(y, z), P_3(x, y, z)\}$ and $\mathcal{B}_{Q_{ex}}(P_3(x, y, z)) =$
812 $\{yz, []\}$, where $[]$ denotes the empty sequence corresponding to the occurrence of $P_3(x, y, z)$
813 in Q_{ex} for which the variables x, y, z are all free in Q_{ex} . Note that the variable z in
814 the other occurrence of $P_3(x, y, z)$ in Q_{ex} is bound to the innermost quantifier. Hence,
815 neither zy nor zyz is in $\mathcal{B}_{Q_{ex}}(P_3(x, y, z))$. Furthermore, let $\text{qps}(Q)$ be the set of the quan-
816 tified predicates obtained by existentially quantifying sequences of bound variables in
817 $\mathcal{B}_{Q'}(Q_{ap})$ from the atomic predicates $Q_{ap} \in \text{aps}(Q')$ in all subqueries Q' of Q . Form-

ally, $\text{qps}(Q) := \bigcup_{Q' \sqsubseteq Q, Q_{ap} \in \text{aps}(Q')} \{\exists \vec{v}. Q_{ap} \mid \vec{v} \in \mathcal{B}_{Q'}(Q_{ap})\}$. For instance, $\text{qps}(Q_{ex}) = \{P_3(x, y, z), \exists z. P_3(x, y, z), \exists yz. P_3(x, y, z), P_2(y, z), \exists z. P_2(y, z), P_1(z)\}$.

A crucial property of our translation that is central for the complexity analysis (Theorem 12) is formalized in the following lemma.

► **Lemma 24.** *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' .

Suppose that $Q' \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \bar{Q}$ 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 $\exists x. Q'_x \sqsubseteq Q' \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \bar{Q}$, and every two subqueries $\exists x. Q_x \sqsubseteq Q_1$ and $\exists x. Q'_x \sqsubseteq Q_2$ such that $Q_1 \wedge Q_2 \sqsubseteq Q' \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \bar{Q}$ have the property that $\exists x. Q_x$ or $\exists x. Q'_x$ is a quantified predicate. Then the free variables in $\bigwedge_{\bar{Q} \in \mathcal{Q}} \bar{Q}$ never clash with the bound variables in Q' , i.e., Line 24 in Figure 14 is never executed. Next we observe that $\mathcal{B}_{\text{sr2ranf}(Q', \mathcal{Q})}(Q_{ap}) \subseteq \mathcal{B}_{Q' \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \bar{Q}}(Q_{ap})$ and then $\text{qps}(\text{sr2ranf}(Q', \mathcal{Q})) \subseteq \text{qps}(Q' \wedge \bigwedge_{\bar{Q} \in \mathcal{Q}} \bar{Q})$. Because Q_{fin}, Q_{inf} have the properties from the beginning of this paragraph 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 $\text{qps}(\hat{Q}_{inf}) = \text{qps}(\text{sr2ranf}(Q_{inf})) \subseteq \text{qps}(Q_{inf})$. ◀

Recall Example 13. The query $\exists u, p. S(p, u, s)$ is in $\text{qps}(Q_{vgt})$, but not in $\overline{\text{qps}}(Q)$. Hence, $\text{qps}(Q_{vgt}) \not\subseteq \overline{\text{qps}}(Q)$, i.e., an analogue of Lemma 24 for Van Gelder and Topor's translation, does not hold.

Let a structure \mathcal{S} be fixed. We observe that every tuple satisfying a RANF query \hat{Q} belongs to the set of tuples satisfying the join over some minimal subset $\mathcal{Q}_{qps} \subseteq \text{qps}(\hat{Q})$ of quantified predicates and also satisfying equalities duplicating some columns from \mathcal{Q}_{qps} . 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 $y \in V'$.

► **Lemma 25.** *Let \hat{Q} be a RANF query. Then $\llbracket \hat{Q} \rrbracket$ satisfies*

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

Proof. The statement is proved by well-founded induction over the inductive definition of $\text{ranf}(\hat{Q})$. ◀

Now we 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 26.** *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*

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

858 **Proof.** Applying Lemma 25 to the RANF query \hat{Q}' yields

$$859 \quad \llbracket \hat{Q}' \rrbracket \subseteq \bigcup_{\substack{\mathcal{Q}_{qps} \subseteq \text{qps}(\hat{Q}'), \text{minimal}(\mathcal{Q}_{qps}), \\ \mathcal{Q}^= \subseteq \{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in \text{fv}(\hat{Q}')\}, \\ \text{fv}(\mathcal{Q}_{qps}) \cup \text{fv}(\mathcal{Q}^=) = \text{fv}(\hat{Q}')}} \left[\bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \wedge \bigwedge_{Q^= \in \mathcal{Q}^=} Q^= \right].$$

860 We observe that $\left| \left[\bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \wedge \bigwedge_{Q^= \in \mathcal{Q}^=} Q^= \right] \right| \leq \left| \left[\bigwedge_{Q_{qp} \in \mathcal{Q}_{qps}} Q_{qp} \right] \right| \leq \prod_{Q_{qp} \in \mathcal{Q}_{qps}} \llbracket Q_{qp} \rrbracket$
 861 where the first inequality follows from the fact that equalities $Q^= \in \mathcal{Q}^=$ can only restrict
 862 a set of tuples and duplicate columns. Because \hat{Q}' is a subquery of \hat{Q} , it follows that
 863 $\text{qps}(\hat{Q}') \subseteq \text{qps}(\hat{Q})$. Lemma 24 yields $\text{qps}(\hat{Q}) \subseteq \overline{\text{qps}}(Q)$. Hence, we derive $\text{qps}(\hat{Q}') \subseteq \overline{\text{qps}}(Q)$.

864 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

$$865 \quad |\text{fv}(\mathcal{Q}_{qps})| \cdot |\text{fv}(\hat{Q}')| \leq |\text{fv}(\hat{Q}')|^2 \leq |\text{av}(\hat{Q})|^2,$$

866 where the first inequality holds because $\text{fv}(\mathcal{Q}_{qps}) \cup \text{fv}(\mathcal{Q}^=) = \text{fv}(\hat{Q}')$ and thus $\text{fv}(\mathcal{Q}_{qps}) \subseteq \text{fv}(\hat{Q}')$
 867 and the second inequality holds because the variables in a subquery \hat{Q}' of \hat{Q} are included in
 868 the set of all variables in \hat{Q} . Hence, the number of subsets $\mathcal{Q}^= \subseteq \{x \approx y \mid x \in \text{fv}(\mathcal{Q}_{qps}) \wedge y \in$
 869 $\text{fv}(\hat{Q}')\}$ is at most $2^{|\text{av}(\hat{Q})|^2}$. ◀

870 Next we bound the query cost of a RANF query $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$ over the structure \mathcal{S} .

871 ► **Lemma 27.** *Let Q be an RC query with pairwise distinct (free and bound) variables and*
 872 *let $\text{rw}(Q) = (\hat{Q}_{fin}, \hat{Q}_{inf})$. Let $\hat{Q} \in \{\hat{Q}_{fin}, \hat{Q}_{inf}\}$. Then*

$$873 \quad \text{cost}^{\mathcal{S}}(\hat{Q}) \leq |\text{sub}(\hat{Q})| \cdot |\text{av}(\hat{Q})| \cdot 2^{|\text{av}(\hat{Q})|^2} \cdot \sum_{\mathcal{Q}_{qps} \subseteq \overline{\text{qps}}(Q), \text{minimal}(\mathcal{Q}_{qps})} \prod_{Q_{qp} \in \mathcal{Q}_{qps}} \llbracket Q_{qp} \rrbracket.$$

874 **Proof.** Recall that $|\text{sub}(\hat{Q})|$ denotes the number of subqueries of the query \hat{Q} and thus
 875 bounds the number of RANF subqueries \hat{Q}' of the query \hat{Q} . For every subquery \hat{Q}' of \hat{Q} , we
 876 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
 877 we use the estimation of $|\llbracket \hat{Q}' \rrbracket|$ by Lemma 26. ◀

878 Finally, we prove Theorem 12.

879 **Proof of Theorem 12.** We derive Theorem 12 from Lemma 27 and the fact that the quant-
 880 ities $|\text{sub}(\hat{Q})|$, $|\text{av}(\hat{Q})|$, and $2^{|\text{av}(\hat{Q})|^2}$ only depend on the query Q and thus they do not
 881 contribute to the asymptotic time complexity of capturing a fixed query Q . ◀

882 E Implementation Details

883 In this section, we provide a detailed description of the implementation of our translation
 884 RC2SQL. Overall, the translation is defined as

$$885 \quad \text{RC2SQL}(Q) := (Q'_{fin}, Q'_{inf})$$

886 where

$$887 \quad \begin{aligned} Q'_{fin} &:= \text{ranf2sql}(\text{optcnt}(Q_{fin})), \\ Q'_{inf} &:= \text{ranf2sql}(\text{optcnt}(Q_{inf})), \\ (Q_{fin}, Q_{inf}) &:= \text{rw}(Q). \end{aligned}$$

Function $\text{rw}(\cdot)$ is defined in Section 4.4 as a composition of the $\text{split}(\cdot)$ and $\text{sr2ranf}(\cdot)$ functions, which are defined in Section 4.3 and Appendix C, respectively. Below we first describe how we resolve the nondeterministic choices in all our algorithms. Then we define functions $\text{optcnt}(\cdot)$ and $\text{ranf2sql}(\cdot)$.

E.1 Instantiating Our Translation

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. Nevertheless, our translation satisfies the correctness and worst-case complexity claims (Section 4.3 and 4.4) for every choice of the training database. The training databases used in our empirical evaluation are obtained using the function dg (Appendix F) with parameters described in Section 5. 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:

- Choosing some $X \in \mathcal{X}$ 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 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 $V \subseteq \text{fv}(\tilde{Q})$ is a subset of its free variables. Observe that the measure $\text{measure}(Q)$ on queries, defined in Figure 9, 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 Figure 2 and 3. 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 $\text{measure}(\tilde{Q})$ 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 $\text{cov}(x, \tilde{Q}, \mathcal{G})$.

E.2 Optimization using Count Aggregations

In this section, we introduce count aggregations and describe a generalization of Claußen et al. [9]’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 valuation α satisfies Q' iff α 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 valuations α satisfying Q_x to get valuations α 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 if $\|Q_x\| + \|Q_y\| + \|Q_{xy}\| \ll \|Q_x \wedge Q_y\|$. Furthermore,

we observe that a valuation α satisfies $Q_x \wedge \neg Q'$ if α 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 28.** *Let $\exists \vec{v}. Q_{\vec{v}} \wedge \bigwedge_{\vec{Q} \in \mathcal{Q}} \neg \vec{Q}$, $\mathcal{Q} \neq \emptyset$, be a RANF query. Let c, c' be fresh variables that do not occur in $\text{fv}(Q_{\vec{v}})$. Then*

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

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

► **Lemma 29.** *Let $\hat{Q} \wedge \neg \exists \vec{v}. Q_{\vec{v}} \wedge \bigwedge_{\vec{Q} \in \mathcal{Q}} \neg \vec{Q}$, $\mathcal{Q} \neq \emptyset$, 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*

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

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, 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 implement the translations from Lemmas 28 and 29 and mini-scoping in a function called `optcnt(\cdot)`. Given a RANF query \hat{Q} , `optcnt(\hat{Q})` is an equivalent RANF query after introducing count aggregations and performing mini-scoping.

► **Example 30.** We show how we 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:

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

E.3 Translating RANF to SQL

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` [29].

We define the function `ranf2ra(\hat{Q})` translating RANF queries \hat{Q} into equivalent RA expressions `ranf2ra(\hat{Q})`. The translation is based on Algorithm 5.4.8 by Abiteboul et al. [1], which we modify as follows. We adjust the way closed RC queries are handled. Chomicki and Toman [8] 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 [1, Algorithm 5.4.8] with queries of the form $[\text{CNT } \vec{v}. Q_{\vec{v}}](c)$.

The `radb` interpreter, abbreviated here by the function `ra2sql(\cdot)`, translates a RA expression into SQL, by simply mapping the RA constructors into their SQL counterparts. The function `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 `ranf(Q_1)`, `ranf(Q_2)`, and $\text{fv}(Q_2) \subsetneq \text{fv}(Q_1)$ is translated into RA expression `ranf2ra(Q_1)` \triangleright `ranf2ra(Q_2)`, where \triangleright denotes the anti-join operator (also referred to as the *generalized* difference operator [1]) and `ranf2ra(Q_1)`, `ranf2ra(Q_2)` are the equivalent relational algebra expressions for Q_1 , Q_2 , respectively. The `radb` interpreter only supports the anti-join operator `ranf2ra(Q_1)` \triangleright `ranf2ra(Q_2)` expressed as `ranf2ra(Q_1)` $-$ (`ranf2ra(Q_1)` \bowtie `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 [25].

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 3.3).

Finally, the function `ranf2sql(\hat{Q})` (Figure 12) is defined as `ranf2sql(\hat{Q}) := ra2sql(ranf2ra(\hat{Q}))`, i.e., as a composition of the two translations from RANF to RA and from RA to SQL.

input : A RC query Q with pairwise distinct (free and bound) variables satisfying CON, CST, VAR, REP, a sequence of distinct variables \vec{v} , $\text{fv}(Q) \subseteq \vec{v}$, sets of tuples $\mathcal{T}_{\vec{v}}^+$, $\mathcal{T}_{\vec{v}}^-$ over \vec{v} such that $\mathcal{T}_{\vec{v}}^+[x] \cap \mathcal{T}_{\vec{v}}^-[x] = \emptyset$, $|\mathcal{T}_{\vec{v}}^+[x]| \geq |\mathcal{T}_{\vec{v}}^-[x]|$, and $|\mathcal{T}_{\vec{v}}^-[x]| \geq |\mathcal{T}_{\vec{v}}^+|$, for every $x \in \vec{v}$, a parameter $\gamma \in \{0, 1\}$.

output : A structure \mathcal{S} such that $\mathcal{T}_{\vec{v}}^+[\vec{\text{fv}}(Q)] \subseteq \llbracket Q \rrbracket$, $\mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$, and $|\llbracket Q' \rrbracket|$, $|\llbracket \neg Q' \rrbracket|$ contain at least $\min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$ tuples each, 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  $r(t_1, \dots, t_{\iota(r)})$  do return  $\{r^{\mathcal{S}} \mapsto \mathcal{T}_{\vec{v}}^+[t_1, \dots, t_{\iota(r)}]\}$ ;
4     case  $x \approx y$  do
5       if there exist  $d, d'$  such that  $d \neq d'$  and  $(d, d') \in \mathcal{T}_{\vec{v}}^+[x, y]$ , or
6          $d = d'$  and  $(d, d') \in \mathcal{T}_{\vec{v}}^-[x, y]$  then
7         fail
8     case  $\neg Q'$  do return  $\text{dg}(Q', \vec{v}, \mathcal{T}_{\vec{v}}^-, \mathcal{T}_{\vec{v}}^+, \gamma)$ ;
9     case  $Q_1 \vee Q_2$  or  $Q_1 \wedge Q_2$  do
10       $(\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2) \leftarrow \{(\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2) \mid |\mathcal{T}_{\vec{v}}^1| = |\mathcal{T}_{\vec{v}}^2| = \min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\} \text{ and}$ 
11         $\mathcal{T}_{\vec{v}}^1[x] \cap \mathcal{T}_{\vec{v}}^2[x] = \emptyset \text{ and } (\mathcal{T}_{\vec{v}}^1[x] \cup \mathcal{T}_{\vec{v}}^2[x]) \cap (\mathcal{T}_{\vec{v}}^+[x] \cup \mathcal{T}_{\vec{v}}^-[x]) = \emptyset, \text{ for all } x \in \vec{v}\}$ ;
12      if  $\gamma = 0$  then
13        return  $\text{dg}(Q_1, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^2, \gamma) \cup \text{dg}(Q_2, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^2, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^1, \gamma)$ 
14      else
15        switch  $Q$  do
16          case  $Q_1 \vee Q_2$  do
17            return  $\text{dg}(Q_1, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^2, \gamma) \cup \text{dg}(Q_2, \vec{v}, \mathcal{T}_{\vec{v}}^1 \cup \mathcal{T}_{\vec{v}}^2, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^+, \gamma)$ 
18          case  $Q_1 \wedge Q_2$  do
19            return  $\text{dg}(Q_1, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^-, \mathcal{T}_{\vec{v}}^1 \cup \mathcal{T}_{\vec{v}}^2, \gamma) \cup \text{dg}(Q_2, \vec{v}, \mathcal{T}_{\vec{v}}^+ \cup \mathcal{T}_{\vec{v}}^2, \mathcal{T}_{\vec{v}}^- \cup \mathcal{T}_{\vec{v}}^1, \gamma)$ 
20      case  $\exists y. Q_y$  do
21         $(\mathcal{T}_{\vec{v}.y}^1, \mathcal{T}_{\vec{v}.y}^2) \leftarrow \{(\mathcal{T}_{\vec{v}.y}^1, \mathcal{T}_{\vec{v}.y}^2) \mid \mathcal{T}_{\vec{v}.y}^1[\vec{v}] = \mathcal{T}_{\vec{v}}^+ \text{ and } \mathcal{T}_{\vec{v}.y}^2[\vec{v}] = \mathcal{T}_{\vec{v}}^- \text{ and}$ 
22           $\mathcal{T}_{\vec{v}}^1[y] \cap \mathcal{T}_{\vec{v}}^2[y] = \emptyset \text{ and } (\mathcal{T}_{\vec{v}}^1[y] \cup \mathcal{T}_{\vec{v}}^2[y]) \cap (\mathcal{T}_{\vec{v}}^+[y] \cup \mathcal{T}_{\vec{v}}^-[y]) = \emptyset\}$ ;
23        return  $\text{dg}(Q_y, \vec{v} \cdot y, \mathcal{T}_{\vec{v}.y}^1, \mathcal{T}_{\vec{v}.y}^2, \gamma)$ 

```

Figure 15 Computing the Data Golf structure.

1012 F Data Golf Benchmark

1013 We devise the *Data Golf* benchmark for generating synthetic databases. Given a query Q ,
 1014 the generated database guarantees that no subquery Q' of Q is satisfied by (almost) all
 1015 possible tuples or (almost) no tuple at all as this might make Q 's evaluation trivial and its
 1016 benchmarking less meaningful. We first make the following assumptions on Q :

1017 CON the bound variable y in every subquery $\exists y. Q_y$ of Q satisfies $\text{con}_{\text{vgt}}(y, Q_y, \mathcal{G})$ (Figure 8)
 1018 for some set \mathcal{G} such that $\{y\} \subsetneq \text{fv}(Q_{qp})$ for every $Q_{qp} \in \mathcal{G}$, to avoid subqueries like
 1019 $\exists y. \neg P_2(x, y)$ and $\exists y. (P_2(x, y) \vee P_1(y))$;

1020 CST Q contains no subquery of the form $x \approx c$ because such a subquery is satisfied by
 1021 exactly one tuple;

1022 VAR Q contains no closed atomic predicate, e.g., $P_1(42)$, because a closed subquery is either
 1023 satisfied by all possible tuples or no tuple at all; and

1024 REP Q contains no repeated predicate symbols, to avoid subqueries like $P_1(x) \wedge \neg P_1(x)$.

Given a sequence of 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 $d'_i = z_j$ 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 loosely inspired by *regex golf* [11], a game whose objective is to write a shortest possible regular expression matching a fixed set of strings and not matching another (disjoint) set of strings. Data Golf is formalized by the function $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ (Figure 15) that computes a structure \mathcal{S} such that $\mathcal{T}_{\vec{v}}^+[\vec{\text{fv}}(Q)] \subseteq \llbracket Q \rrbracket$, $\mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$, and $\llbracket Q' \rrbracket$, $\llbracket \neg Q' \rrbracket$ contain at least $\min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$ tuples each, for every subquery Q' of a query Q , where \vec{v} is a sequence of distinct variables such that $\text{fv}(Q) \subseteq \vec{v}$, $\mathcal{T}_{\vec{v}}^+$ and $\mathcal{T}_{\vec{v}}^-$ are sets of tuples over \vec{v} , and $\gamma \in \{0, 1\}$ is a *strategy*.

The function $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ can fail on an equality between two variables $x \approx y$. In this case, the function $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ does not compute any Data Golf structure. We define the *not-depth* of a subquery $x \approx y$ in Q as the number of subqueries that have the form of a negation among the queries $x \approx y \sqsubseteq \dots \sqsubseteq Q$, i.e., the number of negations on the path between the subquery $x \approx y$ and Q 's main constructor. To prevent failure, we generate the sets $\mathcal{T}_{\vec{v}}^+$, $\mathcal{T}_{\vec{v}}^-$ to only contain tuples with equal values for all variables in equalities with even (odd, respectively) not-depth and pairwise distinct values for all variables in equalities with odd (even, respectively) not-depth. This is not always possible, e.g., for $x \approx y \wedge \neg x \approx y$, in which case no Data Golf structure can be computed.

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.

Finally, we justify why a Data Golf structure \mathcal{S} computed by $\text{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ satisfies $\mathcal{T}_{\vec{v}}^+[\vec{\text{fv}}(Q)] \subseteq \llbracket Q \rrbracket$ and $\mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$. 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 $Q_1 \vee Q_2$, $Q_1 \wedge Q_2$ can be combined using the union operator. The only case that does not follow immediately is that $\mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$ for a query Q of the form $\exists y. Q_y$. We prove this case by contradiction. Without loss of generality we assume that $\vec{\text{fv}}(Q_y) = \vec{\text{fv}}(Q) \cdot y$. Suppose that $\vec{d} \in \mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)]$ and $\vec{d} \in \llbracket Q \rrbracket$. Because $\vec{d} \in \mathcal{T}_{\vec{v}}^-[\vec{\text{fv}}(Q)]$, there exists some d such that $\vec{d} \cdot d \in \mathcal{T}_{\vec{v}}^2[\vec{\text{fv}}(Q_y)]$. 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, we get $\vec{d} \cdot d \notin \llbracket Q_y \rrbracket$ and $\vec{d} \cdot d' \in \mathcal{T}_{\vec{v}}^2[\vec{\text{fv}}(Q_y)]$. Because $\text{con}_{\text{vgt}}(y, Q_y, \mathcal{G})$ holds for some \mathcal{G} satisfying the assumption of CON, the query Q_y is equivalent to

$$Q_y \equiv (Q_y \wedge \text{qps}^\vee(\mathcal{G})) \vee \bigvee_{z \in \text{eqs}(y, \mathcal{G})} (Q_y[y \mapsto z]) \vee Q_y[y/\perp].$$

Recall that $\vec{d} \cdot d' \in \llbracket Q_y \rrbracket$. If the tuple $\vec{d} \cdot d'$ satisfies $Q_y[y \mapsto z]$ for some $z \in \text{eqs}(y, \mathcal{G})$ or $Q_y[y/\perp]$, then $\vec{d} \cdot d \in \llbracket Q_y \rrbracket$ (contradiction) because the variable y does not occur in the queries $Q_y[y \mapsto z]$ and $Q_y[y/\perp]$ and thus its valuation in $\vec{d} \cdot d'$ can be arbitrarily changed. Otherwise, the tuple $\vec{d} \cdot d'$ satisfies some quantified predicate $Q_{qp} \in \text{qps}(\mathcal{G})$ such that $\{y\} \subsetneq \text{fv}(Q_{qp})$ (CON). Hence, the tuples $\vec{d} \cdot d$ and $\vec{d} \cdot d'$ agree on the valuation of a variable $x \in \text{fv}(Q_{qp}) \setminus \{y\}$. Next we observe that the valuations of every variable (e.g., x) in the tuples of the sets $\mathcal{T}_{\vec{v}}^+$, $\mathcal{T}_{\vec{v}}^-$ are pairwise distinct (the conditions $\mathcal{T}_{\vec{v}}^+[x] \cap \mathcal{T}_{\vec{v}}^-[x] = \emptyset$, $|\mathcal{T}_{\vec{v}}^+[x]| \geq |\mathcal{T}_{\vec{v}}^+|$, and $|\mathcal{T}_{\vec{v}}^-[x]| \geq |\mathcal{T}_{\vec{v}}^-|$). Because $\vec{d} \cdot d \in \mathcal{T}_{\vec{v}}^2[\vec{\text{fv}}(Q_y)]$, $\vec{d} \cdot d'$ satisfies the quantified predicate Q_{qp} with $y \in \text{fv}(Q_{qp})$, and the tuples $\vec{d} \cdot d$ and $\vec{d} \cdot d'$ agree on the valuation of x , the valuations of y in these two tuples must be equal, i.e., $d = d'$ (contradiction).

The sets $\mathcal{T}_{\vec{v}}^+$, $\mathcal{T}_{\vec{v}}^-$ only grow in dg 's recursion and the properties CON, CST, VAR, 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 $\llbracket Q' \rrbracket$, $\llbracket \neg Q' \rrbracket$ contain at least $\min\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$ tuples each, for every $Q' \sqsubseteq Q$.