Practical Relational Calculus Query Evaluation

- ² Martin Raszyk ⊠©
- ³ Department of Computer Science, ETH Zürich, Switzerland
- ₄ David Basin ⊠ <a>[□]
- 5 Department of Computer Science, ETH Zürich, Switzerland
- 🛛 Srđan Krstić 🖂 回
- 7 Department of Computer Science, ETH Zürich, Switzerland
- [∗] Dmitriy Traytel ⊠©
- 9 Department of Computer Science, University of Copenhagen, Denmark

¹⁰ — Abstract

The relational calculus (RC) is a concise, declarative query language. However, existing RC query 11 evaluation approaches are inefficient and often deviate from established algorithms based on finite 12 tables used in database management systems. We devise a new translation of an arbitrary RC query 13 into two safe-range queries, for which the finiteness of the query's evaluation result is guaranteed. 14 Assuming an infinite domain, the two queries have the following meaning: The first is closed and 15 characterizes the original query's relative safety, i.e., whether given a fixed database, the original 16 query evaluates to a finite relation. The second safe-range query is equivalent to the original query, if 17 the latter is relatively safe. We compose our translation with other, more standard ones to ultimately 18 obtain two SQL queries. This allows us to use standard database management systems to evaluate 19 arbitrary RC queries. We show that our translation improves the time complexity over existing 20 approaches, which we also empirically confirm in both realistic and synthetic experiments. 21

- 22 2012 ACM Subject Classification Theory of computation \rightarrow Database query languages (principles)
- 23 Keywords and phrases Relational calculus, relative safety, safe-range, query translation
- ²⁴ Digital Object Identifier 10.4230/LIPIcs.ICDT.2022.11

25 Supplementary Material (artifact and extended report): https://github.com/rc2sql/rc2sql

²⁶ **1** Introduction

Codd's theorem states that all domain-independent queries of the relational calculus (RC) 27 can be expressed in relational algebra (RA) [10]. A popular interpretation of this result is that 28 RA suffices to express all interesting queries. This interpretation justifies why SQL evolved as 29 the practical database query language with the RA as its mathematical foundation. SQL is 30 declarative and abstracts over the actual RA expression used to evaluate a query. Yet, SQL's 31 syntax inherits RA's deliberate syntactic limitations, such as union-compatibility, which 32 ensure domain independence. RC does not have such syntactic limitations, which arguably 33 makes it a more attractive declarative query language than both RA and SQL. The main 34 problem of RC is that it is not immediately clear how to evaluate even domain-independent 35 queries, much less how to handle the domain-dependent (i.e., not domain-independent) ones. 36 As a running example, consider a shop in which brands (unary finite relation B of brands) 37 sell products (binary finite relation P relating brands and products) and products are reviewed 38 by users with a score (ternary finite relation S relating products, users, and scores). We 39 consider a brand suspicious if there is a user and a score such that all the brand's products 40 were reviewed by that user with that score. An RC query computing suspicious brands is 41

$$_{42} \qquad Q^{susp} \coloneqq \mathsf{B}(b) \land \exists u, s. \ \forall p. \ \mathsf{P}(b, p) \longrightarrow \mathsf{S}(p, u, s).$$

⁴³ This query is domain-independent and follows closely our informal description. It is not,



11:2 Practical Relational Calculus Query Evaluation

⁴⁴ however, clear how to evaluate it because its second conjunct is domain-dependent as it is ⁴⁵ satisfied for every brand that does not occur in P. Finding suspicious brands using RA or ⁴⁶ SQL is a challenge, which only the best students from an undergraduate database course ⁴⁷ will accomplish. We give away an RA answer next (where – is the set difference operator ⁴⁸ and \triangleright is the anti-join, also known as the *generalized* difference operator [1]):

⁴⁹ $\pi_{brand}((\pi_{user,score}(\mathsf{S}) \times \mathsf{B}) - \pi_{brand,user,score}((\pi_{user,score}(\mathsf{S}) \times \mathsf{P}) \triangleright \mathsf{S})) \cup (\mathsf{B} - \pi_{brand}(\mathsf{P})).$

⁵⁰ The highlighted expressions $\pi_{user,score}(S)$ are called *generators*. They ensure that the left ⁵¹ operands of the anti-join and set difference operators include or have the same columns (i.e., ⁵² are union-compatible) as the corresponding right operands. (Following Codd [10], one could ⁵³ 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 domainindependent 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

⁵⁹ Once we have identified suspicious brands, we may want to obtain the users whose scoring ⁶⁰ made the brands suspicious. In RC, omitting u's quantifier from Q^{susp} achieves just that:

$$Q_{user}^{susp} \coloneqq \mathsf{B}(b) \land \exists s. \forall p. \mathsf{P}(b, p) \longrightarrow \mathsf{S}(p, u, s).$$

61

⁶² In contrast, RA cannot express the same property as it is domain-dependent (hence also not ⁶³ evaluable and thus out of scope for Van Gelder and Topor's translation): Q_{user}^{susp} is satisfied ⁶⁴ for every user if a brand has no products, i.e., it does not occur in P. Yet, Q_{user}^{susp} is satisfied ⁶⁵ for finitely many users on every database instance where P contains at least one row for every ⁶⁶ 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 73 of an infinite domain. To deal with queries that are domain-dependent, our translation 74 produces two RA expressions, instead of a single equivalent one. The first RA expression 75 characterizes the original RC query's relative safety, the decidable question of whether the 76 query evaluates to a finite relation for a given database, which can be the case even for 77 a domain-dependent query, e.g., Q_{user}^{susp} . If the original query is relatively safe on a given 78 database, i.e., produces some finite result, then the second RA expression evaluates to the 79 same finite result. Taken together, the two RA expressions solve the query capturability 80 problem [3]: they allow us to enumerate the original RC query's finite evaluation result, or 81 to learn that it would be infinite using RA operations on the unmodified database. 82

Our translation of an RC query to two RA expressions proceeds in several steps via saferange queries and the relational algebra normal form (Section 3). We focus on the first step of translating an RC query to two safe-range RC queries (Section 4), which fundamentally differs from Van Gelder and Topor's approach and produces better generators like $\pi_{user,score}(S)$. Our 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 [30]. Along the way to SQL, we leverage various ideas from the literature to optimize the overall

M. Raszyk, D. Basin, S. Krstić, and D. Traytel

result (Section 6). For example, we generalize Claußen et al. [9]'s approach to avoid evaluating 91 Cartesian products like $\pi_{user,score}(S) \times P$ in the above translation by using count aggregations. 92 The overall translation allows us to use standard database management systems to evaluate 93 RC queries. We implement our translation and use PostgreSQL to evaluate the translated 94 queries. Using a real Amazon review dataset [23] and our synthetic benchmark that generates 95 hard database instances for random RC queries (Section 5), we evaluate our translation's 96 performance. The evaluation shows that our approach outperforms Van Gelder and Topor's 97 translation (which also uses PostgreSQL for evaluation) and other approaches (Section 6). 98 In summary, the following are our three main contributions: 99 We devise a translation of an arbitrary RC query into a pair of RA expressions as 100 described above. The time complexity of evaluating our translation's results improves 101

- ¹⁰² 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 real data and confirm that our translation's improved asymptotic time complexity carries over into practice.
- ¹⁰⁷ To challenge RC2SQL (and its competitors) in our evaluation, we devise the *Data Golf* ¹⁰⁸ benchmark that generates hard database instances for randomly generated RC queries.

109 2 Related Work

We recall Trakhtenbrot's theorem and the fundamental notions of *capturability* and *data complexity*. Given an RC query over a *finite* domain, Trakhtenbrot [27] showed that it is undecidable whether there exists a (finite) structure satisfying the query. In contrast, the question of whether a fixed structure satisfies the given RC query is decidable [2].

Kifer [16] calls a query class capturable if there is an algorithm that, given a query in 114 the class and a database instance, enumerates the query's evaluation result, i.e., all tuples 115 satisfying the query. Avron and Hirshfeld [3] observe that Kifer's notion is restricted because 116 it requires every query in a capturable class to be domain independent. Hence, they propose 117 an alternative definition that we also use: A query class is capturable if there is an algorithm 118 that, given a query in the class, a (finite or infinite) domain, and a database instance, 119 determines whether the query's evaluation result on the database instance over the domain 120 is finite and enumerates the result in this case. Our work solves Avron and Hirshfeld's 121 capturability problem additionally assuming an infinite domain. 122

Data complexity [29] 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 128 evaluate them under a finite structure [18]. The core question here is whether the evaluation 129 produces the same result as defined by the natural semantics, which typically considers infinite 130 domains. Codd's theorem [10] affirmatively answers this question for domain-independent 131 queries, restricting the structure to the active domain. Ailamazyan et al. [2] show that RC is a 132 capturable query class by extending the active domain with a few additional elements, whose 133 number depends only on the query, and evaluating the query over this finite domain. Natural-134 active collapse results [5] generalize Ailamazyan et al.'s [2] result to extensions of RC (e.g., 135 with order relations) by combining the structure reduction with a translation-based approach. 136 Hull and Su [15] study several semantics of RC that guarantee the finiteness of the query's 137

11:4 Practical Relational Calculus Query Evaluation

evaluation result. In particular, the "output-restricted unlimited interpretation" only restricts
the query's evaluation result to tuples that only contain elements in the active domain, but
the quantified variables still range over the (finite or infinite) underlying domain. Our work
is inspired by all these theoretical landmarks, in particular Hull and Su's work (Section 4.1).

¹⁴² Yet we avoid using (extended) active domains, which make query evaluation impractical.

Query translation. Another strategy is to translate a given query into one that can 143 be evaluated efficiently, for example as a sequence of RA operations. Van Gelder and Topor 144 pioneered this approach [13,14] for RC. A core component of their translation is the choice of 145 generators, which replace the active domain restrictions from structure reduction approaches 146 and thereby improve the time complexity. Extensions to scalar and complex function symbols 147 have also been studied [12,19]. All these approaches focus on syntactic classes of RC, for which 148 domain-independence is given, e.g., the evaluable queries of Van Gelder and Topor [14, Defin-149 ition 5.2]. Our approach is inspired by Van Gelder and Topor's but generalizes it to handle 150 arbitrary RC queries at the cost of assuming an infinite domain. Also, we further improve 151 the time complexity of Van Gelder and Topor's approach by choosing better generators. 152

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].

¹⁶⁰ **3** Preliminaries

¹⁶¹ We introduce the RC syntax and semantics and define relevant classes of RC queries.

¹⁶² 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 $\iota : \mathcal{R} \to \mathbb{N}$ 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:

$$_{167} \qquad Q ::= \bot \mid \top \mid x \approx t \mid r(t_1, \dots, t_{\iota(r)}) \mid \neg Q \mid Q \lor Q \mid Q \land Q \mid \exists x. Q.$$

Here, $r \in \mathcal{R}$ is a predicate symbol, $t, t_1, \ldots, t_{\iota(r)} \in \mathcal{V} \cup \mathcal{C}$ are terms, and $x \in \mathcal{V}$ is a variable. We 168 write $\exists \vec{v}. Q$ for $\exists v_1, \ldots, \exists v_k. Q$ and $\forall \vec{v}. Q$ for $\neg \exists \vec{v}. \neg Q$, where \vec{v} is a variable sequence v_1, \ldots, v_k . 169 If k = 0, then both $\exists \vec{v}. Q$ and $\forall \vec{v}. Q$ denote just Q. Quantifiers have lower precedence than 170 conjunctions and disjunctions, e.g., $\exists x. Q_1 \land Q_2$ means $\exists x. (Q_1 \land Q_2)$. We use \approx to denote 171 the equality of terms in RC to distinguish it from =, which denotes syntactic object identity. 172 We also write $Q_1 \longrightarrow Q_2$ for $\neg Q_1 \lor Q_2$. However, defining $Q_1 \lor Q_2$ as a shorthand for 173 $\neg(\neg Q_1 \land \neg Q_2)$ would complicate later definitions, e.g., the safe-range queries (Section 3.2). 174 We define the subquery partial order \sqsubseteq on queries inductively on the structure of RC 175 queries, e.g., Q_1 is a subquery of the query $Q_1 \land \neg \exists y. Q_2$. One can also view \sqsubseteq as the (reflexive 176 and transitive) subterm relation on the datatype of RC queries. We denote by $\mathsf{sub}(Q)$ the 177 set of subqueries of a query Q and by fv(Q) the set of *free variables* in Q. Furthermore, we 178 denote by fv(Q) the sequence of free variables in Q based on some fixed ordering of variables. 179 We lift this notation to sets of queries in the standard way. A query Q with no free variables, 180 i.e., $fv(Q) = \emptyset$, is called *closed*. Queries of the form $r(t_1, \ldots, t_{\iota(r)})$ and $x \approx c$ are called *atomic* 181 predicates. We define the predicate $ap(\cdot)$ characterizing atomic predicates, i.e., ap(Q) is true 182

 $\begin{array}{ll} (\mathcal{S},\alpha) \not\models \bot; \ (\mathcal{S},\alpha) \models \top; \\ (\mathcal{S},\alpha) \models r(t_1,\ldots,t_{\iota(r)}) & \text{iff} \ (\alpha(t_1),\ldots,\alpha(t_{\iota(r)})) \in r^{\mathcal{S}}; \\ (\mathcal{S},\alpha) \models (Q_1 \lor Q_2) & \text{iff} \ (\mathcal{S},\alpha) \models Q_1 \text{ or } (\mathcal{S},\alpha) \models Q_2; \\ (\mathcal{S},\alpha) \models (Q_1 \land Q_2) & \text{iff} \ (\mathcal{S},\alpha) \models Q_1 \text{ and } (\mathcal{S},\alpha) \models Q_2; \\ \end{array} \right| \left(\begin{array}{c} (\mathcal{S},\alpha) \models (x \approx t) & \text{iff} \ \alpha(x) = \alpha(t); \\ (\mathcal{S},\alpha) \models (\neg Q) & \text{iff} \ (\mathcal{S},\alpha) \not\models Q; \\ (\mathcal{S},\alpha) \models (\mathcal{I},\alpha) & \text{iff} \ (\mathcal{S},\alpha) \models Q_1 \text{ or } (\mathcal{S},\alpha) \models Q_2; \\ & \text{for some } d \in \mathcal{D}. \end{array} \right)$

Figure 1 The semantics of RC.

20

iff Q is an atomic predicate. Queries of the form $\exists \vec{v}. r(t_1, \ldots, t_{\iota(r)})$ and $\exists \vec{v}. x \approx c$ are called 183 quantified predicates. We denote by $\tilde{\exists} x. Q$ the query obtained by existentially quantifying 184 a variable x from a query Q if x is free in Q, i.e., $\exists x. Q \coloneqq \exists x. Q \text{ if } x \in \mathsf{fv}(Q)$ and $\exists x. Q \coloneqq Q$ 185 otherwise. We lift this notation to sets of queries in the standard way. We use $\exists x. Q$ (instead of 186 $\exists x. Q$) when constructing a query to avoid introducing bound variables that never occur in Q. 187 A structure \mathcal{S} over a signature $(\mathcal{C}, \mathcal{R}, \iota)$ consists of a non-empty domain \mathcal{D} and interpret-188 ations $c^{\mathcal{S}} \in \mathcal{D}$ 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 189 r^{S} are *finite*. Note that this assumption does *not* yield a finite structure (as defined in finite 190 model theory [18]) since the domain \mathcal{D} can still be infinite. A (variable) assignment is a map-191 ping $\alpha: \mathcal{V} \to \mathcal{D}$. We additionally define α on constant symbols $\mathbf{c} \in \mathcal{C}$ as $\alpha(\mathbf{c}) = \mathbf{c}^{\mathcal{S}}$. We write 192 $\alpha[x \mapsto d]$ for the assignment that maps x to $d \in \mathcal{D}$ and is otherwise identical to α . We lift this 193 notation to sequences \vec{x} and \vec{d} of pairwise distinct variables and arbitrary domain elements of 194 the same length. The semantics of RC queries for a structure \mathcal{S} and an assignment α is defined 195 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. 196 For a fixed S, only the assignments to Q's free variables influence $\alpha \models Q$, i.e., $\alpha \models Q$ is 197 equivalent to $\alpha' \models Q$, for every variable assignment α' that agrees with α on fv(Q). For closed 198 queries Q, we write $\models Q$ and say that Q holds, since closed queries either hold for all variable 199 assignments or for none of them. We call a finite sequence d of domain elements $d_1, \ldots, d_k \in \mathcal{D}$ 200 a tuple. Given a query Q and a structure \mathcal{S} , we denote the set of satisfying tuples for Q by 201

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

We omit S from $\llbracket Q \rrbracket^S$ if S is fixed. We call the values from $\llbracket Q \rrbracket$ assigned to $x \in \mathsf{fv}(Q)$ column x. The active domain $\mathsf{adom}^S(Q)$ of a query Q and a structure S is a subset of the domain \mathcal{D} containing the interpretations c^S of all constant symbols that occur in Q and the values in the relations r^S interpreting all predicate symbols that occur in Q. Since C and \mathcal{R} are finite and all r^S are finite relations of a finite arity $\iota(r)$, the active domain $\mathsf{adom}^S(Q)$ is also a finite set. We omit S from $\mathsf{adom}^S(Q)$ if 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 $(S, \alpha) \models$ $Q_1 \iff (S, \alpha) \models Q_2$, for every S and α . Queries Q_1 and Q_2 over the same signature are *inf-equivalent*, written $Q_1 \stackrel{\infty}{\equiv} Q_2$, if $(S, \alpha) \models Q_1 \iff (S, \alpha) \models Q_2$, for every S with an *infinite* domain \mathcal{D} and every α . Clearly, equivalent queries are also inf-equivalent.

A query Q is *domain-independent* if $[\![Q]\!]^{S_1} = [\![Q]\!]^{S_2}$ holds for every two structures S_1 and S_2 that agree on the interpretations of constants ($c^{S_1} = c^{S_2}$) and predicates ($r^{S_1} = r^{S_2}$), while 213 214 their domains \mathcal{D}_1 and \mathcal{D}_2 may differ. Agreement on the interpretations implies $\mathsf{adom}^{\mathcal{S}_1}(Q) =$ 215 $\operatorname{\mathsf{adom}}^{S_2}(Q) \subseteq \mathcal{D}_1 \cap \mathcal{D}_2$. It is undecidable whether an RC query is domain-independent [24,28]. 216 We denote by $Q[x \mapsto y]$ the query obtained from the query Q after replacing each free 217 occurrence of the variable x by the variable y (possibly renaming bound variables to avoid 218 capture) and performing constant propagation, i.e., simplifications like $(x \approx x) \equiv \top, Q \wedge \bot \equiv$ 219 \perp , $Q \lor \perp \equiv Q$, etc. We lift this notation to sets of queries in the standard way. Finally, we 220 denote by $Q[x/\perp]$ the query obtained from Q after replacing every atomic predicate or equality 221 containing a free variable x by \perp (except for $x \approx x$) and performing constant propagation. 222 The function $\mathsf{flat}^{\oplus}(Q)$, where $\oplus \in \{\lor, \land\}$, computes a set of queries by "flattening" the 223 operator \oplus : $\mathsf{flat}^{\oplus}(Q) \coloneqq \mathsf{flat}^{\oplus}(Q_1) \cup \mathsf{flat}^{\oplus}(Q_2)$ if $Q = Q_1 \oplus Q_2$ and $\mathsf{flat}^{\oplus}(Q) \coloneqq \{Q\}$ otherwise. 224

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

Figure 2 The *generated* relation.

Figure 3 The *covered* relation.

225 3.2 Safe-Range Queries

The class of *safe-range* queries [1] is a decidable subset of domain-independent RC queries. 226 Its definition is based on the notion of range-restricted variables of a query. A variable is 227 called *range-restricted* if "its possible values all lie within the active domain of the query" [1]. 228 Intuitively, atomic predicates restrict the possible values of a variable that occurs in them as 229 a term. An equality $x \approx y$ can extend the set of range-restricted variables in a conjunction 230 $Q \wedge x \approx y$: If x or y is range-restricted in Q, then both x and y are range-restricted in $Q \wedge x \approx y$. 231 We formalize range-restricted variables using the generated relation $gen(x, Q, \mathcal{G})$, defined 232 in Figure 2. Specifically, $gen(x, Q, \mathcal{G})$ holds if x is a range-restricted variable in Q and every 233 satisfying assignment for Q satisfies some quantified predicate, referred to as *generator*, from 234 \mathcal{G} . Note that, unlike in a similar definition by Van Gelder and Topor [14, Figure 5] that defines 235 the rule gen $(x, \exists y, Q_y, \mathcal{G})$ if $x \neq y$ and gen (x, Q_y, \mathcal{G}) , we modify the rule's conclusion to existen-236 tially quantify the bound variable y from all queries in \mathcal{G} where y occurs: $gen(x, \exists y, Q_y, \exists y, \mathcal{G})$. 237 Hence, $gen(x, Q, \mathcal{G})$ implies $fv(\mathcal{G}) \subseteq fv(Q)$. We now formalize these relationships. 238

▶ Lemma 1. Let Q be a query, $x \in \mathsf{fv}(Q)$, and G be a set of quantified predicates such that gen(x,Q,G). Then (i) for every $Q_{qp} \in G$, we have $x \in \mathsf{fv}(Q_{qp})$ and $\mathsf{fv}(Q_{qp}) \subseteq \mathsf{fv}(Q)$, (ii) for every α such that $\alpha \models Q$, there exists $Q_{qp} \in G$ such that $\alpha \models Q_{qp}$, and (iii) $Q[x/\bot] = \bot$.

▶ Definition 2. We define gen(x, Q) to hold iff there exists a set \mathcal{G} such that gen(x, Q, \mathcal{G}). Let nongens(Q) := { $x \in \mathsf{fv}(Q) \mid \mathsf{gen}(x, Q) \text{ does not hold}$ } be the set of free variables in a query Q that are not range-restricted. A query Q has range-restricted free variables if every free variable of Q is range-restricted, i.e., nongens(Q) = Ø. A query Q has range-restricted bound variables if the bound variable y in every subquery $\exists y. Q_y$ of Q is range-restricted, i.e., gen(y, Q_y) holds. A query is safe-range if it has range-restricted free and range-restricted bound variables.

Relational algebra normal form (RANF) is a class of safe-range queries that can be easily mapped to RA and evaluated using the RA operations for projection, column duplication, selection, set union, binary join, and anti-join. Following a standard textbook [1, Section 5.4], we define the predicate ranf(\cdot) characterizing RANF queries and the translation sr2ranf(\cdot) of a safe-range query into an equivalent RANF query.

253 3.3 Query Cost

To assess the time complexity of evaluating a RANF query Q, we define the *cost* of Q over 254 a structure \mathcal{S} , denoted $\operatorname{cost}^{\mathcal{S}}(Q)$, to be the sum of intermediate result sizes over all RANF 255 subqueries of Q. Formally, $\mathsf{cost}^{\mathcal{S}}(Q) \coloneqq \sum_{Q' \sqsubseteq Q, \mathsf{ranf}(Q')} \left| \llbracket Q' \rrbracket^{\mathcal{S}} \right| \cdot |\mathsf{fv}(Q')|$. This corresponds to 256 evaluating Q following its RANF structure using the RA operations. The complexity of these 257 operations is linear in the combined input and output size (ignoring logarithmic factors due 258 to set operations). The output size (the number of tuples times the number of variables) is 259 counted in $|[Q']]^{S} | \cdot |\mathsf{fv}(Q')|$ and the input size is counted as the output size for the input sub-260 queries. Repeated subqueries are only considered once, which does not affect the asymptotics 261 of query cost. In practice, the evaluation results for common subqueries can be reused. 262

4 Query Translation

Our approach to evaluating an arbitrary RC query Q over a fixed structure 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) $\mathsf{fv}(Q_{fin}) = \mathsf{fv}(Q)$ unless Q_{fin} is syntactically equal to \bot ; $\mathsf{fv}(Q_{inf}) = \emptyset$;

²⁶⁷ (EVAL) $\llbracket Q \rrbracket$ is an infinite set if Q_{inf} holds; otherwise $\llbracket Q \rrbracket = \llbracket Q_{fin} \rrbracket$ is a finite set.

Since the queries Q_{fin} and Q_{inf} are safe-range, they are domain-independent and thus $[\![Q_{fin}]\!]$ is a finite set of tuples. In particular, $[\![Q]\!]$ is a finite set of tuples if Q_{inf} does not hold. Our translation generalizes Hull and Su's case distinction that restricts bound variables [15] to restrict all variables. Moreover, we use Van Gelder and Topor's idea to replace the active domain by a smaller set (generator) specific to each variable [14] while further improving the generators.

273 4.1 Restricting One Variable

2

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

 \tilde{Q}' is a disjunction and x is range-restricted in all but the last disjunct.

The disjunct in which x is not range-restricted has a special form that is central to our translation: it is the conjunction of a query in which x does not occur and a query that is satisfied by infinitely many values of x. 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:

$$\tilde{Q} \equiv (\tilde{Q} \land x \in \mathsf{adom}(\tilde{Q})) \lor \bigvee_{y \in \mathsf{fv}(\tilde{Q}) \setminus \{x\}} (\tilde{Q}[x \mapsto y] \land x \approx y) \lor (\tilde{Q}[x/\bot] \land \neg(x \in \mathsf{adom}(\tilde{Q}) \lor \bigvee_{y \in \mathsf{fv}(\tilde{Q}) \setminus \{x\}} x \approx y)).$$

Here, $x \in \operatorname{adom}(\tilde{Q})$ stands for an RC query with a single free variable x that is satisfied by an assignment α if and only if $\alpha(x) \in \operatorname{adom}^{\mathcal{S}}(\tilde{Q})$. The translation distinguishes the following three cases for a fixed assignment α :

if $\alpha(x) \in \mathsf{adom}^{\mathcal{S}}(\tilde{Q})$ holds, then we do not alter the query \tilde{Q} ;

if $x \approx y$ holds for some free variable $y \in \mathsf{fv}(\tilde{Q}) \setminus \{x\}$, then x can be replaced by y in \tilde{Q} ; otherwise, \tilde{Q} is equivalent to $\tilde{Q}[x/\bot]$, i.e., all atomic predicates with a free occurrence of xcan be replaced by \bot (because $\alpha(x) \notin \mathsf{adom}^S(\tilde{Q})$), all equalities $x \approx y$ and $y \approx x$ for $y \in$ $\mathsf{fv}(\tilde{Q}) \setminus \{x\}$ can be replaced by \bot (because $\alpha(x) \neq \alpha(y)$), and all equalities $x \approx z$ for a bound variable z can be replaced by \bot (because $\alpha(x) \notin \mathsf{adom}^S(\tilde{Q})$ and z is range-restricted in its subquery $\exists z. Q_z$, by assumption, i.e., $\mathsf{gen}(z, Q_z)$ holds and thus, for all α' , we have $\alpha' \models$ $\exists z. Q_z$ if and only if there exists $d \in \mathsf{adom}^S(Q_z) \subseteq \mathsf{adom}^S(\tilde{Q})$ such that $\alpha'[z \mapsto d] \models Q_z)$.

11:8 Practical Relational Calculus Query Evaluation

Note that $\exists \vec{\mathsf{fv}}(Q) \setminus \{x\}$. 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 $\mathsf{adom}^{\mathcal{S}}(Q)$ of the query Q and a structure \mathcal{S} , Van Gelder and Topor [14] replace the condition $x \in \mathsf{adom}(Q)$ in their translation by $\exists \vec{\mathsf{fv}}(\mathcal{G}) \setminus \{x\}$. \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, for all α ,

$$\begin{array}{l} \alpha \models \neg \exists \tilde{\mathsf{fv}}(\mathcal{G}) \setminus \{x\}. \mathcal{G} \Longrightarrow (\alpha \models Q \Longleftrightarrow \alpha \models Q[x/\bot]) \ (\text{VGT}_1) & \text{and} \\ \alpha \models Q[x/\bot] & \Longrightarrow \alpha \models \forall x. Q & (\text{VGT}_2). \end{array}$$

Note that (VGT_2) does not hold for the query $Q \coloneqq \neg B(x)$ and thus a generator set \mathcal{G} of atomic 303 predicates satisfying (VGT₂) only exists for a proper subset of all RC queries. In contrast, we 304 only require that \mathcal{G} satisfies (VGT₁) in our translation. To this end, we define a *covered* relation 305 $cov(x, Q, \mathcal{G})$ (in contrast to Van Gelder and Topor's constrained relation con [14, Figure 5]) 306 such that, for every variable x and query \hat{Q} with range-restricted bound variables, there exists 307 at least one set \mathcal{G} such that $\operatorname{cov}(x, \hat{Q}, \mathcal{G})$ and (VGT_1) holds. Figure 3 shows the definition 308 of this relation. Unlike the generator set \mathcal{G} in $gen(x, Q, \mathcal{G})$, the cover set \mathcal{G} in $cov(x, Q, \mathcal{G})$ 309 may also contain equalities between two variables. Hence, we define a function $qps(\mathcal{G})$ that 310 collects all *generators*, i.e., quantified predicates and a function $eqs(x, \mathcal{G})$ that collects all 311 variables y distinct from x occurring in equalities of the form $x \approx y$. We use $qps^{\vee}(\mathcal{G})$ to 312 denote the query $\bigvee_{Q_{qp} \in qps(\mathcal{G})} Q_{qp}$. We state the soundness and completeness of the relation 313 $\operatorname{cov}(x, Q, \mathcal{G})$ in the next lemma, which follows by induction on the derivation of $\operatorname{cov}(x, \tilde{Q}, \mathcal{G})$. 314

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

$$_{^{318}} \qquad \alpha \models \neg (\mathsf{qps}^{\vee}(\mathcal{G}) \vee \bigvee_{y \in \mathsf{eqs}(x,\mathcal{G})} x \approx y) \Longrightarrow (\alpha \models \tilde{Q} \Longleftrightarrow \alpha \models \tilde{Q}[x/\bot])$$

Finally, to preserve the dependencies between the variable x and the remaining free variables of Q occurring in the quantified predicates from $qps(\mathcal{G})$, we do not project $qps(\mathcal{G})$ on the single variable x, i.e., we restrict x by $qps^{\vee}(\mathcal{G})$ instead of $\exists \vec{fv}(Q) \setminus \{x\}$. $qps(\mathcal{G})$. From Lemma 3, we derive our optimized translation characterized by the following lemma.

▶ Lemma 4. Let \tilde{Q} be a query with range-restricted bound variables, $x \in \mathsf{fv}(\tilde{Q})$, and \mathcal{G} be such that $\mathsf{cov}(x, \tilde{Q}, \mathcal{G})$ holds. Then $x \in \mathsf{fv}(Q_{qp})$ and $\mathsf{fv}(Q_{qp}) \subseteq \mathsf{fv}(\tilde{Q})$, for every $Q_{qp} \in \mathsf{qps}(\mathcal{G})$, and

$$\hat{Q} \equiv (\hat{Q} \land \mathsf{qps}^{\lor}(\mathcal{G})) \lor \bigvee_{y \in \mathsf{eqs}(x,\mathcal{G})} (\hat{Q}[x \mapsto y] \land x \approx y) \lor (\hat{Q}[x/\bot] \land \neg(\mathsf{qps}^{\lor}(\mathcal{G}) \lor \bigvee_{y \in \mathsf{eqs}(x,\mathcal{G})} x \approx y)).$$

Note that x is not guaranteed to be range-restricted in (\bigstar)'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.

330 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., $gen(x, \tilde{Q})$ does not hold. To translate $\exists x. \tilde{Q}$ into an infequivalent 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 (\bigstar) to \tilde{Q} and distribute the existential quantifier binding x over disjunction. Next we observe that

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

input: An RC query Q. **input:** An RC query Q. output: Safe-range query pair (Q_{fin}, Q_{inf}) **output:** A query \hat{Q} with for which (FV) and (EVAL) hold. range-restricted bound variables such that $Q \stackrel{\infty}{\equiv} \tilde{Q}$. 1 function fixfree(Q_{fin}) = $\{(Q_{fix}, Q^{=}) \in \mathcal{Q}_{fin} \mid \operatorname{nongens}(Q_{fix}) \neq \emptyset\};\$ 1 function fixbound(Q, x) = 2 function $\inf(\mathcal{Q}_{fin}, Q) = \{(Q_{\infty}, Q^{=}) \in$ $\{Q_{fix} \in \mathcal{Q} \mid x \in \operatorname{nongens}(Q_{fix})\};\$ $\mathcal{Q}_{fin} \mid \mathsf{disjointvars}(Q_{\not\infty}, Q^{=}) \neq \emptyset \lor$ 2 function $\mathsf{rb}(Q) =$ $\mathsf{fv}(Q_{\infty} \land Q^{=}) \neq \mathsf{fv}(Q) \};$ switch Q do 3 **3 function** split(Q) =case $\neg Q'$ do return $\neg \mathsf{rb}(Q')$; 4 $\mathcal{Q}_{fin} \coloneqq \{ (\mathsf{rb}(Q), \top) \}; \mathcal{Q}_{inf} \coloneqq \emptyset;$ case $Q'_1 \vee Q'_2$ do return 4 5 while fixfree(Q_{fin}) $\neq \emptyset$ do $\mathsf{rb}(Q'_1) \lor \mathsf{rb}(Q'_2);$ 5 $(Q_{fix}, Q^{=}) \leftarrow \mathsf{fixfree}(\mathcal{Q}_{fin});$ 6 6 case $Q'_1 \wedge Q'_2$ do return $x \leftarrow \mathsf{nongens}(Q_{fix});$ $\mathsf{rb}(Q'_1) \wedge \mathsf{rb}(Q'_2);$ 7 $\mathcal{G} \leftarrow \{\mathcal{G} \mid \mathsf{cov}(x, Q_{fix}, \mathcal{G})\};$ case $\exists x. Q_x$ do 8 7 $\mathcal{Q}_{fin} \coloneqq (\mathcal{Q}_{fin} \setminus \{(Q_{fix}, Q^{=})\}) \cup$ $\mathcal{Q} \coloneqq \mathsf{flat}^{\vee}(\mathsf{rb}(Q_x));$ 9 8 $\{(Q_{fix} \land \mathsf{qps}^{\lor}(\mathcal{G}), Q^{=})\} \cup$ **while** fixbound(Q, x) $\neq \emptyset$ **do** 9 $\bigcup_{y \in \mathsf{eqs}(x,\mathcal{G})} \{ (Q_{\mathit{fix}}[x \mapsto y], Q^{=} \land x \approx y) \};$ $Q_{fix} \leftarrow \mathsf{fixbound}(\mathcal{Q}, x);$ 10 $\mathcal{G} \leftarrow \{\mathcal{G} \mid \mathsf{cov}(x, Q_{fix}, \mathcal{G})\};\$ $\mathcal{Q}_{inf} \coloneqq \mathcal{Q}_{inf} \cup \{Q_{fix}[x/\bot]\};$ 11 10 $\mathcal{Q} \coloneqq (\mathcal{Q} \setminus \{Q_{fix}\}) \cup$ while $\inf(\mathcal{Q}_{fin}, Q) \neq \emptyset$ do 12 11 $\{Q_{fix} \land qps^{\vee}(\mathcal{G})\} \cup$ $(Q_{\not\infty}, Q^{=}) \leftarrow \inf(\mathcal{Q}_{fin}, Q);$ 12 $\bigcup_{y \in \mathsf{eqs}(x,\mathcal{G})} \{Q_{\mathit{fix}}[x \mapsto y]\} \cup$ $\mathcal{Q}_{fin} \coloneqq \mathcal{Q}_{fin} \setminus \{(Q_{\not\infty}, Q^{=})\};$ 13 $\mathcal{Q}_{inf} \coloneqq \mathcal{Q}_{inf} \cup \{Q_{\not\infty} \land Q^{=}\};$ $\{Q_{fix}[x/\bot]\};$ $\mathbf{14}$ return $\bigvee_{\tilde{Q}\in\mathcal{Q}} \tilde{\exists} x. \tilde{Q};$ return $(\bigvee_{(Q_{\infty},Q^{=})\in\mathcal{Q}_{fin}}(Q_{\infty}\wedge Q^{=})),$ 13 $\mathbf{15}$ otherwise do return Q; $\mathsf{rb}(\bigvee_{Q_{\infty}\in\mathcal{Q}_{inf}}\exists \vec{\mathsf{fv}}(Q_{\infty}).Q_{\infty}));$ 14 Figure 4 Restricting bound variables. **Figure 5** Restricting free variables.

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

$$\exists x. \, (\tilde{Q}[x/\bot] \land \neg (\mathsf{qps}^{\lor}(\mathcal{G}) \lor \bigvee_{y \in \mathsf{eqs}(x,\mathcal{G})} x \approx y)) \stackrel{\sim}{=} \tilde{Q}[x/\bot]$$

33

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

▶ Lemma 5. Let \tilde{Q} be a query with range-restricted bound variables, $x \in \mathsf{fv}(\tilde{Q})$, and \mathcal{G} be a set of quantified predicates and equalities such that $\mathsf{cov}(x, \tilde{Q}, \mathcal{G})$ holds. Then

$$\exists x. \, \tilde{Q} \stackrel{\infty}{=} (\exists x. \, \tilde{Q} \land \mathsf{qps}^{\vee}(\mathcal{G})) \lor \bigvee_{y \in \mathsf{eqs}(x, \mathcal{G})} (\tilde{Q}[x \mapsto y]) \lor \tilde{Q}[x/\bot]. \tag{\textbf{\textbf{f}}}$$

Our approach for restricting all bound variables recursively applies Lemma 5. Because 346 the set \mathcal{G} such that $cov(x, Q, \mathcal{G})$ holds is not necessarily unique, we introduce the following 347 (general) notation. We denote the non-deterministic choice of an object X from a non-empty 348 set \mathcal{X} as $X \leftarrow \mathcal{X}$. We define the recursive function $\mathsf{rb}(Q)$ in Figure 4, where rb stands for 349 range-restrict bound (variables). The function converts an arbitrary RC query Q into an 350 inf-equivalent query with range-restricted bound variables. We proceed by describing the 351 case $\exists x. Q_x$. First, $\mathsf{rb}(Q_x)$ is recursively applied on Line 8 to establish the precondition of 352 Lemma 5 that the translated query has range-restricted bound variables. Because existential 353 quantification distributes over disjunction, we flatten disjunction in $rb(Q_x)$ and process 354

11:10 Practical Relational Calculus Query Evaluation

the individual disjuncts independently. We apply $(\bigstar \exists)$ 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 $(\bigstar \exists)$ to

 $_{357}$ Q_{fix} the variable x is either range-restricted or does not occur in Q'_{fix} , i.e., $x \notin \text{nongens}(Q'_{fix})$.

This entails the termination of the loop on Lines 9–12.

Example 6. Consider the query $Q_{user}^{susp} := \mathsf{B}(b) \land \exists s. \forall p. \mathsf{P}(b, p) \longrightarrow \mathsf{S}(p, u, s)$ from Section 1. Restricting its bound variables yields the query

$$\mathsf{rb}(Q_{user}^{susp}) = \mathsf{B}(b) \land ((\exists s. (\neg \exists p. \mathsf{P}(b, p) \land \neg \mathsf{S}(p, u, s)) \land (\exists p. \mathsf{S}(p, u, s))) \lor (\neg \exists p. \mathsf{P}(b, p))).$$

The bound variable p is already range-restricted in Q_{user}^{susp} and thus only s must be restric-362 ted. Applying (\bigstar) to restrict s in $\neg \exists p. \mathsf{P}(b,p) \land \neg \mathsf{S}(p,u,s)$, then existentially quantifying 363 s, and distributing the existential over disjunction yields the first disjunct in $\mathsf{rb}(Q_{user}^{susp})$ 364 above and $\exists s. (\neg \exists p. \mathsf{P}(b, p)) \land \neg (\exists p. \mathsf{S}(p, u, s))$ as the second disjunct. Because there exists 365 some value in the infinite domain \mathcal{D} that does not belong to the finite interpretation of 366 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, 367 $\exists s. (\neg \exists p. \mathsf{P}(b, p)) \land \neg (\exists p. \mathsf{S}(p, u, s))$ is inf-equivalent to $\neg \exists p. \mathsf{P}(b, p)$, i.e., the second disjunct in 368 $\mathsf{rb}(Q_{user}^{susp})$. This reasoning justifies applying ($\bigstar \exists$) to restrict s in $\exists s. \neg \exists p. \mathsf{P}(b, p) \land \neg \mathsf{S}(p, u, s)$. 369

4.3 Restricting Free Variables

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

▶ Lemma 7. Let a structure S with an infinite domain D be fixed. Let x be a free variable in a query \tilde{Q} with range-restricted bound variables and let $cov(x, \tilde{Q}, \mathcal{G})$ for a set of quantified predicates and equalities \mathcal{G} . If $\tilde{Q}[x/\bot]$ is not satisfied by any tuple, then

$$[\![\tilde{Q}]\!] = \left[\![(\tilde{Q} \land \mathsf{qps}^{\vee}(\mathcal{G})) \lor \bigvee_{y \in \mathsf{eqs}(x,\mathcal{G})} (\tilde{Q}[x \mapsto y] \land x \approx y) \right]\!].$$

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

Proof. If $\tilde{Q}[x/\perp]$ is not satisfied by any tuple, then (\bigstar) follows from (\bigstar) . If $\tilde{Q}[x/\perp]$ is satisfied by some tuple, then the last disjunct in (\bigstar) 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 $y \in \exp(x, \mathcal{G})$ and x does not appear among the finitely many values interpreting the quantified predicates from $qps(\mathcal{G})$.

We remark that $[\![\tilde{Q}]\!]$ might be an infinite set of tuples even if $\tilde{Q}[x/\perp]$ is never satisfied, for some x. This is because $\tilde{Q}[y/\perp]$ might be satisfied by some tuple, for some y, in which case Lemma 7 (for y) implies that $[\![\tilde{Q}]\!]$ is an infinite set of tuples. Still, (\mathfrak{A}) can be applied to \tilde{Q} for x resulting in an equivalent query that is also satisfied by an infinite set of tuples.

Our approach is implemented by the function $\operatorname{split}(Q)$ defined in Figure 5. In the following, we describe this function and informally justify its correctness, formalized by the input/output specification. In $\operatorname{split}(Q)$, 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

$$Q_{fin} \coloneqq \bigvee_{(Q_{\infty},Q^{=})\in\mathcal{Q}_{fin}}(Q_{\infty}\wedge Q^{=}), \qquad \qquad Q_{inf} \coloneqq \bigvee_{Q_{\infty}\in\mathcal{Q}_{inf}}\exists \vec{\mathsf{fv}}(Q_{\infty}). Q_{\infty}$$

and, for every $(Q_{\infty}, Q^{=}) \in \mathcal{Q}_{fin}, Q^{=}$ is a conjunction of equalities. As long as there exists some $(Q_{fix}, Q^{=}) \in \mathcal{Q}_{fin}$ such that $\operatorname{nongens}(Q_{fix}) \neq \emptyset$, we apply (\bigstar) to Q_{fix} and add the query $Q_{fix}[x/\bot]$ to \mathcal{Q}_{inf} . We remark that if we applied (\bigstar) 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 (\bigstar) to Q_{fix} , nongens (Q'_{fix}) is a proper subset of nongens (Q_{fix}) . This entails the termination of the loop on Lines 5–10. Finally, if $[\![Q_{fix}]\!]$ is an infinite set of tuples, then $[\![Q_{fix} \land Q^=]\!]$ 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 (\bigstar) to Q_{fix} instead of $Q_{fix} \land Q^=$.

After the loop on Lines 5–10 in Figure 5 terminates, for every $(Q_{\infty}, Q^{=}) \in \mathcal{Q}_{fin}, Q_{\infty}$ is a 402 safe-range query and $Q^{=}$ is a conjunction of equalities such that $fv(Q_{\infty} \wedge Q^{=}) = fv(Q)$. How-403 ever, the query $Q_{\infty} \wedge Q^{=}$ need not be safe-range, e.g., if $Q_{\infty} := \mathsf{B}(x)$ and $Q^{=} := (x \approx y \wedge u \approx v)$. 404 Given a set of equalities $Q^{=}$, let classes $(Q^{=})$ be the set of equivalence classes of free variables 405 $\mathsf{fv}(\mathcal{Q}^{=})$ with respect to $\mathcal{Q}^{=}$. For instance, $\mathsf{classes}(\{x \approx y, y \approx z, u \approx v\}) = \{\{x, y, z\}, \{u, v\}\}$. 406 Let disjointvars $(Q_{\infty}, Q^{=}) \coloneqq \bigcup_{V \in \mathsf{classes}(\mathsf{flat}^{(Q^{=})}), V \cap \mathsf{fv}(Q_{\infty}) = \emptyset} V$ be the set of all variables in equi-40 valence classes from classes(flat^(Q=)) that are disjoint from Q_{∞} 's free variables. Then, $Q_{\infty} \wedge$ 408 $Q^{=}$ is safe-range if and only if disjointvars $(Q_{\infty}, Q^{=}) = \emptyset$ (recall the definition of safe-range). 409 Now if disjointvars $(Q_{\infty}, Q^{=}) \neq \emptyset$ and $Q_{\infty} \wedge Q^{=}$ is satisfied by some tuple, then $[\![Q_{\infty} \wedge Q^{=}]\!]$ 410 is an infinite set of tuples because all equivalence classes of variables in disjointvars $(Q_{\infty}, Q^{=}) \neq Q_{\infty}$ 411 \emptyset can be assigned arbitrary values from the infinite domain \mathcal{D} . In our example with 412 $Q_{\infty} := \mathsf{B}(x)$ and $Q^{=} := (x \approx y \wedge u \approx v)$, we have disjoint $\mathsf{vars}(Q_{\infty}, Q^{=}) = \{u, v\} \neq \emptyset$. 413 Moreover, if $\mathsf{fv}(Q_{\infty} \wedge Q^{=}) \neq \mathsf{fv}(Q)$ and $Q_{\infty} \wedge Q^{=}$ is satisfied by some tuple, then this tuple 414 can be extended to infinitely many tuples over fv(Q) by choosing arbitrary values from the 415 infinite domain \mathcal{D} for the variables in the non-empty set $\mathsf{fv}(Q) \setminus \mathsf{fv}(Q_{\infty} \wedge Q^{=})$. Hence, for 416 every $(Q_{\not\infty}, Q^{=}) \in \mathcal{Q}_{fin}$ with disjointvars $(Q_{\not\infty}, Q^{=}) \neq \emptyset$ or $\mathsf{fv}(Q_{\not\infty} \land Q^{=}) \neq \mathsf{fv}(Q)$, we remove 417 $(Q_{\phi}, Q^{=})$ from \mathcal{Q}_{fin} and add $Q_{\phi} \wedge Q^{=}$ to \mathcal{Q}_{inf} . Note that we only remove pairs from \mathcal{Q}_{fin} , 418 hence, the loop on Lines 11–14 terminates. Afterwards, the query Q_{fin} is safe-range. However, 419 the query Q_{inf} need not be safe-range. Indeed, every query $Q_{\infty} \in Q_{inf}$ has range-restricted 420 bound variables, but not all the free variables of Q_{∞} need be range-restricted and thus 421 the query $\exists fv(Q_{\infty}), Q_{\infty}$ need not be safe-range. But the query Q_{inf} is closed and thus the 422 inf-equivalent query $\mathsf{rb}(Q_{inf})$ with range-restricted bound variables is safe-range. 423

▶ Lemma 8. Let Q be an RC query and $\operatorname{split}(Q) = (Q_{fin}, Q_{inf})$. Then the queries Q_{fin} and Q_{inf} are safe-range; $\operatorname{fv}(Q_{fin}) = \operatorname{fv}(Q)$ unless Q_{fin} is syntactically equal to \bot ; and $\operatorname{fv}(Q_{inf}) = \emptyset$.

⁴²⁶ ► Lemma 9. Let a structure S with an infinite domain D be fixed. Let Q be an RC query ⁴²⁷ and split(Q) = (Q_{fin}, Q_{inf}). If \models Q_{inf}, then $\llbracket Q \rrbracket$ is an infinite set. Otherwise, $\llbracket Q \rrbracket = \llbracket Q_{fin} \rrbracket$ ⁴²⁸ is a finite set.

By Lemma 8, Q_{fin} is a safe-range (and thus also domain-independent) query. Hence, for a 429 fixed structure S, the tuples in $[Q_{fin}]$ only contain elements in the active domain $\operatorname{adom}(Q_{fin})$, 430 i.e., $[\![Q_{fin}]\!] = [\![Q_{fin}]\!] \cap \mathsf{adom}(Q_{fin})^{|\mathsf{fv}(Q_{fin})|}$. Our translation does not introduce new constants in 431 Q_{fin} and thus $\mathsf{adom}(Q_{fin}) \subseteq \mathsf{adom}(Q)$. Hence, by Lemma 9, if $\not\models Q_{inf}$, then $\llbracket Q_{fin} \rrbracket$ is equal to 432 the "output-restricted unlimited interpretation" [15] of Q, i.e., $[\![Q_{fin}]\!] = [\![Q]\!] \cap \mathsf{adom}(Q)^{|\mathsf{fv}(Q)|}$. 433 In contrast, if $\models Q_{inf}$, then $\llbracket Q_{fin} \rrbracket = \llbracket Q \rrbracket \cap \mathsf{adom}(Q)^{|\mathsf{fv}(Q)|}$ does not necessarily hold. For 434 instance, for $Q := \neg \mathsf{B}(x)$, our translation yields $\mathsf{split}(Q) = (\bot, \top)$. In this case, we have $Q_{inf} =$ 435 \top and thus $\models Q_{inf}$ because $\neg \mathsf{B}(x)$ is satisfied by infinitely many tuples over an infinite domain. 436 However, if $\mathsf{B}(x)$ is never satisfied, then $\llbracket Q_{fin} \rrbracket = \emptyset$ is not equal to $\llbracket Q \rrbracket \cap \mathsf{adom}(Q)^{\lvert f_V(Q) \rvert}$. 437

Frample 10. Consider the query $Q := B(x) \vee P(x, y)$. The variable y is not rangerestricted in Q and thus split(Q) restricts y by a conjunction of Q with P(x, y). However, if $Q[y/\bot] = B(x)$ is satisfied by some tuple, then $[\![Q]\!]$ contains infinitely many tuples. Hence, split(Q) = ((B(x) ∨ P(x, y)) ∧ P(x, y), $\exists x. B(x)$). Because $Q_{fin} = (B(x) \lor P(x, y)) \land P(x, y)$ is only used if $\not\models Q_{inf}$, i.e., if B(x) is never satisfied, we could simplify Q_{fin} to P(x, y). However, our translation does not implement such heuristic simplifications.

Example 11. Consider the query $Q := B(x) \wedge u \approx v$. The variables u and v are not range-restricted in Q and thus split(Q) chooses one of these variables (e.g., u) and restricts

11:12 Practical Relational Calculus Query Evaluation

⁴⁴⁶ it by splitting Q into $Q_{\infty} = \mathsf{B}(x)$ and $Q^{=} = u \approx v$. Now, all variables are range-restricted ⁴⁴⁷ in Q_{∞} , but the variables in Q_{∞} and $Q^{=}$ are disjoint. Hence, $\llbracket Q \rrbracket$ contains infinitely many ⁴⁴⁸ tuples whenever Q_{∞} is satisfied by some tuple. In contrast, $\llbracket Q \rrbracket = \emptyset$ if Q_{∞} is never satisfied. ⁴⁴⁹ Hence, we have $\mathsf{split}(Q) = (\bot, \exists x. \mathsf{B}(x)).$

⁴⁵⁰ ► Example 12. Consider the query $Q_{user}^{susp} \coloneqq \mathsf{B}(b) \land \exists s. \forall p. \mathsf{P}(b, p) \longrightarrow \mathsf{S}(p, u, s)$ from Sec-⁴⁵¹ tion 1. Restricting its bound variables yields the query $\mathsf{rb}(Q_{user}^{susp}) = \mathsf{B}(b) \land ((\exists s. (\neg \exists p. \mathsf{P}(b, p) \land \neg \mathsf{S}(p, u, s)) \land (\exists p. \mathsf{S}(p, u, s))) \lor (\neg \exists p. \mathsf{P}(b, p)))$ derived in Example 6. Splitting Q_{user}^{susp} yields

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

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

 $\mathsf{rb}(Q_{user}^{susp}) \equiv (\mathsf{rb}(Q_{user}^{susp}) \land (\exists s, p. \, \mathsf{S}(p, u, s))) \lor (\mathsf{B}(b) \land (\neg \exists p. \, \mathsf{P}(b, p)) \land \neg \exists s, p. \, \mathsf{S}(p, u, s)).$

If the subquery $\mathsf{B}(b) \land (\neg \exists p. \mathsf{P}(b, p))$ 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 $\mathsf{S}(p, u, s)$ and thus satisfy the subquery $\neg \exists s, p. \mathsf{S}(p, u, s)$. Hence, $\left[Q_{user}^{susp}\right]^{\mathcal{S}} = \left[\mathsf{rb}(Q_{user}^{susp})\right]^{\mathcal{S}}$ is an infinite set of tuples whenever $\mathsf{B}(b) \land \neg \exists p. \mathsf{P}(b, p)$ is satisfied for some b. In contrast, if $\mathsf{B}(b) \land \neg \exists p. \mathsf{P}(b, p)$ is not satisfied for any b, then Q_{user}^{susp} is equivalent to $\mathsf{rb}(Q_{user}^{susp}) \land (\exists s, p. \mathsf{S}(p, u, s))$ obtained also by applying $(\mathbf{\hat{x}})$ to Q_{user}^{susp} for the free variable u.

▶ Definition 13. Let Q be an RC query and split(Q) = (Q_{fin}, Q_{inf}) . Let $\hat{Q}_{fin} \coloneqq$ sr2ranf (Q_{fin}) and $\hat{Q}_{inf} \coloneqq$ sr2ranf (Q_{inf}) be the equivalent RANF queries. We define rw(Q) \coloneqq $(\hat{Q}_{fin}, \hat{Q}_{inf})$.

464 4.4 Complexity Analysis

In this section, we analyze the time complexity of capturing Q, i.e., checking if [Q] is finite and 465 enumerating [Q] if it is finite. To bound the asymptotic time complexity of capturing a fixed Q, 466 we ignore the (constant) time complexity of computing $\mathsf{rw}(Q) = (\hat{Q}_{fin}, \hat{Q}_{inf})$ and focus on the 467 time complexity of evaluating the RANF queries \hat{Q}_{fin} and \hat{Q}_{inf} , i.e., the query cost of \hat{Q}_{fin} 468 and Q_{inf} . Without loss of generality, we assume that the input query Q has pairwise distinct 469 (free and bound) variables to derive a set of quantified predicates from Q's atomic predicates 470 and formulate our time complexity bound. Nevertheless, the RANF queries \hat{Q}_{fin} and \hat{Q}_{inf} 471 computed by our translation need not have pairwise distinct (free and bound) variables. 472

Let $\operatorname{av}(Q)$ be the set of all (free and bound) variables in a query Q. We define the relation \lesssim_Q on $\operatorname{av}(Q)$ such that $x \lesssim_Q y$ iff the scope of an occurrence of $x \in \operatorname{av}(Q)$ is contained in the scope of an occurrence of $y \in \operatorname{av}(Q)$. Formally, we define $x \lesssim_Q y$ iff $y \in \operatorname{fv}(Q)$ 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 $\operatorname{aps}(Q)$ be the set of all atomic predicates in a query Q. We denote by $\overline{\operatorname{qps}}(Q)$ the set of quantified predicates obtained from $\operatorname{aps}(Q)$ by performing the variable substitution $x \mapsto y$, where x and y are related by equalities in Q and $x \leq_Q y$, and existentially quantifying from a quantified predicate Q_{qp} the innermost bound variable x in Q that is free in Q_{qp} . Let $\operatorname{eqs}^*(Q)$ be the transitive closure of equalities occurring in Q. Formally, we define $\overline{\operatorname{qps}}(Q)$ by: $Q_{ap} \in \overline{\operatorname{qps}}(Q)$ if $Q_{ap} \in \operatorname{aps}(Q)$;

⁴⁸⁴ $Q_{qp}[x \mapsto y] \in \overline{qps}(Q) \text{ if } Q_{qp} \in \overline{qps}(Q), (x, y) \in eqs^*(Q), \text{ and } x \lesssim_Q y;$

$$\exists x. Q_{qp} \in \overline{\mathsf{qps}}(Q) \text{ if } Q_{qp} \in \overline{\mathsf{qps}}(Q), x \in \mathsf{fv}(Q_{qp}) \setminus \mathsf{fv}(Q), \text{ and } x \lesssim_Q y \text{ for all } y \in \mathsf{fv}(Q_{qp}).$$

We bound the complexity of capturing Q by considering subsets Q_{qps} of quantified predicates $\overline{qps}(Q)$ that are *minimal* in the sense that every quantified predicate in Q_{qps} contains a unique free variable that is not free in any other quantified predicate in Q_{qps} . Formally, we define minimal $(Q_{qps}) := \forall Q_{qp} \in Q_{qps}$. $\mathsf{fv}(Q_{qps} \setminus \{Q_{qps}\}) \neq \mathsf{fv}(Q_{qps})$. Every ⁴⁹⁰ minimal subset Q_{qps} of quantified predicates $\overline{qps}(Q)$ contributes the product of the numbers of ⁴⁹¹ tuples satisfying each quantified predicate $Q_{qp} \in Q_{qps}$ to the overall bound (that product is an ⁴⁹² upper bound on the number of tuples satisfying the join over all $Q_{qp} \in Q_{qps}$). Similarly to Ngo ⁴⁹³ et al. [22], we use the notation $\tilde{\mathcal{O}}(\cdot)$ to hide logarithmic factors incurred by set operations.

⁴⁹⁴ ► **Theorem 14.** Let Q be a fixed RC query with pairwise distinct (free and bound) variables. ⁴⁹⁵ The time complexity of capturing Q, i.e., checking if $\llbracket Q \rrbracket$ is finite and enumerating $\llbracket Q \rrbracket$ if it ⁴⁹⁶ is finite, is in $\tilde{\mathcal{O}}\left(\sum_{Q_{qps} \subseteq \overline{qps}(Q), \min[m] \in Q_{qps}} \prod_{Q_{qps} \in Q_{qps}} |\llbracket Q_{qps} \rrbracket\right)$.

We prove Theorem 14 in our extended report [25]. Examples 15 and 16 show that the time complexity from Theorem 14 cannot be achieved by the translation of Van Gelder and Topor [14] or over finite domains. Example 17 shows how equalities affect the bound in Theorem 14.

▶ **Example 15.** Consider the query $Q \coloneqq \mathsf{B}(b) \land \exists u, s. \neg \exists p. \mathsf{P}(b, p) \land \neg \mathsf{S}(p, u, s)$, equivalent 500 to Q^{susp} from Section 1. Then $aps(Q) = \{B(b), P(b, p), S(p, u, s)\}$ and $\overline{qps}(Q) = \{B(b), P(b, p), S(p, u, s)\}$ 501 $P(b,p), \exists p. P(b,p), S(p,u,s), \exists p. S(p,u,s), \exists s, p. S(p,u,s), \exists u, s, p. S(p,u,s)\}$. The translated 502 query Q_{vat} by Van Gelder and Topor [14] restricts the variables r and s by $\exists s, p. S(p, u, s)$ 503 and $\exists u, p. S(p, u, s)$, respectively. For an interpretation of B by $\{(c') \mid c' \in \{1, \ldots, n\}\}$, P by 504 $\{(c',c') \mid c' \in \{1,\ldots,n\}\}, \text{ and } S \text{ by } \{(c,c',c') \mid c \in \{1,\ldots,n\}, c' \in \{1,\ldots,m\}\}, n,m \in \mathbb{N}, d$ 505 computing the join of $\mathsf{P}(b,p)$, $\exists s, p. \mathsf{S}(p,u,s)$, and $\exists u, p. \mathsf{S}(p,u,s)$, which is a Cartesian 506 product, results in a time complexity in $\Omega(n \cdot m^2)$ for Q_{vgt} . In contrast, Theorem 14 yields 507 an asymptotically better time complexity in $\tilde{\mathcal{O}}(n+m+n\cdot m)$ for our translation: 508

$$\mathcal{O}\left(|[[\mathsf{B}(b)]]| + |[[\mathsf{P}(b,p)]]| + |[[\mathsf{S}(p,u,s)]]| + (|[[\mathsf{B}(b)]]| + |[[\mathsf{P}(b,p)]]|) \cdot |[[\mathsf{S}(p,u,s)]]|\right)$$

Example 16. The query $\neg S(x, y, z)$ is satisfied by a finite set of tuples over a finite domain \mathcal{D} (as is every other query over a finite domain). For an interpretation of S by $\{(c, c, c) \mid c \in \mathcal{D}\}$, the equality $|\mathcal{D}| = |[S(x, y, z)]|$ holds and the number of satisfying tuples is

$$|[\![\neg \mathsf{S}(x,y,z)]\!]| = |\mathcal{D}|^3 - |[\![\mathsf{S}(x,y,z)]\!]| = |[\![\mathsf{S}(x,y,z)]\!]|^3 - |[\![\mathsf{S}(x,y,z)]\!]| \in \Omega(|[\![\mathsf{S}(x,y,z)]\!]|^3),$$

which exceeds the bound $\tilde{\mathcal{O}}(|[[S(x, y, z)]]|)$ of Theorem 14. Hence, our infinite domain assumption is crucial for achieving the better complexity bound.

Example 17. Consider the following query over the domain $\mathcal{D} = \mathbb{N}$ of natural numbers:

$$Q \coloneqq \forall u. (u \approx 0 \lor u \approx 1 \lor u \approx 2) \longrightarrow \\ (\exists v. \mathsf{B}(v) \land (u \approx 0 \longrightarrow x \approx v) \land (u \approx 1 \longrightarrow y \approx v) \land (u \approx 2 \longrightarrow z \approx v))$$

Note that this query is equivalent to $Q \equiv B(x) \wedge B(y) \wedge B(z)$ and thus it is satisfied by a finite set of tuples of size $|[[B(x)]]| \cdot |[[B(y)]]| \cdot |[[B(z)]]| = |[[B(x)]]|^3$. The set of atomic predicates of Q is $aps(Q) = \{B(v)\}$ and it must be closed under the equalities occurring in Q to yield a valid bound in Theorem 14. In this case, $\overline{qps}(Q) = \{B(v), \exists v. B(v), B(x), B(y), B(z)\}$ and the bound in Theorem 14 is $|[[B(v)]]| \cdot |[[B(x)]]| \cdot |[[B(y)]]| \cdot |[[B(z)]]| = |[[B(x)]]|^4$. In particular, this bound is not tight, but it still reflects the complexity of evaluating the RANF queries produced by our translation as it does not derive the equivalence $Q \equiv B(x) \wedge B(y) \wedge B(z)$.

525 **5** Data Golf Benchmark

In this section, we devise the *Data Golf* benchmark for generating structures for given RC queries. We will use the benchmark in our empirical evaluation (Section 6). Given an RC query, we seek a structure that results in a nontrivial evaluation result for the overall query and for all its subqueries. Intuitively, the resulting structure makes query evaluation potentially more challenging compared to the case where some subquery results in a trivial (e.g., empty) evaluation result. More specifically, Data Golf has two objectives. The first

input: An RC query Q with pairwise distinct (free and bound) variables satisfying CON, CST, VAR, REP, a sequence of distinct variables \vec{v} , $fv(Q) \subseteq \vec{v}$, sets of tuples $\mathcal{T}_{\vec{v}}^+$ and $\mathcal{T}_{\vec{v}}^-$ over \vec{v} such that $|\mathcal{T}_{\vec{v}}^+[x]| = |\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-[x]| = |\mathcal{T}_{\vec{v}}^-|$, and $\mathcal{T}_{\vec{v}}^+[x] \cap \mathcal{T}_{\vec{v}}^-[x] = \emptyset, \text{ for every } \vec{x} \in \vec{v}, \text{ a parameter } \vec{\gamma} \in \{0,1\}.$ **output:** A structure \mathcal{S} such that $\mathcal{T}^+_{\vec{v}}[\vec{\mathsf{fv}}(Q)] \subseteq \llbracket Q \rrbracket, \mathcal{T}^-_{\vec{v}}[\vec{\mathsf{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$, and $|\llbracket Q' \rrbracket|$ and $|[\![\neg Q']\!]|$ contain at least min $\{|\mathcal{T}_{\vec{v}}^+|, |\mathcal{T}_{\vec{v}}^-|\}$ tuples, for every $Q' \sqsubseteq Q$. 1 function dg $(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma) =$ switch Q do $\mathbf{2}$ case $r(t_1,\ldots,t_{\iota(r)})$ do return $\{r^{\mathcal{S}} \mapsto \mathcal{T}^+_{\vec{v}}[t_1,\ldots,t_{\iota(r)}]\};$ 3 case $x \approx y$ do 4 if there exist d, d' such that $d \neq d'$ and $(d, d') \in \mathcal{T}^+_{\vec{v}}[x, y]$, or 5 d = d' and $(d, d') \in \mathcal{T}_{\vec{v}}^{-}[x, y]$ then fail; case $\neg Q'$ do return dg $(Q', \vec{v}, \mathcal{T}_{\vec{v}}^-, \mathcal{T}_{\vec{v}}^+, \gamma);$ 6 case $Q_1 \vee Q_2$ or $Q_1 \wedge Q_2$ do 7 $\begin{array}{l} (\mathcal{T}^1_{\vec{v}},\mathcal{T}^2_{\vec{v}}) \leftarrow \{(\mathcal{T}^1_{\vec{v}},\mathcal{T}^2_{\vec{v}}) \mid \left|\mathcal{T}^1_{\vec{v}}[x]\right| = \left|\mathcal{T}^2_{\vec{v}}[x]\right| = \left|\mathcal{T}^2_{\vec{v}}\right| = \min\{\left|\mathcal{T}^+_{\vec{v}}\right|, \left|\mathcal{T}^-_{\vec{v}}\right|\},\\ \mathcal{T}^1_{\vec{v}}[x] \cap \mathcal{T}^2_{\vec{v}}[x] = \emptyset, (\mathcal{T}^1_{\vec{v}}[x] \cup \mathcal{T}^2_{\vec{v}}[x]) \cap (\mathcal{T}^+_{\vec{v}}[x] \cup \mathcal{T}^-_{\vec{v}}[x]) = \emptyset, \text{for all } x \in \vec{v}\}; \end{array}$ 8 if $\gamma = 0$ then 9 $|\operatorname{\mathbf{return}} \mathsf{dg}(Q_1, \vec{v}, \mathcal{T}^+_{\vec{v}} \cup \mathcal{T}^1_{\vec{v}}, \mathcal{T}^-_{\vec{v}} \cup \mathcal{T}^2_{\vec{v}}, \gamma) \cup \mathsf{dg}(Q_2, \vec{v}, \mathcal{T}^+_{\vec{v}} \cup \mathcal{T}^2_{\vec{v}}, \mathcal{T}^-_{\vec{v}} \cup \mathcal{T}^1_{\vec{v}}, \gamma);$ 10 else 11 switch Q do 12 case $Q_1 \vee Q_2$ do 13 $\big| \operatorname{\mathbf{return}} \, \mathsf{dg}(Q_1, \vec{v}, \mathcal{T}^+_{\vec{v}} \cup \mathcal{T}^1_{\vec{v}}, \mathcal{T}^-_{\vec{v}} \cup \mathcal{T}^2_{\vec{v}}, \gamma) \cup \mathsf{dg}(Q_2, \vec{v}, \mathcal{T}^1_{\vec{v}} \cup \mathcal{T}^2_{\vec{v}}, \mathcal{T}^-_{\vec{v}} \cup \mathcal{T}^+_{\vec{v}}, \gamma);$ 14 case $Q_1 \wedge Q_2$ do 15 $\big| \operatorname{\mathbf{return}} \, \mathsf{dg}(Q_1, \vec{v}, \mathcal{T}^+_{\vec{v}} \cup \mathcal{T}^-_{\vec{v}}, \mathcal{T}^1_{\vec{v}} \cup \mathcal{T}^2_{\vec{v}}, \gamma) \cup \mathsf{dg}(Q_2, \vec{v}, \mathcal{T}^+_{\vec{v}} \cup \mathcal{T}^2_{\vec{v}}, \mathcal{T}^-_{\vec{v}} \cup \mathcal{T}^1_{\vec{v}}, \gamma);$ 16 case $\exists y. Q_y$ do 17 $(\mathcal{T}^1_{\vec{v}\cdot y}, \mathcal{T}^2_{\vec{v}\cdot y}) \leftarrow \{(\mathcal{T}^1_{\vec{v}\cdot y}, \mathcal{T}^2_{\vec{v}\cdot y}) \mid \mathcal{T}^1_{\vec{v}\cdot y}[\vec{v}] = \mathcal{T}^+_{\vec{v}}, \mathcal{T}^2_{\vec{v}\cdot y}[\vec{v}] = \mathcal{T}^-_{\vec{v}},$ 18

$$\left| \left| \left| \mathcal{T}^{1}_{\vec{v} \cdot y}[y] \right| = \left| \mathcal{T}^{1}_{\vec{v} \cdot y} \right| = \left| \mathcal{T}^{+}_{\vec{v}} \right|, \left| \mathcal{T}^{2}_{\vec{v} \cdot y}[y] \right| = \left| \mathcal{T}^{2}_{\vec{v} \cdot y} \right| = \left| \mathcal{T}^{-}_{\vec{v}} \right|, \mathcal{T}^{1}_{\vec{v} \cdot y}[y] \cap \mathcal{T}^{2}_{\vec{v} \cdot y}[y] = \emptyset \};$$

19 | | return dg
$$(Q_y, \vec{v} \cdot y, \mathcal{T}^1_{\vec{v} \cdot y}, \mathcal{T}^2_{\vec{v} \cdot y}, \gamma)$$

Figure 6 Computing the Data Golf structure.

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}^- over a fixed domain \mathcal{D} , representing assignments of $\mathsf{fv}(Q)$, Data Golf produces a structure \mathcal{S} (represented as a partial mapping from predicate symbols to their interpretations), such that $\mathcal{T}^+ \subseteq \llbracket Q \rrbracket, \mathcal{T}^- \cap \llbracket Q \rrbracket = \emptyset$, and $|\llbracket Q' \rrbracket|$ and $|\llbracket \neg Q' \rrbracket|$ contain at least min $\{|\mathcal{T}^+|, |\mathcal{T}^-|\}$ tuples, for every $Q' \sqsubseteq Q$. To be able to produce such a structure \mathcal{S} , we make the following assumptions on Q:

⁵⁴¹ CON the bound variable y in every subquery $\exists y. Q_y$ of Q satisfies $\operatorname{con}(y, Q_y, \mathcal{G})$ [14, Figure 5] ⁵⁴² for some set \mathcal{G} such that $\operatorname{eqs}(y, \mathcal{G}) = \emptyset$ and, for every $Q_{qp} \in \mathcal{G}$, $\{y\} \subsetneq \operatorname{fv}(Q_{qp})$ holds; ⁵⁴³ this avoids subqueries like $\exists y. \neg \mathsf{P}_2(x, y)$ and $\exists y. (\mathsf{P}_2(x, y) \lor \mathsf{P}_1(y));$

- 544 CST Q contains no subquery of the form $x \approx c$, which is satisfied by exactly one tuple;
- VAR Q contains no closed subqueries, e.g., $P_1(42)$, because a closed subquery is either satisfied by all possible tuples or no tuple at all; and
- ⁵⁴⁷ REP Q contains no repeated predicate symbols; this avoids subqueries like $\mathsf{P}_1(x) \land \neg \mathsf{P}_1(x)$.

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

Data Golf is formalized by the function $dg(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$, defined in Figure 6, where 553 \vec{v} is a sequence of distinct variables such that $\mathsf{fv}(Q) \subseteq \vec{v}, \mathcal{T}^+_{\vec{v}}$ and $\mathcal{T}^-_{\vec{v}}$ are sets of tuples over \vec{v} , and $\gamma \in \{0,1\}$ is a *strategy*. The function $\mathsf{dg}(Q, \vec{v}, \mathcal{T}^+_{\vec{v}}, \mathcal{T}^-_{\vec{v}}, \gamma)$ can fail on an equality 554 555 between two variables $x \approx y$. In this case, the function $dg(Q, \vec{v}, \vec{T}_{\vec{u}}^+, T_{\vec{u}}^-, \gamma)$ does not compute 556 a Data Golf structure. We define the *not-depth* of a subquery $x \approx y$ in Q as the number of 557 subqueries that have the form of a negation among the queries $x \approx y \sqsubseteq \cdots \sqsubseteq Q$, i.e., the 558 number of negations on the path between the subquery $x \approx y$ and Q's main connective. To 559 prevent failure, we generate the sets $\mathcal{T}_{\vec{v}}^+$, $\mathcal{T}_{\vec{v}}^-$ to only contain tuples with equal values for all 560 variables in equalities with even (odd, respectively) not-depth and pairwise distinct values 561 for all variables in equalities with odd (even, respectively) not-depth. This is not always 562 possible, e.g., for $x \approx y \land \neg x \approx y$, in which case no Data Golf structure can be computed. 563 In the case of a conjunction or a disjunction, we add disjoint sets $\mathcal{T}_{\vec{u}}^1, \mathcal{T}_{\vec{v}}^2$ of tuples over \vec{v} 564 to $\mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-$ so that the intermediate results for the subqueries are neither equal nor disjoint. 565 We implement two strategies (parameter γ) to choose these sets $\mathcal{T}_{\vec{v}}^1, \mathcal{T}_{\vec{v}}^2$. 566

Finally, we justify why a Data Golf structure S computed by $dg(Q, \vec{v}, \mathcal{T}^+_{\vec{v}}, \mathcal{T}^-_{\vec{v}}, \gamma)$ satisfies 567 $\mathcal{T}^+_{\vec{v}}[\vec{\mathsf{fv}}(Q)] \subseteq \llbracket Q \rrbracket$ and $\mathcal{T}^-_{\vec{v}}[\vec{\mathsf{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$. We proceed by induction on the query Q. 568 Because of REP, the Data Golf structures for the subqueries Q_1, Q_2 of a binary query $Q_1 \lor Q_2$ 569 or $Q_1 \wedge Q_2$ can be combined using the union operator. The only case that does not follow 570 immediately is that $\mathcal{T}_{\vec{v}}^{-}[\vec{\mathsf{fv}}(Q)] \cap \llbracket Q \rrbracket = \emptyset$ for a query Q of the form $\exists y. Q_y$. We prove this 571 case by contradiction. Without loss of generality we assume that $\vec{\mathsf{fv}}(Q_y) = \vec{\mathsf{fv}}(Q) \cdot y$. Suppose 572 that $d \in \mathcal{T}_{\vec{v}}^{-}[\mathsf{fv}(Q)]$ and $d \in \llbracket Q \rrbracket$. Because $d \in \mathcal{T}_{\vec{v}}^{-}[\mathsf{fv}(Q)]$, there exists some d such that 573 $\vec{d} \cdot d \in \mathcal{T}^2_{\vec{v} \cdot y}[\vec{\mathsf{fv}}(Q_y)]$. Because $\vec{d} \in \llbracket Q \rrbracket$, there exists some d' such that $\vec{d} \cdot d' \in \llbracket Q_y \rrbracket$. By the 574 induction hypothesis, $\vec{d} \cdot d \notin [\![Q_y]\!]$ and $\vec{d} \cdot d' \notin \mathcal{T}^2_{\vec{v} \cdot y}[\vec{\mathsf{fv}}(Q_y)]$. Because $\mathsf{con}(y, Q_y, \mathcal{G})$ holds for 575 some \mathcal{G} satisfying CON, the query Q_y is equivalent to $(Q_y \wedge qps^{\vee}(\mathcal{G})) \vee Q_y[y/\perp]$. We have 576 $\vec{d} \cdot d' \in \llbracket Q_y \rrbracket$. If the tuple $\vec{d} \cdot d'$ satisfies $Q_y[y/\bot]$, then $\vec{d} \cdot d \in \llbracket Q_y \rrbracket$ (contradiction) because 577 the variable y does not occur in the query $Q_y[y/\perp]$ and thus its assignment in $d \cdot d'$ can be 578 arbitrarily changed. Otherwise, the tuple $\vec{d} \cdot d'$ satisfies some quantified predicate $Q_{qp} \in qps(\mathcal{G})$ 579 and (CON) implies $\{y\} \subsetneq \mathsf{fv}(Q_{qp})$. Hence, the tuples $d \cdot d$ and $d \cdot d'$ agree on the assignment of a 580 variable $x \in \mathsf{fv}(Q_{qp}) \setminus \{y\}$. Let $\mathcal{T}_{\vec{v}'}^+$ and $\mathcal{T}_{\vec{v}'}^-$ be the sets in the recursive call of dg on the atomic predicate from Q_{qp} . Because $\vec{d} \cdot \vec{d} \in \mathcal{T}_{\vec{v} \cdot y}^2[\vec{\mathsf{fv}}(Q_y)]$ and $\mathcal{T}_{\vec{v} \cdot y}^2[\vec{\mathsf{fv}}(Q_y)] \subseteq \mathcal{T}_{\vec{v}'}^+[\vec{\mathsf{fv}}(Q_y)] \cup \mathcal{T}_{\vec{v}'}^-[\vec{\mathsf{fv}}(Q_y)]$, 581 582 the tuple $\vec{d} \cdot d$ is in $\mathcal{T}^+_{\vec{v}'}[\vec{\mathsf{fv}}(Q_y)] \cup \mathcal{T}^-_{\vec{v}'}[\vec{\mathsf{fv}}(Q_y)]$. Because $\vec{d} \cdot d'$ satisfies the quantified predicate 583 Q_{qp} , the tuple $\vec{d} \cdot d'$ is in $\mathcal{T}^+_{\vec{v}'}[\vec{\mathsf{fv}}(Q_y)]$. Next we observe that the assignments of every variable 584 (in particular, x) in the tuples from the sets $\mathcal{T}_{\vec{v}'}^+$, $\mathcal{T}_{\vec{v}'}^-$ are pairwise distinct (the conditions 585 $\mathcal{T}^+_{\vec{v}'}[x] \cap \mathcal{T}^-_{\vec{v}'}[x] = \emptyset, \ \left|\mathcal{T}^+_{\vec{v}'}[x]\right| = \left|\mathcal{T}^+_{\vec{v}'}\right|, \text{ and } \left|\mathcal{T}^-_{\vec{v}'}[x]\right| = \left|\mathcal{T}^-_{\vec{v}'}\right|. \text{ Because the tuples } \vec{d} \cdot d \text{ and } \vec{d} \cdot d'$ 586 agree on the assignment of x, they must be equal, i.e., $\vec{d} \cdot d = \vec{d} \cdot d'$ (contradiction). 587

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

Example 18. Consider the query $Q \coloneqq \neg \exists y. \mathsf{P}_2(x, y) \land \neg \mathsf{P}_3(x, y, z)$. This query Q satisfies the assumptions CON, CST, VAR, REP. In particular, $\operatorname{con}(y, \mathsf{P}_2(x, y) \land \neg \mathsf{P}_3(x, y, z), \mathcal{G})$ holds for $\mathcal{G} = \{\mathsf{P}_2(x, y)\}$ with $\{y\} \subsetneq \operatorname{fv}(\mathsf{P}_2(x, y))$. We choose $\vec{v} = (x, z), \mathcal{T}_{\vec{v}}^+ = \{(0, 4), (2, 6)\}$, and $\mathcal{T}_{\vec{v}}^- = \{(8, 12), (10, 14)\}$. The function $\operatorname{dg}(Q, \vec{v}, \mathcal{T}_{\vec{v}}^+, \mathcal{T}_{\vec{v}}^-, \gamma)$ first flips $\mathcal{T}_{\vec{v}}^+$ and $\mathcal{T}_{\vec{v}}^-$ (because Q's main connective is negation) and then extends the tuples in the sets $\mathcal{T}_{\vec{v}}^-$ and $\mathcal{T}_{\vec{v}}^+$ with a value

11:16 Practical Relational Calculus Query Evaluation

⁵⁹⁶ for the bound variable y: $\mathcal{T}^1_{\vec{v}\cdot y} = \{(8, 12, 16), (10, 14, 18)\}$ and $\mathcal{T}^2_{\vec{v}\cdot y} = \{(0, 4, 20), (2, 6, 22)\}.$

For conjunction (a binary operator), two additional sets of tuples are computed: $\overline{\mathcal{T}^1_{\vec{v}\cdot y}} =$

⁵⁹⁸ {(24, 28, 32), (26, 30, 34)} and $\overline{\mathcal{T}_{v\cdot y}^2} = \{(36, 40, 44), (38, 42, 46)\}$. Depending on the strategy ⁵⁹⁹ ($\gamma = 0 \text{ or } \gamma = 1$), one of the following structures is computed: $\mathcal{S}_0 = \{\mathsf{P}_2 \mapsto \{(8, 16), (10, 18), (24, 32), (26, 34)\}, \mathsf{P}_3 \mapsto \mathcal{T}_{xyz}^+\}, \text{ or } \mathcal{S}_1 = \{\mathsf{P}_2 \mapsto \{(8, 16), (10, 18), (0, 20), (2, 22)\}, \mathsf{P}_3 \mapsto \mathcal{T}_{xyz}^+\},$ ⁶⁰¹ where $\mathcal{T}_{xyz}^+ = \{(0, 20, 4), (2, 22, 6), (24, 32, 28), (26, 34, 30)\}.$

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

6 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 14), we further improve its average-case complexity by implementing the following optimizations, described in more detail in our extended report [25, Section E].

We use a sample structure of constant size, called a *training database*, to estimate the query 611 cost when resolving the nondeterministic choices in our algorithms. A good training data-612 base should preserve the relative ordering of queries by their cost over the actual database 613 as much as possible. Nevertheless, our translation satisfies the correctness and worst-case 614 complexity claims (Section 4.3 and 4.4) for every choice of the training database. All our 615 experiments used a Data Golf structure with $|\mathcal{T}^+| = |\mathcal{T}^-| = 2$ as the training database. 616 We use the function optcnt optimizing RANF subqueries of the form $\exists \vec{y}. Q^+ \wedge \bigwedge_{i=1}^k \neg Q_i^-$ 617 using the count aggregation operator. Inspired by Claußen et al. [9], we compare the number of assignments of \vec{y} that satisfy Q^+ and $\bigvee_{i=1}^k (Q^+ \wedge Q_i^-)$, respectively. 618 619

To compute an SQL query from a RANF query, we define the function ranf2sql(·). 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 [30]. 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 $fv(\hat{Q}_2) \subsetneq fv(\hat{Q}_1)$, and we perform common subquery elimination.

To validate our translation's improved asymptotic time complexity, we compare it with 626 the translation by Van Gelder and Topor [14] (VGT), an implementation of the algorithm 627 by Ailamazyan et al. [2] that uses an extended active domain as the generators, and the 628 DDD [20,21], LDD [7], and MonPoly^{REG} [4] tools that support direct RC query evaluation using 629 binary decision diagrams. We could not find a publicly available implementation of Van Gelder 630 and Topor's translation. Therefore, the tool VGT for evaluable RC queries is derived from our 631 implementation by modifying the function $rb(\cdot)$ in Figure 4 to use the con relation [14, Figure 5] 632 instead of $\operatorname{cov}(x, Q, \mathcal{G})$ (Figure 3) and to use the generator $\exists \overline{\mathsf{fv}}(Q) \setminus \{x\}$. $\operatorname{qps}^{\vee}(\mathcal{G})$ instead of 633 $qps^{\vee}(\mathcal{G})$. Evaluable queries Q are always translated into (Q_{fin}, \perp) by $rw(\cdot)$ because all of Q's 634 free variables are range-restricted. We also consider translation variants that omit the count 635 aggregation optimization $optcnt(\cdot)$, marked with a minus (⁻). 636

SQL queries computed by the translations are evaluated using the PostgreSQL database engine. We have also used the MySQL database engine but omit its timings 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

Experiment	SMA	LL, Ev	valuabl	le pseu	dorand	iom qi	ueries	Q, su	b(Q)	= 14	n, n = 500:
RC2SQL	0.3	0.3	0.3	0.2	0.2	0.3	0.3	0.2	0.2	0.3	
$RC2SQL^{-}$	0.3	0.2	150.3	0.3	0.2	0.3	0.3	0.3	5.9	0.2	
VGT	31.5	6.7	4.2	2.5	37.5	9.3	2.4	2.3	11.3	2.7	
VGT^{-}	33.7	4.8	119.9	6.3	11.2	21.9	31.4	11.3	12.3	21.9	
DDD	9.1	2.5	RE	7.1	5.9	RE	5.1	RE	2.2	5.1	
LDD	59.2	24.1	169.1	38.8	53.3	37.4	64.0	ТО	16.0	61.6	
MonPoly ^{REG}	64.2	31.4	143.0	57.6	67.8	54.4	72.4	174.6	33.6	71.3	
Experiment MEDIUM, Evaluable pseudorandom queries Q , $ sub(Q) = 14$, $n = 20000$:											
RC2SQL	2.6	1.4	3.9	2.1	1.5	2.8	3.3	1.6	1.2	2.6	
RC2SQL ⁻	2.0	1.0	ТО	2.0	1.7	2.5	2.3	1.8	ТО	1.8	
VGT	TO	ТО	7.8	3.9	ТО	ТО	5.2	4.7	ТО	4.8	
VGT^{-}	TO	ТО	ТО	ТО	ТО	ТО	ТО	ТО	ТО	ТО	
Experiment	LARG	GE, E	valuab	le pseu	dorano	dom q	ueries	Q, sl	ub(Q)	= 14	I, tool = RC2SQL:
n = 40000	3.5	2.7	8.1	4.0	3.2	5.5	6.7	4.1	1.9	5.8	
n = 80000	7.5	5.4	16.1	8.0	6.1	11.5	14.0	8.1	4.2	11.7	
n=120000	13.2	8.2	24.6	11.5	8.9	16.3	20.9	11.0	7.2	16.7	
Experiment	INFIN	NITE,	Non-ev	valuabl	e pseu	doran	dom o	queries	Q, s	ub(Q)	=7, n=4000:
	Infinite results $(\gamma = 0)$ Finite results $(\gamma = 1)$										
RC2SQL	0.8	0.8	0.8	0.8	0.8	1.0	1.1	0.9	2.4	1.1	
$RC2SQL^{-}$	0.5	0.5	0.4	0.5	0.5	0.6	0.7	0.6	ТО	2.0	
DDD	89.5	49.1	46.9	116.3	50.4	81.7	44.1	45.8	89.8	44.6	
LDD	TO	ТО	ТО	ТО	ТО	TO	ТО	ТО	ТО	ТО	
$MonPoly^{REG}$	ТО	ТО	ТО	ТО	ТО	ТО	ТО	ТО	ТО	ТО	

Figure 7 Experiments SMALL, MEDIUM, LARGE, and INFINITE. We use the following abbreviations: TO = Timeout of 300s, RE = Runtime Error.

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 the SMALL, MEDIUM, and LARGE experiments, we generate ten pseudorandom queries 644 with a fixed size 14 and Data Golf structures \mathcal{S} . The queries satisfy the Data Golf assumptions 645 along with a few additional ones: the queries are not safe-range, have no repeated equalities, 646 disjunction only appears at the top-level, every bound variable actually occurs in its scope, and 647 only pairwise distinct variables appear as terms in predicates. The queries have 2 free variables 648 and every subquery has at most 4 free variables. We control the size of the Data Golf structure 649 S in our experiments using a parameter $n = |\mathcal{T}^+| = |\mathcal{T}^-|$. Because the sets \mathcal{T}^+ and \mathcal{T}^- grow 650 in the recursion on subqueries, relations in a Data Golf structure typically have more than n651 tuples. The values of the parameter n for Data Golf structures are summarized in Figure 7. 652 The INFINITE experiment consists of five pseudorandom queries Q that are *not* evalu-653 able and $\mathsf{rw}(Q) = (Q_{fin}, Q_{inf})$, where $Q_{inf} \neq \bot$. Specifically, the queries are of the form 654 $Q_1 \wedge \forall x, y, Q_2 \longrightarrow Q_3$, where Q_1, Q_2 , and Q_3 are either atomic predicates or equalities. For 655 each query Q, we compare the performance of our tool to tools that directly evaluate Q on 656 structures generated by the two Data Golf strategies (parameter γ), which trigger infinite 657 or finite evaluation results on the considered queries. For infinite results, our tool outputs 658 this fact (by evaluating Q_{inf}), whereas the other tools also output a finite representation 659 of the infinite result. For finite results, all tools produce the same output. 660

Figure 7 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. The

11:18 Practical Relational Calculus Query Evaluation

Query Q^{susp}		Q_{user}^{susp}		Q_{text}^{susp}		$ $ Query $ $ Q^s		usp	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	2.0	2.2	3.0	3.5	6.2	7.1	RC2SQL	2.9	16.2	4.2	21.4	8.9	91.3
$RC2SQL^{-}$	61.7	ТО	63.4	ТО	484.9	ТО	$RC2SQL^{-}$	273.9	TO	270.1	ТО	ТО	TO
VGT	3.9	2.9	-	_	213.2	ТО	VGT	3.5	18.9	_	_	ТО	ТО
VGT^{-}	433.8	ТО	-	_	495.4	ТО	VGT^{-}	ТО	ТО	_	_	ТО	ТО
DDD	7.1	ТО	6.3	ТО	28.8	TO	DDD	93.3	TO	90.1	TO	178.5	TO
LDD	36.3	ТО	34.0	ТО	213.9	ТО	LDD	ТО	ТО	ТО	ТО	ТО	ТО
$MonPoly^{REG}$	49.9	ТО	47.3	ТО	181.2	ТО	MonPoly ^{REG}	TO	ТО	TO	ТО	ТО	TO

Figure 8 Experiment with the queries Q^{susp} , Q^{susp}_{user} , and Q^{susp}_{text} . We use the following abbreviations: GC = Gift Cards dataset, MI = Musical Instruments dataset, TO = Timeout of 600s.

⁶⁶⁴ lowest time for a query is typeset in bold. We do not report the translation time because it ⁶⁶⁵ does not contribute to the time complexity for a fixed query. Still, RC2SQL's translation time ⁶⁶⁶ is at most 0.6 seconds on every query in our experiments. We also omit the rows for tools that ⁶⁶⁷ time out or crash on all queries of an experiment, e.g., Ailamazyan et al. [2]. We conclude ⁶⁶⁸ that our translation RC2SQL significantly outperforms all other tools on all queries and scales ⁶⁶⁹ 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^{susp}_{user} from the introduction and on 670 the more challenging query $Q_{text}^{susp} := \mathsf{B}(b) \land \exists u, s, t. \forall p. \mathsf{P}(b, p) \longrightarrow \mathsf{S}(p, u, s) \lor \mathsf{T}(p, u, t)$ with 671 an additional relation T that relates user's review text (variable t) to a product. The query 672 Q_{text}^{susp} computes all brands for which there is a user, a score, and a review text such that all 673 the brand's products were reviewed by that user with that score or by that user with that 674 text. We use both Data Golf structures (strategy $\gamma = 1$) and real-world structures obtained 675 from the Amazon review dataset [23]. The real-world relations P, S, and T are obtained by 676 projecting the respective tables from the Amazon review dataset for some chosen product 677 categories (abbreviated GC and MI in Figure 8) and the relation B contains all brands from 678 P that have at least three products. Because the tool by Ailamazyan et al., DDD, LDD, and 679 MonPoly^{REG} only support integer data, we injectively remap the string and floating-point 680 values from the Amazon review dataset to integers. 681

Figure 8 shows the empirical evaluation results: execution times on Data Golf structures 682 (left) and execution times on structures derived from the real-world dataset for two specific 683 product categories (right). We remark that VGT cannot handle the query Q_{user}^{susp} as it is not 684 evaluable [14]. Our translation RC2SQL significantly outperforms all other tools (except VGT 685 on Q^{susp} , but RC2SQL still outperforms VGT) on both Data Golf and real-world structures. 686 VGT^- translates Q^{susp} into a RANF query with a higher query cost than $RC2SQL^-$. How-687 ever, the optimization $\mathsf{optcnt}(\cdot)$ manages to rectify this inefficiency and thus VGT exhibits 688 a comparable performance as RC2SQL. Specifically, the factor of $80 \times$ in query cost between 689 VGT^- and $RC2SQL^-$ improves to $1.1 \times$ in query cost between VGT and RC2SQL on a Data 690 Golf structure with n = 20 [25]. Nevertheless, VGT does not finish evaluating the query 691 Q_{text}^{susp} on GC and MI datasets within 10 minutes, unlike RC2SQL. 692

693 7 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, or aggregations.

699		References
700	1	Serge Abiteboul, Richard Hull, and Victor Vianu. Foundations of Databases. Addison-Wesley,
701		1995. URL: http://webdam.inria.fr/Alice/.
702	2	Alfred K. Ailamazyan, Mikhail M. Gilula, Alexei P. Stolboushkin, and Grigorii F. Schwartz.
703		Reduction of a relational model with infinite domains to the case of finite domains. <i>Doklady</i>
704		Akademii Nauk SSSR, 286(2):308-311, 1986. URL: http://mi.mathnet.ru/dan47310.
705	3	Arnon Avron and Yoram Hirshfeld. On first order database query languages. In LICS. July
706		15-18, 1991, Amsterdam, The Netherlands, pages 226–231. IEEE Computer Society, 1991.
707		doi:10.1109/LICS.1991.151647.
708	4	David A. Basin, Felix Klaedtke, Samuel Müller, and Eugen Zalinescu. Monitoring metric
709	-	first-order temporal properties. $L ACM, 62(2):15:1-15:45, 2015, doi:10.1145/2699444.$
710	5	Michael Benedikt and Leonid Libkin. Belational queries over interpreted structures. J. ACM.
711	Ū	47(4):644–680. 2000. doi:10.1145/347476.347477.
712	6	Achim Blumensath and Erich Grädel Finite presentations of infinite structures: Auto-
712	Ŭ	mata and interpretations Theory Comput Syst 37(6):641–674 2004 doi:10.1007/
714		s00224-004-1133-v
715	7	Sagar Chaki Arie Gurfinkel and Ofer Strichman Decision diagrams for linear arithmetic
715	'	In EMCAD 15-18 November 2009 Austin Teras USA pages 53-60 IEEE 2009 doi:
710		10 1109/FMCAD 2009 5351143
710	8	Ian Chomicki and David Toman. Implementing temporal integrity constraints using an active
710	0	DBMS IEEE Trans Knowl Data Eng. 7(4):566-582 1005 doi:10.1109/69.404030
719	0	Long Claußen Alfons Kompor Guide Moarkette and Klaus Poithner Ontimizing queries
720	9	with universal quantification in object oriented and object relational databases. In Matthias
721		Jarka Michael I. Caroy Klaus B. Dittrich Frederick H. Lochovsky, Parieles Loucopoulos
722		and Manfred A Jourfold editors VLDB August 25.20 1007 Athene Greece pages 286-205
723		Morgan Kaufmann 1007 URL: http://www.wldb.org/conf/1007/D286 DDE
724	10	F. E. Codd. Rolational completeness of data base sublanguages. Research Report / RI / IRM
725	10	/ San Lose California B 1087 1072
720	11	Frling Ellingson Boroy golf 2013 https://olf nu/PogoyColf
727	12	Martha Escapar Molano, Richard Hull, and Dean Jacoba, Safaty and translation of calculus
728	12	quories with scalar functions. In Catriel Boori editor PODS May 25-28 1002 Washington
729		DC USA pages 252–264 ACM Pross 1003 doi:10.1145/153850.153000
730	12	Allen Van Colder and Bodney W. Topor. Sofety and correct translation of relational calculus
731	15	formulas. In Mosho V. Vardi aditor, PODS, March 22 25, 1087, San Diago, California, USA
732		pages 212 227 ACM 1087 doi:10.1145/28660.28602
733	1/	Allen Van Calder and Bodney W. Toner, Safety and translation of relational calculus queries
734	14	Alleli Vali Gelder and Rodney W. Topor. Safety and translation of relational calculus queries.
735	15	Pichard Hull and Jianwan Su. Domain independence and the relational colorly.
730	13	Informatica 21(6),512 524 1004 doi:10.1007/PE01212204
/3/	16	Mighael Kifer, On safety domain independence, and contumbility of detabase queries (prelim
738	10	inary report). In Catriel Boori, Joachim W. Schmidt, and Umschwar Daval, editors, Proceedings
739		of the Third International Conference on Date and Knowledge Bases: Improving Hashility and
740		Begrandium and June 28 20 1088 Jerusalem Jerus Langel pages 405 415 Margan Kaufmann 1088
741		doi:10_1016/b079_1_4930_1312_0_E0037_9
742	17	Nils Klarlund and Anders Meller MONA at / Hear Manual BRICS Department of Computer
743	11	Science, University of Aarbus, January 2001, URL: http://www.bricg.db/marc/
/44	10	Loopid Liblin Elements of Einite Model Theory. Touts in Theoretical Commuter Commuter
745	10	An EATCS Sories Springer 2004 UDL, http://www.co.touriet.com/W7Eliblic/art.doi
746		An EATOS Series. Springer, 2004. URL: http://www.cs.toronto.edu/%/E110kin/Imt, dol: 10.1007/079-2-662-07002-1
747	10	10.1001/510-3-002-01003-1. Hong Chou Liu Loffron Yu Yu and Woife Liang Safety domain independence and translation
748	19	of complex value database quories Inf Sci 178(19):9507 9532 2008 doi:10.1016/j.joc
750		2008 02 005
100		

11:20 Practical Relational Calculus Query Evaluation

Jesper B. Møller. DDDLIB: A library for solving quantified difference inequalities. In Andrei
 Voronkov, editor, *CADE*, July 27-30, 2002, Copenhagen, Denmark, volume 2392 of Lecture
 Notes in Computer Science, pages 129–133. Springer, 2002. doi:10.1007/3-540-45620-1_9.

- Jesper B. Møller, Jakob Lichtenberg, Henrik Reif Andersen, and Henrik Hulgaard. Difference decision diagrams. In Jörg Flum and Mario Rodríguez-Artalejo, editors, CSL, September 20-25, 1999, Madrid, Spain, volume 1683 of Lecture Notes in Computer Science, pages 111–125.
 Springer, 1999. doi:10.1007/3-540-48168-0_9.
- Hung Q. Ngo, Christopher Ré, and Atri Rudra. Skew strikes back: new developments in the theory of join algorithms. SIGMOD Rec., 42(4):5–16, 2013. doi:10.1145/2590989.2590991.
- Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying recommendations using distantly labeled reviews and fine-grained aspects. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun
 Wan, editors, *EMNLP, November 3-7, 2019, Hong Kong, China*, pages 188–197. Association
 for Computational Linguistics, 2019. doi:10.18653/v1/D19-1018.
- Robert A. Di Paola. The recursive unsolvability of the decision problem for the class of definite formulas. J. ACM, 16(2):324-327, 1969. doi:10.1145/321510.321524.
- Martin Raszyk, David Basin, Srðan Krstić, and Dmitriy Traytel. Implementation, evaluation,
 and extended report associated with this paper, 2022. https://github.com/rc2sql/rc2sql.
- Peter Z. Revesz. Introduction to Constraint Databases. Texts in Computer Science. Springer,
 2002. doi:10.1007/b97430.
- Boris A Trakhtenbrot. Impossibility of an algorithm for the decision problem in finite classes.
 Doklady Akademii Nauk SSSR, 70(4):569–572, 1950.
- Moshe Y. Vardi. The decision problem for database dependencies. Inf. Process. Lett., 12(5):251–254, 1981. doi:10.1016/0020-0190(81)90025-9.
- Moshe Y. Vardi. The complexity of relational query languages (extended abstract). In
 Harry R. Lewis, Barbara B. Simons, Walter A. Burkhard, and Lawrence H. Landweber,
 editors, STOC, May 5-7, 1982, San Francisco, California, USA, pages 137–146. ACM, 1982.
 doi:10.1145/800070.802186.
- 778 30 Jun Yang. radb, 2019. https://github.com/junyang/radb.