# Analysis and Formalization of Typing in SQL Engines

Wenjia Ye
The University of Hong Kong
Hong Kong, China
yewenjia@connect.hku.hk

Matías Toro
Computer Science Department,
University of Chile
Santiago, Chile
mtoro@dcc.uchile.cl

Claudio Gutierrez
Computer Science Department,
University of Chile & IMFD
Santiago, Chile
cgutierr@dcc.uchile.cl

Bruno C. d. S. Oliveira
The University of Hong Kong
Hong Kong, China
bruno@cs.hku.hk

Éric Tanter
Computer Science Department,
University of Chile & IMFD
Santiago, Chile
etanter@dcc.uchile.cl

## ABSTRACT

Practical SQL engines differ in subtle ways in their handling of typing constraints and implicit type casts. These issues, usually not considered in formal accounts of SQL, directly affect portability between engines. We present a general formal semantics framework for SQL, named RAFT, that explicitly captures the semantics of types, both static and dynamic. RAFT is expressed in terms of abstract operators that provide the necessary leeway to precisely model different SQL engines. We show that this formalism has good theoretical properties (e.g. practical conditions that ensure type safety and soundness) and is of practical use. Indeed, we provide RAFT instantiations for 5 major SQL engines: PostgreSQL, Microsoft SQL Server, MySQL, SQLite, and Oracle Database. To validate the adequacy of the formalism, we implement RAFT in Python, along with these five instances, and test them with thousands of randomly-generated queries. RAFT can serve as precise documentation of typing in existing engines and potentially guide their evolution, as well as provide a formal basis to study type-aware query optimizations, and design provably-correct query translators.

## 1 INTRODUCTION

Query translation between different SQL engines is a common practice arising in different scenarios, like database migration (to reduce costs, maintenance, changes in software, etc.) [7, 21, 23] and prototyping (code in lightweight databases like SQLlite and then port to a more robust database). Today there are many tools addressing this task [1, 2, 19].

Translation between SQL engines could bring many surprises. One of the most notorious is the semantic discrepancy between SQL engines related to typing behaviour, the problem we address in this paper. This problem is mainly due to differences in datatypes, type checking, when to perform the checks, and explicit and implicit type casts that may or may not be performed by the engines. This poses substantial challenges for developers and database administrators. Existing tools, whether paid or open source, often fall short in addressing these differences, tending to prioritize syntax over behavioral disparities.[1]

Table 1 presents a sample of minimal examples of this wide difference in behaviour. As these examples show, database engines have different treatment of types, that obey to different design models, like some being statically typed while others embrace dynamic typing, to different approaches to the overloading of basic arithmetic and comparison operations, etc. For now we only want to stress that there is a relevant problem here. In Section 2 we will give more detailed explanations for these cases.

Addressing and understanding these anomalies is not a simple task. On one hand, SQL standards do not cover many issues related to typing or leave the interpretation rather open.[2] On the other hand, many of the design decision of the engines are hidden under optimization mechanisms that either are not public or complex to find in the code. Last, but not least, the problem has not been addressed by the research literature. Indeed, although there is solid work on the formalization of the semantics of SQL [20, 31], they assume that all comparisons and operations apply to the right types. On the other hand, typing in SQL and syntax errors have been explored (particularly in the Programming Language community from where we took some inspiration [5, 14, 25, 27, 32], but the problem addressed of type constraints and casts which may induce errors at runtime, the one addressed in this paper, has not been

---

[1]Some examples: https://en.wikibooks.org/wiki/Converting_MySQL_to_PostgreSQL http://www.sqlines.com/online https://www.rebasedata.com/convert-bmysql-bto-bpostgres-bonline

[2]A good example is the following: To cast an exact number to an exact numeric type, e.g. from a real to int, it specifies: "If there is a representation of SV [source value] in the data type TD [target value] that does not lose any leading significant digits after rounding or truncating if necessary, then TV is that representation. The choice of whether to round or truncate is implementation-defined." ANSI SQL 1992, Sec. 6.10, Case 3)a)i) )

Table R

| A | B |
|---|---|
| 'Bob' | 10 |
| '1' | 20 |
| '1.1' | 30 |

| | Query | PSQL | MSSQL | Oracle | MySQL | SQLite |
|---|---|---|---|---|---|---|
| E1 | `SELECT 1.1 + 1 FROM R` | 2.1 | 2.1 | 2.1 | 2.1 | 2.1 |
| E2 | `SELECT '1' + 1 FROM R` | 2 | 2 | 2 | 2 | 2 |
| E3 | `SELECT '1.1' + 1 FROM R` | ✗ | ↛ | 2.1 | 2.1 | 2.1 |
| E4 | `SELECT '1.1' + 1.1 FROM R` | 2.1 | 2.1 | 2.1 | 2.1 | 2.1 |
| E5 | `SELECT '1' + '1' FROM R` | ✗ | '11' | 2 | 2 | 2 |
| E6 | `SELECT 'sql' + '2ra' FROM R` | ✗ | 'sql2ra' | ↛ | 2 | 2 |
| E7 | `SELECT 1+A FROM R WHERE B=20` | ✗ | 3 | 3 | 3 | 3 |
| E8 | `SELECT 1+A FROM R WHERE B=10` | ✗ | ↛ | ↛ | 1 | 1 |
| E9 | `SELECT 1 + A FROM (SELECT '2' AS A) B` | ✗ | 3 | 3 | 3 | 3 |
| E10 | `SELECT 1 FROM R WHERE '1' < 2` | 1 | 1 | 1 | 1 | ∅ |
| E11 | `SELECT 1 FROM R WHERE '1.1' < 2` | ✗ | ↛ | 1 | 1 | ∅ |
| E12 | `SELECT '1.1' FROM R INTERSECT SELECT 1.1 FROM R` | 1.1 | 1.1 | ✗ | 1.1 | ∅ |
| E13 | `SELECT '1.1' FROM R INTERSECT SELECT 1 FROM R` | ✗ | ↛ | ✗ | ∅ | ∅ |

**Table 1: Examples of discrepant behaviors of different database engines. The queries are artificially chosen to be minimal examples of typing issues. For simplicity we show only one element of the result. We use ✗ to denote a static error (before execution), and ↛ to denote a runtime error.**

dealt with. In Section 8 of related work we present with more detail these antecedents.

In this paper we address the problem by proposing a general formal framework (called RAFT) that models the intricate behavioral discrepancies across database engines. Our framework (1) formalizes the static semantics of queries via a type system, (2) makes implicit type casts explicit via an elaboration phase that introduces explicit type casts, and (3) formalizes the dynamic semantics of queries via a monadic evaluator. The type system, elaboration, and evaluator are expressed in terms of abstract operators. This framework provides, in theory and in practice, the necessary leeway to precisely model different SQL engines. To the best of our knowledge, no comprehensive modeling of these features exists.

From a practical point of view, we tested the model over five different and well-known SQL engines: PSQL, MSSQL, Oracle, MySQL and SQLite. It shows the flexibility and comprehensiveness of the framework, as it permits to understand different behaviors and thus opens the door for informed decisions when translating queries. The framework is tested with randomly generated queries. At the same time, the framework is a basis to formally discuss future typing-related design decisions for SQL, as it distills the core parameters involved.

From a theory point of view, we prove several formal results. First, a type safety result that proves that well typed queries either reduce to a table or raise a (controlled) type error. In other words, evaluation of well-typed queries does not get stuck. We enhance this result proving that the resulting table can be typed to the same type of the query. Moreover, PSQL, MySQL and SQLite satisfy a theorem that states that if the programmer do not uses explicit casts in well- typed queries, then the queries will evaluate without errors. Finally, regarding cast insertion, we can prove that independently of the engine (the proofs are parameterized by light constraints on abstract operators), translated queries preserve types, and the translation is unique.

In summary, this paper presents a formal and practical foundation for addressing the typing problems and challenges posed by database migration, by elucidating the database engine behavior in the presence of types, bridging the gap between various systems, and offering practical insights for both researchers and practitioners in the realm of database management.

The contributions and organization of the paper are as follows:

1. **Identification of Practical Issues:** In Section 2 we detect, localize and study the problems of semantic discrepancies due to typing, specifically type casts, which would affect query translation between SQL engines.

2. **Formal Semantics Framework:** We describe a formal framework, RAFT, based on a typed core subset of SQL with support for type casts. Section 3 overviews the overall design of RAFT and Section 4 provides the formal semantics of RAFT, which are expressed in terms of abstract operators that provide the necessary leeway to precisely model different SQL engines. We illustrate two different instantiations in Section 5, and provide three others in the supplementary material.

3. **Metatheory:** in Section 6 we present the metatheory of RAFT, demonstrating several important properties, such as the assurance that well-typed queries do not get stuck. For each property, we provide constraints on abstract operators required for these properties to hold. Furthermore, we illustrate which engine satisfy (or not) those properties.

4. **Empirical Validation:** In Section 7 we empirically validate our framework by developing multiple database interpreters. We test these interpreters with thousands of randomly-generated queries and compare their results with those obtained from real database engines. Additionally, we shed light on the impact and challenges of query optimizations on the evaluation process.

Section 8 discusses related work and Section 9 concludes. Additional material (proofs, full definitions, and prototype implementation) can be found in the supplementary material.

## 2 TYPING SEMANTIC DISCREPANCIES

In this work we study two kinds of different behavior between engines: (1) different results due to type conversions, and (2) different behavior related to type errors. Regarding the latter, we identify two kinds of type errors: *static errors*, which happen during compilation/before running the query; and *runtime errors* that occur during the execution of the query. We say that two engines differ in behavior if, for a given query, the results are different (including different kinds of errors).

To illustrate some of the differences, we explore different queries using a core subset of standard SQL. Like several other authors[9, 20, 28], this core comprises booleans, numbers, selections, cross-products, nested queries within `FROM` clauses, and set operations.

We also add support for arithmetic operators and type casts. Notably, features like nulls, aggregation, `EXISTS`, `IN`, and other forms of nested queries are not supported within this core. However, the aforementioned core serves to illustrate the primary distinctions, leaving the extension to these additional features as potential future work. For each query we illustrate their behavior in a mix of statically and dynamically typed engines such as PostgreSQL (PSQL), Microsoft SQL Server (MSSQL), MySQL, SQLite, and Oracle Database (Oracle).

The summary of examples and results can be found in Table 1. Relation R can be denoted using bags as R = {('Bob', 10), ('1', 20), ('1.1', 30)}.

## 2.1 Arithmetic operations

For the first set of examples we focus on the expression being selected rather than the tables or the conditions.

*E1 and E2.* To begin, we present two examples that behave uniformly across the engines. For simplicity, we say that the result is 2.1, to denote a bag of uniform elements {2.1, 2.1, 2.1}.

*E3.* This is the first example that illustrates a difference in behavior. PSQL and MSSQL throw a type error, whereas the other engines return 2.1. The reason for the error is that in PSQL and MSSQL we can implicitly cast a string to an integer if the string is an integer, but not if the string is a real number. To fix this problem in PSQL, we must explicitly *cast* the string to a real number: `SELECT CAST('1.1' as FLOAT) + 1 FROM R`.

*E4.* If we take example 3, and change 1 for a real number, such as 1.1, then the result is now 2.1 for every engine. The difference with respect to the previous example is that the plus operation is defined for both integers and real numbers, and now `'1.1'` can be cast directly to a real number.

*E5.* PSQL reports an static error due to the lack of an addition operator between two strings. MSSQL returns `'11'` as addition is overloaded for string. The other engines return 2 as addition is only defined between numbers, thus implicit casting `'1'` to 1.

*E6.* Both PSQL and Oracle report an error, MSSQL uses string concatenation yielding `'sql2ra'`, and MySQL and SQLite yield 2. The reason for the latter is due to the way these engines cast strings to numbers: they search for a number in the prefix of the string (if nothing is found then 0 is returned).

*E7.* PSQL rejects this query as column A is of type String, and the addition between numbers and strings is not defined. This behavior is more conservative than E2, as now it cannot determine statically if the given string can be cast to number or not. Other engines defer the check to runtime and return 2. To fix this query in PSQL we can explicitly cast column A to integer: `SELECT 1+CAST(A as INT) FROM R WHERE B=20`.

*E8.* PSQL (statically) rejects this query similarly to E7. MSSQL and Oracle now fail dynamically as they cannot convert `'Bob'` to a number. MySQL and SQLite on the other hand do not fail and return 1 as `'Bob'` is cast to 0. Casting A to `INT` in PSQL would make the error dynamic, and a runtime type error would be reported instead, similarly to MYSQL and SQLITE.

*E9.* Contrary to E2, PSQL raises a static error as the nested query hides the actual String returned by the subquery. All other engines yield 3 as conversions are optimistically performed at runtime. To

| Operation | SQLite | Other engines |
|---|---|---|
| `0 < 1` | 1 | t |
| `'0' < 1` | 0 | t |
| `'1' < 0` | 0 | f |
| `'0'+0 < 1` | 1 | t |
| `'0' < CAST(1 AS INT)` | 1 | t |
| `'0' < 1 + 0` | 0 | t |

**Table 2: Behavior of the comparison operator in** SQLite

fix this query in PSQL an explicit cast must be inserted `SELECT 1 + CAST(A AS INT) FROM (SELECT '2' AS A) B`.

## 2.2 Boolean operations

The following examples illustrate difference in behavior related to comparison operators on conditionals.

*E10.* Almost every engine is able to cast `'1'` to `INT`, returning 1. SQLite, on the other hand, returns a empty result as the condition is `false`. This is because SQLite does not perform implicit conversions at the boundaries of comparisons. SQLite has a type hierarchy where every string is bigger than any number.

*E11.* Now, if the left operand is a string representing a real, then PSQL and MSSQL return a type error. This is because the best type of the comparison operator is the one that takes two integers as argument (because 2 is an integer). As there is no direct implicit conversion between a string representing a real and an integer the query is rejected. SQLite still returns an empty result, and MySQL and Oracle returns 1.

*Comparison operator in* SQLite. SQLite warrants special attention to illustrate unique cases related to the comparison operator, as exemplified in Table 2. Notably, operations `'0' < 1` and `'1' < 1` yields 0. This behavior is attributed to the fact that strings are considered larger than integers, as previously explained. However, When an arithmetic operation or a cast is introduced the expressions now yields 1. One might then expect that `'0' < 1+0` would yield 1. Nevertheless, the result is actually 0, since the result of the arithmetic operation remains uncertain.

## 2.3 Set operations

*E12.* PSQL, MSSQL and MySQL, yield 1.1. This is because, during intersection, two conditions are checked: (1) the cardinality of both subqueries must match, and (2) the types of the columns must also be consistent. To achieve the second condition, an implicit conversion from string to real is inserted in the left subquery (the other direction is forbidden). However, Oracle encounters a type error since the column types do not match. On the other hand, SQLite returns an empty result as 1.1 is not the same as `'1.1'`.

*E13.* Now as `'1.1'` cannot be implicitly cast to an integer (the column type of the right subquery), this program is rejected by PSQL, MSSQL and Oracle. Both MySQL and SQLite return an empty result.

Having explored some discrepancies in semantics across different database engines due to types, we now transition to presenting an overview of our solution.
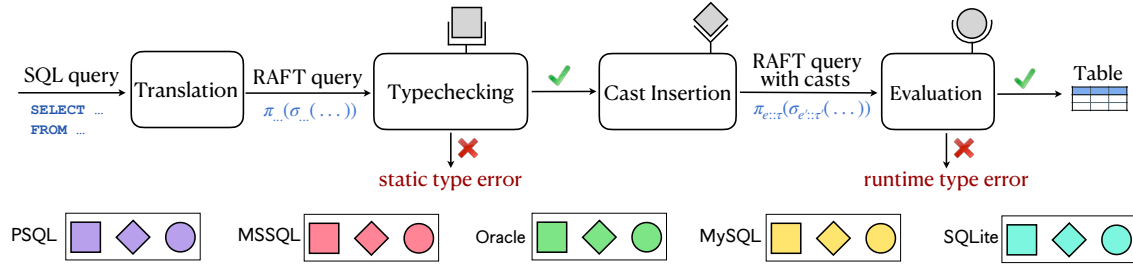
**Figure 1: Overview of RAFT**

## 3 OVERVIEW OF RAFT

In this paper, we introduce RAFT, a general formal semantics framework for SQL that explicitly captures the semantics of types, both static and dynamic. Figure 1 depicts a high-level overview of the framework. The processing of an SQL query involves four sequential phases: translation, typechecking, cast insertion, and evaluation.

*Translation.* In the translation step (Section 4.6), an SQL query is (1) analyzed to rule out syntactically invalid queries such as `SELECT FROM FROM` and (2) transformed to a RAFT query. RAFT is a core typed relational algebra (Section 4), with classical features such as projections, selections, and set operations, and unusual features such as arithmetic operations and casts.

*Typechecking.* RAFT does not assume the absence of type mismatches; users can write queries such as `SELECT R.A from R WHERE 'Bob' = 1`. For this reason, a typechecker validates queries before evaluation (Section 4.4). This typechecker is responsible for identifying mismatches between types and the number of columns of subqueries in set operations, and to ensure proper access to column names in scope. More importantly, it validates the appropriate usage of both implicit and explicit casts, rejecting operations that may result in casts known to fail during evaluation. If the typechecker rejects the query, the process terminates and reports a static type error to the user. However, if the query is deemed well-typed, it proceeds to the third phase.

*Cast Insertion.* The cast insertion phase (Section 4.5) transforms a RAFT query by making all implicit type casts explicit. The purpose of this phase is to circumvent the complexities that would arise from handling implicit casts in the evaluation phase, thereby simplifying the dynamic semantics of RAFT. This approach aligns the dynamic semantics closely with the classical relational algebra semantics, while only necessitating an extension to accommodate explicit casts and errors.

*Evaluation.* The evaluation phase (Section 4.3) interprets the transformed RAFT query. For this purpose, we present a monadic evaluator expressed using denotational semantics, cleanly encompassing the management of explicit type conversions and runtime errors. Should a type error arise during evaluation, the evaluation process is halted, and a runtime type error is reported. In the absence of runtime type errors, the outcome of evaluation is a table.

*Parametrization.* RAFT is meant to precisely capture the actual semantics of different SQL engines. To do so, typechecking, cast insertion, and evaluation are parametrized by engine-specific components (Figure 1). Specifically, these phases are expressed in terms of abstract operators that provide the necessary leeway to precisely

model different SQL engines (Section 4.2), and represented by the different shapes in the figure. These operators include: how to cast a given value to a specific type, how to resolve overloading (*i.e.* the use of an operator with different argument types), determining the types of values and operations, and when to insert explicit casts. We provide instantiations of these operators for 5 major SQL engines: PSQL, MSSQL, Oracle, MySQL and SQLite. Section 5 presents RAFT/PSQL and RAFT/SQLite in full. The instantiations of RAFT for the other three engines are provided in the supplementary material.

*Properties.* Providing a formal framework makes it possible to precisely formulate and prove properties satisfied by all or some engines (Section 6). Some properties can be proven generally at the level of RAFT, assuming some basic conditions of the abstract operators; consequently, any instantiation that respects these conditions is guaranteed to enjoy the said properties. Others can be shown of specific instantiations. In particular, we show that every well-typed RAFT query does not get stuck (such as in queries which select columns that do not exist), either returning a table or raising a runtime type error.

We prove that—with the exception of SQLite—the table resulting from the evaluation of a well-typed query indeed has the predicted type. Additionally, we establish that the absence of explicit casts in the original source SQL query guarantees that no type errors can occur at runtime. We show that this requirements are satisfied by PSQL, MySQL, and SQLite, but not by MSSQL and Oracle.

Finally, we prove that cast insertion is a type-preserving function. This means that cast insertion is deterministic, and that the resulting query is typed to the same type of the source query. We prove that it holds for all the engines considered in this work.

*Implementation.* To validate the adequacy of the formalism, we implemented RAFT in Python along with these five engine instances. We randomly-generated thousands of random queries, corroborating that queries behave the same in RAFT as in the corresponding engines. We validate that, disregarding optimizations, the behavior of queries remains consistent: either both the prototype and the real engine encounter errors, or both evaluate to the same table. Additionally, considering optimizations, we validate that when both the prototype and the engine successfully evaluate a query, the results match.

*Non-goals.* Finally, it is important to note that several *non-goals* we set for the development of RAFT include performance and query optimization, and are left as future work. Our primary focus centers on exploring the semantics and properties of the framework as we discuss next.

**Types and Schemas**

$\tau ::= \mathbb{R} \mid \mathbb{Z} \mid \mathbb{B} \mid \text{String} \mid ?$      (value types)

$T ::= N \mapsto \tau \mid T, N \mapsto \tau$      (relation types)

$\Gamma ::= \emptyset \mid \Gamma, R \mapsto T$      (schema)

**Values and Expressions**

$0_A \in \{+\}$      (arithmetic ops)

$0_C \in \{<, =\}$      (comparison ops)

$0_B \in \{\wedge, \vee\}$      (boolean ops)

$w ::= r \mid n \mid b \mid s$      (simple values)

$v ::= w \mid w :: ?$      (values)

$e ::= N \mid v \mid e\, 0_A\, e \mid e :: \tau$      (general expressions)

$\theta ::= e\, 0_C\, e \mid \theta\, 0_B\, \theta \mid \neg 0_B$      (boolean expressions)

$\beta ::= e$ as $N \mid \beta, e$ as $N$      (aliased expressions)

**Queries**

$0_S \in \{\cup, \cap, \times, -\}$      (set query ops)

$Q ::= R \mid \pi_\beta(Q) \mid \sigma_\theta(Q) \mid Q\, 0_S\, Q$      (queries)

**Figure 2: Syntax of** RAFT. $r, n, b, s$ **denote respectively a real number, an integer number, a Boolean value and a string.** $N$ **is a name.** $R$ **a relation.**

## 4 THE RAFT FORMAL FRAMEWORK

In this section we present the syntax, dynamic semantics, type system, cast insertion and translation of RAFT. A RAFT query is expressed using a core typed relational algebra, supporting projections, selection, set operations, arithmetic and boolean operations, and implicit and explicit casts.

### 4.1 Syntax

The syntax of RAFT is presented in Figure 2. The formalization is inspired by the work of Guagliardo and Libkin [20], except that here we deal with typing instead of assuming a prior (unstudied) typing phase. Types play a central role in this work, because as we illustrated, typing discrepancies are a source of behavioral differences between engines. This section first presents types and schemas, then values and expressions, and finally queries.

*Types and Schemas.* There are two category of types: *value types* $\tau$ that are used to type expressions (and values), and *relation types* $T$ to type queries and tables (relations). A value type $\tau$ is either a real $\mathbb{R}$, an integer $\mathbb{Z}$, a boolean $\mathbb{B}$ or a string String. We use the symbol ? for the unknown type, used by PSQL to type string literals [18], and to model flexible typing in SQLite. A relation type $T$ is an ordered list of pairs of the form $N : \tau$, where $N$ is a column name and $\tau$ its corresponding value type. For instance, $T_1 = \text{Name} \mapsto \text{String}, \text{Age} \mapsto \mathbb{Z}$ and $T_2 = \text{Name} \mapsto \text{String}, \text{Price} \mapsto \mathbb{R}$ are relation types. A schema $\Gamma$ is a list of pairs $R, T$ where $R$ is a relation name and $T$ a relation type. Intuitively, schemas represent types of databases. For instance, the schema of a database of persons and products can be written as $\text{Person} \mapsto T_1, \text{Product} \mapsto T_2$.

*Values and expressions.* We represent tables as sets of *rows*, where each row is a list of *values*[3]. A *simple value* $w$, which represents an atomic piece of data, is either a real number $r$, a natural number $n$, a boolean $b$ or a string $s$. Values $v$ are either a simple value $w$, or a simple value cast to the unknown type $w :: ?$ (the latter is not

---

[3]For simplicity we represent tables as sets, but the supplementary material presents rules using multisets.

used directly by programmers, and it is used exclusively in engines such as SQLite). Expressions, which can be reducible, are classified into three categories: *general expressions*, *boolean expressions*, and *aliased expressions*. A general expression $e$, is used in projections and selections, and is either a column name $N$ such as Name or Age, a value $v$, an arithmetic operation $e\, 0_A\, e$ (for simplicity we only use $+$), or an explicit cast $e :: \tau$. Boolean expressions $\theta$ are used exclusively as conditions in selections, and are either comparison operations $e\, 0_C\, e$, or logical combinations of them. Note that we do not include boolean expressions as general expressions, because some engines do not support selecting boolean values in queries (e.g. in MSSQL and Oracle, `SELECT` 1 < 2 `FROM` R is a syntactically invalid query while it is valid in PSQL, MySQL and SQLite). An aliased expression $\beta$ is a (non-empty) list of pairs $e$ as $N$, where $e$ is an expression and $N$ its alias. Note that this is a richer version than the classical relational algebra aliasing operation, which allows binding of names to expressions, instead of only to queries. For example, in `SELECT` A.C `FROM` (`SELECT` 1+1 `AS` C `FROM` R) A), the expression 1+1 gets a name C that can be used in the outer query.

*Queries.* A query is either a relation $R$, a projection $\pi_\beta(Q)$, a selection $\sigma_\theta(Q)$, or a set operation. Set operations include cross product $(Q \times Q)$, intersection $(Q \cap Q)$, union $(Q \cup Q)$ and difference $(Q - Q)$. Note that projection is non-standard as it is parametrized by a list of aliased expressions, instead of a list of names. This allows RAFT to model SQL queries like `SELECT` 1+1 `AS` C `FROM` R, which projects the arithmetic expression $1 + 1$ from every row of $R$, as $\pi_{(1+1)}$ as $C(R)$.

### 4.2 Abstract operators

As mentioned in Section 3, certain key operators in RAFT are left abstract, as they depend on the specific engine being used.

*Bidirectional Implicit Cast.* Operator $biconv(e_1, \tau_1, e_2, \tau_2) = \tau_3$ determines the optimal implicit type cast for a set operator applied to two columns of (possibly) different types. The *biconv* operator takes as argument two expressions $(e_1, e_2)$ and their corresponding types $(\tau_1, \tau_2)$, returning a type. It either returns $\tau_2$ or $\tau_1$ by testing if $e_1$ can be implicitly cast to $\tau_2$, or if $e_2$ can be implicitly cast to $\tau_1$ respectively.

*Explicit Cast.* Operator $cast(v, \tau) = v'$ attempts to cast value $v$ to a value $v'$ of type $\tau$. This function is primarily used to evaluate casts at runtime.

*Overloading Resolution.* Operator $resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4$ determines which specific operation among a set of overloaded ones should be called based on the provided arguments during invocation. More specifically, given a set of function types $\overline{\tau_3} = \overline{\tau_5 \times \tau_6 \to \tau_7}$, it tries to find the best candidate type for expressions $e_1$ and $e_2$, typed as $\tau_1$ and $\tau_2$ respectively. A function type of the form $\tau_5 \times \tau_6 \to \tau_7$, denotes a function that accepts two arguments, of types $\tau_5$ and $\tau_6$ respectively, and produces a result of type $\tau_7$.

*Explicit Cast Feasibility.* This operator takes one expression and two types, and returns a boolean. For simplicity, it is presented as a relation $e : \tau' \approxtail \tau$, and rules out explicit casts that are known to fail at runtime. It test if it is possible to explicitly cast expression $e$ of type $\tau'$, to an expression of type $\tau$.

*Type of Values.* Operator $ty(v) = \tau$ computes the type of values. *Candidate Types.* Operation $ty(0) = \overline{\tau}$ determines the set of possible

$$\llbracket R \rrbracket_{\mathsf{D}} = \mathbf{ok}\ \mathsf{D}\ (R) \tag{R$R$}$$

$$\llbracket Q_1\ \mathsf{O_S}\ Q_2 \rrbracket_{\mathsf{D}} = \mathrm{do}\{t_1 \leftarrow \llbracket Q_1 \rrbracket_{\mathsf{D}}; t_2 \leftarrow \llbracket Q_2 \rrbracket_{\mathsf{D}}; \mathbf{ok}\ t_1\ \mathsf{O_S}\ t_2\} \tag{R$\mathsf{O_S}$}$$

$$\llbracket \pi_\beta(Q) \rrbracket_{\mathsf{D}} = \mathrm{do}\{t \leftarrow \llbracket Q \rrbracket_{\mathsf{D}}; \lceil\{\llbracket\beta\rrbracket^r_{\eta^r_{\ell(Q)}}\ |r \in t\}\rceil\} \tag{R$\pi$}$$

$$\llbracket \sigma_\theta(Q) \rrbracket_{\mathsf{D}} = \mathrm{do}\{t \leftarrow \llbracket Q \rrbracket_{\mathsf{D}}; \_ \leftarrow \lceil\{\llbracket\theta\rrbracket^r_{\eta^r_{\ell(Q)}}\ |r \in t\}\rceil;$$
$$\mathbf{ok}\ \{r|r \in t \wedge \llbracket\theta\rrbracket^r_{\eta^r_{\ell(Q)}}\}\} \tag{R$\sigma$}$$

$$\llbracket e \text{ as } N \rrbracket_\eta = \llbracket e \rrbracket_\eta \tag{R as}$$

$$\llbracket \beta, e \text{ as } N \rrbracket_\eta = \mathrm{do}\{r \leftarrow \llbracket\beta\rrbracket_\eta; v \leftarrow \llbracket e \rrbracket_\eta; \mathbf{ok}\ r, v\} \tag{R$\beta$}$$

$$\llbracket N \rrbracket_\eta = \mathbf{ok}\ \eta(N) \tag{R$N$}$$

$$\llbracket v \rrbracket_\eta = \mathbf{ok}\ v \tag{R$v$}$$

$$\llbracket v :: \tau \rrbracket_\eta = \begin{cases} \mathbf{ok}\ v' & \text{if}\ \boxed{cast(v, \tau) = v'} \\ \mathbf{error} & \text{otherwise} \end{cases} \tag{R$v$ ::}$$

$$\llbracket e :: \tau \rrbracket_\eta = \mathrm{do}\{v \leftarrow \llbracket e \rrbracket_\eta \quad e \neq v; \llbracket v :: \tau \rrbracket_\eta\} \tag{R$e$ ::}$$

$$\llbracket \theta_1\ \mathsf{O_B}\ \theta_2 \rrbracket_\eta = \mathrm{do}\{b_1 \leftarrow \llbracket\theta_1\rrbracket_\eta; b_2 \leftarrow \llbracket\theta_2\rrbracket_\eta; \mathbf{ok}\ b_1\ \mathsf{O_B}\ b_2\} \tag{R$\mathsf{O_B}$}$$

$$\llbracket \neg\theta \rrbracket_\eta = \mathrm{do}\{v_1 \leftarrow \llbracket\theta_1\rrbracket_\eta; \mathbf{ok}\ \neg b_1\} \tag{R$\neg$}$$

$$\llbracket e_1\ \mathsf{O}\ e_2 \rrbracket_\eta = \mathrm{do}\{v_1 \leftarrow \llbracket e_1 \rrbracket_\eta; v_2 \leftarrow \llbracket e_2 \rrbracket_\eta; \tag{R$\mathsf{O}$}$$
$$\mathbf{ok}\ \boxed{apply(\mathsf{O}, v_1, v_2)} \quad \mathsf{O} \in \mathsf{O_A} \cup \mathsf{O_C}\}$$

**Figure 3: Dynamic Semantics of RAFT.
(abstract operators highlighted in gray)**

function types available for operation $\mathsf{O}$. The output is a list of function types $\bar{t}$ given that many arithmetic and comparison operators are overloaded with more than one type.

*Unknown Type Removing.* Operator $clean(T) = T'$ performs postprocessing for relation type $T$ returning a new relation type $T'$. In some engines, this operation replaces occurrences of the unknown type with another type.

*Annotation Insertion.* Operator $insert(e, \tau, \mathsf{O}) = e'$ attempts to explicitly cast expression $e$ to type $\tau$ in a context within the $\mathsf{O}$ operation. It returns a new expression $e'$ depending on whether the expression was cast or not.

*Value Operation Application.* Operator $apply(\mathsf{O}, v_1, v_2) = v_3$ performs arithmetic or comparison operation $\mathsf{O}$ to values $v_1$ and $v_2$, yielding value $v_3$.

In the next subsections, we use these operators to define the dynamic semantics, the type system and the cast insertion procedure. Examples of instantiation of these abstract operators can be found in Section 5.

### 4.3 Dynamic Semantics

We now present the dynamic semantics of RAFT, described in Figure 3. As standard in programming languages, the *dynamic* semantics specify the way programs are evaluated (a.k.a the runtime), while the *static* semantics refer to typing. Here, the dynamic semantics define the evaluation of expressions, operators and queries. Typing for RAFT is presented in Section 4.4.

We follow standard notations in database and programming languages [3, 20, 30]. The dynamic semantics of RAFT differ from the ones of Guagliardo and Libkin [20] (GL) as follows. First, RAFT

dynamic semantics are defined over a relational algebra, whereas GL uses SQL syntax. Second, RAFT supports implicit and explicit casts, and arithmetic operations, but not null values, and `DISTINCT`, `IN`, and `EXISTS` operators supported by GL. Finally, and more importantly, as some casts may be invalid, the dynamic semantics must deal with the possibility of runtime errors. For instance, in PSQL the query `SELECT CAST('s' as INT) FROM` R $(\sigma_{\text{'s'}::\mathbb{Z}}(R))$ evaluates to an error.

To concisely account for the possibility of runtime errors, we present the dynamic semantics in *monadic* style in order to streamline the handling of errors [37]. The monadic presentation of computations that may fail consists in so-called "optional" values: either an actual value tagged with **ok**, or **error** to denote an error. The sequential composition is given in a do block (e.g. $\mathrm{do}\{A; B; C\}$). Errors are transparently propagated through such sequences: the evaluation returns **ok** $v$ if all steps in a do block (e.g. $A$, $B$, and $C$) evaluate successfully, or **error** if one step in the sequence evaluates to **error**.

There are four categories of evaluations : $\llbracket Q \rrbracket_{\mathsf{D}}$ to reduce queries, $\llbracket e \rrbracket_\eta$ to reduce expressions, $\llbracket\beta\rrbracket_\eta$ to reduce aliased expressions, and $\llbracket\theta\rrbracket_\eta$ to reduce boolean expressions. The evaluation of a query is parametrized by a database $\mathsf{D}$, which maps relation names $R$ to tables $t$. A table $t$ is a set of *rows* $r$, where a row $r$ is an ordered list of values $v$.[4] The other categories of evaluation are parametrized by an environment $\eta$, which maps column names $N$ to values $v$. Intuitively, this environment is used to extract the value associated with a column name for a given row. For instance, consider a relation of persons P(Name $\mapsto$ String, Age $\mapsto$ $\mathbb{Z}$), and its table $\{(\text{'Bob'}, 10), (\text{'Alice'}, 20)\}$. Then, for the first row, $\eta(\text{Name}) = \text{'Bob'}$ and $\eta(\text{Age}) = 10$.

*Queries.* The basic case is Rule (R$R$), which evaluates the name of a relation yielding the table associated to that name in database $\mathsf{D}$. Rule (R$\mathsf{O_S}$) evaluates set operations by first reducing the subqueries and then combining the resulting tables. Rule (R$\pi$) starts by evaluating subquery $Q$. If the result is successful (a table $t$), then for each row $r$ in $t$, we try to project columns as dictated by $\beta$. To do this, we evaluate $\beta$ under environment $\eta^r_{\ell(Q)}$ (similarly to [20]). This environment is formed by matching corresponding column names of $Q$ with the values of $r$. The list of column names of a query $Q$ is obtained using function $\ell(Q)$, and defined inductively in the supplementary material

Finally, as the evaluation of aliased expressions can also produce errors, we lift the set of optional rows S to an optional list of rows using the $\lceil \cdot \rceil$ function defined as:

$$\lceil S \rceil = \begin{cases} \mathbf{error} & \mathbf{error} \in S \\ \mathbf{ok}\ \{r|(\mathbf{ok}\ r) \in S\} & \text{otherwise} \end{cases}$$

Rule (R$\sigma$) follows a similar approach by first reducing the subquery $Q$. If the result is successful (yielding table $t$), we proceed to test the reduction of the condition $\theta$ for each row. We employ a strategy akin to that of (R$\pi$), but in this case, the resulting set of Booleans is unused (binding the resulting in variable "_"). This way, in the third instruction, we can be confident that the evaluation of the conditions does not result in errors. Finally, we filter the rows from table $t$ that satisfy the given condition.

---

[4] The supplementary material uses multi-sets instead.

$$\boxed{T \vdash e : \tau}$$

$$(Tv) \frac{}{T \vdash v :\ \boxed{ty(v)}}$$

$$(TN) \frac{(N \mapsto \tau) \in T}{T \vdash N : \tau} \qquad (T::) \frac{T \vdash e : \tau' \qquad \boxed{e : \tau' \approx\!\!\!\approx \tau}}{T \vdash (e :: \tau) : \tau}$$

$$(T0) \frac{T \vdash e_1 : \tau_1 \qquad T \vdash e_2 : \tau_2 \qquad 0 \in 0_A \cup 0_C}{\boxed{resolve(e_1, \tau_1, e_2, \tau_2, ty(0)) = \tau_3 \times \tau_4 \rightarrow \tau_5}}{T \vdash e_1\ 0\ e_2 : \tau_5}$$

$$(T0_B) \frac{T \vdash \theta_1 : \mathbb{B} \qquad T \vdash \theta_2 : \mathbb{B}}{T \vdash \theta_1\ 0_B\ \theta_2 : \mathbb{B}} \qquad (T\neg) \frac{T \vdash \theta : \mathbb{B}}{T \vdash \neg\theta : \mathbb{B}}$$

$$\boxed{T \vdash \beta : T'}$$

$$(T\beta) \frac{\overline{T \vdash e : \tau} \qquad unique(\overline{N})}{T \vdash \overline{e\ \mathsf{as}\ N} : \overline{N \mapsto \tau}}$$

$$\boxed{\Gamma \vdash Q : T}$$

$$(TR) \frac{\Gamma(R) = T}{\Gamma \vdash R : T} \qquad (T\pi) \frac{\Gamma \vdash Q : T \qquad \boxed{clean(T)} \vdash \beta : T'}{\Gamma \vdash \pi_\beta(Q) : T'}$$

$$(T\sigma) \frac{\begin{array}{c} \Gamma \vdash Q : T \\ T \vdash \theta : \mathbb{B} \end{array}}{\Gamma \vdash \sigma_\theta(Q) : T} \qquad (T\times) \frac{\begin{array}{c} \ell(Q_1) \cap \ell(Q_2) = \emptyset \\ \Gamma \vdash Q_1 : T_1 \qquad \Gamma \vdash Q_2 : T_2 \end{array}}{\Gamma \vdash Q_1 \times Q_2 : T_1, T_2}$$

$$(T0_S) \frac{\begin{array}{c} \Gamma \vdash \pi_{\beta_1}(Q_1) : T_1 \qquad \Gamma \vdash \pi_{\beta_2}(Q_2) : T_2 \\ \boxed{biconv^*(\beta_1, T_1, \beta_2, T_2) = T} \qquad 0_S \in \{\cup, \cap, -\} \end{array}}{\Gamma \vdash \pi_{\beta_1}(Q_1)\ 0_S\ \pi_{\beta_2}(Q_2) : T}$$

**Figure 4: Type System of** RAFT.
**(abstract operators highlighted in gray)**

*Expressions.* Rule (Ras) evaluates a single aliased expression by evaluating the subexpression $e$ and disregarding the name $N$. Rule (R$\beta$) applies when we are evaluating multiple aliased expressions. Initially, it recursively reduces the sublist to a row $r$, and then the head of the list to a value $v$, resulting in a new row $r, v$.

Rule (R$N$) successfully evaluates a column name $N$ to its corresponding value in $\eta$. Rule (R$v$) successfully evaluates values to themselves. Rule (R$v$ ::) attempts to cast a value into a value of a different type using the *cast* function. For instance, in SQLite, the expression `CAST('hi' as INT)` evaluates to 1, whereas in PSQL is not defined. If the function is defined for the given value and type, the resulting value is returned; otherwise, an error is raised. Rule (R$e$::) applies to subexpressions that are not already values. It first reduces the subexpression to a value, and then casts the value using rule (R$v$::).

Rules (R$0_B$), (R$\neg$), and (R$0$) operate in a similar fashion. First, each subexpression is reduced, then the resulting values are combined using the specific operation at hand. For arithmetic and comparison operations, the exact operation is performed using the *apply* operation. For instance, in PSQL, the expression `'0'` $< 1$ evaluates to true, whereas in SQLite, it yields false.

## 4.4 Type System

A type system automatically checks a query or an expression for type errors before its evaluation, ensuring that such errors never occur at runtime. The type system of RAFT is presented in Figure 4, and is composed of three categories of typing rules, for general and aliased expressions, and for queries. Each typing rule attempts to assign a type to a query or expression and is written as an inference rule of the form $\frac{A \ldots}{B}$, where $A \ldots$ are the premises and $B$ is the conclusion. This notation is equivalent to the implication $A \ldots \Rightarrow B$.

*Expressions.* Judgment $T \vdash e : \tau$ denotes that expression $e$ has type $\tau$ under relation type $T$.

Rule (T$N$) assigns the type $\tau$ to the column name $N$ if $N$ is mapped to $\tau$ in the relation type $T$.

Rule (T::) assigns the type $\tau$ to an ascription $\tau :: e$ if the following conditions are met: (1) the expression $e$ must be well-typed for some $\tau'$; (2) the explicit cast of $e$ to $\tau$ is checked to rule out explicit casts that are known to fail at runtime, using the $\cdot : \cdot \approx\!\!\!\approx \cdot$ relation. For example, the attempt to cast `'hi'` to $\mathbb{Z}$ (`'hi'` $:: \mathbb{Z}$) is rejected in PSQL, while casting `'1'` to $\mathbb{Z}$ is accepted in both PSQL and SQLite.

Rule (T$v$) assigns types to values based on the specific database engine. We achieve this by employing the $ty$ operator. For instance, in PSQL, integers are typed as $\mathbb{Z}$ and literal strings as ?. However, in SQLite, both are typed as ?.

Operations are typed using rules (T0) and (T0$_B$) and (T$\neg$). Rule (T0) is responsible for both arithmetic and comparison operations. It typechecks the sub-expressions $e_1$ and $e_2$, which results in obtaining their respective types, $\tau_1$ and $\tau_2$. Next, the rule searches for the best candidate type signature for the given operation, taking into consideration that an operation might be overloaded with multiple types. This selection process is captured by the *resolve* and $ty()$ operations. If a single candidate function type is identified, the expression is typed using the codomain. If more than one candidate type is found, the expression is considered ill-typed. For instance, in the case of the query `SELECT` 1+1 `FROM` P both PSQL and SQLite choose numeric addition ($\mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z}$). However, when dealing with strings (e.g. `SELECT` `'a'` + `'b'` `FROM` P) PSQL rejects the query because it cannot not choose a best candidate operator. On the other hand, SQLite instantiates the query with numeric addition, implicitly casting both string arguments to integers. Finally, rules (T0$_B$) and (T$\neg$) type boolean operations.

Rule (T$\beta$) is used to type an aliased expression $\beta$. For simplicity, we adopt the notation $\overline{A}$ to denote a list of elements $A$. For instance, $\overline{e\ \mathsf{as}\ N}$ denotes the list $e_1$ as $N_1, ..., e_2$ as $N_2$, and $\overline{e}$ and $\overline{N}$ denote $e_1, ..., e_n$ and $N_1, ..., N_n$ respectively. Every expression $e_i$ is typed under $T$, yielding a type $\tau_i$. Subsequently, the type of the aliased expression as a whole is a relation type, where each name $N_i$ is mapped to the type $\tau_i$ of their corresponding subexpression $e_i$, represented as $\overline{N \mapsto \tau}$.

*Queries.* Judgment $\Gamma \vdash Q : T$ represents that query $Q$ has relation type $T$ under schema $\Gamma$. Rule (T$R$) assigns a relation type to relation name $R$ according to schema environment $\Gamma$.

Rule (T$\pi$) first recursively typechecks the subquery $Q$ yielding type $T$. Then the projection is assigned the same type of the aliased expression $\beta$, which is typed under relation type $clean(T)$. Recall, that the *clean* operator performs a transformation on a relation type,

removing unknown occurrences. For instance, in PSQL, `SELECT '1' + 1 FROM P` runs successfully, but `SELECT C + 1 FROM (SELECT '1' AS C) B` is rejected statically. The reason for this is that in the first query `'1'` has type unknown ?, which can be cast to integer, whereas in the second query, the unknown type ? is transformed to String disallowing the cast to integer.

Rule (T$\sigma$) first typechecks the subquery, resulting in a relation type $T$. Then, the condition must be successfully typed as Boolean under the context of $T$ (considering that the condition may reference columns from the subquery). Finally, the selection operation is assigned the same type as the subquery $T$.

Typechecking query union, intersection, and difference (just set operations for simplicity) requires special attention. Usually, it is assumed that the names, cardinality, and types of the columns in the subqueries match. In practice, the column names do not need to match, but cardinality and types must align. To achieve this, many engines perform implicit casts between the columns of the subqueries to align their types. In particular, string literals are analyzed to check the plausibility of casts. For instance in PSQL, (`SELECT '1.1' AS A FROM P`) `INTERSECT` (`SELECT 1.1 AS A FROM P`) runs successfully, resulting in a non-empty result, whereas (`SELECT CAST(Age AS TEXT) AS A FROM P`) `INTERSECT` (`SELECT Age FROM P`) is rejected before execution. To deal with these special cases, and without loss of generality, we require that both subqueries within a set operation must be projections in order to verify column casts. For instance, (`SELECT '1.1' AS A FROM P`) `UNION` (`SELECT 1.1 AS A FROM P`) is encoded as $\pi_{'1.1' \text{ as } A}(P) \cup \pi_{1.1 \text{ as } A}(P)$.

Rule (T0$_S$) deals with set operations over projections of lists of aliased expressions. First, it typechecks both projections. Second, it examines whether the lists of aliased expressions can be implicitly cast between each other, taking their relation types into account. Finally, if this cast is feasible, the target relation type is captured and used to typecheck the whole set operation. To check casts between list of aliased expressions, we use the $biconv^*$ operation defined as follows:

$$biconv^*(\overline{e \text{ as } N, N \mapsto \tau, e' \text{ as } N', N' \mapsto \tau'}) = \overline{N \mapsto biconv(e, \tau, e', \tau')}$$

This operation returns a relation type, where each name $N$ is mapped to the application of $biconv$ over each pair of expressions and their types.[5] This partial function checks casts and, if possible, returns one of the two types as a result. The expressions are provided to the function to rule out casts that are known to fail during evaluation. For instance, PSQL accepts query (`SELECT '1' AS A, 1 AS B FROM R`) `UNION` (`SELECT 2 AS A, '2' AS B FROM R`), as both `'1'` and `'2'` can be cast to integers; but rejects (`SELECT '1.1' AS A FROM R`) `UNION` (`SELECT 1 AS A FROM R`) as `'1.1'` cannot be implicitly cast to an integer (note that implicit casts from integers to strings are not allowed).

Finally, rule (T$\times$) types cross products using the concatenation of the relation types of both subqueries. This is under the condition that the sets of names in each subquery are disjoint.

## 4.5 Cast Insertion

Recall from Figure 1, that to avoid the complexity of dealing with implicit casts during runtime, before execution, we transform each

---

[5]We select $N$, the name of the left subquery, as it is a more commonly adopted practice in various database engines.

---

$$\boxed{T \vdash e : \tau \rightsquigarrow e'} \qquad (\text{E}v)\dfrac{T \vdash v : \tau}{T \vdash v : \tau \rightsquigarrow v :: \tau}$$

$$(\text{E0})\dfrac{\begin{array}{cc} T \vdash e_1 : \tau_1 \rightsquigarrow e'_1 & T \vdash e_2 : \tau_2 \rightsquigarrow e'_2 \\ \text{0} \in \{<, =, +\} & resolve(e_1, \tau_1, e_2, \tau_2, ty(\text{0})) = \tau_3 \times \tau_4 \rightarrow \tau_5 \end{array}}{\begin{array}{c} T \vdash e_1 \text{ 0 } e_2 : \tau_4 \rightsquigarrow \\ insert(\boxed{insert(e'_1, \tau_3, \text{0}) \text{ 0 } insert(e'_2, \tau_4, \text{0}), \tau_5, \text{0}}) \end{array}}$$

$$\boxed{\Gamma \vdash Q : T \rightsquigarrow Q'}$$

$$(\text{T0}_S)\dfrac{\begin{array}{cc} \Gamma \vdash \pi_{\beta_1}(Q_1) : T_1 \rightsquigarrow \pi_{\beta'_1}(Q'_1) & \Gamma \vdash \pi_{\beta_2}(Q_2) : T_2 \rightsquigarrow \pi_{\beta'_2}(Q'_2) \\ biconv^*(\beta_1, T_1, \beta_2, T_2) = T & \text{0}_S \in \{\cup, \cap, -\} \\ e_1 = \boxed{insert^*(\beta'_1, T, \text{0}_S)} & e_2 = \boxed{insert^*(\beta'_2, T, \text{0}_S)} \end{array}}{\Gamma \vdash \pi_{\beta_1}(Q_1) \text{ 0}_S \pi_{\beta_2}(Q_2) : T \rightsquigarrow \pi_{e_1}(Q'_1) \text{ 0}_S \pi_{e_2}(Q'_2)}$$

**Figure 5:** RAFT **Cast Insertion (excerpt).**
**(abstract operators highlighted in gray)**

---

implicit cast to an explicit cast. For instance, in our model, a PSQL query `SELECT '1' + 1 AS A FROM R` is transformed to `SELECT CAST('1' AS INT) + 1 AS A FROM R`, which in RAFT corresponds to the elaboration from $\pi_{'1'+1 \text{ as } A}(R)$ to $\pi_{('1'::\mathbb{Z})+1 \text{ as } A}(R)$.

An excerpt of the elaboration rules are presented in Figure 5; the complete rules can be found in the supplementary material. They are *type directed*, meaning that (1) we only elaborate well-typed terms, and (2) we use type information during elaboration. Like for typing, elaboration rules are defined inductively and grouped in three categories: for general and aliased expressions, and for queries. Most elaboration rules directly follow their corresponding typing rule. Judgment $T \vdash e : \tau \rightsquigarrow e'$ represents that expression $e$ typed as $\tau$ under relation type $T$ is elaborated to $e'$. Judgment $\Gamma \vdash Q : T \rightsquigarrow Q'$ is defined analogously.

Rule (E$v$) inserts an explicit cast to the type of each value. This is especially relevant for engines like SQLite where some values need to be tagged as unknown. For example, the expression `'0' < 1` evaluates to 0 (which represents false in that engine). To facilitate this special case, the expression is elaborated to `'0'` :: ? < 1 :: ?.

Rule (E0) introduces explicit casts based on the operation and the best candidate type. Specifically, we insert casts for both operands to match their expected corresponding domain types and also insert a cast in the result of the operation to the expected codomain type. Certain engines, like SQLite, do not insert explicit casts for comparison operations. To achieve this, we use the *insert* operator.

Rule (E$\pi\beta$) elaborates set operations by inserting casts to the output type of $biconv^*$. This is done using the *insert*$^*$ operation defined as follows:

$$insert^*(\overline{e \text{ as } N, N \mapsto \tau}, \text{0}_S) = \overline{insert(e, \tau, \text{0}_S) \text{ as } N}$$

For instance, for PSQL, query (`SELECT '1' FROM R`) `INTERSECT` (`SELECT 1 FROM R`) is elaborated to (`SELECT CAST('1' AS INT) FROM R`) `INTERSECT` (`SELECT CAST(1, INT) FROM R`).

## 4.6 Translation from SQL to RAFT

Figure 6 presents a subset of the translation rules from SQL to RAFT. These rules are defined inductively and divided into three

$$\boxed{B \leadsto_b \beta}$$

$$\text{(Tr}E)\ \frac{E \ \text{AS}\ alias(E) \leadsto_b e}{E \leadsto_b e}$$

$$\boxed{F \leadsto Q}$$

$$\text{(Tr}Fas)\ \frac{S \leadsto Q \qquad \ell(Q) = N_1, \ldots}{S\ N \leadsto \pi_{N_1 \text{as} N.N_1, \ldots}(Q)}$$

**Figure 6: SQL to** RAFT **translation (excerpt)**

categories: translation rules for expressions, aliased columns, and queries. Most of the rules are straightforward. Rule (Tr$E$) transforms expressions without an alias, as seen in **SELECT** 1+1 **FROM** R. In these cases, we utilize the $alias(\cdot)$ function, which generates an appropriate name for the subexpression. For example, $alias(1 + 1) = $ "1+1", and $alias(\text{NAME}) = \text{NAME}$. Rule (Tr$Fas$) transforms aliased subqueries by creating a projection of the translated subquery, where each column is aliased using the table name plus the column name. The complete rules can be found in the supplementary material.

In addition, we prove that every select query $S$ can be translated to a unique query in RAFT.

THEOREM 4.1. *For any SQL query $S$ there exists a unique relational algebra query $Q$ such that $S \leadsto Q$.*

## 5 INSTANTIATING RAFT

We now demonstrate how to instantiate RAFT for two radically different database engines: PSQL and SQLite. The instantiations for other engines can be found in the supplementary material. An instantiation is achieved by providing specific definitions of the abstract operators. We obtained these definitions by exploring the documentation of each engine, and by conducting black-box analyses interacting with each engine whenever the documentation was lacking in details.

### 5.1 RAFT/PSQL

Figure 7 describes the RAFT/PSQL instantiation. PSQL is characterized by being a strongly-typed SQL engine, meaning that it is more conservative than the rest. In many cases, if it cannot check the feasibility of casts, it rejects the query before its execution.

*Bidirectional implicit cast.* Operator $biconv$ is defined using the auxiliary *implicit cast* relation $e : \tau \leadsto \tau'$. An expression $e$ of type $\tau$ can be cast to $\tau'$ if either type $\tau$ can be implicitly cast to $\tau'$ ($\tau \Rightarrow \tau'$), or if $e$ is a value $v$ of type unknown ? (i.e. a string) and that value can be implicitly cast to some value $v'$ under type $\tau'$ ($icast(v, \tau') = v'$).

Implicit type cast $\tau \Rightarrow \tau'$ is only defined between identical types $\tau \Rightarrow \tau$, or from integers to reals $\mathbb{Z} \Rightarrow \mathbb{R}$. For instance, (**SELECT** 1 **FROM** R) **INTERSECT** (**SELECT** 1.0 **FROM** R) is accepted and evaluates to $\{1.0\}$, since $\mathbb{Z}$ (the type of 1) can be implicitly cast to a $\mathbb{R}$ (the type of 1.0). Implicit value cast $icast(v, \tau) = v'$ is defined for extracting numbers from strings, but not viceversa. Conversions from integer to string are only allowed in explicit casts. For instance, (**SELECT** '-1.0' **FROM** R) **INTERSECT** (**SELECT** -1.0 **FROM** R) is accepted and evaluates to a set containing $-1.0$.

*Explicit cast feasibility.* If the expression tested is not a value, the relation vacuously holds. Conversely, if $e$ is a value $v$, it is eagerly checked whether the explicit cast to $\tau'$ fails or not, using the explicit cast relation.

*Explicit cast.* In addition to what an implicit cast can do, explicit casts support casts from real numbers to integers by removing

$$\boxed{biconv(e_1, \tau_1, e_2, \tau_2) = \tau}$$

$$\frac{e_1 : \tau_1 \leadsto \tau_2}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau_2} \qquad \frac{e_2 : \tau_2 \leadsto \tau_1}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau_1}$$

$$\boxed{e : \tau \leadsto \tau'} \qquad \boxed{icast(v, \tau) = v'}$$

$$\frac{\tau \Rightarrow \tau'}{e : \tau \leadsto \tau'} \qquad icast(\text{'}n\text{'}, \mathbb{R}) = n$$

$$cast(n, \mathbb{R}) = n$$

$$\frac{icast(v, \tau) = v'}{v : ? \leadsto \tau} \qquad icast(v, ty(v)) = v$$

$$\boxed{e : \tau \approx\!\!\!\gg \tau'} \qquad icast(\text{'}v\text{'}, ty(v)) = v \quad v \neq s$$

$$\frac{e \neq v}{e : \tau \approx\!\!\!\gg \tau'} \qquad \boxed{cast(v, \tau) = v'}$$

$$cast(n_1.n_2, \mathbb{Z}) = n_1$$

$$\frac{cast(v, \tau') = v'}{v : \tau \approx\!\!\!\gg \tau'} \qquad \begin{array}{c} cast(v, \text{String}) = \text{'}v\text{'} \quad v \neq s \\ cast(v, \tau) = icast(v, \tau) \end{array}$$

$$\boxed{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau}$$

$$\{m\} = \arg\min_i ((\tau 1 - \tau_{1i}) + (\tau_2 - \tau_{2i}))$$

$$\frac{\tau_1 \Rightarrow \tau_{1m} \qquad \tau_2 \Rightarrow \tau_{2m}}{bestCandidate(\tau_1, \tau_2, \overline{\tau_{1i} \times \tau_{2i} \to \tau_{3i}}) = \tau_{3m} \times \tau_{4m} \to \tau_{3m}}$$

$$\tau - \tau = 0 \quad \mathbb{Z} - \mathbb{R} = 1 \quad \text{String} - \mathbb{Z} = 1 \quad \text{String} - \mathbb{R} = 1 \quad \_ - \_ = 2$$

$$\boxed{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau}$$

$$\frac{\tau_1 \neq ? \qquad\qquad \tau_2 \neq ?}{\quad bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6 \quad}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

$$\frac{\tau_1 \neq ? \qquad\qquad\qquad \tau_2 = ?}{bestCandidate(\tau_1, \tau_1, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6 \quad icast(e_2, \tau_5) = v_2}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

$$\frac{\tau_1 = ? \qquad\qquad\qquad \tau_2 \neq ?}{bestCandidate(\tau_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6 \quad icast(e_1, \tau_4) = v_1}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

$$\frac{\tau_1 = ? \qquad\qquad \tau_2 = ?}{\quad bestCandidate(\text{String}, \text{String}, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6 \quad}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

**Basic definitions**

$$ty(n) = \mathbb{Z} \quad ty(r) = \mathbb{R} \quad ty(s) = ?$$

$$insert(e, \tau, 0) = e :: \tau \qquad clean(T) = T[\text{String}/?]$$

$$ty(0_C) = \{\mathbb{Z} \times \mathbb{Z} \to \mathbb{B}, \mathbb{R} \times \mathbb{R} \to \mathbb{B}, \text{String} \times \text{String} \to \mathbb{B}\}$$

$$ty(+) = \{\mathbb{Z} \times \mathbb{Z} \to \mathbb{Z}, \mathbb{R} \times \mathbb{R} \to \mathbb{R}\}$$

$$apply(0, v_1, v_2) = v_1\ 0\ v_2$$

**Figure 7: The** RAFT/PSQL **Instantiation**

the decimals, and from non-string values to strings by enclosing them in quotes. For instance, `SELECT CAST('1.1' AS DOUBLE) FROM` R is accepted and evaluates to 1.1, but `SELECT CAST('1.1' AS INT)` `FROM` R is rejected by the typechecker.

*Overload resolution.* This rule is divided in four cases. The general case arises when the types of the two expressions are not unknown. Function $bestCandidate(\tau_1, \tau_2, \overline{\tau_3})$ is used to determine the best candidate. We model this function by initially calculating the sum of the type differences between the corresponding types of the expression and the domain of each candidate. The type difference $\tau - \tau'$ quantifies the "cost" of changing type $\tau$ to $\tau'$: it yields a value of 0 if both types are identical, 1 when transitioning from an integer to a real or from a string to a number, and 2 in all other cases. If there are ties, and more than one type is selected, then the function is not defined. Furthermore, both types must satisfy the implicit cast relation with their corresponding type from the domain of the chosen best candidate. For instance, for query `SELECT` 1 + 1.1 `FROM` R, the + operation has two possible candidates: $\mathbb{Z} \times \mathbb{Z} \to \mathbb{Z}$, and $\mathbb{R} \times \mathbb{R} \to \mathbb{R}$. The best candidate in this case is $\mathbb{R} \times \mathbb{R} \to \mathbb{R}$, as both operands can be implicitly cast to reals.

The remaining cases are analogous. When only one type is unknown, the best candidate function is applied using the other known type in both positions. Additionally, the expression of unknown type is implicitly cast to the chosen type to rule out potential errors. Finally, if both types are unknown, they are assumed to be strings when looking for the best candidate. For example, in `SELECT '1'` + 1.1 `FROM` R, the best candidate type is $\mathbb{R} \times \mathbb{R} \to \mathbb{R}$, as the implicit cast from `'1'` to real is possible. On the contrary, query `SELECT` `'1'` + `'1'` `FROM` R is rejected by the typechecker because there is more than one candidate available (int and real versions).

*Basic definitions.* For the type of values, reals are typed as $\mathbb{R}$, integers to $\mathbb{Z}$, and strings literals to ?. The operator for removing unknown types replaces ? occurrences with String, while leaving other types unchanged.[6] The operator for inserting annotations adds explicit casts regardless of the operation's type. The candidate types of a given operator is represented as sets of binary function types. Lastly, the semantics of operations between values are passed to the real implementation without any modifications.

## 5.2 RAFT/SQLite

Figure 8 describes the RAFT/SQLite instantiation. SQLite is one of the most flexible SQL engines: type enforcement is not mandatory, and types on columns are optional. However, it attempts its best effort to perform casts without ever raising a type error at runtime. To capture this kind of flexibility, we model the type of values as ? and define all abstract operators as total functions.

*Bidirectional implicit cast.* Operator *biconv* yields always the unknown type for any pair of types and expressions.

Implicit value casts closely resemble explicit value casts in PSQL. The only distinction lies in the fact that the relation is always defined. Consequently, for implicit casts from strings to numbers, when the string is not a valid number, the result is the same string.

*Explicit cast feasibility.* It is always feasible to cast any expression of type $\tau$ to any type $\tau'$.

[6] Notation $C[A/B]$ denotes type $C$ where occurrences of $B$ have been replaced by $A$.

$$\boxed{icast(v, \tau) = v'} \qquad \dots$$
$$icast(s, \tau) = icast(number(s), \tau) \quad \tau \in \{\mathbb{Z}, \mathbb{R}\}$$
$$icast(s, \tau) = s \quad number(s) \text{ is not def.}$$

$$\boxed{cast(v, \tau) = v'} \qquad \dots$$
$$cast(s, \tau) = cast(number(nprefix(s)), \tau) \quad \tau \in \{\mathbb{Z}, \mathbb{R}\}$$
$$cast(s, \tau) = 0 \quad number(nprefix(s)) \text{ is not defined}, \tau \in \{\mathbb{Z}, \mathbb{R}\}$$

$$\tau - \tau = 0 \quad \mathbb{Z} - \mathbb{R} = 0 \quad \text{String} - \mathbb{R} = 1 \quad ? - \tau = 1 \quad \_ - \_ = 2$$

**Basic definitions**

$$ty(v) = ? \quad dty(w :: ?) = ? \quad dty(r) = \mathbb{R} \quad dty(s) = \text{String}$$
$$insert(e, \tau, 0_C) = e \quad insert(e, \tau, +) = e :: \tau$$
$$ty(0_C) = \{? \times ? \to ?, \mathbb{R} \times \mathbb{R} \to ?, \text{String} \times \text{String} \to ?\}$$
$$ty(+) = \{\mathbb{R} \times \mathbb{R} \to ?\}$$
$$apply(+, v_1, v_2) = v_1 + v_2$$
$$apply(0_C, v_1, v_2) = compare(|v_1'|, |v_2'|) :: ?$$
$$\text{where } resolve(v_1, dty(v_1), v_2, dty(v_2), ty(0_C)) = \tau_1 \times \tau_2 \to \tau_3,$$
$$icast(v_1, \tau_1) = v_1', icast(v_2, \tau_2) = v_2'$$
$$|w :: ?| = w \quad |w| = w$$

**Figure 8: The** RAFT/SQLite **Instantiation (excerpt)**

*Explicit cast.* This operator is defined almost identical to implicit casts, except for the case when the expression is a string. In this case, the cast is performed by extracting the largest numeric prefix from the string and then casting it to the required number type. If there is no numeric prefix, then the cast yields 0. For instance, `SELECT` `CAST('12.3hi' AS INT)`, `CAST('hi') FROM` R evaluates to $(12, 0)$.

*Overload resolution.* There is only one rule for overloading resolution. This rule is similar to PSQL but it yields the first best candidate found. The best candidate function does not perform implicit casts between types. Type difference now yields 0 when going from integer to real, and 2 from unknown to any other type.

*Basic definitions.* Statically, the type of every value is unknown. We introduce a dynamic type operator $dty(v)$ to obtain precise type information during the evaluation of comparisons. The type of a value, initially unknown, may be refined on runtime to a more precise type. In fact 1 :: ? does not behave as 1. Therefore the operator for removing unknown types must act as the identity function. Regarding the operator for inserting annotations, it only inserts explicit casts for arithmetic operations, while for comparison operations, it behaves as the identity function.

The dynamic semantics for comparison is involved. First, the best candidate type is searched, using the dynamic type information of the operands. Then, once the best candidate is found, both operands are implicitly cast to the corresponding types. Finally, the cast values, stripped of (potentially) casts to unknown, are compared using the $compare(w_1, w_2)$ function defined as 1 if $(dty(w_1) = dty(w_2) \wedge w_1 < w_2) \vee dty(w_1) < dty(w_2)$; 0 otherwise. If the dynamic types of the operands are equal, then a regular comparison operation is performed. If the types are different then the types are compared using an arbitrary hierarchy such that $\mathbb{Z} = \mathbb{R} < \text{String}$.

# 6 PROPERTIES

Based on the core relational algebra of RAFT, we can establish metatheoretical results for each engine considered, consisting of multiple lemmas and theorems regarding queries and their evaluation. Specifically, we can articulate the formal distinctions among various engines. This aids us in better comprehending the process of transforming queries from one engine to another.

To state these theorems, we need several new definitions, in particular a way to typecheck rows $r$, tables $t$ and databases $\mathsf{D}$:

$$\frac{\vdash v : \tau}{\vdash v : (N \mapsto \tau)} \qquad \frac{\vdash r : T \qquad \vdash v : \tau}{\vdash r, v : T, (N \mapsto \tau)} \qquad \frac{}{\vdash \cdot : T}$$

$$\frac{\forall r \in t. \vdash r : T}{\vdash t : clean(T)} \qquad \frac{\forall R \in dom(\mathsf{D}). \vdash \mathsf{D}(R) : \Gamma(R)}{\vdash \mathsf{D} : \Gamma}$$

A row is well-typed if every value is typed to its corresponding type in $T$. A table is typed $T$ if every row is typed as $T$. An empty row is typed to any relation type $T$. A database is typed $\Gamma$, if every relation in $\mathsf{D}$ is typed to its corresponding type in $\Gamma$. The proofs require some properties about the implementation of abstract operators:

R1: Every operator must be deterministic.
R2: $biconv(e_1, \tau_1, e_2, \tau_2)$ yields either $\tau_1$ or $\tau_2$.
R3: $resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3})$ must be contained in $\overline{\tau_3}$.
R4: If $cast(v, \tau) = v'$, then the (cleaned) type of $v'$ must be $\tau$.
R5: If $resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6$, then
$[\![e_1 :: \tau_4]\!]_\eta \neq \textbf{error}$ and $[\![e_2 :: \tau_5]\!]_\eta \neq \textbf{error}$.
R6: If $biconv(e_1, \tau_1, e_2, \tau_2) = \tau$, then
$[\![e_1 :: \tau]\!]_\eta \neq \textbf{error}$ and $[\![e_2 :: \tau]\!]_\eta \neq \textbf{error}$.
R7: If $resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6$ then $e_1 : \tau_1 \approxeq \tau_4$ and $e_2 : \tau_2 \approxeq \tau_5$.
R8: If $biconv(e_1, \tau_1, e_2, \tau_2) = \tau$ then $e_1 : \tau_1 \approxeq \tau$ and $e_2 : \tau_2 \approxeq \tau$.
R9: Either $clean(T) = T[\text{String}/?]$ or $clean(T) = T$.

Assuming R1, R2, R3, and R4, any instantiation of RAFT is *type safe*, meaning that well-typed queries either reduce to a table or raise a (controlled) type error. In other words, evaluation of well-typed queries does not get stuck. For instance, an ill-type query, or `SELECT` Foo `FROM` P, where Foo is not defined in P, gets stuck as $[\![\text{Foo}]\!]_\eta$ and does not evaluate.

THEOREM 6.1 (TYPE SAFETY). $\forall T\ Q\ \mathsf{D}, if\ \Gamma \vdash Q : T\ and\ \vdash \mathsf{D} : \Gamma\ then\ (\exists t, [\![Q]\!]_\mathsf{D} = t) \vee ([\![Q]\!]_\mathsf{D} = \textbf{error}).$

Assuming also R4, a stronger theorem called *type soundness* states that, in addition to type safety, the resulting table indeed has the type of the query.

THEOREM 6.2 (TYPE SOUNDNESS). $\forall T\ Q\ \mathsf{D}, if\ \Gamma \vdash Q : T\ and\ \vdash \mathsf{D} : \Gamma\ then\ (\exists t, T'.[\![Q]\!]_\mathsf{D} = t, \vdash t : T'\ and\ clean(T') = clean(T)) \vee ([\![Q]\!]_\mathsf{D} = \textbf{error}).$

The use of $clean(\cdot)$ is exclusively for PSQL, and for cases such as `SELECT CAST`('hi' `as` String) `as` A `from` R. This query of type $A \mapsto$ String evaluates to $\{$'hi', ...$\}$, typed as $A \mapsto ?$ (literal strings are typed ?), but $clean(\text{String}) = clean(?) = \text{String}$.

Type safety is satisfied by every engine we consider, but type soundness is satisfied by all except SQLite. This is because SQLite does not satisfy R3. For instance, in SQLite, `SELECT CAST`(1 `as` `INT`) `AS` A `FROM` R has type $A \mapsto \mathbb{Z}$, but its evaluation is typed $A \mapsto ?$. Also, SQLite permits storing string values in integer columns.

Moreover, PSQL, MySQL and SQLite satisfy a theorem that states that if the programmer does not use any explicit cast in a well-typed query, then the query evaluates without error. To state this theorem we use the cast-free metafunction $\mathsf{CF}(Q)$, which is defined when $Q$ does not have explicit casts of the form $e :: \tau$ (definition in supplementary material).

THEOREM 6.3 (CAST-FREE QUERIES ARE TOTAL). *If* $\mathsf{CF}(Q), \vdash \mathsf{D} : \Gamma$ *and* $\Gamma \vdash Q : T \rightsquigarrow Q'$ *then* $[\![Q']\!]_\mathsf{D} \neq \textbf{error}$.

This theorem requires R1, R5, R6 and R9. Both MSSQL and Oracle do not satisfy this property. Specifically, MSSQL does not satisfy R5 and R6, and Oracle does not satisfy R5. To illustrate why, let us consider table $R = \{($'1'$)\}$, schema $R \mapsto (A \mapsto \text{String})$, and query `SELECT` A + 1 `FROM` R. This query does not typecheck in PSQL, and evaluates successfully in other engines. But with one more row to $R$: $\{($'1'$), ($'hi'$)\}$, the same query evaluates to a runtime error in MSSQL and Oracle.

Regarding cast insertion, translated queries preserve types, and the translation is unique:

THEOREM 6.4 (CAST INSERTION IS A TYPE-PRESERVING FUNCTION). *If* $\Gamma \vdash Q : T$ *then there exists a unique* $Q'$ *such that* $\Gamma \vdash Q : T \rightsquigarrow Q'$ *and* $\Gamma \vdash Q' : T$.

This theorem requires R1, R2, R4, R7, R8 and R9, and is satisfied by every engine we considered.

# 7 VALIDATION

To validate the adequacy of the formalism, we develop PyRAFT, an implementation of RAFT in Python, and create five instances for each engine. We generate multiple random queries and verify if the results from the actual engine match those obtained from the prototype.

PyRAFT is composed of four main modules: a generator for creating random queries, a parser that transforms SQL queries into relational algebra queries, a type checker, and an evaluator responsible for evaluating relational algebra queries. Contrary to the formalism presented in previous sections, the evaluator represents tables as unordered lists instead of sets, closely resembling the bag semantics found in the supplementary material.

For testing purposes, we constructed three tables, each comprised of three columns containing strings, integers, and real numbers. Keeping the tables small accelerates the validation process, as our focus lies not on performance but on behavioral differences to validate the RAFT formalism. Additionally, we generated random expressions consisting of arithmetic and comparison operations.

We generate a total of $100,000$ random SQL queries for each engine, successfully confirming that our design aligns with the behavior of each individual engine. This process is challenging when dealing with engine-specific query optimizations. In particular, sometimes PyRAFT reports an error while the engine returns a table. This discrepancy occurs due to avoidance of executing certain subexpressions or subqueries that are prone to failure. For this reason, we divided the validation in two categories: a *termination-insensitive validation*, and a *termination-sensitive validation*.

The termination-insensitive validation approach involves verifying that, if the evaluation of a query in PyRAFT and in the engine result in tables, then these tables must be equivalent. In PyRAFT,

the query generation is parameterized by the engine due to subtle discrepancies between engines. For instance, MSSQL and Oracle lack a Boolean type and represent Booleans using integers. To avoid floating number precision mismatches, real numbers are represented as decimals in both PyRAFT and the engines.[7] Note that comparing real numbers using a notion of closeness might be feasible, but it presents a greater challenge when these results are then cast to strings. Finally, in MySQL, we had to cast some operands of arithmetic operators to decimal to avoid precision issues. For instance, query `select 'hi' + 5.1 - 9 FROM R` evaluates in PyRAFT to $-3.9$, but to $\{-3.9000000000000004\}$ in MySQL.

The termination sensitive validation approach is a stronger result. It involves verifying that, if the evaluation of a query in PyRAFT yields a table, then evaluation of the query in the engine results in an equivalent table. Furthermore, if a query in PyRAFT reports an error, then the query in the engine also reports an error. It's important to note that sometimes distinguishing between errors resulting from type checking or evaluation solely by inspecting the engine's output might not be feasible. Consequently, if PyRAFT reports a type error (either statically or dynamically), we verify if the real engine throws any kind of errors.

To achieve this stronger validation, we had to perform several simplifications on the generation of queries. This was necessary because the engine's query optimizations prevent some sub-expressions from being evaluated. Among the troublesome optimization we find: (1) expressions within subqueries are not evaluated if the selection columns do not depend on the subquery's columns; (2) unpredictable order of evaluation of logical conjunctions and short-circuit implementation; (3) the engine predicts that intersections are empty by examining only the first column of both tables. A detailed list of simplifications is described in the supplementary material.

## 8 RELATED WORK

Traditionally SQL has been implemented with some sort of either static typing or syntactic checking, though the issue of type errors and type disciplines has received little attention. Nonetheless, there are two lines of works that relates to our problem. In the Database literature, the consideration of corner cases such as NULLs and dynamic generated queries involved typing issues. Also, some engines, like SQLLite have "flexible" type systems. In Programming Languages type system is a central topic, although SQL and databases have received little attention. Both areas have been source of inspiration and techniques for our work.

*Classical Database literature.* There are many works formalizing SQL [9–11, 28]. Guagliardo and Libkin [20] developed a comprehensive formal semantics for SQL whose core we follow here. Following the classic framework in the area, they assume that all comparisons and operations are applied to arguments of the right types. Therefore essentially they do not deal with typing issues. Regarding errors in SQL, based on previous work [4, 33, 39], Taipalus et al. [35] review SQL errors to build a unified error categorization. In further work, Taipalus et al. [34] compared the error messages of the four most popular relational database management systems (MySQL, Oracle, PostgreSQL, and SQL Server) in terms of error message

effectiveness, effects, and usefulness, and error recovery confidence. Our work does not deal with error messages, but instead with detecting errors. Finally, regarding formalization, Benzaken et al. [6] proposed the first mechanically verified compiler (DBCert, using Coq) for SQL queries. Ricciotti and Cheney [31] complemented and deepened the work of Guagliardo and Libkin [20] by making the notions of their semantics and proof precise and formal using Coq. These works assume that every type is matched and therefore there are no implicit casts.

*Flexible typing databases.* A distinctive example is SQLite, the most widely deployed database engine [22], that enjoys flexible typing. Data of any type may be stored in any column of an SQLite table (except an INTEGER PRIMARY KEY column, in which case the data must be integral) and columns can be declared without a data type [13]. SQL queries are viewed as strings and little error checking is done for dynamically generated SQL query strings. Wassermann et al. [38] proposed a static program analysis technique to verify that dynamically generated query strings do not contain type errors. Similar to RAFT, they employ a type system to reject dynamically generated queries. However, their focus does not lie in the formalism of dynamic semantics or casts, and the evaluation uses the grammar of Oracle.

*Strongly-typed queries.* The development of programming language libraries and tools for type-checking queries has been extensively explored [5, 14–17, 25–27, 32]. However, the formalization of implicit and explicit type casts has not been addressed. Additionally, these studies lack a practical exploration of the varied behaviors induced by typing in industry-standard database engines. From a formal perspective, significant progress has been made in the area of type inference for relational algebra [8, 29, 36] and SQL [12, 24]. In RAFT, we presume the existence of a typed schema, and consider the implementation of type inference as an area for future exploration.

## 9 CONCLUSION

In this paper, we identified some discrepancies in behavior regarding the handling of types both statically and dynamically in current SQL engines. This makes difficult to port queries from one engine to another. As a first step to address this problem, we presented a formal framework for a typed relational algebra, named RAFT, with support for implicit and explicit type casts. Through RAFT, we can formally understand the behavior of different database engines; we provided five instantiations and described those of PostgreSQL and SQLite. Our framework, highlights the necessary requirement for any concrete instantiation to satisfy formal properties such as type safety and soundness, among others. We validated our formalization via prototype implementations, comparing thousands of randomly generated queries. The immediate future work we foresee is to study discrepancies under query optimizations performed by many practical engines, and to improve error messages.

## REFERENCES

[1] 2020. Converting MySQL to PostgreSQL. https://en.wikibooks.org/wiki/Converting_MySQL_to_PostgreSQL.
[2] 2023. Converting MySQL to Postgre online. https://en.wikibooks.org/wiki/Converting_MySQL_to_PostgreSQL.
[3] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of Databases.* Addison-Wesley. http://webdam.inria.fr/Alice/

---

[7]In Python, $1.9 - 0.1 = 1.7999999999999998$

[4] Alireza Ahadi, Vahid Behbood, Arto Vihavainen, Julia Prior, and Raymond Lister. 2016. Students' Syntactic Mistakes in Writing Seven Different Types of SQL Queries and Its Application to Predicting Students' Success. In *Proceedings of the 47th ACM Technical Symposium on Computing Science Education* (Memphis, Tennessee, USA) *(SIGCSE '16)*. Association for Computing Machinery, New York, NY, USA, 401–406. https://doi.org/10.1145/2839509.2844640

[5] Lennart Augustsson and Mårten Ågren. 2016. Experience Report: Types for a Relational Algebra Library. *SIGPLAN Not.* 51, 12 (sep 2016), 127–132. https://doi.org/10.1145/3241625.2976016

[6] Véronique Benzaken, Évelyne Contejean, Mohammed Houssem Hachmaoui, Chantal Keller, Louis Mandel, Avraham Shinnar, and Jérôme Siméon. 2022. Translating Canonical SQL to Imperative Code in Coq. *Proc. ACM Program. Lang.* 6, OOPSLA1, Article 83 (apr 2022), 27 pages. https://doi.org/10.1145/3527327

[7] HL Bhandari and R Chitrakar. 2020. Comparison of Data Migration Techniques from SQL Database to NoSQL Database. *J Comput Eng Inf Technol 9* 6 (2020), 2.

[8] Peter Buneman and Atsushi Ohori. 1996. Polymorphism and Type Inference in Database Programming. *ACM Trans. Database Syst.* 21, 1 (mar 1996), 30–76. https://doi.org/10.1145/227604.227609

[9] Stefano Ceri and Georg Gottlob. 1985. Translating SQL Into Relational Algebra: Optimization, Semantics, and Equivalence of SQL Queries. *IEEE Transactions on Software Engineering* SE-11 (1985), 324–345. https://api.semanticscholar.org/CorpusID:22717180

[10] Shumo Chu, Chenglong Wang, Konstantin Weitz, and Alvin Cheung. 2017. Cosette: An Automated Prover for SQL. In *Conference on Innovative Data Systems Research.* https://api.semanticscholar.org/CorpusID:12408033

[11] Shumo Chu, Konstantin Weitz, Alvin Cheung, and Dan Suciu. 2016. HoTTSQL: proving query rewrites with univalent SQL semantics. *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation* (2016). https://api.semanticscholar.org/CorpusID:644867

[12] Dario Colazzo and Carlo Sartiani. 2011. Precision and Complexity of XQuery Type Inference. In *Proceedings of the 13th International ACM SIGPLAN Symposium on Principles and Practices of Declarative Programming* (Odense, Denmark) *(PPDP '11)*. Association for Computing Machinery, New York, NY, USA, 89–100. https://doi.org/10.1145/2003476.2003490

[13] Kevin P. Gaffney, Martin Prammer, Laurence C. Brasfield, D. Richard Hipp, Dan R. Kennedy, and Jignesh M. Patel. 2022. SQLite: Past, Present, and Future. *Proc. VLDB Endow.* 15 (2022), 3535–3547. https://api.semanticscholar.org/CorpusID:252066674

[14] Carl Gould, Zhendong Su, and Premkumar T. Devanbu. 2004. JDBC Checker: A Static Analysis Tool for SQL/JDBC Applications. In *26th International Conference on Software Engineering (ICSE 2004), 23-28 May 2004, Edinburgh, United Kingdom*, Anthony Finkelstein, Jacky Estublier, and David S. Rosenblum (Eds.). IEEE Computer Society, 697–698. https://doi.org/10.1109/ICSE.2004.1317494

[15] The Doobie Development Group. 2019. Typechecking Queries. https://tpolecat.github.io/doobie/docs/06-bChecking.html.

[16] The Kysely Development Group. 2021. Kysely. https://kysely.dev/.

[17] The PgTyped Development Group. 2020. PgTyped. https://github.com/adelsz/pgtyped.

[18] The PostgreSQL Global Development Group. 1996. PostgreSQL Documentation. https://www.postgresql.org/docs/current/typeconv-boverview.html.

[19] The SQLines Development Group. 2010. SQLines. http://www.sqlines.com/online.

[20] Paolo Guagliardo and Leonid Libkin. 2017. A Formal Semantics of SQL Queries, Its Validation, and Applications. *Proc. VLDB Endow.* 11, 1 (sep 2017), 27–39. https://doi.org/10.14778/3151113.3151116

[21] Stephanie W. Haas. 2004. Erik Peter Bansleben . Database Migration : A Literature Review and Case Study. https://api.semanticscholar.org/CorpusID:17518212

[22] D. Richard Hipp. [n.d.]. Most Widely Deployed and Used Database Engine. https://www.sqlite.org/mostdeployed.html

[23] Shafaq Khan, Ajish Kalia, Haleh Mehdipour Dastjerdi, and Nishara Nizamuddin. 2023. Automated Tool for NoSQL to SQL Migration. In *Proceedings of the 7th International Conference on Information Systems Engineering* (, Charleston, SC, USA,) *(ICISE '22)*. Association for Computing Machinery, New York, NY, USA, 20–23. https://doi.org/10.1145/3573926.3573931

[24] Weisheng Lin. 2004. *Type inference in SQL.* Ph.D. Dissertation. Concordia University.

[25] Simon Marlow et al. 2010. Haskell 2010 language report. *Available online http://www. haskell. org/(May 2011)* (2010).

[26] MIT. 2019. TS-SQL-Query. https://ts-bsql-bquery.readthedocs.io/.

[27] Claudia Mónica Necco and J Nuno Olivera. 2005. Toward generic data processing. In *XI Congreso Argentino de Ciencias de la Computación.*

[28] M. Negri, G. Pelagatti, and L. Sbattella. 1991. Formal Semantics of SQL Queries. *ACM Trans. Database Syst.* 16, 3 (sep 1991), 513–534. https://doi.org/10.1145/111197.111212

[29] Atsushi Ohori and Peter Buneman. 1988. Type Inference in a Database Programming Language. In *Proceedings of the 1988 ACM Conference on LISP and Functional Programming* (Snowbird, Utah, USA) *(LFP '88)*. Association for Computing Machinery, New York, NY, USA, 174–183. https://doi.org/10.1145/62678.62700

[30] Benjamin C. Pierce. 2002. *Types and Programming Languages* (1st ed.). The MIT Press.

[31] Wilmer Ricciotti and James Cheney. 2022. A Formalization of SQL with Nulls. *J. Autom. Reason.* 66, 4 (nov 2022), 989–1030. https://doi.org/10.1007/s10817-b022-b09632-b4

[32] Alexandra Silva and Joost Visser. 2006. Strong Types for Relational Databases. In *Proceedings of the 2006 ACM SIGPLAN Workshop on Haskell* (Portland, Oregon, USA) *(Haskell '06)*. Association for Computing Machinery, New York, NY, USA, 25–36. https://doi.org/10.1145/1159842.1159846

[33] John B. Smelcer. 1995. User Errors in Database Query Composition. *Int. J. Hum.-Comput. Stud.* 42, 4 (apr 1995), 353–381. https://doi.org/10.1006/ijhc.1995.1017

[34] Toni Taipalus, Hilkka Grahn, and Hadi Ghanbari. 2021. Error messages in relational database management systems: A comparison of effectiveness, usefulness, and user confidence. *Journal of Systems and Software* 181 (2021), 111034. https://doi.org/10.1016/j.jss.2021.111034

[35] Toni Taipalus, Mikko Siponen, and Tero Vartiainen. 2018. Errors and Complications in SQL Query Formulation. *ACM Trans. Comput. Educ.* 18, 3, Article 15 (aug 2018), 29 pages. https://doi.org/10.1145/3231712

[36] Jan Van den Bussche and Emmanuel Waller. 2002. Polymorphic Type Inference for the Relational Algebra. *J. Comput. System Sci.* 64, 3 (2002), 694–718. https://doi.org/10.1006/jcss.2001.1812

[37] Philip Wadler. 1992. Comprehending Monads. *Mathematical Structures in Computer Science* 2 (1992), 461–493.

[38] Gary Wassermann, Carl Gould, Zhendong Su, and Premkumar Devanbu. 2007. Static Checking of Dynamically Generated Queries in Database Applications. *ACM Trans. Softw. Eng. Methodol.* 16, 4 (sep 2007), 14–es. https://doi.org/10.1145/1276933.1276935

[39] Charles Welty. 1985. Correcting User Errors in SQL. *International Journal of Man-Machine Studies* 22, 4 (1985), 463–477.

$\boxed{biconv(e_1, \tau_1, e_2, \tau_2) = \tau}$

$$\frac{}{biconv(e_1, \tau_1, e_2, \tau_2) = ?}$$

$\boxed{icast(v, \tau) = v'}$

$$
\begin{aligned}
number(`n`) &= n & number(`r`) &= r \\
icast(v, ty(v)) &= v & icast(n_1.n_2, \mathbb{Z}) &= n_1 \\
icast(n, \mathbb{R}) &= n & icast(v, \mathrm{String}) &= `v` \quad v \neq s \\
icast(w :: ?, \tau) &= icast(w, \tau) \\
&\quad\ ... \\
icast(s, \tau) &= icast(number(s), \tau) & \tau &\in \{\mathbb{Z}, \mathbb{R}\} \\
icast(s, \tau) &= s & number(s) &\text{ is not def.}
\end{aligned}
$$

$\boxed{e : \tau \approx\!\!\!\gg \tau'}$ $\qquad \dfrac{}{e : \tau \approx\!\!\!\gg \tau'}$

$\boxed{cast(v, \tau) = v'}$ $\qquad$ ...

$$cast(s, \tau) = cast(number(nprefix(s)), \tau) \quad \tau \in \{\mathbb{Z}, \mathbb{R}\}$$
$$cast(s, \tau) = 0 \quad number(nprefix(s)) \text{ is not defined}, \tau \in \{\mathbb{Z}, \mathbb{R}\}$$

$\boxed{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau}$

$$\frac{\{m\} = \arg\min_i ((\tau1 - \tau_{1i}) + (\tau_2 - \tau_{2i}))}{bestCandidate(\tau_1, \tau_2, \overline{\tau_{1i} \times \tau_{2i} \to \tau_{3i}}) = \tau_{3m} \times \tau_{4m} \to \tau_{3m}}$$

$$\tau - \tau = 0 \quad \mathbb{Z} - \mathbb{R} = 0 \quad \mathrm{String} - \mathbb{R} = 1 \quad ? - \tau = 1 \quad \_ - \_ = 2$$

$\boxed{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$

$$\frac{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

**Basic definitions**

$$ty(v) = ? \quad dty(w :: ?) = ? \quad dty(r) = \mathbb{R} \quad dty(s) = \mathrm{String}$$
$$clean(T) = T$$
$$insert(e, \tau, 0_C) = e \quad insert(e, \tau, +) = e :: \tau$$
$$ty(0_C) = \{? \times ? \to ?, \mathbb{R} \times \mathbb{R} \to ?, \mathrm{String} \times \mathrm{String} \to ?\}$$
$$ty(+) = \{\mathbb{R} \times \mathbb{R} \to ?\}$$
$$apply(+, v_1, v_2) = v_1 + v_2$$
$$apply(0_C, v_1, v_2) = compare(|v_1'|, |v_2'|) :: ?$$
$$\text{where } resolve(v_1, dty(v_1), v_2, dty(v_2), ty(0_C)) = \tau_1 \times \tau_2 \to \tau_3,$$
$$icast(v_1, \tau_1) = v_1', icast(v_2, \tau_2) = v_2'$$
$$|w :: ?| = w \quad |w| = w$$

**Figure 9: Sqlite**

# A ALL ENGINES

We show the full SQLite (9) the other engines: MySQL (Figure 10), MSSQL (Figure 11) and Oracle (Figure 12).

# B DYNAMIC SEMANTICS

In figure 13, we show the semantics where tables are multiple sets. The list of column names of a query $Q$ is obtained using function

---

$\boxed{e : \tau \leadsto \tau'}$

$$\frac{}{e : \tau \leadsto \tau'}$$

$\boxed{biconv(e_1, \tau_1, e_2, \tau_2) = \tau}$

$$\frac{\tau_1 \longleftrightarrow \tau_2 : \tau}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau}$$

$\boxed{\tau_1 \longleftrightarrow \tau_2 : \tau}$

$$\frac{}{\tau \longleftrightarrow \tau : \tau} \qquad \frac{}{\mathrm{String} \longleftrightarrow \tau : \mathrm{String}}$$

$$\frac{}{\tau \longleftrightarrow \mathrm{String} : \mathrm{String}} \qquad \frac{\tau \neq \mathrm{String}}{\mathbb{R} \longleftrightarrow \tau : \mathbb{R}} \qquad \frac{\tau \neq \mathrm{String}}{\tau \longleftrightarrow \mathbb{R} : \mathbb{R}}$$

$\boxed{cast(e, \tau) = v}$

$$
\begin{aligned}
cast(v, ty(v)) &= v \\
cast(v, \mathrm{String}) &= `v` \quad v \neq s \\
cast(n_1.n_2, \mathbb{Z}) &= n_1 \\
cast(n, \mathbb{R}) &= n \\
cast(n, \mathbb{R}) &= n \\
cast(s, \tau) &= icast(number(s), \tau) \quad \tau \in \{\mathbb{Z}, \mathbb{R}\} \\
cast(w :: ?, \tau) &= cast(w, \tau)
\end{aligned}
$$

$\boxed{e : \tau \approx\!\!\!\gg \tau'}$

$$\frac{}{e : \tau \approx\!\!\!\gg \tau'}$$

$\boxed{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau}$

$$\frac{\{m\} = \arg\min_i ((\tau1 - \tau_{1i}) + (\tau_2 - \tau_{2i}))}{bestCandidate(\tau_1, \tau_2, \overline{\tau_{1i} \times \tau_{2i} \to \tau_{3i}}) = \tau_{3m} \times \tau_{4m} \to \tau_{3m}}$$

$$\tau - \tau = 0 \quad \mathbb{Z} - \mathbb{R} = 0 \quad \mathrm{String} - \mathbb{R} = 1 \quad \_ - \_ = 2$$

$\boxed{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$

$$\frac{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

**Basic definitions**

$$ty(r) = \mathbb{R} \quad ty(n) = \mathbb{Z} \quad ty(s) = \mathrm{String}$$
$$clean(T) = T$$
$$insert(e, \tau, 0) = e :: \tau \quad 0 \in \{+, <, =\}$$
$$ty(0_C) = \{\mathbb{R} \times \mathbb{R} \to \mathbb{R}, \mathrm{String} \times \mathrm{String} \to \mathbb{R}\}$$
$$ty(+) = \{\mathbb{R} \times \mathbb{R} \to \mathbb{R}\}$$

**Figure 10: Mysql**

$\ell(Q)$, and defined inductively as follows:

$$
\begin{aligned}
\ell(R) &= \text{list of attributes of base relation } R \\
\ell(Q_1 \times Q_2) &= \ell(Q_1) \mathbin{++} \ell(Q_2) \\
\ell(Q_1 \ 0_S \ Q_2) &= \ell(Q_1) \quad \text{where } 0_S \in \{\cup, \cap, -\}; (*) \\
\ell(\sigma_\theta(Q)) &= \ell(Q) \\
\ell(\pi_\beta(Q)) &= \ell(\beta) \\
\ell(\beta, e \text{ as } N) &= \ell(\beta), \ell(e \text{ as } N) \\
\ell(e \text{ as } N) &= N
\end{aligned}
$$

$\boxed{e : \tau \leadsto \tau'}$

$$\frac{\tau \Rightarrow \tau'}{e : \tau \leadsto \tau'} \qquad \frac{}{e : \text{String} \leadsto \tau}$$

$\boxed{icast(e, \tau) = v}$

$$icast(`n`, ty(v)) = n$$
$$icast(`r`, ty(v)) = r$$
$$icast('n', \mathbb{R}) = n$$
$$icast(v, ty(v)) = v$$
$$icast(n, \mathbb{R}) = n$$

$\boxed{biconv(e_1, \tau_1, e_2, \tau_2) = \tau}$

$$\frac{e_1 : \tau_1 \leadsto \tau_2}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau_2} \qquad \frac{e_2 : \tau_2 \leadsto \tau_1}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau_1}$$

$\boxed{e : \tau \approxeq \tau'}$

$$\frac{e \neq v}{e : \tau \approxeq \tau'} \qquad \frac{cast(v, \tau') = v'}{v : \tau \approxeq \tau'}$$

$\boxed{cast(e, \tau) = v}$

$$cast(n_1.n_2, \mathbb{Z}) = n_1$$
$$cast(v, \text{String}) = `v`$$
$$cast(v, t) = icast(v, t)$$

$\boxed{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau}$

$$\frac{\{m\} = \underset{i}{\arg\min}((\tau 1 - \tau_{1i}) + (\tau_2 - \tau_{2i}))}{bestCandidate(\tau_1, \tau_2, \overline{\tau_{1i} \times \tau_{2i} \to \tau_{3i}}) = \tau_{3m} \times \tau_{4m} \to \tau_{3m}}$$

$$\tau - \tau = 0 \quad \mathbb{Z} - \mathbb{R} = 1 \quad \text{String} - \mathbb{Z} = 1 \quad \text{String} - \mathbb{R} = 1 \quad \_ - \_ = 2$$

$\boxed{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$

$$\frac{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6 \quad e1 : \tau_1 \approxeq \tau_4 \quad e2 : \tau_2 \approxeq \tau_5}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

**Basic definitions**

$$ty(r) = \mathbb{R} \quad ty(n) = \mathbb{Z} \quad ty(s) = \text{String}$$
$$clean(T) = T$$
$$insert(e, \tau, 0) = e :: \tau \quad 0 \in \{+, <, =\}$$
$$ty(0_C) = \{\mathbb{Z} \times \mathbb{Z} \to \mathbb{B}, \mathbb{R} \times \mathbb{R} \to \mathbb{B}, \text{String} \times \text{String} \to \mathbb{B}\}$$
$$ty(+) = \{\mathbb{Z} \times \mathbb{Z} \to \mathbb{Z}, \mathbb{R} \times \mathbb{R} \to \mathbb{R}, \text{String} \times \text{String} \to \text{String}\}$$

**Figure 11: Mssql**

where "++" denotes list concatenation, and "," list appending. Rule (*) follows standard engine usage of using the left schema of a set expression as output schema.

$\boxed{e : \tau \leadsto \tau'}$

$$\frac{\tau \Rightarrow \tau'}{e : \tau \leadsto \tau'}$$

$\boxed{biconv(e_1, \tau_1, e_2, \tau_2) = \tau}$

$$\frac{e_1 : \tau_1 \leadsto \tau_2}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau_2} \qquad \frac{e_2 : \tau_2 \leadsto \tau_1}{biconv(e_1, \tau_1, e_2, \tau_2) = \tau_1}$$

$\boxed{cast(e, \tau) = v}$

$$cast('n', \mathbb{Z}) = n$$
$$cast('r', \mathbb{R}) = r$$
$$cast(v, ty(v)) = v$$
$$cast('n', \mathbb{R}) = n$$
$$cast('n_1.n_2', \mathbb{Z}) = n_1$$
$$cast(n_1.n_2, \mathbb{Z}) = n_1$$
$$cast(n, \mathbb{R}) = n$$
$$cast(v, \text{String}) = `v`$$

$\boxed{e : \tau \approxeq \tau'}$

$$\frac{}{e : \tau \approxeq \tau'}$$

$\boxed{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau}$

$$\frac{\{m\} = \underset{i}{\arg\min}((\tau 1 - \tau_{1i}) + (\tau_2 - \tau_{2i}))}{bestCandidate(\tau_1, \tau_2, \overline{\tau_{1i} \times \tau_{2i} \to \tau_{3i}}) = \tau_{3m} \times \tau_{4m} \to \tau_{3m}}$$

$$\tau - \tau = 0 \quad \mathbb{Z} - \mathbb{R} = 0 \quad \text{String} - \mathbb{R} = 1 \quad \_ - \_ = 2$$

$\boxed{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$

$$\frac{bestCandidate(\tau_1, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}{resolve(e_1, \tau_1, e_2, \tau_2, \overline{\tau_3}) = \tau_4 \times \tau_5 \to \tau_6}$$

**Basic definitions**

$$ty(r) = \mathbb{R} \quad ty(n) = \mathbb{Z} \quad ty(s) = \text{String}$$
$$clean(T) = T$$
$$insert(e, \tau, 0) = e :: \tau \quad 0 \in \{+, <, =\}$$
$$ty(0_C) = \{\mathbb{R} \times \mathbb{R} \to \mathbb{B}, \text{String} \times \text{String} \to \mathbb{B}\}$$
$$ty(+) = \{\mathbb{R} \times \mathbb{R} \to \mathbb{R}\}$$

**Figure 12: Oracle**

## C ELABORATION

In figure 14, we show the full elaboration rules.

LEMMA C.1. *If* $T \vdash e : \tau \leadsto e'$ *then* $\Gamma \vdash e' : \tau$.

PROOF. The proof is derived easily from the definition. □

THEOREM C.2 (ELABORATION PRESERVE TYPE). *If* $\Gamma \vdash Q : T \leadsto Q'$ *then* $\Gamma \vdash Q' : T$.

PROOF. The proof follows by Lemma C.3. □

LEMMA C.3. *If* $T \vdash e : \tau$ *then there exists a unique* $e'$ *such that* $T \vdash e : \tau \leadsto e'$.

$$[\![R]\!]_D = \mathbf{ok}\ D\ (R) \tag{RR}$$

$$[\![Q_1\ 0_S\ Q_2]\!]_D = \mathbf{do} \begin{cases} t_1 \leftarrow [\![Q_1]\!]_D \\ t_2 \leftarrow [\![Q_2]\!]_D \\ \mathbf{ok}\ \varepsilon(t_1\ 0_S\ t_2) \end{cases} \tag{R0$_S$}$$

$$[\![\pi_\beta(Q)]\!]_D = \mathbf{do} \begin{cases} t \leftarrow [\![Q]\!]_D \\ \lceil\{\underbrace{[\![\beta]\!]_{\eta^{\bar{r}}_{\ell(Q)}}, \ldots, [\![\beta]\!]_{\eta^{\bar{r}}_{\ell(Q)}}}_{\text{k times}}\ |\bar{r} \in_k t\}\rceil \end{cases} \tag{R$\pi$}$$

$$[\![\sigma_\theta(Q)]\!]_D = \mathbf{do} \begin{cases} t \leftarrow [\![Q]\!]_D \\ \_ \leftarrow \lceil\{[\![\theta]\!]_{\eta^{\bar{r}}_{\ell(Q)}}\ |\bar{r} \in_k t\}\rceil \\ \mathbf{ok}\ \{\underbrace{\bar{r}, \ldots, \bar{r}}_{\text{k times}}\ |\bar{r} \in_k t \wedge [\![\theta]\!]_{\eta^{\bar{r}}_{\ell(Q)}}\} \end{cases} \tag{R$\sigma$}$$

$$[\![e\ \mathbf{as}\ N]\!]_\eta = [\![e]\!]_\eta \tag{Ras}$$

$$[\![\beta, e\ \mathbf{as}\ N]\!]_\eta = \mathbf{do} \begin{cases} r \leftarrow [\![\beta]\!]_\eta \\ v \leftarrow [\![e]\!]_\eta \\ \mathbf{ok}\ r, v \end{cases} \tag{R$\beta$}$$

$$[\![N]\!]_\eta = \mathbf{ok}\ \eta(N) \tag{RN}$$

$$[\![v]\!]_\eta = \mathbf{ok}\ v \tag{Rv}$$

$$[\![v :: \tau]\!]_\eta = \begin{cases} \mathbf{ok}\ v' & cast(v, \tau) = v' \\ \mathbf{error} & \text{otherwise} \end{cases} \tag{Rv ::}$$

$$[\![e :: \tau]\!]_\eta = \mathbf{do} \begin{cases} v \leftarrow [\![e]\!]_\eta \\ [\![v :: \tau]\!]_\eta \end{cases} \tag{Re ::}$$

$$[\![\theta_1\ 0_B\ \theta_2]\!]_\eta = \mathbf{do} \begin{cases} b_1 \leftarrow [\![\theta_1]\!]_\eta \\ b_2 \leftarrow [\![\theta_2]\!]_\eta \\ \mathbf{ok}\ b_1\ 0_B\ b_2 \end{cases} \tag{R0$_B$}$$

$$[\![\neg\theta]\!]_\eta = \mathbf{do} \begin{cases} v_1 \leftarrow [\![\theta_1]\!]_\eta \\ \mathbf{ok}\ \neg b_1 \end{cases} \tag{R$\neg$}$$

$$[\![e_1\ 0\ e_2]\!]_\eta = \mathbf{do} \begin{cases} v_1 \leftarrow [\![e_1]\!]_\eta \\ v_2 \leftarrow [\![e_2]\!]_\eta \\ \mathbf{ok}\ apply(0, v_1, v_2) \end{cases} \tag{R0}$$

$$0 \in \{+, =, <\}$$

**Figure 13: Dynamic Semantics**

PROOF. We proceed by induction on the typing $T \vdash e : \tau$.

*Case* $(N)$. Trivial as there is the only matching translation rule, so we choose $e' = N'$.

*Case* $(v)$. Trivial as there is the only matching translation rule.

*Case* $(e :: \tau)$. From the induction hypotheses, we know that $T \vdash e : \tau \rightsquigarrow e'$ so the result holds.

*Case* $(0_C, 0_A)$. From the induction hypotheses, we know that $T \vdash e_1 : \tau \rightsquigarrow e'_1$ and $T \vdash e_2 : \tau \rightsquigarrow e'_2$ so the result holds.

□

LEMMA C.4. *If $T \vdash \theta : \mathbb{B}$ then there exists a unique $\theta'$ such that $T \vdash \theta : \mathbb{B} \rightsquigarrow \theta'$.*

PROOF. The proof follows by Lemma C.3. □

LEMMA C.5. *If $T \vdash \beta : T'$ then there exists a unique $\beta'$ such that $T \vdash \beta : T' \rightsquigarrow \beta'$.*

PROOF. The proof follows by Lemma C.3. □

THEOREM C.6. *If $\Gamma \vdash Q : T$ then there exists a unique $Q'$ such that $\Gamma \vdash Q : T \rightsquigarrow Q'$.*

PROOF. We proceed by induction on the typing $\Gamma \vdash Q : T$.

*Case* $(R)$. Trivial as there is the only matching translation rule, so we choose $Q' = Q$.

*Case* $(\sigma_\theta(Q))$. For the induction hypotheses, we know that $\Gamma \vdash Q : T \rightsquigarrow Q'$ and by Lemma C.4, we know that $T \vdash \theta : \mathbb{B} \rightsquigarrow \theta'$. Thus the result holds.

*Case* $(\pi_\beta(Q))$. For the induction hypotheses, we know that $\Gamma \vdash Q : T \rightsquigarrow Q'$ and by Lemma C.5, we know that $T \vdash \beta : T' \rightsquigarrow \beta'$. Thus the result holds.

*Case* $(0_S)$. The proof follows by the induction hypotheses and Lemma C.3.

*Case* $(Q_1 \times Q_2)$. For the induction hypotheses, we know that $\Gamma \vdash Q_1 : T_1 \rightsquigarrow Q'_1$ and $\Gamma \vdash Q_2 : T_2 \rightsquigarrow Q'_2$. Thus the result holds.

□

## D TYPE SOUNDNESS

In the proof we use some different notations but they are equivalent.

*Definition D.1 (Well-formed Mapping).*
$$\frac{\forall N \in dom(T). \vdash \eta(N) : \tau \qquad clean(\tau) = clean(T(N))}{T \vdash \eta}$$

*Definition D.2 (Well-formed Database with Schema).*
$$\frac{\forall R \in dom(\Gamma). \cdot \vdash D(R) : T \qquad clean(T) = clean(\Gamma(R))}{\Gamma \vdash D}$$

$$\frac{\vdash v : \tau}{\vdash v : (N \mapsto \tau)} \qquad \frac{\vdash r : T \qquad \vdash v : \tau}{\vdash r, v : T, (N \mapsto \tau)}$$

$$\frac{\forall r \in t. \vdash r : T}{\vdash t : T} \qquad \frac{}{\vdash \cdot : T}$$

LEMMA D.3. *$\forall v\ \tau_1\ \tau_2$, if $\vdash v : \tau_1$ then $(\exists v', cast(v, \tau_2) = v') \vee (cast(v, \tau) = \mathbf{error})$.*

PROOF. The result follows from the definition of conversion of values and the proof is trivial. □

LEMMA D.4. *$\forall v\ \tau_1\ \tau_2$, if $\vdash v : \tau_1$ then $(\exists v', cast(v, \tau_2) = v' \wedge \vdash v' : \tau_3 \wedge clean(\tau_3) = clean(\tau_2)) \vee (cast(v, \tau) = \mathbf{error})$.*

PROOF. The result follows from the definition of conversion of values and the proof is trivial. □

LEMMA D.5. *$\forall T\ e\ \tau\ \eta$, if $T \vdash e : \tau$ and $T \vdash \eta$ then $(\exists v, [\![e]\!]_\eta = v) \vee ([\![e]\!]_\eta = \mathbf{error})$.*

PROOF. The proof follows by induction on the size of $e$.

$$\boxed{T \vdash e : \tau \rightsquigarrow e'}$$

$$(\text{E::}) \frac{\begin{array}{c} T \vdash e : \tau' \rightsquigarrow e' \\ e : \tau' \gtrapprox \tau \end{array}}{T \vdash e :: \tau : \tau \rightsquigarrow e' :: \tau} \qquad (\text{E}v) \frac{T \vdash v : \tau}{T \vdash v : \tau \rightsquigarrow v :: \tau} \qquad (\text{E}N) \frac{T(N) = \tau}{T \vdash N : \tau \rightsquigarrow N}$$

$$(\text{EO}_\text{B}) \frac{T \vdash \theta_1 : \mathbb{B} \rightsquigarrow \theta_1' \quad T \vdash \theta_2 : \mathbb{B} \rightsquigarrow \theta_2' \quad \text{O}_\text{B} \in \{\wedge, \vee\}}{T \vdash \theta_1 \, \text{O}_\text{B} \, \theta_2 : \mathbb{B} \rightsquigarrow \theta_1' \, \text{O} \, \theta_2'} \qquad (\text{E}\neg) \frac{T \vdash \theta : \mathbb{B} \rightsquigarrow \theta'}{T \vdash \neg\theta : \mathbb{B} \rightsquigarrow \neg\theta'}$$

$$(\text{EO}) \frac{\begin{array}{c} T \vdash e_1 : \tau_1 \rightsquigarrow e_1' \qquad T \vdash e_2 : \tau_2 \rightsquigarrow e_2' \\ \text{O} \in \{<, =, +\} \quad resolve(e_1, \tau_1, e_2, \tau_2, ty(\text{O})) = \tau_3 \times \tau_4 \rightarrow \tau_5 \end{array}}{T \vdash e_1 \, \text{O} \, e_2 : \tau_4 \rightsquigarrow insert(insert(e_1', \tau_3, \text{O}) \, \text{O} \, insert(e_2', \tau_4, \text{O}), \tau_5, \text{O})}$$

$$\boxed{T \vdash \beta : T' \rightsquigarrow \beta'}$$

$$(\text{E}\beta) \frac{\overline{T \vdash e : \tau \rightsquigarrow e'} \quad unique(\overline{N})}{T \vdash \overline{e \text{ as } N} : \overline{N \mapsto \tau} \rightsquigarrow \overline{e' \text{ as } N}}$$

$$\boxed{\Gamma \vdash Q : T \rightsquigarrow Q'}$$

$$(\text{E}R) \frac{\Gamma(R) = T}{\Gamma \vdash R : T \rightsquigarrow R} \qquad (\text{E}\pi) \frac{\Gamma \vdash Q : T \rightsquigarrow Q' \quad clean(T) \vdash \beta : T' \rightsquigarrow \beta'}{\Gamma \vdash \pi_\beta(Q) : T' \rightsquigarrow \pi_{\beta'}(Q')} \qquad (\text{E}\sigma) \frac{\Gamma \vdash Q : T \rightsquigarrow Q' \quad T \vdash \theta : \mathbb{B} \rightsquigarrow \theta'}{\Gamma \vdash \sigma_\theta(Q) : T \rightsquigarrow \sigma_{\theta'}(Q')}$$

$$(\text{T}\pi\beta) \frac{\begin{array}{c} \Gamma \vdash \pi_{\beta_1}(Q_1) : T_1 \rightsquigarrow \pi_{\beta_1'}(Q_1') \quad \Gamma \vdash \pi_{\beta_2}(Q_2) : T_2 \rightsquigarrow \pi_{\beta_2'}(Q_2') \\ biconv(\beta_1, T_1, \beta_2, T_2) = T \qquad \qquad \text{O}_\text{S} \in \{\cup, \cap, -\} \end{array}}{\Gamma \vdash \pi_{\beta_1}(Q_1) \, \text{O}_\text{S} \, \pi_{\beta_2}(Q_2) : T \rightsquigarrow \pi_{insert(\beta_1', T, \text{O}_\text{S})}(Q_1') \, \text{O}_\text{S} \, \pi_{insert(\beta_2', T, \text{O}_\text{S})}(Q_2')}$$

$$(\text{E}\times) \frac{\begin{array}{c} \ell(Q_1) \cap \ell(Q_2) = \emptyset \\ \Gamma \vdash Q_1 : T_1 \rightsquigarrow Q_1' \quad \Gamma \vdash Q_2 : T_2 \rightsquigarrow Q_2' \end{array}}{\Gamma \vdash Q_1 \times Q_2 : T_1, T_2 \rightsquigarrow Q_1' \times Q_2'}$$

**Figure 14: Elaboration rules**

*Case* $(N)$.
we know that

(1) $T \vdash \eta$
(2) $\dfrac{T(N) = \tau}{T \vdash N : \tau}$

we need to show $(\exists v, [\![N]\!]_\eta = v) \vee ([\![N]\!]_\eta = \textbf{error})$
from 1, we know that
$$\frac{\forall N' \in dom(T). \vdash \eta(N') : \tau' \quad clean(\tau') = clean(T(N'))}{T \vdash \eta}$$
from 2, we know that $N \in dom(T)$
so
$[\![N]\!]_\eta = \eta(N)$
thus the result holds.

*Case* $(v)$.
$[\![v]\!]_\eta = v$
thus the result holds.

*Case* $(e :: \tau)$.
we know that

(1) $T \vdash \eta$
(2) $\dfrac{\begin{array}{c} T \vdash e : \tau' \\ e : \tau' \gtrapprox \tau \end{array}}{T \vdash e :: \tau : \tau}$

we need to show
$(\exists v, [\![e :: \tau]\!]_\eta = v) \vee ([\![e :: \tau]\!]_\eta = \textbf{error})$
from the hypothesis, we know
$(\exists v, [\![e]\!]_\eta = v) \vee ([\![e]\!]_\eta = \textbf{error})$

from Lemma D.3, we know that
$(\exists v', cast(v, \tau) = v') \vee (cast(v, \tau) = \textbf{error})$
thus the result holds.

*Case* $(e_1 \, \text{O}_\text{A} \, e_2)$.
we know that

(1) $T \vdash \eta$
(2) $\dfrac{T \vdash e_1 : \tau_1 \quad T \vdash e_2 : \tau_2}{T \vdash e_1 \, \text{O}_\text{A} \, e_2 : \tau_3}$

we need to show
$(\exists v, [\![e_1 \, \text{O}_\text{A} \, e_2]\!]_\eta = v) \vee ([\![e_1 \, \text{O}_\text{A} \, e_2]\!]_\eta = \textbf{error})$
from the hypothesis, we know
$(\exists v_1, [\![e_1]\!]_\eta = v_1) \vee ([\![e_1]\!]_\eta = \textbf{error})$
$(\exists v_2, [\![e_2]\!]_\eta = v_2) \vee ([\![e_2]\!]_\eta = \textbf{error})$
after the real implement of $\text{O}_\text{A}$,
$v_1 \, [\![\text{O}_\text{A}]\!] \, v_2 = v$
thus the result holds.

$\square$

**LEMMA D.6.** $\forall T \, e \, \tau \, \eta$, if $T \vdash e : \tau$ and $T \vdash \eta$ then $(\exists v, [\![e]\!]_\eta = v \wedge \vdash v : \tau' \wedge clean(\tau') = clean(\tau)) \vee ([\![e]\!]_\eta = \textbf{error})$.

**PROOF.** The proof follows by induction on the size of $e$.

*Case* $(N)$.
we know that

(1) $T \vdash \eta$
(2) $\dfrac{T(N) = \tau}{T \vdash N : \tau}$

we need to show $(\exists v, [\![N]\!]_\eta = v \wedge \vdash v : \tau' \wedge clean(\tau') = clean(\tau)) \vee ([\![N]\!]_\eta = \mathbf{error})$

from 1, we know that

$$\frac{\forall N' \in dom(T). \vdash \eta(N') : \tau' \quad clean(\tau') = clean(T(N'))}{T \vdash \eta}$$

from 2, we know that $N \in dom(T)$

so

$[\![N]\!]_\eta = \eta(N)$

$\vdash \eta(N) : \tau'$

thus the result holds.

*Case* $(v)$.

$[\![v]\!]_\eta = v$

thus the result holds.

*Case* $(e :: \tau)$.

we know that

(1) $T \vdash \eta$

$\qquad T \vdash e : \tau'$

(2) $\dfrac{e : \tau' \rightsquigarrow \tau}{T \vdash e :: \tau : \tau}$

we need to show

$(\exists v, [\![e :: \tau]\!]_\eta = v \wedge \vdash v : \tau'' \wedge clean(\tau'') = clean(\tau)) \vee ([\![e :: \tau]\!]_\eta = \mathbf{error})$

from the hypothesis, we know

$(\exists v, [\![e]\!]_\eta = v \wedge \vdash v : \tau_1 \wedge clean(\tau_1) = clean(\tau')) \vee ([\![e]\!]_\eta = \mathbf{error})$

from Lemma D.4, we know that

$(\exists v', cast(v, \tau) = v' \wedge \vdash v' : \tau_2 \wedge clean(\tau_2) = clean(\tau)) \vee (cast(v, \tau) = \mathbf{error})$

thus the result holds.

*Case* $(e_1 \, 0_A \, e_2)$.

we know that

(1) $T \vdash \eta$

(2) $\dfrac{T \vdash e_1 : \tau_1 \quad T \vdash e_2 : \tau_2}{T \vdash e_1 \, 0_A \, e_2 : \tau_3}$

we need to show

$(\exists v, [\![e_1 \, 0_A \, e_2]\!]_\eta = v \wedge \vdash v : \tau_3' \wedge clean(\tau_3') = clean(\tau_3)) \vee ([\![e_1 \, 0_A \, e_2]\!]_\eta = \mathbf{error})$

from the hypothesis, we know

$(\exists v_1, [\![e_1]\!]_\eta = v_1 \wedge \vdash v_1 : \tau_1' \wedge clean(\tau_1') = clean(\tau_1)) \vee ([\![e_1]\!]_\eta = \mathbf{error})$

$(\exists v_2, [\![e_2]\!]_\eta = v_2 \wedge \vdash v_2 : \tau_2' \wedge clean(\tau_2') = clean(\tau_2)) \vee ([\![e_2]\!]_\eta = \mathbf{error})$

after the real implement of $0_A$,

$v_1 \, [\![0_A]\!] \, v_2 = v \wedge \vdash v : \tau_3' \wedge clean(\tau_3') = clean(\tau_3)$

thus the result holds.

□

LEMMA D.7. $\forall T \, \theta \, \tau \, \eta$, *if* $T \vdash \theta : \mathbb{B}$ *and* $T \vdash \eta$ *then* $(\exists b, [\![\theta]\!]_\eta = b) \vee ([\![\theta]\!]_\eta = \mathbf{error})$.

PROOF. The proof follows by Lemma D.6. □

LEMMA D.8. $\forall T \, \beta \, T' \, \eta$, *if* $T \vdash \beta : T'$ *and* $T \vdash \eta$ *then* $(\exists r, [\![\beta]\!]_\eta = r) \vee ([\![\beta]\!]_\eta = \mathbf{error})$.

PROOF. We induction on the typing.

*Case* $(\beta = (e \text{ as } N))$. since $[\![\beta]\!]_\eta = [\![e]\!]_\eta$

by Lemma D.6

there are two cases:

if $[\![e]\!]_\eta = \mathbf{error}$ then $[\![\beta]\!]_\eta = \mathbf{error}$, the result holds.

else:

$\exists v', [\![e]\!]_\eta = v$

thus the result holds.

*Case* $(\beta = \beta', (e \text{ as } N))$. we know

$[\![\beta]\!]_\eta = [\![\beta']\!]_\eta, [\![e]\!]_\eta$

from the induction hypothesis,

$(\exists r \, T_1', [\![\beta']\!]_\eta = r) \vee ([\![\beta']\!]_\eta = \mathbf{error})$

if $[\![\beta']\!]_\eta = \mathbf{error}$ then $[\![\beta]\!]_\eta = \mathbf{error}$

the result holds.

else:

by Lemma D.6, we have two cases:

if $[\![e]\!]_\eta = \mathbf{error}$ then $[\![\beta]\!]_\eta = \mathbf{error}$

thus the result holds

else:

$\exists v \, \tau', [\![e]\!]_\eta = v$

so $[\![\beta]\!]_\eta = [\![\beta']\!]_\eta, [\![e]\!]_\eta = r, v$

thus the result holds.

□

LEMMA D.9. $\forall T \, \beta \, T' \, \eta$, *if* $T \vdash \beta : T'$ *and* $T \vdash \eta$ *then* $(\exists r \, T'', [\![\beta]\!]_\eta = r \wedge \vdash r : T'' \wedge clean(T'') = clean(T')) \vee ([\![\beta]\!]_\eta = \mathbf{error})$

PROOF. We induction on the typing.

*Case* $(\beta = (e \text{ as } N))$. since $[\![\beta]\!]_\eta = [\![e]\!]_\eta$

by Lemma D.6 and the typing rule,

$$\frac{T' \vdash e : \tau}{T' \vdash (e \text{ as } N) : (N \mapsto \tau)}$$

so there are two cases:

if $[\![e]\!]_\eta = \mathbf{error}$ then $[\![\beta]\!]_\eta = \mathbf{error}$, the result holds.

else:

$\exists v' \, \tau', [\![e]\!]_\eta = v \wedge \vdash v : \tau' \wedge clean(\tau') = clean(\tau)$

then $\vdash v : N \mapsto \tau'$ and $clean(N \mapsto \tau') = clean(N \mapsto \tau)$

thus the result holds.

*Case* $(\beta = \beta', (e \text{ as } N))$. we know

$[\![\beta]\!]_\eta = [\![\beta']\!]_\eta, [\![e]\!]_\eta$

$$\frac{T' \vdash \beta' : T_1 \quad T' \vdash e \text{ as } N : N \mapsto \tau \quad N \notin dom(T_1)}{T \vdash (\beta', e \text{ as } N) : T_1, N \mapsto \tau}$$

from the induction hypothesis,

$(\exists r \, T_1', [\![\beta']\!]_\eta = r \wedge \vdash r : T_1' \wedge clean(T_1') = clean(T_1)) \vee ([\![\beta']\!]_\eta = \mathbf{error})$

if $[\![\beta']\!]_\eta = \mathbf{error}$ then $[\![\beta]\!]_\eta = \mathbf{error}$

the result holds.

else:

by Lemma D.6, we have two cases:

if $[\![e]\!]_\eta = \mathbf{error}$ then $[\![\beta]\!]_\eta = \mathbf{error}$

thus the result holds

else: $\exists v \, \tau', [\![e]\!]_\eta = v \wedge \vdash v : \tau' \wedge clean(\tau') = clean(\tau)$

so $[\![\beta]\!]_\eta = [\![\beta']\!]_\eta, [\![e]\!]_\eta = r, v \wedge \vdash r, v : (T_1', N \mapsto \tau')$

thus the result holds.

□

LEMMA D.10. *If $e_1 : \tau_1 \rightsquigarrow \tau_2$ and $e_1 : \tau_2 \rightsquigarrow \tau_1$ then $\tau_1 = \tau_2$.*

PROOF. From $\tau_1 \Rightarrow \tau_2$, we know that if $\tau_1 \Rightarrow \tau_2$ and $\tau_2 \Rightarrow \tau_1$ then $\tau_1 = \tau_2$, so the result holds. □

THEOREM D.11 (TYPE SAFETY). *$\forall T\ Q\ D$, if $\Gamma \vdash Q : T$ and $\Gamma \vdash D$ then $(\exists t, [\![Q]\!]_D = t) \vee ([\![Q]\!]_D = \textbf{error})$.*

PROOF. The proof follows by induction on the size of $Q$.

*Case* $(R)$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma(R) = T}{\Gamma \vdash R : T}$

we need to show $\exists t, [\![R]\!]_D = t$
from 1, we know that

$\dfrac{\forall R' \in dom(\Gamma). \vdash D(R') : T'' \quad clean(T'') = clean(\Gamma(R'))}{\Gamma \vdash D}$

from 2, we know that
$\Gamma(R) = T$
so $R \in dom(\Gamma)$
then there is a $D(R)$
thus the result holds.

*Case* $(\sigma_\theta(Q))$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma \vdash Q : T \quad T \vdash \theta : \mathbb{B}}{\Gamma \vdash \sigma_\theta(Q) : T}$

we need to show
$(\exists t\ T', [\![\sigma_\theta(Q)]\!]_D = t) \vee ([\![\sigma_\theta(Q)]\!]_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t', [\![Q]\!]_D = t') \vee ([\![Q]\!]_D = \textbf{error})$
if $[\![Q]\!]_D = \textbf{error}$ then $([\![\sigma_\theta(Q)]\!]_D = \textbf{error})$
thus the result holds.
else $(\exists t', [\![Q]\!]_D = t')$,
$\forall r \in t'$, from Lemma D.7, we know that
$(\exists b, [\![\theta]\!]_{\eta^r_{\ell(Q)}} = b) \vee ([\![\theta]\!]_{\eta^r_{\ell(Q)}} = \textbf{error})$
if $([\![\theta]\!]_{\eta^r_{\ell(Q)}} = \textbf{error})$ then $([\![\sigma_\theta(Q)]\!]_D = \textbf{error})$
thus the result holds
if $[\![\theta]\!]_{\eta^r_{\ell(Q)}} = true$ then $r \in \tau'$ is picked,
else $r \in \tau'$ is not picked
so the resulting table $t$ is subset of $t'$
thus the result holds.

*Case* $(\pi_\beta(Q))$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma \vdash Q : T_1 \quad T_1' = clean(T_1) \quad T' \vdash \beta : T_1''}{\Gamma \vdash \pi_\beta(Q) : T_1''}$

we need to show
$(\exists t, [\![\pi_\beta(Q)]\!]_D = t) \vee ([\![\pi_\beta(Q)]\!]_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t', [\![Q]\!]_D = t') \vee ([\![Q]\!]_D = \textbf{error})$
if $[\![Q]\!]_D = \textbf{error}$ then $[\![\pi_\beta(Q)]\!]_D = \textbf{error}$
thus the result holds.
else $[\![\pi_\beta(Q)]\!]_D = \{[\![\beta]\!]_{\eta^r_{\ell(Q)}} | r \in [\![Q]\!]_D\}$

then by Lemma D.8
thus the result holds.

*Case* $(Q_1\ 0\ Q_2; 0 \in \{\cap, \cup, -\})$.

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma \vdash \pi_{e_1 asN_1}(Q_1)\ 0\ \pi_{e_2 asN_2}(Q_2) : (N_1 \mapsto \tau) \quad \Gamma \vdash \pi_{\beta_1}(Q_1)\ 0\ \pi_{\beta_2}(Q_2) : T}{\Gamma \vdash \pi_{\beta_1, e_1 asN_1}(Q_1)\ 0\ \pi_{\beta_2, e_2 asN_2}(Q_2) : T, N_1 \mapsto \tau}$

we need to show
$(\exists t, [\![\pi_{\beta_1, e_1 asN_1}(Q_1)\ 0\ \pi_{\beta_2, e_2 asN_2}(Q_2)]\!]_D = t)$
$\vee ([\![\pi_{\beta_1, e_1 asN_1}(Q_1)\ 0\ \pi_{\beta_2, e_2 asN_2}(Q_2)]\!]_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t_1, [\![\pi_{e_1 asN_1}(Q_1)\ 0\ \pi_{e_2 asN_2}(Q_2)]\!]_D = t_1)$
$\vee ([\![\pi_{e_1 asN_1}(Q_1)\ 0\ \pi_{e_2 asN_2}(Q_2)]\!]_D = \textbf{error})$
$(\exists t_2, [\![\pi_{\beta_1}(Q_1)\ 0\ \pi_{\beta_2}(Q_2)]\!]_D = t_2) \vee ([\![\pi_{\beta_1}(Q_1)\ 0\ \pi_{\beta_2}(Q_2)]\!]_D = \textbf{error})$
so
$(\exists t_1, t_2, [\![\pi_{\beta_1, e_1 asN_1}(Q_1)\ 0\ \pi_{\beta_2, e_2 asN_2}(Q_2)]\!]_D = t_1, t_2)$
$\vee ([\![\pi_{\beta_1, e_1 asN_1}(Q_1)\ 0\ \pi_{\beta_2, e_2 asN_2}(Q_2)]\!]_D = \textbf{error})$
thus the result holds.

*Case* $(Q_1 \times Q_2)$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\ell(Q_1) \cap \ell(Q_2) = \Phi \quad \Gamma \vdash Q_1 : T_1 \quad \Gamma \vdash Q_2 : T_2}{\Gamma \vdash Q_1 \times Q_2 : T_1, T_2}$

we need to show
$(\exists t, [\![Q_1 \times Q_2]\!]_D = t) \vee ([\![Q_1 \times Q_2]\!]_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t_1, [\![Q_1]\!]_D = t_1) \vee ([\![Q_1]\!]_D = \textbf{error})$
$(\exists t_2, [\![Q_2]\!]_D = t_2) \vee ([\![Q_2]\!]_D = \textbf{error})$
so
$(\exists t_1, t_2, [\![Q_1 \times Q_2]\!]_D = t_1, t_2) \vee ([\![Q_1 \times Q_2]\!]_D = \textbf{error})$
thus the result holds.

□

THEOREM D.12 (TYPE SOUNDNESS). *$\forall T\ Q\ D$, if $\Gamma \vdash Q : T$ and $\Gamma \vdash D$ then $(\exists t, [\![Q]\!]_D = t$ and $\vdash t : T' \wedge clean(T') = clean(T)) \vee ([\![Q]\!]_D = \textbf{error})$.*

PROOF. The proof follows by induction on the size of $Q$.

*Case* $(R)$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma(R) = T}{\Gamma \vdash R : T}$

we need to show $\exists t, [\![R]\!]_D = t$ and $\cdot \vdash t : T' \wedge clean(T') = clean(T)$
from 1, we know that

$\dfrac{\forall R' \in dom(\Gamma). \vdash D(R') : T'' \quad clean(T'') = clean(\Gamma(R'))}{\Gamma \vdash D}$

from 2, we know that
$\Gamma(R) = T$
so $R \in dom(\Gamma)$
then there is a $D(R)$ and $\vdash D(R) : T'$ and $clean(T') = clean(T)$
so $[\![R]\!]_D = D(R)$ and $\vdash D(R) : T'$
thus the result holds.

Case $(\sigma_\theta(Q))$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma \vdash Q : T \qquad T \vdash \theta : \mathbb{B}}{\Gamma \vdash \sigma_\theta(Q) : T}$

we need to show
$(\exists t\ T', \llbracket \sigma_\theta(Q) \rrbracket_D = t$ and $\vdash t : T' \wedge clean(T') = clean(T)) \vee (\llbracket \sigma_\theta(Q) \rrbracket_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t', \llbracket Q \rrbracket_D = t'$ and $\vdash t' : T''$ and $clean(T'') = clean(T)) \vee (\llbracket Q \rrbracket_D = \textbf{error})$
if $\llbracket Q \rrbracket_D = \textbf{error}$ then $(\llbracket \sigma_\theta(Q) \rrbracket_D = \textbf{error})$
thus the result holds.
else $(\exists t', \llbracket Q \rrbracket_D = t'$ and $\vdash t' : T'')$,
$\forall r \in t'$, from Lemma D.7, we know that
$(\exists b, \llbracket \theta \rrbracket_{\eta^r_{\ell(Q)}} = b) \vee (\llbracket \theta \rrbracket_{\eta^r_{\ell(Q)}} = \textbf{error})$
if $(\llbracket \theta \rrbracket_{\eta^r_{\ell(Q)}} = \textbf{error})$ then $(\llbracket \sigma_\theta(Q) \rrbracket_D = \textbf{error})$
thus the result holds
if $\llbracket \theta \rrbracket_{\eta^r_{\ell(Q)}} = true$ then $r \in \tau'$ is picked,
else $r \in \tau'$ is not picked
so the resulting table $t$ is subset of $t'$
and $\vdash t' : T''$
thus the result holds.

Case $(\pi_\beta(Q))$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\Gamma \vdash Q : T_1 \qquad T'_1 = clean(T_1) \qquad T' \vdash \beta : T''_1}{\Gamma \vdash \pi_\beta(Q) : T''_1}$

we need to show
$(\exists t, \llbracket \pi_\beta(Q) \rrbracket_D = t$ and $\vdash t : T''_2$ and $clean(T''_2) = clean(T''_1)) \vee (\llbracket \pi_\beta(Q) \rrbracket_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t', \llbracket Q \rrbracket_D = t'$ and $\vdash t' : T_2 \wedge clean(T_2) = clean(T_1)) \vee (\llbracket Q \rrbracket_D = \textbf{error})$
if $\llbracket Q \rrbracket_D = \textbf{error}$ then $\llbracket \pi_\beta(Q) \rrbracket_D = \textbf{error}$
thus the result holds.
else $\llbracket \pi_\beta(Q) \rrbracket_D = \{\llbracket \beta \rrbracket_{\eta^r_{\ell(Q)}} | r \in \llbracket Q \rrbracket_D\}$
then by Lemma D.9
so $\vdash \{\llbracket \beta \rrbracket_{\eta^r_{\ell(Q)}} | r \in \llbracket Q \rrbracket_D\} : T''_2 \wedge clean(T''_2) = clean(T''_1)$
thus the result holds.

Case $(Q_1\ O\ Q_2; O \in \{\cap, \cup, -\})$.

(1) $\Gamma \vdash D$

(2) $\dfrac{\begin{array}{c}\Gamma \vdash \pi_{e_1\,as\,N_1}(Q_1)\ O\ \pi_{e_2\,as\,N_2}(Q_2) : (N_1 \mapsto \tau) \\ \Gamma \vdash \pi_{\beta_1}(Q_1)\ O\ \pi_{\beta_2}(Q_2) : T\end{array}}{\Gamma \vdash \pi_{\beta_1, e_1\,as\,N_1}(Q_1)\ O\ \pi_{\beta_2, e_2\,as\,N_2}(Q_2) : T, N_1 \mapsto \tau}$

we need to show
$(\exists t, \llbracket \pi_{\beta_1, e_1\,as\,N_1}(Q_1)\ O\ \pi_{\beta_2, e_2\,as\,N_2}(Q_2) \rrbracket_D = t$ and $\vdash t : N_1 \mapsto \tau', T' \wedge clean(N_1 \mapsto \tau', T') = clean((N_1 \mapsto \tau), T)) \vee (\llbracket \pi_{\beta_1, e_1\,as\,N_1}(Q_1)\ O\ \pi_{\beta_2, e_2\,as\,N_2}(Q_2) \rrbracket_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t_1, \llbracket \pi_{e_1\,as\,N_1}(Q_1)\ O\ \pi_{e_2\,as\,N_2}(Q_2) \rrbracket_D = t_1$ and $\vdash t_1 : (N_1 \mapsto \tau') \wedge clean((N_1 \mapsto \tau')) = clean((N_1 \mapsto \tau))) \vee (\llbracket \pi_{e_1\,as\,N_1}(Q_1)\ O\ \pi_{e_2\,as\,N_2}(Q_2) \rrbracket_D =$

$\textbf{error})$
$(\exists t_2, \llbracket \pi_{\beta_1}(Q_1)\ O\ \pi_{\beta_2}(Q_2) \rrbracket_D = t_2$ and $\vdash t_2 : T' \wedge clean(T') = clean(T)) \vee (\llbracket \pi_{\beta_1}(Q_1)\ O\ \pi_{\beta_2}(Q_2) \rrbracket_D = \textbf{error})$
so
$(\exists t_1, t_2, \llbracket \pi_{\beta_1, e_1\,as\,N_1}(Q_1)\ O\ \pi_{\beta_2, e_2\,as\,N_2}(Q_2) \rrbracket_D = t_1, t_2$ and $\vdash t_1, t_2 : (N_1 \mapsto \tau'), T') \vee (\llbracket \pi_{\beta_1, e_1\,as\,N_1}(Q_1)\ O\ \pi_{\beta_2, e_2\,as\,N_2}(Q_2) \rrbracket_D = \textbf{error})$
thus the result holds.

Case $(Q_1 \times Q_2)$.
we know that

(1) $\Gamma \vdash D$

(2) $\dfrac{\ell(Q_1) \cap \ell(Q_2) = \Phi \qquad \Gamma \vdash Q_1 : T_1 \qquad \Gamma \vdash Q_2 : T_2}{\Gamma \vdash Q_1 \times Q_2 : T_1, T_2}$

we need to show
$(\exists t, \llbracket Q_1 \times Q_2 \rrbracket_D = t$ and $\vdash t : T'_1, T'_2 \wedge clean(T_1, T_2) = clean(T'_1, T'_2)) \vee (\llbracket Q_1 \times Q_2 \rrbracket_D = \textbf{error})$
from the hypothesis, we know that
$(\exists t_1, \llbracket Q_1 \rrbracket_D = t_1$ and $\vdash t_1 : \wedge clean(T'_1) = clean(T_1)) \vee (\llbracket Q_1 \rrbracket_D = \textbf{error})$
$(\exists t_2, \llbracket Q_2 \rrbracket_D = t_2$ and $\vdash t_2 : T'_2 \wedge clean(T'_2) = clean(T_2)) \vee (\llbracket Q_2 \rrbracket_D = \textbf{error})$
so
$(\exists t_1, t_2, \llbracket Q_1 \times Q_2 \rrbracket_D = t_1, t_2$ and $\vdash t_1, t_2 : T'_1, T'_2) \vee (\llbracket Q_1 \times Q_2 \rrbracket_D = \textbf{error})$
thus the result holds.

$\square$

$\boxed{CF(Q), CF(\beta), CF(e), CF(\theta)}$

$$\frac{}{CF(R)} \qquad \frac{CF(\beta) \quad CF(Q)}{CF(\pi_\beta(Q))} \qquad \frac{CF(\theta) \quad CF(Q)}{CF(\sigma_\theta(Q))}$$

$$\frac{CF(Q_1) \quad CF(Q_2)}{CF(Q_1\ O_S\ Q_2)} \qquad \frac{O \in \{+, <, =\} \quad CF(e_1) \quad CF(e_2)}{CF(e_1\ O\ e_2)}$$

$$\frac{CF(\theta_1) \quad CF(\theta_2)}{CF(\theta_1\ O_B\ \theta_2)} \qquad \frac{}{CF(N)} \qquad \frac{}{CF(v)}$$

$$\frac{CF(e)}{CF(e\ as\ N)} \qquad \frac{CF(e) \quad CF(\beta)}{CF(\beta, e\ as\ N)}$$

LEMMA D.13. *If $CF(e)$, $T \vdash \eta$ and $T \vdash e : \tau$ then $\llbracket e \rrbracket_\eta \neq \textbf{error}$.*

PROOF. We prove by induction on the typing.

Case $(N)$.
we know that

(1) $T \vdash \eta$

(2) $\dfrac{T(N) = \tau}{T \vdash N : \tau}$

we need to show $\llbracket N \rrbracket_\eta \neq \textbf{error}$
from 1, we know that
$\dfrac{\forall N' \in dom(T).\ \vdash \eta(N') : T(N')}{} = T \vdash \eta$
from 2, we know that $N \in dom(T)$
so
$\llbracket N \rrbracket_\eta = \eta(N)$
thus the result holds.

*Case* ($v$).
$[\![v]\!]_\eta = v$
thus the result holds.

*Case* ($e_1 \; 0_A \; e_2$).
we know that

(1) $T \vdash \eta$

(2) $\dfrac{T \vdash e_1 : \tau_1 \qquad T \vdash e_2 : \tau_2}{T \vdash e_1 \; 0_A \; e_2 : \tau_3}$

we need to show
$[\![e_1 \; 0_A \; e_2]\!]_\eta \neq \mathbf{error}$
from the hypothesis, we know
$[\![e_1]\!]_\eta \neq \mathbf{error}$
$[\![e_2]\!]_\eta \neq \mathbf{error}$
so $[\![e_1 \; 0_A \; e_2]\!]_\eta \neq \mathbf{error}$

□

LEMMA D.14. *If* $\mathsf{CF}(\theta)$, $T \vdash \eta$ *and* $T \vdash \theta : \mathbb{B}$ *then* $[\![\theta]\!]_\eta \neq \mathbf{error}$.

PROOF. The proof follows by Lemma D.13. □

LEMMA D.15. *If* $\mathsf{CF}(\beta)$, $T \vdash \eta$ *and* $T \vdash \beta : T'$ *then* $[\![\beta]\!]_\eta \neq \mathbf{error}$.

PROOF. We prove by induction on the typing.

*Case* ($e$ as $N$).
$[\![e$ as $N]\!]_\eta = [\![e]\!]_\eta$
from Lemma D.13,
so $[\![e]\!]_\eta \neq \mathbf{error}$
thus the result holds.

*Case* ($\beta, e$ as $N$).
$[\![\beta, e$ as $N]\!]_\eta = [\![\beta]\!]_\eta, [\![e]\!]_\eta$
from induction hypothesis,
$[\![\beta]\!]_\eta \neq \mathbf{error}$
from Lemma D.13,
so $[\![e]\!]_\eta \neq \mathbf{error}$
thus the result holds.

□

THEOREM D.16 (FREE OF RUNTIME ERROR FOR TARGET LANGUAGE).
*If* $\mathsf{CF}(Q)$, $\Gamma \vdash \mathsf{D}$ *and* $\Gamma \vdash Q : T$ *then* $[\![Q]\!]_\mathsf{D} \neq \mathbf{error}$.

PROOF. We prove by induction on the typing.

*Case* ($R$).
we know that

(1) $\Gamma \vdash \mathsf{D}$

(2) $\dfrac{\Gamma(R) = T}{\Gamma \vdash R : T}$

we need to show $[\![R]\!]_\mathsf{D} \neq \mathbf{error}$
from 1, we know that
$\dfrac{\forall R' \in dom(\Gamma). \vdash \mathsf{D}(R') : clean(\Gamma(R'))}{\Gamma \vdash \mathsf{D}}$
from 2, we know that
$\Gamma(R) = T$
so $R \in dom(\Gamma)$
then there is a $\mathsf{D}(R)$
and $[\![R]\!]_\mathsf{D} = \mathsf{D}(R)$
thus the result holds.

*Case* ($\sigma_\theta(Q)$).
we know that

(1) $\Gamma \vdash \mathsf{D}$

(2) $\dfrac{\Gamma \vdash Q : T \qquad T \vdash \theta : \mathbb{B}}{\Gamma \vdash \sigma_\theta(Q) : T}$

we need to show
$([\![\sigma_\theta(Q)]\!]_\mathsf{D} \neq \mathbf{error})$
from the hypothesis, we know that
$([\![Q]\!]_\mathsf{D} \neq \mathbf{error})$
By Lemma D.14,
$\mathbf{error} \notin \{\mathbf{error} | r \in [\![Q]\!]_\mathsf{D} \wedge [\![\theta]\!]_{\eta^r_{\ell(Q)}} = \mathbf{error}\}$
then $[\![\sigma_\theta(Q)]\!]_\mathsf{D} = \{r | r \in [\![Q]\!]_\mathsf{D} \wedge [\![\theta]\!]_{\eta^r_{\ell(Q)}}\}$
thus the result holds.

*Case* ($\pi_\beta(Q)$).
we know that

(1) $\Gamma \vdash \mathsf{D}$

(2) $\dfrac{\Gamma \vdash Q : T_1 \qquad T'_1 = clean(T_1) \qquad T' \vdash \beta : T''_1}{\Gamma \vdash \pi_\beta(Q) : T''_1}$

we need to show
$[\![\pi_\beta(Q)]\!]_\mathsf{D} \neq \mathbf{error}$
from the hypothesis, we know that
$[\![Q]\!]_\mathsf{D} \neq \mathbf{error}$
By Lemma D.14,
$\mathbf{error} \notin \{[\![\beta]\!]_{\eta^r_{\ell(Q)}} | r \in [\![Q]\!]_\mathsf{D}\}$
so $[\![\pi_\beta(Q)]\!]_\mathsf{D} = \{[\![\beta]\!]_{\eta^r_{\ell(Q)}} | r \in [\![Q]\!]_\mathsf{D}\}$
thus the result holds.

*Case* ($Q_1 \; 0 \; Q_2; 0 \in \{\cap, \cup, -, \times\}$).
from the induction hypothesis, we know that
$[\![Q_1]\!]_\mathsf{D} \neq \mathbf{error}$
$[\![Q_2]\!]_\mathsf{D} \neq \mathbf{error}$
so $Q_1 \; 0 \; Q_2 = [\![Q_1]\!]_\mathsf{D} \; 0 \; [\![Q_2]\!]_\mathsf{D}$
thus the result holds.

□

LEMMA D.17. *If* $icast(v, \tau) = v'$ *then* $[\![v :: \tau]\!] \neq \mathbf{error}$.

PROOF. From the implicit conversion definition, it is the subset of explicit conversion, thus the result holds. □

LEMMA D.18. *If* $\vdash v : \tau_1$ *and* $\tau_1 \Rightarrow \tau_2$ *then* $[\![v :: \tau_2]\!] \neq \mathbf{error}$.

PROOF. From the definition of $\tau_1 \Rightarrow \tau_2$, the result holds. □

LEMMA D.19. *If* $resolve(e_1, \tau_1, e_2, \tau_2, ty(0)) = \tau_3 \times \tau_4 \to \tau_5$ *then* $[\![e_1 :: \tau_3]\!] \neq \mathbf{error}$ *and* $[\![e_2 :: \tau_4]\!] \neq \mathbf{error}$.

PROOF. for postgres, the proof follows by Lemma D.17 and Lemma D.18. for mysql and sqlite, it holds because casting is defined for all values and types. □

LEMMA D.20. *If* $\mathsf{CF}(e)$, $T \vdash \eta$ *and* $T \vdash e : \tau \rightsquigarrow e'$ *then* $[\![e']\!]_\eta \neq \mathbf{error}$.

PROOF.

*Case* ($N$).
we know that

(1) $T \vdash \eta$

(2) (E$N$) $\dfrac{T(N) = \tau}{T \vdash N : \tau \rightsquigarrow N}$

we need to show $[\![N]\!]_\eta \neq \textbf{error}$

from 1, we know that

$$\dfrac{\forall N' \in dom(T). \vdash \eta(N') : T(N')}{T \vdash \eta}$$

from 2, we know that $N \in dom(T)$

so

$[\![N]\!]_\eta = \eta(N)$

thus the result holds.

*Case* ($v$).

(E$v$) $\dfrac{T \vdash v : \tau}{T \vdash v : \tau \rightsquigarrow v :: \tau}$

the proof follows by Lemma D.18.

*Case* ($e_1 \ 0 \ e_2$).

we know that

(1) $T \vdash \eta$

(2) $\dfrac{T \vdash e_1 : \tau_1 \rightsquigarrow e_1' \qquad T \vdash e_2 : \tau_2 \rightsquigarrow e_2' \qquad 0 \in \{<, =, +\} \quad resolve(e_1, \tau_1, e_2, \tau_2, ty(0)) = \tau_3 \times \tau_4 \rightarrow \tau_5}{T \vdash e_1 \ 0 \ e_2 : \tau_4 \rightsquigarrow insert(insert(e_1', \tau_3, 0) \ 0 \ insert(e_2', \tau_4, 0), \tau_5, 0)}$

we need to show

$[\![insert(insert(e_1', \tau_3, 0) \ 0 \ insert(e_2', \tau_4, 0), \tau_5, 0)]\!]_\eta \neq \textbf{error}$

from the hypothesis, we know

$[\![e_1]\!]_\eta \neq \textbf{error}$

$[\![e_2]\!]_\eta \neq \textbf{error}$

the result follows by Lemma D.19.

$\square$

LEMMA D.21. *If* $\mathsf{CF}(\theta)$, $T \vdash \eta$ *and* $T \vdash \theta : \mathbb{B} \rightsquigarrow \theta'$ *then* $[\![\theta']\!]_\eta \neq$ **error**.

PROOF. The proof follows by Lemma D.20. $\square$

LEMMA D.22. *If* $\mathsf{CF}(\beta)$, $T \vdash \eta$ *and* $T \vdash \beta : T' \rightsquigarrow \beta'$ *then* $[\![\beta']\!]_\eta \neq$ **error**.

PROOF. The proof follows by Lemma D.20. $\square$

THEOREM D.23 (FREE OF RUNTIME ERROR). *If* $\mathsf{CF}(Q)$, $\Gamma \vdash \mathsf{D}$ *and* $\Gamma \vdash Q : T \rightsquigarrow Q'$ *then* $[\![Q']\!]_\mathsf{D} \neq$ **error**.

PROOF. The proof follows by Lemma D.21 and D.22 and the induction hypothesis. $\square$

# E TRANSLATION

We show the full translation from SQL to Relational Algebra in Figure 15.

LEMMA E.1. *Consider SQL expression $E$. Then there exists a unique $e$ such that $E \rightsquigarrow_e e$.*

PROOF. We proceed by induction on the structure of $E$.

*Case* ($E = N$). Trivial as $N \rightsquigarrow_e N$ is the only matching translation rule, so we choose $e = N$.

*Case* ($E = v$). Similar to previous case, $v \rightsquigarrow_e v$ is the only matching translation rule. So we choose $e = v$.

$$E ::= N \mid v \mid E \ 0_\mathsf{A} \ E \mid \mathsf{CAST}(E \ \mathsf{AS} \ \tau) \quad \text{(expressions)}$$
$$C ::= E \ 0_\mathsf{C} \ E \mid C \ 0_\mathsf{B} \ C \qquad\qquad \text{(boolean expressions)}$$
$$G ::= E \mid E \ \mathsf{AS} \ N \qquad\qquad\qquad \text{(aliased column)}$$
$$B ::= G \mid B, G \qquad\qquad\qquad\quad \text{(columns)}$$
$$S ::= \mathsf{SELECT} \ B \ \mathsf{FROM} \ F \ \mathsf{WHERE} \ C \mid$$
$$\quad\quad \mathsf{SELECT} \ B \ \mathsf{FROM} \ F \mid S \ 0_\mathsf{S} \ S \quad \text{(select queries)}$$
$$F ::= R \mid R \ N \mid S \ N \mid F, F \qquad \text{(from queries)}$$

$\boxed{E \rightsquigarrow_e e, C \rightsquigarrow_e \theta}$

(Tr$N$) $\dfrac{}{N \rightsquigarrow_e N}$ 
$\qquad$ (Tr::) $\dfrac{E \rightsquigarrow_e e}{\mathsf{CAST}(E \ \mathsf{AS} \ \tau) \rightsquigarrow_e e :: \tau}$

(Tr$v$) $\dfrac{}{v \rightsquigarrow_e v}$ 
$\qquad$ (Tr$0_\mathsf{A}$) $\dfrac{E_1 \rightsquigarrow_e e_1 \qquad E_2 \rightsquigarrow_e e_2}{E_1 \ 0_\mathsf{A} \ E_2 \rightsquigarrow_e e_1 \ 0_\mathsf{A} \ e_2}$

(Tr$0_\mathsf{C}$) $\dfrac{E_1 \rightsquigarrow_e e_1 \qquad E_2 \rightsquigarrow_e e_2}{E_1 \ 0_\mathsf{C} \ E_2 \rightsquigarrow_e e_1 \ 0_\mathsf{C} \ e_2}$

(Tr$0_\mathsf{B}$) $\dfrac{C_1 \rightsquigarrow_e \theta_1 \qquad C_2 \rightsquigarrow_e \theta_2}{C_1 \ 0_\mathsf{B} \ C_2 \rightsquigarrow_e \theta_1 \ 0_\mathsf{B} \ \theta_2}$

$\boxed{B \rightsquigarrow_b \beta}$

(Tr$E$) $\dfrac{E \ \mathsf{AS} \ alias(E) \rightsquigarrow_b e}{E \rightsquigarrow_b e}$ 
$\qquad$ (Tr$Eas$) $\dfrac{E \rightsquigarrow_e e}{E \ \mathsf{AS} \ N \rightsquigarrow_b e \ \text{as} \ N}$

(Tr$B$) $\dfrac{B \rightsquigarrow_b \beta \qquad G \rightsquigarrow_b e}{B, G \rightsquigarrow_b \beta, e}$

$\boxed{F \rightsquigarrow Q}$

(Tr$R$) $\dfrac{}{R \rightsquigarrow R}$

(Tr$SW$) $\dfrac{B \rightsquigarrow_b \beta \qquad F \rightsquigarrow Q \qquad C \rightsquigarrow_e \theta}{\mathsf{SELECT} \ B \ \mathsf{FROM} \ F \ \mathsf{WHERE} \ C \rightsquigarrow \pi_\beta(\sigma_\theta(Q))}$

(Tr$S$) $\dfrac{B \rightsquigarrow_b \beta \qquad F \rightsquigarrow Q}{\mathsf{SELECT} \ B \ \mathsf{FROM} \ F \rightsquigarrow \pi_\beta(Q)}$

(Tr$\times$) $\dfrac{F_1 \rightsquigarrow Q_1 \qquad F_2 \rightsquigarrow Q_2}{F_1, F_2 \rightsquigarrow Q_1 \times Q_2}$

(Tr$Fas$) $\dfrac{S \rightsquigarrow Q \qquad \ell(Q) = N_1, \dots}{S \ N \rightsquigarrow \pi_{N_1 \text{as} N.N_1, \dots}(Q)}$

**Figure 15: SQL to RA translation**

*Case* ($E = \mathsf{CAST}(E' \ \mathsf{AS} \ \tau)$). By induction hypothesis on $E'$, we know that $E' \rightsquigarrow_e e$ for some unique $e$. Then by (Tr::), we can conclude that $\mathsf{CAST}(E \ \mathsf{AS} \ \tau) \rightsquigarrow_e e :: \tau$, for unique $e :: \tau$.

*Case* ($E = E_1 \ 0_\mathsf{A} \ E_2$). By induction hypotheses on $E_1$ and $E_2$, we know that $E_1 \rightsquigarrow_e e_1$ and $E_2 \rightsquigarrow_e e_2$, for some unique $e_1$ and $e_2$. Then by (Tr$0_\mathsf{A}$), we can conclude that $E \rightsquigarrow_e e$, for unique $e = e_1 \ 0_\mathsf{A} \ e_2$.

$\square$

LEMMA E.2. *Consider SQL expression $C$. Then there exists a unique $\theta$ such that $C \rightsquigarrow_e \theta$.*

PROOF. We proceed by induction on the structure of $C$.

*Case* ($C = E_1 \ 0_C \ E_2$). By Lemma E.1 on $E_1$ and $E_2$, we know that $E_1 \leadsto_e e_1$ and $E_2 \leadsto_e e_2$, for some unique $e_1$ and $e_2$. Then by (Tr0$_C$), we can conclude that $C \leadsto_e \theta$, for unique $\theta = e_1 \ 0_C \ e_2$.

*Case* ($C = C_1 \ 0_B \ C_2$). By Lemma E.1 on $C_1$ and $C_2$, we know that $C_1 \leadsto_e \theta_1$ and $C_2 \leadsto_e \theta_2$, for some unique $\theta_1$ and $\theta_2$. Then by (Tr0$_B$), we can conclude that $C \leadsto_e \theta$, for unique $\theta = \theta_1 \ 0_B \ \theta_2$.

$\square$

LEMMA E.3. *Consider SQL expression B. Then there exists a unique $\beta$ such that $B \leadsto_e \beta$.*

PROOF. We proceed by induction on the structure of $B$.

*Case* ($B = E$). By Lemma E.1 on $E$, we know that $E \leadsto_e e$, for a unique $e$. Suppose a deterministic total function $alias(\cdot)$ that returns a string given an expression. Then if $alias(E) = N$, then by (Tr$Eas$) we can conclude that $E$ **AS** $N \leadsto_e e$ as $N$, for unique $e$ as $N$. Then by (Tr$E$) we conclude that $E \leadsto_e \beta$, for unique $\beta = e$ as $N$.

*Case* ($B = E$ **AS** $N$). By Lemma E.1 on $E$, we know that $E \leadsto_e e$, for a unique $e$. Then by (Tr$Eas$) we conclude that $E \leadsto_e \beta$, for unique $\beta = e$ as $N$.

*Case* ($B = B', G$). By induction hypothesis on $B'$ and $G$ we know that $B' \leadsto_e \beta'$ and $G \leadsto_e e$ for unique $\beta'$ and $e$. Therefore by (Tr$B$), $B \leadsto_e \beta$, for unique $\beta = \beta', e$.

$\square$

LEMMA E.4. *Consider query $Q$, then $\ell(Q) = \overline{N}$, for a unique $\overline{N}$.*

PROOF. Straightforward induction on the structure of $Q$. $\square$

LEMMA E.5. *Consider SQL expression $F$. Then there exists a unique $Q$ such that $F \leadsto_e Q$.*

PROOF. We proceed by induction on the structure of $F$.

*Case* ($F = R$). Trivial as $R \leadsto_e R$.

*Case* ($F =$ **SELECT** $B$ **FROM** $F'$ **WHERE** $C$). By Lemma E.3, $B \leadsto_e \beta$ for unique $\beta$. By induction hypothesis on $F'$, $F' \leadsto_e Q'$ for unique $Q'$. By Lemma E.2, $C \leadsto_e \theta$ for unique $\theta$. We conclude by (Tr$SW$) that $F \leadsto_e Q$, for unique $Q = \pi_\beta(\sigma_\theta(Q'))$.

*Case* ($F =$ **SELECT** $B$ **FROM** $F'$). By Lemma E.3, $B \leadsto_e \beta$ for unique $\beta$. By induction hypothesis on $F'$, $F' \leadsto_e Q'$ for unique $Q'$. We conclude by (Tr$S$) that $F \leadsto_e Q$, for unique $Q = \pi_\beta(Q')$.

*Case* ($F = F_1, F_2$). By induction hypotheses on $F_1$ and $F_2$, $F_1 \leadsto_e Q_1$ and $F_2 \leadsto_e Q_2$ for unique $Q_1$ and $Q_2$. We conclude by (Tr$\times$) that $F \leadsto_e Q$, for unique $Q = Q_1 \times Q_2$.

*Case* ($F = A \ N$). By induction on $A$ we know that $A \leadsto_e Q'$ for a unique $Q'$. Also, by Lemma E.4, $\ell(Q') = N_1, ...$ for unique $N_1, ...$. Therefore by (Tr$Fas$), $F \leadsto_e Q$, for unique $Q = \pi_{N_1 \text{as} N.N_1,...}(Q')$.

$\square$

THEOREM E.6. *Consider an SQL query denoted as $F$, which takes the form* **SELECT** $B$ **FROM** $F$ **WHERE** $C$ *or* **SELECT** $B$ **FROM** $F$. *Then, there exists a unique relational algebra query, denoted as $Q$, such that $F \leadsto_e Q$.*

PROOF. Direct consequence of Lemma E.5. $\square$

# F  VALIDATION

The termination-sensitive validation method ensures a more robust outcome. It involves confirming that when a query is evaluated in language PyRAFT and yields a table, the same query evaluated in the engine produces an identical table. Additionally, if a query in PyRAFT encounters an error, the engine also reports an error.

To achieve this validation we needed to simplify several aspects of the generation of SQL queries, due to query optimizations, which we describe next.

*Avoiding subqueries in FROM position..* Consider the PSQL query **SELECT** S.c **FROM** (**SELECT CAST**(Name **AS INT**) **AS** c **FROM** P) **AS** S. This query results in an error as the values in the column Name cannot be cast to a number. However, a slight change to the query by evaluating **SELECT** 1 **FROM** (**SELECT CAST**(Name **AS INT**) **AS** c **FROM** P) **AS** S allows the query to run without errors. This change in behavior occurs due to an optimization process, where expressions within subqueries are not evaluated if the selection columns do not depend on the subquery's columns. Therefore, we avoid generating subqueries in the FROM position.

*Avoiding conjunctions..* Consider the PSQL query **SELECT** 1 **FROM** P **WHERE** 5 = **CAST**('hello' **AS FLOAT**) **AND** 2 < 1. This query fails at runtime because the cast from a non-numerical string to an integer is not defined. Now, let's examine the query **SELECT** 1 **FROM** P **WHERE** 5 = **CAST**(Name **AS FLOAT**) **AND** 2 < 1. Here, it does not fail; the condition 2<1 is evaluated before the cast from string to integer. Thanks to the *short-circuit* behavior, the entire condition is deemed false, resulting in an empty result. Given the unpredictable nature of the evaluation order of logical conjunctions in various queries across different database engines, we opt to avoid using conjunctions in the query generator.

*Single column selections..* Consider the following MSSQL query

```
(SELECT 1, Age FROM P) INTERSECT
(SELECT 1, Name FROM P)
```

This query fails because string values in column Name cannot be cast to an integer. However, a slightly different query does not fail:

```
(SELECT 1, Age FROM A) INTERSECT
(SELECT 2, Name FROM A)
```

In this scenario, MSSQL predicts that the intersection is empty by examining only the first column of both tables. It avoids equating the remaining columns, thereby avoiding the evaluation of the doomed-to-fail cast. To circumvent this optimization issue, our approach avoids generating of selections involving more than one column.

*Tailoring for* MSSQL. Consider the following MSSQL query

```
(SELECT Age FROM P WHERE 1 < CAST(Name AS INT))
INTERSECT (SELECT Age FROM P)
```

This query results in an error due to the condition in the first subquery. However, the following query:

```
(SELECT Age FROM P WHERE 1 < CAST(Name AS INT))
INTERSECT (SELECT Age+1 FROM P)
```

Yields an empty result because the engine recognizes that the intersection between columns Age and Age+1 is empty before evaluating

the condition. To avoid this kind of optimization issues, we avoid generating queries with conditionals.

## G EXAMPLES

To illustrate the comparison behavior in SQLite, we provide the following examples:

```
SELECT '0' < 1 FROM R;
SELECT '0' < CAST(1 AS INT) FROM R;
SELECT '0hi' < CAST(1 AS INT) FROM R;
```

The first example evaluates to $0$. Since the (E$v$) elaboration rule inserts an explicit cast on every value, both values are cast to ?. Consequently, the chosen candidate is ? $\times$ ? $\rightarrow$ ?, and the implicit cast leaves them untouched. Finally, as the dynamic type of both operands, with annotations removed, is String and $\mathbb{R}$ respectively, the comparison function yields $1$ as result.

The second example evaluates to $1$. In the process of elaboration, the left expression is cast to the unknown type, while the right expression is cast to an integer type. During the actual evaluation, the most suitable candidate is determined to be $\mathbb{R} \times \mathbb{R} \rightarrow$ ?, which

implies an implicit cast of both values into numerical values $0$ and $1$, respectively. As both casted values share the same dynamic type, a standard comparison is carried out, resulting in $1$.

The third example evaluates to $0$. It follows a pattern similar to the second example, but with the implicit cast to a numeric type failing on the left expression, resulting in the same string as the output. Consequently, since the dynamic types of both operands differ (string and real), a type comparison is carried out, leading to the final result of $0$.

To illustrate the insert function, let us consider three examples in SQLite:

- Expression `'0'` $< 1 :: \mathbb{Z}$ is elaborated to `'0'` $:: ? < (1 :: ?) :: \mathbb{Z}$, which evaluates to 1 (because one of the operand's type is not unknown, see Section 5.2).
- Expression `'0hi'` $< 1 :: \mathbb{Z}$ is elaborated to `'0hi'` $:: ? < (1 :: ?) :: \mathbb{Z}$, which evaluates to 0. If an explicit cast to int is inserted in both operands then the result would be 1 (to denote true), as $[\![$ `'0hi'` $:: \mathbb{Z}]\!]_\eta = 0$.
- Expression `'1'` $+1$ is elaborated to (`'1'` $:: ?) :: \mathbb{Z} + (1 :: ?) :: \mathbb{Z}$, which evaluates to 2.